

Motivation

- **Long-Context Reasoning Bottlenecks:** The scarcity of reliable human annotations and programmatically verifiable reward signals.
- **Non-Verifiable RL:** String match is **brittle and unreliable** for semantic equivalence.
- **SPELL:** A scalable, **label-free** Self-Play framework where a single LLM generates, solves, and verifies its own long-context tasks.

Long-Context RL

- Maximize the KL-regularized expected reward:

$$\max_{\pi_{\theta}} \mathbb{E}_{\mathbf{c}, \mathbf{q} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{c}, \mathbf{q})} [r_{\phi}(\mathbf{c}, \mathbf{q}, \mathbf{y})] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta}(\mathbf{y} | \mathbf{c}, \mathbf{q}) || \pi_{\text{ref}}(\mathbf{y} | \mathbf{c}, \mathbf{q})]$$

Long Documents

- Fully on-policy GRPO (no KL regularization):

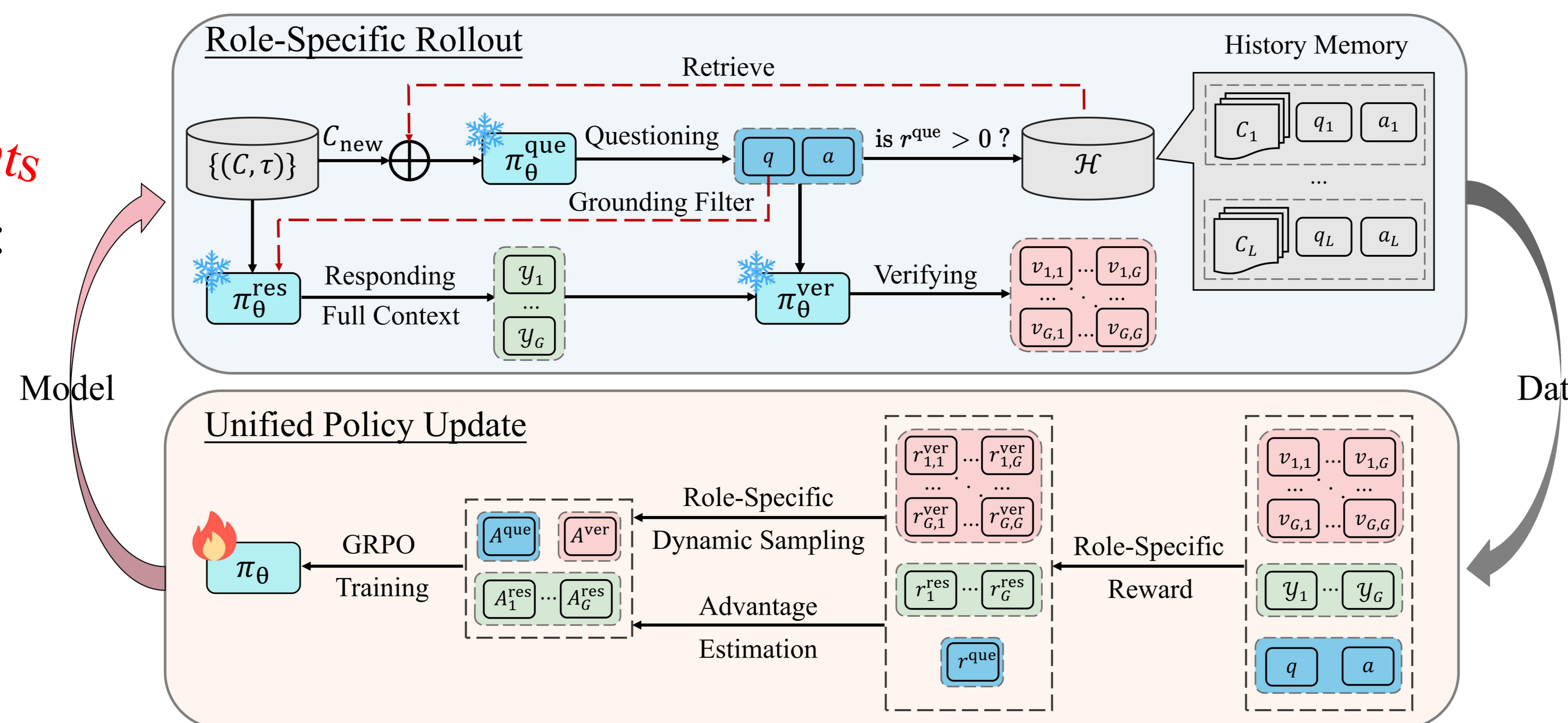
$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\mathbf{c}, \mathbf{q} \sim \mathcal{D}, \{\mathbf{y}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}}$$

$$\left[\frac{1}{\sum_{j=1}^G |\mathbf{y}_j|} \sum_{i=1}^G A_i \sum_{t=1}^{|\mathbf{y}_i|} \frac{\pi_{\theta}(\mathbf{y}_{i,t} | \mathbf{c}, \mathbf{q}, \mathbf{y}_{i,<t})}{\pi_{\theta_{\text{old}}}(\mathbf{y}_{i,t} | \mathbf{c}, \mathbf{q}, \mathbf{y}_{i,<t})} \right]$$

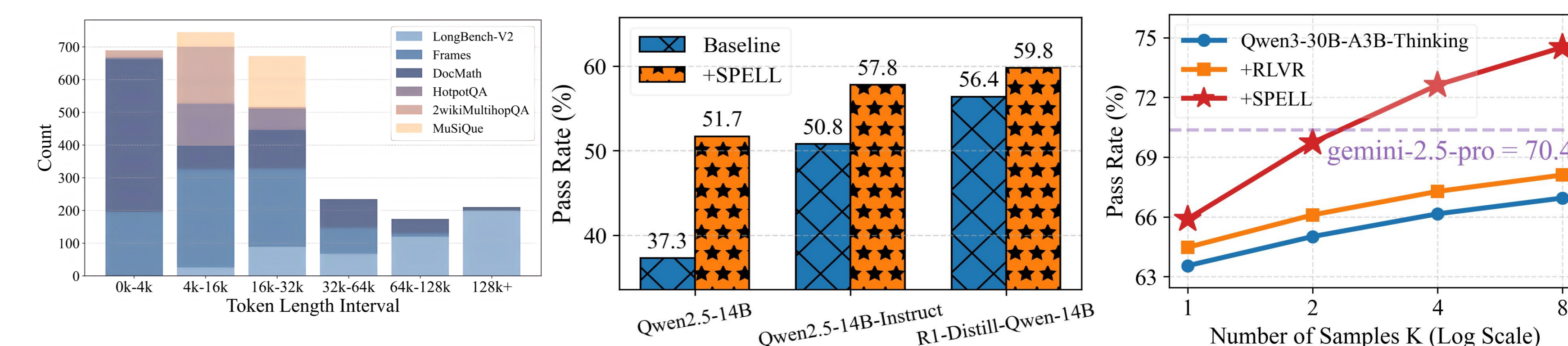
$$A_i = \frac{r_i - \text{mean}(\{r_k\}_{k=1}^G)}{\text{std}(\{r_k\}_{k=1}^G)}$$

Method

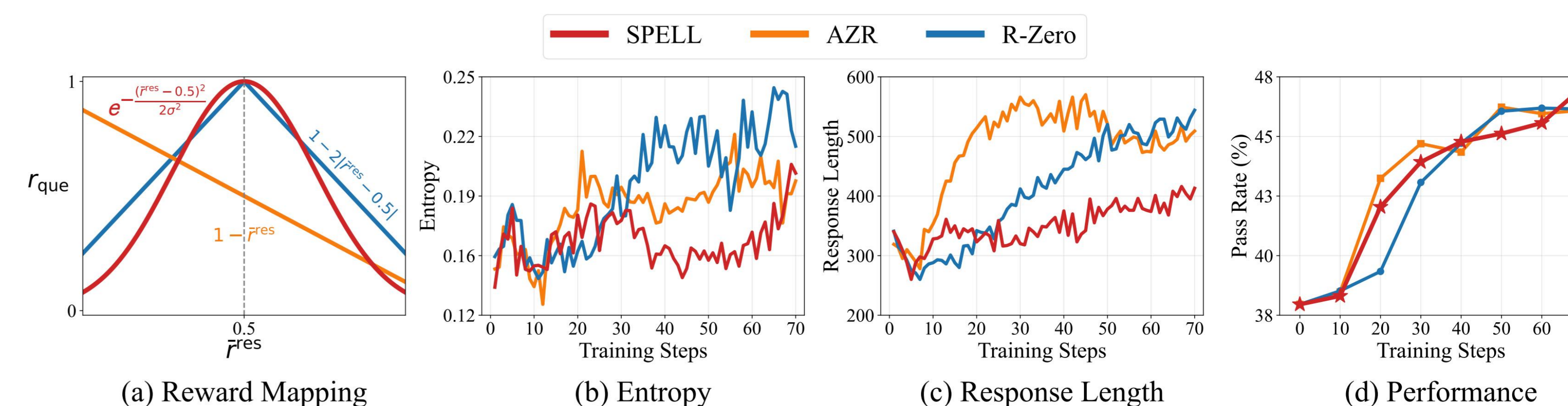
- **Role-Specific Rollout:** A single policy model (π_{θ}) dynamically assumes three complementary roles:
 - Questioner: Generates questions based on documents and history memory.
 - Responder: Attempts to solve the generated questions.
 - Verifier: Evaluates semantic equivalence between the responder's output and the reference answer to produce reward signals.
- **Role-Specific Reward:**
 - Verifier: Trained to improve judgment reliability through **self-consistency**.
 - Responder: Combines a rule-based check and the verifier's consensus.
 - Questioner: Incentivized to generate questions of **intermediate difficulty**.
- **Unified Policy Update:** Jointly optimizes the fully on-policy GRPO objective across all three roles.



Experiments



- **Broad Generalizability:** Consistent improvements across 12 open-source LLMs (4B to 32B parameters), including dense and MoE architectures, Llama and Qwen families.
- **Superior Exploratory Ability:** SPELL achieves an average 7.6-point gain in pass@8 on Qwen3-30B-A3B-Thinking, significantly outperforming traditional RLVR.



- **Stable Training Dynamics:** SPELL outperforms prior Self-Play baselines and maintains a more stable training entropy and more controlled growth in response length.



Contact us

- **Shenzhen Area Loop Institute (深圳河套学院)** is recruiting: **Professor/AP/RAP/RA/PhD**
- Contact Email: xiaojunquan@slai.edu.cn