

---

# Analysis of Reinforcement Learning with Verifiable Rewards for Language Model Reasoning

Luka Meladze

National University of Singapore

Project: LLM Post-training with Reinforcement Learning

## Abstract

1 This report provides an analysis of two papers that explore enhancing Large  
2 Language Models' (LLMs) reasoning capabilities through "pure reinforcement  
3 learning". The first paper, "DeepSeek-R1: Incentivizing Reasoning  
4 Capability in LLMs via Reinforcement Learning," introduces a method  
5 to cultivate advanced reasoning in LLMs from a base model using RL  
6 (Guo et al., 2025). The second, "Understanding R1-Zero-Like Training: A  
7 Critical Perspective," critically explains the underlying assumptions and  
8 methodologies of the first paper, offering a more nuanced view and sug-  
9 gesting refinements (Liu et al., 2025). The paper will highlight details about  
10 "pure reinforcement learning," the differences between Reinforce-  
11 ment Learning from Human Feedback (RLHF) and the technique used  
12 here, which we can term Reinforcement Learning with Verifiable Rewards  
13 (RLVR), the algorithms that drive it and the importance of algorithmic  
14 design and base model selection.

## 15 1 From Subjective Feedback to Verifiable Outcomes: RLHF vs. RLVR

### 16 1.1 Shift in Reinforcement Learning for LLMs

17 The common application of Reinforcement Learning (RL) to LLMs is through Reinforcement  
18 Learning from Human Feedback (RLHF). The primary objective of RLHF is to align a  
19 model with human preferences, which are often subjective and difficult to quantify. In this  
20 paradigm, a separate reward model is trained on a dataset of human-provided comparisons  
21 (e.g., selecting which of two responses is better). The LLM is then fine-tuned using this  
22 reward model as the source of its learning signal.

23 The papers by Guo et al. (2025) and Liu et al. (2025) focus on a different more objective  
24 **Reinforcement Learning with Verifiable Rewards (RLVR)**. This approach is used for  
25 domains where correctness can be programmatically verified.

- 26 • **RLHF**: Optimizes for subjective, human-aligned behavior using a learned reward  
27 model. The goal is often helpfulness, harmlessness, and adherence to a specific  
28 experts response style.
- 29 • **RLVR**: Optimizes for verifiable correctness using a rule-based or verifier-based  
30 reward function. The goal is to improve performance on specific reasoning tasks  
31 like mathematics, coding, and logic puzzles.

32 RLVR sidesteps the costly human data collection and potential reward-hacking pitfalls of a  
33 learned reward model, providing a direct and unambiguous signal for improving a model's  
34 reasoning faculties.

### 35 1.2 Language Model Reasoning as a Markov Decision Process

36 The process of generating a reasoned response can be framed as a Markov Decision Process  
37 (MDP), a foundational concept in RL.

- 
- 38 • **State ( $s_t$ ):** The current state is the concatenation of the input question ( $q$ ) and the  
39 sequence of tokens generated so far ( $o_{<t}$ ).
  - 40 • **Action ( $a_t$ ):** The action is the selection of the next token ( $o_t$ ) from the model’s  
41 vocabulary.
  - 42 • **Policy ( $\pi_\theta$ ):** The LLM itself is the policy. Parameterized by  $\theta$ , it maps a state  $s_t$  to a  
43 probability distribution over all possible actions (tokens).
  - 44 • **Reward ( $r_t$ ):** In the RLVR setting, the reward is sparse. It is zero for all intermediate  
45 tokens and a terminal reward is given only at the end of the sequence. For example,  
46  $R(q, o) = 1.0$  if the final answer in the response  $o$  is correct, and 0.0 otherwise.
- 47 The objective of reinforcement learning is to update the policy’s parameters  $\theta$  to maximize  
48 the expected cumulative reward, effectively teaching the model to generate token sequences  
49 that lead to correct final answers.

## 50 2 The DeepSeek-R1 Methodology: A Blueprint for Reasoning

51 [Guo et al. \(2025\)](#) introduces a novel pipeline for developing reasoning capabilities, centered  
52 around two key models: **DeepSeek-R1-Zero** and the more refined **DeepSeek-R1**.

### 53 2.1 The “Pure RL” Paradigm: DeepSeek-R1-Zero

54 The most significant contribution is the concept of training a reasoning model via “pure  
55 reinforcement learning,” which means starting directly from a base model without an initial  
56 Supervised Fine-Tuning (SFT) phase on reasoning data.

57 **Base Model:** The process begins with DeepSeek-V3-Base, a foundational model not specifi-  
58 cally tuned for complex reasoning.

59 **Training Template:** To guide the model’s output structure, a specific template is enforced.  
60 This is critical for the verifier to parse the output.

61 User: {question}

62 Assistant: <think> reasoning process here </think> <answer> answer here </answer>

63 **Reward Model:** A simple, rule-based reward function provides the learning signal.

- 64 • **Accuracy Reward:** A primary reward is given if the content within the <answer> tag  
65 is correct. This is verified by checking against the ground-truth solution for math  
66 problems or running code against unit tests.
- 67 • **Format Reward:** A small auxiliary reward is given for correctly using the <think>  
68 and </think> tags, encouraging the model to adhere to the desired structure.

69 This pure RL process was shown to be remarkably effective, increasing the AIME 2024  
70 pass@1 score from 15.6% to 71.0%.

### 71 2.2 The GRPO Algorithm: An Efficient Policy Gradient Method

72 The RL updates are driven by **Group Relative Policy Optimization (GRPO)**, a policy  
73 gradient algorithm designed for efficiency at scale.

74 **Training Steps:**

- 75 1. **Group Sampling:** For each question  $q$  in a batch, the policy  $\pi_{\theta_{old}}$  generates a group  
76 of  $G$  different responses  $\{o_1, o_2, \dots, o_G\}$ .
- 77 2. **Reward Calculation:** Each response  $o_i$  is evaluated by the verifier, yielding a set of  
78 rewards  $\{r_1, r_2, \dots, r_G\}$ .

79 3. **Advantage Estimation:** The advantage  $A_i$  for each response is calculated relative to  
80 the other responses in its group. This is the core of GRPO’s efficiency.

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})} \quad (1)$$

81 4. **Policy Update:** The policy parameters  $\theta$  are updated using a clipped surrogate  
82 objective, similar to PPO, to maximize the likelihood of responses with high advan-  
83 tage.

### 84 2.3 The Full Pipeline: The Creation of DeepSeek-R1

85 While DeepSeek-R1-Zero demonstrated powerful reasoning, its outputs were often unpol-  
86 ished and sometimes mixed different languages. To create a more robust and user-friendly  
87 model, **DeepSeek-R1**, a multi-stage pipeline was developed.

- 88 1. **Cold Start (SFT):** The DeepSeek-V3-Base model is first fine-tuned on a small (thou-  
89 sands) set of high-quality, human-readable reasoning examples. This provides a  
90 better starting point for RL.
- 91 2. **Reasoning-Oriented RL:** The GRPO algorithm is applied to this “cold-started”  
92 model to aggressively enhance its core reasoning abilities on math and code.
- 93 3. **Rejection Sampling & SFT:** The powerful RL-trained model is used to generate a  
94 large dataset of correct solutions ( $\sim 600k$ ). This high-quality synthetic data, com-  
95 bined with data for general tasks (writing, Q&A), is used for another round of SFT  
96 to broaden the model’s capabilities.
- 97 4. **RL for All Scenarios:** A final RL stage is performed to align the model with human  
98 preferences for helpfulness and harmlessness, using a hybrid of RLVR for reasoning  
99 and traditional RLHF for general dialogue.

## 100 3 A Critical Perspective: Deconstructing the R1-Zero Paradigm

101 The second paper by [Liu et al. \(2025\)](#) provides a critical re-evaluation of the “R1-Zero-like”  
102 training paradigm, revealing that the narrative of “pure emergence” is more complex.

### 103 3.1 The Illusion of Emergence: Pre-existing Capabilities

104 The analysis reveals that the base models are not the blank slates they might appear to be.

- 105 • **Latent Abilities:** Base models like DeepSeek-V3-Base and especially the Qwen2.5  
106 family already possess significant reasoning capabilities before any RL is applied.
- 107 • **“Aha Moment” is Not Emergent:** The phenomenon of the model appearing to  
108 self-correct mid-thought (the “aha moment”) was found to exist in the base models  
109 themselves. RL amplifies this behavior rather than creating it from nothing.
- 110 • **Pre-training Contamination:** Some models, like Qwen2.5, perform best without any  
111 template, suggesting they were likely pre-trained on concatenated question-answer  
112 text, making them functionally similar to SFT models from the outset.

### 113 3.2 Algorithmic Flaws and Unintended Biases in GRPO

114 Most critically, the paper identifies two significant biases in the GRPO objective function  
115 that can lead to misleading interpretations of model behavior.

- 116 1. **Response-Level Length Bias:** The GRPO loss is normalized by the length of the  
117 response. For an incorrect response (negative advantage), a longer response re-  
118 ceives a smaller penalty. This incentivizes the model to generate increasingly long,  
119 meandering chains of thought for its *incorrect* answers, which can be mistaken for  
120 deeper reasoning.

121 2. **Question-Level Difficulty Bias:** Normalizing the advantage by the standard deviation  
122 of rewards within a group gives disproportionately high weight to questions  
123 that are either very easy or very hard (where reward variance is low). This biases  
124 the learning process.

### 125 3.3 The Fix: GRPO Done Right (Dr. GRPO)

126 The authors propose a simple and effective solution: **Dr. GRPO**, which removes the biasing  
127 normalization terms from the objective. This seemingly small change has a profound impact.  
128 The following code snippet highlights the difference between a typical, biased PPO loss  
129 implementation (as found in many libraries) and the unbiased approach.

```
130 # A common, but biased, way to calculate loss in many RL libraries
131 # This normalizes loss by the length of each individual response,
132 # introducing bias.
133 def biased_loss_calculation(ppo_loss, response_mask):
134     # masked_mean normalizes by the number of non-pad tokens in each
135     # response
136     per_response_loss = masked_mean(ppo_loss, response_mask, dim=-1)
137     return per_response_loss.mean()
138
139 # The unbiased approach proposed by Dr. GRPO
140 # This normalizes by a fixed constant, removing the length bias.
141 MAX_TOKENS = 4096 # A fixed budget
142 def unbiased_loss_calculation(ppo_loss, response_mask):
143     # Here, the loss for each token is summed and normalized by a
144     # constant
145     total_loss = (ppo_loss * response_mask).sum()
146     return total_loss / (BATCH_SIZE * MAX_TOKENS)
```

149 Dr. GRPO achieves comparable or better performance while being significantly more  
150 token-efficient, as it no longer rewards the model for generating lengthy incorrect responses.

## 151 4 Conclusion and Key Takeaways

152 The insights from DeepSeek-R1’s initial breakthrough paper to the critical analysis in the  
153 second paper provides important considerations for the successful application of RL to  
154 LLMs.

- 155 1. **RLVR is a Powerful Tool for Objective Tasks:** For domains with verifiable correctness, RLVR provides a direct, scalable, and powerful method for improving model capabilities without the complexities and extra costs of RLHF.
- 156 2. **Base Models are Not Blank Slates:** The reasoning success of RL is heavily dependent on the pre-existing, latent capabilities of the base model. The notion of “pure emergence” should be treated with skepticism; RL is more accurately described as a process of amplifying and refining these latent skills.
- 157 3. **Algorithmic Design is Critical:** Subtle choices in the RL algorithm’s objective function can introduce significant, unintended biases. An unbiased optimizer like Dr. GRPO is more robust, efficient, and leads to more interpretable model behavior.
- 158 4. **Iterative Refinement and Distillation are State-of-the-Art:** The most effective reasoning models are built through a multi-stage, iterative process that combines the strengths of SFT (for learning patterns) and RL (for exploration and generalization). The reasoning abilities of these large models can then be effectively distilled into smaller, more accessible models.

---

## References

- 170  
171 Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X.,  
172 et al. (2025). Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement  
173 learning. *arXiv preprint arXiv:2501.12948*.
- 174 Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W.S., & Lin, M. (2025). Understanding  
175 r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.