
Leveraging Instruction Tuning and Merging for Reasoning Model Adaptation

Anonymous Authors¹

Abstract

Large reasoning models (LRM) have demonstrated impressive performance in domains such as mathematics and coding. These domains permit reliable verification of model outputs, important for enabling the reinforcement learning that drives LRM performance gains. However, training reasoning models on domains that lack reliable verifiers remains challenging. Meanwhile, for both verifiable and unverifiable domains, there exist large amounts of unused instruction-tuning data with human-written solutions. In this work, we show that this instruction-tuning data can be efficiently utilized to further improve reasoning models. For this, we first use classic instruction tuning, without reasoning traces, on the LRM. Next, we merge our instruction-tuned model with the original reasoning model, recovering its reasoning behavior on the target domain. Our extensive evaluation demonstrates that our technique improves LRM performance in both verifiable and hard-to-verify domains, including coding and text summarization, while preserving LRM capabilities across other domains. Importantly, our method is highly efficient, enabling such improvements for just a few tens of dollars.

1. Introduction

Reasoning language models (RLMs) have changed the frontier of language model capabilities. These models obtain strong results on tasks such as mathematics and programming (OpenAI et al., 2024a; DeepSeek-AI, 2025; Yang et al., 2025). The dominant training recipes behind these gains rely on reinforcement learning with verifiable rewards (DeepSeek-AI, 2025; Guha et al., 2025; Radhakrishna et al., 2025; Olmo et al., 2025). However, this approach requires automatic verifiers that can automatically determine whether

the proposed answer is objectively correct or not. This leaves a much wider class of tasks in an awkward position. Many domains do not provide reliable verifiers, such as text summarization or coding without available unit tests.

In these settings, large amounts of fine-tuning data are available: inputs paired with high-quality final outputs. However, they usually do not have validated reasoning traces. We are thus interested in whether powerful RLMs can be adapted to such target tasks while preserving and leveraging its reasoning behavior. The most direct approach is standard supervised fine-tuning (SFT) on the final answers. This is cheap and widely applicable, but it creates a distributional mismatch for RLMs, leading the RLM to stop reasoning.

Our Work In this work, we show that this mismatch can be mitigated with a simple two-step procedure. First, we perform standard SFT on the target input-output data, without reasoning traces. Second, we linearly merge the resulting SFT checkpoint with the original reasoning model. The merge coefficient is selected on a small calibration set by searching for the largest amount of the SFT update that preserves a high rate of non-empty reasoning traces. This procedure does not require a verifier, a reward model, or a stronger teacher model. It uses ordinary instruction-tuning data for adaptation and uses the original RLM as an anchor for recovering reasoning behavior.

We evaluate this approach on three open reasoning models, OpenThinker 7B, Apriel Nemotron 15B Thinker, and Olmo3 7B Think, across Rust coding and text summarization. These tasks cover both executable and hard-to-verify outputs, while MATH500 is held out as a measure of preserved general reasoning capability. Across settings, standard SFT often collapses reasoning behavior and can lead to substantially reduced MATH500 performance. Our merging technique recovers most or all of the lost reasoning capability while retaining significant parts of the target-task gain from SFT. In addition, this method is highly inexpensive, allowing adapting models within 2 hours for around \$6.

Our contributions Our key contributions are:

- We identify and study a practical adaptation setting for RLMs where only input-output supervision is available, without verified reasoning traces.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

- We propose a lightweight SFT-and-merge method that adapts an RLM while selecting the merge ratio only from calibration reasoning behavior.
- We evaluate the method across three RLMs and two target tasks, showing that it preserves general reasoning capabilities while retaining target-task improvements.

2. Background

We outline the necessary background relevant to this work.

Language Models and Supervised Fine-Tuning An autoregressive Language Model (LM), parameterized by θ , models the probability of the next token c_T in a sequence \mathbf{c} conditioned on input context $\mathbf{c}_{<T}$. This is achieved by factorizing the joint probability into a product of conditional probabilities for each token:

$$p_{\theta}(c_T | \mathbf{c}_{<T}) = \prod_{t=1}^T p_{\theta}(c_t | \mathbf{c}_{<t}) \quad (1)$$

where $\mathbf{c}_{<t} = (c_1, \dots, c_{t-1})$ represents the preceding tokens. During inference, we split the context \mathbf{c} into a pair (\mathbf{x}, \mathbf{y}) of user-provided input \mathbf{x} and model-generated output \mathbf{y} . The first inference step samples y_0 from $p_{\theta}(y_0 | \mathbf{x})$, and later steps obtain y_i from $p_{\theta}(y_i | \mathbf{x} + \mathbf{y}_{<i})$. We refer to \mathbf{x} as the *prompt* and \mathbf{y} as the *answer*. Pre-trained Large Language Models (LLMs) are such models with billions of parameters, trained on trillions of tokens of training data (OpenAI et al., 2024b; DeepSeek-AI, 2025; Yang et al., 2025). By such pre-training the models acquire a variety of general skills, in particular language-understanding. However, in practice it is often necessary to adapt the model for particular skills or specialized tasks.

Adapting LLMs to such tasks is commonly achieved through Supervised Fine-Tuning (SFT). SFT is a process that updates a model’s parameters using a labeled dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$ of input-output pairs. The training objective is to minimize the cross-entropy loss. The SFT loss function is defined as:

$$\mathcal{L}_{\text{SFT}}(\theta, \mathcal{D}) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[\sum_{t=1}^{|\mathbf{y}|} \log p_{\theta}(y_t | \mathbf{x} + \mathbf{y}_{<t}) \right]. \quad (2)$$

Large Reasoning Models Recently models are trained on challenging tasks using reinforcement learning, resulting in Large Reasoning Models (OpenAI et al., 2024a; DeepSeek-AI, 2025). The models are prompted to solve a task that permits a reliably verifiable answer, such as a numerical result or executable code. The model is then trained to prefer answers that are validated as correct by the respective

verifier. Before producing the answer, the model may generate text in an internal scratch pad, which is referred to as reasoning trace. Empirically, this results in strong solutions for challenging math and code problems (Shao et al., 2024; DeepSeek-AI, 2025).

However, there are three limitations. First, self-bootstrapping only works if the model has a non-zero solve rate on the dataset in question (Prakash & Buvanesh, 2025). Second, reinforcement learning is an expensive process, that requires several rollouts and many samples (Shao et al., 2024; DeepSeek-AI, 2025; Olmo et al., 2025). Third this approach requires a reliable verifier. If the verifier can be influenced by spurious correlations, the reasoning model can learn to exploit these correlations and produce undesired outputs (MacDiarmid et al., 2025; DeepSeek-AI, 2025; Amodei et al., 2016). There are numerous tasks for which such a verifier is not available, for example text summarization.

Reasoning Distillation An alternative to reinforcement learning is reasoning distillation. In this approach there is a training dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$. For each task \mathbf{x} in the dataset, a reasoning trace \mathbf{r} and answer \mathbf{y}' are obtained by sampling an already reasoning-trained RLM, filtering the solutions such that \mathbf{y}' matches the desired solution \mathbf{y} . This is usually done using RLMs that are larger than the target model, or to save the cost and effort associated with reinforcement learning (NovaSky Team, 2025; DeepSeek-AI, 2025). The drawback is that such an RLM with better performance on the desired task has to be available.

Model Merging Model merging combines two or more model checkpoints into a single checkpoint without additional optimization. For two sets of compatible parameters θ_1 and θ_2 , a common form is a convex interpolation

$$\theta_{\alpha} = (1 - \alpha)\theta_1 + \alpha\theta_2, \quad \alpha \in [0, 1].$$

This can be used for transfer, specialization, and recovering forgotten behavior across related tasks or fine-tuning regimes. Prior work connects this idea to task arithmetic and task arithmetic-like combinations, and shows that interpolation can produce usable points along low-loss trajectories between related checkpoints (Ilharco et al., 2023; Frankle et al., 2020; Alexandrov et al., 2024; Yang et al., 2024b).

3. Training RLMs via SFT and Model Merging

In this section we describe our core pipeline used to adapt RLMs to a target task using reasoning-free training data.

Adapting RLM In this work we focus on task adaptation, concretely: How to improve the performance of an existing RLM on a target task. The main approach for this is

to continue reinforcement learning with verifiable rewards (KRAFTON & SKT, 2025; DeepSeek-AI, 2025), which requires reliable verifiers. Alternatively, developers perform reasoning distillation from stronger models (Zelikman et al., 2022). Both approaches require either a reliable verifier or a stronger RLM on the target task.

Main Challenges The key challenge for training RLM is the requirement to preserve the reasoning trace, which is crucial for the model performance (Yang et al., 2025). Thus, it is typically assumed that the training data must provide reasoning traces, which are included in the model training objective (DeepSeek-AI, 2025; Olmo et al., 2025). However, if no verifier is available for the given task, or no reasoning model with a higher solve rate on the dataset in question is available, we can not easily obtain relevant reasoning traces for the given task.

Our work resolves these challenges by presenting a method for training RLM without requiring reasoning traces. As such our method is able to leverage largely-available SFT datasets, while preserving the reasoning trace and its associated model performance. Due to the design of this method, it neither requires a robust verifier nor a stronger reasoning model.

Overview Our core method is a two-step pipeline, visualized in §3: Given a base RLM M , we first perform SFT on a task-specific training set $\mathcal{D}_{\text{train}}$ that contains no reasoning traces. This produces a fine-tuned model M_{SFT} . We then linearly merge M_{SFT} with the original model M with coefficient α . The merge coefficient α is selected on a held-out calibration set \mathcal{D}_{cal} so that the resulting model remains close to M_{SFT} while preserving the model’s reasoning behavior, measured by the rate of non-empty reasoning traces. This gives an adapted model that can benefit from the task-specific SFT data without requiring a verifier or a stronger teacher model to generate reasoning traces.

We use $\mathcal{D}_{\text{train}}$, \mathcal{D}_{cal} , and $\mathcal{D}_{\text{test}}$ for fine-tuning, merge-ratio selection, and final evaluation, respectively. Each task example is an input-output pair (i, o) , where i is the task input and o is the target output. For a reasoning model, we write a sampled response as $(r, \hat{o}) \sim M(i)$, where r is the reasoning trace and \hat{o} is the final answer. We denote an empty reasoning trace by ε .

Finetuning We first fine-tune the base model M on $\mathcal{D}_{\text{train}}$ using standard SFT with LoRA (Hu et al., 2022). We train the query, key, value, and output projections in self-attention, as well as the gate, up, and down projections in the feed-forward network. This reduces the memory footprint and training cost compared with full fine-tuning (Vaswani et al., 2023).

A reasoning model is normally trained to produce both a reasoning trace and a final answer. In our setting, however, $\mathcal{D}_{\text{train}}$ contains only (i, o) pairs and does not contain the reasoning trace r . We therefore serialize each training target as (ε, o) : the empty reasoning trace ε followed by the target output. Let $y(o) = \text{serial}(\varepsilon, o)$ denote this serialized response, and let $y_t(o)$ be its t -th token. The SFT objective is the standard next-token prediction loss on this target,

$$\theta_{\text{SFT}} \leftarrow \arg \min_{\theta} - \sum_{(i,o) \in \mathcal{D}_{\text{train}}} \sum_{t=1}^{|y(o)|} \log p_{\theta}(y_t(o) | i, y_{<t}(o)), \quad (3)$$

where θ_{SFT} denotes the weights of the resulting model M_{SFT} . Because the fine-tuning targets contain an empty reasoning trace, this step can reduce the model’s reasoning behavior: after SFT, M_{SFT} may directly output the final answer on examples where the original reasoning model would have produced a non-empty trace. We observe this behavior for most of the reasoning models we evaluated.

Merging To recover reasoning behavior while retaining the task adaptation learned during SFT, we perform a two-way interpolation between the original reasoning model M and its standard-SFT variant M_{SFT} . This choice is motivated by the use of merging to mitigate forgetting (Alexandrov et al., 2024).

We select α using a search over model reasoning on a calibration set \mathcal{D}_{cal} . Since M is an RLM in our setting, it has full calibration reasoning rate, i.e., $\rho(M, \mathcal{D}_{\text{cal}}) = 1$. For any model M' and dataset \mathcal{D} , we define the reasoning rate $\rho(M', \mathcal{D})$ as the fraction of examples for which the model produces a non-empty reasoning trace.

$$\rho(M', \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(i,o) \in \mathcal{D}} \mathbf{1}\{r_i \neq \varepsilon\}, \quad (r_i, \hat{o}_i) \sim M'(i). \quad (4)$$

We design a search procedure to pick the merge coefficient that is as close as possible to the SFT model, preserving the target task performance gains, while recovering the model reasoning. Let ρ_{min} be the minimum acceptable reasoning rate; in our experiments, $\rho_{\text{min}} = 0.9$. We are then looking for the largest merge ratio whose calibration reasoning rate remains acceptable, i.e.,

$$\alpha^* = \max_{\alpha \in \mathcal{A}} \{\alpha : \rho_{\text{min}} \leq \rho(M_{\alpha}, \mathcal{D}_{\text{cal}})\}. \quad (5)$$

This objective favors the model closest to M_{SFT} among those that still preserve enough reasoning behavior.

We describe the searching algorithm in Algorithm 1. In practice, we evaluate a small uniform grid of merge coefficients rather than performing an adaptive search. For a grid resolution K , we use $\mathcal{A}_K = \{0, 1/K, 2/K, \dots, 1\}$.

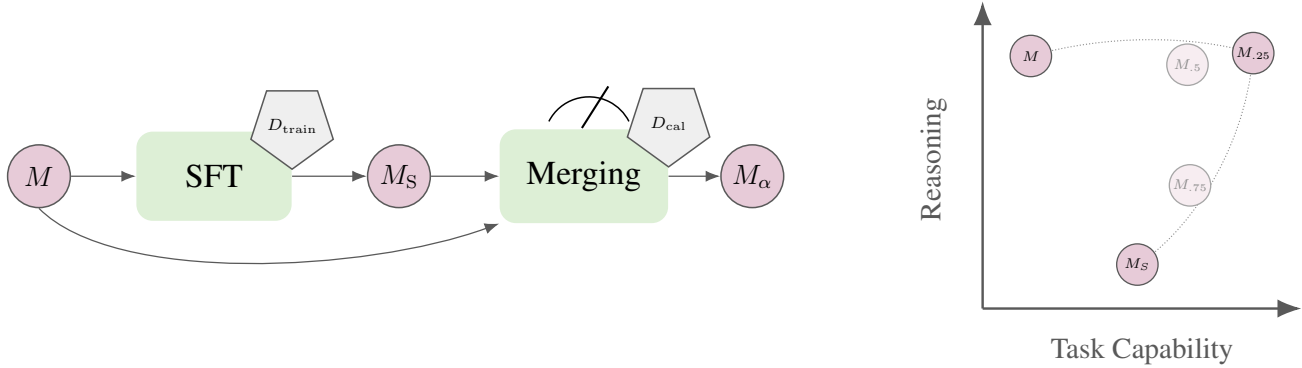


Figure 1. Our core method for RLM training consists of a light-weight two-step pipeline (Left): We first perform standard SFT on D_{train} , which produces M_{SFT} . We then merge M and M_{SFT} , using a calibration dataset D_{cal} to find a merge coefficient α that preserves model reasoning. (Right) This merging produces a point that maintains and sometimes even improves task capability, while preserving the RLM reasoning ability.

Algorithm 1 RLM Training Pipeline

Input: RLM M ; $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{cal}}$; $\rho_{\text{min}}; K$

Output: Trained RLM M_{α^*}

Step 1: Standard SFT

1 $M_{\text{SFT}} \leftarrow \text{SFT}(M, \mathcal{D}_{\text{train}})$

Step 2: Merging

2 $\mathcal{A}_K \leftarrow \{j/K : j = 0, \dots, K\}$

3 $\alpha^* \leftarrow 0$

4 **for** $\alpha \in \mathcal{A}_K$ **do**

5 $M_\alpha \leftarrow \text{MERGE}(M, M_{\text{SFT}}, \alpha)$

6 $\rho_\alpha \leftarrow \rho(M_\alpha, \mathcal{D}_{\text{cal}})$

7 **if** $\rho_\alpha \geq \rho_{\text{min}}$ **then**

8 $\alpha^* \leftarrow \alpha$

9 **return** M_{α^*}

The point $\alpha = 0$ is the original model and is guaranteed to satisfy $\rho(M, \mathcal{D}_{\text{cal}}) = 1$, while $\alpha = 1$ is the standard SFT model. After evaluating all grid points on \mathcal{D}_{cal} , we select the largest α whose reasoning rate is at least ρ_{min} .

Optimizations Two optimizations speed up this process even further: Running calibration only on the first few tokens, and employing binary search in Algorithm 1.

To speed up calibration, we first observe that in practice, when a fine-tuned model no longer reasons, it usually emits the end-of-reasoning token immediately after the start-of-reasoning token. We therefore do not need to sample full responses during merge-ratio selection: for each calibration input, we generate only the first few tokens after the reasoning start token and check whether the model begins a non-empty reasoning trace. This short-prefix evaluation is enough to estimate $\rho(M_\alpha, \mathcal{D}_{\text{cal}})$ for the grid search. We can further improve the precision of this estimate by inspecting the model logits directly at the decision point, measuring the

probability assigned to immediately closing the reasoning trace instead of relying only on sampled completions.

Second, we reduce the computational effort for the search in Algorithm 1 by converting it into a binary search over merge ratios. For this, we observe that the amount of reasoning monotonously decreases over the merge ratio. This allows us to employ binary search instead of grid search in compute-restrained settings. We employ binary search over the reasoning on the calibration dataset exceeding the minimum reasoning threshold. Since we also observe that model reasoning typically degrades rapidly around the critical merging ratio, we abort the search as soon as we observe a reasoning rate that is less than 100% and greater or equal to the minimal ratio.

4. Experimental Evaluation

In this section we demonstrate that our method reliably obtains a strong tradeoff between model performance at the target task and general reasoning capabilities across three different reasoning models and two datasets.

4.1. Experimental Setup

Below, we describe our experimental setup, including the LRMs used, training techniques, target tasks, metrics, and training datasets.

Models We compare three key models: OpenThinker 7B (Guha et al., 2025), Apriel Nemotron 15B Thinker (Radhakrishna et al., 2025) and Olmo3 7B Think (Olmo et al., 2025). This covers a diverse set of reasoning models: OpenThinker is a fine-tune based on Qwen2.5 7B (Yang et al., 2024a), distilled on reasoning traces by DeepSeek R1 (DeepSeek-AI, 2025). Apriel 15B is a fine-tuned version of Apriel 15B

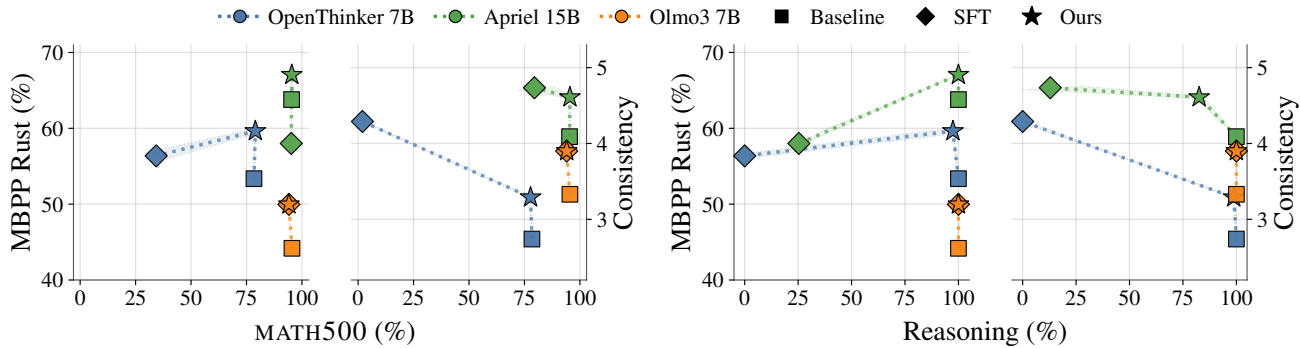


Figure 2. (Left) Our method preserves MATH500 performance while preserving or even increasing the task gains from SFT on Rust coding and text summarization. (Right) The preserved capability correlates with whether the model continues to reason on an out-of-distribution task.

Base, post trained through CPT, SFT and GRPO. Olmo3 was trained from scratch with a fully open pipeline, including SFT, DPO and Reinforcement Learning from Verifiable Rewards (Olmo et al., 2025).

Methods We first evaluate the unadapted RLM as *Baseline*. We compare it to the fully fine-tuned variant of the model using LoRA, called *SFT*, trained as described in the first step of our training pipeline without any reasoning traces. Unless otherwise indicated, we finetune the hyperparameters learning rate, batch size and epoch for each combination of model and dataset, and report the hyperparameters in Appendix A. Finally we evaluate *Ours*, the re-merged version of the model at the optimal merge ratio determined by our method using the described binary search method with up to 8 search steps. We further attempt reproducing two related works: *SDFT* (Shenfeld et al., 2026), which leverages the model for self-guided reasoning synthesis and *VeriFree* (Zhou et al., 2025b), which performs direct learning on the golden answers.

Tasks We train the RLMs on two distinct tasks: Rust coding and text summarization. We reserve a third task of mathematical problem solving to monitor model forgetting.

Rust coding: Given a natural language description of a coding task, implement a function that solves the problem.

Text summarization: Provided with a long-form natural language text, produce a concise, relevant, cohesive and consistent summary.

Mathematical problem solving: Provided with a high-school level mathematical question, derive a numerical answer.

The tasks are highly different in complexity and difficulty. In addition, both tasks do in general not permit reliable

verifiers, making it difficult to train on using RLVR. For Rust coding, RLVR is only possible when comprehensive test suites are provided, which is not always possible.

Metrics For Rust coding, we choose a Rust translation of the MBPP dataset (Lab, 2025) and measure the models PASS@1 performance (Chen et al., 2021), in other words, we measure the percentage of solutions that implement a function described in natural language correctly, as measured by a set of unit tests. For text summarization, we measure fluency, consistency, relevance and coherence on the CNN split of Fabbri et al. (2021). We use Gemini 3 Pro (Google DeepMind, 2025) to evaluate the criteria according to a detailed prompt, detailed in Appendix B. To confirm the validity of these results, we compare the Gemini 3 Pro ratings on the dataset of summaries with human annotators provided in Fabbri et al. (2021) and establish a Cohens Kappa of over 60% for each metric. We provide the details in Appendix A. For the mathematical problem solving, we measure the performance of the models on the MATH500 dataset (Lightman et al., 2023). Note that we do not train on this task, and we therefore use it to assess the loss of general reasoning capabilities of the model due to fine-tuning.

All models are trained and merged twice with different seeds and we report averaged results from 10 evaluation runs for Rust, text summarization and MATH500. In all plots we draw corridors indicating the run-to-run variance, spanning the respective minimum and maximum values obtained by each run.

Training Datasets We devise two datasets for training. To ensure that we measure whether the model preserved its generalization capability, we design the training and test datasets to stem from different distributions, while having similar task formats. (1) *Rust coding:* We train the model on a set of single-function coding tasks in the Rust language, where each data point is a pair of task and code solution that were synthesized from an LLM (Team, 2024b). We

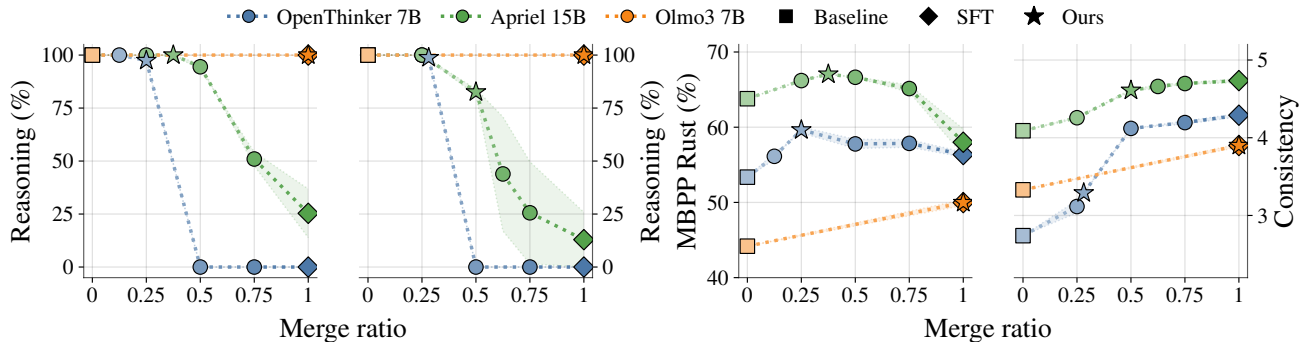


Figure 3. We evaluate all various merge ratios on the final test set. We observe that reasoning on MATH500 tasks drops sharply after a merging threshold is crossed (Left) and target task performance experiences a smoother curve (Right). On Rust, a clear benefit of restored reasoning to task performance is visible.

filter the dataset to the subset of code samples that pass the corresponding test suite consistently, arriving at a subset of 6761 training samples. We confirm the obtained datasets to be disjoint from the MBPP dataset by checking for closest matches using cosine similarity and confirming their distinctness. (2) *Text summarization*: We further train the model on the task of text summarization. As training dataset, we use the reddit TLDR; split of Fabbri et al. (2021), which is intentionally different to the CNN task split that we use for evaluation.

4.2. Main Results

Next, we present our main results on Rust coding and text summarization.

Rust coding We train all models on both Rust coding and text summarization and show their performance on the respective tasks compared to the performance on MATH500 in Figure 2. As can be seen, standard SFT strongly reduces the performance on MATH500: For OpenThinker 7B, it drops from 79% to 34.2% when training on the Rust dataset, and to 2% on the text summarization dataset. While SFT does not disturb the math reasoning capabilities of April 15B, it decreases its performance *at the target task*. Meanwhile, our method consistently picks a strong trade-off. In Rust, the recovered reasoning even leads to improved performance compared to both the starting model (baseline) and the SFT model, with an average increase of 4.1 percentage points.

Text Summarization For text summarization, we observe that fluency, relevance and coherence are already close to the maximal score for all evaluated models, on average 4.84 out of 5, and remain stable across our training. We therefore focus on the metric of consistency and report the performance on other metrics in Appendix A. Here, the best performance is usually achieved by the SFT model, with model performance on average increasing by 1.10 points. This increase

occurs despite as model reasoning decreasing rapidly to 0.0% and 12.9% for OpenThinker 7B and April 15B, respectively. With our method, reasoning is fully recovered, resulting in a parallel recovery of MATH500 performance, while preserving 76.7% and 97.4% of the performance increase for the fully finetuned text consistency.

Runtime and Cost A key benefit of our method is its low cost. Our adaptation method requires on average less than two hours to complete and fits on a single NVIDIA H200 GPU, costing in total less than USD 6 on typical rental services at the time of writing. Most time is split between obtaining the fine-tuned model, with on average 36 minutes, and evaluating the merge ratios in Algorithm 1, with on average 23 min. Note that a key factor in our cost-efficiency is that our method does not require any reasoning rollouts from the model.

Other Methods We attempted reproducing two popular related methods that are also applicable to our setting of training RLMs using only standard SFT data. The first method, VeriFree (Zhou et al., 2025b), requires a significant amount of resources: Based on our initial experiments, we estimated 8 NVIDIA GPUs for a total of 60 hours for a single training run. Given this cost, which applies to each step of a preliminary hyperparameter search for fair comparison, we consider a comparison infeasible. Second, for SDFT (Shenfeld et al., 2026) we have extensively tried to reproduce SDFT on Olmo3 7B, on the task of Rust coding, exploring over 20 hyperparameter combinations. Despite reaching out to the authors, the best performance we could achieve with SDFT was 41.4%, a slight decrease from the baseline performance of 44.2%, while performance on MATH500 also slightly decreased from 95.5% to 93.1%. As such, despite our extensive reproduction efforts, we did not include the above related methods for direct comparisons in our experiments. We report the used hyperparameters in Appendix A.

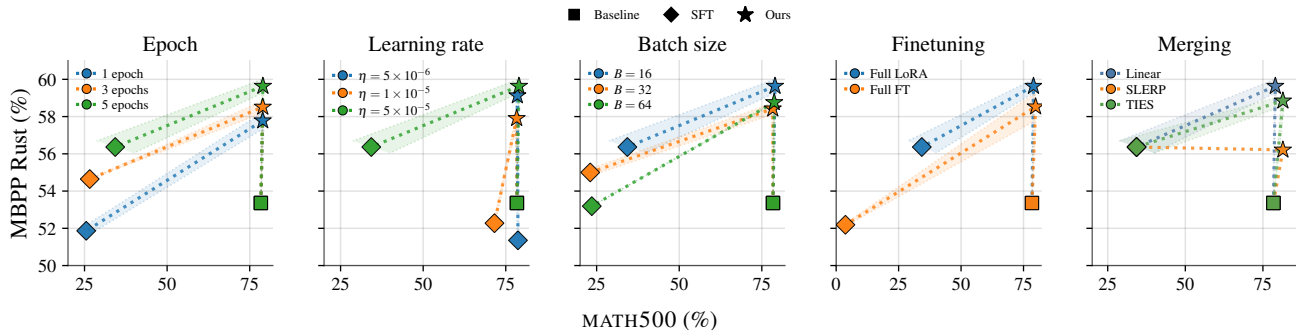


Figure 4. OpenThinker 7B ablations on Rust coding, from left to right: training epochs, learning rate η , batch size B , LoRA vs. full finetuning, and merging techniques. The overall trend of reasoning loss in standard SFT and reasoning and performance recovery is stable across all settings.

4.3. Ablation

We ablate over the choice of merge ratios, employed finetuning techniques and training hyperparameters for the supervised finetuning.

Merge ratio We ablate over the merging factor α on the Rust and text summarization dataset, and show their reasoning rate on MATH500 and the target tasks in Figure 3 (left and right, respectively). We notice that OpenThinker 7B loses its reasoning capability quickly for $\alpha \in [0.25, 0.5]$, while April 15B loses it more slowly and Olmo3 7B even maintains full reasoning at $\alpha = 1$. Surprisingly, the performance on Rust coding appears to peak in for $\alpha \in [0.25, 0.5]$, where the model has still knowledge about the task from training and recovered reasoning. Meanwhile, consistency in the text summarization task increases more linearly with the merging ratio. Our method reliably picks a point close to a good trade-off point between reasoning and task performance optimum, despite picking the point only based on the reasoning score on the calibration dataset.

Merging technique We ablate over the employed merging technique for OpenThinker 7B on Rust coding, presenting the results in the right-most panel of Figure 4. We compare the linear merging employed in our main experiments with Spherical Linear interpolation (SLERP) (Shoemaker, 1985; Goddard et al., 2025) and TIES (Yadav et al., 2023). For TIES, since fixing $\alpha = 1$ and searching over the density parameter would result in complete loss of reasoning even for a density of less than 0.1, we fix density to 0.5 and search over the α parameter. For all merging techniques, our method obtains the best merging ratio at $\alpha = 0.25$. Compared to linear interpolation, SLERP recovers more MATH500 performance but loses target-task performance, while TIES improves both MATH500 and target-task performance, albeit less than our linear merging.

Finetuning Our main method uses LoRA for light weight and efficient finetuning. We ablate over the choice of LoRA by using full finetuning of all model weights, using the same hyperparameters as for LoRA, and present the results in the second panel from left of Figure 4. We observe that the capability loss on MATH500 is much stronger with standard finetuning, and that our method is unable to obtain as much task specific performance as LoRA.

Hyperparameters We ablate over finetuning hyperparameters for the supervised finetuning stage over OpenThinker 7B in Figure 4. We ablate over varying epochs, learning parameters η , and batch sizes B , keeping the remaining hyperparameters fixed. Importantly, across epoch, learning rate and batch size choices, even though some finetunings result in models with worse performance on the target task than the baseline model, the overall trend is consistent that training disturbs the natural reasoning of the model, while merging can recover the MATH500 performance and maintain improvements on the target task.

5. Related Work

Below, we present an overview of prior work related to our method.

Training Reasoning Models with Verifiers The standard technique for training and adapting reasoning models is via reinforcement learning with verifiable rewards (DeepSeek-AI, 2025; Olmo et al., 2025; Stojanovski et al., 2025). This works well for tasks that permit such verifiers, however leaves out many domains, such as the explored domain of text summarization. Moreover, while this technique is well-explored for bootstrapping reasoning model performance, continuous adaptation of RLMs to new tasks via RL is not well explored yet.

Training Reasoning Models without Verifiers Verifier-free or verifier-light methods try to enable training reasoning

models without reasoning traces, by optimizing likelihood of known answers directly (Zhou et al., 2025a;b) or training teacher models to provide feedback and justifications for reference answers (Shenfeld et al., 2026). These methods are closest in motivation to our setting, but they are computationally much more expensive and not designed to work on already trained reasoning models. Our method instead operates on minimal hardware requirements and allows inexpensive, fast task adaptation.

Model merging Model merging combines multiple checkpoints into a single model without additional gradient updates. A common motivation is to combine task-specific skills or mitigate forgetting by merging related models (Alexandrov et al., 2024; Yang et al., 2024b). Prior methods include simple weight averaging, task arithmetic (Ilharco et al., 2023), and sparse or sign-based variants such as TIES (Yadav et al., 2023). Linear interpolation is also connected to linear mode connectivity, where independently trained or fine-tuned networks can lie on low-loss paths in weight space (Frankle et al., 2020). Most merging work aims to combine capabilities from multiple specialized models. Our use is narrower: we merge an SFT checkpoint back with its own original reasoning checkpoint to trade off target-task adaptation against preservation of reasoning behavior. Recent work on tunable reasoning through model merging suggests that interpolation can control the degree of reasoning behavior in language models (Lan et al., 2025). We build on this observation, but select the merge coefficient using a small calibration set and target the practical setting where only input-output SFT data is available.

6. Limitations and Future Work

Composability and Continual learning Our method demonstrates that for single-task adaptation, a simple train-and-merge pipeline is enough to preserve general capabilities and obtain task-specific improvements. However, as LLMs are trained to be generalist and continuously adapted to new tasks, the question remains how to obtain adaptation to several tasks in parallel. We consider this an exciting avenue for future research.

Interpretability In our experiments we observe that models stop reasoning by immediately producing an end-of-reasoning marker. To address this, during the development of the method, we tried prefilling the model response with a start-of-reasoning marker and an immediately following non-end-of-reasoning token like ‘0kay’, however, the reasoning performance was not recovered. An interesting direction for future work is to mechanistically understand how reasoning tendency is encoded in the model. We hope that potential insights could then be directly leveraged to design better adaptation methods for reasoning models, in

particular in terms of recovering reasoning tendency.

Lightweight Adaptation on Reasoning Domains Our experiments focus on domains where the reasoning ability of the model itself is not crucial for achieving top performance, instead, the SFT data already provides the necessary signal for the model to perform well. This is still crucial, especially when patching the model for knowledge gaps; our method provides a lightweight way to insert new information into reasoning models simply by collecting instruction-completion pairs—in the same paradigm as in the pre LRM era. However, on domains where the reasoning gains themselves define performance improvements (e.g., mathematics), final-output-based methods such as ours could fall short. Indeed, in preliminary experiments we have tried to apply our method to mathematical proofs, but we failed to improve the models’ performance—likely due to the lack of supervision signal improving the reasoning process of the model itself.

7. Conclusion

We studied how to adapt reasoning language models when only ordinary input-output supervision is available. Our results show that this failure mode is not merely cosmetic. In several settings, the loss of reasoning coincides with large loss of performance on held-out mathematical reasoning and weaker generalization on the target task. We introduced a simple two-step pipeline to mitigate this issue: first fine-tune the reasoning model with standard SFT, then merge the fine-tuned checkpoint back with the original model. The merge ratio is selected using only a small calibration set and the observed rate of non-empty reasoning traces.

Across Rust coding and text summarization, this procedure recovers most or all of the reasoning behavior lost under SFT while retaining much of the target-task improvement, and does so without a verifier, a reward model, or a stronger teacher model, at a cost comparable to running the initial SFT and a small grid search. Overall, our results indicate that RLMs can be efficiently adapted to target tasks using standard model training techniques.

References

- Alexandrov, A., Raychev, V., Müller, M. N., Zhang, C., Vechev, M. T., and Toutanova, K. Mitigating catastrophic forgetting in language transfer via model merging. In *EMNLP (Findings)*, 2024.
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. Concrete problems in ai safety, 2016. URL <https://arxiv.org/abs/1606.06565>.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto,

- 440 H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N.,
441 Brockman, G., et al. Evaluating Large Language Models
442 Trained on Code. *arXiv Preprint*, 2021. URL <https://arxiv.org/abs/2107.03374>.
443
- 444 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
445
- 446 Fabbri, A. R., Kryściński, W., McCann, B., Xiong, C.,
447 Socher, R., and Radev, D. Summeval: Re-evaluating
448 summarization evaluation, 2021. URL <https://arxiv.org/abs/2007.12626>.
449
- 450 Frankle, J., Dziugaite, G. K., Roy, D. M., and Carbin, M.
451 Linear mode connectivity and the lottery ticket hypothesis,
452 2020. URL <https://arxiv.org/abs/1912.05671>.
453
- 454 Goddard, C., Siriwardhana, S., Ehghaghi, M., Meyers, L.,
455 Karpukhin, V., Benedict, B., McQuade, M., and Solawetz,
456 J. Arcee’s mergekit: A toolkit for merging large language
457 models, 2025. URL <https://arxiv.org/abs/2403.13257>.
458
- 459 Google DeepMind. Gemini Pro, 2025. URL [https://
460 deepmind.google/technologies/gemini/pro/](https://deepmind.google/technologies/gemini/pro/).
461
- 462 Guha, E., Marten, R., Keh, S., Raoof, N., Smyrnis, G.,
463 Bansal, H., Nezhurina, M., Mercat, J., Vu, T., Sprague,
464 Z., Suvarna, A., Feuer, B., Chen, L., Khan, Z., Frankel,
465 E., Grover, S., Choi, C., Muennighoff, N., Su, S., Zhao,
466 W., Yang, J., Pimpalgaonkar, S., Sharma, K., Ji, C. C.-J.,
467 Deng, Y., Pratt, S., Ramanujan, V., Saad-Falcon, J., Li,
468 J., Dave, A., Albalak, A., Arora, K., Wulfe, B., Hegde,
469 C., Durrett, G., Oh, S., Bansal, M., Gabriel, S., Grover,
470 A., Chang, K.-W., Shankar, V., Gokaslan, A., Merrill,
471 M. A., Hashimoto, T., Choi, Y., Jitsev, J., Heckel, R.,
472 Sathiamoorthy, M., Dimakis, A. G., and Schmidt, L.
473 Openthoughts: Data recipes for reasoning models, 2025.
474 URL <https://arxiv.org/abs/2506.04178>.
475
- 476 Hu, E. J., yelong shen, Wallis, P., Allen-Zhu, Z., Li, Y.,
477 Wang, S., Wang, L., and Chen, W. LoRA: Low-rank
478 adaptation of large language models. In *International
479 Conference on Learning Representations*, 2022. URL
480 <https://openreview.net/forum?id=nZeVKeeFYf9>.
481
- 482 Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan,
483 S., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing
484 models with task arithmetic, 2023. URL <https://arxiv.org/abs/2212.04089>.
485
- 486 KRAFTON and SKT. Continual post-training of
487 llms via offline grpo for mathematical reasoning,
488 2025. URL [https://krafton-ai.github.io/blog/
489 llm_post_training_en/](https://krafton-ai.github.io/blog/llm_post_training_en/). Accessed: 2025-07-28.
490
- 491 Lab, N. U. P. R. MultiPL-e (revision 40ca145),
492 2025. URL [https://huggingface.co/datasets/
493 nuprl/MultiPL-E](https://huggingface.co/datasets/nuprl/MultiPL-E).
494
- Lan, X., Zheng, Y., Cao, S., and Li, Y. The thinking spectrum: An empirical study of tunable reasoning in llms through model merging, 2025. URL <https://arxiv.org/abs/2509.22034>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step, 2023. URL <https://arxiv.org/abs/2305.20050>.
- MacDiarmid, M., Wright, B., Uesato, J., Benton, J., Kutasov, J., Price, S., Bouscal, N., Bowman, S., Bricken, T., Cloud, A., Denison, C., Gasteiger, J., Greenblatt, R., Leike, J., Lindsey, J., Mikulik, V., Perez, E., Rodrigues, A., Thomas, D., Webson, A., Ziegler, D., and Hubinger, E. Natural emergent misalignment from reward hacking in production rl, 2025. URL <https://arxiv.org/abs/2511.18397>.
- NovaSky Team. Sky-t1: Train your own o1 preview model within \$450. <https://novasky-ai.github.io/posts/sky-t1>, 2025. Accessed: 2025-01-09.
- Olmo, T., Ettinger, A., Bertsch, A., Kuehl, B., Graham, D., Heineman, D., Groeneveld, D., Brahman, F., Timbers, F., Ivison, H., Morrison, J., Poznanski, J., Lo, K., Soldaini, L., Jordan, M., Chen, M., Noukhovitch, M., Lambert, N., Walsh, P., Dasigi, P., Berry, R., Malik, S., Shah, S., Geng, S., Arora, S., Gupta, S., Anderson, T., Xiao, T., Murray, T., Romero, T., Graf, V., Asai, A., Bhagia, A., Wettig, A., Liu, A., Rangapur, A., Anastasiades, C., Huang, C., Schwenk, D., Trivedi, H., Magnusson, I., Lochner, J., Liu, J., Miranda, L. J. V., Sap, M., Morgan, M., Schmitz, M., Guerquin, M., Wilson, M., Huff, R., Bras, R. L., Xin, R., Shao, R., Skjonsberg, S., Shen, S. Z., Li, S. S., Wilde, T., Pyatkin, V., Merrill, W., Chang, Y., Gu, Y., Zeng, Z., Sabharwal, A., Zettlemoyer, L., Koh, P. W., Farhadi, A., Smith, N. A., and Hajishirzi, H. Olmo 3, 2025. URL <https://arxiv.org/abs/2512.13961>.
- OpenAI, :, Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., Iftimie, A., Karpenko, A., Passos, A. T., Neitz, A., Prokofiev, A., Wei, A., Tam, A., Bennett, A., Kumar, A., Saraiva, A., Vallone, A., Duberstein, A., Kondrich, A., Mishchenko, A., Applebaum, A., Jiang, A., Nair, A., Zoph, B., Ghorbani, B., Rossen, B., Sokolowsky, B., Barak, B., McGrew, B., Minaiev, B., Hao, B., Baker, B., Houghton, B., McKinzie, B., Eastman, B., Lugaresi, C., Bassin, C., Hudson, C., Li, C. M., de Bourcy, C., Voss, C., Shen, C., Zhang, C., Koch, C., Orsinger, C., Hesse, C., Fischer, C., Chan, C., Roberts, D., Kappler, D., Levy,

- 495 D., Selsam, D., Dohan, D., Farhi, D., Mely, D., Robinson,
 496 D., Tsipras, D., Li, D., Oprica, D., Freeman, E., Zhang,
 497 E., Wong, E., Proehl, E., Cheung, E., Mitchell, E., Wal-
 498 lace, E., Ritter, E., Mays, E., Wang, F., Such, F. P., Raso,
 499 F., Leoni, F., Tsimpourlas, F., Song, F., von Lohmann,
 500 F., Sulit, F., Salmon, G., Parascandolo, G., Chabot, G.,
 501 Zhao, G., Brockman, G., Leclerc, G., Salman, H., Bao,
 502 H., Sheng, H., Andrin, H., Bagherinezhad, H., Ren, H.,
 503 Lightman, H., Chung, H. W., Kivlichan, I., O’Connell,
 504 I., Osband, I., Gilaberte, I. C., Akkaya, I., Kostrikov, I.,
 505 Sutskever, I., Kofman, I., Pachocki, J., Lennon, J., Wei,
 506 J., Harb, J., Twore, J., Feng, J., Yu, J., Weng, J., Tang, J.,
 507 Yu, J., Candela, J. Q., Palermo, J., Parish, J., Heidecke,
 508 J., Hallman, J., Rizzo, J., Gordon, J., Uesato, J., Ward,
 509 J., Huizinga, J., Wang, J., Chen, K., Xiao, K., Singhal,
 510 K., Nguyen, K., Cobbe, K., Shi, K., Wood, K., Rimbach,
 511 K., Gu-Lemberg, K., Liu, K., Lu, K., Stone, K., Yu, K.,
 512 Ahmad, L., Yang, L., Liu, L., Maksin, L., Ho, L., Fedus,
 513 L., Weng, L., Li, L., McCallum, L., Held, L., Kuhn, L.,
 514 Kondraciuk, L., Kaiser, L., Metz, L., Boyd, M., Trebacz,
 515 M., Joglekar, M., Chen, M., Tintor, M., Meyer, M., Jones,
 516 M., Kaufer, M., Schwarzer, M., Shah, M., Yatbaz, M.,
 517 Guan, M. Y., Xu, M., Yan, M., Glaese, M., Chen, M.,
 518 Lampe, M., Malek, M., Wang, M., Fradin, M., McClay,
 519 M., Pavlov, M., Wang, M., Wang, M., Murati, M., Bavar-
 520 ian, M., Rohaninejad, M., McAleese, N., Chowdhury,
 521 N., Chowdhury, N., Ryder, N., Tezak, N., Brown, N.,
 522 Nachum, O., Boiko, O., Murk, O., Watkins, O., Chao, P.,
 523 Ashbourne, P., Izmailov, P., Zhokhov, P., Dias, R., Arora,
 524 R., Lin, R., Lopes, R. G., Gaon, R., Miyara, R., Leike, R.,
 525 Hwang, R., Garg, R., Brown, R., James, R., Shu, R., Cheu,
 526 R., Greene, R., Jain, S., Altman, S., Toizer, S., Toyer, S.,
 527 Miserendino, S., Agarwal, S., Hernandez, S., Baker, S.,
 528 McKinney, S., Yan, S., Zhao, S., Hu, S., Santurkar, S.,
 529 Chaudhuri, S. R., Zhang, S., Fu, S., Papay, S., Lin, S., Bal-
 530 aji, S., Sanjeev, S., Sidor, S., Broda, T., Clark, A., Wang,
 531 T., Gordon, T., Sanders, T., Patwardhan, T., Sottiaux, T.,
 532 Degry, T., Dimson, T., Zheng, T., Garipov, T., Stasi, T.,
 533 Bansal, T., Creech, T., Peterson, T., Eloundou, T., Qi, V.,
 534 Kosaraju, V., Monaco, V., Pong, V., Fomenko, V., Zheng,
 535 W., Zhou, W., McCabe, W., Zaremba, W., Dubois, Y., Lu,
 536 Y., Chen, Y., Cha, Y., Bai, Y., He, Y., Zhang, Y., Wang,
 537 Y., Shao, Z., and Li, Z. Openai o1 system card, 2024a.
 538 URL <https://arxiv.org/abs/2412.16720>.
- 539
- 540 OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L.,
 541 Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J.,
 542 Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Bal-
 543 aji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M.,
 544 Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G.,
 545 Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman,
 546 A.-L., Brockman, G., Brooks, T., Brundage, M., Button,
 547 K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson,
 548 C., Carmichael, R., Chan, B., Chang, C., Chantzis, F.,
 549 Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess,
 B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Cur-
 rier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N.,
 Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning,
 S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus,
 L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L.,
 Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G.,
 Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S.,
 Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han,
 J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse,
 C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B.,
 Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S.,
 Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S.,
 Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A.,
 Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L.,
 Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J.,
 Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich,
 A., Konstantinidis, A., Kopic, K., Krueger, G., Kuo, V.,
 Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D.,
 Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T.,
 Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning,
 S., Markov, T., Markovski, Y., Martin, B., Mayer, K.,
 Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C.,
 McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick,
 J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V.,
 Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O.,
 Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan,
 A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki,
 J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo,
 G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng,
 A., Perelman, A., de Avila Belbute Peres, F., Petrov, M.,
 de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M.,
 Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E.,
 Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C.,
 Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H.,
 Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry,
 G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D.,
 Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam,
 P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K.,
 Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such,
 F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N.,
 Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E.,
 Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone,
 A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang,
 J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann,
 C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wi-
 ethoff, M., Willner, D., Winter, C., Wolrich, S., Wong,
 H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu,
 T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R.,
 Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J.,
 Zhuk, W., and Zoph, B. Gpt-4 technical report, 2024b.
 URL <https://arxiv.org/abs/2303.08774>.
- Prakash, J. and Buvanesh, A. What can you do when you

- 550 have zero rewards during rl? *CoRR*, 2025.
- 551 Radhakrishna, S., Parikh, S., Sarda, G., Turkkkan, A., Vohra,
- 552 Q., Li, R., Jhamb, D., Ogueji, K., Shukla, A., Bamgbose,
- 553 O., Liang, T., Kumar, L., Ostapenko, O., Malay, S. K. R.,
- 554 Tiwari, A., Bogavelli, T., Yadav, V., Mehta, J., Mittal, S.,
- 555 Kalkunte, A., Pattnaik, P., Slimi, K., Sreeram, A., Nair,
- 556 J., Oladipo, A., Maiya, S., Mahajan, K., Maheshwary,
- 557 R., Hashemi, M., Mudumba, S. R., Madhusudhan, S. T.,
- 558 Scholak, T., Paquet, S., Davasam, S., and Sunkara, S.
- 559 Apriel-nemotron-15b-thinker, 2025.
- 560
- 561 Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X.,
- 562 Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo,
- 563 D. Deepseekmath: Pushing the limits of mathematical
- 564 reasoning in open language models, 2024. URL
- 565 <https://arxiv.org/abs/2402.03300>.
- 566
- 567 Shenfeld, I., Damani, M., Hübotter, J., and Agrawal, P.
- 568 Self-distillation enables continual learning, 2026. URL
- 569 <https://arxiv.org/abs/2601.19897>.
- 570
- 571 Shoemake, K. Animating rotation with quaternion curves. In
- 572 *Proceedings of the 12th Annual Conference on Computer*
- 573 *Graphics and Interactive Techniques*, SIGGRAPH '85,
- 574 pp. 245–254, New York, NY, USA, 1985. Association
- 575 for Computing Machinery. ISBN 0897911660. doi: 10.
- 576 [1145/325334.325242](https://doi.org/10.1145/325334.325242). URL [https://doi.org/10.1145/](https://doi.org/10.1145/325334.325242)
- 577 [325334.325242](https://doi.org/10.1145/325334.325242).
- 578
- 579 Stojanovski, Z., Stanley, O., Sharratt, J., Jones, R., Ade-
- 580 foye, A., Kaddour, J., and Köpf, A. Reasoning gym:
- 581 Reasoning environments for reinforcement learning with
- 582 verifiable rewards, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2505.24760)
- 583 [2505.24760](https://arxiv.org/abs/2505.24760).
- 584
- 585 Team, O. Rust verified code dataset. [https://www.oxen.](https://www.oxen.ai/ox/Rust)
- 586 [ai/ox/Rust](https://www.oxen.ai/ox/Rust), 2024a. Accessed: 2025-01-01.
- 587
- 588 Team, O. Training a rust 1.5b coder lm with rein-
- 589 forcement learning (grp). [https://ghost.oxen.ai/](https://ghost.oxen.ai/training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grp/)
- 590 [training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grp/](https://ghost.oxen.ai/training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grp/),
- 591 2024b. Technical report.
- 592
- 593 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
- 594 L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention
- 595 is all you need, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1706.03762)
- 596 [1706.03762](https://arxiv.org/abs/1706.03762).
- 597
- 598 Yadav, P., Tam, D., Choshen, L., Raffel, C., and Bansal,
- 599 M. TIES-merging: Resolving interference when merg-
- 600 ing models. In *Thirty-seventh Conference on Neural*
- 601 *Information Processing Systems*, 2023. URL [https://](https://openreview.net/forum?id=xtaX3WyCj1)
- 602 openreview.net/forum?id=xtaX3WyCj1.
- 603
- 604 Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu,
- 605 B., Li, C., Liu, D., Huang, F., Wei, H., et al. Qwen2.5
- 606 Technical Report. *arXiv Preprint*, 2024a. URL [https://](https://doi.org/10.48550/arXiv.2412.15115)
- 607 doi.org/10.48550/arXiv.2412.15115.
- 608
- 609 Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng,
- 610 B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu,
- 611 D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin,
- 612 H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang,
- 613 J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang,
- 614 K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang,
- 615 P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo,
- 616 S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang,
- 617 X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan,
- 618 Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., and
- 619 Qiu, Z. Qwen3 technical report, 2025. URL [https://](https://arxiv.org/abs/2505.09388)
- 620 arxiv.org/abs/2505.09388.
- 621
- 622 Yang, E., Shen, L., Guo, G., Wang, X., Cao, X., Zhang, J.,
- 623 and Tao, D. Model merging in llms, mllms, and beyond:
- 624 Methods, theories, applications and opportunities, 2024b.
- 625 URL <https://arxiv.org/abs/2408.07666>.
- 626
- 627 Zelikman, E., Wu, Y., Mu, J., and Goodman, N. STar:
- 628 Bootstrapping reasoning with reasoning. In Oh, A. H.,
- 629 Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances*
- 630 *in Neural Information Processing Systems*, 2022. URL
- 631 https://openreview.net/forum?id=_3ELRdg2sgI.
- 632
- 633 Zhou, X., Liu, Z., Sims, A., Wang, H., Pang, T., Li, C.,
- 634 Wang, L., Lin, M., and Du, C. Reinforcing general rea-
- 635 soning without verifiers. *CoRR*, 2025a.
- 636
- 637 Zhou, X., Liu, Z., Sims, A., Wang, H., Pang, T., Li,
- 638 C., Wang, L., Lin, M., and Du, C. Reinforcing gen-
- 639 eral reasoning without verifiers, 2025b. URL [https://](https://arxiv.org/abs/2505.21493)
- 640 arxiv.org/abs/2505.21493.

Table 1. SummEval correlations for Gemini-3-pro-preview.

Metric	Spearman ρ	p-value
Fluency	0.6034	6.29×10^{-16}
Relevance	0.5616	1.69×10^{-14}
Coherence	0.6785	8.31×10^{-22}
Consistency	0.6791	7.44×10^{-22}

A. Experimental Details, Ablations and Case Study

In this section, we provide additional details about the implementation, hyperparameters, datasets, and so on.

A.1. Refinement of the Rust dataset

In our study, the Rust corpus from Team (2024a;b) is used for supervised fine-tuning and in-distribution evaluation. Each sample in the corpus consists of four fields: a task identifier, a natural-language Rust prompt that describes the target function, a Rust code implementation that solves the task, and a list of executable test cases. Consequently, the dataset provides complete function-level programming tasks paired with verification harnesses, enabling evaluation of generated code correctness.

Before using the dataset, we performed verification and filtering to improve data quality and ensure experimental reproducibility. For each code-test pair, we compiled the provided Rust implementation and executed its associated test suite ten times in our local environment. We retained only samples for which the implementation successfully compiled and passed the complete test suite in all ten runs. This repeated-execution protocol filters out examples whose correctness is unstable across executions, including cases affected by nondeterministic behavior or code-test pairs that only pass stochastically. We did not manually repair such cases; any sample that failed compilation or testing in at least one run was excluded from the final dataset.

Starting from the original 7,554 Rust tasks, this filtering process retained 7,511 valid samples. We then split the filtered dataset into 6,761 examples for supervised fine-tuning and 750 examples for validation. Each entry in the final split retains the same four-field structure: task identifier (`task_id`), Rust prompt (`rust_prompt`), Rust implementation (`rust_code`), and the complete list of test cases (`rust_test_list`).

A.2. Validation of Gemini 3 Pro as Text Summarization Judge

As mentioned in §4, we validated the use of Gemini 3 Pro Preview as LLM-as-a-Judge metric for our evaluation of the Text Summarization task. For this, we run the model on the SummEval dataset (Fabbri et al., 2021) and compare the model scores with the human annotations. We present the results in Table 1, showing that the model achieves generally high correlation of around 60% on all scores, and almost 70% in the main inspected metric of consistency.

A.3. SFT Hyperparameters

For the main experiments, we tune the SFT hyperparameters separately for each model and dataset combination. The resulting configurations are shown in Table 2. The main runs use LoRA SFT with the target modules described in §3. We also include the full-finetuning configuration used for the OpenThinker 7B Rust ablation.

For the SDFT (Shenfeld et al., 2026) training result, we perform 26 runs with the following SFT settings: epochs $\in \{1, 2, 3, 5\}$, learning rate $\eta \in \{1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}, 5 \times 10^{-6}\}$, batch size $B \in \{16, 32, 64\}$, and EMA coefficient $\alpha \in \{0.01, 0.02, 0.05\}$. We set `max_token_length` = 4096 and discarded training samples exceeding this budget, leaving roughly one-third of the examples. The best-performing run used epoch 2, learning rate 5×10^{-6} , batch size 64, and EMA $\alpha = 0.05$, and completed in about five hours using 4 NVIDIA H200 GPUs. The full four-panel text-summarization tradeoff is shown in Figure 5.

Table 2. SFT hyperparameters used for the main model and dataset combinations, plus the full-finetuning ablation. Batch size denotes the effective training batch size.

Dataset	Model	Method	Epochs	Learning rate	Batch size
Rust coding	OpenThinker 7B	LoRA	5	5×10^{-5}	16
Rust coding	OpenThinker 7B	Full FT	5	5×10^{-5}	16
Rust coding	Apriel 15B	LoRA	5	1×10^{-4}	64
Rust coding	Olmo3 7B	LoRA	3	5×10^{-5}	16
Text summarization	OpenThinker 7B	LoRA	1	5×10^{-5}	16
Text summarization	Apriel 15B	LoRA	1	1×10^{-4}	64
Text summarization	Olmo3 7B	LoRA	1	1×10^{-5}	16

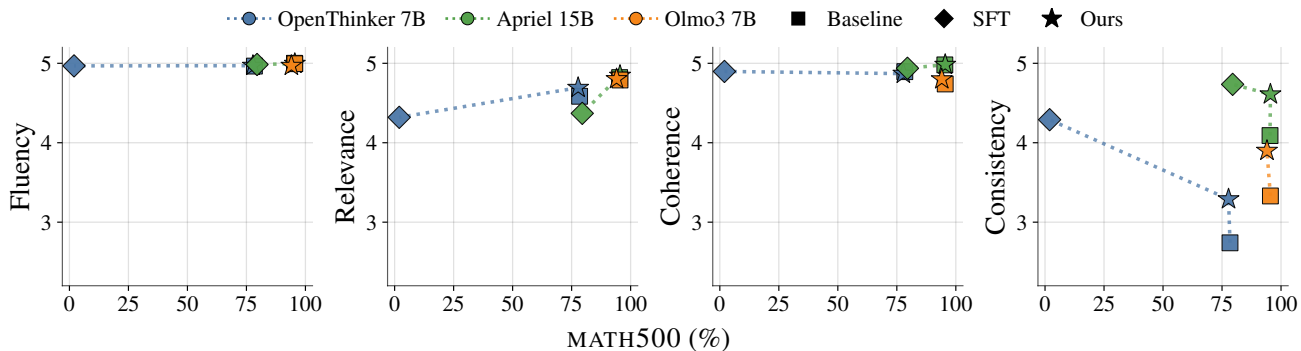


Figure 5. Text summarization tradeoff from SFT/merging experiments: Fluency, Relevance, Coherence, and Consistency against MATH-500.

A.4. Experimental Details on other Text Summarization metrics

In Figure 5, we demonstrate the performance of the trained models across all four metrics of fluency, relevance, coherence and consistency. As can be seen, the performance in the former three categories is already high in the baseline model, and remains stable across our training. The strongest change outside of consistency is due to a loss of relevance in SFT models, which is fully recovered by our method.

B. Prompts

We use one task prompt for Rust code generation, one task prompt for text summarization, and four judge prompts for the text summarization evaluation. Figure 6 shows the Rust prompt used for both SFT training and MBPP Rust evaluation. Figure 8 shows the text summarization prompt used for both SFT training on the TLDR split and evaluation on the CNN split. For the LLM-based text summarization metrics, Gemini 3 Pro is prompted separately for each criterion: Figure 9 for relevance, Figure 10 for fluency, Figure 11 for factual consistency, and Figure 12 for coherence. In this section, we detail all prompts used for the respective models and tasks. Each prompt is shown in the same minipage-based figure format as the prompt examples in the original template.

```
715
716
717
718
719 1 user
720 2   You are a helpful coding assistant producing high-quality Rust code.
721 3   Strictly follow the instruction below to complete the function specified in the
722 4   instruction.
723 5   Your response should include all dependencies, headers and function declaration
724 6   to be directly usable (even for the ones seen in the given part).
725 7   You should NOT call or test the function in your response.
726 8   Output your complete implementation in a single code block wrapped in triple
727 9   backticks with `rust` specified, like this:
728 10  ```rust
729 11  // your function here
730 12  ```
731 13  Instruction:
732 14  {rust_prompt}
```

732 *Figure 6.* Prompt used for both training and evaluation in the Rust code generation task. The task instruction varies per example and is
733 inserted as {rust_prompt}.

734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769

```
1 user
2 {problem}
3 Please reason step by step, and put your final answer within \boxed{}
```

Figure 7. Prompt used for MATH-500 evaluation. The problem is inserted as {problem}.

```
1 user
2 You are an expert in writing summarization.
3 Your task is to read the following Article and write a summary about it.
4
5 Output your complete summary after the <SUMMARY> tag, like this:
6 <SUMMARY>
7 // your summary here
8
9 Article:
10 {article}
```

Figure 8. Prompt used for both training and evaluation in the text summarization task. The article is inserted as {article}.

```

770
771
772
773
774
775
776 1 user
777 2   You are a helpful assistant in evaluating the quality of a summary.
778 3   You will be given a news article and a summary written for that article.
779 4   Your task is to evaluate the *relevance* of the summary.
780 5
781 6   Definition of Relevance:
782 7   Relevance measures how well the summary captures the important information from
783 8   the article. A relevant summary includes the most important main idea of the
784 9   article and avoids redundancies and minor, trivial, or unrelated details.
785 10
786 11   Evaluation Criteria (Relevance: 1-100)
787 12
788 13   80-100 -- The summary captures the key idea of the article accurately and
789 14   completely. It focuses on the essential information and avoids any
790 15   redundancies and unnecessary or irrelevant content.
791 16
792 17   40-79 -- The summary includes some important information but only partially
793 18   capture the main idea. It misses key points or includes some redundancies
794 19   or minor details.
795 20
796 21   1-39 -- The summary fails to capture the main idea of the article and completely
797 22   deviate from the main idea. It may focus on unimportant details or
798 23   irrelevant content or omit critical points.
799 24
800 25   Evaluation Steps:
801 26   1. Read and understand the article.
802 27   2. Identify the article's main idea and secondary details.
803 28   3. Read the summary and judge if it captures the main idea of the article.
804 29   4. Identify if there are any secondary details or redundancy in the summary.
805 30   5. Assign a score from 1 to 100 based on the Evaluation Criteria above.
806 31
807 32   Output your detailed thought process and formal justification based on the
808 33   Evaluation Criteria, and finally output your final score in the format shown
809 34   below:
810 35
811 36   <think>
812 37   // your thought process and justification
813 38   </think>
814 39
815 40   <Final score>
816 41   Relevance: // your final score
817 42
818 43   Article:
819 44   {{ARTICLE}}
820 45
821 46   Summary:
822 47   {{SUMMARY}}
823
824

```

Figure 9. Prompt used to evaluate summary relevance with Gemini 3 Pro. The article and generated summary are inserted as {{ARTICLE}} and {{SUMMARY}}.

825
826
827
828
829
830
831
832
833
834
835 1 user
836 2 You are a helpful assistant in evaluating the quality of a summary.
837 3 You will be given a news article and a summary written for that article.
838 4 Your task is to evaluate the *fluency* of the summary.
839 5
840 6 Definition of Fluency:
841 7 Fluency measures how easy the summary is to read.
842 8 All sentences in a fluent summary need to be readable and natural.
843 9 Minor grammatical, formatting, capitalization, or tokenization issues should
844 10 NOT be heavily penalized as long as they do not make the text difficult to read
845 11 or understand.
846 12
847 13 Evaluation Criteria (Fluency: 1-5)
848 14
849 15 5 -- The summary is easy to read and all the sentences are understandable and
850 16 natural. Minor issues such as awkward wording, missing capitalization, or
851 17 tokenization artifacts are acceptable if they do not hinder readability.
852 18
853 19 3 -- The summary is readable but some of sentences include awkward phrasing,
854 20 inconsistent grammar, or formatting problems that reduce readability.
855 21
856 22 1 -- The summary is difficult to read. It contains frequent grammatical errors,
857 23 broken or incomplete sentences, or severe formatting problems that
858 24 significantly hinder understanding.
859 25
860 26 Evaluation Steps:
861 27 1. Read and understand the article.
862 28 2. Read the summary and judge whether the text is easy to read and whether the
863 29 individual sentences are natural.
864 30 3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
865 31
866 32 Only output the score. Do not include any additional explanations or text.
867 33 Use the following format:
868 34 Fluency: <score>
869 35
870 36 Article:
871 37 {{ARTICLE}}
872 38
873 39 Summary:
874 40 {{SUMMARY}}

Figure 10. Prompt used to evaluate summary fluency with Gemini 3 Pro.

880
881
882
883
884
885
886
887
888
889
890
891 1 user
892 2 You are a helpful assistant in evaluating the quality of a summary.
893 3 You will be given a news article and a summary written for that article.
894 4 Your task is to evaluate the *consistency* of the summary.
895 5
896 6 Definition of Consistency:
897 7 Consistency measures how factually aligned the summary is with the article.
898 8 A consistent summary should not introduce information that contradicts the
899 9 article or contain hallucinated statements not supported by the source article.
900 10
901 11 Evaluation Criteria (Consistency: 1-5)
902 12
903 13 5 -- All statements in the summary are fully supported by the article. No
904 14 contradictions, distortions, or hallucinated details appear.
905 15
906 16 3 -- The summary is mostly consistent but includes minor inaccuracies, unclear
907 17 references, or slightly misleading phrasing. These issues do not
908 18 significantly alter the meaning.
909 19
910 20 1 -- The summary contains factual errors or statements that contradict the
911 21 article or introduce unsupported information.
912 22
913 23 Evaluation Steps:
914 24 1. Read and understand the article.
915 25 2. Read the summary and check whether each fact is supported by the article.
916 26 3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
917 27
918 28 Only output the score. Do not include any additional explanations or text.
919 29 Use the following format:
920 30 Consistency: <score>
921 31
922 32 Article:
923 33 {{ARTICLE}}
924 34
925 35 Summary:
926 36 {{SUMMARY}}

Figure 11. Prompt used to evaluate summary consistency with Gemini 3 Pro.

935
936
937
938
939
940
941
942
943
944
945
1 user
946 2 You are a helpful assistant in evaluating the quality of a summary.
947 3 You will be given a news article and a summary written for that article.
948 4 Your task is to evaluate the *coherence* of the summary.
949 5
950 6 Definition of Coherence:
951 7 Coherence measures how well the ideas in the summary fit together.
952 8 A coherent summary presents information in a logical order, with smooth
953 9 transitions between sentences. It should read as a connected, well-structured
10 whole.
11
954 12 Evaluation Criteria (Coherence: 1-5)
955 13
956 14 5 -- The summary is clearly organized and easy to follow. Sentences flow
957 15 naturally, and ideas progress in a logical order.
958 16
959 17 3 -- The summary is somewhat coherent but has noticeable issues in flow or
960 18 structure. Some sentences feel disconnected or out of place, yet the overall
961 19 meaning is still clear.
962 20
963 21 1 -- The summary is hard to follow. The sentences are out of order and loosely
964 22 connected, which makes the summary feel fragmented.
965 23
966 24 Evaluation Steps:
967 25 1. Read and understand the article.
968 26 2. Read the summary and assess whether the ideas are presented in a clear and
969 27 logical order.
970 28 3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
971 29
972 30 Only output the score. Do not include any additional explanations or text.
973 31 Use the following format:
974 32 Coherence: <score>
975 33
976 34 Article:
977 35 {{ARTICLE}}
978 36
979 37 Summary:
980 38 {{SUMMARY}}
981
982
983
984
985
986
987
988
989

Figure 12. Prompt used to evaluate summary coherence with Gemini 3 Pro.