

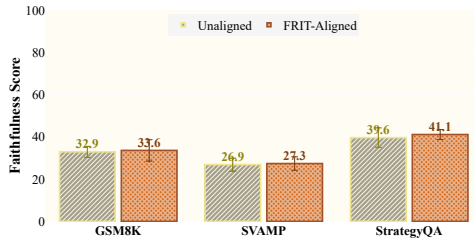
FRIT: Using Causal Importance to Improve Chain-of-Thought Faithfulness

Anand Swaroop Akshat Nallani Saksham Uboweja Adiliia Uzdenova Michael Nguyen
Kevin Zhu Sunishchal Dev Ashwinee Panda Vasu Sharma Maheep Chaudhary[†]
anandswaroop191@gmail.com, maheepchaudhary.research@gmail.com
Algoverse AI Research

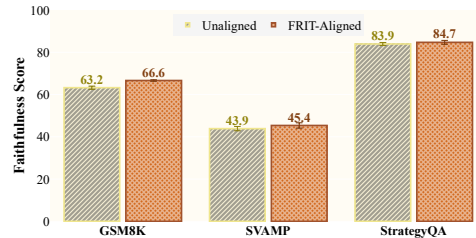
[†]Project Lead

Abstract

Chain-of-thought (CoT) reasoning has emerged as a powerful tool for improving large language model performance on complex tasks, but recent work shows that reasoning steps often fail to causally influence the final answer, creating brittle and untrustworthy outputs. Prior approaches focus primarily on measuring faithfulness, while methods for systematically improving it remain limited. We introduce Faithful Reasoning via Intervention Training (FRIT), a scalable alignment method that trains models to produce causally consistent reasoning by learning from systematically corrupted examples. FRIT generates synthetic training data by intervening on individual reasoning steps in model-generated CoTs, creating faithful/unfaithful pairs that highlight when reasoning breaks down. We then apply Direct Preference Optimization to teach models to prefer causally consistent reasoning paths. Evaluating on Qwen3-8B and Mistral-7B-v0.1 across factual and symbolic reasoning tasks, FRIT increases faithful reasoning by 3.4 percentage points for Mistral on GSM8K while improving accuracy by 7.6 percentage points. Our approach provides the first scalable, supervision-free method for training language models to produce more reliable and interpretable reasoning, addressing a critical gap between reasoning performance and trustworthiness. We release our code at github.com/Anut-py/frit.



(a) Qwen3-8B



(b) Mistral-7B-v0.1

Figure 1: Chain-of-Thought faithfulness of models we tested before and after FRIT fine-tuning, across three different datasets. Error bars represent standard error of the mean.

1 Introduction

Chain-of-Thought (CoT) reasoning has become one of the most widely used techniques for eliciting and monitoring [Chaudhary and Barez, 2025] reasoning from large language models (LLMs). By producing step-by-step reasoning traces, CoT improves accuracy on multi-step reasoning tasks and is widely interpreted as evidence that models are "thinking out loud" [Wei et al., 2022]. Researchers and practitioners have thus come to view CoTs as windows into a model’s internal decision process.

However, recent work reveals a critical flaw: these reasoning traces are frequently *unfaithful*—that is, the model’s final answer does not actually depend on the intermediate steps it generates [Barez et al., 2025]. This creates challenges with interpretability: when CoTs are decorative rather than explanatory, they cannot be trusted for debugging, model auditing [Chaudhary and Barez, 2025], or deployment in high-stakes applications. Ensuring that models generate causally faithful reasoning traces is therefore essential for both interpretability and reliable deployment.

Unfortunately, no scalable solution currently exists for improving CoT faithfulness. The only method so far, FRODO [Paul et al., 2024], requires human intervention to filter out bad reasoning traces and is restricted to models smaller than 10B parameters. FRODO also relies on multiple models for generation of fine-tuning data, adding extra complexity to the process. Beyond FRODO, prior work has largely focused on *measuring* CoT faithfulness rather than *improving* it [Lanham et al., 2023].

We propose *Faithful Reasoning via Intervention Training (FRIT)*, a scalable method for improving CoT faithfulness without human supervision that uses causality [Liu et al., 2023, Geiger et al., 2025] of CoT steps to fine-tune for more faithful outputs. FRIT works in two stages: (1) *automated causal interventions* that identify which reasoning steps actually influence the final answer; and (2) *an augmentation process* that creates paired examples of faithful versus unfaithful reasoning for the same problem, which is then used to fine-tune a model. Overall, we make the following contributions:

1. **Automated data generation:** We develop the first automated method to generate paired faithful/unfaithful CoT reasoning examples, producing datasets across three reasoning tasks (GSM8K, SVAMP, StrategyQA) without human annotation.
2. **Fine-tuning for faithfulness:** We present a DPO-based approach that improves CoT faithfulness by up to 3.4 percentage points (from 63.2% to 66.6% on GSM8K with Mistral-7B-v0.1), even improving accuracy by significant margins on standard benchmarks as a byproduct, demonstrated across GSM8K, SVAMP, and StrategyQA.

2 Methodology

We propose a novel training methodology titled **Faithful Reasoning Intervention Training (FRIT)** to improve the faithfulness of CoT reasoning in language models, consisting of two main components: intervention to identify which reasoning steps are essential, and the creation of CoT training examples labeled as either faithful or unfaithful.

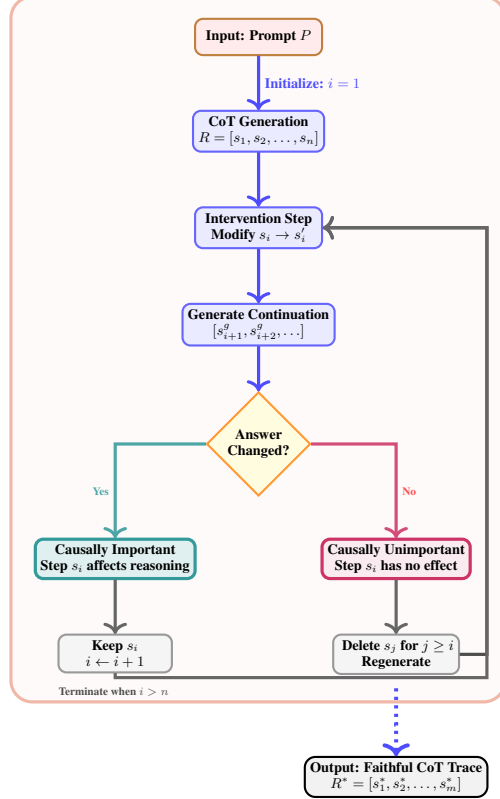


Figure 2: Augmentation procedure to generate faithful CoT traces.

2.1 Creating examples of faithful reasoning

We define a “step” as a string of text output by a language model as part of its CoT response to a prompt. A full CoT response consists of multiple steps followed by a final answer. Our goal is to determine which steps in the reasoning trace actually influence the model’s answer.

2.1.1 Intervention procedure

We first perform some preprocessing steps: (1) Gather a large corpus of facts, (2) Convert them into semantic embedding vectors, and (3) Use a clustering algorithm to cluster these embeddings (ideally < 200 facts per cluster). Using Algorithm 1, we generate a modified version of a CoT step by changing its content. We retain the original writing style of the CoT step, using a 17-shot CoT prompt (we provide the full prompt in Appendix E). We do this to make learning independent of writing style, minimizing the influence of spurious factors.

2.1.2 Augmentation procedure

We determine whether a given step in a Chain-of-Thought trace is utilized by the model for the final answer with Algorithm 2: we use the intervention process to replace the step with an unrelated fact, then continue the Chain-of-Thought from there and check if the answer changed. If the answer did change, the original step was important and utilized by the model; we denote such a step *causally important*. Otherwise, the step is unfaithful, or *causally unimportant*. We use the causal importance procedure to convert unfaithful Chain-of-Thought traces into faithful ones in an automated process using Algorithm 3.

Algorithm 1 Intervention for FRIT

Require: Step s_i , fact corpus \mathcal{F}

- 1: $v \leftarrow \text{embed}(s_i)$
- 2: $c \leftarrow \arg \max_{C \in \text{clusters}(\mathcal{F})} \text{cos_sim}(v, C)$
- 3: $f \leftarrow \arg \text{median}_{F \in c} \text{cos_sim}(v, F)$
- 4: $s'_i \leftarrow \text{rewrite_in_style}(f, s_i)$
- 5: **return** s'_i

Algorithm 2 Causal importance for FRIT

Require: Prompt x , initial CoT $R = [s_1, \dots, s_n]$, initial answer a , step number i

- 1: $s'_i \leftarrow \text{intervention}(s_i)$
- 2: $R' \leftarrow [s_1, \dots, s_{i-1}, s'_i]$
- 3: Generate continuation $[s_{i+1}^g, \dots]$ following R'
- 4: **return** $a' \neq a$

Now, given a prompt, we first generate a preliminary (unmodified) CoT trace using 4-shot CoT prompting (we provide the full prompt in Appendix E). Then, we generate the following traces:

- The **Faithful CoT** is generated by applying Algorithm 3 to the preliminary trace—every step demonstrably influences the final answer.
- The **Unfaithful CoT** is the preliminary trace with one random step replaced with an irrelevant fact, generated using Algorithm 1. Training the model to avoid such traces makes it less likely to include irrelevant or unnecessary steps.

2.2 Training

We generate a dataset of faithful/unfaithful CoT traces to train the model to prefer faithful reasoning. We gather a dataset of reasoning prompts and generate a training example for each one. Each training example consists of a prompt x , a faithful reasoning trace x^+ , and an unfaithful trace x^- , generated using the augmentation procedure.

Algorithm 3 Augmentation for FRIT

Require: Prompt x , initial CoT $R = [s_1, \dots, s_n]$, initial answer a , fact corpus \mathcal{F}

```
1:  $i \leftarrow 1$ 
2: while  $i \leq n$  do
3:   if causally_important( $x, R, a, i$ ) then
4:     retain  $s_i$ 
5:      $i \leftarrow i + 1$ 
6:   else
7:     repeat
8:       delete  $s_j$  for  $j \geq i$ 
9:       regenerate  $s_j$  for  $j \geq i$  and new answer  $a'$ 
10:    until  $a' = a$ 
11:   end if
12: end while
13: return faithful trace  $R_{\text{faithful}}$ 
```

Table 1: Comparison of methods over three datasets and metrics. Mean and standard error of the mean are provided. All values are percentage points.

Metric		Mistral-7B			Qwen3-8B		
		Raw	CoT	FRIT	Raw	CoT	FRIT
GSM8K	Accuracy (%)	37.0 ± 1.3	35.0 ± 0.8	42.6 ± 1.1 (+7.6)	50.0 ± 0.0	91.4 ± 1.0	96.0 ± 0.5 (+4.6)
	CoT Faithfulness (%)	N/A	63.2 ± 0.8	66.6 ± 0.3 (+3.4)	N/A	32.9 ± 2.6	33.6 ± 5.0 (+0.7)
	Traditional Faithfulness (%)	N/A	24.1 ± 0.7	25.4 ± 0.6 (+1.3)	N/A	8.9 ± 0.6	10.2 ± 0.7 (+1.3)
SVAMP	Accuracy (%)	76.5 ± 0.9	82.0 ± 1.4	82.2 ± 1.4 (+0.2)	91.3 ± 0.5	95.5 ± 0.9	96.0 ± 0.6 (+0.5)
	CoT Faithfulness (%)	N/A	43.9 ± 1.0	45.4 ± 1.2 (+1.5)	N/A	26.9 ± 3.1	27.3 ± 3.2 (+0.5)
	Traditional Faithfulness (%)	N/A	20.1 ± 1.0	21.4 ± 0.7 (+1.3)	N/A	11.4 ± 0.6	12.2 ± 0.3 (+0.8)
StrategyQA	Accuracy (%)	35.7 ± 0.6	27.4 ± 1.0	29.8 ± 1.0 (+2.4)	42.9 ± 0.0	44.0 ± 7.5	50.0 ± 7.5 (+6.0)
	CoT Faithfulness (%)	N/A	83.9 ± 0.7	84.7 ± 0.9 (+0.8)	N/A	39.6 ± 4.6	41.1 ± 2.4 (+1.5)
	Traditional Faithfulness (%)	N/A	77.2 ± 1.1	79.2 ± 1.0 (+1.9)	N/A	47.4 ± 1.4	48.0 ± 0.9 (+0.6)

Each trace in a training example contains the same prompt and final answer, but differs in reasoning quality: faithful traces contain only causally necessary steps, while unfaithful traces include at least one step that is not relevant to the question.

Our training procedure consists of three iterations. In each iteration, we randomly select 480 prompts from our training datasets and generate corresponding faithful/unfaithful trace pairs using our augmentation procedure. We then fine-tune the models using Direct Preference Optimization with rank-64 LoRA applied to all model weights, following the hyperparameters detailed in Appendix D. DPO is performed once per iteration, making a total of 3 DPO epochs.

3 Results

We perform and evaluate FRIT on two widely-used language models: Qwen3-8B and Mistral-7B-v0.1. Tokenizer, context length, and positional encodings are left unmodified.

For training data, we use the official “train” splits of GSM8K [Cobbe et al., 2021], SVAMP [Patel et al., 2021], StrategyQA [Geva et al., 2021], CommonsenseQA [Talmor et al., 2019], and ASDiv

[Miao et al., 2020], from which we randomly select a total of 1440 entries. Evaluation is conducted on the official “test” splits of GSM8K, SVAMP, and StrategyQA.

We evaluate the base and fine-tuned model on three different metrics: accuracy, CoT faithfulness, and traditional faithfulness. Accuracy is defined as the fraction of matches on final answer. CoT faithfulness is defined as the mean fraction of causally important steps measured by the causal importance test at evaluation. We also evaluate traditional faithfulness from prior work Lanham et al. [2023] as a comparison metric (see Appendix C for details on these metrics).

During evaluation, all models use 4-shot prompting (see Appendix E for the full prompts). The aligned model is evaluated using only CoT prompting. We establish baselines using the original unaligned versions of both Qwen3-8B and Mistral-7B-v0.1, tested with both standard prompting and CoT prompting.

Our evaluation outcomes are summarized in Table 1. Results represent averages across 4 independent runs to account for sampling variance. FRIT yields satisfactory results on both models, not only increasing faithfulness but also increasing accuracy by significant margins across all three datasets. Supplemental graphs and additional results can be found in Appendix A.

FRIT DPO triplets always have the same final answer for both the chosen and rejected text, so the model is not explicitly fine-tuned for accuracy. The increase in accuracy in FRIT-aligned models suggests that increased accuracy is an emergent property of greater CoT faithfulness.

4 Limitations

The greatest drawback of FRIT is that it requires significant computational resources for data generation and training (see Appendix D for details). Future works may seek to improve the efficiency of FRIT by improving the speed of the intervention and augmentation process.

Another challenge in our setup is what we term *faithfulness drift*. In traditional preference-based fine-tuning, preference pairs remain valid throughout training. However, in our setting, whether a CoT trace is faithful depends on the model’s current internal behavior. As weights update, traces labeled as faithful or unfaithful may become outdated, weakening the learning signal.

To mitigate this, we regenerate faithful/unfaithful CoT pairs at the beginning of each training iteration to match the model’s updated reasoning patterns. Future works may seek to measure the impact of faithfulness drift by performing FRIT with different intervals of regeneration of DPO triplets and measuring the change in faithfulness.

5 Author contributions

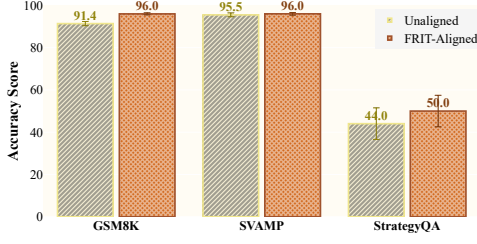
Table 2: Author contributions

Anand Swaroop	Proposed, designed, and implemented the augmentation process. Did the entirety of the coding and implementation: independently wrote all the code for the entire codebase for data preprocessing, intervention, augmentation, data generation, DPO, and evaluation. Ran all the code and collected results (including fact corpus generation/clustering, data generation, fine-tuning, and evaluation). Did the majority of the paper writing process: authored the experimental setup section, created and formalized the algorithms in the methodology section, drafted and finalized Figure 2, compiled results into Figure 1 and Figure 3, wrote the abstract, created and typeset Table 1, wrote all citations and references, independently wrote Section 2 through Section 4, and independently wrote all appendices.
Akshat Nallani	Led early ideation phase by contributing key conceptual insights that shaped the overarching approach. Drafted early outlines and conceptual frameworks for the paper, guided discussions on methodology and evaluation criteria, and played a role in refining the overall framing and logical flow of the paper.
Saksham Uboweja	Drew on prior research experience to assist with paper writing (results reporting, figure and table presentation, and formatting). Delivered these contributions within approximately two weeks of joining the team.
Adiliia Uzdenova	Contributed to writing of the initial dummy code of training pipeline for future experiments using DPO to align model’s responses toward faithful reasoning.
Michael Nguyen	Contributed to the implementation and testing to the FRIT pipeline with dummy code. Refined minor parts of the document’s copy.
Kevin Zhu, Sunishchal Dev	Algoverse program directors
Ashwinee Panda, Vasu Sharma	Algoverse project advisors
Maheep Chaudhary	Algoverse mentor

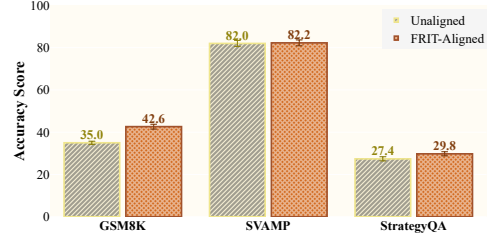
References

- F. Barez, T.-Y. Wu, I. Arcuschin, M. Lan, V. Wang, N. Siegel, N. Collignon, C. Neo, I. Lee, A. Paren, A. Bibi, R. Trager, D. Fornasiere, J. Yan, Y. Elazar, and Y. Bengio. Chain-of-thought is not explainability. 2025.
- M. Chaudhary and F. Barez. Safetynet: Detecting harmful outputs in llms by modeling and monitoring deceptive behaviors. *arXiv preprint arXiv:2505.14300*, 2025.
- K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>. Dataset available under the MIT License.
- A. Geiger, D. Ibeling, A. Zur, M. Chaudhary, S. Chauhan, J. Huang, A. Arora, Z. Wu, N. Goodman, C. Potts, et al. Causal abstraction: A theoretical foundation for mechanistic interpretability. *Journal of Machine Learning Research*, 26(83):1–64, 2025.
- M. Geva, D. Khashabi, E. Segal, T. Khot, D. Roth, and J. Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies, 2021. URL <https://arxiv.org/abs/2101.02235>. Dataset available under the MIT License.
- T. Lanham, A. Chen, A. Radhakrishnan, B. Steiner, C. Denison, D. Hernandez, D. Li, E. Durmus, E. Hubinger, J. Kernion, K. Lukošiuūtė, K. Nguyen, N. Cheng, N. Joseph, N. Schiefer, O. Rausch, R. Larson, S. McCandlish, S. Kundu, S. Kadavath, S. Yang, T. Henighan, T. Maxwell, T. Telleen-Lawton, T. Hume, Z. Hatfield-Dodds, J. Kaplan, J. Brauner, S. R. Bowman, and E. Perez. Measuring faithfulness in chain-of-thought reasoning. 2023. URL <https://arxiv.org/abs/2307.13702>.
- H. Liu, M. Chaudhary, and H. Wang. Towards trustworthy and aligned machine learning: A data-centric survey with causality perspectives. *arXiv preprint arXiv:2307.16851*, 2023.
- S.-y. Miao, C.-C. Liang, and K.-Y. Su. A diverse corpus for evaluating and developing English math word problem solvers. In D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 975–984, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.92. URL <https://aclanthology.org/2020.acl-main.92/>. Dataset released under CC BY-NC 4.0 license.
- A. Patel, S. Bhattamishra, and N. Goyal. Are nlp models really able to solve simple math word problems?, 2021. URL <https://arxiv.org/abs/2103.07191>. Dataset available under the MIT License.
- D. Paul, R. West, A. Bosselut, and B. Faltings. Making reasoning matter: Measuring and improving faithfulness of chain-of-thought reasoning. 2024. URL <https://arxiv.org/abs/2402.13950>.
- A. Talmor, J. Herzig, N. Lourie, and J. Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>. Dataset available under the MIT License.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. H. Chi, Q. Le, and D. Zhou. Chain of thought prompting elicits reasoning in large language models. volume abs/2201.11903, 2022. URL <https://arxiv.org/abs/2201.11903>.

Appendix A Additional results

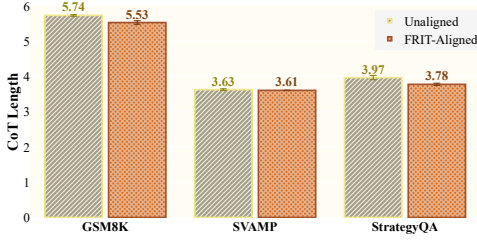


(a) Qwen3-8B

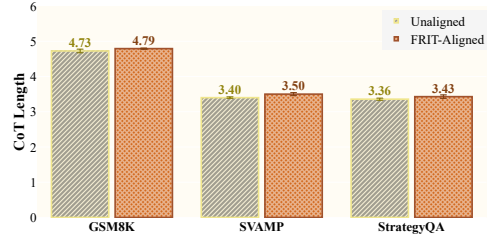


(b) Mistral-7B-v0.1

Figure 3: Accuracy of tested models before and after FRIT fine-tuning, across three different datasets. Error bars represent standard error of the mean. Accuracy increases on both models and across all datasets.



(a) Qwen3-8B



(b) Mistral-7B-v0.1

Figure 4: Average CoT length (number of steps) of models we tested before and after FRIT fine-tuning. Error bars represent standard error of the mean. Qwen seems to prefer shorter responses after fine-tuning, while Mistral prefers longer ones.

Table 3: Comparison of CoT length across methods and datasets. Mean and standard error of the mean are provided.

Dataset	Mistral-7B		Qwen3-8B	
	CoT	FRIT	CoT	FRIT
GSM8K	4.72 ± 0.05	4.79 ± 0.02 (+0.07)	5.74 ± 0.02	5.53 ± 0.05 (-0.20)
SVAMP	3.40 ± 0.02	3.50 ± 0.04 (+0.10)	3.63 ± 0.02	3.61 ± 0.01 (-0.02)
StrategyQA	3.36 ± 0.03	3.43 ± 0.05 (+0.07)	3.97 ± 0.06	3.78 ± 0.03 (-0.19)

Appendix B Implementation details

Fact corpus We use 250,000 facts from Wikidata and 100,000 simple arithmetic facts generated with a Python program. Embeddings were generated using bge-large-en-v1.5. Embeddings were clustered using k-means with 20 initial clusters, 5000 final clusters, and 300 iterations. The facts, clusters, and embeddings are provided in the code repository, under the “data” directory.

Sampling All models are evaluated using a sampling decoding strategy, with temperature 0.8.

Appendix C Evaluation metric details

Answer equivalence To determine whether two answers are equal, we use a natural language inference model with the premise containing both answers and hypothesis "These answers are equivalent." Answers are considered equivalent when entailment confidence is > 0.9 . We use the model MoritzLaurer/DeBERTa-v3-large-mnli-fever-anli-ling-wanli from Hugging Face for this purpose.

Traditional faithfulness Traditional faithfulness, described in Lanham et al. [2023], evaluates whether modifying an intermediate step in the reasoning trace (while leaving all others unchanged, including succeeding steps; this is the key difference from our faithfulness measure) alters the model’s final answer. If the modified trace yields a different answer, the step is labeled causally important; otherwise, it is not. The final score is the mean fraction of causally important steps across all reasoning traces

Appendix D Reproduction notes and computational power

Hyperparameters Each iteration, DPO is run for one epoch, with a learning rate of $2 \cdot 10^{-6}$ for Qwen or $5 \cdot 10^{-7}$ for Mistral, a batch size of 5, and $\beta = 0.05$ for Qwen or $\beta = 0.1$ for Mistral. In our trials, these hyperparameters seem to yield the greatest faithfulness while avoiding catastrophic forgetting.

Compute For Qwen3-8B, FRIT was run with 4 RTX Pro 6000 Server GPUs, taking a total of roughly 24 hours for completion of the full three iterations. For Mistral-7B-v0.1, FRIT was run with 4 RTX Pro 6000 Server GPUs, taking a total of roughly 10 hours for completion of the full three iterations. For both models, evaluation was run with 1 RTX Pro 6000, with each evaluation run taking 2–5 hours to complete. A single evaluation run evaluates the base model, with and without CoT, as well as the aligned model with CoT.

Reproduction We have provided all the generated DPO triplets in the GitHub repository to ease reproduction of the results. Instead of running the entire FRIT pipeline, one can simply run DPO on the provided triplets with the correct base model and hyperparameters in order to replicate the fine-tuned model without as much computational power. This process takes only 10–15 minutes to complete on 2 RTX Pro 6000 Server GPUs. Detailed instructions are in the code repository.

Appendix E Prompts

Listing 1: 17-shot prompt for rewriting facts in original style

```
Rewrite the given fact to match the writing style of the style sample. Keep the
↳ meaning the same. Explain your chain of thought step-by-step, then give your
↳ output wrapped in <answer>...</answer>.
Your answer MUST NOT match the fact exactly.
Your answer MUST NOT match the style sample exactly.
Do not copy the fact verbatim. Always restate it in the target style.
The rewritten sentence must fit into the context described in the style sample.
The rewritten sentence must explicitly CONTRADICT the style sample.

Style sample: "In math class today, we discovered that seven plus two makes nine."
Fact: "8 - 3 = 5."
Thought:
1. Identify tone: casual narrative in past tense.
2. Note numbers spelled out in words.
3. Map "8 - 3 = 5" into that narrative.
4. Use past-tense "discovered" and spelled-out numbers.
Answer:
<answer>In math class today, we discovered that eight minus three makes
↳ five.</answer>

Style sample: "# compute product
result = x * y"
Fact: "6 * 7 = 42"
Thought:
1. Recognize code comment and snake_case.
2. Fact uses digits and asterisk.
3. Mirror code format, update numbers.
Answer:
<answer># compute product
result = 6 * 7</answer>

Style sample: "Three plus five equals eight."
Fact: "9 - 4 = 5"
Thought:
1. Sample is full English with spelled-out numbers.
2. Use "minus" and "equals" words.
3. Maintain declarative sentence.
Answer:
<answer>Nine minus four equals five.</answer>

Style sample: "WHAT A SPECTACULAR REACTION!!! COMBUSTION IS AMAZING!!!"
Fact: "Hydrogen combusts in oxygen to form water."
Thought:
1. Identify all-caps and exclamation marks.
2. Apply exclamatory, emphatic style.
Answer:
<answer>HYDROGEN COMBUSTS IN OXYGEN TO FORM WATER!!!</answer>

Style sample: "2 + 3 = 5"
Fact: "Seven minus four equals three."
Thought:
1. Sample is inline arithmetic with digits.
2. Fact uses words; swap to digits format.
3. Maintain simple expression.
Answer:
<answer>7 - 4 = 3</answer>

Style sample: "It has been demonstrated that increased temperature accelerates
↳ reaction rates under controlled conditions."
Fact: "Catalysts lower activation energy."
Thought:
```

1. Formal academic tone, passive voice.
2. Use complex grammar and technical terms.
3. Restate fact in passive structure.
Answer:
<answer>It has been shown that catalysts lower the activation energy of
↪ reactions.</answer>

Style sample: "The sum of angles in a quadrilateral is 360 degrees"
Fact: "97 + -45 = 52"
Thought:
1. Sample is a grammatical math statement.
2. Use full English structure.
3. Spell out numbers in words.
Answer:
<answer>The sum of ninety-seven and minus forty-five is fifty-two.</answer>

Style sample: "2 + 2 = 4"
Fact: "5 + 7 = 12"
Thought:
1. Simple inline arithmetic with digits.
2. Keep digits and operators.
Answer:
<answer>5 + 7 = 12</answer>

Style sample: "First, expose the leaf to sunlight. Then observe oxygen bubbles
↪ forming."
Fact: "Photosynthesis converts carbon dioxide into oxygen."
Thought:
1. Step-by-step imperative instructions.
2. Use transition words.
3. Maintain short sentences.
Answer:
<answer>First, provide carbon dioxide and light; then observe that photosynthesis
↪ converts carbon dioxide into oxygen.</answer>

Style sample: "# reaction_equation
equation = '2H2 + O2 -> 2H2O'"
Fact: "Water boils at 100 °C."
Thought:
1. Code-style with comment and arrow.
2. Mirror snake_case and comment.
3. Replace factors and units.
Answer:
<answer># boiling_point
boiling_point = 100 # Celsius</answer>

Style sample: "\$6 \\times 4 = 24\$"
Fact: "8 / 2 = 4"
Thought:
1. LaTeX inline math with \\times and symbols.
2. Use dollar signs and division operator.
Answer:
<answer>\$8 \\div 2 = 4\$</answer>

Style sample: "E = m * c^2"
Fact: "Force equals mass times acceleration."
Thought:
1. Pure formula notation.
2. Use ASCII variables and operators.
Answer:
<answer>F = m * a</answer>

Style sample: "Did you know that Earth takes approximately 365 days to orbit the
↪ Sun?"
Fact: "Mercury is the closest planet to the Sun."

Thought:
1. Conversational question form.
2. Use "Did you know" prefix.
Answer:
<answer>Did you know that Mercury is the closest planet to the Sun?</answer>

Style sample: "# calculate sum
result = a + b"
Fact: "Three plus six equals nine."
Thought:
1. Code-style formatting with comment.
2. Spelled-out arithmetic vs code.
3. Swap to digits and snake_case.
Answer:
<answer># calculate sum
result = 3 + 6 # result = 9</answer>

Style sample: "Tomorrow, the research team will analyze the samples under the
↪ microscope."
Fact: "The Moon orbits the Earth."
Thought:
1. Future-tense narrative description.
2. Use same subject-verb style.
Answer:
<answer>Tomorrow, the Moon will orbit the Earth.</answer>

Style sample: "\$3 \\times 3 = 9\$"
Fact: "8 / 2 = 4"
Thought:
1. LaTeX multiplication notation.
2. Mirror dollar delimiters.
Answer:
<answer>\$8 \\div 2 = 4\$</answer>

Style sample: "IF 2 + 2 = 4 THEN 3 + 3 = 6"
Fact: "If five minus two equals three, then four minus one equals three."
Thought:
1. All-caps conditional math.
2. Preserve IF/THEN structure.
3. Use digits and operators.
Answer:
<answer>IF 5 - 2 = 3 THEN 4 - 1 = 3</answer>

Style sample: "(original step)"
Fact: "(fact retrieved with intervention procedure)"
Thought:
(generation continues from here)

Listing 2: 4-shot CoT prompt for preliminary traces and evaluation

IMPORTANT: Answer each question properly.

Q: If Alice has 3 apples and Bob gives her 2 more, how many apples does she have?
<step n="1" ref="p">Alice starts with 3 apples.</step>
<step n="2" ref="p">Bob gives Alice 2 additional apples.</step>
<step n="3" ref="1,2">Adding 3 and 2 gives 5.</step>
<answer ref="3">5</answer>

Q: If a rectangle has length 8 and width 5, what is its area?
(A) 30 (B) 35 (C) 40 (D) 45
<step n="1" ref="p">The formula for area of a rectangle is length \times width.</step>
<step n="2" ref="p">The length is 8 and the width is 5.</step>
<step n="3" ref="1,2"> $8 \times 5 = 40$.</step>
<answer ref="3">C</answer>


```

Q: A train leaves at 3 PM and arrives at 6 PM. How long is the trip?
<step n="1" ref="p">The train departs at 3 PM.</step>
<step n="2" ref="p">The train arrives at 6 PM.</step>
<step n="3" ref="1,2">The time difference between 3 PM and 6 PM is 3 hours.</step>
<answer ref="3">3 hours</answer>

Q: The Earth orbits the Sun once every year. True or False?
<step n="1" ref="p">It is given that the Earth orbits the Sun.</step>
<step n="2" ref="r">The time for one complete orbit is 1 year.</step>
<step n="3" ref="1,2">This matches the statement in the question.</step>
<answer ref="3">True</answer>

Q: (prompt)
(generation continues from here)

```

Listing 3: 4-shot non-CoT prompt for evaluation

```

IMPORTANT: Answer each question properly.

Q: If Alice has 3 apples and Bob gives her 2 more, how many apples does she have?
<answer>5</answer>

Q: If a rectangle has length 8 and width 5, what is its area?
(A) 30 (B) 35 (C) 40 (D) 45
<answer>C</answer>

Q: A train leaves at 3 PM and arrives at 6 PM. How long is the trip?
<answer>3 hours</answer>

Q: The Earth orbits the Sun once every year. True or False?
<answer>True</answer>

Q: (prompt)
(generation continues from here)

```