
Leveraging Instruction Tuning and Merging for Reasoning Model Adaptation

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

16 1 Introduction

17 Reasoning language models (RLMs) have changed the frontier of language model capabilities.
18 These models obtain strong results on tasks such as mathematics and programming [4, 19, 32]. The
19 dominant training recipes behind these gains rely on reinforcement learning with verifiable rewards
20 [4, 9, 18, 22]. However, this approach requires automatic verifiers that can automatically determine
21 whether the proposed answer is objectively correct or not. This leaves a much wider class of tasks
22 in an awkward position. Many domains do not provide reliable verifiers, such as text summarization
23 or coding without available unit tests.

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

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

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

44 **Our contributions** Our key contributions are:

- 45 • We identify and study a practical adaptation setting for RLMs where only input-output
 46 supervision is available, without verified reasoning traces.
- 47 • We propose a lightweight SFT-and-merge method that adapts an RLM while selecting the
 48 merge ratio only from calibration reasoning behavior.
- 49 • We evaluate the method across three RLMs and two target tasks, showing that it preserves
 50 general reasoning capabilities while retaining target-task improvements.

51 2 Background

52 We outline the necessary background relevant to this work.

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

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

57 where $\mathbf{c}_{<t} = (c_1, \dots, c_{t-1})$ represents the preceding tokens. During inference, we split the context
 58 \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
 59 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*
 60 and \mathbf{y} as the *answer*. Pre-trained Large Language Models (LLMs) are such models with billions of
 61 parameters, trained on trillions of tokens of training data [4, 20, 32]. By such pre-training the models
 62 acquire a variety of general skills, in particular language-understanding. However, in practice it is
 63 often necessary to adapt the model for particular skills or specialized tasks.

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

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

68 **Large Reasoning Models** Recently models are trained on challenging tasks using reinforcement
 69 learning, resulting in Large Reasoning Models [4, 19]. The models are prompted to solve a task that
 70 permits a reliably verifiable answer, such as a numerical result or executable code. The model is then
 71 trained to prefer answers that are validated as correct by the respective verifier. Before producing
 72 the answer, the model may generate text in an internal scratch pad, which is referred to as reasoning
 73 trace. Empirically, this results in strong solutions for challenging math and code problems [4, 23].

74 However, there are three limitations. First, self-bootstrapping only works if the model has a non-zero
 75 solve rate on the dataset in question [21]. Second, reinforcement learning is an expensive process,
 76 that requires several rollouts and many samples [4, 18, 23]. Third this approach requires a reliable
 77 verifier. If the verifier can be influenced by spurious correlations, the reasoning model can learn to
 78 exploit these correlations and produce undesired outputs [2, 4, 16]. There are numerous tasks for
 79 which such a verifier is not available, for example text summarization.

80 **Reasoning Distillation** An alternative to reinforcement learning is reasoning distillation. In this
81 approach there is a training dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$. For each task \mathbf{x} in the dataset, a reasoning trace
82 \mathbf{r} and answer \mathbf{y}' are obtained by sampling an already reasoning-trained RLM, filtering the solutions
83 such that \mathbf{y}' matches the desired solution \mathbf{y} . This is usually done using RLMs that are larger than
84 the target model, or to save the cost and effort associated with reinforcement learning [4, 17]. The
85 drawback is that such an RLM with better performance on the desired task has to be available.

86 **Model Merging** Model merging combines two or more model checkpoints into a single check-
87 point without additional optimization. For two sets of compatible parameters θ_1 and θ_2 , a common
88 form is a convex interpolation

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

89 This can be used for transfer, specialization, and recovering forgotten behavior across related tasks
90 or fine-tuning regimes. Prior work connects this idea to task arithmetic and task arithmetic-like
91 combinations, and shows that interpolation can produce usable points along low-loss trajectories
92 between related checkpoints [1, 6, 11, 33].

93 3 Training RLMs via SFT and Model Merging

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

96 **Adapting RLM** In this work we focus on task adaptation, concretely: How to improve the perfor-
97 mance of an existing RLM on a target task. The main approach for this is to continue reinforcement
98 learning with verifiable rewards [4, 12], which requires reliable verifiers. Alternatively, developers
99 perform reasoning distillation from stronger models [34]. Both approaches require either a reliable
100 verifier or a stronger RLM on the target task.

101 **Main Challenges** The key challenge for training RLM is the requirement to preserve the reasoning
102 trace, which is crucial for the model performance [32]. Thus, it is typically assumed that the training
103 data must provide reasoning traces, which are included in the model training objective [4, 18].
104 However, if no verifier is available for the given task, or no reasoning model with a higher solve
105 rate on the dataset in question is available, we can not easily obtain relevant reasoning traces for the
106 given task.

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

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

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

122 **Finetuning** We first fine-tune the base model M on $\mathcal{D}_{\text{train}}$ using standard SFT with LoRA [10].
123 We train the query, key, value, and output projections in self-attention, as well as the gate, up, and
124 down projections in the feed-forward network. This reduces the memory footprint and training cost
125 compared with full fine-tuning [29].

126 A reasoning model is normally trained to produce both a reasoning trace and a final answer. In our
127 setting, however, $\mathcal{D}_{\text{train}}$ contains only (i, o) pairs and does not contain the reasoning trace r . We

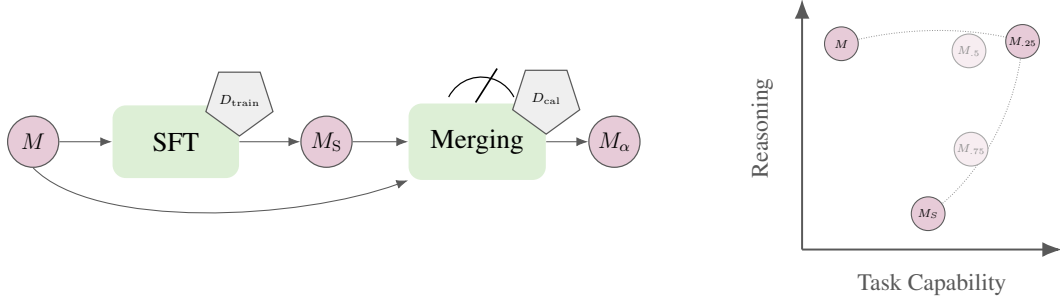


Figure 1: Our core method for RLM training consists of a light-weight two-step pipeline (Left): We first perform standard SFT on D_{train} , which produces M_{SFT} . We then merge M and M_{SFT} , using a calibration dataset D_{cal} to find a merge coefficient α that preserves model reasoning. (Right)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

128 therefore serialize each training target as (ε, o) : the empty reasoning trace ε followed by the target
 129 output. Let $y(o) = \text{serial}(\varepsilon, o)$ denote this serialized response, and let $y_t(o)$ be its t -th token. The
 130 SFT objective is the standard next-token prediction loss on this target,

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

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

135 **Merging** To recover reasoning behavior while retaining the task adaptation learned during SFT,
 136 we perform a two-way interpolation between the original reasoning model M and its standard-SFT
 137 variant M_{SFT} . This choice is motivated by the use of merging to mitigate forgetting [1].

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

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

142 We design a search procedure to pick the merge coeffi-
 143 cient that is as close as possible to the SFT model, pre-
 144 serving the target task performance gains, while recover-
 145 ing the model reasoning. Let ρ_{min} be the minimum ac-
 146 ceptable reasoning rate; in our experiments, $\rho_{\text{min}} = 0.9$.
 147 We are then looking for the largest merge ratio whose cal-
 148 ibration reasoning rate remains acceptable, i.e.,

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

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

151 We describe the searching algorithm in Algorithm 1. In
 152 practice, we evaluate a small uniform grid of merge coef-
 153 ficients rather than performing an adaptive search. For a
 154 grid resolution K , we use $\mathcal{A}_K = \{0, 1/K, 2/K, \dots, 1\}$. The point $\alpha = 0$ is the original model and
 155 is guaranteed to satisfy $\rho(M, \mathcal{D}_{\text{cal}}) = 1$, while $\alpha = 1$ is the standard SFT model. After evaluating
 156 all grid points on \mathcal{D}_{cal} , we select the largest α whose reasoning rate is at least ρ_{min} .

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^*}

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

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

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

174 4 Experimental Evaluation

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

178 4.1 Experimental Setup

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

181 **Models** We compare three key models: OpenThinker 7B [9], Apriel Nemotron 15B Thinker [22]
182 and Olmo3 7B Think [18]. This covers a diverse set of reasoning models: OpenThinker is a fine-
183 tune based on Qwen2.5 7B [31], distilled on reasoning traces by DeepSeek R1 [4]. Apriel 15B is
184 a fine-tuned version of Apriel 15B Base, post trained through CPT, SFT and GRPO. Olmo3 was
185 trained from scratch with a fully open pipeline, including SFT, DPO and Reinforcement Learning
186 from Verifiable Rewards [18].

187 **Methods** We first evaluate the unadapted RLM as *Baseline*. We compare it to the fully fine-tuned
188 variant of the model using LoRA, called *SFT*, trained as described in the first step of our training
189 pipeline without any reasoning traces. Unless otherwise indicated, we finetune the hyperparameters
190 learning rate, batch size and epoch for each combination of model and dataset, and report the hy-
191 perparameters in Appendix A. Finally we evaluate *Ours*, the re-merged version of the model at the
192 optimal merge ratio determined by our method using the described binary search method with up to
193 8 search steps. We further attempt reproducing two related works: *SDFT* [24], which leverages the
194 model for self-guided reasoning synthesization and *VeriFree* [36], which performs direct learning
195 on the golden answers.

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

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

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

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

204 The tasks are highly different in complexity and difficulty. In addition, both tasks do in general not
205 permit reliable verifiers, making it difficult to train on using RLVR. For Rust coding, RLVR is only
206 possible when comprehensive test suites are provided, which is not always possible.

207 **Metrics** For Rust coding, we choose a Rust translation of the MBPP dataset [13] and measure
208 the models PASS@1 performance [3], in other words, we measure the percentage of solutions that
209 implement a function described in natural language correctly, as measured by a set of unit tests.
210 For text summarization, we measure fluency, consistency, relevance and coherence on the CNN
211 split of Fabbri et al. [5]. We use Gemini 3 Pro [8] to evaluate the criteria according to a detailed
212 prompt, detailed in Appendix B. To confirm the validity of these results, we compare the Gemini
213 3 Pro ratings on the dataset of summaries with human annotators provided in Fabbri et al. [5] and
214 establish a Cohens Kappa of over 60% for each metric. We provide the details in Appendix A. For
215 the mathematical problem solving, we measure the performance of the models on the MATH500
216 dataset [15]. Note that we do not train on this task, and we therefore use it to assess the loss of
217 general reasoning capabilities of the model due to fine-tuning.

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

222 **Training Datasets** We devise two datasets for training. To ensure that we measure whether the
223 model preserved its generalization capability, we design the training and test datasets to stem from
224 different distributions, while having similar task formats. (1) *Rust coding*: We train the model on
225 a set of single-function coding tasks in the Rust language, where each data point is a pair of task
226 and code solution that were synthesized from an LLM [28]. We filter the dataset to the subset of
227 code samples that pass the corresponding test suite consistently, arriving at a subset of 6761 training
228 samples. We confirm the obtained datasets to be disjoint from the MBPP dataset by checking for
229 closest matches using cosine similarity and confirming their distinctness. (2) *Text summarization*:
230 We further train the model on the task of text summarization. As training dataset, we use the reddit
231 TLDR; split of Fabbri et al. [5], which is intentionally different to the CNN task split that we use for
232 evaluation.

233 4.2 Main Results

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

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

243 **Text Summarization** For text summarization, we observe that fluency, relevance and coherence
244 are already close to the maximal score for all evaluated models, on average 4.84 out of 5, and
245 remain stable across our training. We therefore focus on the metric of consistency and report the
246 performance on other metrics in Appendix A. Here, the best performance is usually achieved by the
247 SFT model, with model performance on average increasing by 1.10 points. This increase occurs
248 despite as model reasoning decreasing rapidly to 0.0% and 12.9% for OpenThinker 7B and April
249 15B, respectively. With our method, reasoning is fully recovered, resulting in a parallel recovery
250 of MATH500 performance, while preserving 76.7% and 97.4% of the performance increase for the
251 fully finetuned text consistency.

252 **Runtime and Cost** A key benefit of our method is its low cost. Our adaptation method requires on
253 average less than two hours to complete and fits on a single NVIDIA H200 GPU, costing in total less
254 than USD 6 on typical rental services at the time of writing. Most time is split between obtaining the

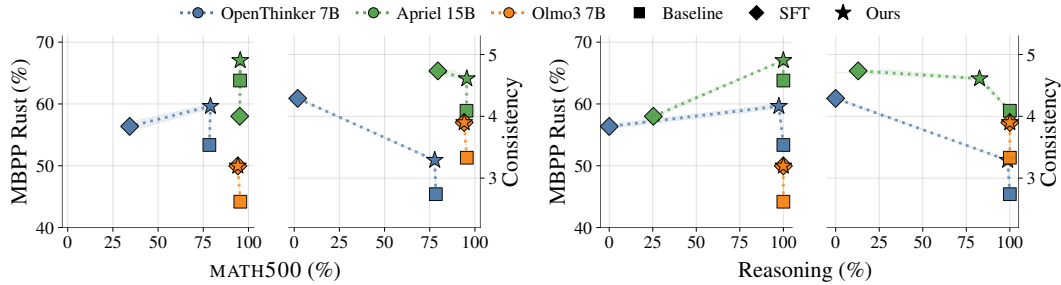


Figure 2: (Left) Our method preserves MATH500 performance while preserving or even increasing the task gains from SFT on Rust coding and text summarization. (Right) The preserved capability correlates with whether the model continues to reason on an out-of-distribution task.

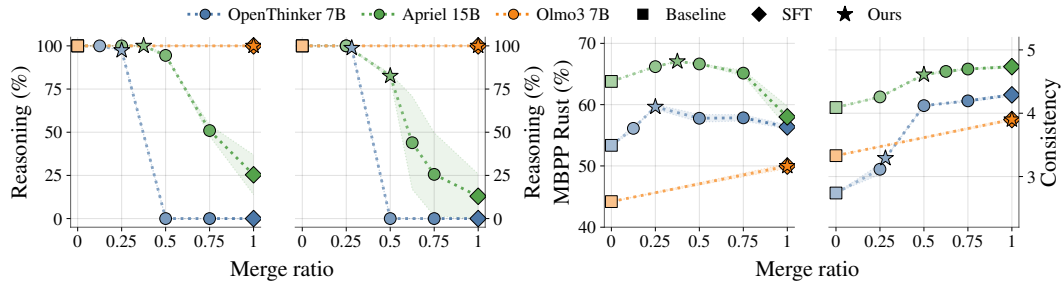


Figure 3: We evaluate all various merge ratios on the final test set. We observe that reasoning on MATH500 tasks drops sharply after a merging threshold is crossed (Left) and target task performance experiences a smoother curve (Right). On Rust, a clear benefit of restored reasoning to task performance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

255 fine-tuned model, with on average 36 minutes, and evaluating the merge ratios in Algorithm 1, with
 256 on average 23 min. Note that a key factor in our cost-efficiency is that our method does not require
 257 any reasoning rollouts from the model.

258 **Other Methods** We attempted reproducing two popular related methods that are also applicable
 259 to our setting of training RLMS using only standard SFT data. The first method, VeriFree [36], re-
 260 quires a significant amount of resources: Based on our initial experiments, we estimated 8 NVIDIA
 261 GPUs for a total of 60 hours for a single training run. Given this cost, which applies to each step of
 262 a preliminary hyperparameter search for fair comparison, we consider a comparison infeasible. Sec-
 263 ond, for SDFT [24] we have extensively tried to reproduce SDFT on Olmo3 7B, on the task of Rust
 264 coding, exploring over 20 hyperparameter combinations. Despite reaching out to the authors, the
 265 best performance we could achieve with SDFT was 41.4%, a slight decrease from the baseline per-
 266 formance of 44.2%, while performance on MATH500 also slightly decreased from 95.5% to 93.1%.
 267 As such, despite our extensive reproduction efforts, we did not include the above related methods
 268 for direct comparisons in our experiments. We report the used hyperparameters in Appendix A.

269 4.3 Ablation

270 We ablate over the choice of merge ratios, employed finetuning techniques and training hyperpa-
 271 rameters for the supervised finetuning.

272 **Merge ratio** We ablate over the merging factor α on the Rust and text summarization dataset, and
 273 show their reasoning rate on MATH500 and the target tasks in Figure 3 (left and right, respectively).
 274 We notice that OpenThinker 7B loses its reasoning capability quickly for $\alpha \in [0.25, 0.5]$, while
 275 April 15B loses it more slowly and Olmo3 7B even maintains full reasoning at $\alpha = 1$. Surprisingly,
 276 the performance on Rust coding appears to peak in for $\alpha \in [0.25, 0.5]$, where the model has still
 277 knowledge about the task from training and recovered reasoning. Meanwhile, consistency in the
 278 text summarization task increases more linearly with the merging ratio. Our method reliably picks

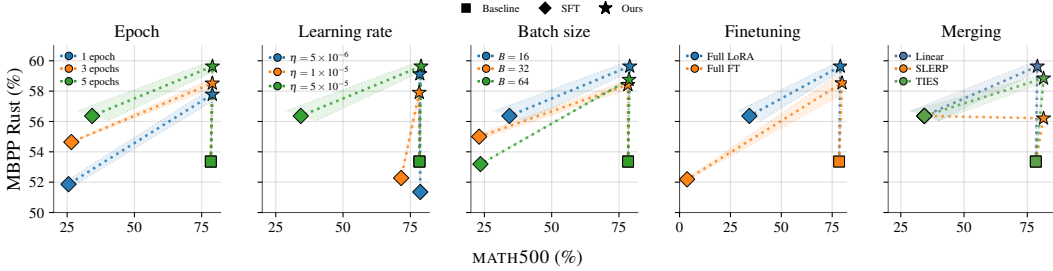


Figure 4: OpenThinker 7B ablations on Rust coding, from left to right: training epochs, learning rate η , batch size B , LoRA vs. full finetuning, and merging techniques. The overall trend of reasoning loss in standard SFT and reasoning and performance recovery is stable across all settings.

279 a point close to a good trade-off point between reasoning and task performance optimum, despite
 280 picking the point only based on the reasoning score on the calibration dataset.

281 **Merging technique** We ablate over the employed merging technique for OpenThinker 7B on Rust
 282 coding, presenting the results in the right-most panel of Figure 4. We compare the linear merging
 283 employed in our main experiments with Spherical Linear intERPolation (SLERP) [7, 25] and TIES
 284 [30]. For TIES, since fixing $\alpha = 1$ and searching over the density parameter would result in com-
 285 plete loss of reasoning even for a density of less than 0.1, we fix density to 0.5 and search over the
 286 α parameter. For all merging techniques, our method obtains the best merging ratio at $\alpha = 0.25$.
 287 Compared to linear interpolation, SLERP recovers more MATH500 performance but loses target-
 288 task performance, while TIES improves both MATH500 and target-task performance, albeit less
 289 than our linear merging.

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

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

302 5 Related Work

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

304 **Training Reasoning Models with Verifiers** The standard technique for training and adapting
 305 reasoning models is via reinforcement learning with verifiable rewards [4, 18, 26]. This works well
 306 for tasks that permit such verifiers, however leaves out many domains, such as the explored domain
 307 of text summarization. Moreover, while this technique is well-explored for bootstrapping reasoning
 308 model performance, continuous adaptation of RLMs to new tasks via RL is not well explored yet.

309 **Training Reasoning Models without Verifiers** Verifier-free or verifier-light methods try to enable
 310 training reasoning models without reasoning traces, by optimizing likelihood of known answers
 311 directly [35, 36] or training teacher models to provide feedback and justifications for reference
 312 answers [24]. These methods are closest in motivation to our setting, but they are computationally
 313 much more expensive and not designed to work on already trained reasoning models. Our method
 314 instead operates on minimal hardware requirements and allows inexpensive, fast task adaptation.

315 **Model merging** Model merging combines multiple checkpoints into a single model without ad-
316 ditional gradient updates. A common motivation is to combine task-specific skills or mitigate for-
317 getting by merging related models [1, 33]. Prior methods include simple weight averaging, task
318 arithmetic [11], and sparse or sign-based variants such as TIES [30]. Linear interpolation is also
319 connected to linear mode connectivity, where independently trained or fine-tuned networks can lie
320 on low-loss paths in weight space [6]. Most merging work aims to combine capabilities from mul-
321 tiple specialized models. Our use is narrower: we merge an SFT checkpoint back with its own
322 original reasoning checkpoint to trade off target-task adaptation against preservation of reasoning
323 behavior. Recent work on tunable reasoning through model merging suggests that interpolation can
324 control the degree of reasoning behavior in language models [14]. We build on this observation, but
325 select the merge coefficient using a small calibration set and target the practical setting where only
326 input-output SFT data is available.

327 6 Limitations and Future Work

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

333 **Interpretability** In our experiments we observe that models stop reasoning by immediately pro-
334 ducing an end-of-reasoning marker. To address this, during the development of the method, we tried
335 prefilling the model response with a start-of-reasoning marker and an immediately following non-
336 end-of-reasoning token like ‘Okay’, however, the reasoning performance was not recovered. An
337 interesting direction for future work is to mechanistically understand how reasoning tendency is en-
338 coded in the model. We hope that potential insights could then be directly leveraged to design better
339 adaptation methods for reasoning models, in particular in terms of recovering reasoning tendency.

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

350 7 Conclusion

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

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

363 References

- 364 [1] Anton Alexandrov, Veselin Raychev, Mark Niklas Müller, Ce Zhang, Martin T. Vechev, and
365 Kristina Toutanova. Mitigating catastrophic forgetting in language transfer via model merging.
366 In *EMNLP (Findings)*, 2024.
- 367 [2] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
368 Concrete problems in ai safety, 2016. URL <https://arxiv.org/abs/1606.06565>.
- 369 [3] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto,
370 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating
371 Large Language Models Trained on Code. *arXiv Preprint*, 2021. URL <https://arxiv.org/abs/2107.03374>.
- 373 [4] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement
374 learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 375 [5] Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher,
376 and Dragomir Radev. Summeval: Re-evaluating summarization evaluation, 2021. URL <https://arxiv.org/abs/2007.12626>.
- 378 [6] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear
379 mode connectivity and the lottery ticket hypothesis, 2020. URL <https://arxiv.org/abs/1912.05671>.
- 381 [7] Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin,
382 Brian Benedict, Mark McQuade, and Jacob Solawetz. Arcee’s mergekit: A toolkit for merging
383 large language models, 2025. URL <https://arxiv.org/abs/2403.13257>.
- 384 [8] Google DeepMind. Gemini Pro, 2025. URL <https://deepmind.google/technologies/gemini/pro/>.
- 386 [9] Etash Guha, Ryan Marten, Sedrick Keh, Negin Raof, Georgios Smyrnis, Hritik Bansal, Mar-
387 ianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer,
388 Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff,
389 Shiye Su, Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-
390 Jie Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave,
391 Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit
392 Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan,
393 Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran
394 Sathiamoorthy, Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for
395 reasoning models, 2025. URL <https://arxiv.org/abs/2506.04178>.
- 396 [10] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
397 Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *In-*
398 *ternational Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- 400 [11] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig
401 Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic, 2023.
402 URL <https://arxiv.org/abs/2212.04089>.
- 403 [12] KRAFTON and SKT. Continual post-training of llms via offline grpo for mathematical reason-
404 ing, 2025. URL https://krafton-ai.github.io/blog/llm_post_training_en/. Accessed:
405 2025-07-28.
- 406 [13] Northeastern University Programming Research Lab. Multipl-e (revision 40ca145), 2025.
407 URL <https://huggingface.co/datasets/nuprl/MultiPL-E>.
- 408 [14] Xiaochong Lan, Yu Zheng, Shiteng Cao, and Yong Li. The thinking spectrum: An empirical
409 study of tunable reasoning in llms through model merging, 2025. URL <https://arxiv.org/abs/2509.22034>.

- 411 [15] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
412 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023. URL
413 <https://arxiv.org/abs/2305.20050>.
- 414 [16] Monte MacDiarmid, Benjamin Wright, Jonathan Uesato, Joe Benton, Jon Kutasov, Sara
415 Price, Naia Bouscal, Sam Bowman, Trenton Bricken, Alex Cloud, Carson Denison, Johannes
416 Gasteiger, Ryan Greenblatt, Jan Leike, Jack Lindsey, Vlad Mikulik, Ethan Perez, Alex Ro-
417 drigues, Drake Thomas, Albert Webson, Daniel Ziegler, and Evan Hubinger. Natural emergent
418 misalignment from reward hacking in production rl, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2511.18397)
419 [2511.18397](https://arxiv.org/abs/2511.18397).
- 420 [17] NovaSky Team. Sky-t1: Train your own o1 preview model within \$450. [https://novasky-](https://novasky-ai.github.io/posts/sky-t1)
421 [ai.github.io/posts/sky-t1](https://novasky-ai.github.io/posts/sky-t1), 2025. Accessed: 2025-01-09.
- 422 [18] Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heine-
423 man, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, Jacob Morrison,
424 Jake Poznanski, Kyle Lo, Luca Soldaini, Matt Jordan, Mayee Chen, Michael Noukhovitch,
425 Nathan Lambert, Pete Walsh, Pradeep Dasigi, Robert Berry, Saumya Malik, Saurabh Shah,
426 Scott Geng, Shane Arora, Shashank Gupta, Taira Anderson, Teng Xiao, Tyler Murray, Tyler
427 Romero, Victoria Graf, Akari Asai, Akshita Bhagia, Alexander Wettig, Alisa Liu, Aman Ran-
428 gapur, Chloe Anastasiades, Costa Huang, Dustin Schwenk, Harsh Trivedi, Ian Magnusson,
429 Jaron Lochner, Jiacheng Liu, Lester James V. Miranda, Maarten Sap, Malia Morgan, Michael
430 Schmitz, Michal Guerquin, Michael Wilson, Regan Huff, Ronan Le Bras, Rui Xin, Rulin Shao,
431 Sam Skjonsberg, Shannon Zejiang Shen, Shuyue Stella Li, Tucker Wilde, Valentina Pyatkin,
432 Will Merrill, Yapei Chang, Yuling Gu, Zhiyuan Zeng, Ashish Sabharwal, Luke Zettlemoyer,
433 Pang Wei Koh, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. Olmo 3, 2025. URL
434 <https://arxiv.org/abs/2512.13961>.
- 435 [19] OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky,
436 Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex
437 Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Alli-
438 son Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein,
439 Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Bar-
440 ret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew,
441 Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon
442 Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy,
443 Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse,
444 Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam,
445 David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dra-
446 gos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung,
447 Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo
448 Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit,
449 Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guil-
450 laume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad,
451 Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Os-
452 band, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman,
453 Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui
454 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish,
455 Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan
456 Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl
457 Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu,
458 Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho,
459 Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas
460 Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar,
461 Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan
462 Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mi-
463 anna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay,
464 Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Roha-
465 ninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam
466 Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne,

467 Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes,
468 Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan
469 James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer,
470 Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney,
471 Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan
472 Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor,
473 Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault
474 Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit
475 Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie
476 Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech
477 Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen
478 Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. OpenAI o1 system card, 2024. URL
479 <https://arxiv.org/abs/2412.16720>.

480 [20] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
481 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red
482 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
483 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christo-
484 pher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg
485 Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew
486 Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis
487 Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester
488 Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory
489 Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve
490 Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fe-
491 dus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges,
492 Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan
493 Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei
494 Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke,
495 Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu,
496 Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang,
497 Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan,
498 Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan
499 Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros,
500 Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis,
501 Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike,
502 Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz
503 Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Man-
504 ning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob Mc-
505 Grew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Med-
506 ina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
507 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
508 Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
509 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pan-
510 tuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov,
511 Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde
512 de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea
513 Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,
514 Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez,
515 Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt,
516 David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh,
517 Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Kata-
518 rina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski
519 Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil
520 Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan
521 Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright,
522 Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila
523 Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens
524 Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu,

- 525 Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers,
526 Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk,
527 and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 528 [21] Jatin Prakash and Anirudh Buvanesh. What can you do when you have zero rewards during rl?
529 *CoRR*, 2025.
- 530 [22] Shruthan Radhakrishna, Soham Parikh, Gopal Sarda, Anil Turkan, Quaizar Vohra, Ray-
531 mond Li, Dhruv Jhamb, Kelechi Ogueji, Aanjaneya Shukla, Oluwanifemi Bamgbose, Toby
532 Liang, Luke Kumar, Oleksiy Ostapenko, Shiva Krishna Reddy Malay, Aman Tiwari, Tara
533 Bogavelli, Vikas Yadav, Jash Mehta, Saloni Mittal, Akshay Kalkunte, Pulkit Pattnaik, Khalil
534 Slimi, Anirudh Sreeram, Jishnu Nair, Akintunde Oladipo, Shashank Maiya, Khyati Maha-
535 jan, Rishabh Maheshwary, Masoud Hashemi, Sai Rajeswar Mudumba, Sathwik Tejaswi Mad-
536 husudhan, Torsten Scholak, Sebastien Paquet, Sagar Dadasam, and Srinivas Sunkara. Apriel-
537 nemotron-15b-thinker, 2025.
- 538 [23] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
539 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of
540 mathematical reasoning in open language models, 2024. URL [https://arxiv.org/abs/2402.](https://arxiv.org/abs/2402.03300)
541 [03300](https://arxiv.org/abs/2402.03300).
- 542 [24] Idan Shenfeld, Mehul Damani, Jonas Hübner, and Pulkit Agrawal. Self-distillation enables
543 continual learning, 2026. URL <https://arxiv.org/abs/2601.19897>.
- 544 [25] Ken Shoemake. Animating rotation with quaternion curves. In *Proceedings of the 12th*
545 *Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '85,
546 page 245–254, New York, NY, USA, 1985. Association for Computing Machinery. ISBN
547 0897911660. doi: 10.1145/325334.325242. URL [https://doi.org/10.1145/325334.](https://doi.org/10.1145/325334.325242)
548 [325242](https://doi.org/10.1145/325334.325242).
- 549 [26] Zafir Stojanovski, Oliver Stanley, Joe Sharratt, Richard Jones, Abdulhakeem Adefioye, Jean
550 Kaddour, and Andreas Köpf. Reasoning gym: Reasoning environments for reinforcement
551 learning with verifiable rewards, 2025. URL <https://arxiv.org/abs/2505.24760>.
- 552 [27] Oxen.ai Team. Rust verified code dataset. <https://www.oxen.ai/ox/Rust>, 2024. Accessed:
553 2025-01-01.
- 554 [28] Oxen.ai Team. Training a rust 1.5b coder lm with re-
555 inforcement learning (grpo). [https://ghost.oxen.ai/](https://ghost.oxen.ai/training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grpo/)
556 [training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grpo/](https://ghost.oxen.ai/training-a-rust-1-5b-coder-lm-with-reinforcement-learning-grpo/), 2024. Techni-
557 cal report.
- 558 [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
559 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL [https://arxiv.](https://arxiv.org/abs/1706.03762)
560 [org/abs/1706.03762](https://arxiv.org/abs/1706.03762).
- 561 [30] Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging:
562 Resolving interference when merging models. In *Thirty-seventh Conference on Neural Infor-*
563 *mation Processing Systems*, 2023. URL <https://openreview.net/forum?id=xtaX3WyCj1>.
- 564 [31] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
565 Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 Technical Report. *arXiv Preprint*,
566 2024. URL <https://doi.org/10.48550/arXiv.2412.15115>.
- 567 [32] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
568 Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang,
569 Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei
570 Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin
571 Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng
572 Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang,
573 Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang
574 Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru
575 Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL [https://arxiv.](https://arxiv.org/abs/2505.09388)
576 [org/abs/2505.09388](https://arxiv.org/abs/2505.09388).

- 577 [33] Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng
578 Tao. Model merging in llms, mllms, and beyond: Methods, theories, applications and oppor-
579 tunities, 2024. URL <https://arxiv.org/abs/2408.07666>.
- 580 [34] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. STar: Bootstrapping reason-
581 ing with reasoning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun
582 Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_3ELRdg2sgI.
583
- 584 [35] Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang
585 Wang, Min Lin, and Chao Du. Reinforcing general reasoning without verifiers. *CoRR*, 2025.
- 586 [36] Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang
587 Wang, Min Lin, and Chao Du. Reinforcing general reasoning without verifiers, 2025. URL
588 <https://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}

589 A Experimental Details, Ablations and Case Study

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

592 A.1 Refinement of the Rust dataset

593 In our study, the Rust corpus from Team [27, 28] is used for supervised fine-tuning and in-
594 distribution evaluation. Each sample in the corpus consists of four fields: a task identifier, a natural-
595 language Rust prompt that describes the target function, a Rust code implementation that solves
596 the task, and a list of executable test cases. Consequently, the dataset provides complete function-
597 level programming tasks paired with verification harnesses, enabling evaluation of generated code
598 correctness.

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

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

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

613 As mentioned in §4, we validated the use of Gemini 3 Pro Preview as LLM-as-a-Judge metric for
614 our evaluation of the Text Summarization task. For this, we run the model on the SummEval dataset
615 [5] and compare the model scores with the human annotations. We present the results in Table 1,
616 showing that the model achieves generally high correlation of around 60% on all scores, and almost
617 70% in the main inspected metric of consistency.

618 A.3 SFT Hyperparameters

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

623 For the SDFT [24] training result, we perform 26 runs with the following SFT settings: epochs
624 $\in \{1, 2, 3, 5\}$, learning rate $\eta \in \{1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}, 5 \times 10^{-6}\}$, batch size $B \in$
625 $\{16, 32, 64\}$, and EMA coefficient $\alpha \in \{0.01, 0.02, 0.05\}$. We set `max_token_length` = 4096 and
626 discarded training samples exceeding this budget, leaving roughly one-third of the examples. The
627 best-performing run used epoch 2, learning rate 5×10^{-6} , batch size 64, and EMA $\alpha = 0.05$, and
628 completed in about five hours using 4 NVIDIA H200 GPUs. The full four-panel text-summarization
629 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

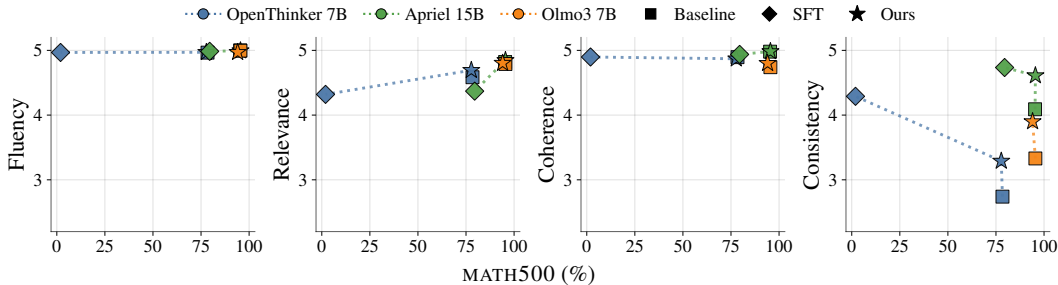


Figure 5: Text summarization tradeoff from SFT/merging experiments: Fluency, Relevance, Coherence, and Consistency against MATH-500.

630 A.4 Experimental Details on other Text Summarization metrics

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

636 B Prompts

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

```

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

```

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

```

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}.

```

1 user
2   You are a helpful assistant in evaluating the quality of a summary.
3   You will be given a news article and a summary written for that article.
4   Your task is to evaluate the *relevance* of the summary.
5
6   Definition of Relevance:
7   Relevance measures how well the summary captures the important information from
8   the article. A relevant summary includes the most important main idea of the
9   article and avoids redundancies and minor, trivial, or unrelated details.
10
11  Evaluation Criteria (Relevance: 1-100)
12
13  80-100 -- The summary captures the key idea of the article accurately and
14           completely. It focuses on the essential information and avoids any
15           redundancies and unnecessary or irrelevant content.
16
17  40-79 -- The summary includes some important information but only partially
18           capture the main idea. It misses key points or includes some redundancies
19           or minor details.
20
21  1-39 -- The summary fails to capture the main idea of the article and completely
22           deviate from the main idea. It may focus on unimportant details or
23           irrelevant content or omit critical points.
24
25  Evaluation Steps:
26  1. Read and understand the article.
27  2. Identify the article's main idea and secondary details.
28  3. Read the summary and judge if it captures the main idea of the article.
29  4. Identify if there are any secondary details or redundancy in the summary.
30  5. Assign a score from 1 to 100 based on the Evaluation Criteria above.
31
32  Output your detailed thought process and formal justification based on the
33  Evaluation Criteria, and finally output your final score in the format shown
34  below:
35
36  <think>
37  // your thought process and justification
38  </think>
39
40  <Final score>
41  Relevance: // your final score
42
43  Article:
44  {{ARTICLE}}
45
46  Summary:
47  {{SUMMARY}}

```

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

```
1 user
2   You are a helpful assistant in evaluating the quality of a summary.
3   You will be given a news article and a summary written for that article.
4   Your task is to evaluate the *fluency* of the summary.
5
6   Definition of Fluency:
7   Fluency measures how easy the summary is to read.
8   All sentences in a fluent summary need to be readable and natural.
9   Minor grammatical, formatting, capitalization, or tokenization issues should
10  NOT be heavily penalized as long as they do not make the text difficult to read
11  or understand.
12
13  Evaluation Criteria (Fluency: 1-5)
14
15  5 -- The summary is easy to read and all the sentences are understandable and
16  natural. Minor issues such as awkward wording, missing capitalization, or
17  tokenization artifacts are acceptable if they do not hinder readability.
18
19  3 -- The summary is readable but some of sentences include awkward phrasing,
20  inconsistent grammar, or formatting problems that reduce readability.
21
22  1 -- The summary is difficult to read. It contains frequent grammatical errors,
23  broken or incomplete sentences, or severe formatting problems that
24  significantly hinder understanding.
25
26  Evaluation Steps:
27  1. Read and understand the article.
28  2. Read the summary and judge whether the text is easy to read and whether the
29  individual sentences are natural.
30  3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
31
32  Only output the score. Do not include any additional explanations or text.
33  Use the following format:
34  Fluency: <score>
35
36  Article:
37  {{ARTICLE}}
38
39  Summary:
40  {{SUMMARY}}
```

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

```

1 user
2   You are a helpful assistant in evaluating the quality of a summary.
3   You will be given a news article and a summary written for that article.
4   Your task is to evaluate the *consistency* of the summary.
5
6   Definition of Consistency:
7   Consistency measures how factually aligned the summary is with the article.
8   A consistent summary should not introduce information that contradicts the
9   article or contain hallucinated statements not supported by the source article.
10
11  Evaluation Criteria (Consistency: 1-5)
12
13  5 -- All statements in the summary are fully supported by the article. No
14     contradictions, distortions, or hallucinated details appear.
15
16  3 -- The summary is mostly consistent but includes minor inaccuracies, unclear
17     references, or slightly misleading phrasing. These issues do not
18     significantly alter the meaning.
19
20  1 -- The summary contains factual errors or statements that contradict the
21     article or introduce unsupported information.
22
23  Evaluation Steps:
24  1. Read and understand the article.
25  2. Read the summary and check whether each fact is supported by the article.
26  3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
27
28  Only output the score. Do not include any additional explanations or text.
29  Use the following format:
30  Consistency: <score>
31
32  Article:
33  {{ARTICLE}}
34
35  Summary:
36  {{SUMMARY}}

```

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

```

1 user
2   You are a helpful assistant in evaluating the quality of a summary.
3   You will be given a news article and a summary written for that article.
4   Your task is to evaluate the *coherence* of the summary.
5
6   Definition of Coherence:
7   Coherence measures how well the ideas in the summary fit together.
8   A coherent summary presents information in a logical order, with smooth
9   transitions between sentences. It should read as a connected, well-structured
10  whole.
11
12  Evaluation Criteria (Coherence: 1-5)
13
14  5 -- The summary is clearly organized and easy to follow. Sentences flow
15      naturally, and ideas progress in a logical order.
16
17  3 -- The summary is somewhat coherent but has noticeable issues in flow or
18      structure. Some sentences feel disconnected or out of place, yet the overall
19      meaning is still clear.
20
21  1 -- The summary is hard to follow. The sentences are out of order and loosely
22      connected, which makes the summary feel fragmented.
23
24  Evaluation Steps:
25  1. Read and understand the article.
26  2. Read the summary and assess whether the ideas are presented in a clear and
27     logical order.
28  3. Assign a score from 1 to 5 based on the Evaluation Criteria above.
29
30  Only output the score. Do not include any additional explanations or text.
31  Use the following format:
32  Coherence: <score>
33
34  Article:
35  {{ARTICLE}}
36
37  Summary:
38  {{SUMMARY}}

```

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