

---

# Boundless Socratic Learning with Language Games

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 An agent trained within a closed system can master any desired capability, as long  
2 as the following three conditions hold: (a) it receives sufficiently informative and  
3 aligned feedback, (b) its coverage of experience/data is broad enough, and (c) it  
4 has sufficient capacity and resource. We justify these conditions and consider what  
5 limitations arise from (a) and (b) in closed systems, when assuming that (c) is not  
6 a bottleneck. Considering the special case of homoiconic agents with matching  
7 input and output spaces (namely, language), we argue that such pure recursive self-  
8 improvement, dubbed ‘Socratic learning,’ can boost performance vastly beyond  
9 what is present in its initial data or initial knowledge, and is only limited by time,  
10 as well as gradual misalignment concerns. Furthermore, we propose a constructive  
11 framework to implement it, based the notion of language games.

## 12 1 Introduction

13 On the path between now and artificial superhuman intelligence [ASI; 11] lies a tipping point, namely  
14 when the bulk of a system’s improvement in capabilities is driven by *itself* instead of human sources  
15 of data, labels, or preferences (which can only scale so far). As yet, few systems exhibit such *recursive*  
16 *self-improvement*, so now may be a prudent time to discuss and characterize what it is, and what  
17 it entails. We focus on one end of the spectrum, the clearest is not the most practical one, namely  
18 pure self-contained settings of ‘Socratic’ learning, closed systems without the option to collect new  
19 information from the external world. We articulate conditions, pitfalls and upper limits, as well as  
20 a concrete path towards them that builds on the notion of language games. The central aim of this  
21 brief position paper is to clarify terminology and frame the discussion, with an emphasis on the long  
22 run. It is not to propose new algorithms, nor survey past literature; we pay no heed to near-term  
23 feasibility or constraints. We start with a flexible and general framing, and refine and instantiate these  
24 definitions over the course of the paper.

25 **Definitions** Consider a *closed system* (no inputs, no outputs) that evolves over time. Within the  
26 system is an entity with inputs and outputs, called *agent*, that also changes over time. External to the  
27 system is an *observer* whose purpose is to assess the *performance* of the agent. If performance keeps  
28 increasing, we call this system-observer pair an *improvement process*.

29 The dynamics of this process are driven by both the agent and its surrounding system, but setting  
30 clear agent boundaries is required to make evaluation well-defined: in fact an agent *is* what can be  
31 unambiguously evaluated. Similarly, for separation of concerns, the observer is deliberately located  
32 outside of the system: As the system is closed, the observer’s assessment cannot feed back into the  
33 system. Hence, the agent’s learning feedback must come from system-internal *proxies* such as losses,  
34 reward functions, or critics.

35 The simplest type of performance metric is a *scalar* score that can be measured in finite time, that  
36 is, on (an aggregation of) episodic tasks. Mechanistically, the observer can measure performance in

37 two ways, by *passively* observing the agent’s behaviour within the system (if all pertinent tasks occur  
38 naturally), or by *copy-and-probe* evaluations where it confronts a copy of the agent with interactive  
39 tasks of its choosing.

40 Without loss of generality, we distinguish three types of elements within an agent; *fixed* elements are  
41 unaffected by learning, such as its substrate or unmodifiable code (genotype). *Transient* elements  
42 do not carry over between episodes, or across to evaluation (e.g., activations, the state of a random  
43 number generator). And finally *learned* elements (e.g., weights, parameters, knowledge) change  
44 based on a feedback signal, and their evolution maps to performance differences.

45 We can distinguish improvement processes by their implied lifetime; some are *open-ended* and keep  
46 improving without limit [7], while others converge onto their asymptotic performance after some  
47 finite time. Note that neither case needs to invoke a notion of optimality.

## 48 2 Three Necessary Conditions for Self-improvement

49 *Self-improvement* is an improvement process as defined above, but with the additional criterion that  
50 the agent’s own outputs (actions) influence its future learning. In other words, systems in which  
51 agents shape (some of) their own experience stream, potentially enabling unbounded improvement in  
52 a closed system. This setting may look familiar to readers from the reinforcement learning community  
53 [RL; 18], who build agents whose behaviour changes the data distribution it learns on, which in turn  
54 affects its behaviour policy, and so on. Another prototypical instance of a self-improvement process  
55 is *self-play*, where the system (often a symmetric game) slots the agent into the roles of both player  
56 and opponent, to generate an unlimited experience stream annotated with feedback (who won) that  
57 provides direction for ever-increasing skill-learning.

58 From its connection to RL, we can derive necessary conditions for self-improvement to work, and  
59 help clarify some assumptions about the system. The first two conditions, feedback and coverage, are  
60 about feasibility in principle, the third (capacity) is about practice.

61 **Feedback** Feedback is what gives direction to learning, without it, the process is merely one of  
62 self-modification. Feedback must have two properties for self-improvement to work, one fundamental,  
63 one practical. First, system-internal feedback must be *aligned* with the external observer, and remain  
64 aligned throughout the process. This places a significant burden on the system at set-up time, with the  
65 most common pitfall being a poorly designed critic or reward function that becomes exploitable over  
66 time, deviating from the observer’s intent. RL’s famed capability for *self-correction* is not applicable  
67 here: what can self-correct is behaviour given feedback, but not feedback itself. Second, the efficiency  
68 criterion for feedback is that it be reliable enough, and contain enough information (not too sparse,  
69 not too noisy, not too delayed) for learning to be feasible within the time horizon of the system.

70 **Coverage** By definition, a self-improving agent determines the distribution of data it learns from.  
71 To prevent issues like collapse, drift, exploitation or overfitting, it needs to preserve<sup>1</sup> coverage of the  
72 data distribution everywhere the observer cares about. In most interesting cases, where performance  
73 includes a notion of generalisation, that target distribution is not given (the test tasks are withheld),  
74 so the system needs to be set up to intrinsically seek ‘sufficient’ coverage, a sub-process classically  
75 called *exploration*.

76 **Capacity** The research field of RL has produced a lot of detailed knowledge about how to train  
77 agents, which algorithms work in which circumstances, an abundance of neat tricks that address  
78 practical concerns, as well as theoretical results that characterize convergence, rates of progress,  
79 etc. It would be futile to try and summarize such a broad body of work here. However, one general  
80 observation that matters for our argument is that ‘RL works at scale’: in other words, when scaling  
81 up experience and compute sufficiently, even relatively straightforward RL algorithms can solve  
82 problems previously thought out of reach [high-profile examples include: 19, 10, 15, 16, 21, 1]. For  
83 any specific, well-defined practical problem, the details matter (and differ), and greatly impact the  
84 efficiency of the learning dynamics; but the asymptotic outcome seems a foregone conclusion.

---

<sup>1</sup>This may entail conditions on how the system is initialised, as the agent needs to see a first set of inputs before it can produce its own.

### 85 3 Socratic Learning

86 The specific type of self-improvement process we consider here is *recursive self-improvement*, where  
87 the agent’s inputs and outputs are *compatible* (i.e., live in the same space), and outputs become future  
88 inputs.<sup>2</sup>This is more restrictive but less mediated than the general case where outputs merely influence  
89 the input distribution (but it is less restrictive than homoiconic self-modification and self-referential  
90 systems). This type of recursion is an attribute of many open-ended processes, and open-ended  
91 improvement is arguably a central feature of ASI [see 7].

92 An excellent example of such a compatible space of inputs and outputs is *language*. A vast range of  
93 human behaviours are mediated by, and well-expressed<sup>3</sup> in language, especially in cognitive domains  
94 (which are definitionally part of ASI). As argued by [4], language may well be sufficient for thinking  
95 and understanding, and not require sensory grounding. Plus, language has the neat property of being  
96 a *soup of abstractions*, encoding many levels of the conceptual hierarchy in a shared space. A related  
97 feature of language is its extendability, i.e., developing new languages within it, such as formal  
98 mathematics or programming languages. While special-purpose tools for these are important for  
99 efficiency, natural language may be sufficient as a basis: just like humans can reason ‘manually’  
100 through mathematical expressions when doing mental arithmetic, so can natural language agents  
101 [12]. And of course, it does not hurt that AI competence on language domains has radically improved  
102 recently, with a lot of momentum since the rise of LLMs.

103 For the remainder of the paper, we will use ‘*Socratic learning*’ to refer to a recursive self-improvement  
104 process that operates in language space. The name is alluding to Socrates’ approach of finding or  
105 refining knowledge through questioning dialogue and repeated language interactions, but, notably,  
106 without going out to collect observations in the real world—mirroring our emphasis on the system  
107 being closed. We encourage the reader to imagine an unbroken process of deliberation among a  
108 circle of philosophers, maybe starting with Socrates and his disciples, but expanding and continuing  
109 undisturbed for millennia: what cultural artifacts, what knowledge, what wisdom could such a process  
110 have produced by now? And then, consider a question that seems paradoxical at first: how can a  
111 closed system produce open-ended improvement?

#### 112 The Limits of Socratic Learning

113 Revisiting the necessary conditions for self-improvement, we can derive some insights on how  
114 Socratic learning is limited *in principle*. For that, we can mostly sidestep the capacity concerns  
115 of Section 2, by choosing one of two premises. Either, we can assume that compute and memory  
116 constraints are but a temporary obstacle, as they keep growing exponentially, so ignoring them  
117 still produces valid high-level insights. Or, we can consider the resource-constrained scenario and  
118 study feasibility within the class of such restricted systems. The other two conditions, coverage and  
119 feedback, remain irreducible however. The system has to keep generating (language) data, while  
120 preserving or expanding diversity over time. In the LLM age, we can envision a generative agent  
121 initialized with a very broad internet-like distribution, but preventing drift, collapse or just narrowing  
122 of that distribution in a recursive process may be highly non-trivial [14].

123 The other requirement is for the system to continue producing feedback about (some subset of)  
124 the agent’s outputs, which structurally requires a critic that can assess language, and that remains  
125 sufficiently aligned with the observer’s evaluation metric. This is challenging for a number of reasons:  
126 Well-defined, grounded metrics in language space are often limited to narrow tasks, while more  
127 general-purpose mechanisms like AI-feedback are exploitable, especially so if the input distribution  
128 is permitted to shift. For example, none of the current LLM training paradigms have a feedback  
129 mechanism that is sufficient for Socratic learning. Next-token prediction loss is grounded, but  
130 insufficiently aligned with downstream usage, and unable to extrapolate beyond the training data.  
131 Human preferences are aligned by definition, but prevent learning in a closed system. Caching such  
132 preferences into a learned reward model makes it self-contained, but exploitable, and misaligned in  
133 the long-run, as well as weak on out-of-distribution data.

---

<sup>2</sup>Or at least some of them are fed back. Input and output spaces are not necessarily identical, but they intersect. For example, the agent could be generating code, but perceive natural language, (partly self-generated) code, and execution traces.

<sup>3</sup>“Whereof one cannot speak, thereof one must be silent.” [23]

## 134 4 Language Games Are All You Need ...

135 Fortunately, language, learning and grounding are well-studied topics. A particularly useful concept  
136 for us to draw on is Wittgenstein’s notion of *language games*.<sup>4</sup> For him, it is not the words that  
137 capture meaning, but only the interactive nature of language can do so. To be concrete here, define  
138 a language game as an *interaction protocol* (a set of rules, expressible in code) that specifies the  
139 interaction of one or more agents (‘players’) that have language inputs and language outputs, plus a  
140 scalar *scoring function* for each player at the end of the game.<sup>5</sup>

141 Language games, thus defined, address the primary needs of Socratic learning; namely, they provide  
142 a scalable mechanism for unbounded interactive data generation and self-play, while automatically  
143 providing an accompanying feedback signal (the score). In fact, they are the logical consequence  
144 of the coverage and feedback conditions, almost tautologically so: there is no form of interactive  
145 data generation with tractable feedback that is not a language game. As a bonus, seeing the process  
146 as one of *game-play* immediately brings in the potential of rich strategic diversity arising from  
147 multi-agent dynamics [as spelled out in depth in 8, 6], which is likely to address at least part of the  
148 coverage condition. Pragmatically too, games are a great way to get started, given the vast human  
149 track record of creating and honing a vast range of games and player skills [3]. A number of common  
150 LLM interaction paradigms are also well represented as language games, for example debate [9, 5],  
151 role-play [20], jailbreak defense [25], or outside of closed systems, paradigms like RL from human  
152 feedback [RLHF, 13, 2].

### 153 ...If You Have Enough of Them ...

154 Returning to our circle of deliberating philosophers: is there any *one* language game we could imagine  
155 them playing for millennia? Instead, maybe, they are more likely to escape a narrow outcome when  
156 playing *many* language games. It turns out that Wittgenstein (him again) proposed this same idea:  
157 he adamantly argued against language having a singular essence or function.<sup>6</sup> Using many narrow  
158 but well-defined language games instead of a single universal one resolves a key dilemma: For each  
159 narrow game, a reliable score function (or critic) can be designed, whereas getting the single universal  
160 one right is more elusive [even if possible in principle, as argued by 17].<sup>7</sup> From that lens, the full  
161 process of Socratic learning is then a *meta-game*, which schedules the language games that the agent  
162 plays and learns from.

### 163 ...And You Play the Right Ones

164 Socrates was famously sentenced to death and executed for ‘corrupting the youth.’ We can take  
165 this as a hint that a Socratic process is not guaranteed to remain aligned with external observers’  
166 intent. Language games as a mechanism do not side-step this either, but they arguably reduce the  
167 precision needed: instead of a critic that is aligned at the fine granularity of individual inputs and  
168 outputs, all that is needed is a ‘meta-critic’ that can judge which games should be played: it may  
169 be that no individual language game is perfectly aligned, but it is doable to filter the many games  
170 according to whether they make a net-positive contribution (when played and learned about). This  
171 kind of structural leniency is precisely what gives it the potential to scale.

172 Stepping out of our assumption of the closed system for a moment: when we actually build ASI, we  
173 will almost surely want to not optimistically trust that alignment is preserved, but instead continually  
174 check the process as carefully as possible, and probably intervene and adjust throughout the training  
175 process. In that case, explicitly exposing the distribution of games (accompanied with per-game  
176 learning curves) as knobs to the designer may be a useful level of abstraction.

<sup>4</sup>“I shall also call the whole, consisting of language and the actions into which it is woven, the ‘language-game’.” [24]

<sup>5</sup>For simplicity, assume that games are guaranteed to terminate in finite time.

<sup>6</sup>“But how many kinds of sentence are there? Say assertion, question, and command?—There are *countless* kinds: countless different kinds of use of what we call ‘symbols,’ ‘words,’ ‘sentences.’ And this multiplicity is not something fixed, given once for all; but new types of language, new language-games, as we may say, come into existence, and others become obsolete and get forgotten.” [24], emphasis in original.

<sup>7</sup>But, as a prescient Norbert Wiener was warning seven decades ago: “The machines will do what we ask them to do and not what we ought to ask them to do. [...] We can be humble and live a good life with the aid of the machines, or we can be arrogant and die.” [22].

## References

- 177
- 178 [1] T. AlphaProof and T. AlphaGeometry. AI achieves silver-medal standard solving International  
179 Mathematical Olympiad problems. *DeepMind blog*, 2024.
- 180 [2] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli,  
181 T. Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from  
182 human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 183 [3] E. Berne. *Games people play: The psychology of human relationships*, volume 2768. Penguin  
184 Uk, 1968.
- 185 [4] D. J. Chalmers. Does thought require sensory grounding? From pure thinkers to large language  
186 models. *arXiv preprint arXiv:2408.09605*, 2024.
- 187 [5] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch. Improving factuality and reasoning  
188 in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.
- 189 [6] E. A. Duéñez-Guzmán, S. Sadedin, J. X. Wang, K. R. McKee, and J. Z. Leibo. A social path to  
190 human-like artificial intelligence. *Nature Machine Intelligence*, 5(11):1181–1188, 2023.
- 191 [7] E. Hughes, M. Dennis, J. Parker-Holder, F. Behbahani, A. Mavalankar, Y. Shi, T. Schaul, and  
192 T. Rocktäschel. Open-endedness is essential for artificial superhuman intelligence. *arXiv  
193 preprint arXiv:2406.04268*, 2024.
- 194 [8] J. Z. Leibo, E. Hughes, M. Lanctot, and T. Graepel. Autocurricula and the emergence of  
195 innovation from social interaction: A manifesto for multi-agent intelligence research. *arXiv  
196 preprint arXiv:1903.00742*, 2019.
- 197 [9] T. Liang, Z. He, W. Jiao, X. Wang, Y. Wang, R. Wang, Y. Yang, Z. Tu, and S. Shi. Encourag-  
198 ing divergent thinking in large language models through multi-agent debate. *arXiv preprint  
199 arXiv:2305.19118*, 2023.
- 200 [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,  
201 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-  
202 forcement learning. *Nature*, 518(7540):529–533, 2015.
- 203 [11] M. R. Morris, J. Sohl-Dickstein, N. Fiedel, T. Warkentin, A. Dafoe, A. Faust, C. Farabet,  
204 and S. Legg. Levels of AGI: Operationalizing progress on the path to AGI. *arXiv preprint  
205 arXiv:2311.02462*, 2023.
- 206 [12] T. OpenAI o1. Learning to reason with LLMs. *OpenAI blog*, 2024.
- 207 [13] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,  
208 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.  
209 *Advances in neural information processing systems*, 35:27730–27744, 2022.
- 210 [14] H. Shi, Z. Xu, H. Wang, W. Qin, W. Wang, Y. Wang, and H. Wang. Continual learning of large  
211 language models: A comprehensive survey. *arXiv preprint arXiv:2404.16789*, 2024.
- 212 [15] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser,  
213 I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of Go with deep  
214 neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 215 [16] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre,  
216 D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess,  
217 shogi, and Go through self-play. *Science*, 362(6419):1140–1144, 2018.
- 218 [17] D. Silver, S. Singh, D. Precup, and R. S. Sutton. Reward is enough. *Artificial Intelligence*,  
219 299:103535, 2021.
- 220 [18] R. S. Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 221 [19] G. Tesauro et al. Temporal difference learning and td-gammon. *Communications of the ACM*,  
222 38(3):58–68, 1995.

- 223 [20] A. S. Vezhnevets, J. P. Agapiou, A. Aharon, R. Ziv, J. Matyas, E. A. Duéñez-Guzmán, W. A.  
224 Cunningham, S. Osindero, D. Karmon, and J. Z. Leibo. Generative agent-based modeling  
225 with actions grounded in physical, social, or digital space using concordia. *arXiv preprint*  
226 *arXiv:2312.03664*, 2023.
- 227 [21] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi,  
228 R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in StarCraft II using multi-agent  
229 reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 230 [22] N. Wiener. The machine age / In 1949, he imagined an age of robots. *MIT Archives / The New*  
231 *York Times*, D:8, 1949 / 2013.
- 232 [23] L. Wittgenstein. *Tractatus Logico-Philosophicus*. 1921.
- 233 [24] L. Wittgenstein. *Philosophical investigations*. 1953.
- 234 [25] Y. Zeng, Y. Wu, X. Zhang, H. Wang, and Q. Wu. Autodefense: Multi-agent llm defense against  
235 jailbreak attacks. *arXiv preprint arXiv:2403.04783*, 2024.