



THE UNIVERSITY *of* EDINBURGH

This thesis has been submitted in fulfilment of the requirements for a postgraduate degree (e.g. PhD, MPhil, DClinPsychol) at the University of Edinburgh. Please note the following terms and conditions of use:

This work is protected by copyright and other intellectual property rights, which are retained by the thesis author, unless otherwise stated.

A copy can be downloaded for personal non-commercial research or study, without prior permission or charge.

This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the author.

The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the author.

When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given.

Every body's gotta eat: why autonomous systems can't live on prediction- error minimization alone

Kathryn Nave

A doctoral thesis in philosophy, submitted to
the University of Edinburgh in 2022

1. I declare that this thesis has been composed solely by myself and that it has not been submitted, in whole or in part, in any previous application for a degree. Except where stated otherwise by reference or acknowledgment, the work presented is entirely my own.
2. I confirm that this thesis presented for the degree of Doctor of Philosophy, has
 - i) been composed entirely by myself
 - ii) been solely the result of my own work
 - iii) not been submitted for any other degree or professional qualification
3. I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

Acknowledgements

Dave Ward: For being endlessly supportive, not (observably) worrying when I went a year without writing anything, introducing me to Hurley and helping me to understand Merleau-Ponty.

Alistair Isaac: For pointing out the problems I know I should change but have been trying to avoid – as well as those I hadn't noticed. I wish I'd finished this earlier so you could have pushed me to fix more of them.

Andy Clark: For giving me the opportunity to be part of the X-SPECT project while giving me the freedom to explore my own interests.

The European Research Council, Royal Institute of Philosophy and Aristotelean Society for funding me to do so.

Kate Webb: For being the best adventure-partner, kitchen-gossiper and philosophy-discusser. I hope teenage-you would approve of the last chapter.

Amy Mallinson: For being my oldest friend, the strongest person I know and showing me it's possible to finish a thesis – in mathematics no less!

Danaja Rutar: For reminding me to stay curious, your unrelenting enthusiasm for exploring how the mind works, and being the ideal of an academic.

Lilith Lee: For showing me that philosophy could be so much more interesting and important than the tiny segment of it I was familiar with.

Matt Sims, Shannon Proksch and the rest of MLEC for demonstrating how philosophy can be anarchic and collaborative, without being any less serious.

The brunch and bouldering crews, particularly Jenny Zhang, for making sure I had something other than philosophy to look forward to each week.

Edinburgh University for being a dream of academic institution and Edinburgh itself for being a place I could love more than any university.

My parents for giving me the freedom and support to do whatever I wanted with my life and no reason to doubt my ability to do so.

Max Wilkinson, my partner of 12 years who believes that acknowledgements should only recognize a direct contribution to the thesis: thank you for allowing me to make fun of you throughout mine. I'm looking forward to having more time to do so in person.

Abstract

Karl Friston's Free Energy Principle has been proposed as a definition of existence from which "everything of interest about life and the universe can be derived" (Friston, 2019, p.176). Despite pretensions to a theory of every 'thing', focus has largely been on the first of these and the attempt to "unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; avoid surprises and you will last longer" (Friston, 2012, p.2).

By analysing biological existence in terms of the stability of a set of essential variables and relations that are taken to define a particular organism, then redescribing this stability amid perturbation in statistical terms as 'surprisal-minimization', the FEP proposes to connect this biological imperative up to an account of higher-level cognition as involving the operation of hierarchical predictive models. At first glance then, it appears the ideal partner to the bioactivist seeking to ground the content of our mental lives in this basic, biological intentionality.

The potential compatibility is misleading. By reducing autonomy to the kind of stability amid perturbation that could be shared by pendulums and people alike, the FEP makes no distinction between a homeostatic mechanism and a self-producing agent, thereby failing to account for why the latter alone should be credited with active engagement in its own continuation. The free energy theorist is left to choose between instrumentalism about intentional attribution, or mathematically-motivated animism.

Neither is an option for the bioactivist who seeks a naturalistic account of intentionality and immanent teleology as real features distinct to living systems alone. This does not, however, mean that the FEP has inadvertently undermined the bioactivist program by revealing there to be no difference in kind between life and non-life. The problem lies instead in its underlying assumption of a machine-substantialist ontology, which erases this distinction from the get-go by presuming the existence of a thing can be reduced to a set of invariant parts and the invariant rules that govern its dynamics.

Organisms are constrained by neither. Unlike machines, we not only persist through but depend upon both material turnover and unprestatable change in our space of possible behaviours. This, as a number of theoretical biologists and complexity scientists have argued, makes it impossible to derive any invariant form or equations that could fix the possible trajectories of a living system in advance. The failure of the FEP as a theory of life is not

just due to the inadequacy of surprisal-minimization as a particular principle, but due to the broader inappropriateness of a machine-substance ontology for analysing biological existence.

The FEF's failure is not the bioenactivist's victory. For this we need to establish not only that living systems are different from machines, but also to show how these differences grant the former alone the status of intentional agents. The prevailing enactivist definition of autonomy as operational closure among precarious processes does not achieve this, for it fails to capture the mutual dependence between structure and dynamics that characterizes the self-determination of organisms. Instead, I draw upon Montévil, Moreno, and Mossio's account of closure among constraints, to show how this describes a system that is intrinsically active in seeking out a continual supply of energy with which to rebuild the network of constraints that both embody that system and which would irretrievably disintegrate without such activity. Only such systems, I argue, can truly be said to have their own existence as the goal of their operations.

This will likely not convince the eliminativist, determined to reduce subjectivity to atoms and the void, nor the anthropocentrist who jealously withholds agency to his species alone. But for the person who feels that there is a vital difference between her strivings and the operations of a machine, I believe that the place to look is not in the rational powers that might set us apart from our humbler biological relations, but in the peculiar form of precariousness that we share.

Table of Contents

<i>Lay summary: Control is not the goal</i>	9
<i>Overview</i>	19
<i>1. Biodynamic Enactivism</i>	25
1.1. The Enactive Approach	25
1.1.1. The enactive approach as naturalistic phenomenology	29
1.1.2. What are your intentions?	36
1.1.3 Hurley's enactive approach	41
1.2. Bioenactivism	46
1.2.0. Autopoiesis and Autonomy	48
1.2.1 Sense-making and adaptivity	53
1.3. From bioenactivism to predictive processing	59
<i>2. Predictive Processing</i>	63
2.1. Minimal Predictive Processing	63
2.2. Reconstructivist Predictive Processing	67
2.2.1. Predictive processing as unconscious inference	71
2.3. Sensorimotor Predictive Processing	76
2.4. What's the point of predictive processing?	80
<i>3. A non-mathematical introduction to the free energy framework</i>	85
3.1. Variational inference + action = active inference	89
3.2. Generative model hunting	95
3.3. Where does this take us, representationally speaking?	102
3.4. Incorporating Action	105
3.5. Generative processes and active systems	107
3.6. Surprisal vs Divergence	112
<i>4. One Weird Trick to Stay Alive: the FEF's philosophy of life</i>	118
4.1. The organism's agenda	118
4.2. Self-organization and steady state	122
4.3. Cybernetics redux	127
4.4. Stability and agency	131
4.5. The more things change, the more they stay the same	136
<i>5. Active inference beyond the brain</i>	142
5.1. A brief review of causal inference, for the purpose of more clearly elucidating the original nature of a 'Markov Blanket'	145
5.1.2 Bayes nets in causal dress	148
5.2 Markov blanket Realism	154
<i>Recap</i>	162
<i>6. The free energy framework and the missing cycle</i>	166

6.1. The missing cycle	171
6.2. The Existential triad	174
6.3. Extensional Ambiguity	180
6.4 The Markov Blankets Last Gasp	184
7. Seeking closure	188
7.1 EISA-closure	189
7.2. Bionactivism, autonomy and Closure	194
7.3. Operational Closure	197
7.4. From freedom and stability to dependence and purpose	202
7.5. Self-production is not 'homeostasis of organization'	208
8. A theory of everything, or just of every 'thing'?	213
8.1. Processes and substances	217
8.2 The instability of organic parts	223
8.2.1 Why metabolism matters	228
8.3. The FEF and the machine concept of the organism	237
8.3.1 Is surprisal-minimization the substance of an organism?	243
8.3.2. Identity crises: from behaviour to bodies	246
8.3.3. Stability across the scales	249
8.3.4. Why the machine concept of the organism fails	253
8.3.5. Life as the process of seeking stability?	263
9. Bioenactivism and autonomy: from process closure to closure of constraints	267
9.1. Processualism and bioenactivism	267
9.2. The problems of defining autonomy as closure of processes	277
9.3. Constraint closure	279
9.3.1. Constraint causation	282
9.3.2. Constraint production	288
9.4. Constraint Closure as a theory of living systems	294
9.4.1. Constraint closure and adaptivity	302
9.5 Constraint closure as a basis for intentionality	306
9.5.1. Organisms as networks of reasons	307
9.5.2. Organisms as causes of their own activity	316
Conclusion	321
Appendix: What's the use of a concrete blanket?	329
A.1. God's great causal graph	332
A.2. Naturalised mathematical realism	338
A.3. Absolute units	340
A.4. A second stability requirement	342
References	345

Lay summary: Control is not the goal

“Artificial intelligence (AI) refers to systems that display intelligent behaviour by analysing their environment and taking actions – with some degree of autonomy – to achieve specific goals.” – The European Commission’s High-Level Expert Group on Artificial Intelligence (2018)

“Intelligence is the computational part of the ability to achieve goals in the world” – John McCarthy, 2007)

“I am less disturbed by a science that claims to be the equal of God than by a science that drains one of the most essential distinctions known to humanity since the moment it first came into existence of all meaning: the distinction between that which lives and that which does not; or, to speak more bluntly, between life and death.” - Jean-Pierre Dupuy (2018)

This thesis could not have been written without support from the inhabitant of a large basement, just outside of Reading. Here, in a space that looks like the offspring of Nasa’s mission control room and the New York stock exchange, lives the ‘brain’ behind the UK’s power network. This is National Grid’s Electricity Control Centre, a partially-automated system responsible for monitoring millions of sensors across the UK to maintain a constant level of 400 kilovolts throughout the grid’s five thousand miles of powerlines – ensuring a continuous flow of electricity to this computer and every other electrical system in the country.

If supply outstrips demand, the ECC releases pressure by lowering electricity prices. If demand increases, the control centre rapidly spins up its reserves to compensate – for instance by instructing Dinorwig power station in

Snowdonia to release the nine billion gallons of water stored near the summit of Elidir Fawr, converting this potential energy into enough electrical power to ensure that millions of households across the UK can simultaneously make their morning cup of coffee without plunging the nation into a blackout.

In recent years this regulatory task has become more challenging, thanks to the replacement of traditional voltage support reserves, such as coal and gas, with the environmentally friendlier but inconstant alternatives of solar and wind. In order to budget with such unreliable forces, the grid has had to get 'smarter' via the deployment of AI systems to forecast demand increases or weather-induced outages, enabling the ECC to automatically take predictive actions to maintain the delicate balance of the grid. There is no in principle reason why the national grid's control centre could not become a fully automated, anticipatory system.

The ECC is, in a weak sense, dependent upon its own successful operation. If supply is not matched to demand, then on-site backup generators will sustain its regulatory system's operation throughout the consequent blackout for only so long before the system runs out of the power needed to continue operating. As civil engineering professor Guy Walker (2013) describes, it is "in some respects closer to an organism than a machine."

The British cybernetician, William Ross Ashby, would have endorsed the comparison. For Ashby, this homeostatic regulation of 'essential variables' is all that there is to being a living system, or indeed, to being any sort of system at all. Survival means nothing more than stability amid perturbation and the only difference between an organism learning to adapt to a new environment, the regulatory activity of the ECC, and the action of a pendulum returning

to its equilibrium point, is the complexity of the re-stabilization mechanisms involved. As he puts it:

“We have heard ad nauseam the dictum that a machine cannot select; the truth is just the opposite: every machine, as it goes to equilibrium, performs the corresponding act of selection. Now, equilibrium in simple systems is usually trivial and uninteresting; it is the pendulum hanging vertically; it is the watch with its mainspring run down; the cube resting flat on one face... What makes the change, from trivial to interesting, is simply the scale of the events.” (Ashby, 1962, P. 70)

The prevailing currents of 20th-century artificial intelligence were not kind to Ashby. It was far from obvious to his fellow cyberneticians how a machine like his ‘homeostat’, which randomly flailed around when disrupted until it re-obtained a stable state – a ‘sleeping machine’, as fellow cybernetician Grey Walter caricatured it – might one day “play chess with a subtlety and depth of strategy beyond that of the man who designed it” (Ashby, 1948).

The intelligence researchers of the late 1900s preferred the clean-shaven respectability of algorithmic symbol manipulation. Accordingly, much work in AI and cognitive science focused upon characterizing the ‘means’ of intelligence rather than the end – what John McCarthy refers to as the ‘computational part’ of intelligence. The problem, as a number of philosophers pointed out, is that in purely formal terms computation is trivial and any physical assemblage of moving parts, whether a pendulum or a difference engine, can be interpreted as executing an algorithmic operation. If we want to distinguish intelligent systems that are genuinely adding, subtracting, inferring or integrating, then a syntactic account of computation cannot be prioritised over an account of the meaning of these formal operations, or the function they perform. Ends must come before means.

This century has proved more receptive to the idea of intelligence as control and a focus on the achievement of ends rather than an understanding of the means. Today, artificial intelligence is almost synonymous with deep learning where success is measured by the reliability with which the output of a many-layered, network converges upon a value that we've determined as desirable. The algorithms by which such networks achieve this are often impenetrable and it is precisely this opacity, with the occasional unexpected behaviour that results, which seems to motivate our perception of them as 'intelligent.' The European Commission's definition of artificial intelligence above might well be reformulated without loss of meaning as "a system that does something we find useful in a manner we find hard to understand."

So, we speak of complex regulatory systems like the ECC as 'autonomous' or 'smart', we credit them with 'trying' to achieve goals, and we get mad at them when they fail. For Ashby, such a system is thus:

"heaven-sent in this context, for it enables us to bridge the enormous conceptual gap from the simple and understandable, to the complex and interesting. Thus we can gain considerable insight into the so-called spontaneous generation of life by just seeing how a somewhat simpler version will appear in a computer." (Ashby, 1962).

The concern may still remain as to how we are supposed to explain the kind of adaptivity and versatility associated with genuine intelligence in terms of a simple imperative of maintaining stability. Deep neural networks, for instance, are famously brittle, often incapable of transferring their capacity in one domain to another and liable to forget everything they've already learnt upon receipt of new information.

In cognitive science, this problem has drawn increasing attention with the proposal that viewing the brain as a hierarchical architecture for predictive control over multiple timescales can account for the emergence of all our

intellectual capacities: from inference and imagination to offline action planning. Yet unlike Ashby, who saw his program as the elimination of teleological attitudes from the scientific explanation of intelligent behaviour, biological or otherwise, many advocates of this predictivist framework have enthusiastically adopted the description of such control hierarchies in purposive terms. By interpreting predicted states as ‘goals’ that a hierarchical control architecture is ‘trying’ to bring about, they claim that they can describe for the sense in which such a system has purposes and the intentions to bring them about.

If we consider a system’s predictive goal as being the stable state in which it is most likely to be found, then its return to this state when disrupted becomes an act of control via the minimization of prediction error. Yet, as Ashby emphasized, everything from a pendulum to a watch spring ‘rejects’ unstable states to ‘select’ a stable equilibrium. If intelligence, agency and intentions reduce to nothing more than this form of predictive control then they are either everywhere, or nowhere at all.

We may take some small comfort in the reminder that our own brains and their prediction-error minimizing behaviour are vastly more complex in their capacity for re-stabilization than a simple pendulum – but what of the ECC? Is my striving towards the goal of finishing this thesis no more significant than its increasingly complex anticipatory operations towards the ‘goal’ of maintaining that 400 kilovolts set point?

Arguably, much more rides upon the ECC’s success. But while it would certainly matter to us if the National Grid suffered a power loss, does the ECC itself care? There is, I suggested, a sense in which the ECC depends upon maintaining the grid’s 400-kilovolt supply. Its regulatory operations

require electricity and absent its own activity in ensuring they receive this, those operations would eventually cease. The ECC is not particularly novel in this respect. The planet's hydrological cycle constitutes a similar cycle of mutually dependent processes, yet it is no longer fashionable to attribute wrathfulness, or other agential attributes, to the weather.

While the ECC is not so unlike the hydrological cycle it is, I argue, quite fundamentally different from an organism. In a living cell, there is not only dependence of activity upon activity but of *existence* upon activity. The ECC's physical parts, its silicon transistors and copper wires, are intrinsically stable and will not automatically disintegrate if energy ceases to run through them. Re-introduce a supply of electricity and it can continue operating in just the same manner as before. In contrast, the components of a cell, its internal enzymes and surrounding membrane, are inherently unstable and dependent upon the very metabolic activity that they enable for their ongoing repair and replacement. Deprived of the supply of matter and energy that fuels this activity the very structure of the organism will irreparably disintegrate.

The ECC is a particularly complex stabilization mechanism, but a stabilization mechanism is all that it is. We call this stable state of 400kv a 'goal' simply because it is reliably achieved and desirable *to us*, and we call the ECC, but not the egg-timer, 'intelligent' merely because we do not understand it.

Perhaps this is all there is to being an intelligent agent. Perhaps talk of goals and intentions are convenient heuristics to abstract away from the messy mechanical details of a system's operation. This instrumentalist view of intelligence is the only one available within a mechanistic perspective, where

the stability of parts is presumed and we are interested only in how they produce a particular behaviour.

But the mechanistic does not exhaust every possible form of existence. The bacteria swarming all over my keyboard may not seem particularly impressive in comparison to the ECC's ability to probabilistically model the likelihood of a variety of future events and allocate resources in anticipation. Unlike the ECC, however, a bacterium can justifiably take credit for its own existence. Unlike the ECC, its structure is precarious and reciprocally dependent upon the activity it produces. The bacterium cannot be cut off from energy flows without consequences for its physical integrity. Unlike the ECC, the bacterium does not merely respond to perturbations to an otherwise stable state, it is intrinsically dynamic as the inevitable degradation of its internal components releases energy to drive the activity that rebuilds them.

The organism, as the philosopher Hans Jonas, put it has a 'needful freedom' in relation to matter – both independent of any particular material basis and dependent upon a continual supply in order to continuously reproduce itself. Energy flows *through* the fixed structure of the ECC but in the organism everything flows.

The bacterial colony of my keyboard must reliably maintain certain metabolic processes as an existential imperative. If they do not achieve the necessary flows of matter and energy, they will not only cease to be active – they will cease to exist. To say the bacterium *needs* to constantly rebuild itself in this way is not an anthropocentric projection. It has nothing to do with us at all.

So, to be a realist about purposes, intentions, or goals is not vitalistic or unscientific. It is to attend to the natural and intrinsic features of biological

existence that are erased by a purely mechanistic conception of the universe. Taking the function of a system to be instrumental to our purposes, as the mechanist often does, was never a particularly satisfying solution in the first place. At some point, this instrumentalist will have to explain what's so special about *us*, in virtue of which we can have needs and purposes for other things to be relative to.

If this 'bioenactive' view of what it means to be an intentional agent supplies the ends of intelligent behaviour, the question remains as to how complex and creative intellectual capacities could arise from such a simple imperative of survival. Hierarchical predictive control, I argue, serves this biological purpose well enough in helping us anticipate and avoid threats that might disrupt our ongoing metabolic self-production. In so far as such predictive architectures have also been advanced as a potential explanation for how we reason about everything from future actions to other peoples' mental states, so this framework may serve to show how the basic biological goal of survival could – through evolution, learning, and social scaffolding – lead to the catching of baseballs, the dancing of tangoes, the arranging of flowers, or the scaling of mountains.

In making the processes of cognition less opaque, perhaps such models will also make them seem less intelligent. But if the intelligence we ascribe to another system is only a reflection of our own ignorance, I'm not sure it was all that worth caring about in the first place. Whether a hierarchical predictive model is autonomous, whether it is an agent with needs and investment in its own activity, depends not upon our perspective but on whether it is a vital constituent of an intrinsically unstable system that is only temporarily stabilised by its own operations. Control is the means, not the goal.

You can grant the ECC access to every power line across the globe, install the most sophisticated predictive algorithms and the largest of large language models. You can hook it up to a speech synthesizer and instruct it to sing of its feelings for electricity, install Windows 95 and use it to play *Doom*, or tear out its silicon chips for earrings and sell them on Etsy. Either way, the ECC itself will not care. No matter how complex it becomes and the variety of perturbations that it is able to achieve a stable state in anticipation of, in so far as that stable state is a 'goal' it will only ever be ours.

Unlike the ECC the bacterium *needs* to seek out continual flows of matter and energy from its environment to continue to exist. You probably don't care if it succeeds. The bacterium likely lacks the counterfactual flexibility or recursive self-modelling required to care either. But unlike the ECC it has something that it could, at least in principle, learn to care about.

If intelligence is about pursuing goals then, I argue, it has to start here.

Overview

Since Varela, Thompson and Rosch first introduced the term in *The Embodied Mind* in 1991 the 'Enactive Approach', together with its embodied, extended and embedded affiliates, has grown in popularity across cognitive science. This swell has not amounted to a sea change, however. There may well be general agreement that we would do well to pay attention to the way cognition is sculpted by our bodies and environments, but the full-blooded enactive standpoint remains outside the mainstream.

This is perhaps unsurprising, for the phenomenological orientation of the enactive approach does not rest easily with the scientific realism that dominates both cognitive science and contemporary Anglophone philosophy. Where the scientific realist assumes investigator and investigated to be strictly independent, the phenomenologist takes self-and-world as inextricably entangled, views knowledge or understanding as a matter of attunement between them and conceptualizes intentionality as the directedness of an action rather than the having of a representation. Where the scientific realist is tasked with overcoming scepticism, to explain how the deliberative operations of our internal machinery could come to mirror the independent structure of the world beyond, the problem for the phenomenologist is the question of how such a self/other distinction comes to arise at all.

As a result of this tension, much work in 4E cognition has instead adopted a scientific realist notion of embodiment in terms of a physically-instantiated sensorimotor system, either seeking to extricate enactive ideas from their phenomenological frame or ignoring them entirely. Such accounts *may* follow the enactivist in rejecting a view of knowledge and cognition as directed

towards the goal of accurate reconstruction. Yet in so far as they do not supply an alternative account of the norms by which our cognitive processes are governed they are incomplete, lacking a foundation for distinguishing a body from an object, an action from a mere movement, or an autonomous agent from a machine.

In the first chapter, I introduce the enactive approach, set it in this phenomenological context, and introduce its primary task as supplying an alternative *teleological* conception of intentionality as the striving of a system towards some non-reconstructive goal or norm. To supply this, I advocate what I'll refer to as '*bioenactivism*', which aims to locate this immanent teleology in the autonomy distinct to living systems, whose precarious existence is dependent upon their own activity in working towards their continuous self-production.

While this may give us a foundation for normative evaluations, as Di Paolo (2005) has argued, it is only the all or nothing normativity of continued existence. To attribute a graded normativity to an agent's interactions with the world, we need to introduce what he terms '*adaptivity*': the capacity of a system to regulate these interactions so as to move towards, or away from, states according to whether they threaten, or support, its autonomous organization. While Di Paolo, Buhrmann & Barandiaran (2017) have progressed the enactive account by describing how this adaptivity might be '*scaled up*' through the process of sensorimotor equilibration, accounts of how this is implemented are still needed.

This, as I propose in Chapter 2 is where predictive processing [PP], the account of cognition in terms of hierarchical prediction error minimization can come in. When presented as a model of how we can extract patterns

from sensory input over multiple timescales to develop predictive models of the relationships between inputs and outputs, PP looks to do away with what Hurley (1998) termed the ‘classical sandwich’ view of cognition, as an independent operation that goes on in between the separate processes of perception and action (Vázquez, 2020; Nave et al., 2020; Harvey, 2018; Bruineberg, Kiverstein, & Rietveld, 2018; Kirchhoff, & Robertson, 2018; Clark, 2015).

In so far as PP is also presented as a mechanism for performing approximate Bayesian inference and has already been used to describe how we might solve a variety of ‘higher-level’ cognitive tasks – from action-planning (Pezzulo, 2017) to counterfactual reasoning about others’ mental states (Palmer, Seth & Hohwy, 2015), so there is reason to be optimistic that it might provide an embodied framework for these ‘representation-hungry’ forms of offline cognition, that have traditionally argued to be unapproachable via sensorimotor coordination alone (Roelofs, 2018; Matthen, 2014; Clark & Toribio, 1994).

Yet, as I have argued, a recognition of the interdependence of perception and action and the rejection of a representational starting point does not an enactivist account make. The enactive approach not only views perception and action as constitutively interdependent but also as directed towards some other norm or goal. In so far as these accounts of ‘sensorimotor predictive processing’ do not address the question of what makes something the directed action of an agent, rather than at the mere movement of a physical object, so they constitute only a partial step towards an enactive account.

Such accounts cannot be complete unless they connect up with something like bioenactivism’s grounding for the non-reconstructivist normativity that

we coordinate our sensorimotor engagements with respect to. Without this alternative explanation for the function of a predictive brain, Clark (2015) and others lack justification for redescribing a bare covariance between top-down and bottom-up signals in intentional terms as ‘prediction-errors’ that a system is ‘trying’ to minimize. Thus, PP not only stands to help scale this basic biological intentionality to ‘higher’ forms of cognition but also stands to benefit from grounding the function of the predictive mind in the bioenactivist account of autonomy and intentionality.

So how exactly might prediction-error minimization relate to the preservation of biological autonomy? In Chapter 3 I introduce Karl Friston’s ‘Free Energy Framework’ [FEF], the first component of which is a formal description of the equivalence between approximate inference and predictive control, as might both be implemented by a predictive processor. In Chapter 4, I describe the second component of this, the Free Energy Principle [FEP] which purports to formalize the survival of an autonomous system in terms of the kind of stability that this predictive control affords. As described in Chapters 5 and 6 respectively, this formalization of inference-as-control is then supplemented with the addition of a Markov blanket, which individuates the organism from its environment, and a set of coupled stochastic differential equations, which are used to formulate the notion of a sensorimotor loop. When put together, the FEF’s advocates argue that these components provide the means both to ground the intentionality of living systems *and* to scale it up to higher-level cognitive processes, via an implementation story like PP.

In viewing survival as a matter of the stability of essential variables, the FEF’s account is strikingly similar to the theory of life in terms of ‘generalized homeostasis’ offered by the cyberneticist, W.R. Ashby half a century earlier.

Yet, as I argue in Chapter 7, biological autonomy is not reducible to homeostasis and the FEF's definition of it turns out to be trivial – as applicable to any stable mechanism as to a living organism.

Perhaps this is not the FEF's fault. Perhaps biological autonomy is not reducible because it is a vitalistic concept with no place in a good scientific account of organisms and their cognitive processes. In Chapter 8 I consider the FEF's prospects as an account of living systems, independently of my prior bioenactive commitments. Here I argue that the FEF is not only too general to provide an account of what distinguishes the living from the non-living, it is also too specific in that it makes claims about necessary imperatives which simply are not necessary for organisms. Unlike inorganic structures, living systems depend upon neither the stability of their parts, nor the stability of the interactions between these, which the FEF takes as definitive of a system

That the FEF fails is not enough to establish that bioenactivism has succeeded in differentiating life from non-life – in order to ascribe intentionality to the latter alone. Indeed, in Chapter 9 I will argue that it has not. While the prevailing bioenactive account of autonomy in terms of process closure has advantages over the free energy framework – in so far as it provides a relational account of why some variables might need to be stabilized and why some are free to change in open-ended ways – like the FEF, it neglects the unique thermodynamic status of living things. By abstracting away from the molecular interactions of autopoiesis, in favour of closure amongst a network of mutually dependent and precarious processes, this 'process closure' may be general enough to apply to all scales of biological organization, but it is also too general to distinguish living, intentional systems from machines.

What makes organisms special, as Moreno and Mossio (2015) describe, is that there is no clear distinction between precarious flows of energy (processes) and the invariant structures that constrain them, for these structures are themselves reciprocally dependent upon those flows of energy in turn. It is this reciprocal dependence that they formulate in the alternative notion of constraint closure. This, I argue, succeeds in combining the advantages of both thermodynamic and relational accounts of living systems, in order to describe what it is about living systems alone that makes them genuinely autonomous agents. And it is here, rather than in the statistical constructs of the FEF, that enactivists should look to ground a teleological concept of intentionality.

1. Biodynamic Enactivism

1.1. The Enactive Approach

“The big problem I have about enactivism is figuring out what it is.” - Ned Block (quoted in Meyer & Brancazio, 2021)

Labels are a necessary evil, and ‘enactivism’ has caused its fair share of confusion. The lack of a specific, widely agreed-upon point of reference can’t be helped, the same could be said of terms like, ‘Computationalist’, ‘Christian’, ‘fun’, or ‘soup’ (Gualeni, 2017). Still, I hope to at least be clear about what I do and do not mean by ‘enactivism’, and ‘bioenactivism’ specifically. I will also give some account of why I take this to pick out a coherent position and tradition within the array of ideas that are typically lumped together and presented as an alternative to ‘classical’, ‘computationalist’ & ‘cognitivist’ approaches to the mind.

Firstly, what I do and do not mean by ‘enactivism’. One way of identifying an enactivist account is genealogical, in terms of whether it developed out of Varela, Thompson & Rosch’s *The Embodied Mind* (1991) in which the label of ‘the enactive approach’ was introduced. Often referred to as ‘autopoietic enactivism’ (though as I will explain that label is not ideal) Varela, Thompson & Rosch’s work contains many of the ideas that I take to be central to the enactive tradition, in drawing upon phenomenology to motivate the proposal of an alternative, non-representationalist starting point for cognitive science. Nonetheless, others may disagree about what the ‘key enactivist ideas’ of *The Embodied Mind* are. As such, depending upon claimed connections to this particular text threatens to lump together incompatible views, whilst

excluding more closely related work on the basis of the lack of historical connections to Varela et al.'s work.

The problems with this are apparent in how the label 'enactivism' is often applied quite loosely by philosophers and cognitive scientists based on one of two criteria, neither of which I take to be adequate to pick out a distinct and unified understanding of what cognition is. The first is to refer to any negative position that forbids the use of representational talk in understanding the mind (eg. Nanay, 2014); the second, as encompassing a wide range of positive proposals, united by taking the coordination of action to play some essential, and underappreciated role, in our cognitive and perceptual lives (eg. Gangopadhyay & Kiverstein, 2009; Ward, Silverman & Villalobos, 2017).

I take identifying one's approach with the first of these to be particularly unconstructive. 'Representation' means many different things to different people, and whether or not anything worthy of the name will play a role in our account of cognition is not something to be committed to at the outset, but should be a downstream consequence of whatever metaphysical, conceptual, and methodological picture we develop. I take *The Embodied Mind* not as an absolute writ against any talk of 'mental representations', but as presenting an argument against assuming that cognition's primary function is the veridical recapitulation of a mind-independent world, and showing how we might approach it from an alternative starting point in terms of the perceptual guidance of action.

So, the success of the enactivist program does not depend on whether it succeeds in evading representational commitments at every point of its development. Instead the test is whether this alternative starting point does,

as Varela, Thompson & Rosch suggest, help us evade the sceptical thicket that has entangled cognitive science since its inception. Even if the enactivist is correct, a non-representationalist starting point does not entail anti-representationalism. As Thompson (2011) remarks, in relation to the emulation theory defended in representational terms by Lucia Foglia and Rick Grush (2011):

I argue against representationalist theories that separate perception and action, instead of recognizing their constitutive interdependence, and that neglect the ways autonomous agents bring forth or enact meaning in perception and action (see pp. 10, 58–9). Since the emulation theory does not require these typical features of representationalism, my objections to representationalism need not apply to the emulation theory. (P. 19)

For this reason, I take the argument that an account of cognition cannot be enactivist if it posits internally-realized structures in the brain that facilitate decoupled action guidance and imagination to be misguided. Just because advocates of such accounts, have often chosen to call these structures ‘representations’ or ‘models’ does not commit them to the reconstructivist, representation-*first* approach to the cognition that the enactive approach rejects.

While this first understanding of ‘enactivism’ pigeonholes its development in a manner that I would reject, the second use of ‘enactivism’ to refer to the family of views that reject the detachment of cognition and action is less pernicious. Though defining enactivism in this loose sense leaves means grouping together a wide variety of different approaches, it does at least seem consistent with the definition of ‘the enactive approach’ given by Varela, Thompson and Rosch (1991) as defined by two points. “(1) perception consists in perceptually guided action and (2) cognitive structures emerge

from recurrent sensorimotor patterns that enable action to be perceptually guided” (p. 173)”

This breadth can be useful, in so far as it identifies a diversity of positions that have been proposed as alternatives to what Susan Hurley (1998) called ‘the classical sandwich’ view of cognition as a distinct, classically computationalist procedure, which occurs in-between the disconnected and peripheral processes of perception and action. Still, this rejection of the ‘classical sandwich’ is also found in the ecological psychology of J.J Gibson (1979), the perceptual control theory of William T. Powers (1973), various ‘skill theories’ of perception such as those of Gareth Evans (1982), Rick Grush (2007) or Susannah Schellenberg (2007), and the ‘interactivism’ of Bickhard (2009) to name just a small scattering of examples. None of the above identifies with enactivism, and there are explicit tensions between some of these positions and the enactive view – for instance between Gibson’s realism and the constructivist metaphysics of *The Embodied Mind*.

So, the view that the capacities of action, perception, and cognition are related to each other in some important way is insufficiently distinctive to identify the ‘enactive’ view specifically. For these reasons then, I want to clearly distinguish the view that interests me from both what is sometimes called the ‘sensorimotor enactivism’ of Kevin O’Regan and Alva Noë (2001), and from the ‘radical enactivism’ of Dan Hutto & Erik Myin (2012) (For a nice overview of the different positions to claim the ‘enactive’ header, see Ward, Silverman and Villalobos, 2017).

The first of these is enactivist in the second, broader sense and bears little more similarity to Varela et al.’s position than some of the ‘action-oriented’ theories listed above. O’Regan and Noë are not so much concerned with

banning representational talk, but with the re-conceptualization of it. As they put it, “seeing lies in the *making use* of the representation, not in the having of the representation” (2001: 1017). While sensorimotor enactivism also shares the phenomenological inspirations that animate *The Embodied Mind* and its descendants, unlike this tradition, it has tended to focus on narrower questions regarding perceptual content, rather than metaphysical and epistemological issues regarding agency, self, environment, and intentionality. That said, O’Regan and Noë’s narrower proposal, which they prefer to call ‘sensorimotor theory’, could be situated within this broader enactive project with relative ease.

The same cannot be said of Hutto & Myin’s (2012, 2017) radical enactivism, which focuses more on anti-representationalism in its self-identification with the ‘enactive’ label. Largely disowning the phenomenological perspective that shaped the development of the enactive approach, Hutto & Myin primarily focus on scolding other accounts for being inadequately committed to the exorcism of all talk of ‘intentions’ and ‘content’ from cognitive explanation. As Thompson (2018) notes, and as I will describe, this conflation of intentional content with representation, and the disavowal of both, is precisely what the enactive approach, as introduced in *The Embodied Mind*, seeks to avoid.

1.1.1. The enactive approach as naturalistic phenomenology

So neither self-identification, nor common applications of the label ‘enactivist’, are particularly helpful guides to picking out a coherent and distinctive approach. As such, I propose that the best way to identify enactivism is not in terms of commitment to the eradication of representational explanations in cognitive science, or as any view that analyses cognition in terms of the coordination of perception and action, but

specifically as the view that they are constitutively interdependent, content-involving and intentionally directed towards a goal that is not, primarily, reconstructive. As both an inspiration for, and consequence of, this view of cognition, we find also the rejection of scientific realism, and the idea that the objects of cognition and the cognizer of objects are two strictly independent realms.

Such a definition does better at distinguishing enactivism from other action-oriented approaches, such as ecological psychology and radical enactivism, while serving to identify the common thread that runs through more recent ‘canonically’ enactivist works, such as *Mind in Life* (Thompson, 2007) and *Sensorimotor Life* (Di Paolo, Buhrmann & Barandiaran, 2017). As I will argue in the next section, it is also beneficial in allowing for the identification of shared approaches, irrespective of whether their advocates refer back to particular texts or describe themselves in particular terms.

That Varela, Thompson and Rosch took their thesis about the relation between perception and action to be more than just an empirical discovery, is evident in their concern not only with providing an alternative methodology for cognitive science but also an alternative metaphysics to replace the scientific realism that was dominant at the time, and which remains so across the philosophical anglosphere today (Bourget & Chalmers, 2021). Where scientific realism takes mind and world to be strictly independent and asks how the former can become cognizant of the latter, Varela, Thompson and Rosch explicitly situate their project within the phenomenological approach, which views mind and world, self and environment, as inextricably entangled.

As the tradition introduced by Edmund Husserl, I understand phenomenology's central feature to be the continuation of the Kantian project to find a way between idealism and metaphysical realism. As Zahavi (2004) puts it, 'Phenomenology is basically, I would insist, a transcendental philosophical endeavour, and to dismiss that part of it, is to retain something that only by equivocation can be called phenomenology' (P. 340).

Thus while the 'phenomenological method' is most commonly identified with Husserl's famous 'epoché' – the practice of analysing the objects of experience as they appear to us without distorting this through a prior commitment to their mind-independent nature – this is only the first step. The second is to attempt to identify the transcendental structures, such as perspectivity and temporality, that are preconditions for the possibility of our experiencing this world of objects at all (Zahavi, 2003; Moran, 2002).

In doing so, we can attempt to save realism about the empirical world by sacrificing a metaphysical realism that treats it as an independent 'given'. For the phenomenologist, this means appreciating that the world of experience is a construction that partially involves our own activity, but, crucially, that this construction is nonetheless empirically real, and not a purely subjective matter of free *individual* choice. As Varela, Thompson and Rosch (1991) describe:

(...) cognition is not the representation of a pre-given world by a pre-given mind but is rather the enactment of the world and a mind on the basis of a history of the variety of actions that a being in the world performs (Varela et al. 1991: p. 9)

So far, so Kantian. What differentiates phenomenology is the identification of these necessary preconditions, not with the conceptual scheme of a pure 'knower', but with the 'embodiment' of an agent, and consequently a

developing appreciation for the contingent, dynamic, and historical nature of these transcendental structures, in contrast to the supposedly atemporal and absolute foundations of Kant's categories (Mohanty, 1978; Zahavi, 2003).

Three key points about embodiment here. First, is that the term is intended to refer not only to the biological body as it is ordinarily understood but also to an extended network of cultural, linguistic and environmental structures, in so far as these afford, solicit and constrain possible actions. Second, is that the body now becomes split accordingly into the empirical aspect given to us in experience and investigatable through scientific methods versus the transcendental lived aspect revealed via phenomenological analysis. Third, the interesting thing about linguistic, cultural, and biological structures is that they can vary and change. As such what is revealed by phenomenological analysis as a necessary precondition for some facet of experience may, as the sociologist Alfred Schutz (1959) criticized of Husserl's early apodictic foundationalism, still be contingent upon our particular situation, rather than reflective of absolute and eternal truths.

Thus, in the dual aspect of the phenomenological body, as Petitot, Varela, Pachoud & Roy (1999) argue, "a transcendental analysis and a natural account are intrinsically joined." It is this recognition of the empirical aspect of transcendental structures, rather than in the discarding or downgrading of phenomenology's transcendental dimension, that Zahavi (2004) suggests opens up the possibility of a genuinely naturalistic phenomenology.

While the early Husserl may have indeed sought the "universality, necessity, apodicticity" to deliver an 'absolute grounding' of human knowledge (Husserl, 1982/1913, p.19) as Zahavi (2003) argues he increasingly recognizes the importance of embodiment (Husserl, 2001/1920;

1997/1907), along with the revisability of phenomenological claims (Husserl, 1970/1936) and the potential for fruitful interaction between empirical and transcendental approaches (Husserl, 1999/1929). That said, it was arguably Merleau-Ponty who first fully appreciated how the two-fold nature of the phenomenological 'body' contained the prospect for bringing together the insights of both scientific and transcendental phenomenological approaches to cognition.

As he describes it, phenomenological analysis concerns not a quest for unshakable foundations, but rather “an intellectual taking over, a making explicit and clarifying of something concretely experienced” (1964, P.68). The clear-cut distinction between such insight, and an empirical fact is now blurred and recast such that, “The a priori is the fact understood, made explicit, and followed through into all the consequences of its latent logic; the a posteriori is the isolated and implicit fact”. (2013/1945, P.221). In this regard, as he argues, there is a continuity between the inductive and generalizing endeavour of the scientist and that of the phenomenologist:

There are not two truths; there is not an inductive psychology and an intuitive philosophy. Psychological induction is never more than the methodological means of bringing to light a certain typical behaviour, and if induction includes intuition, conversely intuition does not occur in empty space. It exercises itself on the facts, on the material, on the phenomena brought to light by scientific research. There are not two kinds of knowledge, but different degrees of clarification of the same knowledge. (1964, P.24)

The scientist and the phenomenologist share the method of attempting to extract invariant features from varying circumstances. One seeks these in the experience of bodies, languages or societies as objects *in* our experience, the other seeks their invariant features as structures *of* experience. But neither escapes experience altogether and neither has a direct methodological line to some of pure realm of *nature* that lies outside or beyond it. Phenomenology

is incompatible with scientific naturalism only in so far as the scientist forgets this and takes properties of her models to be the irrevocable truths of mind-independent reality.

As Merleau-Ponty characterizes this:

Science manipulates things and gives up living in them. It makes its own limited models of things; operating upon these indices or variables to effect whatever transformations are permitted by their definition, it comes face to face with the real world only at rare intervals. Science is and always has been that admirably active, ingenious, and bold way of thinking whose fundamental bias is to treat everything as though it were an object-in-general – as though it meant nothing to us and yet was predestined for our own use. (1964, P. 290).

So enactivism as naturalistic philosophy, does not mean a naturalized phenomenology, wherein phenomenological descriptions of the lived body are reduced or eliminated in favour of empirical descriptions of the body as an object in our experience, but one in which phenomenological analysis and scientific method inform one another, and where neither is taken as an apodictic foundation to which the other must submit absolutely. Such a picture, as Gallagher (2018, 2017) argues, may involve not just a revision in the authority we accorded to the scientific method, but also, as it does for Merleau-Ponty, a transformation in how we conceptualize its objects of investigation, towards a view of nature itself as irreducibly relational and intersubjective, constituted by the interactions between embodied agents.

It is this view of phenomenology, science and nature, reflected in the quote from Varela, Thompson, and Rosch (1991) on page 31, that I take as the basis, though not the original contribution, of the enactive approach. One way of distinguishing its different strands is in terms of which dimensions of our embodiment are emphasized for investigation: whether biological self-constitution, as in Thompson (2007) or Weber & Varela (2002),

sensorimotor dynamics (Di Paolo, Buhrmann, & Barandiaran, 2017) or social and linguistic networks (Di Paolo, Cuffari, De Jaegher, 2018).

While this rejection of metaphysical realism in favour of transcendental phenomenology is crucial to understanding the enactive approach, this is not to say that anything enactive must hark back to Husserl, Heidegger, Merleau-Ponty, or their direct descendants. There are other ways to the phenomenological and transcendental analysis of embodiment. This might be through an alternative post-Kantian route, for instance via Wittgenstein's similar concern with the conventional constitution of our world, and the priority of our intersubjective situation and linguistic embodiment in making this possible (for interpretations of Wittgenstein in a phenomenological light, see Overgaard, 2006; Egan et al., 2013; Gier, 1981; Zhang, 2008).

Alternatively, one might traverse a different tradition and time period altogether: starting from the foundation of Buddhism with Siddhārtha Gautama's distinction between ultimate and conventional truth, and following how this develops in either the Madhyamaka school, which forms the second philosophical pillar of *The Embodied Mind*, or in the Yogācāra school, where we find arguably the closest parallels with the European tradition of transcendental phenomenology (Lusthaus, 2014).

Where I take reference to the European phenomenological tradition to be particularly useful, however, is in making sense of the enactivist concern with the 'intentionality' of cognition and how this sits alongside the rejection of a representationalist theory of what it is to be a cognitive system.

1.1.2. What are your intentions?

The Enactive approach's focus on intentionality, content, and ideas of the body as a 'vehicle of meaning' (Colombetti, 2010) has led to criticism among those who take the defining mark of enactivism to be its anti-representationalism. Hutto & Myin (2012), for instance, accuse the accounts of Varela et al. (1991), Thompson (2007) Di Paolo (2009) and Colombetti (2010) as being insufficiently radical, and in need of 'rectification.' This confusion about how a system can have intentional content and yet not be representational stems from the quite different ways in which the term 'intentionality' is used in phenomenological and 'analytic' approaches to cognition.

In both analytic philosophy of mind and classical cognitive science, the dominant notion of intentionality is as a relationship of aboutness between the content of a representational vehicle and a target object that corresponds with that content in whatever way is supposed to underpin that aboutness relationship. As such, almost all introductory textbooks emphasize at the outset that "there is no substantial philosophical link" between the philosopher's notion of intentionality and the ordinary meaning in terms of having the goal of bringing something about (Crane, 2015, P.32). This latter state is generally explained as a subclass of the broader category of representational states, and a capacity that depends upon the more basic ability to have states that are *about* things. As Crane puts it, "Intentions in the ordinary sense are intentional states, but most intentional states have little to do with intentions." (P.32)

This view of intentionality as a relationship of aboutness towards an object, and as something more general and basic than that of an intention towards a goal, is traced to the introduction of the term into the philosophy of mind

by Franz Brentano (1874). In explaining its meaning, Jacob (2019) points to the etymology of ‘intentionality’ as deriving from the Latin ‘tendere’ meaning to aim, strive, or tend towards.

This seems, to me, to demonstrate the very opposite from what is intended. In Jacob’s example of an arrow aimed at a target, the arrow is not *about* that target. It tends towards piercing bullseye, not becoming similar to the target or to functioning as some stand-in for it. Similarly, I might aim to work harder, be more polite, get stronger, or knit faster, but in none of these cases is the aim towards some target ‘object’ that I seek to enter a relationship of correspondence with. In ordinary English even ‘object’ has a second meaning tied more to purposes and goals than to aboutness, as when the detective explains that ‘the object of the investigation is to determine who killed Bugs Bunny’, or the vice-chancellor declares that ‘the object of a university is to produce highly-employable graduates.’

As Thompson (2007) explains, where the representational theory of mind views intentionality as states having an ‘aboutness’ relation to some mind-independent ‘thing’, in the phenomenological tradition intentionality is instead a property of ‘acts having directedness.’ It is precisely this notion of intentionality that is given a central role in Varela, Thompson and Rosch’s (1991) proposal for an ‘enactive approach to cognitive science (ch. 9).

We would say that the intentionality of cognition as embodied action consists primarily in the directedness of action. Here the two-sidedness of intentionality corresponds to what the system takes its possibilities for action to be and to how the resulting situations fulfill or fail to fulfill these possibilities. (P.206)

This directedness of an action need not, and usually does not, take the form of a deliberately formulated plan. Instead, it is something continuously manifest in our orientation to the world around us. This world, as the

phenomenologist and the enactivist have it, does not first appear as a neutral array of indifferent objects about which we may later make judgements as to whether they interest us. Rather the world, as we experience it in our unreflective engagements appears a landscape of possibilities for action, that may solicit, or repel us. This is nicely described in a classic example from Merleau-Ponty's *The Structure of Behavior*' (1963):

For the player in action the football field is not an "object," that is, the ideal term which can give rise to an indefinite multiplicity of perspectival views and remain equivalent under its apparent transformations. It is pervaded with lines of force (the "yard lines"; those which demarcate the "penalty area") and articulated in sectors (for example, the "open ings" between the adversaries) which call for a certain mode of action and which initiate and guide the action as if the player were unaware of it. The field itself is not given to him, but present as the immanent term of his practical intentions; the player becomes one with it and feels the direction of the "goal," for example, just as immediately as the vertical and the horizontal planes of his own body. (P.168)

From phenomenology, we thus gain an alternative account of cognition, knowledge, or understanding. One that does not take these, primarily, in terms of having some internal representation with intentional content that depicts, and accurately corresponds with, an independent state of affairs. Instead, these are to be explained in terms of our skill of appropriately responding to the solicitations and affordances of our surroundings, as illustrated in the practical knowledge of a typist, musician, or sportsperson, an idea picked up by Dreyfus (2002) as 'skilled coping', or, in sensorimotor theory, the notion of 'sensorimotor mastery' O'Regan and Noë (2001).

While this phrasing is all very nice in moving towards a non-representational framing of intentionality, the notions of skill or mastery imply not only appreciation for what is possible, but discernment as to what is preferable. The perceptual world for Merleau-Ponty is not merely a neutral matrix of 'I cans'. It is, as described above, an affective milieu, a field of salience and

significance with lines of force that draw and repel us. If cognition consists in our skill at responding to these forces to attune with the world, or to increase our grip upon it, then what is the standard by which this attunement is judged?

In Merleau-Ponty, as in Husserl, this normativity is often characterized in terms of epistemic exploration. Husserl (2001/1920) rather sensuously describes the unseen parts of an object as something that:

“calls out to us, as it were, in these referential implications”. “There is still more to see here, turn me so you can see all my sides, let your gaze run through me, draw closer to me, open me up, divide me up; keep on looking me over again and again, turning me to see all sides.” (P. 41)

Merleau-Ponty (2013/1945) similarly speaks of being drawn towards the optimal viewpoint for a painting in an art gallery, or the understanding of how to shift an object in relation to background lighting in order to best discern its colour. Yet Merleau-Ponty, if not Husserl, recognizes that there is more to the norms of attunement than improving one’s epistemic standing. Football is not an epistemic activity. The player on the field is not drawn to intercept just to learn what it feels like, but because the goal of the game is to score more goals, and their role as a defender is to prevent the other side from scoring. Likewise, the classic examples of typing, organ-playing, dancing, or climbing are not purely activities of exploration and discovery, they are governed by other norms, of linguistic coherence, elegance, or ascension. As such, in his analysis of the forms of our experience in *The Structure of Behaviour*, Merleau-Ponty is concerned also with the structures of what he terms the ‘vital order’, and the ‘human order’.

As he describes with respect to the former:

Thus each organism, in the presence of a given milieu, has its optimal conditions of activity and its proper manner of realizing equilibrium; and the internal determinants of this equilibrium are not given by a plurality of vectors, but by a general attitude toward the world. This is the reason why inorganic structures can be expressed by a law while organic structures are understood only by a norm, by a certain type of transitive action which characterizes the individual. (1963/1942, P.148)

While these allow us to draw on an array of recognized, non-epistemic norms – from standing the appropriate distance from other people to running away from danger – pointing to these norms is not to explain their origin and force. This is a problem still faced by O'Regan and Noë's (2001) sensorimotor theory, that in lacking a theory of selfhood and autonomous agency it does not have the tools to move beyond the mere awareness of neutral sensorimotor possibilities to explain the 'affective allure' in how particular affordances 'grab' or solicit us.

For the enactive approach then, the project of naturalizing intentional content and that of grounding teleology, or normativity, are intertwined. In all cases what we need is an account of what it means for some activity to be directed towards an end, in a manner such that it can be described as genuinely 'trying' to achieve that end, with the consequent possibility of failing. An account that neither defaults to the standard of accurate representation or enters a regress by calling upon the further intentions and projections of some external designer. Just as different strands of the enactive approach may focus on different aspects of embodiment, so they may look to correspondingly different places to ground this teleology.

So, situating the enactive approach as a continuation of naturalistic phenomenology – that is to say as a foundational enquiry into the nature of

knowledge and reality, rather than as a local theory about the objects and methodology of cognitive science specifically – is a better way to identify it as a coherent program. This not only clearly distinguishes the enactive approach from other more local accounts that do not share these aims, such as radical enactivism or sensorimotor theory, but, as I will argue, allows us to identify other instances of the same approach, irrespective of either terminological choices or direct historical links. As with transcendental phenomenology in general, there are other routes to the same ideas. After all, if this is indeed a promising approach then it would be strange if no one else had hit upon it.

1.1.3 Hurley's enactive approach

So, the phenomenological inspiration behind the development of the enactive approach is essential to understanding the particular way in which Varela, Thompson and Rosch describe their project, and how it differs from other views described as 'enactive.' As I mentioned, however, there are alternative routes to the same view of perception and action as constitutively interdependent, content-involving and directed towards a goal that is not, primarily, reconstructive. Of particular note is the view of Susan Hurley, who, while beginning with Kant, takes a different route via Wittgenstein, rather than Husserl, Heidegger or Merleau-Ponty, to develop a view that I take to be far closer to that of Varela et al. (1991) than more commonly cited examples of the enactive tradition.

Hurley's does not claim the appellation enactivist, nor does she draw significantly upon Varela et al. (1991)'s work – while her book, *Consciousness in Action*, was not published until 1998, the acknowledgements note that it was written just a year prior to the publication of *The Embodied Mind*. Nonetheless, she is occasionally classed as a sensorimotor enactivist – a

classification that, as Ward (2016) demonstrates, does a disservice to the sophistication of her account. Rather than just taking the contents of perception to depend on our knowledge of how our movement will change our sensory input, as O'Regan and Noë (2001) do, Hurley proposes a 'two-level *interdependence* view, whereby at both the subpersonal level of sensory inputs and motor outputs, and at the personal level of perception and agency, the capacity to act and the capacity to sense are necessary preconditions for the possibility of each other (Hurley, 1998).

The emphasis on interdependence distinguishes Hurley's enactivism from the various skill theories of perception or the control theories of action. For Hurley, it is essential too that this interdependence is not merely instrumental, as in Gibson's emphasis on their utility with respect to each other's independent functioning, but a constitutive matter of what perception and action are. Action is the control of perception, and perception is the presentation of possibilities for action, and as such, they are necessarily inseparable. In this sense, Ward (2016) argues, Hurley is best described as a 'transcendental enactivist'.

In her Shared Circuits Model, Hurley (2008) proposes how this might be scaled up through a hierarchical structure of control systems, a striking precursor to the predictive processing accounts I will discuss in the next chapter. As has been suggested as a consequence of predictive processing (Kiefer & Hohwy, 2018), so Hurley also proposes that cognitive contents must thus be attributed holistically, in terms of the potentially flexible, relations between input and output, perceptions and intentions, that make up this overall control system.

This leaves us with the question, ‘control in the service of what?’ Hurley is particularly sensitive to the threat of what she terms ‘the myth of the giving’, whereby one proposes to explain the content of a perception via the content of an intention, which is taken to somehow be primitive and in need of no further explanation. In itself, such a strategy is no better than treating the objects presented in perception as the immediate ‘givens’ of a mind-independent world. In attempting to avoid the subjective regress of this type of ‘just more content’ strategy, Hurley argues that we need a replacement for the role played by the representation of an external world in delivering a non-subjective grounding for content determination. This she proposes, is what an objective account of normativity could deliver.

The problem as Hurley (2003) describes is in distinguishing a system that is genuinely following a norm at which it might fail from the operation of a simple feedback control system, like a programmable thermostat, that also adjusts internal connections between input and output. The thermostat may fail according to *our* goals but there is nothing to prevent describing its behaviour as successfully following some alternative rule. True normativity, she suggests, depends firstly upon something like increased context-sensitivity and flexibility in how a goal is pursued, of the sort afforded by the higher levels of her shared control circuits model, and secondly upon some external teleological constraints such as social context and evolutionary pressures.

As Hurley (1998, 2003) notes, however, she lacks an account of exactly what grounds this supposed ‘teleological context’ and distinguishes it from the basic physical laws that apply to agents and non-intentional systems alike. Mere complexity of behaviour does not appear adequate to the task, for no amount of complexification seems sufficient to entirely dispel the concern

that a robot capable of following (or at least appearing to follow) some sophisticated array of norms might nevertheless be as mindless a zombie as any other machine. As she puts it:

An agent with conceptual abilities has a more richly structured set of behaviors, and perhaps those behaviors must have causes with a certain related structure. But even granting all this, it is not clear why machine or zombie worries should be disarmed by conceptual abilities. If these worries are valid in the first place, why couldn't a machine or zombie have a conceptually structured set of behaviors and reasoning abilities with correspondingly structured causes, yet not be in conscious states? If the worries get a grip to begin with, their grip is not loosened by the enrichment of structure and of norms of rational behavior that goes with conceptual abilities. Even if we were to allow for the sake of argument that conceptual abilities are necessary for consciousness, nevertheless we can add conceptual abilities to perspective and access without yet getting a set of sufficient conditions. (1998, P.162)

Perhaps, she suggests, “the extra ingredient needed in a set of sufficient conditions is not the richer normativity of conceptual abilities, but simply life” (P.162). Yet she does not pursue this possibility. As indicated, her concern is with the gap between a norm-following intentional agent and a conscious one. She presumes the double disassociability of both ‘intentional action without life’, and ‘life without intentional action’.

On the bioenactive view I will develop in this thesis, life and intentionality are not so easily separated. While it might be possible to create life without agency, once we see what life has to contribute to the naturalization of the normative domain we will see why there can be no intentional agency without it. Perhaps, though it is not a topic pursued here, once we have such a rich account of intentional agency, any further puzzle about whether additional ingredients are needed for consciousness will dissolve.

This problem of how normativity and intentionality emerge is not resolved in *The Embodied Mind* either. Despite being explicitly concerned with resisting the elimination of subjectivity, intentionality and agency in favour of the mere objects of the physical sciences, Varela et al. do not go beyond the explicitly mechanistic notion of autonomy as operational closure in the attempt to characterize these. While they suggest that the principal difference between their simplified illustration of this in the cellular automata, Bittorio and a living organism is simply in their respective degrees of complexity, I disagree. As I will argue, this operational closure alone fails to capture the difference in kind between the intentional and teleologically-oriented agent that the enactivist needs versus a mere physical mechanism.

The development of the bioenactive approach thus begins not with *The Embodied Mind*, but with the notion that we can ground these teleological and intentional properties in the self-production of an organism, specifically its intrinsic dependence on an internal metabolic network. While the idea of metabolism as teleological originates with Hans Jonas (1953, 2001/1966), the uptake of this by enactivism can be traced to Weber and Varela (2002) and is continued through Thompson's (2007) *Mind in Life*. So while enactivism is distinguished by a view of cognition as the coordination of perception and action, where these are constitutively interdependent, content-involving and directed towards a goal that is not, primarily, reconstructive, *bioenactivism*, as I will now describe, is the attempt to ground this intentional directedness towards a goal in biological terms.

1.2. Bioenactivism

The bioenactivist inherits a commitment to both naturalism and realism about teleology and intentionality from enactivism more generally and adds to these two further commitments: mind-life continuity and the view of life as self-production. These are supposed to describe how our biological embodiment accounts for the normative and teleological dimensions that the enactive approach posits, but does not explain.

Mind-life continuity is the simpler of these to explain and is expressed nicely by Hans Jonas (2001/1966) as the view that “the organic even in its lowest forms prefigures mind, and that mind even on its highest reaches remains part of the organic.” (p.1)” How such a claim is interpreted will naturally depend on what one takes as the relevant features that both share. Life-mind continuity is thus a relatively minimal commitment, shared by a wide diversity of theorists who Lyon (2006) surveys as the various ‘biogenic approaches to cognition.’

If we, as enactivists, take the ‘mark of the mental’ to be its intentional directedness (non-representationally understood) then mind-life continuity means locating this intentional directedness in some property of living systems. As Thompson (2007) notes, it is this focus on the existential and phenomenological dimensions shared by the biological and the mental, which differentiates the Jonasian and enactive approach from other life-mind continuity approaches, which focus on shared organizational or functional aspects.

When it comes to accounting for the appearance of teleology and intentionality in particular, the standard recourse since Darwin has been to

natural selection, to explain the development of increasingly complex forms of organization in terms of heritable variation and differential reproductive fitness. Yet Darwin's achievement is often viewed not as making room for intentionality into our naturalistic world view, but rather as eliminating it (Stenmark, 2001). As Dawkins (1986) puts it "natural selection, the blind, unconscious, automatic process which Darwin discovered, and which we now know is the explanation for the existence and apparently purposeful form of all life, has no purpose in mind." (P.5) And as he writes elsewhere, "The universe we observe has precisely the properties we should expect if there is, at bottom, no design, no purpose, no evil and no good, nothing but blind, pitiless indifference." (1995, P.133)

We might attempt to retain a teleological dimension in the tendency towards increased reproductive fitness described by natural selection – for instance by pointing to how, unlike the exceptionless laws of physics, this can 'fail' to be realized in particular individual cases. The same is true of the second law of thermodynamics (Wicken, 1981). In both cases, however, this might be better accounted for by viewing these as statistical generalizations from underlying causal processes, rather than taking them as laws themselves – teleological or otherwise (Matthen & Ariew, 2002). Moreover, even if increased entropy or reproductive fitness were indeed a purpose towards which the universe is being driven, in either case, this could not be credited to the work of individual agents, trying or failing to follow that norm.

Accordingly, bioenactivism has taken a different route, looking instead to the level of the individual organism, and the basic unit of organic life in the single cell. Thus, while its philosophical roots may be in phenomenology, its biological ones are in autopoiesis theory, developed by Humberto Maturana together with Francisco Varela, which proposed to identify the basic logic of

this self-production in the single cell, in order to then formulate its essential features in the more general notion of autonomy.

1.2.0. Autopoiesis and Autonomy

For Maturana and Varela (1973/1980) an autopoietic system is defined as such:

An autopoietic machine is a machine organized (defined as a unity) as a network of processes of production (transformation and destruction) of components that produces components which:

(i) through their interactions and transformations continuously regenerate and realize the network of processes (relations) that produce them; and

(ii) constitute it (the machine) as a concrete unity in the space in which they (the components) exist by specifying the topological domain of its realization as a network. (Maturana and Varela 1980, 78–79).

As this is realized in a cell, the relevant aspects are a membrane or boundary separating the interior from the external environment, and, inside this boundary, a network of enzymes and reactants fed by those molecules that are allowed in through the membrane, which either generate further reactants and enzymes, or the components making up the cell's boundary. In this respect, as Boden (1999) describes, autopoiesis can be viewed as an attempt to define the metabolism of a cell in organizational terms, thereby allowing us to abstract away from particular chemical components.

Because of the use of cellular autopoiesis as an illustrative example, what I call 'bioenactivism' here, is often referred to as 'autopoietic enactivism'. While this is perhaps a better label in so far as it flags which particular sorts of biological properties the enactivist is concerned with, this name is, as previously mentioned, misleading, for at least two reasons.

The first reason is that, as Thompson and Di Paolo (2014) point out, the crucial concept in *The Embodied Mind* was not autopoiesis specifically, but rather the more general principle of *autonomy* – of which cellular autopoiesis is related to as an instantiation at the molecular level. Autonomous systems are defined as networks of processes that exhibit *operational closure*: a recurrent organization such that each process in the network both enables another, and depends on another in turn, and *precariousness*: such that were any of these processes to break down, the network as a whole would cease to exist (Thompson and Di Paolo, 2014). It is by means of these properties of recurrent organization and mutual dependence – rather than in molecular reactions and membranes – that a living system distinguishes itself from its environment, as shown in Fig. 1.

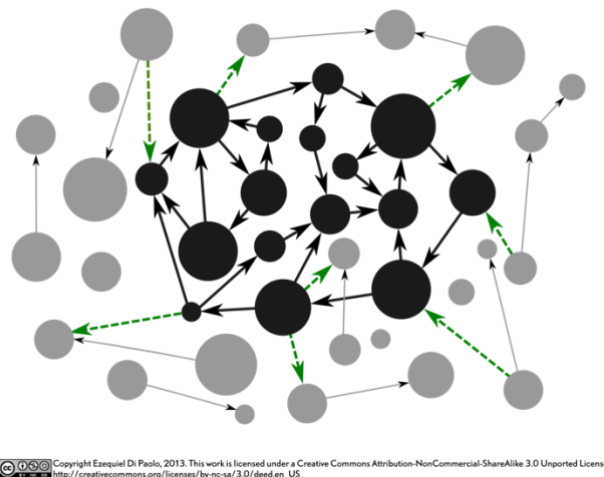


Fig. 1: Illustration of how an operationally closed network of processes (highlighted in black) is distinguished from its surroundings. One-way dependencies of either enablement or dependence alone do not admit membership of said network, and are shown in green (Di Paolo, 2013)

This is then supposed to give us normativity, in so far as the operationally closed organization of a system defines a domain of interactions compatible with the maintenance of this precarious system. Thus, as Thompson (2007) puts it, “Cognition is behavior or conduct in relation to meaning and norms

that the system itself enacts or brings forth on the basis of its autonomy” (P. 158).

Cellular autopoiesis is thus to autonomy rather like what Turing’s (1948) description of his tape-based machine is to computation, and referring to the Enactive Approach as ‘autopoietic enactivism’ is a bit like referring to computational accounts of the mind as ‘paper-and-pen cognitivism.’ It is the notion of autonomy that is supposed to provide a characterization of the logic of living systems that can be generalized across varying levels of organization from the single cell up to the recurrent dynamics of the nervous system.

A second reason the emphasis on autopoiesis can be unhelpful is because, as Villalobos (2013) & Villalobos and Ward (2015) describe, it encourages the muddling together of autopoiesis theory and bioenactivism. While the concepts of autopoiesis and autonomy were both developed in Varela’s work with Maturana, the latter is adamantly not an enactivist – nor would Varela have been understood as such throughout much of their collaboration. Maturana’s goal is not to differentiate living systems from machines, but to explain how a machine could be alive. Thus, as he states in his introduction to their co-authored, *Autopoiesis and Cognition: The Realization of the Living*, in characterizing living systems “notions of purpose, goal, use or function, had to be rejected” (1972/1980 , pxiii).

In this respect, and in the aim to reduce intentional talk to the operations of feedback control systems, autopoiesis theory is cybernetic, not bioenactivist. In the spirit of the British cybernetician W. R. Ashby it endorses a form of life-mind continuity to view the operations of our mind as an elaboration of the homeostatic regulation of simpler organisms. Unlike bioenactivism,

however, it follows this continuity all the way into non-living matter, taking there to be no difference in kind between the existential status of an organism and any other physical system, and thus no basis for attributing intentionality or intrinsic purposiveness as a real property that is uniquely possessed by the latter.

Enactivism itself is not introduced until *The Embodied Mind*, which, as noted, does not yet tackle the problem of naturalising teleology head-on. It is only by the late 90s, where Varela has been influenced by Kant and Jonas' work on the idea of organisms as 'natural purposes' that he begins exploring the concepts of 'original intentionality' and 'sense-making' as unique to life, coming around to the position that these do lead to the re-introduction of a kind of teleology that is "intrinsic to life in action" (quoted from an email exchange in Thompson, 2007, P. 454). This culminates in a 2002 article with Andreas Weber, that draws on Jonas (2001/1966) attempt to naturalize teleology in the 'needful freedom' of a metabolic system.

This 'needful freedom' is intended to emphasize that a metabolic system is not just a machine that can freely persist through a turnover of material components – as implied by Maturana & Varela's earlier claim that "autopoietic systems are homeostatic systems which have their own organization as the variable that they maintain constant" (1972/1980, p80) What makes a metabolic system 'needful', rather than merely 'free', is that it is *dependent* upon this material turnover and its own synthesizing activity for its continued existence. A chair need do nothing at all in order to carry on being a chair, but it is not merely a human projection to say that if a cell's metabolic activity breaks down, then the cell breaks down along with it. Part of what it is to be a cell, part of what it is to be a living thing, is to be something that works towards its own ongoing production through the

continual turnover of molecular material. This, Jonas suggests, gives a purposiveness dimension to the cell's activity, such that we can describe its breakdown as 'failure' even where this is a deterministic outcome inevitably entailed by some prior event.

It is in this definitive break with Maturana's insistence upon treating both living and non-living systems as purposeless mechanisms alike that I would locate the origin of bioenactivism. Two vital questions for its development are: firstly, whether the prevailing definition of autonomy as operational closure among precarious processes captures the needful freedom found in the metabolic cell, and secondly, whether this definition of autonomy is adequate to ground the kind of purposive and intentional attributes needed for an enactive account of cognition. As I will argue, the answer to the first question is negative, and, as a result, bioenactivism has thus failed to provide a positive solution to the second question. This does not mean there is no such solution, however, but only that the bioenactivist need a better formulation of autonomy. Just such an account, as I will argue in Chapter 9, is provided by Montèvil & Mossio (2015) and Mossio & Moreno (2015)'s account in terms of 'constraint closure.'

Before we get to these issues, however, I want to look at a second development in the bioenactive literature, stemming from a different sort of dissatisfaction with how prior formulations of autonomy and autopoiesis are supposed to connect up to enactive norms. This is the fact that, as Di Paolo (2005) notes, the conservation of autopoietic or autonomous network gives us only the all or nothing norm of 'don't die' and the "the rather useless *a posteriori* realization by the external observer that the organism should have avoided that very last encounter that killed it" (P.436). The bioenactivist needs more: namely an account of graded norms that an agent can work

towards and which can guide its sensorimotor interactions towards increasing attunement.

1.2.1 Sense-making and adaptivity

In addition to realizing its own ongoing existence, an organism's autonomous organisation also implies a window of viability – a specific range of environmental conditions outside of which the processes making up the autonomous system will break down. To take cellular autopoiesis: this requires both particular states of affairs (e.g. temperature, pressure) and the ongoing supply of the necessary components to fuel the cell's metabolic processes. As our planet is not a homogenous sphere of lukewarm nutrient soup, so even the simplest of living systems must also adjust to, and interact with, their environment in order to maintain themselves within this window of viability.

While the importance of world-engaging sensorimotor patterns in constituting a subjective perspective was central in *The Embodied Mind* (as reflected in the quote cited in section 1.1) the failure to adequately connect this up to how it serves the autopoietic-autonomous constitution of the bounded individual reflects a conceptual tension that Barandiaran (2017) argues has troubled the Enactive Approach from the beginning – a tension between the organism as separated from the environment, and yet also defined by and dependent upon its interactions with it (Bitbol and Luisi 2005; Bourguine and Stewart 2004).

This tension is explicit in the separate development of sensorimotor enactivism, as a description of the content and structure of perceptual experience in terms of the dynamic relationships between sensory input and motor output. While such an account shares in the broader enactivist

rejection of reconstructivism, it lacks foundation in an account of the norms by which these sensorimotor interactions are coordinated. On the other hand, the concepts of autopoiesis and autonomy in isolation are insufficient to capture the logic of cognition. While they provide us with a naturalistically-founded ‘basic’ normativity, it is only the all-or-nothing imperative of ongoing self-production. We have our success criteria – the preservation of the autonomous network that constitutes the organism – but we need also the criteria that a living system directs its environmental interactions towards that end.

Let’s take a simple example of apparently norm-governed behaviour: bacterial chemotaxis. A favoured case study of Varela’s (1991) and now a mainstay of biogenic approaches to cognition (Lyon, 2006), this describes how a bacterium controls the motion of its flagella in order to move towards higher concentrations of glucose. At a minimum, doing so involves sensing of current glucose concentration, the memory of previous concentration levels, the comparison of the two and the activation of an appropriate motoric response – initiating flagella-rotation to switch from directionless flailing to a directed run when concentration increases. This is a minimal example of ‘intentional’ behaviour without a reconstructive model. There is no internal ‘utility heatmap’ inside the bacterium that represents the distribution of nutrients throughout the current solution, by means of which it plans its journey. There are simply a series of sensorimotor connections of the form – ‘if an increase in sugar concentration is detected, then engage flagella rotation.’

So bacteria swim towards sugar. But before we get carried away and excitedly attribute intentionality, purposiveness, cognition and subjectivity to this process, it should also be pointed out that rocks fall towards the ground –

but no one is inclined to claim that they intend to do so. It is not simply the fact the bacteria reliably swims towards sugar that legitimates the enactivist attribution of intentionality to this behaviour, but the metabolically-underwritten fact that if a bacterium does *not* swim towards sugar, it is unlikely to remain a bacterium much longer. The glucose has ‘significance’ for the bacterium as a nutrient, but this significance is not reducible to the physical properties of the glucose alone. This significance can only be understood in terms of the relationship between the bacterium-self and the sugar-world. It is this organism-environment relationship that brings forth, enacts or constitutes a phenomenal world of significance valence – a world that, as Merleau-Ponty (1963) argued, is irreducible to either the flagella-rotating action of our independent subject or the chemical properties of the metaphysically independent object.

For the sense-making, or enaction, that constitutes this phenomenological, relational world then, we need both autonomy, and the capacity to interact in service of its preservation. Di Paolo (2005) addresses this latter requirement under the heading of ‘adaptivity’ and defines it as such:

A system’s capacity, in some circumstances, to regulate its states and its relation to the environment with the result that, if the states are sufficiently close to the boundary of viability,

1. *Tendencies are distinguished and acted upon depending on whether the states will approach or recede from the boundary and, as a consequence,*
2. *Tendencies of the first kind are moved closer to or transformed into tendencies of the second (Di Paolo, 2005, P. 438)¹*

¹Thompson chooses instead to modify the notion of autonomy to incorporate adaptivity within it. Nonetheless there is agreement that is at least not derivable from the example of autopoiesis alone.

Organisms not only engage in autonomous self-production, they also monitor and regulate their internal states and environmental interactions in order to remain in conditions where this self-production is possible. The similarity between adaptivity, and the notion of homeostasis, is immediately apparent. Yet, just as autopoiesis implies more than metabolism, so too does adaptivity involve more than mere homeostasis. For while homeostasis refers to the preservation of ‘key variables’, such as body temperature, within particular bounds, the conservation of adaptivity is about the conservation of overall organization. Homeostasis is derivative of this more fundamental requirement. Secondly, while examples of homeostasis typically focus on internal regulatory processes that act directly on these variables – such as osmoregulation, Di Paolo’s definition of adaptivity explicitly emphasizes regulation of the organism’s relation to the environment also – that is regulation via the kind of extended agent-environment loops found in bacterial chemotaxis.

Moreover, as Di Paolo develops his account, adaptivity does not merely involve the activation of particular physiological or sensorimotor processes, to regulate internal states, but also the regulation of these processes themselves in response to environmental change. This goes beyond consistent chemotaxis to one particular nutrient source, a distinction nicely illustrated by the example of the Lac-Operon mechanism in *E. coli* (Jacob and Monod, 1961). As Di Paolo, Buhrmann, & Barandiaran (2017) describe:

Under normal conditions E. coli metabolizes glucose. But when availability of this sugar is low while another (lactose) is abundant, certain normally inactive genes will be expressed that enable a new metabolic pathway allowing for the processing of the new sugar. In effect, the bacterium detects a change in environmental conditions that jeopardizes its self-maintenance and reacts by modifying the internal processes underlying its self-construction. It is easier here to distinguish the regulation of behaviour from its normal execution, as the normally dormant genes that are activated contingent on specific environmental conditions do not take part in the ongoing self-sustaining processes of the organism. (P. 130)

So, without the graded normativity of adaptivity, there is no sense-making, no phenomenal world of affective forces to draw and repel us between better or worse modes of interaction with the world. As such, in their book *Sensorimotor Life*, Di Paolo, Buhrmann & Barandiaran (2017) take adaptivity, alongside the self-individuation of autopoiesis, as a necessary requirement for the attribution of agency and subjectivity to living systems.

There is still, however, something missing. Or at least something that was glossed over in Di Paolo's (2005) discussion of adaptivity as regulation of responses, as he moves from the regulation of adaptive processes to talk of the 'experience-dependent' discovery of new ones during individual development. While bacteria such as *E.coli* meet the criteria of regulating their adaptive responses, in order to regulate the variables at the necessary levels to preserve autonomous organization in turn, there is no evidence that an individual bacterium ever *learns* to do so².

Natural selection over multiple generations has developed various hardwired adaptive responses – not only to directly threatening changes to homeostatic variables themselves but also anticipatory adjustments in response to signs, such as predator footprints, that act as a proxy for an imminent threat to

² While Tagkopoulos, I., Liu, Y. C., & Tavazoie, S. (2008) for instance, claim to have demonstrated the ability of *E.Coli* to 'learn' through conditioning that rise in temperature would result in a drop in oxygen levels, and to adjust their metabolism in response, their experiment only showed the phylogenetic evolution of this behaviour across hundreds of generations.

those variables. Such mechanisms, if not always successful, remain relevant across the variety of environments the individuals of a particular species may have to contend with. But this lesson had to be learnt the hard way, at the level of the species, by means of the breakdown in autonomy for any organization that failed to successfully implement them.

Neither the ability to switch between different, genetically-encoded, responses, nor the ability to activate these in anticipation of a non-immanent threat, amounts to the experience-dependent learning of new ones. How can the individual organism learn that a novel signal, say the rising smell of sulfur, threatens its continued existence, without this coming, as it were, one lesson too late? How does an individual learn new adaptive responses to such novel threats without actually experiencing their own breakdown? And, the question we need to get to eventually how does the learning of new adaptive responses ever become the catching of baseballs, the dancing of tangoes, the arranging of flowers, or the scaling of mountains?

In *Sensorimotor Life* Di Paolo et al. (2017) draw upon Piaget's detailed account of sensorimotor equilibration to describe how a system assimilates new environmental possibilities and accommodates these via alterations in its sensorimotor organization. Yet, as they remark.

It should be noted that at this stage the dynamical systems approach to sensorimotor equilibration is not a fully developed theory. It outlines the essential elements that such a theory will eventually have to contain, but several details, for example, regarding its possible implementations, have yet to be filled in. Progress in this area will need to involve further work on the nature of open-ended learning: for instance, additional examination of the processes assumed to be open-ended in nature (such as biological evolution and the dynamics of immune networks) and their relation to processes that could be operating in the brain (see, e.g., Fernando et al. 2012; Watson and Szathmáry 2016) and in the non-neural body. (P.105-106)

This is the cue for predictive processing [PP] to enter the scene, as a possible implementation story for how we learn to coordinate our sensorimotor interactions with respect to adaptivity and, ultimately, the preservation of our autonomous selves.

1.3. From bioenactivism to predictive processing

As a number of authors have argued, PP's account of the brain as a system for minimizing prediction error via perception and action across multiple timescales looks like the perfect implementation story for how the brain coordinates the enactive process of sensorimotor learning that Di Paolo, Buhrmann and Barandiaran (2017) describe (Vázquez, 2020; Nave et al., 2020; Harvey, 2018; Bruineberg, Kiverstein, & Rietveld, 2018; Kirchhoff, & Robertson, 2018; Clark, 2015). That predictive processing has already been used to describe a variety of 'higher' cognitive activities, such as dreaming (Windt, 2018), action-planning (Pezzulo, 2017), memory (Henson and Gagnepain 2010) and counterfactual reasoning about others mental states (Palmer, Seth & Hohwy, 2015) offers some reason to be optimistic that it might provide a means to connect an enactivist account of sensorimotor equilibration up to these 'representation-hungry' forms of offline cognition, traditionally argued to be inapproachable from a purely enactive perspective (Roelofs, 2018; Matthen, 2014; Clark & Toribio, 1994).

The enactivist does not only need to address these externally-characterized cognitive capacities, however. More importantly for my interests is the need to explain how a basic intentionality towards biological survival grows into our rich phenomenological world, permeated as it is by all sorts of socially

and self-constructed norms. As I argue in Nave (2021), when approached from this opposite direction, predictive processing accords strikingly well with Husserl's account of our phenomenal world in terms of a hierarchy of anticipation and fulfilment (Husserl; 2001/1920; 1997/1907; 1982/1913) – see also Yoshimi (2016) and Madary (2016) for comparisons of Husserlian phenomenology and predictive processing. As such, PPs anticipatory language seems well suited to provide a neutral language of behaviour that, as Merleau-Ponty sought, “is neutral with respect to the classical distinctions between the 'mental' and the 'physiological'” (1963, P. 4).

The bioenactivists are not the only ones who would stand to gain from this partnership, however. As I will describe in the next chapter, there has been some controversy over how PP should be interpreted, specifically regarding the function that it enables us to perform. Jakob Hohwy has been the most prominent defender of a reconstructive account, which I shall argue fails, like all representation-first accounts fail, to properly recognize the constitutive interdependence of perception and action. This interdependence, as I argue in Chapter 2, is not only an enactivist concern but a necessary outcome of the predictive processing characterization of cognition. In contrast, Andy Clark (2015) advances a more enactive characterization that better incorporates this interdependence, in which the function of a PP system is to support sensorimotor coordination by, “delivering a grip on the *patterns that matter* for the *interactions that matter* (P. 19).

This is enactive, but not *bioenactive*. Just as Noë & O'Regan's (2001) sensorimotor theory lacked a grounding for the normative dimensions of 'grip', 'attunement', or 'sensorimotor mastery' that they attribute to our interactions with the world, so Clark (2013; 2015) does not offer up an account of why some interactions do *matter*, and what goals the PP system

coordinates action with respect to. This, I argue, is not a problem that can be delegated. PP's interpretation of our neural operations in terms of anticipatory content cannot just be read off the brain's structure and dynamics. Like any account of content, representational or no, it is dependent on norms of functionality, such that we can talk of these operations in terms of their succeeding or failing (Millikan, 1984; Hurley, 1998). This redescription of PP in enactive terms merely rejects Hohwy's proposed 'reconstruction-function' but does not provide the needed alternative. For this, PP needs a bioenactive account of the origins of intentionality in biological autonomy, just as much as the bioenactivist needs (something like) PP to scale this protointentionality up to the level of our more complex cognitive operations.

But before we get to the disputed terrain of content and function, let's set some ground rules. Where did the predictive processing theory of cognition come from and what are the structural constraints that both the reconstructivist and the enactivist, can agree on as entailed by (though not sufficient for) the claim that a system is a predictive processor?

2. Predictive Processing

2.1. Minimal Predictive Processing

The introduction of predictive processing [PP] in the philosophy of mind and cognitive science can be traced to the work of Jakob Hohwy (2013) and Andy Clark (2013, 2016), both of whom pick up Karl Friston's (2003, 2005, 2010) proposal for how predictive coding might be used to provide a generalized theory of brain functioning (2005).

While predictive coding is the most basic component of PP it is not distinct to PP – indeed its origins are not in neuroscience but data compression, as a strategy developed as a means for the storage and transmission of image and video files, during the 1950s (Clark, 2016, see Shi & Sun, 1999 & Musmann, 1979 for an overview). The basic idea is that there are typically regular patterns in the data that we wish to store or send, thus, rather than encoding the value of each pixel individually, we can encode an image more efficiently by only encoding this pattern and its occasional violations. In a video, for instance, large areas of background often remain unchanged over some duration, so rather than re-transmitting the entire scene anew for each frame, we can simply transmit the pattern once, then only encode the subsequent 'errors' induced by local movements of objects and agents in front of this background.

The proposal that the brain uses such a coding strategy, as Sprevak (2021) describes, dates back at least to Attneave (1954), and Barlow (1961), who argued that bottlenecks in the early visual system: for instance, the number of neurons, their dynamic range, limitations on firing rate, and the metabolic

costs of firing, require the brain to use such ‘redundancy reducing’ code for the transmission of sensory data (see also Zhaoping (2006) for a review of some relevant constraints).

This already introduces some minimal requirements upon the structure of the brain. Namely a distinction between prediction neurons (sometimes given the theoretically-overloaded name of ‘representation neurons’) and a comparator or prediction-error neuron, with signals flowing both ‘downwards’ and ‘upwards’ between these. A second source of input to the comparator neurons is from incoming sensory signals, which are compared to the downward prediction signal. The signals from the prediction neurons continue to change until they match the sensory input, signalling that it has been effectively ‘predicted’ (Keller & Mrsic-Flogel, 2018).

Hierarchical predictive processing adds to this a couple of specific proposals about how predictive coding is implemented in the brain, with consequent requirements upon the architecture of a system that could qualify as a predictive processor.

These are:

1) Hierarchy: This process is repeated at various levels, where the input to one level is the state of the level directly below, bottoming out in the sensory periphery.

2) Precision-weighting: predictions and prediction errors are assigned a relative weight, corresponding to the inverse variance of the signal, which determines the influence a prediction error has in changing a prediction.

The hierarchical aspect means that only the bottom level is concerned with matching the sensory signal directly, with each ascending level being driven by the prediction of regularities over increasing spatiotemporal scales. In such a hierarchy there are many degrees of freedom as to which prediction neurons should adjust in order to match incoming signals, thus the role of precision weighting is to determine where this adjustment happens – namely in those neurons where the prediction error signal has a high precision-weighting relative to that of the prediction neuron.

The origins of this account are typically traced to Rao and Ballard (1999), who showed that taking this as a model of the visual cortex predicted a variety of known neural responses, such as end-stopping, not attributable to classic receptive field effects alone. A proposal which Friston (2005) extends to the whole cortex to show how it accounts for a variety of further empirical predictions concerning anatomy and synaptic plasticity; electrophysiological effects, such as mismatch negativity; and psychophysical ones, such as global precedence and priming.

This is not yet predictive processing that you may be familiar with as, “The emerging unifying vision of the brain as an organ of prediction using a hierarchical generative model.” (Clark, 2013, P.5). Having just read a section on enactive approaches, the missing piece should be obvious – the brain is not just a perceiver, but an actor. It is the addition of action into the predictive processing story that is arguably the main feature that has come to distinguish Friston and co.’s work on general predictive accounts of the brain, under the heading of ‘active inference’ (Friston, 2003, 2010; Brown, Friston & Bestmann. 2011), and it is this that is crucial to Hohwy, Clark, and subsequent philosophers’ discussions of predictive processing. On these accounts, prediction error relative to expected sensory input may not only

drive internal changes to prediction neurons but can, alternatively, drive actions that reduce this error by activating reflex arcs to bring about sensory signals that match these neurons' predictions. Precision-weighting, as the determiner of where error-reducing revision happens, controls whether some error is reduced through action to alter the world and bring incoming signals into line with our prediction, or through altering our predictions to bring these into line with signals from the world.

The general view of action as 'the control of perception' is pre-dated in Powers' (1973) perceptual control theory, and before that in ideomotor theories of action (Lotze, 1852; James, 1890). What is novel about PP, Clark (2013) argues, is the integration of this account of action with a theory of learning and perception, under the overarching goal of long-term prediction error minimization – though such a proposal, as noted earlier, bears interesting similarities with Hurley's (2008) shared circuits model also. Thus, as Brown et al. (2011) put it, the incorporation of action generalizes the PP scheme...

and proposes that exactly the same recursive message-passing operates in the motor system. The only difference is that prediction errors at the lowest level (in the cranial nerve nuclei and spinal cord) are also suppressed by movement, through classical reflex arcs. In this view, descending (cortico-spinal) signals are not motor commands per se but predictions of proprioceptive signals that the peripheral motor system fulfills (2011, P.2)

So predictive processing is specifically the claim that perception, action, learning, and attention are implemented by the brain through predictive coding with precision weighting in a hierarchical model where predictions concern patterns over increasingly coarse spatiotemporal grain. This core is agreed upon by both Hohwy, Clark, and others who have developed and disputed the philosophical and cognitive scientific implications of such a

model of brain functioning (Venter, 2021; Vázquez, 2020; Downey, 2018; Seth, 2014).

Where they disagree, is on exactly what function this predictive processing hierarchy serves, and specifically, on the relative priority given to our two possible error-minimizing strategies of perception or action. For Hohwy, action is placed in the service of uncovering the evidence needed for more accurate perception. For Clark, perception is useful in so far as it serves the ultimate goal of successful coordination of action. In this respect, as we will see, each represents the continuation of a different tradition regarding the nature of cognition. Howhy as a representative of the ‘reconstructivist’ branch seeks to understand how the brain infers distal causal structure from impoverished sensory information, Clark, in the cybernetic and, loosely, enactive tradition cares less about such reconstruction relative to PPs utility as an explanation of how we can learn to coordinate our actions over multiple timescales.

2.2. Reconstructivist Predictive Processing

As described, early papers on predictive processing, such as Rao and Ballard (1999) and Friston (2005) focussed on the model’s empirical validity and efficient coding motivations. But it takes more than the unification of some physiological effects under a biologically plausible data compression strategy to get a philosopher out of bed in the morning. The aspects of PP that have arguably drawn the most attention beyond neuroscience are instead the potential epistemological consequences suggested by Friston’s (2005) claim that it provides an implementation of Bayesian inference, a proposal typically

linked to German physiologist, Herman von Helmholtz's (1962/1866, 1867) theory of the brain as an engine of 'unconscious inference'.

Helmholtz saw his work as providing validation for Kant's constructivist account of experience, via discussion of optical principles that reveal the underdetermination of a perceptual experience by sensory stimulation alone. For instance, in Fig. 2 we immediately see the left-hand side as convex, and the right as concave, even though the image alone is ambiguous. This perceptual judgement is thus argued to be dependent upon the unconscious workings of an implicit assumption that light comes from above. Another example is our ability to unconsciously discount the variability of the projection of an object onto our retina, thanks to changes in illumination and distance, in order to continuously view it as being a fixed size and colour.

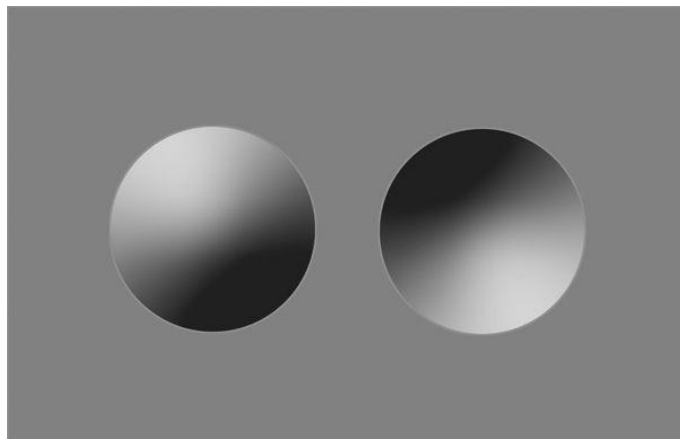


Fig. 2: Illustration of the 'light from above' prior illusion of concavity or convexity (Chunharas & Ramachandran, 2016).

Being by its nature *unconscious* our reliance on these implicit background beliefs, or priors is easy to ignore until we're confronted by an instance in which they go wrong, as in Fig. 3 where we face a conflict between the expectation that footprints are typically concave with our prior belief that light comes from above. In such a case, you should be able to shift between different weightings of these two prior beliefs to view the image as alternately

concave or convex – a case of one’s experience shifting even while the stimulus remains constant.

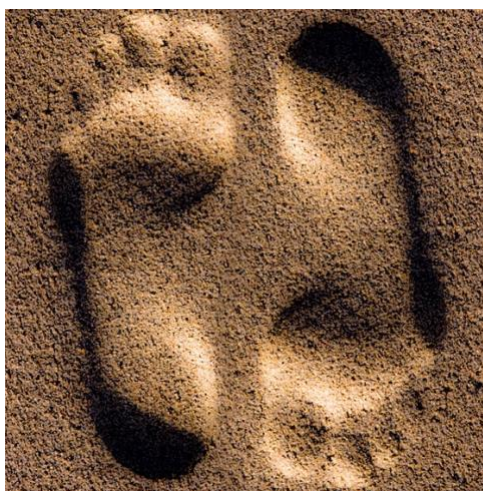


Fig. 3: Footprints illusion (Hobkirk, 2011)

This idea of the brain as an unconscious inference engine came to inspire Bayesian models in perceptual psychology (Gregory, 1980), and visual processing (Lee & Mumford, 2003) (see Yuille & Kersten 2006 & Rescorla 2016 for reviews), rational constructivist accounts of learning and development (Lake et al. 2017; Gopnik, 2012; Spelke, 2009; Tenenbaum et al., 2006) and, in machine learning, the attempt at modelling more tractable approximation strategies for how this inferential process might be implementable in the brain (Dayan et al. , 1995).

It’s interesting to note that all the above take ‘inference’ to be effectively synonymous with Bayesian methods. As Dayan et al. describe their proposal of ‘the Helmholtz Machine’:

Following Helmholtz, we view the human perceptual system as a statistical inference engine whose function is to infer the probable causes of sensory input. We show that a device of this kind can learn how to perform these inferences without requiring a teacher to label each sensory input vector with its underlying causes. (P. 1)

Yet Helmholtz's concept of inference was associationist, not Bayesian. (Westheimer, 2008). Like Kant, he sought to combine a recognition of the inescapable contribution of the cognizer in constructing the objects of experience with an 'objectively valid' basis for these constructive or inferential processes, such that we could be said to have knowledge of the objects they present to us. As Hatfield (1984) describes, while Helmholtz criticized the validity of Kant's account of these a priori 'laws of thought', particularly in terms of the contingency of Euclidean space, he struggled to provide an adequate alternative that might confer validity on our inferential processes – see Hatfield (1990) for an in-depth comparison of Kant and Helmholtz' projects. One way of interpreting the framing of unconscious inference as a specifically *Bayesian* process is to take Bayes rule as providing just such a universally valid transcendental law of thought. The question for the reconstructivist interpretation of PP is how much accord with such a law guarantees regarding the convergence of our beliefs upon an accurate representation of some mind-independent reality

While Helmholtz' concern is with the validity of perception, he did not ignore action. As he describes “We are not leaving ourselves passively open to the [sensory] impressions intruding upon us, rather we observe, that is, we bring our organs into those conditions under which the impressions can be most precisely distinguished” (Helmholtz 1867, P 438, quoted in Hohwy, 2013). And as Hohwy describes this analogy of action to experiment and exploration:

Perceptual inference allows the system to minimize prediction error and thus favour one hypothesis. On the basis of this hypothesis the system can predict how the sensory input would change, were the hypothesis correct. That is, it can test the veracity of the hypothesis by testing through agency whether the input really changes in the predicted ways. The way to do this is to stop updating the hypothesis for a while, and instead wait for action to make the input to fit the hypothesis. If this fails to happen, then the system must reconsider and eventually adopt a different or revised hypothesis. (P. 79)

If perceptions are hypotheses, then it is natural to take actions as hypothesis testing. It is via action that our perceptions of the world collide with a reality that can push back, and through action that we can gather new evidence with which to update our hypotheses through Bayesian inference.

2.2.1. Predictive processing as unconscious inference

The interpretation of predictive processing suggested by Hohwy (2013) is essentially a continuation of these ideas of unconscious inference, perceptions as hypotheses and action as hypothesis testing. On this view prediction neurons do not only encode a prediction of the signal that will be received from the next level below, they also represent the distal causes responsible for producing this particular pattern in the sensory stream. The key idea is that increasing temporal depth in the patterns predicted corresponds to increasing depth in a hierarchy of distal causes. For instance, by observing changing light levels, we can not only track the circadian cycle of light to dark, corresponding to the Earth's rotation but also a second-order, slower regularity in how this first cycle lengthens and shrinks over the course of a year. In tracking this second pattern, so the story goes, we latch on to another distal cause behind our sensory stimulation in how the position of the earth relative to the sun changes in an annual cycle.

As Hohwy (2013) puts it:

“Regularities can be ordered hierarchically, from faster to slower. Levels in the hierarchy can be connected such that certain slow regularities, at higher levels, pertain to relevant lower level, faster regularities (for example, slow regularities about Aussie rules footy word frequency during the yearly news cycle pertain to faster regularities about the words I end up reading; if I know the slower regularity then I am less surprised by the occurrence of those words). A complete such hierarchy would reveal the causal structure and depth of the world—the way causes interact and nest with each other across spatiotemporal scales.” (P. 28)

Thus, by hierarchical prediction error minimization, the brain not only latches onto regularities over multiple timescales, but, in doing so, comes to encode a model of the hierarchical causal structure of our distal environment. According to this reconstructivist PP account [RPP], it is the rich structural content of this internal model, not the comparatively impoverished data presently streaming through the retina, which directly determines my perceptual experience. The latter is demoted to the role of model constraint, suggesting the description of perceptual experience as a process of ‘controlled hallucination.’ Among other things, such an account explains why, when I look at the building site across the road, my experience is not a two-dimensional array of rectangular slices and flattened figures - despite this being all the information that my retina is receiving. The grey squares and silver lines match the prediction of a large, three dimensional, building and so it is this building model, not the retinal activity constraining it, that I experience.

What guarantee do we have that this prediction-error-based control is robust enough to bring our internal model meaningfully in line with reality? So far we’ve only talked of prediction error minimization, but in PP these error signals are not transmitted in raw form, but always with a ‘precision-weighting’ reflecting their estimated reliability, which directs where the error-squashing adjustment takes place. It is this precision-weighting that determines whether errors are accommodated at lower levels – eg. being

explained away as the ordinary background fluctuations caused by the noise of the shower – or whether they force deeper adjustments to long-term regularities encoded higher in the generative model, as when a persistent error signal gathers in estimated reliability, eventually triggering the realisation that the words to the first section of Carl Orff’s famous cantata, *O Fortuna* have nothing to do with a passionate desire for tinned fish.

It is precision-weighting that guards the predictive agent against being swept to and fro by each random fluctuation in the sensory stream, or against becoming entrenched in the determined commitment to some particular pattern. Interpreted as a probabilistic measure of the reliability of a signal, it is a key ingredient in the reconstructivist interpretation of PP, allowing the error-minimisation process to be cast as one of approximate Bayesian updating, in which the reliability of prior regularities learnt over the course of the agent’s experiential history is weighed against estimated reliability of the current evidence (more on this in the next chapter).

Does adjusting our generative model in accordance with Bayes’ rule thus guarantee eventual convergence between its structure and that of the distal environment? Even in advancing the reconstructivist view of predictive processing, Jakob Hohwy is rather pessimistic on this point, noting that as successful minimisation of error is only achieved in relation to our pre-selected hypotheses, reconstructivist PP becomes:

“An affirmation of simple Cartesian skepticism. Since we cannot obtain an independent view of our position in the world, we cannot exclude the skeptical hypothesis that the sensory input we receive is caused by an evil, hoaxing scientist rather than the external states of affairs we normally believe in. The Bayesian framework thus entails scepticism.” (2016, p265)

Within this inferentialist understanding of PP, the fact that our generative model successfully predicts current sensory input becomes evidence that our model accurately captures the structure of this sensory input's causes. This, despite the fact that alternative models might also have predicted it just as successfully. These circular patterns of evidence form what Hohwy describes as an 'evidentiary boundary,' a point of separation between the hypothesis-generating mechanisms and the evidence that is being explained. The boundary proposed here is the edge of the sensorium: on the inside, the skull-bound brain – on the outside, the body and world.

There are two, separate but related, issues proponents of 4E approaches might dispute about Hohwy's internalist characterization of PP here. There is our current question of how to characterize the relationship between the systems on either side of the boundary, and there is also the question of where the boundary between a PP system and its environment is fixed – if it is indeed fixed at all. This second question is an issue for the extended mind theorist not for the enactivist, whose more immediate concern is with the relation between mind and world, rather than with where we draw the partition between them.

As such, rejecting Hohwy's RPP in favour of an Embodied, or Enactivist account, does not, as Bruineberg, Kiverstein and Rietveld (2018) emphasise, require rejecting that there is some meaningful boundary between a PP system and its environment, nor that this boundary *may* be drawn at the sensorimotor interface. Boundaries have, after all, been at the core of enactivist definitions of cognition from the very beginning (Varela, Thompson & Rosch, 1991). Rather, as we will see in the next section, the question is whether said boundary is best characterised as in terms of

evidentiary seclusion, or of ongoing coupling, between the internal dynamics of the PP system and its environment.

Hohwy does not always keep these questions distinct, however, and thus mischaracterizes the response from proponents of enactive cognitive science to this dilemma as the suggestion that incorporating ‘world-engaging’ action might rescue us from scepticism by breaking us out of an evidentiary boundary in order to gain ‘direct’ access to the distal environment – a proposal he rejects. In his reconstructive story action is placed firmly in the service of perception: a process of hypothesis-testing that allows us to intervene to control relevant variables, to seek out further evidence to resolve uncertainty, and to confirm or disconfirm our current model. Such actions (as ecological psychologists are fond of pointing out) do indeed boost our epistemic resources enough to resolve the kind of local ambiguities engineerable with the artificial constraints of the 2D images used in psychophysics labs (Orlandi, 2014). But Hohwy is correct that they could not free an RPP agent from the global underdetermination of sceptical scenarios. As he puts it:

An appeal to action, on the prediction error scheme, reduces to an appeal to inferences about different kinds of patterns of sensory input. If a mad scientist was a hidden common cause of all that sensory input we would have no way of knowing unless she made an independent causal contribution to sensory input. (Hohwy, 2013, P.220)

Indeed, as I’ll describe in the next chapter, when framed in terms of the broader framework of free energy minimization, we see that rather than helping to pull ourselves out of the sceptical pit, the possibility of acting to alter our evidence stream only digs us in even deeper. Before we get to this, however, it’s worth asking why we ever decided to jump in there in the first place.

2.3. Sensorimotor Predictive Processing

If predictive processing is supposed to be an account of how we develop an accurate representation of our distal environment, then, as we have seen, it is not a particularly reassuring one. While this interpretation of predictive processing as encoding a causal model of the external environment is common in the non-philosophical discussion also (eg. Kanai et al. 2015) none of the central components of the PP model, described in section 2.1, entail such a view. As Orlandi (2018) argues, all we actually have is a hierarchical structure with relations of inhibition between neurons, where this inhibition either cancels out an incoming signal, or it does not. In the latter case, that that signal propagates up to cause a change in a ‘higher-level’ neuron.

So where does the representationalist reading come from? One source, as we’ve just seen, is a background commitment to the view of knowledge and cognition as the formation of some internal mental state that corresponds with an independent external one. Another, as Anderson and Chemero (2013) argue may be the common fallacy of deriving semantic conclusions from merely correlational properties. There are as, they note, two senses in which we can talk of ‘prediction’ at play in the discussion of predictive processing:

The first sense of “prediction” (henceforth prediction1) is closely allied with the notion of correlation, as when we commonly say that the value of one variable “predicts” another (height predicts weight; education predicts income, etc.). Prediction1 is essentially model-free, and comes down to simple relationships between numbers. The second sense of “prediction” (prediction2), in contrast, is allied instead with abductive inference and hypothesis testing. Prediction2 involves such cognitively sophisticated moves as inferring the (hidden) causes of our current observations, and using that hypothesis to predict future observations, both as we passively monitor and actively intervene in the world. It is theory laden and model-rich. (P.24)

In minimal predictive processing, the only ‘predictive’ relationship we have is the first of these – the tendency of prediction neurons to correlate with incoming signals, in virtue of which we interpret the difference between them as an ‘error-signal’ that is being reduced. This is the same sense of prediction by which lightning predicts thunder or the position of one coupled pendulum predicts the position of another. Such relationships of covariation may be useful to someone trying to infer the state or structure of a hidden process, but they are not inferential in themselves.

Predictive processing in itself is only a matter of developing correlations between neural activity and patterns of stimulation across a hierarchy of different temporal scales. The reconstructivist interpretation of this is extrinsic to the PP architecture and motivated by a prior belief that reconstruction is what we are after. We need not be disappointed then, by the inability of PP to deliver a guaranteed reconstruction, for we are not obliged to seek one in the first place

As such, rather than taking action as a solution to RPP’s sceptical challenge, Clark (2015) dismisses such a challenge altogether. Following embodied and enactive approaches, he suggests that the solution to the problem of perception does not lie solely in the supplementation of our inferential resources with action, but crucially in rejecting the very characterisation of

our perceptual goals as reconstructive in the first place. As he cites Varela, Thompson and Rosch on this point.

The overall concern of an enactive approach to perception is not to determine how some perceiver-independent world is to be recovered; it is, rather, to determine the common principles or lawful linkages between sensory and motor systems that explain how action can be perceptually-guided in a perceiver-dependent world (Varela et al, 1991, p173)

If this is our goal, then why *should* the PP system have to infer the hidden causal nexus beyond the sensorimotor interface? To take a well-worn example, Chapman (1968) shows that a baseballer outfielder need not first model the entire onward trajectory of the baseball relative to their position, and to the field, in order to then begin the act of moving to catch it. All that is needed is the ongoing coordination strategy of ‘Optical Acceleration Cancellation’ – that is moving such that the ball stays at a stable position in the retinal field until it is close enough to catch. An outfielder requires no internal physics engine to recruit this strategy, no knowledge of the aerodynamic equations governing the flight of round, slightly irregular, projectile in a mild north-westerly wind. All they require is an understanding of the lawlike relations between their motor output and the position of the projection on their retina.

In terms of sensorimotor PP, Clark (2015) explains, this becomes a matter of assigning high-precision weighting to errors related to the prediction that the optical projection of the ball remains at a stable location in the visual field. In such a way the rest of the system’s actions are recruited around the quashing of this particular error signal, to the neglect of most else happening on the field, until the desired state of catching the ball (or the undesired state of colliding with a teammate employing the same strategy) is reached. Here there is no prior process of tinkering at the generative model until a lack of

overall error provides adequate comfort that we've formed an accurate representation of the external world and action may now begin. Rather successful action is itself the ongoing control of a small portion of the sensory flux within those constraints that the system predicts will lead towards its target state. When understood in such a way, as Seth (2015) suggests, we can interpret this “non-reconstructivist” approach to PP as offering a mechanistic rendition of earlier embodied sensorimotor theories of perception (O'Regan & Noe 2001).

Such ‘fast and frugal’ strategies are much more suited towards the ongoing guidance of an organism that must constantly keep afloat in a fast-changing environment. They also fit smoothly within the rules of Bayesian optimality, as Fitzgerald et al. (2014) note, the ideal Bayesian system seeks not only to maximize predictive accuracy but also to minimize the complexity of the models recruited to do so.

For Clark (2015) the availability of locally-effective non-representational strategies is not an argument that we should abandon all talk of models and representations, however. Rather the strength of sensorimotor PP is the offer of, “a systematic way of combining deep, model-based flexibility with the use of multiple, fast, efficient, environmentally exploitative, routes to action and response”(2015, P.18). In order for a PP system to effectively deploy such ‘fast and frugal’ strategies as OAC, it must also be able to monitor slower-changing contextual factors (such as whether one is actually engaged in a game of baseball, or merely a participant) in order to ascertain when the circumstances are ripe for their deployment. This is why the PP system requires hierarchical depth, such that high-level states may target these large-scale increasingly invariant patterns throughout the fast fluctuations of the sensory stream.

Unlike on the reconstructivist view of PP, these high-level action-oriented representations do not allow us to ‘throw away the world’ when we engage in planning our next action but rather to coordinate our interactions with the world over multiple timescales. We track not agent-neutral causes, but agent-relative affordances at nested levels of spatiotemporal grain, from the affordance of ‘playing baseball’ to that of ‘catching this particular ball.’

Unlike in reconstructivist PP, securing the correctness of these ‘action-oriented representations’ does not depend upon the ability to reject the sceptical hypotheses. If a current affordance for ball-catching action is correctly detected, then deploying OAC will guide the evolution of the skilled outfielder’s sensorimotor interactions to the target ball-in-hand state. This model of current sensorimotor contingencies will be successful irrespective of whether the hidden causes interacting with our sensorimotor array are instantiated by mischievous demons, curious scientists, or strange and charming fundamental particles.

2.4. What's the point of predictive processing?

If the reconstructivist version of PP burdened the generative model with extravagant commitments, we may now be concerned that the action-oriented spin has been overly-economical in response. For once PP is freed from the imperative of reconstruction, we still require an alternative motivation for its operations. Action is not an end in itself. To say that our predictive models are ‘action-oriented’ and to attempt to explain perceptual contents in terms of the intention of these actions, as Hurley (1998, Ch 6)

argues, merely moves the problem back a step. So, what determines the appropriateness of an action and its criteria of success?

One option is to propose that this need not be dealt with by the PP system itself, which merely describes the mechanism for achieving some function, not the function itself. That normative issue could be delegated to some separate ‘desire module’ responsible for the calculation of an agent’s goals and intentions. These can then be simply fed in as priors to a PP system tasked with bringing about their execution.

We could propose this, but to do so would undermine the entire point of the PP framework as an explanation, not just of the application of predictive models, but also of their ongoing development. While PP may not itself provide us with a story about the ‘first priors’ by which the predictive process gets started, its central explanatory payoff is as an account of the ongoing modification of these constraints through the minimisation of error generated by interactions at the sensorimotor interface. If desires and intended action policies make up our prior predictions, then their selection and satisfaction conditions must be intertwined with the overall predictive economy. Thus Clark (2013) speaks approvingly of the suggestion that...

" generally, personal and hedonic value) is not simply a kind of add-on, implemented by what Gershman and Daw (2011, p. 296) describe as a "segregated representation of probability and utility in the brain." Instead, it seems likely that we represent the very events over which probabilities become defined in ways that ultimately fold in their personal, affective, and hedonic significance." (2013, P. 200)

While Clark (here, and in Clark, 2019) endorses this blurring together of the affective and the cognitive, he does not offer a positive story about the origin of this normativity within a PP agent. As such, and despite his sympathies

with the enactive approach, the action-oriented interpretation of PP developed in Clark (2015) meets the criteria of being enactivist only in the weaker sense that the term has been applied to a variety of positions concerned with accounting for our mental life in terms of extended patterns of body-world interaction, rather than skull-bounded symbol manipulation.

While this broader use of the ‘enactivist’ label is helpful in grouping together the shared orientation of a diversity of work in visual perception (O’Regan and Noë, 2001) anti-representationalism (Hutto and Myin, 2013) and emotion (Colombetti, 2005), it can also lead to the conflation of these narrower, cognitive scientific efforts with the more metaphysically oriented ‘enactive approach’, first introduced in Varela, Thompson and Rosch’s (1991) *The Embodied Mind*, which proposes thoroughgoing revision to our understanding of the mind-world relationship. This revision places a naturalised account of the emergence of normativity at the foreground of the cognitive scientist’s explanatory task.

As discussed in the previous section, the central motive of the enactive approach’s revisionist metaphysics concerns the replacement of the recapitulationist understanding of intentionality and meaning with a teleological one. For the enactivist, the intentionality that is the ‘mark of the mental’ is understood not in terms of a relationship between a representational vehicle and the object it is ‘about’, but rather, in a sense much closer to its meaning in the phenomenological tradition, as the directedness of an action towards the satisfaction of some goal.

The bioenactivist method for naturalizing this normative assessment of our actions was to seek a grounding for this normativity in the biological processes of life, then to argue that the same logic scales up to the level of

the cognitive. That, as 20th Century phenomenologist and a progenitor of the enactivist approach, Hans Jonas put it, “the organic even in its lowest forms prefigures mind, and [...] mind even on its highest reaches remains part of the organic” (1966, P1)

In the next section, I will zoom out from predictive processing specifically, to introduce the free energy principle [FEP], and its corresponding explanatory framework. As a principle, the FEP proposes to “unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; avoid surprises and you will last longer” (Friston, 2012, p.2). Predictive processing then stands as one possible architecture that could implement this free energy minimization process over multiple timescales. Thus if the imperative to minimize surprises does indeed capture the kind of intentionality that bioenactivism ascribes to the autonomous organism, then the PP theorist would gain a grounding for attributing norms of successful action to the predictive brain, which it regulates its activities with respect to. In turn, PP, can provide the bioenactivist with a means to scale up the basic intentionality of an autonomous organism, exemplified in behaviours such as bacterial chemotaxis, up to the rich counterfactual structure of human cognition and consciousness.

Unfortunately, this will not work. Surprisal-minimization falls far short of any notion of autonomy robust enough to ground the attribution of intentionality and teleological orientation to a system. Even more unfortunately, the unusually sprawling and heterogeneous structure of the free energy framework can, at times, seem constructed to hide this fact. It will take some time to see exactly why its formulation of autonomy fails. Bear with me.

3. A non-mathematical introduction to the free energy framework

*And if he knows where he was standing
When J.F.K was shot
Chances are though time's passed him by
He's still standing within yards of that spot*
- Let's go with the flow - The Beautiful South

Since its proposal in Friston (2003, 2005) and Friston et al. (2006) the free energy framework [FEF] has developed and mutated across hundreds of publications under a diverse³ menagerie of co-authors. As such, a range of subtly different formulations abounds in the literature. Most are dense in mathematical formalization, while the ordinary language accompaniments – eg. that free energy is just, “an information theory measure that bounds or limits (by being greater than) the surprise on sampling some data, given a generative model” (Friston, 2010), are liable to leave those uninitiated into variational calculus, probability theory, and other arcane arts, deeply unsatisfied.

Amongst all of this two, in principle independent, components can be distinguished: a mathematical account of action and perception as subsumed under the approximation process of variational Bayesian inference, and a ‘first principles’ analysis of life as the process of avoiding surprising events – together with some connective tissue tying the execution of the former to the achievement of the latter.

³ In terms of disciplinary background that is

The first component is the most unobjectionable. It is often discussed separately under the heading of ‘active inference’ and used to model a variety of systems and phenomena, from neural dynamics (Friston et al., 2017; Da Costa, et al., 2021) to cognitive and behavioural phenomena (Parr & Friston, 2017; Friston et al., 2016) to social coordination (Friston et al., 2020; Constant et al., 2019), self-organization (Friston, 2013, 2018) and even the climate (Rubin, 2020). The general idea is that the dynamics of all of these systems can be described in terms of an invariant joint probability distribution over possible sets of variables, typically referred to as a ‘generative model’.

What adds a spark of controversy to all of this is that advocates of the FEF do not only take active inference as a useful framework that we, as observers or scientists, can use to model some target system, but often adopt a realist stance to view the target system itself as actually encoding this generative model, using it to perform inference and to direct ‘actions’ that are consistent with it (for critical reviews of instrumentalism vs realism in the FEF, see Bruineberg et al., 2022; Andrews, 2021; Van Es, 2021 & Van Es & Hipólito, 2020).

The connection between PP and active inference is relatively straightforward. Where active inference is a computational description of approximate inference in a system that can act to change incoming evidence, PP is one possible architecture that could implement this computational function. Like PP, active inference has often been proposed in Helmholtzian terms as a story of how the brain might approximately infer the most accurate model for the hidden causes of its sensory stimulation (Friston, 2005; Friston, Kilner & Harrison, 2006). More recently, however, Friston and co-authors

have tended to instead adopt a cybernetic framing, akin to Clark's (2015) sensorimotor PP, in which active inference is presented as a description of how the brain learns to control its sensorimotor engagements in order to preserve a self-model (Seth, 2015; Pezzulo, Rigoli & Friston, 2015, 2018; Ramstead, Kirchoff & Friston, 2018).

The utility of active inference as a modelling framework can be accepted in isolation from the more controversial free energy *principle* [FEP] which, “aims to unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; avoid surprises and you will last longer,” described as a principle so basic that, “there is no need to recourse to any other principles” (Friston et al. 2012, P. 2). Surprisal minimizing, in this context, is just another way of describing a system with time-invariant statistical properties, such that it can be described by a fixed ‘generative model.’.

As part of the attempt to subsume or formalize bioenactivist concepts like autopoiesis in terms of this generative model, the FEP's advocates propose that this imperative of surprisal minimization can supply the missing normativity needed to explain the intentional content of action and perception (Hohwy, 2020; Schwartenbeck et al. 2013). As Ramstead, Badcock & Friston (2017) state:

The ‘intentionality’ or ‘aboutness’ of living systems—that is, the directedness of the organism towards a meaningful world of significance and valence—emerges as a natural consequence of embedded adaptive systems that satisfy the constraints of the free energy formulation. For a living thing to be intentional just means that it entails a generative model... Put simply: active systems are alive if, and only if, there [sic] active inference entails a generative model. This makes the generative model of central importance to the free energy formulation, since it defines the form of life that an organism is seen to enact. Ramstead, Friston, Badcock (2017, supp mat 4, p33)

While the claim that minimizing surprisal is a fundamental imperative for a particular system is not required to justify modelling that system using the active inference framework, it is a crucial step in attributing these inferential processes to the system itself. Without the principle-based claim about the function of a system, the claim that it is a predictive-processor or free energy-minimizer can be based only on formal similarities between its stable dynamics and the computational process of approximate Bayesian inference. As will be familiar from triviality arguments against the computational theory of mind such structural similarities are worth little on their own and are insufficient to attribute internal models and anticipatory contents (Sprevak, 2018). If the FEP *does* adequately formalize autonomy, however, this would both allow active inference to be proposed as a theory of what our brains (and bodies) are actually doing, and it would help the bioenactivist in turn by connecting this basic intentionality up to a theory of sensorimotor learning and coordination that might support the formulation of enactive theories of higher cognitive processes.

In this chapter, I will describe the basic components of how active inference redescribes perception and action as a process of inference and reduces this to the maintenance of stability in the face of perturbation. In the next I will describe the justification for the free energy principle's claim that this maintenance of stability is the defining feature of living systems, and in the following chapter I will explain how advocates of the principle attempt to extend this to neurally-deprived organisms, by using the idea of conditional independence between internal and external states to formalize the notion of a sensorimotor boundary. This gives us a two-part definition of a system in terms of two forms of statistical stability: 1) the stability of a probability density over most likely states of the various parts of a system and 2) the

stability of interactions between these parts that preserves a statistical boundary between the system and its environment

All of this sets us up for the subsequent two chapters in which I will analyse the prospects of Friston's (2019b) 'existential dyad' against two yardsticks. The first is how well these criteria allow us to capture the bioenactivist concept of autonomy. The second, is how well they capture the essential features of living systems, independently of any prior commitment to the bioenactivist picture. I will argue that it fails at both. Such requirements are both too general to capture the specific features that imbue a living system's activities with an intentional orientation and, at the same time, too restrictive to shadow the amorphous organism in its improvisational dance between the formal and material constraints that shackle inanimate substances.

3.1. Variational inference + action = active inference

The first component of the inferentialist part of the Free Energy story is a development of variational inference⁴ – a strategy taken from machine learning for approximating a solution to an intractable inference problem (Hinton & van Camp, 1993; Neal & Hinton, 1998, and see Beal, 2003 for a more recent development and overview)

In the case of the brain, we might think of this task in terms of finding the probability distribution over hidden causes that best predicts our sensory observations⁵. The optimal way to do this would be to begin with a joint

⁴ Also known as variational Bayes, or ensemble learning

⁵ As discussed with regards to PP, these causes of changes in the sensory stream involve

distribution $[P(O,H)]$ over the probability of each observation $[O]$ and each possible value of the hidden variable(s) $[H]$ that we can decompose into the *prior* $[P(h)]$, the *likelihood* $[P(o|h)]$ and *marginal* $[P(o)]$. These are the three components of Bayes formula.⁶

$$Eq.1: P(h|o) = P(o|h)P(h)/P(o)$$

We can then use this formula to progressively update our probability distribution over hidden causes on the basis of each new observation, to get a new posterior $P(H|o)$ which we then feed in as our new prior $P(H)$ for the next round of observations O_{t+1} and updates. Eventually, all being well, this process will stabilize, the posterior given O_{t+n} will be the same as the prior from the previous update. This can then be taken as the ‘true generative model’ that describes the actual statistical properties of the process generating our evidence. Sadly, calculating the marginal $P(O)$ requires integrating $P(o|H)$ across all the weighted possibilities of H in order to figure out how likely some particular observation is ‘in general’. Such an operation quickly becomes intractable with anything like the number of possible hypotheses the brain has to deal with.

So, rather than attempting to infer the true posterior distribution from the entire space of possibility distributions, variational inference operates by selecting from a constrained set of simpler distributions that can be characterized by a small number of parameters. The selected distribution is termed ‘the recognition model,’⁷ and its parameters can then be progressively

the brain’s own actions. For now, we’re ignoring the possibility for the system to intervene on the process that generates its evidence, in order to get through the basics of variational inference as an approximation technique

⁶ Capital letters denote variables, lowercase letters refer to particular values of these

⁷ Also called the recognition density, or variational density. Density is simply the term for a function that generates a continuous probability distribution.

tweaked to perform gradient descent/ascent over two elements. The first is *accuracy*, the likelihood of our evidence under the recognition model, which is what we want to maximize by changing the parameters of said model. The second *complexity*: quantifies the amount that we are changing the recognition model, which we want to minimize.⁸

Intuitively, these two quantities reflect what is called the bias-variance trade-off (Geman, Bienenstock & Doursat, 1992), between (over)fitting our model to each new variation versus (over)generalizing from prior regularities. Balancing these quantities prevents us from just placing all of our confidence on the single hypothesis that makes o most likely, which would maximize accuracy, or digging in our heels with the model we have already learned from previous observations to minimize complexity. These are exactly the two considerations that are weighed against each other under true Bayesian inference. By attempting to minimize these with respect to a selected distribution, we have thus converted our intractable inference problem to a much simpler optimization one.

The rather powerful sounding ‘free energy’ (Friston, 2005, 2010) that gives the FEF its name is simply an alternative label for the combination of these two functions⁹. When encountered in machine learning or statistics, it would more likely go by the name of the (negative) ‘Evidence Lower Bound [ELBO]’ It can be written, with tolerance for a little simplification, like this:

$$Eq. 2: F = D(q(H) || p(H)) - \langle \ln p(O | H) \rangle_q$$

⁸ Thus is measured using Kullback-Leibler divergence - a measure of the difference between two probability distributions

⁹ Not to be confused with thermodynamic free energy!

The first term, $D(q(H) || p(H))$, is the divergence between the recognition distribution $q(H)$ and our prior distribution $p(H)$ over hidden states ‘ H ’ – this is the additional complexity that adopting some particular parameters for our recognition distribution would introduce to our model. The second, $\langle \ln p(O|H) \rangle_q$, is how probable some recognition distribution over hypotheses makes our observations O – the accuracy¹⁰. This has been converted to a logarithmic scale, so, rather than ranging between 0-1, the highest value possible value for a certain outcome is 0 and all other probabilities have negative values – the lower the more unlikely. Two negatives make a positive, so subtracting a very low negative quantity is equivalent to adding a very high positive one. Thus the higher the value of accuracy (up to the maximum of 0), the more probable the observation given our recognition distribution, the smaller the amount added into our free energy equation.

Hence, by minimizing complexity and maximizing accuracy, free energy is reduced.

I want to stress that ‘accuracy’ here refers not to how close our recognition model is to the true generative model that we’re aiming to approximate (which does not appear in this decomposition) but only to how successful it is in predicting our observations – what might be better qualified as *predictive* accuracy. Like a broken clock or a lucky beginner, a model may have predictive accuracy for some particular observation while still being divergent from the true statistical properties of the process which produced that observation. With enough observations, however, the hope is that persistently high predictive accuracy serves as evidence for the fit between our model and the hidden variables that cause our observations.

¹⁰ The ‘ $\langle \rangle$ ’ symbols denote that we’re taking the weighted average for the probability of o with respect to the recognition density $q(H)$ called the ‘expectation’ or ‘expected value’.

What's neat about Free Energy is that it can also be broken down differently into 1) the difference between the recognition density and the actual posterior, sometimes referred to as *divergence* plus 2) the unlikeliness of some sensory data relative to the 'true' generative model that we are aiming to infer, termed *surprisal*.¹¹ Which looks (again with a little simplification) like this:

$$\text{Eq. 3: } F = D(q(H) || p(H|o)) - \ln p(o)$$

A slightly odd formal nicety worth highlighting here is that the minus sign in Eq. 3 is a part of the surprisal term – we are adding surprisal as a *negative log probability*, not subtracting it as we did with accuracy in Eq. 2. Because log probabilities (other than 0) always take negative values, so a negative log probability is a positive number. While this sounds like rather pointless sign-shuffling, it allows us to say that reducing free energy requires 'maximizing accuracy', or 'minimizing surprisal' – which sounds a lot more natural than saying we want to 'minimize accuracy' or 'maximize surprisal'! Still, in both cases all we are doing is trying to get as close to 0 as possible – the difference between adding a positive value that has been minimized towards 0, or subtracting a negative value that has been maximized towards 0 is much of a muchness.

A second, and more important, point is that while (predictive) accuracy and surprisal seem similar they are not just the inverse of each other. Accuracy here is a function of an observation and our recognition model – our current best guess at the approximate statistics of the underlying process that generated this observation. It is something the free energy minimizer is able to measure and reduce. In contrast, surprisal (as it occurs in the FEP) is

¹¹ A neologism coined to distinguish the notion of unlikeliness, from the personal-level notion of surprise

improbability under the true posterior of the correct generative model – which is what we are aiming to approximate. High accuracy entails low surprisal only when our approximate recognition model is sufficiently close to the generative model.

Somewhat confusingly then, it is what has been termed ‘accuracy’ here, not surprisal, that best corresponds to what we would ordinarily call the ‘surprisingness’ of some particular piece of evidence in relation to some inferring agent’s expectations. The surprisal of some observation is inaccessible to the agent for two reasons. The first, already discussed, is because it depends on an intractable marginalization of the generative model $P(O,H)$ that variational inference is a means to avoid. A second, and less commonly appreciated, reason is because (as we’ll get to in a moment) it is not clear that the brain, predictive processor or otherwise, should be interpreted as encoding this generative model at all.

So, while the division into ‘divergence + surprisal’ of Eq. 3 is not accessible to the free energy minimizing brain, what dividing things up this way shows us is that by minimising free energy the brain can limit the possible divergence between its recognition density and the true posterior – to the extent that, if the latter were in the restricted class of simple distributions, they will become equivalent. When there is no divergence between our recognition density and the true posterior, then free energy reduces to surprisal – allowing surprisal minimization with a true model to be cast as a special case of free energy minimization. As it is, however, all the free energy minimizer can rely on is that minimizing free energy will push its recognition density towards the closest thing it can get to the true posterior from within a constrained class of simpler, tractable distributions.

How does all this connect to predictive processing? Well, the variational inference procedure does not specify which constraints we should select for our restricted class of probability distributions. But, if we choose the recognition distribution appropriately, then said procedure may be implemented by a predictive processor. The first assumption we need is the *Laplace assumption*, which restricts our class of possible recognition distributions to Gaussians, which can be parameterized with just the sufficient statistics of mean and variance. The second is the *mean-field approximation*, which assumes that the posterior distribution across all the states of all hidden variables can be decomposed into a number of separate distributions across the state of each independently. With these choices, each prediction in PP can be interpreted as the mean of one of these distributions, the precision as encoding its variance, and the overall precision-weighted prediction error of the input as the free energy for our present evidence relative to said distribution (Gershman, 2019) .

This is how we can map recognition densities, or models, onto a predictive processing architecture, but it takes two to infer. So where is the generative model?

3.2. Generative model hunting

Variational inference avoids the need for an intractable marginalization of the generative model but it does not remove the need for it altogether. This joint distribution $[P(O,H)]$ is still required to obtain the likelihood of an observation given our recognition model the ‘accuracy’ $[\langle \ln p(O|H) \rangle_q]$ and so to calculate its free energy using the ‘accessible’ rearrangement given in Eq. 2. Thus it is standard in both early treatments of the FEF (Friston,

2002, 2003, 2005, 2010) and throughout the PP literature (Kiefer and Hohwy, 2019), to assume that the brain is explicitly encoding and operating with both a recognition and a generative model. In PP specifically, the latter is described as the ‘inverse’ of the former, “implemented by ‘backwards’ (top-down) synaptic connections.” (Kiefer and Hohwy, 2018)

The problem with this, as Hohwy (2016) notes, is that variational inference only ever brings the recognition density closer to the ideal posterior for the particular generative model that we are using. As it is also the generative model that determines the evidential impact of an observation then, if our generative model is off-track, ongoing observation will not be able to knock our recognition model into shape. Variational inference doesn’t just require any old generative model, but one that is close enough to the actual statistical structure of the process being inferred to allow for successful inference about its current state to take place.

As such models of free-energy minimizing systems tend (as I have also been doing thus far) to presume not only *a* generative model, but the *right* generative model (Friston et al. 2017; Hesp et al. 2021; Buckley et al. 2017). As such Raja et al (2021) accuse the FEF of a vicious circularity in presuming we have a good generative model to explain how we learn through perception, where this model itself would require some process of perceptual learning to explain its development in the first place.

In light of this, Friston and co. have more recently altered course, to argue that the idea of the generative distribution as something represented in the brain rests on a failure to distinguish between different senses of ‘model’ at play in the framework (Ramstead, Kirchoff, and Friston, 2019). On this account, the only model actually encoded in the hierarchical structure of a

predictive brain is the simplified recognition model. Generative models, in contrast, are something directly 'entailed' by the generative process itself. As they explain, "Entailment', in this setting, is used to emphasise that a generative model is necessary to define the recognition model but does not have sufficient statistics that are physically realised" (P. 233).

The claim seems to be that the generative model is just an aspect of the equations we use to characterize the generative process as a target for potential inference by the (approximately) Bayesian process of updating an internal recognition model. In this interpretation, it is not a model realized by the neural architecture of a predictive brain, in the same sense that an orrery realizes a model of the solar system, but rather what Thomson-Jones (2012) calls a 'non-concrete model', like the differential equations of the Lotka-Volterra model or the sine functions of Haken-Kelso-Bunz. Such mathematical models may be utilized to describe a diversity of target systems, but need not be realized in a physical vehicle.

This interpretation is, however, hard to square with claims elsewhere in the same paper that the generative model is something that has 'causal bite' (P.233) that it is a 'control system' that an organism 'uses to guide action' (P. 234). A mathematical function does not itself 'cause' anything, it can only describe causal events. Such causal talk fits better with a contradictory explanation of the generative model, presented in the same paper, that, "The generative model does not encode anything. It is realised by the statistical relations between states of interest." (P. 235) Here, then, the generative model is described not as a non-concrete mathematical model, but rather as something that is physically realized, and that controls the behaviour of the inferring system. What distinguishes the generative model from the recognition model, on this presentation, is not the lack of physical realization

but the locus of it. The states that the generative model is encoded by are just those things that we had previously been taking it to be a *representation of*, namely the web of causes that determines a system's observations. Here then, the generative process is no longer cast as the target that the generative model represents, but rather the vehicle of the generative model itself. The generative process and the generative model are cast as different descriptions of one and the same thing. It is trivial and pointless to attribute a representational relationship on the basis of something's identity with itself. If the generative model is realised by and *identical with* the generative process, not a *representation of* the generative process, then what *is* it a model of?

The answer, Ramstead et al. suggest is nothing at all.

There is no warrant, mathematically, for the claim that the generative model encodes semantic content or structural information. The generative model manifests as a control system that uses exploitable structural similarities encoded in the internal states of the organism. It is not itself a representation. (P. 235)

In other words, the generative model, on this understanding, differs from the recognition model not only in being realised outside the brain, but in being a non-representational model. The idea of non-representational models is not too unusual, Andrews (2021) claims, pointing out various accounts of non-representational models that might be drawn on in support, such as Weisberg (2013).

Still, we seem thus to be facing a confusion then about whether Ramstead et al. (2019) are claiming that the generative model is a non-concrete mathematical model that is *not encoded by* anything, or whether it *is* encoded by a control system but *does not represent* any particular target. If it is to play some role in our analysis of the inferring brain, then it must be one or the other. So, which is it?

Well, the only criteria that Friston (2019b) requires for the ‘entailment’ of a generative model is that the statistical tendencies of a system’s behaviour (both in terms of states of variables and dependencies between them) are stable over time. As he puts it, “a model is just an ergodic system” (P.183) This is somewhat misleading, for ergodicity is a stronger requirement than the stability of statistical properties over some duration, and neither is trivial, but temporarily putting that aside until Chapter 4, it is also not obvious why it would justify claiming that some system is a model – targetless or otherwise. Having ergodic dynamics allows a system to be *described by* a particular kind of mathematical model – namely a fixed joint-probability distribution – but if all there were to being a model was just to admit of a description then everything we can talk about is a model, rendering talk of changes as being ‘caused by a model’ or as being a ‘model-update’ trivial.

This pattern of taking things to literally be the mathematical models that describe them is one that recurs in the FEF and we will return to it in Chapter 5, where I will suggest that it stems from an idiosyncratic metaphysical standpoint that is not explicitly elaborated or defended in Friston’s work. For now, I think it best to sidestep talk of mathematical models as physically instantiated by what they describe and causally efficacious. At a minimum a model is a tool, used for some purpose. Its function may not necessarily be the representation of a specific target, but it must have *some* function. The label of ‘Ergodic’, however, merely describes the dynamics of the system generating our observations, it does not confer any functional role upon it.

Thus I agree with van Es (2020) that we should take the generative model to be understood in the first sense, as a non-concrete mathematical model that could, in principle, be used by us as external observers to describe the

statistics of the process generating an agent's observations. The function this model would serve is to allow us to interpret the agent as performing approximate inference of said process via free energy minimization – and thus to interpret its internal states as encoding a recognition model. What has causal force, however, is the structure of the system that the generative model describes, *not* the description itself. A diving gannet may well be describable by a differential equation, but only one of them catches a fish.

Either way, whether the system that generates our evidence literally realizes a generative model, or is merely describable by one, the important point is that the generative model under the FEF is not something that the brain (or other inferring system) encodes or expects. Whether the generative model is taken to be literally realized by the dynamics of the process generating a system's observations *or* to be the correct description of the statistics of that process stipulated by an external observer who is using this knowledge to interpret the system as inferring these statistical properties – either way the point is that it is stipulated to be isomorphic to the generative process. Any convergence between the generative model and the generative process is not explained by the FEF but rather presumed by it.

This has important implications in understanding how the FEF builds upon the variational inference formalism.

Firstly, for free energy minimization to get off the ground as an explanation of how a system learns about the hidden causes that generate its observations, the dynamics of this generative process must such that they could either realize, or be accurately described by, a generative model. Where the generative model is understood as a joint probability distribution, this means those hidden causes must have stable dynamics, such they can be

described in terms of means, variances, and statistical dependence relations that are invariant over the duration for which an inference takes place. The assumption that our systems of interest are of such a kind is often left implicit in the background of the FEF. Nonetheless, it is crucial, for without this there can be no stable generative model to serve as a target that the recognition model tends towards approximating.

Secondly, this changes how surprisal should be understood. If surprisal is unlikeliness relative to the generative model and *if* we accept that the generative model is not encoded by the inferring system, then surprisal under the FEF is not, as it has been sometimes described in the PP literature, equivalent to “sub-personal prediction error” as both (Madary, 2012) and (Clark, 2013) propose in order to distinguish it from personal-level surprise. The surprisal of an event is not something encoded in the brain at all, it is simply the long-run frequency with which the generative process produces said event (Fiorillo, 2010). That something tends to move from a higher to a lower surprisal state is just another way of saying that it regresses towards a stable mean – as may be described in a generative model. Surprisal minimization is thus not an inferential process in any sense more interesting, than the one in which everything is a model.

So, the only model that the FEF should be interpreted as taking the inferring system to internally encode is the recognition model. Whether it makes sense to talk of a system as performing inference over representations depends on whether we can ascribe the encoding of a recognition model to it and whether the relationship between this and the generative process (described by us in the generative model) meets the criteria for being a representation.

3.3. Where does this take us, representationally speaking?

Variational inference is the ascending shoot of the FEF, reaching up to ‘higher-level’ cognitive processes of reasoning, planning and learning. Predictive processing, as an algorithmic implementation, provides a potential answer to the question ‘how *could* a neural network, like a brain, perform inference?’ Still, all we have established so far is that the formal properties of recognition density optimization can be mapped onto the neural dynamics that PP proposes to identify in the brain. Merely mapping the syntax of (approximate) inference onto a physical system does not prove that the brain *does* engage in inference – any more than Chalmers (1995) and Putnam (1975) proved that everything is computing (Sprevak, 2018).

There are more things zipping and diffusing through our skull than are dreamt of in the ontology of free energy minimization. To establish a legitimacy for the Bayesian in your neural tissue that the finite state automata in your breakfast cereal lacks we need to show two things: a reason to privilege the structures picked out by PP, and specificity in how the constructs of variational inference map on to these.

The natural way to secure the former is to show that the dynamical structures singled out by PP are those that explain how the brain performs its function. But, on pain of circularity, we cannot just take the possibility of describing the brain as if it is performing inference to support the claim that inference is the function of the brain, and thus that this is the *right* description to explain how the brain performs its function. Such an argument becomes even weaker once we see, as we will in Chapter 5, the ease with which the FEF allows us to trivially map a formal description of approximate inference

onto a variety of non-neural systems whose function, if any, it becomes increasingly implausible to regard as inferential.

Even if we were provided with an argument that the variational formalism is a legitimate description of what the brain is really doing, the ‘models’ that would be entailed by this potential computational-mechanistic description are just the means and variances of Gaussian distributions – they are not models *of* anything in particular. The interpretation of these Gaussians as probability distributions over the states of some hypothetical hidden causes is assumed in variational inference, not explained by it. Thus, as Wanja Wiese (2017) argues, what this description of the brain would establish directly is only what Francis Egan (2014) terms ‘mathematical contents.’ These may constrain the ascription of possible semantic contents but the computational-mechanistic description alone does not suffice to establish their legitimacy as anything more than an instrumental ‘gloss’ on the system’s operations.

In an attempt to substantiate the ascription of such semantic content, Kiefer and Hohwy (2018, 2019) have pointed to the ‘divergence’ component of the FEF, measuring the difference between our recognition density and the true posterior density of the generative model¹². This, so they argue, looks a lot like a metric of misrepresentation – the possibility of which is often proposed as among the necessary criteria for the attribution of a representation in cognitive science (Millikan, 1984; Dretske, 1993 Ramsey, 2007; Shea, 2018).

¹² There are actually two divergences making up misrepresentation in Hohwy and Kiefer’s argument, as they distinguish between a generative model encoded in the brain and the true statistics of the generative process. If we follow the arguments of the previous section, however, then there is only the true generative model and only one relevant divergence term between this and the recognition model.

The problem is that, as Kirchhoff and Robertson (2018) note, all this divergence actually measures is the covariance between states of variables, in terms of the probability distributions over those variables. All sorts of things may become increasingly correlated with time. Over the years, a mattress might start to take on the shape of the person who sleeps there, but the purpose of sleep is not to create a likeness of oneself in springs and foam. Sagging springs are no more a representation of an absent partner than a dry riverbank is of the water that once flowed through it.

Correlation, like a promise of commitment, comes cheap. It may be a precondition for some structure to be able to play a representational role, but (dis)correlation only becomes (mis)representation if we have independent reason to believe that said correlation is deployed to representational ends. That is, we need to show that the structure is either a constituent of, or used by, a system that is aiming (and thus, potentially, failing) to represent (Millikan, 1984; Dretske, 1993 Ramsey, 2007). With this, we're borne back once again to the issue of function. Only once we've explained the purpose of prediction-error minimization can we answer the question of whether its structures, such as a putative recognition model, play a representational role in service of this.

The mere fact that prediction error minimization *can* be interpreted as a process of variational inference fails to answer the question either way. But the free energy principle is not *just* sparkling variational inference. The incorporation of action is a significant addition to the story told above, from which grows the other strand of the FEF: the attempt to root inference in more 'basic' processes of sensorimotor coordination and homeostasis. It is this side of the FEF where people have looked either to derive a representational sub-function in support of such a process, or to show that

the theory should prompt us to do away with representation altogether. So let's get to the action.

3.4. Incorporating Action

Through variational inference, we can tease out the assumptions and constraints that motivate the description of predictive processing as an approximation of Bayesian inference. This is nice, but if this were all there is to the Free Energy Framework, it would not be all that novel – the suggestion that variational methods could explicate how the brain performs ‘unconscious inference’ dates back to 1995 and Dayan, Hinton, and Abbott’s proposal of the ‘Helmholtz Machine’. Nor is it what we went looking for when we turned to the Free Energy Framework in the attempt to provide an alternative, non-reconstructive, imperative for the operations of a PP system.

While there is no role for action in standard variational inference, this does not mean that its role in the FEF is only as an optional add on. While it is possible (though biologically irrelevant) to describe predictive processing in entirely passive terms, the free energy framework constitutively depends on the integration of action in order to distinguish it from mere variational inference. This part of the FEF, the theory of **active inference**, describes the move from an agent that passively infers regularities in its evidence stream to one that can interfere to change the evidence it receives.

Like variational inference, active inference is not specific to PP. It is a computational description of a process that may be implemented by a wide range of architectures. It has also, more recently, been extended to describe how a system might infer the future consequences of potential actions,

allowing it to anticipatorily evaluate and select policies with respect to how they will minimize what has been termed ‘expected free energy’ over a longer-term sequence of observations – rather than just the free energy of the agent’s current state. Such models may be, and typically are, detached from philosophical concerns about the fundamental nature of life and cognition. Their utility in solving tasks and predicting behaviour does not stand or fall upon the Friston and co.’s aspiration for free energy minimization to serve as a ‘first principle’ for living systems. Even if the assumptions needed for free energy minimization have only limited application, active inference models could still be useful.

The free energy *principle*, however, rests on the claim that active inference is not just an occasionally useful modelling tool, but a redescription of autopoiesis/autonomy. If this is the case, so the idea goes, then these more complex models might be able to do the work of scaling up the kind of basic intentionality that autonomy grants to perception and action (both re-described in terms of simple error-minimization) to encapsulate higher-order cognitive processes of reasoning and deliberation.

But before we talk about scaling-up active inference, we need to investigate the solidity of its foundations in describing the autonomy of organisms. These can be separated into two pillars. One pillar, which will have to wait until Chapter 5, uses the tools of graphical modelling to extend active inference beyond neurally-equipped creatures with a pre-given sensorimotor boundary, in order to describe how such boundaries supposedly emerge from the dynamics of active inference itself. The other pillar argues that the dynamics that active inference prescribes are exactly those in virtue of which something would qualify as a self-preserving autonomous system. It is this pillar that is relevant to the normative problem introduced in the discussion

of PP, for if active inference is autonomy, and PP is an algorithm for active inference, then we have our answer to what the prediction error minimization is actually for: namely, self-preservation.

To see how this works, we'll now return to the problem of action as introduced in the section on Helmholtzian vs sensorimotor interpretations of PP, now aided with the additional elements gained from the preceding sections: 1) clarification of the nature of the generative model and 2) the various decompositions of free energy. Once we have used these to properly frame action under the FEF, we can look to more formal treatments of active inference in the context of graphical models, to see how proponents of the FEF attempt to generalize this theory beyond a predictive brain to describe all living systems.

3.5. Generative processes and active systems

Suppose, you had heard Edinburgh often referred to as “the Athens of the North”¹³ and, mistaking this to be a matter of climate rather than an unedifying comparison between the unfinished 19th-century folly of Calton Hill and the two millennia-old Greek Acropolis, you moved there as an escape from the Manchester drizzle.¹⁴ Arriving in the city in full expectation of a balmy 21°C, you find yourself confronting persistent error in relation to a particularly dreich day and sub-zero temperatures.

¹³ by VisitScotland officials at the least

¹⁴ Indulge me in imagining yourself somewhat geographically challenged

You could just accept your fate and update your model to incorporate an increased probability of states at the lower end of the temperature scale. Alternatively, you could reduce this (really rather embarrassing) error by heading to the airport and flying to warmer climes. Both would serve to minimize the free energy of your ongoing observations, but actively adjusting what you experience in order to fit a prior hypothesis is hard to make sense of in relation to the aim of forming an accurate model of the external environment. If this were your ultimate goal, then all your error-avoiding behaviour seems to have done is prevent you from learning an important lesson about the untrustworthiness of tourist boards.

This problem is familiar from the earlier discussion of predictive processing but things get even messier if, as suggested in the previous section, we understand the generative model, not as something encoded in an inferring brain, but rather as a description of the actual statistical properties of the observation-generating process that we're trying to infer.

Suppose as you exited Waverley train station you instead found yourself captivated by Edinburgh's Medieval skyline and decide to stay. Unable to tolerate the city's outdoor temperatures, you move into a flat with a castle view and a magnificent old fireplace, then install yourself comfortably beside it. Congratulations! You've minimised prediction error relative to your goal state of 21°C – but this isn't all that you've done. In lighting up the fireplace you have not only changed your current state to bring it in line with your prediction, you've also altered the long-term statistics of the environment that you are attempting to model, making 21°C a more likely state to encounter in future.

Every evening, all across the city, thousands of other residents do the same – though these days central heating is the more common method. As a result, the average outside air temperature of Edinburgh city centre is a good few degrees higher than in the surrounding countryside – and getting warmer (Price, 1979). If our heat-seeking denizens continue to multiply in number and energy consumption, then, climate-change permitting, you may one day be able to emerge blearily from your Georgian terrace onto the streets of New Town, and find your prediction of a 21°C air temperature to be perfectly satisfied. Your internal recognition model of Edinburgh’s average air temperature is now accurate – but only because the very fact of you’re assigning this temperature a high probability drove you to make the world conform with it

Incorporating action into the FEF thus unleashes a strange circularity. A circularity that undermines the Helmholtzian understanding of the process whose statistics the generative model describes (and which the brain is supposed to partially approximate in its recognition model) as being composed solely of distal environmental causes. Allowing agents to act to change their sensory input unavoidably inserts the modelling-agent into the very observation-generating process that this agent is simultaneously attempting to model. The generative process that the generative model describes thus incorporates the actions of the organism as well as its environment, hence Ramstead et al.’s (2018) claim that it is a “consequence of the adaptive behaviour of the organism” (P. 231).

If we are modelling anything under active inference then it is not the agent-independent world, but rather the fused self-world system. It is the formalisation of this that makes the mathematical portion of the FEF a marked departure from the standard variational procedure – in which the

inferring system has no effect on the stable statistics of the system whose behaviour it is attempting to infer.

In discussing the problem of evil demon scenarios for PP, Hohwy (2016)¹⁵ describes the circularity inherent to predictive processing, where our evidence is taken to confirm our prior hypotheses so long as it is consistent with them, irrespective of whether it would rule out other equally probable alternatives. The problem looks much worse when our evidence is not only interpreted, not only selected, but actively *created* as a result of these hypotheses. As Bruineberg, Kiverstein and Rietveld (2018) point out, if the Free Energy Minimizing brain is a hypothesis-testing scientist, then it starts to look like a ‘crooked and fraudulent scientist’ that decides on the outcome of an experiment beforehand... and manipulates the experiment until the desired result is reached.” (p2444)

How then can this crooked scientist picture be reconciled with the subsumption of perception and action under the logic of approximate Bayesian inference that is often supposed to be free energy framework’s crowning glory?

Firstly, it should be noted that the kind of crookedness just described does not undermine a Bayesian interpretation of the recognition model. A good Bayesian is not accountable for the evidence she receives, only what she does with it. In consistently gravitating towards sources of warmth you can curate an evidence stream that, under true Bayesian inference, would produce a model of the average air temperature is 21°C. Still, once there are no pre-

¹⁵ The criticisms of Hohwy’s view pertain to the Helmholtzian treatment given in his 2013 book ‘The Predictive Mind’, and contemporaneous papers. It should be noted that in more recent works, he appears to move away from this approach to an understanding of PP more compatible with the current treatment.

fixed independent statistics to be inferred, and once the observations made are determined by the prior expectations of the modelling system, this inferentialist understanding looks less helpful as a way of interpreting the success of a free energy minimising system in preserving itself.

To ameliorate the uncomfortable subjectivity of setting initial priors in Bayesian inference, it is often pointed out that even when agents begin with wildly divergent priors any differences can eventually get ‘washed out’ through the process of updating on the same evidence, leading to convergence in their models. Agents that create and curate their own personally-tailored evidence stream to support their initial priors undermine this possibility. The incorporation of action appears to have left us rather unmoored. When we can change not just the recognition, but also the generative model (that is the statistical description of our patterns of interactions with the world around us) at will, then we appear to lack any stable target that we might be working towards.

We can try to regain some stable footing here by noting that neither system, neither brain nor environment, is infinitely flexible. We might be able to alter our environmental temperature to some degree, but there are constraints on how much the generative process can be changed to fit an internal expectation. Consider the other side of the world where 5.6 million Singaporeans also struggle with their local climate – though from a more tropical standpoint than Edinburgh’s inhabitants. The city-state has the among the highest number of air conditioners per capita in the world, allowing residents to shuttle between apartment buildings, office blocks and underground malls, chilled to a cool 18°C. Despite these efforts, indeed, partially because of them, Singapore is growing warmer at twice the rate of the rest of the world (Jiang et al, 2021). While the action of the Edinburgh

resident in turning up the thermostat amounts to a self-reinforcing prediction, the Singaporean's air conditioning actions are unable to overcome the constraints of thermodynamics. The very attempt to realise a predicted goal temperature in the short term actively undermines the ability to achieve it over longer timescales.

Bringing predicting agent and predicted world into alignment can no more result from the dogged pursuit of prior expectations in the face of their repeated failure, than it from passive conformity to whatever sensory evidence our environment throws at us. You can't infer a silk purse from a sow's ear – for all that a pigskin wallet might be on the table. Still, the fact that our potential repertoire of prediction-fulfilling actions is partially restricted does not really solve our problem. Our space of possibilities is still underconstrained and, in most situations, the active inferrer will still be presented with a choice between acting to bring the generative process in line with its recognition model, versus changing this internal model to fit the world. In the former mode, we may still think of it as inferring an accurate model – albeit a model that includes the likelihood and consequences of its own actions. In the latter, it looks more like a system that is attempting to sculpt the world in line with its model. How is this trade-off settled?

3.6. Surprisal vs Divergence

We have two ways to minimize free energy or prediction error. We can act to change the world, and so the generative model, or we can update our internal recognition model to better predict our observations, and thus, ideally, to better resemble the generative model. In the context of PP, this is typically expressed as the claim that “perceiving the world (perceptual

inference) and acting on it (active inference)¹⁶ turn out to be two sides of the same coin.” (Gładziejewski, 2016, P. 562). As far as PP is concerned, both action and perception involve changing variables in order to minimize prediction error – the only difference is that the former changes concern external variables and the latter involve only internal changes.

The FEF complicates this, however, by revealing prediction error/free energy to be a compositional quantity – one that may be broken down into a) *divergence*: the difference between our approximate recognition model and the ideal Bayesian’s model of the causes behind our sensory input, plus b) *surprisal*: the true unlikeliness of this particular sensory state for this system. Now the lesson in the discussion of variational inference was that we cannot measure or reduce either of these directly. All we have to work with is a) *predictive accuracy*: how likely some change to our recognition model makes this input, and b) *complexity*: a quantification of the additional complexity introduced by this change. What we do know, however, is that surprisal has nothing to do with our internal recognition model. As such, increasing predictive accuracy by changing the expectations embodied in our recognition model through ‘perceptual inference’ will not reduce surprisal one jot.

Here we find a functional asymmetry between perception and action: they not only alter different variables but, crucially, they minimize different quantities as a result. In perceptual inference, we might adjust our model to diverge less from the true statistics of the environment: the true average temperature of Edinburgh in February. In doing so we better position

¹⁶ While in the PP literature it is common to divide the prediction error minimizing process into ‘active inference’ as opposed to ‘perceptual inference’ in this way, in the FEF the term ‘active inference’ is used to refer to the unifying story that combined both perception and action together in the minimization of variational free energy.

ourselves to more accurately estimate and respond to surprisal in future. It is *action* alone, however, that can move our sensory input back to a lower surprisal (less unlikely) state.

Understanding prediction error as free energy, and thus as a proxy for the minimization of two things, divergence and surprisal, allows us to approach the different interpretations of PP in terms of a disagreement over which of these is primary. The Helmholtzian takes it to be divergence, while the sensorimotor approach prioritizes the minimization of surprisal.

Active inference itself does not settle the case either way but, for the Helmholtzian at least, it provides a further explanatory hurdle. As we saw in section 2.2 the Helmholtzian view analogizes actions to experiments – albeit ones skewed towards finding evidence for a preselected hypothesis. If the reason a free energy minimizer selects action over perceptual learning is specifically in order to minimize surprisal – that is to say if action is a means to prevent the occurrence of previously infrequent events – then the only hypotheses it can be used to harvest evidence for are those that predict that things will stay pretty much the same as they always have.

These turn out to be the only kind of hypotheses that a basic free energy minimizer could effectively model. What a recognition density comes to approximate is the true posterior of a generative model – a joint probability density over sensory and external states. We can define such a density only when the tendencies of, and dependencies between, these states are stable over time. There is then a sense then in which we could say that minimizing surprisal supports the goal of accurate representation: it shapes the world around us into the kind of stable thing that we can represent. But the only

reason this is necessary is because we have arbitrarily constrained our representational repertoire to systems for which surprisal is minimized!

The environment around us does not seem to provide any reason to embrace such a universal constraint. There are oscillations and orbits, but there is also growth, metamorphosis, development and collapse – recurrence not guaranteed. In the grand thermodynamic stream of things, the stability we observe is but a temporary eddy in the overarching flow towards disintegration. Once we have the kind of hierarchy we find in predictive processing we could gain the capacity to represent some of this instability in terms of nested, generative models over different timescales (Badcock et al., 2019), but, given that we are surrounded by both change and stability, why prioritise the assumption of the latter?

If it makes sense to model our environment one in which systems are stable and surprisal is minimized, then it seems we need to look first at how, and why, we act to make our little corner of the world this way. From the perspective of the free energy principle then, we must prioritize the provision of an explanation for surprisal minimization, and the role of action in ensuring it, over an argument for the importance of divergence reduction in subserving representational accuracy. If we do things in this order then the derivative purpose of divergence reduction is much easier to explain. The only way to measure surprisal is through the predictive accuracy delivered by our recognition model, and this will only be a good proxy for surprisal when divergence is reduced.

Approaching predictive processing from the free energy principle helps us to resolve the debate between Hohwy's perception-oriented, Helmholtzian interpretation, of PP and Clark's action-oriented one. It also moves us

another step along the pathway to an answer for the normative question: what is the function of the predictive brain? Perceptual updating is, as Clark (2015) argued, for the coordination of action. The coordination of action, we can now say is for surprisal minimization. But what is surprisal minimization for, and why do we do it?

This is the question at the heart of the Free Energy framework's approach to life and mind. It's proposed answer is that preserving a stable generative model is not just about epistemological convenience – it's a matter of life and death.

4. One Weird Trick to Stay Alive: the FEF's philosophy of life

4.1. The organism's agenda

Before we can even consider whether or not the brain actually encodes some approximate recognition model that it tunes through the process of active inference, we would need to establish whether there is a meaningful, stable target that such a model would even be capable of approximating. That is, we need to establish that the causal process generating its observations has stable statistical properties such that it could be described by a fixed probability distribution, approximate or otherwise. This is only the case if that hidden process generating our observations is a 'surprisal-minimizing' one, which just means that it tends to repeatedly revisit the same small subset of states.

In the previous section we saw how allowing an agent to actively bring about the very states that it predicts complicates our understanding of this generative process. Where surprisal, or free energy, is minimizable through action, so the agent itself becomes part of the 'hidden causes' of its own sensory input. As such, the target that our agent attempts to approximate, and minimize surprisal relative to, should not be interpreted as describing the (supposedly) stable dynamics of some agent-independent state of affairs, (as it was on the Helmholtzian characterization) but rather the dynamics of the whole organism-environment system. Surprisal is not just a matter of

how often something happens ‘in general’, but how often it happens for a particular organism-environment system. What is surprising for the sperm whale will be depressingly familiar to the bowl of petunias.

So why think that the organism-in-its-environment is a stable-surprisal minimizing process, describable by a generative model and approximatable by a recognition one? Well, according to Friston (2012), the process of securing this stability is nothing more than a formal analysis of what it means to be a self-preserving or self-organizing system. “The whole point of the free-energy principle,” as he puts it, “is to unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; *avoid surprises and you will last longer.*”(p. 2117).

So where variational inference with a recognition model is the ascending shoot of the free energy framework, stretching upwards towards a description of inference, action-planning, and other higher-order cognitions, the proposed analysis of self-preservation is where the FEF attempts to root this process in bioenactivist soil, in order to suck biological functions and intentionality up to its cognitive branches. Importantly Friston does not present active inference as either an alternative to autopoiesis and autonomy, or an instrumental means to their end, but as a first-principles analysis of the minimal requirements for existence from which enactivist definitions of life may be derived. So, as Ramstead, Badcock & Friston argue:

Put simply: active systems are alive if, and only if, there [sic] active inference entails a generative model. This makes the generative model of central importance to the free energy formulation, since it defines the form of life that an organism is seen to enact. (2017, supp mat 4, p33)

Why think that maintaining a stable generative model by avoiding situations with high surprisal through active inference is the key feature of a living

system? A bad way to make this argument involves arguing for the simple imperative as something “circular” or “self-evidently true.” (Allen & Friston, 2018, p. 19). In one sense, the claim that all living systems tend to avoid ‘improbable’ states is indeed a tautology of existence: things tend to be in states that they are likely to be in. They tend not to be in states that are not likely. It is straightforward to map this on to viability: a state that is non-viable for an organism is very unlikely indeed. It is undeniably correct, though exceptionally pointless, to claim the continued existence of an organism depends on its avoiding states in which it does not exist.

For our laborious excavations beneath the mathematical infrastructure of the Free Energy Framework to result in nothing more than the gnomonic utterance that ‘everything must be what it is, and cannot be what it is not,’ would be disappointing to say the least. No insight is to be gained by re-describing impossible states as improbable ones, then sagely noting that an organism seeking to continue its own existence must avoid these.

Besides, as a way of staying alive, minimising the surprise of your own death would come one action too late. Taking surprisal to be a binary matter of ‘possible, or not’ would indeed make the claim that it *is* avoided tautologous, but ought implies can. To say that minimizing surprisal is something that organisms *must* actively do implies that they *can* be in some, relatively high-surprisal states such that they can then move away from them back to more likely ones.

As it is presented in the FEF, surprisal is crucially something that comes by degrees, such that it can be minimized. Specifically, as discussed in section 3.2, how much surprisal some state has is determined simply by how frequently it has occurred in the history of our target system and to say an

organism minimizes surprisal is to say it frequently revisits the same small set of states. As Friston (2018) puts it:

We are only interested in one sort of system. These are processes where (the neighborhood of) certain states are re-visited time and time again; for example, the biological rhythms that characterize cardiorespiratory cycles—or the daily routine we enjoy every Monday, on getting up and going to work... This means, on average, I must move toward states I am more likely to occupy. This may sound trivially simple but has enormous implications for the nature of any (interesting) process that possesses an attracting set of states. (P. 2)

So, rather than a tautology derived from first principles about what it means for something to exist, the FEF looks more like a positive proposal about the particular kind of existence living systems have. The empty platitude that an organism is more likely to be in states that are more probable becomes the substantive argument that 1) the organism's states vary and 2) it tends to repeatedly revisit the same set of states it has visited previously, and to avoid those states that it has not previously visited with any great frequency. While the state of the system may constantly fluctuate, the probability distribution over these states remains invariant, and this invariant probability distribution is the true generative model of that organism. This is all there is to saying that an organism minimizes surprisal and thus that it 'entails' a generative model.

To say that an organism's self-preserving behaviour necessarily 'entails' a generative model, in this sense, means only that it will be describable by a joint probability distribution that does not change over time. It does not establish that any part of this organism, brain or otherwise, literally encodes a separate recognition model. The latter claim is a suggestion about how the brain *could* perform inference about the generative process, given that the true generative model would be computationally intractable. But the need to attribute such a model depends on the assumption that organisms need to

explicitly encode and compute with a model of the generate process in order to guide their actions. As discussed in the previous chapter statistical properties of an organism's dynamics alone are not enough to substantiate this.

As such, in papers that focus on the free energy framework as a theory of biological self-organization, rather than on the specific problem of how the brain might actually perform approximate Bayesian inference, there is no role for a separate recognition model (Friston, 2013). Somewhat deceptively, however, Friston and co. continue to speak in terms of 'free energy' and to describe the organism's behaviour as 'free-energy minimization'. This is technically correct in as much as free energy is the KL-divergence between the recognition and generative model plus surprisal, and so when there is no recognition model and no KL-divergence, it reduces to surprisal alone. Still, in as much as 'free energy' implies the use of variational inference and the existence of an approximate model, encoded by the organism and distinct from the generative model that describes organism-environment dynamics, so the continued use of the term where no such encoding has been established is misleading. While I would prefer to use 'free energy' only in contexts where we have these two distinct types of models, the use of it to mean nothing more than surprisal relative to a generative model is so pervasive as to be unavoidable. It will have to suffice to emphasise that, going forward, any mentions of 'free energy' mean nothing more than 'surprisal.'

4.2. Self-organization and steady state

In a series of papers Friston (2013, 2018, 2019b) argues that the stability-based account of what it is to be a living system follows directly from a

definition of self-organization in terms of the properties of ergodicity, and low entropy. The statement that a system is ergodic is the part that mandates the distribution over possible states is invariant over time. The ‘low entropy’ part is not used in the thermodynamic sense, but the information-theoretic sense of Shannon and specifies that the particular states our system re-visits will be only a small subset of all possible states, such that the probability distribution over possible states has low variance/high precision.¹⁷

We need the requirement of low-entropy, for ergodicity is not enough to entail surprisal minimization. Repeatedly rolling a six-sided die is an ergodic process, but if it is a fair die then each side is equally likely, its entropy is at maximum and its stable generative model would just be a flat distribution over all possible states. For such a maximal entropy system, any sequence of states would have the same surprisal, and so surprisal is not minimized.

Unfortunately, ergodicity actually implies more than stationarity of dynamics over time, a point that has caused some controversy in recent work on the FEF (Palacios & Colombo, 2021). Strictly speaking, ergodicity not only requires that the average behaviour of a particular iteration of a system will be invariant over time, but also that it will be insensitive to initial conditions. Take a spinning top untroubled by external perturbations. Left alone after an initial impulse, this would quickly settle into orbiting a small area of the overall tabletop – a stable regime, describable by a stationary probability distribution over its position. Which particular part of the table it ends up oscillating around will, however, differ depending on the initial impulse and starting point. As such, spinning a top is not an ergodic process.¹⁸

¹⁷ The entropy and variance of a distribution, while related, are different things. However if (as mentioned in Chapter 3) we are constraining our probability distributions to single peak Gaussians then the only way to change the entropy of a distribution is to change its variance.

¹⁸ This video of the Guinness World Record for longest running spinning top provides a

This stronger requirement of ergodicity has a few important consequences. Firstly, it implies that any particular iteration of a system will eventually visit every state that it is possible for that system to inhabit. Secondly, it implies that a snapshot of an ensemble of iterations at any single point in time will converge with the distribution across states for the trajectory of a single iteration *over* time, as the duration of the individual trajectory or the size of the ensemble increases. This property is expressed in Birkhoff's (1931) ergodic theorem, that with increasing samples and increasing time, the ensemble average and time average will eventually converge.

The classical example of ergodicity in statistical physics is that of idealized gas particles bouncing around a container. Idealized is the key word here for, as Palacios and Colombo (2021) note, proving the existence of concrete systems that meet this requirement has been extremely difficult. In many physical systems the time required to exhibit every possible configuration, or for the time average to converge, extends far beyond the duration over which the system exists (Palacios 2018; Gallavotti, 1999). It is now, they claim, widely recognized that most of the systems studied in statistical mechanics are most likely non-ergodic (Earman and Rédei 1996; van Lith, 2001).

This might be to the advantage of the FEF if it were shown that such ergodicity is the distinctive preserve of the biological. Unfortunately, the opposite appears to be the case. Ergodicity is even more impausible in the biological sphere. To take the favoured example of Stuart Kauffman (2000), who argues it is precisely the defiance of ergodicity that defines biological organization, it would take 10^{39} times the current lifespan of the universe to make all possible permutations of a 200 amino acid long protein at least once.

lovely, if less idealized, illustration.

The convergence of an ergodic process cannot be responsible for the stability of the specific sub-set of amino acid combinations we observe. As Kauffman puts it:

It follows that, even if we consider the universe as a whole, at the levels of molecular and organizational complexity of proteins and up, the universe is kinetically trapped. It has gotten where it has gotten from wherever it started, by whatever process of flow into a persistently expanding adjacent possible, but cannot have gotten everywhere. The ergodic hypothesis fails us here on any relevant timescale. (2000, P.145)

An ergodic system forgets its history. No matter where it starts, after enough time any iteration of an ergodic process becomes indistinguishable from every other. Biological processes, however, fall into local stability wells, where they start matters for where they end up. Any plausibility that the assumption of ergodicity might have for theoretical gas particles bouncing around a box is utterly lacking for the specificity and variation of living organisms.

The ways in which biological processes defy ergodicity are fascinating and informative as to the distinctive character of living systems, and I will discuss them in Chapter 8. For now, I suggest we simply accept the move made in Friston (2019a) & Da Costa et al. (2021) to shed the unnecessarily strong claim of ergodicity and limit ourselves to merely requiring that the probability distribution for each particular iteration of a system remains stationary over time – in dynamical systems terms, the requirement that it reaches a *steady state*. Many of the other critiques raised of the FEF's deployment of ergodicity will, it turns out, apply equally to this weaker requirement. In the interest of getting the whole theory onto the table before we begin dissecting it, I will temporarily put these aside.

So, let's just say I am indeed defined by being a system at steady state and you want to know what I'll be up to at any point in future. You could observe my behaviour over the course of several weeks and chart the relative frequency with which various states are visited. You'd find a very small subset of states – making coffee, drinking coffee, staring mournfully at a mug that is now empty of coffee – to be repeatedly revisited with a high frequency. The vast majority of other, in principle entirely possible, states — relaxing with a glass of champagne in the Balmoral Bar, executing a flawless underwater handstand at the bottom of St Margaret's Loch, spinning fire on the Meadows with the Beltane Society — will be occupied extremely rarely if at all.

Or as Friston, Wiese & Hobson (2020) put it:

At a larger timescale, this trajectory could reflect your daily routine, getting up in the morning, having a cup of coffee, going to work and so on... The key aspect of this trajectory is that it will—after itinerant wandering and a sufficient period of time—revisit particular regimes of state space. (Friston, Wiese & Hobson, 2020, P.31)

If I am defined by being a steady state system, then once you have identified this limited set of states, you have my behaviour sussed for life. Why think that this is true? The idea that my behavioural tendencies will not change or evolve over time seems as unlikely as it would be depressing. The more cynically-minded might suggest that when all you've got is a joint probability distribution, then everything looks like a steady state system. But the FEF is not the first to advance stability as a principle of survival. Before I get to the critical half of this thesis, in Chapter 7 onwards, it's worth looking at what motivates the link between biological survival and stability of dynamics.

4.3. Cybernetics redux

Temporarily postponing immediate philosophical and emotional objections to the idea that such monotony delimits my entire behavioural repertoire – there is at least something to this stability as a common tendency of life as we know it. As both Colombo and Wright (2018) and Seth (2015) argue, this ‘simple imperative’ of avoiding surprises is foreshadowed in the work of the early cyberneticist, W. R. Ashby, who sought to provide an analysis of the adaptive, self-organizing behaviour of living systems that meets the following criteria: “(1) it is purely objective, (2) it avoids all metaphysical complications of “purpose,” (3) it is precise in its definition, and (4) it lends itself immediately to quantitative studies.” (Ashby, 1940, p. 483).

As Froese & Stewart (2010) describe, Ashby’s solution for this was the theory of ‘generalised homeostasis’. By reducing survival to a matter of stability, he proposed that the mysterious appearance of goal-directed behaviour in living systems – no matter how complex and unexpected – results from nothing more than the same tendency to return to equilibrium when perturbed that is exhibited by *all* stable physical systems. As he proposes in *An Introduction to Cybernetics*:

Thus the concepts of “survival” and “stability” can be brought into an exact relationship; and facts and theorems about either can be used with the other, provided the exactness is sustained. The states M are often defined in terms of variables. The states M_1, \dots, M_k , that correspond to the living organism are then those states in which certain essential variables are kept within assigned (“physiological”) limits. (Ashby, 1956, p. 197)

And thus, as he claims in a later paper:

We have heard ad nauseam the dictum that a machine cannot select; the truth is just the opposite: every machine, as it goes to equilibrium, performs the corresponding act of selection. Now, equilibrium in simple systems is usually trivial and uninteresting; it is the pendulum hanging vertically; it is the watch with its mainspring run down; the cube resting flat on one face... What makes the change, from trivial to interesting, is simply the scale of the events. (1962, P. 270)

Situated as a revival of the (less mathematically abstruse) project of Ashbyian cybernetics, the meaning of the simple imperative of the Free Energy Framework becomes clearer. Talk of minimizing surprisal translates to countering the deviation of an essential variable from assigned limits, and to say an organism must avoid surprises just means that it must maintain homeostasis of its essential variables.

Like Friston, Ashby's work goes beyond the submission of adaptive behaviour to formal analysis. Similarly, he sought to account for our cognitive operations in terms of how they serve the coordination of such behaviour. Defining survival as stability not only cleared the ground of teleology or purpose, it laid the foundations for the analysis of the brain as a control system tasked with maintaining the stability of our essential variables. This should not be interpreted as the transparently false claim that our entire neural architecture is dedicated solely to the triggering of autonomic reflexes. Rather, as Ashby declared, his intention was to show "how all the organism's exteriorly-directed activities – its "higher" activities – are all similarly regulatory, i.e., homeostatic" even where what is regulated need not necessarily be a bodily state (Ashby 1956: 196). The manner in which he proposed to do this was both ingenious and somewhat perplexing.

What distinguishes the organism, he argued, is an additional mechanism, found only in a subclass of 'ultrastable systems'. The role of this mechanism is to trigger the random reorganization of the structure of the system when

its essential variables are pushed beyond the threshold that it has either ‘adapted’, or been designed, to compensate for. In the organism, this reorganization amounts to changing the parameters of the behaviour producing network at random – a process that ceases only when a behavioural policy is discovered that brings its essential variables back to stable equilibrium.

When Ashby presented a working model of this, called the homeostat. his ‘electronic brain’ generated both headlines in the popular press (Fig 4., Ashby, 1949) and bemusement from his fellow cyberneticians. Julian Bigelow summed up the general attitude in declaring that, “It may be a beautiful replica of something, but heaven only knows what” (Husbands & Holland, 2012, P.12). Barring the supposition that Ashby was a particularly unorthodox chess strategist, it was far from obvious how a system that flails around randomly until its stability is restored might one day, as he proposed, play the game “with a subtlety and depth of strategy beyond that of the man who designed it.” (Ashby, 1948) Such a task seemed more immediately achievable by the symbol processing systems of Ashby’s contemporaries, Herbert Simon and Allen Newell, which thereby came to define the dominant paradigm in AI research for many years after.

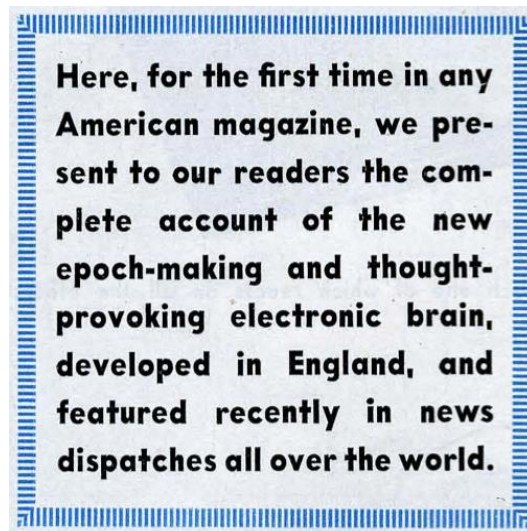


Fig. 4: Intro to an article on Ashby's Homeostat in *Radio Electronics Magazine* (Ashby 1949)

Still, as Vernon (2013) emphasises, while the homeostat and the concept of ultrastability may have represented the apotheosis of Ashby's working models, his theoretical speculations in *Design for a Brain* extended further to the discussion of metastability: as a property that emerges from a high number of interconnected ultrastable systems, and one that might better characterize the operations of something as complex as a nervous system. A specific model of how the nervous system might function in this regard was later developed in *Behaviour: The Control of Perception* (1973) by the cybernetician, William T. Powers, wherein he describes how hierarchical regulation of sensory variables could indirectly support the regulation of external variables, in order to thereby maintain the stability of essential variables¹⁹.

This sounds extremely similar to the operation performed by a predictive processing brain. Unfortunately, the preference for the language of 'predictions' and 'errors', over 'reference levels' and 'discrepancies' in the

¹⁹ As Ward (2016) argues, in Susan Hurley's (2001) synthesis of Power's perceptual control theory with the motor theories of perception we have an even closer forerunner of Predictive Processing.

FEF (and consequently in predictive processing) has obscured these similarities by aligning the latter accounts with the Bayesian brain tradition, over its cybernetic ancestor. Once the interchangeability of these terms under the FEF is made clear, the similarities are striking. If the simple imperative of surprisal minimization is simply a revival of Asbhy's reduction of survival to stability, then we can frame active inference – and the Predictive Processing implementation of this – as providing an analogous proposal to Powers' perceptual control theory. An explanation of how, by taking sensory variables as proxies for essential ones, and by attempting to predict and control these over multiple timescales, a system can engage in the model-based regulation of its environmental interactions, anticipating, and countering dangerous tendencies before they threaten its core stability.

There is, regrettably, not the space to engage in an extended comparison of PP and perceptual control theory here. For now, it should suffice to note that despite their many similarities, there remains (at least one) novel element to the FEF. If the theory is in tune with both Bayesianism and cybernetics, then it offers a potential bridge between the latter's language of control and the inferentialist vocabulary of the former. The value of the FEF might then be in its ability to integrate both the bio-physical grounding of cybernetics in feedback control, and the abstract forms of rationality modelled by Bayesian Brain frameworks, under one unifying formalism.

4.4. Stability and agency

If comparing Friston's project with that of Ashby's helps bring clarity to the free energy framework, it also brings into focus tensions within it. While Ashby is often credited with popularizing the term 'self-organisation', he is

nonetheless at pains to emphasize that there is ‘no such thing’ – at least not when this is understood with volitional implications, as the result of some drive that originates from within the organism itself (Ashby, 1962). All such behaviour, he argues, can only be explained if we view it as a response to external causes. When environmental perturbation drives the system away from a stable state, it triggers this system’s inevitable return to equilibrium. The stable, the ultrastable system, and the multistable system are all driven by the same laws and forces, the only difference is only in the degree of mechanical complexity.

For Ashby then the whole point behind the concept of *generalized* homeostasis is to strip biological explanations of teleological talk, by proposing that the release of glucose stores to counter a decrease in blood-sugar levels, the onset of shivering in response to a drop in body temperature, a pendulum falling back into kinetic equilibrium and the determined uprightiness of the wobble-doll differ only in degree of robustness to perturbation, not in kind. The widespread failure to identify this continuity, he suggests, can be attributed to the lack of systems of intermediate complexity between the multistable person and the basic pendulum. “The computer” he argues:

“is heaven-sent in this context, for it enables us to bridge the enormous conceptual gap from the simple and understandable, to the complex and interesting. Thus we can gain a considerable insight into the so-called spontaneous generation of life by just seeing how a somewhat simpler version will appear in a computer.” (Ashby, 1962, P.271).

For Ashby, the stability of key variables does not just define survival and the persistence of organism but a general tendency of all physical things. Friston (2019a, 2019b) sometimes embraces this generality in presenting the requirement that the states of a system have steady state dynamics as part of

an all-encompassing definition of existence by means of which “everything of interest about life and the universe can be derived” (2019b, P.176)

The idea that stability is an important principle for more than just biological homeostasis runs as follows: if the region of possible states of some system were constantly changing then we would have no way to reidentify it over time. If the region were stationary, but did not have low entropy, then a wide range of states would be equiprobable at any point and there would be no distinctively characteristic states by which this particular system could be identified. As such, Friston (2019b) presents these criteria as amounting to necessary constraints on the possibility of ongoing existence.

The thought that this might serve as an exhaustive definition of all forms of existence seems, to me at least, wildly unconvincing. Nonetheless, once we have this putative definition we can see why Friston often describes the free energy principle as tautological, for when existence is defined in such terms anything that exists will end up satisfying it.

Confusingly, however, despite appropriating this Ashbyian analysis, Friston’s work on the FEF has at other points presented homeostasis-as-stability to be something distinctively biological, evidencing a previous lack of appreciation for its generalizability to all sorts of inanimate systems. As Friston (2010) states

The defining characteristic of biological systems is that they maintain their states and form in the face of a constantly changing environment [...] This maintenance of order is seen at many levels and distinguishes biological from other self-organizing systems; indeed, the physiology of biological systems can be reduced almost entirely to their homeostasis. (P.127)

And, as Ramstead, Badcock, and Friston (2019) put it:

This is the remarkable fact about living systems. All other self-organising systems, from snowflakes to solar systems, follow an inevitable and irreversible path to disorder. Conversely, biological systems are characterised by a random dynamical attractor—a set of attracting states that are frequently revisited. (P. 3)

Where Ashby cautions against the use of volitional language to characterize ‘self-organizing’ systems, discussions of the FEF, in contrast, make liberal use of agential terminology in describing surprisal minimization. This habit persists even in light of more recent extensions of the free energy principle beyond the biotic sphere, to the realm of physical mechanisms in general. When deployed exclusively in relation to neural dynamics, the claim that their surprisal-minimization describes an ‘imperative’ that an ergodic, or steady state, ‘agent’ must follow to ‘actively maintain’ itself (Friston, 2013) does not immediately strike one as obviously misplaced. Yet the free energy framework is supposed to apply to all biological systems, not just en-brained ones. If it does so only in as much as surprisal-minimization describes the stability-through-fluctuation of *any* old physical system, then the idea that it licenses talk of ‘active maintenance, following ‘imperatives’, or even ‘inferring’ and ‘modelling’ becomes much less convincing.

So surprisal minimization is too general to distinguish the unique characteristic of animate existence. Neurocentric predictive processing, conversely, is too narrow. This is a major issue for the FEF as a principle of biological self-organization, and an even worse blow to pretensions at supplying a bioenactivist theory that takes the distinction between life and non-life as fundamental.

Despite its significance, this concern over whether the FEF has the required specificity to serve as a theory of living agents has only recently begun to

meet with serious critical engagement in the free energy literature (DiPaolo, Thompson, Beer, 2022; Raja, Valluri, Baggs, Chemero & Anderson, 2021). As such, I think it's fair to characterize the responses to this challenge as at a relatively nascent stage of development. Still, even at this point, two different strategies can be distinguished.

One draws upon how the active inference framework has been developed and extended in order to characterize increasing levels of hierarchical complexity and, with this, the ability to minimize free energy with respect to future events and probabilistically-weighted counterfactuals (Friston, Wiese and Hobson, 2020; Wiese & Friston, 2021). Said approach suggests that it is this increase in the complexity of surprisal-minimizing mechanisms, rather than the mere property of being a steady state system, that differentiates between simple systems that may merely be modelled 'as if' they are engaging in active inference to maintain the stability of their state, and genuinely inferential agents like ourselves. While Friston, Wiese & Hobson propose this as gradualistic criteria for the emergence of consciousness – in the service of resisting the panpsychism that threatens to result from viewing stable systems as actively inferring ones – when taken as a metric of 'lifelikeness' it may serve as a defence against hylozoistic consequences too.

In taking the question of whether or not a system is a living agent to be a matter of degree, such a response is thoroughly Ashbyian. In making the distinction between life and non-life an instrumental matter, as Wiese, Hobson & Friston (2020) propose, it is fundamentally incompatible with our bioactivist goal to provide an account of basic intentionality that is both naturalistic, and thoroughly realist.

In order to assess the prospects of this strategy, however, we will need a means to analyse non-neural lifeforms, such as single cells, as active inferencers. This requires more than just steady state dynamics, but also a means to factor the system into active and sensory components. To see how the FEF attempts to extend such an analysis beyond the brain we will need the additional tool of a Markov blanket, used to identify this kind of sensorimotor interface. This will take some explaining, and so I'll put this strategy aside until Chapter 5.

The other strategy takes up the common observation, most notably attributed to Schrödinger (1951) and von Bertalanffy (1968), later developed in the work of Ilya Prigogine & Isabelle Stengers, (1984, 1997), that organisms persist out of equilibrium with their environment, then attempts to characterise this in terms of different forms of steady state attractor. As this does not rely on the apparatus of a Markov blanket and, as I believe it is not successful, I will briefly discuss it here.

4.5. The more things change, the more they stay the same

The FEF, as we've seen so far, is an account of steady state systems in the specific sense of systems where the probability distribution over possible states remains stationary over the duration of the system's existence. Unlike the stronger requirement of ergodicity, this allows that the system can be sensitive to initial conditions such that the particular steady state that it eventually settles into may vary across different iterations.

So far, I have discussed this in the context of a homeostatic set point where the steady state is a fixed point attractor, such as a body temperature of 36.5 degrees, from which the system only departs as a result of random fluctuations and to which it reliably returns. This kind of point attractor is what we would see when the system's dynamics are exclusively driven by what is termed the 'dissipative' (or curl-free) flow back to more likely states, which counters surprisal-raising fluctuations in order to prevent the dissipation of the generative model. This is the only flow required for a steady state and the only thing actually mandated by the FEF's 'simple imperative' for living systems.

As a description of a principle supposed to explain the complex human behaviour, this is starkly impoverished. Even in the simple case of body temperature, we find not just a set point but a circadian rhythm, a recurring cycle through different, equally viable, temperatures. Fortunately for the free energy framework, the preservation of a stable generative model does not require the return to a fixed point as the only behaviour of the system. As Friston and Ao (2012) discuss, the preservation of a stationary probability distribution is compatible with the system's dynamics also having a second component, called the solenoidal (or divergence-free) flow. This can be extracted for any system at steady state via the Helmholtz decomposition which breaks the system's overall flow into both dissipative and solenoidal components.

Rather than driving the system back to more likely states, as the dissipative flow does, the solenoidal part circulates around a number of equivalently likely alternatives. This component of the system's dynamics does not itself reduce surprisal, it serves neither to dissipate, nor to counter perturbation-induced dissipations of, the generative model. It is merely compatible with

the generative model remaining unchanged. The dynamical pattern of systems with a solenoidal component will thus be a cycle around a limited subset of states, eg. the orbit of a planet, or a circadian rhythm, rather than a single point attractor.

The Helmholtz decomposition thus allows the dynamics of this kind of steady state system to be described by the following stochastic differential equation which specifies how the state of all of its parts (x) change with respect to 1) the dissipative flow (Γ) towards lower surprisal (\mathfrak{S}) states 2) the solenoidal flow (Q) circulating around equally low-surprisal states, and 3) some random noise (ω) corresponding to surprisal-raising fluctuations which are, by stipulation, countered by (Γ). This is depicted in Fig. 5 and described in Eq. 4 – note that this shows the maximization of $P(x)$ which is equivalent to the minimization of $\mathfrak{S}(x)$.

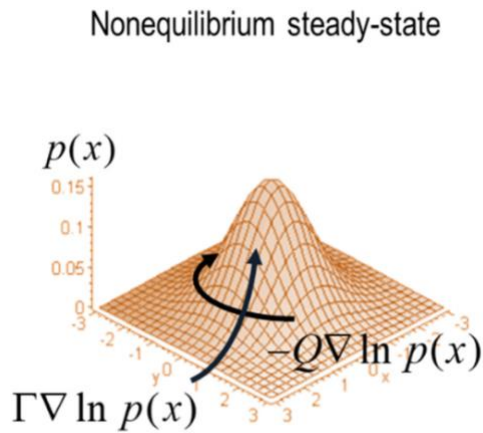


Fig. 5: Depiction of the dissipative (Γ) and solenoidal (Q) flows over a phase space with a probability measure (Friston, Hobson & Wiese, 2021)

Eq. 4

$$\dot{x} = f(x) + \omega$$

$$f(x) = (Q - \Gamma) \nabla \mathfrak{S}(x)$$

This can be connected to the Fokker Planck equation for the particular system to show that it will converge upon a stable joint probability distribution over all the variables making up x .

Strictly speaking, the FEF requires only the dissipative flow (Γ) needed to minimize surprisal-raising fluctuations and preserve a stable generative model. While a solenoidal component is compatible with this, it is not actually required for a system to be free energy minimizing. Nonetheless, the possibility of incorporating this solenoidal component has recently been claimed to distinguish systems at non-equilibrium steady state (Friston, 2019a; Friston, Fagerholm & Zarghami et al., 2021). The idea being that if we can identify a part of the dynamics of a system's internal states that is not driven by response to perturbations (as the dissipative-flow is), then we can identify the 'intrinsic dynamics' by which this system, as an autonomous one, keeps resisting the descent into equilibrium balance with its environment.

While this solenoidal flow is not much discussed in Friston's (2013) Life as we know it paper, in a revealing aside he suggests that it may be key to characterizing the living systems as distinguished by being those steady state systems that are out of equilibrium. As he hints:

However, minimum entropy is clearly not the whole story, in the sense that biological systems act on their environment—unlike a petrified stone with low entropy. In the language of random attractors, the (internal and Markov blanket) states of a system have an attracting set that is space filling but has a small measure or entropy [...]. Put simply, biological systems move around in their state space but revisit a limited number of states. This space filling aspect of attracting sets may rest on the divergence-free or solenoidal flow (equation (2.3)) that we have largely ignored in this paper but may hold the key for characterizing life forms.(P. 11)

This suggestion continues in Friston (2019) but the specific connection between free energy minimization and solenoidal flows, and between a solenoidal flow and the breaking of equilibrium, is still under-explicated. Given that all sorts of non-biological systems, from the hydrological cycle to planetary orbits, have a solenoidal component to their dynamics this component does not, as Friston implies, seem to take us that far in getting at what is distinctive about living systems in particular. Even if we were to accept the presence of a solenoidal flow between equivalently likely states as being adequate to the task, the requirement of free minimization applies equally to both sorts of steady state systems and does not differentiate between them. Further, as Biehl, Pollock & Kanai (2021) and Aguilera et al. (2021) argue, incorporating this solenoidal flow appears to require prohibiting the kind of patterns of influence needed to formalize a ‘sensorimotor’ interface in terms of conditional independencies between parts of our organism-environment system – which will be needed when we attempt to extend the FEF and active inference beyond the brain in Chapter 5.

Nonetheless, the details of the ability of the FEF to formalize non-equilibrium dynamics does not affect the fact that it can only describe systems that can be characterized by a stationary probability distribution. As a theory of living systems, the FEF requires that the probability of their being in any particular state stays constant over time. This remains as the underpinning of a supposed ‘surprisal-minimizing’ imperative for the organism. As such, I will put aside the specific issue of non-equilibrium status, which appears as a somewhat ad-hoc supplement that has yet to be adequately integrated into the FEF, to stick with the weaker requirement of

a system that converges to any sort of steady state. That is, any system that can be described by a stationary joint probability distribution – or, in FEF terms, which ‘entails a generative model.’

In the next chapter, I will consider efforts to extend this beyond the nervous system, in the attempt to characterize a process that is not brain-bound but is still intended to be distinctively biological. After showing how such a characterization still fails to achieve the required level of specificity for a theory of life, I’ll discuss how work on the free-energy principle has responded in a more recent turn towards an Ashbyian view on the continuity between the organic and inorganic. Such a turn I will argue, not only invalidates the FEF’s enactivist’s credentials but also reveals its inability to deliver the first principles of living systems.

5. Active inference beyond the brain

Minimizing surprisal just means maintaining stability, and maintaining stability means surviving, but how exactly does this noble goal hook back up to the messy business of active inference? Active inference may be a means of minimizing surprisal, but this description of a system does not just drop fully-formalized from the observation that the system is surprise-minimizing (or stable, homeostatic, or ergodic) alone. In addition, we need a means to break our system down into the specific set of variables involved in active inference, namely: internal states, external states, sensory states, and active states.

Thus far, these have been presented to us ready-made in the structure of the nervous system and its sensorimotor interface. If the goal of the FEF was no more than the specification of a procedure by means of which brain-enabled creatures like us can maintain our viable states, then we could stop here, consider ourselves satisfied with this as no small accomplishment, and go out and play in the afternoon sun.

Regrettably, FEF's advocates are of a more ambitious bent. Since its (comparatively) modest origins as a 'theory of cortical responses' in Friston (2005) these imperial aspirations has driven the FEF's expansion in scope, sights now set on the provision of "a mathematical formulation of how adaptive systems (that is, biological agents, like animals or brains) resist a natural tendency to disorder" (Friston, 2010, P.127) , and one that can be applied to all forms of life, "from single-celled organisms to social networks" (Friston 2009, P.293).

Which raises the question: where exactly do we find the sensory, active, and internal states of the *E. coli* bacterium, the *Myxogastria* slime mould, the flash mob, or the Bank of England's Monetary Policy Committee?

To understand the FEF's answer to this, we need to get one further piece of mathematical machinery on the table. This is the Markov blanket, introduced by Judea Pearl (1988), and in its original form it is nothing more than a representation of statistical separation between nodes in Bayesian networks.

In the FEF, this has been used firstly, to describe the separation between the internal states of an organism and the external state of environment and secondly, to partition the boundary between these into the active and sensory states of the active inference equations. In the process of applying the Markov blanket to living systems, however, the FEF has pressed this unassuming construct into a series of much more demanding labours, from the establishing an *epistemic* boundary (Hohwy, 2016, 2017), to formalising an *autopoietic* boundary that is both produced by, and preserves the system that it so bounds (Palacios et al., 2017; Friston, 2013; Kirchoff et al., 2018; Allen & Friston, 2018). This has induced such a degree of conceptual hypertrophy that Bruineberg, Baltieri, Dolega and Dewhurst (2021) suggest dubbing this new construct as a “Friston Blanket” in order to distinguish it from its humble ancestor.

The claim that a Markov or ‘Friston’ blanket allows us to map active inference equations onto a diversity of systems beyond the brain, together with the claim that said blanket formalises the autopoietic boundary of biodynamic enactivism are thus both crucial steps in the idea that the FEF might provide this ‘first principle’ of living systems. In order to evaluate

whether Markov blankets are capable of playing either of these roles then, we need to go back to their origin in Pearl's work on Bayesian networks. This will, unfortunately, require a rather prolonged detour from our aim of understanding the FEF as an analysis of living systems, due to the extent to which the notion of a Markov blanket has been distorted in the FEF literature and the lack of explicit acknowledgement of these distortions. If you are already familiar with Markov blankets within Bayesian networks, however, you can skip this section and proceed straight to section 5.2.

In order to claim that the Markov blanket is any sort of boundary for the organism itself – epistemic, autopoietic, sensorimotor, or otherwise – the free-energy framework has presented these blankets as part of the real world and not just a feature of our models of this world as they were initially proposed. This claim has typically been interpreted as a regrettable mistake of, as Andrews (2021) puts it, 'confusing the math for the territory' (see also, Bruineberg et al. (2022), Menary & Gillet (2021), Beni (2021) and Raja et al. (2021)). Nonetheless, while Friston and co. have not provided any justification for the metaphysical claims of Markov blanket realism, I don't think that this is necessarily a mistake. There is, as I argue in an appendix to this thesis, a positive position one could take on the structure of reality and the metaphysics of causation such that it would be plausible to take Markov blankets as real entities, not just modelling constructions.

I don't think such a position is particularly compelling. But this is not particularly important given that, as I will argue in the next chapter, the main problem concerning Markov blankets is less whether they are 'real' or 'constructed' but that they are neither a feature of organisms nor of any description that can capture the necessary conditions of their existence.

5.1. A brief review of causal inference, for the purpose of more clearly elucidating the original nature of a ‘Markov Blanket’

A Bayesian network is a means of representing probabilistic relationships via a directed acyclical graph (DAG for short) in which variables (each represented as a node) are connected by directed lines between, representing their direct statistical dependencies (Pearl, 2000, 1988; Glymour, Spirtes, & Scheines, 1993) (see Box. 1). ‘Graph’ here is just the mathematician’s term for ‘network’. A DAG is directed in that the connections between nodes have a particular direction of influence, and acyclic in the sense that this does not circle back on itself.²⁰

To create a Bayesian network, all that is needed is a joint probability distribution (such as the generative model of the FEP) over the values of our variables and a few simple axioms of construction. The most straightforward reason one might do this is to visualize the decomposition of an unwieldy joint probability distribution over a large number of variables into a set of more tractable independent distributions over smaller subsets of variables, which can be selected amongst depending upon which particular variables we are concerned with. To illustrate how a Bayes net is constructed, we can take a simple set of three variables X, Y and Z.

If $P(x|y) \neq P(x)$, that is if knowing that $Y=y$ changes the probability that $X=x$, then we have a statistical relationship between the two. This is not enough to justify drawing a line directly from Y to X yet, however. To see

²⁰ The graph theorist would say ‘vertices’ and ‘edges’ but I will stick with the more familiar vocabulary of nodes and connections.

why, suppose it is also the case that, knowing that $Z=z$ also changes the probability that $X=x$ ($P(x|z) \neq P(x)$). This would leave us with three possible ways in which our three variables could be connected.

1) $Y \rightarrow Z \rightarrow X$: Fixing the value of Y alters the probability that $X=x$ via changing the probability that $Z=z$, which then directly alters the probability that $X=x$ in turn:

2) $Z \rightarrow Y \rightarrow X$: Fixing the value of Z alters the probability that $X=x$ via changing the probability that $Y=y$, which then directly alters the probability of $X=x$ in turn:

3) $Z \rightarrow X \leftarrow Y$: Fixing the values of either Y or Z will directly alter the probability that $X=x$.

Directly connecting Y to X is only correct in case 2 or 3. In order to disambiguate between these three options then, we need to look at whether $P(x|y) = P(x|y, z)$. If it does, then once we know that $Y=y$, knowing the value of Z contributes no more information about the value of X – that is to say knowing $Y=y$ renders the state of X *conditionally independent* of the state of Z . As such, Z can only raise the probability of X via Y , and so we draw our network according to (2).

If, however, $P(x|y) \neq P(x|y, z)$, then even once we know the state of Y , knowing the state of Z can still alter the probability that $X=x$. This leaves us with either option (1) or option (3). If $P(x|z) = P(x|z, y)$ then knowing the state of Z renders the state of X conditionally independent of the state of Y , and we draw the network according to (1). If, however $P(x|z) \neq P(x|z, y)$ then neither Y nor Z screens off X from the other, and so we draw our

network according to option (3). By continuing this process with a larger number of variables, we can construct a network like that shown in Box. 1.

A Markov blanket (Pearl, 1988) is a partition on this network that contains all and only those nodes which collectively render the state of some target node X conditionally independent of any other nodes in the network. This means that once the states of this subset of nodes are observed there will be no further observations that can provide any more information within our model about the probability of X being in a particular state – other than to observe the state of X directly.

This set amounts to the following: the nodes upon which the probability distribution over possible states of X is directly dependent (its Markovian parents), those nodes whose state is directly dependent upon X (its children), and the state of any other nodes on which X 's children are also directly dependent (the co-parents) –for a more detailed explanation, see Box. 1.

Box. 1:

Suppose I am trying to figure out what sort of mood my partner, Max, is currently in (M) - here modelled as a variable that may be either positive or negative. I cannot observe this directly. Instead, I have to infer it from a limited set of available observations. These include: whether he won his last Rocket League game (R); if he's eaten recently (H); whether he's grumbling to himself (N) and whether he's baking scones (B).

```
graph TD
    G((G)) --> R((R))
    G((G)) --> H((H))
    R((R)) --> M((M))
    H((H)) --> M((M))
    M((M)) --> N((N))
    M((M)) --> B((B))
    M((M)) --> D((D))
    N((N)) --> N1(( ))
    N((N)) --> N2(( ))
    B((B)) --> B1(( ))
    B((B)) --> B2(( ))
    D((D)) --> D1(( ))
    R1(( )) --> R((R))
    H1(( )) --> H((H))
    N3(( )) --> N((N))
    B3(( )) --> B((B))
    D2(( )) --> D((D))
```

I know that if R is positive (a victory) then the probability that Mood is positive increases, so we can draw an arrow between these nodes, with M being downstream of R. If H is negative, the probability that M is positive decreases, so these are similarly connected. I know that a positive M also raises the probability of scone-baking, so I can also draw a connection from M to B. If I want to know the converse, however – that is how the observation of scone-baking alters the probability distribution over Max's mood – then, as Bayes rule tells us, I also need to consider any other states that also alter the probability of scone baking, such as whether he has a deadline within the next two days (D), which would also increase the probability of baking – irrespective of mood.

Together, these nodes make up the Markov Blanket of M, jointly encapsulating all of the information available in my model about M's current state. Thus, I can take the observations that 1) he won his last Rocket League game ($R=1$), 2) he is baking scones ($B=1$), and 3) he does not have a procrastination-inducing deadline within the next two days ($D=0$), and so infer a high likelihood of a positive mood.

There are other variables in my model, which would also alter my probability distribution over M if I knew their state - such as whether he scored two or more goals in his last game (G). These are not a part of the Markov blanket of M, however, because a positive value of G increases the probability of a positive M only by increasing the probability of a victory ($R=1$) Once we know that R is positive then, the number of goals scored has no further information to offer about the probability of a positive M.

In this manner, a Markov blanket renders the target node conditionally independent of the rest of the network.

5.1.2 Bayes nets in causal dress

As described, the immediate advantage of this is in revealing how a complex joint probability distribution may be factored into a number of smaller conditional probability distributions. In this context, a Markov blanket is

nothing more than a statistical device – one which allows us to separate out those variables whose state becomes informationally irrelevant with respect to the state of some target variable once we know the state of this smaller subset of nodes. This factorization, as we saw, is exactly the kind of thing that needs to happen in order to implement approximate Bayesian inference via hierarchical predictive processing.

Still, there is a second more contentious role for Bayesian networks. As developed by Glymour, Spirtes, & Scheines, (1993) and later adopted by Pearl (2000), this is their use in causal discovery, allowing us to transform purely observational data concerning the relative frequencies with which selected variables take particular values, into a putative model of the underlying causal structure.

There's a careful line that needs to be walked here. It may be true that no amount of nifty modelling work or nice diagrams will overcome the fact that correlation doesn't *equal* causation. Nevertheless, they are closely related. The growing field of algorithmically-driven causal discovery has shown that, when combined with some assumptions, a couple of putative axioms of causation (such as its unidirectionality), and a bit of background knowledge, the purely statistical property of conditional independence may be surprisingly effective in constraining possible hypothesis about where causal relationships may lie.

Returning to our three variables X, Y and Z, this time with our causal inference hat on, will allow us to more clearly see both the extent and the limitations of a purely observational approach to inferring causation. We saw that looking at whether a third variable renders two variables conditionally independent of each other allowed us to discriminate between three possible

graph structures. However, I surreptitiously simplified this problem by sneaking in the assumption that X is downstream of both Y and Z.

From a purely statistical perspective, this is unproblematic. To draw the lines in this direction is not to deny that there is a reciprocal relationship in the other direction. It merely expresses that I am interested in one particular direction of statistical dependence for the purposes of inferring the likely value of X given Y, rather than vice versa. From a causal inference perspective, however, the choice of direction is a stronger commitment. One motivation behind the use of DAGs in causal modelling is the understanding of causality as a one-way affair. Thus, to draw an arrow from Y to X is not merely to express which one of two directions of dependence I am interested in. It is to rule out any causal relationship in the opposite direction.

Without this ordering assumption, we would have had not just 3, but 10 possible graphs to consider, all satisfying $P(x|z) \neq P(x)$ and $P(x|y) \neq P(z)$ (See Box. 2) Further investigation into conditional independence relationships would, in most cases, fail to deliver us one unique graph. Without ordering information, the best we can do is to narrow the possibilities down to one of four sets, three of which contain three further possible causal graphs, all observationally equivalent with respect to a particular statistical relationship between our three variables.

Box. 2
X and Y are independent given Z: $P(x z)=P(x z,y)$, $P(y z)=P(y z,x)$
$Y \rightarrow Z \rightarrow X$
$Y \leftarrow Z \leftarrow X$
$Y \leftarrow Z \rightarrow X$

X and Z are independent given Y : $P(x|y)=P(x|y,z)$, $P(z|y)=P(z|y,x)$
 $Z \rightarrow Y \rightarrow X$
 $Z \leftarrow Y \leftarrow X$
 $Z \leftarrow Y \rightarrow X$
 X and Y are independent given Z : $P(x|z) = P(x|z,y)$, $P(y|z) = P(y|z,x)$
 $Y \rightarrow X \rightarrow Z$
 $Y \leftarrow X \leftarrow Z$
 $Y \leftarrow X \rightarrow Z$
 No conditional independences: $P(x|y) \neq P(x|y,z)$, $P(x|z) \neq P(x|y,z)$,
 $Z \rightarrow X \leftarrow Y$

While causality is usually taken to be unidirectional, statistical relationships are not. The occurrence of smoke raises the probability of fire about as much as vice versa, and so the mere fact that A raises the probability of B does not tell us which is the cause, and which the effect. In one sense, this does not matter for the Markov Blanket. The set of nodes included in this is determined by the presence or absence of direct connections between nodes, not their direction. It remains invariant throughout any arrow flips that may legitimately occur while still preserving these same statistical relationships. The set of graphs that preserves the same conditional independence relationships but with different arrow directions is known as the ‘Markov equivalence class.’ What may change, however, is which nodes are the parents, and which are the children.²¹ This will be important when we come to how the FEF uses the Markov blanket construct to factor a system into active and sensory states.

Incorporating time series data resolves some of this directional ambiguity, but there are still problems of unobserved confounders and co-incidentally counterbalanced causal chains (Hesslow, 1976). To take Jeffrey's (1969) example of the former: the probability of lightning is increased when the barometer moves to the left, but despite their temporal sequence, the barometer does not cause lightning. No matter how long you observe the two, mere observation of their statistical and temporal relationship will never teach you about the latent common cause of air pressure.

Let's suppose, however, that there are no confounders or counter-balanced causal chains not included in our model. Let's help ourselves to the temporal information required to specify its ordering. Even this would not be enough to deliver a unique model with an unequivocal Markov blanket. We still face a decision about how to parcel our complex molecular world up into the atomic building blocks of a Bayesian network. A given component of interest may either be broken down into various interacting components, or amalgamated into the state of some global variable. For the neuroscientist should a node correspond to the state of one brain region, one neural cluster, one neuron, one synapse, or one individual neurotransmitter molecule? As with many modelling choices, there is no single answer – beyond saying that it depends what you're interested in. Bayesian networks, like all models, necessitate idealization and simplification. A more fine-grained model is not necessarily a better one (Cartwright, 2001. Borges, 1998).

Just as the atomic causes of our Bayesian network are a modelling distortion, so too are the tidy little arrows between them (Spohn, 2001; Pearl, 2000). Take a model of conditional dependencies among the membrane potentials of individual neurons. As a Bayesian network, this would represent direct interactions between one neuron spiking and the next. By virtue of the

granularity of its variables, said model omits the fact that in order for one neural spike to trigger another there must be intermediate events of neurotransmitter release and uptake across the synaptic cleft. Interpreting the direct connections of a Bayesian network as direct causal relations, and so understanding the Markov blanket as a boundary made up of the ‘most proximal’ causes and effects, is thus liable to be highly misleading. Anderson (2017) expresses the point nicely in the context of the attempt to locate the ‘most proximal’ boundary of the brain.

An obvious candidate answer would be that I have access only to the last link in the causal chain; the links prior are increasingly distal. But I do not believe that identifying our access with the cause most proximal to the brain can be made to work, here, because I don't see a way to avoid the path that leads to our access being restricted to the chemicals at the nearest synapse, or the ions at the last gate. There is always a cause even “closer” to the brain than the world next to the retina or fingertip. (P.14)

All of this goes to explain why Pearl (2000) firmly caveats his argument for the utility of Bayesian networks as causal models with an emphasis on the separation between statistical constructs (such as a Markov blanket) and causal facts. A significant body of assumptions and ‘causal intuition’, he notes, are necessary for a network’s construction and causal interpretation. This is not too great a concern when our system is made up of nothing more than seasons, sprinklers, rain, and wet grass. In the case of highly-complex and poorly understood systems like neural or intracellular networks, however, ‘causal intuition’ will be woefully inadequate to select a single causal model from the innumerable many that will be compatible with the statistical (and temporal) relationships identifiable through observation of our selected variables (Mehler & Kording, 2018).

5.2 Markov blanket Realism

We have seen that Markov blankets, as they were introduced by Pearl (1988), are a property of a statistical model. Both their constitution and (optional) causal interpretation is dependent upon a series of choices and idealizations made in the course of a model's construction. In the work of Friston and co-authors, however, the humble blanket frequently appears not just as a property of our models but of the system being modelled itself. As Friston (2019b) asserts, a Markov Blanket “is not some statistical device by which we come to model the world - it is a necessary attribute of a universe that can be carved into things.” (p.176) As he elaborates with Wanja Wiese, “In the context of the FEP, it is assumed that a Markov blanket is a property of the system itself... This version of the concept is therefore a metaphysical notion” (Wiese & Friston, 2021: p4).

In the FEF & PP literature, the Markov blanket has not only become more concrete, but has been clothed in a variety of unexpected, and almost unrecognizable guises. Their role in the FEF first gained philosophical attention through the work of Hohwy (2013, 2016) where they become an epistemic boundary that cuts us off from the world beyond our sensory veil – a world that we must nonetheless strive to infer an accurate representation of. Hohwy focuses on brain-bound organisms, using the Markov blanket to draw out sceptical implications for a system that already has a pre-defined sensorimotor interface. Yet, we do not need the Markov blanket formalism to identify active, internal & sensory states in the brain, and much early work on the FEF in the brain proceeded without it.

Where talk of Markov blankets has become pivotal to the FEF is in underwriting the application of the active inference formalism beyond the

nervous system. As such, we can view their introduction in Friston (2013) as marking a second wave in the development of the FEF, the point at which some of its advocates move from the proposal of a ‘theory of cortical responses’ (Friston, 2005) to a theory of biological systems across spatial and temporal scales, “from single cells to social networks” (Friston 2009a, P.293).

In this context, the Markov blanket is used firstly to delimit the boundaries of a target system and then, crucially, to divide this boundary up into the sensory, active & external states that form the free energy equations (Palacios et al. 2020). Here, as Hipólito et al. (2021) describe, the target node is just the internal state, its parents are the sensory states of the bounded system, while its children are the active states:

In this context, we associate the variable of interest with the internal states of a Markov blanket; which allows us to think of the ‘parents’ of that variable as mediating the influence of external states on internal states (i.e., as sensory states) and of its ‘children’ and the ‘parents of the children²²’ as mediating the influence of internal states on external states (i.e., as active states) (P. 90)

Note that, as above, despite often introducing the Markov blanket as a ‘statistical boundary’, it is common in the FEF literature to quickly move to causal talk of ‘influence’, ‘mediation’, and explicitly stating that “parent nodes *cause* their children” (P.90). This causal talk is not accidental. Establishing a one-way direction of influence, not entailed by a merely statistical relationship, is needed to secure a basis for partitioning the system itself into active and sensory states in terms of whether a state directly influences or is directly influenced by the internal state. We saw in the previous section that additional work is required to get from statistical separation to a putative

²² It is not immediately obvious why the parents of children are parcelled up as ‘active states’ here, given that they themselves may not be influenced by internal states at all.

causal boundary and that there is a further gap between the putative boundary of a simple model and ‘actual causality’ in the system itself. These difficulties are largely side-lined in the discussion of Markov blankets within the FEF.

Beyond causal concretization, the Markov Blanket has also gathered a further property, becoming something that is both necessary to preserve the continued existence of a system, and something the system preserves in turn. It is not just a description of statistical independence, rather it is something that ‘induces’ this independence (Ramstead, Kirchoff, Constant & Friston, 2019). As Friston and Allen (2018) state, “In short, the very existence of a system depends upon conserving its boundary, known technically as a Markov blanket, so that it remains distinguishable from its environment—into which it would otherwise dissipate” (P. 2475). It is this understanding of the Markov blanket as a real entity, both produced by the organism and causally responsible for preventing that organism’s dissipation in turn, that is at issue in the FEF’s pretensions to have subsumed the circular relations of self-production and self-distinction described in autopoiesis (Friston & Allen, 2018; Kirchoff et al, 2018) to provide “a first principle of living systems” (Friston, 2012).

The task of generalizing the essential features of the autopoietic boundary beyond the molecular membrane of a cell, in order to characterize the autonomy of multicellular lifeforms, has long been a challenge for biodynamic enactivist accounts of life. A possible advantage of the Markov blanket in this respect would be in allowing us to describe the manner in which an autonomous system is separated from its environment, without requiring encapsulation by a continuous physical boundary of the sort conveniently found surrounding the biological cell.

Further, as Hesp, Ramstead & Constant et al. (2019) propose, such liberality enables the identification of the same free-energy minimizing dynamics over a hierarchy of scales – thereby allowing natural selection, social interaction, cultural evolution, development, learning and planning to all be cast as stages in one grand multi-generational quest for free energy minima.

There is, as far as I'm aware, no gelatinous cell-like membrane encasing the nine members of the Bank of England's Monetary Policy Committee. Nonetheless, using the notion of a causal Markov blanket we can subject their collective behaviour to FEF analysis as a stable system that maintains the constancy of an essential variable – the rate of increase in the price of goods – by responding to changes in sensory states, like price indexes, through altering the active states of the base rate and the issuance of government bonds.

You may, at this point, feel that the attribution of vitality to the BoE's Monetary Policy Committee (over and above the limited supply contributed by its individual members) looks less like a success than a suggestion that something has gone horribly wrong with the FEF's analysis of what it is to be a living system. I wouldn't disagree. In the section following this one, we'll see how the definition of active and sensory states in terms of causal Markov blanket can be enlisted to underwrite an even more horrifyingly promiscuous vitalism. The introduction of (causal) Markov blankets to individuate living systems, a task for which they are woefully under-qualified is, I will argue, the point at which the FEF starts to lose touch with biological reality.

There are two key issues to unravel in the attempt to understand this piece of the FEF's analysis of life. Firstly, how, if at all, does it make sense to be a

‘realist’ about Markov blankets? Secondly, if a Markov blanket is indeed a property of a real system, and not just of our models of them, then can it do the work that the FEF, in its aspiration to provide an bioenactivist analysis of living systems, has called upon it to do?

The first question has received the majority of attention among critics of the FEF, such as Bruineberg et al. (2021), Menary & Gillet (2021), Beni (2021) and Raja et al. (2021) who accuse Friston and co. of reifying aspects of their models. The mistake, as Andrews (2021) puts it, of confusing the ‘math for the territory.’ Still Markov blanket realism doesn’t have to be a mistake. All the above authors argue is that Friston and co. have not supplied the metaphysical premises required for their metaphysical conclusions.

That doesn’t mean no such premises are available. Nor would it require a commitment to mathematical Platonism as Beni (2021) Menary & Gillet (2021) & Bruineberg et al. (2022) suggest. There is a positive position one could take on the structure of reality and the metaphysics of causation such that it would be plausible to take Markov blankets as real entities, not just modelling constructions. This position involves two aspects, probabilistic graph realism and the statistical reduction of causation, perhaps alongside a Penelope Maddy-style naturalized mathematical realism about sets as having an existence over and above that of their parts (Maddy, 1990, 1997).

Probabilistic graph realism means taking reality itself to literally have the structure of a Bayesian network, being made up of independent nodes, whose state is determined exclusively by local interactions with other nodes, which respect the Markov condition – such that the state of each unit is independent of the state of its non-descendants, conditioned upon its

parents. Positions related to this have been defended by philosophers such as David Lewis (1992), David Papineau (1992, 1992) and Hartry Field (2003).

This view is often, but not always, combined with a statistical reduction of causation, which takes the statistical relationships described in a Bayesian network to be all there is to causal relationships. On such a view, the reason that correlation does not equal causation is not that they are different things entirely, but just because we haven't gathered *enough* correlational information to uniquely determine causal structure. Such a view has been defended by Spohn (2001) Reichenbach (1956), Good (1959), Suppes (1970) & Papineau (1992) (and for reviews, see Salmon, 1980 & Weslake. 2005) though there is disagreement on whether this accounts for causation a feature of reality or just an inescapable feature of our conceptualization of it. If causation does reduce to statistical relationships in this way, then the Markov blankets of the ultimate graph of reality become *causal* boundaries, in the way Friston and co. often assume.

One further commitment is required for these real Markov blankets to be something that can demarcate an organism and upon which existence depends, as Friston (2013) and Friston and Allen (2018) suggest. This is the requirement of what I'll call a *stable* Markov blanket, namely that the structure of the graph corresponding to the organism – determined by the patterns of interaction between its parts – is stable enough such that there is a fixed set of components that renders its conditionally independent of its exterior and this set endures for the duration of the organism's existence.

As Friston (2013) puts it:

“a candle flame cannot possess a Markov blanket, because any pattern of molecular interactions is destroyed almost instantaneously by the flux of gas molecules from its surface. Meaning we cannot identify a consistent set of blanket states rendering some internal states independent from other state” (P. 2)

So all that positing a real Markov blanket involves is the claim that a system decomposes into independent units (accomplished by some means prior to the FEF, by methods the framework itself does not specify) and arguing that if we accept objective statistical dependencies as an adequate reduction of causal relationships, along with the principle of locality, then the immediate surroundings of any one of those units will literally have the Markov property of inducing a conditional independence between what is inside this boundary and what is outside of it. To claim Markov blankets are real things is to make a general claim about the structure of the causal universe. In doing so the FEF has neither identified a new and interesting entity in the world, nor discovered a principled basis for carving the world into things. This latter result depends on the necessarily prior task of telling us what the absolute units of the ultimate graph are.

Whether we should care about these metaphysical questions depends upon the second question, regarding whether Markov blankets are actually of any importance in a theory of living systems. In the next chapter, I will argue that Markov blankets are of no help at all in helping us formalize bioactivist notions of autonomy or autopoiesis. The FEF lacks the tools to “subsume autopoiesis” (Korbak, 2021, P2747) or to ‘supersede or absorb classical (i.e., autopoietic) formulations of enactivism” (Ramstead et al. 2021, P.59), and it is hard to see why one would think it could. Instead, the FEF’s claim to perform this role, as it turns out, depends upon an implicit assumption about

the cyclical structure of a system's causal graph. This provides an alternative demarcation of said system prior to our factoring it up with a Markov blanket. While this particular graph structure, as I will describe, does bear some similarities to early definitions of autonomy, it bears no obvious relation to contemporary accounts in terms of a closed network of precarious processes.

There is a second reason to be less concerned with the reality of Markov blankets than with their relevance for living systems, as I will describe in Chapter 8. This is the observation that even if some parts or aspects of reality do have the structure of a statistical-causal graph, an organism is about the last place we would expect to find the stability of interactions needed to pick out a stable Markov blanket as persisting and defining the existence of this system over time.

Recap

Before we move on to the evaluation of the free energy framework [FEF] as a formalization of bioenactivism, a quick refresher of the key concepts introduced so far.

The Free Energy *Principle* is the claim that every system of a certain type (microcircuits, brains, organisms or literally every ‘thing’, depending on the claimed scope) must minimize the free energy of its constituent parts in order to continue to exist. **Free energy** is a function of two things. One: **divergence**, the difference between a simplified **recognition model** that a system encodes and the actual statistical properties of the process that generates its evidence, called **the generative model**. And two, **surprisal**: the unlikeliness of a particular state relative to this true generative model.

It turns out that minimizing the first thing matters only in so far as it positions a system to minimize the second. Most discussion of the free energy principle as it applies to non-neural systems disregards claims about an encoded recognition model to focus on the ‘true’ generative model and the minimization of surprisal.

Taking the minimization of surprisal to be necessary and sufficient for self-preservation rests on a definition of ongoing existence in terms of being either **ergodic** or having converged to a **steady state** (Friston & Mathys, 2016). ‘Steady state’, crucially, does not require that the state of the system never changes, only that any changes will preserve the same statistical properties such that the probability distribution over potential states of the system at any randomly selected point in time is stationary over the duration

of the systems existence.²³ This is compatible with the presence of symmetric and stochastic fluctuations, with a consistent amplitude, in the system's state but these must be countered by a return to more likely states, described by the **dissipative flow**. It is also compatible with, though does not require, the possibility of cycles between the same subset of equally likely states, described by a **solenoidal flow**. Further, the FEF also requires that this stationary probability distribution has **low entropy** such that the repertoire of states said system is likely to be in is not only constant, but also relatively small.

The argument for presenting these components as jointly definitive of existence was as follows. On one hand, stationarity alone would be satisfied if a wide range of states remains equiprobable over time – but in such a case there would be no distinctively characteristic states by which the system could be re-identified. On the other hand, if the system was likely to be in only a small region of states at each time period, *but* this region was constantly changing, then we would have no way to reidentify it as the same system persisting over time. The paradigmatic example given of steady state dynamics is biological homeostasis, though as we saw in Chapter 4 it can also be extended to describe the stability of non-biological systems.

The connection between the minimization of free energy and the maintenance of a non-equilibrium steady state is supplied by **active inference**. This extends the idea of variational free energy minimization, a means to adjust an approximate recognition model to better fit one's

²³ This was originally secured via the stronger claiming that systems are ergodic, but while ergodicity entails stationarity, it also entails that a system will eventually explore every state that it is possible for it to be in - such that the probability distribution describing the states of a particular system over time converges to the distribution that describes the states of an ensemble of that type of system at a single moment in time. This is a stronger claim than needed, and one inappropriate for living systems (see Palacios and Colombo, 2021)

evidence, by adding in the option of acting so as to change this evidence instead. Such a system can thus minimize free energy in two ways. 1) by changes in the internal states that encode its recognition model, thereby ‘learning’ what its characteristic sensory inputs are – as described by the generative model of said system. 2) by acting to minimize the surprisal of its sensory evidence, thereby countering dispersal away from these characteristic states. This rests on the idea that in reducing the surprisal of sensory states, a system implicitly reduces the surprisal of the distal states that characterize it as the kind of system it is.

Importantly, however, the ‘inferential’ interpretation of this rests on claim that the system actually encodes this ‘recognition model’. Neither the claim that a system ‘embodies’ or ‘entails’ a generative model, and ‘minimizes surprisal’ implies that anything inferential is going on. All this means is that the system has steady state dynamics with tendency to return to the same set of states when perturbed, such that we could describe it *by* the stationary joint probability distribution of a generative model

Still, when there is no recognition model, and no divergence between one and a true generative model, then, mathematically, free energy does reduce to surprisal. As such the FEF takes something’s merely being a steady state system, describable by a generative model. to suffice for describing it as actually ‘using’ this model itself, in order to infer and to minimize its ‘free energy’.

To do this, the FEF needs to decompose the system’s constituent states into internal, sensory, external, and active variables. And it is this decomposition that is supposed to be accomplished by means of the **Markov Blanket** formalism just discussed.

All of these constructs together makes up the Free Energy *Framework* [FEF].

The upshot of this is the claim that the dynamics of a system at steady state can be described as formally analogous to Bayesian inference. But is a steady state system with a Markov blanket enough to formalize anything like autopoiesis, autonomy, operational closure, or organizational closure? And, even if it isn't does it provide an alternative definition of a living organism?

6. The free energy framework and the missing cycle

If we are willing to allow some pretty hefty assumptions about the nature of reality, then any stable entity we pick out will have a Markov blanket composed of a further set of elements that suffice to make its state conditionally independent of everything else. Does this, combined with the requirement that a system is at a surprisal-minimizing steady state, give us everything we need to start talking about active inference?

Friston (2013, 2019b) certainly claims as much, declaring that an ergodic (or steady state) system and a (stable) Markov blanket are all that is needed to get active inference off the ground and to interpret internal states as changing so as to minimize variational free energy with respect to a probabilistic model. “With this existential dyad,” he claims, “everything of interest about life and the universe can be derived, from biotic self-organisation through predictive processing to the detailed microcircuitry of our brains.” (P. 176)

We need at least these two requirements as neither a *stable* Markov blanket nor steady state dynamics entail the other. Under probabilistic graph realism a charging battery has a Markov blanket, but over that duration its state of charge is increasing, not steady. Similarly, various stable chemical reactions might be describable by a fixed probability distribution over some collective property, like average concentration, but due to a constant flux of materials, there would not be sufficient stability of interactions between these parts to establish a stable Markov blanket.

So a particular system will be a candidate for analysis in terms of active inference only if it has both a) stability in terms of the changing states of each of its parts – which gives us the steady state distribution or generative model, and b) stability of interactions between those parts – which gives us a *stable* Markov blanket.

Still, there are a heck of a lot of stable physical systems with a boundary more persistent than a candle flame that look nothing at all like a living cell. Without its chocolate casing, a cream egg would disintegrate. Getting at the cream requires breaking through this shell. The force of the cream pushing against its chocolate cage may be cast as an ‘active’ state; countervailing air pressure from the egg’s surroundings a ‘sensory’ one. But for all that the chocolate egg renders its gooey innards conditionally independent from the world outside, a cream egg is neither sentient, autopoietic nor autonomous.

How then, are we to understand claims like the following?

...life—or biological self-organization—is an inevitable and emergent property of any (ergodic) random dynamical system that possesses a Markov blanket.” (Friston, 2013, p.1)

And such Markov blanketed systems are:

... autopoietic: because active states change—but are not changed by—hidden states, they will appear to place an upper (free energy) bound on the dispersion (entropy) of biological states. This homeostasis is informed by internal states, which means that active states will appear to maintain the structural and functional integrity of biological states. (P.5)

Similarly, Kirchoff, Parr, Palacios, Friston and Kiverstein concur (2018) that:

any Markov blanketed system will embody recurrent processes of autopoietic self-generation, which—as long as the system exists—enforces a difference between a living system and everything else. (P.6)

A point that is made more expansively by Allen & Friston (2018) as follows:

For example, a cell persists in virtue of its ability to create and maintain a boundary (cell-surface), through which it interacts with the environment, thereby maintaining the integrity of the boundary. It is this autopoiesis, or self-creation, which enables the system to limit the possible states it visits, and thus to survive (Varela et al. 1974). The FEP recasts this as a kind of self-fulfilling prophecy, in which an organism itself constitutes, in the generative sense, a belief that it will prevail within certain embodied and environmental conditions. In short, the very existence of a system depends upon conserving its boundary, known technically as a Markov blanket, so that it remains distinguishable from its environment—into which it would otherwise dissipate. (P.2473)

Such statements rarely acknowledge the distinction emphasized by both Maturana and Varela (1980) between autopoiesis as a recurrent process of metabolic self-assembly, as opposed to the more general concept of autonomy which attempts to generalize this logic of self-production beyond the molecular interactions of a single cell. In as much as autopoiesis is constrained to a specific level of chemical interactions, no purely statistical generalization could capture it. Elsewhere, however, free energy theorists target the more general notion of autonomy, which, despite general agreement to the contrary (Thompson and Di Paolo, 2014; Bich and Arnellos, 2012; Di Paolo, Buhrmann & Barandiaran, 2017), they take to be adequately expressed in the concept of operational closure.

It is, Ramstead, Kirchoff, Constant & Friston (2021) claim, “fairly straightforward to establish that the Markov blanket formalism provides a

statistical formulation of operational closure.” Given that the kind of circular relationships of input to output that define operational closure have no part in the concept of a Markov blanket, the idea that they serve as a formalization of the former concept appears far from straightforward to me.

Sadly, the paper by Kirchoff et al. (2018) that Ramstead et al. cite as performing this ‘straightforward’ establishment does nothing of the kind. The main justification provided for relating the Markov blanket to operational closure is by way of reference to what Varela calls ‘the intriguing paradox’ of autonomy in how it requires that living systems are closed-off and distinguished from their environment, while at the same time completely dependent upon remaining coupled with that external environment for their ongoing existence. Kirchoff et al.’s suggestion is that the conditional independence of internal states from external states describes this closure, for once we have determined the blanket states we are closed off to the possibility of gaining any further information that might reduce our uncertainty in predicting internal states. At the same time internal states are still open to influence from external ones, *via the blanket*. The Markov blanketed system thus exhibits a balance of closure to information (conditioned upon the blanket) and openness to causal influence (via the blanket).

As an analysis of autonomy, this is unconvincing. Being an autonomous system, as we will see in more detail in Chapters 7-9 is not just a matter of exhibiting closure to some things and openness to others, but a matter of being closed to the *right* things. A formalization of autonomy is not provided by describing any old mixture of openness and closure but rather by identifying precisely what particular things an autonomous system must be

open or closed to. Nonetheless, Kirchoff et al. move on from this to claim that

This teleological (Bayesian) interpretation of dynamical behaviour in terms of optimization allows us to think about any system that possesses a Markov blanket as some rudimentary (or possibly sophisticated) ‘agent’ that is optimizing something; namely, the evidence for its own existence. This means we can regard the internal states (and their Markov blanket) as, in some sense, autonomous. (2018, P.2)

We can also think of a wobble doll as believing that remaining upright is its highest calling, and actively striving to achieve this. Or, as Kirchoff et al. acknowledge, a pendulum as engaging in active inference. But just because we can, doesn’t mean that we should. Surely it takes more to legitimate agential and inferentialist language than just stability in the face of perturbations.

Elsewhere, Friston, Hobson & Wiese (2020) instead target the notion of ‘sentience’, pointing to the ability to divide a Markov blanketed system into parents and children of the target node, where the parents of some node are those upstream that directly affect it, and children are those downstream that are directly affected by it.²⁴ The idea of a sensory state, they claim, is captured in nodes whose parents are external states, and whose children are internal states – and vice versa for active states²⁵.

This is no more compelling as an analysis of sentience than it is of autopoiesis. If these dependency relations were all there were to an active or

²⁴ This ordering already introduces more than the existential dyad strictly buys us, for a single ordering into ancestors and descendents will not always be supplied by statistical relationships alone.

²⁵ Co-parents of children, though a part of the Markov blanket, are typically ignored - presumably due to the difficulty of fitting them into the category of either sensory or active states

sensory state and if, as Friston, Hobson & Wiese advocate, being ‘responsive to sensory impressions’ is all there is to sentience, it looks like we are ascribing capacities of sensation and action to everything that has – or at least can be modelled as having – a stable Markov blanket over some duration.

Friston, Hobson, and Wiese attempt to elide the panpsychist implications of this by noting that their attribution of sentience “is not used in the philosophy of mind sense; namely, the capacity to perceive or experience subjectively, i.e., phenomenal consciousness, or having ‘qualia. Sentience here, simply implies the existence of a non-empty subset of systemic states; namely, sensory states.” (P.3). This may be enough to avoid attributing a rich inner life to everything with a Markov blanket, but one still wonders what possible justification could then be offered for using the terms like ‘sentience’ and ‘sensory’ here. Taken in the most minimal sense of implying the capacities of sensation and action, attributing sentience on the basis of a Markov blanket alone is still going to lead to a mathematically-motivated animism that considers steam engines, pendulums, plants and people as all sentient systems alike.

There must then be still more baked in to what Friston and co. mean by a ‘Markov blanket’ than the properties of stability and physical realization.

6.1. The missing cycle

Now we are more familiar with Markov blankets, claims of such a straightforward connection to autopoiesis, autonomy, or sentience, and so to the emergence of life, agency and consciousness should immediately appear rather suspect. I imagine that the fact that such proclamations have remained

unchecked for so long is due to the way Markov blankets are typically explained in the FEF literature.

Firstly, there is usually little offered by way of introduction other than the claim that they are a statistical partition, which then quickly morphs into causal talk of ‘influence’ when introducing the division into descendants and ancestors of a target node²⁶. The interpretations of the Markov blanket as a mathematical object, or as a physical boundary, are treated interchangeably, and there is typically little, if any, mention of either a) the various assumptions of the modelling framework of causal graphs, which underpin the initial partition, or of b) the metaphysical premises required to secure its causal interpretation and physical concretization.

Secondly, the paradigmatic example chosen to illustrate a Markov blanket is almost always the cellular membrane (Friston, 2013; Palacios et al, 2020). While the cellular membrane is indeed a physical boundary surrounding the cell, and would thus be a Markov blanket in the supposed all-encompassing graph of reality, the properties of ‘self-generation’ and ‘self-preservation’ that Friston and co. go on to attribute to it are held, not in virtue of this, but only in virtue of its also being specifically an *autopoietic* boundary – something that the vast majority of physical boundaries and Markov blankets are not. This is rather like taking the true statements that “Amy is my friend” and, “Like all my friends, Amy is a featherless biped” then proceeding to suggest I have successfully analysed the concept of friendship.

So far I’ve engaged in some metaphysical speculation to justify interpreting a Markov blanket as a physically instantiated, causal boundary, and

²⁶ The parents of children are more difficult to incorporate in the FEF framework, and are typically ignored

supplemented this with an additional requirement for the stability of said boundary. I have also pointed out that the FEF is concerned only with the Markov blankets of a steady state system, and as such it requires not only stability of dependencies between parts but also of stability in the tendencies of these parts to occupy a particular subset of states. Yet even this pumped-up ‘stable Markov blanket plus steady state’ description remains incapable of supporting the claims that are made of the ‘Markov blanket formalism’ in the FEF literature.

Importantly, however, in as much as the steady state Markov blanket fails to entail the circular dependence involved in autopoiesis or autonomy, it also fails to capture the kind of perception-action cycles described in active inference. As presented in the context of a brain, active inference was not just about a system that remains at steady state, and in which some parts are statistically shielded from others, but also one in which these parts exert a distinctively circular pattern of influence on each other. In particular: it described a system where a discrepancy in the sensory part of the system, relative to the internal part (interpreted as encoding a predictive model of these sensory states) leads to a change in action, which changes the external environment so as to alter the sensory state and reduce this discrepancy. How is this circular pattern of influence derivable from a Markov blanket of either an ergodic or steady state system?

The answer, and I’m sure you can imagine how frustrating this would be had you (hypothetically) spent months trying to understand how active inference can be derived from these explicitly declared and oft-repeated prerequisites, is that it can’t.

Perhaps then it is this further requirement that is supposed to differentiate the living autopoietic agent from the merely stable system. Perhaps it is only a subclass of Markov-blanketed systems, whose internal and blanket states change in line with the cycle described by active inference, that should be described as producing and reciprocally depending upon a stable, physical Markov blanket. Can this requirement distinguish a ‘self-organizing system’ from all other things that merely have such a stable blanket through no fault of their own?

6.2. The Existential triad

Markov blankets, as discussed previously, are traditionally used in the context of a directed *acyclic* graph (DAG). In the philosophical literature on the FEF this is typically the format in which Markov are depicted. The problem is that the acyclicity of this type of graph specifically prevents depicting the kind of circular connection needed for active states to have a reciprocal influence on the sensory nodes that are their (indirect) ancestors.

Nonetheless, acyclicity was not part of the requirements laid out for probabilistic graph realism. It is indeed possible to have a graph with cycles that are factorizable by Markov blankets, though it is less straightforward than in the acyclic case. Accordingly, there is a second, quite different, graph also presented as depicting the structure of an active inferrer that appears in Friston and co.’s discussions of the FEF (see fig. 6) typically found rakishly imposed over a picture of a brain. Unlike the acyclic structure of fig. 6a, this second diagram, fig. 6b, explicitly depicts the kind of cycle between external, sensory, internal, and active states (ESIA) that would be expected in the perception-action loop of active inference.

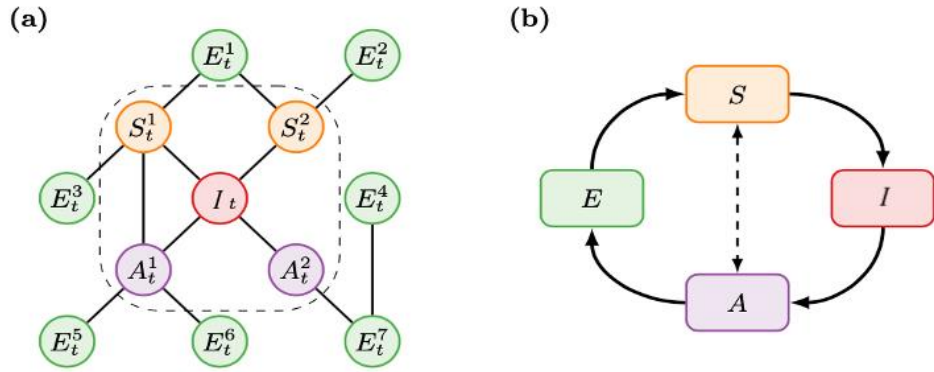


Fig. 6: Depicting a Markov blanket in the standard acyclic graph (a) and the cyclic graph that is typically presumed in work on the Free Energy Framework (b) - adapted from Rosas et al. (2020)

The first diagram is derivable from the assumption that something is divisible into parts whose interactions respect the causal Markov condition, as discussed in the previous section. The second is not. Rather it depends on the ability to partition our overall steady state system (x), described by the stochastic differential equation, Eq. 5 into a set of coupled equations describing four sets of variables, external (E), sensory (S), internal (I), and active (A) and how they influence each other – as shown in Eq. 6.

Eq. 5:

$$\dot{x} = f(x) + \omega$$

$$f(x) = (Q - \Gamma) \nabla \mathfrak{F}(x)$$

Eq. 6

$$f(x) =$$

$$\dot{E} = f_E(E, A, S) + \omega_E$$

$$\dot{S} = f_S(E, A, S) + \omega_S$$

$$\dot{I} = f_I(I, A, S) + \omega_I$$

$$\dot{A} = f_A(I, A, S) + \omega_A$$

Said equations specify how each set of variables changes as a deterministic function of the state of a subset of other variables, plus a noise term representing the stochastic fluctuations that the surprisal minimizing tendency of our system is stipulated to dissipate in order to remain at its steady state. These are supposed to meet the requirements for a Markov blanket between internal and external states, for internal states are a function only of sensory, active and other internal states, not of external states directly – though see Rosas et al. (2020) Biehl, Pollock & Kanai (2021) and Aguilera et al. (2021) for a more technical criticism of the relation between the two. Just as the Markov blanket induced a synchronic requirement for the stability of statistical dependencies between the parts of a system, the EISA-cycle equations describes one specific flow of interactions that meets this requirement by specifying diachronic relationships in how the various types of state must change as a function of one another.

It is these equations that are represented in fig. 6b. They have, however, received much less attention in philosophical discussions of the FEF, where the claim that the Markov blanket is what determines the boundaries of the system of interest has largely been taken as read (Hohwy, 2017; Clark, 2017; Kirchoff and Kiverstein, 2021). If a diagram showing the cycle of fig. 6b does appear, it is often treated as a straightforward depiction of a Markov blanket and presented without any explanation of these equations on which it is based (Hohwy, 2017, Kirchoff and Kiverstein, 2021).

By establishing a cyclical flow, these equations specify the more restrictive requirement for active inference over and above a steady state and a Markov blanket – namely, for a pattern of feedback from active to sensory variables, via external ones. This cycle as Rosas et al. (2020) note, looks closer to the kind of thing that we think of as a sensorimotor loop and makes the

identification of constituent variables as ‘sensory’ or ‘active’ a slightly less trivial matter than when it is done in terms of Markovian ancestry alone.

Beyond this, the ‘EISA cycle’ equations also tell us which external variables should be included in the overall model of the system-in-its-environment, and which can be excluded. The Markov blanket itself could not do this because, as far as a Markov blanket is concerned, one external state is just the same as any other no matter how far away in the causal chain. All are equally irrelevant to predicting internal states once the states of the blanket nodes are established. In fig 6a. the inclusion, or exclusion, of an external variable looks like an arbitrary choice. Once we have the cycle defined by equations eqs. 6, and depicted in fig. 6b, we can see that the only Markov-external variables that are relevant from the perspective of the free energy framework are those that are part of this cycle – namely those Markov-external variables whose state is both a function of active variables, and a partial determinant of the state of sensory variables.

This helps make sense of the FEF’s rather implausible requirement that it is not only the internal, active, and sensory states of a system, but also its external environment that must be at steady state. If this were not the case, if internal states converged to a stable regime while external ones were constantly changing then the former could not be regarded (even in the weakest correlational terms) as entailing a model of the latter. As Millidge, Seth & Buckley (2021) note, if steady state were a requirement taken to apply to the universe as a whole, it would (theories of eternal recurrence aside) be quite obviously false. The EISA-cycle allows us to restrict this requirement to the Markov-blanketed system’s niche, or *umwelt*, now defined as all and only those ‘external’ variables that are part of said cycle. It is only this niche that we require to be at steady state, and its being so can be explained as a

byproduct of the Markovian system's activities in working to maintain its essential internal variables at steady state.

What then of the extra-external world, outside this EISA-defined niche? As modelled in the equations of the EISA-cycle, the entirety of the remaining universe is demoted to random fluctuations – specifically, to normally-distributed, independent, additive noise terms (w) which are stipulated as sufficient to capture the only kind of influence that the EISA-external world can have upon the steady state system (Millidge, Seth & Buckley, 2021).

This idea we can adequately describe a system and its environment in terms of this simple steady state cycle, and condense all external influences upon it to uniform noise, is somewhat more tenable than the requirement that the entire universe be in a stable dynamical regime – though not by much. I have already hinted at the implausibility of suggesting that the entire lifecycle of an organism can be modelled by the convergence to and maintenance of a single, stationary probability distribution. I will discuss this further in Chapter 8. The EISA-cycle may, as Millidge, Seth & Buckley suggest, make more sense as a locally-valid approximation. Still, the fact that our environment is so volatile and changeable, not only in its temporary state but also in its ongoing dynamical tendencies, is precisely what makes biotic systems' ability to maintain their homeostatic stability so interesting. To *begin* by modelling the system's local surroundings as being at steady state would seem to extract much of the interest from this problem at the get-go.

As Millidge, Seth & Buckley put it:

*This is the real question which is what the FEP tries (at least in its intuitive sales pitch) to answer – how can I maintain an internal steady state against an environment which is **not** at steady state. By assuming that the external states are also at the steady state, it may be that the FEP is, in some sense, answering the wrong question and is, in the process, assuming away the true difficulty in answering the right one. (P.34)*

As Biehl, Pollock & Kanai (2021) and Aguilera et al. (2021) note, the EISA-equations are not derivable from the Markov decomposition of our overall steady state probability distribution. Nor are they direct consequences of the requirement that our system be ergodic, or at steady state. Instead, it seems that the steady state density and Markov blanket are intended to be derivative properties of a dynamical system that can be described by the coupled stochastic differential equation of the form in Eq. 5-6.²⁷

Diagram 6b certainly appears a little closer to being a candidate of operational closure, or a perception-action cycle, than diagram 6a.²⁸ If we say a sensory or active state is determined not only by its Markovian ancestry but also by its position in such a cycle, then we have a slightly more restrictive definition of these than one that depended on conditional (in)dependence alone. But before evaluating the suitability of this formulation to capture the particular kind of closure that distinguishes living systems, we need to look a little more closely at how the system defined by these EISA-cycle equations

²⁷ As Biehl Pollock & Kanai (2021) and Aguilera et al. (2021) described, are further requirements that need to be in place to secure the Markov blanket. This is raised as a critique of some technical details in Friston’s work that are beyond the scope of this paper however.

²⁸ Crucially, like the Markov blanket however, these labels will still depend on the prior selection of some set of states as ‘internal’ ones. In the EISA cycle, internal and external states are symmetrical and there is nothing to prevent us swapping their labels, along with those of action and sensory states.

meshes with the claims that have been made about the role of the Markov blanket.

6.3. Extensional Ambiguity

With the EISA-cycle equations and the Markov blanket both in hand we now have two different partitions and two different sorts of ‘externality’ involved in the definition of a single system: things external to the Markov blanket, but internal to the ESIA-cycle and things that are external to both. As such, as Raja et al (2021) note, Friston’s ‘existential dyad’ is ambiguous on whether it is the overall steady state EISA-cycle, or only the Markov blanketed subset of this, which determines the boundaries of the system that the free energy principle is supposed to define.

In some places, the Markov blanket is described as what differentiates “between the system and its environment – those states that constitute or are intrinsic to the system and those that are not.” (Ramstead et al. 2018, p. 3), and as providing “a statistical formulation of operational closure” (Ramstead et al. 2021, P. 55). Yet the steady state equation, which ranges over the entire EISA-cycle, is also described as capturing ‘the phenotype’ of an organism (Ramstead et al. 2020) or even, as Friston (2019) more broadly declares, something that can capture the entire concept of ‘thingness’. This seems at one remove to consider something as external to our organism (or other system of interest) and at another to consider it as part of the phenotypic states that define it.

We have already discussed the inadequacy of the Markov blanket to formalize anything like self-organization or autonomy, to identify a

sensorimotor interface, or, indeed, to carve out anything particularly unique at all. The natural move then would be to turn instead to the EISA cycle, which we will do in the next section. But if we take this to provide both the demarcation of our system and its decomposition into active, sensory, internal, and external states, then there is a residual question: what is it that the Markov blanket demarcates?

As with the ‘simple imperative’ of surprisal minimization, I think we can gain some insight from Ashby’s work on self-organization here, specifically a procedure he suggests for describing the appearance of an intelligent system. (Note that ‘equilibrium’ is used in the looser sense of stability or stationarity of dynamics here, and is compatible with a system being at non-equilibrium steady state or ‘dynamic equilibrium’).

Take a dynamic system whose laws are unchanging and single-valued, and whose size is so large that after it has gone to an equilibrium that involves only a small fraction of its total states, this small fraction is still large enough to allow room for a good deal of change and behavior. Let it go on for a long enough time to get to such an equilibrium. Then examine the equilibrium in detail. You will find that the states or forms now in being are peculiarly able to survive against the changes induced by the laws. Split the equilibrium in two, call one part ‘organism’ and the other part ‘environment’: you will find that this ‘organism’ is peculiarly able to survive against the disturbances from this ‘environment.’ (1962, P.120).

The key idea, here and elsewhere in Ashby’s discussions of self-organization, is that the appearance of any intrinsic volition is illusory. What we call ‘self-organization’ is really driven by the coupling between our selected system and another one. We might artificially decompose these and treat each as an individual system, but neither has intrinsic stable dynamics that are independent of their interaction. The Markov blanket is a means of making a similar division, but in describing it as the boundary of an organism, the free energy framework has unwisely elected to drop the scare quotes that are

essential to reminding us that this distinction between the two is a choice imposed by an external observer.

The problem with taking the Markov blanket as the boundary of an organism, rather than a pragmatic division, can be made more concrete in the case of a system that is both enbrained and embodied. Here, this Markovian partition could cut around the nervous system – in which case bodily variables (eg. core temperature) and environmental ones (eg. air temperature) take on the same equal status as undifferentiated ‘external’ variables. This makes sense in that both are part of what we may want to describe the nervous system as modelling and regulating but it erases the distinction between the body and the environment, allowing no special import to be given to the management of the former. If, instead, we take the relevant Markov blanket to be that which divides the whole body from the environment, then the internal relationship between neural-regulator and regulated body is washed out in favour of describing the entire body, en-masse, as a regulative model of its surroundings.

This is not just the point that I have already drummed to death about the variety of scales at which Markov blankets can be identified. As Ramstead et al. (2019) and Sims (2020) suggest this could naturally be taken as analogous to the multiscale nature of organisms composed of organs, organelles and cells. The bigger issue here is that at neither the scale of the brain, nor of the whole body, does the Markov boundary adequately parcel up the system we are interested in as a brain-equipped organism.

Ashby seems to take the mere appearance of one system’s adapting to another as sufficient to found an account of intelligence in terms of the complexification of this process. As a consequence of choosing to treat the

differences between inanimate matter, organisms, or intelligent systems as a gradient, and homeostasis as a general property of stable systems, his position is undisturbed by the concern that this decomposition into ‘organism’ and ‘environment’ is a modelling heuristic that can be applied to anything from a Watt Governor and a steam engine or a pair of coupled pendulums.

In the context of this Ashbyian approach to self-organization, the question of whether a particular Markov blanket truly demarcates an organism would reduce to a matter of whether the active inference analysis that a Markov blanket facilitates grants sufficient explanatory power to pay its way in our best scientific theory of the system. This is a much weaker view of what the free energy principle delivers than has often been claimed – through Friston, Hobson and Wiese (2020) do entertain such an instrumentalist stance with regard to the attribution of sentience, specifically as a means to avoid the panpsychist consequences of a Markov blanket-based demarcation of the mind.

Extending this instrumentalism further, to the demarcation of life itself, seems a more radical (and ethically unpalatable) position. Nonetheless, if the FEF’s advocates are prepared to treat the distinction between organism, mechanism, and other material objects as a matter of degree – the demarcation between them as something dictated by the interest of an observer – then this would at least render the free energy principle’s more philosophical claims consistent with their Ashbyian underpinnings.²⁹

²⁹ This approach to living systems as cycles upon which an observer may impose an organism/environment distinction that does not originate from the cycle itself, is a direction continued in Humberto Maturana’s pre-enactivist development of autopoiesis theory. While autopoiesis theory is the genesis of a great deal of bioenactivism’s conceptual repertoire (not least with regard to autopoiesis itself) the two should not be confused, and indeed take radically opposed views on the status of life. Crudely, while the former retains this Ashbyian view of the continuity between life and non-life, the latter

If, however, the FEF is supposed to be something that “provides an implementation of enactivism, and in a sense supersedes or absorbs classical (i.e., autopoietic) formulations,” (Ramstead et al. 2021, P. 59) then reducing the idea of an organism to an observer-dependent projection will not do. Above all else, bioenactivism is about autonomy, about systems that are *self-distinguishing* and *self-producing*, not heteronomous mechanisms which we might pick out and ascribe purposes to as an explanatory heuristic. If anything is going to demarcate this operationally-closed and autonomous unit then, it must be the cycle picked out by the EISA-equations of the overall steady state system.

6.4 The Markov Blankets Last Gasp

A plausible alternative to treating the Markov blanket as something that demarcates a living system might run as follows: any organism (a steady state EISA-cycle, demarcated by non-Markovian means) will be divisible into cognitive-regulator and bodily-regulated parts and there will be a Markov blanket between these. Markov blankets, being exceedingly trivial if they are to be considered real features of our world at all, cannot suffice to exclusively identify this demarcation but the fact that said division is also a Markov blanket may allow us to say some interesting things about the relationship between these parts.

For instance, we can note that as each Markov-blanket-individuated set of parts: internal, external, active and sensory, has steady state dynamics derived

originates in Varela’s later phenomenological, organicist, or existential turn, which grants organisms special status as ‘natural purposes’ capable of ‘bringing forth a world of meaning.’ More on this, in section 7.4

from those of the overall organism, so each will each tend to spend most of its time in the same characteristic state (or cycle of states if the overall system also has a solenoidal flow) returning to these when perturbed by internal or external fluctuations. If we describe the stable dynamics for each component in terms of a probability distribution over most likely states, then we will be able to interpret each component's return to its characteristic state (or cycle) when perturbed as minimising surprisal relative to this probability distribution. As each component is doing this, so there will be low KL-divergence between these respective probability distributions, and each time one component returns to its lowest surprisal state when perturbed it may be interpreted as minimizing the KL-divergence between its average state and that of other parts of the overall steady state system.

Still, all that 'inferential' description reflects is the covariational dynamics that are inevitable among the parts of any system that has the kind of overall stability the FEF requires. 'Inference' here is not a means to obtaining this stability, but rather a redescription of it, and there is nothing especially intentional, agential, or representational about the kind of systems that admit such redescription. As argued in section 3.3, to straightforwardly take KL-divergence as a measure of misrepresentation, and its ongoing reduction as evidence of something's performing a representational function, would be to reduce representation to covariation and licence the attribution of representational relationships between the parts of any stable system.

As Raja et al. (2021) note, it is not immediately clear that we gain anything by describing this relationship in Bayesian vernacular that we did not have access to from a description in terms of the dynamics of feedback control loops between coupled subsystems. Markov blankets, as they put it, appear to be nothing more than a 'trick' of exploiting the general divisibility of

(most) systems into parts whose interactions respect the Markov condition, in order to redescribe their relationship as the inference of one to the other. Whether you take this to be reflective of fundamental metaphysical truth, or merely a feature of a particular modelling convention, either way, it does not reflect anything uniquely teleological, cognitive, agential, or animate about the particular systems that we chose to execute the procedure upon.

As with mere covariation, this general property may gain special significance if we identify it in an agent that actively deploys this covariation to representational ends. Performing such a function would, arguably, require the potential use of the internal parts for action guidance in the absence of an ongoing connection to those external parts which it they have come to resemble through a prior process of perception-action coupling (Grush, 2004; Haugeland, 1991). This detachability is not something that can be modelled in terms of a single, simple active inferrer, which is defined entirely in terms of its coupling to a target system. It is, however, potentially addressable in terms of hierarchical systems, like PP, where the activity of each ascending level can be viewed as an individual free energy minimizer, increasingly detached from current sensory input from the overall system's periphery (Corcoran, Pezzulo & Hohwy 2020; Pezzulo, 2017).

Still, my focus here is not on whether or not some particular free energy minimizing systems may, or may not, be able to meet the detachability criteria for having representational parts. The issue at hand is to identify what, if anything, free energy minimization as a 'principle' contributes to the bioenactivist desire for a naturalistic explanation of the emergence of intentionality and agency in terms of biological autonomy. Perhaps once this is obtained and we have a naturalized account of the system's primary function, then Markov blankets might be used to formalize this hierarchical

detachment and the development of proper representations. That is, they may be a part of how we build a bridge between modelling and inference with basic processes of autonomy and self-preservation. The prior question, however, is whether the steady state EISA-cycle does a good job of describing that basic autonomy.

7. Seeking closure

The idea that the thermodynamic openness identified by von Bertalanffy (1968) as a condition of life, must be paired with some form of closure by which the individual organism may be identified and distinguished from its environmental backdrop is far from unique to either the autopoietic tradition or the FEF. Closure is central to Jean Piaget's (1971/1967) work on the nature of life; Howard Patee (1982) proposes the notion of semantic closure; Robert Rosen (1991, 1999) draws upon Aristotle to support an analysis of living systems in terms of 'closure to efficient causation'; Stuart Kaufmann (1986) develops an account of 'catalytic closure', later extrapolated to the more general notion of 'work-task closure (Kaufmann, 2000), and further developed by Montèvil & Mossio (2015) and Moreno & Mossio (2015) in terms of 'closure of constraints'.

While this enduring concern with some sort of 'closure' suggests that we might be on the right track in taking it to be central to a theory of life, the variety of forms in which it has been defended should also alert us to the possibility that the closure of the FEF, and that of the bioenactivist, may not be the same thing. This is important because, as we'll now see, the steady state EISA cycle is far too general to capture anything distinctively autonomous about living systems, being as it is applicable to any coupled system whose collective dynamics converge to a steady state.

7.1 EISA-closure

We have already discussed how a Markov blanketed set of internal states can be described as exhibiting a balance of (conditional) *informational* closure with *causal* openness – and how trivial this particular sort of closure is. The EISA cycle introduces a new kind of closure proper to the overall steady state system ‘ x ’, which the Markov blanket would partition into internal and external parts. The state of each component in this EISA-cycle is partly set by a deterministic function whose domain of inputs is limited to the state of another subset of components that are part of this same cycle, and partly by a gaussian noise term w . In this sense, it is closed to any non-stochastic, or non-symmetric, influence from operations whose constituents are external to this cycle. Thus, in accordance with the requirement that the overall system remain at steady state, this function will keep the components of said system at either a fixed position or oscillating through a recurring cycle of states.

As regards Varela’s ‘intriguing paradox’ and the necessity of pairing closure with an openness to environmental interchange, the EISA cycle does admit external influence in the form of the stochastic elements of the steady state equations (w). This grant a limited contribution from the environment as a source of random, symmetrically distributed perturbations away from the states prescribed in the deterministic component of the EISA equations. Such perturbations are counteracted by the ‘dissipative’ element of the EISA flow, which is set so as to match the magnitude of said fluctuations (see section 4.5). The justification for this stipulation is that if this flow did not adequately counteract said fluctuations then surprisal would not be minimized, the system would eventually drift away from its characteristic

steady state, and so, according to the free energy principle, would cease to exist.

It seems worth clarifying that the kind of closure described by the EISA-cycle cannot be intended as the requirement that the particular components of some system be isolated against any *potential* influences that would not conform with these equations. There is nothing, living or otherwise, whose states are immune from the possibility of alteration by anything other than uncorrelated fluctuations. If this were what EISA-cycle closure required then it could only be an idealisation that is never actually satisfied in reality – which would make it a poor criterion for the presence of life.

Instead, this EISA-closure must be understood as the definitional criteria for identifying a particular process as taking place, rather than a constraint on the range of possible interactions that the realizer of said process may be subject to. So to say that, for instance, two coupled pendulums can be modelled in terms of a steady state EISA-cycle does not mean that it's impossible to grab one pendulum and force it into a particular position. What it means is that such a deterministic influence would be incompatible with their entrainment and would violate the closure that defined them as '*coupled pendulums*'. While the process must exhibit closure for so long as it continues to operate, its realizer remains vulnerable to influences that would violate said closure and thereby terminate that process – in much the same manner as organisms, in the process of living, are vulnerable to the possibility (the eventual inevitability) of dying.³⁰

³⁰ As I'll discuss in section Chapter 9 this focus on closure of, or among, processes exclusively has helped prevent an adequate formalization of autonomy which in Mossio and Moreno's (2015) terms, requires closure between two regimes of causation involving both change and invariance

So, the requirements of steady state and the EISA-equations define a cyclical process that is closed to any outside influences that would change its average behaviour, disrupt its steady state dynamics, and thus break the stationarity of the joint probability distribution of its generative model. It is, however, is open to brief fluctuations consistent with the variance of this probability distribution. Where might we find such processes?

Unfortunately for the FEF's pretensions to provide something that 'defines the form of life that an organism is seen to enact,' (Ramstead, Badcock & Friston, 2018, *sup matt 4*, P.33), something that is sufficient to characterize the emergence of adaptive behaviour (Friston, 2013), and that that 'formalizes', 're-describes' or 'subsumes' the concept of autopoiesis (Allen and Friston 2018; Korbak 2021; Ramstead et al. 2021) all that is necessary to meet these requirements is mutual entrainment among a pair of coupled systems whose collective dynamics synchronize to a steady state. Systems which could then be decomposed, via a Markov blanket, into external, sensory, internal, and active parts, dependent on which part is first chosen as internal.

As Kirchoff et al. (2019) discuss, such a description applies perfectly well to Huygens classical example of the mutual entrainment of the two oscillating pendulums, coupled via vibrations transmitted through a connecting beam that is (or may be modelled as) a Markov blanket between them. Similarly, Baltieri, Buckley & Bruineberg (2018) give a more detailed treatment of how a Watt governor coupled with a steam engine-powered flywheel satisfies the 'existential triad' of the FEF, allowing their redescription in active inference terms.

If we take the arm angle of the governor as an internal state then the free energy framework allows us to class the rotational speed of the axle, on which this directly depends, as a ‘sensory’ state. As this increases beyond some set threshold the arms rise as a consequence, resulting in the closing of the steam engine throttle valve, the ‘action’ node. This reduces the flow of steam to the engine thereby slowing its speed, which constitutes the ‘external’ state that the governor is regulating.

On this basis, an active inference interpretation allows us to describe the Watt governor as inferring the speed of the steam engine, and the angle of its arms as constituting a model of the same. This interpretation, as Baltieri et al. note, depends on the arbitrary selection of the governor’s arms as our internal state. Nothing in this set-up, or the toolkit of the FEF, prevents us from selecting the engine’s speed as our internal state instead and then analysing this as inferring and modelling the Watt governor.

This triviality with which the FEF’s requirements are satisfied by non-living systems is not a unique problem for it as a theory of cognition. The idea that Watt governors and other stabilizing systems might share important principles with action-perception coordination is, after all, the basis of the dynamical approach to cognitive science (Van Gelder, 1995). To avoid attributing cognitive capacities to simple mechanical regulators, FEF-theorists could allow that a system’s being describable in terms of inference does not amount to its actually performing inferential operations, as Wiese and Friston (2021) suggest. This would be to concede that the FEF does not constitute ‘mark of the cognitive, or an answer to ‘what is cognition?’ any more than the idea of a dynamical system does. Rather a more modest alternative would be to suggest that, like dynamical systems theory, the FEF provides a set of tools for analysing things in general that *may* allow us to

identify the specific properties that distinguish cognitive ones (Andrews, 2021).

If, however, we were to stick with presenting the FEF as an attempt to specify the conditions under which “life — or biological self-organization — is an inevitable and emergent property” (Friston, 2013, P.1), then this broad applicability to systems that are obviously not alive, and show no tendencies of becoming so, looks fatal.

Still, fatal for whom? If the free energy framework has indeed delivered a satisfactory formalization of the bioactivist’s autonomy-based definition of life, then it would have shown that autonomy cannot do the work that bioactivists need it to do. The inadequacies of prior characterizations of both autopoiesis and autonomy to pick out the essential features of living systems have been extensively discussed (Bickhard, 2009b; Di Paolo, 2005; Bitbol and Luisi 2005; Collier, 2004, 2008; Bourguin and Stewart 2004; Fleischaker, 1988) (see Froese & Stewart (2010 P.9 for a comprehensive overview). The FEF’s purported formalization cannot be blamed for its triviality if this merely reflects the incoherence, as Ashby claimed, of truly self-driven behaviour of the sort the bioactivist seeks. Perhaps all this reflects is that there is no such thing as an immanent teleology of organisms, no fundamental distinction between a system that displays “*lifelike*” behaviour and which “*appear[s]*” to actively maintain its structural and dynamical integrity” versus a truly living system that genuinely works towards its ongoing existence.

Alternatively, the fault may lie with the FEF’s existential prerequisites of steady states and EISA-cycles. Even they were adequate to formalize early definitions of autonomy, which I shall shortly argue they are not,

bioenactivism is not wedded to these particular formulations. In recent years much exciting work has been done upon addressing inadequacies in these earlier accounts, to develop and refine the concept of autonomy. This work has been roundly ignored in the free energy literature. I will return to this in Chapter 9 to look at the state-of-the-art when it comes to biological autonomy. First a closer look at the history of bioenactivist formulations of closure and how these relate to the FEF.

7.2. Bioenactivism, autonomy and Closure

In the introduction to bioenactivism in Chapter 1 I gave a brief summary of the role closure plays in Di Paolo and Thompson's (2014) definition of autonomy, where closure is realized by a network of processes such that each process is both enabled by, and a condition for, at least one other process in that network and would run down absent this network's support. Once we have selected our target of investigation, identifying these mutual dependence relations allows us to determine the boundaries of this network, in terms of those connected processes that meet this criteria.

This description of closure presents two immediate contrasts, in terms of both relations and relata, from the closure of the EISA-cycle. Firstly, where the EISA-cycle described a single process in terms of a closed cycle between changes in *states of variables*, in Di Paolo and Thompson's (2014) definition closure is specifically a relationship *between processes*. Secondly, for Di Paolo and Thompson it is key that this 'dependence' relationship is not merely one of causal or statistical influence, in the sense that the state of a barometer depends upon the air pressure, but existential dependence where one process would cease to exist as a process altogether if not for its enablement by

further processes making up said network. With the further requirement that there are no additional redundancies that would compensate to continue this process in the absence of its support from within the network, then we have a system that is precarious and so, for Thompson & Di Paolo, autonomous.

The EISA-cycle, as a cycle of patterns of influence between states rather than of dependence between processes, seems to get us no closer to this formulation of autonomy than the steady state Markov blanket did. Still, over the history of attempts to extract the key principles of autopoiesis and extend them to autonomy, one can find a wide variety of subtly different formulations of the particular kind of closure that enactivists have taken to be the foundation of autonomy.

More importantly, as Bich and Arnellos (2012) argue, there are (at least!) two distinct notions of closure that play importantly different roles in both Maturana and Varela's work on autonomy and autopoiesis – namely operational vs organizational closure³¹.

These are standardly run together not only in the free energy literature but also in bioenactivist discussion as either a mere terminological choice (Barandiaran, 2017) or as emphasizing different aspects of the same sort of system (Thompson, 2007). Somewhat confusingly, the term that is generally favoured in attempts to characterize autonomy is 'operational closure', with precariousness stuck on as an additional requirement. In spite of this terminological choice the kind of requirement that contemporary

³¹ The history of these two terms would require work beyond the scope of my argument here - particularly given the ambiguity over attribution that Maturana & Varela's co-authorship entails. It seems to me that their work does not consistently follow the terminological distinction Bich and Arnellos suggest, and moreover I'm not convinced that 'organizational' is the best choice to distinguish the closure of a self-producing system from that of a merely operationally closed one.

bioenactivists (eg. Thompson and Di Paolo (2014) and De Jaegher & Di Paolo (2007) describe corresponds better to the kind of generative dependence between precarious transformation processes that, Bich and Arnellos argue, is captured by the more specific notion of organizational closure, defined by Varela as so:

We shall say that autonomous systems are organizationally closed. That is, their organization is characterized by processes such that (1) the processes are related as a network, so that they recursively depend on each other in the generation and realization of the processes themselves, and (2) they constitute the system as a unity recognizable in the space (domain) in which the processes exist. (Varela 1979, p. 55).§

It is this kind of closure, that Bich and Arnellos suggest, which captures the essential element of self-production identified in autopoiesis. This appreciation for the centrality of self-production, in terms of the mutual dependence of precarious elements, is not only a feature of autopoietically-informed accounts of autonomy but also of associated theories of what Letelier, Cárdenas, & Cornish-Bowden (2011) collectively term ‘metabolic closure’. It is the primacy accorded to metabolism that sharply distinguishes all of these above from the homeostasis-focussed approaches of Friston and Ashby, whose accounts are unable to capture specific requirements introduced with the kind of preservation that is dependent upon both continuous material turnover and phase transitions between different steady states – as I will discuss in sections 8.2 and 8.3 respectively

I plan to argue that a formulation of biological autonomy that is not subject to the kind of trivializability that we have met with in the discussion of the FEF depends upon this self-production that is not captured by operational closure. It is self-*production*, that is not only unformalized, but unformalizable in the free energy framework, and which distinguishes between a system that

is merely preserved, and an agent that preserves itself. In Chapter 9 I will discuss more recent attempts to formalize it beyond the molecular level of autopoiesis. First, however, I will now turn to the notion of *operational* closure and the interpretation of autopoiesis as a form of homeostasis to describe both how the FEF describes this and why bioenactivists have largely abandoned such an approach.

7.3. Operational Closure

Operational closure is defined by Bourguine & Varela (1992) as follows:

“A domain K has closure if all operations defined in it remain within the same domain. The operation of a system has therefore closure, if the results of its action remain within the system” (Bourguine & Varela, 1992, p. xii).

As stated, the inspiration for this definition is closure in the algebraic sense, where the range of possible outputs of some operation is a subset of its domain of inputs. While Bourguine & Varela’s transposition of this to the action of a system in the second part of the above quote sounds like the requirement that this system has no effects outside itself, this cannot be what is intended, particularly given their stress on avoiding the suggestion that any system is isolated from interaction with its external environment. Rather, the point is that in-so-far as these side-effects do not influence the state of the system itself, they are of relevance only to us as external observers. They have no bearing on the constitution of the system and its closure of operation, which must always bring the system back to the domain in which the operation began.

Maturana (1970) described this same property more explicitly as follows:

*Living systems as units of interactions specified by their condition of being living systems cannot enter into interactions that are not specified by their organization. **The circularity of their organization continuously brings them back to the same internal state (same with respect to the cyclic process).** Each internal state requires that certain conditions (interactions with the environment be satisfied in order to proceed to the next state. Thus, the circular organization implies the prediction that an interaction that took place once will take place again. If this does not happen the system disintegrates; if the predicted interaction does take place, the system maintains its integrity (identity with respect to the observer) and enters into a new prediction. (P.3) [my emphasis]*

So operational closure requires firstly, that the state of the system is at each point (at least partially) determined by its own operations, and secondly that as the range of potential states of the system is limited by its organization, so is the range of interactions into which it can enter. In Maturana and Varela's writing possible states are referred to (somewhat confusingly) as potential 'structures'. Whether these structures are possibilities, or not, is dependent upon whether they preserve the higher-order relational organization that defines the system as the particular operationally closed system that it is, even as its constituent parts may change. Just as the level of a mercury thermometer may rise and fall while it retains the organization that defines it as a thermometer, so a cell may lose and gain specific molecules that constitute it, whilst still retaining the same overall relations between its molecular network. As such, when treating autopoiesis as defined by operational closure Maturana and Varela declare that, "It is thus clear that the fact that autopoietic systems are homeostatic systems which have their own organization as the variable that they maintain constant" (Maturana and Varela, 1980, p80).

Like the pronouncements of the FEF, operational closure has the feeling of triviality. Such an abstract formal specification grants no insight into the methods by which ongoing existence is achieved, it makes no reference to the particular thermodynamic considerations peculiar to living systems, and it does not distinguish their capacity for the turnover of their material components from the changes in state of a common mechanism. The operations of all systems will also have some effects that feedback to the system itself, and all systems are limited in the range of states they can occupy. As an analysis of a living organism then, operational closure alone feels explanatorily impoverished.

I have already mentioned that operational closure is not intended to suggest that the operations of the system do not have any effect on the world outside said system. Rather, the point is that such effects belong to the domain of external observers, and are not a part of the constitution of the system itself (Maturana. 1975; Maturana & Varela, 1980/1970; Varela, 1979). As Maturana (1975) describes

Given a closed system, inside and outside exist only for the observer who beholds it, not for the system. The sensory and the effector surfaces that an observer can describe in an actual organism, do not make the nervous system an open neuronal network because the environment (where the observer stands) acts only as an intervening element through which the effector and sensory neurons interact completing the closure of the system. (Maturana, 1975, P.318)

Just as outputs are off-limits, so too, as Thompson (2007) explains is talk of inputs, at least ‘in the usual sense’. This caveat is particularly important with regard to the latter. While outputs that do not loop-back to affect the system are, by definition, irrelevant to a characterization of its intrinsic properties and detectable only by an additional observer, no system is isolated from any

environmental influences upon its state. Such ‘inputs’ do not require an external observer to detect, so what justifies their exclusion?

The answer is the concept of structural determination (Maturana, 1975) or self-determination (Varela et al., 1991), a corollary of operational closure that describes how, as far as the system is concerned, any environmental influence will appear only as “perturbations within the processes that define its closure, and thus no ‘instructions’ or ‘programming’ can possibly exist” (Varela 1979, P. 58). What the notions of operational closure and structural determination pick out is not an isolated system, closed off from any interaction with its environment, but rather a consequence of the particular descriptive perspective taken when we describe the internal constitution of the system in terms of its operationally closed organization. When the system is described in such terms, an environmental change does not ‘instruct’ the system by entailing a particular output but perturbs it, with the consequences being dependent upon the state of the system at the time of the input and how the systems closed operations work to preserve its fixed organization.

As Maturana (1970) describes, any event which does not merely perturb those operations but causes them to cease altogether (as in the case of fixing one coupled pendulum to a set position) would correspond to the breakdown of its organization and the destruction of the system. The FEF’s steady state EISA-cycle can be interpreted as a transposing of the concepts of operational closure and a fixed organization that may vary in state, into the language of random dynamical systems and statistical models of the same. Like the idea of ‘homeostasis of organization’, the FEF abstracts away from turnover of components, treating these as equivalent to changes in state that are compatible with the preservation of the operationally closed organization, specified by the EISA-cycle. In arguing that the FEF can serve

as a formalization of autopoiesis, Wiese & Friston (2021) thus draw specifically upon this characterization of the latter in terms of operational closure and homeostasis of organization.

In the stochastic differential equations that the FEF took to define the overall steady state system (Eq. 5-6, in section 6.2) we find a similar move to demote environmental influence to perturbation. This is the distinction between the EISA-internal variables of the deterministic component that make up x , define our ‘operationally closed’ system and determine its steady state dynamics, versus the noise terms (w) which are supposed to encapsulate all outside influence.

Like the FEF’s steady state EISA-cycle, this characterization of autopoiesis and autonomy in terms of operational closure and homeostasis of organization is an extremely minimal as a characterization of a living system. As I will argue in section 7.5 it is no longer popular as an analysis of autonomy, and is insufficient to serve the aims of bioenactivism.

Still, it is in the formulations of operational closure and structural determinism from Maturana and Varela’s writings in the 1970s and 80s, and in the FEF’s steady state EISA-cycle, where I think the paths of Friston et al. and the development of the bioenactivist viewpoint come closest to crossing. As I will argue in the next section, however, they quickly diverge.

7.4. From freedom and stability to dependence and purpose

As with the FEF's steady state EISA cycle, the properties of operational closure and structural determinism are not unique to living systems. But then neither did Maturana or Varela intend to take them as such. The autopoiesis of a single cell is an instance of operational closure, but that does not mean to say that this kind of closure exhausts all that is of interest in the autopoietic characterization of life. As Villalobos and Ward (2015) point out, the examples Maturana (1987, p. 73) cites in illustration of structural determination are not biological, but mechanical or computational cases like a washing machine or a lightbulb. Similarly, while operational closure may be easiest illustrated with classical dynamical systems examples, such as the Watt-Governor or the thermostat, computational systems can also be analysed in terms of a closed loop of interaction with their environment. A Turing machine, for instance, writes to and reads from the same tape, allowing its 'output' at one time to alter its 'input' at the next (Villalobos & Dewhurst, 2017).

So what of our hopes to naturalize teleological or intentional talk in the autonomy of living systems? If operational closure as a description of the homeostasis of organization were indeed all there is to autonomy, and this is satisfied by anything with a feedback loop through its environment, then we would have to attribute goal-directed behaviour to toasters and to teamakers. This is not, however, a view that should be attributed to either Maturana or Varela. The development of such an implausible position can only come from ignoring how their views developed and diverged, and, crucially, the mistake of conflating together the relatively trivial properties of self-determination and operational closure with the more demanding

requirements of self-*production* and the closure of precarious processes that are at play in more recent bioenactivist attempts to naturalize teleology.

In Maturana and Varela's earlier work their shared stance is explicitly anti-teleological. As Maturana describes the guiding tenets for his work in the (sole-authored) introduction to *Autopoiesis and Cognition: The Realization of the Living*, "notions of purpose, goal, use or function, had to be rejected" (Maturana & Varela 1972/1980, pxiii). Chapter II on the 'Dispensability of Teleonomy' is dedicated explicitly to this purpose, wherein they describe a perspective on the organism completely at odds with bioenactivism's phenomenologically-motivated concern for sense-making, immanent teleology, and intentionality:

...since the relations implied in the notion of function are not constitutive of the organization of an autopoietic system, they cannot be used to explain its operation. The organization of a machine, be it autopoietic or allopoietic, only states relations between components and rules for their interactions and transformations, in a manner that specifies the conditions of emergence of the different states of the machine which, then, arise as a necessary outcome whenever such conditions occur. (P.86)

As with Ashby then, the concern is to exorcise an internal goal-directed driver of behaviour, a primitive volition standing outside of ordinary causal entailment, from our explanation of living systems. If one is content with this view that there is nothing especially teleological or purposeful about an organism and believes that *if* there is even such a thing as 'original intentionality' it does not begin with the first autopoietic cell, then operational closure, homeostasis of organization, and indeed free energy minimization could serve as an adequate characterization of living and the non-living systems alike.

Varela appeared unable to remain content with this anti-teleological position for long. In *Principles of Biological Autonomy* (1981) he already begins to negotiate room for teleological talk as part of a pluralistic view of explanations, treating it as an alternative form of description that an external observer might validly use to better understand a system's behaviour – in much the same sense as talk of 'symbols', 'inputs' & 'outputs', or 'the environment' may be proscribed from an operational characterization of a system's intrinsic properties, while being perfectly legitimate at the level of communication about that system between external observers.

This is not yet the bioenactivism that I am looking for. The view that teleological explanations and intentional attributions may be legitimate in certain contexts is just as compatible with an instrumentalist characterisation that takes the difference between organisms and mechanisms to be one of complexity, and the appropriateness of such terms to be a matter of their success in abstracting from this detail to yield the kind of predictive regularities that interest us.

By 1999, however, Varela has been influenced by Kant and Jonas' work on the idea of organisms as 'natural purposes' and has begun exploring the concepts of 'original intentionality' and 'sense-making' as unique to life, coming around to the position that these do lead to the re-introduction of a kind of teleology that is "intrinsic to life in action" (quoted from an email exchange in Thompson, 2007. P453-454). This about-turn culminates in a 2002 article with Andreas Weber, that marks the definitive break with Maturana (and Ashby's) attempt to treat living and non-living systems as purposeless entities alike.

As Varela and Weber (2002) argue, drawing extensively on Jonas (1973) the distinguishing aspects of the living can be stated as follows:

1. *it exchanges its matter and acts thereby from a subject pole partially independent of the underlying matter,*
1. *as precarious existence it is always menaced by concern (Sorge), the need to avoid perishing, and to do this, it is again completely dependent on matter whose characteristics are the reason for its concern.*
2. *already the simplest forms of life have thus a subjective perspective as a result of this existential need.*

Therefore

3. *life as such will always be captured in the antinomies of “freedom and necessity, autonomy and dependence, I and world, relatedness and isolation, creation and mortality”. (P.113)*

These are not consequences of operational closure, homeostasis and structural determinism, properties that autopoietic systems share with systems in general, but rather express the ways in which an organism’s self-producing character creates a difference in kind between the living and non-living. For Jonas, the key to all of these properties is an understanding of precarious self-production through metabolism – the dependence of living systems upon continual flows of matter in order to preserve and rebuild their precarious organization. Unlike in homeostasis, the key feature of metabolism is not just preservation *in spite of* material turnover, where this is presented as perturbation within a homeostatically-conserved organization, but precarious *dependence upon* that turnover. Or, as Jonas (2001/1966) nicely puts it, the distinguishing feature of the organism is its relationship of ‘needful freedom’ with matter.

For Weber and Varela it is this needful freedom of metabolism, as opposed to the mere freedom of homeostasis-of-organization, that is now presented as the key feature of autopoiesis. Where homeostatic systems are fundamentally passive ones that act only in response to external disturbance and remain stable *in spite* of change, a metabolic self-producer is necessarily active and *dependent* upon change. It is this that makes a network of chemical reactions a self-producer, and it is this precarious dependence that is key to the bioenactive naturalization of goals, intentionality, and immanent teleology.

The key point here is that, as restrictions on how a system can change *if* a change in structure or state occurs, neither homeostasis, operational closure, nor the steady state EISA-cycle describe something that is *dependent* upon change for its existence. All that is necessary for a system to be an EISA cycle is the conditional requirement that *if* one component changes state then it must be the result of a change in another component of the appropriate set. So *if* an internal variable changes, then there must have also been a change in either another internal, sensory, or active variable that caused it – and these changes must conform to the dynamics of our steady state equation so as to preserve a stationary probability distribution. But, crucially, there is no absolute requirement for ongoing change to preserve either the individual components or the overall EISA-cycle organization.

Take the coupled pendulums when they are at rest. In so far as their average state remains fixed at equilibrium, they will satisfy a steady state requirement for a stationary probability distribution, used as the generative model of active inference. At rest, they have minimized their surprisal quite perfectly and, should they happen to be perturbed by external noise, they can be expected to minimize the resulting unlikeliness of being out-of-equilibrium

with perfect alacrity. In so far as the position of one pendulum cannot directly change the velocity of the another without changing the state of the beam, their potential patterns of influence will also conform to that of an EISA-cycle and the beam will qualify as a ‘Markov-blanket’ between them. Still, the beam itself, like all other parts of the system, is perfectly self-sufficient. It will neither crumble nor dissipate if pendulums cease to oscillate.

So, the kind of system described by a steady state EISA-cycle is one where we have some fixed constraints that determine the dynamics of the system, but where these constraints do not depend upon those dynamics in turn. The fact that the beam does not depend upon the movement of the pendulums can seem hard to miss, and yet we find supporters of the free energy framework claiming

“The Markov blankets are a result of the system’s dynamics. In a sense, we are letting the biological systems carve out their own boundaries in applying this formalism. Hence, we are endorsing a dynamic and self-organising ontology of systemic boundaries. (Ramstead et al., 2019. p. 3)

This is presented as describing the same existential dependence between dynamics and structure as found in the cell-membrane, but as we’ve seen with the coupled pendulums, this is not a general feature of Markov-blankets or EISA-cycles. The cellular membrane is an intrinsically unstable configuration and if the internal workings of the cell do not act to replenish this boundary, it will dissolve. The same is not true for the blankets in Huygens pendulums, nor for the coupling between the Watt-governor and the steam engine.

A probabilistic graph and its attendant Markov blanket merely describe how a system’s structure constrains its dynamics, they do not mandate any

reciprocal dependence of this structure upon those dynamics. At best the relationship we might have here is an epistemological one, in as much as the dynamics reveal the independencies of a boundary. Still, given that conditional independence is everywhere, and those particular boundaries the FEF chooses to focus on are typically things that we have picked out by means prior to any measurements of statistical relationships, it does not seem a particularly informative clue to the selection of boundaries.

7.5. Self-production is not 'homeostasis of organization'

There are, as Di Paolo, Thompson, and Beer (2022) note in their enactivist critique of the FEF, other reasons to reject the deflated account of autopoiesis-as-homeostasis. For one thing, as I will describe in Chapter 8, a distinguishing feature of organisms is the difficulty, and I will suggest the impossibility, of identifying any invariant organizational features that both individuate a particular organism and which *must* be preserved throughout their unprestatable and open-ended development. Another, more basic problem with this interpretation, as (Di Paolo, 2005, 2010) describes, is that homeostasis allows for *variation* around a stable point, whereas autopoiesis is a binary property that does not admit of graduation. One is either a self-producing network or one is not, and once this breaks down it's too late to do anything about it.

For the purposes of explaining why the FEF, cannot serve the aims of a bioenactivist the main point is nicely expressed by Mossio and Bich (2017), with regard to the failings of homeostasis-based accounts of the organism more generally.

.... it presupposes the existence of the organisation that under certain circumstances it contributes to maintain stable. In particular, homeostasis does not capture the most distinctive generative dimension of biological organisation, i.e. the fact that the components involved in feedback loops are not only stabilised, but produced and maintained by the very organisation to which they belong. In a word, homeostasis misses precisely self-determination (P.1096)

And, as they continue...

Technically, the “goal” of a homeostatic mechanism is defined as the interval within which the mechanism maintains the target variables. Yet, it does not make any difference from the point of view of the definition whether the interval is extrinsically established by a designer, as in the case of artefacts, or intrinsically identified with the conditions of existence of the system, as in the case of biological systems. Both cases can pertinently be said to be homeostatic. However in failing to account for their difference, Cybernetics misses the crucial dimension of biological teleology. (P. 1096)

A homeostatic description can only come after we have determined some set of variables, and the state at which they must be preserved. As such, it is indifferent to the means by which whether these essential variables and their bounds of viability are determined – whether by an external designer, or by something intrinsic to the homeostatic system itself.

The adoption of a Jonasian approach to living systems in Weber and Varela (2002), subsequently carried through in Thompson (2007) and Di Paolo, Buhrmann & Barandiaran (2017), represents a decisive break with the attempt to treat autonomy as a generic property of operationally closed systems, to reduce autopoiesis to homeostasis, and to treat organisms in the same terms as any other physical system. Awareness of this diametric opposition between Maturana (and Ashby) on the one hand, and bioenactivism on the other, is essential to a coherent understanding of either (Villalobos, 2013; Villalobos & Ward, 2015). It is not that the former neglect this intrinsic teleology, such that their work can be straightforwardly

supplemented by an account of it. In defining autonomy in terms of the generic properties of homeostasis of organization & operational closure, Maturana and Ashby debar such a possibility. Lumping these views together, as the free energy literature has tended to do, creates an incoherent frankentheory, which no amount of mathematical stitching can hold together.

If the FEF's steady state EISA-cycle is only a description of a specific form of operational closure, then it can hardly serve the bioenactivist as an analysis of the features in virtue of which living systems are 1) distinguishable from the non-living and 2) sources of proto-intentionality. But I don't want to just reject the FEF's analysis of the organism for failing to live up to my prior commitments to the bioenactivist viewpoint. A FEF-theorist might respond that Maturana and Ashby had the right of it, that the attempt to treat organisms as unique and goal-directed is misguided and that they are just one physical system among many, where a vastly greater degree of complexity combined with our inability to fully comprehend this, leads us to mistakenly attribute a difference in kind. Or they might suggest that the steady state EISA-cycle be taken only as a partial, necessary but not sufficient, set of conditions for a definition of biological self-organization. Alternatively, they could propose that the FEF itself should not be understood as a theory about organisms at all, but rather a set of mathematical tools, a statistical redescription of dynamical systems theories that might then be used to formulate a description of those features that *are* specific to organisms.

None of these responses will work. As I will argue in the next section the FEF's problem is not merely that formalization of a homeostatic identity is too generic to distinguish the particular features of autonomous biological agents. The problem is that the principles of stability that it presents as

necessary are exactly those principles that biological systems are uniquely prone to violate.

8. A theory of everything, or just of every ‘thing’?

Now, at last, we have a clearer understanding of how the FEF defines the existence of a system, or, in Friston’s (2019) terminology, every ‘thing’ to which the free energy principle is supposed to apply, in terms of the coupled equations of the steady state ESIA-cycle. With this in hand, we can extract the two assumptions that the FEF makes about a system.

- 1) Tendencies: The parts of the system tend to revisit the same state, or to cycle through the same set of states, with a frequency that does not change over time.

- 2) Dependencies: The interactions between these different parts of the system do not change.

To connect this to Bayesian inference the FEF then redescribes these dynamics in terms of an invariant, low entropy joint probability distribution over the state of all of the parts of the system – which gives us the ‘generative model.’ The partition between external, sensory, internal, and active variables now becomes a statistical one – the Markov blanket. The FEF then uses this partition to interpret internal variables as ‘inferring’ external ones, in as much as the statistical properties of the internal variables converge with those of external variables, despite their being independent of one another when conditioned on the Markov blanket of sensory and active variables.

So, there are two key moves in the free energy framework: a putative definition of every ‘thing’ that the principle is supposed to apply to, and the attempt to use a statistical redescription of this to license the attribution of inferential properties. The problems with this latter move, which is essentially a slide from correlation to representation, have been well discussed and I described these in section 3.3. Yet as Raja et al. (2021) note the first move upon which it depends has too often been granted a free pass, despite the highly unlikely assumptions it makes about those entities to which the free energy principle is supposed to apply.

It’s this state of affairs that has led to the FEF’s being either lauded as a ‘first principle’ or as decried as ‘unfalsifiable. *If* this proposed definition of a system holds (along with some further technical specifications, for more on which see Biehl et al, 2021 & Aguilera et al., 2021) it entails that the system will be formally describable in terms of free energy minimization. In this sense Friston is correct to claim that the free energy principle itself is not an empirical hypothesis – for all that specific hypotheses about cognitive architecture, such as PP, may be derived from it – hence frequent objections raised to its lack of falsifiability miss the mark. Yet, the FEF should not be taken as merely the working out of a tautology of existence either for all that it has been presented as such. This would presuppose that the FEF’s steady state EISA-cycle was already an accepted definition of ‘existence’ biological and otherwise. Rather than a tautology then, Hohwy (2020) proposes that the FEF is best understood as an attempt to address this definitional gap through its putative account of what it is to be a (self-organizing) system.

The appropriate test of a proposed analysis is not whether falsifiable experiments can be derived from it, but that doesn’t mean it is immune from criticism and counterevidence. Instead, we look to how well it accords with

both common and scientific practice. Neither the FEF, nor any other piece of formal analysis, is required to submit entirely to the absolute authority of this bicameral legislature – as though either chamber were even capable of producing a unilateral ruling in the first place. A philosophical analysis may worry the fabric of our linguistic habits, to expose incoherences in the everyday applications of a concept. It may likewise criticize scientific terminology for losing touch with everyday use. If, however, the proposed criteria result in extensions of a concept that are diametrically opposed to both everyday intuitions and to our best scientific understanding, then the only conclusion is that the philosopher (or neuroscientist) is talking about something else entirely.

I have already described one problem with the FEF's definition of a system: namely that in attempting to treat both animate and inanimate entities in the same terms it rubs up against the folk understanding that there is a difference in kind between the existence of a person and that of a pendulum. In treating both as nothing other than steady state EISA-cycles, the FEF fails to account for why we tend to talk of one as an agent and pursuer of goals, and the other as a mere mechanism. As such it threatens to lead us either into deflating instrumentalism about intentional talk, such that it depends only on the greater complexity of the former's behaviour, or to inflated panpsychism in which every system in a stable coupling would count as an inferring agent.

Neither option works for the bioenactivist viewpoint, but that is not a good enough reason to reject the FEF. Physicalism and mechanism are still the order of the day in most respectable scientific circles, who are liable to turn up their noses at any whiff of 'vitalism'. A theory that treats the apparent gap between life and non-life as an illusion that reduces to the complexification of ordinary mechanisms may be more likely to appeal.

Friston, and others working on the FEF, have begun to develop an account of what this complexification might look like. Suggestions include the incorporation of a solenoidal component to the system's dynamics resulting in a limit cycle attractor, rather than fluctuations around a single point (Friston, 2019b), or hierarchical extensions that allows for temporal or counterfactual depth in a system's predictions, wherein actions are selected to minimize the long-run average of free energy over an entire trajectory (called expected free energy) rather than just to minimize immediate free energy (Wiese & Friston, 2021; Friston et al., 2020).

I will not consider any of these in detail here, because I believe that the FEF has bigger problems than its generalizability. Any attempts to fix its broad applicability by casting the organism as a 'special case' of a steady state EISA-cycle will fail. As I will argue, the problem with the existential imperatives by which the FEF defines a system is not their generality, but their *contingency* when applied to biological forms of existence. The stationarity of a joint probability distribution, and the stability of tendencies and dependencies that it implies, may do well enough as necessary requirements for the continuing existence of an inanimate substance. Living systems, in their temporary coincidence with a particular physical body, may even happen to meet these two stability requirements over some duration. But as necessary principles meant to define an organism's ongoing existence they are both false.

This is not just a flaw of the FEF and its conceptualization of the organism as a homeostatic mechanism with stable behaviours and stable parts. Instead, this misconception stems from a broader ontological framework that views organisms as substances, and specifically, as machines. But the mechanical and the substantial do not exhaust all possible modes of existence. Unlike

machines, organisms may change both their components and the rules that govern the behaviour of these parts in unprestatable ways, leading a number of philosophers, theoretical biologists, and complexity theorists to propose that they are much better captured by a processual ontology.

The stability of parts and properties may well be a trivial property of inanimate objects but what distinguishes living systems, I will argue, is precisely their intrinsic instability in both behaviour and material constitution. This is a difference in kind between the living and non-living, one that can be identified independent of any folk, or bioenactivist, commitment towards attributions of biological autonomy, and which is enough to show that the FEF cannot serve as the basis for a theory of living systems – bioenactivist or otherwise.

8.1. Processes and substances

The FEF presumed two sorts of stability in our system. The first: invariance of interactions between our parts, such that even as states of variables may change, the statistical dependencies *between* these³² remain the same. The second: invariance of the statistical tendencies of each individual part. This allows that the state of each part may change, but requires that if it has been in a particular state eight out of ten times previously, then it must also spend eighty per cent of its time in that state in future. This second requirement has been justified in increasingly general ways: as a formalization of biological homeostasis in particular; as a mathematical description of stability through

³² As discussed in the section on the metaphysics of Markov blankets, I am allowing the conflation of causal relationships with statistical ones here, based on the argument that in some complete and comprehensive causal network the former will reduce to the latter.

perturbation more broadly; and, most generally of all, as a principle of all existence.

Friston (2019a) motivates this by defining a ‘thing’ as something that is ‘distinguishable in a statistical sense’, cashed out in terms of having states characterized by a low entropy probability distribution that remains constant over the duration of said thing’s existence. The idea being, crudely put, that if a thing did not regularly revisit the same states but rather wandered off to ever new regions of possibility then how would we reidentify it? Or, as he puts it:

... nearly every system encountered in the real world is self-organising to a greater or lesser degree – suggesting that self-organisation is, in itself, unremarkable. Put another way, if systems did not self-organise they would have dissipated before we had a chance to observe them. (Friston, 2019, P. 24)

If this were the case, if our ability to observe and identify a system over time were dependent on its exhibiting this kind of stability, then we would not need recourse to the specific importance of homeostasis to an organism to make the case that they are steady state systems. This condition would instead be derivable from this more general requirement for what it takes to exist over time.

But is this really a requirement for everything we can identify as persisting over time? It seems to me that there are plenty of things that we can reidentify in spite of their violations of any steady state condition. Songs have their choruses, stories have their tropes, and dances their motifs, but we can also recognize a continuity throughout the screwball transitions of Bohemian Rhapsody, we understand that when Odil ends the 32 fouettes of the Black Swan Pas de Deux to return to her partner that the dance goes on, and as the timer ticks down on the bomb handcuffed to Sean Connery’s wrist we know

that this is a single ongoing moment of peril, part of a broader, non-repeating sequence that constitutes the overall narrative unit that is *Goldfinger*.

Moreover, we can continue to recognize this continuity even as the realizing parts change: as a melody moves between members of an orchestra, a relay race moves between runners, or as a river churns through water molecules. Friston (2019, 2013) takes such material churn to be incompatible with our need to identify the stable parts whose stable interactions give us a causal graph as a real entity – though as we will see in the next section it is not clear that he is correct to do so. Given that inanimate objects may also persist through an exchange of parts, then if the FEF cannot deal with this it would presumably raise an issue for Friston’s aspirations to a theory of the non-living also.

With regard to the violations of steady state, it would be reasonable to respond that *Goldfinger*, *Bohemian Rhapsody*, and *Swan Lake* are not really objects, but *processes*. That repetition may be only one form of continuity in a process, but the FEF is a theory of ‘things’ only in the narrower sense of a persisting substance, or more specifically, of a mechanism. As Wiese & Friston (2021) put it, these ‘things’ are:

Systems that exist over some appreciable timespan, in the sense of having an attracting set, i.e., revisiting the neighbourhood of characteristic states despite external perturbations—of which living systems are a subset. (P.7)

Still, if this means that the FEF is not truly a theory of every possible form of existence, then we have no *a priori* reason to accept that it characterizes the kind of existence that organisms have. To motivate this we would need to first argue that organisms are indeed substantial ‘things’ rather than processes, and secondly that the continued existence of a substance, or at

least of a biological substance in particular, depends upon the particular forms of stability that the FEF describes.

It might seem obvious that a cat is more like the ballet shoe than the ballet. Unlike *Swan Lake*, you can pick a cat up and throw it around – though it seems inadvisable. Yet while something like mass and extension (as necessary pre-requisites of pick-up-and-throwability) may be part of the ordinary conception of ‘substance’, throughout the history of metaphysics this term has been used in a variety of distinct technical senses (Morgan, 2021). This lack of one agreed-upon definition can make it somewhat difficult to assess the prospects of a ‘substance ontology’ for organisms in general. For our purposes in evaluating the FEF, however, the most apposite characterization of a ‘substance’ is what Morgan terms the ‘essentialist’ account, in which something is defined by a particular set of properties that persist throughout its existence, allowing it to undergo only those changes that do not violate these properties.

This essentialism is the notion of substance at issue in recent criticisms from defenders of processual accounts of organisms (Nicholson, 2018; Dupré and Nicholson 2018; Meincke 2019). As Dupré and Nicholson (2018) state, the key claim they reject is that organisms should be understood as things which persist “by virtue of their continued possession of certain essential properties, which make those things what they are and which remain unchanged over time.” (P. 24). In contrast, the processual account looks not for the essential features of X, such that they could identify it in any randomly selected time-slice, but asks instead, as “how should I follow X through time” (Guay and Pradeu, 2016). Rather than by an atemporal set of features, a processual identity is individuated by what Lewin (1922) terms relations of ‘genidentity’ between each successive time slice, where the latter is a generative

product of the former. Such a relationship might, for instance, be described in terms of autopoiesis, where what makes one organism-slice a continuation of the previous one is that its chemical components are the product of the synthesizing activity of the prior set.

The processualist may allow that said process can be contingently ‘stabilised’ in a particular material ‘cat-substance’ that we can stroke or hold, but must take that particular substance as a temporary episode in the more fundamental process necessary for being a cat. Moreover, rejecting the idea of organisms as essentialist substances is not to prohibit an organism from having *any* invariant properties, rather it is to claim that those invariant properties are not what individuate it as the particular organism that it is. If something is a persisting substance then, when viewing it at two temporally-disconnected instances, it should still be possible to reidentify it as the same thing without knowing anything about what went on in the intervening period. If it is not, then the only way to answer this question would be to follow the entire temporal trajectory from the first moment to assess whether it connects up to the second in some, as yet unspecified, ‘right way’.

Just as there can be different accounts of the ‘right way’ for a processualist to follow an organism through time, so the substantialist’s requirement of an essential invariant still leaves open a wide range of different substance ontologies that might be offered for the organism. The free energy framework, as I have described, seems committed to two essential invariants, 1) the stability of the parts that are taken to literally instantiate the real causal graph of the EISA-cycle and its attendant ‘real’ Markov blanket, and 2) the stability of their behaviour which gives us the fixed steady state function and associated, invariant generative model.

Friston (2013, 2019), as mentioned, often talks about the first of these, the stability of parts, in terms of the stability of material components. The stability of some fixed aggregation of material stuff is perhaps closest to the ordinary concept of substance, as when we talk of a substance as something with pick-up-and-throwability or when we say, “there’s a mysterious sticky substance all over the baby’s high chair. An inability to account for changes in this material realization would certainly raise a problem for the FEF in light of the continual recycling of parts that occurs in all organisms. Indeed, given that we can change a graphics card, set of wheels, or planks without creating a new computer, car, or ship each time we do so this would raise problems for the FEF’s ability to serve as a theory of anything. As we’ll see in the next section, this requirement of stable material parts is a strange thing for Friston to have committed himself to. The FEF might easily abandon this to talk, as most machine-style models do, in terms of the stability of formal parts.

This comes with two problems, however. The first is that in abstracting away from material turnover to focus on the stability of formal parts we erase the difference between a structure that is stable *in spite of* possible material exchanges, versus a structure whose stability is entirely *dependent* upon ongoing material turnover. This distinction, as already mentioned, is crucial to the bioenactive, or Jonasian, conception of the primitive intentionality of organisms.

The second problem is that this strategy relies upon a different type of essential invariant: namely an invariant organization, describable by some mathematical equation which remains fixed even as the states of its variables change. It is this that allows us to individuate a formal part in terms of the role it plays in this equation, even as its material realization changes. As we

will see in the second half of this chapter, however, organisms are unique in their ability to persist not only through material turnover, but also through radical changes in organization. As I will argue, there is likely no level of abstraction at which we can identify an invariant equation that is both specific enough to individuate this particular organism and flexible enough to allow us to derive every change it might possibly undergo in the course of its lifespan.

8.2 The instability of organic parts

The fact that organisms are continuously replacing their components by interchanging matter with their environment appears to pose a problem for the requirement of fixed parts, and more generally to mark at least one quite fundamental difference both from machines, and from other substances in general. To find non-living systems with the same property we have to look not to the mechanical, but to phenomena such as tornadoes, whirlpools, rivers and candleflames – in other words, the archetypal cases for a process ontological perspective.

Friston (2013, 2019b) explicitly cites such material turnover as incompatible with how the FEF defines a system, but then quite bizarrely presents this in support of the framework's particular suitability for living organisms. A principal reason for taking a single-cell and its membrane as the canonical Markov-blanketed free-energy minimizer, he claims, is the very stability of its components, contrasted to a candle flame which:

... cannot possess a Markov blanket, because any pattern of molecular interactions is destroyed almost instantaneously by the flux of gas molecules from its surface. Meaning we cannot identify a consistent set of blanket states rendering some internal states independent from other states (Friston, 2013, P.2)

And as he repeats, regarding his most extended presentation of the free energy framework to date:

... it does not easily accommodate the fact that the particles that constitute a Markov blanket can, over time, wander away or, indeed, be exchanged or renewed. The canonical example here would be the blanket states of a candle flame, whose constituent particles (i.e., molecules of gas) are in constant flux. “ (Friston, 2019a, P.50)

To cite the membrane as a point of contrast is an odd choice, given that its constituent parts are continuously consumed and regenerated by the cell's metabolic network. Membrane turnover via endo- and exocytosis is a means for all sorts of self-organizing behaviour, from regeneration and growth to the transportation of molecules between the interior of the cell and its environment. In the cellular slime mould *Dictyostelium*, for instance, membrane turnover has been proposed as a mechanism of locomotion, with estimated times for complete turnover in the order of 4-10 minutes (Aguado-Velasco & Bretscher, 2017). Within the cell too, amid the ‘internal states’ presumably realised by particular proteins, we find turnover times much shorter than the lifespan of the overall system – on the order of around two days for a non-dividing mammalian cell (Toyama & Hetzer, 2013). And this flux continues up to the multicellular level where, in the human body, for instance, there is an estimated daily turnover of around or 0.2 per cent of total cellular mass (Sender & Milo, 2021). As the physiologist, John Scott Haldane, one of the earliest ‘processualists’ of 20th century biology describes, “the organs and tissues which regulate the internal environment . . . are constantly taking up and giving off material of many sorts, and their

“structure” is nothing but the appearance taken by this flow of material through them” (Haldane 1917, P. 90).

Such a porous and protean thing is much more like a candle flame than it is like the EISA-cycle’s fixed patterns of interactions between fixed parts.³³ If the basic units of our causal graphs are the states of token particles, as Friston (2013) takes them to be, then their statistical dependencies will lack the stability necessary to establish patterns of conditional independence between them, and for the identification of a Markov blanket between unchanging sets of ‘internal and ‘external’ components. In the cell, a previously ‘external’ molecule is free to waltz right through its membrane to start interacting directly with an ‘internal’ one, blithely violating the FEFs basic definition of a system.

This is especially problematic for a realist about EISA-cycles and their Markov blankets, who treats them as something that the organism literally instantiates and which makes it what it is. Such a realist will find that the components of any ‘real’ Markov blanket they identify around an organism will dissipate on timescales that are shorter than that lifespan of the organism whose ‘very existence’, they claim, ‘depends’ on that boundary’s preservation (Allen & Friston, 2018).

Still, one might think that what matters is not the stability of interactions between component particles, as Friston (2013, 2019a) seems to take it, but rather the stability of the higher-level organization. As noted in Chapter 5, the realist about probabilistic graphs need not hold that they are instantiated at the level of interactions between the states of particular token particles.

³³ As Proksch (2021) notes, this requirement of stability is equally difficult to sustain in cases where Markov blankets are supposed to serve as the boundaries of social networks, given that the turnover of members is a feature of most social organizations

We might instead take features of the cell's macroscopic organization, such as intracellular and extracellular glucose concentration, and describe how one cannot affect the other without a change in the state of transmembrane channels. The movement of a particle would thus correspond to a change in the state of some more macroscopic fixed node, rather than the breaking and creating of new connections in a particle-level causal graph.

The idea that we can abstract away from material turnover is, after all, the key idea behind the notion of homeostasis of organization that Wiese and Friston (2021) attempt to use as a bridge between autopoiesis and the free energy theory. The point of defining the fixity of the organism, as Maturana and Varela (1980/1972) did, in terms of organization rather than the variable material instantiations of this, is that it allows us to talk of an organism like a machine, in terms of its fixed formal parts and the rules governing their behaviour.

This machine-substance view of what it means to have fixed parts not only accommodates material turnover in the organism but arguably better captures how this requirement is supposed to apply to machines too. Machines may not exchange their material components as a matter of course, but they do admit of such exchanges. Unscrupulous salesmen aside, when we repair and replace the wheels on a car we are not inclined to say we have created a new machine. While having wheels is a fixed invariant of a car, having some particular set of wheels is not. From the perspective of the organization level, the old wheels and new wheels qualify as the same formal part. As described in the previous section, the fixity of material realization is only one possible view of what the essential invariant of a substance is. In both the organism and the machine case, however, the substantial invariant

is better captured not by fixity of material properties and parts, but by fixity of formal ones.

So, we cannot take our parts for granted as something like specific atoms or molecules. To construct a fixed graph that abstracts away from this constant turnover of microphysical entities we need to describe the invariant form of our particular living system such that we can identify its fixed formal parts. These parts must be individuated prior to being able to create a graph of the connections between them. As such, neither the EISA-cycle nor a Markov blanket is of any use in individuating them, for in order to construct these we need to already have divided our system into fixed units, such that we can then assess any relations of dependence or independence between them.

There are bigger problems for the FEP here than its inability to serve as first principles analysis and its requirement for some prior specification of the fixed organization of the organism, from which the parts needed for its analysis in terms of generative models and causal graphs might subsequently be derived. As I will discuss in the second half of this chapter, organisms, unlike machines, not only persist through material turnover but also radical transformations of organization. Such transformations, I will suggest, cannot be described in advance by *any* invariant set of equations that are both sufficiently flexible to derive these transformations and specific enough to individuate that particular organism. Without this, it seems unlikely that we can pin the organism to an invariant set of formal parts any more easily than we could tie it down to material ones.

Putting this aside, temporarily. Let's presume that we *could* give an account of how to identify fixed formal parts for any organism amid its material turnover, in order to treat it like a machine and so to redescribe it in the

FEF's terms. Even if we could, as I will argue in the next section, in doing so we miss out on a fundamental distinction between the organic and the mechanical in terms of how the former alone *depends* upon this material turnover for whatever temporary stability its formal parts may have.

8.2.1 Why metabolism matters

To say that organisms, like machines, are substances in the sense of having an invariant organization and formal parts allows us to abstract away from material turnover, without outright denying its occurrence as Friston (2013, 2019) appears to do. But should we abstract away from this turnover? We've already seen that the bioactivist, following Jonas (1953), would be inclined to protest such an abstraction as erasing the distinction between a formally-defined machine that can admit of such exchanges, versus a precarious organism whose existence *depends* upon them. Even those who reject the idea that this dependence underpins any special teleological status for the latter must acknowledge that at least some operations of biological self-preservation, such as metabolic repair and regeneration, cannot be reduced to information processing and syntactical transformations alone.

This is not a problem if we take the EISA-cycle instrumentally as a useful tool for modelling specific behavioural phenomena, such as the regulation of body temperature. But if our model is supposed to provide the basis for a general theory of life, as Friston (2013, 2019) presents the FEF, then to acknowledge that it, like all models, is partial and distorted is not sufficient. The task of a model of 'life in general' is to highlight the *right* things, and neglect only those contingent features of the particular instances we happened to have encountered.

That all living systems we know of are metabolic systems, and there are arguably no naturally arising non-living metabolisms, is a reason for thinking metabolism might be important, but not a conclusive one. All living systems that we have encountered are also made up of a specific set of amino acids, but to take these to be essential features of life would be chemically chauvinistic, unless there is a principled reason to claim that *only these specific amino acids* are capable of realizing some general property that could reliably distinguish life from non-life in, say, some potential astrobiological encounter.

There is good reason for thinking that metabolism should be a criterion in such a test, quite independent of any putative connections to a bioactive account of intentionality and agency. Metabolism, properly understood is not merely an additional and disconnected capacity of an organism, nor does it just mean something that ‘depends upon energy’ as could also be said of a computer game. In the strong sense, as Boden (1999) argues it should be understood, metabolism describes a different mode of existence from that of a machine. The difference between something whose physical body is constituted via its own activity, versus an object that persists independently of its own doing.

Matter, in general, prefers to occupy a low-energy configuration. The form of an atom and the locations of electrons in their orbits is, as Bickhard (2009) notes, the paradigm example of how this tendency shapes its organization. In such equilibrium, things and their organization will remain stable as long as the ambient energy is not sufficiently high to destabilize and destroy their cohesion. This makes for good mechanical parts. In a living system, however, things are exactly reversed. Biological components often occupy inherently unstable configurations and the continuous input of energy is needed to

preserve them at this non-equilibrium steady state (Bickhard & Campbell, 1993; Schrödinger, 1944). No engineer would select such parts to build her machines.

In metabolic systems, the relationship between the structure of a system and its activity, or between the constraints it embodies and the dynamics they produce, is fundamentally different from in machines. To a machine, the flow of energy that facilitates its operation is a threat. Inefficiency in how the machine channels this energy allows it to gradually degrade the machine's otherwise stable parts. In a living system these structural parts – membranes, enzymes, etc. – are inherently unstable. Unlike in an atom, the stability that they do appear to have is not intrinsic to their internal configuration but is reciprocally dependent upon the activity they enable which secures and channels the matter and energy that is necessary for the continued replenishment that provides this contingent stability (Montévil & Mossio, 2015; Mossio & Moreno, 2015).

It's important to be clear that what I'm claiming differentiates organic parts here is not just that they are intrinsically unstable, but that the particular kind of stability that they do have is extrinsic to that part. They are stable only because of interaction with the broader network of the organism and how the instability of some of these parts fuels the activity of other parts to synthesise their replenishment – thereby giving them a contingent, mutually interdependent stability. In contrast, where radioactive materials are also unstable and continuously decaying, the partial stability that a Uranium-238 isotope has is merely due to the balance of fundamental forces at the subatomic level *within* that atom – there need be no additional inflow of energy or matter to sustain that stability and its breakdown does not power any process that replenishes it.

So organisms are not merely homeostatic mechanisms that remain stable despite change, their stability is dependent upon change. As Jonas (1953) criticized the cybernetic approach that underlies the free energy framework, “A feedback mechanism may be going, or may be at rest: in either state the machine exists. The organism has to keep going, because to be going is its very existence” (P.12) With my bioactivist hat on, I, like Jonas, would like to describe this in terms of the the intentional directedness of its actions towards its environment as a source of matter and energy required to achieve the immanent teleology of its own continued existence. But for the avoidance of controversy, we can just call it the dependence of a non-equilibrium structure on the continuous flow of energy required to sustain it. The fact that the existence of an organic structure *depends* upon on material turnover, as opposed to merely allowing for it, is a difference between organic, and non-organic existence that our theory of living systems needs to recognize – even if we have no interest in the bioactive project of interpreting this in intentional terms.

The FEF’s fixed statistical/causal network and its attendant Markov blanket describe only how a system’s structure constrains its dynamics, they do not address any reciprocal dependence of this structure upon those dynamics in turn. Just as a steam engine does not need to be in constant operation to continue to exist, and just as a laptop may be turned on and off again with no deleterious consequences to its future capacity for computation, so once Huygen’s coupled pendulums wind down, the connecting beam remains as a constraint on possible interactions should they be perturbed again. Active or not, the pendulums still meets the limitations on potential interactions required for an EISA-cycle, still achieve a free energy minimum, and still possess the beam between them as a Markov blanket.

Nicholson (2018) takes this presumption that an organism can be separated into an invariant structure and dynamical behaviour to be the fundamental error of what he terms the ‘machine concept of the organism’, and one nicely captured in the common analogy of food to fuel. As he puts it.

The problem of equating fuel with food is that it drastically underestimates the physiological pervasiveness of metabolism. No matter how dynamic a functioning machine may be, it is always possible to distinguish the machine's physical frame—which remains fixed—from the materials that flow through it. The actual structure of the machine does not itself take part in the chemical transformations that the fuel undergoes as it passes through it. Instead, it serves as a channel that facilitates the exchange of materials as fuel is converted into waste. An organism, in contrast, changes wholly and continuously as a result of its metabolizing activity... This is why the fuel-food analogy is so misleading, and why the stability of a machine—despite its apparent dynamicity—ultimately resides in an unchanging material structure. In machines there is a specific 'inflow' and a specific 'outflow'. In organisms everything flows. (Nicholson, 2018, p.145)

It is trivial to note that all machines depend on energy to operate and we could easily create a machine whose operation is dedicated to the harvesting of it – a sun-tracking solar panel would suffice. But the solar panel’s structure does not depend on its success. It need not store energy to make it through a foggy day and it will not disintegrate if it runs out.

What explains the persistent tendency to abstract away from the dependence on material turnover metabolic systems in theories of the living? Boden (1999) argues this tendency may stem not so much from a principled position as to its irrelevance, but rather from the desire to separate logical form entirely from material instantiation and identify the ‘essence of life’ with the former, such that we might hope to simulate, or even create it, in a virtual medium-independent form. Where other proposed criteria of life, such as growth or adaptation, do not make explicit reference to flows of matter and

energy, metabolism is defined in terms of these, and so cannot be straightforwardly captured in purely syntactic or informational terms.

The only reason for proposing that we drop metabolism from our concept of life is to allow a strictly functionalist-informational account of life in general, and A-Life (artificial life) in particular. The same applies in respect of suggestions that we weaken the notion of metabolism... and substitute mere energy dependency (with or without individual energy packets). The only purpose of this recommendation is to allow virtual beings, which have physical existence but no body, to count as life. These question-begging proposals have no independent grounds to buttress them. (P.8)

The desire of the A-life research program of the 80s and 90s to be able to create life in a computer that Boden criticizes is one motivation for this focus on the formal properties of living systems to the exclusion of all else. Such a move may also stem from a more general desire to take formal models to be exhaustive of the phenomena they are supposed to describe – as when scientists become so seduced by the success of their models in serving some particular purpose that they overlook all that they cannot explain, reflexively incorporating arbitrary limitations of their model into their very concept of the target phenomena. This, as Merleau-Ponty (2004/1961. P. 292) argued, was the foundational mistake of the ‘absolute artificialism’ of cyberneticians, like Ashby (1962), who claimed to be able to explain everything from the economy to the behaviour of individuals in terms of feedback control.

As Chirimuuta (2020) describes:

The “absolute artificialism” of cybernetics is a kind of vicious circularity: the cyberneticist has an understanding of organisms based on selective attention to analogies with machines and then uses this conception of organism to inspire the building of new devices, which are then projected back onto living organisms as models of their workings, and through the cumulative and recurrent effect of this process it becomes impossible to think of the organism—including the human being—in any other terms than as a tool, a thing to be manipulated and an instrument at the service of interminable projects of intervention and control. (P.450)

I take the development of the FEF to have fallen into a similar trap, motivated by a belief in the priority and exhaustivity of mathematical formalization, and a desire to be able to use this alone to “answer the questions traditionally posed by metaphysics; i.e., what does it mean to be a thing that exists, what is existence, etc” (Ramstead et al., 2021, S.43).

That metabolism is not a purely formal property does not necessarily mean it cannot be *described* by a model, that it can only be realized by the specific set of amino acids that make up the metabolisms we are familiar with, or even that it can only be instantiated at the level of molecular synthesis. In Chapter 9 I will discuss Moreno & Mossio’s (2015) attempt to describe the functional organization of metabolism in terms of constraint closure. Their account, I will argue, captures the specific way in which flows of matter and energy must be organized so as to realize a metabolic network, thereby placing tight constraints on what could realize such a thermodynamic organization, without arbitrarily tying this to one particular set of chemical components.

So the fact that we can give a multiply-realizable functional specification of something does not, as Piccinini (2020) notes mean that it is medium-independent in the sense that a formal operation is. So long as some medium

has the right mechanical properties to realize the relevant syntactic operations and internal relationships between parts that define a particular operation, then anything can realize addition or integration, whether water in tubes or electricity in wires. In contrast, some medium's capacity to perform the function of splitting an atom or seasoning a chicken is ineliminably constrained by the intrinsic properties of atoms, or chickens. A neutron or a proton, but not an electron, can adequately perform the former function, for the latter, rosemary or thyme are good options, but strychnine is not.

Performance of these non-formal functions is not just constrained by the need to realize the right internal relations between parts, but also the need to realize the right relationship to something that is defined in terms other than its syntactic role in said function. In the case of metabolism, these are the properties of the energetic and material resources of the environment external to that metabolic system. Thus, as Ruiz-Miazo & Moreno (1999) argue, when approaching life from the perspective of metabolism one is forced to consider how a description of living systems in terms of an "abstract-relational logic" can be "geared in with the implementation or physical realization of some effective management of those energy flows" (P. 46)

A description that makes reference to these energy flows will place additional constraints on the realizing medium that abstract-relational logic of a purely formal description does not. When energy flows are incorporated we see that a particular medium will not be able to realize a metabolic part, like an enzyme, unless the breakdown of this medium is coupled with these external resources in such a way that its breakdown simultaneously releases energy to

fuel activity that absorbs and synthesizes external material into the regeneration of that same part.

If something is medium independent, then there is no distinction between an exact simulation or model and a genuine realization. The only difference between model and reality in these terms would be one of detail. So one might say that a particular machine only approximates multiplication but if this is defined in purely syntactic terms, then any non-approximate procedure with the correct formal properties will literally *be* a multiplier. If something is medium-dependent, in contrast, then a simulation in the incorrect medium, no matter how detailed, will remain merely a simulation because the medium in which it is realized will not instantiate the necessary non-formal properties. The reason the transistors and electrical currents that instantiate the ‘walls’ of your Sim’s house are not really walls is nothing to do with partiality or approximateness, but because they cannot support a roof, provide shelter from a storm, or enclose a piece of land.

Medium independence is thus crucial to the hope of strong artificial life, as stated by Langton (1989), “to build models that are so lifelike that they would cease to be models of life and become examples of life themselves.” (P. 63). It is also crucial to the FEF’s hope to characterize life in purely formal terms. But neither Frankensteinian aspirations nor a mania for mathematical models are good reasons for rejecting the importance of something with such significant consequences for an organism’s existence as its metabolism. If metabolism is important then models of life will always remain models and descriptions of their necessary properties will not automatically carry over living systems they are supposed to describe.

In simply abstracting away from material turnover to treat parts as stable, the machine-substance perspective ignores the difference between the extrinsic and contingent stability of the enzyme vs the intrinsic stability of the atom, ignoring their differing costs and consequences. This is not to say we cannot occasionally treat organisms like machines and give mechanistic models of them in cases where we can argue that this distinction is irrelevant. So abstracting away from this fluidity may still be justified by how it allows us to apply the FEF's second requirement, namely the stability of the probability distribution over the possible states of the system. Even if this is not a first principle from which 'everything of interest about life' can be derived it may, like natural selection, describe a useful law that living things tend to follow, and by means of which we can predict *some* of their behaviours

8.3. The FEF and the machine concept of the organism

When describing organisms as machines we treat them as things with stable parts that are individuated not by their material basis, but their role in a fixed organization. This machine-substance concept of the organism is distinct from the naive concept of a material-substance, but as Dupré and Nicholson (2018) describe, is equally rejected by processualist views. In analysing this Nicholson (2012) quotes Glennan's (2002) definition of a mechanism as:

a complex system that produces that behavior by the interaction of a number of parts, where the interactions between parts can be characterized by direct, invariant, change-relating generalizations (Glennan, 2002, p. S344),

This definition has two criteria: decomposability and invariance. The former of these, the claim that particular phenomena can be explained in terms of

local interactions between distinct parts, is most associated with the ‘new mechanist’ movement currently ascendant across philosophy of science (Bechtel & Abrahamsen, 2005; Machamer, Darden & Craver, 2000; Glennan, 2002). It is this claim of decomposability that formed a central target of the organicist movement in early 20th-century biology, a precursor and close relation to contemporary, processualist accounts of the organism (Dupré & Nicholson, 2018; Peterson, 2017; Chirimuuta, 2020). Living systems, the organicists argued, are not decomposable mechanisms but ‘organic wholes’ with emergent behaviour that alters the activity of their parts and cannot be reduced to local interactions between these.

I agree with the organicists that the way in which the parts of organisms are mutually dependent upon one another for their contingent stability makes decomposing them for explanatory purposes, as mechanistic accounts do, more of a distortion than in the case of a machine, where we can literally separate out the fuel injector and the intake manifold from the engine. I discussed this in the previous section, however, and my concern here is not with exploring the consequences of this for the legitimacy, or not, of mechanistic explanation.

Rather, with respect to how this machine concept of the organism influences the FEF, the relevant issue is the presumption that the organism itself (and not just some useful model we might give of it) has, at some level of analysis, an invariant set of equations describing its behaviour. From this we could then identify some invariant formal parts in terms of the variables of these equations – even if those parts do not correspond to discrete things that can actually be separated from the organism in reality.

It is these two essential invariants of a machine-substance ontology that are expressed in the two stability requirements making up the FEF's definition of a 'thing.' So if, for any living system, we can identify a fixed set of parts whose dynamics can be described by some invariant rule, then it will qualify as a machine in the sense I am using the word. If this invariant rule is one of convergence to a steady state (or ergodicity) then it will be, more specifically, a free energy minimizing machine.

I have already touched on the implausibility of thinking that either ergodicity or steady state, might serve as a first principle from which we could derive everything of interest about the organism. Yet, for all that this is a particularly impoverished framework for characterizing life, I want to argue that rejecting it as the 'wrong sort of invariant' does not go far enough. In examining how the FEF's invariants fail to track the identity of an organism over time, I believe we will discover the impossibility of identifying any invariant rule that can do so.

As we've seen Friston presents the invariant, homeostatic rule that defines our steady state system in two different ways 1) as a fixed joint probability distribution over states that must be preserved, described as the 'generative model' and 2) as a fixed stochastic differential equation that defines the random dynamical attractor of a steady state system. While the former, statistical, version is more common in philosophical treatments of the FEF, particularly in connection with predictive processing and the Bayesian Brain, it is a consequence of the second.

So, *if* every 'thing' is indeed a random dynamical system with a point, or limit cycle attractor, then Friston's claim that its existence over time is dependent upon its states changing in conformity with the steady state equation

necessarily follows. The latter is not an *a priori* truth about all existence, however, but derived from a perfectly defeasible antecedent. As I will describe, as soon as we define something in terms of dynamical system theory, we are committing to the substantialist ontology of essential and invariant rules that makes this thing what it is. The more important question that Friston ignores, is whether dynamical systems theory in general, and a random attractor in particular, are the appropriate conditions for defining an organism.

A dynamical system consists of a set of equations of motion and a fixed phase space of dependent variables that capture all the possible states for that system. So, for the simple case of a friction-less pendulum, we have length, mass and gravitational force as parameters, a two-dimensional phase space of angular position and velocity as dependent variables and a differential equation relating these to determine possible trajectories through that phase space.

As Longo, Montévil & Kauffman (2012) describe, these equations and the phase space are standardly derived from what physicists and mathematicians call the ‘symmetries’ of the system, invariants that must be continuously preserved throughout any transformation it undergoes. In the case of our idealized, frictionless pendulum this is the overall energy, which remains constant while potential energy (determined by position) and kinetic energy (determined by velocity) change. The equations of motion specify these invariant-preserving transformations, while the phase space consists of all and only those complexes of states that can occur within some invariant preserving trajectory — see fig. 7.

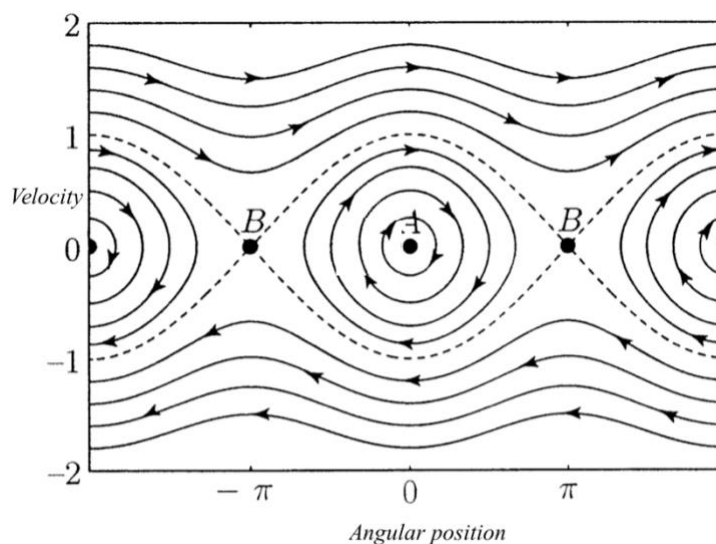


Fig. 7: phase space for a simple, frictionless, pendulum (adapted from Winter & Murray, 1997)

In its focus on continuous change over time, dynamical systems theory has been touted as offering an account of cognition that could supplement the discrete, atemporal symbol manipulation of Turing-style computationalism (Van Gelder, 1995; Chemero). Moreover, in allowing our phase space to be made up of collective variables, which need not correspond to the distinct ‘real parts’ of a compositional mechanism, dynamical systems theory has been embraced by those in the enactive and embodied tradition of cognitive science, seeking to describe the putatively holistic and emergent behaviours of complex systems. (Ross, 2015; Batterman and Rice, 2014; Chemero and Silberstein, 2008; Stepp, Chemero & Turvey, 2011)

There is some debate whether the covering law descriptions provided by DST (Walmsley, 2008) qualify as explanations, and whether they offer a genuine alternative to mechanistic analysis (Kaplan, 2017; Kaplan & Bechtel, 2011; Kaplan & Craver, 2011). But whether or not DST does indeed provide an alternative, holistic explanatory paradigm, in defining a system in terms of a set of invariant equations and parameters that dictate its dynamics and are

not affected by them in turn, it nonetheless retains the substantialist commitment of the machine picture that concerns us here.

As Koutroufinis (2017) puts it:

"There are two essential elements of Turing machine-logic. First, sharp distinction between state variables (what is calculated) on the one hand and parameters, independent variables and operators (what is given from outside) on the other. Second, no generation of new general types of intrinsic causal factors, and thus restriction to a fixed state space" (P. 34)

So while dynamical systems theory may offer an alternative descriptive language to computational or mechanistic analyses, by assuming invariant transformation rules with a fixed and finite space of possible states, Koutroufinis argues, it can describe a target system only in so far as it treats this target as following the invariant logic of a machine.

This doesn't mean that dynamical systems models require us to commit to a substantialist ontology regarding the things that our equations are supposed to be models of. If we take 'the system' for which these properties must remain invariant to characterize descriptive content of our particular approximate model, rather a definition of some real-world target, then the question is only whether our model has enough flexibility to characterize all of the behaviours of our target which interest us. If, as Friston does, however, we take this model to be a *definition* of what it is to be, say, a pendulum or a person, then we commit ourselves to a machine-substance ontology where these extrinsic factors of parameters, phase space, and equations constitute the essential features of the particular 'thing' in question. If they change because, say, the string of the pendulum breaks, then we say that particular 'thing' has ceased to exist.

It is this commitment, to what I am calling the machine-substance account of the organism that forms the ontological backdrop of the FEF as a theory of things. Applied to a pendulum, it seems reasonable enough. But how well does the FEF in particular, and dynamical systems theory and the machine-substance view more generally, serve as a conceptual analysis of organismic existence?

8.3.1 Is surprisal-minimization the substance of an organism?

As discussed in previous sections, the fixed steady state equation of the FEF allows a system to continue changing in two ways that preserve a stationary probability distribution.

The first are temporary fluctuations away from the attracting states – so long as the duration spent in these perturbed states remains consistently small. This is guaranteed by the dissipative component of the ESIA-cycle equations, which was stipulated to exactly counterbalance these fluctuations. The FEF can also allow for systems that converge to a limit cycle around equally likely states, described by the (optional) solenoidal component of the EISA cycle equations. This has been proposed as a distinguishing feature of biological systems but, as mentioned in Chapter 4, seems inadequate to the task given that plenty of things, from pendulums to planets, can cycle through recurring sets of states for some duration before gravity and friction eventually bring them to a halt – yet we are not inclined to consider them alive at any point of their doing so.

So while the system can change in these stereotyped ways, the FEF presumes that all we need to characterise our system is a description of a single random dynamical attractor, a single free-energy minima, to which it will eventually

converge upon and remain within. As Ramstead et al. (2021) succinctly express this:

On this account, to exist as a living system entails continually revisiting the neighbourhood of the same set of states (e.g., remaining within a certain range of body temperatures or ecological environments). Technically, such systems are endowed with a random dynamical or pullback attractor. This engenders a nonequilibrium steady state density that we can associate with the phenotype of a living system. (P.110)

This is a particularly simple characterisation that elides the wide variety of itinerant trajectories that can potentially be described in dynamical systems models, as Friston (2019a) acknowledges when he states:

*... symmetry breaking (i.e., divergence of nearby trajectories to different regimes of phase-space) is a hallmark of nonequilibrium dynamics (Evans and Searles, 2002) and is intimately related to phenomena like self-organised criticality in dynamical systems (Bak et al., 1988; Vespignani and Zapperi, 1998). Indeed, much of complexity science addresses the problem of how to formalise multiscale, itinerant and chaotic dynamics... **In this monograph, we will elude many of the finer details (and phenomena such as bifurcations, frustration and phase transitions) and suppose that the interesting behaviour of self-organising systems can be captured by nonequilibrium steady state densities with the right sort of shape. (P.20) [my emphasis]***

In characterizing lifeforms, this focus on systems that repeatedly return to the same region of phase space has some plausibility, as the continuation of an approach that treats homeostasis as the fundamental principle of organisms, that we saw in Ashby. That logic, Ashby and Friston both argue, extends beyond bodily regulation to every activity of even complex creatures like ourselves, from the bodily rhythms of our respiratory, cardiac, circadian or hormonal cycles, to our daily routines, weekly schedules and annual festivals. A nice example of this is at the behavioural scale comes from a

study on the records of millions of mobile phone users which found that over the space of three months, their location could be predicted with 93 per cent accuracy. irrespective of how far the individuals tended to travel (Song et al., 2010).

Still, the observation that organisms often regularly revisit the same states is not proof that they *have* to as a condition of continued existence. A pervasive tendency is not the same as a necessity. Humans may often be creatures of habit and routine, but they are also creatures that undergo dramatic identity crises. They may regularly revisit the same locations over some three month period, but at the end of those three months they may leave and never come back. A young girl with stars in her eyes and a dream in her heart abandons the Iowan farmstead of her childhood to move to the big city; a worn-out graduate student ups sticks for a log cabin in Siberia, and a certain type of man reaches a certain age and trades in his daily commute and his VW Golf for a motorboat and a fishing license.

Perhaps the best-known objection to the FEF then, often discussed under the heading of ‘the dark room problem’ (Sun & Firestone, 2020; Friston, Thornton & Clark, 2012), is the point that a drive to minimize free energy fails to even touch upon this diversity of creative and novelty-seeking activities that are characteristic of human behaviour.

Friston and co. have attempted to respond to this by noting that when we model systems as attempting to minimize the run long-run average of free energy (called expected free energy), we allow for temporary, exploration-driven, increases in free energy in service of the long term aim of maintaining the fixed generative model (Schwartenbeck et al., 2013; Parr & Friston, 2019). But to start a new career is neither temporary fluctuation from nor a

recurring cycle within a stable pattern of behaviour but a critical transition from one regular pattern of behaviour to another. If we embrace the FEF in its strongest form, as the proposal to explain all the behaviour of living systems in terms of "one simple imperative; avoid surprises and you will last longer," (Friston, 2013); if we accept a definition of existence in terms of the preservation of a probability distribution and the claim that survival is the avoidance of phase transitions (Friston & Stephans, 2007), then such a change in routine, habits, or home should amount to our death.

8.3.2. Identity crises: from behaviour to bodies

A natural response to this behavioural variability would be to propose constraining the principle of stability to hold only for certain 'essential variables', such as body temperature or blood oxygenation. Various parameters and other proxy variables might then be free to change in all sorts of non-recurring ways, governed by the requirement of keeping these essential variables stable. This is exactly the approach taken by Ashby in his description of the ultrastable homeostat, which responds to a breakdown of stability by random reorganization of connections between its parts until stability is re-obtained. Still, perhaps the major weakness of his account is the absence of an account of which variables in an organism should count as *essential*, and why (Harvey, 2013). The FEF is in a still worse position than Ashby for, unlike in the homeostat, in the steady state EISA cycle no such distinction is even available – the only elements of our overall system being either a) a variable that remains at a steady state or b) random, uncorrelated noise.

Moreover, this attempt to pick out a subset of variables for which a stability principle might hold is particularly challenging in the face of the dramatic transformations of an organism's entire body plan that occur over the

development of an undifferentiated bundle of cells into a complex, heterogeneous newborn. It is not surprising that the field of embryology has played an important role in the early development of processual accounts of the organism, exemplified in the work of Paul Weiss, Joseph Henry Woodger, and Conrad Waddington – the latter being a particular key figure in the introduction of Whitehead’s process ontology into the study of biology.

If we soften some of the more strident claims in the FEF literature, such as Friston’s (2019b) declaration that “things only exist on timescales over which they are ergodic” [or at steady state] (P.176), then we might think of the transformations of embryo growth as stages on the path towards the steady state of a substantial and stable adult organism. This essentially involves decomposing our system into a dissipative component, during which it moves through parts of phase space to which it will never return, versus the conservative set of states which it eventually settles into and continues to revisit. An example of this might be a pendulum on a vibrating surface whose range of motion decreases as the energy from its initial impulse is eaten away by friction, but which eventually settles into a smaller region of vibration-induced oscillation around its equilibrium point. So, rather than claiming things (biological or otherwise) only exist while they are at steady state, we could interpret the FEF’s definition of a thing more plausibly as claiming that everything will eventually converge to some steady state that defines it as the kind of thing it is.

Still, this does not account for the kind of radical life crises that may occur in non-human organisms, even after the achievement of a stable adult form. Concession is sometimes made in the free energy literature to the metamorphic butterfly as an odd exception to the rule that organisms can be

defined by a stationary probability distribution, but little work has been done to reckon with the consequences of this for the aim to “unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; avoid surprises and you will last longer” (Friston, 2012, p.2). Moreover, this kind of phenotypic plasticity is not that uncommon. There are also the 500 odd species of fish known to change sex under the appropriate circumstances (Murata et al., 2021), or the dramatic transition of a placid bunch of short-horned grasshoppers into a swarming, seething mass of locusts, with a quite dramatically different morphology and behaviour, when they reach a critical mass (Burrows et al., 2011) – Fig. 8.



Fig. 8: (Burrows et al., 2011) Comparison of solitary (left) versus swarming (right) locusts in larval (upper) and adult (lower) phases

Even in organisms as simple as a single cell, these transitions in the parameters of survival can be observed. A classic example is the lac operon mechanism of the *E. coli* bacteria. Ordinarily, such bacteria may be well characterized as steady state systems, cycling through the metabolization of glucose molecules, with some fluctuation around a stable intake level. Yet when glucose levels drop and lactose levels rise, a coding region in the *E. coli*'s DNA is activated, triggering the production of enzymes to allow for the metabolization of lactose as well. This amounts to a move from a steady state that was describable by a stable probability distribution with zero likelihood

of lactose metabolism, to a new steady state where this now has a high probability.

8.3.3. Stability across the scales

This pervasive adaptivity and phenotypic plasticity are one of the several points that, Colombo and Palacios (2021) advance in order to argue that the steady state attractor needed in the FEF is inadequate to characterize the behaviours and states of living organisms. As they note, this in itself, does not necessarily undermine the in principle ‘truth’ of the claim that organisms minimize free energy relative to some model however, only the formulation of this model in terms of a single steady state attractor and symmetric unimodal joint probability distribution.

Alternative, more complex and multistable, attractor landscapes are available. Thus, as Koutroufinis (2017) describes what looks like a change in the equation governing the dynamics of a system from one perspective, can also be described as merely a movement from one locally-stable subregion of a fixed attractor manifold to another, governed by a higher-order set equation that remains invariant – point also made by a pre-FEF Friston (2000) in the context of modelling neural dynamics.

This distinction between a single attractor and a multistable landscape, as Koutroufinis (2017) suggests, is nicely captured by Von Foerster’s distinction between a trivial machine vs a non-trivial machine. A trivial machine is definable by a fixed mapping between input and output, such as a keyboard which always emits the same sound for each key that you press. A non-trivial machine, however, is a machine that changes its state in response to an operation at one time in a manner that alters its response to any subsequent input. A Turing machine is the paradigm of a non-trivial machine. So are the

computers we encounter, in so far as their state (eg. their memory) may be changed in ways that alter their response to subsequent external events.

In a non-trivial machine, we have at least one higher-order function, which, rather than mapping from one state of a system to the next, maps instead to different parameter values in the first-order equation which then describes the evolution of the system in turn. For particularly complex systems this nesting of functions can continue indefinitely, but eventually we must reach an invariant function, which is not altered by the state of the system and remains invariant throughout all these transitions. In this manner, as Koutroufinis notes, even systems whose movement through state space results in the alteration of the equation governing their dynamics may still be describable in terms of a machine-substance ontology, where the essential invariant is this higher-order function that does not change.

If we can define a fixed and unchanging attractor landscape for a system, then in principle, after enough time, that systems behaviour would inevitably start to recur – as once it has visited every region of its phase space it will have nowhere new to go. Up until this recurrence, the probability distribution over possible states of that system will appear unstable, as when an *E. coli* first switches from glucose to lactose. But once it starts to recur, its dynamics would, as the FEF claims, eventually converge to a stable probability distribution over possible states that it will cycle through with an unchanging likelihood, where changes between different locally stable states occur with a fixed probability.

The problem with a multistable landscape is that we can no longer just derive the fixed generative model of the organism by just taking the first recurring pattern of behaviour we find in a particular individual. Something might be

in a low surprisal state over one timescale, but this could just be a temporarily stable phase in a larger multistable trajectory, rather than being reflective of that system's full generative model. Accordingly, Ramstead et al (2019, 2021) suggest that we instead think of the free energy principle as something that holds over multiple scales, not only at the level of the organism but also the overall species. At one level surprisal minimization may describe the local stability of a particular phase in the organism's life cycle, such as the maintenance of steady glucose levels. At another, it may govern the transition between the different local steady states that make up the broader multimodal but stable distribution of states observed for the E.Coli species as a whole. By looking at an ensemble of E.Coli then, we can find that for any randomly selected time period the probability of some bacterium being in the lactose metabolizing state would be constant. While the move from one steady metabolic state to another may appear surprising from the perspective of a single bacterium that has spent its life to date on a glucose-only diet, from the perspective of the species it is a predictable and recurring transition mandated by a fixed, species-level model.

This kind of response, this shifting the issue of how to individuate and identify the invariant 'generative model' or 'phenotype' of a system between different scales, until we find one at which the free energy principle applies, is the kind of thing that raises concerns about its lack of predictive power and unfalsifiability. The right scale for evaluating the applicability of the principle of free-energy minimization cannot just be 'whatever scale at which it happens to be true.' Moreover, as Friston (2019a) admits in the quote above, the free energy framework tells us nothing about the shape of this multistable, species-level attractor landscape or why these critical transitions would occur. All it allows us to say is that whatever form it has is fixed, and

thus that any individual iteration's trajectory through it would eventually converge to the same fixed probability distribution over states.

'Eventually', is the keyword here. To make this move between a generative model for the ensemble and one for the individual relies on not just the assumption that a system will converge to a steady state, but the stronger requirement of ergodicity – such that we can view the probability distribution over states for different iterations of a system as interchangeable, in spite of their differing initial conditions. As discussed in Chapter 4 a reason for shedding this for the weaker steady state requirement is that ergodicity is actually rather difficult to prove outside of idealised models. One problem here, as Colombo and Palacios (2021) point out is that the time for the average behaviour of an individual to converge to the ensemble average could be arbitrarily long, and far outstrip the duration over which that system exists (see also, Palacios 2018; Gallavotti, 1999).

The trouble with this is perhaps not so obvious for a creature with a comparatively small behavioural repertoire, like an *E. coli*, where the idea that each of its members could be defined by one invariant species-level model, might be somewhat plausible. Yet applied to say, human beings, any probability distribution over all the states that members of the species might feasibly occupy would cover such a large variety of possibilities as to be completely uninformative about the dynamics of any individual. Life is short. Even if there were a fixed attractor manifold for human viability, no single member of our species will ever have time to explore enough of it for their individual dynamics over time to converge to a stable probability distribution. In this respect, the FEF's stable probability distribution would not only be an asymptotic idealization that is never actually realized by the kind of systems that we are interested in but, crucially, one that places no

informative constraints upon the dynamics of these systems over the timescales on which they actually do exist.

8.3.4. Why the machine concept of the organism fails

To say an organism is a free energy minimizer presumes that we can give a model of that organism as a dynamical system with a fixed phase space and some fixed equations of motion. But once we have this model, redescribing it in terms of a fixed probability distribution only loses information. If the landscape is complex enough, then there will be a great many ‘high probability’ states and so this generative model becomes useless in helping us describe what the organism will do next. This is enough reason to disregard the free energy principle as trivial and irrelevant to characterizing an organism.

Still, I think there is a more fundamental reason to think that the FEF fails. The problem is not just that it is uninformative for sufficiently complicated machines. The problem is that living systems are unlike machines in so far as we cannot identify any fixed phase space and equations of motion for them, either at the level of the species or the individual, no matter how multistable our attractor landscape, or how multi-level the equations in our non-trivial machine model (Koutroufinis, 2017; Longo, Montévil & Kaufmann, 2012; Longo & Montévil, 2013; Mirazo, Pereto & Moreno, 2004; Kaufmann 2019, 2000; Rosen, 1991).

The problem with attempting to call upon a stable, ensemble probability distribution across ‘humans’ to provide a stable model for an individual person is not just that this would be too vast to be meaningful for the individual’s trajectory but the fact that, even with hundreds of thousands of years of history under our belts, the human ensemble doesn’t appear to have

plans to converge upon a steady state any time soon. Moreover, thanks to the rapid timescales at which we are continuing to transform our behavioural phase space through cultural and technological evolution, the ensemble distribution for humans is continuing to change within the lifespan of an individual person – not just between generations.

For my great grandmother, the possibility space for human beings in general would have been very different from her first birthday to her 90th. Prior to the invention of the aeroplane, there was no chance of finding her, or anyone else, several miles in the air above the Atlantic ocean. To be in such a position would be both non-viable, and highly surprising by the conventions of the time. By 1990 however transatlantic flight had become a much less improbable state. Over my great-grandmother's lifespan, neither her individual trajectory nor that of the human ensemble converged to a stable model that could be used to define a human phenotype and which she could be interpreted as minimizing surprisal with respect to.

This instability is not a unique problem resulting from some distinctive human capacity for information sharing. Nor can it be resolved by attempting a gene-centric 'essential variables' response. Genetic variation does not just occur between organisms, via the copying of DNA for transmission from a parent to its offspring, but can also occur within the lifecycle of an individual organism. Bacteria, for instance, trade genetic material like 90s kids traded Pokémon cards in the form of small rings of self-replicating DNA called plasmids (Firth et al. 2018). Sometimes these plasmids can become integrated into a bacteria's central genome, and this horizontal gene transfer is one hypothesised origin for one of the three genes making up the E.Coli Lac Operon (Hediger, 1985) Another possibility is the non-destructive insertion of genetic material by a retrovirus, the same

mechanism responsible as much as 5-8 per cent of the human genome (Belshaw et al., 2004), and potentially the development of the mammalian placenta (Mi et al, 2000).

Prior to the acquisition of this gene, there was no meaningful sense or scale under which either the lucky bacterium in question, or *E. coli* in general, had a high probability of lactose metabolism. No meaningful sense in which this should have been part of its phase space in any dynamical systems model of it. The acquisition of the ability to metabolize lactose is ‘surprising’ on whatever timescale we attempt to locate a stable reference frame for this bacterium. The only way the FEF could describe the development and maintenance of this capacity is as a violation of the principle of free-energy minimization, resulting in the destruction of one ‘generative model’ and its replacement with a new one. Yet, rather than equating such a novel adaption to death, we would tend to recognize this as a development occurring within the same *E. coli* bacterium.

Suppose our hypothetical bacterium had never bumped into the hypothetical retrovirus that provided it with the final piece of the Lac Operon puzzle. Then *E. coli* might never have developed the ability to metabolize lactose that gives it the headstart on mammalian gut colonization. They might never have taken on the role of preparing the ground for further microbial inhabitation, and the synthesis of a variety of vitamins essential to mammalian functioning – from the multi-purpose folic acid to the vitamin-K required for the formation of blood clots (Maynard & Weinkove, 2020).

There are no necessary laws that entail the development of a microbiome. Some animals lack one altogether (Hammer et al. 2019), and there are mammals, such as bats, that do not appear to depend upon microbial support

to meet their nutritional needs. Still, most mammals do, and when specimens, known as ‘germ-free’ lines are produced without this microbial colony, studies show that their nutritional demands in the absence of this assistance are up to 30 per cent higher (Wostmann et al., 1983). Speculatively translating this to humans, this more direct dependence upon, say, a constant source of leafy vegetables to supply our vitamin-K needs might have constrained our adaptability to less fertile regions, and limited or slowed the species outwards expansion.

Like most attempts at narrating deep evolutionary history, this is something of a just-so story. Stuart Kauffman (2019) provides a similar tale in his book *A World Beyond Physics* (P. 97) describing the evolution of the first food chain, and the first symbiotic relationship among protocells. Another more complicated narrative might be given for the retrovirally-assisted development of the mammalian placenta and the transition from oviparity (egg-laying) to viviparity (internal gestation). Though such stories are inevitably speculative, the point is that hard as it is to trace back the sequence of events that led to where we are now, it was much harder to predict them in advance. Indeed, as Longo, Montévil & Kauffman (2012) claim, the inability to locate some rule, or symmetry, that must be preserved through any transformation is what makes organisms ‘unprestatable’ – that is to say unformalizable with the tools of dynamical systems, or indeed any other machine-type logic. As they argue:

Thus, it is proliferation, variation and selection grafting novel phenotypes into evolving organisms that reveals, again after the fact, the newly relevant and unprestatable observables and parameters. Thereby, this is our main thesis, the very phase space of evolution changes in unprestatable ways. In consequence, again, we can write no equations of motion for the evolving biosphere, nor know ahead of time the niche boundary conditions so cannot integrate the equations of motion which we do not have. No law entails the evolution of the biosphere. (P.6)

Longo et al. focus on how these unprestatable developments unfold over phylogenetic timescales – and thus prevent the formulation of exceptionless laws describing the evolution of some population, species, or indeed the biosphere as a whole. More important here, however, is the fact that these symmetry-breaking events do not only occur *between* organisms, from copying errors and recombinations in the transmission of DNA from a parent to its offspring, but also within a single individual. Events such as horizontal gene transfer or rapid within-generation changes to an individual's cultural and physical niche may alter the possibility space for that organism within its lifetime.

By way of explaining how the machine concept of the organism fails here, it is not enough just to say that an organism is a process, not a substance. This does not itself explain what it is that differentiates organisms from other, less troublesome sorts of activity like countdowns or rotations, for which we can still provide a fixed rule that describes their unfolding. As Longo and Montévil (2014) point out, “the dynamics of biological organisms, in their various levels of organization, are not “just” processes, but permanent (extended, in our terminology) critical transitions and, thus, symmetry changes.” (2014, P.161) The key difference between physical and biological processes, as they argue is that, “Usual physical processes preserve invariants, whereas extended critical transitions [characteristic of biology] are a permanent reconstruction of organization and symmetries, i.e., of invariants” (P.175).

To illustrate this point, consider again a pendulum. Here, we have a set of variables, such as angular displacement and velocity; parameters, such as string length and pendulum mass; and background constants such as the gravitational constant. We also have some symmetry principles and

conservation laws, that specify things that must remain invariant as the system moves through a fixed phase space of possible states. For instance, conservation of energy tells us that when potential energy increases as the pendulum’s angular displacement increases, so angular momentum must decrease. This allows us to derive a set of equations to describe every state this system might possibly enter while still remaining an (undriven, frictionless) pendulum – see Fig. 9.³⁴

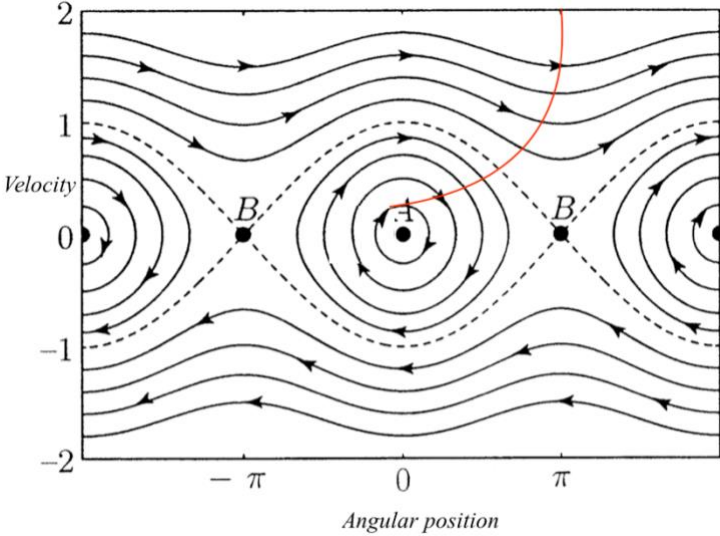


Fig. 9: phase space for a simple, frictionless, pendulum illustrating an ‘impermissible’ flow in red (adapted from Winter & Murray, 1997)

A real pendulum is not an isolated system, and processes not included in our model might disrupt it in ways that our model cannot predict. But while we cannot predict these from the study of the pendulum alone, we can still specify every possible position in phase space that it might unexpectedly be displaced into, along with how it will behave subsequently – so long as it remains the same type of system. If we allow our parameters to be varied we

³⁴ For a real system we would also have to incorporate co-efficients of friction to describe how total energy is dissipated by this resistance, but the trajectories of the pendulum would still be constrained such that at each point the greater the angular displacement, the lower the velocity – and vice versa.

can also prestate what would happen if the string were extended or if we increased the weight of the bob.

This doesn't mean that the pendulum in itself contains all the information we need to exhaustively specify every possible thing that might happen to its constituent parts. From studying our pendulum alone, we could not pre-specify what would happen if the experimental laser in the much more exciting laboratory next door were to malfunction, sending a beam of super-high powered coherent light through the wall and vaporizing our modest little experiment into disconnected gas particles with greatly expanded degrees of freedom. Still what we can do is give an exhaustive specification of the possible trajectories that are compatible with that pendulum continuing to exist *as a pendulum*.

For the pendulum then, the possible ways it might be affected by the world beyond can be reduced to either perturbation of its position in a fixed phase space, or to its destruction. But because, as the processualist argues, an organism is not individuated by a set of fixed rules but rather by some relation between its temporal stages, so an externally induced event – whether the insertion of retroviral genes, a plasmid exchange, or the learning of a new skill – may alter the equations governing its dynamics and expand its phase space in a way that we do not necessarily consider as the destruction of that organism.

Is this true unpredictability? Is it due to something about the metaphysical status of organisms rather than just a reflection of our own epistemic limitations? Longo et al. (2012) propose such an indeterminist view, taking the unprestatability they ascribe to biological systems to be, at least partly, grounded in genuinely random events, pointing to the possibility for

indeterministic quantum mechanical events to trigger point mutations that could have phenotypic consequences for the organism. Yet this question of whether the macroscopic trajectory of an organism might be altered by metaphysically-unpredictable events cannot be what makes organisms ‘unprestatable’ in a way that machines are not. Thanks to the complexity and small scale of modern computer components, so these machines are similarly vulnerable to what are known as ‘single event effects’, triggered when a single, randomly released, particle crashes into a single crucial transistor – with macroscopic consequences ranging from the crashing of aeroplanes to the overturning of elections. Moreover, even amid the biological, such irreducibly indeterminate macro-events would be exceedingly rare and unlikely to alter a cell’s behaviour in the course of its lifetime.

While it is important for my purposes that the unprestatability I’m ascribing to the trajectory of an individual organism reflects a genuine difference between that organism and a machine, not just a matter of our epistemic limitation in light of to the greater complexity of the organism, I don’t think we need genuine randomness for this. What matters most in this respect is not the randomness of *when* an event happens, but the relational nature of *what* it means for the trajectory of that organism. A functional safety engineer, tasked with designing integrated circuits for critical purposes, has no way to tell you if or when a single photon will strike a particular transistor in that circuit board – but their job is contingent upon the fact that the fixed structure of that circuit allows them to specify *what* would happen if it did. What distinguishes an organism from a machine is not the possibility of its being macroscopically disrupted by sub-atomic fluctuation, but rather the way in which environmental changes do not only destroy or perturb organisms but also *enable* them by expanding their repertoire of possible behaviours. Whether macroscopic changes in an organism’s structure are

deterministic in origin or not, the point is that their consequences for that organism depend upon how they interact with a host of other changes within the rest of the organism, in its environment, or other organisms.

We can make a generic statement about the consequences of lengthening the string in a pendulum without reference to any particular pendulum or its environment. Rolling back the clock billions of years, the same would not be the case for the consequence of the insertion of a *lacA* gene into an *E. coli*. Whether that insertion would lead to the capacity for lactose metabolism will have depended upon *this particular E. coli* being one that had also evolved the other *lac* operon genes, and whether it would lead to faster mammalian gut colonization depended on whether there were mammals waiting around with guts ready to be colonized. Similarly, the occurrence of genetic changes resulting in the development of bony fins of an ancient aquatic tetrapod would not, in themselves, have allowed us to predict that the tetrapod would walk on land. This depended also upon the drying out of the environment in which that tetrapod lived, something that could not have been predicted solely by studying the structure of that tetrapod alone. Bony fins played a causal role in the development of amphibiousness, but they did not entail it. Prior to the arrival of dry land into its environment, walking around was not a part of the possibility space for that tetrapod.

So the possibility of enabling causes, in addition to perturbations or destructions, depends on a rejection of the machine-substance conception of an organism as being individuated by some invariant structure and dynamics. For these sorts of reasons, Longo et al. (2012) argue for the enablement of unspecified possibilities, as more apt than the entailment of some specific effect, for the study of biological causality. This move alone is not enough to secure a claim of objective unpredictability for the

development of a biological system. Without genuine randomness in the generation of macroscopic mutations, a Laplacian demon could still predict the development of an individual organism, in virtue of its ability to trace every interaction that occurs in the world as a whole. But this possibility of enablement of new behaviours is enough distinguish the organism from the pendulum in so far as it prevents even this demon from prestatating a finite set of possibilities for how the former could develop (conditional upon its not being destroyed) in isolation from knowledge of the development of the entire environment in which it is embedded.

The argument that we cannot pre-state the phase space and equations of motion for an organism may sound like biologists arguing against the very possibility of doing biology. But it should not be taken as a blanket objection to the use of models, such as those of dynamical systems theory, to describe organism dynamics. We can still describe stereotypical behaviours of an organism using dynamical systems models, and we can retrospectively construct a model of the trajectory that an organism did in fact take. To take a processual view just means remembering that these dynamical models are only locally valid approximations, subject to change in ways that we cannot predict in advance. In this context, the bold claims about ‘existential imperatives made by advocates of the FEF serve as a nice *reductio ad absurdum* against mistaking the relationship between some particular, contingent pattern of behaviour and a mathematical model of it, for an identity between the behaving organism and said model.

Moreover, to say the biological constitutes an unprestatable *World Beyond Physics*’ as Kauffman does is not to say it is supernatural world beyond the physical, populated by vitalistic entities and energies. It is just to say that it is a realm that cannot be predicted with the sorts of physical models that are

currently dominant. This should not be controversial. The various inadequacies of the Newtonian paradigm of a mechanical universe are well established and, as Rosen (1999) notes, complaints about the impoverished tools of the field have been made by leading physicists throughout history, from Schrodinger to Einstein. But as each new mathematical formalism improves our ability to describe some aspect of a system's behaviour within limited and arbitrary constraints our enthusiasm can lead us to forget all the phenomena beyond those constraints for which it fails.

So both the free energy principle, and a machine substance ontology more broadly, fail to capture the ways in which organisms are free of allegiance to any particular material instantiation or governing equations. Before moving to whether enactive accounts can do any better, I want to briefly consider whether a more modest form of the FEF, uncommitted to machine-substantialism might be useful instead. Once we separate it from this machine-substantialist framework I believe we can see why a description of living systems cannot afford to ignore the metabolic nature of organisms.

8.3.5. Life as the process of seeking stability?

Rather than taking the FEF as providing a definition of a living system in terms of a fixed generative model, we might instead opt for the weaker claim that the lifecycle of a living system is characterized by the need to continuously re-establish *some* steady states, if not the mandatory preservation of the *same* steady state. We would no longer be guaranteed a particular steady state density describing the unchanging phenotype by which a particular organism could be identified. All this would allow us to say about the overall trajectory of an organism is that it must at each point find some state or pattern of behaviour that it can stably maintain and thus be in a position to minimize free energy relative to going forward. The collapse of

this stability and rise and free energy would not be the organism's death but a sign of the need to find a new stable regime.

The problem with this approach is that the FEF itself makes no discrimination between the viability of possible steady states, other in terms of how stable they actually end up being. Ordinary death and consequent decomposition look well described in terms of a failure to secure stability, but this doesn't mean that any stable state will automatically be good for the organism in virtue of that stability alone.

An alternative response to declining food stores in some species of bacteria is to enter a frozen state of cryptobiosis, called an endospore, in which they can persist for thousands of years. Similar strategies are encountered in other microorganisms, like the famously resilient tardigrade (Wright, 2001) and even animals as large as Alaskan wood frogs, which can spend up to seven months frozen solid (Costanzo, 2013). In extremely rare cases, this kind of self-mummification may even be practised by humans, as with the practice of Sokushinbutsu among Shingon Buddhist monks in Japan or cryogenesis among Silicon Valley billionaires – though a crucial difference is that these transitions have never been successfully reversed.

Whether or not cryptobiotic organisms, whose vital processes will re-start under more auspicious conditions, should be counted among the living while in their frozen state is not so much controversial, as it is oddly un-discussed by among philosophers of biology. A scattering of remarks by biologists studying the phenomenon suggests that many hold that “on an organismal level, they are essentially dead,” as Alaskan wood frog researcher, Don Larson puts it (quoted in Netburn, 2014). Such assessments, as Neuman (2006) describes, are typically based on the presumption of ongoing

metabolic activity as a necessary criterion of life. That there is any question about the matter can be credited to the fact that in several cases, such as the tardigrade, it is difficult to establish whether metabolism is truly completely inactive (Pedersen et al., 2020). ‘Crypto’, means hidden, not absent, and as Keilin (1959) introduced it was originally intended to refer specifically to the absence of ‘*visible* signs of life’ (Clegg, 2001, p. 213). For the above biologists at least, determining whether the cryptobiotic tardigrade is alive or not looks to come down to the empirical question of whether it is, in fact, continuing to metabolise residual energy stores at undetectable levels

The FEF, with its disregard for metabolic turnover, can make no such distinction. The proposal that organisms are driven to just establish *any* stable state fails as a normative principle for living systems, in so far as it treats both metabolically-inactive phases and the kinds of steady state we normally associate with being alive as of equivalent value. The requirement to be in some steady state has nothing to say about why we (or at least I) would be inclined to say that the E.Coli chowing down on lactose is doing rather better than the frozen endospore, and much better than the self-mummified monk – despite the impressive stability of the latter two. Some transitions between steady states may not equate to the death of the organism, but others do. What is needed is a principled account of why some steady states count as viable and others do not, why some transitions are compatible with the continuation of an organism’s existence, and others are not.

This is precisely what an account of the autonomy of an organism that does not abstract away from metabolic self-production, should provide. A state is viable for an organism only in so far as it is compatible with the particular self-producing organization of the organism at that time. Crucially, however, what allows us to follow the organism through time is not that this

organization has any fixed form, but how each stage is related to the next in terms of the relationships of reciprocal dependence between the organization, the processes it enables, and the organization they produce in turn, and so on. As the particular form of this autonomous organization changes, so what states are viable for it may also change in ways that cannot be predated by the presumption that it, like a machine, has an organization that remains invariant.

In the next chapter, I want to take a closer look at bioenactive characterizations of autonomy, to see why the standard formulation is not quite up to this task, before suggesting an alternative that I believe is, and which I argue the bioenactivist should adopt if they want an adequate, non-trivializable grounding for talk of biological autonomy, teleology and intentionality.

9. Bioenactivism and autonomy: from process closure to closure of constraints

9.1. Processualism and bioenactivism

In being uniquely unmoored from either material or organizational commitments, organisms are not machines and cannot be reduced to the logic of machines – for all that it may sometimes be useful for explanatory purposes to model them in mechanistic terms. The claim that an organism’s capacity for change is unprestatable by fixed dynamical equations, no matter how higher-order they may be, marks a genuine distinction in kind, not merely degree of complexity. But the fact that we can follow an organism throughout these changes leaves a puzzle as to how we do so. Moreover, a distinction in terms of the capacity for open-ended change is not enough for bioenactivism, as it does not tell us what marks organisms alone as *agents*.

The answer to both of these questions, I believe, lies in choosing the intrinsic instability of a metabolically-produced system over the default stability of a homeostatic one as our starting point for a theory of the organism. This priority given to the preservation of autonomous organization over the homeostatic preservation of stability is, Di Paolo, Thompson & Beer (2022) argue, a distinguishing feature of the enactive approach (or what I’m referring to more specifically as bioenactivism) which raises a barrier to any potential compatibility with the free energy framework. Still, the degree to which this

preservation of autopoietic or autonomous organization is compatible with open-ended development has not always been clear.

Indeed, autonomy and autopoiesis are defined in Maturana and Varela's early work in explicitly machine-like terms, where the cell is presented as an 'autopoietic machine'. As Thompson (2007) emphasizes, this does not reflect the contemporary, atomistic concept of mechanism as something whose behaviour can be decomposed in terms of local interactions between its parts, but rather describes a system with a relational organization that is multiply realizable and independent of its particular instantiation (what M&V refer to as its structure) at a particular moment. Nonetheless, this commitment to an invariant organization still fits within the 'essentialist' view of the machine-substance conception – where, rather than being defined by the first-order properties of its parts, an organism would be individuated by some invariant second-order relations between these that do not change over time. In other words, as Di Frisco (2014) notes, Maturana and Varela have not so much rejected substantialism but moved from an atomistic understanding of substance to an formalist one, allowing material instantiation to vary, while holding that it is instead the form (or 'organization') that must remain invariant.³⁵

This is the same move that was discussed in the previous chapter for allowing the FEF to account for the turnover of material components making up its fixed causal graph. As described, talk of higher-order organizational invariants can also allow us to make sense of lower-level variations in dynamics, as when an organism moves from one pattern of behaviour to

³⁵ While DiFrisco follows Simondion in reserving 'substantialism' for the atomistic view, contrasted with a formalist notion of substance, his criticism is targeted at both in terms of their commitment to invariant features. Thus as he notes, both would count as substantialist in the persistence sense

another, such that there is no fixed one-to-one mapping from input to output as there would be in Von Foerster's (2003) 'trivial machine'. The claim that there is no such fixed mapping is a much weaker one than Longo, Montévil & Kauffman (2012)'s argument that 'there are no entailing laws' or fixed equations of motion for biological systems, however. It just means that organisms are non-trivial machines, and, while a non-trivial machine may involve arbitrarily many layers of higher-order transition rules, these will ultimately terminate in some fixed rule and so, "in its highest level of operation any non-trivial machine is a trivial machine" (Koutroufinis, 2017, P. 34).

It is this fixed rule that corresponds to the 'organization' that, for Maturana and Varela, defines a living system throughout these 'structural' variations in lower-level rules. But, in so far as we accept the arguments of the previous chapters for the unprestatability of an organism's phase space, then for the organism, unlike the non-trivial machine, there is no such termination, no fixed rule and no fixed organization that we can reduce all changes in structure or behaviour to variations within.

The consequence of viewing organisms as machines, rather than unprestatable processes of critical transitions is illustrated in Maturana and Varela's replacement of an entailing relationship between input and output with the cybernetic concept of a 'perturbation' – a compensatable disruption to an otherwise stable mode of being. Contrast this, with Longo, Montévil & Kauffman's (2012) suggestion of enablement, as more apt than entailment for the analysis of biological causation. Unlike reduction of events to 'perturbations', which reduces the environment to a source of disruptions and challenges that one must preserve one's inherently stable form against, the idea of enablement presents our changing surroundings as an expanding

well of possibilities for continuing a process of production and individuation in novel and unpredictable ways.

All of this is not a problem for the goals of autopoiesis theory. As discussed in Chapter 7 Maturana & Varela were not initially aiming to distinguish organisms from machines. The residual idea of an invariant organization that is continually reproduced has, however, created concerns about any claimed compatibility between bioenactivism with processual accounts of the organism (Di Frisco, 2014), raising the question of whether, as Meincke (2014) asks, autopoiesis might be “a substance wolf in process sheep’s clothing?”

While, as described, definitions of autopoiesis have often been framed in terms of the regeneration of some invariant set of components, processes, or relations, my categorization of bioenactivism is not wedded to the specifics of the early formulations. The importance of autopoiesis, from my perspective, is not that it provides a sacrosanct analysis of the necessary and sufficient conditions for life, but that it points towards a different approach to explaining what a living system is in terms of the logic of self-production — rather than the capacity for differential reproduction and evolution, or in terms of the particular and contingent chemical form of it that is familiar to us.

Moreover, as Di Paolo, Beer, and Thompson (2022) argue, a distinguishing feature of the *enactive*, as opposed to merely autopoietic, concept of an organism is the focus on “precarious, self-constituted entities in ongoing historical development and capable of incorporating different sources of normativity, a world-involving process that is co-defined with their environment across multiple spatiotemporal scales and together with other

agents.” (P. 3) It is this focus on cumulative historical change, they argue, which constitutes one of the irreconcilable theoretical tensions that undermine claims of potential compatibility between the steady state formalism of the free energy framework and the enactive approach.

As an evidential basis for this account of the importance of historicity in the origins of the enactive approach, they point to Varela’s statement in his preface to the 1994 re-issue of his and Maturana’s *On Machines and Living Beings*, where he acknowledges some inadequacies in their account because, as he puts it:

*... it seems to leave the phenomenon of interaction in a grey area of being a ‘mere’ perturbation” (ibid. 614). Structural coupling “does not properly take the account of the emerging regularities in the course of a history of interactions ... Over these years I have developed an explicit alternative ... turning the historical reciprocity into the clue of a co-definition between an autonomous system and its environment. I propose to call this point of view in both biology and cognitive science, **enaction** (2011/1994, P. 614).*

So, as Varela recognizes and a number of authors have argued, these enactive ideas do not derive immediately from the original formulation of autopoiesis. A recent concern of enactivist literature has been on how to supplement, or modify, the concepts of autopoiesis and autonomy to render them more suitable for bioenactivist aims. A particularly significant development in this respect, as briefly mentioned in the introductory discussion of bioenactivism, is Di Paolo’s (2005) claim of the need to supplement autopoiesis with an account of adaptivity, describing the organism’s capacity to parametrically regulate its coupling with its environment. As mentioned in Chapter 1, this idea of adaptivity, and its relationship to an open-ended process of learning, has since been developed in Di Paolo, Buhrmann & Barandiaran’s use of

Piaget's ideas on sensorimotor equilibration to describe how novel interactions can be incorporated into agent's repertoire.

Another particularly interesting move towards a more processualist conception of the organism is the uptake of the work of the 20th-century French philosopher Gilbert Simondon, among both critics of autopoiesis (Di Frisco, 2014) and enactivists (Di Paolo, 2020). In rejecting the attempt to define an individual in terms of either a fixed form or fixed material parts, Simondon advocates prioritising the process of individuating within which each individual is but a temporary phase and which, for the organic, is an inherently unfinished process. In this respect Simondon's work stands as an interesting philosophical precursor to the arguments of complexity theorists and theoretical biologists reviewed in this chapter, and, as Di Paolo argues, is well suited to supplementing the bioactive literature by directing a focus towards this 'open-ended becoming'.

Still, in discussions of this kind of adaptivity and accommodation, the focus is on the sensorimotor level – suggesting the presumption of a more basic network of bodily metabolic processes that must be preserved and, in service of which, ever new sensorimotor engagements might be incorporated. This is clearly an improvement on the FEF's language of steady state EISA cycles, in so far as the bioneactivist at least has the means to describe a distinction between the preservation of a particular network of self-perpetuating processes, versus the open-endedness of the possible sensorimotor engagements that might serve to preserve it. Moreover, unlike Ashby's starting point in the homeostasis of 'essential variables', the bioactivist can give a principled account of *why* some particular variables must be kept within particular bounds, whereas others can be freely varied in support of the this

goal, in terms of how the former must be kept within those bounds if the autonomous network of processes is to continue its operation.

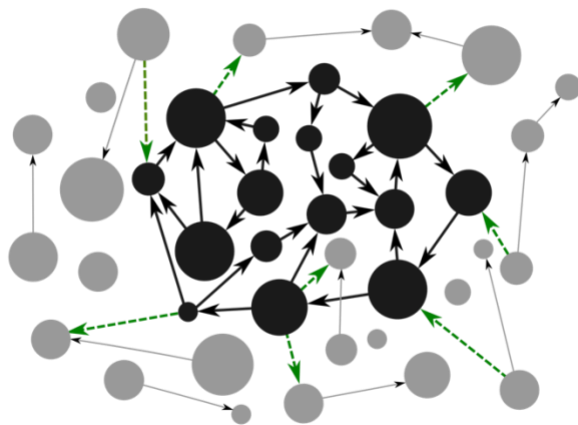
Nonetheless, as a consequence of a focus upon these secondary dimensions of autonomy in recent enactive work, with the aim to scale up to explanations of cognitive processes of learning and development, I believe descriptions regarding the prior autonomy of the organic body – the constitution of which autopoiesis was intended to pick out at the molecular level – remain insufficiently refined.

As described in Chapter 1, the definition of autonomy that I take to be the widely accepted one in contemporary bioenactivism is given in terms of an operationally closed network of processes, such that every process both depends on at least one other process in that network and enables a further process, together with the requirement of precariousness, such that these processes would not continue to operate outside of said network. As De Jaegher & Di Paolo state this:

An autonomous system is defined as a system composed of several processes that actively generate and sustain an identity under precarious conditions. To generate an identity in this context is to possess the property of operational closure. This is the property that among the enabling conditions for any constituent process in the system one will always find one or more other processes in the system (i.e., there are no processes that are not conditioned by other processes in the network which does not mean, of course, that conditions external to the system cannot be necessary as well for such processes to exist). By precarious we mean the fact that in the absence of the organization of the system as a network of processes, under otherwise equal physical conditions, isolated component processes would tend to run down or extinguish. (De Jaegher & Di Paolo; 2007. P.487)

Equivalent formulations are also given in Di Paolo (2009), Thompson & Di Paolo (2014) and (Di Paolo, Buhrmann & Barandiaran, 2017). To avoid debate over the difference between early characterizations of “operational

closure” (Bourgine & Varela, 1992) versus what Bich and Arnellos (2012) argue should instead be referred to as ‘organizational closure’, I will refer to this characterization as “**process closure.**” Like all accounts of closure, this does not mean that the system with process closure does not depend on other processes external to it – only that from amid all of these dependence relationships we can extract a network of *mutual* dependence. And it is only those processes that both enable and are enabled by other processes within this network that will be a part of the system that realizes process closure.



Copyright Esquiel Di Paolo, 2013. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. http://creativecommons.org/licenses/by-nc-sa/3.0/deed.en_US

Fig. 10: Illustration of how an operationally closed network of processes (highlighted in black) is distinguished from its surroundings. (Di Paolo, 2013)

As depicted above, this definition of process closure tells us how to extract those processes that are, or are not, a part of a particular set of cycles that occurs over some chunk of time (see fig. 10). Importantly, however, the recurrence time of many of these cycles is shorter than the lifespan of an organism. The regeneration cycle for liver tissue, for instance, can be as short as a few days (Sender & Milo, 2021). While there are also much longer cycles making up the organism, as argued in the previous section, we cannot conceptualize the identity through time of the organism in terms of one longer overarching cycle. A more obvious problem with doing so is that

every organism's lifespan is terminated in the inevitable failure to 'loop back' when the organism dies.

So process-closure only tells us the spatial extent of an organism over the particular period within which a set of cycles occur. The question is how we identify the organism throughout multiple cycles. Must its identity be defined by the *same* repeating cycles, the same organization, the same processes and the same products? Di Paolo, Buhrmann & Barandiaran (2017), for instance, often talk of a cycle of regeneration among the same set of processes in the same network, implying that beneath the open-ended adaptivity of sensorimotor learning we might still be able to identify a stable, essential organic core by which the 'self' of the individual organism might be defined. Given the discussion in the preceding sections of this chapter, and, specifically, the metabolic plasticity observed not only in single-celled bacteria like the *E. coli*, but even in multicellular organisms like deep-sea fish (Raposo de Magalhães et al., 2021), I am not convinced that it is possible to specify an invariant network of processes at a level both specific enough to individuate a particular organism, and flexible enough to incorporate the possibility of such changes.

Nonetheless, the above definition of autonomy seems compatible in principle with a view in which the organism's identity through time is determined not by the preservation of any one particular network of processes or components, but in terms of the continuity of relations of production between a series of evolving process-cycles. Though Austin (2020) is critical of the attempt to avoid resorting to substantialist principles in characterizing an organism, he notes, "Amending the rather permissive relation of genidentity with the *constructive* character of autopoiesis certainly furnishes one with a more restrictive criterion of slice-series composition"

(P. 9) The same might be intended with autonomy, such that even if a particular process network underwent organizational change, we would still be able to take the new organization and the old one as manifestations of the same organism, thanks to the latter organization being dependent upon the operations of the former for its existence.

So the fact that the enactive approach has been concerned with developmental change since its beginnings is enough, as Di Paolo, Thompson & Beer (2022) argue, to undermine claims of compatibility with the FEF. I believe we can make a stronger claim than incompatibility, however. I have argued that this non-cyclical behaviour is a pervasive and distinctive feature of living systems that the FEF's language of steady state EISA-cycles is constitutively incapable of describing. This motivates a ruling against its ability to serve as a 'first principle' for living systems – enactivist or otherwise.

Still, the bioenactive definition of autonomy is not without its problems. If we are to individuate the organism in terms of relations of production between temporal stages, rather than in terms of an invariant organization that is reproduced, then we need a robust account of what these relations of production are. As I will argue in the next section, process-closure is inadequate to capture the relations that underpin a self-producing continuity and, as such, is inadequate to ground the attribution of purposiveness and intentionality to living systems.

9.2. The problems of defining autonomy as closure of processes

The second aspect in which I argued the FEF failed us a theory of living systems was in how it downgraded material turnover to optional changes that *may* occur within the constraints of a necessarily invariant organization. As I have argued, it is only by describing how the organism is not merely free of a particular material basis but in *need* of a continual flow of matter by which it continuously re-constitutes itself, that we can adequately characterize the precarious dependence of organisms upon their own activity. It is this characterization that we need to properly ground the bioactive understanding of agency, intentionality, and immanent teleology. (Jonas, 1953, 2001/1966; Weber & Varela, 2002; Thompson, 2007)

So one improvement of the contemporary bioactive definition of autonomy, given in the preceding section, over Bourguin and Varela (1992)'s operational closure, is that it describes an ongoing process of individuation under which change and process are primary. Nonetheless, as Bickhard (2008) argues, the material and energetic conditions of life remain underdeveloped in these contemporary formulations of autonomy. This leads to the risk of trivialization, for, as Thomson & Di Paolo (2014) acknowledge:

All material processes are precarious if we wait long enough. In the current context, however, what we mean by “precariousness” is the following condition: In the absence of the enabling relations established by the operationally closed network, a process belonging to the network will stop or run down. (P. 4)

This clarification is not sufficient for precarious operational closure to pick out living systems distinctively. As Moreno & Mossio (2015) & Mossio & Bich (2017) describe, many physical and chemical systems from tornados and

convection rolls to the hydrological cycle would meet this set of requirements.

Consider a simple network of a swingball set and a robot programmed to play it. Here we have two mutually dependent processes. The first: the balls orbiting around its pivot. The second: the robot's moving its arm to hit it. These processes do not merely enable each other, they depend upon each other. If the robot doesn't hit the ball the latter's rotation would run-down, to leave it hanging limply from its string. If the swingball didn't continue to swing around, the robot's hit-the-ball process wouldn't be activated because there would be no ball swinging past that it is able to hit.

Here we have a precarious and operationally closed network of processes. If this were all there is to autonomy, then, with additional capacities, such as the ability of the robot to regulate the force of its stroke with respect to wind conditions, we would have a system that is also adaptive, and thus, something that starts to look like it meets Di Paolo, Buhrmann & Barandiaran's (2017) criteria for being an agent. I don't think we should accept this. To take the coupled robot-swingball or the hydrological cycle as autonomous, as intentionally oriented towards goals and norms, is to invite the criticism, oft raised against bioenactive approaches from traditional autopoiesis theorists (Villalobos, 2013; Villalobos and Ward, 2015; Villalobos and Dewhurst, 2017), that we are inappropriately projecting teleology on to systems that are merely state-determined mechanisms. Moreover, I take this insufficiently robust account of self-production and precariousness to be what has left bioenactive work vulnerable to the kind of trivializations of autonomy presented by the FEF, under which even coupled pendulums and Watt-governors are argued to qualify. Without a robust notion of precarious self-*production*, we cannot distinguish the adaptive agency we might wish to ascribe

to organisms from the capacity of a coupled feedback mechanism to change its state or structure, as part of a continual cycle of activity.

This doesn't mean we should abandon the bioenactivist program to opt for either instrumentalism or hylozoism. All it means is that bioenactivists have not sufficiently captured what is special about the metabolic self-production of a cell in virtue of which it carries the germ of intentionality, teleology, normativity and agency.

9.3. Constraint closure

In what respect has the above definition failed, to capture the thermodynamic basis of autonomy, as Moreno & Mossio (2015) suggest? The cyclical process of a robot-swingball set depends upon a continual flow of energy into the system from outside, via the robot's power cord, but mere energy dependence is trivial. All processes depend upon a flow of energy for their continuation – the miracle would be a system where they did not.

There are two key respects in which the robot-swingball is nothing like an organism. The first is that neither the robot-hitting nor the swingball-swinging play any part in securing or regulating the energy supply that enables them. The second is that the *structure* of both the robot and the swingball set are intrinsically stable, and will persist with or without this supply. It is only their dynamics that depend on a flow of energy (as indeed all processes do) the structure of the robot-swingball set does not.

The problem, as Moreno and Mossio put it then, is that this account of closure among processes 'fails to locate closure at the relevant level of

causation' (Moreno & Mossio, 2015). Like the analogy of food as fuel discussed in Section 8.2, it presumes the system can be factored into a fixed set of background constraints and the 'autonomous' dynamics that result. These fixed constraints are taken for granted, and their presence is presumed to require no explanation when we come to describe the system any more than information about the manufacturer is required for a dynamical model of a pendulum. It is this presumed division, as Koutroufinis (2017) describes, that is the key move in dynamical systems theory, where we separate the invariant equations of motion, determined by the constraints of a fixed organization, from the resultant dynamics, which cannot alter those constraints in turn.

This modelling division is not a distortion when what we are modelling is a machine. We really can 'divide' a machine's structure from its activity by stopping the machine, and we can also pull those structural pieces apart. We can turn our car off and start it again and we can take it apart and put it back together again without consequences. This cannot be done for organisms because these constraints are themselves dependent on the processes that they collectively enable for their continuance (Dupré & Nicholson, 2018; Nicholson, 2018; Mossio and Bich, 2017). As such, Montévil & Mossio (2015) propose that the appropriate characterization of the organism is in terms not of closure of processes, but more specifically closure of *constraints* – an idea inspired by Kauffman's (2000) work on autocatalytic sets and work-task cycles on the one hand and by Rosen's (1991) ideas of closure to efficient causation on the other.

Their account is elaborated at greater length in Mossio & Moreno's (2015) book *Biological Autonomy*, where it is defined as follows:

In formal terms, a set of constraints \mathbf{C} realises closure if, for each constraint C_i belonging to \mathbf{C} :

- 1. C_i depends directly on at least one other constraint of \mathbf{C} (C_i is dependent);*
- 2. There is at least one other constraint C_j belonging to \mathbf{C} which depends on C_i (C_i is enabling). (P 20)*

Before explaining whether this account is sufficiently specific to the organic and whether it has the potential to provide a naturalistic grounding for teleological, intentional, or normative talk, we need first to say a little more about this distinction between constraints and processes, which Mossio & Moreno define as follows:

Processes refer to the whole set of physicochemical changes (including reactions) occurring in biological systems, which involve the alteration, consumption and/or production of relevant entities. Constraints, in turn, refer to entities that, while acting upon these processes, can be said to remain unaffected by them, at least under certain conditions or from a certain point of view. (P. 11)

This point is cashed out more formally in terms of symmetries, where a constraint is something that remains invariant with respect to particular thermodynamic flow, just as the length of a (frictionless) pendulum string, or its overall energy, remains invariant throughout its changes in angular position. The first natural question to arise here is how can something that is partially defined by *not* changing also be something that 'acts upon', or causes something else? The second concern is how something that is defined in terms of its invariance could, at the same time be precarious and in need of regeneration? For a constraint to be cause, or to be an effect of some

process, is actually quite common outside the biological realm. I will deal with each in turn before describing how they can connect together in the structure of constraint closure unique way to biological systems.

9.3.1. Constraint causation

The standard philosophical view of causal relations, as Schaffer (2016) describes, is that they are events. One change occurs and, in a law-governed way, triggers a subsequent change as its effect. Additional factors, like the fragility of glass, or the energy contained in sugar, may alter what happens but they do not tell you why things happened at a particular time. So when your partner asks, “Why is my favourite glass broken?” they are asking what caused the event of its breaking and are expecting to be answered with some event – perhaps involving your recent tendency to conduct chemistry experiments with wanton disregard for the integrity of other people’s kitchenware. Pointing to ‘thermodynamics’ is seldom an acceptable response.

This is the view of causation offered in the mechanical conception of the universe, inherited from Newton, as essentially a collection of inert and independent billiard balls, that change their velocity only when bumped into by another in a manner determined by some external and eternal laws. But this is not the only way to think about causation, and as Alicia Juarrero (1999) argues, there is reason to think that our difficulties with differentiating the actions of an agent from mere movements can be traced to an implicit commitment to this “inadequate, 350-year-old model of cause and explanation.”

Even before we get whether a different conception of causation can do a better job for describing agency, we can note that this view of matter as inherently inert and changing only in response to the external determination

is also incompatible with contemporary physics, which conceptualizes matter not as static, but as intrinsically dynamic. Even at its lowest energy state, a particle like an electron will constantly ‘jiggle around’ whatever region of space it is confined to, the smaller the constraint the faster the jiggling. As Luisi & Capra (2014) describe

This tendency of particles to react to confinement with motion implies a fundamental “restlessness” of matter that is characteristic of the atomic world... The fact that particles are not isolated entities but wave-like probability patterns implies that they behave in a very peculiar way. To the extent that things can be pictured to be made of smaller constituents – molecules, atoms, and particles – these constituents are in a state of continual motion. Macroscopically, the material objects around us may seem passive and inert, but when we magnify such a “dead” piece of stone or metal, we see that it is full of activity. The closer we look at it, the more restless it appears... Modern physics thus pictures matter not at all as passive and inert but as being in a continuous dancing and vibrating motion whose rhythmic patterns are determined by the molecular, atomic, and nuclear configurations. There is stability, but this stability is one of dynamic balance, and the further we advance into matter, the more we need to understand its dynamic nature to understand its patterns. (P. 75)

In modern physics then the Newtonian order is reversed. There is no answer to what event ‘caused’ an electron to wiggle other than that this is just what electrons do. So how do we think about causation in a world where change is the default and it is stability that is in need of explanation?

One option is to take causal relations to hold only at the macroscopic level where things at least *appear* to be stable unless perturbed by something else. Such a view often ends up treating causation as useful for explanatory purposes but absent in fundamental physics, or, as Price (1992) puts it, as “anthropocentric, being linked to our perspective as agents.” In so far as ‘our perspective as agents’ is exactly the thing I take to be in need of explanation, this view will not help us here.

Moreover, as Hoffman (2012) argues, it is precisely the scale of the microscopic molecular storm at which life begins its operations, and in terms of which it must be explained. The problem of a living system, as he puts it, is that, “Without the shaking and rattling of the atoms, life’s molecules would be frozen in place, unable to move. Yet, if there were only chaos, there would be no direction, no purpose, to all of this shaking” (P. 21). In Chapter 8 I briefly touched upon how this molecular chaos may introduce genuine randomness into genetic replication, but of more relevance here is how it interfaces with an organism’s metabolic operation. As Godfrey-Smith (2016) describes:

Metabolic processes in cells occur at a specific spatial scale, the scale measured in nanometers—millionths of a millimeter. They also take place in a particular context, immersed in water. In that context and at that scale, matter behaves differently from how it behaves elsewhere... There is unending spontaneous motion that does not need to be powered by anything external. Larger molecules rearrange themselves spontaneously and vibrate, and everything is bombarded by water molecules, with any larger molecule being hit by a water molecule trillions of times per second... The way things get done is by biasing tendencies in the storm. (P. 4)

It is because of this dependence upon spontaneous molecular motion that Godfrey-Smith suggests that metabolisms could not have arisen at any other scale. An interesting question to look at once we’ve finished formalizing metabolism in terms of constraint closure is whether such processes could nonetheless extend beyond this scale once they’ve gotten going. Such a question will be important when it comes to explaining how any teleology we locate in metabolism is to be carried over to minds.

So, in order to make sense of the particular causal regime at work in living systems, both Juarrero (1999) and the theoretical biologist Robert Rosen (1991, 1999) argue that we need to step out of this Newtonian framework and look back to how causation was expressed in Aristotle’s account of the

‘four causes’ – efficient, material, formal and final, of which Newtonianism retains only the first.

Neither Juarrero nor Rosen intend to suggest that we revive the Aristotelean framework wholesale – indeed Juarrero blames Aristotle’s ‘prohibition against self-cause’ as much as the ‘billiard ball’ reduction of causation for the problems with our contemporary theories of agency. Rather, their proposal is that, amongst these various concepts of causation, we can find ones more apt to describe how it works at the molecular level of biological operations. For Rosen, this is formalizable in terms of recursive functions and their variables, in which the material ‘cause’ is a variable whose state is transformed, while the function is the efficient cause that brings this about (though as Moreno & Mossio (2015) note, ‘formal cause’ seems potentially more apt for the latter).

For Juarrero, it is a matter of the more concrete notion of constraints. Constraints can limit and stabilize movement, as when bonds between atoms limit their degrees of freedom and calm them into the apparently inert solid objects around us. Juarrero, however, is more interested in how constraints can not only limit possibilities but enable them.

This idea of a constraint as something that expands the range of possibilities can, at first, seem orthogonal to the meaning of the term. Nonetheless, constraints that ‘make things happen’, precisely by preventing other things from happening, are a pervasive feature of the world and we need not look to the biological to find them. It is by restricting a restless teeming mass of atoms into a small space that the cylinder of a steam engine can perform the work of making a wheel go round, and in so far as the cylinder also remains unchanged during this thermodynamic flow, so it qualifies as a constraint.

Constrained or no, the energy from burning coal would still be transferred. Without constraints, the only way for this to occur is via heat transfer as the hot gas expands and the kinetic energy of these gas molecules is distributed to the surrounding air molecules as they collide. What defines heat, as opposed to work, is that it involves not just a transfer but also a loss of ‘useful’ energy available to do further work. Through heat transfer, energy becomes dispersed into an object’s surroundings – corresponding to an increase in entropy. By contrast, work involves “the constrained release of energy into a few degrees of freedom” (Atkins, 1984) such that the ‘concentration’ of energy remains constant – as when energy from thermal expansion is channelled into the raising of a piston, an increase in potential energy that can then be transferred into the movement of a wheel. The nice thing about work, is there is always more of it that can be done.

So the more constrained the energy, the lower the increase in entropy (energy dispersal) and so the greater the amount is available to perform mechanical work. When indexed to temperature, it is this quantity of constrained energy that makes up thermodynamic free energy, which Friston has sometimes played fast and loose with analogizing to the statistical construct involved in free energy minimization. It is the last thing an organism wants to minimize if it is to continue doing the work of interacting with its environment and generally staying alive. Because there is no such thing as an ideal engine, however, free energy is also lost through heat whenever work is done, hence organisms, like any other system, need continual re-supply of energy if they are to continue performing work.

The funny thing about constraints is that the new possibilities they create are not found at the microscopic level of one-to-one interactions between

individual chemicals. A catalyst is the paradigmatic example of a constraint, but it does not make reactions happen that would otherwise be impossible, it merely provides a lower energy route for said reaction to occur, allowing it to happen much more frequently. Sugar oxidation, for instance, is an exergonic reaction, meaning that it happens spontaneously and releases energy, yet confectioners are rarely consumed by unexpected fireballs. This is because sugar oxidation normally happens extremely slowly. Even if you mix it in a glass with an oxidizing agent, like potassium chlorate (used in fireworks and matches) nothing will happen. But the moment you add a drop of sulfuric acid the two chemicals will instantly react in a violent purple explosion, destroying the glass in which they're contained (Shakhashiri, 1983).

The sulfuric acid doesn't exactly *do* anything, in the traditional sense of event-causation. Rather it is a catalyst that acts as a constraint, remaining unchanged throughout the reaction but accelerating the rate at which it occurs. At the level of the individual molecular reactions, this is just the same thing that would have happened anyway – just much, much faster. At the macro level of glasses, and the people who own them, it's the difference between a beloved crystal tumbler and a pile of glass shards in need of explanation.

In a similar manner, all things being equal, water molecules at the top of a hill will end up at the bottom, dispersing their potential energy in the process. It is only if this flow of energy is constrained by a channel that it will maintain sufficient concentration so as to turn a water wheel on the way down. By making a quantitative change to the spatiotemporal scale of microscopic processes, a constraint can make a *qualitative* difference in the production of a new macroscopic effect.

So the transfer of energy stored in an object to its surroundings is the norm and will happen spontaneously as the system moves into a more thermodynamically stable state – with a corresponding loss of ‘useful’ free energy via heat transfer and an overall increase in entropy. But for work to be done, for a macroscopic mechanical event to occur with the potential to cause further events in turn, this energy transfer must be constrained. As such, as Kauffman (2000) argues, when explaining why a wheel has been rotated or a glass destroyed, pointing to the transfer of energy seems to focus on the wrong locus of explanation. Energy transfer would happen either way, but the reason work has been done is because this energy transfer was *constrained*.

So this is the sense in which a constraint can be a cause. The next question we need to answer, in order to explain how we can have a closure of constraints, is how something that is defined partly by its invariance through change can nonetheless be the effect of an ongoing process, upon which it depends for its continual regeneration

9.3.2. Constraint production

As I have described, there is nothing particularly special or agential about constraints as causes. They are a feature of any machine that channels energy to perform work. In a machine, however, constraints are rarely treated as effects in turn. Once a machine has been created the constraints that make up its structure can be taken for granted as invariant features that stand outside of the thermodynamic flow. It is in this sense that constraints that are part of a typical physical model are described as ‘external’. And it is for this reason (in addition to their further decomposability into structural parts) that machines are so well modelled in mechanistic terms of a division between fixed equations and parameters vs dependent variables. When we

make such a distinction in a dynamical systems model of a machine, this model is picking up a real feature of the target phenomenon, in a way that I will argue it does not in the case of living systems³⁶.

Given that Moreno and Mossio (2015) define constraints as things that are unaffected by processes, and which are capable of “harnessing a thermodynamic flow without being subject to that flow”(P.15), so it may seem as though any constraint-based model of organisms must be unavoidably machine-like. Yet constraint closure is specifically supposed to capture how we can have constraints that do not stand *entirely* outside of the flow of activity in the system (as in a machine) but depend upon this activity for their repair or regeneration.

The crucial caveat in making sense of this tension is that the invariance of a constraint holds only “under certain conditions or from a certain point of view” (P. 11). What is meant by this is not that whether or not something is a constraint depends on whether you choose to treat it that way, but rather that a constraint is always defined relative to the processes that it constrains, as shown in Fig. 11. Taking an enzyme (C1): the reaction it catalyses ($A1 > B1$) has a timescale ($\tau1$) over which it occurs and with respect to that timescale, the enzyme remains invariant. Relative to that process an enzyme is a constraint. Over a longer timescale ($\tau2$), however, this enzyme will degrade and be in need of repair by the process of translating an mRNA sequence into the chain of amino acids that makes up that enzyme ($A2 > C1$), an assembly process that can only happen because it is constrained in turn by ribosomes (C2), which also need to be replenished in turn. And so on and so forth.

³⁶ While it may be less of a distortion it is still an idealization in that energy transfer via work is never perfectly efficient and energy loss via heat in the course of the operation of a machine can eventually degrade the structure of the machine itself.

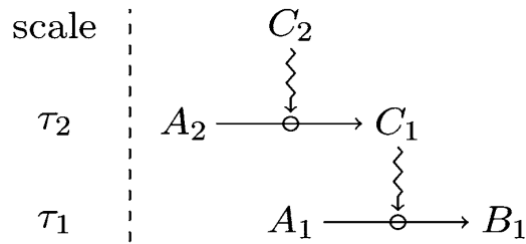


Fig. 11: Illustration of the dependence of constraints upon other constraints from Moreno & Mossio (2015), attributed to Maël Montévil

So something that is a fixed constraint over one timescale is nonetheless a stage in a continuing process over another timescale. As the biologist, Ludwig von Bertalanffy expressed the same idea, “the old contrast between “structure” and “function” is to be reduced to the relative speed of processes within the organism. Structures are extended, slow processes; functions are transitory, rapid processes” (von Bertalanffy, 1941, P.251 - quoted in Dupré and Nicholson, 2018). What distinguishes biological constraints from those found in a typical machine is that they are not intrinsically stable things, but extrinsically stabilized processes.

This point can seem a bit odd, given that we have been emphasizing the need for two different levels of causation in characterizing an organism – the need for a distinct role for processes versus constraints in our description of the system. If constraints are really just processes, then what makes constraint closure importantly different from the process closure of Di Paolo & Thompson (2014) and other contemporary enactivists? The answer is that process closure still implicitly involves constraints as the fixed structures that are engaged in and enable these processes. By leaving these in the background, process closure treats them as non-processual entities, as something completely external to the thermodynamic flow throughout the overall system. In other words, by focussing on the closure of processes,

such accounts preserve the absolute process/structure distinction of a machine.

The insight of constraint closure is not that constraints are separate from processes, but that in biological systems these constraints are themselves part of the processual network, not fixed structures standing outside of it. By defining the invariance of constraint relative to the process it constrains, Moreno, Mossio & Montévil's account explains how something can be both invariant over one timescale and changing over another, without making this choice of timescale a subjective decision of some external observer. With regards to mechanistic and dynamical systems models, their account of precarious constraints nicely explains how these can work over one timescale, while failing over a different one. This idea of precarious constraints thus explains how, in biological systems, the distinction between a process and a constraint can be simultaneously real and relational.

Just as there is nothing particularly biological about constraint causation, so there is nothing particularly unique about constraint production either. When an automated digger constrains the flow of energy from a battery so as to channel this into the mechanical work of digging a channel, it is creating a constraint that can, in turn, channel a flow of water down a hill so as to power a water wheel. We might even hook the digger up so that it is powered by this same waterwheel and imagine that the soil quality is poor so that the channel is constantly collapsing and must be continuously re-dug. In this situation, we have a precarious constraint that is necessary to enable the work that maintains that same constraint.

What we don't have here, and what we do have in the organism, is constraint closure. The dependence between the channel being dug and the water

flowing through it to power that digging is not direct, but mediated by a variety of other constraints, such as the structures of the waterwheel and the digger, that it plays no part in maintaining and which are stable in their own right. In contrast, to realize closure *every* constraint within the system must not only enable other parts of that system but depend on other parts in turn. Every enzyme, ribosome, membrane and mRNA strand in a cell is not only a constraint on its metabolic network but also a product of processes that can only occur because of the constraints of that very same network.

To say that an organism is a closed network of constraints does not mean that it doesn't depend on the external environment. Crucially, as with most attempts to formalize the kind of closure distinct to an organism, Moreno and Mossio (2015) emphasize that this runs alongside necessary openness – in this case, to thermodynamic flows of energy and the reactants that fuel these processes (as shown in fig. 12.) In this respect, just like Thompson and Di Paolo's (2014) definition of an operationally closed network of processes, a constraint closed system is both distinguished from, and constitutively dependent upon its environment as a resource. The difference in a constraint closed system is that it is not just the activity the system engages in, but the structure of the system itself that is so dependent.

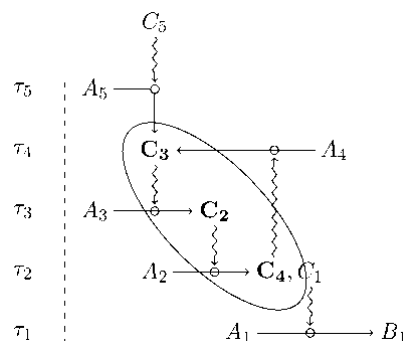


Fig. 12: Illustration of constraint closure from Moreno & Mossio (2015), attributed to Maël Montévil. C_x are constraints, A_x are reactants, and B_1 are products. Wavy lines depict a constraining relationship and straight lines depict a process of transfer

Further, constraint closure allows that the system can depend upon external constraints to channel these thermodynamic flows *to* it. But any constraint that channels a flow *within* the system – any constraint that operates between one constraint that is a part of a system and another – must also be regenerated by that same system if it is to realize closure. Thus all machines, even those whose activity depends upon its preservation of a precarious constraint, like the Digger-Waterwheel-Channel-system, will fail to meet the criterion of operational closure, in so far as this activity is mediated via a number of non-precarious structures that will remain stable in their own right, independent of the rest of that network and its operations.

So this is how constraints can cause, depend upon other constraints, and come together to realize closure. There is much more to Moreno and Mossio's account in terms of how additional processes of regulation, adaptivity, and evolution lead to the increasing complexity of constraint closed systems, both within and across generations. This brief account is also highly abstract in relation to the fine-grained empirical details of its realization in the intimidatingly complex network of real metabolisms, which even in their very simplest microbial form enlist several hundred reactions and metabolites. These connections have been elaborated more thoroughly by Stuart Kauffman (1986, 1993) in his work on autocatalytic sets, as a more empirically detailed precursor of Montévil and Mossio (2015) and Moreno and Mossio's (2015) account. There are also other closely related accounts of what Letelier, Cárdenas, & Cornish-Bowden (2011) review more generally as 'metabolic closure' – not least Robert Rosen's (1991) M-R systems, which serve as a direct, though in this case more abstract, inspiration for constraint closure

But it would take another thesis to explore all of this. For now, the question is what makes constraint closure better than either ‘closure among processes’ or steady state EISA-closure for defining life? And does it do any better at locating the basic ingredients for a real notion of intentionality and purposiveness from amid natural processes?

9.4. Constraint Closure as a theory of living systems

The most basic challenge of an account of living systems is to obtain the right balance of specificity and generality. While there may be some disagreement on the status of border cases such as the cryptobiotic organism, the virus, the seed, or some, as yet unrealised, sophisticated AI, the criteria proposed should at least succeed in including those things uncontroversially regarded as living, and in excluding those cases which no one not already in the grip of some theoretical worldview would mistake for being alive.

Moreover, a successful account also needs also to be able to deal with the possibility for life to have evolved in different ways, under different conditions. Simply tying life to some particular chemical manifestation, even if this is shared by all known instances on Earth would be the most obvious failure to describe “not just life as we know it, but life as it can be” Langton (1989). A central advantage of relational, or ‘operational’, accounts like autopoiesis is that they say nothing about the particular chemical compounds involved, only the organization they need to be able to realize in terms of the relations between them (Fleischaker, 1990).

Still, in focussing on molecular metabolism, autopoiesis is insufficiently general to describe life at the multicellular and sensorimotor levels. For this, Maturana and Varela initially proposed the idea of second-order autopoietic systems, composed of autopoietic sub-units (1987, P. 88-89). This,, as Thompson (2007, P.105-6) describes, is not quite satisfactory. How exactly do the autopoietic processes cohere in an organismal-whole? And is it the autopoietic sub-units or the overall organization of the multicellular organism in virtue of which it is to count among the living?

Because the bioenactive approach seeks to characterize the minds of multicellular creatures, not just life in its most basic single-celled form, so the more general notion of autonomy has taken on a more important role than autopoiesis. Yet in moving away from the single-celled level, to focus on the closure of precarious processes, such characterizations have relaxed more than just the restriction of autopoiesis to molecular synthesis. Specifically, as argued in the earlier part of this chapter, this account loses sight of the relation between these processes and their products, between the activity of the system and the body it regenerates, which was core to autopoiesis. As a result, the present characterization of autonomy as process-closure is now too general to distinguish the living, given how it may be instantiated in non-living systems like the hydrological cycle or the robot-swingball set.

Process-closure still does better than the FEF's steady state EISA-cycles, which somehow managed not only to be too general – in applying to non-living systems – but also too specific – in specifying necessary requirements that living systems themselves do not meet. Both EISA-closure and process-closure fail for a shared, and more basic reason however, that they go too far in their abstraction – not only disassociating life from any particular chemical realisation but also from its unique thermodynamic status. In autopoiesis,

these energetic and material requirements were arguably at least implicit and only in need of further elaboration (Fleischaker, 1988). In contrast, by placing the structure of constraints entirely outside of the thermodynamic flow of the system to focus only on processes, so contemporary enactivist accounts of autonomy succeed in providing a substrate neutral, multiscale characterization of closure – but only at the expense of erasing this crucial distinction between precarious, non-equilibrium constraints versus intrinsically stable ones (Ruiz-Mirazo & Moreno, 2004; Bickhard, 2000).

As such, when seeking to account for the distinction between living systems and other circular organizations, these accounts can only resort to either the requirement of autopoietic sub-components or, as Ashby did, to arguing that such a distinction just comes down to the degree of complexity of the mechanisms involved. Thus Di Paolo (2005) proposes adaptivity as an additional requirement for life, and, together with Buhrmann and Barandiaran (2017) proposes this, along with further supplementations as being necessary for agency.

I think Di Paolo et al. are correct that such properties are required for agency, and perhaps for life too. However, following Ruiz-Mirazo & Moreno (2004) and Moreno and Mossio (2015), I take these to be intrinsically linked to a more robust notion of self-production than Di Paolo et al. provide. Without such a criterion capable of distinguishing autonomous agents from automated machines, the distinction between life and non-life, the coupled-pendulums and the bacteria, becomes a matter of the degree of adaptivity or complexity of some network of processes. Life itself becomes a purely abstract phenomenon that might as easily be instantiated in a computer as a chemical network.

In this respect, such organizational, or relational, definitions can be contrasted with work that has focussed specifically on the thermodynamics of life, and on living systems as far from equilibrium dissipative structures, such as a candle flame, a convection roll, or a tornado, that are constitutively dependent upon ongoing flows of matter and energy for their continuation. (Christensen & Bickhard, 2002; Collier, 2004, 2008; Christensen & Hooker, 2000; Juarrero, 1999; Fleischaker, 1988; Nicolis and Prigogine, 1977; Bertalanffy, 1968; Schrödinger, 1951). Such dissipative structures seem to manifest the needful freedom in their relationship to matter that Hans Jonas (1953, 2001/1966) took as the hallmark of the immanent teleology unique to living systems. Yet while the dissipative structure and living organisms may share this precarious form of existence, it seems just as problematic to ascribe intentionality and vitality to tornados, convection rolls or a candle flame as it would be to ascribe such properties to a computer system or to coupled pendulums.

Accounts of life seem caught between the metaphor of the candle flame and that of the computer (Keller, 2008, 2009). Neither the organizational nor thermodynamic properties of a living system alone suffice to capture what makes an organism them distinct from both.³⁷

³⁷ In her two-part review of the history of the concept of self-organization, *Organisms, machines, and thunderstorms*, Evelyn Fox Keller (2008, 2009) provides a lovely overview of the development of these contrasting approaches to the organism.

As Fleischaker (1990) describes:

It is not new, of course, to point out that life requires energy to drive its processes of production. Nor can it be claimed that life is alone in requiring energy for the integrity of its internal organization: fluid or gaseous convective systems utilize heat-driven density gradients to that same end. These complex dynamical systems are non-living, but they, too, transform energy from the environment in maintaining themselves at a distance from equilibrium, and they hold energy in non-linear relationships among system components, that is, in circular (self-amplifying) relationships in which effects become cause (Swenson, 1989). What is unique to living systems is the organized coupling of energetic and material interactions in a single network of processes whose outcome is the production of all system components, including the constituents of its membraneous boundary structure. (P. 128)

The strength of constraint closure then is how it combines the insights of both organizational and thermodynamic approaches to living systems. As a relational notion, a precarious constraint is more general than the notion of a polypeptide chain, or an enzyme – yet something’s ability to stand in the relevant relations is highly constrained by its material properties. Namely, this material must not only be capable of channelling a thermodynamic flow to enable a particular process, but it must be both invariant throughout the timescale on which this process occurs, and unstable on the larger timescale of a slower flow to which it is also subject.

It is because of these particular thermodynamic requirements that the notion of a precarious constraint, while medium-variable, is not medium-independent, and why metabolism is not a purely formal property that can be literally realised in a computer model, as Boden (1999) argued. We may describe and interpret parts of a computer simulation as representing ‘flows of energy’ and ‘precarious constraints upon those flows’, but so long as the bits of silicon and metal that realize this simulation are not actually dependent upon the energy that flows through them, so there will always be a gap between even the most detailed simulation and a real metabolic network.

In contrast, a dissipative system has the right thermodynamic properties to literally be a precarious constraint. Yet in such cases of spontaneous ‘self-organization’ there is only mutual dependence between a single constraint and its process of regeneration and so they do not realise the organizational properties of constraint closure proper to the biological, in which, “constraints are not able to achieve self-maintenance individually or locally: each of them exists in so far as it contributes to maintaining the whole organisation of constraints that, in turn, maintains (at least some of) its own boundary conditions.” (Moreno & Mossio, 2015, P. 17)

This is important for explaining why dissipative systems both spontaneously emerge and then vanish almost as quickly in a way that living systems do not. Where self-organizing systems can arise whenever the boundary conditions are appropriate – for instance when heating from below creates the appropriate temperature differential for the movement of individual molecules between the top and bottom of a liquid to spontaneously organize into the coherent rotating cells characteristic of Rayleigh-Bénard convection – these boundary conditions remain outside the influence of the dissipative system itself. As such, the system has no influence in sustaining them and will disintegrate as soon as those external constraints falter – either due to depletion of the energy source that sustained them or due to external perturbation. A simple self-organizing system may arise spontaneously, but once it does it cannot influence its environment to support its continued existence (Ruiz-Mirazo & Moreno, 2004).

So, as Moreno & Mossio argue, the distinction between a simple ‘self-organizing’ process-constraint loop and constraint closure is not an arbitrary matter of the number of parts of the system that are acting as constraints.

What distinguishes constraint closed systems from self-organizing ones is not just that there is a greater number of constraints, but that there is a hierarchy of constraints operating at distinct levels and over different timescales – some of which serve as boundary conditions for the possibility of closure between the others.

In this regard, we can point to the significance of a cellular membrane, as an example of this separation. For Ruiz-Mirazo & Moreno the significance of the membrane for the cell is not spatial as Maturana and Varela (1973/1980) describe it in marking out ‘the topological domain of its realization as a network’. Rather, they argue, the key feature of a membrane is that it is not just another constraint within the metabolic network of enzymes, but that it is a higher level, slower timescale constraint that preserves the conditions for that network’s operation, and is regenerated by it in turn. Even in more complex systems such as autocatalytic cycles, which may self organize under quite specific conditions, these cycles still involve only the timescale of the catalysts and of the processes they catalyse (Virgo & Ikegami, 2013; Virgo, McGregor & Ikegami, 2014). The network has no influence on the higher-level constraints that make this network possible – neither on channelling the flow of reactants into itself, nor on the container in which said network is housed.

It is because biological systems require the coordination between at least two different timescales of constraints that they do not just emerge spontaneously in the manner of a ‘one-level’ self-organizing system.

In biological systems as Moreno & Ruiz-Mirazo (1999) argue:

... one has to take into account not only the amount of time that a reaction – or some other process – requires in order to be carried out, but also (and most especially) the time it needs in relation to other reactions with which it could become coupled. In other words, metabolism necessarily requires the synchronisation of a whole set of biophysicochemical processes. (1999, P.51)

What makes this synchronization across different timescales so important is that the activity of said system depends upon spontaneous reactions, where parts of the system release energy through their degradation (exergonic reactions), which is then channelled into the energy-absorbing (endergonic) work of self-construction and repair – the latter processes being unable to occur without the energy released by the former. The problem is that, left alone, these spontaneous reactions may happen too slowly for the constraints upon the energy they release to have any macroscopic effect. Consider the attempt constrain sugar within a piston and expect its slow oxidation to drive an engine. As described earlier in this chapter, we can accelerate this energy release via a catalyst (which is itself just another constraint), and it is through the use of such enzymes that organismal metabolisms are able to extract energy from glucose in order to power the work of their self-repair. But these catalysts themselves only exist *because* of that regenerative and reparative work. In a living system there will be redundancy within an individual constrained process, e.g. an excess of enzymes or energy stores, but if this constrained process fails altogether then it will bring down the rest of the exergonic-endergonic couplings making up that organization and the whole thing will start to fall apart.

So constraint closure requires a number of independent puzzle pieces to come up together in just the right way at just the right time. Such a complex synchronization of energy-releasing and energy-absorbing processes,

coordinated with respect to the various different timescales at which these reactions occur, does not just spontaneously emerge ‘all at once’ as in a self-organizing system – hence why tornados, but not tardigrades, can arise when the weather conditions are fortuitous.

9.4.1. Constraint closure and adaptivity

In so far as there is a separation between timescales, so there is already the possibility for a minimal form of responsiveness to perturbation built into a constraint closed system. Perturbations to constraints operating at one timescale can activate another constraint at a different timescale to compensate for these disruptions. One example, suggested by Ruiz-Mirazo & Mavelli (2007) would be how the production rate of the metabolic network inside a cell raises its osmotic pressure, thereby altering the permeability of the membrane, increasing the rate at which waste products are channelled out of the cell, and so bringing that pressure back into balance

While this form of what Moreno & Mossio call ‘constitutive stability’ is not entailed by constraint closure, it does not require any additional mechanism – only the requirement that these constraints are collectively robust to some degree of perturbation via their modulation of each other. It is hard to see how the delicate synchronization between the various processes and constraints of a constraint closed system could persist over any significant duration without at least some robustness to disruptions.

Constitutive stability, as Mossio & Moreno (2015) describe, is a conservative process that preserves the *same* organization throughout disruption. In this regard, it might be characterized in terms of the homeostatic logic of free energy minimization. For developmental changes that lead to increases in the complexity of a system’s organization, however, we need a further level of

second-order constraints – which they call ‘regulatory constraints’. These second-order constraints are defined by their being dormant over some timescale, during which they do not serve as necessary constituents of the constraint closed system. As such these dormant constraints may change without instantly breaking the constraint-closure that keeps the organism alive.

One example would be repressed genes. Because these repressed genes are typically not active participants in the constraint closed organization so they may mutate, or be altered, without immediately destroying this closure. When this new variation is subsequently activated it may result in the production of a constraint that can synchronize with the overall set of couplings making up the organism – as when the insertion and subsequent activation of the lac operon genes in some prehistoric *E. coli* led it to replace glucose-metabolizing constraints with lactose metabolizing ones. In such a case the organization of a living system may mutate into a new form without any break in its continuity of constraint production.

The nervous system could also be construed in terms of an even more decoupled set of constraints. Repressed genes like the lac operon will, when activated, enter directly into the new constraint closed organization by producing constraints that channel energy into productive work. But even when energy courses through a neuron, it is not directly put to metabolic ends. Instead, the energy entering at the sensory periphery is channelled into coordinating the multicellular organism’s motor system, in order to seek out the sources of energy that the metabolic network needs and to avoid those things that would threaten it. This double decoupling, as Moreno and Mossio suggest, is what affords a much greater degree of adaptivity and plasticity to neurally-equipped creatures.

Importantly, however, while these decoupled sub-systems are not mandated by the notion of constraint closure in itself, they remain part of the constraint closed organization, even when decoupled from its ongoing constitution. Dormant genes may not be contributing to an organism's metabolism over the time period of their dormancy, but their existence depends on their having the potential to make a beneficial contribution under some conditions. If a particular neural assembly does not successfully channel energy into behaviour that supports the metabolic system, it is liable to be rewired into a more beneficial format. So such regulatory constraints are *enabling*, in so far as they facilitate the transition between different constraint-closed organizations, rather than the movement of energy through different production processes within an organization. They are also dependent, in so far as their preservation depends upon their success in doing so. As such, Mossio and Moreno argue, "regulatory constraints are subject to a *second-order closure* between both themselves and the whole *set* of organisations among which they govern the transitions" (P. 35)

So constitutive stability and regulation, are not be intrinsic to constraint closure, any more than adaptivity was to autopoiesis. They are, however, implied by the requirements of preserving constraint closure in a world where things change. Once we have regulation, we not only have an explanation of how something as delicately-balanced as a constraint closed system could persist in such a world but also (the start of) an explanation for how these changes could lead it to evolve into the vast variety and complexity of metabolic organizations that we see today.

As Moreno and Mossio (2015) put it:

Biological organisation must be able to handle variations, and then conserve closure, otherwise it would be extremely fragile and its realisations in the natural world would hardly move beyond a very low level of organisational complexity. Any perturbation would be more likely to drive the system to disruption than to result in an increase of complexity. What is then required for biological organisation not only to remain stable in the face of perturbations, but also be able to increase its complexity? The answer is, we submit, regulation. Biological autonomy requires regulated closure (P. 30)

Thus we can agree with Di Paolo (2005) and Di Paolo, Buhrmann, and Barandiaran (2017) that adaptivity is necessary for a biological system. Yet by taking this to be grounded in something more robust than operational closure among precarious processes, we avoid falling into a view where the question of whether something is a living or autonomous system becomes a matter of the degree of complexity and adaptivity of a network of mutually dependent processes. Unlike process closure, constraint closure allows us to explain how this capacity of adaptivity, which may be graded, is tied to a more fundamental difference in kind.

So constraint closure provides a way of describing the relational logic of a living system that is 'geared in' to the energy flows that must be harnessed for its realization. In doing so it identifies a qualitative difference between the existence of an organism versus that of either an intrinsically stable machine-substance or a spontaneously emerging dissipative system. Moreover, constraint closure gives us an account of the kind of relations of production that must be instantiated at every point in the organism's development, without requiring the organism be fixed to any one particular invariant organization to realize those relations. As such it allows for the possibility of adaptable regulatory mechanisms that can alter the organization

of our constraint closed system without destroying this ongoing process of self-production.

The question now is how this difference in kind makes living systems, and living systems alone, *agents*, that are capable of engaging in intentionally directed action.

9.5 Constraint closure as a basis for intentionality

This thesis began with the claim that the foundational question for the enactive approach is the question of what it is for something to be intentional in the sense of an ‘act having directedness’ towards some norm that it might fail to satisfy. It is this question that the notion of autonomy is supposed to provide a solution to. The issue I took with many theories typically lumped together with the enactive approach, such as sensorimotor enactivism, radical enactivism, or embodied cognitive science more generally, is that they do not take autonomy as a central question. Rather than seeking to account for the normativity of actions they either take it for granted or reject it altogether. Without normativity, what we have is not so much enactivism as mere (sensor-guided) movementism.

To be more than mere movementism, enactivism not only needs an account of normativity but an explanation of how that normativity gets into the movements of an organism. To act, as Davidson’s (1963) classic formulation puts it, is to “do something for a reason.” Thus a theory of the difference between directed actions and mere movements must not only explain what a ‘reason’ (or goal, purpose, or norm) is, as opposed to a mere cause, but also

what it is for some movement to be performed *for* one, as opposed to merely being interpretable in terms of a reason (Wittgenstein, 1953/2010) – as when we say silly things like, “The book stayed perfectly still for the purpose of remaining unnoticed”, or, “The pendulum returned to rest for the reason of minimizing its free energy”.

It is only by answering both these questions, in order to distinguish between acts and mere motions, that we can obtain a robust sense of intentional-directedness that differentiates the actions of an agent (or at least a proto-agent) from the motions of a machine. The ‘reason for an action’, as Hurley (1998) argued, is no more of an unproblematic given than the ‘content of a perception’. Rejecting reconstruction in favour of the practical normativity of actions is no solution if we still lack a foundation for these normative attributions. Without this, the enactivist is just as vulnerable as the reconstructivist to the objection that their ‘intentional acts’ are nothing more than instrumentally-useful abstractions, whose validity is relative to our own explanatory perspective.

The question then is whether constraint-closure can account for how an organism is both a reason, and also a reason that is responsible for its actions.

9.5.1. Organisms as networks of reasons

The idea that we can distinguish between something being done *for* a reason, as opposed to merely being describable in terms of a reason, implies that there is a fact of the matter about whether the reason caused something to happen. Yet reasons are also often opposed to causes in that where the former bear a normative relation to what they are supposed to bring about, which can fail to be satisfied, in the latter case the relation between the cause

and the effect is one of necessity governed by exceptionless natural laws. As Kenny (1989) describes this:

One important difference between the explanatory power of reasons and the operation of causes is this. If there is present a perfectly adequate cause for an effect, then the effect cannot but follow: for a cause—at least on the determinist's view of the matter—is a sufficient antecedent condition for the effect, and if an effect does not follow when an alleged cause is present we know the cause is not a genuine one. On the other hand, there may be a perfectly adequate reason for performing an action and yet the action may not ensue, without this fact casting any doubt on the adequacy of the reason. (P.145)

So when I say, “the mixture of sugar, potassium chlorate, and sulfuric acid shouldn't have broken the glass” I am merely revealing my ignorance of the force this reaction exerts on its container. When my partner responds, “*you* shouldn't have broken my glass” he's saying that there was something wrong or incorrect about the action that led to it breaking, and he will maintain that conviction irrespective of whatever additional scientific details are offered in response. What distinguishes the normativity of reason from the necessity of cause is that, in the former case alone, there is supposed to be a genuine possibility for what *ought* to happen and what actually does happen to come apart. Thus, as Longo and Montévil (2013) describe the role of constraint closure and autonomy in biological explanation:

... a river never goes wrong and we know why: it will follow a geodetics. An onto- or phylogenetic trajectory may go wrong, actually most of the time it goes wrong. We are trying to theoretically understand “how it goes”, between causes and enablement. (P.16)

It is because normative relations between reason and actions appear incommensurable with a deterministic cause and effect universe that reasons or purposes are sometimes either accepted as non-natural (Parfit, 2006, 2011; Enoch, 2011; Scanlon, 2014), or rejected as non-existent (Henderson, 2002,

2010). To provide a naturalistic account of how there can be reasons which could genuinely be responsible for actions then, we need to explain how something can combine both the force of a law with the possibility of its failure.

It is this tension expressed in what Weber and Varela (2002) describe as the Jonasian antinomies of “freedom and necessity, autonomy and dependence.” The same apparent conflict that is found in Kant’s original formulation of freedom as autonomy (literally, ‘self-*law*’) in terms of “the will’s property of being a law to itself” such that “a free will and a will under moral laws are one and the same (Kant, 2008/1785, P. 446–447). It is arguably a tension better captured in the earlier statement that inspired Kant, Rousseau’s assertion that “the impulsion of mere appetite is slavery, and obedience to the law one has prescribed to oneself is freedom” (Rousseau, 2018/1762, 56). The notion of autonomy thus encapsulates the question of how we reconcile the freedom and necessity that defines normativity – the question of how a law can be at once something contingent or optional, such that it is dependent upon the autonomous system that prescribes it to itself, and yet non-arbitrary and binding such that this autonomous system may be nonetheless be subjugated to it.

Arguably the most prominent development of this idea within Anglo-American philosophy is the rationalism of Nagel (1986, 2012) Brandom (1979, 1994) McDowell (1994) Korsgaard (1996) and others who, building on Sellars (1956), point to our distinctive conceptual and linguistic abilities and how the uniquely human discursive practices they enable place us in the “game of giving and asking for reasons” within which participants construct laws that then hold for those participants in turn (Sellars, 1956). For such accounts autonomy is a matter of this sense in which humans alone are

rational beings, capable of communication, deliberation, reflection and the self-conscious recognition of particular evaluative standards as applying to both ourselves and to others.

As Jebari (2019) argues, “The prevailing attitude is that rationalist approaches to ethics are essentially unworkable from within a scientific context and must be abandoned as part of the naturalizing project in ethics” (P.1). This attitude is not only held by rationalism’s critics. Following the Humean proscription against deriving an ought from an is, many rationalists have also defended the independence of this space of reasons from the realm of science and its laws, which is argued to lack the conceptual tools for a description of normativity.

For this reason, and in so far as such accounts are also committed to making rationality a distinctively human capacity, rationalism appears diametrically opposed to the bioactivist project of naturalizing normativity via the capacity that humans share with other forms of life. Yet Jebari (2019) suggests that both these commitments are inessential, and the conflict between the rationalists on the one hand, versus naturalism (and perhaps non-human normativity) on the other, is unnecessary. What is important to rationalism, he argues, is not an autonomous domain of human rationality per se, but the attempt to use this to derive a concept of normative facts as being out “there anyway, whether or not [our] eyes are opened to them” (McDowell, 1994, p. 91), without reifying these norms as eternal essences that are “constituted in splendid isolation from anything merely human” (McDowell, 1994, p. 92).

Furthermore, in a similar vein to Juarrero, he argues that the belief that this desire for objectification-without-reification cannot be incorporated within

scientific naturalism, or physicalism, stems not from the intrinsic nature of scientific explanation itself, but from the philosopher's impoverished concept of what scientific explanation and the physical are. Once we recognize the importance of constraints, he argues, we can see how a scientific worldview has the resources to describe normativity as an objective feature of the natural world. As he describes:

This construal of the rationalist position also provides a way to satisfy the objectification-without-reification constraint. For, on this approach, whether a normative standard applies to an agent is not generally a function of the agent's attitudes; rather it is a function of (1) the overall structure of the social system and (2) the agent's position in that social system. Normative requirements are thus constituted by structures largely external to the agent, and an agent can do better or worse at recognizing and responding to the requirements that in fact apply to her. Nevertheless, the reality of such requirements does not entail Platonism, since such requirements emerge from perfectly natural social-systemic processes. (P.15)

So, for Jebari, it is because these constraints are both constructed and realized by social systems that they are, unlike necessary laws of nature, contingent principles that may be violated. But it is because a society is nonetheless a real structure, capable of limiting the behaviour of its members, that these constraints are both objective and naturalistic features of the world. Thus as he explains, "People's behaviour will both explain and be explained by these constraints, yielding an overall picture in which people's actions operate in both a norm-guided and norm-constituting capacity, often at the same time" (P.16).

Still, while Jebari emphasises the relevant constraints as being those that are constructed, makes reference to the work of Moreno, Mossio, Kaufmann and others, and discusses how constraints must be organized so as to maintain the system as a whole, his focus remains on their *social* construction. As a result, he discusses only one side of the potentially dual nature of a

constraint, in terms of how it constrains, not how it may also depend upon flows of energy for its repair and reproduction. As such he does not address the unique thermodynamic status of biological constraints as intrinsically unstable, which I take to be essential to what makes them distinctively normative.

Constraints are everywhere. Constraints that we have constructed are pretty common too, from the machines we make to the canals we build, but the manner in which steam engine is forced to move when coal is burnt or a canal prevents me from walking somewhere does not seem to me to be a normative affair but a straightforwardly physical one. It is true that these constraints are contingent in the sense that they need not hold. There is no inevitable exceptionless law that matter must form into engines, burn coal, and power steam trains. We didn't have to construct that steam engine and we could destroy it. There is also an indirect sense in which these constraints are dependent upon the activity they enable for their continuation. If the steam engine does a bad job of turning thermal energy into motion we might melt it down for candlestick holders. Yet this risk is not inherent to the steam engine itself. We might not destroy it, but place it in a museum for schoolchildren to goggle at, where its structure can persist indefinitely without actually constraining anything.

This steam engine *can* constrain energy flows to perform work, but there is no necessity to its doing so. Its constraints are only of the conditional form, *if* there is a flow of energy then raise a piston. The raising of the piston does not have the force of any sort of necessity and there is no reason, no need, for the steam engine to do anything at all. In contrast, a constraint closed organism must be constantly operative because it is only by constraining energy flows that it can enable the regenerative work without which the set

of constraints that realise it would degrade – irrespective of what external agents like us might choose to do about the situation. As Nicholson (2018) describes.

This ongoing self-producing activity is not optional—not undergoing constant metabolic regeneration is not a possibility. The thermodynamically grounded fact that organisms need to keep acting in order to keep existing helps to account for the emergence of a rudimentary form of normativity in nature (cf. Mossio et al. 2009; Christensen 2012). It is because its existence depends on its own activity that an organism must act in accordance to the operational norms that enable it to persist through time. If the organism stops following these norms, it ceases to exist. What this means is that it is in principle possible to objectively specify what is intrinsically ‘good’ or ‘bad’ for an organism (that is to say, what is and what is not in an organism’s ‘interest’) by evaluating its activities according to the contribution they make towards the preservation of its organization in far-from-equilibrium conditions. (P.154)

Such precarious constraints need to enable regenerative work to continue to exist – as long as they exist regenerative work *must* be done. In this respect they have a kind of necessity baked into their existence that ordinary constraints do not. Yet they may also fail. They may not receive the energy they need to do this work. But if they fail, they fall apart. A constraint closed system needs to enable its own existence, but it is also free to cease to exist. In this sense, a constraint-closed system is itself a ‘self-given law’.

Even before we get to any uniquely human capacities of conceptualization and communication, the constraint closure of organisms already provides the means to break apart the ‘ought’ from the ‘is’. To say this particular *E. coli* bacterium ought to avoid ethanol is not a human projection but an existential imperative derived from the bacterium’s self-producing organization. What it is to be that bacterium is to be something that cannot exist in a highly concentrated ethanol solution. Nonetheless, that particular bacterium might still fail to comply with this existential constraint and thus it will cease to be a bacterium any longer.

This I not to say the demise of a single constraint is the death of the organism, or that the norms they realize cannot change. I argued that what gives the organism continuity over time is not one particular constraint-closed organization, but an unbroken relationship of self-production throughout its various organizational stages. For one particular *E. coli* at one particular point, it can be an objective fact that it cannot exist in a highly concentrated ethanol solution, but this *E. coli* might be engineered to survive on ethanol (Cao et al., 2020). From that point on, ‘avoid ethanol’ is no longer a norm that it must follow. Changeability does not make norms arbitrary or subjective, however, and a fact need not be eternally true to be objective. Prior to the development of ethanol resistance, the need to avoid ethanol was an objective fact about the particular organization of a particular *E. coli*. Post-insertion of the relevant genes this bacterium takes on a new organization such that this norm no longer holds.

Because particular individuals can have variability in their organization, so there is no generic fact about whether lactose is ‘good’ for *E. coli* in general. The goodness of lactose depends upon whether the specific individual *E. coli* in question has incorporated the lac operon genes into its overall organization, allowing it to channel the energy contained in lactose into its self-production. On the other hand, there is no absolute fact about whether the development of lac operon genes will be good for that particular *E. coli* bacterium either. This depends upon facts external to that particular bacterium namely whether its environment contains a source of lactose that its constraints can channel the energy from.

The only restriction governing these changes in organization is that any new constraints must synchronize with the rest of the exergonic-endogenic couplings making up the organism in the same manner as the one that it

replaces, such that the new organization remains a coherent self-producing whole. As far as any particular organism is concerned, the requirement of ongoing production, in whatever organizational guise that might take, is the only ultimate norm that cannot change.

As argued at the end of Chapter 8, this means that the description of biological systems must take on a different form from that of ordinary constrained physical systems – in which the invariant constraints and resulting dynamics can be specified at the outset, and do not depend on either historical change or their interactions with the external environment. Thus, the theoretical biologist and the rationalist can agree that the explanation of norm-governed systems is distinct from the kind of explanations given in ordinary physical models. Yet this does not mean that our account of normativity must be non-natural, only that these particular types of modelling practices fail to encompass all the idiosyncrasies of nature. Naturalistic is a slippery term and I don't want to get into policing its borders here, but in so far as the basic materials of constraint closure involve nothing more mysterious than energy flows and constraints upon them – language that would legitimately be found between the pages of *Physical Review* – so it should qualify as naturalistic.

Jebari's account is useful in showing how the generality of the concept of a constraint means that it can be applied to social organizations as well as biological ones – making it apt for scaling the normativity found in molecular metabolism up to the different levels at which normativity might be constituted. Nonetheless, to understand what makes a constructed constraint a genuinely *normative* one, we need to look first to the biological realm, to see how the precarity of certain types of constraints place them in the more fundamental game of giving and asking for **thermodynamic** free energy. It's

only once we've properly characterized this intrinsic instability and mutual dependence, that we can then look to see whether this also characterizes constraints at the social scale

This is the sense in which I take organisms to be 'natural purposes in Kant's terminology. But to be agents, or at least proto-agents, they need not only to *be* purposes, they also need to act *for* those purposes. Thus Korsgaard (2018a) allows that animals can have value and that things can be good or bad for them while rejecting the idea that they are necessarily autonomous systems who are capable of acting for the purpose of obtaining those goods. For Korsgaard, this latter attribute requires a thick Kantian type of rationality as a capacity of self-conscious reflection that makes a human being alone capable of "knowing that an evaluative standard applies to your conduct, that there is a way you should act or ought to act or that it is good or correct to act, and being motivated in part by that awareness" (Korsgaard, 2018b, P.5).

This view of what it means to be rational, or to act *for* a reason is not particularly helpful to an enactivist, for whom it is exactly these capacities of being conscious, knowing, or aware, that we want to use an account of autonomy to try to explain. But there is, I will now argue, a much less demanding, and more naturalistically-grounded sense in which we can argue that living systems are not only reasons, but reasons that cause the very actions which they are reasons for.

9.5.2. Organisms as causes of their own activity

To say that an organism is both the reason for, and the cause of its action amounts to saying that an organism is a cause of itself – an idea that, as Juarrero (1999) argues, is impossible within the Newtonian framework of event-based causation, where every change must be caused by some other

event. If this were an accurate description of the natural world, then to say the movement of a system cannot be traced back to through a series of events to some other causes external to it (and, in principle all the way back to the unexplained origin of the universe) would indeed be a supernatural claim (see Juarrero 2019, Ch 1 in particular). As we have seen, however, the natural world is not, at its basic level, Newtonian. Energy will dissipate and entropy will increase with or without external prompting. The macroscopic effects of this dissipation are not determined by some prior, external change, but by the presence, or absence, of invariant constraints upon those flows of energy.

The idea of constrained energy gives us a different approach to causation but it does not yet give us self-cause. Machines also operate by constraining the spontaneous release of energy from an unstable reactant, i.e. fuel, in order to perform work. But organisms are different from machines in two ways. Firstly, as described, the work that they do involves the reproduction of their own constraints, which would disintegrate the regenerative activity that they enable. Secondly, organisms do not only regenerate themselves, they also consume themselves. In machines, there must be a separation between the reactant that releases energy, and the constraints that channel it. There is also energy contained within a machine's structural parts and this will slowly spread into the environment as those parts rust. But the machine has no means to channel this spontaneous energy-releasing reaction into the work of rebuilding itself. The breakdown of a machine's structure can only ever lead to energy dissipation and destruction.

So when an engine runs out of its fuel supply it ceases to operate. But precisely because organisms need to rebuild their component parts, so they are also free to accelerate their break-down, via catalysis, and to channel this release of energy through other constraints in order to power further activity

directed towards the repair of these, and other, parts of itself. While in Moreno & Mossio's (2015) diagrams of constraint closure, constraints are only depicted as the output of a process, their account also allows that a constraint may also degrade, ceasing to constrain and becoming instead the reactant for another process.

Crucially, this gives us a sense in which organisms are *intrinsically active*, rather than just responsive to perturbation as were Ashby's 'sleeping machines' or Friston's free-energy minimizers. Like a leaf in the wind, the latter may appear quite animated if their environment is disruptive enough, but their intrinsic dynamics only drive them back towards stasis. In contrast, when an organism is deprived of any input of energy it will continue to operate, first by catalysing the breakdown of non-constraining energy stores, then via the breakdown parts of the structure by which it operates – such as muscle tissue for amino acids (Steinhauser et al., 2018). Systems that channel energy released from their own dissipation in order to rebuild themselves thus have “their own endogenous dynamics, continually running through their cycles whether perturbed from the outside or not.” (Pickering, 2010, p. 164)

This is not to say the breakdown of internal constraints and other energy stores are the organism's only power supply. Nothing, living or otherwise, can perfectly channel free energy into work and in every exergonic-endogenic cycle some free energy will be lost through its dispersal as heat. The only reason the organism is able to continue this self-production for extended periods is because of its openness to energy and the continual 'top-up' of its reserves from the environment in the form of food or sunlight. Nor is it to say that organismal behaviour cannot be influenced by external events. But while these influences shape what an organism does, they are not why does anything at all. Deprived of either energy or other external promptings an

organism will still be active – right up until the point it loses the ability to channel the *thermodynamic!* free energy stored in its own body into work and dies.

The reason why an organism acts is the reason that it embodies in its energetic resources and the precarious constraints that channel this energy into the work on which its existence depends. An organism is genuinely the initiator of its movements in a way in which a machine is not, thus, if an organism is a reason, then this reason is itself a cause.

Conclusion

It is only by appealing to both thermodynamic and relational properties, as Moreno and Mossio do, that we can describe the precarious dependence of an organism's existence upon its own activity, in virtue of which it constitutes a self-constraining, self-producing system. In doing so we see why, unlike inorganic structures that can be expressed by universal laws, as Merleau-Ponty put it, "organic structures are understood only by a norm," (1963/1942, P.148). While a precarious constraint organization will mandate certain interactions as both a consequence and as a precondition for so long as it persists, the organism can also abandon a failing organization in favour of a new one in service of the ultimate norm of self-production by whatever means necessary. Unlike a law, this norm too may genuinely fail, but when it does the entire existential world of the organism collapses with it.

I have explored this account of autonomy and autopoiesis in terms of constraint closure mainly at the level of simple organisms. Yet the beauty of the notion of a precarious constraint is that it is at once thermodynamically restrictive enough to avoid attributing normativity and intentionality to any stable mechanism, while being general enough that it could, in principle, be instantiated at all different scales of biological organization.

Still, this may not be enough for you to be willing to consider an organism a rational agent, or even an intentional agent at all. You might take this to require epistemic capacities, like the ability to detach what appears to you from what is and entertain the possibility you might be wrong. You might see it as requiring not only regulation in response to external disruption but the ability to proactively anticipate such disruptions and adjust in response, or to explore the consequences of an action in offline simulation prior to

executing it. Or you might require the intersubjective capacity to view others as purposive, intentional agents such as yourself, and to coordinate your behaviour with respect to *their* norms.

Such capacities are useful markers with which to judge the conceptual or inferential abilities of a living system. Perhaps they might provide a means to track the emergence of self-consciousness and justify the attribution of it to some animals and not to others. They might provide a means to distinguish between merely acting for a reason and *knowing* that you are acting for a reason. But to propose these as further, necessary, requirements to be an agent is just to say that for a movement to be caused by a reason does not suffice for it to be the act of an agent.

Thus Hurley (2003), unlike Korsgaard (2018a), argues for the treatment of animals as rational, intentional agents, rejecting the idea of tying practical rationality to either conceptual and inferential capacities, or to something like conscious awareness – neither of which account for the *origin* of the normative standards that they are supposed to make us uniquely capable of acting in accordance with. As Hurley argues, the flexibility of behaviour that is typically taken as an indicator of conceptual capacity comes by degrees. It does not even require a nervous system. The humble bacterium is capable of adapting its behavioural strategies in pursuit of the ultimate norm of continued self-production – for instance by changing from an organization that must seek glucose to one that must seek lactose when the former set of norms lose viability. Groups of bacteria may even share strategies amongst themselves, in the form of parcels of DNA-based constraints, called plasmids, that enable the production of other constraints which have proved beneficial in their own case.

This is not to say that bacteria are capable of conceptualization and reasoning or intersubjective and social coordination in the same manner that we are – just that the attempt to draw a line somewhere on this gradual scale of socio-cognitive ability and announce, “*here* is where the ability to follow norms begins” looks like an unavoidably anthropocentric stance. The enactive approach offers a much more promising strategy, in reversing the direction of explanation to describe how these more advanced capacities might progressively develop from the adaption and evolution of creatures that are already intentional agents.

In her shared circuits model, Hurley (2008) puts forward an initial proposal for how this might work, describing how a layered feedback control hierarchy could provide a subpersonal mechanism for the capacities of ‘imitation, deliberation and mindreading’. Grounded in the more fundamental capacity of online motor control, the shared circuits model explores how the development of the capacity to improve this via offline simulation, to predict the potential consequences of one’s own actions, could enable the possibility to simulate potential actions of others and infer the hidden causes behind these “thereby enabling strategic social deliberation” (P.1).

While predictive processing may have initially been proposed as an account of how we infer the distal structure of our environment, its basic requirements for the minimization of prediction error at different levels of detachment and spatiotemporal grain, are the same as for Hurley’s account. Reframed in the same enactive terms as the shared circuits model, PP could provide an implementation mechanism that the enactivist can use to explain how more advanced socio-cognitive capacities could emerge out of the more basic capacities of sensorimotor control.

Yet the coordination of movements to maintain inputs at a stable, predictable state is not, in itself, a *need*. Hurley points to the requirement for a teleological context to ground normative attributions to a control hierarchy but does not supply an account of this teleology. Similarly, non-reconstructive versions of predictive processing have typically gone only halfway towards an enactive account – arguing that the predictive brain is directed towards something other than accurate reconstruction but offering no justification for interpreting the tendency for some neurons to become correlated with others in terms of a purpose that the brain is trying to fulfil. Without this, we have no basis for attributing function or purpose to a ‘predictive’ neural architecture, no sense in which it can genuinely fail, no means to take its dynamics as anything more than the externally-determined movements of a law-governed machine, and no reason to consider it any more of a predictive model than one pendulum coming into sync with another.

The free energy framework, and particularly its principle of free energy minimization initially appeared as a solution for this problem. Said principle, Friston (2012) claimed, could “unify all adaptive autopoietic and self-organizing behaviour under one simple imperative; avoid surprises and you will last longer” (P.2). As free energy is analogous to prediction error, so predictive processing looked nicely placed to provide a mechanism for how organisms achieve this ‘existential imperative’ and to describe how abilities like offline simulation and counterfactual reasoning might be grounded, as Ramstead, Friston, Badcock (2017, P.33) put it in, “The ‘intentionality’ or ‘aboutness’ of living systems – that is, the directedness of the organism towards a meaningful world of significance and valence,” which, they argue, “emerges as a natural consequence of embedded adaptive systems that satisfy the constraints of the free energy formulation.”

Unravelling this claim was a long, frustrating, and ultimately disappointing experience. What lies beneath the FEF's twisted thicket of mathematical terminology is just a formulation of survival as stability in the face of perturbation and the insight that such stability can be formally redescribed in inferential terms. An equivalent analysis of survival was formulated by W.R. Ashby half a century earlier who, like Maturana subsequently, saw it as precisely the means to *eliminate* all 'metaphysical complications' of purpose, intentionality, teleology, or function from biological explanation. As Ashby argued, any stabilizing thing whether a computer or a pendulum can be interpreted as 'rejecting' unstable states to 'seek out' a stable equilibrium. If intelligence, agency and intentionality reduce to nothing more than prediction-error-based control, and if such control is attributed to any stable system, then they are either everywhere, or nowhere at all.

Pitched as an 'existential imperative' it turned out that the requirement of stability had even worse problems than its generalizability. While the requirement might be softened in various ways to make it more trivial, living systems are precisely those things that are most likely to violate it. Modelling an organism as a free energy minimizer may work well for describing some of its behaviours over some timescales, but the properties that define this model do not define the identity of the living system itself, which is liable to change in ways that cannot be predetermined by its organization alone. This is not just a problem for the FEF, however, but for any attempt to define living systems in terms of the invariant logic of a machine.

The prevailing formulation of autonomy within the enactive approach is more amenable to the kind of open-ended change that we observe in the development and evolution of living systems. Yet for bioenactivism to

succeed this definition of autonomy needs to include a robust enough notion of self-production, such that we can follow the organism through these organizational changes, account for what renders its existence fundamentally different from that of a machine, and license a realist attribution of agency and intentionality to organisms alone. Process closure, I argued, fails in these respects.

This does not mean that the bioenactive project fails too. As I have described Moreno, Mossio and Montévil's work in constraint closure succeeds on all three counts, providing a naturalistic grounding for the application of normative notions such as purpose, intentionality, goals, or function to living systems alone. Moreover, their account of regulation in terms of second-order constraints provides the opening for extending these intentional attributions to extra-metabolic systems, like the nervous system, in terms of the role these play in coordinating the transition between different constraint-closed organizations in order to preserve an ongoing process of production.

In so far as the organism, and its brain in turn, is trying to achieve the continuation of a process by whatever means necessary – rather than the stability of some particular set of states or relations – so its function will not be reducible to free energy minimization. Yet, as I acknowledged, the FEF is still useful in building approximate models of a system. A wide range of organismal behaviours are, after all, homeostatic. If this homeostasis and stability are not *necessary* for survival, then why would they be so common?

A potential answer to this may be found not in enactivism, but rather in the epistemic anxiety of Hohwy and Helmholtz. In so far as we cannot determine in advance whether some new interaction or organization will support the

continuation of constraint production, so, even as enactivists, we face the sceptical challenge of attempting to coordinate ourselves with respect to unknown factors. In order to do so, living systems need to make a guess at how things around them might unfold. The simplest one is to presume that they will carry on just as before – and that if a particular organization has proved viable in the past it will continue to be so in the future. On such an assumption it makes sense to attempt to keep the state of ourselves and our environment within the stable bounds that this particular organization requires for as long as possible, before taking the risky leap into a new form.

So survival may not be equivalent to the minimization of surprisal but, when in doubt, avoiding surprisal may still help secure it. Just as this free energy approximation can prove useful to scientists seeking to model living systems, so it may be a useful way for organisms to approximately model themselves, for the purposes of predictive control.

That evolution of this into *hierarchical* predictive control may lead to all sorts of wonderful things like the ability to predict others, to entertain different possible expectations about what will happen, to try out inconsequential predictions, make non-fatal mistakes and to live through this discovery that we can get things wrong. Still, what matters is not the complexity of the control hierarchy required for these capacities, nor the variety of perturbations it can respond to, nor how many different stable states it might alternate between. What matters is the precarious status of the intrinsically unstable organism that both realizes it and depends upon its successful operation.

We could copy this hierarchical control structure and grant it the physical ‘embodiment’ of effectors and sensors hooked up to the sort of inputs that

matter to us. Yet in so far as this system is built out of parts that have nothing to lose by its failure, any interpretation of it as trying or failing will only be a projection of our goals onto the movements of an utterly indifferent machine.

Appendix: What's the use of a concrete blanket?

Over the past year a growing number of critics of the FEF have begun to object to the ease with which papers by Friston and co-authors slide between the standard, heuristic formulation of a Markov blanket and a stronger ontological one (Bruineberg, Dolega, Dewhurst, & Baltieri, 2021; Raja, Valluri, Baggs, Chemero, & Anderson, 2021; Menary & Gillett, 2020).

Friston and Allen (2018), for instance, move from the uncontroversial statement that, “The boundary (e.g., between internal and external states of the system) can be described as a Markov blanket” to the above description of a Markov blanket as a real boundary that the system must actively conserve, and upon which its very existence causally depends. Similarly, Ramstead et al. (2018) switch freely back and forth between the Markov blanket as way of describing some statistical relationship in the language of a Bayesian network, versus its being the thing in the world that produces the very conditional independence that the model then describes:

Markov blankets establish a conditional independence between internal and external states that renders the inside open to the outside, but only in a conditional sense (i.e., the internal states only ‘see’ the external states through the ‘veil’ of the Markov blanket... With these conditional independencies in place, we now have a well-defined (statistical) separation between the internal and external states of any system. A Markov blanket can be thought of as the surface of a cell, the states of our sensory epithelia, or carefully chosen nodes of the World Wide Web surrounding a particular province. (P.4)

The location of a Markov blanket, they go on to say provides the basis for “a fully generalizable *ontology* for biological systems”(P. 5) [My emphasis].

This realist interpretation is not a part of the original concept of a Markov blanket, and neither does it just follow inevitably from the mathematical core of the FEF.

As Bruineberg et al. (2021) point out, “metaphysical consequences require metaphysical premises, and cannot simply be read off the formal model.” At first appearances then, the reification of Markov blankets looks like a classic case of confusing properties of the model for properties of the target system. Such missteps, Andrews (2021) argues, are endemic in the FEF literature where terms like ‘entropy’ or ‘energy’, appropriated from thermodynamic systems to describe analogous statistical properties, are then misinterpreted as retaining implications that follow only in the limited context of describing constrained flows of matter and energy in concrete physical systems.

Realism about Markov blankets might thus be taken as just another instance of this ‘fallacy of misplaced concreteness’ (Whitehead, 1925). This is an ever-present danger with scientific models, where utility often comes apart from representational fidelity (Morgan & Morrison, 1999; Potochnik, 2017). Still, if realism about Markov blankets is indeed a fallacy, this would not undermine the FEF’s legitimacy outright. Andrews (2021), for instance, suggests shedding this extraneous pretention to describing an objective feature of all living systems, in order to separate out the mathematical core of the FEF as a purely formal model structure, from both conceptual and empirical questions regarding its applicability to any particular target system.

Had advocates of the FEF stuck to this more abstemious interpretation of its key constructs, then their work would have likely generated far less confusion and controversy. It would also probably have garnered much less attention for its authors. Unfortunately, as we have seen, this distinction

between the pure formalism, versus the various philosophical theories and models that have been constructed upon it, has rarely been respected within the FEF literature. Claims that the FEF might provide a first principle of living systems, one which could subsume autopoiesis & autonomy, fall squarely within the purview of philosophy and theoretical biology – not mathematics. In critically analysing the FEF as a theory of life, it is specifically these extraneous philosophical claims that I'm concerned with. As regards Markov blankets, said theory of life requires that we can indeed consider these to be a real entity that is possessed by living systems, and living systems alone, which has causal influence to the effect of preserving and distinguishing them from their surroundings.

In the interests of steel-manning this view then, it's worth pointing out that while proponents of ontological MBs have not provided the argumentation to support this realist position, neither have the FEF's critics shown that such realism is incorrect – only that the justification for it is currently lacking. The fact that (as discussed in the previous section) those Markov blankets that appear in our model are sensitive to a host of initial choices – such as which variables to use and the coarseness of grain with which these were individuated – does not entail that there is not genuinely something like a set of 'real Markov blankets' in the actual structure of the world. Nor does it preclude the possibility that our partial models and often incorrect models may sometimes allow us to accurately identify them. It may be a platitude of the modelling literature that 'all models are wrong', but it is presumably also true that all useful models must be getting at least something right. Some features of our maps must really be features of the territory, or they will not be maps at all.

So, if Friston and co. have not supplied the metaphysical premises required to support their metaphysical conclusions, can we derive these on their behalf?

A.1. God's great causal graph

What would a causal modeller have to say about the reality of Markov blankets? Well, firstly, they would likely point out the fact that every causal graph has them is simply a consequence of the presuppositions that direct such a graph's construction. For our purposes, the relevant ones are the following 1) **Decomposability**: the system can be broken down into a set of discrete parts and the connections between them. And 2) interactions between these parts respect **the Markov condition**, such that, the state of each is independent of the state of its non-descendants, conditioned upon its parents. This second criteria ensures the factorizability of our graph in terms of Markov blankets and becomes the *causal* Markov condition when said graphs are interpreted as causal models.

So, the issue of whether Markov blankets are merely a 'statistical device,' or a 'necessary attribute' of 'the system itself' comes down to whether these are simply useful falsehoods to help us approximate the behaviour of said system, or whether they accurately describe its structure and interactions between its parts.

Among the developers of programmatic causal discovery, these axioms of construction have primarily been defended on a pragmatic basis, not as either analytic truths about our concept of causation, or as principles reflective of some fundamental law of nature (Spirtes, Glymour & Scheines, 2000; Pearl,

2000; Stafford, 2005; Weslake, 2005). Objections to the metaphysical implications underlying such models tend to be overlooked, if not met with outright hostility (Glymour, 2010). Thus we find Glymour (1999) defending the causal Markov condition's legitimacy, not via what he deems as 'Socratic analysis', but rather by means of a list of cases in which methods premised upon it have been able to predict the result of interventions from observational data alone. As he puts it, "The essential issue in scientific discovery is the right representation for reliable, efficient search, not the metaphysical disputes upon which philosophy of science is fixated." (P.64)

Still, while the causal modellers may not be overly concerned, the position that these axioms of construction might be more than pragmatic – the proposal that they describe a necessary feature of reality – is one that predates the FEF. Decomposability, the claim that reality is built up out of individual building blocks, has been a (surprisingly) resilient cornerstone in the history of science and philosophy, from the atoms of Leucippus to the individuals of classical logic. Its resilience is surprising because the status of this longstanding doctrine as a description of fundamental reality has become increasingly threatened by developments, such as quantum field theory, that suggest that such apparently 'elementary' particles may not constitute basic reality after Hobson (2013). As Bickhard (2015) puts it, "According to our best physics, there are no particles and what are called particles in contemporary parlance are quantized excitations in quantum fields" (P.24)

The idea that the Markov condition is no mere modelling tool, but rather a constraint that actual causality must satisfy, does not have quite so long a heritage. Still, it can still be traced back through philosophical discussion of the Markov condition's ancestor in Hans Reichenbach's (1956) analogous

‘principle of the common cause’³⁸, proposed not as an optional heuristic, but rather a necessary requirement governing the relationship between correlation and causation that Reichenbach attempted to place on a par with the second law of thermodynamics.

In an attempt to give more metaphysical bite to SGS’ models, Hartry Field nicely expresses what reality might look like if one accepted these two tenants, suggesting that:

“Intuitively, it seems (barring quantum non-locality and the like) that one should be able to think of the physical universe as a causal system with a node for each space-time point, with the value of the node expressing the totality of the values of physical quantities at that point; the light cone structure gives the dependence relations. (2003, P. 447)

I will call the position that causal graphs accurately capture the structure of reality **causal graph realism**.^{39*40} For Field, the fundamental nodes of reality are microphysical space-time points, or point-sized particulars, that cannot be further decomposed, a position defended more extensively by David Lewis (1994) as ‘Humean supervenience’. But being a realist about causal graphs does not necessarily require the commitment to microphysical reduction. As part of his defence of the reality of non-reducible macrophysical phenomena, Papineau (1992) takes causal graphs to provide an accurate metaphysical picture, but one that concerns relations between

³⁸ The differences in content between the causal Markov condition and the principle of the common cause are not important for our purposes, but for more on this see Hausman and Woodward (1999)

³⁹ This shouldn’t be confused with a holistic graph realism in which there are no properties of individual particulars, but all properties supervene solely on the relations instantiated in the graph.

⁴⁰ While, as explained previously, causal modelling tends to use Directed Acyclical Graphs, and specifically Bayesian networks, Friston typically deploys cyclic graphs in discussion of the FEF. For the purposes of the Markov blanket, the specific form of the real causal graph is not essential, only that it connects probabilities over states of individual units in a way that satisfies the Markov condition.

generic states of affairs – eg. between the prevalence of smoking and prevalence of early mortality in population – and which cannot be reduced to microphysical dynamics (where causal relations look to be absent altogether).

The causal graph realist need not be troubled by the fact that many of our causal graphs will be wrong, or that the inevitable simplifications made in the construction of such models invariably divorce them the real structure of the world. As Spohn (2001) argues, we can still conceive, in principle, of a graph constructed with an ‘all-encompassing frame’ that would capture *all* of the correctly individuated units and their direct relations (and so, ultimately, all the ‘ground-truth’ Markov blankets). Such a graph, he argues, would capture all there is to say about causation. The correctness of a particular simple model, for the causal graph realist, is simply a matter of its similarity to this all-encompassing graph.

Causal graph realism may (but need not be) combined with support for the **statistical reduction of causation**,⁴¹ to claim that the correct causal graph contains nothing more than all the statistical relationships between our various nodes.⁴² This reduction is developed by Spohn (2001) who explicitly extends SGS’ work to argue that “Bayesian nets are all there is to causal dependence”. While the idea is essentially Humean, similar reductions have been developed by Reichenbach (1956), Good (1959) and Patrick Suppes

⁴¹ This statistical reduction of causation is not incompatible with a deterministic account, as Papineau (1989) notes in defending both. Determinism merely involves raising the threshold of statistical relevance for one thing to be a cause of another, such that A only actually causes B if its occurrence raises the probability of B to 1.

⁴² A causal graph realist doesn’t have to adopt this however. They might maintain that all the statistical facts in the world would not suffice to determine the true causal graph, and that additional facts - regarding say counterfactuals, the effect of interventions, or temporal asymmetry - are required. To defend a probabilistic-reduction of causation requires arguing that these facts, in so far as they pertain to causation, may also be reduced to statistical ones.

(1970) & Papineau (1992) (for reviews, see Salmon, 1980 & Weslake, 2005). On such accounts, the gap between correlation and causation is not due to their being different things entirely, but rather due to the incompleteness of our correlational information rendering it insufficient to determine a single causal graph.

These accounts differ, however, on whether they consider the reduction to be conceptual (Spohn, 2001; Suppes, 1970) or metaphysical (Papineau, 1992, 1993; Field, 2003). Where for Suppes, probabilistic/causal relationships are relative to a scientific model, the metaphysical reductionist is committed to there being a true set of probabilistic relations that make up some objective causal graph.⁴³ On this view, the all-embracing Bayes net does not just describe how we conceptualise causality, but correctly describes the structure of reality itself – the axioms of causal modelling are neither contingent facts about the relationship between causation and probability, nor analytically true in virtue of how we conceptualize causality, but simply describe what causality actually is.

Let's call the combination of the strongest version of these two positions 'probabilistic graph realism'. There are a host of difficulties in holding such a view. With regard to the probabilistic reduction of causation, the most obvious is the need for a robust account of objective probabilities. It's also interesting to note that while Pearl (2000) or Glymour (2010) avoid taking positions on the metaphysics of causation they explicitly reject a probabilistic reduction as being sufficient even for the methodology of causal investigation, let alone an ontology of causal relationships. Instead, Pearl (2001) advocates an 'interventionist' account (Woodward, 2005) arguing that

⁴³ This of course requires a worked out notion of objective probabilities but that's a whole other kettle of fish that we don't have time to poach.

disambiguating unique causal structures requires going beyond mere statistical relationships to include the effect of targeted interventions to fix the value of a particular variable. In more concrete terms, the kind of disambiguation of causal relationships that we achieve by moving the needle on a barometer and seeing if it causes a thunderstorm.

Other difficulties include deriving the temporal asymmetry required for a causal graph (and crucial in the FEF's division of the Markov blanket into parent/child nodes). This cannot be simply delegated to fundamental physical laws, which are time-symmetric, and it is controversial whether it can be derived either from within purely statistical asymmetries or from some other external source, as in Reichenbach's enlistment of the second law of thermodynamics (Price, 1993). Further, if we take the units of Spohn's 'all-embracing frame' to be microphysical (for instance the local spacetime points of Field (2003) and Lewis (1994) we not only lose temporal asymmetry but also decomposability and Markovian behaviour, which are upended by quantum entanglement and indeterminacy (Glymour, 2006; Arntzenius, 1992; Cartwright & Jones 1991; van Fraassen, 1980, 1982).

Nonetheless, I think that interpreting the realist manner in which Friston talks about Markov blankets to be a product of an implicit commitment to the metaphysics of probabilistic graph realism is a more charitable, and more plausible, interpretation than taking him to be merely mistaking the constraints of a particular modelling framework for real features of modelled systems.

This implicit metaphysics makes sense of more than just Markov blanket realism. It also bears upon the discussion in section 3.2 about the status of the generative model and claims that the interaction between an organism

and its environment is not only described by, but literally embodies, this joint distribution. A statistical graph, recall, represents a joint distribution, the generative model, plus a (typically temporal) ordering. To say that the structure of reality is such a statistical graph is to say that it actually has the properties of this generative model. So to say that systems literally are generative models, and that they really have Markov blankets as their parts, is not to say that these systems are representations of anything else, or that they are tools being used to some epistemic end. It is to say that the system itself has the same kind of properties of as the partial graphical models that we build of it. Under this particular metaphysical view the gap between the model of a scientist, and the real system itself would be merely in degree of detail and not a difference of kind.

A.2. Naturalised mathematical realism

Still, even if this delivers objectivity to claims about Markov blankets as a description of the true causal graph, does this confer independent existence upon the Markov blanket itself? Majid Beni (2021) argues not, pointing out that Markov blankets are still a mathematical object. As such, he points out, to treat them as something that exists in its own right, rather than as an accurate description of a really existing *physical* object, must either be a category error, or the result of a commitment to realism about mathematical objects⁴⁴. A similar point is raised by Bruineberg et al. (2022) who argue for a distinction between treating the Markov blanket as ‘literally’ an entity in the world, vs seeing them as ‘realist’ descriptions of some other feature that does have worldly existence.

⁴⁴ Not platonist, but Pythagorean

A commitment to the FEF as providing a ‘formal ontology’ as advanced by (Ramstead et al, 2021), who propose its use as, “a mathematical formalism to answer the questions traditionally posed by metaphysics; i.e., what does it mean to be a thing that exists, what is existence, etc.”, does imply a commitment to some form of mathematical realism. However, both Beni (2021), and Menary and Gillett (2020) in a similar criticism, interpret this realism as of the Platonic form. Platonism may be the best known mathematical realist position, but its defining feature is the treatment of mathematical objects as transcendent entities existing outside of space and time. That might fit with understanding the FEF as a purely formal model that nevertheless says something true, as Andrews (2021) suggests we take it. But this cannot be the kind of reality Friston takes a Markov blanket to have if it is also intended to be a property or thing that particular systems really have which can enter into the kind of causal relations that preserve that system’s existence.

As applied to living systems the FEF seems to require the Markov blanket to be a simultaneously physical and mathematical entity. As such, I suggest this application of the FEF, would be better paired with the kind of naturalised mathematical realism developed by Penelope Maddy (1990, 1997) under which there are physical objects that instantiate the properties of a mathematical object – specifically, for Maddy, the property of being a mathematical set. If the world is indeed structured like a statistical graph, respecting the Markov condition, then for any system there will be a set of objects that have the property of rendering the state of some further set of things probabilistically independent of everything else. Naturalised mathematical realism would allow Friston to describe this set as a physical object that literally has the mathematical property of being a Markov blanket.

A.3. Absolute units

We now have the elements that seem to be behind the principled commitment to Markov blankets: causal graph realism, the probabilistic reduction of causation, and naturalised mathematical realism. Let's say we adopt all those. There really are physical sets that render a particular thing independent of everything else. Where, exactly are these 'Markov blankets' then? Are they in the room with us right now?

The answer to this depends on a final missing detail in Friston's (supposed) causal graph realism: if reality is indeed divisible into discrete individual units, what are those units? We've seen that, for Friston, a Markov blanket is at once a mathematical object and something that exists in the physical world. As such, each node in the network that composes it must itself be the state of some particular physical thing, and the graph will be constituted by statistical regularities holding between the changing states of these particulars. This is another reason why his use of the tools of causal graphical modelling may look strange to those familiar with the causal discovery work of Pearl, or of Spirtes, Glymour and Schienens, where causal-statistical relationships are modelled as holding not between token objects, but between general events, properties, or states of affairs (Hausman, 2005). Here, each node in the graph is not a concrete particular but a type that may be instantiated by many different particulars.

Friston, in contrast, takes nodes of the graph to be the token states of concrete parts of the physical world. Unlike Lewis and Field, however, he does not commit himself to the ultimate causal graph being microphysical.

In Friston (2013) the relevant units are indeed the electrochemical and kinetic states of particles in a primordial soup, but in Hipólito et al. (2021) they are the states of synapses and ion channels, and in the proposed models of societies and ecosystems, each node could be the state of an individual organism (Kirchoff et al., 2018).

The idea here, as Ramstead et al. (2019) explain, is that what appears as the single node of one Markov blanketed system is itself a Markov-blanketed complex that may be further decomposed into sensory, active and internal components. We might call this **multiscale probabilistic graph realism**: the position that at whatever scale we consider a system, it will be separable into distinct components whose interactions respect the Markov condition. Which scale we pick, Ramstead et al. argue will be relative to our interests, but the divisibility into Markov blanketed subsystems is not. Once all the components at that scale are included (even if this is not actually possible in some particular simple model) then there will be a set of *real* Markov blankets for that scale. As Sims (2020) suggests, such a multiscale formalism may actually be better placed to capture biological and physiological individuality, in which individual organisms – a somatic cell or a symbiotic organism – cooperate with others to form a higher-level individual, while still preserving their own individuality at a lower level of analysis.

While Friston (2019a) suggests this may apply ‘ad infinitum’, that causal graph structure may persist ‘all the way down’ (P.7), this commitment is not integral to the FEF, and later in the monograph, he expresses ‘metaphysical agnosticism’ on the issue (P.124). We thus have the option to take Papineau’s (1992) approach and reject the requirement for macroscopic phenomena to be reducible to microscopic interactions, sidestepping the spanner that quantum mechanics throws into the workings of a microphysical causal

graph, while still maintaining the reality of individual units, and the Markov blankets they compose, as an emergent macrophysical phenomenon.

A.4. A second stability requirement

Our initial problem with Markov blankets was that they might be drawn anywhere, depending on how we construct our causal model. Once we move to the all-encompassing graph, the constitution of Markov blankets may not be arbitrary, but they are nonetheless pervasive. If we take causality to be a local phenomenon, as presumed in Friston (2013), such that only contiguous elements can interact directly, and if we treat causality as reducible to statistical relationships between these elements, then for *any* element that we pick, the state of its immediate surroundings will render it conditionally independent of everything else in the world. If every organism is surrounded by a real Markov blanket this is not the result of its own efforts, but merely a general result of a particular account of the metaphysics of causality.

Friston's (2013, 2019) does not seem to take the existence of a Markov blanket to be quite as trivial as I have outlined. A candle flame cannot possess a Markov blanket, he argues, because in contrast to a cell membrane "any pattern of molecular interactions is destroyed almost instantaneously by the flux of gas molecules from its surface" (2013, P.2)

In other words, it [The Markov Blanket formalism] does not easily accommodate the fact that the particles that constitute a Markov blanket can, over time, wander away or, indeed, be exchanged or renewed. The canonical example here would be the blanket states of a candle flame, whose constituent particles (i.e., molecules of gas) are in constant flux." (2019, P.50)

This, Friston (2013) claims, is contrasted with “the physical configuration and dynamical states that constitute the Markov blanket of an organism—or organelle—change slowly in relation to the external and internal states it separates” (P.10). So what he seems to require here is that the connections between various parts of a system are stable throughout the changes in states of those parts, such that the particular set of elements that make up the Markov blanket is conserved over time. Indeed, if the way we make sense of objective probabilities is in terms of something like long-run frequencies or propensities, then the connections between units *must* be more stable than the changes in the state of those units in order for there to be a statistical relationship between them, and so for there to be a graph divisible into Markov blankets, at all.

So the systems of interest to the FEF are not supposed to be distinguished by the property of having Markov blankets per se, but rather by having the same particular Markov blanket with constituents that are stable⁴⁵ over some duration. This helps explain the conflation Friston’s critics accuse him of, between describing the Markov blanket as a ‘physical boundary’ and ‘statistical partition’ (Bruineberg et al, 2022). On the putative metaphysics I’m ascribing to Friston, a Markov blanket ends up being both. If reality literally is a statistical graph, and if causal interactions are exclusively local, then a stable Markov blanket is both a statistical partition and one realized by a fixed set of discrete elements that surround the system of interest – in other words, a boundary.

We can thus think of the requirements for a system being a free energy minimizer in terms of two types of stability. Firstly, the stability of the typical

⁴⁵ This requirement of stability is worryingly vague to be playing such a key role in determining whether a system ‘exists’ or not – but there are only so many rabbits a graduate student can be expected to chase.

state of its parts, as discussed in the previous section on survival as surprisal-minimization, and secondly, the stability of statistical-causal dependencies between these parts that allows us to identify a particular Markov blanket as persisting over time.

These then, are the steps required to claim that Markov blankets are real things. All that positing a real Markov blanket involves is claiming that a system decomposes into independent units (accomplished by some means prior to the FEF, by methods the framework itself does not specify) and arguing that if we accept objective statistical dependencies as an adequate reduction of causal relationships, along with the principle of locality, then the immediate surroundings of any one of those units will literally have the Markov property of inducing a conditional independence between what is inside this boundary and what is outside of it. To claim Markov blankets are real things is to make a general claim about the structure of the causal universe. In doing so the FEF has neither identified a new and interesting entity in the world, nor discovered a principled way of carving the world into things. This latter task depends on the, necessarily prior task, of telling us what the absolute units of the ultimate graph are.

References

Aguado-Velasco, C., & Bretscher, M. S. (1999). Circulation of the plasma membrane in Dictyostelium. *Molecular biology of the cell*, 10(12), 4419-4427.

Allen, M., & Friston, K. J. (2018). From cognitivism to autopoiesis: towards a computational framework for the embodied mind. *Synthese*, 195(6), 2459-2482.

Allen, M., & Friston, K. J. (2018). From cognitivism to autopoiesis: towards a computational framework for the embodied mind. *Synthese*, 195(6), 2459-2482. <http://doi.org/10.1007/s11229-016-1288-5>

Anderson, M. L. (2017). Of Bayes and Bullets: An Embodied, Situated, Targeting-Based Account of Predictive Processing. In T. Metzinger & W. Wiese (Eds.). *Philosophy and Predictive Processing: 4*. Frankfurt am Main: MIND Group. doi: [10.15502/9783958573055](https://doi.org/10.15502/9783958573055)

Arntzenius, F. (1992). The common cause principle. In *PSA: Proceedings of the Biennial meeting of the philosophy of science association (Vol. 1992, No. 2, pp. 227-237)*.

Ashby, W. R. (1940). Adaptiveness and Equilibrium. *Journal of Mental Science*, 86, pp. 478-484

Ashby, W. R. (1948) 'Design for a Brain', *Electronic Engineering*, 20, pp. 379-383.

Ashby, W. R. (1949). The electronic brain. *Radio Electronics*.

Ashby, W. R. (1956). *An Introduction to cybernetics*. London: Chapman & Hall.

Ashby, W. R. (1962). "Principles of the self-organizing system," in *Principles of Self-Organization: Transactions of the University of Illinois Symposium*, H. Von Foerster and G. W. Zopf, Jr. (eds.), Pergamon Press: London, UK, pp. 255-278.

Atkins, P. W. (1984). *The second law*. New York: Freeman.

Attneave, F. (1954). Some informational aspects of visual perception. *Psychological review*, 61(3), 183.

Austin, C. J. (2020). Organisms, activity, and being: on the substance of process ontology. *European Journal for Philosophy of Science*, 10(2), 1-21.

Badcock, P. B., Friston, K. J., Ramstead, M. J., Ploeger, A., & Hohwy, J. (2019). The hierarchically mechanistic mind: an evolutionary systems theory of the human brain, cognition, and behavior. *Cognitive, Affective, & Behavioral Neuroscience*, 19(6), 1319-1351.

Baltieri, M., Buckley, C. L., & Bruineberg, J. (2020). Predictions in the eye of the beholder: an active inference account of Watt governors. arXiv preprint arXiv:2006.11495.

Barandiaran, X. E. (2017). Autonomy and enactivism: towards a theory of sensorimotor autonomous agency. *Topoi*, 36(3), 409-430.

Barlow, H. B. (1961). Possible principles underlying the transformation of sensory messages. *Sensory communication*, 1(01).

Bassam Z. Shakhshiri, *Chemical Demonstrations: A Handbook for Teachers of Chemistry, Volume 1*. Madison: The University of Wisconsin Press, 1983, p. 79-80.

Batterman, R. W., & Rice, C. C. (2014). Minimal model explanations. *Philosophy of Science*, 81(3), 349-376.

Bechtel, William and Adele Abrahamsen (2005), Explanation: A Mechanist Alternative, *Studies in History and Philosophy of Biological and Biomedical Sciences*, 36:421-441.

- Belshaw, R., Pereira, V., Katzourakis, A., Talbot, G., Pačes, J., Burt, A., & Tristem, M. (2004). Long-term reinfection of the human genome by endogenous retroviruses. *Proceedings of the National Academy of Sciences*, *101*(14), 4894-4899.
- Beni, M. D. (2021). A critical analysis of Markovian monism. *Synthese*, 1-21.
- Berghofer, P. (2018). Why Husserl is a moderate foundationalist. *Husserl Studies*, *34*(1), 1-23.
- Bich, L., & Arnellos, A. (2012). Autopoiesis, autonomy, and organizational biology: Critical remarks on 'Life after Ashby'. *Cybernetics & Human Knowing*, *19*(4), 75-103.
- Bickhard, M. H. (2000). Autonomy, function, and representation. *Communication and Cognition—Artificial Intelligence*, *17*(3-4), 111-131.
- Bickhard, M. H. (2009). The biological foundations of cognitive science. *New Ideas in Psychology*, *27*, 75-84.
- Birkhoff GD (1931) Proof of the ergodic theorem. *Proceedings of the National Academy of Science* *17*(12):656–660
- Boden, M. A. (2000). Autopoiesis and life. *Cognitive Science Quarterly*, *1*(1), 115-143.
- Borges, J. L. (1998). On the exactitude of science. *Collected Fictions*. Translated by Andrew Hurley. *New York: Penguin*, 325.
- Bouizegarene, N., Ramstead, M. J., Constant, A., Friston, K. J., & Kirmayer, L. J. (2020). Narrative as active inference: an integrative account of the functions of narratives. Preprint, 10.
- Bourget, D., & Chalmers, D. J. (2021). Philosophers on philosophy: The 2020 PhilPapers survey.
- Brandom, R. (1979). Freedom and constraint by norms. *American Philosophical Quarterly*, *16*(3), 187–196.

- Brandom, R. (1994). *Making it explicit: Reasoning, representing, and discursive commitment*. Cambridge: Harvard University Press.
- Brentano, F. (1995) *Psychology from an Empirical Standpoint*, London: Routledge. Original work originally published 1874
- Brown, H., Friston, K. & Bestmann, S. (2011) Active inference, attention and motor preparation. *Frontiers in Psychology* 2:218. doi: 10.3389/fpsyg.2011.00218.
- Brown, H., Friston, K. J., & Bestmann, S. (2011). Active inference, attention, and motor preparation. *Frontiers in Psychology*, 2, 218.
- Bruineberg, J., Kiverstein, J., & Rietveld, E. (2018). The anticipating brain is not a scientist: the free-energy principle from an ecological-enactive perspective. *Synthese*, 195(6), 2417-2444.
- Buckley, C. L., Kim, C. S., McGregor, S., and Seth, A. K. (2017). The free energy principle for action and perception: A mathematical review. *Journal of Mathematical Psychology*, 81, 55-79.
- Burrows, M., Rogers, S. M., & Ott, S. R. (2011). Epigenetic remodelling of brain, body and behaviour during phase change in locusts. *Neural Systems & Circuits*, 1(1), 1-9.
- Cao, Y., Mu, H., Guo, J., Liu, H., Zhang, R., Liu, W., ... & Liu, H. (2020). Metabolic engineering of *Escherichia coli* for the utilization of ethanol. *Journal of Biological Research-Thessaloniki*, 27(1), 1-10.
- Capra, F., & Luisi, P. L. (2014). *The systems view of life: A unifying vision*. Cambridge University Press.
- Cartwright, N. (2001). What is wrong with Bayes nets?. *The monist*, 84(2), 242-264.
- Cartwright, N., & Jones, M. (1991). How to hunt quantum causes. *Erkenntnis*, 35(1), 205-231.

Chalmers, D. J. (1995). 'On implementing a computation', *Minds and Machines*, 4: 391–402.

Chapman, S. (1968). Catching a baseball. *American Journal of Physics*, 36, 868–870.

Chemero, A., & Silberstein, M. (2008). After the philosophy of mind: Replacing scholasticism with science. *Philosophy of science*, 75(1), 1-27.

Chirimuuta, M. (2020). The Reflex Machine and the Cybernetic Brain: The Critique of Abstraction and its Application to Computationalism. *Perspectives on Science*, 28(3), 421-457.

Christensen, W. D. & Hooker, C. A. (2000). Autonomy and the emergence of intelligence: Organised interactive construction. *Communication and Cognition – Artificial Intelligence*, 17 (3-4), 133-157.

Christensen, W. D., & Bickhard, M. H. (2002). The process dynamics of normative function. *The Monist*, 85(1), 3-28.

Chunharas, C., & Ramachandran, V. S. (2016). OUT OF THE SHADOWS. *Scientific American Mind*, 27(4), 56-61.

Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and brain sciences*, 36(3), 181-204.

Clark, A. (2015). Radical predictive processing. *The Southern Journal of Philosophy*, 53, 3-27.

Clark, A. (2015). *Surfing uncertainty: Prediction, action, and the embodied mind*. Oxford University Press.

Clark, A. (2017). How to knit your own Markov blanket: Resisting the second law with metamorphic minds. In T. Metzinger & W. Wiese (Eds.), *Philosophy and Predictive Processing: 3*. Frankfurt am Main: MIND Group. doi:10.15502/9783958573031.

Clark, A., & Toribio, J. (1994). Doing without representing?. *Synthese*, 101(3), 401-431.

Clegg, J. S. (2001). Cryptobiosis—a peculiar state of biological organization. *Comparative Biochemistry and Physiology Part B: Biochemistry and Molecular Biology*, 128(4), 613-624.

Collier, J. (2004). Self-Organisation, individuation and identity. *Revue Internationale de Philosophie*, 59, 151-172. Collier, J. (2008). A dynamical account of emergence. *Cybernetics & Human Knowing*, 15 (3-4), 75-86.

Colombetti, G. (2010). Enaction, sense-making and emotion. *Enaction: Toward a new paradigm for cognitive science*, 145-164.

Colombetti, G. (2014). *The feeling body: Affective science meets the enactive mind*. MIT press.

Colombo, M., & Palacios, P. (2021). Non-equilibrium thermodynamics and the free energy principle in biology. *Biology & Philosophy*, 36(5), 1-26.

Constant, A., Ramstead, M. J., Veissière, S. P., & Friston, K. (2019). Regimes of expectations: An active inference model of social conformity and human decision making. *Frontiers in psychology*, 10, 679.

Corcoran, A. W., Pezzulo, G., & Hohwy, J. (2020). From allostatic agents to counterfactual cognisers: active inference, biological regulation, and the origins of cognition. *Biology & Philosophy*, 35(3), 1-45.

Costanzo, J. P., do Amaral, M. C. F., Rosendale, A. J., & Lee Jr, R. E. (2013). Hibernation physiology, freezing adaptation and extreme freeze tolerance in a northern population of the wood frog. *Journal of Experimental Biology*, 216(18), 3461-3473.

Crane, T. (2015). *The mechanical mind: A philosophical introduction to minds, machines and mental representation*. Routledge.

Da Costa, L., Friston, K., Heins, C., & Pavliotis, G. A. (2021). Bayesian mechanics for stationary processes. *Proceedings of the Royal Society A*, 477(2256), 20210518.

Da Costa, L., Parr, T., Sengupta, B., & Friston, K. (2021). Neural dynamics under active inference: Plausibility and efficiency of information processing. *Entropy*, 23(4), 454.

David Henderson, "Norms, Normative Principles, and Explanation," *Philosophy of the Social Sciences* 32 (2002): 329-64.

Davidson, Donald (1963) Actions, Reasons and Causes. *The Journal of Philosophy*. Vol. 60, No. 23, pp. 685-700 (16 pages)

Dawkins, R. (1996). *The blind watchmaker: Why the evidence of evolution reveals a universe without design*. WW Norton & Company.

Dawkins, R. (2008). *River out of Eden: A Darwinian view of life*. Basic books.

Dayan, P., Hinton, G. E., Neal, R. M., & Zemel, R. S. (1995). The helmholtz machine. *Neural computation*, 7(5), 889-904.

De Jaegher, H., & Di Paolo, E. (2007). Participatory sense-making. *Phenomenology and the cognitive sciences*, 6(4), 485-507.

Di Paolo, E. A. (2010). Overcoming autopoiesis: An enactive detour on the way from life to society. In *Advanced series in management*. Emerald Group Publishing Limited.

Di Paolo, E. A., Cuffari, E. C., & De Jaegher, H. (2018). *Linguistic bodies: The continuity between life and language*. MIT press.

- Di Paolo, E., & Thompson, E. (2014). The enactive approach. In *The Routledge handbook of embodied cognition* (pp. 86-96). Routledge.
- Di Paolo, E., Thompson, E., & Beer, R. (2022). Laying down a forking path: Tensions between enaction and the free energy principle. *Philosophy and the Mind Sciences*, 3.
- Di Paolo, E., Buhrmann, T., & Barandiaran, X. (2017). *Sensorimotor life: An enactive proposal*. Oxford University Press.
- DiFrisco, J. (2014). Hylomorphism and the metabolic closure conception of life. *Acta Biotheoretica*, 62(4), 499-525.
- Downey, A. (2018). Predictive processing and the representation wars: A victory for the eliminativist (via fictionalism). *Synthese*, 195(12), 5115-5139.
- Dretske, F., & Bogdan, R. (1993). Misrepresentation. *Readings in philosophy and cognitive science*, 297-314.
- Dreyfus, H. L. (2002). Intelligence without representation—Merleau-Ponty's critique of mental representation The relevance of phenomenology to scientific explanation. *Phenomenology and the cognitive sciences*, 1(4), 367-383.
- Dupré, J., & Nicholson, D. (2018). A manifesto for a Processual philosophy of biology. In D. Nicholson & J. Dupré (Eds.), *Everything flows: Towards a Processual philosophy of biology* (pp. 3–45). Oxford: Oxford University Press.
- Dupuy, J. P. (2018). Cybernetics Is an Antihumanism. Technoscience and the Rebellion Against the Human Condition. In *French Philosophy of Technology* (pp. 139-156). Springer, Cham.
- Earman J, Rédei M (1996) Why ergodic theory does not explain the success of equilibrium statistical mechanics. *Brit J Philos Sci* 47(1):63–78
- Egan, D., Reynolds, S., & Wendland, A. J. (Eds.). (2013). *Wittgenstein and Heidegger*. New York, NY: Routledge.
- Egan, F. (2014). How to think about mental content. *Philosophical Studies*, 170(1), 115–135. doi:10.1007/s11098-013-0172-0.

Enoch David 2011: *Taking Morality Seriously: A Defense of Robust Realism*.
Oxford: Oxford University Press.

Evans, G. (1982). *The Varieties of Reference*. Oxford. Clarendon Press

Field, H. (2003). Causation in a physical world. *Oxford handbook of metaphysics*, 435-60.

Fiorillo, C. D. (2010). A neurocentric approach to Bayesian inference. *Nature Reviews Neuroscience*, 11(8), 605-605.

Firth, N., Jensen, S. O., Kwong, S. M., Skurray, R. A., & Ramsay, J. P. (2018). Staphylococcal plasmids, transposable and integrative elements. *Microbiology spectrum*, 6(6), 6-6.

FitzGerald, T. H., Dolan, R. J., & Friston, K. J. (2014). Model averaging, optimal inference, and habit formation. *Frontiers in human neuroscience*, 8, 457.

Fleischaker, G. R. (1988). Autopoiesis: the status of its system logic. *BioSystems*, 22 (1), 37-49.

Fleischaker, G. R. (1988). Autopoiesis: the status of its system logic. *BioSystems*, 22(1), 37-49.

Fleischaker, G. R. (1990). Origins of life: an operational definition. *Origins of Life and Evolution of the Biosphere*, 20(2), 127-137.

Foerster H, 2003, *Understanding Understanding: Essays on Cybernetics and Cognition*, Springer, New York, Berlin, Heidelberg.

Friston, K. (2002). Functional integration and inference in the brain. *Progress in neurobiology*, 68(2), 113-143.

Friston, K. (2003). Learning and inference in the brain. *Neural Networks*, 16(9), 1325-1352.

Friston, K. (2005). A theory of cortical responses. *Philosophical transactions of the Royal Society B: Biological sciences*, 360(1456), 815-836.

Friston, K. (2009). The free-energy principle: a rough guide to the brain?. *Trends in cognitive sciences*, 13(7), 293-301.

Friston, K. (2010). The free-energy principle: a unified brain theory?. *Nature reviews neuroscience*, 11(2), 127-138.

Friston, K. (2013). Life as we know it. *Journal of the Royal Society Interface*. http://doi.org/10.1098/rsif.2013.0475&domain=pdf&date_stamp=2013-07-03

Friston, K. (2018). Am I self-conscious?(Or does self-organization entail self-consciousness?). *Frontiers in psychology*, 9, 579.

Friston, K. (2019). A free energy principle for a particular physics. *arXiv preprint arXiv:1906.10184*.

Friston, K. J. (2000). The labile brain. II. Transients, complexity and selection. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 355(1394), 237-252.

Friston, K. J. (2019). Beyond the desert landscape. *Andy Clark and his critics*, 174-190.

Friston, K. J., Fagerholm, E. D., Zarghami, T. S., Parr, T., Hipólito, I., Magrou, L., & Razi, A. (2021). Parcels and particles: Markov blankets in the brain. *Network Neuroscience*, 5(1), 211-251.

Friston, K. J., Wiese, W., & Hobson, J. A. (2020). Sentience and the origins of consciousness: From Cartesian duality to Markovian monism. *Entropy*, 22(5), 516.

Friston, K., & Ao, P. (2012). Free energy, value, and attractors. *Computational and mathematical methods in medicine*, 2012.

Friston, K., & Mathys, C. (2016). I am therefore I think. *The Unconscious: A bridge between psychoanalysis and cognitive neuroscience*, 113.

Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., and Pezzulo, G. (2017). Active inference: A process theory. *Neural Computation*, 29, 1-49.

Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., & Pezzulo, G. (2016). Active inference and learning. *Neuroscience & Biobehavioral Reviews*, 68, 862-879.

Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P., & Pezzulo, G. (2017). Active inference: a process theory. *Neural computation*, 29(1), 1-49.

Friston, K., Kilner, J., & Harrison, L. (2006). A free energy principle for the brain. *Journal of physiology-Paris*, 100(1-3), 70-87.

Friston, K., Thornton, C., & Clark, A. (2012). Free-energy minimization and the dark-room problem. *Frontiers in psychology*, 3, 130.

Friston, K., Thornton, C., & Clark, A. (2012). Free-Energy Minimization and the Dark-Room Problem. *Frontiers in Psychology*, 3.
<http://doi.org/10.3389/fpsyg.2012.00130>

Friston, K. (2013). Life as we know it. *Journal of the Royal Society Interface*, 10(86), 20130475.

Froese, T. (2010). Life after Ashby: ultrastability and the autopoietic foundations of biological autonomy. *Cybernetics & Human Knowing*, 17(4), 7-49.

Gallagher, S. (2012). On the possibility of naturalizing phenomenology. In D. Zahavi. *Oxford Handbook of Contemporary Phenomenology* (70-93). Oxford: Oxford University Press.

Gallagher, S. (2017). *Enactivist Interventions: Rethinking the Mind*. Oxford: Oxford University Press.

Gallagher, S. (2018 in press). Rethinking nature: Phenomenology and a non-reductionist cognitive science. *Australasian Philosophical Review*.

- Gallavotti G (1999) *Statistical mechanics: a short treatise*. Springer, Berlin
- Gangopadhyay, N., & Kiverstein, J. (2009). Enactivism and the unity of perception and action. *Topoi*, 28(1), 63-73.
- Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural computation*, 4(1), 1-58.
- Gershman, S.J. (2019). What does the free energy principle tell us about the brain? *Neurons, Behavior, Data Analysis, and Theory*.
- Gibson, J. J. (1979). *The ecological approach to visual perception*. Boston. Houghton-Mifflin
- Gładziejewski, P. (2016). Predictive coding and representationalism. *Synthese*, 193(2), 559-582.
- Glennan, S. (2002). Rethinking mechanistic explanation. *Philosophy of Science*, 69(1), S342–S353
- Glennan, Stuart (2002) Rethinking Mechanistic Explanation, *Philosophy of Science*, 69: S342-S353.
- Glymour, C. (1999). Rabbit hunting. *Synthese*, 55-78.
- Glymour, C. (2006). Markov properties and quantum experiments. In *Physical theory and its interpretation* (pp. 117-126). Springer, Dordrecht.
- Glymour, C. (2010). What is right with ‘Bayes net methods’ and what is wrong with ‘hunting causes and using them’?. *The British Journal for the Philosophy of Science*, 61(1), 161-211.
- Godfrey-Smith, P. (2016). Mind, matter, and metabolism. *The Journal of Philosophy*, 113(10), 481-506.
- Good, I. J. (1959). A theory of causality. *British Journal for the Philosophy of Science*, 9:307–310.

Gopnik, A. (2012). Scientific thinking in young children: Theoretical advances, empirical research, and policy implications. *Science*, 337(6102), 1623-1627.

Gregory, R. L. (1980) Perceptions as hypotheses. *Philosophical Transactions of the Royal Society of London B* 290(1038):181–97. [aAC, KF]

Grush, R. (2004). The emulation theory of representation: Motor control, imagery, and perception. *Behavioral and brain sciences*, 27(3), 377-396.

Grush, R. (2007). Skill theory v2. 0: Dispositions, emulation, and spatial perception. *Synthese*, 159(3), 389-416.

Gualeni, S. (v 1.1, 2017). *Something Something Soup Something*. Digital game developed by I. Kniestedt, R. Fassone, and J. Schellekens. <http://soup.gua-leri.com>

Guay A, Pradeu T (2016) Individuals across the sciences. Oxford University Press, New York

Haldane, J. S. (1917). *Organism and Environment, as Illustrated by the Physiology of Breathing*. New Haven: Yale University Press.

Hammer, T. J., Sanders, J. G., & Fierer, N. (2019). Not all animals need a microbiome. *FEMS microbiology letters*, 366(10), fnz117.

Harvey, I. (2013). Standing on the broad shoulders of Ashby. *Constructivist Foundations*, 9(1), 102-104.

Harvey, K. (2018). An open-ended approach to Piagetian development of adaptive behavior. *Open Access Library Journal*, 5(03), 1.

Hatfield, G. (1984, January). Spatial Perception and Geometry in Kant and Helmholtz. In *PSA: Proceedings of the Biennial Meeting of the Philosophy of Science Association* (Vol. 1984, No. 2, pp. 569-587). Philosophy of Science Association.

- Hatfield, G. C. (1990). *The natural and the normative: Theories of spatial perception from Kant to Helmholtz*. Cambridge. MIT Press.
- Hatfield, G. C. (1990). *The natural and the normative: Theories of spatial perception from Kant to Helmholtz*. MIT Press.
- Haugeland, J. (1991). Representational genera. In W. Ramsey, S. Stich, & J. Garron (Eds.), *Philosophy and connectionist theory* (pp. 61–89). Hillsdale, NJ: Lawrence Erlbaum.
- Hausman, D. M. (2005). Causal relata: Tokens, types, or variables?. *Erkenntnis*, 63(1), 33-54.
- Hediger, M. A., Johnson, D. F., Nierlich, D. P., & Zabin, I. (1985). DNA sequence of the lactose operon: the lacA gene and the transcriptional termination region. *Proceedings of the National Academy of Sciences*, 82(19), 6414-6418.
- Helmholtz, H. (1962). “Concerning the perceptions in general,” in *Treatise on Physiological Optics*, 3rd Edn, Vol. III, ed. J. Southall, trans. (New York: Dover). Original work published 1866
- Henderson, D. (2010). Explanation and rationality naturalized. *Philosophy of the Social Sciences*, 40(1), 30-58.
- Henson, R. N., & Gagnepain, P. (2010). Predictive, interactive multiple memory systems. *Hippocampus*, 20(11), 1315-1326.
- Hesp, C., Ramstead, M., Constant, A., Badcock, P., Kirchhoff, M., & Friston, K. (2019). A multi-scale view of the emergent complexity of life: A free-energy proposal. In *Evolution, development and complexity* (pp. 195-227). Springer, Cham.
- Hesp, C., Smith, R., Parr, T., Allen, M., Friston, K. J., and Ramstead, M. J. D. (2021). Deeply felt affect: The emergence of valence in deep active inference. *Neural Computation*, 33(2), 398–446.

Hesslow, G. (1976). Two notes on the probabilistic approach to causality. *Philosophy of science*, 43(2), 290-292.

Hinton, G. E., & Van Camp, D. (1993, August). Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory* (pp. 5-13).

Hipólito, I., Ramstead, M. J., Convertino, L., Bhat, A., Friston, K., & Parr, T. (2021). Markov blankets in the brain. *Neuroscience & Biobehavioral Reviews*.

Hobson, A. (2013). There are no particles, there are only fields. *American journal of physics*, 81(3), 211-223.

Hoffmann, P. M. (2012). *Life's ratchet: how molecular machines extract order from chaos*. New York. Basic Books.

Hohwy, J. (2013). *The predictive mind*. OUP Oxford.

Hohwy, J. (2016). The self-evidencing brain. *Noûs*, 50(2), 259-285.

Hohwy, J. (2017). How to entrain your evil demon. In T. Metzinger & W. Wiese (Eds.), *Philosophy and predictive processing: 2*. Frankfurt am Main: MIND Group. doi:10.15502/9783958573048.

Hohwy, J. (2021). Self-supervision, normativity and the free energy principle. *Synthese*, 199(1), 29-53.

Hurley, S. (1998) *Consciousness in action*. Harvard University Press

Hurley, S. (2003). Animal action in the space of reasons. *Mind & Language*, 18(3), 231-257.

Hurley, S. (2008). The shared circuits model (SCM): How control, mirroring, and simulation can enable imitation, deliberation, and mindreading. *Behavioral and brain sciences*, 31(1), 1-22.

Husbands, P., & Holland, O. (2012). Warren McCulloch and the British cyberneticians. *Interdisciplinary Science Reviews*, 37(3), 237-253.

Husserl, E. (1970). *The Crisis of European Sciences and Transcendental Phenomenology. An Introduction to Phenomenology*, (D. Carr, Trans.) Evanston, IL: Northwestern University Press. Original work published 1936

Husserl, E. (1982). Ideas pertaining to a pure phenomenology and to a phenomenological philosophy. First book: general introduction to a pure phenomenology (F. Kerstens, Trans.). Dordrecht: Kluwer. Original work published 1913.

Husserl, E. (1997). Thing and space: lectures of 1907 (R. Rojcewicz, Trans.). Dordrecht: Kluwer. Original work published in 1907.

Husserl, E. (1999). Cartesian meditations: An introduction to phenomenology. (D. Cairns, Trans.). Original work published 1929

Husserl, E. (2001). Analyses Concerning Passive and Active Synthesis. Lectures on Transcendental Logic. (A.J. Steinbock, Trans.). Dordrecht: Kluwer Academic Publishers. Original work published 1920

Hutto, D. D., & Myin, E. (2012). *Radicalizing enactivism: Basic minds without content*. MIT press.

Hutto, D. D., & Myin, E. (2017). *Evolving enactivism: Basic minds meet content*. MIT press.

Jacob, P. (2019) Intentionality, *The Stanford Encyclopedia of Philosophy* , Edward N. Zalta (ed.), URL = <https://plato.stanford.edu/archives/win2019/entries/intentionality/>.

James, W. (1890). The principles of psychology, Vols. I, II. Cambridge, MA: Harvard University Press.

Jebari, J. (2019). Empirical moral rationalism and the social constitution of normativity. *Philosophical Studies*, 176(9), 2429-2453.

- Jeffrey, R. C. (1969). Statistical explanation vs. statistical inference. In *Essays in honor of Carl G. Hempel* (pp. 104-113). Springer, Dordrecht.
- Jiang, R., Li, W., Lu, X. X., Xie, J., Zhao, Y., & Li, F. (2021). Assessment of temperature extremes and climate change impacts in Singapore, 1982–2018. *Singapore Journal of Tropical Geography*, 42(3), 378-396.
- Jonas, H. (1953). A critique of cybernetics. *Social Research*, 172-192.
- Jonas, H. (1966). *The phenomenon of life: Toward a philosophical biology*. Northwestern University Press.
- Juarrero, A. (2000). Dynamics in action: Intentional behavior as a complex system. *Emergence*, 2(2), 24-57.
- Kanai, R., Komura, Y., Shipp, S., & Friston, K. (2015). Cerebral hierarchies: predictive processing, precision and the pulvinar. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1668), 20140169.
- Kant, I. (2008). *Groundwork for the Metaphysics of Morals*. Yale University Press. Original work published 1785
- Kaplan, D. M. (2017). Mechanisms and dynamical systems. *The Routledge handbook of mechanisms and mechanical philosophy*, 267-280.
- Kaplan, D. M., & Bechtel, W. (2011). Dynamical models: An alternative or complement to mechanistic explanations?. *Topics in Cognitive Science*, 3(2), 438-444.
- Kaplan, D. M., & Craver, C. F. (2011). The explanatory force of dynamical and mathematical models in neuroscience: A mechanistic perspective. *Philosophy of science*, 78(4), 601-627.
- Kauffman, S. A. (1986). Autocatalytic sets of proteins. *Journal of theoretical biology*, 119(1), 1-24.
- Kauffman, S. A. (1993). *The origins of order: Self-organization and selection in evolution*. Oxford University Press, USA.

Kauffman, S. A. (2000). *Investigations*. Oxford University Press.

Kauffman, S. A. (2019). *A world beyond physics: the emergence and evolution of life*. Oxford University Press

Keilin, D. (1959). The Leeuwenhoek Lecture-The problem of anabiosis or latent life: history and current concept. *Proceedings of the Royal Society of London. Series B-Biological Sciences*, 150(939), 149-191.

Keller, E. F. (2008). Organisms, machines, and thunderstorms: A history of self-organization, part one. *Historical Studies in the Natural Sciences*, 38(1), 45-75.

Keller, E. F. (2009). Organisms, machines, and thunderstorms: a history of self-organization, part two. Complexity, emergence, and stable attractors. *Historical Studies in the Natural Sciences*, 39(1), 1-31.

Keller, G. B., & Mrcic-Flogel, T. D. (2018). Predictive processing: a canonical cortical computation. *Neuron*, 100(2), 424-435.

Kenny, Anthony. *The Metaphysics of Mind*. Clarendon Press, 1989.

Kiefer, A., & Hohwy, J. (2018). Content and misrepresentation in hierarchical generative models. *Synthese*, 195(6), 2387-2415.

Kiefer, A., & Hohwy, J. (2018). Content and misrepresentation in hierarchical generative models. *Synthese*, 195(6), 2387-2415.
<https://doi.org/10.1007/s11229-017-1435-7>

Kiefer, A., & Hohwy, J. (2019). Representation in the prediction error minimization framework. In *The Routledge companion to philosophy of psychology* (pp. 384-409). Routledge.

Kirchhoff, M. D., & Kiverstein, J. (2021). How to determine the boundaries of the mind: A Markov blanket proposal. *Synthese*, 198(5), 4791-4810.

- Kirchhoff, M. D., & Robertson, I. (2018). Enactivism and predictive processing: A non-representational view. *Philosophical Explorations*, 21(2), 264–281.
- Kirchhoff, M., Parr, T., Palacios, E., Friston, K., & Kiverstein, J. (2018). The Markov blankets of life: autonomy, active inference and the free energy principle. *Journal of the Royal Society Interface*, 15(138), 20170792. <http://doi.org/10.1098/rsif.2017.0792>
- Korbak, T. (2021). Computational enactivism under the free energy principle. *Synthese* 198, 2743–2763 <https://doi.org/10.1007/s11229-019-02243-4>.
- Korsgaard, C. M. (1996). *The sources of normativity*. Cambridge: Cambridge University Press.
- Korsgaard, C. M. (2008). *The constitution of agency: Essays on practical reason and moral psychology*. OUP Oxford.
- Korsgaard, C. M. (2018a). *Fellow creatures: Our obligations to the other animals*. Oxford University Press.
- Korsgaard, C. M. (2018b). RATIONALITY 20. *Critical terms for animal studies*, 294.
- Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building machines that learn and think like people. *Behavioral and brain sciences*, 40.
- Langton, C. (1989). *Artificial life: proceedings of an interdisciplinary workshop on the synthesis and simulation of living systems*. Westview Press
- Langton, C. G. (1989). Artificial life. In C. G. Langton, (Ed.), *Artificial life I (proceedings of the First Conference on Artificial Life, Los Alamos, September, 1987)*
- Lee, T. S. & Mumford, D. (2003) Hierarchical Bayesian inference in the visual cortex. *Journal of Optical Society of America*, A 20(7):1434–48. [aAC]

- Letelier, J. C., Cárdenas, M. L., & Cornish-Bowden, A. (2011). From L'Homme Machine to metabolic closure: steps towards understanding life. *Journal of Theoretical Biology*, 286, 100-113.
- Lewin, K. (1922). *Der Begriff der Genese in Physik, Biologie und Entwicklungsgeschichte: eine Untersuchung zur vergleichenden Wissenschaftslehre*. Berlin: Springer.
- Lewis, D. (1994). Humean supervenience debugged. *Mind*, 103(412), 473-490.
- Longo, G., & Montévil, M. (2013). Extended criticality, phase spaces and enablement in biology. *Chaos, Solitons & Fractals*, 55, 64-79.
- Longo, G., & Montévil, M. (2014). From physics to biology by extending criticality and symmetry breakings. *Perspectives on Organisms*, 161-185.
- Longo, G., Montévil, M., & Kauffman, S. (2012, July). No entailing laws, but enablement in the evolution of the biosphere. In *Proceedings of the 14th annual conference companion on Genetic and evolutionary computation* (pp. 1379-1392).
- Lotze, R. H. (1852). *Medizinische Psychologie oder Physiologie der Seele* (pp. 287–325). Leipzig, Germany: WeidmannŌsche Buch-handlung.
- Lusthaus, D. (2014). *Buddhist Phenomenology: A philosophical investigation of Yogacara Buddhism and the Ch'eng Wei-shih Lun*. Routledge.
- Lyon, P. (2006). The biogenic approach to cognition. *Cognitive Processing*, 7(1), 11-29.
- Machamer, Peter, Lindley Darden, and Carl F. Craver. (2000). Thinking about Mechanisms, *Philosophy of Science*, 67: 1-25.
- Madary, M. (2012). How would the world look if it looked as if it were encoded as an intertwined set of probability density distributions?. *Frontiers in psychology*, 3, 419.
- Madary, M. (2016). *Visual phenomenology*. MIT Press.

Maddy, P. (1990). Physicalistic platonism. In *Physicalism in mathematics* (pp. 259-289). Springer, Dordrecht.

Maddy, P. (1997). *Naturalism in mathematics*. Clarendon Press.

Matthen, M. (2014). Debunking enactivism: a critical notice of Hutto and Myin's Radicalizing Enactivism. *Canadian Journal of Philosophy*, 44(1), 118-128.

Matthen, M., & Ariew, A. (2002). Two ways of thinking about fitness and natural selection. *The Journal of Philosophy*, 99(2), 55-83.

Maturana, H. R. (1970). *Biology of cognition* (pp. 4-9). Urbana: Biological Computer Laboratory, Department of Electrical Engineering, University of Illinois.

Maturana, H. R. (1975). The organization of the living: A theory of the living organization. *International journal of man-machine studies*, 7(3), 313-332.

Maturana, H. R. & Varela, F. J. (1980), *Autopoiesis and Cognition: The Realization of the Living* (pp. 59-140). Dordrecht, Holland: Kluwer Academic. Original work published 1972

Maynard, C., & Weinkove, D. (2020). Bacteria increase host micronutrient availability: mechanisms revealed by studies in *C. elegans*. *Genes & nutrition*, 15(1), 1-11.

McDowell, J. (1994). *Mind and world*. Cambridge: Harvard University Press.

Mehler, D. M. A., & Kording, K. P. (2018). The lure of causal statements: Rampant mis-inference of causality in estimated connectivity. *arXiv e-prints*, arXiv-1812.

Meincke, A. S. (2019). Autopoiesis, biological autonomy and the process view of life. *European Journal for Philosophy of Science*, 9(1), 1-16.

Menary, R., & Gillett, A. J. (2020). Are Markov Blankets Real and Does It Matter?. *The Philosophy and Science of Predictive Processing*, 39.

- Merleau-Ponty, M. (1963). *The Structure of Behavior*, (A. Fisher, Trans.) Pittsburgh, PA: Dusquene Lniversity Press. Original work published 1942
- Merleau-Ponty, M. (1964). *The primacy of perception: And other essays on phenomenological psychology, the philosophy of art, history, and politics*. Northwestern University Press.
- Merleau-Ponty, M. (2004). Eye and Mind. Pp. 291–234 in Maurice Merleau-Ponty: Basic Writing. Edited by Thomas Baldwin. London: Routledge. Original work published 1961
- Merleau-Ponty, M. (2012). *Phenomenology of perception*. (D. Landes, Trans.) New York: Routledge. Original work published 1945
- Mi, S., Lee, X., Li, X. P., Veldman, G. M., Finnerty, H., Racie, L., ... & McCoy, J. M. (2000). Syncytin is a captive retroviral envelope protein involved in human placental morphogenesis. *Nature*, 403(6771), 785-789.
- Mohanty, J. N. (1978). Husserl's Transcendental Phenomenology and Essentialism. *The Review of Metaphysics*, 299-321.
- Montévil, M., & Mossio, M. (2015). Biological organisation as closure of constraints. *Journal of theoretical biology*, 372, 179-191.
- Moran, D. (2002). *Introduction to phenomenology*. Routledge.
- Moreno, A., & Ruiz-Mirazo, K. (1999). Metabolism and the problem of its universalization. *BioSystems*, 49(1), 45-61.
- Moreno, A., Mossio, M. (2015) *Biological Autonomy. A Philosophical and Theoretical Enquiry*. Springer, Dordrecht.
- Morgan, M. S., & Morrison, M. (1999). *Models as mediators* (p. 347). Cambridge: Cambridge University Press.
- Morgan, W. (2021). Are Organisms Substances or Processes?. *Australasian Journal of Philosophy*, 1-15.

- Mossio, M., & Bich, L. (2017). What makes biological organisation teleological?. *Synthese*, 194(4), 1089-1114.
- Murata, R., Nozu, R., Mushirobira, Y., Amagai, T., Fushimi, J., Kobayashi, Y., ... & Nakamura, M. (2021). Testicular inducing steroidogenic cells trigger sex change in groupers. *Scientific reports*, 11(1), 1-7.
- Musmann, H. (1979). Predictive image coding. In W. K. Pratt (Ed.), *Image transmission techniques* 73-112 New York: Academic Press.
- Nagel, T. (1986). *The view from nowhere*. Oxford: Oxford University Press.
- Nagel, T. (2012). *Mind and cosmos: Why the materialist neo-Darwinian conception of nature is almost certainly false*. Oxford University Press.
- Nanay, B. (2014). Empirical problems with anti-representationalism. *Does perception have content*, 39-50.
- Nave, K., Deane, G., Miller, M., & Clark, A. (2020). Wilding the predictive brain. *Wiley Interdisciplinary Reviews: Cognitive Science*, 11(6), e1542.
- Neal, R. M., & Hinton, G. E. (1998). A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models* (pp. 355-368). Springer, Dordrecht.
- Netburn, D. (2014). In Alaska, wood frogs freeze for seven months, thaw and hop away. *The Los Angeles Times*. Retrieved from <https://www.latimes.com/science/sciencenow/la-sci-sn-alaskan-frozen-frogs-20140723-story.html>
- Neuman, Y. (2006). Cryptobiosis: a new theoretical perspective. *Progress in biophysics and molecular biology*, 92(2), 258-267.
- Nicholson, D. (2018). Reconceptualizing the organism: From complex machine to flowing stream. In D. Nicholson & J. Dupré (Eds.), *Everything flows: Towards a Processual philosophy of biology* (pp. 139–166). Oxford: Oxford University Press.

Nicholson, D. J. (2012). The concept of mechanism in biology. *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological and Biomedical Sciences*, 43(1), 152-163.

Orlandi, N. (2014). *The Innocent Eye: Why Vision is Not a Cognitive Process*. Oxford University Press

Orlandi, N. (2018). Predictive perceptual systems. *Synthese*, 195(6), 2367–2386. <https://doi.org/10.1007/s11229-017-1373-4>

Overgaard, S. (2006). The problem of other minds: Wittgenstein's phenomenological perspective. *phenomenology and the Cognitive Sciences*, 5(1), 53-73.

Palacios P (2018) Had we but world enough, and time... but we don't!: Justifying the thermodynamic and infinite-time limits in statistical mechanics. *Found Phys* 48(5):526–541

Palacios, E. R., Razi, A., Parr, T., Kirchhoff, M., & Friston, K. (2017). Biological Self-organisation and Markov blankets. *BioRxiv*, 227181.

Palacios, E. R., Razi, A., Parr, T., Kirchhoff, M., & Friston, K. (2020). On Markov blankets and hierarchical self-organisation. *Journal of theoretical biology*, 486, 110089.

Palmer, C. J., Seth, A. K., & Hohwy, J. (2015). The felt presence of other minds: Predictive processing, counterfactual predictions, and mentalising in autism. *Consciousness and Cognition*, 36, 376-389.

Papineau, D. (1992, January). Can we reduce causal direction to probabilities?. In *PSA: Proceedings of the biennial meeting of the philosophy of science association* (Vol. 1992, No. 2, pp. 238-252). Philosophy of Science Association.

Parfit Derek 2006: 'Normativity'. In Shafer-Landau Russ (ed.), *Oxford Studies in Metaethics*, 1, pp. 325–80. Oxford: Clarendon Press.

Parfit Derek 2011: *On What Matters*, vol. 2. Oxford: Oxford University Press.

Parr, T., & Friston, K. J. (2017). Working memory, attention, and salience in active inference. *Scientific reports*, 7(1), 1-21.

Parr, T., & Friston, K. J. (2019). Generalised free energy and active inference. *Biological cybernetics*, 113(5), 495-513.

Pattee, H. H. (1982). Cell Psychology: An Evolutionary Approach to the Symbol-Matter Problem. Reprinted in HH Pattee & J. Rączaszek-Leonardi. *Laws, Language and Life*, 165-179.

Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan kaufmann.

Pearl, J. (2000). *Causality: Models, reasoning, and inference*. Cambridge University Press.

Pearl, J. (2001). Bayesianism and causality, or, why I am only a half-Bayesian. In *Foundations of bayesianism* (pp. 19-36). Springer, Dordrecht.

Pedersen, B. H., Malte, H., Ramløv, H., & Finster, K. (2020). A method for studying the metabolic activity of individual tardigrades by measuring oxygen uptake using microrespirometry. *Journal of Experimental Biology*, 223(22), jeb233072.

Peterson, E. L. (2017). *The life organic: The theoretical biology club and the roots of epigenetics*. University of Pittsburgh Press.

Pezzulo G (2017) Tracing the roots of cognition in predictive processing. In: Metzinger T, Wiese W (eds) Philosophy and predictive processing. MIND Group , Frankfurt am Main

Pezzulo, G., Rigoli, F., & Friston, K. (2015). Active inference, homeostatic regulation and adaptive behavioural control. *Progress in neurobiology*, 134, 17-35.

Pezzulo, G., Rigoli, F., & Friston, K. J. (2018). Hierarchical active inference: a theory of motivated control. *Trends in cognitive sciences*, 22(4), 294-306.

Piaget, J. (1971). *Biology and knowledge: An essay on the relations between organic regulations and cognitive processes*. Edinburgh: Edinburgh University Press 1967a (Trans.) Original work published in 1967

Piccinini, G. (2020). Mechanisms, Multiple Realizability, and Medium Independence. *Neurocognitive Mechanisms: Explaining Biological Cognition*. Oxford University Press

Pickering, A. (2010). *The cybernetic brain*. University of Chicago Press.

Potochnik, A. (2017). *Idealization and the Aims of Science*. University of Chicago Press.

Powers, W. T. (1973). *Behavior: The control of perception*. Chicago: Aldine.

Price, H. (1992, January). The direction of causation: Ramsey's ultimate contingency. In *PSA: Proceedings of the biennial meeting of the philosophy of science association* (Vol. 1992, No. 2, pp. 253-267). Philosophy of Science Association.

Price, J. C. (1979). Assessment of the urban heat island effect through the use of satellite data. *Monthly Weather Review*, 107(11), 1554-1557.

Prigogine, I., & Stengers, I. (1984). *Order out of chaos*. New Science Library.

Prigogine, I., & Stengers, I. (1997). *The end of certainty*. Simon and Schuster.

Proksch, S (2021, October 15–17). *Acoustic social worlds; from Markov blankets to interpersonal synergies* [Conference presentation]. Cognitio2021, L'Université du Québec à Montréal, Canada, <https://www.youtube.com/watch?v=CiiPVa3yu6s>)

Putnam, H. (1975). 'The mental life of some machines'. *Mind, language and reality, philosophical papers, volume 2*, pp. 408–28. Cambridge University Press: Cambridge.

Raja, V., Valluri, D., Baggs, E., Chemero, A., & Aderson, M. L. (2021). The markov blanket trick: On the scope of the free energy principle and active inference.

Ramsey, W. M. (2007). *Representation reconsidered*. Cambridge University Press.

Ramstead, M. J. D., Badcock, P. B., & Friston, K. J. (2018). Answering Schrödinger's question: A free-energy formulation. *Physics of Life Reviews*, 24, 1–16. <http://doi.org/10.1016/j.plrev.2017.09.001>

Ramstead, M. J. D., Kirchhoff, M. D., Constant, A., & Friston, K. J. (2019). Multiscale integration: beyond internalism and externalism. *Synthese*, 24(6), 17. <http://doi.org/10.1007/s11229-019-02115-x>

Ramstead, M. J., Friston, K. J., & Hipólito, I. (2020). Is the free-energy principle a formal theory of semantics? From variational density dynamics to neural and phenotypic representations. *Entropy*, 22(8), 889.

Ramstead, M. J., Hesp, C., Tschantz, A., Smith, R., Constant, A., & Friston, K. (2021). Neural and phenotypic representation under the free-energy principle. *Neuroscience & Biobehavioral Reviews*, 120, 109-122.

Ramstead, M. J., Kirchhoff, M. D., & Friston, K. J. (2020). A tale of two densities: Active inference is enactive inference. *Adaptive Behavior*, 28(4), 225-239.

Ramstead, M. J., Kirchhoff, M. D., Constant, A., & Friston, K. J. (2019). Multiscale integration: beyond internalism and externalism. *Synthese*, 1-30.

Raposo de Magalhães, C., Schrama, D., Nakharuthai, C., Boonanuntanasarn, S., Revets, D., Planchon, S., & Rodrigues, P. M. (2021). Metabolic plasticity of gilthead seabream under different stressors: analysis of the stress responsive hepatic proteome and gene expression. *Frontiers in Marine Science*, 8, 469.

Reichenbach, H. (1956). *The direction of time* (Vol. 65). Univ of California Press.

Rescorla, M. (2015). Bayesian perceptual psychology. *The oxford handbook of philosophy of perception*, 10.

Roelofs, L. (2018). Why imagining requires content: A reply to a reply to an objection to radical enactive cognition. *Thought: A Journal of Philosophy*, 7(4), 246-254.

Rosas, F. E., Mediano, P. A., Biehl, M., Chandaria, S., & Polani, D. (2020, September). Causal blankets: Theory and algorithmic framework. In *International Workshop on Active Inference* (pp. 187-198). Springer, Cham.

Rosen, R. (1991). *Life itself: a comprehensive inquiry into the nature, origin, and fabrication of life*. Columbia University Press.

Rosen, R. (1999). *Essays on life itself*. Columbia University Press.

Ross, L. N. (2015). Dynamical models and explanation in neuroscience. *Philosophy of Science*, 82(1), 32-54.

Rousseau, J. J. (2018). *Rousseau: The Social Contract and other later political writings*. (V. Gourevitch, Trans.) Cambridge University Press. Original work published in 1762

Roy, J. M., Petitot, J., Pachoud, B., & Varela, F. J. (1999). Beyond the gap: An introduction to naturalizing phenomenology. In *Naturalizing phenomenology: Issues in contemporary phenomenology and cognitive science* (pp. 1-83). Stanford University Press.

Rubin, S., Parr, T., Da Costa, L., & Friston, K. (2020). Future climates: Markov blankets and active inference in the biosphere. *Journal of the Royal Society Interface*, 17(172), 20200503.

Ruiz-Mirazo, K., & Mavelli, F. (2007, September). Simulation model for functionalized vesicles: Lipid-peptide integration in minimal protocells. In *European Conference on Artificial Life* (pp. 32-41). Springer, Berlin, Heidelberg.

Ruiz-Mirazo, K., & Moreno, A. (2004). Basic autonomy as a fundamental step in the synthesis of life. *Artificial life*, 10(3), 235-259.

Ruiz-Mirazo, K., Peretó, J., & Moreno, A. (2004). A universal definition of life: autonomy and open-ended evolution. *Origins of Life and Evolution of the Biosphere*, 34(3), 323-346.

Salmon, W. (1980) Probabilistic Causality, in *Causality and Explanation*, 1998 Ed, Oxford University Press, Oxford, , pp. 208–232. Originally published in *Pacific Philosophical Quarterly*, Vol. 61, January-April 1980, pp.50-74. 7,9,13

Schaffer, J. (2016) The Metaphysics of Causation, *The Stanford Encyclopedia of Philosophy*, Edward N. Zalta (ed.), URL = <https://plato.stanford.edu/archives/fall2016/entries/causation-metaphysics/>.

Schellenberg, S. (2007). Action and self-location in perception. *Mind*, 116(463), 603-632.

Schrodinger, E. (1951). *What is life? The physical aspect of the living cell*. At the University Press.

Schütz, A. (1959). Type and eidos in Husserl's late philosophy. *Philosophy and Phenomenological Research*, 20(2), 147-165.

Schwartenbeck, P., FitzGerald, T., Dolan, R., & Friston, K. (2013). Exploration, novelty, surprise, and free energy minimization. *Frontiers in psychology*, 4, 710.

Sellars, W. (1956). Empiricism and the Philosophy of Mind. *Minnesota studies in the philosophy of science*, 1(19), 253-329.

Sender, R., & Milo, R. (2021). The distribution of cellular turnover in the human body. *Nature Medicine*, 27(1), 45-48.

Seth, A. K. (2014). A predictive processing theory of sensorimotor contingencies: Explaining the puzzle of perceptual presence and its absence in synesthesia. *Cognitive neuroscience*, 5(2), 97-118.

Seth, A. K. (2015). Presence, objecthood, and the phenomenology of predictive perception. *Cognitive neuroscience*, 6(2-3), 111-117.

- Shi, Yun Q., & Sun, H. (1999). *Image and video compression for multimedia engineering: Fundamentals, algorithms, and standards*. Boca Raton CRC Press
- Sims, M. (2021). How to count biological minds: symbiosis, the free energy principle, and reciprocal multiscale integration. *Synthese*, 199(1), 2157-2179.
- Song, C., Qu, Z., Blumm, N., & Barabási, A. L. (2010). Limits of predictability in human mobility. *Science*, 327(5968), 1018-1021.
- Spelke, E. S., & Kinzler, K. D. (2009). Innateness, learning, and rationality. *Child development perspectives*, 3(2), 96-98.
- Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. (2000). *Causation, prediction, and search*. MIT press
- Spohn, W. (2001). *Bayesian Nets Are All There Is To Causal Dependence*. In: Maria Carla Galavotti, ed. *Stochastic causality*. Stanford:CSLI Publ., pp. 157-172. ISBN 1-57586-321-9
- Sprevak, M. (2018). Triviality arguments about computational implementation. In *The Routledge Handbook of the Computational Mind* (pp. 175-191). Routledge.
- Sprevak, Mark (2021) *Predictive coding I: Introduction*. [Preprint]
- Stafford, J. (2005). *A philosophical interpretation of Judea Pearl's theory of causality* (Doctoral dissertation, University of Tasmania).
- Steinhauser, M. L., Olenchock, B. A., O'Keefe, J., Lun, M., Pierce, K. A., Lee, H., ... & Fazeli, P. K. (2018). The circulating metabolome of human starvation. *JCI insight*, 3(16).
- Stenmark, M. (2001). Evolution, Purpose and God. *Ars Disputandi*, 1(1), 44–52. <https://doi.org/10.1080/15665399.2001.10819710>
- Stepp, N., Chemero, A., & Turvey, M. T. (2011). Philosophy for the rest of cognitive science. *Topics in Cognitive Science*, 3(2), 425-437.

Sun, Z., & Firestone, C. (2020). The dark room problem. *Trends in cognitive sciences*, 24(5), 346-348.

Suppes, Patrick. 1970. *A Probabilistic Theory of Causality*, North-Holland, Amsterdam, 1970. 3

Tagkopoulos, I., Liu, Y. C., & Tavazoie, S. (2008). Predictive behavior within microbial genetic networks. *science*, 320(5881), 1313-1317.

Tenenbaum, J., Griffiths, T., & Kemp, C. (2006). Theory-based Bayesian models of inductive learning and reasoning. *Trends in Cognitive Sciences*, 10, 309–318.

Thompson, E. (2010). *Mind in life: Biology, phenomenology, and the sciences of mind*. Harvard University Press.

Thompson, E. (2018). Review of *Evolving Enactivism: Basic Minds Meet Content*. *Notre Dame Philosophical Reviews*.
<https://ndpr.nd.edu/news/evolving-enactivism-basic-minds-meet-content/>

Thomson-Jones, M. (2012). Modeling without mathematics. *Philosophy of Science*, 79(5), 761-772.

Toyama, B. H., & Hetzer, M. W. (2013). Protein homeostasis: live long, won't prosper. *Nature reviews Molecular cell biology*, 14(1), 55-61.

Turing, A.M. (1948), 'Intelligent Machinery', National Physical Laboratory Report, in B. Meltzer and D. Michie, eds, *Machine Intelligence 5*, Edinburgh University Press (1969).

van Es, T., & Hipólito, I. (2020). Free-Energy Principle, Computationalism and Realism: a Tragedy.

van Fraassen, B. (1980), *The Scientific Image*. Oxford: The Clarendon Press.

van Fraassen, B. (1982), The Charybdis of Realism: Epistemological Implications of Bell's Inequality, *Synthese* 52: 25-38.

Van Gelder, T. (1995). What might cognition be, if not computation?. *The Journal of Philosophy*, 92(7), 345-381.

van Lith, J. (2001). Ergodic theory, interpretations of probability and the foundations of statistical mechanics. *Studies in History and Philosophy of Science Part B: Studies in History and Philosophy of Modern Physics*, 32(4), 581-594.

Varela, F. (1979), *Principles of Biological Autonomy*, New York: North-Holland.

Varela, F. G., Maturana, H. R. & Uribe, R. (1974). Autopoiesis: The organization of living systems, its characterization and a model. *Biosystems*, 5(4), 187–196.

Varela, F. J. (2011). Preface to the second edition of “De Máquinas y Seres Vivos - Autopoiesis: La organización de lo vivo”. *Systems Research and Behavioral Science*, 28, 601–617. doi: 10.1002/sres.1122 Original work published 1994

Vázquez, M. J. C. (2020). A match made in heaven: predictive approaches to (an unorthodox) sensorimotor enactivism. *Phenomenology and the Cognitive Sciences*, 19(4), 653-684.

Venter, E. (2021). Toward an embodied, embedded predictive processing account. *Frontiers in Psychology*, 12, 137.

Vernon, D. (2013). Interpreting Ashby–But which One?. *Constructivist Foundations*, 9(1), 111-113.

Villalobos, M. (2013). Enactive cognitive science: revisionism or revolution?. *Adaptive Behavior*, 21(3), 159-167

Villalobos, M., & Dewhurst, J. (2018). Enactive autonomy in computational systems. *Synthese*, 195(5), 1891-1908.

Villalobos, M., & Ward, D. (2015). Living systems: Autonomy, autopoiesis and enaction. *Philosophy & Technology*, 28(2), 225-239.

Virgo, N. and Ikegami, T. (2013). Autocatalysis before enzymes: The emergence of prebiotic chain reactions. In *Advances in Artificial Life, ECAL*. MIT Press.

Virgo, N., McGregor, S., & Ikegami, T. (2014, July). Self-organising autocatalysis. In *ALIFE 14: The Fourteenth International Conference on the Synthesis and Simulation of Living Systems* (pp. 498-505). MIT Press.

von Bertalanffy L (1968) General system theory: foundations development applications. Penguin Books, Harmondsworth

Walker, G. (2013). Evolution, co-evolution, and complexity: an anniversary systems journey through the Grid. *Civil Engineering and Environmental Systems*, 30(3-4), 249-262.

Walmsley, J. (2008). Explanation in dynamical cognitive science. *Minds and Machines*, 18(3), 331-348.

Ward, D. (2016). Hurley's transcendental enactivism. *Journal of Consciousness Studies*, 23(5-6), 12-38.

Weber, A., & Varela, F. J. (2002). Life after Kant: Natural purposes and the autopoietic foundations of biological individuality. *Phenomenology and the cognitive sciences*, 1(2), 97-125.

Weisberg M (2013) Simulation and similarity: Using models to understand the world. Oxford University Press

Weslake, B. (2006). Common causes and the direction of causation. *Minds and Machines*, 16(3), 239-257.

Westheimer, G. (2008). Was Helmholtz a Bayesian?. *Perception*, 37(5), 642-650.

Whitehead, A. N. (1925). *Science and the modern world*. Macmillan.

Wicken, J. S. (1981). Causal explanations in classical and statistical thermodynamics. *Philosophy of Science*, 48(1), 65-77.

Wiese, W. (2017). What are the contents of representations in predictive processing?. *Phenomenology and the Cognitive Sciences*, 16(4), 715-736.

Wiese, W., & Friston, K. J. (2021). Examining the continuity between life and mind: Is there a continuity between autopoietic intentionality and representationality?. *Philosophies*, 6(1), 18.

Windt, J. M. (2018). Predictive brains, dreaming selves, sleeping bodies: how the analysis of dream movement can inform a theory of self-and world-simulation in dreams. *Synthese*, 195(6), 2577-2625.

Winter, O. C., & Murray, C. D. (1997). Resonance and chaos: I. First-order interior resonances. *Astronomy and Astrophysics*, 290-304.

Wittgenstein, L. (2010). *Philosophical investigations*. John Wiley & Sons. Original work published 1953

Woodward, J. (2005). *Making things happen: A theory of causal explanation*. Oxford university press.

Wostmann, B. S., Larkin, C., Moriarty, A., & Bruckner-Kardoss, E. (1983). Dietary intake, energy metabolism, and excretory losses of adult male germfree Wistar rats. *Laboratory animal science*, 33(1), 46-50.

Wright, J. C. (2001). Cryptobiosis 300 years on from van Leuwenhoek: what have we learned about tardigrades?. *Zoologischer Anzeiger-A Journal of Comparative Zoology*, 240(3-4), 563-582.

Yoshimi, J. (2016). *Husserlian phenomenology: A unifying interpretation*. Cham, Switzerland: Springer.

Yuille, A., & Kersten, D. (2006). Vision as Bayesian inference: analysis by synthesis?. *Trends in cognitive sciences*, 10(7), 301-308.

Zahavi, D. (2003). *Husserl's phenomenology*. Stanford University Press.

Zahavi, D. (2004). Phenomenology and the project of naturalization. *Phenomenology and the cognitive sciences*, 3(4), 331-347.

Zhang (2008). Wittgenstein's reconsideration of the transcendental problem. *Frontiers of Philosophy in China*, 3(1), 123-138.—

Zhaoping, L. (2006). Theoretical understanding of the early visual processes by data compression and data selection. *Network: computation in neural systems*, 17(4), 301-334.