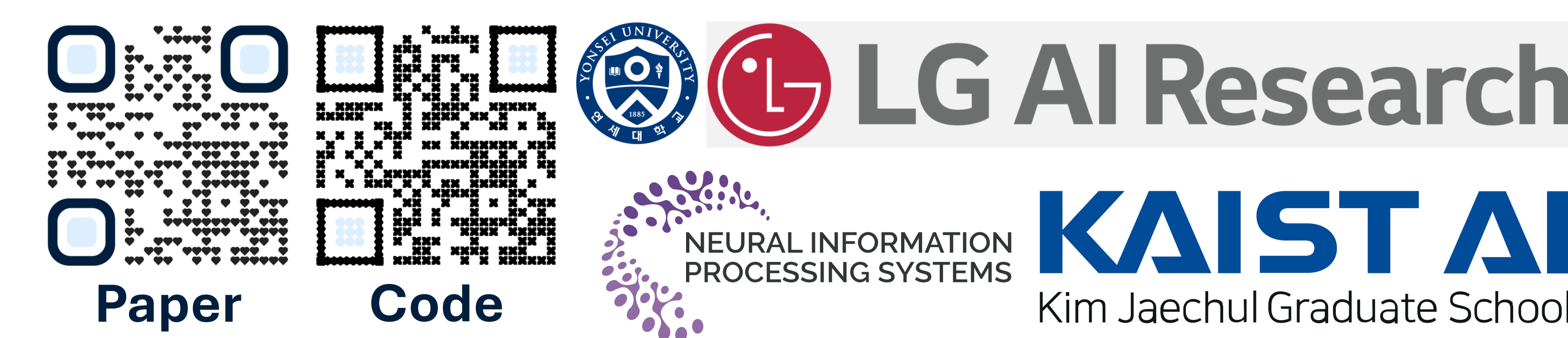


AdaSTaR: Adaptive Data Sampling for Training Self-Taught Reasoners

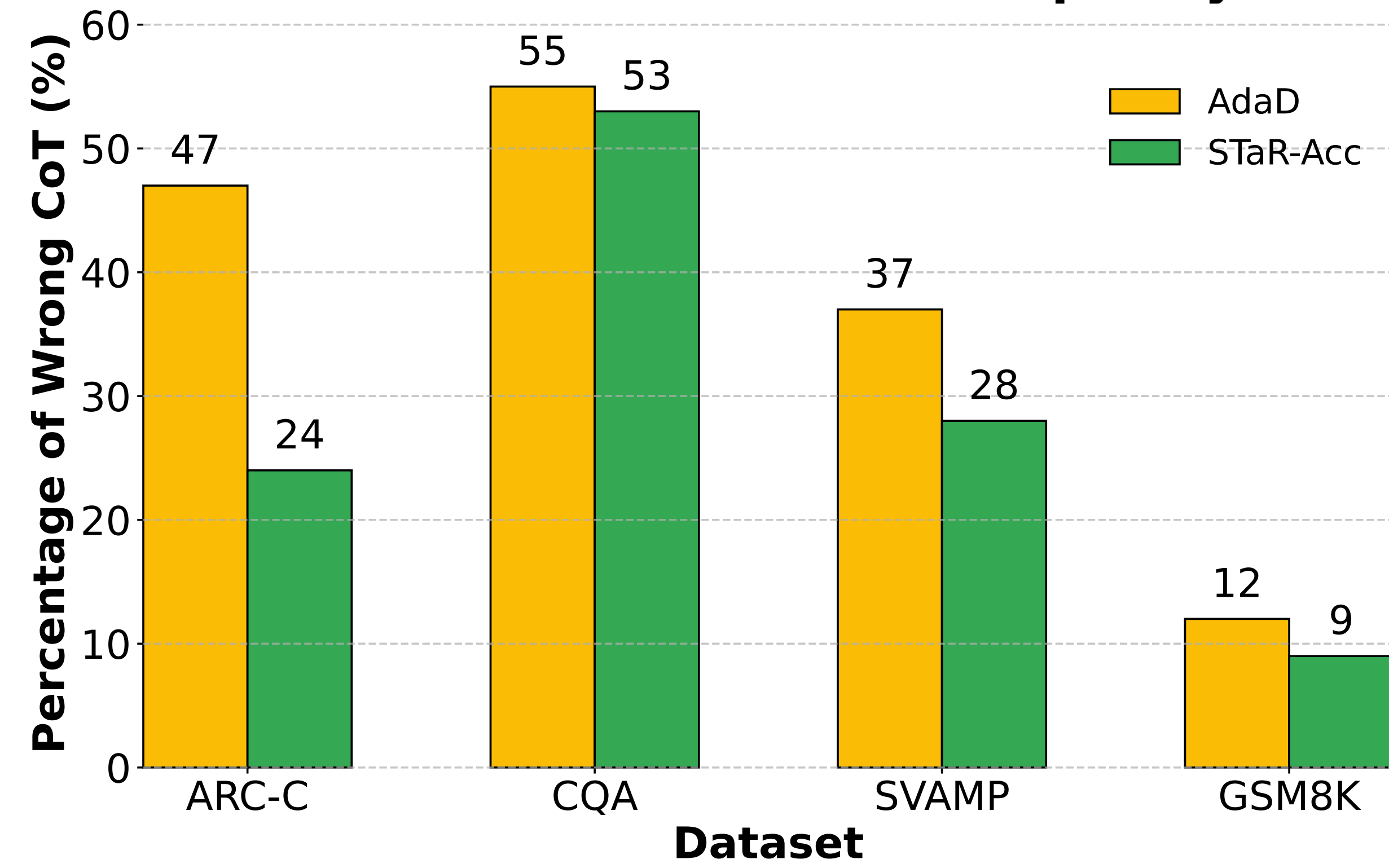
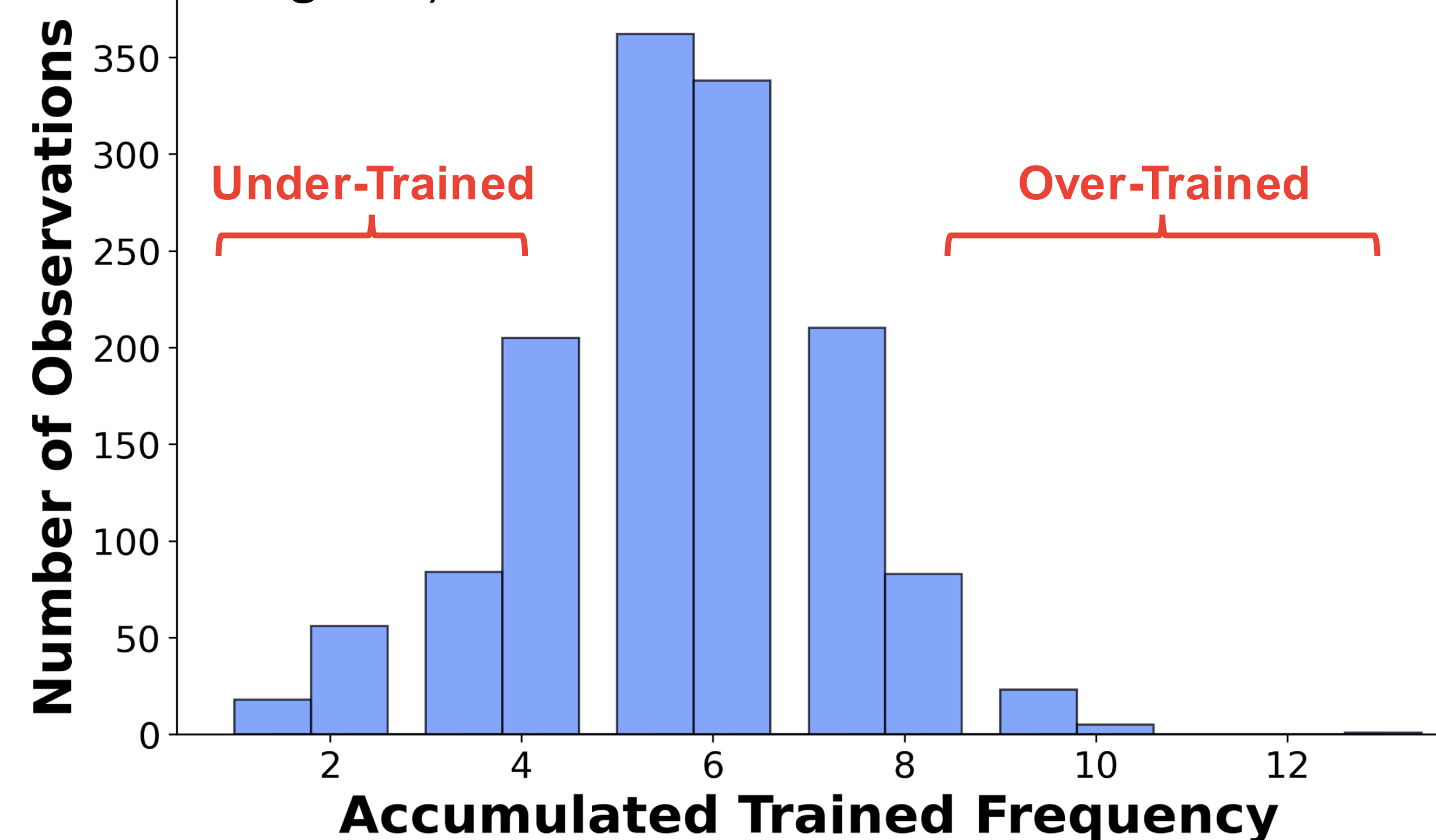
Woosung Koh, Wonbeen Oh, Jaein Jang, MinHyung Lee, Hyeongjin Kim, Ah Yeon Kim, Joonkee Kim, Junghyun Lee, Taehyeon Kim, Se-Young Yun

 **Still using *sample inefficient* Self-Taught Reasoners (Rejection Fine-tuning)?**

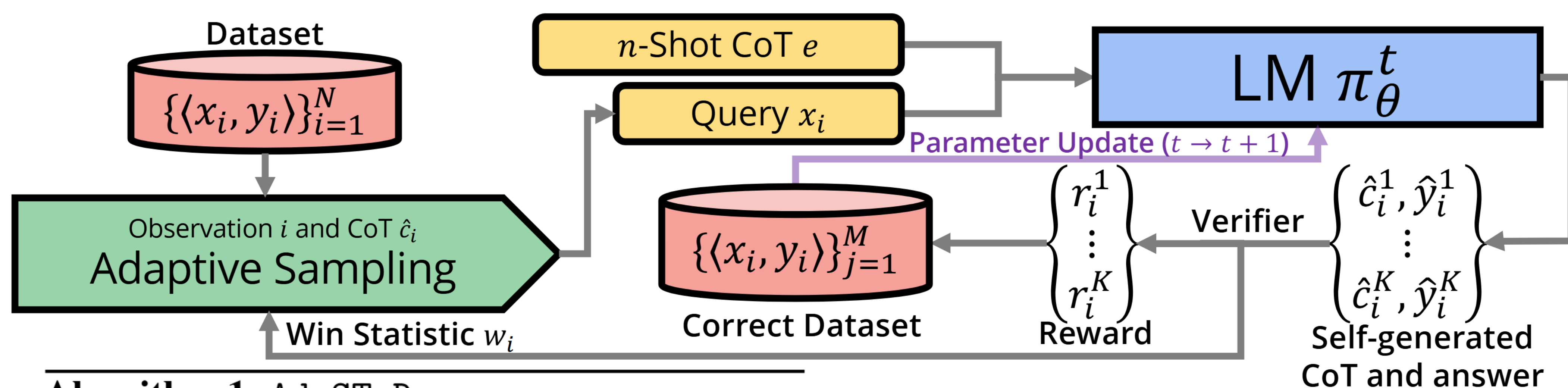


Existing Problem

- **STaR (or RFT)** allows LMs to improve iteratively by training on their own correct CoT solutions.
- Standard STaR relies on **random data sampling**.
- **Inefficiency:** Random sampling causes the model to over-train on easy examples (solved repeatedly) while under-training on challenging ones.
- **False Positives:** Simply prioritizing "hard" problems increases false positives (correct final answer, but flawed reasoning/CoT).



Method: Adaptive Sampling for Diversity and Curriculum



Algorithm 1: AdaSTaR

```

Input:  $\mathcal{D}, \pi_{\theta}^{t=1}, e$ 
/* AdaD (§3.1; lines 1-14) */
1  $\tilde{t} \leftarrow \text{dict}\{i : \tilde{t}_i = 0\}_{i=1}^N$ ;
2  $w \leftarrow \text{dict}\{i : w_i = 0\}_{i=1}^N$ ;
3 init  $\text{HieMinHeap}(\mathcal{D}, \tilde{t}, w)$ ;
4 for iteration  $t = 1, \dots$  do
5    $\mathcal{D}_+^t \leftarrow \emptyset, m \leftarrow 0$ ;
6    $w^{tmp} \leftarrow \text{dict}\{i : w_i^{tmp} = 0\}_{i=1}^N$ ;
7   while  $|\mathcal{D}_+^t| < \beta^t$  do
8      $i \leftarrow \text{HieMinHeap.peek\_next}$ ;
9      $m \leftarrow m + 1$ ;
10    for sample  $k = 1, \dots, K$  do
11       $\langle \hat{c}_i, \hat{y}_i \rangle \leftarrow \pi_{\theta}^t(e, x_i)$ ;
12       $w_i^{tmp} \leftarrow \frac{k-1}{k} w_i^{tmp} + \frac{1}{k} \mathbb{I}[\hat{y}_i = y_i]$ ;
13      if  $\hat{y}_i = y_i$  then
14         $\mathcal{D}_+^t \leftarrow \mathcal{D}_+^t \cup \{ \langle x_i, \hat{c}_i, \hat{y}_i \rangle \}$ ;
/* AdaC (§3.2; lines 15-19) */
15  $\alpha, \pi_{\theta}^{t+1} \leftarrow \text{Train}(\pi_{\theta}^t, \mathcal{D}_+^t)$ ;
16 for  $1, \dots, \lfloor m\alpha^2 \rfloor$  do
17    $i \leftarrow \text{HieMinHeap.pop}$ ;
18    $\tilde{t}_i \leftarrow t, w_i \leftarrow w_i^{tmp}$ ;
19    $\text{HieMinHeap.push}(i, \tilde{t}_i, w_i)$ ;
    
```

- **Pseudocode: Color-coded**
 - **Black:** Existing STaR
 - **Green:** AdaSTaR
- **Diversity (AdaD):** Encourage more diverse training samples (\downarrow standard deviation)
 - \tilde{t}_i : Last sampled iteration
 - w_i : Win statistic (accuracy of i at \tilde{t}_i)
- **Curriculum (AdaC):** Sample easier observations earlier and harder later
 - α : Train accuracy

✓ TL;DR: Main Contribution

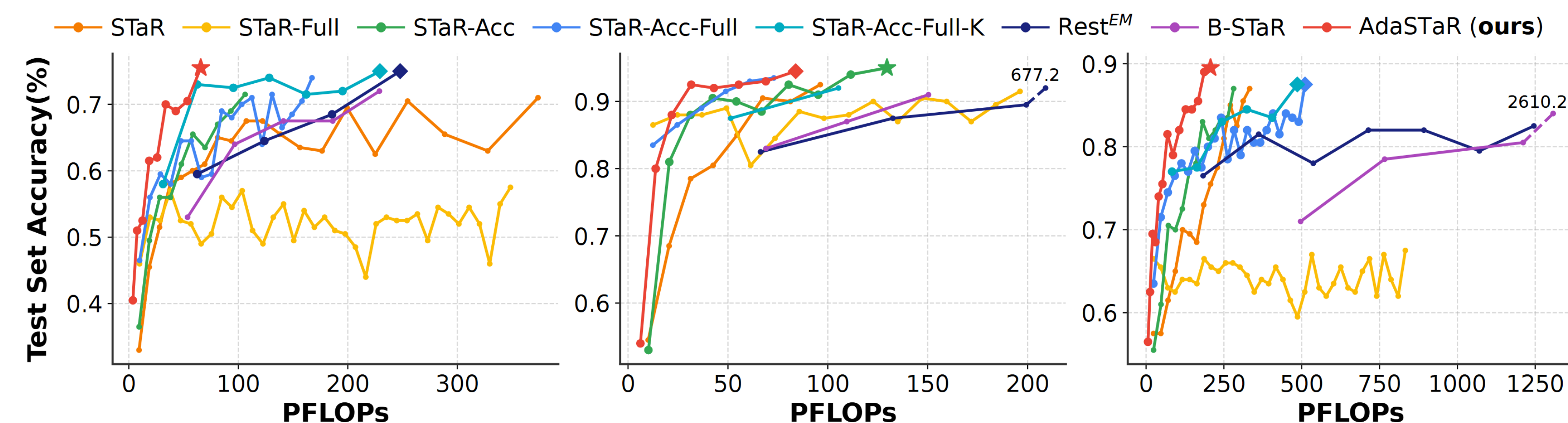
1. AdaSTaR solves the **under-training** and **over-training** of examples (data points) during rejection fine-tuning (RFT) self-improvement post-training.
2. This is achieved using a **compute-overhead free** sampling algorithm.
3. This algorithm results in improved post-training **performance**, and **significantly reduces training FLOPs**.

Empirical Study

- Achieves best test accuracy in **6/6 benchmarks** and reduced training FLOPs by an average of **58.6%**.
- **Pre-trained Base Models:** Llama 3.2 3B, Qwen 2.5 3B, Gemma 7B.

Evaluation Metric	ARC-C			CQA			CLadder 1.5		
	Acc. (\uparrow)	t	PFL0Ps (\downarrow)	Acc. (\uparrow)	t	PFL0Ps (\downarrow)	Acc. (\uparrow)	t	PFL0Ps (\downarrow)
SFT	61.4	1.0 e	7.0	71.8	1.0 e	24.0	31.0	7.0 e	382.3
SFT + 8-CoT	59.0	1.5 e	10.5	71.6	2.5 e	60.1	43.6	3.0 e	163.9
SFT + 5-SC	63.8	4.5 e	31.6	76.4	2.5 e	60.1	45.2	8.0 e	437.0
STaR	71.6	13 it	351.4	72.2	25 it	2877.8	53.4	25 it	8427.3
STaR-Full	69.8	27 it	739.4	72.2	12 it	1502.7	53.8	19 it	6523.7
STaR-Acc	73.2	18 it	639.8	74.6	19 it	1745.3	94.2	28 it	9663.0
STaR-Acc-Full	71.8	5 it	135.8	76.0	10 it	1158.3	94.2	15 it	4465.4
STaR-Acc-Full-K	71.4	3 it	302.2	73.0	4 it	1760.9	80.0	6 it	6382.3
ReST ^{EM}	70.8	4 it	637.1	72.8	2 it	1548.4	53.4	5 it	10498.3
B-STaR	67.8	2 it	222.8	68.4	2 it	800.9	52.8	4 it	3937.3
AdaSTaR (ours)	73.8	10 it	174.4 (\downarrow 72.7%)	78.0	20 it	779.3 (\downarrow 32.7%)	95.6	23 it	3610.0 (\downarrow 19.2%)

Evaluation Metric	ANLI			GSM8K			SVAMP		
	Acc. (\uparrow)	t	PFL0Ps (\downarrow)	Acc. (\uparrow)	t	PFL0Ps (\downarrow)	Acc. (\uparrow)	t	PFL0Ps (\downarrow)
SFT	64.2	4 e	262.9	61.0	2.5 e	177.3	57.0	5.5 e	21.7
SFT + 8-CoT	65.2	5 e	328.7	68.0	1 e	70.9	61.5	7.5 e	29.6
SFT + 5-SC	49.2	2 e	131.5	67.2	2.5 e	177.3	61.5	5.5 e	21.7
STaR	61.0	23 it	4195.3	76.0	4 it	409.2	71.0	20 it	373.8
STaR-Full	57.6	13 it	2604.6	72.6	4 it	684.8	57.5	37 it	348.5
STaR-Acc	64.8	22 it	3528.4	77.0	3 it	305.2	71.5	10 it	106.2
STaR-Acc-Full	64.6	5 it	986.0	74.6	2 it	333.0	74.0	18 it	167.3
STaR-Acc-Full-K	58.8	4 it	2528.4	77.0	2 it	1456.5	75.0	7 it	229.3
ReST ^{EM}	63.0	9 it	10938.5	77.0	2 it	2229.1	75.0	4 it	247.8
B-STaR	59.4	10 it	6373.4	73.6	3 it	2120.2	72.0	5 it	228.9
AdaSTaR (ours)	66.8	21 it	1340.9 (\downarrow 62.0%)	77.0	2 it	19.3 (\downarrow 93.7%)	75.5	9 it	65.7 (\downarrow 71.3%)



- Star ★ and diamond ♦ represents the **best** and **second best** accuracy.
- Training curve visualized up to the peak performance; afterwards it overfits.
- Dashed line represents methods that use large FLOPs that lead outside of the figure; numeric annotation denotes FLOPs at peak performance.