



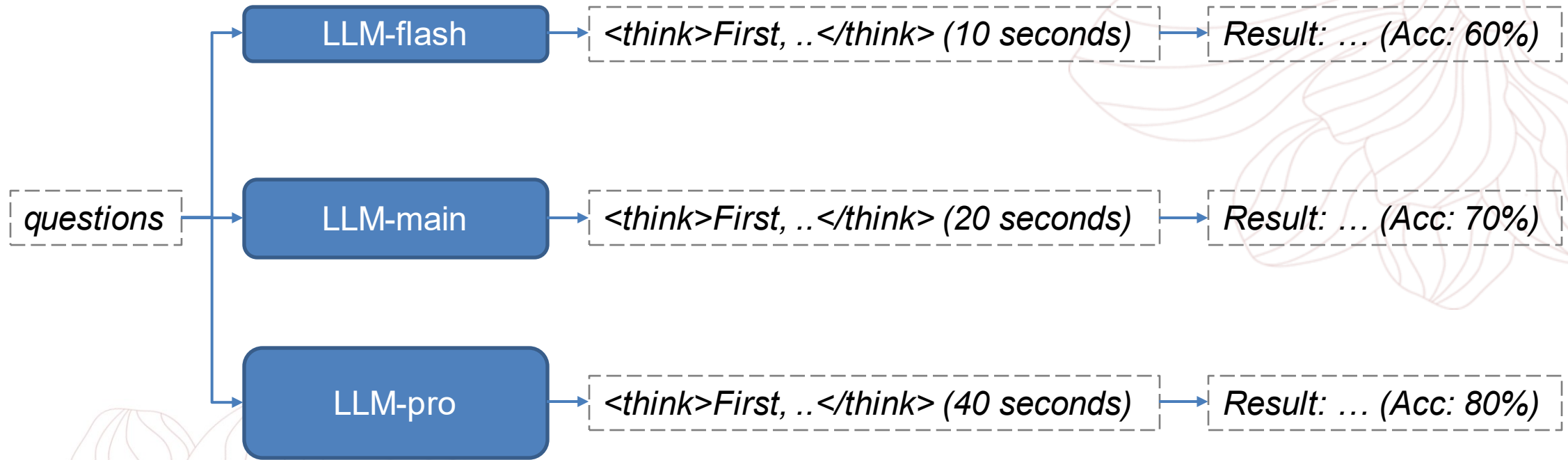
Think Silently, Think Fast: Dynamic Latent Compression of LLM Reasoning Chains

NeurIPS 2025

CoLaR-latent-reasoning.github.io

Wenhui Tan, Jiaze Li, Jianzhong Ju, Zhenbo Luo, Ruihua Song, Jian Luan

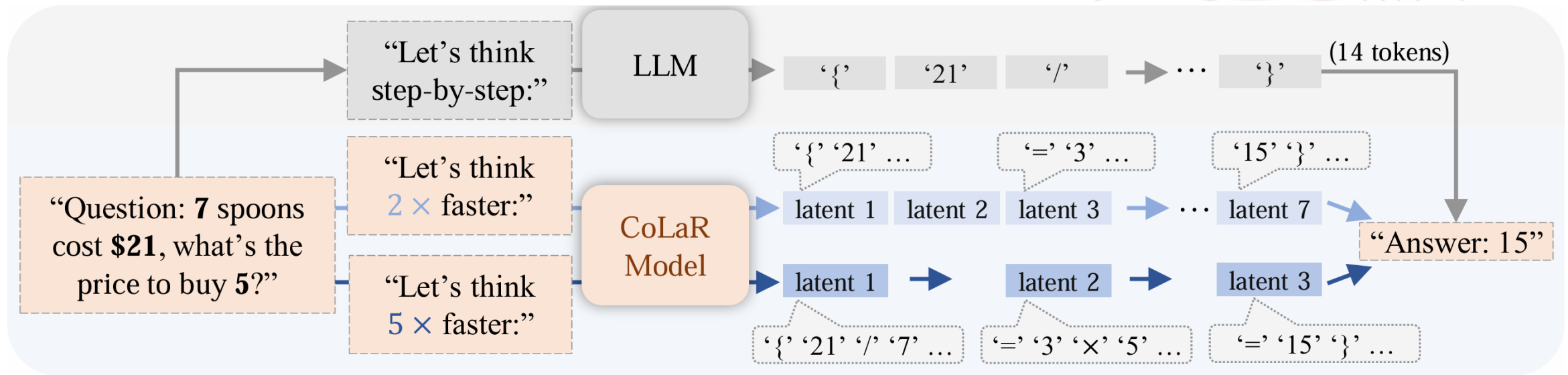
Background: LLM reasoning



Can we use **one omnivorous model** rather than separated models?

Word tokens are too expensive!
Must LLMs think in **textual space**?

Overview: Let's think *silently* and *fast*!



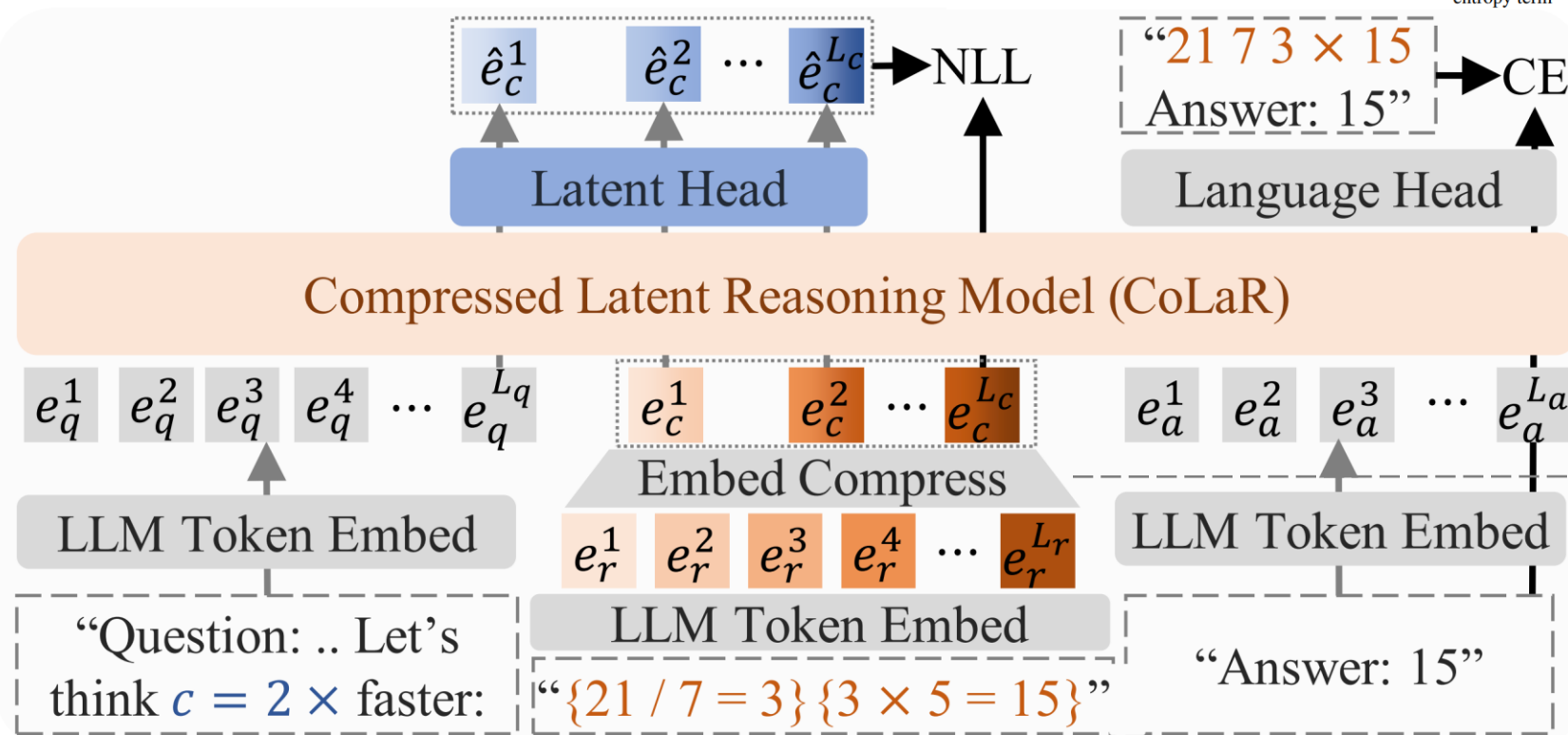
- **Think latent-by-latent**, where **one latent** compresses semantics from **multiple word tokens**
- **Dynamic** and controllable **compression factor** by prompting the thinking speed

Method: SFT training

- The model should **compress & predict & understand** latents.

$$\mathcal{L}_{\text{latent}}(i) = -\log p(e_c^i | \hat{\mu}_c^i, \hat{\sigma}_c^i) = \frac{(e_c^i - \hat{\mu}_c^i)^2}{2\hat{\sigma}_c^i} + \log \hat{\sigma}_c^i$$

$$\mathcal{L}_{\text{latent}}(i) = \underbrace{\mathbb{E}_\epsilon [(\hat{\mu}_c^i + \hat{\sigma}_c^i \epsilon - e_c^i)^2]}_{\text{MSE term}} - \alpha \underbrace{\left(\frac{1}{2} \log(2\pi e (\hat{\sigma}_c^i)^2) \right)}_{\text{entropy term}}$$

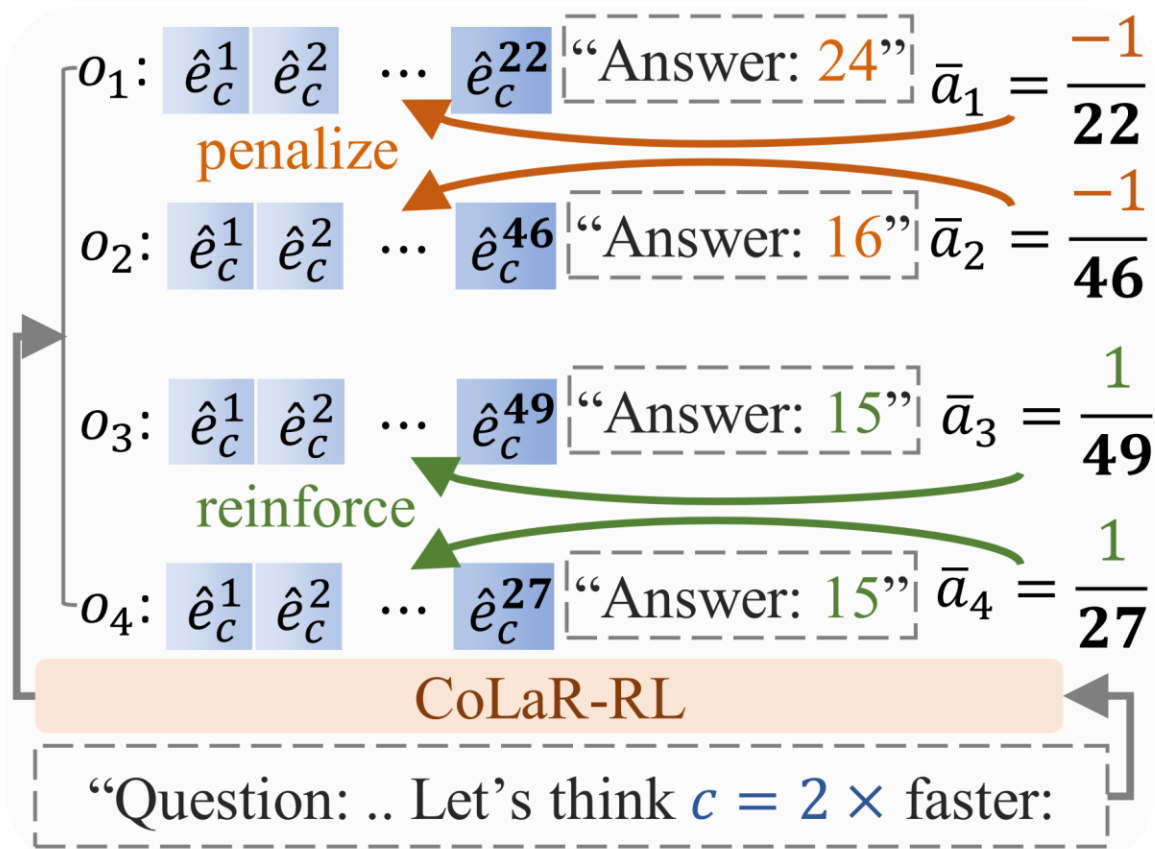


$$\mathcal{L}_{\text{comp}} = -\frac{1}{L_a + L_c} \sum_{i=1}^{L_a + L_c} \log p([t_c, t_a]^i | [e_c, e_a]^{1:i-1}, e_q),$$

- $e \sim \mathcal{N}(0, \sigma_e)$
- Random two embeddings could be highly uncorrelated (high dimensionality)
- Mean Pooling -> Add & divide by \sqrt{c}

Method: RL training

- Encourage model to **explore correct** reasoning pathways, and **exploit shorter** ones.



$$\mathcal{L}_{\text{GRPO}} = -\frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{\text{old}}}(o_i|q)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{\text{old}}}(o_i|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) \right)$$

$$A_i = \frac{r_i - \text{mean}(r_1, r_2, \dots, r_G)}{\text{std}(r_1, r_2, \dots, r_G)}$$

Experiments: CoLaR vs. SOTA

Table 1: Experiment results of baseline methods and CoLaR on four grade-school math reasoning datasets. We test the methods for five times with different random seeds to report the averaged number and 95% confidence interval (\pm) on accuracy (Acc. %) and reasoning chain length (# L). CoLaR- c denotes a same CoLaR model tested with different compression factors c . For ablation methods (marked in gray), suffixes DL, OC, MP and NLL denote CoLaR with a Deterministic Latent head, training withOut Compressed reasoning chain in cross entropy labels, using Mean Pooling to compress embeddings, and training with NLL loss, respectively.

	GSM8k-Aug		GSM-Hard		SVAMP		MultiArith		Average	
	Acc.	# L	Acc.	# L	Acc.	# L	Acc.	# L	Acc.	# L
CoT	49.4 \pm .72	25.6 \pm .11	11.9 \pm .16	34.2 \pm .11	59.8 \pm .29	12.1 \pm .03	93.2 \pm .49	13.7 \pm .09	53.6	21.4
iCoT	19.8 \pm .23	0.00 \pm .00	3.87 \pm .16	0.00 \pm .00	36.4 \pm .51	0.00 \pm .00	38.2 \pm .66	0.00 \pm .00	24.6	0.00
Coconut	23.1 \pm .28	6.00 \pm .00	5.49 \pm .33	6.00 \pm .00	40.7 \pm .65	6.00 \pm .00	41.1 \pm .24	6.00 \pm .00	27.6	6.00
Distill	13.3 \pm .62	6.00 \pm .00	2.97 \pm .24	6.00 \pm .00	21.7 \pm .73	6.00 \pm .00	19.2 \pm .83	6.00 \pm .00	14.3	6.00
CoLaR-5	26.8 \pm .17	5.57 \pm .02	5.87 \pm .10	6.53 \pm .01	48.4 \pm .45	2.95 \pm .02	86.4 \pm .35	3.21 \pm .01	41.7	4.57
- DL	26.7 \pm .11	5.74 \pm .01	5.53 \pm .11	8.20 \pm .04	48.3 \pm .05	2.90 \pm .01	84.5 \pm .19	3.22 \pm .01	41.3	5.02
- OC	24.8 \pm .27	5.14 \pm .12	6.46 \pm .11	5.49 \pm .06	46.5 \pm .18	2.85 \pm .01	85.9 \pm .22	3.13 \pm .01	40.1	4.15
- MP	20.6 \pm .22	5.61 \pm .02	4.20 \pm .07	6.18 \pm .02	47.7 \pm .41	2.96 \pm .01	80.7 \pm .59	3.20 \pm .01	38.3	4.49
- NLL	20.3 \pm .64	5.99 \pm .06	4.52 \pm .39	16.6 \pm .25	43.9 \pm .43	3.06 \pm .03	81.6 \pm .23	3.20 \pm .02	37.6	8.01
CoLaR-2	40.1 \pm .20	12.7 \pm .02	9.08 \pm .03	14.0 \pm .07	54.9 \pm .20	6.11 \pm .01	91.3 \pm .12	7.35 \pm .01	48.8	10.0
- DL	39.7 \pm .18	12.8 \pm .01	8.84 \pm .06	17.2 \pm .09	54.3 \pm .23	6.10 \pm .01	90.1 \pm .17	7.46 \pm .01	48.2	10.9
- OC	39.1 \pm .33	12.3 \pm .04	8.96 \pm .01	16.9 \pm .13	54.7 \pm .18	6.08 \pm .02	90.1 \pm .25	7.36 \pm .01	48.2	10.6
- MP	36.9 \pm .30	12.4 \pm .02	8.46 \pm .19	12.0 \pm .05	54.1 \pm .42	6.14 \pm .01	86.8 \pm .20	7.43 \pm .01	46.6	9.49
- NLL	32.3 \pm .51	12.2 \pm .04	7.57 \pm .16	16.6 \pm .25	51.0 \pm .24	5.50 \pm .03	88.3 \pm .41	7.09 \pm .02	44.8	10.3

- **14.1% \uparrow** performance compared to Coconut with shorter reasoning chains
- Reduces reasoning chain length by **53.3% \downarrow** with only a **4.8% \downarrow** performance degradation compared to CoT

Table 2: Experimental results on the challenging MATH dataset. We evaluate our proposed method CoLaR on two base models and three settings: -DL denotes using a Deterministic Latent head, -NLL denotes CoLaR trained with NLL Loss as $\mathcal{L}_{\text{latent}}$, which is our main method, and -/w GRPO denotes the post-trained CoLaR-NLL with GRPO reinforcement learning process. We calculate the performance gain between CoLaR-NLL and CoLaR-NLL-RL to highlight the effectiveness of reinforcement learning. Compression factor c and $\# L_{\text{max}}$ are set to 2 and 128, respectively.

	DeepSeek-R1-Distill-Qwen-1.5B		Llama-3.2-1B-Instruct	
	Acc.	# L	Acc.	# L
CoT	23.5 \pm .29	209 \pm 1.6	9.71 \pm .33	210 \pm 1.4
CoLaR-DL	9.04 \pm .12	99.4 \pm .25	3.07 \pm .28	134 \pm .46
CoLaR-NLL	8.94 \pm .21	56.8 \pm .14	5.28 \pm .16	83.1 \pm .52
CoLaR-NLL-RL	14.3 \pm .25 (5.36% \uparrow)	9.79 \pm .40 (82.8% \downarrow)	7.08 \pm .07 (1.80% \uparrow)	16.1 \pm .14 (80.6% \downarrow)
- w/o average	13.8 \pm .14	128 \pm .00	0.00 \pm .00	128.0 \pm .00

- **5.36% \uparrow** accuracy while reducing the length of reasoning chain significantly by **82.8% \downarrow**

Experiments: Case study

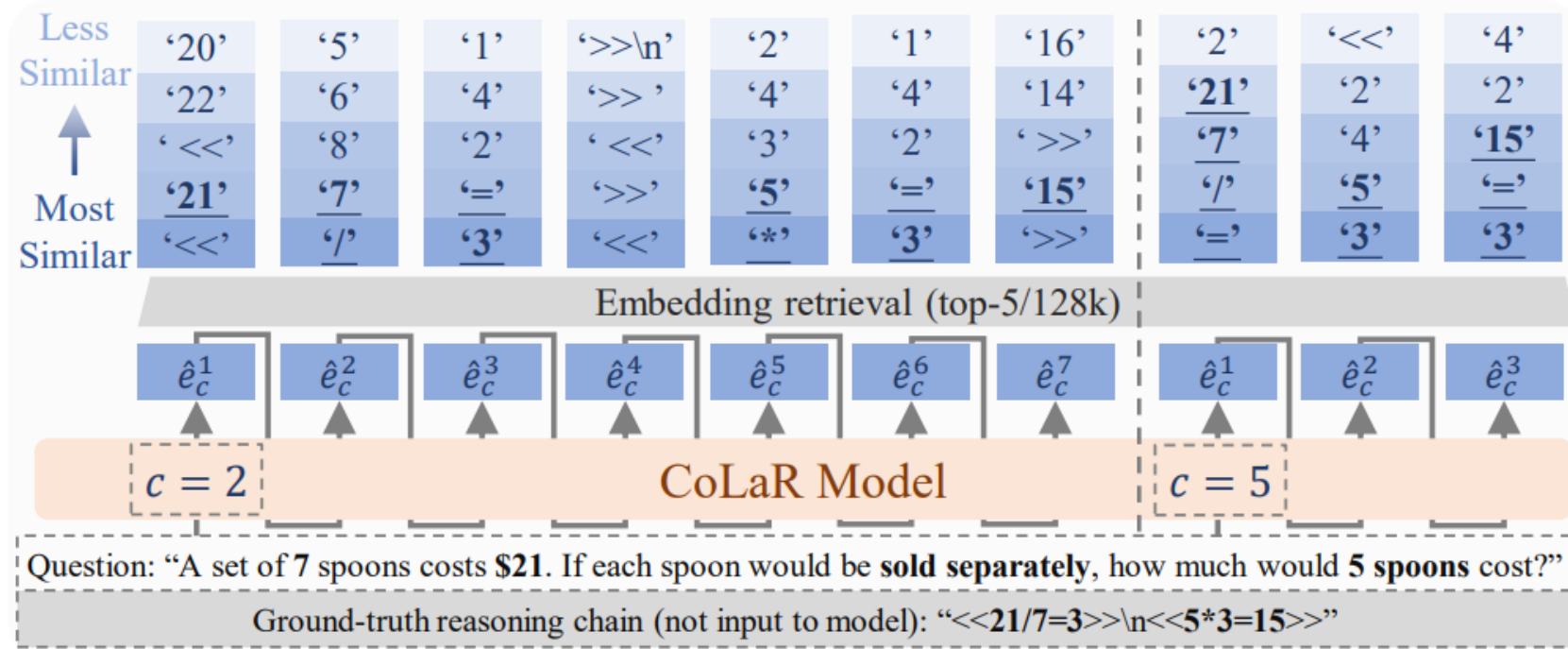


Figure 5: A case study on the GSM8k validation set. We set the compression factor c to 2 and 5, which produce two latent reasoning chains in length 7 and 3, respectively. We then retrieve tokens with the predicted latents by embedding cosine similarity, and underscore those informative tokens.

Higher compression factor captures more tokens while ignoring less informative tokens (like "<<")

Experiments: Analyses on compression factor c

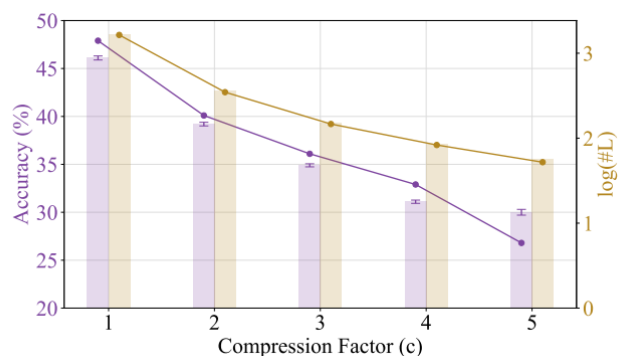


Figure 3: Accuracy and reasoning chain length (# L) of CoLaR on GSM8k dataset when trained with random $c \in [1, 5]$ (the **lines**) or trained solely on specific c (the **bars**).

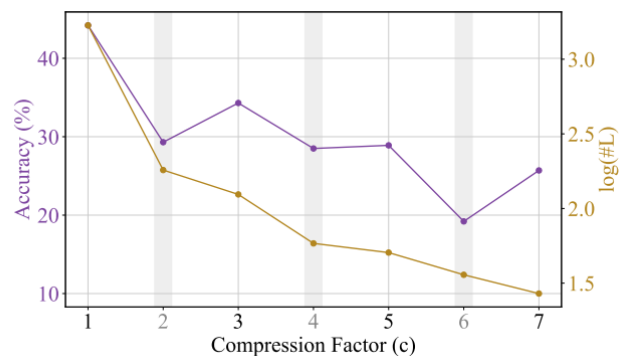
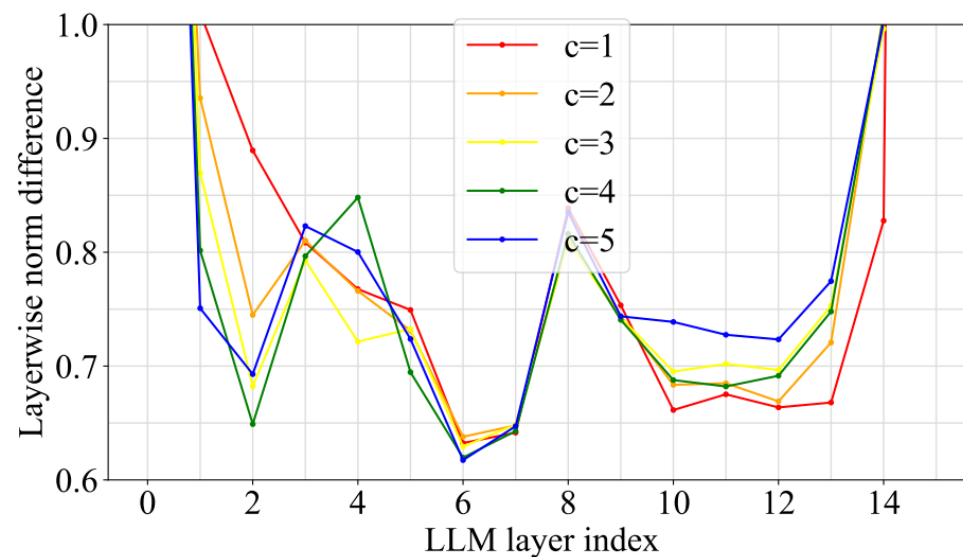
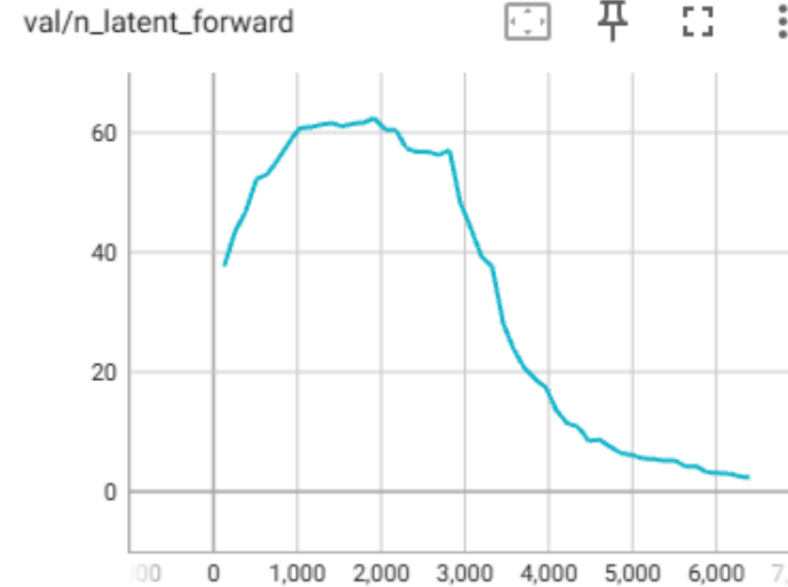
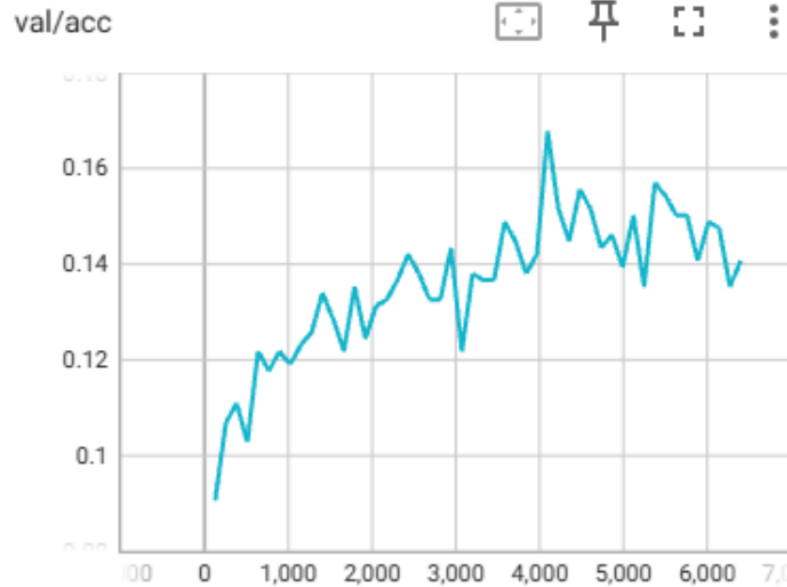


Figure 4: Accuracy and reasoning chain length (# L) of CoLaR on GSM8k dataset when trained with $c \in \{1, 3, 5, 7\}$ and tested with extra $c \in \{2, 4, 6\}$ (under gray bars).



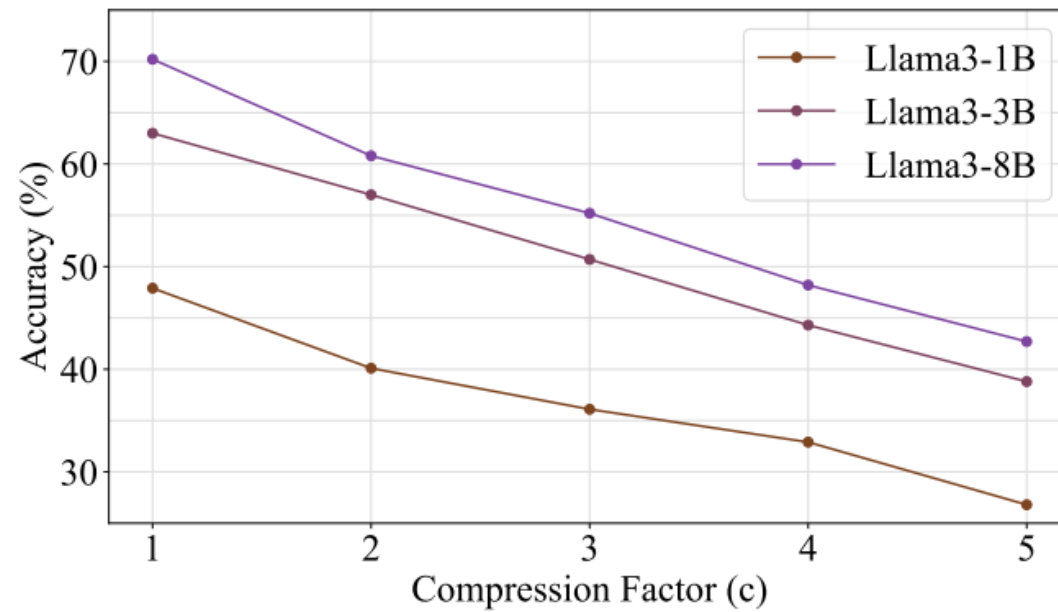
Experiments: RL training curve analyses



Three stage:

- Exploration: acc and reasoning chain length both increases
- Exploitation: acc fluctuates while reasoning chain length decreases
- Overfitting and early-stopping

Experiments: Model size scaling



Conclusion

Main contributions:

- Novel framework: Compressed latent reasoning with controllable test-time compression factors
- Training pipeline: First work demonstrating the effectiveness of reinforcement learning on latent reasoning

Limitations:

- No significant performance gain applying RL on **simple** math reasoning datasets
- Not surpassing explicit CoT method on Acc.

Future work:

- Multimodality domain
- Adaptive test-time compression factor