

Mathematical approaches to the study of agents

✉ Manuel Baltieri^{1,2,†}, ✉ Keisuke Suzuki³

¹ Araya Inc., Tokyo, Japan

² University of Sussex, Brighton, UK

³ Center for Human Nature, Artificial Intelligence, and Neuroscience (CHAIN),
Hokkaido University, Sapporo, Japan

The definition of life remains one of science's most profound challenges, with contemporary approaches usually focusing on two research programmes: Darwinian evolution and self-maintenance in chemical systems. While evolution has been successfully abstracted and mathematically modelled, the concept of a self-sustaining system has so far resisted a comparable level of formalisation. This paper tackles this challenge by reframing the concept of self-sustaining system within a more abstract framework to study *agents*: goal-directed systems acting in an environment. We build on an existing conceptual framework comprising three requirements for agents: individuality, normativity (or goal-directedness), and interactional asymmetry. We then provide a systematic analysis, under a unified notation, of several mathematical approaches aiming to formalise these requirements, including the free energy principle, integrated information theory and dynamical systems. Unlike this conceptual framework, which commits to an intrinsic perspective on agency, we commit to a less ontologically committed *as-if* stance. Using this, we discuss links between identity and normativity, and a way to understand actions as if they were produced by causal interventions. Taken together, our systematic analysis clarifies the limitations of current proposals and reveals how they can work synergistically within a unified, mathematical account of agency across natural and artificial domains.

Keywords: agency, individuality, normativity, asymmetry

1. Introduction

Artificial Life (ALife) formulates an approach to the study of living systems based on the study of "life as it could be" [1]. ALife can be seen as a form of "comparative biology", that extends research in biology and origins of life to artificial universes, allowing us to study the properties of life, individuality and evolution from a more general, substrate-independent perspective [2]. This line of research is tightly connected to inquiries about Terran life, i.e. life as it is here on Earth, and its origins. To this point, the working definition of life provided by NASA, "a self-sustaining chemical system capable of undergoing Darwinian evolution" [3], elegantly captures the two central pillars that drive modern inquiries into the origins of living systems. The second part of this definition, "[C]apable of undergoing Darwinian evolution", has been the subject of extensive and successful formal abstraction, particularly in evolutionary biology [4], computer science [5] and artificial life [6]. Computational systems such as Avida [7, 8] and Tierra [9] have demonstrated that core principles of evolution including selection, replication, and mutation can be instantiated in silico, yielding complex emergent phenomena that mirror biological evolution without recapitulating the intricacies

[†]Correspondence e-mail: manuel_baltieri@araya.org



of terrestrial biochemistry. The implementation and study of modern, complex artificial worlds like Lenia [10, 11] further show the power of this approach, revealing a universe of digital creatures with their own evolving dynamics (see also [12] for a historical review of this and other related approaches in studies of cellular automata).

On the other hand, the first part, “a self-sustaining chemical system [...]”, has received comparatively less attention from both a mathematical and a computational perspectives. While we have established mathematical formalisms and powerful computational models describing the dynamics and structure of evolution over time given an entity of interest, i.e. some kind of self-sustaining (chemical) system, we mostly lack corresponding mathematical and computational apparatuses to characterise the spectrum of possible entities of interest. The main goal of this paper is to provide a unifying mathematical presentation of different proposals for the study of such entities of interest. We will, henceforth, address them as *agents*, and informally describe them, at least for now, as *goal-directed systems acting in an environment*. Computational perspectives consisting of the application of the formal ideas we will cover in the next few sections to the analysis of existing datasets or the simulation of artificial systems based on these ideas, will not be discussed here. This is for various reasons: some of the frameworks below cannot be applied at face value to the study of any realistic system, some can be approximated but the “cost” of these approximations is not clear, while others can be applied but a comparative analysis of their achievements is beyond the scope of the current work.

Agents are described in different ways across different fields, and at times simply explained away. In much of physics for instance, the world is explained by dynamical laws that contain no explicit variables that can be attributed to agents or actions. On the other hand, thermodynamic and complex-systems views link agency to irreversibility and information use [13], but risk conflating generic non-equilibrium systems (like a steam engine [14, 15]) with living ones. By contrast, fields like computer science pragmatically assumes the existence of agents: software agents act in digital environments [16], agent-based models explore macro-patterns from micro-rules [17], while AI extends the notion to learning systems, including contemporary “agentic” models and reinforcement learning’s agent-environment loop [18, 19]. Control theory likewise presupposes goals and actions, formalising controller-plant-environment loops that steer systems toward targets and thereby bridge physics and AI using pragmatic notions of goal-directedness [20]. Beyond these, philosophy [21, 22, 23, 24], psychology [25, 26, 27, 28], biology [29, 30, 31], sociology [32], economics [33], and law [34] mostly proceed as if agents simply exist, even when their ontological status is debated, studying processes such as sense of agency, organismic organisation, and social or institutional agency. These literatures thus split between traditions that assume agents (natural/social sciences) and those that tend to deny them or explain them away (much of physics/chemistry).

In this work we adopt a pragmatic stance that seeks to identify the *minimal* requirements for a possible definition of agents [35], similar in spirit to the quest to identify the minimal requirements of cognition, or rather minimally cognitive behaviour [36, 37]. As mentioned above, we start by informally addressing agents as goal-directed systems acting in an environment. *Goal-directedness* will provide a vast generalisation and abstraction to the idea of self-sustainment, *system* will address relevant parts of the notion of individuality and include entities beyond chemical systems, *action* will be used to try and differentiate an agent’s output from mere physical happenings, and *environment* as a notion complementary to the agent that is necessary to frame or carve out an identity for agents as distinct entities.

To operationalise the core components of a mathematical theory of agency, we will then take inspiration from an existing conceptual framework proposed in [38]. We will use, in particular, their three requirements of individuality, normativity (or goal-directedness) and interactional asymmetry



to organise and discuss the role of various mathematical proposals found in the literature that are relevant for a definition of agents. Similarly to the work done in artificial life for evolution, we will then focus on formal methods that abstract the concept of agent, for instance, going beyond the particular chemical or physical instantiations of systems subject to Darwinian evolution such as those part of Terran Life. We will then consider different formal proposals through the lenses of this unified conceptual framework, introducing a unifying mathematical notation to see how they can be seen as complementing each other. This approach will provide a more cohesive picture of the differences and complementary roles of these proposals, and in particular will highlight two important points that we will elaborate at the end of the work based on an *as-if* approach to agents previously introduced in [39, 40] and adopted here. Firstly, we will discuss a possible *individuality-normativity correspondence* motivated by recent work showing how individuality is deeply linked to goal-directedness. Secondly, we will consider the exact formal meaning currently attributed to actions by interventionist notions of causality [41], where the agent is seen as a causal individual acting on the environment, which are a central part of several proposals to formalise interactional asymmetry.

In Section 2 we will review the conceptual framework for a theory of agency introduced in [38], based on the ideas of individuality, normativity, and interactional asymmetry. In Sections 3 to 5 we will then proceed to classify and present under a unifying mathematical language, some of the most prominent existing proposals for mathematical theories of agents that capture at least one of the main aspects highlighted in Section 2. Finally, in Section 6 we will discuss some of the assumptions in our classification and some ideas raised by our unification attempt, highlighting especially the strong overlap of individuality and normativity and the causal status of interactional asymmetry.

2. A conceptual framework for agents

Our starting point for a conceptual framework of agents is the set of requirements for a definition of agents proposed in [38]. We introduce them here to show how they have effectively driven, or can at least be said to capture, most of the research on mathematical formalisations of different aspects of agency discussed in the next few sections. This will allow us to reframe existing proposals in a unifying language and systematically compare different classes of approaches, showing how they need not be mutually exclusive. While most of them will turn out to have implications and consequences beyond our somewhat rigid classification, we believe it is nonetheless important and useful to first and foremost characterise their main goal, background and repercussions. The three requirements advocated by [38] are:

- *Individuality*, which gives the basic unit of inquiry, the agent, as a system with some spatio-temporal persistent properties, i.e. a notion of identity. The starting point is that an agent must be a distinct entity with its own boundaries (physical, functional, etc.) that distinguishes it from its environment. In other words, given a (full) system, there is a way to factorise it into an agent and an environment. There could be multiple parallel or nested agents and environments, but the minimum requirement is to have at least one agent and one environment, which can then be further factorised if needed.
- *Normativity*, or goal-directedness, which states that an agent's interactions with the environment ought to be regulated according to specific norms or goals. Goals can correspond to ways in which an agent comes to be, or keeps on persisting, but crucially are not limited to those cases. Importantly, agents can seemingly fail to achieve goals, and still be agents. Arguably, failure is the main characterising feature of agents which puts them in contrast to non-agents which, on



the other hand, cannot fail to behave according to some goal.

- *Interactional asymmetry*, stating that an agent must be the active source of its activity, its actions, adapting its coupling with the environment rather than being merely a passive receiver of external inputs. Agency involves the system's capacity to adapt the parameters and constraints that govern its interaction with the environment.

Note that in [38] there is a strong emphasis on these requirements being *intrinsic* to the agent. We will not consider this to be a requirement here, and discuss our alternative, as-if stance in Section 6 based on [39, 40]. In the next three sections, we will look at how these three requirements have been explicitly, and at times implicitly, the main driving force behind different proposals to mathematically define an agent. Some proposals focus on a single requirement, while others attempt to encompass two, or even three at times. We do however believe that all the frameworks overviewed here have a clear, main goal in mind that well aligns with one of these three requirements (but see the individuality-normativity correspondence discussion in Section 6), which will be used to map their contributions and connections to other proposals.

In this work uppercase letters (e.g. X_t, Y_t) are used to denote indexed, usually by time $t \in \mathcal{T}$ unless otherwise stated, random variables, and lowercase (e.g. x_t, y_t) their realisations. Deterministic variables are included in this convention as degenerate random variables, i.e. X_t such that $\Pr(X_t = x_t) = 1$ for some constant x_t . Calligraphic letters (e.g. \mathcal{X}, \mathcal{Y}) denote the sets over which they take values and, by abuse of notation, the dynamical systems whose state spaces are these sets when clear from context. We use P (as in $P(\mathcal{X}), P(\mathcal{Y})$) to denote the collection of all distributions over those sets. We use the shorthand notation $p(x_t | y_t)$ to denote the conditional probability mass or density function of X_t given $Y_t = y_t$, and assume that equalities of the form $p(x_t | y_t, z) = p(x_t | y_t)$ hold for all realisations that occur with non-zero probability. When the argument of $\Pr(\cdot)$ is a random variable rather than an event, we use $\Pr(X_t)$ to denote the probability distribution (PMF/PDF) of the random variable X_t . We use the following abbreviations: $X_{a:b} = (X_a, \dots, X_b)$, $X_{:b} = X_{0:b}$, $X_{a:} = X_{a:\infty}$, and when considered as a whole $X_{0:\infty}$, we usually drop the index altogether, $X = X_{0:\infty}$.

3. Relational methods to formalise individuality

This first class of frameworks attempts to formally define an agent starting from the relation it ought to have with its environment, a mode of coupling between the agent and its environment. We associate this class of proposals with the first requirement, individuality: for an individual to exist, it must be distinct from its environment but at the same time interact with it, since complete independence would make it a quite uninteresting entity and undermine the concepts of action and agent as a whole.

The *dynamical systems approach* pioneered in [42, 43] is an attempt to formalise this idea, moving the explanatory focus from the agent's internal mechanisms to the continual mutual interaction between the agent and, importantly, its environment. The core theoretical commitment is to model the agent with states S taking values in \mathcal{S} , and the environment with states E taking values in \mathcal{E} , as two open dynamical systems coupled on the same interface (O, A) of observations and actions taking values in $\mathcal{O} \times \mathcal{A}$, whose states are governed by equations of the form:

$$\begin{aligned} \text{upd}_S : \mathcal{S} \times \mathcal{O} &\rightarrow \mathcal{S}, & \text{out}_S : \mathcal{S} &\rightarrow \mathcal{A}, \\ \text{upd}_E : \mathcal{E} \times \mathcal{A} &\rightarrow \mathcal{E}, & \text{out}_E : \mathcal{E} &\rightarrow \mathcal{O}. \end{aligned} \tag{1}$$

These correspond to classical Moore machines, and can be generalised to systems without the factorisation between updates and output mappings (Mealy machines) as in, e.g. [44, 45, 46, 47, 48]. For



simplicity, in the remainder of this work we will refer to these open dynamical systems as *systems*, and only explicitly bring out different definitions if necessary. We will also overload the notation and just refer to them using their state space when obvious from context, i.e. we will talk about an agent \mathcal{S} and an environment \mathcal{E} , assuming that each system is equipped with update and output maps of the above form unless otherwise stated. For a more general, structural treatment of systems we refer to [49, 50]. The systems above are assumed to be deterministic, and in the original work [42, 43] also in continuous time. Here, for simplicity, we introduce them in their discrete time counterpart to simplify the exposition, since generalisations to other kinds of systems (including those in continuous time), don't affect the overall arguments [49]. Specific extensions of this approach to different kinds of systems, including for instance multiset dynamical systems (forming "organisations") and possibilistic dynamical systems, could follow from existing lines of work such as chemical organisation theory [51] and viability theory [52] but the details won't be explored here.

The central insight of this approach is that these two systems can be coupled, i.e. connected on a given interface to form a single autonomous ¹ dynamical system, \mathcal{X} , whose state space is an appropriate composition of the agent and environment state spaces. The relevant behaviour of an agent is not something generated by the agent alone, but is a property of the trajectories of this entire coupled system. Agency is not located in a controller issuing commands, but distributed across the entire feedback loop. This is captured by the notion of *adaptive fit*, which provides the criterion for defining agency [42]. A system \mathcal{S} is considered an agent if it's adaptively fit: if the trajectories of the coupled system \mathcal{X} remain within a specific constraint, or viability volume, \mathcal{C} in the total state space \mathcal{X} , i.e. as long as agent and environment interact in a way that preserves certain properties represented by \mathcal{C} . For a biological organism, this constraint volume could mean survival, i.e. maintaining the integrity of its autopoietic (self-producing) processes [29]. For an agent more in general, the constraint is defined by the goal it needs to accomplish (e.g. for a robotic vacuum, the floor remaining clean). Agency is thus the capacity to interact with the environment to ensure that the coupled system's trajectory satisfies the constraint \mathcal{C} . This makes agency a relational property, defined not by what an agent is, but by its ability to maintain a specific kind of relationship with its environment over time.

This perspective is explored in more detail using gliders in the Game of Life as a toy model [54, 55, 56, 57, 58, 59]. By treating the Game of Life update rule as an artificial physics, [56] formalises a glider's self-maintaining organisation as a closed process dependency network. This provides a bottom-up, relational definition of the glider's identity. Related work [55] then characterises the agent-environment relationship by exhaustively mapping the interaction graph between a glider and its surroundings. This graph defines a macrodynamic function, $\text{upd}_{\mathcal{G}} : \mathcal{G} \times \mathcal{D} \rightarrow \mathcal{G}$, where \mathcal{G} is the set of glider states and \mathcal{D} is a set of non-lethal perturbations, see also [60]. The concept of structural coupling is then formalised by finding the mutually consistent trajectories of the agent and the environment, yielding the set of all possible glider lives.

The *free energy principle* [61, 62, 63, 64] defines an agent, or actually anything [61, 64, 15], as a system that persists over time by maintaining a statistical boundary that separates it from its environment. This requires the full system to be described as a stochastic differential equation:

$$dX_t = f(X_t)dt + \sigma(X_t)dW_t, \quad (2)$$

which possesses a non-equilibrium steady-state density, $p_{ss}(x_t)$. Here, W_t is a standard Brownian motion in \mathbb{R}^z , and the flow $f \in \mathbb{R}^n$ and $\sigma \in \mathbb{R}^{n \times z}$ are smooth vector and matrix fields, respectively.²

¹NB: not necessarily closed since we may want, as observers, to look at least at its outputs [53].

²Note that proponents of the free energy principle often emphasise "generalised coordinates" as a label for state



On this view, the defining feature of an agent is the statistical boundary known as a Markov blanket, claimed to induce a factorisation of the states of the full system, X , into four subsets. These subsets include internal states, S , of the agent, external states, E , of the environment, and blanket states, $B = (O, A)$, an interface that mediates between the internal and external states, composed of sensory states (observations), O , that are influenced by external states but not internal states, and active states (actions), A , that are influenced by internal states but not external states. It is however unclear at this stage what conditions are required for this factorisation to hold over time: it is quite straightforward to define some notion of conditional independence within a time slice, i.e. considering S_t, E_t, O_t, A_t ³, but rather non-trivial to consider an equivalent notion over time, i.e. for $S = S_{0:\infty}, E = E_{0:\infty}, O = O_{0:\infty}, A = A_{0:\infty}$ [66, 67, 68, 69, 70].

According to its proponents, an appropriate structure of influences, often addressed as sparse coupling in the system's flow [61], $f(X_t)$, ought to guarantee (but see criticisms above) a crucial conditional independence at non-equilibrium steady-state: internal and external states are statistically independent when conditioned on the blanket states. This is formally expressed as a conditional independence, $\Pr(S_t, E_t | O_t, A_t) = \Pr(S_t | O_t, A_t) \Pr(E_t | O_t, A_t)$, or more compactly $(S_t \perp E_t) | O_t, A_t$, which holds within a time slice (i.e. for a time-synchronous blanket) but need not hold once we consider the full history of observations and actions (i.e. for a time-unrolled blanket) [70, 69]. This statistical insulation is the formal definition of what constitutes an agent, it is the boundary that carves the agent out from the rest of the system. It is often claimed that a consequence of this boundary is that internal states must infer the state of the external world, based on classical results from cybernetics [71] and control theory [72]. However these claims remain hard to verify given the role of different possible Markov blankets (within a time slice, or over trajectories) and of another fundamental assumption of the free energy principle, the existence of a synchronisation map, $\zeta : \mathcal{S} \rightarrow \mathcal{E}$. This map captures an explicit relation between an agent and its environment, such that for two given functions $g_S : \mathcal{O} \times \mathcal{A} \rightarrow \mathcal{S}$ and $g_E : \mathcal{O} \times \mathcal{A} \rightarrow \mathcal{E}$, within a time slice:

$$\begin{aligned} g_S(O_t, A_t) &:= \mathbb{E}_{\Pr(S_t | O_t, A_t)}[S_t] \\ g_E(O_t, A_t) &:= \mathbb{E}_{\Pr(E_t | O_t, A_t)}[E_t] \end{aligned} \tag{3}$$

we have

$$\zeta(S_t) := g_E(g_S^{-1}(S_t)). \tag{4}$$

The synchronisation map is a function that links an agent's (expected) internal states to the (expected) external states they supposedly infer. The existence of this map is not guaranteed, as it depends on the condition that any two blanket states mapping to the same expected internal state must also map to the same expected external state, see [73] for a discussion of the linear case. This map is what allows an agent to ostensibly infer its environment, since an agent is not, in general, directly influenced by states E . It can only register its effects through sensory states O . The presence of this map implies that the internal states S parametrise a collection of variational densities, $Q_S(E) := \{Q_S(E_t)\}$, which serves as a probabilistic model of external states given blanket states. For a definition using the synchronisation map used above see [74, Eq. 85]. The persistence of this entire factorised system at non-equilibrium steady-state then is thought to imply that its dynamics must follow a specific

augmentation in continuous time for, with Brownian motion in higher dimension to reduce non-Markovian processes to Markovian ones [62].

³Note that it is quite common to stagger time indexes to explicitly show the flow of time in a sensorimotor loop [65], but here we don't do that to simplify the exposition.



path, one that minimises a functional over trajectories called variational free energy, $F(O, A, S)$ that is defined as:

$$F(O, A, S) = \mathbb{E}_{Q_S(E)}[\ln Q_S(E) - \ln \Pr(E, O, A, S)]. \quad (5)$$

This can be expressed as an upper bound on the surprisal (negative log-probability) of the agent's own states, actions and observations $F(O, A, S) = -\ln \Pr(O, A, S) + D_{\text{KL}}[Q_S(E) || \Pr(E | O, A, S)]$ and according to this proposal, agents are systems that minimise free energy, since that is what ensures that their surprisals remain within acceptable bounds over time.

The *Bayesian interpretation map* framework [74, 44] aims to overcome some of the possible limitations of the free energy principle, highlighted for instance in [68, 66, 67, 69, 75, 76, 77]. Following Dennett's intentional stance [78], this framework provides a formalisation of what it means to interpret a physical system as having beliefs and for its dynamics to be interpreted as a process of Bayesian belief updating [74]. Furthermore, it also provides criteria to interpret a physical system as taking actions according to its beliefs so to achieve a goal [44]. The core of the framework is the interpretation map, a function $\psi : \mathcal{S} \rightarrow P(\mathcal{H})$ proposed by an observer to link the physical states S_t of an agent⁴, taking values in \mathcal{S} , to abstract belief states, i.e. probability distributions, over hidden variables of a model H_t , taking values in \mathcal{H} . The agent's physical state transitions are governed by a machine kernel $\text{upd}_{\mathcal{S}} : \mathcal{O} \times \mathcal{S} \rightarrow P(\mathcal{S})$, which gives the probability of transitioning to a new state S_{t+1} given the current state S_t and a sensory input O_t . An interpretation includes both an interpretation map and a model kernel $\kappa : \mathcal{H} \times \mathcal{A} \rightarrow P(\mathcal{H} \times \mathcal{O})$ which represents the agent's (potentially incorrect) model of how its actions A_t influence the hidden state and produce sensory inputs. Note that this kernel combines update and output maps into one, unlike the dynamical systems approach discussed earlier. Interestingly, the probabilistic system \mathcal{H} can, but need not be the environment, \mathcal{E} . In many interesting cases it however will be at least a model, to some extent veridical, of the environment [46] or of the whole agent/environment combined system [53].

An interpretation is said to be consistent if the physical dynamics of the agent and the belief dynamics align perfectly. This is captured by a consistency equation, which, unpacked from its category-theoretic form, states that for any possible next physical state S_{t+1} , the belief it represents, $\psi(H_{t+1} | S_{t+1})$, must be equal to the Bayesian posterior belief, calculated from the prior belief $\psi(H_t | S_t)$ and the evidence O_t . Formally, this means that

$$\psi(H_{t+1} | S_{t+1}) = \frac{\sum_{H_t \in \mathcal{H}} \kappa(H_{t+1}, O_t | H_t, \pi(S_t)) \psi(H_t | S_t)}{\sum_{H_t, H_{t+1} \in \mathcal{H}} \kappa(H_{t+1}, O_t | H_t, \pi(S_t)) \psi(H_t | S_t)}, \quad (6)$$

where $\pi(S_t)$ gives the action taken in state S_t . This framework has been used to show deep connections to control theory and concepts of agency in that field, proving that any system satisfying the assumptions of the classic internal model principle [53] and of the good regulator [46] can be interpreted as performing Bayesian updates on a model built by an observer, which can sometimes correspond to the environment.

Cartesian frames [79] provide on the other hand a mathematical framework for modelling agency that is designed to overcome the limitations of traditional models that treat the agent-environment boundary as a fixed, primitive reality. Cartesian frames constitute a way to carve an agent-centric perspective of choice onto an objective, physical set of all possible outcomes of a particular full system, allowing for an analysis of agency across different levels of description. The fundamental building

⁴More precisely, a system that can be interpreted to be an agent.



block of this framework is the Cartesian frame, a Chu space $\mathcal{J} = (\mathcal{S}, \mathcal{E}, \cdot)$ over a set \mathcal{X} , a structure that generalises a decision matrix. Here, \mathcal{S} is the set of all possible ways the agent can choose to be, e.g. its states, \mathcal{E} is the environment's set of possible states, and the outcome function $\cdot : \mathcal{S} \times \mathcal{E} \rightarrow \mathcal{X}$ maps agent-environment pairs to a set of possible worlds of the full system \mathcal{X} . An agent is defined as a system \mathcal{S} that can be obtained from the factorisation of a larger system \mathcal{X} whose state can be composed with the environment \mathcal{E} to produce \mathcal{X} itself. An agent can be carved out to be a separate system from the environment, but its identity and interface with the environment are subject to constant change. Note that here there is no explicit notion of dynamics, although one could understand the outcome function as mapping between trajectories, therefore folding time into the sets \mathcal{S} and \mathcal{E} , nor of an agent's observations \mathcal{O} or its actions \mathcal{A} , which are instead derived as follows. Observations are defined as:

$$\mathcal{O} = \{\mathcal{X}' \subseteq \mathcal{X} \mid \forall s_0, s_1 \in \mathcal{S}, \exists s \in \mathcal{S}, s \in \text{if } (\mathcal{X}', s_0, s_1)\} \quad (7)$$

where $\text{if } (\mathcal{X}', s_0, s_1)$ denotes the set of all $s \in \mathcal{S}$ such that for all $e \in \mathcal{E}$, $(s \cdot e \in \mathcal{X}') \rightarrow (s \cdot e = s_0 \cdot e)$ and $(s \cdot e \notin \mathcal{X}') \rightarrow (s \cdot e = s_1 \cdot e)$. Agents in this setting observe events, which are either true or false, not variables in full generality. We will say that given a frame \mathcal{J} , its observables, \mathcal{O} , are the set of all \mathcal{X}' such that \mathcal{J} 's agent, S , can observe whether $x' \in \mathcal{X}'$ is true. Actions are defined as:

$$\mathcal{A} = \{\mathcal{X}' \subseteq \mathcal{X} \mid \exists s \in \mathcal{S}, \forall e \in \mathcal{E}, s \cdot e \in \mathcal{X}'\} \cap \{\mathcal{X}'' \subseteq \mathcal{X} \mid \exists s \in \mathcal{S}, \forall e \in \mathcal{E}, s \cdot e \notin \mathcal{X}''\} \quad (8)$$

where \mathcal{X}' gives a set of world states an agent can ensure, and \mathcal{X}'' a set of world states an agent can prevent. The framework's power comes from its use of morphisms or maps between frames. A morphism from frame $\mathcal{J} = (\mathcal{S}, \mathcal{E}, \cdot)$ to $\mathcal{J}' = (\mathcal{S}', \mathcal{E}', *)$ is a pair of functions $(g : \mathcal{S} \rightarrow \mathcal{S}', h : \mathcal{E} \rightarrow \mathcal{E}')$ that satisfies the adjointness condition: $s \cdot h(e') = g(s) * e$ for all acts $s \in \mathcal{S}$ and environment states $e' \in \mathcal{E}'$. These maps allow for formal operations like moving between different levels of description (coarsening or refining the set of possible worlds \mathcal{X}), and composing or decomposing agents by modelling *subagency* of different kinds, additive to describe agents after some commitment to particular subsets of choices or states, and multiplicative to define agents within a larger agent.

4. Prediction-based methods to formalise normativity

This second class of frameworks posits that agents are systems that possess some kind of predictive power, about themselves, their environment or both. On this view, normativity is formalised by the idea that when agents achieve a goal, they are, or can be interpreted as, necessarily predicting features that are important to achieve their goal [71, 61, 53, 46]. These goals can involve the agent's survival, but in general can be about other things entirely if we specify an appropriate viability volume [42, 52, 53, 46] or condition on a different model [80, 81, 82, 83, 84], and they can include self-destructive behaviour, as in the case of autothysis [39]. This approach doesn't necessarily clearly specify a way to describe systems that at times may *fail* to achieve a goal, but are nonetheless systems that could be labelled as agents, see also Section 6. At this stage we also remain agnostic as to whether this ability to predict is intrinsic for the agent or in the eye of the beholder [40, 14] and return to this briefly in Section 6. We then note that some of the proposals in the previous section, chiefly the free energy principle, make strong claims about agents and their ability to predict features of their environment. This is related to another point we will discuss at the end, a correspondence between individuality and normativity.



The *information theory of individuality* [85] provides a quantitative way to identify individuals at any scale by defining them as aggregates that maintain temporal integrity by propagating information from their past into their future. Building on earlier work on information-theoretic autonomy [86], this approach analyses the total predictive information that the present state of an agent, S_t , and the state of the environment, E_t , jointly provide about the agent's next state, S_{t+1} , captured by the mutual information $I(S_{t+1}; S_t, E_t)$. The framework's key insight is that this total information can be decomposed in two distinct ways using the chain rule for mutual information:

$$I(S_{t+1}; S_t, E_t) = \underbrace{I(S_{t+1}; S_t)}_{\text{organismal individuality}} + \underbrace{I(S_{t+1}; E_t | S_t)}_{\text{environmentally determined individuality}} = I(S_{t+1}; E_t) + \underbrace{I(S_{t+1}; S_t | E_t)}_{\text{colonial individuality}} \quad (9)$$

These decompositions define a taxonomy of individuality. Organismal individuality, originally proposed in [86] as a measure of autonomy for systems that control their environment, measures the information that an agent's own past provides about its future, ignoring the environment. It captures highly autonomous individuals that are largely in control of their own dynamics. Colonial individuality, see again [86], instead measures autonomy for agents that are driven by their environment. More specifically, the additional information the system's past provides about its future, once the influence of the environment's past is already taken into account.⁵ This ought to capture more loosely-coupled individuals, like microbial colonies, whose persistence depends heavily on ongoing interaction. To see this, [85] employs tools from partial information decomposition to more carefully unpack the influence of environment and agent. We do not however venture into this level of detail here. Environmental determined individuality is on the other hand defined by the information flow from the environment, which identifies entities whose structure is largely imposed by external forces.

The *agent description* of [82] provides a formalisation of Dennett's intentional stance [78] by framing the question of agency as a Bayesian model comparison problem. In this setup, an observer is presented with a candidate agent's⁶ behaviour, i.e. a trajectory of outputs (actions), $A_{:t}$, given inputs (observations) for a candidate agent, $O_{:t}$, and must decide which of two explanatory models is more probable. The model M_d describes the system as a device, a reactive machine whose behaviour is explained by a direct, mechanistic input-output mapping, which is associated to Dennett's physical stance. The likelihood of the trajectory under this model is a weighted average over all possible device descriptions:

$$M_d(A_{:t} | O_{:t}) = \sum_{d \in \mathcal{M}_d} d(A_{:t} | O_{:t}) \omega_d, \quad (10)$$

where $d : (\mathcal{O} \times \mathcal{A})^* \rightarrow P(\mathcal{A})$ produces a probability distribution of outputs given an interaction history of inputs $O_{:t}$ and outputs $A_{:t}$, and ω_d is a prior over devices. The model M_s identifies instead the system as an agent whose behaviour appears as the result of optimising a goal or utility function $u \in \mathcal{U}$, proposed to correspond to Dennett's intentional stance. The likelihood is derived using inverse reinforcement learning:

$$M_s(A_{:t} | O_{:t}) = \sum_{u \in \mathcal{U}} \pi_u(A_{:t} | O_{:t}) \omega_u, \quad (11)$$

⁵For Gaussian variables, this also corresponds to Granger autonomy, or G-autonomy, as proposed in [87], see [88].

⁶A candidate agent is a system with the potential to be labelled as an agent when certain conditions are met.



where $\pi_u : (\mathcal{O} \times \mathcal{A})^* \rightarrow P(\mathcal{A})$ is the optimal (or near-optimal) policy for that utility function, and ω_s is a prior over devices. The framework then uses Bayes' rule to calculate the posterior probability that the system is an agent:

$$p(\text{agent} \mid A_{:t}, O_{:t}) = \frac{M_s(A_{:t} \mid O_{:t})}{M_d(A_{:t} \mid O_{:t}) + M_s(A_{:t} \mid O_{:t})} \quad (12)$$

assuming equal priors. The prior weights (ω_d, ω_u) , which can be based on the Kolmogorov complexity of describing the device or the utility function, provide an intuition for this definition by thinking about the compression aspect that arises from the principles of algorithmic information theory. In particular, a goal-based description is the better description if the complexity of describing the goal is lower than the complexity of describing the intricate mechanistic policy that achieves it. This framework thus quantifies the intuition that agency is a useful concept when it provides a more compressed, predictive explanation for a system's behaviour, and while reliant on Bayesian model comparison, extensions can make use of different objectives, see for instance [89].

Complete local integration [30, 90, 91] seeks to describe an important property of agents within a dynamical system by quantifying the internal coherence of the parts that constitute them. The proposal can be described by the idea that a spatiotemporal pattern, or ι -entity, is a system whose every part (which is itself another spatiotemporal pattern) makes every other part more probable. This is formalised through the measure of specific local integration. For a given spatiotemporal pattern, i.e. a set of fixed values of a candidate agent's states $S_{\mathcal{W}} = s_{\mathcal{W}}$ indexed over time and space by $\mathcal{W} \subseteq \mathcal{T} \times \mathcal{R}$, and a partition ξ of the set of variables \mathcal{W} into blocks b , the specific local integration is defined as:

$$\text{mi}_{\xi}(S_{\mathcal{W}}) = \log \frac{\Pr(S_{\mathcal{W}})}{\prod_{b \in \xi} \Pr(S_b)}. \quad (13)$$

A positive specific local integration indicates that the whole pattern $S_{\mathcal{W}}$ is more probable than cases where its parts (defined by the partition ξ) are independent. To identify truly integrated entities that can serve as candidate agents, the framework introduces then complete local integration, which is the minimum specific local integration over all possible non-unit partitions of the pattern:

$$\iota(S_{\mathcal{W}}) = \min_{\xi \in \mathcal{L}(\mathcal{W}) \setminus 1_{\mathcal{W}}} \text{mi}_{\xi}(S_{\mathcal{W}}), \quad (14)$$

where $\mathcal{L}(\mathcal{W})$ is the partition lattice of \mathcal{W} , i.e. the set of partitions of \mathcal{W} partially ordered by refinement, and $1_{\mathcal{W}}$ is the unit partition.

A spatiotemporal pattern is then defined as an ι -entity if its complete local integration is greater than zero, $\iota(S_{\mathcal{W}}) > 0$. The claim here is thus that candidate agents are systems with positive complete local integration. This means that a pattern is integrated with respect to every possible way of decomposing it. The importance of positive complete local integration is formally established through a *disintegration theorem*, demonstrating that ι -entities are the fundamental, maximally disintegrated components of a system's trajectory. Positive complete local integration also implies that, from an information-theoretic perspective, the surprisal (information content) of the whole entity is lower than the sum of the surprisal of the parts (cf. the minimisation of surprisal in the free energy principle [61, 62, 64]). This tells us that the whole is spatially more predictable (lower surprisal) than the parts, and temporally more predictable because the presence of an entity at a certain point in time makes the existence of a corresponding entity at a subsequent moment more probable. From a



coding-theory perspective, the fact that a ι -entity has positive integration across all possible partitions means that the most efficient, compressed description of the pattern is always the one that treats it as a single, integrated whole. Any attempt to describe it as a composition of independent parts results in a less efficient, longer description, meaning that an ι -entity captures, in some sense, a pocket of predictability that arises from a system's dynamics.

Rosen's theory of *anticipatory systems* [92, 93, 94] usually refers to the study of biological systems, but similarly to what we do here, focuses on aspects of living systems not directly related to their ability to undergo (Darwinian) evolution. Because of this, we will treat this as a framework with elements that attempt to capture aspects of agency, in particular claiming that agents are anticipatory systems. Anticipatory systems are closed systems characterised by a modelling relation expressed as a commuting diagram between a candidate agent \mathcal{S} , often addressed as *natural system* in this literature, and a system, \mathcal{H} , usually called *formal system* which serves as its epistemic, predictive model or abstraction, and is given by an observer. The maps representing state transitions are called causal entailment for the physical laws in the agent, e.g. its dynamics, $c : \mathcal{S} \rightarrow \mathcal{S}$, and inferential entailment for the logical rules of the formal system, e.g. the dynamics, $i : \mathcal{H} \rightarrow \mathcal{H}$. The mappings between the systems are the encoding, $\varepsilon : \mathcal{S} \rightarrow \mathcal{H}$, which represents measurement of the natural system, and the decoding, $\delta : \mathcal{H} \rightarrow \mathcal{S}$, which represents an interpretation of properties of a formal model in the natural world. A modelling relation between \mathcal{S} and \mathcal{H} exists if the following diagram commutes [95, 4.14] in the appropriate setup⁷:

$$\begin{array}{ccc} \mathcal{S} & \xrightarrow{\varepsilon} & \mathcal{H} \\ c \downarrow & & \downarrow i \\ \mathcal{S} & \xleftarrow{\delta} & \mathcal{H} \end{array} \quad (15)$$

meaning that the evolution in the natural system \mathcal{S} must be equivalent to the path of encoding, inferring, and decoding, i.e. $c = \delta \circ i \circ \varepsilon$, and

$$i = \varepsilon(c), \quad (16)$$

a condition stating that the updates of the formal model must be related to the updates of the natural system, or agent.⁸ The system \mathcal{H} can be seen as a predictive model because from the perspective of \mathcal{S} , one can first encode its state via ε , then process it as part of the model using i and finally decode the result with δ , checking that it is equivalent to its own evolution [94]. We note the resemblance with some of the arguments formalised in [53], in particular when the system \mathcal{H} is a model comprising both the agent and its environment, however the exact connection remains unclear and won't be explored here.

5. Causality-based methods to formalise interactional asymmetry

The third class of frameworks seeks to operationalise a potential interaction asymmetry between agent and environment in terms of agent causation, by arguing that an agent's internal states are the genuine causes of its actions and their effects in the world. These frameworks aim to identify and quantify this causal power, overcoming some of the philosophical challenges of defining causality. We note that for

⁷We can imagine these objects to be sets in the category of sets and relations for instance, **Rel**, with arrows given by relations or multi-valued functions, since this is the starting point of the categorical setup of anticipatory systems [94, 96].

⁸Different attempts to formalise this using a categorical setup, as a functor between categories of natural systems and of formal models, can be found in, for instance, [95, 97, 94], but the details won't be covered here.



these proposals too, it is unclear whether causality ought to be portrayed as an intrinsic feature of agents (it often is presented as such) and return to this in Section 6.

According to the framework of *actual causation* [98], based on integrated information theory (IIT) [99, 100], an agent can be defined as a physical system that causes its own actions [98, 101]. The existence of an agent is assumed to be determined by the system's internal cause-effect power, defined as the degree of causal irreducibility, and quantified by integrated information. As part of the setup, we start with a full system with states X_K indexed over time and state $K \subseteq \mathcal{T} \times \mathcal{V}$, corresponding to a multivariate stochastic process. Full time slices of variables are denoted as X_t . Integrated information is defined for an occurrence [98], a subset of a full time slice of variables together with the values those variables take, $X'_t = x'_t \subseteq X_t = x_t$, as following:

$$\phi(X'_t) = \min_{t \pm 1} \left(\max_Q \left(\min_{\xi} \left(\text{Dist} \left[r(Q_{t \pm 1} | X'_t); \xi(r(Q_{t \pm 1} | X'_t)) \right] \right) \right) \right) \quad (17)$$

We unpack this equation in more detail here. We define another occurrence, $Q_t = q_t$, with variables Q chosen to maximise $\phi(X'_t)$ via \max_Q , and then label the complements of our two occurrences as $\bar{X}'_t = \bar{x}'_t$ and $\bar{Q}_t = \bar{q}_t$. To study the constraints that variables X'_t impose on Q_{t+1}, Q_{t-1} (NB: future and past of Q_t), we then define effect and cause repertoires respectively [98]. Using the do-operator from do-calculus [41] to simulate a notion of physical intervention that deletes certain functions from a model by replacing them with a constant $X = x$, we define the *effect repertoire* as $r(Q_{t+1} | X'_t = x'_t) = \prod_i 1/|\bar{\mathcal{X}}'| \sum_{\bar{x}'_t \in \bar{\mathcal{X}}'} \Pr(Q_{i,t+1} | \text{do}(X'_t = x'_t, \bar{X}'_t = \bar{x}'_t))$. This quantity "causally marginalises" variables \bar{X}'_t by imposing a uniform distribution over all $\bar{x}_t \in \bar{\mathcal{X}}_t$ using $1/|\bar{\mathcal{X}}'|$, leaving only possible dependencies on the occurrence $X'_t = x'_t$. On the other hand, the *cause repertoire* is given by: $r(Q_{t-1} | X'_t = x'_t) = 1/L \prod_i \sum_{\bar{q}_{t-1} \in \bar{\mathcal{Q}}} \frac{\Pr(X'_{i,t} = x'_{i,t} | \text{do}(Q_{t-1}, \bar{Q}_{t-1} = \bar{q}_{t-1}))}{\sum_{x'_{t-1} \in \mathcal{X}} \Pr(X'_{i,t} = x'_{i,t} | \text{do}(X_{t-1} = x_{t-1}))}$ with L a normalising constant.

This quantity causally marginalises instead variables \bar{Q}_{t-1} to remove common inputs from non- Q_{t-1} variables. Next we introduce ξ for a partition of occurrences in the cause and effect repertoires, see [98] for details.⁹ We have a minimisation, \min_{ξ} , which is used to select the one that makes the least difference to cause and effect repertoires. This gives the so-called minimum information partition. To calculate this difference, we need a distance $\text{Dist}[-; -]$ ¹⁰, which captures the degree to which a system is integrated. A system is said to be integrated if its causal dynamics cannot be decomposed into independent parts. Formally, this means that larger distances give larger ϕ , capturing the idea that a system's joint cause-effect structure cannot be factorised into separate substructures without loss of information. Finally, the outer $\min_{t \pm 1}$ says that if no information is specified by $X'_t = x'_t$ about either $Q_{t-1} = q_{t-1}$ or $Q_{t+1} = q_{t+1}$, this will return 0, making $\phi(X_t = x_t) = 0$. Occurrences for which integrated information is positive, i.e. $\phi(X'_t = x'_t) > 0$, are called mechanisms. This formalism has been used to define causal strength, R , a measure of irreducible information effects applied to the study of the causes, in a subset of an agent's states, S'_{t-1} , of an agent's actions, A_t [102]. This measure is given by:

$$R(A_t) = \max_{S'} \left(\min_{\xi} \left(\log_2 \left(\frac{r(S'_{t-1} | A_t)}{\xi(r(S'_{t-1} | A_t))} \right) \right) \right)$$

where $r(S'_{t-1} | A_t)$ is the cause repertoire defined above.

⁹See also the description given for complete local integration [30, 90, 91] which is in part inspired by this work.

¹⁰For the specific choice of a measure $\text{Dist}[-; -]$ in different versions of IIT, see [100].



On the other hand, IIT based approaches have also been used to define an agent's causal boundary as a local maximum of Φ , i.e. as a macro-level entity with greater causal power than its underlying micro-constituents [103, 104, 105]. The system-level integrated information, Φ , is a measure of integrated information capturing the cause-effect structure $\mathcal{C}(X'_t)$ among mechanisms (occurrences such that $\phi(X'_t = x'_t) > 0$) [101]:

$$\Phi(X'_t) = \min_{\xi} \left(\text{Dist} \left[\mathcal{C}(X'_t); \mathcal{C}(\xi(X'_t)) \right] \right). \quad (18)$$

For large systems, a highly integrated state is achieved by being poised near a critical point (a phase transition), where integration diverges. Based on this, [103, 104, 105] introduce a notion of asymmetry to define autonomous agents as systems that must be the (sub)system which causally drives their interaction with the environment, rather than being determined by it. This is shown for example in [103], where a measure of autonomy is given by the agent's ability to dynamically modulate its causal coupling with the environment. Here, the agent, whose only states are those of its interface (O, A) (no internal states), can in some cases be the most integrated system, above and beyond the full system \mathcal{X} including the environment \mathcal{E} , i.e. $\Phi_{OA} > \Phi_{OAE}$, and thus can be seen as setting its boundary to demarcate itself from the environment. In other situations, the full agent-environment system becomes the most integrated unit ($\Phi_{OAE} > \Phi_{OA}$), with the agent recruiting (parts of) the environment to change its identity and causal presence on the remaining parts of the environment.

Mechanised Causal Graphs [106] provide another proposal of causality-based methods. The core idea is that "agents are systems that would adapt their policy if they were aware that their decisions influenced the world in a different way" [106]. The starting point is a mechanised structural causal game, a tuple $\mathcal{M} = \langle N, P_N, G, \Pr(P_N) \rangle$. As in the standard definition of structural causal models, N is a set of endogenous variables indexed over state \mathcal{V} , $\{N_v\}_{v \in \mathcal{V}}$, P_N is a set of exogenous variables indexed by endogenous variables, $\{N_v = g_{N_v}(N, P_{N_v})\}$ is a set of structural equations indexed by endogenous variables, and $\Pr(P_N)$ is a distribution over the exogenous variables. A mechanised structural causal game extends a structural causal model in two ways. The first extension comes from the definition of mechanised structural causal models and is given by a factorisation of endogenous variables into object-level and mechanism variables, $N = X \cup \tilde{X}$. For each object-level variable X , a corresponding mechanism \tilde{X} parametrises the relative structural equation. The second extension is the assumption that endogenous variables can be factorised into chance (environment's state), decision (actions) and utility variables, which combined with the first extension give $X = E \cup A \cup U$ and $\tilde{X} = \tilde{E} \cup \tilde{A} \cup \tilde{U}$.

The crucial idea behind this framework is that when an object-level variable is a decision, A , its associated mechanism, \tilde{A} , is formally treated as an agent's decision rule or policy. The mechanism is the reasoning that produces the decision, which is then assumed to be the action. This framework thus provides a way to build mechanised structural causal games and identify decision nodes and their associated mechanisms. To do this, it distinguishes between two fundamental types of interventions [41]: standard object-level interventions, $\text{do}(X = x)$, which set a variable's value, and mechanism interventions, $\text{do}(\tilde{X} = \tilde{x})$, which alter the function g^X itself, and must satisfy the following condition: given an object-level variable X_v and a related mechanism \tilde{X}_v :

$$\Pr(X_v | \text{Pa}^{X_v}, \text{do}(\tilde{X}_v = \tilde{x}_v)) = \Pr(X_v | \text{do}(\tilde{X}_v = \tilde{x}_v)). \quad (19)$$

For decisions, this means that one can distinguish between interventions on the action itself, $\text{do}(A = a)$, and interventions on the agent's policy, $\text{do}(\tilde{A} = \tilde{a})$, which force the agent to adapt. The latter



corresponds to a change in the world's dynamics that a candidate agent can be made aware of and adapt to. Agency thus is discovered by performing interventions on mechanism variables, in particular when these are deemed to be policies. To determine this, [106] defines a specific causal structure, the terminal mechanism edge. This edge links a utility's mechanism to a decision's mechanism, $\tilde{U} \rightarrow \tilde{A}$, and is defined by two conditions tested via mechanism interventions, interventions that sever an object-level node from its consequences. The first condition determines if a node is a utility node. The rule is that an agent's policy, \tilde{A} , must remain sensitive to the utility's mechanism, \tilde{U} , even when the utility's own consequences are nullified. For any structural intervention on the children of U , $\text{do}(\text{Ch}^U)$, there must exist a mechanism intervention $\text{do}(\tilde{U} = \tilde{u})$ such that:

$$\Pr(\tilde{A} \mid \text{do}(\tilde{u}, \text{Ch}^U)) \neq \Pr(\tilde{A} \mid \text{do}(\text{Ch}^U)). \quad (20)$$

This identifies U as a utility node because the agent values it for its own sake, not for its downstream effects. The second condition tests if a node is a decision node. To be one, the agent's policy, \tilde{A} , must become insensitive to the utility's mechanism when the decision's own consequences are nullified. This is the case if, for any structural intervention on the children of A , $\text{do}(\text{Ch}^A)$:

$$\Pr(\tilde{A} \mid \text{do}(\tilde{u}, \text{Ch}^A)) = \Pr(\tilde{A} \mid \text{do}(\text{Ch}^A)). \quad (21)$$

This identifies A as a decision node because the policy only adapts for the sake of achieving the decision's consequences. Notice how this approach remains agnostic with respect to time and to a full factorisation of a system that includes an agent's states and its observations, which could perhaps be represented by some mechanism variables, but whose status is otherwise unclear.

Using different measures of causality and a different notion of intervention from that of Pearl [41], *semantic information* is defined as the information that an agent has about its environment that is causally necessary for the agent to maintain its own existence over time [107]. This information ought to be causally necessary for an agent because often an agent needs to take actions in the environment to ensure its existence, and a mechanism to disentangle relevant from irrelevant information is thought to be a good candidate explanation for how this can be achieved efficiently [108]. The theory begins by defining a (candidate) agent as a far-from-equilibrium system \mathcal{S} that actively works to maintain its own existence by keeping itself in a low-entropy state. This degree of existence is quantified by a viability function, V . The primary function proposed is the negative Shannon entropy, H , of the agent's marginal state distribution, $\Pr(S_t)$, at a time t , $V(\Pr(S_t)) = -H(S_t)$. The core claim is that an agent maintains its viability by using semantic information, which is formally defined as the portion of syntactic (Shannon) information the agent has about its environment, \mathcal{E} , that is causally necessary for its self-maintenance. The trajectory of the agent under its actual dynamics is compared to *intervened* trajectories where correlations between the agent and environment have been scrambled. Here, intervening does not imply the presence of a surgical Pearl-style $\text{do}(\cdot)$ that sets variables to fixed values [41]. Instead, the framework simply perturbs the information channel from environment to agent by coarse-graining what the agent can distinguish about \mathcal{E} , performing a "pre-garbling" of the channel, while leaving the environment's own dynamics otherwise unchanged. Concretely, one defines kernels such that an agent has only access to a degraded version of E_t , i.e. $\mu(E_t)$. This is achieved by replacing a channel in a fully observable setup $\mathcal{E} \rightarrow \mathcal{S}$ with a channel $\mathcal{E} \xrightarrow{\mu} \mathcal{E}' \rightarrow \mathcal{S}$ where we can assume that μ is a surjective map and $\mathcal{E}' \subseteq \mathcal{E}$. By varying μ , one obtains a family of partial interventions that scramble correlations between \mathcal{S} and \mathcal{E} to controlled degrees. If μ is a bijection, the intervened and actual processes coincide up to a permutation, if μ is constant we get a "fully



“scrambled” setup where all environment states look the same. The value of information, ΔV_t , is then defined as the loss in viability that results from such an intervention at time t :

$$\Delta V_t = V(\Pr(S_t)) - V(\hat{\Pr}^\mu(S_t)), \quad (22)$$

where $\Pr(S_t)$ is the actual distribution, while the intervened one is given by $\hat{\Pr}^\mu(S_t) := \sum_{S_0, E_0, E_t} \Pr(S_t, S_t | S_0, E_0) \hat{\Pr}^\mu(S_0 | E_0) \Pr(E_0)$ with $\hat{\Pr}^\mu(S_t | E_t) := \Pr(S_t | \mu(E_t))$ ¹¹. To isolate the information meaningful for the agent at time $t = 0$, the theory defines an optimal intervention, $\hat{\Pr}^{\text{opt}}(S_t)$. This is the intervention that preserves the minimum amount of initial mutual information, $I(S_0; E_0)$, required to achieve the exact same viability as the unintervened system at a time τ :

$$\hat{\Pr}^{\text{opt}}(S_t) \in \arg \min_{\hat{\Pr}^\mu} I_{\hat{\Pr}^\mu}(S_0; E_0) \quad \text{s.t. } \Delta V_\tau = 0. \quad (23)$$

All information that survives this scrambling process is assumed to be causally necessary and therefore semantic. The amount of semantic information is the quantity of syntactic information, Σ , preserved under this optimal intervention. For information stored in the initial state, this corresponds to:

$$\Sigma_{\text{stored}} = I_{\hat{\Pr}^{\text{opt}}}(S_0; E_0). \quad (24)$$

This framework considers also a notion of observed semantic information that can be acquired over time, characterising how agents dynamically interact with an environment. For this, this framework defines a different kind of intervention, one that preserves the minimum amount of transfer entropy over one time step, $\text{TE}(E_{:t-1} \rightarrow S_t)$ while still achieving the same viability as the unintervened system:

$$\hat{\Pr}^{\text{opt}}(S_{:t}, E_{:t}) \in \arg \min_{\hat{\Pr}^\mu} \sum_{t=0}^{\tau} \text{TE}_{\hat{\Pr}^\mu}(E_{:t-1} \rightarrow S_t) \quad \text{s.t. } \Delta V_\tau = 0. \quad (25)$$

where $\hat{\Pr}^\mu(S_{:t}, E_{:t})$ is given by iterating over $\hat{\Pr}^\mu(S_t, E_t | S_{t-1}, E_{t-1}) = \Pr(E_t | S_t, S_{t-1}, E_{t-1}) \hat{\Pr}^\mu(S_t | S_{t-1}, E_{t-1})$ with $\hat{\Pr}^\mu(S_t | S_{t-1}, E_{t-1}) := \Pr(S_t | S_{t-1}, \mu(E_{t-1}))$, and with $\hat{\Pr}^\mu(S_0, E_0) := \Pr(S_0, E_0)$. This leads to a notion of observed semantic information which is derived by intervening on the dynamic flow of information, as measured by the amount of transfer entropy that remains under this optimal intervention:

$$\Sigma_{\text{observed}} = \sum_{t=1}^{\tau} \text{TE}_{\hat{\Pr}^{\text{opt}}(S_{:t}, E_{:t})}(E_{:t-1} \rightarrow S_t) = I_{\hat{\Pr}^{\text{opt}}}(S_t; E_{:t-1} | S_{:t-1}). \quad (26)$$

This is the minimal rate of information flow that is causally necessary for the agent to regulate its coupling with the environment, and to take appropriate actions to maintain its existence. It quantifies how an agent’s perceptions are meaningful and essential for guiding its behaviour over time. These measures lead to a rigorous definition of an agent: a physical system is an autonomous agent to the extent that it possesses a high value of semantic information, Σ , both in terms of stored and observed semantic information. This notion is derived from the system’s own dynamics and it’s relative to a goal, in this case the imperative to maintain a low-entropy existence, but could in principle be defined for different goals.

¹¹The authors would like to thank Artemy Kolchinsky for his assistance in clarifying this.



6. Discussion

Our work provides a characterisation of different formal approaches to the study of agents through the lenses of a unifying conceptual framework [38]. This framework relies on three requirements: individuality, normativity or goal-directedness and interactional asymmetry, which we believe to capture desirable aspects of a theory of agents. This framing allows us to consider the strengths of each proposal, and what they appear to be missing given these three different aspects. Under this light, none of the proposals analysed here meets all the three criteria.

6.1. The necessity of individuality, normativity and interactional asymmetry

While our characterisation aimed to classify each theory as mainly fulfilling only one of these requirements, some theories can be seen as implementing two of them, especially in the sense we will see below for individuality and normativity, but none appear to capture them all. This can either mean that none of these frameworks currently constitute a complete theory of agency, that these theories fail to formalise one or more of these requirements, or that some of the requirements described here are in fact not necessary. The first case is stressed by various proposals, and does not in and by itself constitute a problem since different frameworks may in fact have different goals. As for the second case, we highlight for instance that if we strictly followed [38], most of the frameworks discussed here would fail to provide an account of *intrinsic* agency (more on this below). Furthermore, while different proposals seek to characterise normativity in agents in terms of their ability to achieve a goal, none of them seem to clearly discuss the idea that systems can fail to achieve a goal, and still be called agents [38]. If we consider a viability volume \mathcal{C} , as in [42], an agent may for instance start outside of it and *eventually* be found inside within it, or may be momentarily be pushed out by some unfavourable disturbance just to find its way back inside as soon as possible. All of this while still being labelled as an agent. This is a rather important, potential shortcoming that we believe can be addressed in future work, perhaps with the use of relevant tools from, e.g. temporal logic [109]. For the latter, the major candidate is perhaps interactional asymmetry.

This is because it is hard to imagine a way to talk about agents that doesn't include a formalisation of some notion of individuality, defining the unit of investigation, and goal-directedness, giving a way to describe what distinguishes agents from *mere* physical systems. On the other hand, it appears to be less of a stretch to propose a definition of minimal agency that posits that the actions of an agent are simply the outputs of a system labelled as an agent following criteria other than interactional asymmetry. It is also the case because, while this requirement seems to make sense on the surface, it is in fact not formulated very precisely in the original work, as highlighted by [40]. When described using causal language, it is then perhaps more suggestive of a definition of more complex agents that may require some notion of causality for cognition [108] to be properly formalised. This would then lead to a notion of agency that is quite different from that of minimal agency we investigated here as part of a definition of, among others, simple(r) living systems.

To further consider the role of interactional asymmetry, consider for instance a typical example of niche construction, ant gardens [110], cooperative nest structures in which ants appear to cultivate specific epiphytic plants whose growth, in turn, reinforces and sustains the colony's habitat. On a short timescale, ants seem to be acting, by "farming" specific epiphytic plants. They select seeds, embed them in nest structures, and maintain conditions that allow the plants to grow, producing a modified environment that benefits the colony across generations. On this view, ants appear as agents and plants as part of the environment. However, on longer timescales, plants themselves can



be seen as acting, by “farming” ants: they cultivate their own colony by producing chemical cues that induce ants to collect and plant their seeds, and provide rewards that shape the behaviour of both individual ants and the colony they belong to. In this picture, plants would be agents and ants part of the environment. Is there a (more) correct way to carve out agent and environment from this setup? Neither perspective appears fundamentally wrong, and certainly it is possible to consider even more scenarios, where there are no agents at all because these systems are too simple to be labelled as agents, or where both systems are agents of some, maybe different, kind. Does this situation exemplify a case where interaction asymmetry is not helpful or even required? Maybe, but not necessarily. We think however that this is a call to consider the role of observers more explicitly and carefully. The apparent agency in our example depends in fact on the observer’s framing: zoom in, and ants are gardeners, zoom out, and plants are orchestrating ant behaviour.

6.2. As-if agency

One of the core features of the formalisation of agency proposed in [38] is that agents are systems whose individuality, normativity and asymmetry are *intrinsic* defining properties of what it means to be an agent. Identity should not be imposed by an external observer, as with an artifact, but ought to be an intrinsic and ongoing achievement of the system itself. An agent’s norms should also not be arbitrary or merely assigned by a designer, but rather generated from within the system itself, based on its need to maintain its own organisation. Similarly, interactional asymmetry ought to be induced by the internal causal structure of an agent, and not be dependent on external factors. Among the approaches we analysed and reframed in this work, several however explicitly assume that agency is, at least in part, an observer-dependent property [90, 82, 61, 74]. Other proposals could also be interpreted as making a similar assumption, albeit more implicitly. For instance by choosing a good state space over which information theoretic measures can be applied [85], by building a model space over which structural causal models are well defined and meaningful for questions related to agency [106], or by requiring knowledge of other systems over which to calculate the maximum degree of causal power [111]. To fill this explanatory gap between intrinsic and extrinsic, observer-dependent proposals, following [39, 40] we adopt a more pragmatic “as-if” approach to the study of agents based on Dennett’s intentional stance. This means that an adequate account of an agent should, in our opinion, include the requirement that, *for an observer*, a candidate agent ought to look as if it’s updating beliefs like an agent, acting on the environment like an agent, and following goals like an agent. In this respect, intrinsic aspects of agency could still be deemed constitutive of a full definition, but would just further be refining this notion. In other words: no intrinsic aspect of agency should appear in systems that cannot be interpreted to be agents by an observer. Following this, we consider as-if agency a necessary, but perhaps not sufficient (if intrinsic aspects ought to be included) condition for agency.

This as-if stance has two consequences in the present work. Firstly, the fact that at least in an as-if sense, individuality and normativity are deeply connected. To see this, we first note that at a high level the relational methods used to describe individuality and the predictive methods seem to overlap considerably. For instance, the theory of individuality [85] and complete local integration [30, 90, 91], which we described as predictive and relevant for normativity, were originally introduced to define individuality. On the other hand, relational methods such as the free energy principle [61, 62, 64] and the Bayesian interpretation map [74, 44] seem to suggest the necessity of a model that is in some sense predictive for an individual. This hints at a possible *identity-normativity correspondence*, first discussed in [38] for the case of goals that are conducive to self-maintenance, and therefore



individuality. However, as mentioned in Section 4, survival need not be the only possible goal for an agent. Several frameworks discussed here already consider forms of normativity beyond survival, either by allowing to specify viability volumes in more general state-spaces [42, 52, 107, 53, 46], or by expressing measures of goal-directedness that take into account the choice of a particular (observer's) model of the agent [80, 81, 82, 83, 84, 89]. In principle, they can thus even include a characterisation of goals that work against survival, as is perhaps the case for self-destructive behaviour [39].

A more general identity-normativity correspondence captures, informally, the idea that identity, specified as a relation between agent and environment given by a particular factorisation of a full system, implies that it is possible to attribute a goal to an agent in terms of relevant properties of the full coupled agent-environment system that the agent can predict and achieve. Conversely, this also means that specific goals give rise to specific ways to carve particular identities out of a (dynamical) system. For instance, imagine we see somebody cycling. As observers, we can attribute different goals to this person, or at least be uncertain as to what exactly their goal is [112]: going to work, exercising, going for groceries, etc. In the same situation, we can also focus on a particular goal, reaching point B from point A, and ask what the right factorisation of the system is [113] for that particular goal: is the agent the human cycling? Is it the person + the bicycle? Or is it the company this person works for that makes them want to go to point B? The correspondence between individuality and goals of an agent can be many-to-many, but the point is that it is not arbitrary: we shouldn't, according to some notion of rationality [112, 40], interpret a person cycling as an agent with *any* kind of goal. Similarly, we also shouldn't interpret reaching point B from point A as the goal of *any* factorisation, e.g. a single cell in the human body, with all other cells as part of the environment. Formally, this possible correspondence has been discussed for some simple classes of systems in [74, 44, 53, 46, 114, 115]. As explained in [46, 53] this correspondence is connected to the old: "Every good regulator of a system must be a model of that system" [71], roughly claiming to provide a mathematical statement about how systems that achieve some goal must do so by modelling their environment. The mathematical grounding of this statement is however unclear, and later only partially formalised with the internal model principle in control theory [72], under rather strong assumptions as explained in [53]. As suggested by [116], this statement also forms the basis of frameworks such as the free energy principle [61, 62, 64], stating roughly the same idea but for a larger class of systems than the ones covered by [71]. We believe, however, that [46] makes this claim more formal and mathematically precise, and while limited for now to a class of systems arguably smaller than the one used by the free energy principle, within this class the claim applies to systems under quite general conditions.

The second consequence is a more coherent causal treatment of actions. In general, determining whether an agent's actions can be formally regarded as the consequence of causal interventions can be conceptually and formally challenging. This is because there is a difference between an action that *is* produced by an intervention, and an action that can be treated *as-if* it was produced by an intervention. In Pearl's structural causal models for instance, the first perspective implicitly assumes that an agent's actions are "local surgeries" performed by interventions. This means that one has to assume that, somehow, manipulations of a model via the do-operator are implemented by the agent itself [117]. Imagine a model \mathcal{M} that includes a binary random variable for a sprinkler, X_3 , which can either be ON or OFF. On this view, an action that turns on the sprinkler is an intervention "do($X_3 = \text{ON}$)", [where] we delete the equation $X_3 = f_3(X_1, U_3)$ from the theory of Eq. (4), and replace it with $X_3 = \text{ON}$ " [117]. If taken at face value, this however confuses levels of analysis, and reifies interventions as the output of an agent that can change (delete equations, and replace them) the very model it is embedded in, instead of describing them as the output of a scientist updating the model with $\mathcal{A} = \{\text{interventions on model } \mathcal{M}\}$. Conversely, the second perspective aims to describe



an agent’s policy, i.e. the map from states/beliefs to actions, as a soft or policy intervention [118], an intervention where not all causal dependencies are severed (to allow for dependencies on the agent’s previous states/beliefs and actions) [119]. This perspective simply re-describes existing knowledge of a policy using ideas from causal reasoning. It is however non-trivial: on the one hand, this less ontologically committed perspective allows us to use a series of tools from the study of causality to better interpret and understand an agent’s policies and its consequences. On the other, it also allows us to consider cases where actions cannot in fact be seen as some kind of intervention. For instance, in off-line/off-policy reinforcement learning, where observations are in part originating from systems other than the agent itself (another agent, or nature [108]), some policies may be better described without the language of interventions [120, 118].

6.3. Limitations and future work

An aspect we mostly overlooked in this work is that mathematical theories of agents ought to capture agency that can appear at different levels, or scales of abstraction [121]. This ties directly with theories of *emergence* describing intuitions and formal characterisations of how certain properties and functions can appear at some scale while not being well defined at some other. The idea of emergence appears in several fields [122] and while in some cases it could serve for a definition of individuality that is relevant for agents [86, 87, 121, 88, 123], it doesn’t necessarily speak to individuality as induced by the relation between an agent and its environment, in the way we phrased individuality in Section 3 for coupled systems. On the other hand, one could interpret emergence to be more closely aligned to the prediction-based methods used to describe normativity in Section 4, assuming that different kinds of coarse-grainings, which are necessary to define some form of emergence, can be derived from objectives that describe the predictive power of some meso or macro scale of a system (by assuming, or not, some level of supervenience) [121, 88, 123]. This idea of scales is also related to multi-scale agency [124], investigating how parts can form a whole that is an agent, and in particular how agents could be made of parts which are themselves agents (since the case where the parts are not agents is already discussed by theories of emergence). This scenario is considered, for instance, by multiplicative subagency in Cartesian frames [79] (see Section 3) and colonial individuality in the information theory of individuality [85] (see Section 4). Often, this line of work leads to discussions relevant for complex forms of agency where parts with seemingly different goals come together to act as an individual, for instance in a multi-cellular organism (which differs considerably from a group of cells on a Petri dish). This approach relies however on a pre-existing notion of an individual agent (which is the focus of this work), where the parts that form this organisation are themselves agents. We thus leave explorations of these and other related ideas to future work.

7. Conclusion

In this work, we started from the standard definition of life provided by NASA, “a self-sustaining chemical system capable of undergoing Darwinian evolution” [3] and sought to explore how the first half of this definition, “a self-sustaining chemical system”, can be abstracted and mathematically formalised in the study of life as it could be, ALife, and other related fields. We suggested the study of agents as a step in this direction, and initially operationalised agents as goal-directed systems acting in an environment. Building on a conceptual framework that identifies individuality, normativity or goal-directedness, and interactional asymmetry as key requirements of a theory of agents, we then analysed and reframed different mathematical approaches under a shared formal notation, enabling a more



direct comparison of different proposals, their strengths and possible missing parts. We then saw how adopting an as-if stance on agency resolves a few conceptual tensions and reveals a deep link between individuality and normativity: an agent's identity strongly depends on its goal-directed behaviour and vice-versa, a relation shown so far formally for a specific class of systems but likely to generalise. Within this as-if perspective, actions can be interpreted as produced by interventions that may capture the causal asymmetry between agent and environment. This perspective allows us to sidestep questions about who or what, if not the agent itself, ought to formally implement interventions in a more ontologically committed interpretation of actions in the language of causal reasoning. Overall, this synthesis exposes complementary strengths and gaps among existing frameworks, discussing mathematical and philosophical perspectives of agency and outlining a path toward an integrated, mathematically explicit theory of agents based on a common conceptual framing.

Acknowledgments

This manuscript is based on M.B.'s lectures for the 2023 CHAIN Winter School on "Minimal Cognition and Agency" in Sapporo, Japan, organised by K.S., and on a lecture at the "Life As We Don't Yet Know It," Breakthrough Discuss 2025 in Oxford, UK, in the session on "Forms of Non-Terrestrial Life," by M.B. The current material reflects a revised, improved and edited version with K.S., that makes use of feedback received in these and other events. The authors would also like to thank in particular Martin Biehl for feedback on a previous version of the manuscript. M.B. was supported by JST, Moonshot R&D, Grant Number JPMJMS2012. K.S. was supported by JST, CREST Grant Number JPMJCR21P4 and JSPS KAKENHI Grant Numbers 24H01534, Japan.

References

- [1] Christopher Langton. *Artificial Life: Proceedings Of An Interdisciplinary Workshop On The Synthesis And Simulation Of Living Systems*. Routledge, New York, 1989.
- [2] John Maynard Smith. Byte-sized evolution. *Nature*, 355(6363), 1992.
- [3] Steven A. Benner. Defining Life. *Astrobiology*, 10(10):1021–1030, December 2010.
- [4] Martin A. Nowak. *Evolutionary Dynamics: Exploring the Equations of Life*. Harvard university press, 2006.
- [5] David E. Goldberg. *The Design of Innovation*, volume 7 of *Genetic Algorithms and Evolutionary Computation*. Springer US, Boston, MA, 2002.
- [6] Christoph Adami. *Introduction to Artificial Life*. Springer Science & Business Media, 1998.
- [7] Richard E. Lenski, Charles Ofria, Robert T. Pennock, and Christoph Adami. The evolutionary origin of complex features. *Nature*, 423(6936):139–144, May 2003.
- [8] Charles Ofria and Claus O. Wilke. Avida: A Software Platform for Research in Computational Evolutionary Biology. *Artificial Life*, 10(2):191–229, 2004.
- [9] Thomas Ray. An Approach to the Synthesis of Life. In *Proceedings of Artificial Life II* (p. 371). Reading, MA: Addison-Wesley, 1991.
- [10] Bert Wang-Chak Chan. Lenia: Biology of Artificial Life. *Complex Systems*, 28(3), 2019.
- [11] Erwan Plantec, Gautier Hamon, Mayalen Etcheverry, Bert Wang-Chak Chan, Pierre-Yves Oudeyer, and Clément Moulin-Frier. Flow-Lenia: Emergent Evolutionary Dynamics in Mass Conservative Continuous Cellular Automata. *Artificial Life*, 31(2):228–248, May 2025.



- [12] Hiroki Sayama and Christopher L. Nehaniv. Self-Reproduction and Evolution in Cellular Automata: 25 Years After Evoloops. *Artificial Life*, 31(1):81–95, February 2024.
- [13] Carlo Rovelli. Agency in Physics, July 2020.
- [14] Manuel Baltieri, Christopher L. Buckley, and Jelle Bruineberg. Predictions in the eye of the beholder: An active inference account of Watt governors. In *The 2020 Conference on Artificial Life*, pages 121–129, Online, 2020. MIT Press.
- [15] Manuel Baltieri and Ryota Kanai. Steam-engine naturalism, October 2025.
- [16] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall Press, 3rd edition, 2009.
- [17] Sameera Abar, Georgios K. Theodoropoulos, Pierre Lemarinier, and Gregory M.P. O'Hare. Agent Based Modelling and Simulation tools: A review of the state-of-art software. *Computer Science Review*, 24:13–33, May 2017.
- [18] Deepak Bhaskar Acharya, Karthigeyan Kuppan, and B. Divya. Agentic AI: Autonomous Intelligence for Complex Goals—A Comprehensive Survey. *IEEE Access*, 13:18912–18936, 2025.
- [19] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning Series. The MIT Press, Cambridge, Massachusetts, second edition edition, 2018.
- [20] Karl J. Åström and Richard M. Murray. *Feedback Systems: An Introduction for Scientists and Engineers*. Princeton Univ. Press, Princeton, NJ, 2008.
- [21] Gertrude Elizabeth Margaret Anscombe. *Intention*. Harvard University Press, 1963.
- [22] Donald Davidson. Actions, reasons, and causes. *The Journal of Philosophy*, 60(23):685–700, 1963.
- [23] Markus Schlosser. Agency. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2019 edition, 2019.
- [24] Luca Ferrero. An introduction to the philosophy of agency. In *The Routledge Handbook of Philosophy of Agency*, pages 1–18. Routledge, 2022.
- [25] Patrick Haggard. Sense of agency in the human brain. *Nature Reviews Neuroscience*, 18(4):196–207, April 2017.
- [26] Keisuke Suzuki, Peter Lush, Anil K. Seth, and Warrick Roseboom. Intentional Binding Without Intentional Action. *Psychological Science*, 30(6):842–853, June 2019.
- [27] Wen Wen and Hiroshi Imamizu. The sense of agency in perception, behaviour and human-machine interactions. *Nature Reviews Psychology*, 1(4):211–222, April 2022.
- [28] Michael Tomasello. *The Evolution of Agency: Behavioral Organization from Lizards to Humans*. MIT Press, 2022.
- [29] Humberto R Maturana and Francisco J Varela. *Autopoiesis and Cognition: The Realization of the Living*. Springer Science & Business Media, 1980.
- [30] Martin Biehl, Takashi Ikegami, and Daniel Polani. Towards information based spatiotemporal patterns as a foundation for agent representation in dynamical systems. In *ALIFE 2016, the Fifteenth International Conference on the Synthesis and Simulation of Living Systems*, pages 722–729. MIT Press, July 2016.
- [31] Phillip Ball. Organisms as agents of evolution. Technical report, John Templeton Foundation: West Conshohocken, PA, USA, 2023.



- [32] Mustafa Emirbayer and Ann Mische. What Is Agency? *American Journal of Sociology*, 103(4):962–1023, January 1998.
- [33] Stephen A. Ross. The Economic Theory of Agency: The Principal’s Problem. *The American Economic Review*, 63(2):134–139, 1973.
- [34] Roderick Munday. *Agency: Law and Principles*. OUP Oxford, 2010.
- [35] Manuel Baltieri, Hiroyuki Iizuka, Olaf Witkowski, Lana Sinapayen, and Keisuke Suzuki. Hybrid Life: Integrating biological, artificial, and cognitive systems. *WIREs Cognitive Science*, n/a(n/a):e1662, 2023.
- [36] Pamela Lyon. Of what is “minimal cognition” the half-baked version? *Adaptive Behavior*, 28(6):407–424, December 2020.
- [37] Randall D Beer. Some historical context for minimal cognition. *Adaptive Behavior*, 29(1):89–92, February 2021.
- [38] Xabier E. Barandiaran, Ezequiel Alejandro Di Paolo, and Marieke Rohde. Defining Agency: Individuality, Normativity, Asymmetry, and Spatio-temporality in Action. *Adaptive Behavior*, 17(5):367–386, October 2009.
- [39] Simon McGregor. A More Basic Version of Agency? As If! In *Lecture Notes in Computer Science*, pages 183–194, 2016.
- [40] Simon McGregor. The Bayesian stance: Equations for ‘as-if’ sensorimotor agency. *Adaptive Behavior*, 25(2):72–82, 2017.
- [41] Judea Pearl. *Causality - Models, Reasoning, and Inference*. Cambridge University Press, 2009.
- [42] Randall D. Beer. A dynamical systems perspective on agent-environment interaction. *Artificial intelligence*, 72(1-2):173–215, 1995.
- [43] Randall D. Beer. The dynamics of adaptive behavior: A research program. *Robotics and Autonomous Systems*, 20(2-4):257–289, June 1997.
- [44] Martin Biehl and Nathaniel Virgo. Interpreting systems as solving POMDPs: A step towards a formal understanding of agency. In *International Workshop on Active Inference*, pages 16–31. Springer, 2022.
- [45] Nathaniel Virgo. Unifilar machines and the adjoint structure of Bayesian models. In *Electronic Proceedings in Theoretical Computer Science*, volume 397, pages 299–317, 2023.
- [46] Nathaniel Virgo, Martin Biehl, Manuel Baltieri, and Matteo Capucci. A “good regulator theorem” for embodied agents, August 2025.
- [47] Fernando Rosas, Alexander Boyd, and Manuel Baltieri. AI in a vat: Fundamental limits of efficient world modelling for agent sandboxing and interpretability. In *Proceedings of the Second Reinforcement Learning Conference*, April 2025.
- [48] Alexander Boyd, Franz Nowak, David Hyland, Manuel Baltieri, and Fernando E. Rosas. From monoliths to modules: Decomposing transducers for efficient world modelling, December 2025.
- [49] David Jaz Myers. *Categorical Systems Theory*. 2021.
- [50] Matteo Capucci. Notes on categorical systems theory. Technical report, 2024.
- [51] Francis Heylighen, Shima Beigi, and Tomas Veloz. Chemical Organization Theory as a General Modeling Framework for Self-Sustaining Systems. *Systems*, 12(4):111, April 2024.



- [52] Jean-Pierre Aubin, Alexandre M. Bayen, and Patrick Saint-Pierre. *Viability Theory: New Directions*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.
- [53] Manuel Baltieri, Martin Biehl, Matteo Capucci, and Nathaniel Virgo. A Bayesian Interpretation of the Internal Model Principle, March 2025.
- [54] Randall D. Beer. Autopoiesis and Cognition in the Game of Life. *Artificial Life*, 10(3):309–326, June 2004.
- [55] Randall D. Beer. The cognitive domain of a glider in the game of life. *Artificial Life*, 20(2):183–206, 2014.
- [56] Randall D. Beer. Characterizing Autopoiesis in the Game of Life. *Artificial Life*, 21(1):1–19, February 2015.
- [57] Randall D. Beer. Bittorio revisited: Structural coupling in the Game of Life. *Adaptive Behavior*, 28(4):197–212, August 2020.
- [58] Randall D. Beer. An integrated perspective on the constitutive and interactive dimensions of autonomy. In *Artificial Life Conference Proceedings 32*, pages 202–209. MIT Press, 2020.
- [59] Randall D. Beer. An investigation into the origin of autopoiesis. *Artificial Life*, 26(1):5–22, 2020.
- [60] Martin Biehl and Nathaniel Virgo. The game of life in a glider’s frame of reference. 2023.
- [61] Karl J. Friston. A free energy principle for a particular physics, June 2019.
- [62] Thomas Parr, Giovanni Pezzulo, and Karl J. Friston. *Active Inference: The Free Energy Principle in Mind, Brain, and Behavior*. MIT Press, March 2022.
- [63] Chris Fields, Karl Friston, James F. Glazebrook, and Michael Levin. A free energy principle for generic quantum systems. *Progress in Biophysics and Molecular Biology*, 173:36–59, September 2022.
- [64] Karl J. Friston, Lancelot Da Costa, Noor Sajid, Conor Heins, Kai Ueltzhöffer, Grigorios A. Pavliotis, and Thomas Parr. The free energy principle made simpler but not too simple. *Physics Reports*, 1024:1–29, June 2023.
- [65] Naftali Tishby and Daniel Polani. Information Theory of Decisions and Actions. In Vassilis Cutsuridis, Amir Hussain, and John G. Taylor, editors, *Perception-Action Cycle*, pages 601–636. Springer New York, New York, NY, 2011.
- [66] Jelle Bruineberg, Krzysztof Dołęga, Joe Dewhurst, and Manuel Baltieri. The Emperor’s New Markov Blankets. *Behavioral and Brain Sciences*, 45:e183, 2022.
- [67] Jelle Bruineberg, Krzysztof Dołęga, Joe Dewhurst, and Manuel Baltieri. The Emperor Is Naked: Replies to commentaries on the target article. *Behavioral and Brain Sciences*, 45, 2022.
- [68] Martin Biehl, Felix A. Pollock, and Ryota Kanai. A Technical Critique of Some Parts of the Free Energy Principle. *Entropy*, 23(3):293, February 2021.
- [69] Nathaniel Virgo, Fernando E. Rosas, and Martin Biehl. Embracing sensorimotor history: Time-synchronous and time-unrolled Markov blankets in the free-energy principle. *Behavioral and Brain Sciences*, 45:e215, 2022.
- [70] Fernando E. Rosas, Pedro A.M. Mediano, Martin Biehl, Shamil Chandaria, and Daniel Polani. Causal blankets: Theory and algorithmic framework. In *International Workshop on Active Inference*, pages 187–198. Springer, 2020.



- [71] Roger C. Conant and W. Ross Ashby. Every good regulator of a system must be a model of that system. *International Journal of Systems Science*, 1(2):89–97, 1970.
- [72] W. M. Wonham. Towards an Abstract Internal Model Principle. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-6(11):735–740, November 1976.
- [73] Lancelot Da Costa, Karl J. Friston, Conor Heins, and Grigoris A. Pavlou. Bayesian mechanics for stationary processes. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 477(2256):20210518, December 2021.
- [74] Nathaniel Virgo, Martin Biehl, and Simon McGregor. Interpreting Dynamical Systems as Bayesian Reasoners. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, December 2021.
- [75] Miguel Aguilera, Beren Millidge, Alexander Tschantz, and Christopher L. Buckley. How particular is the physics of the free energy principle? *Physics of Life Reviews*, 40:24–50, March 2022.
- [76] Ezequiel Alejandro Di Paolo, Evan Thompson, and Randall D Beer. Laying down a forking path: Tensions between enaction and the free energy principle. *Philosophy and the Mind Sciences*, 2022.
- [77] Vicente Raja, Dinesh Valluri, Edward Baggs, Anthony Chemero, and Michael L. Anderson. The Markov blanket trick: On the scope of the free energy principle and active inference. *Physics of Life Reviews*, 39:49–72, December 2021.
- [78] Daniel C. Dennett. *The Intentional Stance*. A Bradford Book. MIT Press, 1989.
- [79] Scott Garrabrant, Daniel A. Herrmann, and Josiah Lopez-Wild. Cartesian Frames, September 2021.
- [80] Karl J. Friston, Christopher Thornton, and Andy Clark. Free-energy minimization and the dark-room problem. *Frontiers in Psychology*, 3(MAY), 2012.
- [81] Andrew Barto, Marco Mirolli, and Gianluca Baldassarre. Novelty or surprise? *Frontiers in Psychology*, 4(December):907, 2013.
- [82] Laurent Orseau, Simon McGregor McGill, and Shane Legg. Agents and devices: A relative definition of agency, 2018.
- [83] Manuel Baltieri and Christopher L. Buckley. The dark room problem in predictive processing and active inference, a legacy of cognitivism? In *Artificial Life Conference Proceedings 31*, pages 40–47. MIT Press, 2019.
- [84] Manuel Baltieri. *Active Inference: Building a New Bridge between Control Theory and Embodied Cognitive Science*. PhD thesis, University of Sussex, 2019.
- [85] David Krakauer, Nils Bertschinger, Eckehard Olbrich, Jessica C. Flack, and Nihat Ay. The information theory of individuality. *Theory in Biosciences*, 139(2):209–223, June 2020.
- [86] Nils Bertschinger, Eckehard Olbrich, Nihat Ay, and Jürgen Jost. Autonomy: An information theoretic perspective. *Biosystems*, 91(2):331–345, February 2008.
- [87] Anil K. Seth. Measuring Autonomy and Emergence via Granger Causality. *Artificial Life*, 16(2):179–196, April 2010.
- [88] Lionel Barnett and Anil K. Seth. Dynamical independence: Discovering emergent macroscopic processes in complex dynamical systems. *Physical Review E*, 108(1):014304, July 2023.
- [89] Matt MacDermott, James Fox, Francesco Belardinelli, and Tom Everitt. Measuring Goal-Directedness. In *Advances in Neural Information Processing Systems*, volume 27, 2024.



- [90] Martin Biehl. *Formal Approaches to a Definition of Agents*. PhD thesis, University of Hertfordshire, 2018.
- [91] Martin Biehl, Takashi Ikegami, and Daniel Polani. Specific and Complete Local Integration of Patterns in Bayesian Networks. *Entropy*, 19(5):230, May 2017.
- [92] R Rosen. On anticipatory systems: I. When can a system contain a predictive model of another? *Journal of Social and Biological Systems*, 1(2):155–162, April 1978.
- [93] R Rosen. On anticipatory systems: II. The nature of the modelling relation between systems. *Journal of Social and Biological Systems*, 1(2):163–180, April 1978.
- [94] Aloisius H. Louie. Mathematical Foundations of Anticipatory Systems. In Roberto Poli, editor, *Handbook of Anticipation*, pages 1–29. Springer International Publishing, Cham, 2017.
- [95] Aloisius H. Louie. *More Than Life Itself: A Synthetic Continuation in Relational Biology*. DE GRUYTER, December 2009.
- [96] Fosco Loregian. Rosen’s no-go theorem for regular categories, May 2021.
- [97] Elias Zafiris. Rosen’s modelling relations via categorical adjunctions. *International Journal of General Systems*, 41(5):439–474, July 2012.
- [98] Larissa Albantakis, William Marshall, Erik Hoel, and Giulio Tononi. What Caused What? A Quantitative Account of Actual Causation Using Dynamical Causal Networks. *Entropy*, 21(5):459, May 2019.
- [99] Masafumi Oizumi, Larissa Albantakis, and Giulio Tononi. From the phenomenology to the mechanisms of consciousness: Integrated information theory 3.0. *PLoS Computational Biology*, 10(5):e1003588, 2014.
- [100] Larissa Albantakis, Leonardo Barbosa, Graham Findlay, Matteo Grasso, Andrew M. Haun, William Marshall, William G. P. Mayner, Alireza Zaeemzadeh, Melanie Boly, Bjørn E. Juel, Shuntaro Sasai, Keiko Fujii, Isaac David, Jeremiah Hendren, Jonathan P. Lang, and Giulio Tononi. Integrated information theory (IIT) 4.0: Formulating the properties of phenomenal existence in physical terms. *PLOS Computational Biology*, 19(10):e1011465, October 2023.
- [101] Larissa Albantakis, Francesco Massari, Maggie Beheler-Amass, and Giulio Tononi. A Macro Agent and Its Actions. In *Top-down Causation and Emergence*, pages 135–155. Springer, 2021.
- [102] Bjørn Erik Juel, Renzo Comolatti, Giulio Tononi, and Larissa Albantakis. When is an action caused from within? Quantifying the causal chain leading to actions in simulated agents. In *Artificial Life Conference Proceedings 31*, 2019.
- [103] Miguel Aguilera, Carlos Alquézar, and Manuel G. Bedia. Agency and Integrated Information in a Minimal Sensorimotor Model. In *The 2018 Conference on Artificial Life*, pages 396–403, Tokyo, Japan, 2018. MIT Press.
- [104] Miguel Aguilera and Ezequiel Alejandro Di Paolo. Integrated Information and Autonomy in the Thermodynamic Limit. In *The 2018 Conference on Artificial Life*, pages 113–120, Tokyo, Japan, 2018. MIT Press.
- [105] Miguel Aguilera and Ezequiel Alejandro Di Paolo. Integrated information in the thermodynamic limit. *Neural Networks*, 114:136–146, June 2019.
- [106] Zachary Kenton, Ramana Kumar, Sebastian Farquhar, Jonathan Richens, Matt MacDermott, and Tom Everitt. Discovering agents. *Artificial Intelligence*, 322:103963, September 2023.



- [107] Artemy Kolchinsky and David H. Wolpert. Semantic information, autonomous agency and non-equilibrium statistical physics. *Interface Focus*, 8(6):20180041, December 2018.
- [108] Filippo Torresan and Manuel Baltieri. Disentangled representations for causal cognition. *Physics of Life Reviews*, 51:343–381, December 2024.
- [109] Jonathan Richens, David Abel, Alexis Bellot, and Tom Everitt. General agents contain world models, 2025.
- [110] Diane W. Davidson. Ecological Studies of Neotropical Ant Gardens. *Ecology*, 69(4):1138–1152, 1988.
- [111] H.H. Morch. Is Consciousness Intrinsic? A Problem for the Integrated Information Theory. *Journal of Consciousness Studies*, 26(1-2):133–162, January 2019.
- [112] Chris L. Baker, Rebecca Saxe, and Joshua B. Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, December 2009.
- [113] Andy Clark and David Chalmers. The extended mind. *Analysis*, 58(1):7–19, 1998.
- [114] Simon McGregor, timorl, and Nathaniel Virgo. Formalising the intentional stance 1: Attributing goals and beliefs to stochastic processes, January 2025.
- [115] Simon McGregor, timorl, and Nathaniel Virgo. Formalising the intentional stance 2: A coinductive approach, January 2025.
- [116] Anil K Seth. The Cybernetic Bayesian Brain. In Wanja Wiese and Thomas K Metzinger, editors, *Open MIND*, pages 9–24. Frankfurt am Main, Germany: MIND Group, 2015.
- [117] Judea Pearl. Causation, Action, and Counterfactuals. In Maria Luisa Dalla Chiara, Kees Doets, Daniele Mundici, and Johan Van Benthem, editors, *Logic and Scientific Methods*, pages 355–375. Springer Netherlands, Dordrecht, 1997.
- [118] Elias Bareinboim, Junzhe Zhang, Sanghack Lee, and Snu Ac Kr. An Introduction to Causal Reinforcement Learning, 2021.
- [119] Zhihong Deng, Jing Jiang, Guodong Long, and Chengqi Zhang. Causal Reinforcement Learning: A Survey, November 2023.
- [120] Oliver Schulte and Pascal Poupart. When Should Reinforcement Learning Use Causal Reasoning? *Transactions on Machine Learning Research*, 2025.
- [121] Erik P. Hoel. Agent Above, Atom Below: How Agents Causally Emerge from Their Underlying Microphysics. In Anthony Aguirre, Brendan Foster, and Zeeya Merali, editors, *Wandering Towards a Goal*, pages 63–76. Springer International Publishing, Cham, 2018.
- [122] Ross H. McKenzie. Emergence: From physics to biology, sociology, and computer science, August 2025.
- [123] Fernando E. Rosas, Bernhard C. Geiger, Andrea I. Luppi, Anil K. Seth, Daniel Polani, Michael Gastpar, and Pedro A. M. Mediano. Software in the natural world: A computational approach to hierarchical emergence, June 2024.
- [124] Michael Levin. The Computational Boundary of a “Self”: Developmental Bioelectricity Drives Multicellularity and Scale-Free Cognition. *Frontiers in Psychology*, 10, December 2019.