

# What The F\*ck Is Artificial General Intelligence?

Michael Timothy Bennett<sup>1[0000-0001-6895-8782]</sup>

The Australian National University  
 michael.bennett@anu.edu.au

**Abstract.** Artificial general intelligence (AGI) is an established field of research. Yet some have questioned if the term still has meaning. AGI has been subject to so much hype and speculation it has become something of a Rorschach test. Melanie Mitchell argues the debate will only be settled through long term, scientific investigation. To that end here is a short, accessible and provocative overview of AGI. I compare definitions of intelligence, settling on intelligence in terms of adaptation and AGI as an artificial scientist. Taking my cue from Sutton’s Bitter Lesson I describe two foundational tools used to build adaptive systems: search and approximation. I compare pros, cons, hybrids and architectures like o3, AlphaGo, AERA, NARS and Hyperon. I then discuss overall meta-approaches to making systems behave more intelligently. I divide them into scale-maxing, simp-maxing, w-maxing based on the Bitter Lesson, Ockham’s and Bennett’s Razors. These maximise resources, simplicity of form, and the weakness of constraints on functionality. I discuss examples including AIXI, the free energy principle and The Embiggening of language models. I conclude that though scale-maxed approximation dominates, AGI will be a fusion of tools and meta-approaches. The Embiggening was enabled by improvements in hardware. Now the bottlenecks are sample and energy efficiency.

**Keywords:** artificial general intelligence.

## 1 Introduction

Picture a machine endowed with human intellect. In its most simplistic form, that is Artificial General Intelligence (AGI) [1]. AGI is also a well established and rigorous field of research [2]. However public perception of AGI is plagued by wild speculation and hype. Some see it as Skynet waiting to pounce [3, 1]. Others, like Melanie Mitchell, question if the term still has any meaning [4]. Speculation and hype have reduced it to a Rorschach test. As Mitchell points out, the debate will not be settled not by media but by rigorous, scientific research. Here I present a short and accessible survey to that end. It is framed in intentionally provocative terms, to spark debate<sup>1</sup>.

---

<sup>1</sup> There is precedent for the use of profanity in a paper title [5]. However this paper will be published in the 2025 AGI proceedings under the title “What Is Artificial General Intelligence” because Anton threw a tantrum. The real name of the paper remains What the F\*ck Is Artificial General Intelligence. Please cite it as that. I’d like to dedicate this footnote to Anton’s pearl clutching. Good job Anton.

## 2 Intelligence

I'll begin by defining intelligence and AGI. There are a number of positions [6, 2, 7–12]. Some peg AGI to **human-level performance** across a broad range of tasks [13, 1]. This is intuitive, but anthropocentric and hard to quantify<sup>2</sup>. Chollet argues intelligence is a measure of the ability to generalise and **acquire new skills**. He argues AGI can do this at least as well as a human [11]. He attempts to quantify the ability to acquire new skills, which can encompass the aforementioned anthropocentric definition. His formalism resembles Legg-Hutter intelligence. Legg and Hutter argued intelligence is an ability to **satisfy goals in a wide range of environments** [10]<sup>3</sup>. Chollet's definition descends from Legg-Hutter. It is based on Ockham's Razor. They both use Kolmogorov complexity. They both equate simplicity with generality. They both seek to quantify intelligence, and they are both highly subjective because they treat intelligence as a property of software interacting with the world through an interpreter [15–18].

That a problem. Why? Because if I develop an AI for some purpose, then  $I$  decide whether it has fulfilled that purpose, and  $I$  am part of the agent's environment. The environment is where *objective* success or failure is decided. Assume  $\mathcal{C}$  is a space of software programs, and  $\Gamma$  is a space of behaviours. Imagine  $f_1 \in \mathcal{C}$  is AI software,  $f_2 : \mathcal{C} \rightarrow \Gamma$  is the hardware on which it runs, and  $f_3 : \Gamma \rightarrow \{0, 1\}$  is the environment (including me) where success is decided. Success is a matter of  $f_3(f_2(f_1))$ . The behaviour of  $f_3(f_2(f_1))$  can be changed by changing  $f_2$  or  $f_3$  [12]. It is pointless to make claims about  $f_3(f_2(f_1))$  based on  $f_1$  alone.  $f_1$  and  $f_2$  are like mind and body. Every choice of embodiment biases the system in some way. Each movement it makes constrains the space of possibilities, much like a constraint expressed in a formal language. Complexity is a property of how a body interprets information [18]. The choice of Universal Turing Machine can make any software agent optimal according to Legg-Hutter intelligence [17].

The idea of AI as a software mind is called **computational dualism** [12]<sup>4</sup>. It is a reference to the work of Descartes, who in 1637 argued the pineal gland mediates between mind and body. AI researchers have exchanged the pineal gland for a Turing machine. So what is the alternative? Wang defines intelligence as **adaptation with limited resources** [6]. This leaves room for us to avoid dualism, and it implies the ability to satisfy goals in a wide range of environments anyway [12].

An attempt was made to resolve computational dualism and formalise intelligence as objective adaptability. It does so by formalising software, hardware and

---

<sup>2</sup> More accurately, I might say it is too easy to quantify. There are so many ways we might quantify it. It is ambiguous and too weak a criteria to be much use for anything but padding out Sam Altman's Twitter feed.

<sup>3</sup> This treats intelligence as implicitly separable from goals, endorsing the orthogonality thesis [14].

<sup>4</sup> Computational dualism is grounded in lengthy formal definitions and derivations given elsewhere [19, 18, 12, 20, 21]. For this survey that level of formality would be counterproductive.

environment together [12]. It formalises intelligence as a measure of the **ability to complete a wide range of tasks** [21]. This dispenses with the separation of goals and intelligence in favour of a whole-of-system model that treats the purpose of a system as what it does. One's body implies a set of goals and subgoals. Body, environment and goals together form a task, by which I mean a purpose and a means of fulfilling it. If  $A$  completes a superset of tasks that  $B$  completes, then  $A$  is more adaptable than  $B$ . This encompasses both **sample and energy efficiency**. It is how fast a system can adapt and how much energy it needs to do so. This is the definition I will use for this survey. I'll consider an AGI to be a system that adapts at least as generally as a human scientist [22]. An **artificial scientist** can prioritise, plan and perform useful experiments. This requires autonomy, agency, motives, an ability to learn cause and effect and the ability to balance exploring to acquire knowledge with acting to profit from it [23, 9, 8, 24, 25].

Artificial intelligence (AI) and machine learning (ML) are typically divided up into buckets like supervised learning, reinforcement learning, regression, classification, planning and so on. These are not useful categories for AGI, because an artificial scientist must be able to do all of these things. Instead, I will take my cue from Sutton's Bitter Lesson. It acknowledges that generally applicable tools can be used to learn any behaviour [26], if we scale up resources (compute, memory, data etc).

### 3 Tools

*Search:* Informally, by search I mean systems that take structure, and then *construct* a solution within the confines of that structure. For example, take a map and then plan a route by trying first every combination of one, then two, then three and more turns until you find the smallest sequence of turns that end at your destination. Search typically refers to algorithms like A\* used for symbolic reasoning and planning problems [27]. These involve describing a problem and goal as a set of rules, and then constructing possible courses of action until one is found that obeys all the rules. Heuristics are used to construct a solution faster<sup>5</sup>. In theory any problem can be framed as a search problem<sup>6</sup>. Search has **advantages**. It produces verifiably correct and interpretable answers. It excels at

---

<sup>5</sup> An example of a heuristic is a function that takes a sequence of turns and tells you how far the end is from your destination.

<sup>6</sup> At first glance it might be tempting to object, and say search can only be applied for a well-defined problem with unambiguous rules and goals. However search can easily be applied to poorly defined problems. In the absence of well defined rules and goals, search can be used to infer rules and goals from observed data. It can do this by treating observed data or subsets thereof as rules, and searching the space of all possible criteria until one is found that explains some or all of observed data[28, 29]. Everything can be reduced to a search problem for much the same reason that every imperative program (instruction like “do this”) can be translated into an equivalent declarative program (an assertion like “the pen is red”)[30]; it is simply a matter of framing.

planning [31] and is typically used in map software. It can prove theorems [32]. In the 90s, search defeated the world chess champion [33]. Search can be used to learn, by iterating through possible hypotheses or models until one is found that conforms to observed data. However it also has **disadvantages**. Iterating through large state spaces is expensive. Hand crafted constraints can be added to reduce the search space, but that is not very scalable. Search tends to be sequential<sup>7</sup>, making it ill suited to take full advantage modern hardware, which was originally designed to parallelise graphical rendering and physics simulation in games. Only later was this hardware adapted for AI [34]. Parallel search algorithms exist but there is a lot of room for improvement [35–38]. The consequence is that search is only really practical at a higher levels of abstraction, where problems are represented using a small number of abstract symbols or well defined parts.

*Approximation:* Sutton’s Bitter Lesson described the alternative to search as learning. However search can be used to learn [28], so to avoid confusion I’ll use the term approximation instead. In any case most of modern machine learning is approximate. Typically it involves taking a model that can map inputs to outputs, then changing its parameters to so that the relation between inputs and outputs approximates training data. For example, convolutional neural networks can be taught to classify the contents of images [39]. Transformers trained on large corpus of text can generate human-like responses [40]. Approximation is imprecise, but for that very reason it is great at dealing with noisy data and large state spaces. It is easy to parallelise and scale on current hardware. There are drawbacks. Approximation is unreliable<sup>8</sup>, because it is by definition only approximate. It is not easily interpretable [41]. Most importantly, current methods are extremely sample and energy inefficient [42, 21]. This makes them less adaptable. Sample inefficiency doesn’t just mean a model is slower to learn. It means the model does not cope well with anything outside the norm. In layman’s terms, an approximation is mid. If two models are trained on the same amount of data, then the more sample efficient model will deal as well or better with edge cases<sup>9</sup>.

*Hybrids and Architectures:* Hybrids are those systems which don’t fit neatly into search or approximation. For example collectives of living cells self-organise and adapt. They can traverse a morphospace during development or regeneration

---

<sup>7</sup> Tends to in present day implementations. Doesn’t need to be.

<sup>8</sup> Imprecision implies a sort of unreliability.

<sup>9</sup> The proofs are given elsewhere [28, 29, 43, 19], but suppose for a second I learn faster than you (more sample efficient). We both learn how to fix a table from the same three examples of someone fixing a table. We can now both fix the table with 100% accuracy. However table fixing is an example of fixing in general. If I am more sample efficient, I must now be able to fix something you cannot. An edge case, like a chair. Now, there’s an upper bound for when we have both learned everything, but given that is impossible given finite resources it will always be the case that the more sample efficient learner will deal better with edge cases, all else being equal.

[44], which is like search. Animals mimic and thus approximate behaviour. It is difficult to argue biological self-organisation falls neatly into either search or approximation. Also, our current methods for search and approximation have complementary strengths and weaknesses. They can be combined to get the best of both worlds. Hybrids are inherently more general because they're not tied to one playbook. Need precision? Search. Got a mess of unstructured data? Approximate. By fusing their strengths, hybrids promise robustness where single-track systems choke [45].

Perhaps the simplest example of a hybrid is **AlphaGo** [46]. It vanquished Go's world champion using a combination of search and approximation. Search enabled AlphaGo to explore potential sequences of moves within the game's constraints. Deep neural networks then approximated how likely sequences were to win. Intuitively, think of these as 'how to play' and 'how to win' respectively. This synergy allowed AlphaGo to surpass human champions, demonstrating the potential of hybrid approaches in mastering complex, strategic tasks.

Search tends to be applied in the context of high level symbolic abstractions that depend on human interpretation. For example, the word 'cat' is just a sound until someone interprets it. It must be decided why and how a particular problem is represented using a particular language or set of symbols. This is the symbol grounding problem [47]. **Neuro-symbolic hybrids** attempt to address it [48]. These systems typically employ neural networks to interpret raw input, converting them into symbolic representations that encapsulate meaning. Search can then be applied to these representations to enable tasks such as planning or logical inference. However it should be noted that the complexity of a problem depends on how it is represented [18], and not all choices of symbolic representation are equal. Another hybrid approach is **structured reinforcement learning**. It leverages approximation to reduce high-dimensional raw sensory data to a more manageable symbolic format. Convolutional autoencoders are used to compress high-dimensional data into concise, symbolic forms that try to capture what is *relevant* in the input, enabling more effective adaptation to dynamic environments [49]. More recent examples include OpenAI's o3 and DeepMind's AlphaGeometry. **o3** employs chain-of-thought reasoning, blending approximation with a structured processes for complex problem-solving [50]. **AlphaGeometry** combines neural networks with symbolic reasoning to solve geometry problems [51]. These systems exemplify the shift towards hybrid approaches for more capable AI.

Finally, there are comprehensive frameworks designed to be generally intelligent. Cognitive architectures and autonomous machines constructed from modules that each serve a different purpose. Perception, memory, and reasoning modules. System 1 and system 2. For example scaffolding can be applied to neural networks to facilitate persistent identity and memory [52]. Pioneering examples include cognitive architectures like SOAR [53] and ACT-R [54]. More recent examples include Hyperon, Autocatalytic Endogenous Reflective Architecture (AERA) and the Non-Axiomatic Reasoning System (NARS).

- **Hyperon** is a modular, distributed system integrating probabilistic logic networks, neural networks, and a knowledge metagraph for holistic cognition [23, 55, 7]. It is highly distributed, modular, scalable and self-organising. This makes it a versatile AGI platform that can integrate new technology as it develops. For example it appears Hyperon will soon incorporate a discrete form of active inference [56, 57].
- **AERA** self-programmes, reflecting on its own symbolic structures while learning statistically. It emphasises analogy, causality, autonomy and growth, with predictive modelling supporting proactive adaptation [58, 8, 59–61].
- **NARS** rejects rigid axioms for a fluid, adaptive logic. Operating under the Assumption of Insufficient Knowledge and Resources (AIKR), NARS reasons with incomplete, uncertain data via a non-axiomatic framework. It integrates symbolic reasoning with probabilistic inference, using a custom inheritance-based logic (NAL) to derive conclusions from limited evidence. Designed for real-time adaptability, NARS learns incrementally, refining its knowledge base as new inputs arrive [9, 62].

Hybrid systems have obvious advantages. They can be more efficient, interpretable, they can integrate human priors effectively and above all they allow for autonomy. It can also be difficult to harmonise disparate methodologies. A lack of robust theoretical guidance risks hybrids being ad hoc rather than principled. This brings us to the final piece of the puzzle.

## 4 Meta-Approaches

If we’re to build an artificial scientist, we need a clear idea of what we’re optimising for. What constitutes a ‘good’ hypothesis? We need a theory that predicts whether one model adapts better than another. There are many such cases. These aren’t algorithms so much as they are philosophies with teeth. Guiding principles like ‘always choose the least specific solution’ or ‘delegate control instead of micromanaging tasks’. I call them meta-approaches. Examples include the Free Energy Principle [56], Universal Artificial Intelligence (UAI) [63], the Minimum Description Length Principle [64], the scaling hypothesis [26], Stack Theory or Pancomputational Enactivism [12, 21] and even organisational principles like the military doctrine of Mission Command [65]. To simplify matters I put meta-approaches into buckets based on what they share in common. Scale-maxing, simp-maxing, and w-maxing. These maximise resources, simplicity of form and versatility of function respectively. They can be understood in terms of Sutton’s Bitter Lesson, Ockham’s Razor and Bennett’s Razor respectively [26, 66, 28]. I’ll begin with the elephant in the room.

*Scale-Maxed:* As Sutton observed we seem to be able to just crank up compute, data and model size and get something that looks like intelligence. Scale-maxed approximation has defined recent history. I call this period ‘The Embiggening’. Language and vision models just got bigger as a bottleneck in compute was removed by technology originally developed for games [34]. GPT-3? 175 billion parameters, 45TB of text, and it’s churning out essays, code, and creepy love letters [67]. AlphaFold 2? Threw a data center at protein folding and cracked a puzzle that had biologists weeping for decades [68]. However performance gains diminish with scale [69]. The energy bill is a nightmare [42]. Worst of all is the sample inefficiency. Today, scale-maxed approximations like GPT-4 struggle with novelty and always will, because novelty is by definition that of which we have few examples. Confront a large language model (LLM) with the genuinely unusual, and it’ll flail like a toddler in a calculus exam. The Bitter Lesson says scale will eventually work. It will eventually identify all other meta-approaches. Fine. Scale will eventually work, but *eventually* is doing all the work in that sentence.

*Simp-Maxed:* Simplicity maximisation (simp-maxing) assumes the most accurate predictions are made by the simplest models. Simpler models can be written as shorter programs, so AI researchers have for a long time equated intelligence with compression [70]. There are many such cases. Regularization is the most common (e.g. dropout [71]). Likewise the Minimum Description Length principle (MDL) lends itself well to selecting hypotheses at high levels of symbolic abstraction [64].

Then there is UAI. The AIXI UAI is a mathematical formalism of superintelligence [63]. It equates simplicity with compressibility, and bases its decisions on the most compressed representations of its history. The length of such an smallest self-extracting archive of a dataset is the Kolmogorov Complexity of the dataset

[72, 73]. Solomonoff induction assigns probabilities to programs based on Kolmogorov complexity [74, 75]<sup>10</sup>. AIXI uses Solomonoff induction to identify the best models of its environment<sup>11</sup>. It can then choose the best possible actions based on those models. Conversely, just as we can say the optimal agent is the one that identifies the simplest models, we can use this Kolmogorov Complexity to measure intelligence<sup>12</sup>. We can check to see if an agent learns one of the simplest models, or how close to simplest its model might be. That is Legg-Hutter intelligence: a measure of intelligence [10]. Chollet's measure, mentioned in the introduction, is similarly based on Kolmogorov Complexity [11]. Unfortunately Kolmogorov Complexity is incomputable, but working *approximations* of both AIXI and Legg-Hutter intelligence exist [78, 79]. These represent the universal upper bounds on intelligence.

Except they don't. AIXI is a case of computational dualism [12]. Complexity is a property of form, not function. Kolmogorov complexity hinges on your choice of Turing machine [17]<sup>13</sup>. In an interactive setting, there is an interpreter between the software mind and the world it inhabits. Complexity need not have any bearing on reality at all. However, there is a correlation between simplicity of form and generalisation of function. There are reasons for this correlation [18]. First, a bounded system can contain only a finite amount of information [80]. Second, a goal-directed process like natural selection can select for systems that make accurate predictions. Third, to make accurate predictions using a finite vocabulary of representations of varying complexity, simpler forms must express more generalisable, *weaker* constraints on functionality [18], which brings us to our third meta-approach.

*W-Maxed:* Computational dualism frames intelligence as a disembodied software policy interacting with the world through an interpreter. The alternative is cognition as a process taking place within the environment, as a part of the environment [81]. This is called *enactive* cognition [82, 83]. A formalisation of enactive cognition must formalise the system as a whole, not parts. This is a challenge. One formalism, does this by moving the problem of interpretation outside the environment [12], basically saying “whatever determines the laws of physics is the interpreter” and assuming it is unknowable<sup>14</sup>. This is useful to examine not just complexity of form and generalisation of function, but the relationship between the two in all possible environments. It was subsequently shown that generalisation stems from weakening the constraints on functionality to be as loose as possible while still satisfying the requirements of the system [28, 84]. Weakness, as it is called, is a measure of function as opposed to form.

<sup>10</sup> Now available for Boolean circuits [76].

<sup>11</sup> Now available for incomputable environments [77].

<sup>12</sup> Oh sure, Kolmogorov Complexity is incomputable. Nobody cares. I can approximate it arbitrarily.

<sup>13</sup> A choice of Turing Machine is a choice of interpreter is a choice of descriptive language.

<sup>14</sup> This is called Stack Theory, because it frames everything as an infinite stack of abstraction layers and assumes the bottom layer is unknowable or does not exist[19].

As such, maximising the weakness of constraints on function (w-maxing) is not mutually exclusive with simp-maxing. Both can take place, and optimising a finite set of representations to express the weakest possible collective constraint on functionality will *cause* simple forms to express weak constraints [18].

Given an abstraction layer (a language), w-maxing involves identifying policies or hypotheses that are as non-specific or ‘weak’ as possible, whilst still satisfying basic requirements. In experiments involving binary arithmetic, w-maxing alone yielded 110–500% improvement in generalisation rate over simp-maxing alone. W-maxing also involves delegating control to lower levels of abstraction. This reflects the biological polycomputational architecture self-organisation [85–87, 44, 88]. Biological systems are comparatively more adaptable than artificial intelligence because biology distributes control and delegates it to lower levels of abstraction [21]. In computer science terms, this is like programming in C instead of Python, or on a field programmable gate array instead of C. More efficient, bespoke implementations are possible when adaptation extends down to smaller scales and lower levels of abstraction. Early stage examples of computing systems that delegate control in this manner include soft-robotics [89] and self-organising systems of nano particles [90, 91]. Because it does not separate software or hardware, it simultaneously optimises for both sample and energy efficiency [21]. But to optimise hardware as above, one must search an infinite space of embodiments. This is a process of trial and error, or optimisation through selection. It could be thought of like a biological self-organising system searching the morphospace during development and regeneration [44]. Because of the enactive frame, w-maxing has been used to explain causal reasoning, language and consciousness in terms that apply to both AI and biological self-organising systems [92–94, 19].

## 5 Conclusion

Recent history has been dominated by scale-maxed approximation. I like to call this period The Embiggening. The rapid improvements suggest we were bottlenecked by compute and data. Now there are diminishing returns [69]. Models are expensive. There are far more problems for which we have little data than problems for which we have a lot of training data. Reliable precision and energy efficiency are increasingly important considerations. Perhaps scaling is no longer the easiest way forward. Better opportunities lie in w-maxing and simp-maxing. Newer models are hybrids like o3, not pure approximations as when GPT-3 was released. There has also been a great deal of discussion around the economic potential of autonomous agents [52]. This is where architectures like AERA, NARS and Hyperon stand to shine. Yes scale-maxed approximations dominate, but a fusion is required for an artificial scientist.

## References

1. Russell, S.: Artificial Intelligence and the Problem of Control, pp. 19–24. Springer Nature (2022)

2. Goertzel, B.: Artificial general intelligence: Concept, state of the art. *Journal of Artificial General Intelligence* **5**(1), 1–48 (2014)
3. Bostrom, N.: Superintelligence: Paths, Dangers, Strategies. Oxford University Press, Oxford, UK (2014)
4. Mitchell, M.: Debates on the nature of artificial general intelligence. *Science* **383**(6689), eado7069 (2024). <https://doi.org/10.1126/science.ado7069>, <https://www.science.org/doi/abs/10.1126/science.ado7069>
5. Krauth, S.J., Coulibaly, J.T., Knopp, S., Traoré, M., N'Goran, E.K., Utzinger, J.: An in-depth analysis of a piece of shit: distribution of Schistosoma mansoni and hookworm eggs in human stool. *PLoS Neglected Tropical Diseases* **6**(12), e1969 (12 2012). <https://doi.org/10.1371/journal.pntd.0001969>, <https://doi.org/10.1371/journal.pntd.0001969>
6. Wang, P.: On defining artificial intelligence. *Journal of Artificial General Intelligence* **10**(2), 1–37 (2019)
7. Goertzel, B.: Generative ai vs. agi: The cognitive strengths and weaknesses of modern llms (2023), <https://arxiv.org/abs/2309.10371>
8. Thorisson, K.R.: A New Constructivist AI: From Manual Methods to Self-Constructive Systems, pp. 145–171. Atlantis Press, Paris (2012)
9. Wang, P.: Rigid Flexibility: The Logic of Intelligence. Applied Logic Series, Springer Nature (2006)
10. Legg, S., Hutter, M.: Universal intelligence: A definition of machine intelligence. *Minds and Machines* pp. 391–444 (2007)
11. Chollet, F.: On the measure of intelligence (2019)
12. Bennett, M.T.: Computational dualism and objective superintelligence. In: Artificial General Intelligence. Springer Nature (2024)
13. Sternberg, R.J.: Toward a triarchic theory of human intelligence. *Behavioral and Brain Sciences* **7**(2), 269–287 (1984). <https://doi.org/10.1017/S0140525X00044629>
14. Bostrom, N.: The superintelligent will: Motivation and instrumental rationality in advanced artificial agents. *Minds and Machines* **22**(2), 71–85 (May 2012). <https://doi.org/10.1007/s11023-012-9281-3>, <https://doi.org/10.1007/s11023-012-9281-3>
15. Orseau, L.: Asymptotic non-learnability of universal agents with neural networks. In: Bach, J., Goertzel, B., Iklé, M. (eds.) Artificial General Intelligence: 5th International Conference, AGI 2012. pp. 234–243. Springer Nature, Berlin, Heidelberg (2012)
16. Orseau, L., Ring, M.: Space-time embedded intelligence. In: Bach, J., Goertzel, B., Iklé, M. (eds.) Artificial General Intelligence. pp. 209–218. Springer Berlin Heidelberg, Berlin, Heidelberg (2012)
17. Leike, J., Hutter, M.: Bad universal priors and notions of optimality. *Proceedings of The 28th Conference on Learning Theory, in Proceedings of Machine Learning Research* pp. 1244–1259 (2015)
18. Bennett, M.T.: Is complexity an illusion? In: Artificial General Intelligence. Springer Nature (2024)
19. Bennett, M.T.: How To Build Conscious Machines. Ph.D. thesis, School of Computing, The Australian National University (2025), [github.com/ViscousLemming/Technical-Appendices](https://github.com/ViscousLemming/Technical-Appendices)
20. Bennett, M.T.: Compression, the fermi paradox and artificial super-intelligence. In: Artificial General Intelligence. pp. 41–44. Springer Nature (2022)
21. Bennett, M.T.: Are biological systems more intelligent than artificial intelligence? (2025), forthcoming

22. Bennett, M.T., Maruyama, Y.: The artificial scientist: Logicist, emergentist, and universalist approaches to artificial general intelligence. In: Goertzel, B., Iklé, M., Potapov, A. (eds.) *Artificial General Intelligence*. pp. 45–54. Springer Nature, Cham (2022)
23. Goertzel, B.: The general theory of general intelligence: A pragmatic patternist perspective. Tech. rep., Singularity Net (2021)
24. Thorisson, K.R., Nivel, E., Steunebrink, B., Helgason, H.P., Pezzulo, G., Sanz, R., Schmidhuber, J., Dindo, H., Rodriguez, M., Chella, A., Jonsson, G.K., Ognibene, D., Corbato-Hernandez, C.: Autonomous acquisition of situated natural communication. *Intl. J. Comp. Sci.& Info. Sys.* (2014)
25. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press, MA (2018)
26. Sutton, R.: The bitter lesson. University of Texas at Austin (2019)
27. Russell, S., Norvig, P.: *Artificial intelligence: A modern approach*, global edition 4th. Pearson, London (2021)
28. Bennett, M.T.: The optimal choice of hypothesis is the weakest, not the shortest. In: *Artificial General Intelligence*. Springer Nature (2023)
29. Bennett, M.T.: A formal theory of optimal learning with experimental results. *Proceedings of the Thirty-fourth International Joint Conference on Artificial Intelligence* (2025)
30. Howard, W.A.: The Formulae-as-Types Notion of Construction. In: Seldin, J., Hindley, J. (eds.) *To H.B. Curry: Essays on Combinatory Logic, Lambda Calculus and Formalism*, pp. 479–490. Academic Press, Cambridge MA (1980)
31. Kautz, H., Selman, B.: Planning as satisfiability. In: IN ECAI-92. pp. 359–363. Wiley, New York (1992)
32. Newell, A., Simon, H.: The logic theory machine—a complex information processing system. *IRE Transactions on Information Theory* **2**(3), 61–79 (1956)
33. Campbell, M., Hoane, A., hsiung Hsu, F.: Deep blue. *Artificial Intelligence* (2002)
34. Kirk, D.: Nvidia cuda software and gpu parallel computing architecture. In: *Proceedings of the 6th International Symposium on Memory Management*. p. 103–104. ISMM '07, Association for Computing Machinery, New York, NY, USA (2007). <https://doi.org/10.1145/1296907.1296909>
35. Schulte, C., Carlsson, M.: Chapter 14 - finite domain constraint programming systems. In: Rossi, F., van Beek, P., Walsh, T. (eds.) *Handbook of Constraint Programming*. Foundations of Artificial Intelligence, Elsevier (2006)
36. Edelkamp, S., Schrödl, S.: Chapter 9 - distributed search. In: Edelkamp, S., Schrödl, S. (eds.) *Heuristic Search*, pp. 369–427. Morgan Kaufmann, San Francisco (2012)
37. Zhou, Y., Zeng, J.: Massively parallel a\* search on a gpu. *Proceedings of the AAAI Conference on Artificial Intelligence* (2015)
38. Oswald, J.T., Rozek, B.: Parallel verification of natural deduction proof graphs. *Electronic Proceedings in Theoretical Computer Science* **396**, 36–51 (Nov 2023). <https://doi.org/10.4204/eptcs.396.4>, <http://dx.doi.org/10.4204/EPTCS.396.4>
39. Krizhevsky et al., A.: Imagenet classification with deep convolutional neural networks. *Commun. ACM* (2017)
40. Vaswani et al., A.: Attention is all you need. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. NIPS'17, Curran, NY (2017)
41. Ribeiro, M.T., Singh, S., Guestrin, C.: "why should i trust you?": Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD Interna-*

- tional Conference on Knowledge Discovery and Data Mining. p. 1135–1144. KDD '16, Association for Computing Machinery, New York, NY, USA (2016)
42. Strubell, E., Ganesh, A., McCallum, A.: Energy and policy considerations for deep learning in NLP. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Florence, Italy (2019)
  43. Bennett, M.T.: Optimal policy is weakest policy. Artificial General Intelligence (2025)
  44. McMillen, P., Levin, M.: Collective intelligence: A unifying concept for integrating biology across scales and substrates. Communications Biology (2024)
  45. Bennett, M.T., Maruyama, Y.: Philosophical specification of empathetic ethical artificial intelligence. IEEE Transactions on Cognitive and Developmental Systems **14**(2), 292–300 (2022)
  46. Silver et al., D.: Mastering the game of go with deep neural networks and tree search. Nature **529**(7587), 484–489 (2016)
  47. Harnad, S.: The symbol grounding problem. Physica D: Nonlinear Phenomena **42**(1), 335–346 (1990)
  48. Garcez, A., Gori, M., Lamb, L.C., Serafini, L., Spranger, M., Tran, S.N.: Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. FLAP (2019)
  49. Garnelo, M., Arulkumaran, K., Shanahan, M.: Towards deep symbolic reinforcement learning (2016)
  50. OpenAI: Openai o3-mini system card (2025)
  51. Trinh, T.H., et al.: Solving olympiad geometry without human demonstrations. Nature (2024)
  52. Perrier, E., Bennett, M.T.: Position: Stop acting like language model agents are normal agents (2025), <https://arxiv.org/abs/2502.10420>
  53. Laird, J.E.: The Soar Cognitive Architecture. MIT Press, MA (2012)
  54. Anderson, J.R., Bothell, D., Byrne, M.D., Douglass, S., Lebiere, C., Qin, Y.: An integrated theory of the mind. Psychological Review (2004), because apparently six authors are needed to figure out how your brain works
  55. et al., B.G.: Opencog hyperon: A framework for agi at the human level and beyond. Tech. rep., OpenCog Foundation (2023)
  56. Friston, K.: The free-energy principle: A unified brain theory? Nature Reviews Neuroscience **11**(2), 127–138 (2010)
  57. Goertzel, B.: Actpc-chem: Discrete active predictive coding for goal-guided algorithmic chemistry as a potential cognitive kernel for hyperon and primus-based agi (2024)
  58. Nivel et al., E.: Autocatalytic endogenous reflective architecture. Tech. rep., Reykjavik University, School of Computer Science (2013)
  59. Thórisson, K.R.: Seed-programmed autonomous general learning. In: Proceedings of the First International Workshop on Self-Supervised Learning. Proceedings of Machine Learning Research, vol. 131, pp. 32–61. PMLR (27–28 Feb 2020)
  60. Sheikhlar, A., Thorisson, K.R.: Causal generalization via goal-driven analogy. In: Thorisson, K.R., Isaev, P., Sheikhlar, A. (eds.) Artificial General Intelligence. pp. 165–175. Springer Nature Switzerland, Cham (2024)
  61. Eberding, L.M., Thompson, J., Thorisson, K.R.: Argument-driven planning and autonomous explanation generation. In: Thorisson, K.R., Isaev, P., Sheikhlar, A. (eds.) Artificial General Intelligence. pp. 73–83. Springer Nature Switzerland, Cham (2024)

62. Hammer, P., Lofthouse, T.: ‘opennars for applications’: Architecture and control. In: Goertzel, B., Panov, A.I., Potapov, A., Yampolskiy, R. (eds.) *Artificial General Intelligence*. pp. 193–204. Springer Nature, Cham (2020)
63. Hutter, M., Quarel, D., Catt, E.: *An Introduction to Universal Artificial Intelligence*. Chapman and Hall/CRC, 1st edn. (2024). <https://doi.org/10.1201/9781003460299>
64. Rissanen, J.: Modeling by shortest data description. *Automatica* (1978)
65. Ingesson, T.: *The Politics of Combat: The Political and Strategic Impact of Tactical-Level Subcultures, 1939-1995*. Doctoral thesis (monograph), Department of Political Science (2016)
66. Sober, E.: *Ockham’s Razors: A User’s Manual*. Cambridge Uni. Press (2015). <https://doi.org/10.1017/CBO9781107705937>
67. Roose, K.: Why a conversation with bing’s chatbot left me deeply unsettled. *The New York Times* (February 2023), <https://www.nytimes.com/2023/02/16/technology/why-a-conversation-with-bings-chatbot-left-me-deeply-unsettled.html>
68. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., Bridgland, A., Meyer, C., Kohl, S.A.A., Ballard, A.J., Cowie, A., Romera-Paredes, B., Nikolov, S., Jain, R., Adler, J., Back, T., Petersen, S., Reiman, D., Clancy, E., Zielinski, M., Steinegger, M., Pacholska, M., Berghammer, T., Bodenstein, S., Silver, D., Vinyals, O., Senior, A.W., Kavukcuoglu, K., Kohli, P., Hassabis, D.: Highly accurate protein structure prediction with alphafold. *Nature* (2021)
69. Kaplan, J., McCandlish, S., Henighan, T., Brown, T.B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., Amodei, D.: Scaling laws for neural language models (2020)
70. Chaitin, G.J.: On the length of programs for computing finite binary sequences. *J. ACM* (1966)
71. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**(56), 1929–1958 (2014), <http://jmlr.org/papers/v15/srivastava14a.html>
72. Kolmogorov, A.: On tables of random numbers. *Sankhya: The Indian Journal of Statistics* **A**, 369–376 (1963)
73. Kolmogorov, A.N.: Three approaches to the quantitative definition of information \*. *International Journal of Computer Mathematics* **2**(1-4), 157–168 (1968). <https://doi.org/10.1080/00207166808803030>
74. Solomonoff, R.: A formal theory of inductive inference. part i. *Information and Control* **7**(1), 1–22 (1964)
75. Solomonoff, R.: A formal theory of inductive inference. part ii. *Information and Control* **7**(2), 224–254 (1964)
76. Wyeth, C., Sturtivant, C.: A circuit complexity formulation of algorithmic information theory. *Physica D: Nonlinear Phenomena* **456**, 133925 (2023). <https://doi.org/https://doi.org/10.1016/j.physd.2023.133925>, <https://www.sciencedirect.com/science/article/pii/S0167278923002798>
77. Oswald, J.T., Ferguson, T.M., Bringsjord, S.: Extension to legg and hutter’s universal intelligence measure for uncomputable environments. In: Thórisson, K.R., Isaev, P., Sheikhlar, A. (eds.) *Proceedings of the 17th International Conference on Artificial General Intelligence. Lecture Notes in Computer Science*, vol. 14951, pp. 134–144. Springer, Cham (2024). [https://doi.org/10.1007/978-3-031-65572-2\\_15](https://doi.org/10.1007/978-3-031-65572-2_15)
78. Hutter, M.: *Universal Algorithmic Intelligence: A Mathematical Top→Down Approach*, pp. 227–290. Springer Berlin Heidelberg, Berlin, Heidelberg (2007)

79. Legg, S., Veness, J.: An approximation of the universal intelligence measure. In: Algorithmic Probability and Friends (2011)
80. Bekenstein, J.D.: Universal upper bound on the entropy-to-energy ratio for bounded systems. *Phys. Rev. D* **23**, 287–298 (Jan 1981)
81. Dreyfus, H.L.: What Computers Can't Do: A Critique of Artificial Reason. Harper & Row (1972)
82. Thompson, E.: Mind in Life: Biology, Phenomenology, and the Sciences of Mind. Harvard University Press, Cambridge MA (2007)
83. Vervaeke, J., Lilliacrap, T., Richards, B.: Relevance realization and the emerging framework in cognitive science. *J. Log. Comput.* (2012)
84. Bennett, M.T.: Computable Artificial General Intelligence. Preprint (2022)
85. Bongard, J., Levin, M.: There's plenty of room right here: Biological systems as evolved, overloaded, multi-scale machines. *Biomimetics* **8**(1) (2023)
86. Gershenson, C.: Self-organizing systems: what, how, and why? *npj Complexity* (2025)
87. Fields, C., Levin, M.: Scale-free biology: Integrating evolutionary and developmental thinking. *BioEssays* **42** (06 2020)
88. Solé, R., Seoane, L.F.: Evolution of brains and computers: The roads not taken. *Entropy* **24**(5), 665 (2022)
89. Man, K., Damasio, A.R.: Homeostasis and soft robotics in the design of feeling machines. *Nature Machine Intelligence* **1**, 446 – 452 (2019), <https://api.semanticscholar.org/CorpusID:208089594>
90. Borghi, F., Nieuw, T.R., Galli, D.E., Milani, P.: Brain-like hardware, do we need it? *Frontiers in Neuroscience* **18** (2024)
91. Paroli, B., Martini, G., Potenza, M., Siano, M., Mirigliano, M., Milani, P.: Solving classification tasks by a receptron based on nonlinear optical speckle fields. *Neural Networks* **166**, 634–644 (2023)
92. Bennett, M.T.: Emergent causality and the foundation of consciousness. In: Artificial General Intelligence. Springer Nature (2023)
93. Bennett, M.T.: On the computation of meaning, language models and incomprehensible horrors. In: Artificial General Intelligence. Springer Nature (2023)
94. Bennett, M.T., Welsh, S., Ciaunica, A.: Why Is Anything Conscious? Preprint (2024)