

# ON INHERENT LIMITATIONS OF GPT/LLM ARCHITECTURE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this paper, we show that reasoning/proving issues of GPT/LLM are an inherent logical consequence of the architecture. Namely, they are due to a schema of its prediction mechanism of the next token in a sequence, and randomization involved into the process. After natural formalization of the problem into a domain of finite graphs,  $G(\omega)$ , we prove the following general theorem:

*For almost all proofs, any learning algorithm of inference, that uses randomization in  $G(\omega)$ , and necessitates veracity of inference, is almost surely literal learning.*

In the context, "literal learning" stands for one which is either vacuous, i.e.  $\forall x [P(x) \implies Q(x)]$  where  $P(x)$  is false for every  $x$ , or create a random inference from a false assumption (hallucination), or it essentially memorizes the inferences from training/synthetic data.

## 1 NOTATION AND TERMS

GPT/LLM stands for algorithmic representation of transformer with attention viewed as a main inference mechanism for *LLM*.

$G(\omega)$  stands for the infinite set of finite graphs, and the first-order model on the set as a domain, where the connectivity of nodes  $n_1 \sim n_2$  stands for the first countable ordinal.

Formal definitions for the language of the first-order theory can be found in Appendix B. It contains all the necessary information on the 0-1 law.

In this setting, a fault in  $G(\omega)$  is an erroneous proof – that is, a chain of thought containing a false implication ( $n_1 \not\sim n_2$ ), or a false assumption (node  $n_1$  represents falsehood).

"Randomization" on a probability space means that the events inference mechanism admits variability in the next token selection (such as random seed initialization, temperature, beams search, etc.).

A "Learning algorithm" is a machine learning algorithm that, after being trained on data, produces output based on the training.

"Literal learning" stands for one which is either vacuous, i.e.  $\forall x [P(x) \implies Q(x)]$  where  $P(x)$  is false for every  $x$  or creates a random inference from a false assumption (hallucination), or it essentially memorizes the inferences from training/synthetic data.

First-order logic terms can be found in B.

0-1 Law for graphs is in B.1.

## 2 INTRODUCTION

On one hand, GPT architecture for LLMs demonstrated significant progress in a generative manifestation of summarizations, chat, and representations of materials. On the other hand, the architecture displayed multiple negative effects such as hallucinations, falsehoods, degrading generalization, performance degradation, and alike (e.g., (Yadlowsky et al., 2023a)). In rigorous contexts (where one

requires a consistent mathematical reasoning or a formal proof), the results are consistently discom-  
forting ((Chen et al., 2023), (Hagendorff et al., 2022), (Dziri & et. al., 2023)).

Recently a few authors pointed out various limitations (cf. e.g., (Liu et al., 2023), (Mikhaylovskiy  
& Churilov, 2023), and (Asher et al., 2023a)). Nonetheless, there have been suggested possible  
remedies ((Sel et al., 2023), (L. et al., 2023), and (Z. et al., 2022)).

In this paper, we show that these issues are an inherent *logical* consequence of the GPT architec-  
ture. As a result, multiple phenomena of transformer inference limitations can be explained from  
purely logical view; in particular, some results of (Dziri et al., 2023) can be obtained that way. It is  
shown that some limitations addressed in the paper (e.g., problem of increasingly large parallelism  
requirement) can be relieved with changing a type of attention.

In general, it appears that there is a latent belief in contemporary literature that *all* limitations of  
technology can be resolved within the governing transformers’ model. The goal of this paper is to  
prove that the architecture is inherently limited in case of inference that required rigor; thus, these  
imitations are fundamental and innate.

A crucial observation is a scheme of transformer prediction mechanism of next token in a sequence.  
Using a natural formalization of the problem into the domain of (standard) finite graphs  $G(\omega)$ , we  
prove the following theorem:

*For almost all proofs, any learning algorithm of inference, that uses randomization in  $G(\omega)$ <sup>1</sup>, and  
necessitates veracity of inference, is almost surely literal learning.*

We provide a few proofs for this statement. For a quick look at a formal proof, refer to C.2. To  
develop an intuition for the phenomenon, there is an informal proof (1) where consideration is  
around 0-1 law.

In this form of the inherent limitation, there are a few basic assumptions that need to be addressed.  
Namely, we work in the first-order model  $G(\omega)$  (appendix B.1), where connectivity between nodes  $a$   
and  $b$ , expressible by the first order formula, represents the validity of the implication  $a \implies b$  (we  
are going to use  $a \rightarrow b$  for the implication as well, interchangeably). For our purposes, the graph  
do not have to be directed since only a countable enumeration of the nodes is necessary. Moreover,  
any algorithmic randomization on the nodes allows us to view  $G(\omega)$  as a suitable probability space  
(C.2). Finally, literal (or vacuous) learning is defined as such that, almost surely, the proof chain,  
generated by an algorithm, is equivalent (in  $G(\omega)$ ) to one in a training dataset.

A few corollaries follow. For instance, if its formulation is somewhat original, it is easy to notice the  
issue of solving mathematical problems with LLMs in the case of even low complexity task. Since  
its solution is unlikely to be found in a holistic form in a training dataset, a correct proof is not to be  
expected.

It is because, in a rigorous context, GPT has exponentially decreasing odds of finding a valid proof  
of the result unless it simply “repeats” a known proof, perhaps with trivial modifications (Corollary  
3). Another corollary is that degradation of performance is of exponential rate by length of a proof.  
In other words, an attempt to prove a complex enough statement virtually has no chance of being  
successful.

In a rigorous context of generating a proof, GPT virtually has no chance to find a valid proof of  
the result unless it simply “repeats” a known proof, perhaps with trivial modifications (Corollary  
3). Another corollary is that the degradation of performance has an exponential rate by the length  
of a proof. In other words, an attempt to prove a complex enough statement with GPT/LLM has  
virtually no chance to be successful.

In a novel rigorous context (i.e., when GPT-based architecture is looking to prove a new result, for  
instance, a hypothesis), that is virtually impossible even for a long enough fragment of the proof.  
The probability of success becomes infinitesimal quickly for either a fragment of possible proof or  
a weaker non-trivial statement. That also was empirically shown for data mixtures in (Yadlowsky  
et al., 2023b).

This has been recently confirmed experimentally in ((Hubinger et al., 2024)). Moreover, the logical  
view approach enables us to discover the same limitation patterns for LLM and auto-regressive next-

---

<sup>1</sup>defined in B.1

108 token predictors even though the latter are universal ((Malach, 2023)).

109 The logical view approach is more effective in generalization by varying a domain obeying the  
110 generic 0-1 law.

111 For instance, the results in ((Dziri et al., 2023)) can be obtained as a partial case in the proof of our  
112 main result. The 0-1 law variant for polynomial decision problems ((Blass et al., 1998)) is used for  
113 a more instructive proof.

114 Tools like FunSearch ((Romera-Paredes et al., 2023)) contribute to searching for solution specifica-  
115 tions, instead of providing an actual inference.

117 Recent improvements in LLMs such as a 1 million-long context window would make vacuous infer-  
118 ences quite probable in a standard context, with tantamount consequences to the dangling pointers in  
119 software. In the novel context, it is almost certain that the resulting proof will be incorrect, present a  
120 hallucination, or both. Thus, the expectations that ChatGPT-4.0 or a similar model (e.g., the agents)  
121 would soon be able to reason and plan like a person seem unfounded.

122 For the alignment problem, we only note that our approach is fundamentally different from that  
123 of (Wolf & et. al., 2023) since our methodology is ultimately based on considerations within the  
124 first-order logic of appropriate Random Graph theory while theirs is purely statistical.

125 Note also, that in terms of paper (Nasr & et. al., 2023), leaving the adversarial context of it, we essen-  
126 tially proved that in a rigorous context, given sufficient complexity, LLMs are able only to memorize  
127 the existing proofs in the training dataset. Thus, one cannot expect these models to produce a novel  
128 non-trivial proof. In terms of (Nasr & et. al., 2023), discoverable and extractable memorizations  
129 coincide, given sufficient complexity of a statement  $P$ .

130 In other words, given sufficient complexity of a statement  $P$ , the prompt "Please prove statement  
131  $P$ " would generate a memorized proof if one exists in the training dataset, present an incorrect  
132 proof, or hallucinate (cf. (Chen et al., 2023) as in (Asher et al., 2023b), and (Mikhaylovskiy &  
133 Churilov, 2023)). Similarly, an attempt to fix the LLM (e.g., GTP-4) bugs with LLM critique tech-  
134 nique ((McAleese & et. al., 2024)) will have only a limited scope of applicability.

135 In the paper (Dohare et al., 2024), the authors note that the deep-learning system's performance  
136 degrades during extended training on new data. A method, proposed in the paper, requires random-  
137 ization to establish elasticity; so it is likely that not just LLMs but also traditional  $FFNs$  admit  
138 a version of the main result on inherent limitation. Moreover, an attempt to enrich a model with  
139 synthetic data will go only so far as well as an emerging representation of underlying abstraction  
140 (cf. (Jin & Rinard, 2024) where the context is not rigorous).

141 Additionally, the main theorem is a formal statement of statistical nature, applicable to a more gen-  
142 eral context than GPT/LLM, albeit GPT/LLM is a target example. The future of the GPT/LLM  
143 may bring additional mechanisms into discussion, e.g., RAG integration, advanced planning, and  
144 domain-specific tuning. One can claim that these additional mechanisms will not be sufficient to  
145 overcome the limitations of the core foundation of transformer-based generation for rigorous prov-  
146 ing. However, we are not making this claim since the algorithmic representation of these mecha-  
147 nisms may vary on prompt tuning, RL, adapters, and most importantly, RLHF. Thus, post-training  
148 techniques are out of the scope of the main result. Note that the result does not contradict Malach  
149 (2023) since the scope of the two approaches is entirely different. The formalism of that work is  
150 focused on approximating any Turing function while this paper presents a first-order logic view on  
151 rigorous reasoning at large.

### 152 153 154 3 MOTIVATION/PRELIMINARIES

155 To demonstrate the model that proves the inherent limitation phenomena, we need to formalize  
156 logically the mechanism employed by GPT/LLM to predict the next token in a sequence. It turns out  
157 that randomization used by  $GPT$  architecture (namely "temperature") is the main reason. However,  
158 it isn't easy to devise an alternative when dealing with training on a large text corpus. Note that the  
159 architecture becomes too "predictable/plain" if we choose the most probable pattern in the list of  
160 candidates for the next token. There, we would need a certain randomization to become "creative".  
161

As we will see, that necessary hack is sufficient to preclude the architecture from ever succeeding in a rigorous context, in a formal setting when we need to infer our next supposition with strict regard to its veracity. A good example would be generating proof for a theorem. Note that, despite infamous issues with Generative Learning with rigor in this context, there have been a few feasible attempts to attack this problem that way (e.g., (Saparov et al., 2023)); moreover, there was a claim that we may be able to "recover" from an ostensibly systematic LLM model faltering in mathematical settings ((Shi et al., 2022)). As is known, these attempts were largely unsuccessful. Our results offer a logical explanation of why.

To that end, we introduce formalism to make the subject rigorous enough to have a logical view. Namely, we present a simple first-order theory on the language of (random) graphs where one can state that the generative inference that admits randomization on implications will almost necessarily lead to logical faults (i.e., with probability 1). This result is based on a 0-1 law (and its variations) in (random) graphs theory.

### 3.1 DEFINING THE CONTEXT

We can assume that one can enumerate all the inferences using  $[n]$ , since there is only a countable number of (finite) proofs on a countable number of entities. For a graph in  $G(\omega)$ , define property  $\mathcal{A} := \{\exists \text{ nodes } e_1, \dots, e_j \text{ forming chain } (e_0 \rightarrow e_1 \dots e_i \dots e_j \rightarrow e_m) \text{ for inference } e_0 \rightarrow e_m\}$ . This definition is well-formed since  $\mathcal{A}$  is expressed as a first-order sentence in the first-order logic theory for  $G(\omega)$ , and the axiom of foundation<sup>2</sup>. For chains above, we need to verify that these are first-order expressions. A suitable framework for this is that of least fixed point extension (cf. (Grohe, 2017)). Namely, if the " $\sim$ " is a connectivity relation, then a chain  $C(e_0, e_t)$  where  $e_0$  and  $e_t$  represent a proof starting node  $e_0$  and terminal node  $e_t$  respectively, can be expressed as follows:

$$C(e_0, e_t) \leftarrow ((e_0 = e_t) \vee \exists e_i (C(e_0, e_i) \wedge e_i \sim e_t)) \quad (1)$$

Then, by 0-1 Lemma in B.1, we have two possibilities, namely:  $p = \lim_{n \rightarrow \infty} \mathbb{P}(G_n(\omega) \in \mathcal{A}) = 0$  or it is equal to 1. If it is zero, then no valid proof can be found within the context in the first place. Therefore,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n(\omega) \in \mathcal{A}) = 1$ . By Lemma 0 in B.1, it follows that  $G(\omega) \models \mathcal{A}$ . However, it also means that our inference follows a literal graph representation from the original (i.e., from the given training set). Similar consideration is possible for a novel vs. not novel context. Thus,  $p \neq 1$ . In this case, we create a hierarchy in  $G(\omega)$  as follows.

Consider chain  $(e_0 \xrightarrow{\psi_1} e_1 \dots e_i \dots e_j \xrightarrow{\psi_k} e_m)$  and formula  $\psi := \psi_1 \wedge \dots \wedge \psi_k$ . Clearly,  $\psi$  is true in  $G(\omega)$  for any inference of  $e_m$ . But that means that we again have a "literal" learning. Otherwise, since  $p$  is not 1, we will have a "fault" for sufficiently large  $n$ .

In (Blass et al., 1998), the authors proved a version of the zero-one law for binary sequences and, within the context, a decision problem. Our formal proof is a generalization to a class of algorithms in which logical inference admits a standard graph representation. Namely, we just proved the following:

**Theorem 1** *For almost all proofs, any learning algorithm of inference, based on randomization in  $G(\omega)$ , that necessitates veracity of inference, is almost surely literal learning.*

The complete formal proof is in C.2. Before, we established that there is a natural model for the inference and pointed out the limitations associated with it. In other words, the algorithm almost surely fail unless it is vacuous. ■

**Theorem 2 (reformulation of the theorem 1)** *Given the graph model of inference for machine learning, the only algorithm based on randomization, that also necessitates veracity of inference, is almost surely "literal" learning. In other words, for a sufficiently long proof, any algorithm that randomly deviates from the training data will fail with a probability of 1.* ■

<sup>2</sup>In a second-order logic, one can quantify over sets of domain elements; in the first-order logic, one can quantify over elements only.

**Corollary 3** Within a rigorous inference context, almost surely, no randomization of the prediction scheme of proof patterns can discover new (unknown) non-trivial valid statements. This can be easily explained: since any degradation is inherited in the foundational graph, the subsequent inferences on the trained data tend to deviate from already shortened erroneous paths thus multiplying the faults. ■

**Example of an inference problem that exceeds the current capabilities of generative learning**

*Elementary example.* Proving the statement: "for every number  $2^n$  for any natural  $n$ , there exists a number  $k$  such that  $2^n * k$  does not have zeroes in its digital representation".

*Non-elementary example.* Dedekind numbers sequence.

*Benchmark Examples* Multiple benchmark examples can be found in Glazer (2024).

**Theorem 4 (Generalizations of main result)** Any algorithm of learning enforcing veracity, admitting a 0-1 domain cast as random graphs, is almost surely vacuous.

**Proof** This is the context where proof of theorem 1 is fully applicable. ■

Within the view adopted herein, there is an interesting example of a decidable theory that admits a 0-1 random graph domain yet its classifier comparison is not expressible in its first-order logic. Therefore, it is an example of a.s. learning algorithm with randomization which is ultimately decidable but vacuous and does not support any notion of expressible first-order classifier comparison; thus, there is no feasible notion of fairness for classification tasks.

## 3.2 ELEMENTARY PROBLEMS

### 3.2.1 QUESTION

Question: Can one cut a scalene triangle into two congruent scalene triangles? Answer: Microsoft Bing Copilot: "Certainly! Let's explore how we can cut a scalene triangle into two congruent triangles". Then Copilot generates two methods to create the cut: angle bisector method, and perpendicular bisector method which would work only for isosceles triangles, completely ignoring the fact that the original triangle is scalene. The fault is that the bisectors will not divide the opposite side into two equal segments. So, the subsequent application of angle-angle-side and side-side-side postulates is invalid. However, Bing Copilot "insists" and suggests the question: "Can you cut any triangle into two congruent triangles?" The predictable Copilot's answer is now that any triangle can be, while referring to the very answer to the previous question as a given (one can only note that it looks "logical"). Needless to say the process would be easily repeated with all sorts of fallacious geometrical statements. If the user points out an occasional contradiction, the Copilot produces a loop or changes the subject.

Claude (Antropic): This was a different experiment where the author tried to "teach" Claude to solve the aforementioned elementary problem. It took a few trials before Claude arrived at a plausible reformulation of the problem. A very positive result was, despite an inability to present a complete rigorous proof, Claude came up with a plan for how to obtain the proof. However, after a few unsuccessful attempts to implement the plan, and a few homework sessions later, we agreed that reaching the point is beyond Claude's capabilities yet. With the three assistants, our experience with Claude was the most pleasant and sensible.

Google Gemini (Bard) The result is similar to Claude's. After a few clarifications and direct clues, Gemini produced the following in bold: "Therefore, I cannot confidently claim to have proven the statement about the impossibility of cutting a scalene triangle into two congruent scalene triangles". Note that it is, formally, a weaker statement than what is necessary to solve the problem asked.

Similar results were obtained for other chatbots, e.g., Perplexity.

Incidentally, the paper (Trinh et al., 2024) depicts good results on solving geometry problems of Olympiad's level. First, we have to note that, because the first-order theory of Euclid geometry is

```

270 import Mathlib.Data.Nat.PrimeFin
271 import Mathlib.Data.Nat.Factors
272 import Mathlib.Data.Set.Finite
273
274 theorem exists_infinite_primes (n : ℕ) : ∃ p, n ≤ p ∧ Nat.Prime p :=
275 let p := Nat.minFac (Nat.factorial n + 1)
276
277   have f1 : Nat.factorial n + 1 ≠ 1 := ne_of_gt <| Nat.succ_lt_succ <| Nat.factorial_pos _
278   have pp : Nat.Prime p := Nat.minFac_prime f1
279   -- have ppc := Nat.Prime ↑pp
280   -- have ppc := minFac_to_nat pp
281   have np : n ≤ p :=
282     |le_of_not_ge fun h =>
283       |have h1 : p | Nat.factorial n := Nat.dvd_factorial (Nat.minFac_pos _) h
284       |have h2 : p | 1 := (Nat.dvd_add_iff_right h1).2 (Nat.minFac_dvd _)
285       |pp.not_dvd_one h2
286     ⟨p, np, pp⟩

```

elementary in the logical sense (decidable), the task is achievable by a universal algorithm since we can work in the decidable first-order theory of  $\mathbb{R}$ .

Since the output is natural language (rather than in a code for an automated prover, unlike in the approach of (Zheng et al., 2022)), it isn’t easy to assess the solution’s performance. Because this transformer is trained on synthetic data, and proofs are relatively short by nature of the problems involved, due to our main result, likely, the solution does not exceed a threshold of vacuous/literal learning overall. A more formalized approach is presented in (Krueger et al., 2021).

In (Nezhurina et al., 2024) are more examples of basic reasoning breakdown for foundational industrial models.

### 3.3 NON-ELEMENTARY PROBLEMS

Above is an example of ”hallucinating”: a formal proof of the infinitude of prime numbers in Lean 3 or 4.

Here is the critical fragment of the proof where randomization played a key role:

```

302 { by_contradiction ,
303     have h1 : p | fact N := dvd_fact (min_fac_pos M) a ,
304     have h2 : p | 1 := (nat.dvd_add_iff_right h1).
305     mpr (min_fac_dvd M) ,
306     exact prime.not_dvd_one pp h2 } ,
307 { exact pp }

```

This latest fragment renders the proof unusable. One correct version is placed above, which is unlikely to be found elsewhere (since we use an explicit ”Nat.” prefix).

In (Nguen & Sarah, 2022), the authors describe multiple patterns of software development that reflect erroneous or sub-optimal code generated by Copilot. This leads to an elevated code churn and downward pressure on code quality in *GitHub*. Another survey, (Kabir et al., 2024), shows that coding questions generate up to 50% of errors. A similar study is conducted in (Macmillan-Scott & Musolesi, 2024).

### 3.4 DISCUSSION

#### Discussion - Primary

These results easily explain the phenomenon of ”hallucinations” and brittleness of the GPT models in a rigorous context. It also means that LLMs is unlikely to discover any new mathematical result of sufficient strength.

324 **Discussion-Datasets** In (Gendron et al., 2023), is shown that the baseline dataset construction for  
325 rigorous learning needs to be a formal exercise. Consider the task of equation completion in which  
326 one has to predict a missing symbol. Since this is perfectly aligned with the main premise of LLM  
327 based on transformers, one can expect that the success rate for this task will be quite high. As is  
328 shown in the paper, this is not the case. An associated (and well-known) phenomenon of a plateau  
329 of performance and subsequent degradation in an exponential fashion manifests in the same way as  
330 for the generic sequence case. Similarly, few authors summarize a few problems in the answers in  
331 contemporary systems associated with a low P/R w.r.t citation usage from the underlying sources.  
332 These experimental results are not for the rigorous context.

333 There is a widespread belief that because the training set contains "everything", any result, including  
334 novel ones, can be proven using symbolic inference from the corpus. However, it is just not the case.  
335 It is well-known that any mathematical problem of significance requires one or multiple critical  
336 insights that are just not to be found. These are not combinations of known results (or tactics), but  
337 rather completely new, albeit inevitable, ideas. For instance, for some long-standing problems, new  
338 fields of mathematics had to be created, representing a new body of knowledge. Thus, generalizing  
339 the LLM solution for these targets is a task of yet another level of complexity for which the method  
340 is not suited. Moreover, as we show below, it is guaranteed to fail. The inevitable conclusion is that  
341 the apt inference model has to be more deferential to logic.

342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

## 378 A MAIN RESULT

379  
380 In this section, we assume a natural representation of a proof by a path from a node  $e_0$  to the node  $e_t$   
381 in graph  $G(\omega)$  in appendix B.1 which is, in a suitable enumeration, corresponds to a premise/target  
382 statements  $e_0, e_t$  accordingly.

383 **lemma 1.** (*Accumulating errors Lemma*). *Assuming independence of faults in  $G(\omega)$  representing*  
384 *proofs, the probability of no fault proof tends to zero exponentially over its length.*

385  
386 *Proof.* We can assume that faults are independent since the semantics of a formal inference are out  
387 of scope <sup>3</sup>. Let labeled graph  $G$  represent proofs (chains) of enumerated statements (nodes) where  
388 each label is a probability that the chain ending with the corresponding node contains an error. Then  
389 we have:

$$390 \mathbb{P}(\text{no fault proof}) \leq \exp(-\mathbb{E}(\text{number of faults})) \quad (2)$$

391 Since the right side of the equation tends to zero, we have:

$$392 \lim_{n \rightarrow \infty} \mathbb{P}(\text{no fault proof of length } n) = 0. \quad (3)$$

393  
394 This proves the statement of the lemma. ■

395  
396 **example 1.** Since GPT has no semantic notion on the entities involved, we can assume the lemma  
397 is fully applied to the GPT rigorous context.

398  
399 **Corollary 4** Assume an algorithm admitting model  $G(\omega)$  for inference and using randomization.  
400 Then, the rate of (correctness) decay is exponential over the proof length.

401  
402 **Proof** The proof is similar to a usual consideration for a set of independent events in a clas-  
403 sic probability space generating a fault. The key observation is that, once a fault in the chain of  
404 inference occurs, it is thus erroneous in the chain subsequently. This proves that the correctness rate  
405 decays exponentially over the length of proof. ■

406  
407 **(Heuristic note)** In a few papers, this phenomenon has been shown experimentally. More-  
408 over, it has a few incarnations. These are hallucinations (when there are no references, supporting  
409 an inference), erroneous statements (falsehood, incorrect generalization, non sequitur, etc.), and a  
410 general misalignment. Exponential decay is also noted in a few papers; our result (Corollary 4)  
411 shows for all non-trivial (complex enough) tasks, including performance degradation on synthetic  
412 data in an autophagous loop.

413  
414 **theorem 1.** (*Inherent GPT/LLM Limitation*).

415 *Any algorithm of inference, based on randomization on  $G(\omega)$ , that necessitates veracity of inference,*  
416 *is almost surely literal learning.*

417  
418 *Proof.* (Informal) We give two proofs of the statement. To develop a theoretical intuition, we start  
419 with the one below. The second one, more instructive and rigorous, is in C.2.

420 Note that we can assume that one can enumerate all the inferences using  $[n]$ , since there is only a  
421 countable number of (finite) proofs on a countable number of entities (statements). Without loss of  
422 generality, for that representation of entities, we can assume that node (vertices)  $e_i$  implies  $e_j$  only  
423 if they are connected; we do not need to impose any order on the nodes.

424  
425 For a generative model, that would be enumeration for a proof generated for a particular prompt,  
426 say, prove that  $e_k$  implies  $e_l$ . Moreover, we can assume that a generic proof is an actual chain of  
427 thought, i.e., we have a finite sequence of distinct nodes, connected via regular paths, with possible  
428 cycles which would reflect the equivalency of the statements. The underlying training graph for this  
429 is not necessarily connected, but the model output has to contain a path from the premise to the  
430 desired conclusion.

---

431 <sup>3</sup>GPT algorithm does not follow the syntax of the first-order theory – instead, it uses randomizing and  
inferred statistics. It has no notion of non-statistical meaning.

In appendix B, we define the first-order language of graphs used in an associated model,  $G(\omega)$ .

Thus, the path (or "chain of thought") is just a sequence of tuples  $(s_k \sim s_l)$  where sign " $\sim$ " represents adjacency for vertices  $s_k$  and  $s_l$  and there is a path

$$e_s \sim e_1 \wedge e_1 \sim e_2 \wedge \dots \wedge e_k \sim e_t \quad (4)$$

Now, there are two possibilities:

1. The path (4) exists in the training set (not a novel context).
2. The path (4) does not exist in the training set (a novel context).

Consider the first-order formula  $\phi(\cdot) = e_s \sim e_1 \wedge e_1 \sim e_2 \wedge \dots \wedge e_k \sim e_t$ . Then, again, we have two possibilities. Namely, by 0-1 law (B.1), we have:

$$\lim_{n \rightarrow \infty} \mathbb{P}(G(\omega) \models \phi) = 0 \text{ or } 1. \quad (5)$$

Thus, we have four possibilities, namely:

1. The limit (5) is equal to zero and the path (4) does not exist in the training set.
2. The limit (5) is equal to zero and the path (4) exists in the training set.
3. The limit (5) is equal to one and the path (4) does not exist in the training set.
4. The limit (5) is equal to one and the path (4) exist in the training set.

For each of these, we also need to consider the cases of model temperature, normalized to probability  $p$ , equal to zero or one, or between zero and one. We can assume the following for these cases:

For the case **1**, it is nearly obvious that, within the context, no valid proof can be found almost for sure in the first place either if we try literal learning, falling into a novel context, or varying the probability  $p$  between zero and one - we apply Accumulating errors Lemma since GPT is an accumulating errors algorithm. The latter manifests as a phenomenon of accumulating errors for sufficiently complex (lengthy) proofs.

The case **2** is more interesting. Despite having proof in the training set and a chance of literal learning, we use probability  $p$  other than one. As a result, we are having the phenomenon of accumulating errors described above.

The case **3** is the most interesting – we are in a novel context – and may follow fragments of the proof, somehow creating the final proof as an assembly. Note, we chose  $p$  equal to one. It means that we are trying to assemble the required proof in pieces. The problem is equivalent to finding paths among potentially connected pieces. However, we can simply apply B.1 and note that since  $p$  is equal to 1, we have:

$$\lim_{n \rightarrow \infty} \mathbb{P}(G(n) \models \phi) = 1 \Leftrightarrow G(\omega) \models \phi. \quad (6)$$

Therefore, for ever-growing complexity and length of proofs, we have to follow ever-growing fragments of proof literally which means we have them in the training set. That is literal learning or we have a contradiction with the assumption of this case.

The case **4** is literal learning, by definition.

The conclusion is that *almost for sure*, only literal learning, has a chance of generating an error-free proof.

$\lim_{n \rightarrow \infty} \mathbb{P}(G(n) \in \mathcal{A}) = 0$  or it is equal to 1. If it is zero, then no valid proof can be found within the context in the first place. Therefore,  $\lim_{n \rightarrow \infty} \mathbb{P}(G(n) \in \mathcal{A}) = 1$ . By Lemma 0, in the first-order theory

486 for the language of random graphs, it follows that  $G(\omega) \models \mathcal{A}$ . For  $p = 1$ , it is possible. However, it  
487 also means that our inference follows a literal graph representation from the original (i.e., from the  
488 given training set). Thus,  $p \neq 1$ . In this case, we create a hierarchy in  $G(\omega)$  as follows.  
489

490 Consider chain  $(e_0 \xrightarrow{\psi_1} e_1 \dots e_i \dots e_j \xrightarrow{\psi_k} e_m)$  and formula  $\psi := \psi_1 \wedge \dots \wedge \psi_k$ . Clearly,  $\psi$  is true  
491 in  $G_\omega(p)$  for any inference of  $e_m$ . But that means that we again have a "literal" learning. Otherwise,  
492 since  $p$  is not 1, we will have a "fault" for sufficiently large  $n$ . ■

493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

## B FORMAL DEFINITION FOR LANGUAGE $L$

Let  $L$  be a language (an extension of a basic formal logical language,  $L_0$ ).

**Definition and base notations** The set of  $L$ -terms is the smallest set  $L_t$  such that contain all constant symbols of  $L$ , all variables, and if  $t_1, t_2, \dots, t_n$  are in  $L_t$  then for any n-ary function symbol  $f$ ,  $f(t_1, t_2, \dots, t_n)$  is also in  $L_t$ . Set  $L_a$  of atomic formulas are represented by the properties:

- (1) if  $t_1$  and  $t_2$  are terms then  $t_1 = t_2$  is in  $L_a$ , and
- (2) the corresponding n-ary function symbols are also in  $L_a$ .

In other words, the set of all formulas in  $L$  (expressions, sentences - herein, we use these interchangeably) is the smallest set containing all atomic formulas and closed under logical connectives  $\vee, \wedge, \neg, \rightarrow, \leftrightarrow$ , quantifiers  $\exists, \forall$ , equality symbol " $=$ ", parenthesis "(" and ")", and variables. For our purposes herein and simplicity, it is sufficient to consider that theory in language  $L$  is a set of sentences in first-order logic over  $L$ . We also assume first-order logic with equality; in other words, only normal models are employed. Thus, the models, considered herein (e.g., Erdős-Rényi or finite graph model for random graphs, are normal).

The main language in this paper is that of graphs.<sup>4</sup> We denote  $\mathbb{G}_L$  the first-order theory over language of graphs  $L$ . One convenient (and usual) laxity talking about expressions and formulas in  $L$  is using  $L$  and  $\mathbb{G}_L$  interchangeably.

### B.1 0-1 LAW FOR GRAPHS $L$

We introduce a few known formulations for the 0-1 law for finite graphs.

**0-1 Lemma 0** For any first-order formula  $\phi$  and graph  $G$  in  $\mathbb{G}_L$  (with the equivalent notation  $G(\omega)$  which is intuitively more suitable), let

$$G_{n,\phi} = \frac{|\{G \models \phi : |G| = n \text{ and } G \text{ is a graph}\}|}{|\{G : |G| = n \text{ and } G \text{ is a graph}\}|} \quad (7)$$

Then  $\lim_{n \rightarrow \infty} G_{n,\phi}$  is 0 or 1.

*Proof.* Refer to, e.g., (Fagin, 1976).

This can be reformulated as

**0-1 Lemma** For any property  $\mathcal{A}$  that can be described by a first-order expression  $\phi$  and  $G_n = \{G : |G| = n \text{ and } G \text{ is a graph}\}$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(G_n \in \mathcal{A}) \in \{0, 1\} \quad (8)$$

To wit (assuming notations for  $G(\omega)$ , a set of all finite graphs, and its associated domain  $G(\omega)$ , up to isomorphism):

**Lemma 0, reformulation** For any graph  $G_n \in G(\omega)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n) = 0$  or 1. The equivalent statement is as follows: for any first-order expression  $\phi$  in theory of  $\mathbb{G}_L$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n \models \phi) = 0$  or 1.

We can also say that  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n \models \phi) = 1 \Leftrightarrow G(\omega) \models \phi$ .

**Lemma 0, reformulation** For any *random* graph  $G_n \in G(\omega)$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n) = 0$  or 1. The equivalent statement is as follows:  $\forall$  0-1 probability  $p$  and a first-order expression in theory of Random Graphs,  $\phi$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n \models \phi) = 0$  or 1. We can also say that  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n \models \phi) = 1 \Leftrightarrow G(\omega) \models \phi$ .

**Proof** Standard considerations similar to the previous lemma. ■

One useful representation for the same results is as follows. Given a first-order property  $\mathcal{A}$  of a

<sup>4</sup>i.e. graph is a pair  $G = (G, E)$  for non-empty set  $G$  of nodes (vertices) and a binary relation  $E$  on  $G$  (the edges). For our purposes, we can assume that  $G$  is symmetric and unordered:  $E(a, b) \rightarrow E(b, a)$ , and  $E(a, a)$  is false. We denote  $G(\omega)$  the class of finite graphs and, loosely, the associated first-order logic model, described in the following section B.1.

594 random graph  $G_n$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}(G_n \in \mathcal{A}) \in \{0, 1\}$ . Equivalent notation will be  $G(n, \omega)$  or just  $G(n)$   
 595 when the context is clear.  
 596  
 597

## 598 C CORE PROOFS, PROOF OF THE MAIN THEOREM 599

600 **Lemma on GPT.** GPT prediction schema is (a.s.) an accumulating error algorithm unless it acts as  
 601 a vacuous learning.  
 602  
 603  
 604

### 605 C.1 FIRST PROOF (FORMAL) 606

607 **Proof** Suppose the GPT prediction is not a vacuous/literal learning. Consider formula  $\phi = \bigwedge_i \phi_i$ .  
 608 where  $\phi_i$  are respected edges on a proof sequence paths,  $e_i \xrightarrow{\phi_i} e_j$ , in any enumeration of the nodes  
 609 in the training dataset. In our first-order theory of graphs, this is a first-order expression. Moreover,  
 610 since the theory obeys 0-1 law, for inference graph  $G$ , by Lemma B.1,  $\lim_{n \rightarrow \infty} G_{n, \phi}$  is zero or one.  
 611 Since GPT algorithm is not vacuous/literal learning, we have  $\lim_{n \rightarrow \infty} G_{n, \phi} = 0$ . That means for  
 612 any  $\epsilon$  there exists  $n_0$  such that  $n \geq n_0$  implies  $\lim_{n \rightarrow \infty} G_{n, \phi_n} < \epsilon$ . Viewed as a graph in  $G(\omega)$  and  
 613 GPT randomization with temperature  $p$  selected for inference, the graph satisfies conditions of 1.  
 614 Thus, starting from  $n_0$ , GPT must be an (a.s.) catastrophic / accumulating error algorithm, with the  
 615 veracity of proof exponentially tending to zero. That is to say, it has to be almost surely a vacuous  
 616 learning to generate valid proof.  
 617  
 618  
 619

### 620 C.2 SECOND PROOF (FORMAL) 621

622 **Main Theorem** For almost all proofs, any learning algorithm of inference, based on randomization  
 623 in  $G(\omega)$ , that necessitates veracity of inference, is almost surely literal learning.  
 624

625 **Proof** More instructive than informal considerations is the following proof in which we partially  
 626 follow a version of the 0-1 law in (Blass et al., 1998).

627 The probability space for the GPT algorithm can be viewed as follows. Consider a probability  
 628 distribution over infinite binary strings. Let  $\Psi$  be a set of infinite sequences representing proofs  
 629 (since any string can be encoded by a binary string, in a suitable enumeration (or embedding) and,  
 630 given a proposition, its proofs of any length can be encoded into an infinite binary string).

631 Let  $\Psi$  be a set of infinite sequences  $\phi = \langle \phi_n : n \geq 1 \rangle \in \Psi$ . In this context, we can view the  
 632 set as one of independent trials. The resulting probability distribution over  $\Psi$  is naturally equipped  
 633 with the product measure (cf. (Feller, 1968)). Moreover, we can consider every proof over strings  
 634 semantically. Therefore, for any generative algorithm  $\mathfrak{A}$ , if, given a sequence  $\{e_0 \rightarrow e_1 \dots e_k \rightarrow$   
 635  $e_t\}$ , representing the proof  $\{e_0 \rightarrow e_t\}$ , we have  $\mathfrak{A}(\phi_n) = e_n$ , we say that the algorithm succeeds  
 636 proving  $\phi_n$ ; otherwise, we say it fails. The corresponding notation for any  $\phi \in \Phi$ , if  $\mathfrak{A}$  succeeds, is  
 637  $\mathfrak{A} \models \phi$ ; if  $\mathfrak{A}$  fails, we write  $\mathfrak{A} \not\models \phi$ .

638 Thus, let us introduce the notation:  $p_n(\mathfrak{A}) = \mathbb{P}(\mathfrak{A} \text{ fails on the } n\text{-th step } \phi_n \text{ of } \phi)$  or  $\mathbb{P}(\mathfrak{A} \not\models \phi_n)$   
 639 where  $\phi$  ranges over  $\Psi$ .

640 The following two cases are possible:

641 **Case 1.** There exists an algorithm,  $\mathfrak{A}$  s.t.  $\sum_{n=0}^{\infty} p_n(\mathfrak{A}) < \infty$ . By the (first) Borel-Cantelli lemma  
 642

643 (Feller, 1968),  $\mathbb{P}(\text{there are infinitely many } n \text{ s.t. } \mathfrak{A} \text{ fails on } \phi_n) = 0$ . Thus, for almost all  $\phi \in \Psi$ ,  $\mathfrak{A}$   
 644 succeeds on all but finitely many  $\phi_n$ . Therefore, for almost all  $\phi$ , there exists an algorithm  $\mathfrak{A}' = \mathfrak{A} +$   
 645 finite lookup that succeeds on  $\phi$ . The algorithm  $\mathfrak{A}$  stays the same for all  $\phi$  and only the finite lookup  
 646 depends on  $\phi$ . It means that, for almost all sequences  $\phi \in \Psi$ ,

$$647 \mathbb{P}(\mathfrak{A} \models \phi) = 1. \quad (9)$$

The question becomes whether such an algorithm  $\mathfrak{A}$  can be GPT. We will show below that the assumption it is GPT meets a contradiction. Namely, from (9) we have:

$$\forall \epsilon > 0 \exists n_0 > 0 \text{ s.t. } \forall n > n_0 \mathbb{P}(\mathfrak{A} \models \phi_n) > 1 - \epsilon. \quad (10)$$

On the other hand, from the Accumulating error lemma inequality (2), we see that  $\mathbb{P}(\mathfrak{A} \not\models \phi) > 1 - \exp(-\rho)$  where  $\rho = \mathbb{P}(\mathbb{E}(\# \text{ faults}))$ . Thus, setting  $\epsilon = 1 - \exp(-\rho)$  leads to contradiction with (10). This leaves only two possibilities for the algorithm  $\mathfrak{A}$  to succeed (since we have  $\mathbb{P}(\mathfrak{A} \models \phi) = 1$  for all  $\phi$ ).

In the first instance,  $\mathfrak{A}$  may arrive at nodes representing the false statements, but the inferences would be true (vacuous truths). The proof is still invalid, overall. The second instance is literal learning; that is, the algorithm would generate (potentially, piece-by-piece) a known proof discoverable in the training data.

**Case 2.** For every algorithm  $\mathfrak{A}$ ,  $\sum_{n=0}^{\infty} p_n(\mathfrak{A}) = \infty$ . Again, as in (2), we can assume that  $\phi_n$  are independent events. By the (second) Borel-Cantelli lemma (e.g., (Feller, 1968)), the probability that there exists an infinite number of  $n$  that  $\mathfrak{A}$  fails on  $\phi_n$  is 1. Hence, for every  $\mathfrak{A}$  there exists  $n$  s.t.  $\mathbb{P}(\mathfrak{A} \models \phi_n) = 0$ . Since there are only countably many algorithms, for almost all  $\phi \in \Phi$ , we have:

$$\mathbb{P}(\exists \mathfrak{A}, \mathfrak{A} \models \phi) = 0. \quad (11)$$

Qualitatively, this means that in this case, almost surely, no algorithm using randomization with exponential correctness decay can succeed in generating a proof for the statement. ■

**Main Theorem, Reformulation** *For almost all proofs, any learning algorithm of inference, based on randomization in  $G(\omega)$ , does not generate a valid proof unless it is vacuous.* ■

## REFERENCES

- Nicholas Asher, Swarnadeep Bhar, and Akshay Chaturvedi et. al. Autocorrelations decay in texts and applicability limits of language models. *arxiv*, 2306(12213):1–13, 2023a. URL <https://arxiv.org/pdf/2306.12213>.
- Nicholas Asher, Swarnadeep Bhar, and Akshay Chaturvedi et. al. Limits for learning with language models. *arxiv*, 2306(12213):1–13, 2023b. URL <https://arxiv.org/pdf/2306.12213>.
- Andreas Blass, Yuri Gurevich, Vladik Kreinovich, and Luc Longpré. A variation on the zero-one law. In *Information Processing Letters*, volume 67, pp. 29–30, January 1998. URL <https://www.microsoft.com/en-us/research/publication/132-variation-zero-one-law/>.
- Lingjiao Chen, Matei Zaharia, and James Zou. How is chatgpt’s behavior changing over time? *arxiv*, 2307(09009):1–23, 2023.
- Shibhansh Dohare, J. Fernando Hernandez-Garcia, and Qingfeng Lan et. al. Loss of plasticity in deep continual learning. In *Nature*, volume 632, pp. 768–774, August 2024. URL <https://doi.org/10.1038/s41586-024-07711-7>.
- Nouha Dziri and Ximing Lu et. al. Faith and fate: Limits of transformers on compositionality. *arxiv*, 2305(18654):1–37, 2023.
- Nouha Dziri, Ximing Lu, and Melanie Sclar et. al. Faith and fate: Limits of transformers on compositionality. *arxiv*, 2305(18654):1–40, 2023.
- Ronald Fagin. Probabilities on finite models. *The Journal of Symbolic Logic*, 41(1):50–57, 1976.
- William Feller. *An Introduction to Probability Theory and Its Applications*. Wiley, Princeton, NJ, 1968.
- Gaël Gendron, Qiming Bao, and Michael Witbrock et. al. Large language models are not strong abstract reasoners. *arxiv*, 2305(19555):1–50, 2023. URL <https://arxiv.org/pdf/2305.19555>.

- 702 Eric Glazer. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai. *arxiv*,  
703 2411(04872):1–26, 2024. URL <https://arxiv.org/pdf/2411.04872>.
- 704
- 705 Martin Grohe. *Descriptive Complexity, Canonisation, and Definable Graph Structure Theory*. As-  
706 sociation for Symbolic Logic, Storrs, CT, 2017.
- 707 Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Thinking fast and slow in large language  
708 models. *arxiv*, 2212(0900905206):1–30, 2022.
- 709
- 710 Evan Hubinger, Carson Denison, and Jesse Mu et. al. Sleeper agents: Training deceptive llms that  
711 persist through safety training. *arxiv*, 2401(05566):1–71, 2024.
- 712 Charles Jin and Martin Rinard. Emergent representations of program semantics in language models  
713 trained on programs. *arxiv*, 2305(11169):1–25, 2024.
- 714
- 715 Samia Kabir, David N. Udo-Imeh, Bonan Kou, and Tianyi Zhang. Is stack overflow obsolete? an  
716 empirical study of the characteristics of chatgpt answers to stack overflow questions. *CHI’24*,  
717 2024. URL <https://dl.acm.org/doi/pdf/10.1145/3613904.3642596>.
- 718 Ryan Krueger, Jesse Michael Han, and Daniel Selsam. Automatically building diagrams for  
719 olympiad geometry problems. *arxiv*, 2012(02590):1–22, 2021. URL <https://arxiv.org/pdf/2012.02590>.
- 720
- 721
- 722 Chen L., Zaharia M, and Zou J. Frugalgpt: How to use large language models while reducing cost  
723 and improving performance. *arxiv*, 2305(05176):1–13, 2023.
- 724 Wentao Liu, Hanglei Hu, and Jie Zhou et. al. Mathematical language models: A survey. *arxiv*, 2312  
725 (07622):1–21, 2023.
- 726
- 727 O. Macmillan-Scott and M. Musolesi. (ir)rationality and cognitive biases in large language mod-  
728 els. *R. Soc. Open Sci.*, 11(6):240255, 2024. URL [https://doi.org/10.1098/rsos.](https://doi.org/10.1098/rsos.240255)  
729 240255.
- 730 Eran Malach. Auto-regressive next-token predictors are universal learners. *arxiv*, 2309(06979):  
731 1–22, 2023.
- 732
- 733 Nat McAleese and Rai (Michael Pokorny) Juan Felipe Cerón Uribe et. al. Llmcritics  
734 help catch llm bugs. *Open AI*, pp. 1–23, 2024. URL [https://cdn.openai.com/](https://cdn.openai.com/llm-critics-help-catch-llm-bugs-paper.pdf)  
735 [llm-critics-help-catch-llm-bugs-paper.pdf](https://cdn.openai.com/llm-critics-help-catch-llm-bugs-paper.pdf).
- 736 Nikolay Mikhaylovskiy and Ilya Churilov. Autocorrelations decay in texts and applicability limits of  
737 language models. *arxiv*, 2305(06615):1–20, 2023. URL [https://arxiv.org/pdf/2305.](https://arxiv.org/pdf/2305.06615)  
738 06615.
- 739
- 740 Milad Nasr and Nicholas Carlini and Jonathan Hayase et. al. Scalable extraction of training data  
741 from (production) language models. *arxiv*, 2311(17035):1–64, 2023.
- 742 Marianna Nezhurina, Lucia Cipolina-Kun1, Mehdi Cherti, and Jenia Jitsev. Alice in wonderland:  
743 Simple tasks showing complete reasoning breakdown in state-of-the-art large language models.  
744 *arxiv*, 2406(02061):1–45, 2024.
- 745
- 746 Nhan Nguen and Nadi Sarah. An empirical evaluation of github copilot’s code suggestions. In  
747 *MSR ’22: Proceedings of the 19th International Conference on Mining Software Repositories*,  
748 volume 19, pp. 1–5, May 2022. URL <https://doi.org/10.1145/3524842.3528470>.
- 749 Bernardino Romera-Paredes, Mohammadamin Barekatin, and Alexander Novikov et. al. Mathe-  
750 matical discoveries from program search with large language models. *Nature*, pp. 1–33, 2023.
- 751 Abulhair Saparov, Richard Yuanzhe Pang, and Vishakh Padmakumar et. al. Testing the general  
752 deductive reasoning capacity of large language models using ood examples. *arxiv*, 2305(15269):  
753 1–23, 2023.
- 754
- 755 Bilgehan Sel, Ahmad Al-Tawaha, and Vanshaj Khattar et. al. Algorithm of thoughts: Enhancing  
exploration of ideas in large language models. *arxiv*, 2308(10379):1–29, 2023.

756 Yundi Shi, Piji Li, Changchun Yin, and Zhaoyang Han et. al. Promptattack: Prompt-based attack  
757 for language models via gradient search. *arxiv*, 2209(01882):1–12, 2022.  
758

759 Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, and He He et.al. Solving olympiad geometry without  
760 human demonstrations. *nature*, 625(1):476–482, 2024. URL [https://www.nature.com/  
761 articles/s41586-023-06747-5](https://www.nature.com/articles/s41586-023-06747-5).

762 Yotam Wolf and Noam Wies et. al. Fundamental limitations of alignment in large language models.  
763 *arxiv*, 2304(11082):1–29, 2023.  
764

765 Steve Yadlowsky, Lyric Doshi, and Nilesh Tripuraneni. Pretraining data mixtures enable narrow  
766 model selection capabilities in transformer models. *arxiv*, 2311(00871):1–13, 2023a.

767 Steve Yadlowsky, Lyric Doshi, and Nilesh Tripuraneni. Pretraining data mixtures enable nar-  
768 row model selection capabilities in transformer models. *arxiv*, 2311(0087):1–13, 2023b. URL  
769 <https://arxiv.org/pdf/2311.00871>.  
770

771 Zhang Z., Zhang A., and Li M. et. al. Automatic chain of thought prompting in large language  
772 models. *arxiv*, 2210(03493):1–25, 2022.

773 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: A cross-system benchmark for  
774 formal olympiad-level mathematics. *arxiv*, ICLR 2022(2109):1–11, 2022. URL [https://  
775 arxiv.org/pdf/2109.00110](https://arxiv.org/pdf/2109.00110).  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809