

---

# Bitter Lesson of the ARC-AGI Challenge: Intelligence may look very different in machines and humans

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The Abstraction and Reasoning Corpus (ARC) and the associated ARC-AGI chal-  
2 lenge are benchmarks for evaluating core reasoning skills. Recently, OpenAI’s  
3 new model *o3* solved ARC-AGI, reigniting debate on whether such achievements  
4 reflect true reasoning or mere pattern matching. Rich Sutton’s “Bitter Lesson”  
5 suggests that general methods at scale outperform specialized ones; in ARC, top  
6 solutions heavily rely on data augmentation. We examine this tension and propose  
7 that new concepts or metaphors may be needed to describe machine reasoning.  
8 This position paper argues that reasoning and intelligence in machines can differ  
9 fundamentally from human cognition: recognizing that intelligence is task-specific  
10 and diverse across animals, humans, and machines could lead to more appropriate  
11 benchmarks.

## 12 1 Introduction

13 The Abstraction and Reasoning Corpus (ARC) and the associated ARC-AGI challenge Chollet [2019]  
14 are benchmarks for core reasoning and abstraction in AI. These tasks require minimal prior knowledge  
15 and emphasize generalisation. Recently, OpenAI’s new model *o3* solved many ARC-AGI puzzles,  
16 reigniting debate on whether such achievements reflect genuine reasoning or mere pattern-matching.  
17 Rich Sutton’s “Bitter Lesson” Sutton [2019] suggests that general methods at scale outperform  
18 specialized ones; indeed, many ARC solutions rely on data augmentation and test-time heuristics.

19 **This position paper argues that reasoning and intelligence in machines may differ fundamentally**  
20 **from those in humans. Recognizing intelligence as task-specific and diverse among animals,**  
21 **humans, and machines may lead to more appropriate benchmarks.**

## 22 2 Defining Artificial General Intelligence (AGI)

23 ARC emphasises adaptability to novel problems, aligning with the definition of intelligence as  
24 efficiency of skill acquisition across tasks Chollet [2019]. ARC tasks present input-output puzzles  
25 on small colored grids (each cell one of ten colors); the goal is to produce a correct output grid  
26 from limited examples, demonstrating generalisation. These puzzles rely on basic core priors (object  
27 permanence, goal-directedness, counting, basic geometry) that humans naturally possess. ARC  
28 avoids reliance on acquired cultural knowledge (e.g., language) to ensure a fair comparison between  
29 human and machine intelligence.

30 OpenAI’s *o3* model recently achieved around 87% accuracy on the ARC test, exceeding the 85%  
31 threshold often associated with artificial general intelligence (AGI). However, much of the evaluation  
32 set had previously been accessible through APIs, raising concerns about data contamination. As a  
33 result, it is still uncertain whether ARC-AGI has been solved sincerely.

34 Goodhart’s Law states that once a measure becomes a target, it ceases to be a reliable measure. If  
35 ARC is optimized through data augmentation or heuristics, its ability to assess general reasoning may  
36 diminish. To address this, we propose testing robustness by introducing task variants: for example,  
37 modifying input-output examples or creating concept variants (e.g., ConceptARC tasks Mitchell et al.  
38 [2023]) to ensure methods capture underlying principles rather than shortcuts.

### 39 **3 Machines, reasoning, and the need for neologisms**

40 A key question is what constitutes genuine reasoning in machines. Are we simply projecting human  
41 metaphors onto machine behaviour? Terms like “reasoning”, “understanding”, and “abstraction” are  
42 laden with human-centric assumptions. Machines may operate by entirely different principles. To  
43 describe these emergent behaviours, new terms (neologisms) may be needed: for example, we could  
44 reconceptualize “artificial reasoning” as “mechanical abstraction” or “synthetic inference”.

### 45 **4 Generalisability to other datasets and other domains**

46 AI systems that perform well on ARC tasks may still rely on task-specific shortcuts. To truly evaluate  
47 generalisation, such models should be tested on related tasks in different domains. For example, a  
48 system that solves an object-centric ARC puzzle could be tested on a maze navigation task requiring  
49 the same objectness reasoning (e.g., the PUZZLES benchmark Estermann et al. [2024]). Likewise,  
50 solving an ARC permutation task should translate to solving word-based permutation puzzles. By  
51 evaluating models on such analogous tasks, we can assess whether the underlying reasoning principles  
52 generalise.

### 53 **5 A call for broader agreement on intelligence**

54 “Intelligence” remains a “suitcase word” with many conflicting definitions. AI researchers often  
55 use functional definitions (e.g., goal achievement or optimisation), while cognitive scientists and  
56 neuroscientists may emphasise creativity, emotion, or embodied cognition. This ambiguity leads  
57 to talking past each other about AI’s capabilities. We suggest either broadening the definition  
58 (e.g., intelligence as the ability to adapt to complex environments) or making it more precise by  
59 categorising aspects (such as problem-solving, social, and creative intelligence). Interdisciplinary  
60 dialogue among AI researchers, psychologists, and philosophers is critical to develop clear, context-  
61 dependent definitions.

### 62 **6 Artificial general intelligence is not well defined**

63 Artificial General Intelligence (AGI) itself has no universally accepted definition Mitchell [2024a].  
64 ARC-AGI is explicitly a research tool (not a definitive test) focusing on abstraction and reasoning.  
65 Critics note that focusing only on “cognitive tasks” ignores embodiment: a system may excel at  
66 language but not perform physical tasks (“brain in a vat”). Some have suggested rebranding “artificial  
67 intelligence” to “actual intelligence” Mitchell [2019] to reduce anthropomorphism and acknowledge  
68 that machines, animals, and humans may each exhibit different forms of intelligence.

### 69 **7 There may be a spectrum of reasoning in LLMs**

70 Reasoning and intelligence may lie on a spectrum. It is possible that reasoning in LLMs lies  
71 somewhere along this spectrum. Reasoning might look very different in machines and humans. The  
72 word reasoning comes loaded with many assumptions: it is a metaphor and unfortunately these  
73 metaphors are used repeatedly in the field of AI Mitchell [2024b].

### 74 **8 Alternate View**

75 While this paper argues that intelligence and reasoning in machines and humans are fundamentally  
76 different and need distinct definitions, alternative perspectives expand upon this view. An alternative

77 view asserts that human intelligence and reasoning should remain the gold standard for defining and  
78 evaluating machine intelligence. This position is often implicit in benchmarks like ARC, which are  
79 designed to test AI systems on tasks that are meant to reflect human cognitive abilities. Proponents  
80 argue that the ultimate goal of AI research should be to replicate and surpass human cognitive skills  
81 in a way that aligns closely with human definitions of intelligence, as this ensures practical utility in  
82 domains like healthcare, education, and governance.

83 While anthropocentric definitions have practical utility, they risk constraining our understanding of  
84 intelligence to human-specific parameters. This ignores the potential for machines to exhibit forms of  
85 intelligence that humans cannot. Such a narrow focus may lead to inefficient or misguided approaches  
86 to AI development, as it prioritizes mimicking human cognition rather than leveraging the unique  
87 strengths of machines. *By redefining intelligence to include machine-specific capabilities, we can*  
88 *create benchmarks and evaluation criteria that account for the distinct ways in which machines*  
89 *operate, fostering innovation beyond anthropocentric constraints.* ARC itself could be adapted  
90 to include tasks that explore non-human-centric reasoning, thus broadening the scope of what we  
91 consider “intelligent”.

92 The current situation with AI benchmarks is depicted in Figure 1, where a human is asking a seal, a  
93 penguin, a dog and a shark to climb up a tree. This ignores the fact that intelligence is task-specific:  
94 organisms have evolved over millions of years to solve tasks that are suited for the kind of environment  
95 they find themselves in. This frequently involves tradeoffs: a penguin will trade-off flight and efficient  
96 locomotion on land for swimming efficiently under water. A human (and other primates) have evolved  
97 to walk efficiently on land at the expense of not being able to fly.

98 *Recognizing that intelligence is task-specific and comes in various forms and guises may help us*  
99 *appreciate that animals, humans and machines have different kinds of intelligence, and help us design*  
100 *better benchmarks for machines.*



Figure 1: Our current benchmarks for AI are very anthropocentric. It is similar to a situation illustrated here, where a human is asking a seal, a penguin, a dog and a shark to all climb a tree. The implication is that each organism has its own strengths which have been honed through millions of years of evolution to solve certain tasks that are specific to the environment it lives in. Image generated using DALL-E.

## 101 **9 Discussion**

102 This paper examines the conflict between ARC’s objectives and the employed methods and reflects  
103 on whether intelligence and reasoning might ultimately look very different in machines and humans.  
104 Some key takeaways are:

- 105 1. Goodhart’s law: As ARC becomes a target for optimization, its utility as a benchmark may  
106 diminish unless new task variations are introduced.
- 107 2. Neologisms and new metaphors: We must reconsider the language we use to describe  
108 machine behaviours. Terms like “reasoning” and “understanding” may need redefinition to  
109 account for the unique characteristics of machine intelligence.
- 110 3. Unified theory: A unified approach to meaning and reasoning could help bridge the gap  
111 between human cognition and machine abstraction.
- 112 4. Reasoning and intelligence may be very different in machines and humans.

113 At the risk of stating the obvious, let us go over how we humans solve these visual puzzles. We solve  
114 them in part using our visual apparatus. Specifically, light reflected from these puzzles impinges on  
115 our visual system, where part of the initial processing, such as edge detection, occurs in the eye itself.  
116 The remaining processing is carried out in the brain, which integrates this sensory input with a vast  
117 amount of prior knowledge, including concepts such as objectness, gravity, etc.

118 For machines, the processing happens in a fundamentally different way. In the case of large language  
119 models, the input is ingested as tokens in a continuous one-dimensional stream, devoid of any direct  
120 visual experience. However, if the system is a convolutional neural network, the processing may  
121 mimic certain aspects of the human visual system, such as feature detection, but it still lacks the  
122 inherent embodiment and experiential learning process of humans.

123 Machines are also not embodied in the physical world, meaning that they cannot acquire the same  
124 training data or learn through interaction and sensory experiences as humans. Their learning is  
125 based entirely on the data they are trained on, which is inherently limited compared to the diverse  
126 and continuous experiences humans undergo. Hence, to expect that reasoning and intelligence in  
127 machines will mirror those of humans when processing the abstraction and reasoning corpus is a  
128 fallacy. Differences in sensory processing, embodiment, and prior knowledge result in different  
129 pathways of reasoning and intelligence in humans and machines. In other words, there can be many  
130 different ways to solve the same problem. *Humans and machines can solve similar problems in*  
131 *fundamentally different ways.*

132 The computer scientist Edsger Dijkstra, when asked whether computers can think like humans,  
133 famously responded with a counter-question: “Do we think submarines swim like fish?” This analogy  
134 highlights how the word “swim” is a loaded and problematic metaphor, as it imposes human-centric  
135 attributes on non-human entities. Similarly, words like “thinking,” “intelligence,” and “reasoning”  
136 are metaphors we frequently use, and they are what the computer scientist Marvin Minsky referred to  
137 as “suitcase words” (terms packed with a range of meanings and assumptions).

138 For example, while humans have always taken inspiration from birds to achieve flight, human flight  
139 looks nothing like avian flight. Airplanes achieve heavier-than-air flight using entirely different  
140 principles, such as fixed wings and engines, compared to the flapping of bird wings.

141 This shows that there can be multiple ways to solve the same problem, each fundamentally different  
142 from the other. For instance, the problem of heavier-than-air flight through an atmosphere—has  
143 been addressed in ways that differ profoundly between humans and birds. Similarly, the problems of  
144 intelligence and reasoning, however they are defined, may also be solved in fundamentally different  
145 ways. If broad intelligence is defined as solving problems in a complex novel environment, then  
146 machines and humans may come up with fundamentally different ways of solving these problems.

147 *We might need to get used to the idea that machines and humans may have very different kinds of*  
148 *intelligence.*

149 The way machines and humans approach problem-solving can be fundamentally different, yet both  
150 can lead to effective solutions. Intelligence and reasoning, however defined, need not be constrained  
151 to human-like thought processes. Machines might arrive at solutions that are equally effective, yet  
152 completely alien to human cognition.

153 *Instead of insisting that machines must mimic human cognition to be considered intelligent, we*  
154 *should embrace the possibility that intelligence comes in diverse forms, shaped by the constraints*  
155 *and capabilities of the system in which it operates.*

156 **References**

- 157 François Chollet. On the measure of intelligence. 11 2019. doi: 10.48550/arxiv.1911.01547. URL  
158 <https://arxiv.org/abs/1911.01547v2>.
- 159 Rich Sutton. The bitter lesson, 2019. URL [http://www.incompleteideas.net/IncIdeas/  
160 BitterLesson.html](http://www.incompleteideas.net/IncIdeas/BitterLesson.html).
- 161 Melanie Mitchell, Alessandro B. Palmarini, and Arseny Moskvichev. Comparing humans, gpt-4,  
162 and gpt-4v on abstraction and reasoning tasks. 11 2023. URL [https://arxiv.org/abs/2311.  
163 09247v2](https://arxiv.org/abs/2311.09247v2).
- 164 Benjamin Estermann, Luca A Lanzendörfer, Yannick Niedermayr, Roger Wattenhofer, and Eth Zürich.  
165 Puzzles: A benchmark for neural algorithmic reasoning. 6 2024. URL [https://arxiv.org/  
166 abs/2407.00401v1](https://arxiv.org/abs/2407.00401v1).
- 167 Melanie Mitchell. Debates on the nature of artificial general intelligence. *Science*, 383:eado7069,  
168 3 2024a. ISSN 10959203. doi: 10.1126/SCIENCE.ADO7069. URL [https://www.science.  
169 org/doi/10.1126/science.ado7069](https://www.science.org/doi/10.1126/science.ado7069).
- 170 Melanie Mitchell. *Artificial Intelligence: A Guide for Thinking Humans*. Penguin Publishers, 2019.
- 171 Melanie Mitchell. The metaphors of artificial intelligence. *Science*, 386:eadt6140, 11 2024b. ISSN  
172 10959203. doi: 10.1126/SCIENCE.ADT6140. URL [https://www.science.org/doi/10.  
173 1126/science.adt6140](https://www.science.org/doi/10.1126/science.adt6140).