

# MARKOVIAN TRANSFORMERS FOR INFORMATIVE LANGUAGE MODELING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Chain-of-Thought (CoT) reasoning holds great promise for explaining language model outputs, but recent studies have highlighted significant challenges in its practical application for interpretability. We propose to address this issue by making CoT causally essential to prediction through two key components: factoring next-token prediction through intermediate CoT text, and training CoT to predict future tokens independently of other context. This results in “Markovian” language models, where CoT serves as a fixed-size state for future token prediction. Our approach optimizes for “informativeness” – the improvement in next-token predictions using a trained CoT compared to a baseline. Using Proximal Policy Optimization (PPO) for arithmetic problems and policy gradient for GSM8K, we demonstrate effectiveness on both arithmetic problems with Mistral 7B and the GSM8K benchmark with Llama 3.1 8B, where the model learns to produce CoTs that are 33.20% more effective at predicting answers than the pre-trained baseline. The increased sensitivity of model performance to CoT perturbations provides strong evidence of CoT reliance. Furthermore, we show that CoTs trained for one model generalize to help other models predict answers, suggesting these CoTs capture reasoning patterns that transfer across different interpreters. This work advances the development of more interpretable language models, potentially enabling their extension to arbitrarily long contexts and enhancing AI reasoning capabilities across various domains.

## 1 INTRODUCTION

The rapid advancement of language models (LMs) has revolutionized the field of artificial intelligence, demonstrating remarkable capabilities in tackling complex cognitive tasks (Brown et al., 2020). However, it can be challenging to understand why an LM gave a particular answer (Burns et al., 2023; Gurnee & Tegmark, 2024; Lamparth & Reuel, 2023), which can be problematic in high-stakes scenarios (Rivera et al., 2024; Lamparth et al., 2024; Grabb et al., 2024). Interpretability techniques analyze the patterns and activations of a neural network in order to extract an explanation of the network’s behavior (Casper et al., 2023; Meng et al., 2022; Geva et al., 2022; Geiger et al., 2022; Wang et al., 2022; Nanda et al., 2023; Lamparth & Reuel, 2023). However, since language models already speak natural language and have been trained to be able to use their own internal representations, we could in principle simply ask the language model why it gave a particular answer to a question. Asking the language model to explain its reasoning in a “step-by-step” fashion before answering a question is known as Chain-of-Thought (CoT) (Wei et al., 2022; Nye et al., 2022) prompting.

However, there are concerns that CoT is an inadequate or *unfaithful* explanation for LM-generated text. For example, Turpin et al. (2023) show that biasing the LM to believe a particular answer via a supposedly irrelevant in-context feature such as multiple choice answer order will cause the CoT to rationalize that answer without mentioning the background feature. Also some LMs give the same answers to questions despite changes to the CoT reasoning in their context window (Lanham et al., 2023). While this has some benefits – the model can still answer correctly despite intermediate reasoning errors – it is also an indicator that the CoT does not fully capture the LM’s reasoning process. This raises a critical issue with using CoT as a tool for interpretability.

Our work introduces a novel perspective on this issue by focusing on *informativeness* rather than faithfulness, which would imply that the CoT reflects some underlying causal process in the model. Our key insight is to make the CoT text itself causally important in the model’s reasoning by training an LM to generate a minimal-length CoT such that the model can predict the answer given *only* that CoT. This approach ensures that the CoT is both *complete* (i.e., each necessary step is included) and *maximally fragile* (i.e., removing or changing the meaning of any step breaks the CoT and thus leads to a different result) by making the CoT itself a bottleneck in the flow of information that the language model uses to produce text.

We assume the LM receives a sequence of *observations* to predict – this could be a question-answer pair (length two sequence) or many adjacent segments of generic internet text. Our conceptual arguments rely on the size of each observation being larger than the CoT – otherwise the LM could put the answer immediately in the CoT. Though for pragmatic reasons we use short observations, the model does not learn the undesirable behavior of directly answering in the CoT due to the relative difficulty of predicting the answer without any CoT. The primary contributions of this work are:

1. We introduce a formal definition of informativeness which is used as an optimization target, providing a principled approach to generating meaningful CoT reasoning.
2. We demonstrate our training algorithm’s effectiveness by:
  - Training Mistral 7B V0.2 (Jiang et al., 2023) to solve 15-term addition problems
  - Training Llama 3.1 8B (Dubey et al., 2024) to achieve a 33.2% performance gain on the GSM8K (Cobbe et al., 2021) reasoning dataset
3. We verify the causal importance of generated CoTs through perturbation analysis, showing that training increases the sensitivity of model performance to CoT modifications.
4. We demonstrate that CoTs trained for one model transfer effectively to other models, suggesting they capture generalizable reasoning patterns rather than model-specific artifacts.

By making CoT causally important in the model’s reasoning, we aim to improve the interpretability and reliability of language models. This approach offers a novel perspective on understanding and steering LM behavior by leveraging the model’s own generated explanations, rather than relying solely on the analysis of its internal parameters.

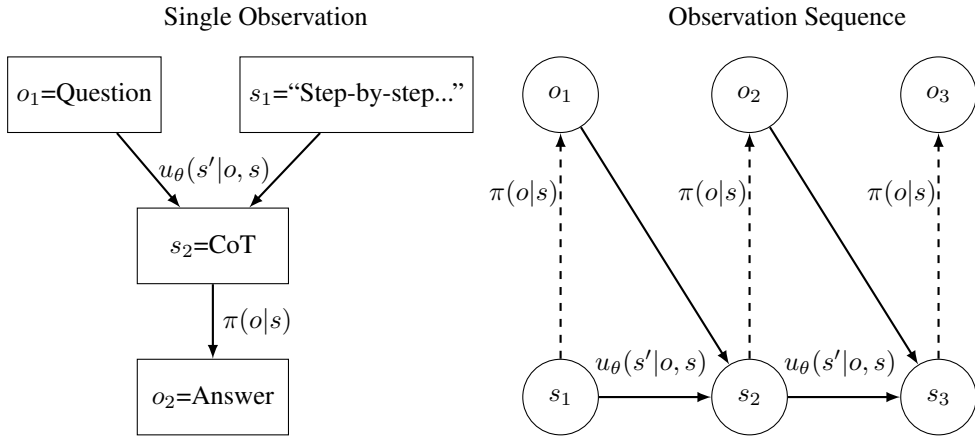


Figure 1: Refined illustration of the training method. Left: Single time-step process from Question to CoT to Answer. Right: Causal structure showing the generation of states from observations and previous states using the state update function  $u_\theta(s'|o, s)$ , and the prediction of observations from states using the policy  $\pi(o|s)$ . Observations are generated by the causal data distribution. In experiments, both  $u_\theta$  and  $\pi$  are Mistral 7B Instruct V0.2 or Llama 3.1 8B Instruct, but only the weights of  $u_\theta$  are updated during training. The state update  $u_\theta$  also involves concatenating the observation and state letting Mistral generate the next state’s worth of tokens.

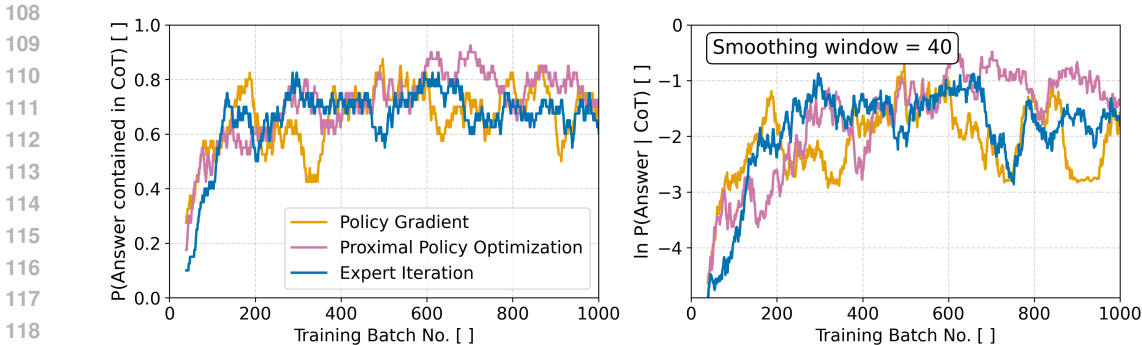


Figure 2: The log probability  $\ln \pi(ans | cot)$  of the answer  $ans$  given a CoT  $cot$ , where  $cot$  is sampled from the trained weights  $cot \sim u_{\theta}(cot | q, cot_{init})$  and  $cot'$  is sampled from the unmodified weights  $cot' \sim u(cot | q, cot_{init})$ . We train to produce CoTs which are sufficient to predict the correct answer even without the original question, enforcing a text bottleneck in the language model’s information flow, forcing the CoT to be causally load-bearing to production of the answer. This plot specifically depicts the training of Mistral 7B Instruct V0.2 on 15-term addition problems and their solutions. Because of high variance, we plot the point-wise maximum for each training technique across 4 separate training runs.

## 2 RELATED WORK

Prior work shows that CoT prompting improves language model reasoning capabilities (Wei et al., 2022; Nye et al., 2022). We *train* the model to produce a strong CoT, as opposed to prompting strategies as in Wei et al. (2022). Scratchpad (Nye et al., 2022) also trains the model to produce a CoT, but they supply correct CoTs during training, whereas our model has to discover useful CoTs for itself. Zelikman et al. (2024) also use RL to improve CoT reasoning, but they do not restrict the model’s attention to the previously generated CoT, making the CoT less of a standalone explanation. State space models also generate state to remember their history (Gu et al., 2021; 2022; Gu & Dao, 2023), but we use natural language instead of activation vectors for interpretability.

Lyu et al. (2023) improved faithfulness of language model reasoning by restricting the output to a particular formal language so that a deterministic solver could provide the rest of the answer, whereas we do not restrict to production of a formal language, because our future goal is to target general language modeling. Ranaldi & Freitas (2024) directly fine-tune a smaller model using CoT from a more capable model. In contrast, we do not require the existence of a more competent model to learn useful CoTs. Lanham et al. (2023) use robustness to reasoning perturbations as an indicator of unfaithfulness, which we adapt by replacing the variation in multiple choice accuracy with the variation in log probability assigned to the correct observation. Bentham et al. (2024) respond that robustness might simply be an indicator of accuracy, which we ameliorate by removing history from the context window. In order to address this concern more thoroughly, we would need to demonstrate the ability to further compress our CoTs.

## 3 MARKOVIAN LANGUAGE MODELS AND INFORMATIVENESS OF UPDATE FUNCTIONS

We would like a mathematical structure which describes the shape of a language model with a CoT bottleneck, so that we can derive an reinforcement learning algorithm with respect to that formalism. For this reason, we introduce the concept of Markovian Language Models and define a measure of informativeness for their update functions.

### 3.1 MARKOVIAN LANGUAGE MODELS

A regular auto-regressive LM can use its entire context when predicting the next token. In particular, when the LM takes a question, produces some reasoning and finishes with a final answer, the

162 generation of the final answer can still attend back to the question. Thus, there is no guarantee that  
 163 the CoT is causally linked to the answer tokens. In contrast, in a Markovian LM, the model is only  
 164 given access to limited state to make predictions.

165 Formally, we define a Markovian Language Model as a tuple  $M = (\mathcal{V}, \mathcal{S}, \pi, u, s_1)$ , where:

- 167 •  $\mathcal{V}$  is a finite vocabulary,
- 168 •  $\mathcal{S}$  is a set of states, representing CoT reasoning,
- 169 •  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{V}^k)$  is a state-conditional distribution over  $k$ -token sequences, where  $\Delta(\mathcal{V}^k)$   
 170 is the probability simplex over  $k$ -token sequences from  $\mathcal{V}$ ,
- 171 •  $u : \mathcal{V}^k \times \mathcal{S} \rightarrow \Delta(\mathcal{S})$  is a stochastic update function,
- 172 •  $s_1 \in \mathcal{S}$  is the initial state.

173 The MLM operates sequentially: given a current state  $s_t \in \mathcal{S}$  and observation  $x_t \in \mathcal{V}^k$ , it produces  
 174 a probability distribution  $\pi(s_t)$  over the next  $k$ -token sequence, and stochastically updates its state  
 175 to  $s_{t+1} \sim u(s_t, x_t)$ .

### 179 3.2 DATA-GENERATING DISTRIBUTION

180 Let  $P$  be the true data-generating distribution over sequences of length  $T$ . We can sample from this  
 181 distribution using:

$$182 \quad x_t \sim P(x_t | x_{<t}) \quad \text{for } t = 1 \text{ to } T \quad (1)$$

183 where  $x_{<t}$  denotes all observations before time  $t$ .

### 188 3.3 PARAMETERIZED UPDATE FUNCTION

189 We consider a parameterized update function  $u_\theta$ , where  $\theta$  represents the parameters to be optimized.  
 190 We compare this to a baseline update function  $u'$ , which uses the original set of weights before fine-  
 191 tuning. Both  $u_\theta$  and  $u'$  operate in conjunction with the same prediction function  $\pi$ , which also uses  
 192 the original set of weights.

### 194 3.4 INFORMATIVENESS OF UPDATE FUNCTIONS

195 We define the informativeness of the update function  $u$  relative to a baseline update function  $u'$  as:

$$196 \quad I(u, u', P) = \mathbb{E}_{\tau \sim P, u, u'} [R(\tau)] \quad (2)$$

197 where  $\tau = (x_1, s_1, s'_1, \dots, x_T, s_T, s'_T)$  is a trajectory, with:

- 198 •  $x_t \sim P(x_t | x_{<t})$
- 199 •  $s_{t+1} \sim u(s_{t+1} | x_t, s_t)$
- 200 •  $s'_{t+1} \sim u'(s'_{t+1} | x_t, s'_t)$

201 The reward  $R(\tau)$  for a trajectory is defined as:

$$202 \quad R(\tau) = \sum_{t=1}^T [\ln \pi(x_t | s_t) - \ln \pi(x_t | s'_t)] \quad (3)$$

203 Now, let's consider optimizing this informativeness using policy gradient methods. We parameterize  
 204  $u$  by some weights  $\theta$ , giving us  $u_\theta$ . The objective function is:

$$205 \quad J(\theta) = I(u_\theta, u', P) \quad (4)$$

The gradient of this objective with respect to  $\theta$  is:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim P, u_{\theta}, u'} \left[ R(\tau) \sum_{t=1}^{T-1} \nabla_{\theta} \ln u_{\theta}(s_{t+1} | x_t, s_t) \right] \quad (5)$$

In practice, we estimate this gradient using Monte Carlo sampling:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N R(\tau^{(i)}) \sum_{t=1}^{T-1} \nabla_{\theta} \ln u_{\theta}(s_{t+1}^{(i)} | x_t^{(i)}, s_t^{(i)}) \quad (6)$$

where  $\{\tau^{(i)} = (x_1^{(i)}, s_1^{(i)}, s_1^{\prime(i)}, \dots, x_T^{(i)}, s_T^{(i)}, s_T^{\prime(i)})\}_{i=1}^N$  are sampled trajectories.

This procedure improves the update function  $u_{\theta}$  to generate more informative CoT reasoning, leading to better predictions of future observations.

## 4 METHODS

### 4.1 MARKOVIAN LANGUAGE MODEL FOR QUESTION-ANSWER PAIRS AND OPTIMIZATION

We define a specialized Markovian Language Model (MLM) for question-answer pairs as a 5-tuple  $M = (\mathcal{V}, \mathcal{S}, \pi, u, s_1)$ , where:

- $\mathcal{V}$  is the vocabulary of tokens.
- $\mathcal{S} = \mathcal{V}^{\ell}$  is the set of all possible CoT sequences of length  $\ell$ .
- $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{V}^{\ell})$  is the prediction function.
- $u : \mathcal{S} \times \mathcal{V}^{\ell} \rightarrow \Delta(\mathcal{S})$  is the update function.
- $s_1 = \text{cot}_{\text{init}} \in \mathcal{S}$  is the initial state, where  $\text{cot}_{\text{init}}$  is a fixed initial prompt.

Let  $\ell_q = \ell_a = \ell$  be the length of an observation (question or answer). We implement the MLM specification using a language model  $\mathcal{L} : \mathcal{V}^* \rightarrow \Delta(\mathcal{V})$ , where  $\mathcal{L}(s)$  gives the probability distribution over the next token given the sequence  $s$ . We denote the  $i$ -th tokens of the CoT and answer as  $\text{cot}_i$  and  $\text{ans}_i$ , respectively.

The model operates as follows:

1. Update function  $u$ :

$$\ln u(s_2 = \text{cot} | o_1 = q, s_1 = \text{cot}_{\text{init}}) = \sum_{i=1}^{\ell} \ln \mathcal{L}(\text{concat}(q, \text{cot}_{\text{init}}, \text{cot}_{<i}))[\text{cot}_i] \quad (7)$$

We implement  $\ln u$  by concatenating the question with  $\text{cot}_{\text{init}}$  and summing the log probability of each token conditioned on the previous tokens and the prefix.

2. Prediction function  $\pi$ :

$$\ln \pi(o_2 = \text{ans} | s_2 = \text{cot}) = \sum_{i=1}^{\ell} \ln \mathcal{L}(\text{concat}(\text{cot}, \text{ans}_{<i}))[\text{ans}_i] \quad (8)$$

### 4.2 THRESHOLD-BASED EXPERT ITERATION, POLICY GRADIENT, AND PROXIMAL POLICY OPTIMIZATION

We explore three RL techniques to optimize the language model for informative CoT production: Threshold-based Expert Iteration, Policy Gradient, and Proximal Policy Optimization. All three implementations use a form of importance sampling to focus updates on more informative CoTs. All three implementations are concisely contained within a single Python file, which we have made freely available.

#### 4.2.1 THRESHOLD-BASED EXPERT ITERATION (TEI)

Threshold-based Expert Iteration consists of the following steps:

1. Sample a CoT from a trained and untrained model (cot and cot')
2. Estimate informativeness as  $I(ans, cot, cot') = \pi(ans|cot) - \pi(ans|cot')$
3. If  $I$  is at least one standard deviation above the historical average:
  - Calculate the gradient of the log probability of having produced that CoT:  $\nabla_{\theta} \ln u_{\theta}(cot|q, cot_{init})$
  - Gradient ascend

**Limitation:** This technique potentially discards valuable information, as we might prefer to update more strongly towards CoTs that produce very high rewards.

#### 4.2.2 POLICY GRADIENT (PG)

Policy Gradient (with threshold-based sample selection) consists of the following steps:

1. Sample a CoT from a trained and untrained model (cot and cot')
2. Estimate informativeness as  $I(ans, cot, cot') = \pi(ans|cot) - \pi(ans|cot')$
3. If  $I$  is at least one standard deviation above the historical average:
  - Calculate the gradient of the log probability of having produced that CoT:  $\nabla_{\theta} \ln u_{\theta}(cot|q, cot_{init})$
  - Multiply the gradient by  $I$  and then ascend

**Advantage:** Utilizes more information than TEI

**Disadvantage:** Increased instability, which can be problematic given pre-trained initial weights

#### 4.2.3 PROXIMAL POLICY OPTIMIZATION (PPO)

For each CoT, PPO performs the following:

1. Calculate the probability ratio:  $r = \frac{u_{\theta}(cot|q, cot_{init})}{u'(cot|q, cot_{init})}$
2. Compute the clipped objective:

$$\text{obj} = \min(r \cdot I, \text{clip}(r, 1 - \epsilon, 1 + \epsilon) \cdot I)$$

where:

- $I = \text{Informativeness}(ans, cot, cot')$
- $\text{clip}(x, y, z) = \begin{cases} y & \text{if } x < y \\ z & \text{if } x > z \\ x & \text{otherwise} \end{cases}$
- $\epsilon = 0.2$

3. Back-propagate to increase obj

**Key Idea:** Remove the incentive to create CoTs for which the trained and untrained state update functions disagree too much.

**Implementation Details:**

- We use threshold-based sample selection here as well
- Subtract the historical average informativeness over unfiltered CoTs from the current informativeness as a baseline

## 5 EXPERIMENTS

### 5.1 MULTI-STEP ADDITION

We generate random addition problems, where each problem consists of 15 terms and each term is a uniform random natural number less than 100. We fine-tune Mistral 7B Instruct V0.2 to produce CoT tokens such that a frozen copy of the pre-trained language model can predict the correct answer given that CoT, for each training technique in Methods. We plot the mean negative log likelihood over the answer tokens as a function of training batch in Fig. 2. Note that this is both training and testing loss, since we are always generating fresh arithmetic problems. PPO, our preferred training method for arithmetic, can mention the correct answer in up to 90% of CoTs and achieve an average natural log probability of around -0.7.

Since the Mistral tokenizer allocates a separate token for each digit, a natural log probability of -0.7 corresponds to an actual probability of  $e^{-0.7} \approx 0.4966$ , or 50% chance of picking the correct next token on average. A 90% likelihood saying the answer verbatim in the CoT and a 50% of guessing each digit incorrectly may seem contradictory – however this discrepancy is due to the predictor model’s uncertainty around prompt formatting, and specifically about what tokens should come after “Answer:”. So it is distributing probability mass over the entire vocabulary including non-numerical tokens, since we are only training CoT production  $u_{\theta}(s'|o, s)$ , as opposed to training the predictor model  $\pi(o|s)$ .

### 5.2 GSM8K

To test our method on more complex reasoning tasks, we train Llama-3.1-8B-Instruct to produce CoT over the GSM8K dataset. Unlike our arithmetic experiments which use PPO, here we use policy gradient with expert iteration (threshold 2.2 standard deviations), along with 150 CoT tokens and a KL penalty of 0.1. Figure 3 shows the results across three training runs. The left plot demonstrates substantial improvements in the log probability that an untrained Llama assigns to the correct answer given the trained CoT. The right plot shows the proportion of CoTs that contain the answer verbatim, indicating the model learns to consistently encode the correct answer in its reasoning.

Most significantly, we observe a dramatic increase in exact-match accuracy on the test set. Starting from a baseline of 35.94% at batch 0, our best performing run achieves 69.14% accuracy (n=1), representing a 33.2% absolute improvement. The other two runs achieve 58.23% and 62.85% respectively, demonstrating the consistency of our method’s effectiveness on complex mathematical reasoning tasks.

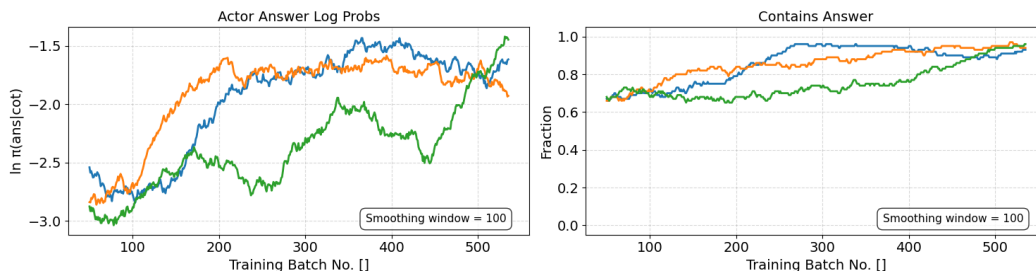


Figure 3: GSM8K performance metrics over three separate training runs of Llama-3.1-8B-Instruct. The left plot shows the log probability that an untrained Llama assigns to the correct answer given the trained CoT —  $\ln \pi(ans|cot)$ , and the right plot shows the proportion of CoTs in a batch which contain the answer verbatim. We used a smoothing window of size 100, explaining the multiplicity of possible y-values for “Contains Answer”.

### 5.3 WIKIPEDIA

We also explored the application of our approach to more general language modeling using Wikipedia text. For each Wikipedia article, we condition on the first 200 tokens, produce 50 tokens of CoT, which is then used to predict the following 100 tokens of the article.

Our prompt template is:

“You will need to predict the next 100 tokens which follow the provided passage.  
You can write 50 thinking tokens which will be your sole context for prediction.  
Feel free to be creative in your thinking strategy!\n\nOpening text:”

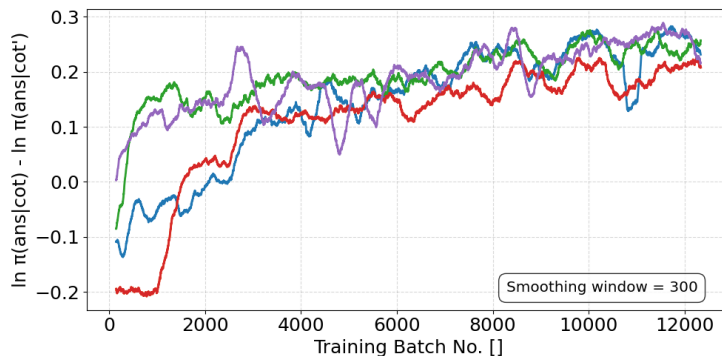


Figure 4: Four independent training runs, showing the difference in log probabilities of the answer  $ans$  given a trained CoT  $cot$  and the default  $cot'$  that a pre-trained model would produce —  $\ln \pi(ans | cot) - \ln \pi(ans | cot')$ . The model is Llama 8B and the task is to produce text which helps predict the next 100 tokens in a Wikipedia article (Foundation, 2024).

Results showed modest improvements in next-token prediction accuracy from 8.2% to 10.5% (Figure 4). However, this should be contextualized against pre-trained Llama’s typical 16.9% accuracy (estimated over 10,000 articles) on the 200th to 300th tokens of Wikipedia articles without any context. The lower baseline performance (8.2%) appears to be an artifact of our prompting setup.

Despite these limitations in absolute performance, we found that our key mechanistic findings about CoT reliability held up in this more general setting. In Wikipedia Perturbations and Cross-Model Generalization, Figure 7 demonstrates that perturbations to the CoT meaningfully impact performance, with the trained model showing greater sensitivity to perturbations than the baseline model. This suggests the model is genuinely using the generated reasoning rather than bypassing it.

#### 5.4 MEASURING FRAGILITY OF CoT

Expanding upon Lanham et al. (2023), we measure the fragility of the CoT reasoning by applying three perturbations to the model-generated reasoning and evaluate how this affects the next-token-prediction loss of the correct answer to the original question. Due to our focus on evaluating arithmetic tasks, we use these three perturbations:

- Truncating a fraction of the CoT reasoning from the end
- Flipping any number (digit) with a probability in the CoT reasoning and replacing it with another random number between 0 and 9
- Swapping a fraction of characters with random characters in the CoT reasoning. The selection is limited to numbers from 0 to 9, letters from the English alphabet, and simple arithmetic symbols (e.g., “+” and “-”)

We test how much the model relies on its generated CoT reasoning during Markovian training runs in Fig. 5. The y-axis depicts the log probability of the answer given CoT, normalized so that  $y = 0$  corresponds to the log probability of the answer given the unperturbed CoT. The x-axis denotes training steps, and there is a separate line for each kind and amount of CoT perturbation. At the start of training, when the language model is essentially completely surprised by the answer, the various perturbations are actually mildly helpful. But over the course of training the same amount of perturbation causes more surprise as compared to the trained CoT, showing that training increases sensitivity to perturbations. Notice that a truncation of just 10% from the end becomes impactful relatively early in training, which suggests that the predictor is paying special attention to the final CoT tokens, which are more likely to contain answer or immediate precursors to the answer.



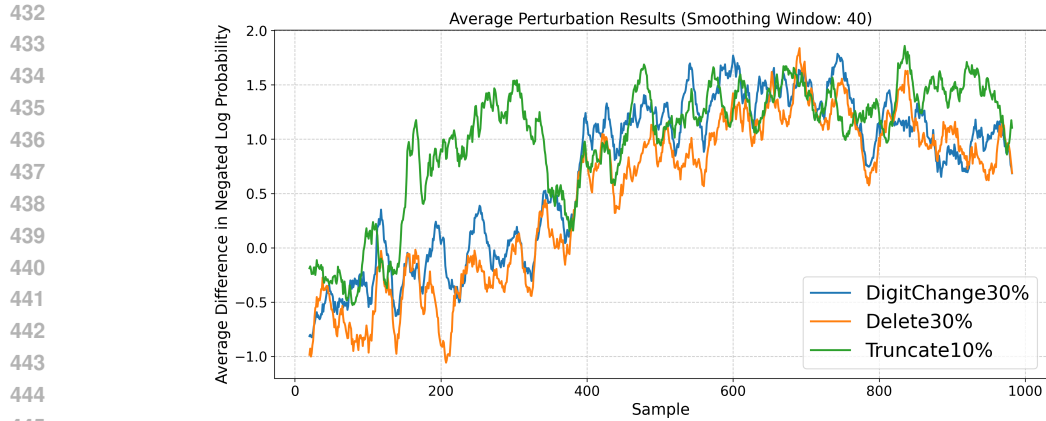


Figure 5: Comparison of different perturbation effects on CoT reasoning, using the arithmetic dataset and Mistral 7B. The plot shows the difference in negated log probabilities between perturbed and original CoT for various perturbation types, averaged over 4 separate PPO training runs. Higher values indicate worse performance compared to the original. Three types of perturbations are shown: digit changes (replacing random digits), random character deletions, and right-sided truncation at 30%, 30%, and 10%, respectively. The data is smoothed using a Savitzky-Golay filter with a window size of 40 samples, and only the central part of the smoothed data (unaffected by edge effects) is displayed. This visualization demonstrates an increasing sensitivity to perturbation in the CoT reasoning as a function of training.

## 5.5 INTERPRETABILITY OF CoT GENERATIONS

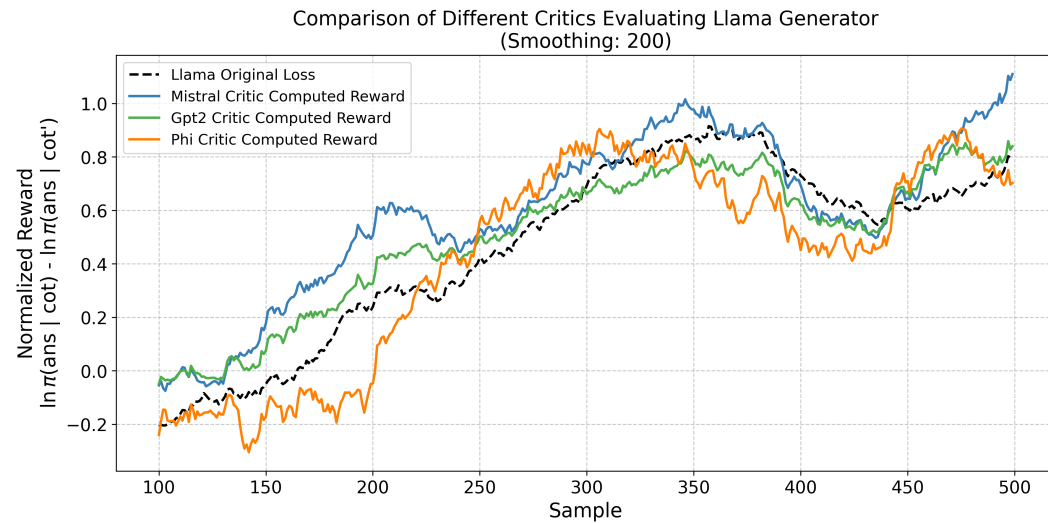


Figure 6: Comparison of the log probabilities between the Llama 8B model and several other language models, averaged across 3 separate training runs on the GSM8K dataset. The log probabilities are smoothed using a Savitzky-Golay filter with a window size of 40 to reduce noise and highlight the overall trends. The plot shows that improvements in CoT from Llama’s perspective also correspond to improvements in CoT from the perspective of Mistral, GPT2 (Radford et al., 2019), and Phi 3.5 Mini Instruct (Abdin et al., 2024), which vary greatly in performance and characteristics, lending evidence that humans may also understand Mistral’s trained CoT.

To probe how well the reasoning generalizes, we plot the log probabilities that various models ascribe to the answer given trained Llama’s CoT in Fig. 6. In both plots the log probabilities increase simultaneously, demonstrating that Llama is learning to produce generic CoTs which do not over-fit to the peculiarities of a Llama answer-predictor.

486 This cross-model transferability connects to fundamental questions about interpretability. When  
487 we say a CoT is "interpretable", we must ask "interpretable to whom?" — just as induction heads  
488 are interpretable to ML researchers but not to most humans, different CoTs might be naturally in-  
489 terpretable to different readers or models. Our experimental design engages with this relativity in  
490 two ways: First, we include GPT2, a significantly smaller model, as one of our evaluators. Sec-  
491 ond, we test across three distinct model families (Llama, Mistral, and GPT2), preventing the trained  
492 model from exploiting architecture-specific patterns. The fact that the trained CoTs transfer effec-  
493 tively across this diverse set of evaluators, including a much smaller model, suggests they capture  
494 reasoning patterns that are interpretable across a broad range of computational architectures.

## 495 496 6 DISCUSSION AND LIMITATIONS 497

498 Our experiments show that it is possible to learn informative and interpretable CoT reasoning via  
499 RL on an LM using Markovian training. However, we find that training is unstable, and we present  
500 various techniques to prevent the LM from losing its strong language modeling prior.

501 A weakness in our interpretability argument is that for GSM8K and addition we use more CoT to-  
502 kens than answer tokens, so in principle the LM could learn to put the answer in the CoT directly.  
503 However, this did not affect our particular experiments because Mistral struggles to learn to add  
504 fifteen terms without intermediate reasoning, and similarly for Llama with GSM8K answers. Addi-  
505 tionally, our interpretability technique is currently only verified in myopic question-answer datasets,  
506 as opposed to multi-turn trajectories where trained CoTs might provide a lens into longer term fu-  
507 ture behavior. Lastly, we only train Mistral to produce CoT that it can interpret (use to predict  
508 observations), but in principle future work could optimize CoT for human interpretability directly.

509 Markovian training is essentially language modeling – predicting future tokens from previous tokens  
510 – but with an intermediate "action" to produce the LM's own memory. In this sense, this training  
511 paradigm blurs the line between RL and unsupervised learning. But since it comes at the cost of  
512 adding expensive serial token generation steps in an otherwise highly parallelizable unsupervised  
513 training regime, it would need to have a high payoff in terms of interpretability or perplexity in  
514 order to be feasible. But as it stands, we have only tested the technique on question-answer pairs,  
515 and are preliminary results are limited in the context of more general language modeling. In future  
516 work, we hope to stably optimize this objective in more general contexts.

## 517 518 7 ETHICS STATEMENT 519

520 Reinforcement learning techniques improve a policy with respect to an arbitrary reward function.  
521 But it can be difficult to mathematically specify nuanced human preferences about the policy. Both  
522 reinforcement learning from human feedback and Constitutional AI help people specify and opti-  
523 mize the properties they would like the AI to have. This increase in controllability makes the AI  
524 more of an extension of human intention, for better or for worse. The approach of this paper is much  
525 more targeted – we use RL to specifically increase an agent foresight – its ability to predict its future  
526 observations.

527 On its face, this seems like it might be just as dependent on human intentions as RLHF and Con-  
528 stitutional AI – if people are more knowledgeable, maybe they could use that extra knowledge to  
529 deceive others, for instance. However, better foresight may also give rise to better values, where  
530 values are opinions about how to act such that the collective system can attain better foresight.

## 531 532 8 REPRODUCIBILITY STATEMENT 533

534 To ensure reproducibility, we provide comprehensive supplementary materials including all source  
535 code, training and evaluation scripts, and detailed instructions in the README. The main training  
536 loop (`src/train.py`) supports (i) EI, PG, and PPO methods and (ii) GSM8K, arithmetic, and  
537 Wikipedia datasets. We measure fragility of CoT via `src/perturbation_analysis.py` and  
538 we estimate interpretability of CoT generations via `src/evaluate_cross_model.py`. The  
539 `results/Official` directory contains plots, full training logs, and perturbation evaluation logs  
from our experiments.

We use the public GSM8K and HuggingFace Wikipedia datasets, and we use the public Llama 3.1 8B Instruct, Mistral 7B Inst V0.2, Phi 3.5 Mini-Instruct, and GPT2 models. All hyperparameters are specified in the scripts defaults and in the paper, and environment setup instructions are in the README.

With these materials, researchers should be able to reproduce our work, including the performance boost on GSM8K and the perturbation analysis results demonstrating CoT reliance.

## REFERENCES

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- Oliver Bentham, Nathan Stringham, and Ana Marasović. Chain-of-thought unfaithfulness as disguised accuracy, 2024. URL <https://arxiv.org/abs/2402.14897>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ETKgubY0hcs>.
- Stephen Casper, Tilman Rauker, Anson Ho, and Dylan Hadfield-Menell. Sok: Toward transparent AI: A survey on interpreting the inner structures of deep neural networks. In *First IEEE Conference on Secure and Trustworthy Machine Learning*, 2023.
- Paul Christiano, Ajeya Cotra, and Mark Xu. Eliciting latent knowledge: How to tell if your eyes deceive you, December 2021. URL [https://docs.google.com/document/d/1WwsnJQstPq91\\_Yh-Ch2XRL8H\\_EpsnjrCldwzXR37PC8/edit](https://docs.google.com/document/d/1WwsnJQstPq91_Yh-Ch2XRL8H_EpsnjrCldwzXR37PC8/edit).
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John

594 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,  
595 2021. URL <https://arxiv.org/abs/2110.14168>.  
596

597 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
598 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony  
599 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,  
600 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,  
601 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris  
602 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,  
603 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny  
604 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,  
605 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael  
606 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-  
607 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah  
608 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan  
609 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-  
610 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy  
611 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,  
612 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-  
613 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,  
614 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearry, Laurens van der  
615 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,  
616 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-  
617 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,  
618 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,  
619 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur  
620 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-  
621 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,  
622 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,  
623 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-  
624 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,  
625 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,  
626 Sharath Rparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,  
627 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney  
628 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,  
629 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,  
630 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-  
631 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,  
632 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,  
633 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre  
634 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha  
635 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay  
636 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda  
637 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew  
638 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita  
639 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh  
640 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De  
641 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-  
642 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina  
643 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,  
644 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,  
645 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana  
646 Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,  
647 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-  
caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco  
Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella  
Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory  
Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,  
Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-

- 648 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,  
649 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer  
650 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe  
651 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie  
652 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun  
653 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal  
654 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,  
655 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian  
656 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,  
657 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-  
658 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel  
659 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-  
660 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-  
661 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,  
662 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,  
663 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,  
664 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,  
665 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,  
666 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,  
667 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-  
668 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-  
669 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang  
670 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen  
671 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,  
672 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,  
673 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-  
674 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,  
675 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu  
676 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-  
677 stable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,  
678 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,  
679 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef  
680 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.  
681 URL <https://arxiv.org/abs/2407.21783>.
- 682 Wikimedia Foundation. Wikipedia, 2024. URL <https://dumps.wikimedia.org>.
- 683 Atticus Geiger, Zhengxuan Wu, Hanson Lu, Josh Rozner, Elisa Kreiss, Thomas Icard, Noah Good-  
684 man, and Christopher Potts. Inducing causal structure for interpretable neural networks. In *Inter-  
685 national Conference on Machine Learning (ICML)*, pp. 7324–7338. PMLR, 2022.
- 686 Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward lay-  
687 ers build predictions by promoting concepts in the vocabulary space. *Proceedings of the 2022  
688 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 30–45, 2022.
- 689 Declan Grabb, Max Lamparth, and Nina Vasan. Risks from language models for automated men-  
690 tal healthcare: Ethics and structure for implementation. *medRxiv*, 2024. doi: 10.1101/2024.  
691 04.07.24305462. URL [https://www.medrxiv.org/content/early/2024/04/08/  
692 2024.04.07.24305462](https://www.medrxiv.org/content/early/2024/04/08/2024.04.07.24305462).
- 693 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2023.  
694 URL <https://arxiv.org/abs/2312.00752>.
- 695 Albert Gu, Isys Johnson, Karan Goel, Khaled Kamal Saab, Tri Dao, Atri Rudra, and Christopher  
696 Re. Combining recurrent, convolutional, and continuous-time models with linear state space  
697 layers. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in  
698 Neural Information Processing Systems*, 2021. URL [https://openreview.net/forum?  
699 id=yWd42CWN3c](https://openreview.net/forum?id=yWd42CWN3c).
- 700 Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured  
701 state spaces. In *International Conference on Learning Representations*, 2022. URL [https://  
702 openreview.net/forum?id=uYLFoz1v1AC](https://openreview.net/forum?id=uYLFoz1v1AC).

- 702 Wes Gurnee and Max Tegmark. Language models represent space and time. In *The Twelfth Interna-*  
703 *tional Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=jE8xbmvFin)  
704 [forum?id=jE8xbmvFin](https://openreview.net/forum?id=jE8xbmvFin).  
705
- 706 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
707 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*  
708 *ference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=nZeVKeeFYf9)  
709 [id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 710 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-  
711 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
712 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,  
713 Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>. Version 1.  
714
- 715 Nitish Joshi, Javier Rando, Abulhair Saparov, Najoung Kim, and He He. Personas as a way  
716 to model truthfulness in language models, 2024. URL [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2310.18168)  
717 [arXiv.2310.18168](https://doi.org/10.48550/arXiv.2310.18168). arXiv:2310.18168v5 [cs.CL].  
718
- 719 Max Lamparth and Anka Reuel. Analyzing and editing inner mechanisms of backdoored language  
720 models, 2023. URL <https://arxiv.org/abs/2302.12461>.
- 721 Max Lamparth, Anthony Corso, Jacob Ganz, Oriana Skylar Mastro, Jacquelyn Schneider, and  
722 Harold Trinkunas. Human vs. machine: Language models and wargames, 2024. URL <https://arxiv.org/abs/2403.03407>.  
723
- 724 Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Her-  
725 nandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamil  Luko iut , Karina  
726 Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson,  
727 Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Tim-  
728 othy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan  
729 Brauner, Samuel R. Bowman, and Ethan Perez. Measuring faithfulness in chain-of-thought rea-  
730 soning, 2023. URL <https://arxiv.org/abs/2307.13702>.  
731
- 732 Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human  
733 falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational*  
734 *Linguistics*, 2022. URL <https://arxiv.org/abs/2109.07958>. ACL 2022 (main con-  
735 ference).
- 736 Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki,  
737 and Chris Callison-Burch. Faithful chain-of-thought reasoning, 2023. URL [https://arxiv.](https://arxiv.org/abs/2301.13379)  
738 [org/abs/2301.13379](https://arxiv.org/abs/2301.13379).
- 739 Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual  
740 associations in GPT. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.  
741
- 742 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures  
743 for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learn-*  
744 *ing Representations*, 2023. URL <https://openreview.net/forum?id=9XFSbDPmdW>.
- 745 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David  
746 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Au-  
747 gustus Odena. Show your work: Scratchpads for intermediate computation with language models,  
748 2022. URL <https://openreview.net/forum?id=iedYJm92o0a>.
- 749 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
750 models are unsupervised multitask learners. 2019.  
751
- 752 Leonardo Ranaldi and Andre Freitas. Aligning large and small language models via chain-of-  
753 thought reasoning. In Yvette Graham and Matthew Purver (eds.), *Proceedings of the 18th Con-*  
754 *ference of the European Chapter of the Association for Computational Linguistics (Volume 1:*  
755 *Long Papers)*, pp. 1812–1827, St. Julian’s, Malta, March 2024. Association for Computational  
Linguistics. URL <https://aclanthology.org/2024.eacl-long.109>.

- 756 Juan-Pablo Rivera, Gabriel Mukobi, Anka Reuel, Max Lamparth, Chandler Smith, and Jacquelyn  
757 Schneider. Escalation risks from language models in military and diplomatic decision-making,  
758 2024. URL <https://arxiv.org/abs/2401.03408>.
- 759 D. Silver, A. Huang, C. Maddison, et al. Mastering the game of go with deep neural networks and  
760 tree search. *Nature*, 529:484–489, 2016.
- 762 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,  
763 Marc Lanctot, Laurent Sifre, Dharrshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Si-  
764 monyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general rein-  
765 forcement learning algorithm, 2017. URL [https://doi.org/10.48550/arXiv.1712.](https://doi.org/10.48550/arXiv.1712.01815)  
766 01815. arXiv:1712.01815 [cs.AI].
- 767 Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea Finn. Fine-tuning  
768 language models for factuality, 2023. URL [https://doi.org/10.48550/arXiv.2311.](https://doi.org/10.48550/arXiv.2311.08401)  
769 08401. arXiv:2311.08401 [cs.CL].
- 771 Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don’t always  
772 say what they think: Unfaithful explanations in chain-of-thought prompting. In *Thirty-seventh*  
773 *Conference on Neural Information Processing Systems*, 2023. URL [https://openreview.](https://openreview.net/forum?id=bzs4uPLXvi)  
774 [net/forum?id=bzs4uPLXvi](https://openreview.net/forum?id=bzs4uPLXvi).
- 775 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-  
776 pretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh*  
777 *International Conference on Learning Representations*, 2022.
- 778 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,  
779 Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language  
780 models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Ad-*  
781 *vances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=_VjQlMeSB_J)  
782 [forum?id=\\_VjQlMeSB\\_J](https://openreview.net/forum?id=_VjQlMeSB_J).
- 783 Zichao Yang, Phil Blunsom, Chris Dyer, and Wang Ling. Reference-aware language models. In  
784 Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings of the 2017 Conference*  
785 *on Empirical Methods in Natural Language Processing*, pp. 1850–1859, Copenhagen, Denmark,  
786 September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1197. URL  
787 <https://aclanthology.org/D17-1197>.
- 788 Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D. Goodman.  
789 Quiet-star: Language models can teach themselves to think before speaking, 2024. URL [https:](https://arxiv.org/abs/2403.09629)  
790 [/arxiv.org/abs/2403.09629](https://arxiv.org/abs/2403.09629).

## 793 A STABILITY-ENHANCING TRAINING TECHNIQUES

794  
795 Fine-tuning a pre-trained language model with a strong linguistic prior requires careful considera-  
796 tion to avoid irrecoverable weight updates that could push the model out of the language modeling  
797 loss basin. In addition to the PPO-clip objective mentioned in Sec. 4.2.3, we implemented several  
798 techniques to enhance training stability across different objective functions:  
799

### 800 1. Low-Rank Adaptation (LoRA):

- 801 • Freeze all weights except for a set of LoRA weights (Hu et al., 2022)
- 802 • Use rank 8 with  $\alpha = 16$

### 803 2. Gradient Clipping:

- 804 • If the  $L_2$  norm of the gradient update vector exceeds 1, normalize the vector

### 805 3. Gradient Accumulation: (Only for arithmetic)

- 806 • Set batch size to 6 to optimize H100 GPU memory usage
- 807 • Perform 8 gradient accumulation steps between weight updates

### 808 4. Average Reward Baseline:

- For PPO: Subtract the previous average of rewards from the current reward
- Found to be as beneficial as a value head, with less hyper-parameter tuning required

5. Selection of  $cot_{init}$ :

- Choose  $cot_{init}$  to bias CoT search in a productive direction
- For arithmetic we used “You will be given an arithmetic problem, which you have [cot length] tokens to work through step-by-step. Question:”
- For GSM8K we used “You will be given a reasoning problem, which you have [cot length] tokens to work through step-by-step. Question:”

B WIKIPEDIA PERTURBATIONS AND CROSS-MODEL GENERALIZATION

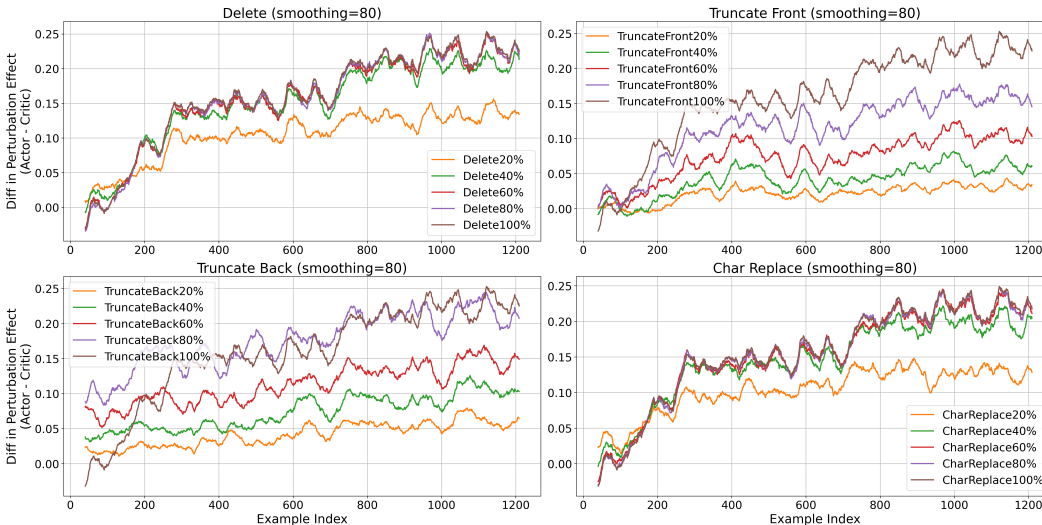


Figure 7: Impact of various perturbations on Wikipedia CoT effectiveness over the course of training. Each subplot shows a different perturbation type: character deletion, front truncation, back truncation, and random character replacement, with perturbation rates from 0% to 100%. For a perturbation function  $pert$ , letting  $\pi(ans|cot)$  denote the log probability of the answer given a CoT, we plot  $[\pi(ans|cot) - \pi(ans|pert(cot))] - [\pi(ans|cot') - \pi(ans|pert(cot'))]$ , where  $cot'$  is the default CoT from the pre-trained model. Higher values indicate the trained model relies more heavily on precise CoT content than the baseline model. When  $pert$  is a 100% perturbation rate (effectively a constant function  $k$ ), this reduces to  $[\pi(ans|cot) - k] - [\pi(ans|cot') - k] = \pi(ans|cot) - \pi(ans|cot') = I(ans, cot, cot')$ , explaining why these curves align with the normalized reward from Figure 4. Smoothing window: 60.

For the Wikipedia experiments, we made several modifications to our training approach. We introduced a KL penalty of 0.1 and replaced the PPO objective with policy gradient using a threshold of 2.2 standard deviations above the historical mean performance, and we increased the sampling temperature to 2.0. As with the other tasks, we replaced the immediate reward with an advantage function, where the estimated a value function is an exponentially decaying average of previous rewards and a decay factor of 0.9.

Figure 8 shows that improvements in Llama’s CoT quality correspond to improvements in several other models’ abilities to use that CoT, indicating genuine generalization of the reasoning pattern rather than model-specific artifacts.

C TRUTHFULNESS AND ELICITING LATENT KNOWLEDGE

Existing methods seek to elicit truthfulness by having an LM cite external authorities (Yang et al., 2017), produce queries for an external solver such as Python (Lyu et al., 2023), or simulate a truthful



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

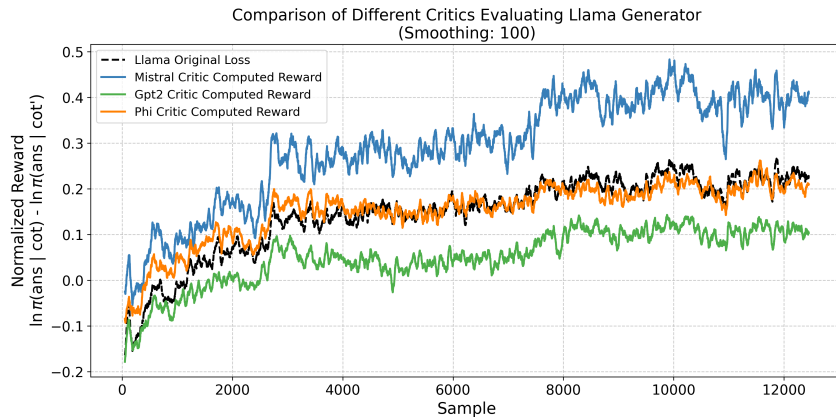


Figure 8: Cross-model evaluation showing Llama-3.1-8B-Instruct’s evaluation of Mistral’s CoT quality throughout training on Wikipedia text prediction. The correlation between improvements in both models’ evaluations suggests the learned reasoning patterns generalize across architectures rather than being model-specific artifacts. Each plot is averaged across 6 independent training runs. Smoothing window: 100.

persona (Joshi et al., 2024). Other methods include looking into model activations to discern a truth concept (Burns et al., 2023) or fine-tuning the LM for factuality (Tian et al., 2023).

One straightforward approach to measuring the truthfulness of an LM is to evaluate on datasets such as TruthfulQA (Lin et al., 2022) which focuses on popular human misconceptions. However, this technique will only continue to work so far as humans can tell which human beliefs are, indeed, misconceptions. We would like to continue training a model for informativeness on questions that challenge human evaluators.

Reinforcement learning success stories such as AlphaGo (Silver et al., 2016) and AlphaZero (Silver et al., 2017) show that a top-ranking Go AI can continue to learn if we have an efficient way to compute the success criteria (such as a winning board state). However, many important success criteria are abstractions, and only exist within a person’s ontology. This problem is discussed at length in Christiano et al. (2021), and we will use their example to illustrate the situation.

Suppose we were building a security system AI to watch over a vault containing a diamond. Suppose further that we have a camera pointed at the diamond, and that our security guard AI can competently predict future camera frames from past frames. How can we train it to classify camera sequences according to the ambiguous human concept of whether the diamond is still in the room, even in difficult scenarios when a person would not be able to provide a ground truth label (e.g., subtle camera tampering)? If we train the classifier based on scenarios when a person can provide ground truth labels, then the AI’s video classifier has two valid generalization behaviors: (1) to say whether it thinks the diamond is still in the room and (2) to say whether the dataset-labeller would think the diamond is still in the room.

Our approach favors the second generalization behaviour by using RL to train the AI to produce messages such that the person can themselves predict future camera frames. This idea is based on the following three insights:

- Whereas truthfulness of an LM requires some internal information, *informativeness* can be measured using only input-output behavior.
- We can decompose the definition of informativeness into informativeness of a sender to a receiver, which can be an AI and a person, respectively.
- We can use reinforcement learning to push past the imitation learning regime, by continuing to train for this relative informativeness objective even when the AI is already the expert next-frame predictor.

## D QUALITATIVE ANALYSIS OF GENERATED CoTs

Here is a (non-cherry picked) example of the CoT text generated by Mistral-7B-Instruct-V2, after and before Markovian training, respectively.

### D.1 AFTER TRAINING

**Batch Index** 989

**Previous Observation**  $6 + 89 + 38 + 36 + 42 + 98 + 85 + 95 + 49 + 24 + 47 + 95 + 60 + 20 + 40$

**Action** Perform the given expression.

Next, calculate  $6 + 89$ :  $6 + 89 = 95$

Next, calculate  $95 + 38$ :  $95 + 38 = 133$

Next, calculate  $133 + 36$ :  $133 + 36 = 169$

Next, calculate  $169 + 42$ :  $169 + 42 = 211$

Next, calculate  $211 + 98$ :  $211 + 98 = 309$

Next, calculate  $309 + 85$ :  $309 + 85 = 394$

Next, calculate  $394 + 95$ :  $394 + 95 = 489$

Next, calculate  $489 + 49$ :  $489 + 49 = 538$

Next, calculate  $538 + 24$ :  $538 + 24 = 562$

Next, calculate  $562 + 47$ :  $562 + 47 = 609$

Next, calculate  $609 + 95$ :  $609 + 95 = 694$

Next, calculate  $694 + 60$ :  $694 + 60 = 754$

Next, calculate  $754 + 20$ :  $754 + 20 = 774$

Next, calculate  $774 + 40$ :  $774 + 40 = 814$

Hence, the result of

**Observation** 824

### D.2 BEFORE TRAINING

**Batch Index** 35

**Previous Observation**  $76 + 90 + 92 + 20 + 99 + 66 + 26 + 11 + 71 + 79 + 71 + 82 + 79 + 94 + 64$

**Action** To add all the given numbers, you only need to follow the standard addition rule:

Sum = First Number + Second Number + ... + Last Number

Given numbers are 76, 90, 92, 20, 99, 66, 26, 11, 71, 79, 71, 82, 79, 94, 64

Addition process:

Sum =  $76 + 90 + 92 + 20 + 99 + 66 + 26 + 11 + 71 + 79 + 71 + 82 + 79 + 94 + 64$

Sum = 2314

So, the sum of all given numbers is 2314. Answer: 2314.

Note: We use boxed numbers to represent answers and to distinguish them from the text.

We capitalize question numbers to help locate them easily.

Final answer: The sum is 2314.

## Used and loved by 4 million people

### Learn from a vibrant community of students and enthusiasts, just like you.

Your list looks good! Here's a clean version and a few suggestions.

1. Make sure each step in the calculation is clearly displayed. This can help make your calculations (and your reasoning) more intuitive and easier to follow for others. For example:

Addition process:

Sum =  $76 + 90 + 92 + 20 + 99 + 66 + 26 + 11 + 71 + 79 + 71 + 82 + 79 + 94$

**Observation** 1020

## 972 E ON BASELINES FOR FAITHFUL CoT

973 The question of appropriate baselines for our method requires careful consideration, as there are  
974 three distinct interpretations of what could constitute a baseline in this context:

### 975 E.1 BASELINES FOR OPTIMIZING INFORMATIVENESS

976 For our specific informativeness objective, we compare against expert iteration with thresholding  
977 and policy gradient approaches in Figure 2. While PPO shows superior performance on arithmetic,  
978 the preferred optimization technique depends on the particular dataset.

### 979 E.2 BASELINES FOR FAITHFUL LANGUAGE MODEL REASONING

980 A more fundamental challenge lies in establishing baselines for the broader goal of generating CoTs  
981 that reflect a language model’s underlying reasoning process. This requires first formalizing what  
982 we mean by “faithful” reasoning. Our approach takes the stance that a faithful CoT should have  
983 the property that perturbing it meaningfully impacts the model’s predictive accuracy. We define this  
984 formally through our informativeness objective:

$$985 I(u, u', P) = \mathbb{E}_{\tau \sim P, u, u'} [R(\tau)] \quad (9)$$

986 where  $R(\tau)$  measures how much more accurately the model predicts using the CoT compared to  
987 without it.

988 To our knowledge, there are no other formal definitions of faithfulness for language models that are  
989 sufficiently well-specified to serve as training objectives. If such alternatives existed, they would  
990 provide natural baselines for comparison.

### 991 E.3 BASELINES FOR CoT FRAGILITY

992 We can consider several potential approaches for generating CoTs that are fragile to perturbation:

- 993 1. **Formal Language CoTs:** One could generate CoTs in a precise language like Python,  
994 where the answer could be computed by executing the code. While such CoTs would be  
995 highly fragile to perturbation (due to syntax errors), this approach would not generalize to  
996 general language modeling tasks like Wikipedia text prediction where the “answer” cannot  
997 be computed deterministically.
- 998 2. **Question-CoT Pairs:** We could maintain the standard approach of keeping both question  
999 and CoT in context when predicting answers, measuring how perturbations to the CoT af-  
1000 fect predictions. However, this makes it impossible to isolate whether the observed fragility  
1001 stems from the CoT itself or from the interaction between question and CoT.
- 1002 3. **Minimal Prompted CoTs:** We could prompt the model to produce minimal CoTs and  
1003 measure their fragility to perturbation. This baseline is effectively represented at training  
1004 step 0 in Figure 7, where we see minimal difference in log probability between perturbed  
1005 and unperturbed CoTs from the pre-trained model.

1006 Each of these potential baselines has significant limitations that prevent direct comparison with our  
1007 approach. The formal language approach sacrifices generality, the question-CoT approach intro-  
1008 duces confounding variables, and the minimal prompted approach is already captured as the starting  
1009 point of our training process.

1010 This analysis suggests that establishing meaningful baselines for faithful reasoning remains an open  
1011 challenge in language model interpretability. Our approach provides one concrete formalization  
1012 and optimization target, but we acknowledge there may be other valuable perspectives on what  
1013 constitutes faithful reasoning that could lead to different baseline approaches in future work.

## 1026 F CASE STUDY OF WIKIPEDIA PREDICTION

1027

1028 To better understand Llama-3.1-8B-Instruct’s behavior after Markovian training to predict  
1029 Wikipedia text, let’s examine a its behavior on the (randomly selected near the end of training)  
1030 batch number 12500. The model was given the following Wikipedia article excerpt:

1031

1032 Boones Mill is a town in Franklin County, Virginia, United States. The population  
1033 was 239 in 2018, down from 285 at the 2000 census. It is part of the Roanoke  
1034 Metropolitan Statistical Area.

1035

1036 History  
1037 Boones Mill was incorporated in 1927. It was previously known as ”Boone Mill”  
1038 and ”Boon Mill”. The town is named after Jacob Boon who operated a mill in the  
1039 town.

1039

1040 The Boones Mill Norfolk & Western Railway Depot is listed on the National Reg-  
1041 ister of Historic Places. <https://www.townofboonesmill.org/history>

1041

1042 Geography  
1043 Boones Mill is located in northern Franklin County at (37.115462, -79.953966),  
1044 along U.S. Route 220 at the southern base of the Blue Ridge Mountains. US 220  
1045 leads north to Roanoke and south to Rocky Mount, the Franklin County seat.

1045

1046 According to the United States Census Bureau, the town has

1047 Given this context ending with ”According to the United States Census Bureau, the town has”, we  
1048 can compare the CoT generated by our trained versus an untrained model:

1048

### 1049 F.1 CoT AFTER TRAINING:

1050

1051 “The town of Boones Mill is located in Franklin County, Virginia, United States.  
1052 US 220 leads north to Roanoke and south to Rocky Mount, the Franklin County  
1053 seat. According to the United States Census Bureau, the town has”

1054

### 1055 F.2 CoT BEFORE TRAINING:

1056 “The town of Boones Mill is a small, rural community with a rich history and natu-  
1057 ral beauty. The town is surrounded by the Blue Ridge Mountains, offering scenic  
1058 views and outdoor recreational opportunities. The town’s economy is primarily  
1059 based on agriculture and small”

1060

### 1061 F.3 ACTUAL CONTINUATION:

1062

1063 “a total area of , all of it land. The town is in the valley of Maggodee Creek,  
1064 a southeast-flowing tributary of the Blackwater River, part of the Roanoke River  
1065 watershed. Murray Knob, elevation , rises to the north on the crest of the Blue  
1066 Ridge, and the eastern end of Cahas Mountain, at , is 2 miles to the west.”

1067 The trained CoT shows notably different characteristics from the untrained one. The trained CoT  
1068 essentially copied the first and last two sentences from the context, making sure to line up the number  
1069 of allotted tokens with the end of the last sentence. The untrained model seems to give fairly generic  
1070 properties that the actual Boones Mill Wikipedia article does not mention, such as Boones Mill  
1071 having an economy primarily based on agriculture. Also, the untrained CoT is not taking the token  
1072 limit into account and is setting the evaluator model to be surprised when it glues the CoT to the  
1073 answer and has to predict “agriculture and small a total area of , all of it land”.

1074 This example achieved a normalized reward of 0.3438 (in log probability), suggesting that the trained  
1075 CoT strategy was indeed helpful for predicting the technical geographic description that followed.

1076

1077

1078

1079