
Transformers Have the Potential to Achieve AGI

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 As large language models (LLMs) based on the Transformer architecture continue
2 to achieve impressive performance across diverse tasks, this paper explores whether
3 Transformers can ultimately achieve artificial general intelligence (AGI). We argue
4 that Transformers have significant potential to achieve AGI, supported by the
5 following insights and arguments. (1) A Transformer is expressive enough to
6 simulate a programmable computer equipped with random number generators and,
7 in particular, to execute programs for meta-tasks such as algorithm design. (2) By
8 the extended Church-Turing thesis, if some realistic intelligence system (say, a
9 human with pencil and paper) achieves AGI, then in principle a single Transformer
10 can replicate this capability; Besides, we suggest that Transformers are well-suited
11 to approximate human intelligence, because they effectively integrate knowledge
12 and functions represented in network form (e.g. pattern recognition) with logic
13 reasoning abilities. (3) We argue that Transformers offer a promising practical
14 approximation of Hutter’s AIXI agent, which is an ideal construction to achieve
15 AGI but is uncomputable.

16 1 Introduction

17 Large language models (LLMs) [1–4] have demonstrated remarkable capabilities across a broad
18 range of challenging tasks. For example, OpenAI’s o-series [5] model achieves 71.7% accuracy on
19 the software engineering benchmark SWE-bench [6], 87.7% on the graduate-level question answering
20 task GPQA [7], and 96.7% on a competition-level mathematics reasoning task [8]. Notably, these
21 results surpass human-expert performance. As LLMs evolve, their capabilities are expected to
22 advance further.

23 These successes are grounded in the Transformer architecture [9], which has proven to be highly
24 effective across a wide range of domains, extending beyond natural language processing to areas
25 such as computer vision [10] and decision-making [11]. Given the impressive achievements of
26 Transformers in tackling challenging tasks across various domains, a fundamental question arises:

27 *Question 1: Can Transformers ultimately achieve artificial general intelligence (AGI)?*

28 To answer this question, we must first establish a rigorous definition of intelligence. Intelligence is
29 multifaceted, encompassing abilities such as creativity, problem-solving, pattern recognition, and
30 reasoning. However, formulating a single, comprehensive definition that captures all these aspects is
31 challenging. As pointed out in [12], most, if not all, aspects of intelligence can be framed in terms of
32 goal-driven behavior, or more precisely, as the maximization of some (often unknown) utility (reward)
33 function. This aligns with the “reward is enough” hypothesis [13], which suggests that the pursuit
34 of maximizing reward alone is sufficient to drive behaviors that exhibit a wide range of capabilities,
35 many of which are traditionally studied in both natural and artificial intelligence.

36 In this paper, we follow the definition that intelligence can be broadly categorized into two types of
37 reasoning abilities:

- Learning an unknown utility function (inductive reasoning): This involves drawing generalizations from specific observations, where the conclusions are probable but not certain. This type of reasoning is extensively explored in the context of inverse reinforcement learning [14, 15]. Examples of inductive reasoning include pattern recognition, natural language processing, prediction, and scientific research, where repeated observations lead to hypotheses or theories.
- Maximizing a known utility function (deductive reasoning): In this case, the solution depends entirely on the explicit, provided information. Successful applications include AlphaGo [16], MuZero [17], AlphaProof [18], OpenAI-o1 [5], and DeepSeek-R1 [19].

In this paper, we argue in favor of Question 1, **supporting the potential of Transformers to achieve AGI** with the following insights and arguments.

1. A single Transformer can simulate a probabilistic programmable computer. Prior works (e.g. [20]) have shown that Transformers (with chain-of-thoughts) can efficiently simulate deterministic Turing machines (DTMs). We extend this result to the potentially more powerful probabilistic Turing machines (PTMs), proving that Transformers can efficiently simulate PTMs as well (Theorem 2).

At first glance, Theorem 2 may suggest adherence to a one-model-one-task paradigm, where different tasks require different transformers. This misaligns with the current practice of training a single general-purpose transformer to perform various tasks. In fact, Theorem 2 provides deeper insights: as also observed in related work (e.g. [21]), it implies that a *single* Transformer can simulate a probabilistic universal Turing machine (UTM), a formalization of a general-purpose programmable computer equipped with random number generators.

Furthermore, while Transformers do not follow the one-model-one-task paradigm, they appear to adhere to a one-prompt-one-task paradigm, where different tasks require different PTMs (or equivalently, programs) to be specified in the prompt or pre-injected during training. We argue that this is *not* the case. Specifically, beyond algorithms for specific tasks, a PTM T can also serve as a program for meta-tasks, such as designing other algorithms (meta-algorithms), or even higher-order tasks, such as meta-meta-algorithms.

2. Implication of the extended Church-Turing thesis. The *extended Church-Turing thesis* (ECT) [22, 23], an extension of the Church-Turing thesis in the modern computer science literature from a complexity-theoretic perspective, asserts that the PTM model is not only as expressive as but also as efficient as any realistic physical device (say, a human brain, a society, or a future neural network). Specifically, any function that can be computed by a realistic finite physical system can also be computed by a PTM with at most a polynomial slowdown. Consequently, if some realistic intelligence system (say, a human brain with pencil and paper) achieves AGI, then in principle, a single Transformer can achieve AGI as well (Thesis 1).

In particular, Thesis 1 suggests that a single Transformer has the potential to achieve human-level intelligence. Moreover, we suggest that Transformers are particularly well-suited as approximations of human intelligence, because they effectively integrate knowledge and functions represented in network form with logical reasoning abilities, and thus can leverage benefits from both connectionism AI and symbolism AI.

3. Algorithmic approximations of general intelligence. Besides mimicking the human reasoning process, another line of research, inspired by algorithmic information theory, seeks to reach or even outperform human intelligence by establishing a formal theory of general intelligence. Several constructions have been proposed to address meta-tasks, including:

Levin's universal search algorithm: Many deductive reasoning tasks, such as theorem proving, planning, and general NP-complete problems, can be effectively modeled as search problems. Levin's universal search is an algorithm that can solve all search problems as quickly as the fastest algorithm for each, up to a large constant factor [24, 25]. The basic idea is to run all programs p in parallel with relative computation time $2^{-\ell(p)}$; i.e. a time fraction $2^{-\ell(p)}$ is dedicated to executing p . Here, we describe programs as Boolean strings using a prefix-free encoding, where $\ell(p)$ denotes the length of the description of p . Note that the sum of all these time fractions satisfies $\sum_p 2^{-\ell(p)} \leq 1$.

Solomonoff's universal induction: Every inductive reasoning task, such as continuing a number of series in an IQ test, classification in machine learning, stock-market forecasting, or scientific research, can be described as a sequence prediction problem, more precisely, predicting future data from past observations [12]. Solomonoff's universal induction [26, 27] is an optimal approach for all

92 sequence prediction problems, where the data is sampled from a computable probability distribution,
93 or equivalently generated by a realistic physical system according to the physical version of Church-
94 Turing thesis. The basic idea is to do Bayesian prediction, using Solomonoff prior as the prior belief,
95 which assigns higher probabilities to simpler hypotheses with shorter descriptions, aligning with
96 Occam’s razor.

97 *Hutter’s AIXI agent*: AIXI is a theoretical agent that achieves AGI [12]. AIXI is somehow a
98 combination of Solomonoff’s universal induction and Levin’s universal search. Specifically, AIXI
99 replaces the unknown environment in the Bellman equation with a generalized Solomonoff prior and
100 then invokes $M_{p^*}^e$ [12], an enhancement of Levin’s universal search, to solve the Bellman equation.
101 Like Solomonoff induction, AIXI tends to hypothesize the environment as shortest possible programs,
102 in line with Occam’s razor.

103 While these constructions are theoretically optimal, they are often intractable in practice or even
104 uncomputable. However, we argue that Transformers provide a promising tractable approximation of
105 these universal constructions. Specifically,

106 *Universal search*: By Theorem 2, a single Transformer, by simulating Levin search, can theoretically
107 solve all search problems as efficiently as the fastest algorithm for each problem. To enhance
108 tractability, Transformers can leverage prior knowledge embedded during training to assign the
109 relative computation time proportion in a more adaptive and efficient way. In addition, Transformers
110 can continually refine search strategies by learning from past experiences.

111 *Universal induction*: Recent works [28–32] have demonstrated that Transformers align with Occam’s
112 razor, the core principle of Solomonoff induction. Specifically, Transformers tend to output sequences
113 generated by shorter programs (a.k.a. with lower Kolmogorov complexity). The alignment with
114 Occam’s razor enables Transformers to generalize effectively across diverse tasks and data modalities,
115 making them good approximations of general-purpose predictors. Furthermore, [32] put forth and
116 explore a hypothesis that Transformers approximate Solomonoff induction better than any other extant
117 sequence prediction method, highlighting their potential as practical implementations of universal
118 induction.

119 *AIXI agent*: we suggest that Transformers have the potential to offer a practical approximation of AIXI
120 for the following reasons: as we just discussed, (i) Transformers have the potential to approximately
121 implement Solomonoff universal induction; (ii) Transformers has potential to implement universal
122 search in practice, enabling efficient solutions for a wide range of deductive reasoning tasks; and (iii)
123 Transformers integrate prior knowledge effectively, leveraging human experience to enhance their
124 practical applicability. Moreover, for practicality, we propose a framework that approximates the
125 AIXI agent using two complementary Transformers in Section 4.3.

126 2 Transformers can efficiently simulate probabilistic TMs

127 We assume that the reader is familiar with the definitions of the Transformer architecture, Turing
128 machine (TM), probabilistic Turing machine (PTM), and universal Turing machine (UTM). For the
129 convenience of readers, we present a background of TM in Appendix A. There is a line of theoretical
130 works [33–37, 20, 38, 39, 21] studying the expressive power of Transformer with chain of thought
131 (CoT) by connecting them with Turing machines. It turns out that decoder-only Transformers with t
132 CoT steps can simulate t DTM steps.

133 **Theorem 1** ([20]). *Let T be a deterministic Turing machine that, on input x of length n , runs for at
134 most $t(n)$ steps. There is a constant-depth decoder-only Transformer that, on input x , takes $t(n)$ CoT
135 steps and then outputs $T(x)$.*

136 In this paper, we extend Theorem 1 to PTMs. In particular, it implies that Transformers with
137 polynomial CoT steps can solve all problems in BPP, the class of decision problems solvable by
138 a PTM in polynomial time, which is *strictly* larger than the P class, which consists of all decision
139 problems solvable by a DTM in polynomial time, unless BPP = P. The proof of Theorem 2 can be
140 found in Appendix B.

141 **Theorem 2.** *Let T be a probabilistic Turing machine that, on input x of length n , runs for at most
142 $t(n)$ steps. There is a constant-depth decoder-only transformer that, on input x , takes at most $2t(n)$
143 CoT steps and returns the same (randomized) output as T .*

144 At first glance, Theorem 2 appears to follow the one-model-one-task paradigm: different tasks require
145 different Transformers. This misaligns with the current practice of training a single general-purpose
146 transformer to perform various tasks. In fact, Theorem 2 provides deeper insights: as also observed
147 in related works (e.g. [21]), Theorem 2 implies that a *single* Transformer can simulate a UTM, or
148 equivalently a programmable computer equipped with random number generators.

149 Though Transformers do not have to follow the one-model-one-task paradigm, they appear to follow
150 the one-prompt-one-task paradigm: for different tasks, different PTMs should be loaded into the
151 prompt or pre-injected during training. We argue that this is not the case. Specifically, beyond
152 algorithms for specific tasks, the PTM T can also be taken as a program that performs meta-tasks,
153 such as a meta-algorithm that designs algorithms, or even meta-meta-algorithms. For example, the
154 Transformer can implement the following meta procedure: it takes a problem description and an input
155 x as input, and then

- 156 1. run some prescribed meta-algorithm to initialize or update a program p ;
- 157 2. run p on input x , and obtain $p(x)$;
- 158 3. evaluate $p(x)$. If not good enough, then go to Step 1.

159 3 The extended Church-Turing thesis and its implication

160 The physical version of *Church-Turing thesis* (CT), as known as Deutsch-Wolfram thesis, asserts [40–
161 42] that every finite physical system (say, a modern personal computer, a human brain, a society, or a
162 future neural network) can be simulated to any specified degree of accuracy by a PTM. Furthermore,
163 there is also a strengthening, referred to as the *Extended Church-Turing thesis* (ECT), of the physical
164 Church-Turing thesis in the modern computer science literature [22, 23] from a complexity-theoretic
165 perspective, asserting that the probabilistic Turing machine model is also as efficient as any computing
166 device can be. That is, if a function is computable by some hardware device in time $T(n)$ for the
167 input of size n , then it is computable by a PTM in time $O(T(n)^k)$ for some constant k .

168 **Remark 1.** *The physical version of CT and ECT are very different from the original version proposed*
169 *by Church and Turing [43, 44], which asserts that every algorithmic process can be carried out by*
170 *a PTM. Specifically, if a task can solved by a human being with paper and pencil by following a*
171 *finite number of exact instructions, then the original CT asserts that it can also be solved by a PTM.*
172 *Notably, no insight, intuition, or ingenuity is demanded on the part of the human being carrying out*
173 *the method, which is very different from the physical version. The original CT is something between*
174 *a theorem and a definition. And the physical version and ECT are neither mathematical theorems nor*
175 *definitions. If they are true, then the truth is a consequence of the laws of physics [45].*

176 By combining Theorem 2 and ECT, we obtain the following thesis:

177 **Thesis 1.** *If some realistic intelligence system (say, a human brain with pencil and paper) achieves*
178 *AGI, then a single Transformer can also achieve AGI with at most a polynomial slowdown.*

179 3.1 Algorithmically description of human reasoning

180 If we accept that (a) the extended Church-Turing thesis applies to the human’s reasoning process,
181 meaning that the reasoning process of humans can be efficiently simulated by a PTM, and (b) a
182 human brain, or a group of human brains (say, a research community) with paper and pencil can
183 achieve AGI, then by Theorem 2, we should also accept that in principle a single Transformer or a
184 group of Transformers can also achieve AGI as well. The related challenge lies in algorithmically
185 describing the human reasoning process, including cognitive functions like intuition or creativity.

186 **Question 2:** *How to algorithmically describe human reasoning process?*

187 There are two kinds of general approaches to this challenge: connectionism and symbolism.

- 188 • **Connectionism: Simulating the Brain at the Physical Level.** Connectionism posits that human
189 reasoning arises from the emergent properties of biologically inspired neural networks. By
190 modeling the brain’s physical and biological substrates—specifically, the interactions of neurons
191 through synaptic connections—this approach seeks to replicate cognitive processes via distributed,
192 parallel computation. Modern artificial neural networks, such as deep learning architectures,

193 exemplify this paradigm. These systems learn hierarchical representations from data, mirroring
194 how the brain processes sensory input and abstracts patterns [46, 47].

195 For instance, Transformer architectures [9] model sequential reasoning by leveraging temporal
196 dependencies and attention mechanisms, achieving expert-level performance in complex math-
197 ematical reasoning and code generation tasks [5, 48]. Connectionist models excel at pattern
198 recognition and probabilistic reasoning but often lack explicit symbolic representations, leading
199 to critiques about their interpretability and inability to handle structured, rule-based logic [49].
200 Recent advances in neuro-symbolic integration, however, aim to bridge this gap by combining
201 neural networks with symbolic reasoning modules [50, 51].

202 • **Symbolicism: Abstracting General Principles of Human Thought.** Symbolicism adopts a
203 top-down perspective, seeking to formalize the universal principles and logical structures that
204 underpin human reasoning. Rooted in classical AI and influenced by philosophy, linguistics, and
205 formal logic, this approach abstracts cognition into discrete symbols and rules, independent of
206 biological implementation. Unlike connectionism, which emulates neural substrates, symbolism
207 prioritizes computational-level explanations of thought—asking what problems cognition solves
208 and why, rather than how the brain physically solves them [52, 53].

209 At its core, symbolism assumes that reasoning can be modeled as manipulation of explicit
210 representations through deterministic or probabilistic rules. For example: Occam’s razor, a
211 heuristic for inductive reasoning, is formalized in algorithmic frameworks like Bayesian model
212 selection [54], where simpler hypotheses are assigned higher prior probabilities. Deductive
213 reasoning is captured by logic-based systems (e.g., Prolog, theorem provers) that apply syllogistic
214 rules (e.g., modus ponens) to derive conclusions from premises [55].

215 Here, we argue that an effective solution requires a combination of these two approaches, since (i)
216 abstracting general principles offers a more tractable and generalizable framework for intelligence
217 and (ii) part of knowledge and functions, such as pattern recognition and cognitive functions, may
218 have no representation more concise than a huge, analog neural network [56], thus are not suitable to
219 be represented as logic or symbolic.

220 In particular, since Transformers can effectively integrate knowledge and functions represented in
221 network form (since they are neural networks) with logical reasoning abilities (Theorem 2), and thus
222 can leverage benefits from both connectionism and symbolism, we suggest that Transformers are
223 particularly well-suited as approximations of human intelligence.

224 4 Algorithmic approximations of general intelligence

225 Besides mimicking human reasoning process, another line of research, inspired by algorithmic infor-
226 mation theory [57], aims to achieve or even surpass human-level intelligence by establishing a formal
227 theory of general intelligence, such as Levin’s universal search algorithm [24, 25], Solomonoff’s
228 universal induction [26, 27], and Hutter’s AIXI agent [12]. While these constructions are theoretically
229 optimal, they are often intractable in practice and even uncomputable. However, we argue that
230 Transformers provide a promising and tractable approximation of these universal constructions.

231 4.1 Levin’s universal search

232 Many deductive reasoning tasks, such as theorem proving, planning, and general NP-complete
233 problems, can be effectively modeled as search problems.

234 **Search problems.** Let $\phi : \{0, 1\}^* \rightarrow \{0, 1\}^*$ be a function where $\phi(\cdot)$ can be computed quickly (say,
235 in polynomial time). The search problem is defined as: given y , find an x such that $\phi(x) = y$.

236 For example, in the Boolean satisfiability problem (SAT), the function $\phi : \{0, 1\}^* \rightarrow \{0, 1\}$ can
237 be defined as a verifier that checks whether a given assignment satisfies the Boolean formula in
238 conjunctive normal form.

239 **Levin search.** The algorithm is simple to describe: just run and verify the output of all algorithms p
240 in parallel with relative computation time $2^{-\ell(p)}$; i.e. a time fraction $2^{-\ell(p)}$ is devoted to executing
241 p [24, 25]. Here, programs are described as Boolean strings in a prefix-free encoding, where $\ell(p)$
242 denotes the length of the description of p . Note that $\sum_p 2^{-\ell(p)} \leq 1$.

243 **Theorem 3** ([24, 25, 12]). *The computation time of Levin search is upper bounded by $\min_p\{2^{\ell(p)} \cdot$*
 244 *time_p⁺(y), where time_p⁺(y) is the runtime of p(y) plus the time to verify the correctness of the result*
 245 *(ϕ(x) = y) by a known implementation for ϕ.*

246 By Theorem 2, we conclude that in principle, a single Transformer can solve all search problems as
 247 quickly as the fastest algorithm for each, up to a constant factor.

248 We note that Levin’s universal search—which optimally allocates computational effort across can-
 249 didate solvers according to their algorithmic probability (Theorem 3)—may provide a theoretical
 250 foundation for the emerging paradigm of inference-time scaling in LLMs [58, 59, 5, 48]. This
 251 framework structures LLM reasoning into two synergistic phases: generating diverse candidate
 252 solutions (or algorithms) and efficiently prioritizing their execution and evaluation, mirroring Levin’s
 253 time-optimal balance between exploration and exploitation. By prescribing a focus on programs with
 254 minimal description length (i.e., favoring simpler, valid solutions), Levin’s principles offer guidance
 255 for designing compute-efficient strategies.

256 Though Levin search is theoretically optimal for all search problems, the large constant overhead
 257 $2^{\ell(p)}$ renders it impractical. A line of research [60–64] has explored adaptations of Levin’s search
 258 that leverage past experience to improve its efficiency. We note that the key lies in generating highly
 259 successful algorithms p with the shortest description length, ensuring they are prioritized during the
 260 search process. Such knowledge can be acquired from experience. For instance, a Transformer could
 261 maintain a parameterized model (e.g., a neural network or program) within its context and employ
 262 bootstrap methods—such as search-and-learn processes [65]—to iteratively refine its performance. By
 263 repeatedly solving increasingly challenging instances and updating the model based on successfully
 264 solved examples, the system could incrementally improve its problem-solving efficiency.

265 4.2 Solomonoff’s universal induction

266 Every inductive reasoning task, such as continuing a number of series in an IQ test, classification in
 267 machine learning, stock-market forecasting, or scientific research, can be described as a sequence
 268 prediction problem, more precisely, predicting future data from past observations [12]. Without loss
 269 of generality and for simplicity, we assume the data $x_i \in \{0, 1\}$ is binary.

270 We first introduce some notations and definitions. Given a subset S of $\{0, 1\}^*$, let $[S]$ denote the set
 271 obtained from S by deleting all elements that have a prefix in S . A *monotone Turing machine* is a
 272 Turing machine with one unidirectional ready-only input tape, one unidirectional write-only output
 273 tape, and some bidirectional work tapes. We say a tape is unidirectional if its head can only move
 274 from left to right, and bidirectional if its head can move in both directions.

275 **Definition 1** (Measure). *We say a function $\mu : \{0, 1\}^* \rightarrow [0, 1]$ is a measure if $\mu(\emptyset) = 1$ and*
 276 *$\mu(x) = \mu(x1) + \mu(x0)$. Here, \emptyset denotes the empty string.*

277 *A measure μ defines a random process generating an infinitely long binary sequence: start with*
 278 *an empty string and repeatedly select the next bit $x_n \in \{0, 1\}$ according to the probability $\mu(x_n |$*
 279 *$x_{<n} := \mu(x_{<n}x_n)/\mu(x_{<n})$ conditioned on the past data $x_{<n} := x_1x_2 \dots x_n$.*

280 We say μ is estimable if there exists a TM that, given $x \in \{0, 1\}^*$ and a precision ϵ , computes an
 281 ϵ -approximation of $\mu(x)$. By the physical version of the Church-Turing thesis, any μ implemented on
 282 a finite, realistic physical device is estimable. Moreover, by Theorem 4.5.2 in [57], for any estimable
 283 μ , there is a monotone TM T that takes an infinitely long uniformly random binary string as input and
 284 generates an infinitely long binary sequence according to μ . Let $K(\mu)$ denote the shortest description
 285 of such a T .

286 **Sequence prediction problem.** Having observed the past data $x_{<n} := x_1x_2 \dots x_{n-1}$, the task
 287 is to predict the next bit x_n . More precisely, let μ denote the unknown underlying mechanism
 288 generating the sequence $x_1x_2 \dots$. The task is to estimate the conditional probability $\mu(x_n | x_{<n}) :=$
 289 $\mu(x_{<n}x_n)/\mu(x_{<n})$.

290 **Solomonoff’s universal induction.** Bayesian prediction provides a framework for sequence predic-
 291 tion problems, which repeatedly employs Bayes’ rule to update its beliefs about each hypothesis based
 292 on newly observed data. The primary challenge is how to select the prior beliefs. Solomonoff [26, 27]
 293 addressed this challenge by introducing a universal prior, rooted in the simplicity of hypotheses. His
 294 approach leverages the fact that simpler hypotheses, represented by shorter programs, are more likely

295 to generalize well—a concept aligned with Occam’s Razor. Solomonoff showed that the Bayesian
 296 prediction with the Solomonoff prior as the prior belief is an optimal way for the sequence prediction
 297 problem, provided that the underlying μ is estimable.

298 **Definition 2** (Solomonoff prior [26, 27]). *Let U be a monotone UTM. The Solomonoff prior is defined
 299 as*

$$M_U(x) := \sum_{\lfloor p \in \{0,1\}^*: U(p) = x \star \rfloor} 2^{-\ell(p)}.$$

300 Here, $U(p) = x\star$ means x is a prefix of $U(p)$. Intuitively, $M_U(x)$ is the probability that the output
 301 starts with x when the input is an infinite-long uniformly random binary string.

302 Solomonoff’s universal induction (SI) is simple to describe: use $M_U(x_n \mid x_{<n}) =$
 303 $M_U(x_n)/M_U(x_{<n})$ as an estimate of the true conditional probability $\mu(x_n \mid x_{<n})$.

304 **Theorem 4** (Solomonoff central theorem [26, 27]). *For any estimable μ , we have*

$$\sum_{n=1}^{+\infty} \sum_{x_n \in \{0,1\}} \mu(x_{<n}) (M(x_t \mid x_{<t}) - \mu(x_n \mid x_{<n}))^2 \leq \ln 2 \cdot K(\mu) + O(1).$$

305 For any estimable μ , the upper bound $\ln 2 \cdot K(\mu)$ is finite, so the difference $M(x_t \mid x_{<t}) - \mu(x_n \mid x_{<n})$
 306 tends to zero as $n \rightarrow \infty$ with μ -probability 1. Consequently, $M(x_t \mid x_{<t})$ converges rapidly to the
 307 true underlying generating process.

308 Unfortunately, Solomonoff prior $M_U(x)$ is inestimable: there is no TM that, given $x \in \{0,1\}^*$ and a
 309 precision ϵ , can compute an ϵ -approximation of $M_U(x)$ in finite time. To address this uncomputability
 310 issue, several approximations have been proposed [66–68, 30].

311 In particular, observing that Transformers are naturally suited for sequence prediction tasks, a line
 312 of work [28–32] has explored whether the Transformer model can approximate SI. Specifically,
 313 [29, 28] showed that transformers can do Bayesian inference. [30] proved that Transformers can
 314 approximate SI by training solely on UTM data, and demonstrated that increasing model size leads to
 315 improved performance. [32] proposed and investigated the hypothesis that Transformers approximate
 316 Solomonoff induction better than any other extant sequence prediction method. This hypothesis was
 317 further supported by [31, 69]. Specifically, they showed that like Solomonoff induction, transformers
 318 also align with Occam’s Razor: transformers prefer generating data with low Kolmogorov complexity.
 319 Occam’s razor provides transformers with good generalization on many different problems and
 320 modalities of data, and makes them powerful general-purpose predictors.

321 4.3 Hutter’s AIXI agent

322 We first briefly introduce Hutter’s AIXI agent, which is claimed to be universal in that it is independent
 323 of the true environment (model-free) and is able to solve any solvable problem and learn any learnable
 324 task. The main idea of AIXI is simple to describe: just replace the unknown environmental distribution
 325 in the Bellman equations with a suitably generalized Solomonoff prior [12].

326 **Setting.** The agent and the environment interact chronologically as follows: in each cycle k , the agent
 327 performs an action $y_k \in \mathcal{Y}$ (output), and then receives a perception $x_k \in \mathcal{X}$ from the environment.
 328 The perception x_k consists of a regular part o_k and a reward r_k . Given the history $y_1 x_1 \cdots x_{k-1} y_k$,
 329 the probability that the environment produces perception x_k is denoted $\mu(x_k \mid y_1 x_1 \cdots x_{k-1} y_k)$.
 330 Here, we make no assumptions about μ other than it is estimable. In particular, μ is allowed to depend
 331 on the complete history $y_1 x_1 \cdots x_{k-1} y_k$.

332 We use p to denote the agent’s policy, which can be described as a monotone Turing machine that
 333 takes $x_1 x_2 \cdots$ as input and outputs $y_1 y_2 \cdots$. As the optimal policy can always be chosen to be
 334 deterministic, we assume p is a deterministic monotone TM. In addition, we say $y_{1:k} = p(x_{<k})$
 335 if $y_i = p(y_1 x_1 y_2 x_2 \cdots x_{i-1})$ for $i \leq k$. We also use $\mu(x_{k:m} \mid y_{1:m} x_{<k})$ as an abbreviation for
 336 $\prod_{i=k}^m \mu(x_i \mid x_{<i}, y_{\leq i})$. We define the value of policy p in environment μ as

$$V_\mu^p := \sum_{x_{1:m}} (r_1 + \cdots + r_m) \mu(x_{1:m} \mid y_{1:m})_{|y_{1:m}=p(x_{<m})}$$

337 where m is the lifespan of the agent.

338 The goal of the agent is to maximize the total reward $\sum_{i=1}^m r_i$. Formally, the agent aims to find a
 339 policy p^μ that maximizes V_μ^p .

340 **The AIXI agent.** If the environment μ is known, then the optimal policy is

$$y_k := \arg \max_{y_k} \sum_{x_k} \cdots \max_{y_m} \sum_{x_m} \left(\sum_{i=k}^m r_i \right) \mu(x_{k:m} \mid y_{1:m}, x_{<k})$$

341 with total reward

$$\max_{y_1} \sum_{x_1} \cdots \max_{y_m} \sum_{x_m} (r_1 + \cdots + r_m) \mu(x_{1:m} \mid y_{1:m}) := V_\mu^*$$

342 The AIXI agent replaces the true but unknown μ with a generalized Solomonoff prior. Specifically,
343 the AIXI policy is

$$y_k := \arg \max_{y_k} \sum_{x_k} \cdots \max_{y_m} \sum_{x_m} \left(\sum_{i=k}^m r_i \right) \xi(x_{k:m} \mid y_{1:m}, x_{<k})$$

344 where

$$\xi(x_{1:k} \mid y_{1:k}) := \sum_{\substack{\text{monotone TM } q: q(y_{1:k}) = x_{1:k}}} 2^{-\ell(q)}.$$

345 Intuitively, the agent continually updates its belief about hypotheses of the unknown environment μ
346 by Bayes' rule. Similar to Solomonoff universal induction, environments with lower Kolmogorov
347 complexity are preferred, in line with Occam's razor. [12] shows that AIXI's environment model
348 converges rapidly to the true environment, and its policy is Pareto-optimal and self-optimizing. Here,
349 we say a policy Pareto-optimal if there is no other agent that performs at least as well as AIXI in
350 all environments while performing strictly better in at least one environment, and self-optimizing if
351 $\frac{1}{m} V_\mu^{\text{AIXI}} \rightarrow \frac{1}{m} V_\mu^*$ for horizon $m \rightarrow +\infty$ for all estimable μ .

352 Unfortunately, like Solomonoff's universal induction, AIXI is uncomputable. To address this issue,
353 several computable approximations have been proposed [12, 70–72, 67, 73–77]. One such approxima-
354 tion is AIXIt ℓ , which performs at least as well as any other agent bounded time t and length ℓ . Some
355 approximations focus on restricted environment classes and have been successfully implemented
356 [72]. [77] studied how to inject knowledge into the AIXI agent.

357 We suggest that Transformers have the potential to approximate AIXI for the following reasons: (i)
358 as discussed above, Transformers might serve as a good approximation of SI, and provide a good
359 estimation of ξ ; (ii) with an estimation of ξ , Transformers can solve the Bellman equation using an
360 enhanced Levin search or an enhanced $M_{p^*}^\ell$ algorithm [12], which solves all well-defined problems as
361 quickly as the fastest algorithm for each; and (iii) Transformers effectively integrate prior knowledge,
362 leveraging human experience to further enhance their practical applicability.

363 For practicality, we propose a framework that approximates AIXI using two complementary Trans-
364 formers: a Environment Modeler (induction) and a Action Planner (deduction).

365 **• Transformer I (induction component): Environment Modeler.** This component aims to
366 approximate SI, providing a good estimation of ξ . We list some potential methods to improve its
367 practical inductive reasoning ability.

- 368 **– Mixing synthetic UTM data with real-world data:** Training exclusively on synthetic UTM data
369 is already sufficient to enable Transformers to converge to SI, as such data spans a universal
370 distribution over computable sequences [30]. Nevertheless, by incorporating real-world data,
371 we can leverage prior knowledge to make this component much more practical. Specifically,
372 adding real-world data biases the model toward assigning higher probability to programs that
373 are relevant to real-world environments. This leads to faster learning and fewer errors when
374 the model is deployed in such environments, while still preserving the universality [30].
- 375 **– Data and model scaling are essential:** As shown in [30], increasing model size leads to better
376 approximation of SI. This supports the need for massive data and extensive pretraining.
- 377 **– Leveraging environment feedback:** After deployment, the model can continually refine its
378 prior through periodic fine-tuning on new feedbacks. While ICL is theoretically sufficient,
379 practical Transformers are limited by finite context lengths, so parameter updates through
380 finetuning are crucial to consolidate the information from the long-term feedback.

381 **• Transformer II (deduction component): Action planner.** This component solves the explicit
382 decision-making or planning task based on the prediction from Transformer I. This is not a

383 machine learning problem by itself; so unlike the inductive component, massive training data are
384 not required in principle.

385 Importantly, our deduction framework for the AIXI agent offers a promising direction for improving
386 LLMs’ test-time scaling behavior. While current methods employ direct reward maximization
387 (e.g., RL) to drive reasoning—focusing narrowly on final outcomes—our approach integrates
388 Levin search to systematically evaluate reasoning steps based on their algorithmic complexity.
389 This principled methodology achieves an optimal balance between: exploration (discovering
390 diverse reasoning pathways) and exploitation (optimizing high probability steps toward solu-
391 tions). By formalizing this trade-off, our framework may provide a foundation for developing
392 next-generation reasoning models capable of more sophisticated and scalable problem-solving.

393 Although our deduction component does not require training in principle, its practical performance
394 can be significantly improved using the learning-to-optimize techniques [78]. For example, by
395 leveraging human experience or prior problem-solving data, one can train a model that takes a
396 problem description as input and outputs a learned prior over candidate programs. This learned
397 prior is expected to assign higher weights to more promising algorithms, effectively replacing the
398 default weighting in Levin’s universal search.

399 Due to the space constraint, the discussion about the data requirement is deferred to Appendix D

400 5 Alternative Views

401 This section discusses alternative views arguing that Transformers are not a sufficient path to AGI.

402 *Alternative view 1: Transformers miss essential capabilities for intelligent beings, such as under-
403 standing and reasoning about the physical world. Specifically, Transformers cannot anchor their
404 understanding in reality: They cannot perform actions in the real world or learn through embodied ex-
405 periences, and they lack the capability for hierarchical planning, a crucial element for understanding
406 and interacting with the world at multiple levels of abstraction (e.g. [79]).*

407 We acknowledge that current Transformers lack the capabilities to interact with the physical world
408 directly and employ embodied learning. However, we do not think this represents an inherent
409 limitation of Transformers. With minor enhancements, Transformers could be embedded within
410 agent models. Specifically, such agents could utilize Transformers as an approximation of universal
411 induction (Section 4.2) to learn about the unknown environment, and subsequently apply them as an
412 approximation of universal search (Section 4.1) to perform deductive reasoning.

413 Besides, we argue that current Transformers really do understand. In our definition of intelligence
414 (Section 1), understanding can be equated to inductive reasoning, i.e., the ability to uncover the
415 underlying general mechanism from specific observations. As we argued in Section 4.2, Transformers
416 provide a promising practical approximation of Solomonoff’s universal induction, which is a universal
417 and optimal way to do inductive reasoning.

418 In addition, as we argued in Section 2, Transformers can execute any meta-process, such as algorithm
419 design, when an algorithmic description is provided. In particular, as argued in Section 4.1, Trans-
420 formers provide a promising practical approximation of Levin’s universal search algorithm, enabling
421 them to efficiently perform various deductive reasoning tasks, including planning and theorem proof.

422 *Alternative view 2: Transformers are limited by their expenditure of bounded compute per input
423 instance, e.g. the finite context window and finite precision, thus cannot simulate a UTM, whose tapes
424 are infinitely long (e.g. [79, 31, 80]).*

425 First, no finite physical system, such as a human brain or a personal computer, can solve problems of
426 infinite size. This limitation naturally extends to the simulation of a UTM, which assumes infinitely
427 long tapes. Therefore, when discussing whether a Transformer can simulate a UTM, the correct
428 interpretation should follow the framework of the logical circuit model [81], specifically: “a uniform
429 family of Transformers can simulate a UTM.” In other words, for any arbitrarily large tape length ℓ ,
430 there exists an efficiently constructible Transformer capable of simulating a UTM with tape length ℓ .

431 In this context, while an exact simulation of a UTM is impossible, a sufficiently large Transformer
432 can approximate its behavior to an arbitrarily high degree of accuracy. This scalability ensures
433 that Transformers, much like circuits, can address increasingly complex problems within practical
434 computational limits.

435 **References**

436 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
437 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4
438 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

439 [2] Team Gemini, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
440 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
441 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

442 [3] Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024.

443 [4] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
444 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd
445 of models. *arXiv preprint arXiv:2407.21783*, 2024.

446 [5] OpenAI. Learning to reason with LLMs. OpenAI Blog, Feb 2024. <https://openai.com/index/learning-to-reason-with-llms>.

447 [6] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
448 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint
449 arXiv:2310.06770*, 2023.

450 [7] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
451 Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a
452 benchmark. *arXiv preprint arXiv:2311.12022*, 2023.

453 [8] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
454 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.
455 *arXiv preprint arXiv:2103.03874*, 2021.

456 [9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
457 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information
458 Processing Systems*, pages 5998–6008, 2017.

459 [10] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at
460 scale. *arXiv preprint arXiv:2010.11929*, 2020.

461 [11] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter
462 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning
463 via sequence modeling. *Advances in Neural Information Processing Systems*, 34:15084–15097,
464 2021.

465 [12] Marcus Hutter. *Universal artificial intelligence: Sequential decisions based on algorithmic
466 probability*. Springer Science & Business Media, 2005.

467 [13] David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. Reward is enough. *Artificial
468 Intelligence*, 299:103535, 2021.

469 [14] Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *International
470 Conference on Machine Learning*, volume 1, page 2, 2000.

471 [15] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J Russell, and Anca Dragan. Inverse
472 reward design. *Advances in neural information processing systems*, 30, 2017.

473 [16] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driess-
474 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mas-
475 tering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489,
476 2016.

477 [17] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-
478 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering
479 atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.

480

481 [18] DeepMind AlphaProof and AlphaGeometry Teams. AI achieves silver-medal standard solving
 482 international mathematical olympiad problems, 2024.

483 [19] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
 484 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
 485 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

486 [20] William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of
 487 thought. In *International Conference on Learning Representations*, 2024.

488 [21] Ruizhong Qiu, Zhe Xu, Wenxuan Bao, and Hanghang Tong. Ask, and it shall be given: Turing
 489 completeness of prompting. *CoRR*, abs/2411.01992, 2024.

490 [22] Andrew Chi-Chih Yao. Classical physics and the Church-Turing thesis. *Journal of the ACM*,
 491 50(1):100–105, 2003.

492 [23] Dorit Aharonov and Umesh V. Vazirani. *Is Quantum Mechanics Falsifiable? A Computational
 493 Perspective on the Foundations of Quantum Mechanics*, pages 329–349. 2013.

494 [24] Leonid Anatolevich Levin. Universal sequential search problems. *Problemy peredachi informatsii*,
 495 9(3):115–116, 1973.

496 [25] Leonid A Levin. Randomness conservation inequalities; information and independence in
 497 mathematical theories. *Information and Control*, 61(1):15–37, 1984.

498 [26] R.J. Solomonoff. A formal theory of inductive inference. part i and ii. *Information and Control*,
 499 7(1):1–22, 1964.

500 [27] Ray Solomonoff. Complexity-based induction systems: comparisons and convergence theorems.
 501 *IEEE Transactions on Information Theory*, 24(4):422–432, 1978.

502 [28] Samuel Müller, Noah Hollmann, Sebastian Pineda-Arango, Josif Grabocka, and Frank Hutter.
 503 Transformers can do bayesian inference. In *International Conference on Learning Representations*, 2022.

504 [29] Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. Tabpfn: A
 505 transformer that solves small tabular classification problems in a second. In *International
 506 Conference on Learning Representations*, 2023.

507 [30] Jordi Grau-Moya, Tim Genewein, Marcus Hutter, Laurent Orseau, Grégoire Delétang, Elliot
 508 Catt, Anian Ruoss, Li Kevin Wenliang, Christopher Mattern, Matthew Aitchison, and Joel
 509 Veness. Learning universal predictors. In *International Conference on Machine Learning*, 2024.

510 [31] Micah Goldblum, Marc Anton Finzi, Keefer Rowan, and Andrew Gordon Wilson. Position:
 511 The no free lunch theorem, kolmogorov complexity, and the role of inductive biases in machine
 512 learning. In *International Conference on Machine Learning*, 2024.

513 [32] Nathan Young and Michael Witbrock. Transformers as approximations of solomonoff induction.
 514 *CoRR*, abs/2408.12065, 2024.

515 [33] Jorge Pérez, Javier Marinkovic, and Pablo Barceló. On the turing completeness of modern
 516 neural network architectures. In *International Conference on Learning Representations*, 2019.

517 [34] Satwik Bhattacharya, Arkil Patel, and Navin Goyal. On the computational power of transformers
 518 and its implications in sequence modeling. In *Conference on Computational Natural Language
 519 Learning*, pages 455–475, 2020.

520 [35] Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. *Journal of
 521 Machine Learning Research*, 22(75):1–35, 2021.

522 [36] Dale Schuurmans. Memory augmented large language models are computationally universal.
 523 *arXiv preprint arXiv:2301.04589*, 2023.

525 [37] Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris
 526 Papailiopoulos. Looped transformers as programmable computers. In *International Conference*
 527 *on Machine Learning*, pages 11398–11442, 2023.

528 [38] Kaiying Hou, David Brandfonbrener, Sham M. Kakade, Samy Jelassi, and Eran Malach.
 529 Universal length generalization with turing programs. *CoRR*, abs/2407.03310, 2024.

530 [39] Zhiyuan Liu, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers
 531 to solve inherently serial problems. In *International Conference on Learning Representations*,
 532 2024.

533 [40] Stephen Wolfram. Undecidability and intractability in theoretical physics. *Physical Review*
 534 *Letters*, 54 8:735–738, 1985.

535 [41] David Deutsch. Quantum theory, the Church-Turing principle and the universal quantum
 536 computer. *Proceedings of the Royal Society of London. A. Mathematical and Physical Sciences*,
 537 400:117 – 97, 1985.

538 [42] B. Jack Copeland and Oron Shagrir. The Church-Turing thesis: logical limit or breachable
 539 barrier? *Commun. ACM*, 62(1):66–74, 2018.

540 [43] Alonzo Church. An unsolvable problem of elementary number theory. *American Journal of*
 541 *Mathematics*, 58:345, 1936.

542 [44] Alan M Turing. Computability and λ -definability. *The Journal of Symbolic Logic*, 2(4):153–163,
 543 1937.

544 [45] Stanford Encyclopedia of Philosophy. The Church-Turing thesis. <https://plato.stanford.edu/entries/church-turing/>, 2023.

546 [46] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep
 547 belief nets. *Neural computation*, 18(7):1527–1554, 2006.

548 [47] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,
 549 2015.

550 [48] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
 551 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
 552 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

553 [49] Gary Marcus. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631*, 2018.

554 [50] Tarek R Besold, Artur d’Avila Garcez, Sebastian Bader, Howard Bowman, Pedro Domingos,
 555 Pascal Hitzler, Kai-Uwe Kühnberger, Luis C Lamb, Priscila Machado Vieira Lima, Leo
 556 de Penning, et al. Neural-symbolic learning and reasoning: A survey and interpretation 1. In
 557 *Neuro-Symbolic Artificial Intelligence: The State of the Art*, pages 1–51. IOS press, 2021.

558 [51] Bikram Pratim Bhuyan, Amar Ramdane-Cherif, Ravi Tomar, and TP Singh. Neuro-symbolic
 559 artificial intelligence: a survey. *Neural Computing and Applications*, pages 1–36, 2024.

560 [52] Zenon W Pylyshyn. Computing in cognitive science. *Foundations of cognitive science*, pages
 561 51–91, 1989.

562 [53] Allen Newell and Herbert A Simon. Computer science as empirical inquiry: Symbols and
 563 search. In *ACM Turing award lectures*, page 1975. 2007.

564 [54] William H Jefferys and James O Berger. Ockham’s razor and bayesian analysis. *American*
 565 *scientist*, 80(1):64–72, 1992.

566 [55] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Pearson, 2016.

567 [56] Paul Graham. How to do philosophy. <http://www.paulgraham.com/philosophy.html>,
 568 2007.

569 [57] Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*,
 570 volume 3. Springer, 2008.

571 [58] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré,
 572 and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated
 573 sampling. *arXiv preprint arXiv:2407.21787*, 2024.

574 [59] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute opti-
 575 mally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*,
 576 2024.

577 [60] Ray Solomonoff. The application of algorithmic probability to problems in artificial intelligence.
 578 In *Machine Intelligence and Pattern Recognition*, volume 4, pages 473–491. 1986.

579 [61] Jürgen Schmidhuber, Jieyu Zhao, and Marco Wiering. Shifting inductive bias with success-
 580 story algorithm, adaptive levin search, and incremental self-improvement. *Machine Learning*,
 581 28:105–130, 1997.

582 [62] Jürgen Schmidhuber. Discovering neural nets with low kolmogorov complexity and high
 583 generalization capability. *Neural Networks*, 10(5):857–873, 1997.

584 [63] Jürgen Schmidhuber. Bias-optimal incremental problem solving. *Advances in Neural Infor-
 585 mation Processing Systems*, 15, 2002.

586 [64] Jürgen Schmidhuber. Optimal ordered problem solver. *Machine Learning*, 54:211–254, 2004.

587 [65] Shahab Jabbari Arfaee, Sandra Zilles, and Robert C Holte. Learning heuristic functions for
 588 large state spaces. *Artificial Intelligence*, 175(16-17):2075–2098, 2011.

589 [66] Jürgen Schmidhuber. The speed prior: A new simplicity measure yielding near-optimal
 590 computable predictions. In *Annual Conference on Computational Learning Theory*, volume
 591 2375, pages 216–228, 2002.

592 [67] Joel Veness, Peter Sunehag, and Marcus Hutter. On ensemble techniques for AIXI approxima-
 593 tion. In *Artificial General Intelligence*, volume 7716, pages 341–351, 2012.

594 [68] Daniel Filan, Jan Leike, and Marcus Hutter. Loss bounds and time complexity for speed
 595 priors. In *International Conference on Artificial Intelligence and Statistics*, volume 51, pages
 596 1394–1402, 2016.

597 [69] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christo-
 598 pher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, Mar-
 599 cus Hutter, and Joel Veness. Language modeling is compression. In *International Conference
 600 on Learning Representations*, 2024.

601 [70] Sergey Pankov. A computational approximation to the AIXI model. In *Artificial General
 602 Intelligence*, volume 171, pages 256–267.

603 [71] Joel Veness, Kee Siong Ng, Marcus Hutter, and David Silver. Reinforcement learning via AIXI
 604 approximation. In *Proceedings National Conference on Artificial Intelligence*, pages 605–611,
 605 2010.

606 [72] Joel Veness, Kee Siong Ng, Marcus Hutter, William T. B. Uther, and David Silver. A monte-carlo
 607 AIXI approximation. *J. Artif. Intell. Res.*, 40:95–142, 2011.

608 [73] Joel Veness, Martha White, Michael Bowling, and András György. Partition tree weighting. In
 609 *Data Compression Conference*, pages 321–330, 2013.

610 [74] Marc G. Bellemare, Joel Veness, and Michael Bowling. Bayesian learning of recursively
 611 factored environments. In *International Conference on Machine Learning*, volume 28, pages
 612 1211–1219, 2013.

613 [75] Marc G. Bellemare, Joel Veness, and Erik Talvitie. Skip context tree switching. In *ICML*,
 614 volume 32, pages 1458–1466, 2014.

615 [76] Samuel Yang-Zhao, Tianyu Wang, and Kee Siong Ng. A direct approximation of AIXI using
616 logical state abstractions. In *Advances in Neural Information Processing Systems*, 2022.

617 [77] Samuel Yang-Zhao, Kee Siong Ng, and Marcus Hutter. Dynamic knowledge injection for AIXI
618 agents. In *Proceedings National Conference on Artificial Intelligence*, pages 16388–16397,
619 2024.

620 [78] Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. Machine learning for combinatorial
621 optimization: a methodological tour d’horizon. *European Journal of Operational Research*,
622 290(2):405–421, 2021.

623 [79] Yann Lecun. Meta ai, open source, limits of llms, agi & the future of ai, 2024. <https://www.youtube.com/watch?v=5t1vTLU7s40>.

625 [80] Shriyash Kaustubh Upadhyay and Etan Jacob Ginsberg. Turing complete transformers: Two
626 transformers are more powerful than one. 2023.

627 [81] Sanjeev Arora and Boaz Barak. *Computational complexity: a modern approach*. Cambridge
628 University Press, 2009.

629 [82] Fred C Hennie and Richard Edwin Stearns. Two-tape simulation of multitape turing machines.
630 *Journal of the ACM*, 13(4):533–546, 1966.

631 [83] Nicholas Pippenger and Michael J Fischer. Relations among complexity measures. *Journal of*
632 *the ACM*, 26(2):361–381, 1979.

633 [84] Claus-Peter Schnorr. The network complexity and the turing machine complexity of finite
634 functions. *Acta Informatica*, 7:95–107, 1976.

635 [85] Paul Benioff. The computer as a physical system: A microscopic quantum mechanical hamil-
636 tonian model of computers as represented by turing machines. *Journal of Statistical Physics*,
637 22:563–591, 1980.

638 [86] Andrew Chi-Chih Yao. Quantum circuit complexity. In *IEEE Annual Symposium on Foundations*
639 *of Computer Science*, 1993.

640 [87] Ethan S. Bernstein and Umesh V. Vazirani. Quantum complexity theory. *ACM Symposium on*
641 *Theory of Computing*, 1993.

642 [88] John R Searle. The rediscovery of the mind. *A Bradford Book*, 1992.

643 [89] Samuel Guttenplan and Samuel D Guttenplan. *A Companion to the Philosophy of Mind*.
644 Blackwell Oxford, 1994.

645 [90] Roger Penrose. *Shadows of the Mind*, volume 4. Oxford University Press Oxford, 1994.

646 [91] Stuart Hameroff and Roger Penrose. Consciousness in the universe: A review of the ‘orch
647 or’theory. *Physics of life reviews*, 11(1):39–78, 2014.

648 **A Background on Turing Machines**

649 **Turing machines.** Turing machines (TMs) are a mathematical model of computation. A k -tape TM
 650 is defined as a tuple $\langle \Sigma, \perp, Q, q_{start}, F, \delta \rangle$ where (i) Σ is a finite tape alphabet including a blank
 651 symbol \perp , (ii) Q is the finite set of states containing initial state q_{start} , (iii) $F \subseteq Q$ is a set of halting
 652 states, and (iv) δ is a transition function $(Q \setminus F) \times \Sigma^k \rightarrow Q \times (\Sigma \times \{L, S, R\})^k$.

653 Throughout this paper, we will assume $\Sigma = \{0, 1, \perp\}$ for simplicity and with loss of generality.

654 **Probabilistic Turing Machines.** A probabilistic Turing Machine (PTM) is a Turing machine with an
 655 additional read-only coin tape full of independent and uniformly random coins.

656 The PTM model is potentially more powerful than the deterministic Turing machine (DTM) model.
 657 An example of a computational problem that can be solved in polynomial time by a PTM but still not
 658 known how by a DTM is the polynomial identity testing problem (PIT) (see, e.g. [81]). In fact, it is a
 659 central question in complexity theory, well-known as the $\text{BPP} = ? \text{P}$ problem, whether any decision
 660 problem solvable by a polynomial-time PTM can also be solved by a polynomial-time DTM.

661 We say that a (deterministic or probabilistic) TM T is *oblivious* if the tape head movements of
 662 T running on input x depend only on the input length $|x|$. That is, T makes the same sequence
 663 of head movements for all inputs x of the same length. [82, 83] proved that: for every multitape
 664 DTM T running in $O(t(n))$ time, there is an equivalent oblivious two-tape DTM T' that runs in
 665 $O(t(n) \log t(n))$ time. Furthermore, as observed by [84], this result also holds for all relative Turing
 666 machines, including PTMs. Specifically, for any multitape PTM T running in $O(t(n))$ time, there is
 667 an equivalent oblivious two-tape PTM T' running in $O(t(n) \log t(n))$ with an additional ready-only
 668 coin tape. So, w.l.o.g., in this paper unless otherwise specified, whenever we refer to DTMs or PTMs,
 669 we refer to two-tape oblivious DTMs or PTMs respectively.

670 **Universal Turing Machines.** A universal Turing machine (UTM) is a TM that can simulate
 671 the execution of every other (deterministic or probabilistic) TM T given T 's description as input.
 672 Specifically, we encode PTMs as Boolean strings in a prefix-free way. A UTM is a PTM U that takes
 673 the concatenations of the encoding of a PTM T and an input x , and outputs the (possibly randomized)
 674 $T(x)$. UTMs capture the notion of a “general-purpose programmable computer”, which is a single
 675 machine that can be adapted to any arbitrary task provided an appropriate program is loaded. We
 676 remark that the parameters of a UTM, such as alphabet size, number of states, and number of tapes
 677 are fixed, though the TM being simulated could have much more parameters.

678 **B Proof of Theorem 2**

679 *Proof.* We first adapt T by introducing a lazy sampling of the coin tape. The coin tape is initially
 680 empty, filled with blank symbols \perp , and will be assigned random coins on the fly during execution.
 681 At one step, if T reads a blank symbol \perp from the coin tape, it first tosses a fair coin and writes the
 682 result on the coin tape at the current head position. Note that the adapted T , denoted by T' , runs for
 683 at most $2t(n)$ steps, since each original step may include an additional coin-tossing operation.

684 Next, we demonstrate how a transformer can simulate T' with at most $2t(n)$ CoT steps. We adapt the
 685 proof of Theorem 2 in [20]. For the i -th step of T' , let $h_i^\tau \in \mathbb{Z}$ and $\gamma_i^\tau \in \Sigma$ denote the head position
 686 and the content on tape τ , and let $q_i \in Q$ denote the state. Let $\Delta := Q \times \Sigma^2 \times \{L, S, R\}^3$, and $\delta_i \in \Delta$
 687 denote the log at the i -th step, indicating the state entered, symbols written, and directions moved.
 688 The crucial observation is that the tape contents at the current head positions can be reconstructed
 689 from the input x and the previous logs $\delta_0, \delta_1, \dots, \delta_{i-1}$.

690 As shown in [20], a Transformer can first obtain all arguments $(q_{i-1}, \gamma_i^1, \gamma_i^2, \gamma_i^3)$ for the transition func-
 691 tion. Suppose tape 3 is the coin tape. If $\gamma_i^3 \neq \perp$, which means that the i -th step of T' will be determin-
 692 istic rather than the coin-tossing operation, then the Transformer computes $\delta_i = \delta(q_{i-1}, \gamma_i^1, \gamma_i^2, \gamma_i^3)$
 693 with a feedforward net outputting the one-hot encoding of δ_i . If $\gamma_i^3 = \perp$, which means that the i -th step
 694 is a coin-tossing operation, then the Transformer outputs the equally weighted linear combination of
 695 the one-hot encodings of $(q_{i-1}, \gamma_i^1, \gamma_i^2, 0, S, S, S)$ and $(q_{i-1}, \gamma_i^1, \gamma_i^2, 1, S, S, S)$. The vector outputted
 696 by the feedforward net is then processed by the final token classification head, which is a softmax
 697 function. One can check that the Transformer exactly simulates the i -th step of T' . \square

698 **C Two remarks on Thesis 1**

699 **Remark 2.** *There is ongoing debate as to whether quantum computers falsify ECT. In particular, it*
700 *is a central problem in quantum computing, well-known as the BQP =?BPP problem, whether all*
701 *decision problems solvable by polynomial-time quantum computers can also be solved by polynomial-*
702 *time PTMs. If ECT is falsified by quantum computers, then a quantum variant of Transformer that*
703 *can simulate universal quantum computers [85, 41, 86, 87] might be necessary to achieve AGI.*

704 **Remark 3.** *It is widely accepted that a human brain can be modeled as a complex computational*
705 *system (say, a huge neural network) following classical physical laws, and thus can be simulated by*
706 *a PTM [88, 89]. However, this traditional view of the brain as a classical system was challenged by*
707 *[90, 91]: they argued that the brain utilizes quantum mechanical effects (e.g., quantum coherence*
708 *or entanglement) for reasoning and recognition, and human consciousness is even non-algorithmic,*
709 *though still lack empirical validation.*

710 **D Data requirement in Section 4.3**

711 We discuss the types and amounts of data required in our framework to approximate AIXI using two
712 Transformers.

- 713 • **Format Flexibility:** Data can be multi-modal (e.g. text, images) and multi-task (understanding,
714 reasoning, etc.), aligning with current LLM practices. In fact, such diversity is essential to achieve
715 increasingly universal models [30]. To further enrich the diversity, we suggest incorporating the
716 synthetic UTM data. This types of data complements human-generated dataset by providing
717 coverage over algorithmic patterns that may be underrepresented in the existing dataset. Moreover,
718 the UTM data also provide a theoretical guarantees of convergence to SI.
- 719 • **Quantity:** Extensive data remain crucial for teaching induction. For deduction, however, the
720 model can self-generate data (e.g., proposing solutions in math tasks and refining them via
721 verification feedback).