

MICHAEL TIMOTHY BENNETT

TECHNICAL APPENDICES

VERSION 3.0

Math and Experiments

THIS CHAPTER covers all the math and experiments published in ¹ and refined the following year in ². Some improvements have been made to the proofs, and many details and examples have been fleshed out. It is effectively a reference manual for the rest of the thesis. The rest of the thesis provides context and verbal argument. I give the definitions for my formalism of everything, a list of proofs concerning those definitions, and finally experimental results.

¹ Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a; Michael Timothy Bennett. Emergent causality and the foundation of consciousness. In *Artificial General Intelligence*. Springer Nature, 2023b; and Michael Timothy Bennett. On the computation of meaning, language models and incomprehensible horrors. In *Artificial General Intelligence*. Springer Nature, 2023c

² Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a; and Michael Timothy Bennett. Is complexity an illusion? In *Artificial General Intelligence*. Springer, 2024b

Foundational Definitions

DEFINITIONS 1 TO 4 are taken from ³:

Definition 1 (environment)

- We assume a set Φ whose elements we call **states**.
- A **declarative program** is $f \subseteq \Phi$, and we write P for the set of all declarative programs (the powerset of Φ).
- By a **truth** or **fact** about a state ϕ , we mean $f \in P$ such that $\phi \in f$.
- By an **aspect of a state** ϕ we mean a set l of facts about ϕ s.t. $\phi \in \bigcap l$. By an **aspect of the environment** we mean an aspect l of any state, s.t. $\bigcap l \neq \emptyset$. We say an aspect of the environment is **realised**⁴ by state ϕ if it is an aspect of ϕ .

Definition 2 (abstraction layer)

By abstraction layer⁵ we mean:

- We single out a subset $\mathbf{v} \subseteq P$ which we call **the vocabulary** of an abstraction layer. The vocabulary is finite unless explicitly stated otherwise. If $\mathbf{v} = P$, then we say that there is no abstraction.
- $L_{\mathbf{v}} = \{l \subseteq \mathbf{v} : \bigcap l \neq \emptyset\}$ is a set of aspects in \mathbf{v} . We call $L_{\mathbf{v}}$ a formal language, and $l \in L_{\mathbf{v}}$ a **statement**.
- We say a statement is **true** given a state iff it is an aspect realised by that state.
- A **completion** of a statement x is a statement y which is a superset of x . If y is true, then x is true.
- The **extension of a statement** $x \in L_{\mathbf{v}}$ is $E_x = \{y \in L_{\mathbf{v}} : x \subseteq y\}$. E_x is the set of all completions of x .
- The **extension of a set of statements** $X \subseteq L_{\mathbf{v}}$ is $E_X = \bigcup_{x \in X} E_x$.
- We say x and y are **equivalent** iff $E_x = E_y$.

Definition 3 (v-task)

For a chosen \mathbf{v} , a task α is a pair $\langle I_{\alpha}, O_{\alpha} \rangle$ where⁶:

- $I_{\alpha} \subseteq L_{\mathbf{v}}$ is a set whose elements we call **inputs** of α .
- $O_{\alpha} \subseteq E_{I_{\alpha}}$ is a set whose elements we call **correct outputs** of α .

I_{α} has the extension $E_{I_{\alpha}}$ we call **outputs**, and O_{α} are outputs deemed correct. $\Gamma_{\mathbf{v}}$ is the set of **all tasks** given \mathbf{v} .

³ Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a

⁴ Realised meaning it is made real, or brought into existence.

⁵ (NOTATION) E with a subscript is the extension of the subscript. For example, E_l is the extension of l .

(INTUITIVE SUMMARY) $L_{\mathbf{v}}$ is everything which can be realised in this abstraction layer. The extension E_x of a statement x is the set of all statements whose existence implies x , and so it is like the sub-table of x 's truth table for which x is true.

⁶ (NOTATION) If $\omega \in \Gamma_{\mathbf{v}}$, then we will use subscript ω to signify parts of ω , meaning one should assume $\omega = \langle I_{\omega}, O_{\omega} \rangle$ even if that isn't written.

(INTUITIVE SUMMARY) To reiterate and summarise the above:

- An **input** is a possibly incomplete description of a world.
- An **output** is a completion of an input [see def. 2].
- A **correct output** is a correct completion of an input.

(GENERATIONAL HIERARCHY) A \mathbf{v} -task α is a **child** of \mathbf{v} -task ω if $I_\alpha \subset I_\omega$ and $O_\alpha \subseteq O_\omega$. This is written as $\alpha \sqsubset \omega$. If $\alpha \sqsubset \omega$ then ω is then a **parent** of α . \sqsubset implies a “lattice” or generational hierarchy of tasks. Formally, the level of a task α in this hierarchy is the largest k such there is a sequence $\langle \alpha_0, \alpha_1, \dots, \alpha_k \rangle$ of k tasks such that $\alpha_0 = \alpha$ and $\alpha_i \sqsubset \alpha_{i+1}$ for all $i \in (0, k)$. A child is always “lower level” than its parents⁷.

⁷ (FURTHER INTUITIVE SUMMARY) A \mathbf{v} -task is a formal, **behavioural** description of an aspect of the environment. For example, a self-organising biological system could be described as a task α enumerating all behaviour in which it remains alive. It begins alive in circumstances given by inputs I_α , and remains alive in circumstances given by outputs O_α , and is dead in circumstances given by $E_{I_\alpha} - O_\alpha$. Likewise, we could describe the game chess played from the perspective of white. We could say Φ contains a state corresponding to each and every move of each and every possible game of chess, I_α contains every possible sequence of moves in which the game has not ended and it remains *possible* for white to win, and O_α contains every possible sequence ending in a move that means white *has* won. Tasks are *behavioural* descriptions of systems in the philosophical sense of the word, and we will next relate these ideas to machine functionalism.

Learning and Inference Definitions

MACHINE LEARNING is typically divided into learning and inference. To learn a task, an organism learns a policy (strictly speaking, an organism *is* a policy). By constraining how we complete inputs, a policy allows us to do abductive inference (much as a constraint in a SAT solver would). Inference requires a **policy** and learning a policy requires a **proxy**, the definitions of which follow.

Definition 4 (inference)

- A **v-task policy** is a statement $\pi \in L_v$. It constrains how we complete inputs.
- π is a **correct policy** iff the correct outputs O_α of α are exactly the completions π' of π such that π' is also a completion of an input.
- The set of all correct policies for a task α is denoted Π_α .⁸

Assume v-task ω and a policy $\pi \in L_v$. Inference⁹ proceeds as follows:

1. we are presented with an input $i \in I_\omega$, and
2. we must select an output $e \in E_i \cap E_\pi$.
3. If $e \in O_\omega$, then e is correct and the task “complete”. $\pi \in \Pi_\omega$ implies $e \in O_\omega$, but $e \in O_\omega$ doesn’t imply $\pi \in \Pi_\omega$ (an incorrect policy can imply a correct output).

Definition 5 (learning) Learning is a collection of definitions that describe the process by which a policy is constructed by any system¹⁰.

- A **proxy** $<$ is a binary relation on statements, and the set of all proxies is Q .
- $<_w$ is the **weakness proxy**¹¹. For statements l_1, l_2 we have $l_1 <_w l_2$ iff $|E_{l_1}| < |E_{l_2}|$.
- $<_d$ is the **description length or simplicity proxy**¹². We have $l_1 <_d l_2$ iff $|l_1| > |l_2|$.

(GENERALISATION) A statement l **generalises** to a v-task α iff $l \in \Pi_\alpha$. We speak of **learning** ω from α iff, given a proxy $<$, $\pi \in \Pi_\alpha$ maximises $<$ relative to all other policies in Π_α , and $\pi \in \Pi_\omega$.

(PROBABILITY OF GENERALISATION) We assume a uniform distribution over Γ_v . If l_1 and l_2 are policies, we say it is less probable that l_1 generalizes than that l_2 generalizes, written $l_1 <_g l_2$, iff, when a task α is chosen at random from Γ_v (using a uniform distribution) then the probability that l_1 generalizes to α is less than the probability that l_2 generalizes to α .

⁸ To repeat the above definition in set builder notation:

$$\Pi_\alpha = \{\pi \in L_v : E_{l_\alpha} \cap E_\pi = O_\alpha\}$$

⁹ (INTUITIVE SUMMARY) To reiterate and summarise the above:

- A **policy** constrains how we complete inputs.
- A **correct policy** is one that constrains us to correct outputs.

¹⁰ (INTUITIVE SUMMARY) Learning is an activity undertaken by an adaptive system, and a task has been **learned** by a system that embodies a correct policy. Humans typically learn from **examples**. An example of a task is a correct output and input.

¹¹ By the weakness of a statement, we mean the cardinality of its extension. By the weakness of an extension we mean its cardinality.

¹² When we speak of simplicity with regards to a statement, we mean its cardinality. The complexity of an extension is the **simplest** statement of which it is an extension.

(SAMPLE EFFICIENCY) Suppose¹³ \mathbf{app} is the set of *all pairs of policies*. In the following formula, when we write “ $(l_1 <_g l_2)$ ”, we mean the number 1 if $l_1 <_g l_2$ or 0 otherwise. For example, if $x <_g y$ and $y \not<_g z$ then the expression $(x <_g y) + (x <_g y) + (y <_g z)$ would equal $1 + 1 + 0 = 2$. Proxy $<_a$ is more sample efficient than $<_b$ iff

$$\left(\sum_{(l_1, l_2) \in \mathbf{app}} |(l_1 <_g l_2) - (l_1 <_a l_2)| - |(l_1 <_g l_2) - (l_1 <_b l_2)| \right) < 0$$

(OPTIMAL PROXY) There is no proxy more sample efficient than $<_w$, so we call $<_w$ optimal. This formalises the idea that “explanations should be no more specific than necessary” (see Bennett’s Razor).

¹³ (FURTHER INTUITIVE SUMMARY) A collection of examples is a child task, so learning is an attempt to generalise from a child, to one of its parents. The lower level the child from which an agent generalises to parent, the ‘faster’ it learns, the more sample efficient the proxy.

Philosophical and Biological Analogues

Definition 6 (λ -tasks)

$\Gamma = \bigcup_{v \in V} \Gamma_v$ is the set of all tasks in all vocabularies. An uninstantiated or “ λ -task” is a function $\lambda : V \rightarrow \Gamma$ that takes a vocabulary and returns a particular task in that vocabulary, such that if $\lambda(v) = \alpha$ and $\lambda(w) = \omega$, and $v \subset w$, then $\Pi_\alpha \subset \Pi_\omega$. Λ is the set of all λ -tasks. A λ -task is a v -task that is yet to be instantiated in a particular vocabulary (it lets us represent a version of the same task in different vocabularies).

Definition 7 (utility of intelligence)

Every task $\gamma \in \Gamma$ has a “utility of intelligence”¹⁴ computed as $\epsilon : \Gamma \rightarrow \mathbb{N}$ such that $\epsilon(\gamma) = \max_{m \in \Pi_\gamma} (|E_m| - |O_\gamma|)$. Maximisation of utility means

maximising $\bullet \stackrel{\epsilon}{<} \bullet$ that returns true iff $\epsilon(\alpha) < \epsilon(\omega)$.

¹⁴ Assuming we accept that intelligence is the ability to generalise, then we can measure the utility of selecting policies in accord with Bennett’s Razor by measuring the weakness of the weakest policy for a task. Tasks with weaker policies make more use of intelligence.

Definition 8 (organism)

We can describe the circumstances of an organism¹⁵ \mathbf{o} as $\langle \mathbf{v}_\mathbf{o}, \mu_\mathbf{o}, \mathbf{p}_\mathbf{o}, <_\mathbf{o} \rangle$ where:

- $O_{\mu_\mathbf{o}}$ contains every output which qualifies as “fit” according to natural selection.
- $\mathbf{p}_\mathbf{o}$ is the set of policies an organism knows, s.t. $\mathbf{p}_\mathbf{o} \subset \mathbf{p}_{n.s.} \cup \mathbf{p}_{h_{< t_\mathbf{o}}}$ and:
 - $\mathbf{p}_{n.s.} \subset L_{\mathbf{v}_\mathbf{o}}$ is **reflexes** hard coded from birth by natural selection.
 - $\mathbf{p}_{h_{< t_\mathbf{o}}} = \bigcup_{\zeta \in h_{< t_\mathbf{o}}} \Pi_\zeta$ is the set of policies it is possible to **learn** from a history of past interactions represented by a task $h_{< t_\mathbf{o}}$.
 - If $\Pi_{h_{< t_\mathbf{o}}} \not\subset (\mathbf{p}_\mathbf{o} - \mathbf{p}_{n.s.})$ then the organism has **selective memory**. It can “forget” outputs, possibly to productive ends if they contradict otherwise good policies.
- $<_\mathbf{o}$ is a binary relation over $\Gamma_{\mathbf{v}_\mathbf{o}}$ we call **preferences**.

¹⁵ (INTUITIVE SUMMARY) Strictly speaking an organism \mathbf{o} would be a policy, but we can describe the circumstances of its existence as a task μ that describes all “fit” behaviour for that organism. We can also identify policies the organism “knows”, because these are implied by the policy that is the organism. Likewise, we can represent lossy memory by having the organism “know” fewer policies than are implied by its history of interactions. Finally, preferences are the particular “protosymbol” the organism will use to “interpret” an input in later definitions.

Semiosis

Original definitions are to be found in ¹⁶.

Definition 9 (protosymbol system)

Assume an organism \mathbf{o} . For each policy $p \in \mathbf{p}_{\mathbf{o}}$ there exists a set $\mathbf{s}_p = \{\alpha \in \Gamma_{\mathbf{v}_{\mathbf{o}}} : p \in \Pi_{\alpha}\}$ of all tasks for which p is a correct policy. The union of all such sets is

$$\mathbf{s}_{\mathbf{o}} = \bigcup_{p \in \mathbf{p}_{\mathbf{o}}} \{\alpha \in \Gamma_{\mathbf{v}_{\mathbf{o}}} : p \in \Pi_{\alpha}\}$$

We call $\mathbf{s}_{\mathbf{o}}$ a “protosymbol system”. A \mathbf{v} -task $\alpha \in \mathbf{s}_{\mathbf{o}}$ is called a “protosymbol”, and is “more abstract” if it is higher in the generational hierarchy.

¹⁶ Michael Timothy Bennett. On the computation of meaning, language models and incomprehensible horrors. In *Artificial General Intelligence*. Springer Nature, 2023c

Definition 10 (interpretation)

Interpretation is inference, with the additional step of choosing a policy according to preference. Interpretation is an activity undertaken by an organism \mathbf{o} . It proceeds as follows:

1. Assume true input $i \in L_{\mathbf{v}_{\mathbf{o}}}$ (meaning $i = i_t$ in an EGRL system at time t).
2. We say that i **signifies** a protosymbol $\alpha \in \mathbf{s}_{\mathbf{o}}$ if $i \in I_{\alpha}$.
3. $\mathbf{s}_{\mathbf{o}}^s = \{\alpha \in \mathbf{s}_{\mathbf{o}} : i \in I_{\alpha}\}$ is the set of all protosymbols which i signifies.
4. If $\mathbf{s}_{\mathbf{o}}^s \neq \emptyset$ then i **means something** to the organism in the sense that there is “value” ascribed to symbols in $\mathbf{s}_{\mathbf{o}}^s$ compelling the organism to act.
5. If i means something, then \mathbf{o} chooses $\alpha \in \mathbf{s}_{\mathbf{o}}^s$ that maximises its preferences $<_{\mathbf{o}}$.
6. The organism then infers an output $o \in E_s \cap E_{\Pi_{\alpha}}$.

Definition 11 (affect)

Suppose we have two organisms, \mathbf{a} (Alice) and \mathbf{b} (Bob). Suppose \mathbf{a} interprets $i \in L_{\mathbf{v}_{\mathbf{o}}}$ as an output o , then:

- a **statement** $v \subset i$ affects \mathbf{a} if \mathbf{a} would have interpreted $e = i - v$ as a different output $g \neq o$.
- an **organism** \mathbf{b} has affected \mathbf{a} by making an output k if, as a consequence of k , there exists $v \subset s$ which affects \mathbf{a} .

Definition 12 (multiscale competency architecture)

The **causal chain** of the multiscale competency architecture (MCA) is a sequence of uninstantiated tasks $\langle \lambda^0, \lambda^1 \dots \lambda^n \rangle$ s.t $\lambda^{i+1} \sqsubset \lambda^i$. Assume a **function** $\mathbf{j} : 2^{L_{\mathbf{v}}} \rightarrow 2^P$ that takes the extension E_{π} of a policy $\pi \in L_{\mathbf{v}}$, and returns a **vocabulary** $\mathbf{v}' = \{f \in P : \exists o \in E_{\pi}(\cap o = f)\}$. The **state** of the

MCA is a sequence of policies $\langle \pi^0, \pi^1 \dots \pi^n \rangle$ and a sequence of vocabularies $\langle \mathbf{v}^0, \mathbf{v}^1 \dots \mathbf{v}^n \rangle$ such that $\mathbf{v}^{i+1} = \mathbf{f}(E_{\pi^i})$ and $\pi^i \in \Pi_{\lambda^i(\mathbf{v}^i)}$. The MCA is **over-constrained** where there exists $i < n$ s.t. and $\Pi_{\lambda^i(\mathbf{v}^i)} = \emptyset$, and **multiscale causal learning** (MCL) occurs when the MCA is under-constrained and the proxy for learning is weakness. Note that the vocabulary is different at each level in the MCA, which means each has its own generational hierarchy of tasks. By a higher level of abstraction in the MCA, we mean a later task in the causal chain.

Causal definitions

ORIGINAL DEFINITIONS are to be found in ¹⁷.

Definition 13 (intervention)

Intuitively, if *int* and *obs* are “events” which have happened, then we say that *int* has **caused** *obs* if *obs* would not have happened in the absence of *int* (counterfactual). In our formalism, an **event** is a statement in L_v , and an event **happens** or is **observed** iff it is a true statement. If $obs \in L_v$ is sensorimotor activity we interpret as an “observed event”, and $int \in L_v$ is in **intervention** (by an organism or other agency, in the sense described by Pearl ¹⁸) to cause that event, then $obs \subset int$ (because *int* could not be said to cause *obs* unless $obs \subset int$).

Definition 14 (causal identity) If ¹⁹ $obs \in L_v$ is an observed event, and $int \in L_v$ is in intervention causing *obs*, then intuitively $c \subseteq int - obs$ “identifies” or “names” the intervening agency if $c \neq \emptyset$. We call *c* a **causal identity** corresponding to *int* and *obs*. Suppose INT and OBS are sets of statements, and we assume OBS contains observed events and INT interventions, then a causal identity corresponding to INT and OBS is $c \neq \emptyset$ s.t. $\forall i \in INT (c \subset int)$ and $\forall obs \in OBS (c \cap obs = \emptyset)$ (we can attempt to construct a causal identity for any INT and OBS). If a policy is a causal identity, then the associated task is to classify interventions.

Definition 15 (purpose, goal or intent)

We consider a policy *c* which is a causal identity corresponding to INT and OBS to be the **intent**, **purpose** or **goal** ascribed to the interventions. *c* is what the interventions share in common, meaning the “name” or “identity” of behaviour is the “intent”, “goal” or “purpose” of behaviour (philosophically my position is that there are no such things as representations, only things in and of themselves). Just as an intervention caused an observation, the particular intent which motivated the agency undertaking the intervention is what caused it (to correctly infer intent, one must infer a causal identity that implies subsequent interventions).

¹⁷ Michael Timothy Bennett. Emergent causality and the foundation of consciousness. In *Artificial General Intelligence*. Springer Nature, 2023b

¹⁸ Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., New York, 1st edition, 2018

¹⁹ (EXAMPLE) Suppose we have organisms *a* (Alice) and *b* (Bob). The inputs Alice has experienced so far $I_{h < t_a}$ can be divided into those in which Bob affected Alice I_a^b and those in which Bob did not $I_a^{-b} = I_{h < t_a} - I_a^b$. By affecting Alice, Bob has intervened in Alice’s experience. Alice can construct a causal identity *b* for Bob corresponding to interventions $INT = I_a^b$ and observations $OBS = I_a^{-b}$. The objects which “exist” in Alice’s experience are those for which she constructs a causal identity, so this is how Bob comes to exist as a distinct “object” which Alice experiences, rather than in parts of other objects).

Meaning

ORIGINAL DEFINITIONS are to be found in ²⁰

Definition 16 (meaning)

The *meaning* an organism ascribes to an input is that protosymbol it uses to interpret. While each organism has a unique sensorimotor vocabulary and thus protosymbol system, but we assume members of a species construct **roughly equivalent** meanings (quality, intent etc). We use $\omega \approx \alpha$ to indicate that ω and α are roughly equivalent. In accord with definition 15, one organism “intends” to affect another if completion of the protosymbol (a task) the former uses to interpret hinges on how the latter’s behaviour is affected.

Assume an organism α (Alice) in input i_α and organism \mathfrak{b} (Bob) in $i_\mathfrak{b}$. α means $\alpha \in \mathfrak{s}_\alpha$ by deciding $u \in E_{i_\alpha}$ iff α intends in deciding u :

1. that \mathfrak{b} interpret $i_\mathfrak{b}$ using $\omega \approx \alpha$,
2. \mathfrak{b} “recognize” this intention by being affected by u such that the input $i_\mathfrak{b} = j$ in which \mathfrak{b} finds itself change to $i_\mathfrak{b} = k \neq j$,
3. and (1) occur on the basis of (2), meaning had α not decided u then \mathfrak{b} would have interpreted $i_\mathfrak{b}$ using $\zeta \not\approx \alpha$.

To communicate meaning organisms must:

1. be able to affect one another.
2. have similar experiences, so $\mathfrak{s}_\mathfrak{b}$ and \mathfrak{s}_α contain roughly equivalent symbols.
3. have similar preferences.

²⁰ Michael Timothy Bennett. On the computation of meaning, language models and incomprehensible horrors. In *Artificial General Intelligence*. Springer Nature, 2023c

Reafference and consciousness

THE ORIGINAL definitions of these can be found in ²¹. A first order self correctly classifies the system's interventions, just as the neuroscientific phenomenon ²² known as "reafference" classifies interventions in living organisms. A first order self allows a system to reason about its own future interventions, plan spatial navigation and so forth.

Definition 17 (preconditions)

If \mathbf{o} is an organism, and c is a causal identity:

- the **representation** precondition is met iff $c \in L_{\mathbf{v}_\mathbf{o}}$, and
- the **incentive** precondition is met if \mathbf{o} must learn c in order to remain "fit".

Definition 18 (first order self)

If c is the lowest level causal identity corresponding to INT and OBS, and INT is every intervention an organism could make (not just past interventions, but all potential future interventions), then we consider c to be the system's **first order self**. If $c \in \mathbf{p}_\mathbf{o}$ then an organism has constructed a first order self. A first order self for an organism \mathbf{o} is denoted \mathbf{o}^1 . An organism has at most one first order self.

Definition 19 (chain notation)

Suppose we have two organisms, \mathbf{a} (Alice) and \mathbf{b} (Bob). $c_\mathbf{a}^\mathbf{b}$ denotes a causal identity for \mathbf{b} constructed by \mathbf{a} (what Alice thinks Bob intends). Subscript denotes the organism who constructs the causal identity, while superscript denotes the object. The superscript can be extended to denote chains of predicted causal identity. For example, $c_\mathbf{a}^{\mathbf{b}\mathbf{a}} \subset c_\mathbf{a}^\mathbf{b}$ denotes \mathbf{a} 's prediction of \mathbf{b} 's prediction of \mathbf{a}^1 (what Alice thinks Bob thinks Alice intends). The superscript of $c_\mathbf{a}^*$ can be extended indefinitely to indicate recursive predictions, however the extent recursion is possible is determined by \mathbf{a} 's vocabulary $\mathbf{v}_\mathbf{a}$. Finally, Bob need not be an organism. Bob can be anything for which Alice constructs a causal identity.

Definition 20 (n^{th} order self)

An n^{th} order self for \mathbf{a} is $\mathbf{a}^n = c_\mathbf{a}^{*\mathbf{a}}$ where $*$ is replaced by a chain, and n denotes the number of reflections. For example, a second order self $\mathbf{a}^2 = c_\mathbf{a}^{\mathbf{b}\mathbf{a}}$, and a third order self $\mathbf{a}^3 = c_\mathbf{a}^{\mathbf{b}\mathbf{a}\mathbf{b}\mathbf{a}}$. We use \mathbf{a}^2 to refer to any second order self, and chain notation to refer to a specific second order self, for example

²¹ Michael Timothy Bennett, Sean Welsh, and Anna Ciaunica. *Why Is Anything Conscious?* Preprint, 2024

²² Andrew B. Barron and Colin Klein. What insects can tell us about the origins of consciousness. *Proceedings of the National Academy of Sciences*, 2016

c_a^{ba} . The union of two n^{th} order selves is also considered to be an n^{th} order self, for example $a^3 = c_a^{baba} \cup c_a^{da\bar{d}a}$, and the weaker or higher level a self is in the generational hierarchy, the more selves there are of which it is part.

Definition 21 (levels of consciousness)

I argue the following stages and the prerequisites:

1. *an organism that acts but does not learn, meaning p_o is fixed from birth.*
2. *an organism that learns, but $o^1 \notin p_o$ either because $o^1 \notin L_{v_o}$ (failing the “representation precondition”) or because the organism is not incentivised to construct o^1 (failing the “incentive precondition”).*
3. *reafference and phenomenal or core consciousness are achieved when $o^1 \in p_o$ is learned by an organism as a consequence of attraction to and repulsion from statements in L_{v_o} .*
4. (a) *access or self reflexive consciousness is achieved when $o^2 \in p_o$.*
 (b) *hard consciousness is achieved when a phenomenally conscious organism learns a second order self (an organism is consciously aware of the contents of second order selves, which must have quality if learned through phenomenal conscious).*
5. *meta self reflexive consciousness (human level hard consciousness) is achieved when $o^3 \in p_o$.*

Illustrative Examples Using General Reinforcement Learning

FOR THE SAKE OF INTUITION I will now provide some simple examples of how these definitions might be applied in the context of general reinforcement learning. In general reinforcement learning (GRL) $\mathcal{A}, \mathcal{O}, \mathcal{R}$ denote the (finite) action, observation and reward set respectively. The following definition of G.R.L. was taken from ²³.

- $\mathcal{E} := \mathcal{O} \times \mathcal{R}$ is the set of **percepts**.
- $\Delta\mathcal{X}$ is the set of distributions over \mathcal{X} .
- $\mathcal{H} := (\mathcal{A} \times \mathcal{E})^*$ is the set of **histories**, where $*$ signifies that histories can be of any length (for example, were we to type $(\mathcal{A} \times \mathcal{E})^3$ would just have all histories of length 3).
- $\pi : \mathcal{H} \rightarrow \Delta\mathcal{A}$ is a function that outputs distributions over actions given a history. π is called a **policy**.
- $\mu : \mathcal{H} \times \mathcal{A} \rightarrow \Delta\mathcal{E}$ is an **environment**. An environment is a distribution over precepts *given* a history, and an action.
- $h_{<t} := a_1e_1...a_{t-1}e_{t-1} = a_1o_1r_1...a_{t-1}o_{t-1}r_{t-1}$ denotes a history of interactions up to time $t - 1$ inclusive.
- Given an environment μ and policy π , the conjunction $\mu\pi$ denotes the induced probability measure on histories, where $h_{<t}$ is assumed to be sampled from $\mu\pi$ (intuitively, an environment does the opposite of a policy, predicting the outcome of an action rather than which action comes next).
- GRL is concerned with how well policies perform (and by extension agents). We measure this performance by the future expected sum of discounted rewards, called value function. The agent selects a policy that maximises the value function.

Note that though it is not explicitly stated, the above definition assumes reward is between 0 and 1.

²³ Elliot Catt, Jordi Grau-Moya, Marcus Hutter, Matthew Aitchison, Tim Genewein, Gregoire Deletang, Li Kevin Wenliang, and Joel Veness. Self-predictive universal AI. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=psXVkk09No>

GRL and Pancomputational Enactivism

HERE I REPRESENT THE DIFFERENT INGREDIENTS of GRL using \mathbf{v} -tasks and an abstraction layer. This is just for illustrative purposes, to give some intuition about how this all works. It is one possible way to translate GRL into the context of my formalism, but is by no means the only way. The most important difference to note is that there is no need to append new outputs to the history. This is because it is an *embodied* formalism. The “input” is not just a stream from sensors, it is the state of the computer; including memory and the surrounding environment. The input at every step therefore includes memory implicitly. These are examples of how these definitions might be applied. I am not asserting that an EGRL system as described below is conscious. It leaves many details unresolved.

Definition 22 (EGRL)

Assume a vocabulary \mathbf{v} , a function $r : L_{\mathbf{v}} \rightarrow \{1, 0\}$ and an integer $len > 0$. The first two imply a \mathbf{v} -task μ such that $O_{\mu} = \{l \in L_{\mathbf{v}} : r(l) = 1\}$ and $\Pi_{\mu} \neq \emptyset$ ²⁴. O_{μ} is every correct output and I_{μ} is every input in which an output can be made. Together these imply several other objects.

- We assume sequence of timesteps from 0 to len inclusive, and that there is a current **timestep** t . At each timestep there is a new “current state” of the environment, which means we get a new input.
- We single out a \mathbf{v} -task $h_{<t} \in \Gamma_{\mathbf{v}}$ of outputs made leading up to time t is a \mathbf{v} -task. We call $h_{<t}$ the **history** until time t . As reward is implicit in the task definition, r is applied to $O_{h_{<t}}$ to obtain a pair of \mathbf{v} -tasks $(h_{<t}^1, h_{<t}^0) \in \Gamma_{\mathbf{v}} \times \Gamma_{\mathbf{v}}$ s.t. $O_{h_{<t}^1} \subseteq O_{\mu}$ and $O_{h_{<t}^0} \cap O_{\mu} = \emptyset$ (meaning $h_{<t}^1$ is a history of correct outputs preceding t , and $h_{<t}^0$ is a history of incorrect outputs).
- A **pair proxy** $<^{\mathbf{pp}}$ is a binary relation on pairs in $\mathbf{pp} = \{(m, n) \in \Pi_{h_{<t}^1} \times \Pi_{h_{<t}^0} : E_m \cap E_n = \emptyset\}$ that classify correct and incorrect outputs. For such pairs, **learning** means choosing $(\pi, n) \in \mathbf{pp}$ that maximises $<^{\mathbf{pp}}$, and **inference** means choosing $o \in E_{i_t} \cap E_{\pi}$.
- Given $(\pi, n), (j, k) \in \mathbf{pp}$ “weakness” pair proxy is $(j, k) <_w^{\mathbf{pp}} (\pi, n)$ iff $|E_j \cup E_k| < |E_{\pi} \cup E_n|$. Intuitively and with some abuse of notation, this is because (vocabulary permitting) there can exist a solitary statement j which has the equivalent extension, meaning $E_{\pi} \cup E_n = E_j$. j is a weakest policy implied by all past outputs, rather than only correct outputs. Whether or not $j \in \mathbf{v}$, it remains the optimal choice according to $<_w$, which means (π, n) is also optimal because it has the equivalent extension.

²⁴ This is because we’re trying to offer a computable and objectively optimal alternative to AIXI, which means we want lossless interpolation. If $\Pi_{\mu} \neq \emptyset$, then we’d need to consider lossy approximations by “selective forgetting” of past outputs inconsistent with otherwise good hypotheses.

For convenience going forward I’ll refer to $\langle \mathbf{v}, \mu, len \rangle$ as an EGRL problem. It is solved by an EGRL system, the pseudocode of which is given below in algorithm 1. To run an EGRL system and solve an EGRL we require $\langle \mathbf{v}, h_{<t}, <^{\mathbf{pp}}, r \rangle$, where $<^{\mathbf{pp}}$ is a pair proxy and $h_{<t}$ is a history (a \mathbf{v} -task).

Algorithm 1: EGRL SYSTEM

```

1: If  $I$  is a set,  $I.choose()$  returns a member of  $I$ , and  $I.add(i) := I \cup \{i\}$ 
2: The  $.choose()$  method assumes a total order or distribution
3:  $\mathfrak{r}_{\mathfrak{h}_{<t}}$  is an int variable for the reward accrued before time  $t$ 
4:
5: function EGRL_LOOP( $\mu, <^{\mathfrak{pp}}, r, len$ )
6:    $i_t \leftarrow I_\mu.choose()$ 
7:    $t, \mathfrak{r}_{\mathfrak{h}_{<t}} \leftarrow 0$ 
8:    $\mathfrak{h}_{<t} \leftarrow \langle \{i_t\}, \emptyset \rangle$ 
9:   while  $t < len$  do
10:     $\pi \leftarrow \text{LEARNING}(\mathfrak{h}_{<t}, <^{\mathfrak{pp}}, i_t, r) \triangleright \text{See LEARNING METHODS}$ 
11:     $o_t \leftarrow \text{INFERENCE}(\pi, i_t)$ 
12:     $\mathfrak{r}_{\mathfrak{h}_{<t}} \leftarrow \mathfrak{r}_{\mathfrak{h}_{<t}} + r(o_t) \triangleright \text{feedback}$ 
13:     $t, i_t, \mathfrak{h}_{<t} \leftarrow \text{NEXT\_TIMESTEP}(t, I_\mu, \mathfrak{h}_{<t})$ 
14:   end while
15:   return  $\mathfrak{r}_{\mathfrak{h}_{<t}}$ 
16: end function
17:
18: function INFERENCE( $\pi, i_t$ )
19:    $E_{\pi \cap i_t} \leftarrow E_\pi \cap E_{i_t}$ 
20:   return  $E_{\pi \cap i_t}.choose()$ 
21: end function
22:
23: function NEXT_TIMESTEP( $t, I_\mu, \mathfrak{h}_{<t}$ )
24:    $t \leftarrow t + 1$ 
25:    $O_{\mathfrak{h}_{<t}} \leftarrow O_{\mathfrak{h}_{<t}}.add(o_t)$ 
26:    $I_{\mathfrak{h}_{<t}} \leftarrow I_{\mathfrak{h}_{<t}}.add(i_t)$ 
27:   return  $t, I_\mu.choose(), \mathfrak{h}_{<t}$ 
28: end function
29:

```

Algorithm 2: LEARNING METHODS

```

1: For a  $\mathbf{v}$ -task  $\alpha$ ,  $\alpha.expand() := \{\omega \in \Gamma_{\mathbf{v}} : \omega \sqsubset \alpha, \neg \exists \zeta \in \Gamma_{\mathbf{v}}(\omega \sqsubset \zeta \sqsubset \alpha)\}$  gets the set of all  $\alpha$ 's immediate children (i.e. children which are only one generation away from  $\alpha$ , meaning we exclude grandchildren etc)
2:  $\mathbf{q}$  is a prioritised queue. If  $(n, obj)$  is in  $\mathbf{q}$ , then  $n$  is inverse priority (smaller  $n$  is higher priority).  $\mathbf{q}.put(n, obj)$  places  $(n, obj)$  in the queue, and  $\mathbf{q}.get()$  returns the highest priority item.  $\mathbf{q}.notempty()$  returns true iff  $\mathbf{q}$  is not empty.
3: If  $\mathcal{V}$  is a set then  $\mathcal{V}.add(c) := \mathcal{V} \cup \{c\}$ .
4:
5: function CLASSIFICATION( $\mathbf{h}_{<t}, r$ ) ▷ used by LEARNING()
6:    $O_{\mathbf{h}_{<t}^1} \leftarrow \{o \in O_{\mathbf{h}_{<t}} : r(o) = 1\}$  ▷ relatively "attractive" memories
7:    $O_{\mathbf{h}_{<t}^0} \leftarrow O_{\mathbf{h}_{<t}} - O_{\mathbf{h}_{<t}^1}$ 
8:    $I_{\mathbf{h}_{<t}^1} \leftarrow \{i \in I_{\mathbf{h}_{<t}} : \exists o \in O_{\mathbf{h}_{<t}^1} (i \subseteq o)\}$ 
9:    $I_{\mathbf{h}_{<t}^0} \leftarrow \{i \in I_{\mathbf{h}_{<t}} : \exists o \in O_{\mathbf{h}_{<t}^0} (i \subseteq o)\}$ 
10:   $\mathbf{h}_{<t}^0 \leftarrow \langle I_{\mathbf{h}_{<t}^0}, O_{\mathbf{h}_{<t}^0} \rangle$ 
11:   $\mathbf{h}_{<t}^1 \leftarrow \langle I_{\mathbf{h}_{<t}^1}, O_{\mathbf{h}_{<t}^1} \rangle$ 
12:  return  $\mathbf{h}_{<t}^1, \mathbf{h}_{<t}^0$ 
13: end function
14:
15: Example 1 : A simple example showing how learning works without noise.
16:
17: function LEARNING( $\mathbf{h}_{<t}, <^{pp}, i_t, r$ )
18:    $\mathbf{h}_{<t}^1, \mathbf{h}_{<t}^0 = \text{CLASSIFICATION}(\mathbf{h}_{<t}, r)$ 
19:    $N_{\mathbf{h}_{<t}^1} \leftarrow \{j \in \Pi_{\mathbf{h}_{<t}^1} : E_j \cap E_{i_t} \neq \emptyset\}$  ▷ discard policies that don't apply
20:    $\Pi_{\mathbf{h}_{<t}^{(1,0)}} \leftarrow \{(j, k) \in N_{\mathbf{h}_{<t}^1} \times \Pi_{\mathbf{h}_{<t}^0} : E_j \cap E_k = \emptyset\}$  ▷ remove ambiguity
21:   return  $\pi$  s.t.  $\exists(\pi, n) \in \Pi_{\mathbf{h}_{<t}^{(1,0)}}$  that maximises  $<^{pp}$ 
22: end function

```

1: *Example 2 : Learning is possible even with noise (the only drawback being that it is harder to understand) by selectively “forgetting” outputs when no valid policy. Given only noise, there would be a policy for each output (rote learning means $|E_\pi| \rightarrow 1$).*

2:

3: **function** LEARNING($h_{<t}, <^{pp}, i_t, r$)

4: $q \leftarrow q.put(0, h_{<t})$

5: $\mathcal{V} \leftarrow \{h_{<t}\}$ ▷ set of visited nodes

6: **while** $q.notempty()$ **do**

7: $n, h \leftarrow q.get()$

8: $h_{<t}^1, h_{<t}^0 = \text{CLASSIFICATION}(h, r)$

9: $N_{h_{<t}^1} \leftarrow \{j \in \Pi_{h_{<t}^1} : E_j \cap E_{i_t} \neq \emptyset\}$

10: $\Pi_{h_{<t}^{(1,0)}} \leftarrow \{(j, k) \in N_{h_{<t}^1} \times \Pi_{h_{<t}^0} : E_j \cap E_k = \emptyset\}$

11: **if** $\Pi_{h_{<t}^{(1,0)}} \neq \emptyset$ **then**

12: **return** π s.t. $\exists (\pi, n) \in \Pi_{h_{<t}^{(1,0)}}$ that maximises $<^{pp}$

13: **end if**

14: **for all** $c \in h.expand()$ **do**

15: **if** $c \notin \mathcal{V}$ **then**

16: $q.put(n+1, c)$

17: $\mathcal{V} \leftarrow \mathcal{V}.add(c)$

18: **end if**

19: **end for**

20: **end while**

21: **return** \emptyset ▷ try uninformed exploration when there is no valid policy

22: **end function**

Algorithm 3: LEARNING METHODS
CONTINUED...

Definition 23 (self-organising EGRL system)

A self-organising EGRL system \mathbf{o} is $\langle \mathbf{v}_\mathbf{o}, \mu_\mathbf{o}, \mathbf{p}_\mathbf{o}, <_\mathbf{o} \rangle$ accompanied by an E.G.R.L system $\langle \mathbf{v}_\mathbf{o}, \mathbf{h}_{<t_\mathbf{o}}, <_\mathbf{o}^{\mathbf{p}\mathbf{p}}, r_\mathbf{o} \rangle$ and an E.G.R.L problem $\langle \mathbf{v}_\mathbf{o}, \mu_\mathbf{o}, len \rangle$:

- $O_{\mu_\mathbf{o}}$ contains every output which qualifies as “fit” according to natural selection.
- $\mathbf{p}_\mathbf{o}$ is the set of policies the self-organising EGRL system knows, s.t. $\mathbf{p}_\mathbf{o} \subset \mathbf{p}_{n.s.} \cup \mathbf{p}_{\mathbf{h}_{<t_\mathbf{o}}}$ and:
 - $\mathbf{p}_{n.s.} \subset L_{\mathbf{v}_\mathbf{o}}$ is **reflexes** hard coded from birth by natural selection.
 - $\mathbf{p}_{\mathbf{h}_{<t_\mathbf{o}}} = \bigcup_{\zeta \in \mathbf{h}_{<t_\mathbf{o}}} \Pi_\zeta$ is the set of policies it is possible to **learn**.
 - If $(\mathbf{p}_\mathbf{o} - \mathbf{p}_{n.s.}) \not\subset \Pi_{\mathbf{h}_{<t_\mathbf{o}}}$ then the self-organising EGRL system has **selective memory**²⁵.
- $<_\mathbf{o}$ is a binary relation over $\Gamma_{\mathbf{v}_\mathbf{o}}$ we call **preferences**.
- \mathbf{o} is an EGRL system iff maximising $<_\mathbf{o}$ maximises r and $\mathbf{p}_\mathbf{o} = \mathbf{p}_{\mathbf{h}_{<t_\mathbf{o}}}$.

²⁵ It can “forget” outputs that contradict otherwise good policies.

Example Definitions ASI and AGI using EGRL

THE FOLLOWING are based on ²⁶. As before these are meant to be illustrative, rather than final definitions:

Definition 24 (bagī [バギー])

“Bennett’s artificial general intelligence” (bagī) is an EGRL system $\langle \mathbf{v}, \mathbf{h}_{<t}, <_w^{\text{pp}}, r \rangle$ that uses weakness as the pair proxy²⁷. From ²⁸ we have that bagī is the optimal system, meaning it is “most intelligent” given embodiment, in that it maximises the ability to complete a wide range of tasks using the vocabulary and history.

Definition 25 (unagi [鰻])

The “unformed artificial general intelligence” (unagi) formalises the upper bound given in ²⁹. It is a system that searches through the space of finite vocabularies $V = \{v \subset P : v \text{ is finite}\}$ to identify the bagī for which the utility ϵ is maximised³⁰. We assume an “unformed GRL problem” (UGRL), which is $\langle \lambda_\mu, \text{len} \rangle$ where λ_μ is an uninstantiated task to be completed and len is the time horizon (it is like an EGRL problem $\langle \mathbf{v}, \mu, \text{len} \rangle$, but λ_μ replaces μ and \mathbf{v}). We also assume a “U.G.R.L system” given by $\langle \lambda_{\mathbf{h}_{<t}}, <_w^{\text{pp}}, r \rangle$ where $<_w^{\text{pp}}$ is the *weakness* pair proxy and $\lambda_{\mathbf{h}_{<t}}$ is the system’s history ($\lambda_{\mathbf{h}_{<t}}$ is an uninstantiated task replacing \mathbf{v} and $\mathbf{h}_{<t}$ from the EGRL system definition)³¹. The unagi is then a search of vocabularies and then policies using those vocabularies. It takes $\langle \lambda_{\mathbf{h}_{<t}}, <_w^{\text{pp}}, r \rangle$ and returns a policy π and vocabulary \mathbf{v} such that:

$$\exists(\pi, \mathbf{v}) \in V \times V \text{ s.t. } \pi \subset \mathbf{v}, \epsilon(\lambda_{\mathbf{h}_{<t}}(\mathbf{v})) \text{ and } <_w^{\text{pp}} \text{ are maximised}$$

Because $<^{\text{pp}}$ is the weakness proxy, an unagi finds the bagī is best suited to an uninstantiated task, and then constructs that bagī (it identifies the best tool for the problem, and then builds it).

²⁶ Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a

²⁷ This assumes a particular vocabulary, which acts as a vehicle for cognition (bagī or バギー means “buggy” in Japanese).

²⁸ Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a

²⁹ Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a

³⁰ Its namesake is the Japanese water eel, slippery (having no particular embodiment). I used to get unagi lunches for free as an undergrad student, and this theoretical unagi system requires a “free lunch” in the sense of the “no free lunch” theorem.

³¹ To have a disembodied history we must embrace dualism, which means unagi is impossible to build within the confines of the environment.

List of proofs

VERSIONS OF THE FOLLOWING THREE proofs were published in ³².

Theorem 1 (sufficiency)

Assume $\alpha \sqsubset \omega$. The weakness proxy sufficient to maximise the probability that a parent ω is learned from a child α ³³.

Proof 1 You're given the definition of \mathbf{v} -task α from which you infer a hypothesis $\mathbf{h} \in \Pi_\alpha$. To learn ω , you need $\mathbf{h} \in \Pi_\omega$:

1. For every $\mathbf{h} \in \Pi_\alpha$ there exists a \mathbf{v} -task $\gamma_{\mathbf{h}} \in \Gamma_{\mathbf{v}}$ s.t. $O_{\gamma_{\mathbf{h}}} = E_{\mathbf{h}}$, meaning \mathbf{h} permits only correct outputs for that task regardless of input. We'll call the highest level task $\gamma_{\mathbf{h}}$ s.t. $O_{\gamma_{\mathbf{h}}} = E_{\mathbf{h}}$ the **policy task** of \mathbf{h} .
2. ω is either the policy task of a policy in Π_α , or a child thereof.
3. If a policy \mathbf{h} is correct for a parent of ω , then it is also correct for ω . Hence we should choose \mathbf{h} that has a policy task with the largest number of children. As tasks are uniformly distributed, that will maximise the probability that ω is $\gamma_{\mathbf{h}}$ or a child thereof.
4. For the purpose of this proof, we say one task is **meaningfully distinct**³⁴ from another if it has different correct outputs.
5. No two policies in Π_α have the same policy task³⁵. This is because all the policies in Π_α are derived from the same set inputs, I_α .
6. The set of statements which might be outputs addressing inputs in I_ω and not I_α , is $\overline{E_{I_\alpha}} = \{l \in L_{\mathbf{v}} : l \notin E_{I_\alpha}\}$ ³⁶.
7. For any given $\mathbf{h} \in \Pi_\alpha$, the extension $E_{\mathbf{h}}$ of \mathbf{h} is the set of outputs \mathbf{h} implies. The subset of $E_{\mathbf{h}}$ which fall outside the scope of what is required for the known task α is $\overline{E_{I_\alpha}} \cap E_{\mathbf{h}}$ ³⁷.
8. $|\overline{E_{I_\alpha}} \cap E_{\mathbf{h}}|$ increases monotonically with $|E_{\mathbf{h}}|$, because for all $e \in E_{\mathbf{h}}$ is $e \notin \overline{E_{I_\alpha}}$ iff $e \in E_{I_\alpha}$.
9. $2^{|\overline{E_{I_\alpha}} \cap E_{\mathbf{h}}|}$ is the number of meaningfully distinct **parents** of α to which \mathbf{h} generalises. It increases monotonically with the weakness of \mathbf{h} .
10. Given \mathbf{v} -tasks are uniformly distributed and $\Pi_\alpha \cap \Pi_\omega \neq \emptyset$, the probability that $\mathbf{h} \in \Pi_\alpha$ generalises to ω is

$$p(\mathbf{h} \in \Pi_\omega \mid \mathbf{h} \in \Pi_\alpha, \alpha \sqsubset \omega) = \frac{2^{|\overline{E_{I_\alpha}} \cap E_{\mathbf{h}}|}}{2^{|\overline{E_{I_\alpha}}|}}$$

$p(\mathbf{h} \in \Pi_\omega \mid \mathbf{h} \in \Pi_\alpha, \alpha \sqsubset \omega)$ is maximised when $|E_{\mathbf{h}}|$ is maximised. Recall from definition 4 that $<_w$ is the **weakness proxy**. For statements l_1, l_2 we have $l_1 <_w l_2$ iff $|E_{l_1}| < |E_{l_2}|$. \mathbf{h} that maximises $<_w$ will also maximise $p(\mathbf{h} \in \Pi_\omega \mid \mathbf{h} \in \Pi_\alpha, \alpha \sqsubset \omega)$. Hence the weakness proxy maximises the probability that³⁸ a parent ω is learned from a child α . \square

³² Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a

³³ Assume there exist correct policies for ω , because otherwise there would be no point in trying to learn it.

³⁴ This is because switching from β to ζ s.t. $I_\beta \neq I_\zeta$ and $O_\beta = O_\zeta$ would be to pursue the same goal in different circumstances. This is because inputs are *subsets* of outputs, so both sets of inputs are implied by the outputs. O_ζ implies I_β and O_β implies I_ζ

³⁵ Every policy task for policies of α is meaningfully distinct from the others.

³⁶ This is because E_{I_α} contains every statement which is a correct output or an incorrect output, and $\overline{E_{I_\alpha}}$ contains every statement which could possibly be in I_ω , E_{I_ω} and thus O_ω .

³⁷ This is because E_{I_α} is the set of all conceivable outputs by which one might attempt to complete α , and so the set of all outputs that can't be made when undertaking α is $\overline{E_{I_\alpha}}$ because those outputs occur given inputs that aren't part of I_α .

³⁸ Subsequently it also maximises the sample efficiency with which a parent ω is learned from a child α .

Theorem 2 (necessity)

To maximise the probability of learning ω from α , it is necessary to weakness as a proxy.

Proof 2 Let α and ω be defined exactly as they were in the proof of prop. 1.

1. If $\mathbf{h} \in \Pi_\alpha$ and $E_{I_\omega} \cap E_{\mathbf{h}} = O_\omega$, then it must be the case that $O_\omega \subseteq E_{\mathbf{h}}$.
 2. If $|E_{\mathbf{h}}| < |O_\omega|$ then generalisation cannot occur, because that would mean that $O_\omega \not\subseteq E_{\mathbf{h}}$.
 3. Therefore generalisation is only possible if $|E_{\mathbf{h}}| \geq |O_\omega|$, meaning a sufficiently weak hypothesis is necessary to generalise from child to parent.
 4. For any two hypotheses \mathbf{h}_1 and \mathbf{h}_2 , if $|E_{\mathbf{h}_1}| < |E_{\mathbf{h}_2}|$ then the probability $p(|E_{\mathbf{h}_1}| \geq |O_\omega|) < p(|E_{\mathbf{h}_2}| \geq |O_\omega|)$ because tasks are uniformly distributed.
 5. Hence the probability that $|E_m| \geq |O_\omega|$ is maximised when $|E_m|$ is maximised.
- To maximise the probability of learning ω from α , it is necessary to select the weakest hypothesis.

To select the weakest hypothesis, it is necessary to use the weakness proxy. \square

Theorem 3 (simplicity sub-optimality)

Description length is neither a necessary nor sufficient proxy for the purposes of maximising the probability that induction generalises.

Proof 3 In propositions 1 and 2 we proved that weakness is a necessary and sufficient choice of proxy to maximise the probability of generalisation. It follows that either maximising $\frac{1}{|m|}$ (minimising description length) maximises $|E_m|$ (weakness), or minimisation of description length is unnecessary to maximise the probability of generalisation. Assume the former, and we'll construct a counterexample with $\mathbf{v} = \{a, b, c, d, e, f, g, h, j, k, z\}$ s.t. $L_{\mathbf{v}} = \{\{a, b, c, d, j, k, z\}, \{e, b, c, d, k\}, \{a, f, c, d, j\}, \{e, b, g, d, j, k, z\}, \{a, f, c, h, j, k\}, \{e, f, g, h, j, k\}\}$ and a task α where

- $I_\alpha = \{\{a, b\}, \{e, b\}\}$
- $O_\alpha = \{\{a, b, c, d, j, k, z\}, \{e, b, g, d, j, k, z\}\}$
- $\Pi_\alpha = \{\{z\}, \{j, k\}\}$

Weakness as a proxy selects $\{j, k\}$, while description length as a proxy selects $\{z\}$. This demonstrates the minimising description length does not necessarily maximise weakness, and maximising weakness does not minimise description length. As weakness is necessary and sufficient to maximise the probability of generalisation, it follows that minimising description length is neither. \square

THE FOLLOWING TWO proofs were published in ³⁹.

Theorem 4 (subjectivity) *If there is no abstraction, complexity can always be minimized without improving sample efficiency, regardless of the task.*

Proof 4 *In accord with definition 2, the absence of abstraction means the vocabulary is the set of all declarative programs, meaning $\mathbf{v} = P$. It follows that for every $l \in L_{\mathbf{v}}$ there exists $f \in \mathbf{v}$ such that $\bigcap l = f$. Statements l and $\{f\}$ are equivalent iff $E_l = E_{\{f\}}$, which is exactly the case here because $\bigcap l = f$. Theorems 1 and 2 show that maximising weakness is necessary and sufficient to maximise the probability of generalisation, which means weakness maximises sample efficiency (is the optimal proxy). This means sample efficiency is determined by the cardinality of extension. For every correct policy l of every task in $\Gamma_{\mathbf{v}}$ there exists $f \in \mathbf{v}$ s.t. $E_l = E_{\{f\}}$. Policy complexity can be minimised regardless weakness, because the simplest representation of every extension is simplicity 1. \square*

Theorem 5 (confounding) *If the vocabulary is finite, then policy weakness can confound⁴⁰ sample efficiency with policy simplicity.*

Proof 5 *We already have that policy weakness causes sample efficiency, in that it is necessary and sufficient to maximise it in order to maximise sample efficiency. Continuing from proof 1, in a finite vocabulary, there may not exist $f \in \mathbf{v}$ s.t. $E_l = E_{\{f\}}$, which means the complexity of all extensions will not be the same. If we choose any vocabulary in which weaker aspects take simpler forms, then simplicity will be correlated with weakness and so will also be correlated with sample efficiency. This means we would choose \mathbf{v} s.t. for all $a, b \in L_{\mathbf{v}}$, the simpler statement has the larger extension, meaning $a <_w b \leftrightarrow a <_d b$. For example, suppose $P = \{a, b, c, \dots\}$, $a = \{1, 2, 4\}$, $b = \{1, 3, 4\}$, $\mathbf{v} = \{a, b\}$, $L_{\mathbf{v}} = \{\{a\}, \{b\}, \{a, b\}\}$, then it follows $\{a, b\} <_w \{a\}$, $\{a, b\} <_w \{b\}$, $\{a, b\} <_d \{a\}$, $\{a, b\} <_d \{b\}$. \square*

³⁹ Michael Timothy Bennett. Is complexity an illusion? In *Artificial General Intelligence*. Springer, 2024b

⁴⁰ A confounds B and C when for example $A = \text{"badlyinjured"}$ causes $B = \text{"died"}$ and $C = \text{"pickedupbyambulance"}$, and it looks like C causes B because $p(B | C) > p(B | \neg C)$, and yet it may be that $p(B | C, A) < p(B | \neg C, A)$.

THE FOLLOWING THREE proofs were published in ⁴¹, ⁴² and ⁴³ respectively. Further proofs and discussions are available in a reply to one of the papers, written by Gabriel Simmons ⁴⁴. It gives an intuitive and succinct description of the formalism, and an example of a task that doesn't have a correct policy. The comment pertains to one paper in particular and so is missing some of the broader context, but it raises some interesting questions that I answer in other publications, or seek to answer in this thesis.

Theorem 6 (upper bound) *The most 'intelligent' choice of policy and vocabulary given uninstantiated task λ_ρ is π and \mathbf{v} s.t. \mathbf{v} maximises utility for $\lambda_\rho(\mathbf{v})$, $\pi \in \Pi_{\lambda_\rho(\mathbf{v})}$ and π maximises weakness.*

Proof 6 *We have equated intelligence with sample efficient generalisation. According to theorems 1 and 2 the weakest correct policies have the highest probability of generalising. Given an uninstantiated task λ_ρ , utility measures the weakness of the weakest correct policies. We can use this to compare vocabularies. By choosing a vocabulary \mathbf{v} which maximises utility for $\lambda_\rho(\mathbf{v})$, we instantiate λ_ρ in a vocabulary that maximises the weakness of correct policies for λ_ρ . Then, using weakness proxy, we can select a policy that has the highest possible probability of generalising, and thus maximise sample efficiency. \square*

Theorem 7 (lower levels constrain higher levels) *The greater the utility $\epsilon(\lambda^{i+1}(\mathbf{v}^{i+1}))$, the weaker the policy π^i s.t. $\mathfrak{f}(E_{\pi^i}) = \mathbf{v}^{i+1}$ must be⁴⁵.*

Proof 7 *If $\mathbf{a} \subset \mathbf{b}$ then $\epsilon(\lambda^{i+1}(\mathbf{a})) < \epsilon(\lambda^{i+1}(\mathbf{b}))$, meaning if \mathbf{b} is the vocabulary at $i+1$, then it will be possible to construct weaker policies than if \mathbf{a} is the vocabulary (intuitively, a larger vocabulary enables a wider range of policies). We consider two policies π_a^i and π_b^i which could be the policy π^i at i . If $\mathbf{a} = \mathfrak{f}(E_{\pi_a^i}) \subset \mathfrak{f}(E_{\pi_b^i}) = \mathbf{b}$, then $\pi_a^i <_w \pi_b^i$, meaning an enlarged vocabulary at $i+1$ implies a weaker policy at i . \square*

Theorem 8 (n^{th} order self convergence) *An organism that uses weakness as its proxy will learn an n^{th} order self if the incentive and representation preconditions are met for that order of self.*

Proof 8 *Assume we have an organism \mathbf{o} that learns using "weakness" as a proxy. A $\mathbf{v}_\mathbf{o}$ -task $\mathfrak{h}_{<t_\mathbf{o}}$ represents the history of \mathbf{o} (meaning $\mathfrak{h}_{<t_\mathbf{o}} \sqsubset \mu_\mathbf{o}$ and $\mathfrak{h}_{<t_\mathbf{o}}$ is an ostensive definition of $\mu_\mathbf{o}$, by virtue of the fact that \mathbf{o} remains alive). The organism explores the environment, intervening to maintain homeostasis. As it does so, more and more inputs and outputs are included in $\mathfrak{h}_{<t_\mathbf{o}}$. It follows that:*

⁴¹ Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a

⁴² Michael Timothy Bennett. Multiscale causal learning. 2024c

⁴³ Michael Timothy Bennett, Sean Welsh, and Anna Ciaunica. *Why Is Anything Conscious?* Preprint, 2024

⁴⁴ Gabriel Simmons. Comment on is complexity an illusion?, 2024. URL <https://arxiv.org/abs/2411.08897>

⁴⁵ Intuitively, this just means adaptability at higher levels implies adaptability at lower levels.

1. From the representation precondition we have that there exists a n^{th} order self $\mathbf{o}^n \in L_{\mathbf{v}_\mathbf{o}}$.
2. To remain fit, \mathbf{o} must “generalise” to $\mu_\mathbf{o}$ from $\mathbf{h}_{<t_\mathbf{o}}$. According to the incentive precondition, generalisation to $\mu_\mathbf{o}$ requires \mathbf{o} learn the n^{th} order self, which is when $\mathbf{o}^n \in \mathbf{p}_\mathbf{o}$.
3. From ⁴⁶ we have proof that weakness is the optimal choice of proxy to maximise the probability of generalisation from child to parent is the weakest policy. It follows that \mathbf{o} will generalise from $\mathbf{h}_{<t_\mathbf{o}}$ to $\mu_\mathbf{o}$ given the smallest history of interventions with which it is possible to do so (meaning the smallest possible ostensive definition, or cardinality $|\mathcal{O}_\alpha|$).

Were we to assume learning under the above conditions does not construct an n^{th} order self for \mathbf{o} , then one of the three statements above would be false and we would have a contradiction. It follows that the proposition must be true. \square

⁴⁶ Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a

THE FOLLOWING are supplementary proofs that exist only here.

Theorem 9 (pair proxy optimality)

The optimal choice of pair proxy for EGRL is the weakness pair proxy.

Proof 9 For every pair (a, b) of statements s.t. $a \cap b = \emptyset$ there exists a statement c such that $E_c = E_a \cup E_b$. The completions of c are completions of either a or b . An extension is a truth table, so this means that a pair (a, b) is equivalent to a single policy c . From 2 and 1 we already have that the optimal policy for a single policy is weakness, and we can the weakness of c in accord with $<_w$ and definition 5 is exactly the weakness of (a, b) in accord with $<_w^{pp}$ and definition 22. Because a pair of policies is equivalent to a policy the optimal pair proxy is equivalent to the optimal proxy, and so the optimal pair proxy is $<_w^{pp}$.

Theorem 10 (do operator equivalence) *We can substitute the do operator for a statement I .*

Proof 10 In the context of the conditional probability of Y , for every intervention $do[X = x]$ there exists a variable A which is equivalent given X , meaning $p(Y \mid A = \text{true}, X = x) = p(Y \mid do[X = x])$ and $p(Y \mid A = \text{false}, X = x) = p(Y \mid X = x)$ (the equivalence between variables and the do operator was pointed out by ⁴⁷, however we will prove it here for illustrative purposes). Likewise, for each possible value y of Y , there are three possibilities to consider pertaining to X and the do operator. Each of these cases has a probability of occurring represented by the rightmost column.

⁴⁷ A. P. Dawid. Influence diagrams for causal modelling and inference. *International Statistical Review / Revue Internationale de Statistique*, 70(2):161–189, 2002. ISSN 03067734, 17515823. URL <http://www.jstor.org/stable/1403901>

$do[X = x]$	$X = x$	$p(Y \mid \dots)$
true	true	$p(Y \mid do[X = x]) = p(Y)$
false	true	$p(Y \mid X = x)$
false	false	$p(Y \mid X \neq x)$

In all cases we either have $do[X = x]$, or we do not. Therefore we can substitute the application of the do operator effecting $X = x$ for $A \in \{\text{true}, \text{false}\}$. A is not a random variable, but something to which one assigns a true or false value to indicate whether $X = x$ is the result of one's intervention. We then assert the following conditional probabilities:

$A = \text{true}$	$X = x$	$p(Y \mid \dots)$
true	true	$p(Y \mid A = \text{true}, X = x) = p(Y)$
true	false	$p(Y \mid A = \text{true}, X \neq x) = 0$
false	true	$p(Y \mid A = \text{false}, X = x) = p(Y \mid X = x)$
false	false	$p(Y \mid A = \text{false}, X \neq x) = p(Y \mid X \neq x)$

As $p(Y \mid A = \text{true}, X = x) = p(Y \mid \text{do}[X = x])$ and $p(Y \mid A = \text{false}, X = x) = p(Y \mid X = x)$ for all possible cases ($p(Y \mid A = \text{true}, X \neq x)$ being an impossibility). From definition 1 we have that 2^Φ contains every declarative program. 2^Φ is sufficient to represent A , because there exists $a \subset 2^\Phi$ s.t. a is true iff $A = \text{true}$, so a is equivalent to $\text{do}[A]$. \square

Experiments

CONTEXT FOR AND A DETAILED ANALYSIS of these experiments (see “experiments.py”) is given in ⁴⁸. What is given here is merely a brief technical report on the two experiments. In these two experiments, a toy program computes policies for 8-bit string prediction tasks (binary addition and multiplication)⁴⁹.

(ABSTRACTION LAYER) There were 256 states, one for every possible 8-bit string. The statements in L were expressions regarding those 8 bits that could be written in propositional logic (\neg , \wedge and \vee).

(TASK) A task was specified by choosing $O \subset L$ such that all $d \in O$ conformed to the rules of either binary addition (for the first experiment) or multiplication (for the second experiment) with 4-bits of input, followed by 4-bits of output.

EACH OF THE TWO experiments (addition and multiplication) involved repeated trials (sampling results). The parameters of each trial were “operation” (a function), and an even integer “number_of_trials” between 4 and 14 which determined the cardinality of the set O_k (defined below). Each trial was divided into training and testing phases.

(TRAINING PHASE)

1. A task T_n was generated:
 - (a) First, every possible 4-bit input for the chosen binary operation was used to generate an 8-bit string. These 16 strings then formed O_n .
 - (b) A bit between 0 and 7 was then chosen, and I_n created by cloning O_n and deleting the chosen bit from every string (meaning I_n was composed of 16 different 7-bit strings, each of which could be found in an 8-bit string in O_n).
2. A child-task $T_k = \langle I_k, O_k \rangle$ was sampled from the parent task T_n . Recall, $|O_k|$ was determined as a parameter of the trial.
3. From T_k two policies (formerly known as models) were generated; a weakest c_w , and a MDL c_{mdl} .

(TESTING PHASE) For each policy $c \in \{c_w, c_{mdl}\}$:

1. The extension E_c of c was then generated.
2. A prediction O_{recon} was then constructed s.t. $O_{recon} = \{e \in E_c : \exists s \in I_n (s \subset e)\}$.

⁴⁸ Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a

⁴⁹ Notation here varies slightly from the formal notation due to the limitations of what can be written in Python (not latex), and because it the experiments coincided with an earlier iteration of the formalism.

3. O_{recon} was then compared to the ground truth O_n , and results recorded.

Between 75 and 256 trials were run for each value of the parameter $|O_k|$. Fewer trials were run for larger values of $|O_k|$ due to restricted availability of hardware. The results of these trails were then averaged for each value of $|O_k|$.

(MEASUREMENTS) Generalisation was deemed to have occurred where $O_{recon} = O_n$. The number of trials in which generalisation occurred was measured, and divided by n to obtain the rate of generalisation for c_w and c_{mdl} . Error was computed as a Wald 95% confidence interval.

Even where $O_{recon} \neq O_n$, the extent to which policies generalised could be ascertained. $\frac{|O_{recon} \cap O_n|}{|O_n|}$ was measured and averaged for each value of $|O_k|$, and the standard error computed.

$ O_k $	c_w				c_{mdl}			
	Rate	$\pm 95\%$	AvgExt	StdErr	Rate	$\pm 95\%$	AvgExt	StdErr
6	.11	.039	.75	.008	.10	.037	.48	.012
10	.27	.064	.91	.006	.13	.048	.69	.009
14	.68	.106	.98	.005	.24	.097	.91	.006

Table 1: Binary addition.

$ O_k $	c_w				c_{mdl}			
	Rate	$\pm 95\%$	AvgExt	StdErr	Rate	$\pm 95\%$	AvgExt	StdErr
6	.05	.026	.74	.009	.01	.011	.58	.011
10	.16	.045	.86	.006	.08	.034	.78	.008
14	.46	.061	.96	.003	.21	.050	.93	.003

Table 2: Binary multiplication.

Further Examples

WE ARE ONLY CONCERNED WITH COMPUTERS THAT CAN ACTUALLY EXIST. Because memory is always finite, a computer that exists is a finite state machine (F.S.M.), not a Turing Machine. Instead of conventional models of computation we're using tasks, which are more general. To show how tasks relate to established notions of computation, we show how a F.S.M. can be converted into a task and back again. Let F be the set of all finite state machines, and recall from definition 6 that the set of all tasks in all vocabularies is Γ .

A F.S.M. is $\langle Q, \sigma, \delta, q_0, A \rangle$:

- Q is a set of F.S.M. states (we use "F.S.M. state" to distinguish Q from Φ).
- Σ is the input alphabet (a finite non-empty set of symbols);
- $\delta : Q \times \Sigma \rightarrow Q$ is the state transition function.
- $q_0 \in Q$ is the initial F.S.M. state.
- $A \subset Q$ is a set of correct or "accepting" F.S.M. states.

Recall that statements and the declarative programs they contain describe any and all aspects of a physicalist environment. This means we can represent the behaviour of a F.S.M. f as a v-task ρ . To do so we define:

- a sensorimotor vocabulary \mathfrak{v} s.t. there is a unique declarative program in \mathfrak{v} for each and every member of Q and Σ , every subset of Q and Σ , and every possible sequence of members of these sets or any part thereof.
- E_{I_ρ} as the set of all statements describing those sequences of input symbols and state transitions that include δ_ρ . Because Σ is finite, there are a finite number of infinite length input sequences, and so we can describe even these infinite length sequences using our finite vocabulary \mathfrak{v} .
- $O_\rho \subset E_{I_\rho}$ is the set of all statements describing sequences ending in an "accepting" state, and $I_\rho = E_{I_\rho} - O_\rho$.
 - Each member $\sigma \in \Sigma$ becomes a member of \mathfrak{v} . We call these **declarative equivalents**. We also include in \mathfrak{v} a declarative program $\sigma_{a < b}$ for every pair $(\sigma_a, \sigma_b) \in \Sigma$ such that $\phi \in \sigma_{a < b}$ iff $\phi \in \sigma_a$ and $\phi \in \sigma_b$ and σ_b comes "after" σ_a . This indicates relative ordering. We call these **declarative orderings**.

- If $\sigma_a, \sigma_b, \sigma_{a < b} \in l \in L_v$, then if l is true then it is a fact that σ_a comes before σ_b .
- We also then include in v declarative equivalents and orderings for the members of Q . This means we can construct a statement l that describes an ordered sequence of states and input symbols.
- To avoid infinite sequences we construct a set K of declarative programs, one for every possible **non repeating** (as in not self similar) sequence of symbols and states, and an additional declarative program to indicate self similarity. Because there are finitely many symbols and states, there are finitely many possible sequences which are not self similar. Doing so lets us represent every possible sequence, so we include these programs in v (meaning $K \subset v$).
- We must also ensure there is a correct policy. An easy way to do this is to include a declarative program $\delta_\rho \in v$ s.t. $\delta_\rho \in l$ only when a sequence described by l is exactly what δ implies (it is the declarative equivalent of the imperative δ).
- We then construct a set $\Delta \subset L_v$ s.t. $\Delta = \{l \in L_v : \exists k \in L_v (l \subseteq k, \delta_\rho \in k)\}$, which means every sequence or part thereof which could be output by δ .
- $I_\rho \subset \Delta$ is then the set of all sequences that do **not** include the declarative equivalent of a member of A .
- $O_\rho \subset \Delta$ is then the set of all sequences that **do** include the declarative equivalent of a member of A .

f is then described in perfect detail, and every $f \in F$ is mapped to a unique task $\rho \in \Gamma$, so we have an injection $i : F \rightarrow \Gamma$. The process described above is reversible. We construct Q from E_{I_ρ} , Σ and q_0 from I_ρ , A from O_ρ , and then δ is just the imperative equivalent of δ_ρ .

Bibliography

Andrew B. Barron and Colin Klein. What insects can tell us about the origins of consciousness. *Proceedings of the National Academy of Sciences*, 2016.

Michael Timothy Bennett. The optimal choice of hypothesis is the weakest, not the shortest. In *Artificial General Intelligence*. Springer Nature, 2023a.

Michael Timothy Bennett. Emergent causality and the foundation of consciousness. In *Artificial General Intelligence*. Springer Nature, 2023b.

Michael Timothy Bennett. On the computation of meaning, language models and incomprehensible horrors. In *Artificial General Intelligence*. Springer Nature, 2023c.

Michael Timothy Bennett. Computational dualism and objective superintelligence. In *Artificial General Intelligence*. Springer, 2024a.

Michael Timothy Bennett. Is complexity an illusion? In *Artificial General Intelligence*. Springer, 2024b.

Michael Timothy Bennett. Multiscale causal learning. 2024c.

Michael Timothy Bennett, Sean Welsh, and Anna Ciaunica. *Why Is Anything Conscious?* Preprint, 2024.

Elliot Catt, Jordi Grau-Moya, Marcus Hutter, Matthew Aitchison, Tim Genewein, Gregoire Deletang, Li Kevin Wenliang, and Joel Veness. Self-predictive universal AI. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=psXVvK09No>.

A. P. Dawid. Influence diagrams for causal modelling and inference. *International Statistical Review / Revue Internationale de Statistique*, 70(2):161–189, 2002. ISSN 03067734, 17515823. URL <http://www.jstor.org/stable/1403901>.

Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., New York, 1st edition, 2018.

Gabriel Simmons. Comment on is complexity an illusion?, 2024. URL <https://arxiv.org/abs/2411.08897>.