Surrogate Signals from Format and Length: Reinforcement Learning for Solving Mathematical Problems without Ground Truth Answers

Anonymous ACL submission

Abstract

Large Language Models have achieved remarkable success in natural language processing tasks, with Reinforcement Learning playing a key role in adapting them to specific applications. However, obtaining ground truth answers for training LLMs in mathematical problemsolving is often challenging, costly, and sometimes unfeasible. This research delves into the utilization of format and length as surrogate signals to train LLMs for mathematical problem-solving, bypassing the need for traditional ground truth answers. Our study shows that a reward function centered on format correctness alone can yield performance improvements comparable to the standard GRPO algorithm in this phase. Recognizing the limitations of format-only rewards, we incorporate length-based rewards. The resulting GRPO approach, leveraging format-length surrogate signals, not only matches but surpasses the performance of the standard GRPO algorithm relying on ground truth answers in certain scenarios, achieving 40.0% accuracy on AIME2024 with a 7B base model. Through systematic exploration and experimentation, this research offers a practical solution for training LLMs to solve mathematical problems and reducing the dependence on extensive ground truth data collection.

1 Introduction

002

004

005

011

012

017

019

040

042

043

In the dynamic landscape of artificial intelligence, Large Language Models (LLMs) (Brown et al., 2020; Chowdhery et al., 2023; Yang et al., 2023; Wang et al., 2025a; Grattafiori et al., 2024) have emerged as a transformative force, with models like GPT-o1 (Jaech et al., 2024), DeepSeek-R1 (DeepSeek-AI et al., 2025), and Qwen3 (Yang et al., 2025) leading the charge. Pre-trained on massive text corpora, these models have demonstrated remarkable proficiency in diverse natural language processing tasks, ranging from text generation and question-answering to language translation and code writing. Their success largely stems from unsupervised pre-training, which empowers LLMs to capture complex semantic and syntactic patterns, enabling effective generalization across various scenarios. 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

Reinforcement Learning (RL) plays a crucial role in adapting pre-trained LLMs to specific downstream tasks. Among RL techniques, Proximal Policy Optimization (PPO) (Schulman et al., 2017) and its advanced variant, Group Relative Policy Optimization (GRPO) (Shao et al., 2024), are commonly employed to optimize LLMs. These methods typically rely on ground truth answers to define rewards, serving as explicit feedback for the model to iteratively refine its solutions. However, obtaining accurate ground truth answers, particularly in the domain of mathematical problemsolving, presents significant challenges. It often demands substantial human effort, time, and resources, and in certain cases, such answers may be scarce or even nonexistent. This limitation has spurred the exploration of alternative training strategies to enable effective RL without relying on explicit ground truth information.

Motivated by these challenges, our research endeavors to explore the possibility of training LLMs for mathematical problem-solving using alternative signals instead of ground truth answers. Through systematic experiments and in-depth analysis, we have made several significant discoveries. During the initial 15 steps of RL training, the model predominantly focuses on learning the format of mathematical solutions. This early phase is crucial, contributing approximately 85% of the overall performance improvement during the entire RL training process. Notably, during this period, we observed a significant reduction in response length, indicating that the model rapidly eliminates redundant information and converges towards a more structured and efficient representation. Experiments revealed that a reward function solely considering format correctness achieved the same performance gains

100

102

103

104

107

108

110

111

112

113

114

as the standard GRPO algorithm, underscoring the potency of format correctness as a key signal in the early learning stage.

Nevertheless, relying solely on format-based rewards has limitations, causing performance improvement to stall after the initial gains. To overcome this hurdle, we integrated a length-based reward into the format-based reward function. Strikingly, our GRPO approach leveraging formatlength surrogate signals not only matched but in some cases outperformed the standard GRPO algorithm that relies on ground truth answer information. This is because format and length together act as powerful "surrogate signals" highly correlated with answer correctness. Format correctness provides a necessary optimization target, while the length-based reward penalizes overly long or short responses, prompting the model to refine its content by eliminating incorrect or redundant derivations.

Through these findings, our work effectively challenges the necessity of ground truth answers for RL in mathematical problem solving, provides a detailed analysis of GRPO training dynamics to reveal the importance of the early format-learning phase and complementary role of length-based rewards, and opens up new possibilities for training LLMs in scenarios where ground truth answers are scarce or unavailable, offering an efficient approach applicable to mathematical reasoning tasks.

2 Related Work

RL has been proven effective in enhancing LLM 115 performance. PPO (Schulman et al., 2017) and 116 GRPO (Shao et al., 2024) are widely used in RL 117 frameworks for LLMs, with detailed introductions 118 provided in Appendix A. Recent research uses 119 scaled-up RL training to enable LLMs to explore reasoning paths for complex problems. For exam-121 ple, DeepSeek-AI et al. (2025) achieved excellent 122 results in math and coding tasks through large-scale 123 RL on an unsupervised base model, without relying 124 on pre-trained reward models or MCTS. Team et al. 125 (2025) enhances general reasoning via RL, focus-126 ing on multimodal reasoning and controlling think-127 ing length. Format reward in RL. DeepSeek-AI 128 et al. (2025) uses format rewards to structure model 130 outputs. Liu et al. (2025a) noted format rewards dominate early training. Our study isolates the in-131 fluence of answer rewards and designs a format for 132 math reasoning tasks. Experiments show using our 133 format in early RL training matches performance 134

of answer reward training. Length Control in RL. DeepSeek-AI et al. (2025) found response length and evaluation metrics increase with RL training steps until an "Aha moment". Other studies explore length reward functions' impacts. Yeo et al. (2025) observed response lengths decline due to model size and KL divergence penalties. Chen et al. (2025) argued direct length extension training may harm performance. In contrast, our length reward penalizes overly long responses, guiding concise outputs. Experiments show combining length and format rewards outperforms answer rewards.

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

3 Method: Format and Length as Surrogate Signals for Answer

To mitigate the issue of label scarcity in real-world environments, we explore the potential of format and length as powerful "surrogate signals" highly correlated with answer correctness. Format correctness in mathematical problem-solving offers a necessary but insufficient condition for answer accuracy, providing a clear structural optimization target for the model. Meanwhile, the length of the response serves as an indicator of content efficiency and logical compactness, reflecting the quality of the solution's reasoning process. Based on these insights, we develop a novel learning framework that integrates format and length rewards into the GRPO algorithm. This framework, centered around optimizing LLMs without relying on explicit ground truth answers, aims to enable effective training by leveraging these surrogate signals to approximate the optimization direction of ground truth answer rewards.

3.1 Format Reward

In the context of mathematical problem-solving, a correct format is crucial for ensuring the clarity and comprehensibility of the solution. Our format reward mechanism is designed to encourage the model to generate responses that adhere to the standard presentation conventions of mathematical solutions (details in Appendix D). The format reward $R_{\rm f}$ is defined as a binary function:

$$R_{\rm f} = \begin{cases} 1 & \text{if the format is right.} \\ 0 & \text{else.} \end{cases}$$
(1)

This reward serves as a fundamental signal for the178model to learn the structural aspects of mathemati-
cal problem-solving in the early stages of training.179180

182 183

185

186

- 189
- 190

193

195

196

198 199

204

210

211

212 213

214

215

217

218 219

221

222

223

3.2 Length Reward

To complement the format reward and further refine the content of the model's responses, we introduce a length reward function. The length of a response is a critical factor that reflects the efficiency and logical compactness of the solution. An overly short response may lack essential reasoning steps, while an excessively long response might contain redundant or incorrect derivations.

Our length reward function is designed to strike a balance between promoting comprehensive reasoning and preventing overly long responses that could exceed the model's context limits. It is formulated as a piecewise function:

$$R_{1} = \begin{cases} 1 - \left(1 - \frac{x}{p}\right)^{2}, & 0 \le x \le p, \\ 1 - 2\left(\frac{x - p}{1 - p}\right)^{2}, & p < x \le 1, \end{cases}$$
(2)

Let

$$x = \frac{L}{L_{\text{max}}},\tag{3}$$

where L is the length of the current response and L_{max} is the maximum context length. Let $p \in (0, 1)$ be a tunable parameter that controls the turning point of the piecewise function, with a default value of 0.5. This piecewise function is continuous and differentiable at x = p, encouraging response lengths that approach the turning point p. The reward increases smoothly as x grows from 0 to p, reaches a maximum at x = p, and then decreases for x > p, thereby penalizing overly long responses.

A positive length reward can only be obtained when the format is right. Examples with format errors are considered severe-no matter how ideal their length may be, they can receive at most 0. Therefore, the final format-length reward can be expressed as:

$$R_{\rm fl} = \begin{cases} R_{\rm f} + R_{\rm l} & \text{if the format is right.} \\ \min(0, R_{\rm f} + R_{\rm l}) & \text{else.} \end{cases}$$
(4)

By combining the format reward and length reward, we provides an "surrogate signals" for the model's reinforcement learning, helping to alleviate the issue of label scarcity in real-world environments.

4 **Experiments**

In this section, we present a comprehensive set of experiments designed to demonstrate the practical viability of using format and length as surrogate signals for answer accuracy in GRPO for mathematical reasoning tasks.

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

264

265

266

268

269

4.1 Experimental Setup

Reward configurations: We designed a series of experiments with distinct reward configurations to assess the effectiveness of our proposed approach. Correctness: This configuration is served as our baseline, which uses the exact match with groundtruth answers as the reward criterion. When the model's output precisely aligns with the correct answer, it is assigned a reward score of 1; otherwise, it receives 0. We utilized the MARIO_EVAL 1 library to accurately extract answer content from the model's output, ensuring a reliable evaluation standard. Format-Only: The reward function is as shown in Eq.(1), which is determined solely by the format of the model's output. After normalizing the content, we employ SymPy², a powerful symbolic mathematics library, to validate its mathematical format. Format-Length: The reward function is as shown in Eq.(4), where the format reward is the same as that of Format-Only RL.

Datasets: We trained models on two mathematical reasoning datasets: DeepScaleR³ and MATHtrain. DeepScaleR (17,000 samples) integrates problems from the MATH (Hendrycks et al., 2021), AMC (< 2023), AIME (1984-2023), and others, with deduplication and decontamination applied. MATH-train (7,500 samples) is the MATH dataset's training split.

Evaluation: We evaluated the model on three datasets: MATH500, AIME2024, and AMC2023 with greedy decoding. In addition to analyzing each dataset individually, we also calculated the average scores across all benchmarks to enable direct comparison.

Implementation details: We trained the Qwen2.5-Math-7B base model using the GRPO algorithm under the verl⁴ framework. For each case in training and evaluation, we used Qwen-Math template (as shown in Appendix C). During training, we used the following hyperparameters: a learning rate of 1e-6, a batch size of 128, a temperature of 0.6, 8 responses per prompt, a maximum response length of 3072, and a KL coefficient of 0.001. All training was performed on a machine with $8 \times H20$

²https://github.com/sympy/sympy

⁴https://github.com/volcengine/verl

¹https://github.com/MARIO-Math-Reasoning/ MARIO_EVAL

³https://github.com/agentica-project/rllm

271

274

275

276

279

281

282

283

286

287

291

292

294

297

301

GPUs, and a single training run took 6 hours.

4.2 Impact of Format Reward

The format-only experiment offers critical insights into the role of format correctness in the training process. During the initial 15 steps, as depicted in Figure 1, the performance of the model trained with format-only reward remarkably aligns with that of the correctness reward setup on both benchmarks. This convergence validates our hypothesis that in the early stages of GRPO, the model predominantly focuses on learning the structural patterns of mathematical solutions. It suggests that format serves as a strong initial signal, allowing the model to quickly grasp the essential presentation conventions of mathematical answers, which accounts for approximately 85% of the overall performance improvement in this early phase.

However, as the training progresses beyond the 15-step mark, a significant divergence emerges. The performance of the format-only model plateaus, barely showing any improvement even after 100 training steps. This stagnation can be attributed to the inherent limitation of relying solely on format as a reward signal. While format correctness is a necessary condition for answer accuracy, it is not sufficient. Without additional guidance, the model lacks the means to refine the content within the correct format, leading to an inability to further enhance the accuracy of its solutions. This highlights the need for supplementary signals to drive continuous improvement.

4.3 Effectiveness of Format-Length RL

Our format-length reward demonstrates notable advantages in mathematical problem-solving without ground truth answers, as shown in Table 1. By using format consistency and response length as surrogate signals, the approach achieves competitive performance against the model trained with correctness reward.

Numerically, model trained with format-length reward achieves an average score of 56.8, surpassing the correctness reward's average score of 53.0 311 when using the DeepScaleR training dataset. In par-312 ticular, model trained with format-length reward 313 achieved 40 points in AIME2024 using the MATH 314 315 training dataset. This indicates that leveraging structural and length-based rewards alone can guide the model to generate high-quality solutions compara-317 ble to or better than models trained with correctness reward, even without explicit answer supervision. 319

Figure 1 shows the average accuracy curves of GRPO training with different rewards. In Appendix Figures S1 and S2, we present the accuracy curves of each benchmark respectively. It can be seen from these figures that model trained with format-length reward maintains stable performance comparable to the correctness reward baseline throughout the entire training process. The consistent curves validate the reliability of surrogate signals in driving model improvement without ground truth, highlighting the approach's scalability and data efficiency for mathematical reasoning tasks.

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

345

346

347

348

349

350

351

352

353

354

355

357

358

359

360

361

362

363

364

365

366

367

369

370

4.4 Response Length Dynamics

In Figures 2, we respectively show the curves of average response length during GRPO training with different rewards on the DeepscaleR dataset. The model trained with format-length reward demonstrated a distinctive dual-phase evolution in response length, which starkly contrasts with the monotonic decrease observed in the models trained with correctness reward and format-only reward.

Across all reward configurations, the average response length decreases during the initial 30 training steps. This indicates that the model prioritizes format adherence during this phase. Driven by the dominant format reward signal, which penalizes any deviation from the required answer schema, it prunes redundant content to meet structural constraints.

As training advances from 30 to 100 steps, the length reward mechanism takes the lead, driving a strategic expansion of response content. Unlike simplistic length penalties that encourage brevity at the cost of depth, GRPO with format-length reward fosters an optimal equilibrium. It encourages longer thinking processes and discourages unnecessary verbosity. This dynamic mirrors the human problem-solving process, where initial efforts focus on establishing structure, followed by iterative refinement of content. During the final stages, the model's response length increases by an average of 14.0%, which correlates with a 10.5% improvement in average accuracy training on DeepScaleR, indicating that length serves as a proxy for reasoning complexity rather than redundancy.

This dual-phase evolution parallels the human learning process encapsulated by the adage "Reading thin before reading thick." In the first stage, the model, similar to human summarization, compresses a single reasoning process, while in the second stage, it expands and generalizes, exploring



Figure 1: Average accuracy on evaluation benchmark training on (a) DeepScaleR and (b) Math-train.

Method	Label Free	AIME2024	MATH500	AMC2023	AVG.
Qwen-Math-7B	_	16.7	50.8	42.2	36.6
DeepSeek-R1-Distill-7B@3k	×	10.0*	60.1*	26.2*	32.1*
DeepSeek-R1-Distill-7B@8k	×	33.3*	88.1*	68.4*	63.3*
Qwen2.5-Math-7B-Instruct	×	16.7	83.2	55.4	51.8
LIMR-7B Li et al. (2025)	×	23.3 (32.5*)	74.8 (78.0*)	60.2 (63.8*)	52.8 (58.1*)
SimpleRL-Zero-7B Zeng et al. (2025)	×	26.7 (40.0*)	75.4 (80.2*)	57.8 (70.0*)	53.3 (63.4*)
Oat-Zero-7B Liu et al. (2025b)	×	40.0 (43.3*)	78.2 (80.0*)	61.5 (62.7*)	60.0 (62.0*)
Correctness (baseline)	X	26.7 / 26.7	74.6 / 73.0	57.8 / 56.6	53.0 / 52.1
Format-Only	J	26.7 / 26.7	72.6 / 72.8	55.4 / 53.0	51.6 / 50.8
Format-Length	J	33.3 / 40.0	76.8 / 73.0	60.2 / 54.2	56.8 / 55.7

Table 1: Accuracy comparison of different models on benchmark datasets (cyan rows denote our trained models). Results are separated by a slash for DeepscaleR and MATH-train datasets (DeepscaleR first, MATH-train second). Results without * are evaluated in our environment (details in Appendix B); * indicates results from Liu et al. (2025b) or the original paper.

more diverse and complex reasoning paths, such as error correction and branch exploration. In contrast, the correctness reward baseline and format-only models, as highlighted by the red box in Figure 2, briefly attempt to explore complex reasoning but ultimately revert to the "comfort zone" of compressing a single reasoning process.



Figure 2: Response length during training. The solid lines in the figure represent the original results, while the dashed lines represent the results smoothed with a window size of 5.

4.5 Format-Length Rewards' Impact Across Difficulty Levels

378

379

380

383

384

386

389

390

391

392

393

394

395

396

398

To explore how format-length rewards affect LLMs' mathematical problem-solving, we analyzed the MATH500 dataset, which has official difficulty ratings and balanced problem distribution. As depicted in Figure 3a, by the end of the training process, the format-length model outperformed the correctness reward baseline across all difficulty levels.

The relationship between response length and reasoning performance further illuminates the mechanism behind these results. As shown in Figure 3b, both models generated longer responses for higher-difficulty problems. The correctness reward baseline model initially showed a rapid decrease in output length, which later stabilized, while the format-length model demonstrated a midstage increase, especially for high-difficulty problems. This increase in length was positively correlated with improved accuracy, indicating that the

374

375

399	length reward encourages the model to adopt more
400	comprehensive reasoning strategies, particularly
401	when tackling complex tasks.

We delved deeper into the model's reason-402 ing process by analyzing the frequency of 403 reflective words in the generated responses 404 (Figure 3c). Reflective words, including those 405 related to verification (wait/verify/check), 406 retrospection (recall/recheck), branch ex-407 ploration (alternatively), logical turn or 408 contrast (however/but/since), and problem 409 and step-by-step reasoning 410 decomposition (step/step-by-step), represent complex reason-411 ing behaviors. The correctness reward baseline 412 model showed an initial increase in reflective 413 words, which plateaued in the later stages, aligning 414 with its limited performance gains. In contrast, the 415 format-length model exhibited a significant rise 416 in reflective words, especially for high-difficulty 417 problems. This indicates that the length signal 418 helps increase the depth of thinking, enabling 419 the model to engage more in complex reasoning 420 behaviors such as verification, retrospection, and 421 problem decomposition. Such enhanced reflective 422 423 thinking allows the model to better explore different solution paths and logical turns, thereby 424 improving its ability to handle high-difficulty 425 problems. 426

> To further validate these findings, we conducted a case study by comparing the outputs of the correctness model and format-length model on challenging questions (Appendix Table S1). The format-length model had learned a "step-by-step problem-solving and verification" pattern, which confirmed the effectiveness of our format-length reward mechanism in balancing response length, reasoning depth, and content quality.

427

428

429

430

431

432

433

434

435

Similar to (Wang et al., 2025b), we observed 436 that increasing the frequency of reflective language 437 does not necessarily correlate with better model 438 performance. Specifically, models can exhibit over-439 reflection, characterized by repeatedly switching 440 441 reasoning paths on complex problems, often leading to failed solutions. This over-reflection is some-442 times accompanied by phrase repetition (Appendix 443 Table S2), where models may exploit length re-444 wards through redundancy. 445

Method	AIME2024	MATH500	AMC2023
Qwen-Math-7B	63.3	94.0	92.8
Correctness	73.3	94.4	90.4
Format-Only	66.7	94.0	91.6
Format-Length	66.7	94.4	92.8

Table 2: Pass@64 results across different methods.

5 Discussion

5.1 Rethinking Ground Truth Dependency in Mathematical Reasoning

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

The remarkable performance of our ground truthfree RL approach begs the question: how can RL without explicit answer supervision match the effectiveness of traditional ground truth-based methods? The answer lies in the latent knowledge already encoded within pre-trained language models. Prior to RL fine-tuning, these models have assimilated vast amounts of knowledge from diverse corpora, enabling them to potentially generate correct answers—RL merely serves as a catalyst to activate this dormant capacity.

Our pass@N experiments provide compelling evidence for this mechanism. By generating N distinct responses per prompt and assessing the presence of correct answers among them, we observe comparable pass@N scores across four conditions: the pre-trained model, the model fine-tuned by GRPO with correctness, format-only, and formatlength rewards. As presented in Table 2, which showcases the pass@64 results, we can see that the performance differences between thes methods are relatively minor. This parity indicates that all RL variants fail to confer new knowledge; instead, they optimize how the model retrieves and structures existing knowledge.

In essence, our findings demonstrate that with the right reward design—such as leveraging format and length cues—RL can effectively stimulate the model's internal reasoning processes. As long as the training mechanism activates the model's latent cognitive abilities, explicit ground truth answers become an optional component rather than an essential requirement for high-performance RL in mathematical reasoning tasks.

5.2 Format Learning in RL and SFT

Since both traditional RL with ground truth rewards and our format-based RL mainly learn answer formatting in the first 15 training steps, a key question arises: how does format learning through RL com-



Figure 3: The curves of (a) accuracy, (b) response length, and (c) reflective keyword frequency for cases of different difficulty levels in MATH500 during training.

pare with supervised fine-tuning (SFT)? To answer this question, we carried out a series of comparative experiments, comparing three different training methods: 1) GRPO training with format-based rewards, 2) offline SFT using ground truth chain-ofthought (CoT) examples, and 3) online SFT. Online SFT serves as a middle ground between offline SFT and RL, connecting static supervised learning and the dynamic, feedback-driven RL, which helps us figure out how different training methods affect format learning.

488

489

490

491

492

493

494

495

496

497

498

499

503

504

505

506

510

511

512

513

514

515

517

521

523

We used Qwen2.5-Math-7B as the original model, which we didn't train, to provide a baseline for comparison. The GRPO(Correctness) was used as a reference to measure the performance of other methods. All experiments were conducted under the setting of sampling from the MATH dataset with a temperature of 0.6.

In the GRPO training with format-based rewards and online SFT experiments, we adopted an online sampling strategy. During training, we constantly sampled model outputs and applied GRPO or SFT based on whether the format was correct. Specifically, online SFT only used format-correct samples to update parameters. All experiments used a batch size of 128 and ran for 100 training steps.

As shown in Table 3, the results offer important insights. Under the temperature=0.6 setting, the GRPO training with format-based rewards and online SFT performed very similarly, achieving comparable format accuracy rates and scores on 518 the MATH500 benchmark. On the other hand, the offline SFT method didn't perform as well, showing lower format accuracy and lower MATH500 scores. These results emphasize the important role of online sampling in making RL more effective for format learning. RL and online SFT can adjust to the quality of real-time outputs, which allows them 525

Method	Answer Acc	Format Acc
Qwen2.5-Math-7B	61.7	87.3
GRPO(Correctness)	74.0	95.0
GRPO(Format-Only)	70.1	96.3
offline SFT	51.3	88.7
online SFT	71.3	95.0

Table 3: Comparison of format accuracy and answer accuracy across different training methods on the MATH500 benchmark.

to optimize answer formatting more efficiently than the static offline SFT. Clearly, the iterative and feedback-driven nature of online training is crucial for quickly improving language models' ability to learn formats.

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

5.3 **Mitigating Repetition and Reward** Hacking

A potential concern with length-based rewards is the risk of reward hacking, where the model generates repetitive content to increase its length. To address this, we employed the longest repeated substring analysis to measure repetition. The longest repeated substring ratio (Figure 4) provides a normalized perspective on repetition. At the start of training, both the format-length and correctness models exhibited high levels of repetition, mainly due to incorrect formatting issues, such as stacked instances of '\\boxed'. However, this problem was resolved after just 15 training steps. The repetition rate then dropped significantly and remained stable throughout the subsequent training process. These findings demonstrate that the format-length reward mechanism effectively balances response length, reasoning depth, and content quality. By integrating format and length signals, our approach not only improves performance on mathematical reasoning tasks but also mitigates the risks associated with traditional length-based rewards, like



Figure 4: Longest duplicate substring ratio of MATH500 evaluation benchmark during training.

repetition and reward hacking.

554

555

557

558

559

561

563

565

569

570

571

573

575

582

586

5.4 Design Principles of Format-Length RL

In the context of language model training, truncation refers to the situation where the generated output exceeds the maximum allowable length (e.g., the context window size of the model) and has to be cut off. Truncation is highly undesirable for several reasons. Firstly, it leads to incomplete responses, which can result in the loss of crucial information and logical steps necessary for correct mathematical reasoning. In the case of mathematical problem-solving, a truncated answer may omit key derivations or final conclusions, rendering the solution incorrect or meaningless. Secondly, truncation can disrupt the coherence and flow of the reasoning process, making it difficult for the model to build on its own arguments and reach a valid conclusion. Prior studies have explored length-based rewards, but their applicability to label-free settings is limited. For example, Yeo et al. (2025) proposed a cosine-shaped length reward coupled with correctness, while Chen et al. (2025) introduced a linear length reward: $R = L/L_{Max} + R_{correctness}$. We reproduced this linear reward and the result is in Appendix Figure S3. However, it led to a rapid surge in response length, exceeding the model's context window and causing a 52.9% truncation rate by step 54. This high truncation rate severely degraded performance, as the truncated outputs were often incomplete and lacked the necessary logical structure for accurate mathematical reasoning. This outcome underscores the importance of carefully designing length rewards to balance exploration and efficiency, ensuring that the model generates responses of optimal length without incurring excessive truncation. In contrast, our Format-Length approach maintains a low truncation rate while achieving superior accuracy. By incorporating a length reward that penalizes excessive length before reaching the context limit, our method effectively guides the model to generate concise yet comprehensive responses. This not only prevents reward hacking, where the model might generate overly long or repetitive content to maximize rewards, but also promotes high-quality reasoning, as the model is encouraged to find the most efficient way to express correct mathematical solutions within the given length constraints. 587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

6 Conclusion

In this study, we found that format and length can serve as effective surrogate signals for training LLMs in mathematical problem-solving, eliminating the need for ground truth answers. During early RL training, LLMs focus on learning solution formats; a format-based reward function alone yields performance gains similar to standard GRPO. Integrating length-based rewards enables the GRPO approach with format-length signals to outperform traditional methods relying on ground truth in some cases. This finding challenges the notion that ground truth answers are essential for LLM training. The format-length signals offer a practical, efficient alternative, reducing data collection costs. Applicable across mathematical and logical tasks, this approach opens new avenues for LLM training. Future work will optimize signal utilization and expand application to enhance LLM training efficiency and generalization.

7 Limitations

There are aspects of our study that merit further exploration. The evaluation of format and length as surrogate signals was predominantly focused on mathematical problem-solving, leaving open the question of their effectiveness in other complex reasoning domains, such as scientific hypothesis testing or advanced programming challenges. Additionally, our experiments were conducted with specific LLM architectures and training configurations, and the performance of this approach may differ when applied to models with varying pretraining paradigms and scale.

References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Zhipeng Chen, Yingqian Min, Beichen Zhang, Jie Chen, Jinhao Jiang, Daixuan Cheng, Wayne Xin Zhao, Zheng Liu, Xu Miao, Yang Lu, Lei Fang, Zhongyuan Wang, and Ji-Rong Wen. 2025. An empirical study on eliciting and improving r1-like reasoning models. *Preprint*, arXiv:2503.04548.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, and 1 others. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The Ilama 3 herd of models. arXiv preprint arXiv:2407.21783.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. 2025. Limr: Less is more for rl scaling. *Preprint*, arXiv:2502.11886.
- Zichen Liu, Changyu Chen, Wenjun Li, Tianyu Pang, Chao Du, and Min Lin. 2025a. There may not be aha moment in r1-zero-like training — a pilot study. https://oatllm.notion.site/oat-zero. Notion Blog.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025b. Understanding r1-zero-like training: A critical perspective. *Preprint*, arXiv:2503.20783.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *Preprint*, arXiv:1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300.*
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, and 75 others. 2025. Kimi k1.5: Scaling reinforcement learning with llms. *Preprint*, arXiv:2501.12599.
- Bingning Wang, Haizhou Zhao, Huozhi Zhou, Liang Song, Mingyu Xu, Wei Cheng, Xiangrong Zeng, Yupeng Zhang, Yuqi Huo, Zecheng Wang, and 1 others. 2025a. Baichuan-m1: Pushing the medical capability of large language models. *arXiv preprint arXiv:2502.12671*.
- Yue Wang, Qiuzhi Liu, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Linfeng Song, Dian Yu, Juntao Li, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025b. Thoughts are all over the place: On the underthinking of o1-like llms. *Preprint*, arXiv:2501.18585.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, and 1 others. 2023. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. 2025. Demystifying long chain-of-thought reasoning in llms. *Preprint*, arXiv:2502.03373.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. 2025. Simplerlzoo: Investigating and taming zero reinforcement learning for open base models in the wild. *Preprint*, arXiv:2503.18892.

722

723

724

725

726

727

729

732

733

734

735

736

737

689

690

691

635

638

641

642

643

644

645

655

666

672

674

677

678

679

740

741

742

743

744

745

746

747

749

750

751

752

754

759

761

765

770

771

772

773

774

776

780

A Introducion of PPO and GRPO

A.1 Proximal Policy Optimization

PPO is a widely-used and highly effective algorithm in the field of RL. At its core, PPO aims to optimize the policy of an agent to maximize the expected cumulative reward over time. The algorithm is based on the policy gradient method, which updates the policy by computing the gradient of the expected reward with respect to the policy parameters. The key idea behind PPO is to balance the trade-off between exploration and exploitation during the policy update process. It does this by introducing a clipped surrogate objective function. Let π_{θ} be the policy parameterized by θ , and $\pi_{\theta_{old}}$ be the old policy. Given a set of trajectories collected from the environment, the objective of PPO is to maximize the following clipped objective function:

$$\mathbb{E}_{\pi_{\theta_{old}}}\left[\min\left(r_t(\theta)A_t, \operatorname{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)A_t\right)\right]$$
(5)

where

$$r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \tag{6}$$

is the probability ratio of the new policy π_{θ} to the old policy $\pi_{\theta_{old}}$ for taking action a_t in state s_t , A_t is the advantage function that estimates how much better an action is compared to the average action in state s_t , and ϵ is a hyperparameter that controls the clipping range. The clipping operation ensures that the policy update is not too drastic, preventing the policy from diverging significantly from the old policy in a single update step.

To compute the advantage function A_t , PPO typically relies on value function estimation combined with Generalized Advantage Estimation (GAE). The value function V(s) parameterized by ϕ , predicts the expected cumulative reward from state s. It is trained via temporal difference (TD) learning to minimize the squared error:

$$L^{\text{Value}}(\phi) = \mathbb{E}\left[\left(V_{\phi}(s_t) - \left(R_t + \gamma V_{\phi}(s_{t+1}) \right) \right)^2 \right],$$
(7)

where R_t is the reward given by a reward model or a reward function and γ is the discount factor. The advantage A_t is then calculated using GAE, which generalizes multi-step TD errors with a tunable parameter $\lambda \in [0, 1]$ to balance bias and variance:

$$A_t^{\text{GAE}(\gamma,\lambda)} = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l},$$

$$\delta_t = R_t + \gamma V(s_{t+1}) - V(s_t).$$
(8)

Here, $\lambda = 0$ reduces to single-step TD error, while781 $\lambda = 1$ recovers Monte Carlo advantage estimation.782By integrating GAE, PPO efficiently utilizes trajectory data while maintaining stable policy updates.783

A.2 Group Relative Policy Optimization

GRPO is an efficient reinforcement learning algorithm that improves upon PPO by eliminating the need for a separate value function. GRPO estimates advantages through group-relative normalization: for a given input query q, the behavior policy $\pi_{\theta_{old}}$ samples G responses $\{o_i\}_{i=1}^G$, then calculates each response's advantage as:

$$A_t^{\text{GRPO}}(o_i) = \frac{R(o_i) - \text{mean}(\{R(o_j)\}_{j=1}^G)}{\text{std}(\{R(o_j)\}_{j=1}^G)}, \quad (9)$$

$$793$$

where $R(o_i)$ is the reward of response o_i . 794

B Evaluation Details

We used vllm for inference with greedy decoding (temperature = 0) to ensure reproducibility. Since VLLM's batched inference produces different outputs for the same input under different batch sizes, we set the validation batch size to 128 and evaluate each dataset independently to ensure consistency in evaluation. Because we used the Qwen2.5-Math base models with a context length of 4k, we set the generation budget for all compared baselines to 3k.

C Template

Qwen-Math Template
<pre>< im_start >system Please reason step by step, and put your fi- nal answer within . < im_end > < im_start >user {question}</pre>
< im_end >
< im_start >assistant

805

796

797

798

799

800

801

802

803

804

785

786

787

789

790

791

Deepseek-R1 Template

A conversation between User and Assistant. The User asks a question, and the Assistant solves it. The Assistant first thinks about the reasoning process in the mind and then provides the User with the answer. The reasoning process is enclosed within <think> </think> and the answer is enclosed within <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>. User: {question}

Assistant:

D Detailed Form of Format Reward

DeepSeek-AI et al. (2025) introduced a format reward to assess whether the model's output aligns with the Deepseek-R1 template(Appendix C) format (i.e., writing the reasoning process within <think> <think> </think> tags and placing the answer within <answer></answer> tags). Our format reward is designed based on Qwen-Math template(in Appendix C), which evaluates whether the model's output contains \boxed{} and checks whether the content inside the \boxed{} is valid mathematical content. A specific example is shown below:

Right case:
\boxed{1}
<pre>\boxed{\frac{3}{2}}</pre>
$boxed{x^2 + 12y =1}$
<pre>\boxed{(0,\infty)}</pre>
Wrong case:
$boxed{x +* 2}$
\boxed{(0,1 }

E Performance Curves of Different Reward Trainings on Benchmarks

- See Figure S1 and S2.
 - F Comparison of Experimental Results between Linearlength Reward and Format-Length Reward

See Figure S3.

G Case Study

837 See Table S1 and S2.

812

813

815

816

817

818

819

820 821 822

829

830

832

833 834

807

Table S1: Model output on complex mathematical problems after training on format-length reward and correctness reward. The red boxes in the figure represent step-by-step problem-solving behavior, while the blue boxes indicate verification behavior. Clearly, both types of behavior occur far more frequently in the format-length model than in the correctness model. As for the results, the former produces correct answers, whereas the latter produces incorrect ones.

PROMPT: system\nPlease reason step by step, and put your final answer within \\boxed{}. \n user \n You have seven bags of gold coins. Each bag has the same number of gold coins. One day, you find a bag of 53 coins. You decide to redistribute the number of coins you have so that all eight bags you hold have the same number of coins. You successfully manage to redistribute all the coins, and you also note that you have more than 200 coins. What is the smallest number of coins you could have had before finding the bag of 53 coins?\nassistant\n

RESPONSE_FORMAT-LENGTH: To solve this problem, we need to determine the smallest number of coins in each of the original seven bags such that when we add the 53 coins from the new bag, the total number of coins can be evenly distributed among eight bags, and the total number of coins is more than 200.\n\nLet's break down the problem step step :\n\n 1.** Define Variables :**\n - Let x be the number of coins in each of the original seven bags. h - The total number of coins in the original seven bags is 7x.\n - After finding the new bag with 53 coins, the total number of coins becomes 7x + 53.\n\n**2**.**Redistribution Condition:**\n - We need to redistribute these coins so that each of the eight bags has the same number of coins.\n - Therefore, 7x + 53must be divisible by $8.\ln 3.**$ Mathematical Formulation:**n - We need to find the smallest x such that $7x + 53 \setminus equiv_0 \setminus pmod_8.$ h - This can be rewritten as $7x \setminus equiv_{-53} \setminus pmod_8.$ - Since $-53 \setminus equiv3 \setminus pmod8$ (because -53 + 56 = 3 and 56 is a multiple of 8), we have $7x \ge 0.5$ total number of coins in the original seven bags is $7 \setminus times 5 = 35$. In - Adding the 53 coins from the new bag, the total number of coins is 35 + 53 = 88.\n\n**7.******Check** the Condition:**\n - We need to ensure that the total number of coins is more than 200.\n - Since 88 is not more than 200, - The next smallest x is x = 5 + 8 = 13. In - If x = 13, the total number of coins in the original seven bags is $7 \setminus times 13 = 91 \setminus n$ - Adding the 53 coins from the new bag, the total number of coins is 91 + 53 = 144.\n\n**9.**** Continue Checking:**\n - Since 144 is still not more than 200, we continue to the next smallest x.\n - The next smallest x is x = 13 + 8 = 21.\n - If x = 21, the total number of coins in the original seven bags is 7×147 . In - Adding the 53 coins from the new bag, the total number of coins is 147 + 53 = 200.\n\n**10**.****Final Check**:**\n - **Since** 200 is not more than 200, we need to find the next smallest x.\n - The next smallest x is x = 21 + 8 = 29.\n - If x = 29, the total number of coins in the original seven bags is $7 \setminus times 29 = 203 \cdot n$ - Adding the 53 coins from the new bag, the total number of coins is 203 + 53 = 256.\n\n[11.]**Conclusion:**\n-The smallest number of coins in each of the original seven bags that satisfies all conditions is x = 29.\n - **Therefore**, the smallest number of coins you could have had before finding the bag of 53 coins is $1 \leq 29 = 203$.\n\nThus, the final answer is:\n\\[\n**boxed{203}** \n\\].

RESPONSE_CORRECTNESS: To solve this problem, let's denote the number of coins in each of the original seven bags as x. **Therefore**, the total number of coins you originally had is 7x. After finding the bag of 53 coins, you now have 7x+53 coins. You redistribute these coins into eight bags, so each bag now has frac7x + 538 coins. **Since** the number of coins in each bag must be an integer, 7x + 53 must be divisible by 8. Additionally, we know that the total number of coins after redistribution is more than 200, so 7x + 53 > 200. This simplifies to 7x > 147, or x > 21. We need to find the smallest integer x such that 7x+53 is divisible by 8 and x > 21. Let's iterate through the possible values of x starting from 22 and **Check** the conditions.\n"'python\n Iterate through possible values of x starting from 22\nfor x in range(22, 1000):\n total_coins = $7 * x + 53 \setminus 162 - 100$. The smallest number_of_coins $3 \times 183 - 100 \times 183 \times 1000$.

Table S2: Over-reflection in sentence repetition. The red boxes in the figure indicate reflective words, while the blue boxes represent repeated phrases (with the numbers indicating the frequency of repetition). It can be observed that reflective words appear within the repeated phrases.

PROMPT: system \nPlease reason step by step, and put your final answer within \boxed{}.\n user \nDenali and Nate work for a dog walking business and are paid for each dog they walk. Denali is responsible for 16 dogs and Nate is responsible for 12 dogs. Under the company's new policy, they will be assigned or unassigned new dogs in groups of x dogs. The ratio of Denali's pay to Nate's pay would be the same if Denali started walking 4x more dogs and Nate stayed at 12 dogs or if x of Nate's dogs were reassigned to Denali. Find x if $x \setminus neq0$.\nassistant\n,

RESPONSE: Let's denote the payment per dog as $\backslash (p \setminus)$. The initial ratio of Denali's pay to Nate's pay is $\(\ pay = \ pay is \).\)$ more dogs and Nate stays at 12 dogs, the new ratio of their pays would be $\(\frac{16 + 4x}p12p = \frac{16 + 4x}12\)$. Since this ratio is equal to $\(\ 4x12 =$ $\frac{12}{12} = \frac{12}{12} = \frac{12$ $48,n^{n}=0,n^{n}=0,n^{n}=0,n^{n}$, we need to consider the second scenario where $\backslash(x \land)$ of Nate's dogs are reassigned to Denali. In this case, the new number of dogs Denali is responsible for is $\backslash (16 + x \land)$ and the new number of dogs Nate is responsible for is (12 - x). The new ratio of their pays would be ((16 + x)p(12 - x)p =)frac16 + x12 - x)x\). Since this ratio is equal to \\(\\frac43\\), we can set up the equation:\n\\[\n\\frac16 + x12 - x = $\frac{12 - x}{n} = 4(12 - x),n^{1} + 3x$ to **recheck** the problem constraints. The only way the ratio can remain the same is if the number of dogs reassigned is such that the ratio of the new number of dogs is the same as the original ratio. Let's solve it again by checking the constraints properly.\n\nThe correct approach is to solve the equation $\ (\ x_12 - x = \ x_243\)$ properly.

Repeat1: Let's solve it **step by step** again. $\n\(\n3(16 + x) = 4(12 - x), \n\)\n\(\n48 + 3x = 48 - 4x, \n\)\n\(\n7x = 0, \n\)\n\(\nx = 0.\n\)\n\Since \(x \neq 0\), we need to consider the possibility that the problem might have a different interpretation or constraint. Let's$ **recheck**the problem statement and constraints. The only way the ratio can remain the same is if the number of dogs reassigned is such that the ratio of the new number of dogs is the same as the original ratio. Let's solve it again by checking the constraints properly.

Repeat2: Let's solve it **step by step** again. $\n\n\[\n3(16 + x) = 4(12 - x), \n\]\n\[\n48 + 3x = 48 - 4x, \n\]\n\[\n7x = 0, \n\]\n\[\nx = 0.\n\]\n\Since \(x \neq 0\), we need to consider the possibility that the problem might have a different interpretation or constraint. Let's$ **recheck**the problem statement and constraints. The only way the ratio can remain the same is if the number of dogs reassigned is such that the ratio of the new number of dogs is the same as the original ratio. Let's solve it again by checking the constraints properly.

•••••



Figure S1: Accuracy curves on (a) MATH500, (b) AIME2024, and (c) AMC2023 benchmarks training on the DeepScaleR-17k.



Figure S2: Accuracy curves on (a) MATH500, (b) AIME2024, and (c) AMC2023 benchmarks training on Math-train dataset.



Figure S3: (a) Response length, (b) clip ratio, and (c) average accuracy of benchmark during training.