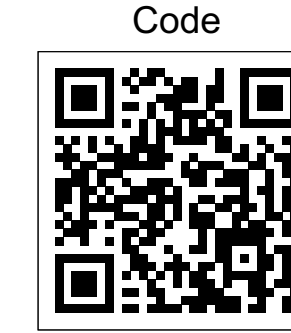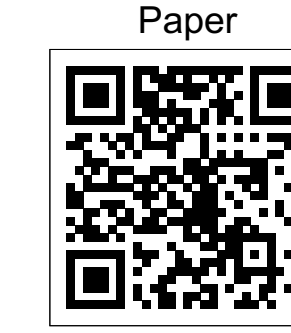# Segment Policy Optimization: Effective Segment-Level Credit Assignment in RL for Large Language Models
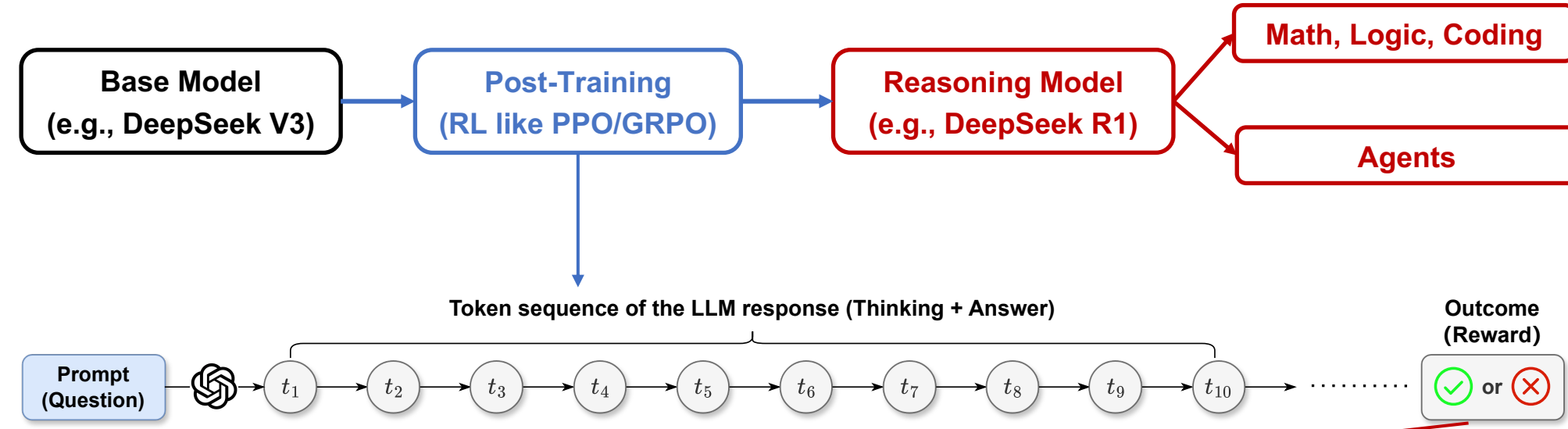
Yiran Guo[1], Lijie Xu[1*], Jie Liu[1*], Dan Ye[1], Shuang Qiu[2]

[1]Institute of Software, Chinese Academy of Sciences   [2]City University of Hong Kong

ISCAS 中国科学院软件研究所 Institute of Software Chinese Academy of Sciences

CityU 香港城市大學 City University of Hong Kong

NEURAL INFORMATION PROCESSING SYSTEMS

Paper   Code   Covered by Jiqizhixin

## Key problem and challenges of RL for LLM

**LLM Training Pipeline**



Base Model (e.g., DeepSeek V3) → Post-Training (RL like PPO/GRPO) → Reasoning Model (e.g., DeepSeek R1) → Math, Logic, Coding / Agents

Token sequence of the LLM response (Thinking + Answer)

Prompt (Question) → $t_1$ $t_2$ $t_3$ $t_4$ $t_5$ $t_6$ $t_7$ $t_8$ $t_9$ $t_{10}$ ... → Outcome (Reward) ✓ or ✗

**Key Problem: Credit assignment**
How to attribute the final evaluation result (reward signal) of the entire sequence (LLM response) to the specific decision-making actions (i.e., tokens) within that sequence?
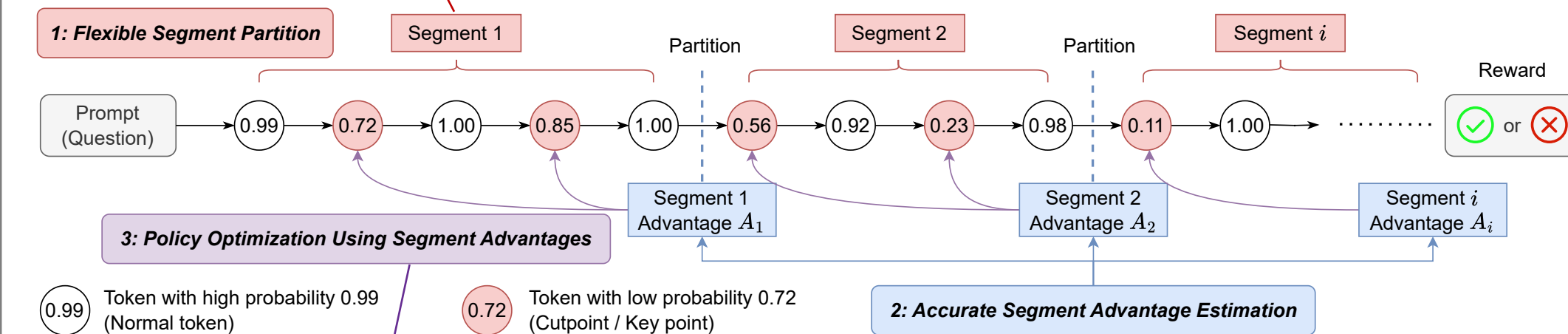
**Challenges**
Different from traditional RL, the LLM reward is extremely sparse (clear success or failure feedback is only available at the end of the entire sequence).

## Existing methods (Token-level vs. Trajectory-level)



Prompt (Question) → $t_1$ $t_2$ $t_3$ $t_4$ $t_5$ ... → $r$ ✓

Critic model (Value network)

$V_1$ $V_2$ $V_3$ $V_4$ $V_5$ → GAE (Generalized Advantage Estimation)

$A_1$ $A_2$ $A_3$ $A_4$ $A_5$

Update LLM policy

**Token-level / fine-grained methods (e.g., PPO)**

**Trajectory-level / coarse-grained methods (e.g., GRPO)**

Sample N responses (e.g., N = 3)

Prompt (Question) → $t_{11}$ $t_{12}$ $t_{13}$ $t_{14}$ $t_{15}$ → $r_1$ ✓ → $A_1$
Prompt (Question) → $t_{21}$ $t_{22}$ $t_{23}$ $t_{24}$ $t_{25}$ → $r_2$ ✗ → $A_2$
Prompt (Question) → $t_{31}$ $t_{32}$ $t_{33}$ $t_{34}$ $t_{35}$ → $r_3$ ✓ → $A_3$

Reward / Group Computation $\frac{r_i - mean}{std}$ / Advantages

Update LLM policy

| Method Granularity | Example | Mechanism | Drawbacks |
|---|---|---|---|
| Token-Level (Fine-Grained) | PPO | Uses a critic model to estimate advantage ($A_i$) for every token ($t_i$). | Inaccurate critic: Hard to train the critic model, leading to unreliable value $V$ predictions. High overhead: Critic model training/prediction. |
| Trajectory-Level (Coarse-Grained) | GRPO | Uses a single advantage signal ($A_i$) from the final reward for the entire sequence (each $t_{ij}$ gets the same $A_i$). | Imprecise credit: Cannot reward partial progress or penalize redundancy. Overfitting: Performance on validation sets saturates early. |

## Our SPO framework (Segment-level)

**Key idea:** Mid-grained (segment-level) advantage estimation can unify and overcome the limitations of token-level and trajectory-level methods. (See Feature 1 and 2; Segment → a number of contiguous tokens)

**Feature 1: Flexible segment partition (without requiring semantic completeness)**
→ Unify token-level and trajectory-level methods (flexible adjustment of granularity from token-level to trajectory-level)

### Our SPO framework with three components



*1: Flexible Segment Partition*
Segment 1, Partition, Segment 2, Partition, Segment i

Prompt (Question) → 0.99 0.72 1.00 0.85 1.00 0.56 0.92 0.23 0.98 0.11 1.00 → Reward ✓ or ✗

*3: Policy Optimization Using Segment Advantages*
Segment 1 Advantage $A_1$, Segment 2 Advantage $A_2$, Segment i Advantage $A_i$

0.99 Token with high probability 0.99 (Normal token)
0.72 Token with low probability 0.72 (Cutpoint / Key point)

*2: Accurate Segment Advantage Estimation*

**Feature 3: Probability-mask optimization strategy (enhance the credit assignment)**
→ Selectively assign the segment advantages to critical (low-probability) tokens instead of all tokens within a segment.
→ Critical tokens (*cutpoints*) represent positions where the model's reasoning trajectory could diverge.

**Feature 2: Accurate segment advantage estimation (based on Monte Carlo (MC) sampling)**
→ SPO vs. PPO (SPO eliminates the need for an additional, unstable critic model.)
→ SPO vs. GRPO (SPO can reward partial progress for unsuccessful responses and penalize redundancy or unnecessary portions within successful responses.)

Compute the segment advantage $\hat{A}_1^{seg} = \hat{V}(s_{t_2}) - \hat{V}(s_{t_1}) = -\frac{1}{3}$

$\hat{A}(n) = 1 - \frac{5}{9} = \frac{4}{9}$   $\hat{V}(\text{Pa}(n)) = (1 + \frac{2}{3} + 0)/3 = \frac{5}{9}$

### Two algorithms based on SPO

**SPO-chain for short CoT**



(a) Chain-based
$\hat{V}(s_{t_1}) = \frac{2}{3}$
Prompt (Question)
Samples e.g., $N = 3$ $\tau^{(1)}$ $\tau^{(2)}$ $\tau^{(3)}$

Compute the segment advantage $\hat{A}_1^{seg} = \hat{V}(s_{t_2}) - \hat{V}(s_{t_1}) = -\frac{1}{3}$
$\hat{V}(s_{t_2}) = \frac{1}{3}$

✓ ✗ ✗   ✓ ✗ ✗

**SPO-tree for long CoT**

(b) Tree-based
Compute the segment advantage $\hat{A}(n) = 1 - \frac{5}{9} = \frac{4}{9}$
$|\text{Ch}(n)| = 3$
Prompt (Question)
$\hat{V}(\text{Pa}(n)) = (1 + \frac{2}{3} + 0)/3 = \frac{5}{9}$

Segment Group 0   One segment has M tokens
$\frac{2}{3}$   1

Segment Group 1   Segment Group 2   Segment Group 3
✓ ✗ ✓   ✓ ✓ ✗   ✗ ✗ ✗

**Features of SPO-chain:**
(1) **Adaptive cutpoint-based segment partition strategy:** Let each segment contain a number of cutpoints, avoiding unnecessary segment (where all tokens within the segment have probabilities close to 1).
(2) **Chain-based Monte Carlo sampling:** For each segment, SPO independently samples N trajectories to estimate the value and segment advantage.

**Features of SPO-tree:**
(1) **Fixed token count segment partition strategy:** Each segment contains a fixed number of tokens (to well support tree-based sampling; unlikely all tokens within a long segment have probabilities close to 1).
(2) **Tree-based segment advantage estimation:** Directly yield segment advantages after tree-based trajectory generation, avoiding costly resampling and significantly improving efficiency for long CoT.

## Example

**Segment partition example (Each segment contains 5 cutpoints, i.e., red tokens)**

Prompt: [MATH_TASK] Problem: Gloria wants to buy the $129,000 mountain cabin that her friend Alfonso is selling. She only has $150 in cash. She intends to raise the remaining amount by selling her mature trees for lumber. She has 20 cypress trees, 600 pine trees, and 24 maple trees. She will get $100 for each cypress tree, $300 for a maple tree, and $200 per pine tree. After paying Alfonso for the cabin, how much money will Gloria have left? Solution:

Segment 1   Blue vertical line → segment boundary

Red tokens → cutpoints (tokens with probability < 0.9)

The total value of the cypress trees that she will | sell is 20 * $100 = $2000.
The total value of the maple trees that she will sell is | 600 * $300 = $18000.
The total value of the pine trees that she will sell is 24 * $200 = $4800.
The total value of all the trees that she will sell is $2000 + $18000 + $4800 = $24800.
| After selling all the trees, she will have $24800 - $150 = $23300 left.
#### 23300

Problem: The total value of the maple trees should be 24 * $300, but the model outputs 600 * $300 here and shows low confidence at the digit 6.
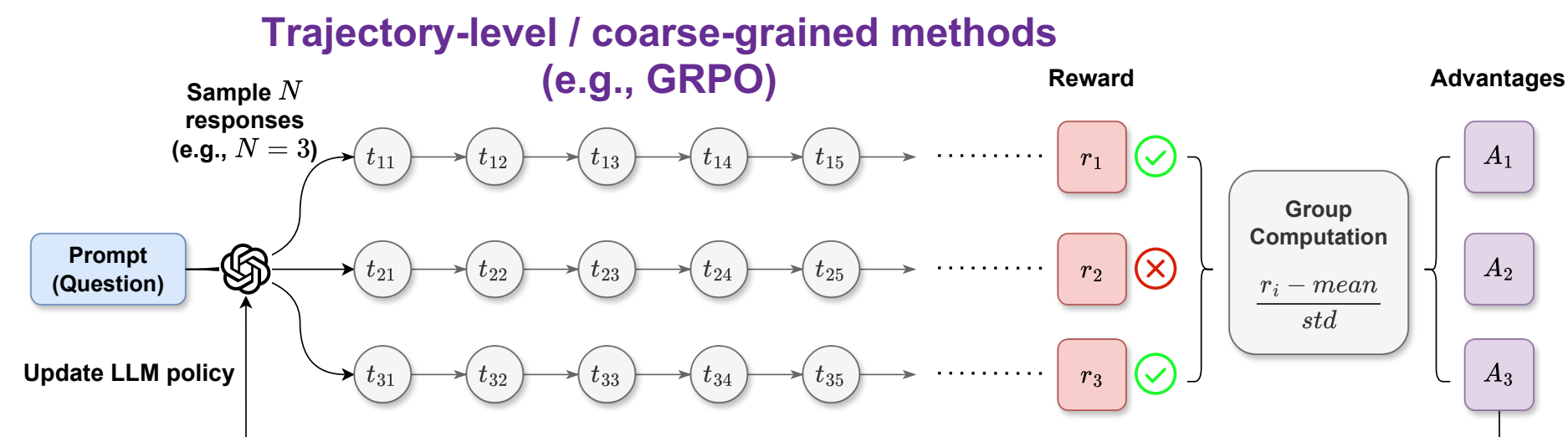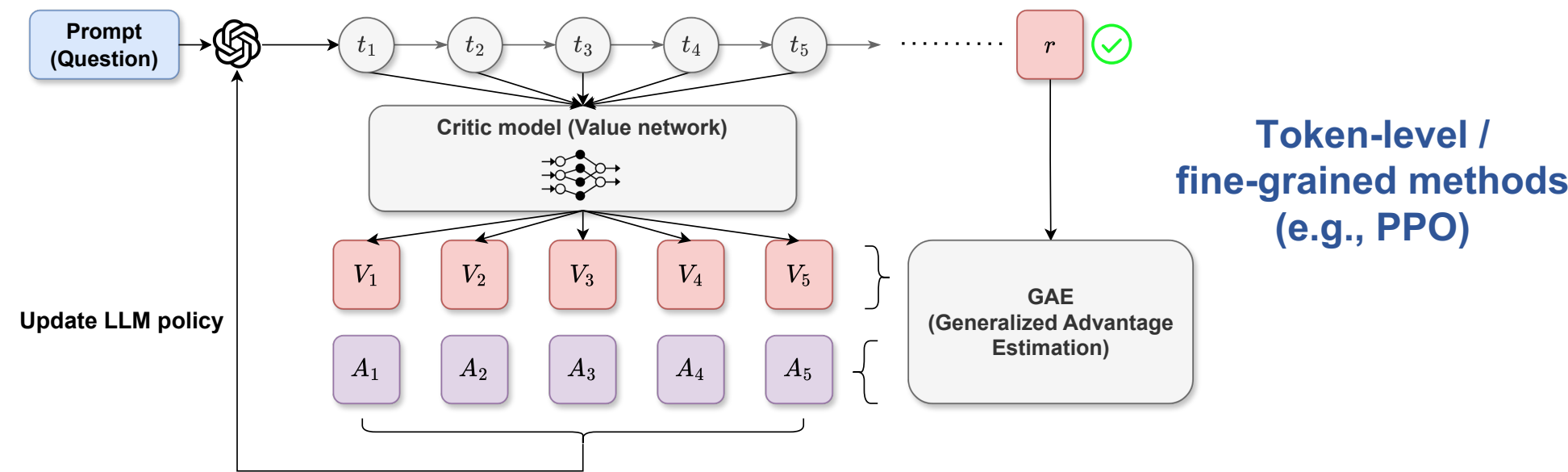
Problem: The calculation should be $24800 - $12900, but the model instead outputs $24800 - $150 and is uncertain at the digit 5.

**Implication: Cutpoints are the locations where the model's reasoning trajectory can shift, and they are the main drivers behind "segment advantage".**

**Cutpoint-based segment partition strategy (generate effective segments; See Features of SPO-chain)**

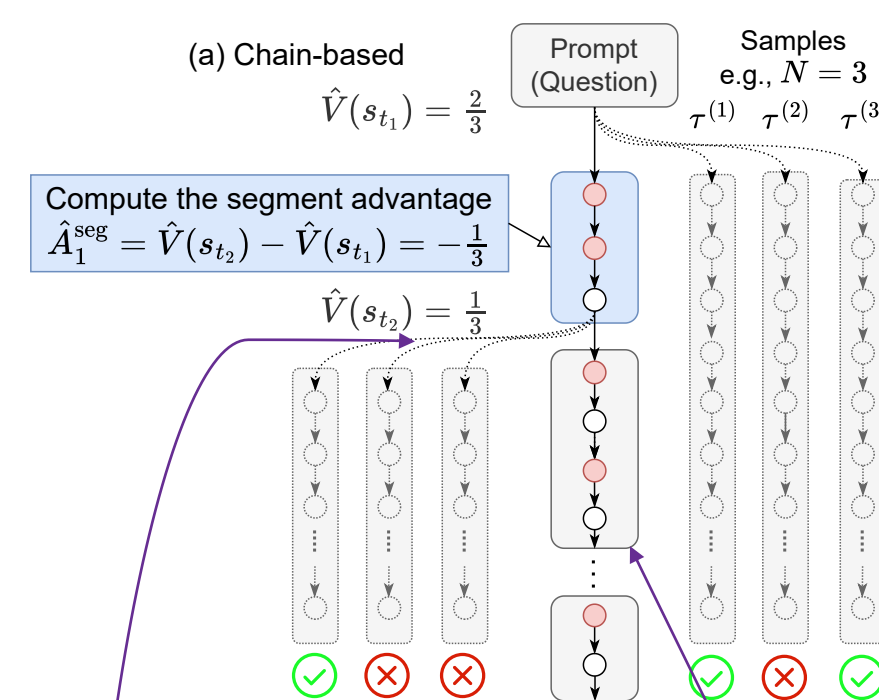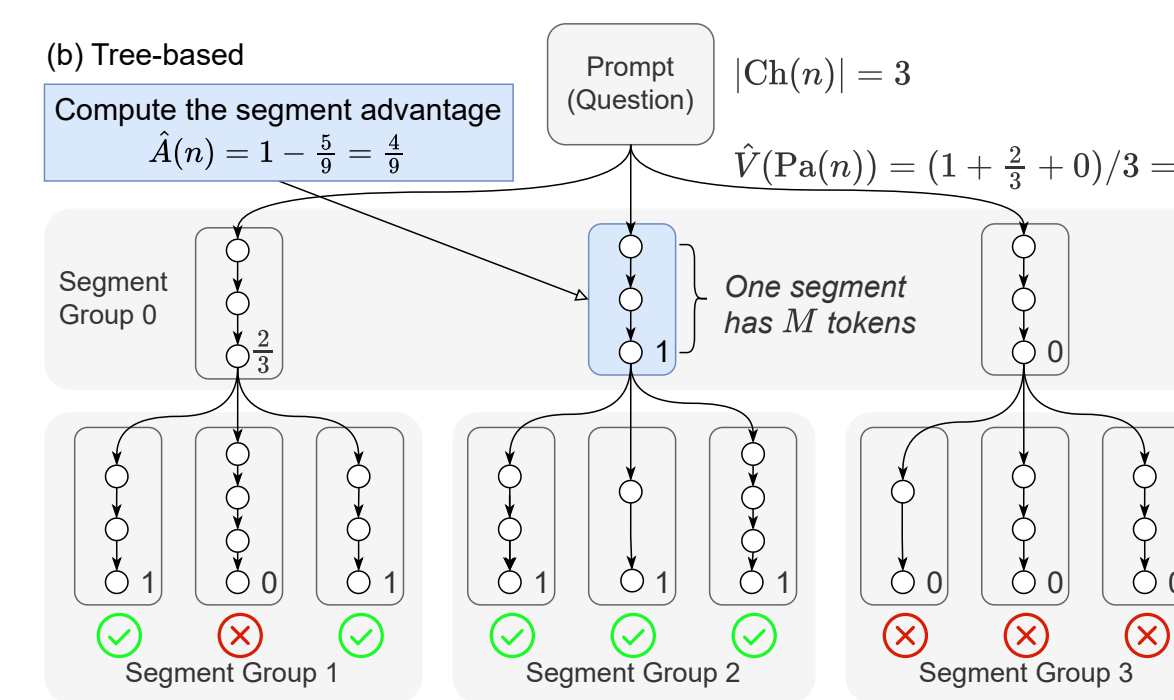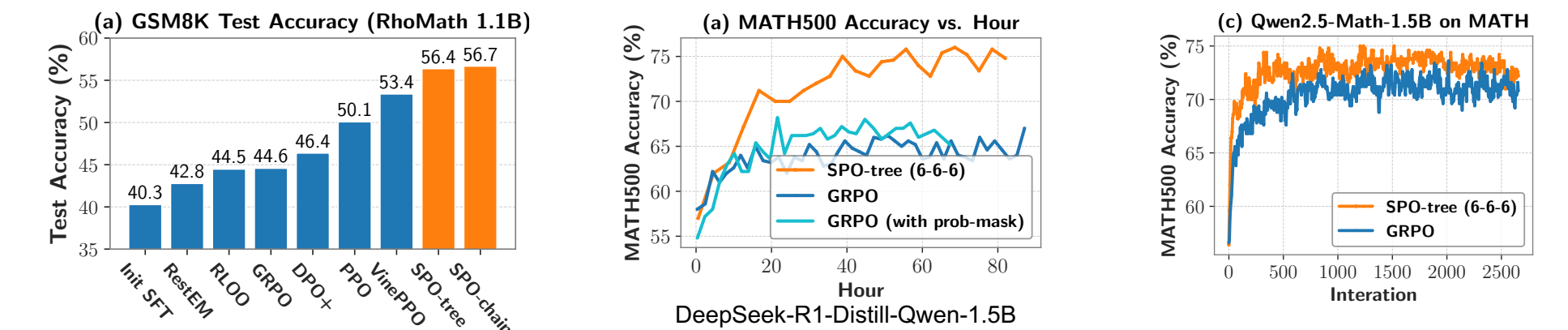**Probability-mask optimization strategy (optimize critical tokens to enhance credit assignment; See Feature 3 of SPO)**

## Experiments

**SPO-chain and SPO-tree outperform existing methods on GSM8K and MATH.**



(a) GSM8K Test Accuracy (RhoMath 1.1B)
Init SFT 40.3, RFT 42.8, RLOO 44.5, GRPO 44.6, DPO+ 46.4, PPO 50.1, VinePPO 53.4, SPO-chain 56.4, SPO-tree 56.7

(b) Episode Generation Time on GSM8K — VinePPO, SPO-chain (int5)

(c) Validation Accuracy on GSM8K — SPO-chain (int5), GRPO

**SPO also outperforms GRPO on logic task (Knights and Knaves Puzzle)**

(a) Qwen2.5-0.5B-Instruct on GSM8K — SPO-tree (6-6-6), GRPO

Accuracy comparison with various context sizes with DeepSeek-R1-Distill-Qwen-1.5B

*Different partition granularities*
(a) Partition Granularity Variations — SPO-chain (int2), SPO-chain (int5), SPO-chain (int100)

*Cutpoint-based vs. Other partition strategies*
(b) Partition Strategy Variations — SPO-chain (int5), VinePPO, Fixed-token-count

*Effectiveness of probability-mask strategy*
(c) Probability-Mask Strategy Ablation — SPO-chain (int5), SPO-chain (int5) (no prob-mask), GRPO, GRPO (with prob-mask)