

# SoftCoT: Soft Chain-of-Thought for Efficient Reasoning with LLMs

Yige Xu, Xu Guo, Zhiwei Zeng, Chunyan Miao

Nanyang Technological University, Singapore

# Research Background

---

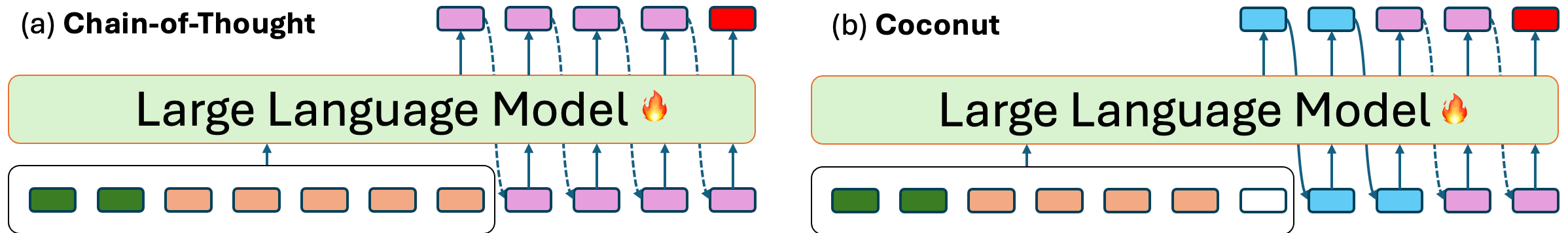
- Chain-of-Thought reasoning has become one of the basic ability of LLMs.
- Three primary concerns:
  - Consistency and Stability: CoT can vary significantly with minor changes in prompts. [1,2]
  - Robustness: CoT's effectiveness depends on the quality of intermediate thoughts. [3]
  - Efficiency: CoT often requires substantial computational resources. [4]

SoftCoT



# Continuous Space Reasoning

- Generate soft thought tokens according to the hidden of last-token last-layer
- Facilitates the reasoning chain generation
- Optimal latent-space exploration
  - Coconut [3], CCoT [5]



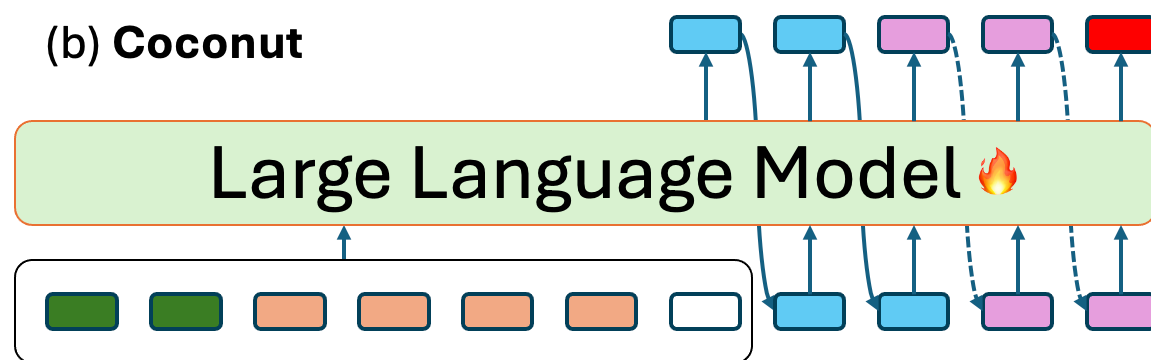
# Motivation

Current latent-space reasoning approaches consider latent-space reasoning as a new task and fine-tunes the whole LLM [3,5], which results in ...

- Catastrophic forgetting problem on SOTA LLMs
- Auto-regressively generate the soft thought tokens

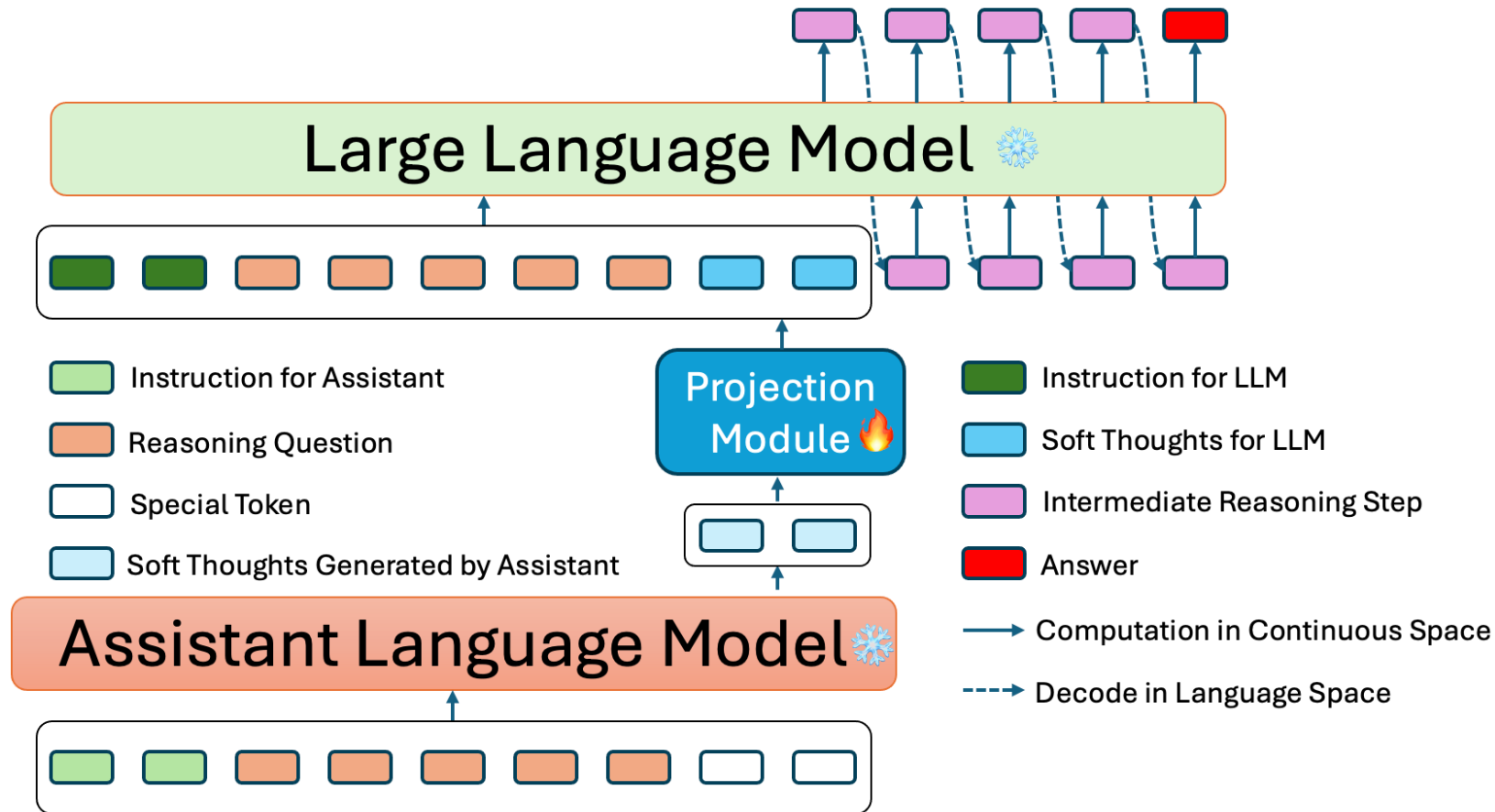
Can we *freeze the LLM* for mitigating the catastrophic forgetting problem?

Challenge: the fixed LLM struggle to generate learnable soft thought tokens.



How to generate the learnable soft thought tokens?

# SoftCoT: Overall Architecture



Background

Motivation

Methodology

Results

Analysis

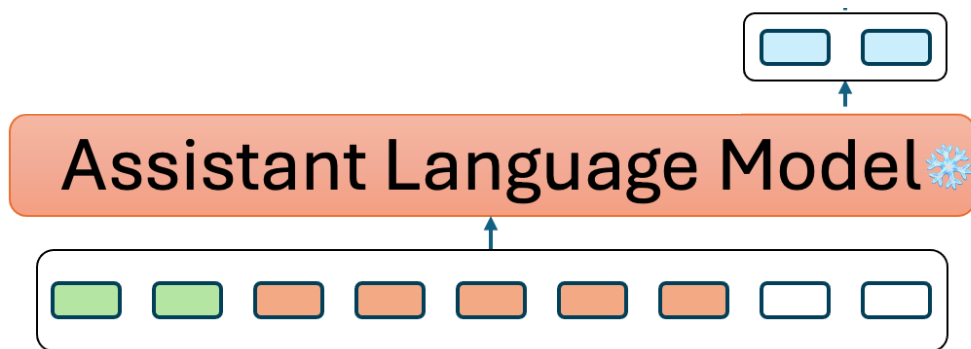
# SoftCoT: Soft Thought Tokens Generataion

- Soft Thought Tokens Generation
  - Use auxiliary assistant model to produce the soft thoughts

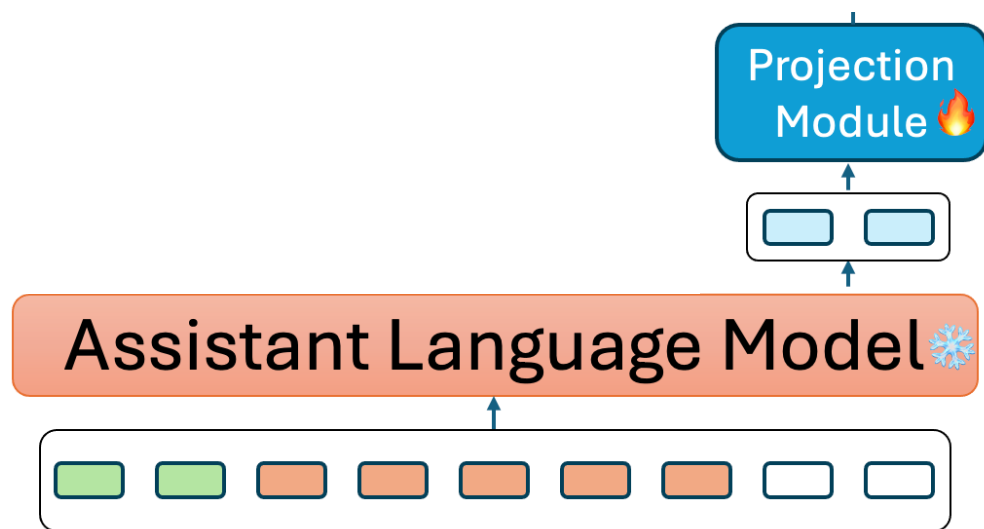
$$\mathbf{x}_{\text{assist}} = \text{concat}[\mathcal{I}_{\text{assist}}, \mathcal{Q}, [\text{UNK}]_{1:N}]$$

$$\mathbf{h}^{\text{assist}} = \text{Assistant}(\mathbf{x}_{\text{assist}}),$$

$$\mathbf{t}_{\text{assist}} = \mathbf{h}_{|\mathcal{I}|+|\mathcal{Q}|+1:|\mathcal{I}|+|\mathcal{Q}|+N}^{\text{assist}}$$



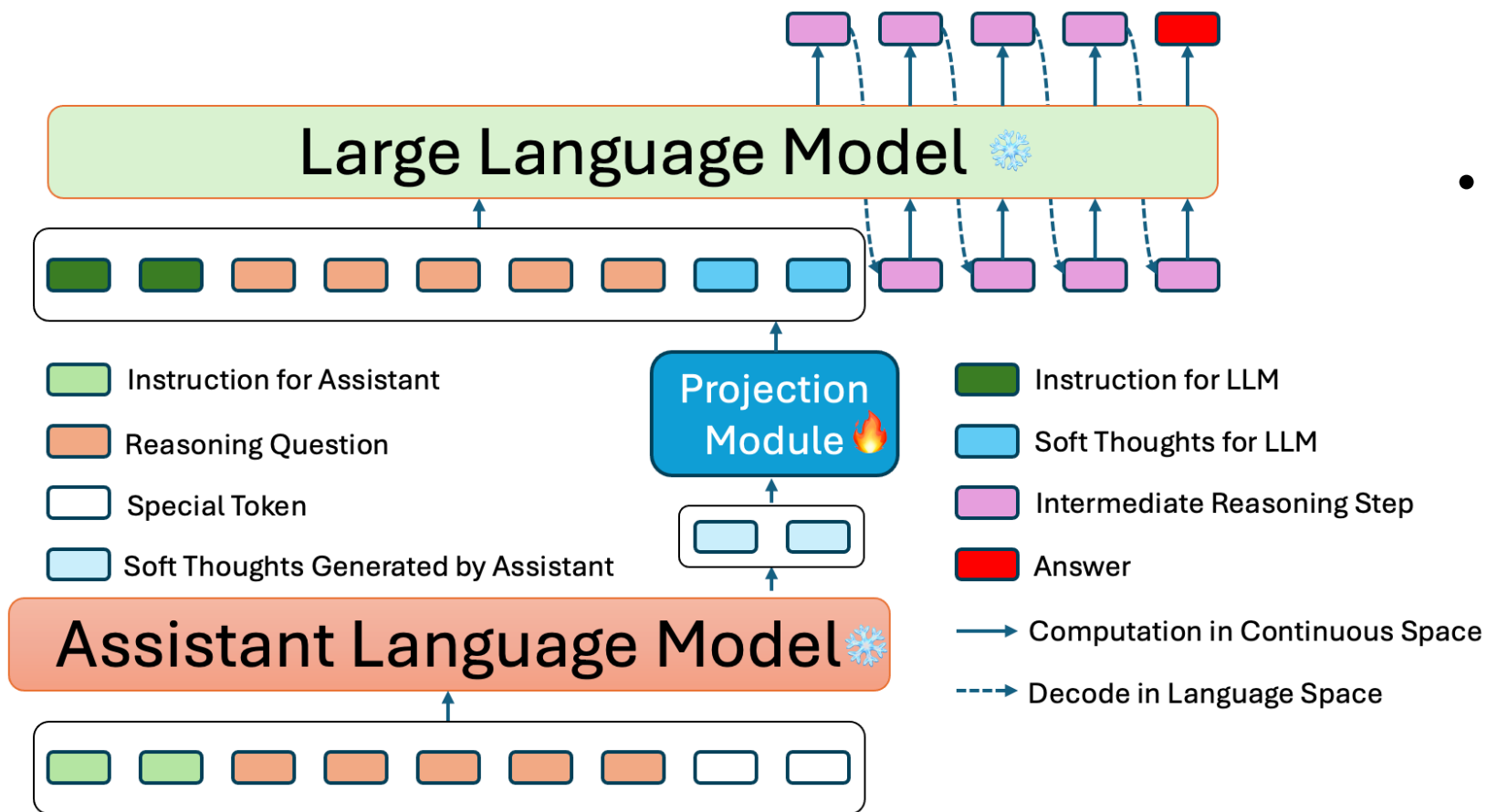
# SoftCoT: Soft Thought Tokens Projection



- Soft Thought Tokens Projection
  - Maps the assistant-generated soft thoughts from the assistant model's embedding space to the LLM's embedding space.
  - Only the parameters in the projection module are trainable.

$$\mathcal{T}_{\text{soft}} = \text{Linear}_{\theta}(\mathbf{t}_{\text{assist}}),$$

# SoftCoT: LLM Reasoning



- LLM Reasoning with SoftCoT
  - Apply the soft thoughts to aid LLMs in CoT reasonings.

$$\mathbf{x}_{\text{LLM}} = \text{concat}[\mathcal{I}_{\text{LLM}}, \mathcal{Q}, \mathcal{T}_{\text{soft}}],$$

$$\bar{\mathcal{R}} = \text{LLM}(\mathbf{x}_{\text{LLM}}),$$

$$\bar{\mathcal{A}} = \text{LLM}(\mathbf{x}_{\text{LLM}}, \bar{\mathcal{R}}),$$

$$\hat{\mathcal{A}} = \mathcal{E}(\bar{\mathcal{A}}),$$

Background

Motivation

Methodology

Results

Analysis



# Comparison with baselines

Model	GSM8K	ASDiv-Aug	AQuA	StrategyQA	DU	Avg.
	Mathematical			Commonsense	Symbolic	
<i>GPT-2</i>						
Coconut (Hao et al., 2024)	34.10 <sub>±1.50</sub> *	38.92 <sub>±0.00</sub> <sup>†</sup>	22.83 <sub>±0.00</sub> <sup>†</sup>	-	-	-
<i>LLaMA-3.1-8B-Instruct</i>						
Zero-Shot CoT	79.61 <sub>±0.81</sub>	86.78 <sub>±0.63</sub>	54.65 <sub>±2.43</sub>	65.63 <sub>±3.31</sub>	54.40 <sub>±2.40</sub>	68.21
Zero-Shot CoT-Unk	79.95 <sub>±0.59</sub>	86.90 <sub>±0.41</sub>	55.28 <sub>±1.88</sub>	66.16 <sub>±2.70</sub>	54.16 <sub>±1.46</sub>	68.49
Zero-Shot Assist-CoT	80.76 <sub>±1.53</sub>	86.96 <sub>±0.46</sub>	55.83 <sub>±2.98</sub>	66.55 <sub>±3.99</sub>	58.24 <sub>±3.56</sub>	69.67
LoRA Fine-Tuning	75.66 <sub>±0.00</sub>	86.67 <sub>±0.00</sub>	52.36 <sub>±0.00</sub>	-	-	-
Coconut (Hao et al., 2024) <sup>†</sup>	76.12 <sub>±0.00</sub>	86.80 <sub>±0.00</sub>	53.15 <sub>±0.00</sub>	-	-	-
<b>SoftCoT (Ours)</b>	<b>81.03<sub>±0.42</sub></b>	<b>87.19<sub>±0.40</sub></b>	<b>56.30<sub>±1.67</sub></b>	<b>69.04<sub>±1.23</sub></b>	<b>59.04<sub>±1.93</sub></b>	<b>70.52</b>

- Supervised LoRA Fine-Tuning performs worse than zero-shot CoT, which make Coconut not applicable to SOTA LLMs
- Assistant model is effective to facilitate CoT reasoning
- SoftCoT consistently benefits from the supervised training

Background

Motivation

Methodology

Results

Analysis

# Generalization to Other LLM Backbones

Model	GSM8K	ASDiv-Aug	AQuA	StrategyQA	DU	Avg.
	Mathematical			Commonsense	Symbolic	
Zero-Shot CoT	83.70 $\pm$ 0.78	87.19 $\pm$ 0.28	64.53 $\pm$ 3.27	49.65 $\pm$ 3.18	66.40 $\pm$ 2.26	70.29
Zero-Shot CoT-Unk	84.12 $\pm$ 0.71	86.94 $\pm$ 0.89	64.72 $\pm$ 2.06	50.74 $\pm$ 1.90	66.48 $\pm$ 1.43	70.60
Zero-Shot Assist-CoT	84.85 $\pm$ 0.35	88.63 $\pm$ 1.05	64.96 $\pm$ 2.83	52.71 $\pm$ 2.65	67.04 $\pm$ 2.84	71.64
LoRA Fine-Tuning	81.80 $\pm$ 0.00	86.80 $\pm$ 0.00	62.60 $\pm$ 0.00	-	-	-
Coconut (Hao et al., 2024)	82.49 $\pm$ 0.00	86.90 $\pm$ 0.00	63.39 $\pm$ 0.00	-	-	-
<b>SoftCoT (Ours)</b>	<b>85.81<math>\pm</math>1.82</b>	<b>88.90<math>\pm</math>1.01</b>	<b>72.44<math>\pm</math>2.19</b>	<b>60.61<math>\pm</math>1.55</b>	<b>67.52<math>\pm</math>2.92</b>	<b>75.06</b>

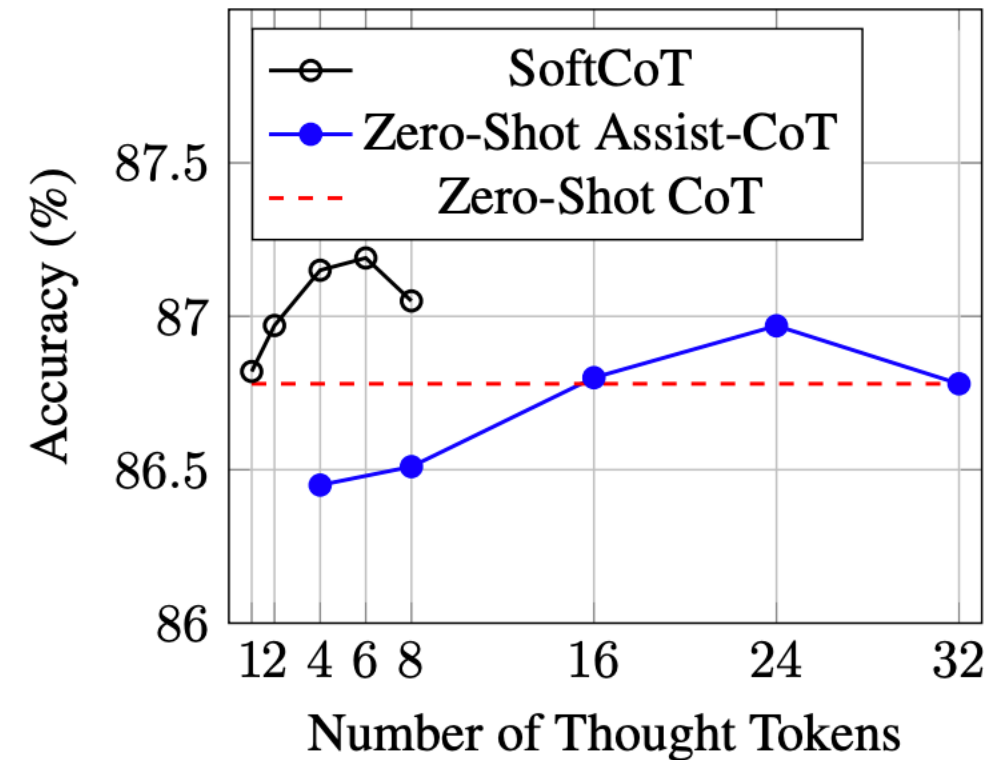
Results on Qwen2.5-7B-Instruct

- SoftCoT is effective across different LLM architectures



# Model Analysis – Number of Thought Tokens

- Soft thoughts reduce the required number of thought tokens



# Model Analysis – Size of Assistant Model

Method	0.5B	1.5B	7B
Zero-Shot CoT	83.70	83.70	83.70
Zero-Shot Assist-CoT	84.78	84.85	84.90
SoftCoT	<b>85.76</b>	<b>85.81</b>	<b>85.84</b>

Table 5: Performance on GSM8K with different sizes of assistant model on Qwen2.5 series.

- The scale of the assistant model has limited impact on the accuracy of the final answer

# Model Analysis – Self-Consistency

Model	GSM8K		ASDiv-Aug		AQuA		StrategyQA		DU	
	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$	$N = 1$	$N = 10$
Zero-Shot CoT	79.61 $\pm$ 0.81	90.36 $\pm$ 0.40	86.78 $\pm$ 0.63	89.23 $\pm$ 0.17	54.65 $\pm$ 2.43	63.23 $\pm$ 0.86	65.63 $\pm$ 3.31	70.13 $\pm$ 0.47	54.40 $\pm$ 2.40	65.76 $\pm$ 1.54
Zero-Shot Assist-CoT	80.76 $\pm$ 1.53	90.43 $\pm$ 0.69	86.96 $\pm$ 0.46	89.48 $\pm$ 0.36	55.83 $\pm$ 2.98	63.62 $\pm$ 0.99	66.55 $\pm$ 3.99	70.48 $\pm$ 0.68	58.24 $\pm$ 3.56	65.84 $\pm$ 1.93
<b>SoftCoT (Ours)</b>	<b>81.03<math>\pm</math>0.42</b>	<b>90.63<math>\pm</math>0.39</b>	<b>87.19<math>\pm</math>0.40</b>	<b>89.75<math>\pm</math>0.29</b>	<b>56.30<math>\pm</math>1.67</b>	<b>65.51<math>\pm</math>0.72</b>	<b>69.04<math>\pm</math>1.23</b>	<b>71.14<math>\pm</math>0.10</b>	<b>59.04<math>\pm</math>1.93</b>	<b>67.36<math>\pm</math>1.12</b>

Table 4: Self Consistency for SoftCoT on LLaMA-3.1-8B-Instruct. “ $N$ ” indicates the number of reasoning chains.

- SoftCoT introduces an independent improvement mechanism, which can be effectively combined with self-consistency for enhanced reasoning performance



# Takeaway messages

---

- We address the need for efficient CoT reasoning on continuous space within SOTA LLMs
  - Freezing the backbone LLM to mitigate the catastrophic forgetting problem.
  - Creating a learnable projection module to map the assistant-generated soft thoughts from the assistant model's embedding space to the LLM's embedding space.
- SoftCoT has demonstrated that
  - it enables reasoning on continuous space and has a better downstream performance than baselines.
  - it can be scaled to multiple LLM architectures
  - it can be scaled to existing test-time scaling methods such as self-consistency.

# References

---

- [1] Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, Yoav Goldberg. ***Measuring and Improving Consistency in Pretrained Language Models***. **TACL 2021**.
- [2] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, Denny Zhou. ***Self-Consistency Improves Chain of Thought Reasoning in Language Models***. **ICLR 2023**.
- [3] Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, Yuandong Tian. ***Training Large Language Models to Reason in a Continuous Latent Space***. **arXiv preprint: 2412.06769**.
- [4] Zhenglin Wang, Jialong Wu, Yilong Lai, Congzhi Zhang, Deyu Zhou. ***SEED: Accelerating Reasoning Tree Construction via Scheduled Speculative Decoding***. **COLING 2025**.
- [5] Jeffrey Cheng, Benjamin Van Durme. ***Compressed Chain of Thought: Efficient Reasoning Through Dense Representations***. **arXiv preprint: 2412.13171**.