

---

# THE PITFALLS OF NEXT-TOKEN PREDICTION

---

Gregor Bachmann<sup>\*1</sup> Vaishnavh Nagarajan<sup>\*2</sup>

## Abstract

Can a mere next-token predictor faithfully model human intelligence? We crystallize this intuitive concern, which is fragmented in the literature. As a starting point, we argue that the two often-conflated phases of next-token prediction — autoregressive inference and teacher-forced training — must be treated distinctly. The popular criticism that errors can compound during autoregressive inference, crucially assumes that teacher-forcing has learned an accurate next-token predictor. This assumption sidesteps a more deep-rooted problem we expose: in certain classes of tasks, teacher-forcing can simply fail to learn an accurate next-token predictor in the first place. We describe a general mechanism of how teacher-forcing can fail, and design a minimal planning task where both the Transformer and the Mamba architecture empirically fail in that manner — remarkably, despite the task being straightforward to learn. We provide preliminary evidence that this failure can be resolved when training to predict multiple tokens in advance. We hope this finding can ground future debates and inspire explorations beyond the next-token prediction paradigm. We make our code available under <https://github.com/gregorbachmann/Next-Token-Failures>

## 1. INTRODUCTION

Long after its inception in the seminal work of Shannon (1948; 1951), next-token prediction has made its way into becoming a core part of the modern language model. But despite its long list of achievements, there is a small but growing belief that a next-token predicting model is merely an impressive *improv* artist that cannot truly model human thought. Humans, when navigating the world, meticulously imagine, curate and backtrack plans in their heads before

executing them. Such strategies are unfortunately not explicitly built into the backbone of the present-day language model. This criticism has been floating around as an informal viewpoint (LeCun, 2024; Bubeck et al., 2023). Our paper is aimed at crystallizing this intuitive criticism of next-token prediction, and developing the core arguments of this debate.

Let us start by making more precise, what it means to say that human-generated language, or problem-solving, does not follow next-token prediction. When formalizing this, we hit an immediate roadblock: isn't every sequence generation task possible autoregressively? Put differently, an optimist would say, every distribution over a sequence of tokens can be captured by an appropriately sophisticated next-token predictor simulating the chain rule of probability i.e.,  $\mathbb{P}(r_1, r_2, \dots) = \prod_i \mathbb{P}(r_i | r_1 \dots r_{i-1})$ . Thus, the autoregressivity in our systems is not antithetical to learning human language, after all.

Although this argument is compelling, a pessimist would worry, realistically, even with minor imperfections in the next-token predictor, the accuracy may break down spectacularly for long sequences (Kääriäinen, 2006; Ross & Bagnell, 2010; LeCun, 2024; Dziri et al., 2023). Say, even if every next-token error is as little as 0.01, the probability of encountering an erroneous token exponentially compounds along the way, and by the end of 200 tokens, blows up to 0.86.

This is a simple and powerful observation. Yet, this does not completely capture the intuition that next-token predictors may be poor planners. Crucially, this argument does not carefully distinguish between the two types of next-token prediction: inference-time autoregression (where the model consumes its own previous outputs as inputs), and training-time *teacher-forcing* (Williams & Zipser, 1989) (where the model is taught to predict token-by-token consuming all previous ground truth tokens as inputs). Framed this way, the compounding of errors only pinpoints a superficial failure to execute a plan during *inference*. It leaves open the possibility that we may have still learned a near-perfect next-token predictor; perhaps, with an appropriate post-hoc wrapper that verifies and backtracks, we can elicit the right plan without compounding errors.

Drawing this distinction allows us to articulate a much more

---

<sup>\*</sup>Equal contribution <sup>1</sup>ETH Zürich, Switzerland <sup>2</sup>Google Research, US. Correspondence to: Gregor Bachmann <gregorb@ethz.ch>, Vaishnavh Nagarajan <vaishnavh@google.com>.

concerning possibility: is it safe to assume that next-token based learning (teacher-forcing) always learns an accurate next-token predictor? We identify this is not always the case. Consider a task where we expect the model to witness a problem statement  $\mathbf{p} = (p_1, p_2, \dots)$  and produce the ground truth response tokens  $(r_1, r_2, \dots)$ . Teacher-forcing trains the model to produce each token  $r_i$  by not only providing the problem statement  $\mathbf{p}$  but also by revealing part of the ground truth  $r_1, \dots, r_{i-1}$ . Depending on the task, we first argue that this can induce shortcuts that use the revealed prefix of the ground truth answer to spuriously fit future answer tokens. We call this the *Clever Hans cheat*.<sup>1</sup> Next, while the later tokens ( $r_i$  for large  $i$ ) become easy to fit by the Clever Hans cheat, in contrast, the earlier answer tokens (say,  $r_0, r_1$  etc.) become harder to learn. This is because they no longer come with any supervision about the full answer — part of the supervision is lost to the Clever Hans cheat. We argue that these two flaws would together arise in “lookahead tasks”: tasks that require implicitly planning a later token in advance of an earlier token. In such tasks, teacher-forcing would result in a highly inaccurate next-token predictor that would fail to generalize to unseen problems  $\mathbf{p}$ , even those sampled in-distribution.

Empirically, we demonstrate that the above mechanism leads to complete in-distribution failure in a path-finding setup on a graph, that we propose as a minimal lookahead task. We design our setup in a way that it is demonstrably straightforward to solve that the failure of any model is remarkable. Yet, we observe failure for both the Transformer (Vaswani et al., 2017) and the Mamba architecture, a structured state space model (Gu & Dao, 2023). We also find that a form of *teacherless* training that predicts multiple future tokens (Monea et al., 2023) is (in some settings) able to circumvent this failure. Thus, we pinpoint a precise and easy-to-learn scenario where, rather than properties that are criticized in existing literature — like convolution or recurrence or autoregressive inference (see 6), — it is next-token prediction during training that is at fault.

We hope that these findings inspire and set future debates around next-token prediction on solid ground. In particular, we believe that the failure of the next-token prediction objective on our straightforward task casts a shadow over its promise on more complex tasks (such as say, learning to write stories). We also hope that this minimal example of failure and the positive results on teacherless training can motivate alternative paradigms of training.

<sup>1</sup>*Clever Hans* (Pfungst & Rahn, 1911) was a famous show horse that could solve simple arithmetic tasks by repeatedly tapping with his hoof until he reached the correct count. It turns out however, that Clever Hans did not really solve the problem, but merely stopped tapping upon detecting certain (involuntary) facial cues from his coach. Clever Hans’ answers were wrong when the coach was absent.

We summarize our contributions below.

1. We consolidate existing critiques against next-token prediction and crystallize new core points of contention (§6 and §3, §4).
2. We identify that the next-token prediction debate must not conflate autoregressive inference with teacher-forcing. Both lead to vastly different failures (§3, §A).
3. We conceptually argue that in lookahead tasks, next-token prediction during training (i.e., teacher-forcing) can give rise to problematic learning mechanisms that are detrimental to even in-distribution performance (§4).
4. We design a minimal lookahead task (§4.1). We empirically demonstrate the failure of teacher-forcing for the Transformer and Mamba architectures, despite the task being easy to learn (§5).
5. We identify that a teacherless form of training that predicts multiple future tokens at once — proposed in Monea et al. (2023) for orthogonal inference-time efficiency goals — shows promise in circumventing these training-time failures in some settings (§5, Eq 7). This further demonstrates the limits of next-token prediction.

## 2. THE TWO MODES OF NEXT-TOKEN PREDICTION

Consider a set of tokens  $\mathcal{V}$ . Let  $\mathcal{D}$  be a ground truth distribution over sequences that consist of a prefix  $\mathbf{p}$  and a response  $\mathbf{r}$ , denoted as  $\mathbf{s} = \mathbf{p}, \mathbf{r}$  where  $\mathbf{p} = (p_1, p_2, \dots) \in \mathcal{V}^{L_{\text{pref}}}$  and  $\mathbf{r} = (r_1, r_2, \dots) \in \mathcal{V}^{L_{\text{resp}}}$ . We assume sequences of fixed length merely for simplicity.

For any sequence  $\mathbf{s}$ , let  $\mathbf{s}_{<i}$  denote the first  $i - 1$  tokens of  $\mathbf{s}$ , and  $\mathbf{s}_{i<}$  the tokens following the  $i$ th token. Note that  $\mathbf{s}_{<1}$  is the empty prefix. With an abuse of notation, let  $\mathbb{P}_{\mathcal{D}}(s_i | \mathbf{s}_{<i})$  denote the ground truth probability mass on  $s_i$  being the  $i$ th token given the prefix  $\mathbf{s}_{<i}$ . Consider a next-token-predicting language model  $\text{LM}_{\theta}$  (with parameters  $\theta$ ) such that  $\text{LM}_{\theta}(\hat{s}_i = s_i; \mathbf{s}_{<i})$  is the probability that the model assigns to the  $i$ th output  $\hat{s}_i$  taking the value  $s_i$ , given as input the sequence  $\mathbf{s}_{<i}$ . Note that the next-token predictor only defines the probability for a single future token given an input, but not the joint probability of multiple future tokens. This joint probability is axiomatically defined analogous to the chain rule of probability:

$$\text{LM}_{\theta}(\hat{\mathbf{r}} = \mathbf{r} ; \mathbf{p}) := \prod_{i=1}^{L_{\text{resp}}} \text{LM}_{\theta}(\hat{r}_i = r_i; \mathbf{p}, \mathbf{r}_{<i}) \quad (1)$$

where  $\hat{\mathbf{r}} = \mathbf{r}$  denotes an exact token-by-token match.

To train the above model, two distinct types of next-token prediction are used. First, during inference, for a given

prefix, we autoregressively sample from the model token-by-token, providing as input the prefix and all previously-generated tokens. Formally,

**Definition 1. (Inference-time next-token prediction via autoregression)** Autoregressive inference is a form of inference-time next-token prediction in that to generate a response  $\hat{\mathbf{r}}$ , we iterate over  $i = 1, \dots, L_{resp}$ , to sample the next token  $\hat{r}_i$  with the distribution given by  $LM_\theta(\hat{r}_i; \mathbf{p}, \mathbf{r}_{<i})$ . We denote this as  $\hat{\mathbf{r}} \stackrel{\text{ag}}{\sim} LM_\theta(\cdot; \mathbf{p})$ .

There is also a second phase of next-token prediction, one that is applied during the training process, called *teacher-forcing*. Here, instead of feeding the model its own output back as input, the model is fed with prefixes of the *ground truth* response  $\mathbf{r}_{<i}$ . Meanwhile, the model is assigned as supervisory target,  $r_i$ , the next ground truth token. Then, the model maximizes a sum of next-token log-probabilities:

**Definition 2. (Training-time next-token prediction via teacher-forcing)** Teacher-forced training is a form of training-time next-token prediction in that we find parameters  $\theta$  that maximize the next-token log-probability sum:

$$\begin{aligned} \mathcal{J}_{next-token}(\theta) &= \mathbb{E}_{(\mathbf{p}, \mathbf{r}) \sim \mathcal{D}} [\log LM_\theta(\hat{\mathbf{r}} = \mathbf{r}; \mathbf{p})] \\ &= \mathbb{E}_{\mathcal{D}} \left[ \sum_{i=1}^{L_{resp}} \log LM_\theta(\hat{r}_i = r_i; \mathbf{p}, \mathbf{r}_{<i}) \right] \quad (2) \end{aligned}$$

The key property of the objective is that we extract the model’s output, *allowing the model access to the ground truth response preceding the current token*. This property will be crucial to the failure we describe in §4.

### 3. FAILURE DUE TO AUTO-REGRESSIVE INFERENCE

A broad criticism against next-token predictors is that intuitively these models are not explicitly designed to plan ahead, and during inference, they do not know how to recover from their own errors. This discourse has been fragmented in literature. Furthermore, the umbrella term “next-token prediction” is used interchangeably with “autoregressive architecture”. Our goal is to analyze these intuitions more systematically, and be careful about distinguishing between the two phases of next-token prediction: teacher-forcing and autoregression. A key insight we will arrive at is that existing arguments capture only a part of the intuitive concern that next-token predictors may not be able to plan.

**The chain-rule-of-probability defense:** We first outline what is arguably the most tempting defense for next-token prediction: the chain rule of probability always promises us a next-token predictor that can fit our distribution.

**Fact 1. (Every sequence distribution can be represented by a next-token predictor)** By the chain rule of probability

we have  $\mathbb{P}_{\mathcal{D}}(\mathbf{r} \mid \mathbf{p}) = \prod_{i=1}^{L_{resp}} \mathbb{P}_{\mathcal{D}}(r_i \mid \mathbf{p}, \mathbf{r}_{<i})$ . Therefore, define a next-token predictor  $LM$  such that for every valid value of  $i$ ,  $\mathbf{p}$ , and  $\mathbf{r}$ , we have  $LM(\hat{r}_i = r_i; \mathbf{p}, \mathbf{r}_{<i}) := \mathbb{P}_{\mathcal{D}}(r_i \mid \mathbf{p}, \mathbf{r}_{<i})$ . Then, sampling  $\mathbf{r} \sim \mathcal{D} \mid \mathbf{p}$ , is equivalent to autoregressively sampling  $\mathbf{r} \stackrel{\text{ag}}{\sim} LM(\cdot; \mathbf{p})$ .

The cleverness of this argument lies in the fact that it can apply to *any* imaginable distribution. Thus, as long as the next-token predictor is sufficiently expressive (by scaling up the context, memory and compute), it can model both natural language and problem-solving. Thus, it may seem that next-token predictors are not antithetical to planning-based tasks, after all.

**The snowballing errors criticism:** A skeptic would however raise the following concern. Regardless of the abundance of computational resources, realistic models are typically not perfect next-token predictors. There may always be a slight chance of error in each step, and once an error is committed, there is no explicit backtracking mechanism to rescue the model. The argument then goes that, the probability of errors in each token, however miniscule, would exponentially snowball along the way. By the end of a long sequence of tokens, the accuracy (of having produced an error-free response) becomes trivial. This has been formalized in various contexts, from that of autoregressive models LeCun (2024), to that of the limits of Transformers in compositional tasks Dziri et al. (2023), and in a different form, in much earlier work in imitation learning and structured prediction Kääriäinen (2006); Ross & Bagnell (2010) (see §6). We present a minimal formalization of this below:

**Failure 1. (Snowballing error due to autoregressive inference)** Consider a model  $LM_\theta$ , prefix  $\mathbf{p}$  and a unique ground truth response  $\mathbf{r}$  such that the next-token error obeys

$$\forall i \leq L_{resp}, LM_\theta(\hat{r}_i \neq r_i; \mathbf{p}, \mathbf{r}_{<i}) \approx \epsilon. \quad (3)$$

Then, for  $\hat{\mathbf{r}} \stackrel{\text{ag}}{\sim} LM_\theta(\cdot; \mathbf{p})$  the probability that the generated response exactly matches the ground truth  $\mathbf{r}$  obeys

$$\mathbb{P}(\hat{\mathbf{r}} = \mathbf{r}) \approx (1 - \epsilon)^{L_{resp}}.$$

We argue that the snowball failure mode only indicates how an autoregressive model can fail to *execute* a plan during inference-time. It does not preclude the possibility that the model may have *learned* a good plan that it simply fails to execute during inference. Concretely, it may still be possible that, at each step, the model has high accuracy of predicting a next token that is consistent with a good plan (as assumed in Eq 3). Depending on the setting, one can potentially exploit this accuracy to elicit a good plan during inference. For instance, one may be able to use a post-hoc wrapper that first verifies whether an error has taken place, then

backtracks and re-executes a different action. One may even simulate backtracking using more elaborate techniques such as chain/tree/graph of thought (Wei et al., 2022; Yao et al., 2023a; Besta et al., 2023; Yao et al., 2023b), or using the model to give itself feedback (Madaan et al., 2023; Huang et al., 2022; Shinn et al., 2023) to elicit the plan that the model has learned.

Thus, the snowball failure mode captures what is primarily a shortcoming of an autoregressive architecture. Likewise, the chain-rule-of-probability defense captures only the expressive power of an autoregressive architecture. Neither of these arguments address the possibility that learning with next-token prediction may itself have shortcomings in learning how to plan. In this sense, we argue that existing arguments capture only a part of the intuitive concern that next-token predictors fare poorly at planning.

#### 4. FAILURE DUE TO TEACHER-FORCING

Can a model trained to predict the next token, fail to predict the next token with high accuracy during test-time? Mathematically, this would mean showing that a model trained with the teacher-forcing objective of Eq 2 has high next-token prediction error on the very distribution it was trained on (thus breaking the assumption in Eq 3 of the snowballing failure mode). Consequently, no post-hoc wrapper can salvage a plan out of the model. The goal of this section is to conceptually argue that this failure can happen for lookahead tasks: tasks that implicitly require computing a future token in advance before an earlier token.

As a running example for our argument, we design a path-finding problem on a simple class of graphs. We view this example as a minimal setting that captures the core essence of what it means to solve problems with lookahead, without irrelevant confounding factors. This task is also demonstrably straightforward to solve, as we will see, thus making any observed failures remarkable. Thus we view this running example as a template for an intuitive argument that can be made about teacher-forced models on more general and harder problems that require lookahead.

##### 4.1. Path-Finding on Path-Star Graphs: A Minimal and Easy Lookahead Task

Consider a path-finding problem on a directed graph  $G$  with a set of nodes  $\{v_{\text{start}}, v_{\text{goal}}, v_1, v_2, \dots\}$ . The graph is a “path-star” graph with  $v_{\text{start}}$  as the central node, with at least 2 paths (each of length  $l \geq 2$  edges) emanating from it, with a unique path ending in  $v_{\text{goal}}$ . The task is to find a path from  $v_{\text{start}}$  to  $v_{\text{goal}}$ . Correspondingly, we assume that the distribution  $\mathcal{D}$  is over sequences where the prefix  $p$  represents a (randomly generated) graph, and the response represents the path from the start to the goal. In particular,

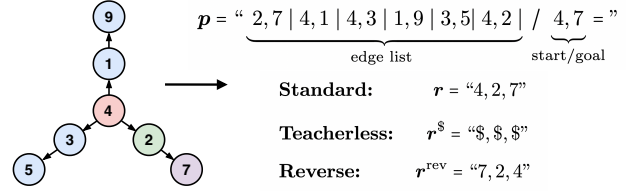


Figure 1. Illustration of a path-star graph. The prefix  $p$  represents the adjacency list and the (central) start and goal node. The target is represented by  $r$ . Under “standard” teacher-forcing, we condition the model on prefixes of  $r$  to predict  $r$ . But in §5 we explore alternatives where we train without a teacher (condition on  $r^S$  and predict  $r$ ) or train with a reversal (condition on and predict  $r^{\text{rev}}$ ).

we sample a graph  $G$  which is represented as an adjacency list as  $\text{adj}(G) = e_1, e_2, \dots$  where each edge  $e = (v, v')$  is represented such that  $v'$  farther away from  $v_{\text{start}}$  than  $v$ . We then set the prefix as  $p = (\text{adj}(G), v_{\text{start}}, v_{\text{goal}})$  so the model knows what the graph, and the desired start and goal states are. The ground truth response  $r$  corresponds to the sequence of vertices  $r = v_{\text{start}}, \dots, v_{\text{goal}}$  on the start-to-goal path. We visualize this construction in Fig. 1.

**The straightforward lookahead solution.** Ideally, we want the model to learn a mapping from the input  $p$  consisting *only* of  $(\text{adj}(G), v_{\text{start}}, v_{\text{goal}})$  to an output that is the full path  $r$ . Two such solutions are possible. One idea is to plan by examining all the paths emanating from  $v_{\text{start}}$  and choosing the one that ends at  $v_{\text{goal}}$ . But a second, straightforward solution exists: the model simply needs to look ahead at the sequence “right-to-left” and observe that it corresponds to the one unique path starting from  $v_{\text{goal}}$  and ending at  $v_{\text{start}}$ . After internally computing the path from  $v_{\text{goal}}$  and reversing it, the model can emit its response.

##### 4.2. Outline of failure mechanism

While we will use the path-star example as a running example, we make our claim more generally for problems that require lookahead (such as story-writing, as we will discuss later). In such classes of problems, we claim that teacher-forcing prevents learning the true mechanisms, causing failure. Intuitively, in teacher-forcing, we decompose the learning of  $p \rightarrow r$  into multiple problems, one for each token  $r_i$ . Specifically, we make the model learn a mapping from the input  $(p, r_{<i})$  — not just  $p$  — to the output  $r_i$ . The additional information  $r_{<i}$  in the input, we argue, is problematic and destroys the core challenge in what the model has to learn. Specifically, our argument puts forth two debilitating mechanisms that would together emerge under teacher-forcing (explained over the next two subsections). While, we will empirically verify these mechanisms for path-star graphs in §5, we also provide a discussion of how our ideas apply to a text-based scenario at the end of



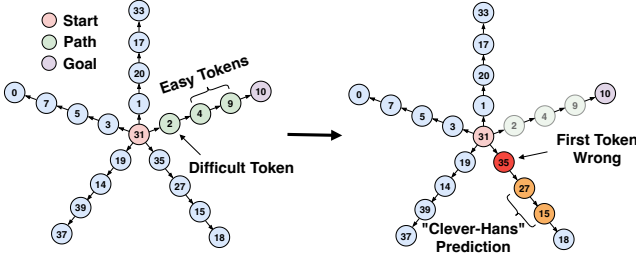


Figure 2. Illustration of the failure of teacher-forcing on a path-star graph. The left image marks the “easy tokens” which can be fit by the Clever Hans cheat (Failure 2a), while the “difficult token” cannot be learned (Failure 2b) due to lost supervision. The right image shows how the model would behave during autoregressive inference, under the absence of the “teacher”.

this section.

### 4.3. The Clever Hans Cheat

First, and most importantly, by revealing parts of the answer to the model as input, we allow the model to fit the data by *cheating* i.e., by using trivial mechanisms that use the extra information in  $\mathbf{r}_{<i}$  to produce  $r_i$ . Such cheats must especially be abundant for the later tokens (large  $i$ ) for which a larger prefix is revealed.

To illustrate this in our path-star example, without loss of generality, consider a ground truth path that is of the form  $\mathbf{r} = v_{\text{start}}, v_1, v_2, \dots, v_{\text{goal}}$ . With a slight abuse of the indexing notation, let  $\mathbf{r}_{<i} = v_{\text{start}}, v_1, \dots, v_{i-1}$  be the prefix of length  $i$  (so we index from 0 instead of 1). Observe that nodes from  $v_2$  onwards, until before  $v_{\text{goal}}$ , have precisely one edge going “away” from  $v_{\text{start}}$ . Thus, consider when the model is given as input,  $(\mathbf{p}, \mathbf{r}_{<i})$  where  $\mathbf{p} = (\text{adj}(G), v_{\text{start}}, v_{\text{goal}})$ , to fit the target  $v_i$ . The model first merely needs to scan the adjacency list  $\text{adj}(G)$  within  $\mathbf{p}$  for the one edge containing  $v_{i-1}$  in the first position. Then, the model only has to predict the other node on that edge as  $v_i$ . Note though, this cheat cannot work on fitting the target  $v_1$  given the input  $\mathbf{r}_{<1} = v_{\text{start}}$  since  $v_{\text{start}}$  has many outward edges — we will address this node in the next section. We illustrate this difference between  $\mathbf{r}_1$  and the remaining tokens as the “easy” vs. “difficult” tokens in Fig. 2.

Crucially, the above cheating mechanism for fitting the easy tokens does not require any lookahead. It is simple, and implementable by an induction head-like module (Olsson et al., 2022). Owing to this simplicity, we hypothesize that these tokens will be quickly fit and ignored during training. This destroys useful signal for the model to efficiently learn the underlying “right-to-left” solution: the solution that requires looking at all tokens in  $\mathbf{r}$ , and then learning that they are simply the unique path from  $v_{\text{goal}}$  spelled in reverse.

We emphasize two key aspects of this cheating behavior.

First, these shortcuts are unlike previously-identified shortcuts (see §6) that map from the original input prefixes  $\mathbf{p}$  to the ground truth  $\mathbf{r}$ . The behavior we identify is unique to the mapping from the teacher-forced prefix  $\mathbf{p}, \mathbf{r}_{<i}$  to  $r_i$ . We name this behavior as *Clever Hans cheating*. Another notable point is that this does not come from a dearth of samples: even if we had infinite training data at our disposal, the model can still fit the easy tokens of all that data by Clever Hans cheating.

### 4.4. The Indecipherable Token

Perhaps, not all is lost. While the later tokens may be fit using the Clever Hans cheat, we may still have some of the earlier tokens (for small  $i$ ), for which such cheats may be unavailable. The supervision from these tokens may eventually coerce the model into learning the true solution. For example, in the path-star task, the model still needs to learn to predict the first node  $v_1$ , where it is not possible to fit the training data by the Clever Hans cheat. If not memorize this token on the training data, the most general way to fit this token is by actually solving the underlying task.

However, we argue that it is significantly harder for the model to learn the correct solution now. Consider the point in training when the Clever Hans cheat is perfected. At this point, the model is deprived of information about much of the full solution which was once present as supervisory targets. The model is simply left with the task of mapping the input  $\mathbf{p}$  to an *incomplete* solution (e.g., the first vertex  $v_1$  in the path-star graph). Recovering the plan in this scenario must first of all be relatively harder from a statistical point of view due to the incomplete supervision. But more importantly, learning this task may become computationally harder, or simply, intractable (Wies et al., 2023).

We provide an informal intuition of intractability for the path-star problem, but this intuition should extend to more general problems as well. Intuitively, our learner has to find an end-to-end algorithm that composes multiple subroutines over one another. For instance, recall that the straightforward solution consists of  $l$  steps: start from the current vertex as  $v_{\text{goal}}$ , and find the preceding vertex in the graph in each subsequent step. Each vertex in this path can be thought of as “intermediate supervision” to learn a corresponding “find-the-adjacent-vertex” subroutine from a space of candidate subroutines.<sup>2</sup> Even if we conservatively assume that there is only a constant-sized space of  $C$  candidate subroutines,

<sup>2</sup>As an illustration of what these candidate subroutines could be, imagine that the model can implement an induction head (Olsson et al., 2022)  $\text{Ind}_k(\mathbf{p}, v)$  that finds  $v$  in the adjacency list of  $\mathbf{p}$ , and outputs the token that precedes it by  $k$  positions. Then the candidate space could be parameterized by  $k$  as  $\{\text{Ind}_k(\mathbf{p}, v) | k = 1, 2, \dots\}$ . For our specific tokenization, the correct subroutine at each of the  $l$  steps is the induction head for which  $k = 2$ .

the end-to-end search space is an exponential space of  $C^l$  algorithms composing  $l$  subroutines.

Now, after the Clever Hans cheat is in effect, the only supervision for this search is a single-token loss of the form  $-\log \text{LM}_\theta(\hat{r}_1 = r_1; \mathbf{p})$ . However, this loss is an “all-or-nothing” loss. Crucially, by the *discrete* nature of the task, even if the learner gets one subroutine incorrect, the final answer  $\hat{r}_1$  would likely be incorrect on all inputs. For instance, imagine that the first subroutine is incorrectly learned and its output takes us to an arbitrary location on the graph. Then, even if all subsequent subroutines were correctly learned (i.e., they are “find-the-adjacent-vertex” subroutines), the final output would be arbitrary. Thus, we have  $\hat{r}_1 = r_1$  precisely for the algorithm where all  $l$  subroutines are correct, and  $\hat{r}_1 \neq r_1$  for any other choice of the algorithm. For such an all-or-nothing loss surface, the end-to-end learner must necessarily brute-force search through the exponential space of algorithms.<sup>3</sup>

We encapsulate the overall claim more generally via the following conditions:

**Proposition 3.** *Consider learning a task that satisfies the following conditions:*

1. *It requires composing  $l$  discrete-output subroutines over one another.*
2. *The  $k$  leading response tokens are sensitive in that even if one subroutine is altered, the first  $k$  tokens are each completely altered.*
3. *The subroutine search space consists of at least  $C$  candidate subroutines.*

*Then learning the task with only supervision from the first  $k$  ground truth tokens requires exponential time of  $\Omega(C^l)$ .*

Indeed, literature on chain-of-thought has identified many instances of similar “multi-hop reasoning tasks” that cannot be solved end-to-end — both empirically (Wei et al., 2022) and theoretically (Wies et al., 2023) — unless there is “complete supervision”. We elaborate on this in §6.

In § 5, we will verify experiments to demonstrate that our models indeed fail to learn the Indecipherable Token as a result of Clever Hans cheating, and that conversely, they succeed in conditions where the Clever Hans cheat is prevented.

<sup>3</sup>The only condition under which the learner may identify the right algorithm efficiently is when the model has accidentally seen similar problems during pretraining and thus has a useful prior. For example, the prior may be that it assigns high probability to the correct subroutine(s). Or, in the specific path-star example, the model may assign high probability to all  $l$  subroutines being identical. In such a case, one can construct slight variations of the tasks that defy this prior, to demonstrate intractability.

## 4.5. Beyond the path-star setting

Framing our argument more generally, and informally, we argue that teacher-forcing can suffer the following failures in order, especially in tasks that require advance lookahead.

**Failure 2a. (Clever Hans cheating due to teacher-forcing)** *Although there is a true mechanism that can recover each  $r_i$  from the original prefix  $\mathbf{p}$ , there can be multiple other mechanisms that can recover a token  $r_i$  from the teacher-forced prefix  $(\mathbf{p}, \mathbf{r}_{<i})$ . These mechanisms can be simpler to learn thus disincentivizing the model from learning the true mechanism.*

**Failure 2b. (Indecipherable token due to lost supervision)** *After the Clever Hans cheat is perfected during training, the model is deprived of a part of the supervision (especially,  $r_i$  for larger  $i$ ). This makes it harder and potentially even intractable for the model to learn the true mechanism from the remaining tokens alone.*

As we demonstrate in the next section, the above failures can cause the model to fail on the very distribution it was trained on. This breakdown of planning abilities emerges right from training, and is orthogonal to the Snowballing Failure that is primarily an inference-time issue (See §A).

While the path-star problem provides a concrete and empirically verifiable setting of this failure (where it is easy to reason and test the mechanisms), it can help us speculate how such failures could occur in more complex and nebulous tasks. Intuitively we expect this failure to occur when there are right-to-left dependencies: a token that appears later must be planned before an earlier-appearing token. We provide an example below.

**Story-writing.** Consider training with teacher-forcing on novels. Imagine that the (sub)plots take the form of a conflict, followed by a backstory, followed by a resolution of the conflict, utilizing the backstory (illustrated in detail in §B). Crucially, although the story explicitly reads as  $v_{\text{conflict}}, v_{\text{backstory}}, v_{\text{resolution}}$ , the model must learn to decide on  $v_{\text{backstory}}$  before all else.

However, we hypothesize that a teacher-forced model trained would fail to learn this story-writing plan. First, the teacher-forced model would utilize the Clever Hans cheat: it would learn the deductive skills required to fit the last segment  $v_{\text{resolution}}$  using the preceding segments revealed by the teacher in the input,  $v_{\text{conflict}}, v_{\text{backstory}}$ . With the supervision from  $v_{\text{resolution}}$  lost, the model has no explicit indication of how  $v_{\text{conflict}}$  and  $v_{\text{backstory}}$  depend on each other. These would become Indecipherable Tokens. We hypothesize that the resulting model would learn to generate uninteresting stories, interjecting arbitrary conflicts and backstories on a whim, subsequently forcing contrived resolutions upon them. While this hypothesis is not straightforward to empirically test for, we provide a more detailed

conceptual illustration in §B.

## 5. EXPERIMENTAL VERIFICATION

In this section, we demonstrate our hypothesized failure modes in practice on the graph path-finding task. We perform our experiments in both Transformers and Mamba to demonstrate that these failures are general to teacher-forced models. We begin by establishing that our teacher-forced models fit the training data but fail in-distribution. Next, we design metrics to quantify the extent to which the two hypothesized mechanisms occur (Failures 2a, 2b). Finally, we design alternative objectives to intervene and remove each of the two failure modes, to test whether the performance improves. We report additional experiments in §D.1 quantifying the Snowballing Failure 1. We describe our experimental setting more precisely below.

**Dataset.** We denote by  $G_{d,l}(N)$  for  $d, l, N \in \mathbb{N}$ , a path-star graph consisting of a center node  $v_{\text{start}}$  with degree  $d \in \mathbb{N}$ , meaning there are  $d$  different paths emerging from the center node, each consisting of  $l - 1$  nodes (excluding the start node). Node values are uniformly sampled from  $\{0, \dots, N - 1\}$  where  $N$  can be larger than the actual number of nodes in the path-star graph. In every graph, we use the center node as the starting node  $v_{\text{start}}$  and then pick as  $v_{\text{goal}}$ , the last node of one of the paths chosen uniformly at random. The order of the edges in the adjacency list is randomized. We describe the tokenization in §E.1.

For each experiment, we generate the training and test graphs from the same distribution  $\mathcal{D}$ , all with the *same* topology of  $G_{d,l}(N)$  with fixed  $d, l$  and  $N$ . Thus, any failure we demonstrate is an *in-distribution* failure, and does not arise from the inability to generalize to different problem lengths (Anil et al., 2022). We note that while the graphs are all of the same topology, this is not a trivial memorization problem for the model, since the graphs are labeled differently, and the adjacency list randomized — the model *has* to learn a general algorithm. Throughout the experiments, we fix the number of samples to  $200k$  and fix the number of node values to  $N = 100$  across topologies to enable diverse instantiations of the topology for training and testing.

**Models.** We evaluate models from two architectural families to highlight that the failures are not tied to a particular architecture but stem from the next-token prediction objective. For Transformers, we use from-scratch GPT-Mini, and pretrained GPT-2 large (Radford et al., 2019). For recurrent models, we use from-scratch Mamba (Gu & Dao, 2023). We optimize using *AdamW* (Loshchilov & Hutter, 2019) until perfect training accuracy. To rule out grokking behaviour (Power et al., 2022), we trained the cheaper models for as long as 500 epochs. More details are in §E.2.

### 5.1. Observations.

**Verifying in-distribution failure.** For a given distribution, we evaluate all our teacher-forced models by autoregressively generating solutions, and comparing that solution with the true one for an exact-match:

$$\text{Acc}_{\text{ag}}(\text{LM}_{\theta}) := \mathbb{P}(\hat{r} = r), \quad p, r \sim \mathcal{D}, \quad \hat{r} \stackrel{\text{ag}}{\sim} \text{LM}_{\theta}. \quad (4)$$

We report  $\text{Acc}_{\text{ag}}(\text{LM}_{\theta})$  for path-star graphs of varying topologies in Fig. 3 and Table 2. We observe that all models (even when pre-trained) struggle to learn the task accurately. The accuracy values are precisely limited to the value achievable if the model uniformly guesses a path starting from  $v_{\text{start}}$  i.e.,  $\approx \frac{1}{d}$ , thus establishing complete in-distribution failure. This is so even when trained to fit sample sizes up to  $200k$  to 100% accuracy, and despite the fact that the training and test graphs have identical topology. Next, we quantitatively demonstrate how this stark failure arises from our two hypothesized mechanisms (Failure 2a, 2b).

**Verifying Failure 2a (The Clever Hans cheat)** We had hypothesized that the teacher-forcing model would use cheating to fit the training tokens (the ones that follow  $r_1$  in each instance). Specifically, to predict node  $v_i$  in the true path, the model can exploit the ground truth node  $v_{i-1}$  that is revealed as input. Rather than learning to plan, the model would simply predict the node that is outwardly adjacent to  $v_{i-1}$ . To quantify whether this behavior emerges, we “teacher-force” the model with a uniform random neighbor  $v'_1$  of  $v_{\text{start}}$ . We then look for whether the model indiscriminately applies the learned Clever Hans cheat here: does the model simply follow along the path that emanates from the neighbor  $v'_1$ , not necessarily ending in  $v_{\text{goal}}$ ?

Formally, let  $\text{Unif}(\mathcal{N}(v_{\text{start}}))$  denote a uniform distribution over the set of adjacent nodes of  $v_{\text{start}}$ . For any node  $v$  in the graph, denote by  $\text{path}(v)$  the path emanating from  $v$  and going outwards, away from the start node. Notice that except for  $v = v_{\text{start}}$ , this path is unique. We thus measure

$$\text{Acc}_{\text{cheat}}(\text{LM}_{\theta}) := \mathbb{P}(\hat{r}_{1<} = \text{path}(v'_1)) \quad (5)$$

$$\text{where } p, r \sim \mathcal{D}, \quad \hat{r}_{1<} \stackrel{\text{ag}}{\sim} \text{LM}_{\theta}(\cdot; p, v_{\text{start}}, v'_1) \\ v'_1 \sim \text{Unif}(\mathcal{N}(v_{\text{start}})).$$

Empirically, we find that  $\text{Acc}_{\text{cheat}}(\text{LM}_{\theta})$  on a held-out test set is  $\approx 100\%$  almost across the board (except for graphs with very high degree where training is challenging). The exact values are in §D.1, Table 1. This establishes that to fit the training data, the teacher-forced model has exploited the Clever Hans cheat.

**Verifying Failure 2b (The Indecipherable Token)** Recall that the Clever Hans cheat only applies to all but the first node  $v_1$  after  $v_{\text{start}}$  lying on the path. After the Clever Hans

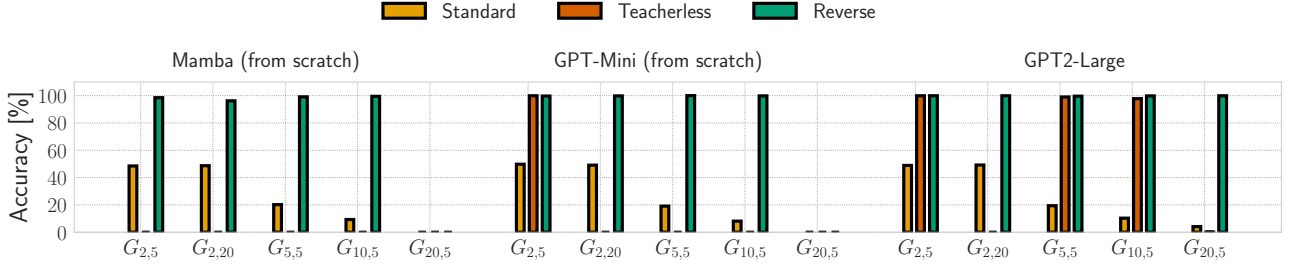


Figure 3. For different architectures, we report the accuracy of the standard teacher-forced model ( $\text{Acc}_{\text{ag}}$ , Eq 4), teacherless-trained model’s accuracy ( $\text{Acc}_{\text{\$}}$ , Eq 8) and accuracy of the model trained with reversed targets ( $\text{Acc}_{\text{rev}}$ , Eq 9) evaluated on path-finding a range of graphs (with degree in the first subscript, and path length in the second).

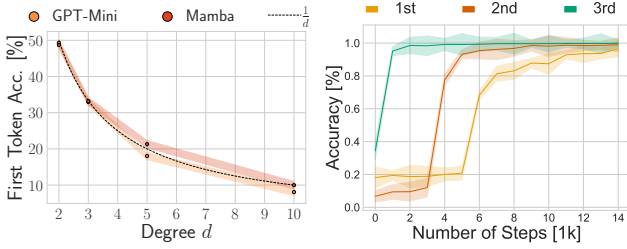


Figure 4.  $\text{Acc}_{1\text{st}}(\text{LM}_{\theta})$  (in percent %, Eq 6) for path-star graphs of various degrees  $d \in \{2, 3, 5, 10\}$  for fixed path length  $l = 5$  (left). Individual token accuracies (for  $v_1, v_2, v_3$ ) for the graph  $G_{5,5}$  under teacherless training (Eq 7) with GPT2-large (right).

cheat fits the rest of the path during training, we hypothesized that node  $v_1$  may become impossible to learn since the model is deprived of all information about the subsequent targets. To quantify this behavior, we evaluate how well the model is able to predict the difficult first node,  $v_1$ :

$$\text{Acc}_{1\text{st}}(\text{LM}_{\theta}) = \mathbb{P}(\hat{r}_1 = r_1), \quad \mathbf{p}, \mathbf{r} \sim \mathcal{D}, \hat{\mathbf{r}} \stackrel{\text{ag}}{\sim} \text{LM}_{\theta}(\cdot; \mathbf{p}). \quad (6)$$

which we estimate using the held-out test set. As shown in Fig. 4 the model achieves a low  $\text{Acc}_{1\text{st}}(\text{LM}_{\theta})$ , approximately  $1/d$ . Thus, the model indeed fails to identify that  $v_1$  is the one on the path to  $v_{\text{goal}}$ . It instead resorts to emitting one of the  $d$  neighbors of  $v_{\text{start}}$  at random.

**Removing the Clever Hans cheat via teacherless training (Monea et al., 2023).** We now consider a training setup where we prevent Clever Hans cheating (Failure 2a) and examine how learning differs. Concretely, consider modifying teacher-forcing by replacing the *input*  $\mathbf{r}$  (which reveals the ground truth) with an uninformative input  $\mathbf{r}^{\text{\$}}$ , consisting of the same special (“lookahead”) token  $\text{\$}$  repeated  $l$  times. For supervision in the loss, we still use the original target  $\mathbf{r}$ . Thus, the model cannot fit the targets by looking at the prefixes  $\mathbf{r}_{<i}$  and by predicting the next token  $v_i$  via cheating. Instead, the model only has access to the graph description

in  $\mathbf{p}$  to lookahead and fit all the targets  $v_i$  for  $i = 1, \dots, l$ . Formally, we maximize:

$$\mathcal{J}_{\text{t-less}}(\theta) = \mathbb{E}_{\mathcal{D}} \left[ \sum_{i=1}^{L_{\text{resp}}} \log \text{LM}_{\theta}(\hat{r}_i = r_i; \mathbf{p}, \mathbf{r}_{<i}^{\text{\$}}) \right]. \quad (7)$$

We denote a model trained in this fashion by  $\text{LM}_{\theta}^{\text{\$}}$  and perform inference simply by conditioning on  $\text{\$}$  tokens i.e.

$$\hat{\mathbf{r}} \stackrel{\text{\$}}{\sim} \text{LM}_{\theta}^{\text{\$}}(\cdot; \mathbf{p}) \text{ where } \hat{r}_i \sim \text{LM}_{\theta}^{\text{\$}}(\cdot; \mathbf{p}, \mathbf{r}_{<i}^{\text{\$}}).$$

and accordingly evaluate the accuracy:

$$\text{Acc}_{\text{\$}}(\text{LM}_{\theta}^{\text{\$}}) = \mathbb{P}(\hat{\mathbf{r}} = \mathbf{r}) \quad \mathbf{p}, \mathbf{r} \sim \mathcal{D}, \quad \hat{\mathbf{r}} \stackrel{\text{\$}}{\sim} \text{LM}_{\theta}^{\text{\$}}(\cdot; \mathbf{p}). \quad (8)$$

This training and inference setup was proposed in Monea et al. (2023) for the orthogonal purpose of improving the computational efficiency of inference. Our goal however is to evaluate whether forcing the model to lookahead can prevent the Clever Hans cheat from being picked up, and thereby allow the model to generalize successfully. We report the accuracy of these teacherless models in Fig. 3 and Table 3. Unfortunately, in most cases, the teacherless objective is too hard for the models to even fit the training data, likely because there is no simple cheat to employ here. However, surprisingly, on some of the easier graphs, the models not only fit the training data, but generalize well to test data. This positive result (even if in limited settings) verifies two hypotheses. First, the Clever Hans cheat is indeed a cause of failure in the original teacher-forced model. Secondly, and remarkably, with the cheat gone, these models are able to fit the first node which had once been indecipherable under teacher-forcing. This verifies our hypothesis that the Clever Hans cheat absorbs away supervision that is critical to learn the first token. Shortly at the end of this section, we provide more intuition about how exactly teacherless models are able to solve this task with the Clever Hans cheat out of their way.



**Removing the Indecipherable Token failure via path reversal.** Back in the teacher-forcing setup, we make a slight change: we train the model to predict the reversal of the true path  $\mathbf{r}$ . Indeed, prior works (Lee et al., 2023; Shen et al., 2023) have proposed reversal in the context of addition tasks as a way of explicitly guiding the next-token predictor to learn a simpler algorithm. Likewise, in our “reversed” path-finding task, the model now needs to predict  $v_{\text{goal}}$  first and make its way to  $v_{\text{start}}$ ; the hope is that since there is only one unique path emanating from  $v_{\text{goal}}$ , there is no planning required. Thus we should never run into an Indecipherable Token. Every next node can be learned as the node that is inwardly adjacent to the previous node.

Notationally, we let  $\text{LM}_{\theta}^{\text{rev}}$  be the model trained to maximize  $\mathcal{J}_{\text{next-token}}$  with the targets (and the teacher-forced inputs) set to  $\mathbf{r}^{\text{rev}} = r_{L_{\text{resp}}}, \dots, r_1$ , the reversal of  $\mathbf{r}$ . We then measure the autoregressive accuracy by comparing against  $\mathbf{r}^{\text{rev}}$ :

$$\text{Acc}_{\text{rev}}(\text{LM}_{\theta}^{\text{rev}}) = \mathbb{P}(\hat{\mathbf{r}} = \mathbf{r}^{\text{rev}}), \quad \mathbf{p}, \mathbf{r} \sim \mathcal{D}, \hat{\mathbf{r}} \stackrel{\text{ag}}{\sim} \text{LM}_{\theta}^{\text{rev}}(\cdot; \mathbf{p}) \quad (9)$$

We display the results in Fig. 3 and Table 4. As expected, we observe that reversing significantly boosts learning, allowing even models trained from scratch to solve the task. This verifies that for the standard model, indecipherability of the first token was indeed a roadblock to successful learning.

## 5.2. Why the failure of teacher-forcing is remarkable.

We conclude by emphasizing that the success of the reversed training (and also, the occasional success of teacherless training) make the in-distribution failure of teacher-forcing particularly surprising. Recall that, when viewed left-to-right, our problem requires complex planning — evaluating multiple paths and selecting the right one — but when viewed right-to-left, the problem is straightforward. The experiments on the reversed formulation confirm that the right-to-left solution is not only expressible by our architectures, but also learnable via gradient descent. Evidently, the left-to-right teacher-forced model is unable to view the problem any differently and falls into the traps outlined in §4.

**Intuition for teacherless training.** We hypothesize that even teacherless training allows the model to implicitly learn the right-to-left view. Concretely, the teacherless model cannot use the trivial Clever Hans cheat to fit the data, since the ground truth prefixes are not available during training. Nor is it explicitly prescribed to fit the target right-to-left. Instead the model is tasked with using only the graph description in  $\mathbf{p}$  to fit all the target nodes  $\mathbf{r}$  (implicitly requiring a lookahead beyond just the next token). In this paradigm, the model would first fit the target token that is simplest to deduce using only information available in the prefix  $\mathbf{p}$ : this is the penultimate vertex  $r_{l-1}$  which is the unique token that

precedes the goal (and can be discovered using a simple scan of the prefix). Once the model figures this out, the model can similarly work backwards to fit each node  $r_{i-1}$  using the previously-fit  $r_i$ . (We note that this solution is still a fairly difficult one to implement in the teacherless setting, which surprisingly some models nevertheless learn — we describe this in §C.)

Our hypothesis is borne out in Fig. 4 where we see that the teacherless model automatically learns right-to-left: the later tokens achieve higher accuracy earlier. Thus, the teacherless objective provides one possible way for future work to build alternatives to next-token prediction that force the models to look ahead, without falling into the next-token traps of the Clever Hans cheat or the Indecipherable Token failure.

## 6. RELATED WORK

We provide an elaborate survey of the various arguments that have been made in favor of or against next-token prediction. We hope this can help consolidate a debate that has been scattered over the literature. Part of the survey is deferred to §F.

**Arguments in support of next-token prediction.** Shannon (1948; 1951); Alabdulmohsin et al. (2024) demonstrate that language has enough redundancy to be conducive for next-token prediction. Empirically, Shlegeris et al. (2022) find that modern language models are surprisingly better than humans at next-token prediction on the text dataset, OpenWebText (Gokaslan & Cohen, 2019). But this does not preclude the possibility that next-token predictors may still be poor at planning. Furthermore, the above result may be confounded by the ability of language models to store more general knowledge than humans.

On the theoretical side, Merrill & Sabharwal (2023b); Feng et al. (2023) show that autoregressive Transformers that generate chains of thought have a markedly larger *expressive* power. Most relevant to us is the positive *learnability* results of Malach (2023); Wies et al. (2023) which argue that complex multi-hop tasks that are otherwise unlearnable, become learnable via next-token prediction when there is a preceding chain-of-thought supervision for each hop. Our negative result does not contradict this. In our path-finding problem, learning the first token requires an implicit chain of thought (the reversed path) that we do not provide *before* the first token.

**Arguments against next-token prediction.** The most well-formulated criticism of next-token prediction is what we term as the snowballing error, scattered in literature, both in recent work Dziri et al. (2023); LeCun (2024) and much earlier in Kääriäinen (2006); Ross & Bagnell (2010). The earlier works capture an additional notion of snowballing, wherein, once an erroneous sub-optimal action is commit-

ted, the model is more likely to commit more sub-optimal actions since it has wandered into territories that it was not trained on. Implicitly, the error here is not evaluated as an exact match of the response (i.e.,  $r \neq \hat{r}$ ) but as a cumulative notion of error over all steps (e.g.,  $\sum \mathbf{1}[r_i \neq \hat{r}_i]$ ). In this setting, there is an additional cause of failure called *exposure bias*: the teacher-forced model has only been trained on correct trajectories, and has not learned how to recover from poor trajectories. Nevertheless, all existing criticisms assumes that teacher-forcing has learned an accurate next-token predictor in the first place, which is the very assumption our counterexample challenges. Indeed, in our setting the error happens “instantaneously” at the beginning of inference.

Our main counterexample can be seen as a way of formalizing an emerging, informal intuition that is often worded as “autoregressive next-token predictors are ill-suited for planning tasks”. Indeed, Momennejad et al. (2023); Valmeekam et al. (2023) report failures on several planning tasks framed as word problems (including path-finding in Momennejad et al. (2023)) and Bubeck et al. (2023) on various arithmetic, summarization and poem/story generation problems (which they speculate is a limit of autoregressivity). McCoy et al. (2023) argue that, for such tasks, the performance of the model must greatly depend on its frequency during pretraining. However, we demonstrate that even when trained on many samples from a distribution, the next-token predictor can fail on the very distribution.

A closely-related criticism (Bubeck et al., 2023; Dawid & LeCun, 2023; LeCun, 2024; Du et al., 2023a) is that to model human thinking, we need to model two types of thinking as outlined in Kahneman (2011): a fast (System 1) thinking process that is also guided by a slower (System 2) thinking process. Theoretically, Lin et al. (2021) show that there are formal languages for which expressing some next-tokens may require super-polynomial time or parameter count during *inference*. Du et al. (2023a) informally note that some next tokens can be hard to *learn* as they require a global understanding of what will be uttered in the future.<sup>4</sup>

Our work extends and clarifies this discourse by introducing the Clever Hans cheat and the Indecipherable Token failure. Next, we empirically report our failure modes in both the Transformer (Vaswani et al., 2017) and the Mamba structured state space model (Gu & Dao, 2023). This establishes that what we witness is indeed a failure of next-token pre-

diction (and not of the Transformer architecture as some existing criticisms are framed). Importantly, existing literature pins these failures broadly on the next-token prediction paradigm and interchangeably, on the inability of the autoregressive architecture to backtrack. We emphasize the need to differentiate between the two types of next-token prediction (teacher-forcing and autoregressive inference) as they lead to distinct planning-related failures and require distinct solutions.

**Going beyond next-token prediction.** Inference-time techniques like chain-of-thought and its variants (Wei et al., 2022; Yao et al., 2023a; Besta et al., 2023; Yao et al., 2023b) or those that elicit feedback from the model (Madaan et al., 2023; Huang et al., 2022; Shinn et al., 2023) can be thought of as going beyond conventional form of inference by allowing the model to think more before producing its final answer. However, the backbone in these models are still trained by standard teacher-forcing. While other techniques (Burtsev et al., 2020; Xue et al., 2023; Goyal et al., 2023) train the model to explicitly think more, even these boil down to next-token prediction during training.

Other works have explored architectures and objectives that train the backbone to go beyond next-token prediction. This includes non-autoregressive models (Gu et al., 2018), energy-based models (Dawid & LeCun, 2023), diffusion models (Gong et al., 2023), and variants of Transformers learning to predict multiple tokens at the same go (Qi et al.; Monea et al., 2023) or injecting “lookahead” data (Du et al., 2023a). Teacherless training — proposed as “parallel speculative sampling (PaSS)” in Monea et al. (2023) — provides an arguably simple such approach that involves a trivial modification to teacher-forcing. Note that while research in parallel decoding too is concerned with predicting multiple future tokens (Stern et al., 2018), the goal is purely inference-time efficiency — which is also the setting under which Monea et al. (2023) propose teacherless training.

One may argue that reinforcement learning-based training (Ranzato et al., 2016; Wu et al., 2016; Bahdanau et al., 2017; Paulus et al., 2018; Ziegler et al., 2019; Stiennon et al., 2020; Ouyang et al., 2022) is another way to build backbones that go beyond teacher-forcing. However, it is worth noting that the gradients in these techniques boil down to teacher-forcing on the model’s own generated answer. Furthermore, if we desire that the model be able to generate a solution that can plan ahead of time, it is unclear how a model can go from a complete inability to plan (that may assign near-zero probability to the true plan in an exponential space of solutions), to discovering the correct plan simply through preference-based feedback (see (Havrilla et al., 2024) for related empirical evidence).

Another line of work — spanning language (Bengio et al., 2015; Goyal et al., 2016), imitation learning (Ross et al.,

<sup>4</sup>The Indecipherable Token failure is related to the “locally unlearnable” hypothesis of Du et al. (2023a), but is not the same. The (first) Indecipherable Token in the path-star problem is locally learnable by the teacherless model, even without access to the full target sequence (which is only presented after this token). Thus, the inability of the teacher-forced model to learn the Indecipherable Token is crucially related to the fact that the supervisory targets that come after this token are lost to the Clever Hans cheat.

2011; Ross & Bagnell, 2010; 2014) and structured prediction (Daumé III et al., 2009; Chang et al., 2015) — has been aimed at addressing the Snowball Failure, under the assumption that the model has otherwise learned an accurate next-step predictor. Broadly, the idea is to train the model on a mixture of the ground truth sequences and the model-generated sequences themselves, as a way to ensure that the test-time and training-time distributions are as similar as possible. These techniques however do not address the failure to learn a good next-step predictor in the first place.

As for reversal-based training, Lee et al. (2023); Shen et al. (2023) observe that addition tasks become much simpler when the digits are reversed. Their argument is that this explicitly assists the model to learn a simpler algorithm. When it comes to natural language however, Papadopoulos et al. (2024) find that reversing hurts the model’s perplexity.

#### End-to-end reasoning and chain-of-thought supervision.

In our path-star graph, learning the Indecipherable Token (the first node  $v_1$ ) can be thought of as a task whose end target is  $v_1$ , but whose implicit intermediate targets (or “chain-of-thought”) correspond to the unique path starting from  $v_{\text{goal}}$  headed towards  $v_1$  (although this is only provided as supervision after the first token). In this terminology, we can rephrase our claim as the model failing to learn the end target once the intermediate targets are lost to the Clever Hans Cheat.

Such limits of end-to-end learning have been echoed in literature on learning with chain-of-thought-type supervision. Recent theoretical works have shown broad classes of tasks (e.g., any function efficiently computed by a Turing machine) where prepending CoT to the end target allows efficiently learning tasks; yet, there are “multi-hop reasoning” tasks that are unlearnable end-to-end (i.e., without intermediate supervision) either due to computational hardness (Wies et al., 2023) or representational limits (Malach, 2023)). Earlier theoretical works Shalev-Shwartz et al. (2017); Shalev-Shwartz & Shashua (2016) have similarly proven negative results for end-to-end learning in similar settings. Similar empirical arguments have been made in neural network literature (Gülçehre & Bengio, 2016; Glasmachers, 2017) and also more recently, in language models on complex reasoning and math problems (Nye et al., 2021; Ling et al., 2017; Cobbe et al., 2021; Piekos et al., 2021; Zelikman et al., 2022; Recchia, 2021; Cobbe et al., 2021).

It is worth noting though that the above lines of work are concerned with chain-of-thought that is present before the end target; in our setup, this supervision is presented only *after* the end target. Surprisingly, some of our teacherless models manage to utilize even such hindsight chain-of-thought. This success is not fully explained by existing positive results about chain-of-thought supervision, such as Wies et al. (2023); Malach (2023), where supervision is provided *be-*

*fore* the end target.

**Shortcut-learning in language models.** A line of work has empirically and theoretically analyzed how Transformer-based language models learn superficial shortcuts to (partially) solve tasks such as learning multiplication (Dziri et al., 2023), automata (Liu et al., 2023), recursion (Young & You, 2023), reading comprehension (Lai et al., 2021) and multiple-choice questions (Ranaldi & Zanzotto, 2023)

However, these shortcuts must *not* be confused with the Clever Hans cheating induced by teacher-forcing. First, these aforementioned shortcuts exist independent of teacher-forcing: these are correlations between the prefix (such as the initial digits of two multiplicands) and the final answer (the initial digits of the product) in the underlying training distribution. But Clever Hans cheats arise only upon teacher-forcing as they are correlations between the prefixes of the answer itself to the rest of the answer. Second, the above shortcuts only fail out-of-distribution (such as when the number of multiplied digits is increased, where the failure is in length generalization (Anil et al., 2022)). In contrast, the Clever Hans cheat is much more severe as it causes in-distribution failure. Thirdly, the aforementioned empirical observations are specific to Transformers, and the theoretical arguments rely crucially on properties of the Transformer (such as its non-recurrence and convolution, or its self-attention modules). Our argument however only relies on the teacher-forcing objective with no reliance on the Transformer architecture, and is demonstrated even for the recurrent Mamba architecture.

Please see § F for more related works.

## 7. LIMITATIONS

We note that our arguments are empirical and conceptual. We have not provided a formal proof for our arguments. We have also not demonstrated failure for very large models such as Llama2 (Touvron et al., 2023) or Mistral (Jiang et al., 2023). Next, beyond the minimal path-finding setting, we have not demonstrated or characterized the range of problems where teacher-forcing-induced failure may occur. We only intuitively believe it should extend to other problem-solving tasks and creative-writing tasks that require lookahead. It is also unclear if it generalizes to run-of-the-mill text-generation tasks.

## 8. CONCLUSION

Next-token prediction lies at the heart of modern language models which have demonstrated tremendous empirical success in a range of general tasks. Theoretically too, we know by the chain rule of probability that, next-token predictors can express any imaginable distribution over tokens. Thus,

it is tempting to view next-token prediction as a formidable approach to modeling language and intelligence. Our work crystallizes the core arguments around why this optimism may be misplaced.

First, we emphasize not to conflate the two modes of next-token prediction: autoregressive inference and teacher-forced training. While existing criticisms primarily challenge autoregressive inference, they assume that teacher-forcing learns a good next-token predictor. We challenge this very assumption, finding that even in a straightforward task, there is failure due to teacher-forcing — not due to autoregressive inference or the architecture. This casts a shadow over more complex tasks. For instance, as we speculate in § B, can a model trained to predict the next token of tens of thousands of fiction novels, learn to generate plot twists?

An immediate way to circumvent this, as our reversal experiments suggest, is to train with chain-of-thought supervision, echoing Malach (2023); Wies et al. (2023). However, it is unclear how that is possible in more unstructured tasks like story-writing. To that end, our minimal counterexample and the idea of teacherless training (Monea et al., 2023) may inspire alternative paradigms to next-token prediction in practice. Overall, we hope our analyses provide a solid ground to pursue future debates on next-token prediction.

## 9. IMPACT

Our results outline the limits of a foundational technique that lies at the heart of modern AI systems. Naturally, there are many potential downstream societal consequences that would apply at large to such foundational work, none we feel must be specifically highlighted here.

**Acknowledgments:** We would like to thank Colin Raffel, Tiago Pimentel, and Surbhi Goel for their extensive feedback on a draft of this preprint, especially for pointers to some key references.

## REFERENCES

- Alabdulmohsin, I., Tran, V. Q., and Dehghani, M. Fractal patterns may unravel the intelligence in next-token prediction, 2024.
- Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.2, knowledge manipulation. *CoRR*, abs/2309.14402, 2023. doi: 10.48550/ARXIV.2309.14402. URL <https://doi.org/10.48550/arXiv.2309.14402>.
- Anil, C., Wu, Y., Andreassen, A., Lewkowycz, A., Misra, V., Ramasesh, V. V., Slone, A., Gur-Ari, G., Dyer, E., and Neyshabur, B. Exploring length generalization in large language models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- Arkoudas, K. Chatgpt is no stochastic parrot. but it also claims that 1 is greater than 1. *Philosophy & Technology*, 36(3):54, 2023.
- Artetxe, M., Du, J., Goyal, N., Zettlemoyer, L., and Stoyanov, V. On the role of bidirectionality in language model pre-training. In Goldberg, Y., Kozareva, Z., and Zhang, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pp. 3973–3985. Association for Computational Linguistics, 2022.
- Bahdanau, D., Brakel, P., Xu, K., Goyal, A., Lowe, R., Pineau, J., Courville, A. C., and Bengio, Y. An actor-critic algorithm for sequence prediction. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. On the dangers of stochastic parrots: Can language models be too big? In Elish, M. C., Isaac, W., and Zemel, R. S. (eds.), *FACtT ’21: 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event / Toronto, Canada, March 3-10, 2021*, pp. 610–623. ACM, 2021.
- Bengio, S., Vinyals, O., Jaitly, N., and Shazeer, N. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 1171–1179, 2015.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Gianinazzi, L., Gajda, J., Lehmann, T., Podstawski, M., Niewiadomski, H., Nyczyk, P., and Hoefler, T. Graph of thoughts: Solving elaborate problems with large language models. *CoRR*, abs/2308.09687, 2023. doi: 10.48550/ARXIV.2308.09687. URL <https://doi.org/10.48550/arXiv.2308.09687>.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S. M., Nori, H., Palangi, H., Ribeiro, M. T., and Zhang, Y. Sparks of artificial general intelligence: Early experiments with GPT-4. *CoRR*, abs/2303.12712, 2023. doi: 10.48550/ARXIV.2303.12712. URL <https://doi.org/10.48550/arXiv.2303.12712>.
- Burtsev, M. S., Kuratov, Y., Peganov, A., and Sapunov, G. V. Memory transformer. *arXiv preprint arXiv:2006.11527*, 2020.

- Chang, K., Krishnamurthy, A., Agarwal, A., III, H. D., and Langford, J. Learning to search better than your teacher. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 2058–2066. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/changb15.html>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Daumé III, H., Langford, J., and Marcu, D. Search-based structured prediction. *Mach. Learn.*, 75(3):297–325, 2009.
- Dawid, A. and LeCun, Y. Introduction to latent variable energy-based models: A path towards autonomous machine intelligence. *CoRR*, abs/2306.02572, 2023. doi: 10.48550/ARXIV.2306.02572. URL <https://doi.org/10.48550/arXiv.2306.02572>.
- Du, L., Mei, H., and Eisner, J. Autoregressive modeling with lookahead attention. *CoRR*, abs/2305.12272, 2023a.
- Du, L., Torroba Hennigen, L., Pimentel, T., Meister, C., Eisner, J., and Cotterell, R. A measure-theoretic characterization of tight language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9744–9770, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.543. URL <https://aclanthology.org/2023.acl-long.543>.
- Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin, B. Y., West, P., Bhagavatula, C., Bras, R. L., Hwang, J. D., Sanyal, S., Welleck, S., Ren, X., Ettinger, A., Harchaoui, Z., and Choi, Y. Faith and fate: Limits of transformers on compositionality. *CoRR*, abs/2305.18654, 2023.
- Feng, G., Zhang, B., Gu, Y., Ye, H., He, D., and Wang, L. Towards revealing the mystery behind chain of thought: a theoretical perspective. *CoRR*, abs/2305.15408, 2023.
- Glasmachers, T. Limits of end-to-end learning. In Zhang, M. and Noh, Y. (eds.), *Proceedings of The 9th Asian Conference on Machine Learning, ACML 2017, Seoul, Korea, November 15-17, 2017*, volume 77 of *Proceedings of Machine Learning Research*, pp. 17–32. PMLR, 2017. URL <http://proceedings.mlr.press/v77/glasmachers17a.html>.
- Gokaslan, A. and Cohen, V. Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- Gong, S., Li, M., Feng, J., Wu, Z., and Kong, L. Diffuseq: Sequence to sequence text generation with diffusion models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=jQj-rLVXsj>.
- Goyal, A., Lamb, A., Zhang, Y., Zhang, S., Courville, A. C., and Bengio, Y. Professor forcing: A new algorithm for training recurrent networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 4601–4609, 2016.
- Goyal, S., Ji, Z., Rawat, A. S., Menon, A. K., Kumar, S., and Nagarajan, V. Think before you speak: Training language models with pause tokens. *CoRR*, abs/2310.02226, 2023.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces, 2023.
- Gu, J., Bradbury, J., Xiong, C., Li, V. O. K., and Socher, R. Non-autoregressive neural machine translation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=B1l8BtlCb>.
- Gülçehre, Ç. and Bengio, Y. Knowledge matters: Importance of prior information for optimization. *J. Mach. Learn. Res.*, 17:8:1–8:32, 2016. URL <http://jmlr.org/papers/v17/gulchere16a.html>.
- Gurnee, W., Nanda, N., Pauly, M., Harvey, K., Troitskii, D., and Bertsimas, D. Finding neurons in a haystack: Case studies with sparse probing. *CoRR*, abs/2305.01610, 2023. doi: 10.48550/ARXIV.2305.01610. URL <https://doi.org/10.48550/arXiv.2305.01610>.
- Havrilla, A., Du, Y., Raparthy, S. C., Nalmpantis, C., Dwivedi-Yu, J., Zhuravinskyi, M., Hambro, E., Sukhbaatar, S., and Raileanu, R. Teaching large language models to reason with reinforcement learning, 2024.
- Huang, W., Xia, F., Xiao, T., Chan, H., Liang, J., Florence, P., Zeng, A., Tompson, J., Mordatch, I., Chebotar, Y., Sermanet, P., Jackson, T., Brown, N., Luu, L., Levine, S., Hausman, K., and Ichter, B. Inner monologue: Embodied reasoning through planning with language models. In Liu, K., Kulic, D., and Ichnowski, J. (eds.), *Conference on Robot Learning, CoRL 2022, 14-18 December 2022, Auckland, New Zealand*, volume 205 of *Proceedings of Machine Learning Research*, pp. 1769–1782. PMLR, 2022.



- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023.
- Kääriäinen, M. Lower bounds for reductions. In *Atomic Learning Workshop*, 2006.
- Kahneman, D. *Thinking, fast and slow*. Farrar, Straus and Giroux, 2011.
- Lai, Y., Zhang, C., Feng, Y., Huang, Q., and Zhao, D. Why machine reading comprehension models learn shortcuts? In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pp. 989–1002. Association for Computational Linguistics, 2021.
- LeCun, Y. Do large language models need sensory grounding for meaning and understanding? University Lecture, 2024.
- Lee, N., Sreenivasan, K., Lee, J. D., Lee, K., and Papailiopoulos, D. Teaching arithmetic to small transformers. *CoRR*, abs/2307.03381, 2023.
- Li, Y., Huang, Y., Ildiz, M. E., Rawat, A. S., and Oymak, S. Mechanics of next token prediction with self-attention. In *27th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2024.
- Lin, C., Jaech, A., Li, X., Gormley, M. R., and Eisner, J. Limitations of autoregressive models and their alternatives. In Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tür, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., and Zhou, Y. (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pp. 5147–5173. Association for Computational Linguistics, 2021.
- Ling, W., Yogatama, D., Dyer, C., and Blunsom, P. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In Barzilay, R. and Kan, M. (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pp. 158–167. Association for Computational Linguistics, 2017.
- Liu, B., Ash, J. T., Goel, S., Krishnamurthy, A., and Zhang, C. Transformers learn shortcuts to automata. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=De4FYqjFueZ>.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Lv, A., Zhang, K., Xie, S., Tu, Q., Chen, Y., Wen, J., and Yan, R. Are we falling in a middle-intelligence trap? an analysis and mitigation of the reversal curse. *CoRR*, abs/2311.07468, 2023.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., Welleck, S., Majumder, B. P., Gupta, S., Yazdanbakhsh, A., and Clark, P. Self-refine: Iterative refinement with self-feedback. *CoRR*, abs/2303.17651, 2023. doi: 10.48550/ARXIV.2303.17651. URL <https://doi.org/10.48550/arXiv.2303.17651>.
- Malach, E. Auto-regressive next-token predictors are universal learners. *CoRR*, abs/2309.06979, 2023. doi: 10.48550/ARXIV.2309.06979. URL <https://doi.org/10.48550/arXiv.2309.06979>.
- McCoy, R. T., Yao, S., Friedman, D., Hardy, M., and Griffiths, T. L. Embers of autoregression: Understanding large language models through the problem they are trained to solve. *CoRR*, abs/2309.13638, 2023.
- Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in GPT. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- Merrill, W. and Sabharwal, A. The parallelism tradeoff: Limitations of log-precision transformers, 2023a.
- Merrill, W. and Sabharwal, A. The expressive power of transformers with chain of thought. *CoRR*, abs/2310.07923, 2023b. doi: 10.48550/ARXIV.2310.07923. URL <https://doi.org/10.48550/arXiv.2310.07923>.
- Momennejad, I., Hasanbeig, H., Frujeri, F. V., Sharma, H., Ness, R. O., Jojic, N., Palangi, H., and Larson, J. Evaluating cognitive maps and planning in large language models with cogeal. *CoRR*, abs/2309.15129, 2023. doi: 10.48550/ARXIV.2309.15129. URL <https://doi.org/10.48550/arXiv.2309.15129>.
- Monea, G., Joulin, A., and Grave, E. Pass: Parallel speculative sampling. *CoRR*, abs/2311.13581, 2023. doi: 10.48550/ARXIV.2311.13581. URL <https://doi.org/10.48550/arXiv.2311.13581>.

- Nye, M. I., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., Dohan, D., Lewkowycz, A., Bosma, M., Luan, D., Sutton, C., and Odena, A. Show your work: Scratchpads for intermediate computation with language models. *CoRR*, abs/2112.00114, 2021. URL <https://arxiv.org/abs/2112.00114>.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., Drain, D., Ganguli, D., Hatfield-Dodds, Z., Hernandez, D., Johnston, S., Jones, A., Kernion, J., Lovitt, L., Ndousse, K., Amodei, D., Brown, T., Clark, J., Kaplan, J., McCandlish, S., and Olah, C. In-context learning and induction heads. abs/2209.11895, 2022. doi: 10.48550/ARXIV.2209.11895. URL <https://doi.org/10.48550/arXiv.2209.11895>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P. F., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- Pal, K., Sun, J., Yuan, A., Wallace, B. C., and Bau, D. Future lens: Anticipating subsequent tokens from a single hidden state. In Jiang, J., Reitter, D., and Deng, S. (eds.), *Proceedings of the 27th Conference on Computational Natural Language Learning, CoNLL 2023, Singapore, December 6-7, 2023*, pp. 548–560. Association for Computational Linguistics, 2023.
- Papadopoulos, V., Wenger, J., and Hongler, C. Arrows of time for large language models, 2024.
- Paulus, R., Xiong, C., and Socher, R. A deep reinforced model for abstractive summarization. In *6th International Conference on Learning Representations, ICLR 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- Pfau, J., Infanger, A., Sheshadri, A., Panda, A., Michael, J., and Huebner, C. Eliciting language model behaviors using reverse language models. In *Socially Responsible Language Modelling Research*, 2023. URL <https://openreview.net/forum?id=m6xyTie61H>.
- Pfungst, O. and Rahn, C. L. *Clever Hans (the horse of Mr. Von Osten) a contribution to experimental animal and human psychology*. New York, H. Holt and company, 1911. URL <https://www.biodiversitylibrary.org/item/116908>. <https://www.biodiversitylibrary.org/bibliography/56164>.
- Piekos, P., Malinowski, M., and Michalewski, H. Measuring and improving bert’s mathematical abilities by predicting the order of reasoning. In Zong, C., Xia, F., Li, W., and Navigli, R. (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021*, pp. 383–394. Association for Computational Linguistics, 2021.
- Power, A., Burda, Y., Edwards, H., Babuschkin, I., and Misra, V. Grokking: Generalization beyond overfitting on small algorithmic datasets, 2022.
- Qi, W., Yan, Y., Gong, Y., Liu, D., Duan, N., Chen, J., Zhang, R., and Zhou, M. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. In Cohn, T., He, Y., and Liu, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pp. 2401–2410.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. 2019.
- Ranaldi, L. and Zanzotto, F. M. Hans, are you clever? clever hans effect analysis of neural systems, 2023.
- Ranzato, M., Chopra, S., Auli, M., and Zaremba, W. Sequence level training with recurrent neural networks. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1511.06732>.
- Recchia, G. Teaching autoregressive language models complex tasks by demonstration. *CoRR*, abs/2109.02102, 2021. URL <https://arxiv.org/abs/2109.02102>.
- Ross, S. and Bagnell, D. Efficient reductions for imitation learning. In Teh, Y. W. and Titterton, D. M. (eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, volume 9 of *JMLR Proceedings*, pp. 661–668. JMLR.org, 2010.
- Ross, S. and Bagnell, J. A. Reinforcement and imitation learning via interactive no-regret learning. abs/1406.5979, 2014. URL <http://arxiv.org/abs/1406.5979>.
- Ross, S., Gordon, G. J., and Bagnell, D. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011*, volume 15 of *JMLR Proceedings*, pp. 627–635. JMLR.org, 2011.

- Shalev-Shwartz, S. and Shashua, A. On the sample complexity of end-to-end training vs. semantic abstraction training. *CoRR*, abs/1604.06915, 2016. URL <http://arxiv.org/abs/1604.06915>.
- Shalev-Shwartz, S., Shamir, O., and Shammah, S. Failures of gradient-based deep learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3067–3075. PMLR, 2017.
- Shannon, C. E. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 1948. doi: 10.1002/j.1538-7305.1948.tb01338.x.
- Shannon, C. E. Prediction and entropy of printed english. *The Bell System Technical Journal*, 30(1):50–64, 1951. doi: 10.1002/j.1538-7305.1951.tb01366.x.
- Shen, R., Bubeck, S., Eldan, R., Lee, Y. T., Li, Y., and Zhang, Y. Positional description matters for transformers arithmetic. *CoRR*, abs/2311.14737, 2023. doi: 10.48550/ARXIV.2311.14737. URL <https://doi.org/10.48550/arXiv.2311.14737>.
- Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning, 2023.
- Shlegeris, B., Roger, F., Chan, L., and McLean, E. Language models are better than humans at next-token prediction. *CoRR*, abs/2212.11281, 2022. doi: 10.48550/ARXIV.2212.11281. URL <https://doi.org/10.48550/arXiv.2212.11281>.
- Springer, J. M., Kotha, S., Fried, D., Neubig, G., and Raghunathan, A. Repetition improves language model embeddings, 2024.
- Stern, M., Shazeer, N., and Uszkoreit, J. Blockwise parallel decoding for deep autoregressive models. In Bengio, S., Wallach, H. M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 10107–10116, 2018.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D. M., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize from human feedback. *CoRR*, abs/2009.01325, 2020. URL <https://arxiv.org/abs/2009.01325>.
- Thrapoulidis, C. Implicit bias of next-token prediction, 2024.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models, 2023.
- Valmeekam, K., Marquez, M., Olmo, A., Sreedharan, S., and Kambhampati, S. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change, 2023.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- Welleck, S., Kulikov, I., Kim, J., Pang, R. Y., and Cho, K. Consistency of a recurrent language model with respect to incomplete decoding. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, November 2020.
- Wies, N., Levine, Y., and Shashua, A. Sub-task decomposition enables learning in sequence to sequence tasks. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=BrJATVZDWEH>.
- Williams, R. J. and Zipser, D. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280, 1989. doi: 10.1162/neco.1989.1.2.270.

- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, L., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016. URL <http://arxiv.org/abs/1609.08144>.
- Xue, F., Likhoshesterov, V., Arnab, A., Houlsby, N., Dehghani, M., and You, Y. Adaptive computation with elastic input sequence. In *International Conference on Machine Learning, ICML 2023, Proceedings of Machine Learning Research*. PMLR, 2023.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *CoRR*, abs/2305.10601, 2023a.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. R., and Cao, Y. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023b. URL [https://openreview.net/pdf?id=WE\\_vluYUL-X](https://openreview.net/pdf?id=WE_vluYUL-X).
- Young, T. and You, Y. On the inconsistencies of conditionals learned by masked language models. *CoRR*, abs/2301.00068, 2023. URL <https://doi.org/10.48550/arXiv.2301.00068>.
- Zelikman, E., Wu, Y., Mu, J., and Goodman, N. D. Star: Bootstrapping reasoning with reasoning. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/639a9a172c044fbb64175b5fad42e9a5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/639a9a172c044fbb64175b5fad42e9a5-Abstract-Conference.html).
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P. F., and Irving, G. Fine-tuning language models from human preferences. *CoRR*, abs/1909.08593, 2019. URL <http://arxiv.org/abs/1909.08593>.

## A. TEACHER-FORCING FAILURE AND SNOWBALLING FAILURE ARE DISTINCT

We emphasize that, while both the Clever Hans failure mode and the Snowball mode are both indicative of the inability to plan, these failure modes are also orthogonal to each other, and demand different solutions. We make this a bit more formal:

**Proposition 4.** *In the path finding problem of §4.1, there exists a next-token predictor that experiences Failures 2a, 2b due to teacher-forcing, but not the snowballing error Failure 1 due to autoregressive inference. Conversely, there exists a next-token predictor that experiences the latter failure but not the former.*

*Proof.* Consider the model learned via teacher-forcing on the graph problem. During inference, we saw that it suffers a debilitating error right in the first step (with accuracy of  $1/d$  for degree  $d$  of the start node). Thus, during inference the error that is experienced is not from an accumulation over length. In fact, if only the first node is set correctly during inference, a model with the perfect Clever Hans cheat, would achieve 100% accuracy rate. Such a model does not experience snowballing errors.

On the other hand, consider a model, that in each step predicts the correct next vertex with a high accuracy of  $1 - \epsilon$  for small  $\epsilon$ . Such a model clearly has learned the correct plan, albeit with minor errors in each token. These errors however can snowball during inference. Thus, this model has no failure due to teacher-forcing, but will fail during autoregressive inference, if the path length is long.  $\square$

**Differing solutions.** Based on the above simple illustration, we note that the two failures need different solution approaches. Specifically, while snowballing errors may be fixable via “backtracking-and-planning” wrappers, teacher-forcing failures is a pathology that cannot be solved post-hoc.

## B. AN ILLUSTRATION VIA STORY-TELLING

Can a teacher-forced model merely trained on thousands of stories learn to write plot twists? Indeed, Bubeck et al. (2023) report instances where models can fail to accomplish tasks involving creative-writing (e.g., poems). We speculatively extend our discussion in §4 to reason about this scenario. Consider for example, teacher-forcing on the following story that follows an often-used plot outline:

- **Event 1 (Setup):** Alex and Bob, who are friends, are trying to defeat the Evil King.
- **Event 2 (Conflict):** One day, surprisingly, Bob turns against Alex, and tries to thwart Alex’s plans, *albeit unsuccessfully*.
- **Event 3:** Alex thinks Bob is evil too, defeats Bob first.
- **Event 4 (Backstory):** Losing the battle, Bob reveals he is a double-agent. In his final words, Bob explains he was ordered to defeat Alex.
- **Event 5 (Resolution):** To preserve the King’s trust, Bob obeyed the command, but also *deliberately* failed at it. Bob then relays critical information he extracted from the King’s inner circles.
- **Event 6:** Alex uses Bob’s insider information to defeat the King.

Evidently, this story requires a plan: Event 5 is a key plot resolution that the narrator must have planned before methodically generating parts of the setup in Event 1 (introducing Bob as a friend) and the conflict in Event 2 (Bob’s turning against Alex, and failing at it). While training, the model must thus treat the story as a whole, and tease apart these dependencies between the events, some of which may be anti-chronological (akin to how, in the path-star graph, the model must learn that the problem is straightforwardly solvable when viewed from right-to-left).

However, we hypothesize that a teacher-forced model would take a rigid chronological (left-to-right) view. First, it would use the Clever Hans cheat to easily fit the plot resolution in Event 5: the model would use the facts of Event 4 and Event 2 (revealed as input) to fit the content of Bob’s final words. Thus, the content of Event 5 would no longer be available as supervision to guide how the model fits Event 1 and Event 2. When the model tries to fit these earlier events, these events would become Indecipherable Tokens — the model would simply learn to fit them as arbitrary events. Thus, we conjecture that a model trained via teacher-forcing merely on raw, unannotated texts of stories — however many stories they may be — would not learn to plan its stories, and would instead create arbitrary twists and turns during inference, and improvise upon that.



## C. MORE ON TEACHERLESS TRAINING

Recall that our hypothesis is that the teacherless model automatically learns to fit the targets in the reverse order, since the path *from* the  $v_{\text{goal}}$  is unique. This is indeed what we find in Fig 4, where the accuracies of the later tokens become higher earlier.

Note though that this is a fairly difficult computation to implement. First, when the model predicts  $v_i$ , it must require the identify of  $v_{i+1}$ . However, this identity is not fed as input to the model, in the absence of the teacher. Thus the model must have computed  $v_{i+1}$  and crucially, stored that in one of its internal representations. Then, by induction, when predicting the first node  $v_1$ , the model must know the identity of *all* the other nodes in the path. In other words, the model must have (a) computed and (b) stored the whole solution in its hidden representations before it outputs the first token. This is a substantial type of lookahead that some of our models are able to achieve under teacherless training.

## D. MORE EXPERIMENTAL RESULTS

### D.1. Snowball Failure

To explicitly measure to what degree the model falls victim to the snowball effect, we train *GPT-Mini* on graphs of various path lengths  $l$ . In order to remove the failure stemming from the difficult first token, we teacher-force the model for the first token and then check how accurate the generations are for subsequent tokens. More concretely, we evaluate

$$\text{Acc}_{\text{sb}}(\text{LM}_\theta) = \mathbb{P}(\hat{r}_{1<} = r_{1<}) \quad (10)$$

$$\text{where } \mathbf{p}, \mathbf{r} \sim \mathcal{D}, \quad \hat{r}_{1<}^{\text{ag}} \approx \text{LM}_\theta(\cdot; \mathbf{p}, r_1)$$

If  $\text{Acc}_{\text{sb}}(\text{LM}_\theta) \approx 1$ , then *Failure 1* is not prominent in our task. If  $\text{Acc}_{\text{sb}}(\text{LM}_\theta) \ll 1$ , then clearly teacher-forcing is responsible for suppressing errors in generation, strongly hinting at the fact that *Failure 1* is at play. We display the results in Fig. 5 (left). We observe that the accuracy  $\text{Acc}_{\text{sb}}$  is barely affected even for graphs with very long paths  $L = 40$ .

As another metric, we proceed token by token during inference, and evaluate the probability of correctly predicting all tokens up to the current one. We report this for  $G_{2,40}$  in Fig. 5 (right). Similarly, while the success probability does decay for larger length (at an exponential rate), it remains very high due to the failure events being so unlikely. We thus conclude that *Failure 1* is not as prominent in this setting.

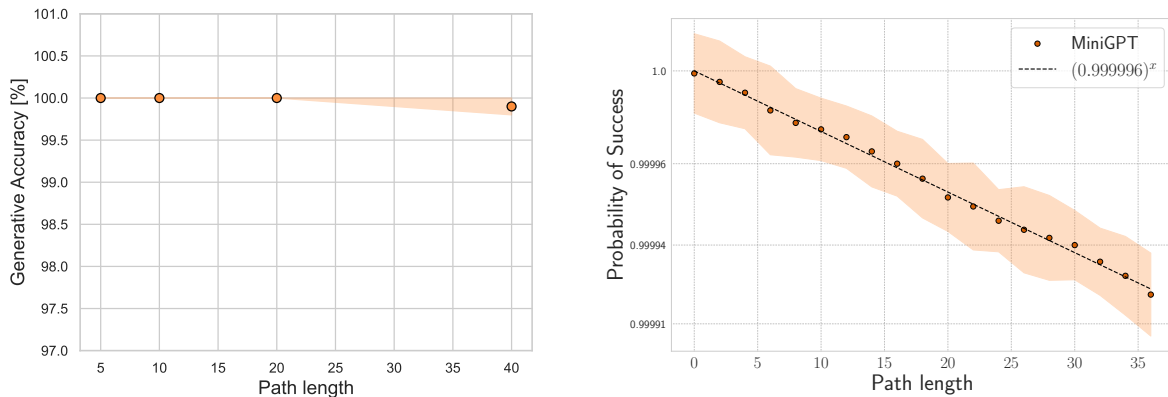


Figure 5. Accuracy of  $\text{LM}_\theta$  when conditioned on the first difficult token (left) for graphs of various length. Probability of correct prediction of  $\text{LM}_\theta$  as a function of current token position on  $G_{2,40}$ , as we walk towards the goal.

### D.2. Clever Hans Cheating Accuracies

In Table 1 we display the Clever Hans cheating accuracies  $\text{Acc}_{\text{cheat}}(\text{LM}_\theta)$ . We observe that in almost all cases, all the models achieve nearly perfect cheating accuracies. The only exception is the high-degree graph  $G_{20,5}$  where all models struggle to even fit the training data.

	$G_{2,5}$	$G_{2,20}$	$G_{5,5}$	$G_{10,5}$	$G_{20,5}$
GPT-MINI	99.7	100	100	99.8	0.0
GPT2-LARGE	99.8	99.7	100	99.8	0.0
MAMBA	97.6	98.3	99.5	95.9	0.0

Table 1. Evaluating Clever Hans cheating accuracies  $\text{Acc}_{\text{cheat}}(\text{LM}_\theta)$  (in percent %) for different types of graphs.

### D.3. More Detailed Accuracies

We report more detailed accuracy values per model in the following tables. We display standard accuracy  $\text{Acc}_{\text{ag}}(\text{LM}_\theta)$  in Table. 2, teacherless accuracy  $\text{Acc}_{\S}(\text{LM}_\theta)$  in Table. 3 and reverse accuracy  $\text{Acc}_{\text{rev}}(\text{LM}_\theta)$  in Table. 4. In general we observe that solving the task with standard next-token prediction is very tough and performance is limited to  $\frac{1}{d}$  where  $d$  is the degree of the graph  $G_{d,l}$ .

	$G_{2,5}$	$G_{2,20}$	$G_{5,5}$	$G_{10,5}$	$G_{20,5}$
GPT-MINI	49.8	49.1	19.1	8.1	0.0
GPT2-LARGE	48.9	49.2	19.4	10.3	3.5
MAMBA	48.5	48.7	20.2	9.3	0.0

Table 2. Autoregressive accuracies  $\text{Acc}_{\text{ag}}(\text{LM}_\theta)$  (in percent %) for different types of graphs.

Teacherless training on the other hand works very well with GPT2-Large, allowing it to solve most graph tasks perfectly. From-scratch models however also struggle to learn the task in this fashion (except for GPT-Mini on the simplest graph,  $G_{2,5}$ ).

	$G_{2,5}$	$G_{2,10}$	$G_{2,20}$	$G_{5,5}$	$G_{10,5}$	$G_{20,5}$
GPT-MINI	99.9	0.0	0.0	0.0	0.0	0.0
GPT2-L	99.9	98.8	0.0	99.0	97.8	0.0
MAMBA	0.0	0.0	0.0	0.0	0.0	0.0

Table 3. Autoregressive accuracy  $\text{Acc}_{\S}$  when using a teacherless response.

Finally, reversing the sequence significantly simplifies the problem for all the models, allowing near perfect accuracies across all graphs.

	$G_{2,5}$	$G_{2,20}$	$G_{5,5}$	$G_{10,5}$	$G_{20,5}$
GPT-MINI	99.7	99.8	100	99.8	0.0
GPT2-LARGE	99.9	99.9	99.6	99.8	99.9
MAMBA	98.5	96.2	99.1	99.5	0.0

Table 4. Autoregressive accuracy  $\text{Acc}_{\text{rev}}$  when reversing the response  $r$ .

## E. OTHER EXPERIMENTAL DETAILS

### E.1. Tokenization

We tokenize the graph in the following manner: (1) we first tokenize the randomly shuffled edge list as “ $|v_1 v_2|v_3 v_4|...$ ” where the first vertex in each edge is the one closest to  $v_{\text{start}}$ , (2) then append start and goal node as “ $/v_{\text{start}} v_{\text{goal}} =$ ” and (3) then append the full path repeating start and goal node, “ $v_{\text{start}} v_{i_1} \dots v_{i_{l-1}} v_{\text{goal}}$ ”. Note that (1) and (2) make up the prefix  $p$ , which the model does not learn to predict. Then, (3) is the target sequence that the model aims to learn. The vocabulary size is thus given by  $N + 3$ , where we add entries for the special tokens “|”, “/” and “=”. When using the pre-trained models GPT2 we use the tokenizer that was employed for pre-training, in this case the *Byte-Pair tokenizer* (Radford et al., 2019).

### E.2. Models

When training Transformer models from scratch, we use a small model consisting of  $n_{\text{layers}} = 12$  blocks with embedding dimension  $e_{\text{dim}} = 384$ ,  $n_{\text{heads}} = 6$  attention heads and MLP expansion factor  $e = 4$ , coined *GPT-Mini*. For pre-trained models, we consider GPT2-Large with  $n_{\text{layers}} = 36$ ,  $e_{\text{dim}} = 1280$ ,  $n_{\text{heads}} = 20$  and expansion factor  $e = 4$  (Radford et al., 2019). To further evaluate purely recurrent models, we perform experiments with the recent Mamba model (Gu & Dao, 2023). We train the Mamba models from scratch with 12 layers and embedding dimension 784. We train all the models with the *AdamW* optimizer (Loshchilov & Hutter, 2019). For models trained from scratch we use a learning rate of  $\eta = 0.0005$  while for pre-trained models we use a smaller one of  $\eta = 0.0001$ . In both cases we use weight decay of strength 0.01. Models from scratch are trained for up to 500 epochs in order to ensure convergence. Pre-trained models require less training time and we usually fit the training data perfectly after 10 epochs.

## F. MORE RELATED WORK

**Other arguments about next-token prediction** We survey related arguments of next-token prediction, orthogonal to our main discussion regarding planning. Allen-Zhu & Li (2023); Lv et al. (2023) report that language models that are trained on  $A \text{ equals } B$  are unable to infer  $B \text{ equals } A$ , which Allen-Zhu & Li (2023) suggest is due to autoregressive left-right training. Du et al. (2023b); Welleck et al. (2020) formalize the limitation that autoregressive models may potentially assign non-zero probability to infinite-length strings, thus leading to non-terminating inference. Li et al. (2024) provide a Transformer-specific analysis of how self-attention affects the optimization geometry of next-token prediction. Thrampoulidis (2024) provide an analysis of the implicit bias of optimization with next-token prediction for linear models.

**Other limitations of Transformers** Merrill & Sabharwal (2023a) identify limitations of the representative power of Transformer architecture when the arithmetic precision is logarithmic in the number of input tokens. Bender et al. (2021) criticize GPT-like language models as simply parroting out training data with minor stochasticity, while Arkoudas (2023) report that such models struggle with reasoning, even if not a stochastic parrot. Young & You (2023) study masked language (T5, BERT) models (not causally-trained) and argue there are inconsistencies in the probabilities that they assign. E.g., when conditioned on ‘white’, the probability of ‘rice’ may be higher ‘bread’ but the probability of ‘white bread’ and ‘white rice’ are the opposite. Artetxe et al. (2022) empirically analyze the effect of bidirectional attention and bidirectional supervision (as in masked language modeling) during pretraining on the ability of the model to do various things, including next-token prediction. Springer et al. (2024) argue that autoregressive Transformers compute sub-optimal embeddings that can be improved by repeating the input text twice.

Finally, we note that (Ranaldi & Zanzotto, 2023) use the term Clever Hans effect to denote how models can pick up spurious correlations between the position of a choice in a multiple-choice question, and the correctness of the answer. We note that the above correlation is inherent to the distribution, and independent of teacher-forcing. We distinguish this from the Clever Hans *cheating* which happens under the guidance of teacher-forcing.

**Going beyond next-token prediction-based training.** Finally, we note that some works (Gurnee et al., 2023; Meng et al., 2022; Pal et al., 2023) aim to recover future tokens that an already-trained model may predict based on the internal layers of the current token. Note that the success of this does not imply that the model necessarily plans well. This only means that it is possible to recover what the already-trained model wants to generate in the future (which may simply be a bad plan). Pfau et al. (2023) train a language model to predict in reverse with the orthogonal goal of finding prefixes that elicit certain behaviors.