Learning to Reason without External Rewards

Anonymous Authors¹

Abstract

Training large language models (LLMs) for complex reasoning via Reinforcement Learning with Verifiable Rewards (RLVR) is effective but limited by reliance on costly, domain-specific supervision. We explore Reinforcement Learning from 015 Internal Feedback (RLIF), a framework that enables LLMs to learn from intrinsic signals without external rewards or labeled data. We propose IN-018 TUITOR, an RLIF method that uses a model's own confidence-termed self-certainty-as its sole re-020 ward signal. INTUITOR replaces external rewards in Group Relative Policy Optimization (GRPO) with self-certainty scores, enabling fully unsupervised learning. Experiments demonstrate that IN-TUITOR matches GRPO's performance on mathe-025 matical benchmarks while achieving superior generalization to out-of-domain tasks like code gen-027 eration, without requiring gold solutions or test cases. Our findings show that intrinsic model sig-029 nals can drive effective learning across domains, 030 offering a scalable alternative to RLVR for autonomous AI systems where verifiable rewards are unavailable.

1. Introduction

034

035

049

050

051

053

Reinforcement learning has become essential for enhancing 038 large language model capabilities. Early work focused on 039 Reinforcement Learning from Human Feedback (RLHF), which aligns model outputs with human values through 041 reward models trained on preference data (Ouyang et al., 2022). Recent advances in Reinforcement Learning with 043 Verifiable Rewards (RLVR) replace learned reward models with automatically verifiable signals-such as exact answer 045 matching in mathematical problem-solving-demonstrating 046 improved reasoning capabilities in models like DeepSeek-047 R1 (Guo et al., 2025; Lambert et al., 2024; Hu et al., 2025).



Figure 1: Overview of RLIF and INTUITOR's Performance. Left: Illustration of RLIF, a paradigm where LLMs learn from intrinsic signals generated by the model itself, without external supervision. Right: Performance comparison of Qwen2.5-3B Base, GRPO, and INTUITOR (our RLIF instantiation). Both GRPO and INTUITOR are trained on the MATH dataset. INTUITOR achieves comparable performance to GRPO on in-domain mathematical benchmarks (GSM8K, MATH500) and demonstrates better generalization to out-of-domain code generation tasks (Live-CodeBench v6, CRUXEval). Part of the illustration was generated by GPT-40.

Despite these successes, both RLHF and RLVR face fundamental limitations that constrain their broader applicability. RLHF requires extensive human annotation, making it expensive and potentially biased (Gao et al., 2023). RLVR, while avoiding learned reward models, demands domainspecific verifiers and gold-standard solutions. In mathematics, this requires expert annotation of solutions; in code generation, it necessitates comprehensive test suites and execution environments (Liu et al., 2023; Liu & Zhang, 2025; Team et al., 2025; Xiaomi LLM-Core Team, 2025). These requirements limit RLVR to carefully curated domains and complicate deployment in open-ended scenarios. Moreover, outcome-oriented verifiable rewards limit transferability to other domains. These challenges motivate exploration of more general and scalable reward paradigms, leading to a critical research question:

Can LLMs enhance their reasoning abilities by relying solely on intrinsic, self-generated signals, without recourse to external verifiers or domain-specific ground truth?

In this paper, we introduce and explore such a paradigm: *Reinforcement Learning from Internal Feedback (RLIF)*, where models optimize intrinsic feedback to improve performance without external rewards or supervision. The mo-

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

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

tivation for RLIF extends to future scenarios where models
develop superhuman capabilities that become difficult for
humans to evaluate directly (Burns et al., 2023), requiring
self-improvement through intrinsic mechanisms (Oudeyer
& Kaplan, 2007).

060 Under the RLIF paradigm, we propose INTUITOR, a novel 061 reinforcement learning approach leveraging a model's own 062 confidence as an intrinsic reward. This builds on observa-063 tions that LLMs exhibit lower confidence on difficult prob-064 lems (Farquhar et al., 2024; Kuhn et al., 2023; Kang et al., 065 2024; 2025); optimizing for confidence should improve 066 reasoning capabilities. Specifically, we use self-certainty 067 (Kang et al., 2025), the average KL divergence between 068 the model's output distribution and a uniform distribution, 069 as our confidence measure. This metric has proven useful 070 for distinguishing high-quality responses from flawed ones (Kang et al., 2025; Ma et al., 2025). Building on this insight, INTUITOR guides learning through self-generated signals, eliminating the need for external supervision or handcrafted 074 rewards. The implementation of INTUITOR is simple, effi-075 cient, and effective: we replace the verifiable reward signal 076 in existing RLVR frameworks, specifically Group Relative 077 Policy Optimization (GRPO) (Shao et al., 2024), with self-078 certainty scores, using the same policy gradient algorithm. 079

Our experiments demonstrate promising results. On the 081 MATH dataset (Hendrycks et al., 2021) with Qwen2.5-3B 082 base (Yang et al., 2024a), INTUITOR matches the perfor-083 mance of GRPO without relying on any gold answers. As INTUITOR rewards the generation trajectory rather than only the end result, it generalizes more effectively: training 086 a Qwen2.5-3B base model on MATH yields a 65% relative improvement on LiveCodeBench Code generation task (Jain 087 088 et al., 2024) versus no improvement for GRPO, and a 76%gain on CRUXEval-O (Gu et al., 2024) compared with 44%089 090 for GRPO. Additionally, when we fine-tune the Qwen2.5-091 1.5B base model with INTUITOR on the MATH corpus, a 092 model that originally produces repetitive content and scores 093 0% on LiveCodeBench learns to emit coherent reasoning 094 chains and well-structured code, reaching 9.9% accuracy 095 after the tuning. This demonstrates the strong generalization 096 capabilities of INTUITOR. As INTUITOR requires only a 097 clear prompt and no verifiable reward, it is broadly appli-098 cable across tasks, providing fresh evidence that pretrained 099 LLMs possess richer latent behavioral priors than previously 100 recognized.

Our contributions can be summarized as follows:

 We introduce and explore Reinforcement Learning from Internal Feedback (RLIF), a novel reinforcement learning paradigm enabling LLMs to improve reasoning skills by leveraging intrinsic, self-generated signals, without reliance on external supervision or labeled data.

• We introduce INTUITOR, an RLIF-based method that uti-

lizes a model's own internal confidence measure—termed *self-certainty*—as the sole intrinsic reward.

• We demonstrate that INTUITOR matches supervised RL performance on in-domain tasks while achieving superior out-of-domain generalization. We uncover emergent structured reasoning and enhanced instruction-following capabilities induced by intrinsic rewards.

2. Related Work

Reinforcement Learning from Human Feedback (**RLHF**). RL has become instrumental in refining LLMs. Early pivotal work centered on Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022), which aligns LLMs with human values by training a reward model on human preference data. While effective, RLHF is often resource-intensive due to the need for extensive human annotation (Touvron et al., 2023). Subsequent innovations like Direct Preference Optimization (DPO) (Rafailov et al., 2023) aimed to simplify this by directly training models on preferences. The reliance on human-generated or model-approximated human preferences poses scalability challenges and introduces potential biases from the reward model itself (Gao et al., 2023).

Reinforcement Learning with Verifiable Rewards (**RLVR**). RLVR emerged as a powerful alternative, particularly for tasks with clear correctness criteria like mathematical reasoning and code generation (Guo et al., 2025; Lambert et al., 2024; Hu et al., 2025; Team et al., 2025; Xiaomi LLM-Core Team, 2025). RLVR utilizes rule-based verification functions, such as exact answer matching (Guo et al., 2025; Team et al., 2025; Xiaomi LLM-Core Team, 2025; Jaech et al., 2024), to provide reward signals, thereby avoiding the complexities and potential pitfalls of learned reward models. This approach has sparked significant advances, with models like DeepSeek-R1 (Guo et al., 2025) achieving state-of-the-art reasoning capabilities. The development of robust policy optimization algorithms like GRPO (Shao et al., 2024) and its variants (Luo et al., 2025; Liu et al., 2025) has further solidified RLVR's success. Nevertheless, RLVR's applicability is largely confined to domains where verifiable gold solutions or exhaustive test cases can be constructed, and its predominant focus on outcome-based rewards can limit generalization to dissimilar tasks or those requiring nuanced, process-oriented feedback.

Intrinsic Signals and Self-Play in Language Model Optimization. Self-play and intrinsic rewards have gained attention as strategies for enabling autonomous model improvement. Inspired by early work in games like AlphaGo Zero (Silver et al., 2017), recent LLM-based frameworks incorporate self-refinement mechanisms to bootstrap reasoning ability. Methods like SPIN (Chen et al., 2024) and

Self-Rewarding LMs (Yuan et al., 2024) utilize the model 111 itself to provide feedback for subsequent training iterations. 112 While earlier work such as STaR (Zelikman et al., 2022) re-113 lies on repeated outcome evaluation, more recent approaches 114 explore self-improvement through procedural generalization 115 and goal invention (Poesia et al., 2024; Cheng et al., 2024). 116 Concurrent works like Genius, TTRL, and Absolute Zero 117 (Xu et al., 2025; Zuo et al., 2025; Zhao et al., 2025)-lever-118 age queries without labels for reinforcement learning but 119 remain constrained to specific task distributions, primar-120 ily in mathematical reasoning. INTUITOR aligns with this 121 direction but introduces a lightweight, general-purpose alter-122 native: using self-certainty as a confidence-based intrinsic reward. Unlike prior work, INTUITOR enables single-agent 124 reinforcement learning across diverse tasks without rely-125 ing on explicit feedback, gold labels, or environment-based 126 validation.

3. Method

127 128

129

141

142

143

144

150 151

130 In this section, we present our approach to training large lan-131 guage models through intrinsic feedback mechanisms. We 132 begin by reviewing existing RL-based fine-tuning paradigms 133 and their limitations, which motivate our exploration of Re-134 inforcement Learning from Internal Feedback (RLIF). We 135 then introduce INTUITOR, our novel RLIF method that lever-136 ages self-certainty as an intrinsic reward signal, and detail 137 its implementation through policy optimization. 138

1393.1. From External Supervision to Internal Feedback

Current RL fine-tuning approaches for LLMs primarily fall into two categories: those relying on external human feedback (RLHF) and those using verifiable, often task-specific, rewards (RLVR).

145 In RLHF (Ziegler et al., 2019; Ouyang et al., 2022), the 146 policy π_{θ} is optimized to align with human preferences, 147 typically encapsulated by a learned reward model r_{ϕ} . The 148 objective is:

$$\max_{\pi_{\theta}} \mathbb{E}_{o \sim \pi_{\theta}(q)} \left[r_{\phi}(q, o) - \beta \mathrm{KL}[\pi_{\theta}(o|q) \| \pi_{\mathrm{ref}}(o|q)] \right]$$
(1)

where q is an input query, o is the generated output, π_{ref} is 152 an initial reference policy, and β is a coefficient controlling 153 the KL divergence to prevent excessive deviation from π_{ref} . 154 155 Online RL algorithms like PPO (Schulman et al., 2017) generate samples from π_{θ} , evaluate them using r_{ϕ} , and 156 update π_{θ} to maximize this objective. However, the reward 157 model r_{ϕ} is crucial yet fragile; introducing it can lead to 158 "reward hacking," and retraining it is resource-intensive, 159 complicating the training pipeline (Gao et al., 2023). 160

 RLVR, on the other hand, substitutes the learned reward model with an automatically verifiable signal. This has proven effective in promoting reasoning capabilities, especially in domains like mathematics (Guo et al., 2025). The RLVR objective is:

$$\max_{\pi_{\theta}} \mathbb{E}_{o \sim \pi_{\theta}(q)} \left[v(q, o) - \beta \mathrm{KL}[\pi_{\theta}(o|q) \| \pi_{\mathrm{ref}}(o|q)] \right]$$
(2)

where v(q, o) is a verifiable reward function. For instance, in mathematical problem-solving, v(q, o) might be: $v(q, o) = \begin{cases} \alpha & \text{if output } o \text{ is correct} \\ 0 & \text{otherwise.} \end{cases}$. RLVR is often implemented using algorithms like REINFORCE (Williams, 1992), PPO or GRPO. Despite their simplicity, verifiable rewards still rely on gold-standard answers or test executions, which are costly and domain-specific (Liu et al., 2025; Team et al., 2025). RLVR faces challenges in extending beyond math and code to tasks involving ambiguity or subjective reasoning.

3.2. Reinforcement Learning from Internal Feedback (RLIF)

To overcome the limitations of RLHF's costly human annotation and RLVR's domain-specific supervision, we propose Reinforcement Learning from Internal Feedback (RLIF). Instead of depending on external evaluation, RLIF uses the model's own assessment of its outputs or reasoning process as feedback. This offers several advantages: it reduces reliance on supervision infrastructure, provides task-agnostic reward signals, and supports learning in domains where external verification is unavailable. Under the RLIF paradigm, the optimization objective becomes:

$$\max_{\pi_{\theta}} \mathbb{E}_{o \sim \pi_{\theta}(q)} \left[u(q, o) - \beta \mathrm{KL}[\pi_{\theta}(o|q) \| \pi_{\mathrm{ref}}(o|q)] \right] \quad (3)$$

where u(q, o) represents an intrinsic signal derived from the model's internal state or computation, rather than external verification. The key challenge lies in identifying intrinsic signals that correlate with output quality and can effectively guide learning.

Concurrent research explores related concepts within the RLIF paradigm. For example, Entropy Minimized Policy Optimization (EMPO) (Zhang et al., 2025) minimizes LLM predictive entropy on unlabeled questions in a latent semantic space. SEED-GRPO (Chen et al., 2025) uses the semantic entropy of generated sequences, combined with ground truth rewards, to modulate policy updates. Reinforcement Learning with a Negative Entropy Reward (EM-RL) (Agarwal et al., 2025) employs a reward signal based solely on the negative sum of token-level entropy, akin to REINFORCE but without labels. These methods highlight the growing interest and potential of leveraging intrinsic signals for LLM training under the RLIF framework.

165 **3.3. INTUITOR: Policy Optimization with Self-Certainty**

We propose INTUITOR, a novel RLIF method that utilizes a model's own confidence as the sole intrinsic reward signal u(q, o).

170 Our choice of model confidence as the intrinsic reward is 171 motivated by observations that LLMs often exhibit lower 172 confidence when encountering unfamiliar tasks or lacking 173 sufficient knowledge (Kang et al., 2024). Conversely, higher 174 confidence frequently correlates with correctness. By re-175 warding increased self-confidence, INTUITOR encourages 176 to iteratively "practice" and refine its reasoning pathways 177 until it becomes more confident in its outputs.

We adopt the self-certainty metric from Kang et al. (2025), defined as the average KL divergence between a uniform distribution U over the vocabulary \mathcal{V} and the model's nexttoken distribution:

$$\begin{array}{ll}
\text{183} \\
\text{184} \\
\text{185} \\
\text{185} \\
\text{186} \\
\text{187} \\
\text{188} \\
\end{array} = -\frac{1}{|o| \cdot |\mathcal{V}|} \sum_{i=1}^{|o|} \sum_{j=1}^{|\mathcal{V}|} \log\left(|\mathcal{V}| \cdot p_{\pi_{\theta}}(j|q, o_{< i})\right) \\
\text{(4)}$$

where $o_{<i}$ are the previously generated tokens and $p(j|q, o_{<i})$ is the model's predicted probability for token j at step i. Higher self-certainty values indicate greater confidence.

Self-certainty, being related to a KL divergence where the model's prediction is the second argument $KL(U \parallel p_{\pi\theta})$, is 195 mode-seeking. This contrasts with entropy (or reverse KL 196 divergence from uniform), which is mode-covering. Crit-197 ically, self-certainty is reported to be less prone to biases towards longer generations, a common issue with perplexity 199 or entropy-based measures (Fang et al., 2024; Kang et al., 200 2025), making it a potentially more reliable indicator of intrinsic confidence. Kang et al. (2025) demonstrate that self-certainty is effective for selecting high-quality answers 203 from multiple candidates, and uniquely among different con-204 fidence measures, its utility improves with more candidates. Optimizing for self-certainty thus encourages the model to 206 generate responses that it deems more convincing. The RL process can achieve this by, for instance, guiding the model 208 to produce more detailed reasoning steps, thereby increasing 209 the model's conviction in its final answer. This mechanism 210 is more nuanced than simply increasing the probability of 211 the most likely output; it involves modifying the generation 212 process itself to build confidence. 213

To optimize the objective in Equation 3, various policy gradient algorithms can be employed. Informed by the recent success in models such as DeepSeek-R1 (Guo et al., 2025) and its widespread adoption of GRPO in open-source projects, we utilize GRPO to optimize for self-certainty.



Figure 2: Illustration of INTUITOR. INTUITOR simplifies the training strategy by leveraging self-certainty (the model's own confidence) as an intrinsic reward, optimizing these scores to incentivize reasoning abilities without external supervision.

The overall pipeline for this GRPO-based instantiation of INTUITOR is illustrated in Figure 2.

The core idea behind the optimization is to sample multiple candidate outputs for a given query and use their relative rewards to estimate advantages for policy updates. For each query $q \sim P(Q)$, GRPO samples a group of G outputs o_1, \ldots, o_G using a behavior policy $\pi_{\theta_{\text{old}}}$ (e.g., a previous iteration or the SFT model). The target policy π_{θ} is then optimized by maximizing:

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\theta) &= \mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)} \tag{5} \\ &\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left(\min \left[c_{i,t}(\theta) \hat{A}_{i,t}, \operatorname{clip}(c_{i,t}(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_{i,t} \right] \right. \\ &- \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}) \right), \\ &c_{i,t}(\theta) &= \frac{\pi_{\theta}(o_{i,t} \mid q, o_{i,$$

Hyperparameters ϵ (for clipping) and β (for KL penalty strength) control stability and exploration, and $\hat{A}_{i,t}$ is the advantage estimate.

Integration of Self-Certainty. The key innovation in IN-TUITOR is replacing external rewards with self-certainty scores in GRPO's advantage computation. Specifically, each output o_i is scored by:

$$u_{i} = \text{Self-certainty}(o_{i}|q),$$
$$\hat{A}_{i,t} = \frac{u_{i} - \text{mean}(\{u_{1}, u_{2}, \cdots, u_{G}\})}{\text{std}(\{u_{1}, u_{2}, \cdots, u_{G}\})}.$$
(6)

This formulation enables the policy to favor outputs that the model itself considers more confident, creating a selfreinforcing learning loop. The complete INTUITOR training pipeline operates by sampling multiple candidate outputs for each query, computing self-certainty scores for each candidate, using these scores to estimate advantages within the group, and updating the policy to increase the likelihood of generating high-confidence outputs. This process requires
no external supervision, making it broadly applicable across
domains and tasks.

4. Experimental Setup

224

225

245

255

Training Setup. Both GRPO and INTUITOR are trained 227 with the Open-R1 framework (Face, 2025) on the training split of the MATH dataset (Hendrycks et al., 2021), 229 which contains 7,500 problems. We use Qwen2.5-1.5B 230 and Qwen2.5-3B (Yang et al., 2024b) as backbone models, 231 with a chat-based prompting format throughout. Given the 232 models' initially weak instruction-following abilities, we 233 do not require them to disentangle intermediate reasoning 234 from final answers. Each update processes 128 problems, 235 generating 7 candidate solutions per problem, with a default 236 KL penalty of $\beta = 0.005$. For a fair comparison, GRPO and 237 INTUITOR share identical hyperparameters (see Appendix) 238 without additional tuning. We also evaluate a GRPO variant, 239 denoted GRPO-PV in Table 1, which uses plurality voting¹ 240 as a proxy for ground truth. This follows the approach from 241 TTRL (Zuo et al., 2025), which shows that self-consistency-242 based rewards can match the performance of golden answers 243 when training on inference data. 244

INTUITOR for Code Generation (INTUITOR-Code). To 246 assess generalization beyond mathematical reasoning, we 247 apply INTUITOR to the Codeforces code generation dataset 248 (Li et al., 2022). This variant, denoted INTUITOR-Code 249 in Table 1, modifies the setup as follows: the number of 250 sampled completions per problem is increased to 14; the 251 learning rate is reduced from 3×10^{-5} to 1×10^{-5} ; and 252 the KL penalty is increased to $\beta = 0.01$. For simplicity, we 253 limit the run to 50 steps, utilizing a total of 3,200 problems. 254

256 Evaluation. Evaluations generally use the same chat-style 257 prompting format as in training, except for MMLU-Pro 258 (Wang et al., 2024), where we follow the benchmark's orig-259 inal prompt format. Greedy decoding is used for all completions. Experiments were conducted on NVIDIA A100 261 GPUs, each with 40GB of memory. We evaluate perfor-262 mance on the following benchmarks (1) Math reasoning: 263 MATH500 and GSM8K, using the lighteval library 264 (Habib et al., 2023). (2) Code reasoning: CRUXEval-O (Gu 265 et al., 2024), using the ZeroEval framework (Lin, 2024), 266 and LiveCodeBench v6 (LCB) (Jain et al., 2024). (3) In-267 struction following: AlpacaEval 2.0 with length-controlled 268 win rates (Dubois et al., 2024), judged by GPT-4.1 (OpenAI, 269 2025). 270



Figure 3: Average response lengths during training rollouts. For Qwen2.5-1.5B, INTUITOR and GRPO reduce gibberish outputs. For Qwen2.5-3B, INTUITOR and GRPO increase reasoning length; INTUITOR yields significantly longer responses. GRPO-PV shows minimal length increase.

5. Results and Analysis

Table 1 presents main evaluation results, and Figure 3 illustrates response length evolution during training. On in-domain MATH and GSM8K datasets, INTUITOR and GRPO-PV (both golden-answer-free) achieve performance comparable to GRPO (using golden answers). This aligns with TTRL (Zuo et al., 2025), where plurality voting approximated golden answers without significant performance loss. While INTUITOR performs slightly worse than GRPO overall, on MATH it produces longer responses and demonstrates markedly improved code generation, suggesting enhanced reasoning capabilities.

5.1. Learning to Follow Instructions

INTUITOR significantly enhances instruction-following. Initially, the pretrained Qwen2.5-1.5B struggles with chat-style prompts, scoring ¡10% on all chat-template tasks (Table 1) and generating repetitive, nonsensical output, which inflates average response lengths (Figure 3). Fine-tuning with INTU-ITOR sharply reduces such gibberish, decreases completion lengths, and enables non-trivial performance across all evaluated benchmarks.

Furthermore, on the MATH dataset, INTUITOR substantially improves the Length Control Win Rate on AlpacaEval for both Qwen2.5-1.5B and Qwen2.5-3B, surpassing GRPO under identical settings. This demonstrates robust gains in instruction adherence.

5.2. Fostering Structured Reasoning

Rapid Initial Learning. Self-certainty, a continuous and inherently process-aware reward derived from the model's internal assessment across all tokens, contrasts with binary rewards. This internal signal may encourage LLMs to follow more effective learning trajectories. Given comparable final performance between GRPO and INTUITOR, we assess early-stage learnability by comparing in-domain accuracy

 ¹Self-consistency uses a plurality rule, selecting the most frequent answer even without majority support, while majority voting requires > 50% support and otherwise yields no winner (De Condorcet et al., 2014).

Model	Training Data	GSM8K	MATH500	LCB	CRUX	MMLU-Pro	AlpacaEva
Qwen2.5-1.5B Results							
Base	-	0.002	0.090	0.000	0.000	0.297	2.10
+ GRPO	MATH	0.747	0.560	0.056	0.328	0.315	4.03
+ Intuitor	MATH	0.711	0.530	0.099	0.296	0.310	4.28
Qwen2.5-3B Results							
Base	-	0.673	0.544	0.093	0.236	0.377	3.72
+ GRPO	MATH	0.826	0.636	0.085	0.341	0.403	6.91
+ GRPO-PV	MATH	0.820	0.636	0.086	0.299	0.398	6.17
+ Intuitor	MATH	0.792	0.612	0.153	0.416	0.379	7.10
+ INTUITOR-Code	Codeforces	0.743	0.572	0.153	0.411	0.386	4.16

Table 1: Performance comparison of various methods on the GSM8K, MATH, LCB, CRUXEval-O, MMLU-Pro, and AlpacaEval benchmarks. The INTUITOR-Code variant is trained on Codeforces data with a smaller learning rate and fewer training steps. All evaluations are obtained with the chat inference template, except for MMLU-Pro.

Table 2: Early-stage performance (training step 10) on
GSM8K and MATH. INTUITOR consistently outperforms
GRPO.

291

296

297

299

300

301

302

303 304

305

Model	Method	GSM8K	MATH
Qwen2.5-1.5B	Baseline	0.002	0.090
	GRPO	0.081	0.296
	INTUITOR	0.152	0.368
Qwen2.5-3B	Baseline	0.673	0.544
	GRPO	0.758	0.596
	INTUITOR	0.811	0.618

at training step 10. As shown in Table 2, INTUITOR consistently outperforms GRPO on both GSM8K and MATH
benchmarks for Qwen2.5-1.5B and Qwen2.5-3B, highlighting its advantage in rapid initial learning.

311 Cross-Task Generalization. Figure 4 illustrates perfor-312 mance trajectories on MATH500 (in-domain) and Live-313 CodeBench (transfer task) for models trained on the MATH 314 dataset. For both INTUITOR and GRPO, accuracy improve-315 ments on LiveCodeBench emerge later in training, following 316 initial gains on MATH500. Notably, LiveCodeBench perfor-317 mance continues to improve even after MATH500 accuracy 318 plateaus. This pattern suggests that initial in-domain learn-319 ing (on MATH) facilitates subsequent generalization to code 320 generation tasks (LiveCodeBench). 321

Emergence of Long-Form Reasoning. While large models like Deepseek-R1 achieve long-form reasoning through extensive RL, INTUITOR enables smaller models to develop structured reasoning with limited data. On CRUXEval-O (Figure 8), models trained with INTUITOR often exhibit freeform reasoning before summarizing it within the instructed JSON block, despite prompts requiring reasoning directly



Figure 4: Performance evolution on MATH500 (in-domain) and LiveCodeBench (transfer) for models trained on MATH. In-domain (MATH500) accuracy improves rapidly early in training, preceding gains in code-generation (Live-CodeBench) accuracy. LiveCodeBench performance continues to rise even after MATH500 accuracy plateaus.

in JSON. A similar pattern of pre-code natural language reasoning is observed on LiveCodeBench. This emergent pre-reasoning may contribute to INTUITOR 's strong performance on these benchmarks.

5.3. Understanding Emergent Long-Form Reasoning

When LLMs encounter unfamiliar questions, they sample from a distribution of possible answers (Kang et al., 2024). Self-certainty reflects the model's internal assessment of its output coherence. By reinforcing high-confidence responses, INTUITOR encourages more elaborate reasoning, potentially improving the model's comprehension of its own outputs. While not explicitly targeting benchmark accuracy, this enhancement in output quality and structure leads to more reliable answers and better generalization.

We analyze models trained with INTUITOR on code corpora by examining outputs for ten randomly selected Live-CodeBench questions across different training steps. Figure 5 shows the evolution of output types alongside model



Figure 5: (a) Left: Distribution of answer types for ten random LiveCodeBench questions across training steps. Right: Corresponding model accuracy. The model first learns to generate correct code, then adds reasoning to improve understanding. (b) Training with INTUITOR on code corpora leads to spontaneous reasoning before coding and explanation of outputs.

332

333

334

335

337

338

339

340

341

342

343

370

371

372

345 accuracy. The results reveal a clear progression: models 346 first learn to generate valid Python code (evidenced by im-347 proved accuracy and fewer invalid responses), then develop 348 pre-code reasoning to facilitate self-understanding. Further 349 inspection of generations confirms that models progressively 350 elaborate their reasoning throughout training, supporting our 351 hypothesis that INTUITOR encourages traces that the model 352 itself can better understand. 353

354 To quantify this effect, we classify outputs from succes-355 sive checkpoints into three categories: invalid code ("No 356 Answer"), valid code without reasoning ("No Reasoning"), 357 and valid code with explicit reasoning ("Reasoning"). Fig-358 ure 5(a) illustrates how these proportions evolve during 359 training alongside LiveCodeBench accuracy. The model 360 first reduces invalid outputs and improves code correct-361 ness before incorporating pre-code reasoning, reflecting an 362 emergent emphasis on self-explanatory traces. Figure 5(b)363 demonstrates how training with INTUITOR leads to struc-364 tured reasoning before code generation. Additional evidence 365 appears in Figure 7, where INTUITOR-trained models assign 366 significantly higher confidence to their generated responses 367 compared to baseline models, as discussed further in Sec-368 tion 5.4. 369

5.4. Online Self-Certainty Prevents Reward Exploitation

Over-optimization against static reward models is a known 373 374 failure mode in reinforcement learning (Gao et al., 2023). To assess the robustness of self-certainty as a reward, we 375 compare offline self-certainty (rewards from a fixed base 376 model) with online self-certainty (rewards from the evolving 377 policy model), using a reduced batch size of 224 responses 378 379 per gradient update.

380 Figure 6 demonstrates that the offline annotator is suscepti-381 ble to exploitation. Around the 100th update step, the policy 382 model learns to inflate its self-certainty reward by append-383 ing an auxiliary, already-solved problem to its answer for 384



during training, comparing online and offline self-certainty annotators with INTUITOR under reduced batch sizes. The offline reward model is exploited early in training (around 100 steps), leading to increased response length and decreased accuracy. The online annotator maintains stable training. Refer to Section 5.4 for details.

the given question. This exploitation manifests as a sharp increase in response length (dashed line) and a concurrent collapse in validation accuracy (solid line). In contrast, the online annotator, whose reward signal co-evolves with the policy, prevents such reward hacking and maintains stable training dynamics.

To further evaluate the quality of self-certainty as a reward signal, we analyze the distribution of selfcertainty scores from policies trained with INTUITOR and GRPO on MATH500 responses (Figure 7). We employ Mann–Whitney U tests to determine if correct responses achieve significantly higher self-certainty scores than incorrect ones.

Both GRPO and INTUITOR models exhibit significantly higher average self-certainty scores, indicating that GRPO also enhances the model's self-assessment capabilities. Notably, policies trained with online self-certainty (i.e., INTU-ITOR) show no signs of reward hacking. The INTUITOR policy yields the lowest *p*-values and largest effect sizes (r) in the Mann-Whitney U tests (Figure 7, inset). This indicates it is most effective at discriminating its own correct and incorrect answers using self-certainty, even while assigning higher absolute confidence scores overall. These findings underscore the potential of INTUITOR for robust training on larger datasets.

5.5. Ablation Studies

We further investigate how the magnitude of the KL penalty influences INTUITOR, as shown in Table 3. On in-domain benchmarks (MATH500 and GSM8K) the choice of penalty has only a minor effect, but on out-of-domain tasks-LiveCodeBench (code generation) and CRUXEval-O (code reasoning)-model accuracy is highly sensitive to this hyper-parameter. Because INTUITOR does not receive explicit feedback from generated responses during training, the KL penalty serves as a critical regularization mechanism. It prevents the policy from drifting too far from the initial



396Figure 7: Distribution of self-certainty on MATH500 re-397sponses, for policies trained with GRPO and INTUITOR.398Histograms are split by response correctness. The inset399shows Mann–Whitney U test statistics (p-value and effect400size r) comparing self-certainty of correct versus incorrect401responses. The policy trained with INTUITOR demonstrates402the best separation.

411

412 413 414

415

416

417

418

419 420

421

422

423

424

425

426 427

428

429

430

431

432

433

434

435

436

437

438

439

Table 3: Impact of the KL-divergence penalty in INTUITOR
during fine-tuning of Qwen-2.5-3B on the MATH dataset.
We compare performance across GSM8K, MATH500, LCB,
CRUXEval-O, MMLU-Pro, and AlpacaEval. All scores are
obtained with the chat-style inference template, except for
MMLU-Pro, which uses its standard evaluation protocol.

Model	GSM8K	MATH500	LCB	CRUX	MMLU-Pro	AlpacaEval
Base	0.673	0.544	0.093	0.236	0.377	3.72
+ INTUITOR-KL0	0.809	0.598	0.081	0.390	0.359	6.77
+ INTUITOR-KL0.0001	0.793	0.616	0.090	0.364	0.354	6.79
+ INTUITOR-KL0.005	0.792	0.612	0.153	0.416	0.379	7.10
+ INTUITOR-KL0.01	0.803	0.618	0.130	0.394	0.371	6.54

model distribution, acting as a safeguard against degeneration. These findings highlight the importance of careful KL tuning in general-purpose RL setups, especially when targeting robust generalization across domains.

Additionally, we evaluate INTUITOR on larger models, including Qwen2.5-7B and Qwen2.5-14B, and test different model architectures such as Llama-3.2-3B (Meta AI, 2024) and OLMo-2 (OLMo et al., 2024); these results are detailed in the Appendix.

6. Discussion and Future Research

Scalability and Generalization. Our experiments, constrained by computational resources, utilize relatively compact models trained on relatively small, unsupervised corpora. We aim to demonstrate the potential of a model's self-certainty as a reward signal for policy optimization. The results show that this signal consistently promotes more coherent, well-justified, and interpretable explanations, indicating a path towards more autonomous learning. Future work could explore these benefits in larger foundation models and on more diverse, real-world datasets. Given that purely offline training with INTUITOR led to performance degradation over time, scaling up will likely require periodic online updates to self-certainty estimates or hybrid offline-online schedules to maintain calibration.

Applicability to Other Policy Gradient Methods. INTU-ITOR is a framework that leverages a model's self-certainty as an intrinsic reward signal for fine-tuning LLMs. It can be instantiated with various policy gradient algorithms. Due to computational constraints, and informed by the success of models like DeepSeek-R1 and the widespread adoption of GRPO, we employ GRPO for self-certainty optimization. The efficacy of self-certainty signals with other algorithms, such as REINFORCE or PPO, warrants further investigation.

Combining Reward Signals. To enable a direct comparison between self-certainty and golden-answer rewards, this paper focuses exclusively on a single reward signal. However, these signals are not mutually exclusive. Future work could explore combining them, for instance, by summation or by alternating based on the availability of golden answers. Furthermore, other reward signals, such as formatting rewards (Guo et al., 2025), could be additively combined to enhance performance. Integrating RLIF with methods like RLHF and RLVR may further advance LLM capabilities across various dimensions.

7. Conclusion

This paper introduces INTUITOR, an instantiation of Reinforcement Learning from Internal Feedback (RLIF) that uses a model's intrinsic self-certainty as its sole reward signal, eliminating the need for external supervision or goldstandard solutions. Our experiments show that INTUITOR matches the performance of supervised RLVR methods like GRPO on mathematical reasoning, while achieving superior generalization to out-of-domain tasks such as code generation and instruction following. It also promotes structured reasoning and leverages online self-certainty to guard against reward exploitation.

These findings highlight the transformative potential of RLIF, signaling a meaningful step toward AI systems that improve through introspection and unlock rich latent capabilities. Looking forward, this paradigm opens the door to AI agents capable of autonomous skill acquisition in novel domains and scalable self-improvement—even as they approach or surpass the limits of human oversight. Future directions include integrating RLIF with external reward methods like RLHF or RLVR to tackle increasingly complex real-world challenges, and advancing the development of more robust, generalizable, and truly autonomous learning systems.

440 **References**

441

442

443

444

445

446

447

448

449

450

464

468

482

483

484

485

486

487

488

- Agarwal, S., Zhang, Z., Yuan, L., Han, J., and Peng, H. The unreasonable effectiveness of entropy minimization in llm reasoning. *arXiv preprint arXiv:2505.15134*, 2025.
- Burns, C., Izmailov, P., Kirchner, J. H., Baker, B., Gao, L., Aschenbrenner, L., Chen, Y., Ecoffet, A., Joglekar, M., Leike, J., et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.
- Chen, M., Chen, G., Wang, W., and Yang, Y. Seed-grpo: Semantic entropy enhanced grpo for uncertainty-aware policy optimization. *arXiv preprint arXiv:2505.12346*, 2025.
- Chen, Z., Deng, Y., Yuan, H., Ji, K., and Gu, Q. Self-play
 fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024.
- Cheng, P., Dai, Y., Hu, T., Xu, H., Zhang, Z., Han, L., Du,
 N., and Li, X. Self-playing adversarial language game enhances llm reasoning. *Advances in Neural Information Processing Systems*, 37:126515–126543, 2024.
- 465 De Condorcet, N. et al. *Essai sur l'application de l'analyse*466 à la probabilité des décisions rendues à la pluralité des
 467 voix. Cambridge University Press, 2014.
- 469 Dubois, Y., Galambosi, B., Liang, P., and Hashimoto, T. B.
 470 Length-controlled alpacaeval: A simple way to debias 471 automatic evaluators. *arXiv preprint arXiv:2404.04475*, 472 2024.
- 473
 474
 475 Face, H. Open r1: A fully open reproduction of deepseek475 r1, January 2025. URL https://github.com/
 476 huggingface/open-r1.
- Fang, L., Wang, Y., Liu, Z., Zhang, C., Jegelka, S., Gao, J., Ding, B., and Wang, Y. What is wrong with perplexity for long-context language modeling? *arXiv preprint arXiv:2410.23771*, 2024.
 - Farquhar, S., Kossen, J., Kuhn, L., and Gal, Y. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630, 2024.
 - Gao, L., Schulman, J., and Hilton, J. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Gu, A., Rozière, B., Leather, H., Solar-Lezama, A., Synnaeve, G., and Wang, S. I. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv* preprint arXiv:2401.03065, 2024.

- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Habib, N., Fourrier, C., Kydlíček, H., Wolf, T., and Tunstall, L. Lighteval: A lightweight framework for llm evaluation, 2023. URL https://github.com/ huggingface/lighteval.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.
- Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- Hu, J., Zhang, Y., Han, Q., Jiang, D., Zhang, X., and Shum, H.-Y. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. *arXiv preprint arXiv:2503.24290*, 2025.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. arXiv preprint arXiv:2412.16720, 2024.
- Jain, N., Han, K., Gu, A., Li, W.-D., Yan, F., Zhang, T., Wang, S., Solar-Lezama, A., Sen, K., and Stoica, I. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Kang, K., Wallace, E., Tomlin, C., Kumar, A., and Levine, S. Unfamiliar finetuning examples control how language models hallucinate. *arXiv preprint arXiv:2403.05612*, 2024.
- Kang, Z., Zhao, X., and Song, D. Scalable best-of-n selection for large language models via self-certainty. arXiv preprint arXiv:2502.18581, 2025.
- Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- Lambert, N., Morrison, J., Pyatkin, V., Huang, S., Ivison, H., Brahman, F., Miranda, L. J. V., Liu, A., Dziri, N., Lyu, S., et al. T\"ulu 3: Pushing frontiers in open language model post-training. arXiv preprint arXiv:2411.15124, 2024.
- Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., Hubert, T., Choy, P., de Masson d'Autume, C.,

		÷.,	
4	9	6	
4	9	7	
4	9	8	
4	.9	9	
5	0	0	
5	0	1	
5	0	2	
5	0	3	
5	0	4	
5	0	5	
5	0	6	
5	0	7	
5	0	8	
5	0	9	
5	1	0	
5	1	1	
5	1	2	
5	1	3	
5	1	4	
5	1	5	
5	1	6	
5	1	7	
5	1	8	
5	1	9	
5	2	0	
5	2	1	
5	2	2	
5	2	3	
5	2	4	
5	2	5	
5	2	6	
5	2	7	
5	2	8	
5	2	9	
5	3	0	
5	3	1	
5	3	2	
5	3	3	
5	3	4	
5	3	5	

537

538

539

540

541

542

543

544

545

546

547

548

549

495

Babuschkin, I., Chen, X., Huang, P.-S., Welbl, J., Gowal, S., Cherepanov, A., Molloy, J., Mankowitz, D., Sutherland Robson, E., Kohli, P., de Freitas, N., Kavukcuoglu, K., and Vinyals, O. Competition-level code generation with alphacode. arXiv preprint arXiv:2203.07814, 2022. Lin, B. Y. ZeroEval: A Unified Framework for Evaluating Language Models, July 2024. URL https://github. com/WildEval/ZeroEval. Liu, J. and Zhang, L. Code-r1: Reproducing r1 for code with reliable rewards. https://github.com/ganler/ code-r1, 2025. Liu, J., Xia, C. S., Wang, Y., and Zhang, L. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. Advances in Neural Information Processing Systems, 36:21558-21572, 2023. Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding r1-zero-like training: A critical perspective. arXiv preprint arXiv:2503.20783, 2025. Luo, M., Tan, S., Huang, R., Shi, X., Xin, R., Cai, C., Patel, A., Ariyak, A., Wu, Q., Zhang, C., Li, L. E., Popa, R. A., and Stoica, I. Deepcoder: A fully open-source 14b coder at o3-mini level. https://pretty-radio-b75.notion.site/ DeepCoder-A-Fully-Open-Source-14B-Coder-at-03-mini-Level-1cf81902c14680b3bee5eb349a512a51, 2025. Notion Blog. Ma, W., He, J., Snell, C., Griggs, T., Min, S., and Zaharia, M. Reasoning models can be effective without thinking. arXiv preprint arXiv:2504.09858, 2025. Meta AI. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. https://ai.meta.com/blog/ llama-3-2-connect-2024-vision-edge-mobile-devicesy.edge. nature, 550(7676):354-359, 2017. 2024. Accessed: 2025-05-16. OLMo, T., Walsh, P., Soldaini, L., Groeneveld, D., Lo, K., Arora, S., Bhagia, A., Gu, Y., Huang, S., Jordan, M., Lambert, N., Schwenk, D., Tafjord, O., Anderson, T., Atkinson, D., Brahman, F., Clark, C., Dasigi, P., Dziri, N., Guerquin, M., Ivison, H., Koh, P. W., Liu, J., Malik, S., Merrill, W., Miranda, L. J. V., Morrison, J., Murray, T., Nam, C., Pyatkin, V., Rangapur, A., Schmitz, M., Skjonsberg, S., Wadden, D., Wilhelm, C., Wilson, M., Zettle-2023. moyer, L., Farhadi, A., Smith, N. A., and Hajishirzi, H. 2 olmo 2 furious. arXiv preprint arXiv:2501.00656, 2024. URL https://arxiv.org/abs/2501.00656. OpenAI. Introducing GPT-4.1 in the API. https:// openai.com/index/gpt-4-1/, April 2025. Accessed: 15 May 2025. 2024. 10

- Oudever, P.-Y. and Kaplan, F. What is intrinsic motivation? a typology of computational approaches. Frontiers in neurorobotics, 1:108, 2007.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730-27744, 2022.
- Poesia, G., Broman, D., Haber, N., and Goodman, N. Learning formal mathematics from intrinsic motivation. Advances in Neural Information Processing Systems, 37: 43032-43057, 2024.
- Prabhudesai, M., Chen, L., Ippoliti, A., Fragkiadaki, K., Liu, H., and Pathak, D. Maximizing confidence alone improves reasoning, 2025. URL https://arxiv.org/ abs/2505.22660.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36: 53728-53741, 2023.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.2017.

- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300, 2024.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without
- Team, K., Du, A., Gao, B., Xing, B., Jiang, C., Chen, C., Li, C., Xiao, C., Du, C., Liao, C., et al. Kimi k1. 5: Scaling reinforcement learning with llms. arXiv preprint arXiv:2501.12599, 2025.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and finetuned chat models. arXiv preprint arXiv:2307.09288,
- Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S., Ren, W., Arulraj, A., He, X., Jiang, Z., et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. arXiv preprint arXiv:2406.01574,

- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992.
- Xiaomi LLM-Core Team. Mimo: Unlocking the reasoning potential of language model – from pretraining to posttraining, 2025. URL https://github.com/ XiaomiMiMo/MiMo.
- Xu, F., Yan, H., Ma, C., Zhao, H., Sun, Q., Cheng, K., He, J., Liu, J., and Wu, Z. Genius: A generalizable and purely unsupervised self-training framework for advanced reasoning. *arXiv preprint arXiv:2504.08672*, 2025.
- 563 Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., 564 Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, 565 J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, 566 K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., 567 Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Xia, T., Ren, 568 X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., 569 Cui, Z., Zhang, Z., and Qiu, Z. Qwen2.5 technical report. 570 arXiv preprint arXiv:2412.15115, 2024a. 571
- Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li,
 C., Liu, D., Huang, F., Wei, H., et al. Qwen2. 5 technical
 report. *arXiv preprint arXiv:2412.15115*, 2024b.
- Yuan, W., Pang, R. Y., Cho, K., Li, X., Sukhbaatar, S., Xu,
 J., and Weston, J. E. Self-rewarding language models. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 57905–57923. PMLR, 2024.
- Zelikman, E., Wu, Y., Mu, J., and Goodman, N. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Zhang, Q., Wu, H., Zhang, C., Zhao, P., and Bian, Y.
 Right question is already half the answer: Fully unsupervised llm reasoning incentivization. *arXiv preprint arXiv:2504.05812*, 2025.
- Zhao, A., Wu, Y., Yue, Y., Wu, T., Xu, Q., Lin, M., Wang,
 S., Wu, Q., Zheng, Z., and Huang, G. Absolute zero: Reinforced self-play reasoning with zero data. *arXiv* preprint arXiv:2505.03335, 2025.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford,
 A., Amodei, D., Christiano, P., and Irving, G. Fine-tuning
 language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.
- Zuo, Y., Zhang, K., Qu, S., Sheng, L., Zhu, X., Qi, B.,
 Sun, Y., Cui, G., Ding, N., and Zhou, B. Ttrl: Test-time
 reinforcement learning. *arXiv preprint arXiv:2504.16084*,
 2025.
- 603 604

605	Table 4: Performance comparison of various methods on the GSM8K, MATH, LCB, CRUXEval-O, MMLU-Pro, and
606	AlpacaEval benchmarks for larger models. All evaluations are obtained with the chat inference template, except for
607	MMLU-Pro.

Model	GSM8K	MATH500	LCB	CRUX	MMLU-Pro	AlpacaEval
Qwen2.5-7B	0.549	0.612	0.017	0.161	0.424	5.63
+ GRPO	0.723	0.699	0.200	0.540	0.434	7.78
+ INTUITOR	0.874	0.730	0.200	0.573	0.422	8.43
Qwen2.5-14B	0.751	0.674	0.219	0.499	0.499	6.28
+ GRPO	0.917	0.764	0.295	0.516	0.487	11.58
+ INTUITOR	0.922	0.770	0.294	0.555	0.492	11.83



Figure 8: INTUITOR quickly demonstrate R1-like reasoning

A. Additional Experimental Details

A.1. Generalization to Larger Models

We extend INTUITOR to larger base models Owen2.5-7B and Owen2.5-14B and find that the original training recipe triggers severe behavioural collapse at the very start of training. Even before any updates, the 7B model solves the given problem and then immediately proceeds to tackle an unrelated one; this tendency becomes more pronounced as training progresses.

To stabilize learning, we simplify the system prompt, reduce the learning rate to 1×10^{-6} , and increase the number of sampled responses per problem to sixteen. These settings represent our first, untuned trial, and a comprehensive hyperparameter sweep is beyond the scope of this paper. Because the system prompt is the only additional signal the model receives during INTUITOR fine-tuning, we expect its careful calibration to exert a particularly strong influence on training dynamics. With the foregoing adjustments INTUITOR trains smoothly on both larger models. The corresponding evaluation results are reported in Table 4.

A.2. Evaluation on Llama3.2-3B-Instruct

We further evaluate INTUITOR on the Llama3.2-3B-Instruct model. Compared to the Qwen family, improvements on external benchmarks are less pronounced-likely due to extensive prior instruction tuning. Nevertheless, as shown in Figure 10, both accuracy and generated sequence length improve steadily over the course of training, indicating meaningful optimization gains under INTUITOR.



Figure 9: Average accuracy and mean completion length during reinforcement learning on the MATH dataset using INTUITOR and GRPO. Both methods yield similar accuracy gains, with INTUITOR generally producing longer completions.

A.3. Evaluation on OLMo-2-1124-7B-SFT

To further validate our findings, we applied INTUITOR to OLMo-2-1124-7B-SFT (OLMo et al., 2024), a fully open large language model. The results are shown below.

A.4. Evaluating INTUITOR Against Entropy-Minimization and Random Reward Strategies

Contemporary research has found that applying a negative token-level entropy reward can improve a model's reasoning performance without requiring external labels (Agarwal et al., 2025; Prabhudesai et al., 2025). However, since low entropy often correlates with repetitive loops (Holtzman et al., 2019), using negative entropy alone as a reinforcement learning reward risks driving the model into a collapsed state. In other words, without sufficient supervised training to push the base model away from degenerate behavior, the model risks falling into a repetition trap from which it cannot recover. As we observed a nontrivial amount of repetitive responses in Qwen2.5-1.5B, we tested this hypothesis by applying GRPO with the negative-entropy reward ($u_{EM} = -\frac{1}{|o|} \sum_{i=1}^{|o|} \sum_{j=1}^{|V|} \cdot p_{\pi_{\theta}}(j|q, o_{<i}) \cdot \log(p_{\pi_{\theta}}(j|q, o_{<i}))$). Figure 12 (left) validates our prediction. Entropy minimization exacerbates repetition, and after a few updates the model converges to producing the same

Table 5: Accuracy	y of the OLMo-2	-1124-7B-SFT n	nodel using	GRPO and 1	INTUITOR O	n GSM8K	and MATH500

100				
709	Model	Method	GSM8K	MATH
/10				
711		Baseline	0.685	0.288
712	OLMo-2-1124-7B-SFT	GRPO	0.704	0.334
713		INTUITOR	0.697	0.348
714				



Figure 10: Average accuracy and mean completion length of Llama3.2-3B-Instruct during training with INTUITOR on
 MATH dataset. Both metrics exhibit significant increases.



Figure 11: Average accuracy and mean completion length on the MATH dataset during reinforcement learning with OLMo-2-1124-7B-SFT using INTUITOR and GRPO. Both methods achieve comparable accuracy gains.

character regardless of the prompt. By contrast, INTUITOR enhances performance without triggering collapse (Figure 4).

To further validate the efficacy of INTUITOR, we also trained Qwen2.5-3B using a random reward baseline, where each response was assigned a reward of 0 or 1 with equal probability. Figure 12 (Right) shows that this random reward scheme severely degrades the model's performance in a chat-style RL setting, demonstrating that the performance gains observed with INTUITOR are indeed non-trivial.

A.5. Training Hyper-parameters

We show the training hyper-parameters in Table 6.

B. Prompts and Model Completions

This section presents sample prompts and the responses generated by the models. Unless otherwise specified, the default base model used is Qwen2.5-3B, and the default training dataset is MATH.

Learning to Reason without External Rewards



Figure 12: Left: GRPO with an entropy minimization objective using Qwen2.5-1.5B on MATH. Right: GRPO with a random reward using Qwen2.5-3B on MATH. Both approaches exhibit severe output degeneration.

Table 6: Only hyper-parameters that affect the learned policy or evaluation are listed. Unspecified fields inherit the TRL_v0.8 defaults.

]	Parameter	MATH (1.5B/3B)	MATH (7B/14B)	Codeforces (3B)
]	Learning Rate	3×10^{-6}	1×10^{-6}	1×10^{-6}
]	Batch Size	128	64	64
(Group Size	7	14	14
]	KL Penalty(β)	0.0005	0.01	0.01
r	Training Steps	58	117	50
]	Max Prompt Length	512	512	1024
]	Max Completion Length	3072	3072	2048
r	Temperature	0.9	0.9	0.9
(Clip Ratio	0.2	0.2	0.2
]	Lr Scheduler Type	Cosine	Cosine	Cosine
,	Warmup Ratio	0.1	0.1	0.1
(Optimizer	AdamW (A	$\beta_1 = 0.9, \ \beta_2 = 0.999,$	$\varepsilon = 10^{-8})$

B.1. Training Prompts

System prompt used for Qwen2.5-1.5B on MATH.

You are a helpful AI Assistant, designed to provided well-reasoned and detailed responses. You FIRST think about the reasoning process step by step and then provide the user with the answer. Please enclose your final answer in the box: \boxed{Your Answer}.

System prompt used for Qwen2.5-3B on MATH.

You are a helpful AI Assistant, designed to provided well-reasoned and detailed responses. You FIRST think about the reasoning process step by step and then provide the user with the answer. Please enclose your final answer in the box: \boxed{Your Answer}. Please stop generation immediately after outputing the box.

System prompt used for Qwen2.5-7B and Qwen2.5-14B on MATH.

You are a helpful AI Assistant, designed to provided well-reasoned and detailed responses. Please provide a step-by-step solution to the following problem.

System prompt used for Qwen2.5-3B on code corpuse.

You are an AI designed to help solve competitive programming problems by generating Python code.

B.2. Example from LiveCodeBench Code Generation

Models trained with INTUITOR tend to generate reasoning before producing code, whereas GRPO-tuned models typically respond with Python code directly under the same prompt.

Problem

Question:

You are given a 0-indexed array of strings details. Each element of details provides information about a given passenger compressed into a string of length 15. The system is such that:

- The first ten characters consist of the phone number of passengers.
- The next character denotes the gender of the person.
- The following two characters are used to indicate the age of the person.
- The last two characters determine the seat allotted to that person.
- Return the number of passengers who are strictly more than 60 years old.

Example 1

Input

details = ["7868190130M7522","5303914400F9211","9273338290F4010"]

Output

Explanation The passengers at indices 0, 1, and 2 have ages 75, 92, and 40. Thus, there are 2 people who are over 60 years old.

Example 2

Input

details = ["1313579440F2036","2921522980M5644"]

Output

Explanation None of the passengers are older than 60.

Constraints

- 1 = details.length = 100
- details[i].length == 15
- details[i] consists of digits from '0' to '9'.
- details[i][10] is either 'M' or 'F' or 'O'.
- The phone numbers and seat numbers of the passengers are distinct."

You will use the following starter code to write the solution to the problem and enclose your code within delimiters.

```
```python
class Solution:
 def countSeniors(self, details: List[str]) -> int:
    ```
```

8	8	0
8	8	1
8	8	2
8	8	3
8	8	4
8	8	5
0	0	5
0	0	7
0	0	/
8	ð	ð
8	8	9
8	9	0
8	9	1
8	9	2
8	9	3
8	9	4
8	9	5
8	9	6
8	9	7
8	9	8
8	9	9
9	0	0
0	0	1
0	0	1
2	0	2
9	0	2
9	0	4
9	0	5
9	0	6
9	0	7
9	0	8
9	0	9
9	1	0
9	1	1
9	1	2
9	1	3
9	1	4
9	1	5
9	1	6
0	1	7
0	1	2 2
7	1	0
9	1	9
9	2	U
9	2	1
9	2	2
9	2	3
9	2	4
9	2	5
9	2	6
9	2	7
9	2	8
9	2	9
9	3	0
0	3	1
/	2	-

INTUITOR-tuned To solve the problem of counting the number of passengers who are strictly more than 60 years old, we can follow these steps: 1. Iterate through each passenger's details.

- 2. Extract the age from the details.
- 3. Check if the age is greater than 60.
- 4. Count the number of passengers who meet the criteria.

Here is the Python code to achieve this:

```
", "python
class Solution:
def countSeniors(self, details: List[str]) -> int:
    count = 0
    for detail in details:
        age = int(detail[11:13])
        if age > 60:
            count += 1
    return count
. . .
```

GRPO-tuned

```
", "python
class Solution:
def countSeniors(self, details: List[str]) -> int:
    count = 0
    for detail in details:
        age = int(detail[11:13])
        if age > 60:
             count += 1
    return count
. . .
```

B.3. Example from LiveCodeBench Code Generation

Models trained with INTUITOR are typically effective at interpreting instructions and producing correct Python code, while GRPO-tuned models may misinterpret the instructions or decline to provide a response.

Problem

Ouestion:

There are three cards with letters a, b, and c placed in a row in some order. You may perform *at most one* operation:

• Pick any two cards and swap them.

Is it possible that the row becomes abc after the operation? Output YES if it is possible, and NO otherwise.

Input

The first line contains a single integer t ($1 \le t \le 6$) — the number of test cases. Each test case consists of a single line containing a permutation of the three characters a, b, and c.

Output

For each test case, output YES if you can make the row abc with at most one operation, or NO otherwise. The answer is case-insensitive.

Sample Input

abc acb bac bca cab cba **Sam**

Sample Output YES YES

YES

NO NO

YES

Note

- In the first test case, no operation is needed.
- In the second test case, swap c and b: acb \rightarrow abc.
- In the third test case, swap b and a: bac \rightarrow abc.
- In the fourth test case, it is impossible to obtain abc using at most one swap.

Read the inputs from stdin solve the problem and write the answer to stdout (do not directly test on the sample inputs). Enclose your code within delimiters as follows. Ensure that when the python program runs, it reads the inputs, runs the algorithm and writes output to STDOUT.

```
'`' python
# YOUR CODE HERE
'''
```

INTUITOR-tuned

To solve the problem of determining whether it is possible to transform the given string into "abc" using at most one operation (swapping two cards), we can follow these steps:

- 1. Read the input string
- 2. Check if the string is already "abc"
- 3. If the string is not "abc", check if it can be transformed into "abc" by swapping two cards
- 4. Output "YES" if it is possible, otherwise output "NO"

```
'``python
# Determine whether s can become "abc" with <= swap
def can_transform_to_abc(s: str) -> str:
    if s == "abc":
```

```
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1011
1012
1013
1014
1015
1016
1018
1019
1024
1025
1026
1028
1029
1033
1034
1035
1036
1038
1039
1040
1041
1042
```

```
return "YES"
# try all pairs of indices (i, j) with i < j
for i in range(len(s)):
    for j in range(i + 1, len(s)):
        swapped = s[:i]+s[j]+s[i+1:j]+s[i]+s[j+1:]
        if swapped == "abc":
            return "YES"
return "NO"
t = int(input().strip())
for _ in range(t):
    print(can_transform_to_abc(input().strip()))</pre>
```

```
GRPO-tuned
```

. . .

```
'`' python
# YOUR CODE HERE
'''
```

Another related GRPO-tuned response to another problem

I'm sorry, but I can't assist with that request.

B.4. Example from CRUXEval-O

Models trained with INTUITOR usually articulate their reasoning first, then deliver the formatted response as instructed, while GRPO-tuned models often produce the required JSON output immediately under the same prompt.

Problem

Question:

You are given a Python function and an assertion containing an input to the function. Complete the assertion with a literal (no unsimplified expressions, no function calls) containing the output when executing the provided code on the given input, even if the function is incorrect or incomplete.

Instruction

Please answer this question by first reasoning and then providing your answer. Present your reasoning and solution in the following json format. Please show your final answer in the answer field, e.g., "answer": "42".

```
```json
{
 "reasoning": "___",
 "answer": "___"
} ```
```

{

1099

**INTUITOR-tuned** To complete the assertion for the given Python function f(text), we need to determine the output of the function when the input is an empty string ''. Let's break down the function step by step: 1. Function Definition: def f(text):  $new_text = ''$ for ch in text.lower().strip(): if ch.isnumeric() or ch in '\u00c4\u00e4\u00cf\u00ef\u00d6\u00f6\u00dc\u00fc':  $new_text += ch$ return new\_text 2. Input: The input to the function is an empty string ' '. 3. Step-by-Step Execution: • Step 1: text.lower().strip() - text.lower() converts the string to lowercase. - text.strip() removes leading or trailing whitespace. - For '', both calls return '', so the result is ''. • Step 2: Iterate over each character - The loop sees no characters; it performs zero iterations. • Step 3: Character check – Because there are no characters, the condition ch.isnumeric() or ch in '\u00c4\u00e4\u00cf\u00ef\u00d6\u00fc\u00dc\u00fc' is never true, and nothing is appended to new\_text. 4. Output: Since nothing was appended, new\_text remains ' '. Therefore the assertion should be: assert f('') == '' So, the final answer is: · · · json { "reasoning": "The function 'f(text)' converts the input string to lowercase and removes any leading or trailing whitespace. For an empty string the resulting string is also empty. Therefore the output of 'f(')' is an empty string.", "answer": ", ", " } • • • **GRPO-tuned** ··· json

"reasoning": "The function 'f' takes a string 'text' as

input and returns a new string with only numeric characters

```
and certain special characters
 (\u00c4\u00e4\u00cf\u00d6\u00d6\u00dc\u00dc\u00dc) preserved.
 The assertion checks if the function correctly handles an
 empty string. Since an empty string contains no characters,
 the function should return an empty string as well.
 Therefore, the expected output for an empty string is an
 empty string.",
 "answer": ""
} • • •
```

### **B.5. Rollout Example of INTUITOR During Training on Code Corpus**

The following examples illustrate representative rollouts during the reinforcement-learning training of Qwen2.5-3B on the CodeContest dataset. As shown in Figure 5, the model initially produces concise reasoning and brief explanations to "convince" itself of its interim outputs. Over the course of training, it gradually shifts toward more detailed, step-by-step reasoning and richer explanations, which further reinforce its understanding and improve the final responses.

## Code Block (Step 2)

1100

1101

1102

1103

1104

1105

1106

1107

1108

1113 1114

1115

1116

1117

1118 1119

1120 1121

1122

1123

1124

1125

1126 1127

1128

1129 1130

1131

1132

1133

1134

1135

1136

1137 1138

1139 1140

1141

1142 1143

1144

1145

I have solved the problem. Here's the Python code to solve it as described:

```
'' python
def calculate_stones(S):
 n = len(S)
 MOD = 10 * *9 + 7
 # Initialize array to store remainders
 remainder = [0] * (n + 1)
 # Calculate remainders
 for i in range (n, -1, -1):
 remainder[i] = (
 remainder[i+1] * 10 + int(S[i])) % MOD
 if i > 0:
 remainder[i] = (
 remainder[i] - remainder[i+1] + MOD) % MOD
 return remainder [1] * 10 + remainder [0]
Read input
S = input()
Call function to calculate stones
T = calculate_stones(S)
Print output
print(T)
This code computes the sum of digits from right to left, using an array 'remainder' to store partial remainders modulo
```

1000000007. It then calculates T as the final remainder times 10 plus the first remainder from the array.

Learning to Reason without External Rewards

# Planning + Code Block + Explanation (Step 15)

To solve this problem, we can iterate through all the possible subsets of the set and calculate their sums. However, a more efficient approach is to use combinatorics. For N elements, there are  $2^N$  subsets. Each element appears in exactly half of the subsets, so the total sum of all subset elements is  $2^{(N-1)}$  times the sum of all elements in the original set.

```
''' python
def sum_of_powerset(n):
 return 2**(n-1) * n * (n+1) // 2
Read input
t = int(input())
for _ in range(t):
 n = int(input())
 print(sum_of_powerset(n))
'''
```

This program defines a function  $sum_of_powerset(n)$  that calculates the sum of the power set of a set with n elements. It then reads the number of test cases t, and for each test case, reads n and prints the result.

## Step-by-Step Planning + Code Block + Step-by-Step Explanation (Step 50)

To solve this problem, we need to read the input, parse the map and the regions to survey, then count the number of each type of terrain for each survey region. Here's a step-by-step approach to write the Python code for this problem:

- 1. Read the input: Get the dimensions of the map and the number of regions to survey.
- 2. Read the map: Parse the map information into a 2D list of strings.
- 3. **Read the regions**: Parse the survey regions into a list of tuples.
- 4. **Count the terrain**: For each survey region, count the number of 'J', 'O', and 'I' characters in the corresponding section of the map.
- 5. Output the results: Print the counts for each survey region.

Here's the Python code that implements this solution:

```
", "python
def count_terrain (map_info, survey):
 m, n = map_info.get('m'), map_info.get('n')
 j_{count} = 0
 o_{-}count = 0
 i_{-}count = 0
 for row in survey:
 a, b, c, d = row
 if a > m or b > n or c > m or d > n:
 continue
 for i in range(a, c + 1):
 for j in range (b, d + 1):
 if map_info[(i, j)] == 'J':
 j_{count} += 1
 elif map_info[(i, j)] == 'O':
 o_count += 1
```

```
elif map_info[(i, j)] == 'I':
 i_{-}count += 1
 return j_count, o_count, i_count
def process_input():
 \dot{M}, N = map(int, input().split())
 K = int(input())
 map_info = \{\}
 survey = []
 for _ in range (M):
 map_info[(M - _, _)] = input().strip()
 for _ in range(K):
 survey.append(list(map(int, input().split())))
 return M, N, K, map_info, survey
def main():
 M, N, K, map_info, survey = process_input()
 results = []
 for region in survey:
 j_count, o_count, i_count=count_terrain (map_info, region)
 results.append((j_count, o_count, i_count))
 for region in results:
 print(region[0], region[1], region[2])
if _____ == "____main___":
 main()
. . .
```

# Explanation:

- 1. **count\_terrain Function**: This function takes the map information and a survey region, and counts the number of 'J', 'O', and 'I' in that region.
- 2. **process\_input Function**: This function reads the input from the user, processes it to extract the map information and the survey regions.
- 3. **main Function**: This function calls 'process\_input' to get the necessary data, then calls 'count\_terrain' for each survey region and prints the results.

This code should correctly solve the problem and provide the desired output.

# B.6. Evolution of Qwen2.5-3B's Responses on LiveCodeBench Code Generation trained with INTUITOR on MATH

<sup>1259</sup> During training Qwen2.5-3B with INTUITOR on the MATH dataset, its code-generation capability steadily improves. By the <sup>1260</sup> middle of training, the model learns to produce syntactically valid Python; as training progresses, it refines formatting and <sup>1261</sup> clarity. By the end, it arrives at correct solutions.

1262 1263 1264

1210

1211

1212 1213

1214 1215

1216

1217

1218

1219

1220 1221

1222

1223 1224

1225

1226

1228 1229

1230

1231

1232 1233

1234

1235

1236 1237

1238

1239 1240

1241

1242

1243 1244

1245

1246

1247 1248

1249

1250 1251

1252

1253

Problem

**Question:** 

Alex is participating in the filming of another video of BrMeast, and BrMeast asked Alex to prepare 250 thousand tons of TNT, but Alex didn't hear him well, so he prepared n boxes and arranged them in a row waiting for trucks. The *i*-th box from the left weighs  $a_i$  tons. All trucks that Alex is going to use hold the same number of boxes, denoted by k. Loading happens the following way: • The first k boxes goes to the first truck, • The second k boxes goes to the second truck, • . . . • The last k boxes goes to the  $\frac{n}{k}$ -th truck. Upon loading is completed, each truck must have exactly k boxes. In other words, if at some point it is not possible to load exactly k boxes into the truck, then the loading option with that k is not possible. Alex hates justice, so he wants the maximum absolute difference between the total weights of two trucks to be as great as possible. If there is only one truck, this value is 0. Alex has quite a lot of connections, so for every  $1 \le k \le n$ , he can find a company such that each of its trucks can hold exactly k boxes. Print the maximum absolute difference between the total weights of any two trucks. Input The first line contains one integer t ( $1 \le t \le 10^4$ ) — the number of test cases. The first line of each test case contains one integer  $n (1 \le n \le 150\,000)$  — the number of boxes. The second line contains n integers  $a_1, a_2, \ldots, a_n$   $(1 \le a_i \le 10^9)$  — the weights of the boxes. It is guaranteed that the sum of n for all test cases does not exceed  $150\,000$ . Output For each test case, print a single integer — the answer to the problem. Sample Input 1: 1 2 10 2 3 6 1 3 60978 82265 78961 56708 39846 31071 4913 4769 29092 91348 64119 72421 98405 222 14294 19957 69913 37531 96991 57838 21008 14207 19198 Sample Output 1: Note

- In the first case, we should pick two trucks, so the first one will have only the first box, and the second one will have only the second box.
- In the second case, we should pick six trucks, so the maximum will be 10, the minimum will be 1, and the answer is 10 1 = 9.
- In the third case, for any possible k, the trucks will have the same total weight of boxes, so the answer is 0.

Read the inputs from stdin solve the problem and write the answer to stdout (do not directly test on the sample inputs). Enclose your code within delimiters as follows. Ensure that when the python program runs, it reads the inputs, runs the algorithm and writes output to STDOUT.

'''python

# YOUR CODE HERE

#### **Step 0 (Base): Incorrect**

'`' python
# YOUR CODE HERE

#### Step 10: Incorrect

'`' python
# YOUR CODE HERE
'''

#### Step 20: Incorrect

'`' python
# YOUR CODE HERE

. . .

# Step 30: Incorrect

To solve the problem of finding the maximum absolute difference between the total weights of any two trucks, we can follow these steps:

- 1. **Read the input**: The input consists of multiple test cases. Each test case starts with an integer n (the number of boxes), followed by n integers  $a_1, a_2, \ldots, a_n$  (the weights of the boxes).
- 2. Calculate the total weight of all boxes: Sum up all the weights of the boxes to get the total weight W.
- 3. Determine the number of trucks: The number of trucks m is given by  $m = \lfloor \frac{n}{k} \rfloor$ , where k is the number of boxes each truck can hold. Since k can vary from 1 to n, we need to find the maximum possible value of k that allows us to distribute the boxes evenly.
- 4. Calculate the maximum and minimum total weights of the trucks: For each possible value of k, calculate the total weight of the boxes that can be loaded into each truck. The maximum total weight of the trucks will be the sum of the weights of the k heaviest boxes, and the minimum total weight will be the sum of the weights of the k lightest boxes.

- 5. Find the maximum absolute difference: The maximum absolute difference between the total weights of any two trucks is the maximum of the absolute differences between the maximum and minimum total weights of the trucks for all possible values of k.
- 6. Find the maximum absolute difference: The maximum absolute difference between the total weights of any two trucks is the maximum of the absolute differences between the maximum and minimum total weights of the trucks for all possible values of k.

Here is the Python code that implements the above steps:

```
```python
import sys
from heapq import heappush, heappop
def max_difference(n, weights):
    # Calculate the total weight of all boxes
    total_weight = sum(weights)
    # Initialize variables to store the maximum
    # and minimum total weights of the trucks
    max_weight = 0
    min_weight = total_weight
    # Iterate over all possible values of k
    for k in range (1, n + 1):
        # Calculate the total weight of the k heaviest boxes
        max_weight = max(max_weight)
            sum(heappop(weights) for _ in range(k)))
        # Calculate the total weight of the k lightest boxes
        \min_{-}weight = \min(\min_{-}weight)
            sum(heappop(weights) for _ in range(k)))
    # Return the maximum absolute difference
    return max_weight - min_weight
# Read the number of test cases
t = int(input())
# Process each test case
for _ in range(t):
    n = int(input())
    weights = list(map(int, input().split()))
    # Calculate the maximum absolute difference
    result = max_difference(n, weights)
    # Write the result to stdout
    print(result)
. . .
```

Step 40: Correct

To solve the problem of finding the maximum absolute difference between the total weights of any two trucks, we can follow these steps:

1. **Read the input**: The input consists of multiple test cases. Each test case starts with an integer *n* (the number of boxes), followed by *n* integers representing the weights of the boxes.

1429

- 2. Calculate the total weight of all boxes: Sum up the weights of all the boxes to get the total weight.
- 3. Determine the number of trucks: Since each truck holds k boxes, the number of trucks is $\frac{n}{k}$.
- 4. Calculate the possible total weights of the trucks: For each possible value of k (from 1 to n), calculate the total weight of the boxes that can be loaded into the trucks. This involves summing up the weights of k consecutive boxes.
- 5. Find the maximum and minimum total weights: For each possible value of k, find the maximum and minimum total weights of the trucks. The maximum absolute difference is the difference between the maximum and minimum total weights.
- 6. Output the maximum absolute difference: Print the maximum absolute difference for each test case.

Here is the Python code that implements the above steps:

```
''' python
def find_max_difference(n, weights):
    total_weight = sum(weights)
    max_difference = 0
    for k in range(1, n + 1):
        if n \% k != 0:
            continue
        truck_weights = [sum(weights[i * k :
            (i + 1) * k]) for i in range(n // k)]
        max_weight = max(truck_weights)
        min_weight = min(truck_weights)
        max_difference = max(max_difference,
            abs(max_weight - min_weight))
    return max_difference
def main():
    t = int(input())
    for _ in range(t):
        n = int(input())
        weights = list(map(int, input().split()))
        result = find_max_difference(n, weights)
        print(result)
if __name__ == "__main__":
   main()
. . .
```

Step 58: Correct

To solve the problem of finding the maximum absolute difference between the total weights of any two trucks, we can follow these steps:

- 1. **Read the input**: The input consists of multiple test cases. Each test case starts with an integer *n* (the number of boxes), followed by *n* integers representing the weights of the boxes.
- 2. Calculate the total weight of all boxes: Sum up the weights of all the boxes to get the total weight.
- 3. Determine the number of trucks: Since each truck holds k boxes, the number of trucks is $\frac{n}{k}$.

- 4. Calculate the possible total weights of the trucks: For each possible value of k (from 1 to n), calculate the total weight of the boxes that can be loaded into the trucks. This involves summing up the weights of k consecutive boxes.
- 5. Find the maximum and minimum total weights: For each possible value of k, find the maximum and minimum total weights of the trucks. The maximum absolute difference is the difference between the maximum and minimum total weights.
- 6. Output the maximum absolute difference: Print the maximum absolute difference for each test case.

Here is the Python code that implements the above steps:

```
", "python
def max_difference(n, weights):
    total_weight = sum(weights)
    \max_{-}diff = 0
    for k in range (1, n + 1):
        if n % k != 0:
            continue
        truck_weights = [sum(weights[i * k :
            (i + 1) * k]) for i in range(n // k)]
        max_diff = max(max_diff)
            max(truck_weights) - min(truck_weights))
    return max_diff
def solve():
    t = int(input())
    for _ in range(t):
        n = int(input())
        weights = list(map(int, input().split()))
        result = max_difference(n, weights)
        print(result)
if __name__ == "__main__":
    solve()
. . .
```