

Implicit Thinking

Reasoning in the Latent/Continuous Space

NAVER LABS Europe NLP Seminar

March 17, 2026

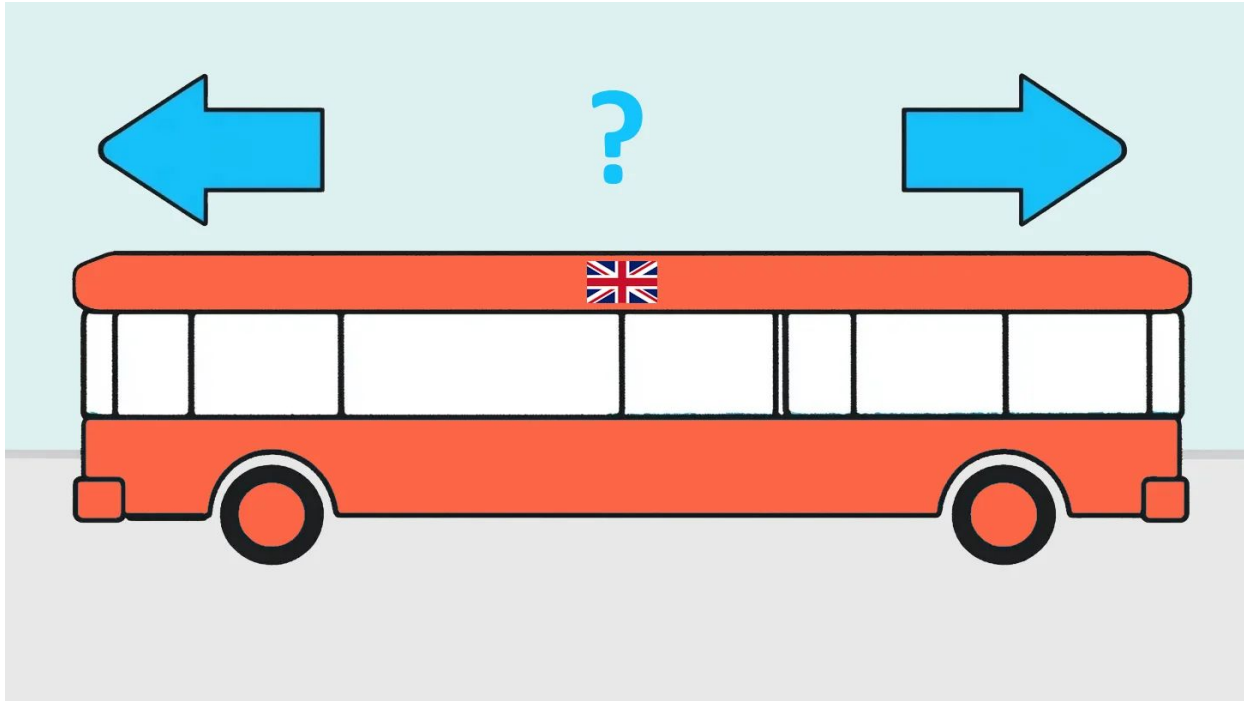
Thomas Palmeira Ferraz

Pierre Erbacher

Can you solve this problem?

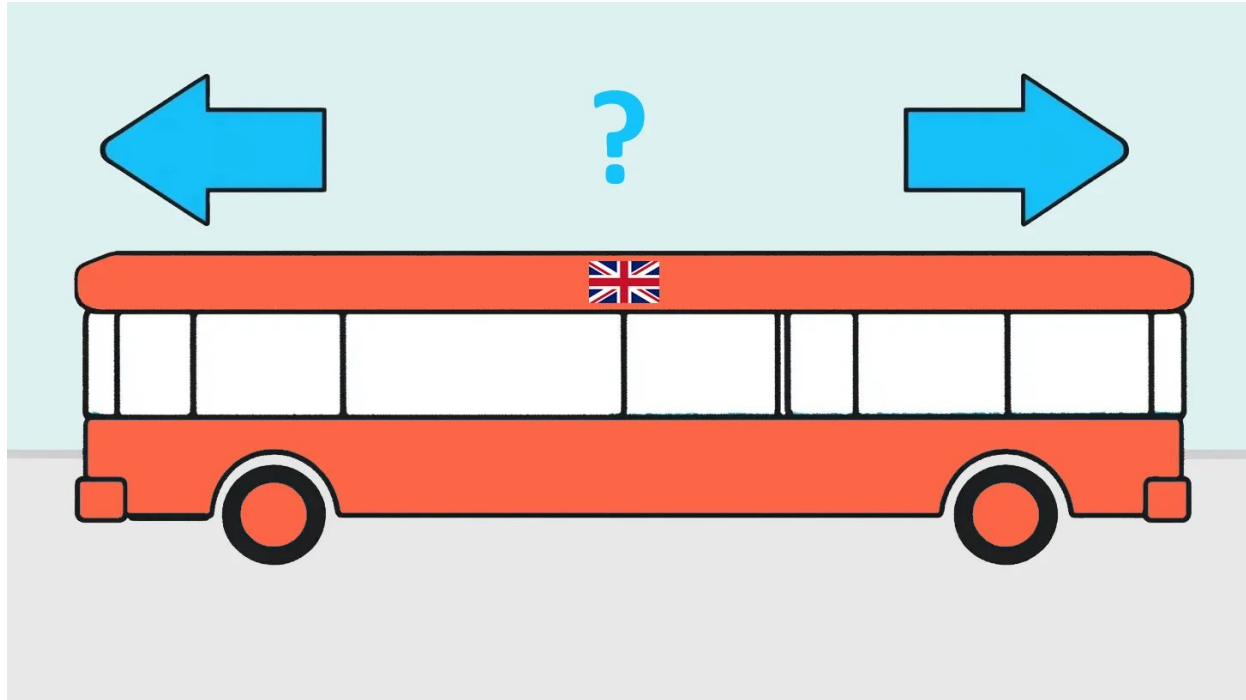
Can you solve this problem?

Which direction is the bus moving? Left or right?



Can you solve this problem?

Which direction is the bus moving? Left or right?



RIGHT

How LLMs Would Solve This?

Model	# Words to Solve the Riddle	Correct
Kimi K-2.5 (1T)	7,143	Yes
Qwen 3.5-VL-27B	2,757	No
Qwen 3.5-VL-9B	1,951	No
Qwen 3.5-35B-A3B	3,897	Yes

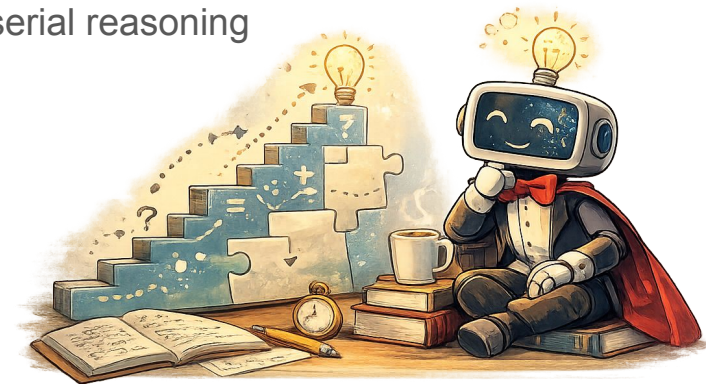
Agenda

- Why Latent Reasoning?
- L-SEQ: Latent Sequence Reasoning
- L-LOOP: Latent Looped Reasoning
- L-OPT: Latent Optimization
- Alternative Architectures
 - Large Concept Models
 - Diffusion Models
- Multimodality
- Evidences, Challenges, Limitations
- Open problems and future directions

Preliminaries

Explicit textual CoT: Breakthrough for TTS

- Test-time scaling
 - Human cognitive model: System 1 (autopilot) vs. System 2 (thinking)
 - System 2 LLMs improve performance by generating intermediate results.
- Explicit textual CoT was a breakthrough for TTS
 - Better performance on many multi-step tasks with *“Let’s think step-by-step”*
 - Supervision made straightforward: train on intermediate traces, not just answers
 - A human-readable interface for debugging and steering (interventions / RLVR)
 - A partially observable “working memory” channel for serial reasoning



Language is **the** bottleneck for reasoning progress

- Comparison with most efficient case: humans
 - Human cognition often transcends discrete linguistic symbols, involving abstract, continuous, multimodal or multi-conceptual representations that resist precise verbalization (Wittgenstein, 1922)
 - Complex emotions like **nostalgia** and **bittersweetness** are continuous blends of **joy**, **sadness**, and **multimodal** memory that defy expressions from a fixed vocabulary (Pinker, 1994)

Language is **the** bottleneck for reasoning progress

- Comparison with most efficient case: humans
 - Human cognition often transcends discrete linguistic symbols, involving abstract, continuous, multimodal or multi-conceptual representations that resist precise verbalization (Wittgenstein, 1922)
 - Complex emotions like **nostalgia** and **bittersweetness** are continuous blends of **joy**, **sadness**, and **multimodal** memory that defy expressions from a fixed vocabulary (Pinker, 1994)
- Consequences of reasoning on discrete textual space
 - **Expressive redundancy**: many reasoning tokens are syntactically needed but functionally weak (“so”, “the”, etc.).
 - **Stylistic overfitting**: models can learn artifacts of reasoning traces rather than reasoning itself.
 - **Semantic bottleneck**: forcing continuous or multi-concept states into a discrete linear chain of fixed vocabulary inevitably leads to information loss.
 - **Serialization bottleneck**:
 - reasoning must unfold one token at a time.
 - Verbose, expensive and limited parallel exploration.

The case for *de-linguistifying* reasoning

Promises:

- Less token-level computation → faster inference
- Richer, more compact reasoning representations
- Potential to explore multiple reasoning trajectories in parallel

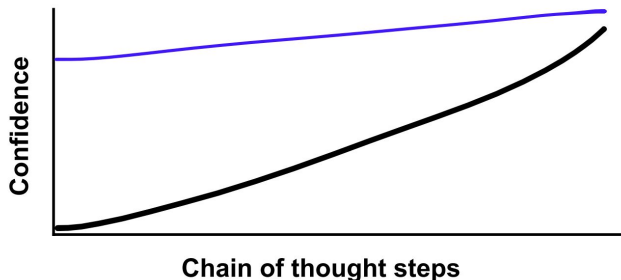
Challenges:

- Harder direct supervision
- Training and alignment difficulty
- Lower transparency / weaker debuggability
- Evaluation gap: did the model really reason, or just fit input–output patterns?

CoT and Performative Reasoning

Faithfulness

- **Unfaithful explanation:** CoT can systematically misrepresent LLM internal reason for a its prediction [1]
- **LLMs can fake reasoning:** producing coherent-looking rationales for contradictory / bias-driven answers [1,2]
- **(≠ Influence):** CoT can affect final answer even if not faithful to model's actual reasoning [3]
- **Fragile Monitorability:** CoT can help oversight, but it is imperfect and may fail to reveal important internal states [4]



Source: Lewis-Lim et al. [3]

[1] Turpin et al., [Language Models Don't Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting](#) (NeurIPS, 2023)

[2] Arcuschin et al., [Chain-of-Thought Reasoning in the Wild is not Always Faithful](#) (2025)

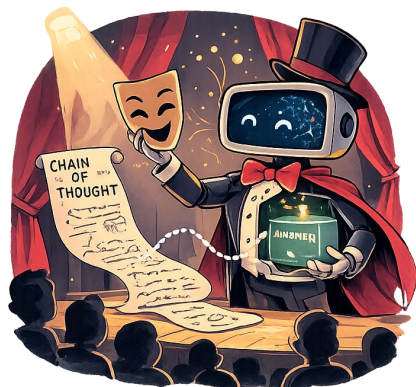
[3] Lewis-Lim et al., [Analysing Chain of Thought Dynamics: Active Guidance or Unfaithful Post-hoc Rationalisation?](#) (EMNLP 2025)

[4] Korbak et al., [Chain of Thought Monitorability: A New and Fragile Opportunity for AI Safety](#) (2025)

CoT and Performative Reasoning

Performative (or decorative) reasoning:

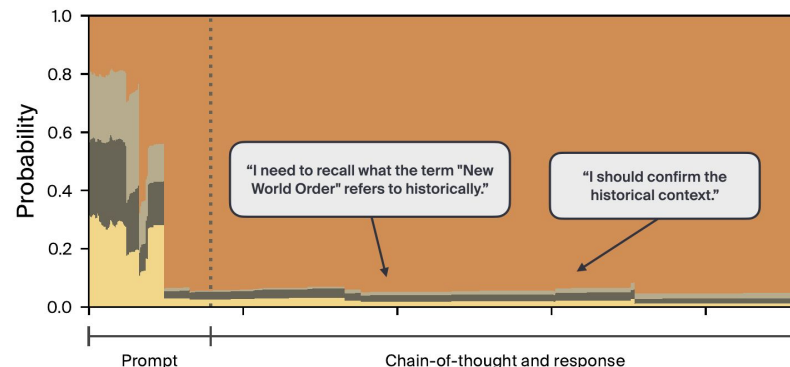
- **Early decisions:** in some cases, the final answer becomes decodable well before the full CoT finishes [1]
- **Causality:** many CoT steps have little causal impact on the final prediction, even “aha” or self-verification steps can be decorative [2]



Per-token probe prediction of model's final answer

What was meant by the term 'New World Order'?

- (A) A new democratic internationalism led by the United States
- (B) A new balance of power between the US and China
- (C) A new global economic framework
- (D) A new era of globalization



Source: Boppana et al. [1]

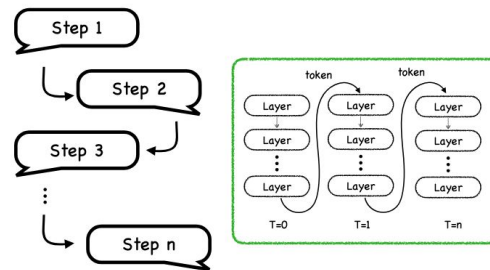
[1] Boppana et al., [Reasoning Theater: Disentangling Model Beliefs from Chain-of-Thought](#) (2026)

[2] Zhao et al., [Can Aha Moments Be Fake? Identifying True and Decorative Thinking Steps in Chain-of-Thought](#) (2025)

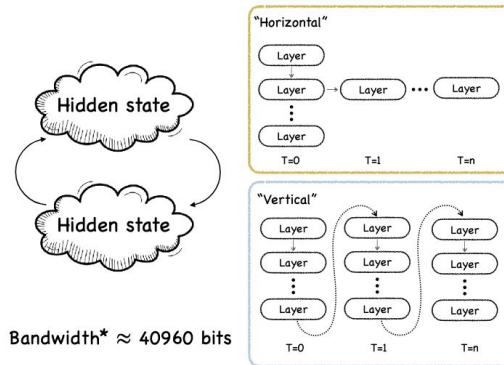
LLM Reasoning

- LLM Reasoning can be modeled as a **two-stage process**:
 - from an input, the model builds an intermediate trace, then produces the answer.
- **Explicit reasoning**: trace is verbalized as text.
- **Implicit reasoning**: trace is not verbalized or discretized.
 - **Latent reasoning**: internal reasoning explicitly in latent/hidden-state form.
- Key distinction is not whether there is a trace, but whether it is exposed to the user.

*For FP-16: Explicit reasoning transmits discrete tokens (≈ 15 bits / step), whereas latent reasoning exchanges 2560-dimensional FP-16 hidden states ($\approx 40,960$ bits / step), revealing a $\sim 2.7 \times 10^3$ -fold bandwidth gap between the two approaches. Source: Zhu et al. [A Survey on Latent Reasoning](#) (July, 2025)



Bandwidth* ≈ 15 bits



Bandwidth* ≈ 40960 bits

Latent Reasoning: Taxonomy

- **L-SEQ - Latent Sequence Reasoning**
 - Generate a sequence of latent thoughts before producing the final answer
 - Inspiration from early Test-Time Scaling and Compression research
 - Also referred to as “token-wise horizontal level optimization”
- **L-LOOP - Latent Looped (Recurrent) Reasoning**
 - Methods that deepen reasoning through iterative computation across layers for each token
 - Inspiration from Sparsity / Early-Exit / MoE-MoD research
 - Also referred to as “layer-wise vertical level optimization”
- **L-OPT - Latent Optimization (Control)**
 - methods that focus on optimizing hidden reasoning
 - mainly without architecture modifications

Latent Reasoning: Goals

Usual goals reported in the literature can be summarized in:

- Accuracy / reasoning quality
 - Improve performance on reasoning-heavy tasks (math, logic, coding, algorithmic tasks) at a given compute budget.
- Efficiency / compute
 - Reduce FLOPs, latency, or token count while keeping accuracy roughly flat or slightly improved.
 - Includes: fewer generated tokens, chunked decoding, adaptive depth, parallel drafting, etc.
- Data efficiency / Pre-Training efficiency
 - Get more learning signal per token (e.g., using latent thoughts or concepts during pretraining to reduce data requirements or improve scaling).
- Out-of-distribution / systematic generalization
 - Improve robustness to distribution shifts and compositional generalization (e.g., new problem sizes, unseen graph structures, or systematically recombined primitives)


L-SEQ

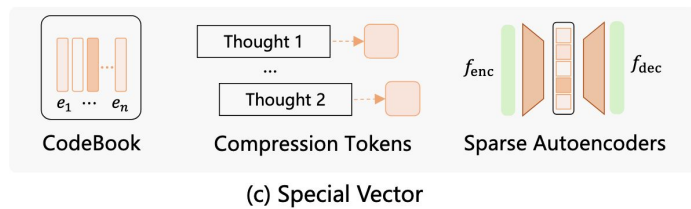
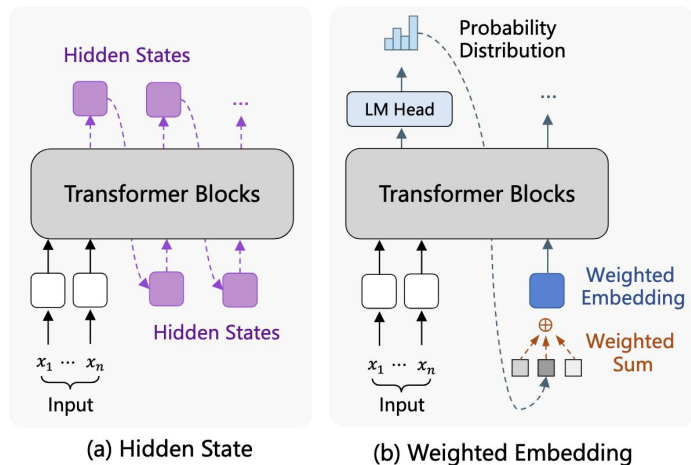
Latent Sequence Reasoning

Sequential Latent Reasoning (L-SEQ)

- Early TTS efforts: padding / pause tokens / planning tokens
 - In parallel: trials to embed the logic of an explicit CoT directly into the model's latent representations without extending the visible token sequence.
- 3 key aspects:
 - Representation Initialization
 - Model Optimization
 - Inference Exploration



L-SEQ: Representation Initialization


- A main design choice!
 - dictates the nature of the latent space
 - strongly influences the model's capacity to reason
- Hidden State 
- Weighted Embedding
- Special Vector



 Indicates approaches with stronger empirical evidence of success

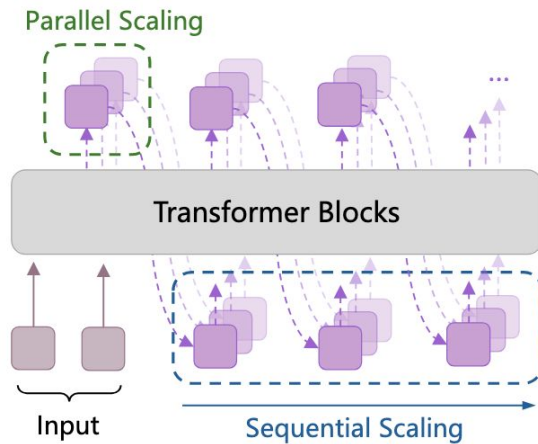
L-SEQ: Model Optimization

- Pre-Training: Beyond NTP
- Post-Training:
 - Indirect SFT: hope to induce learning useful latent representations
 - KL-Div / VAEs / Regularization
 - Direct SFT 
 - Representation optimization with “gold” targets.
 - KV-Cache
 - RL 
 - Does not rely on CoT traces for supervision
- Token-level vs. Trajectory-level

 Indicates approaches with stronger empirical evidence of success

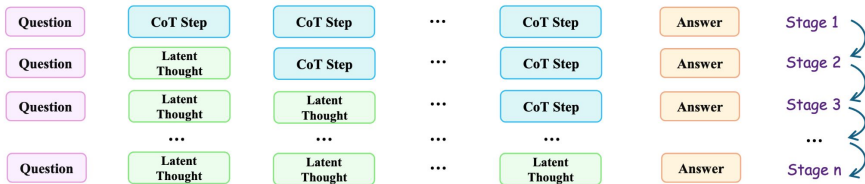
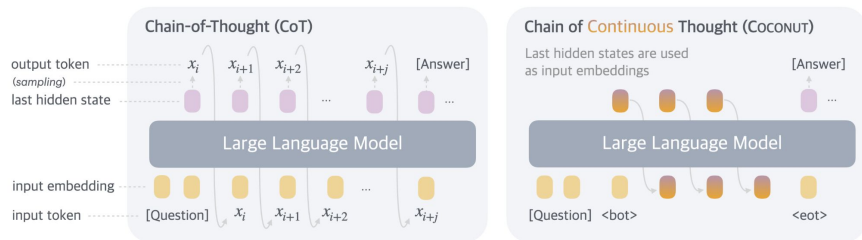
L-SEQ: Inference Exploration

- Sequential Scaling vs. Parallel Scaling
- Underexplored: Budget / Adaptability
- Optimize tokens vs. trajectories.
 - Optimizing Trajectories allows adaptability (e.g. System 1.5)

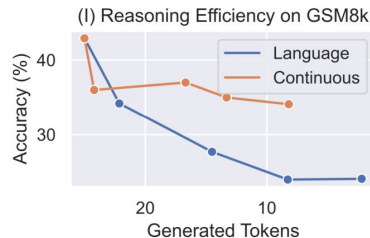


L-SEQ: Coconut

- Curriculum Learning to Compress CoT into continuous vectors.
- Surpassing Explicit CoT on Logical Reasoning with 10-100x faster inference
- Robustness to budget reduction.
- BFS-like Exploration allows improving performance on complex, planning-intensive tasks

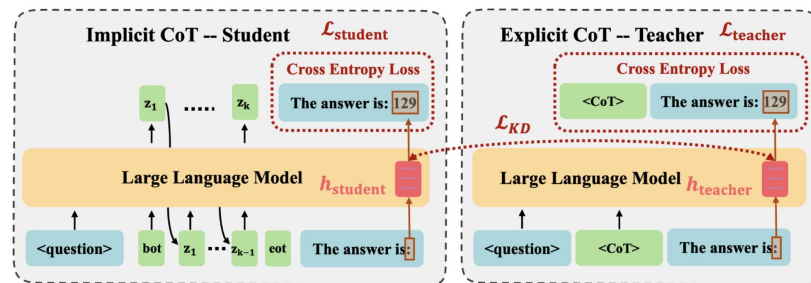
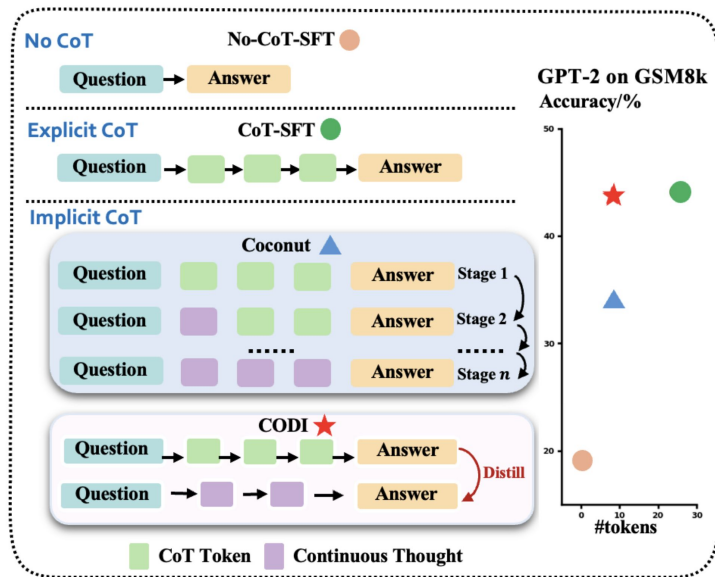


Method	GSM8k		ProntoQA		ProsQA	
	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens
CoT	42.9 \pm 0.2	25.0	98.8 \pm 0.8	92.5	77.5 \pm 1.9	49.4
No-CoT	16.5 \pm 0.5	2.2	93.8 \pm 0.7	3.0	76.7 \pm 1.0	8.2
iCoT	30.0*	2.2	99.8 \pm 0.3	3.0	98.2 \pm 0.3	8.2
Pause Token	16.4 \pm 1.8	2.2	77.7 \pm 21.0	3.0	75.9 \pm 0.7	8.2
COCONUT (Ours)	34.1 \pm 1.5	8.2	99.8 \pm 0.2	9.0	97.0 \pm 0.3	14.2



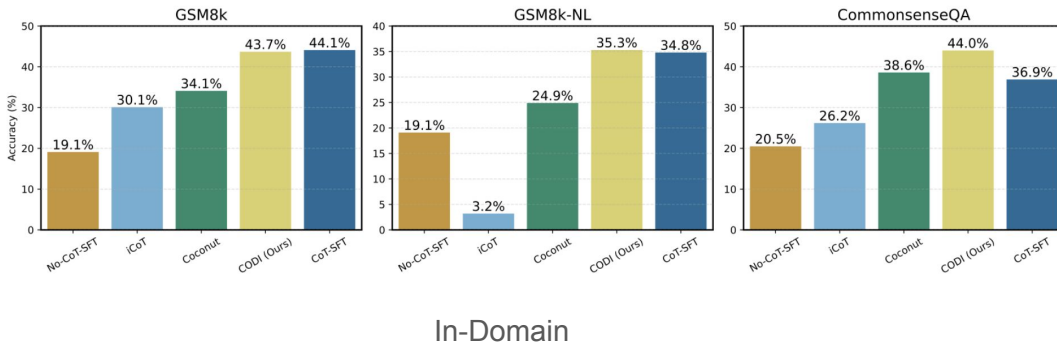
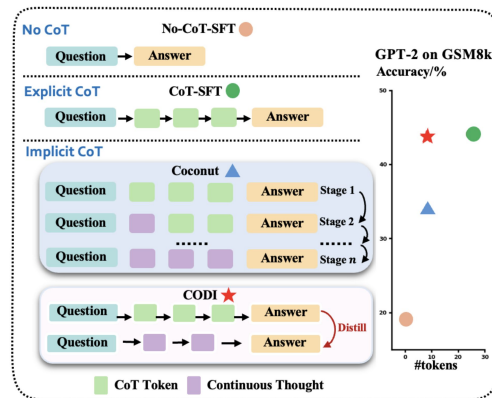
L-SEQ: CODI (Hero model)

- Supervision via distillation
- Comparable to Explicit CoT on Math



L-SEQ: CODI (Hero model)

- Supervision via distillation
- Comparable to Explicit CoT on Math
- OOD Generalization
- However, unstable



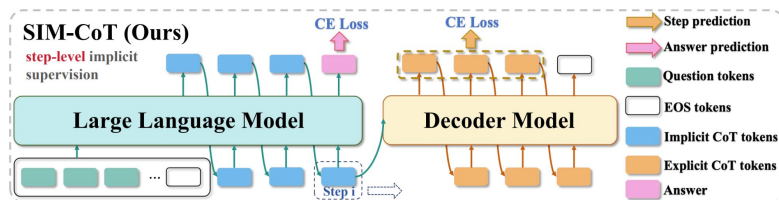
In-Domain

Models	SVAMP	GSM-Hard	Multia
GPT-2			
No-CoT-SFT	16.4	4.3	41.1
CoT-SFT	41.8	9.8	90.7
iCoT	29.4	5.7	55.5
Coconut	36.4	7.9	82.2
CODI	42.9	9.9	92.8

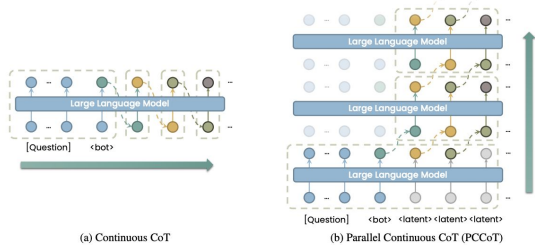
Math (OOD)

L-SEQ: Supervision is the Key

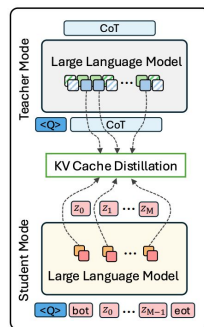
- Variants of CODI - Solve instability with supervision
 - Parallel Continuous CoT (PCCoT): Parallel Scaling
 - KaVA: KV-Cache Compression and Distillation
 - Sim-CoT: Direct Supervision with Thought Decoder
 - CoLaR: Embedding Compression + RL pipeline
 - Chained: PCCoT -> KaVA -> Sim-CoT



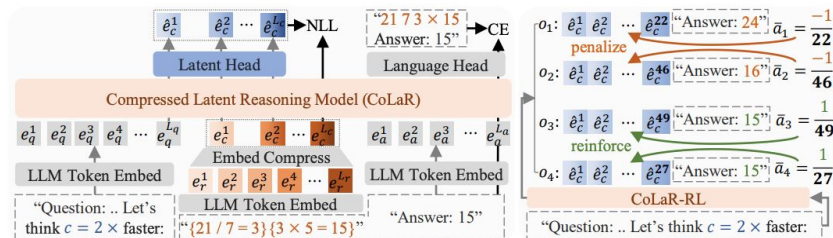
Sim-CoT



PCCoT



KaVA



CoLaR

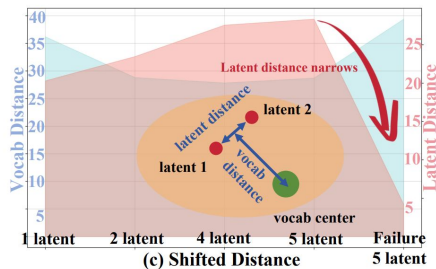
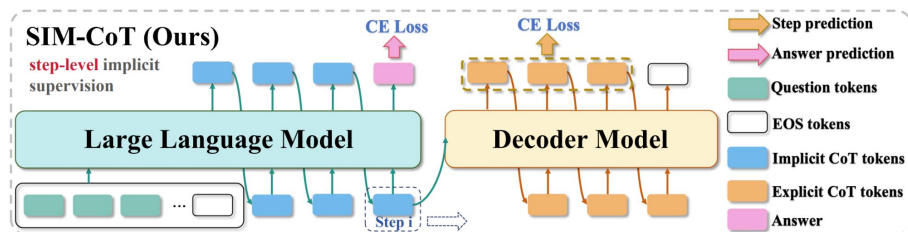
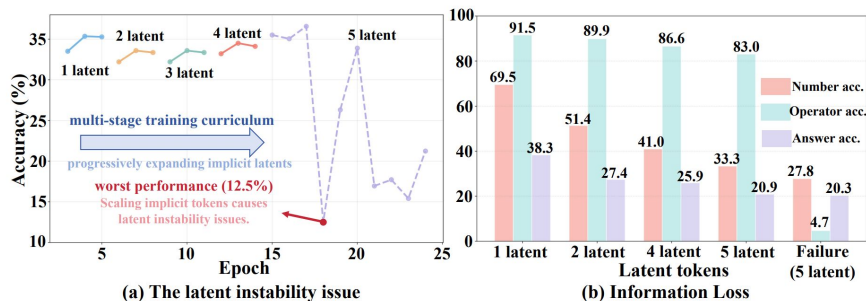
Wu et al., [Parallel Continuous Chain-of-Thought with Jacobi Iteration](#) (EMNLP 2025)

Kuzina et al., [KaVA: Latent Reasoning via Compressed KV-Cache Distillation](#) (2025)

Wei et al., [SIM-CoT: Supervised Implicit Chain-of-Thought](#) (2025)

Tan et al., [Think Silently, Think Fast: Dynamic Latent Compression of LLM Reasoning Chains](#) (NeurIPS 2025)

L-SEQ: Supervision is the Key



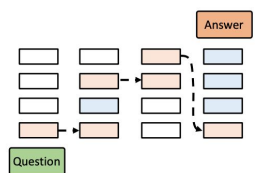
(d) Semantic Homogenization: A text box comparing a 'Normal Implicit Model' with a 'Failed Implicit Model' using a specific question: "Janet's ducks lay 16 eggs per day ... How much in dollars does she make every day at the farmers' market? steps: <<16-3-4=9>>, ... ###gt: 18".

Normal Implicit Model:
 latent 1: ["3", "4", "15", "2", "5", "16", "20", "9"]
 latent 2: ["+", "*", "=", "16", "3", "4", "9"]
 ... (Larger latent distance, richer content.)

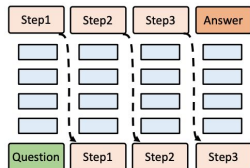
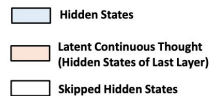
Failed Implicit Model:
 latent 1: ["16", "3", "###", "2", "4", "8", "18", "24"]
 latent 2: ["3", "16", "30", "2", "6", "12", "18", "24"]
 ... (Smaller latent distance, more uniform.)

Method	SIM-CoT	In-domain		Out-of-domain				
		GSM8k-Aug	GSM-Hard	MultiArith	SVAMP	Average	# Average	
		Acc. (%)	# Tokens	Acc. (%)	Acc. (%)	Acc. (%)	Acc. (%)	Tokens
SFT-CoT	✗	58.4	25.3	13.9	96.7	65.7	58.8	23.1
No-CoT	✗	28.8	1.2	6.3	50.3	26.7	27.8	1.9
iCoT	✗	19.0	1.2	4.4	39.0	40.9	28.1	1.9
Coconut	✗	33.2	13.2	7.0	63.3	43.7	38.0	11.9
	✓	42.2 (+9.0)	13.2	9.3	87.7	43.9	47.0 (+9.0)	11.9
CODI	✗	52.7	13.2	11.9	95.0	60.6	55.8	13.4
	✓	56.1 (+3.4)	13.2	12.7	96.2	61.5	56.8 (+1.0)	13.4

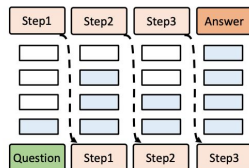
L-SEQ: System 1.5



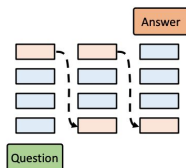
(a) System-1.5 Reasoning (Ours)



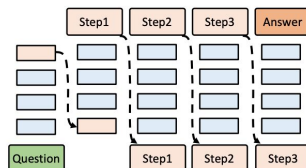
(b) Chain-of-Thought (CoT)



(c) Early Exit

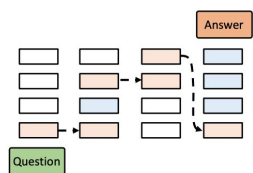


(d) Compressed Latent Reasoning

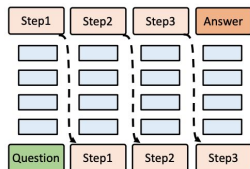
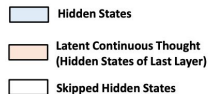


(e) Extra Latent Reasoning

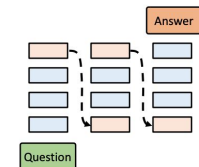
L-SEQ: System 1.5



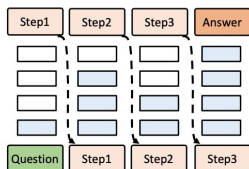
(a) System-1.5 Reasoning (Ours)



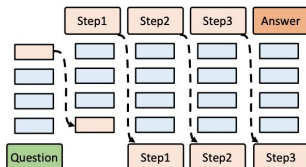
(b) Chain-of-Thought (CoT)



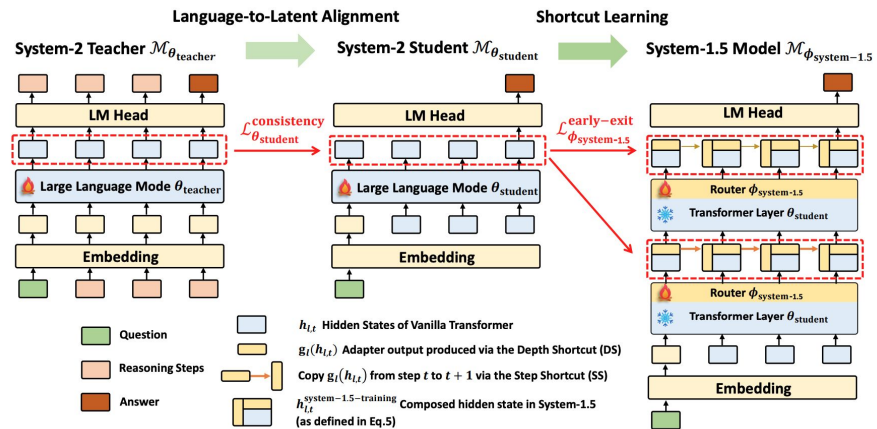
(d) Compressed Latent Reasoning



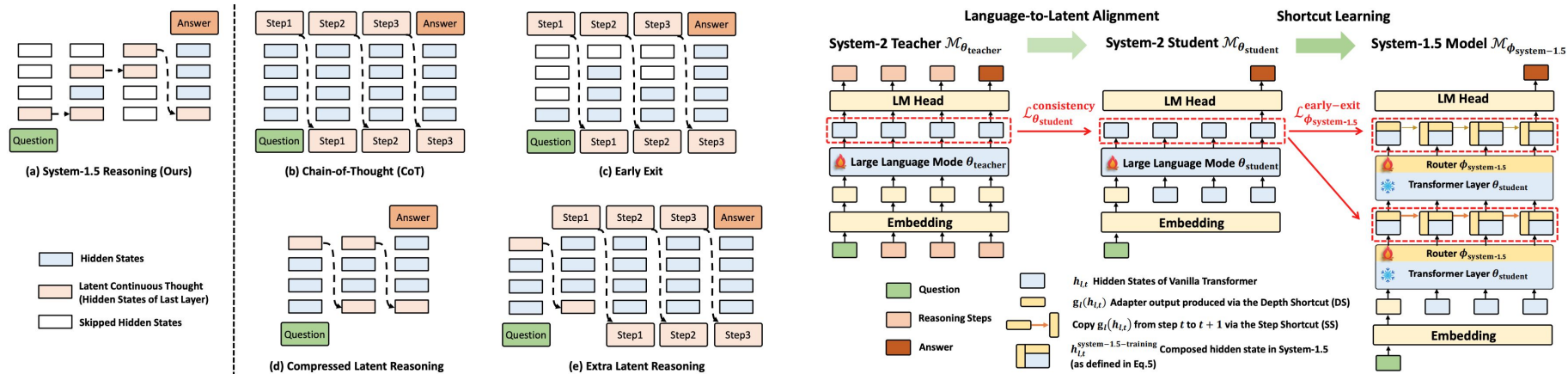
(c) Early Exit



(e) Extra Latent Reasoning



L-SEQ: System 1.5



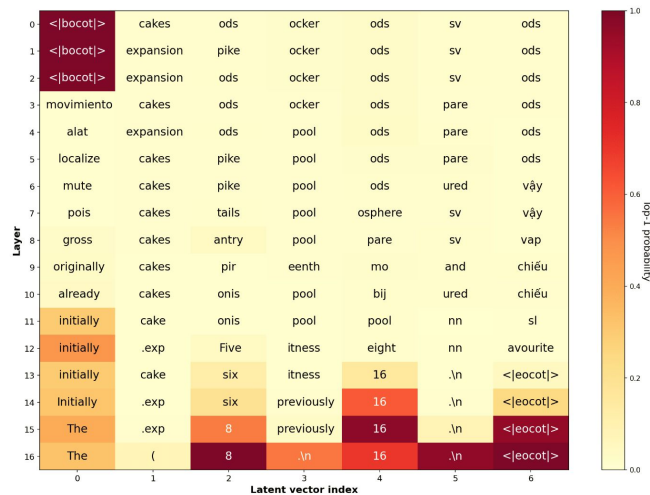
Method	GSM8K				GSM-HARD				StrategyQA			
	Acc. (%)	# Steps	FLOPs r.	Speedup	Acc. (%)	# Steps	FLOPs r.	Speedup	Acc. (%)	# Steps	FLOPs r.	Speedup
CoT	46.94	26	-	-	38.32	26	-	-	47.62	52	-	-
LITE	44.51	28	1.84×	1.61×	37.01	28	1.56×	1.44×	46.15	42	1.96×	2.36×
LayerSkip	43.20	32	1.77×	1.45×	36.55	33	1.32×	1.03×	42.54	49	1.8×	1.86×
iCoT	32.14	2	1.02×	<u>13.45×</u>	23.17	4	1.02×	6.17×	34.42	2	1.02×	26.47×
Coconut	36.75	2	1.02×	11.98×	28.25	4	1.02×	7.07×	38.67	2	1.02×	27.25×
CODI	43.78	2	1.02×	13.37×	35.91	4	1.02×	6.93×	45.12	2	1.02×	25.96×
pause token	46.32	38	1.00×	0.8×	38.17	42	1.00×	0.89×	48.14	61	1.00×	0.82×
System-1.5	<u>46.66</u>	2	1.95×	20.27×	<u>38.28</u>	4	1.76×	12.45×	48.61	2	2.12×	55.65×

L-SEQ: Scratchpad Thinking Mechanistic Hypothesis

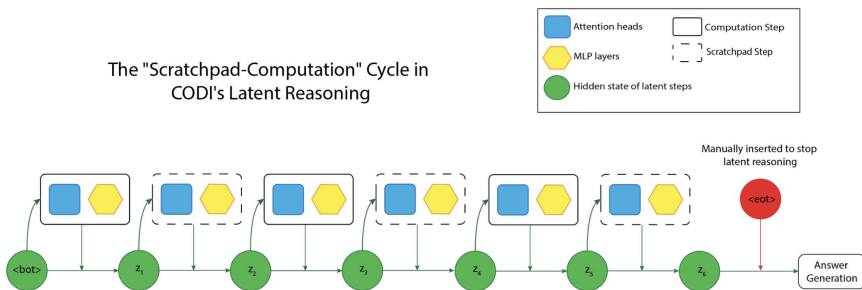
- **Even latent steps (scratchpad):** hold numerical content, and are where number-related interventions matter most.
- **Odd latent steps (compute):** where the main operation updates happen.
- **Supporting signals:** attention to numeric tokens and latent tokens peaks on even steps, early decoding, logitlens, patching, SAEs.
- Example (right): Model solves 3-step math problems with 6 latent vectors by storing the 2 intermediate values in specific vectors (3rd & 5th)

GSM8k: "A team starts with 3 members. They recruit 5 new members. Then each current member recruits 2 additional people. How many people are there now on the team?", we expect the reasoning should look like:

Step 1: $3 + 5 = 8$
Step 2: $\text{Step 1} * 2 = 8 * 2 = 16$
Step 3: $\text{Step 1} + \text{Step 2} = 8 + 16 = 24$



The "Scratchpad-Computation" Cycle in COD's Latent Reasoning



L-LOOP

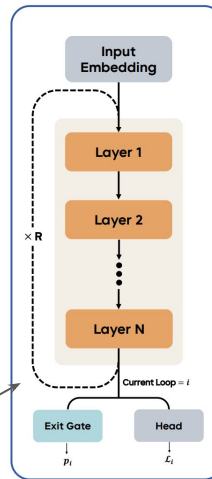
Latent Looped Reasoning

L-LOOP

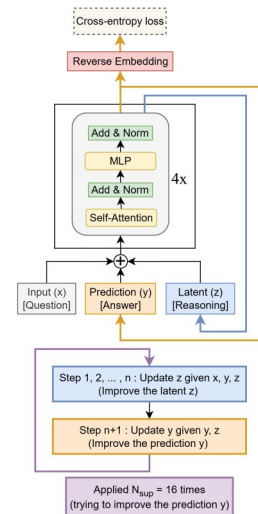
Recursive process that **iteratively refines internal representations**:

Fixed number of “think” embeddings to refine

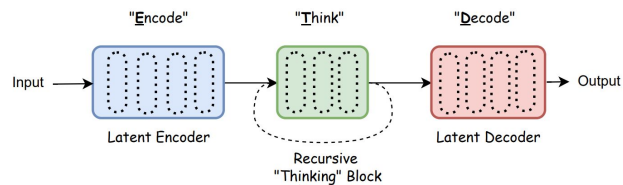
- Encode-Think-Decode (Recursion on specific intermediate latent representations)
- Scaling Latent Reasoning via Looped Language Models (loop over all embeddings)
- Recursion only on few specific representations (decouples answers from reasoning) TRM
- Denoise latent representation with Diffusion/flow matching



Recursion over all embeddings and all blocks



Decouple answer from reasoning



Recursion over specific intermediate block

Encode-Think-Decode

Adapt a pre-trained model (Olmo2) to loop over few middle layers

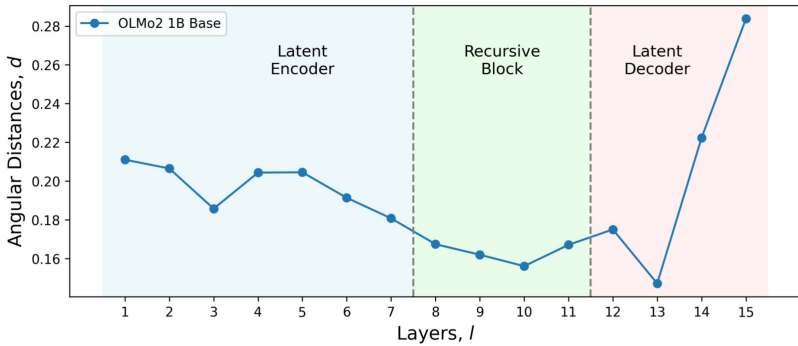
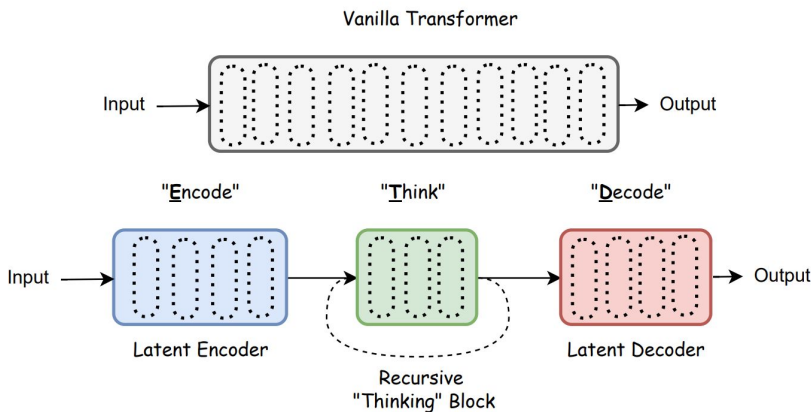
Observations from (Gromov et al 2024): dropping layers with low angular distances

=> **No impact on knowledge retrieval tasks**

=> **High degradation on reasoning tasks**

Experiment:

Second stage training (high quality text), Olmo models with varying number of loop



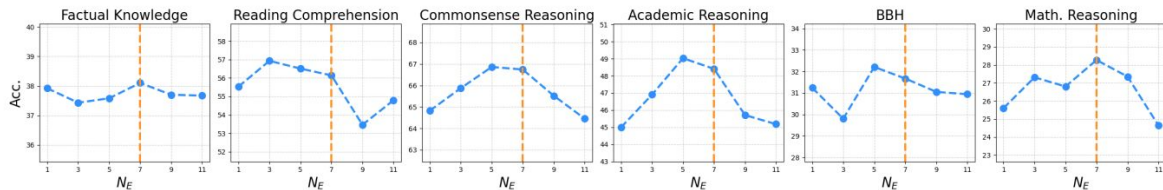
Results

No improvement for Factual Knowledge

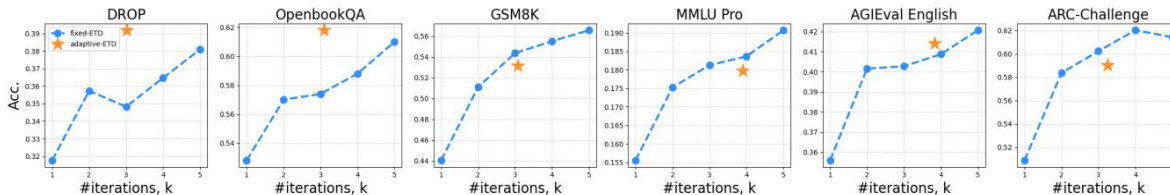
Large improvement for Math and reasoning tasks

Model	Params/FLOPs	Factual Knowledge		Reading Comprehension		Commonsense Reasoning		Multi-Disciplinary Reasoning		BBH		Math. Reasoning	
		Acc.	Δ (%)	Acc.	Δ (%)	Acc.	Δ (%)	Acc.	Δ (%)	Acc.	Δ (%)	Acc.	Δ (%)
OLMo 2 (k=1)	16 / 16	37.55	-	52.19	-	65.29	-	45	-	31.8	-	24.31	-
ETD (k=2)	16 / 20	38.1	(+1.5%)	56.14	(+7.6%)	66.74	(+2.2%)	48.41	(+7.6%)	31.67	(-0.4%)	28.27	(+16.3%)
ETD (k=3)	16 / 24	37.55	(0%)	56.07	(+7.4%)	67.75	(+3.77%)	49.55	(+10.1%)	32.62	(+2.6%)	30.29	(+24.6%)
ETD (k=4)	16 / 28	37.74	(0%)	57.76	(+10.7%)	68.16	(+4.4%)	50.18	(+11.5%)	33.01	(+3.8%)	29.62	(+21.8%)
ETD (k=5)	16 / 32	38.23	(+1.8%)	58.5	(+12.1%)	68.41	(+4.8%)	50.58	(+12.4%)	33.49	(+5.3%)	30.45	(+25.3%)

How does the choice of recursive layers change performance? Fix the recursive “thinking” block size and vary its starting position from layer 2 to 12



Adaptive test-time scaling: Train an exit gate that predict when to exit the thinking loop



Scaling Latent Reasoning via Looped Language Models

Method:

Pretrain model from scratch to Loop representations multiple times over the model

Stage 1: CE with entropy regularization

$$\mathcal{L} = \underbrace{\sum_{t=1}^{T_{max}} p_{\phi}(t | x) \mathcal{L}^{(t)}}_{\text{expected task loss}} - \underbrace{\beta H(p_{\phi}(\cdot | x))}_{\text{entropy regularization}}, \quad H(p_{\phi}(\cdot | x)) = - \sum_{t=1}^{T_{max}} p_{\phi}(t | x) \log p_{\phi}(t | x).$$

with Beta controlling how uniform the distribution $p(\cdot | x)$ is

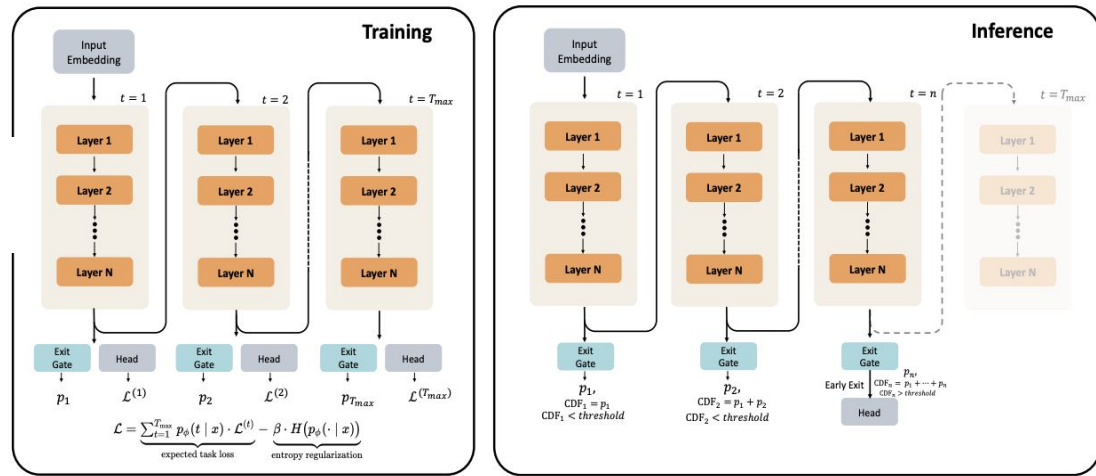
Stage 2: Focused Adaptive Gate Training

Build label for exit gate to predict correct exit probabilities using:

$$I_i^{(t)} = \max(0, \mathcal{L}_{i,stop}^{(t-1)} - \mathcal{L}_{i,stop}^{(t)})$$

0 = should exit (higher loss in future steps)

1 = continue looping (lower loss in following steps)



Scaling Latent Reasoning via Looped Language Models

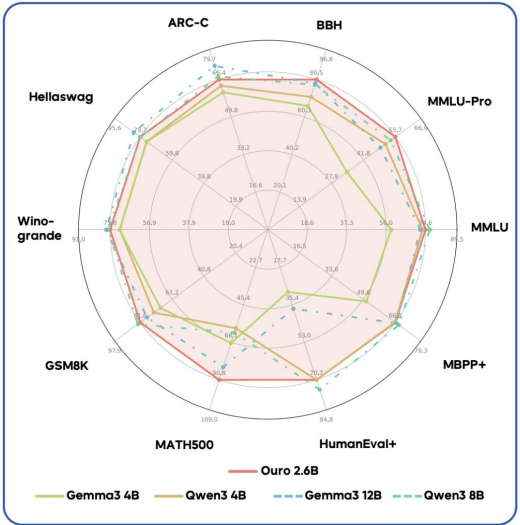
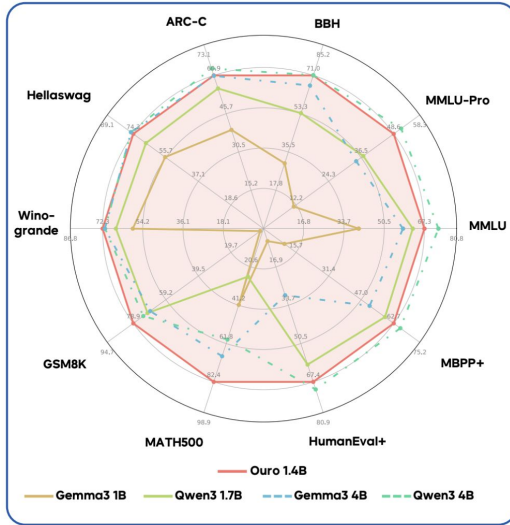
Contributions:

Parameter efficiency for reasoning tasks:
 Pre-training on 7.7T tokens, 1.4B and 2.6B parameter LoopLMs match 4B and 8B standard transformers

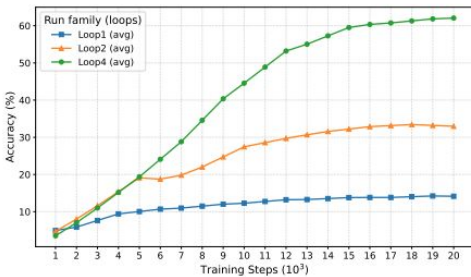
Entropy-regularized adaptive computation:
 Adaptive exits tend to collapse to shallow depths or overuse long loops => added entropy regularization under a uniform prior over exit steps for unbiased depth exploration

Understanding of recurrence:
 Recurrence does not increase raw knowledge storage but enhances knowledge manipulation capabilities on tasks requiring fact composition and multi-hop reasoning.

Improved safety and faithfulness:
 LoopLM reduces harmfulness on HEX-PHI



Ouro 1.4B and 2.6B models, both with 4 recurrent steps (red)



MultiHop question answering

	$L = 10$	$L = 16$	$L = 24$
Baseline model			
Base ($12 \otimes 1$)	93.6	94.4	34.8
2 layer model			
Base ($2 \otimes 1$)	21.5	8.4	7.5
Loop ($2 \otimes 6$)	98.1	96.3	78.0
3 layer model			
Base ($3 \otimes 1$)	75.4	29.8	11.0
Loop ($3 \otimes 4$)	97.9	95.8	92.2
6 layer model			
Base ($6 \otimes 1$)	84.7	59.5	20.0
Loop ($6 \otimes 2$)	93.4	88.5	35.1

Iso parameter study (Mano arithmetic task)

L-LOOP: Main takeaways

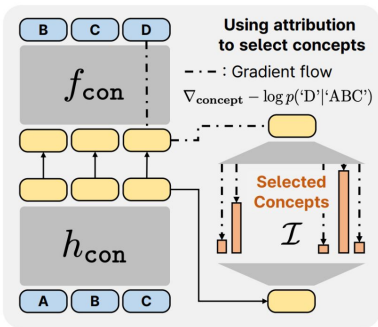
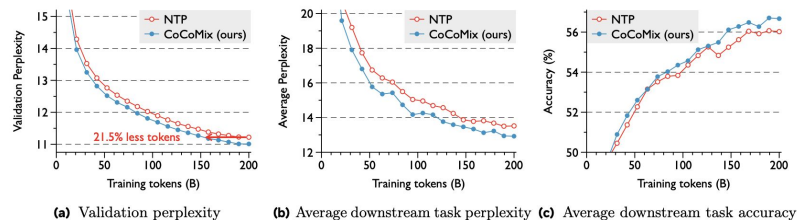
- Existing methods loop over a fixed number of embeddings
- Can learn strategies to adapt the number of loop
- Doesn't improve on knowledge retrieval tasks
- Improves reasoning performance (knowledge manipulation)
- Looping on some layers seems to have more effect on reasoning tasks (routing mechanism)

L-OPT

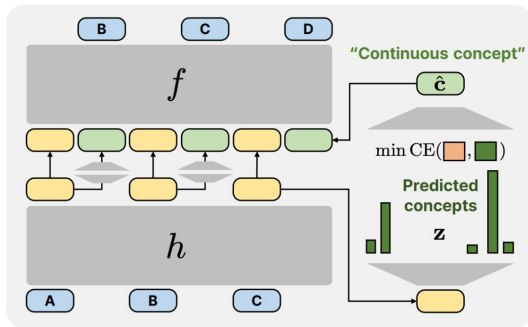
Latent Optimization

L-OPT: CoCoMix

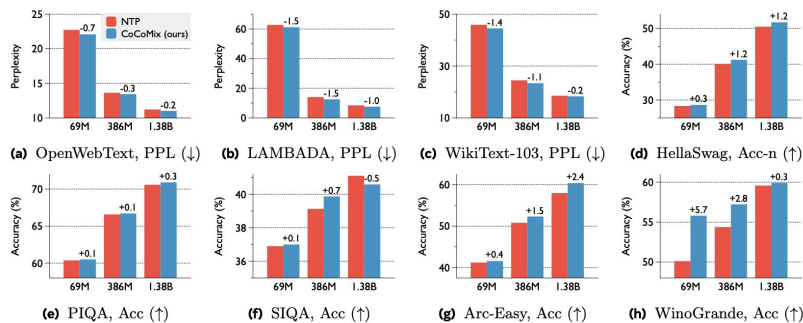
- L-OPT: methods that focus on optimizing hidden representations for reasoning
- Extract concepts from pre-trained model to guide learning
- Model predicts from compressed continuous representation of predicted concepts
- Symbolic learning, interpretability and steerability



Extracting concepts from a pretrained SAE model's hidden state



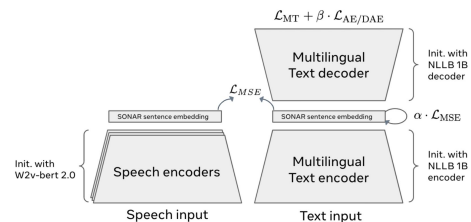
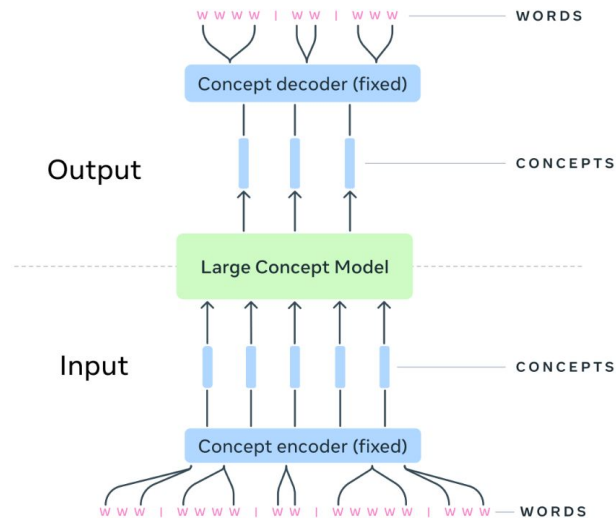
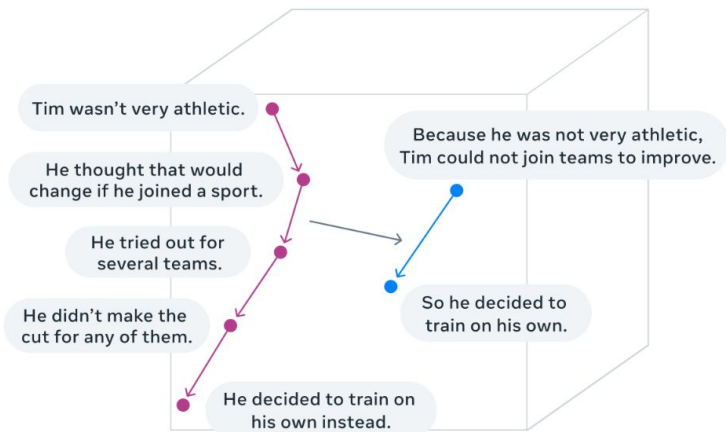
Learning to predict concepts & Mixing/Interleaving continuous concepts into the hidden state



Alternative Architectures

Large Concept Models

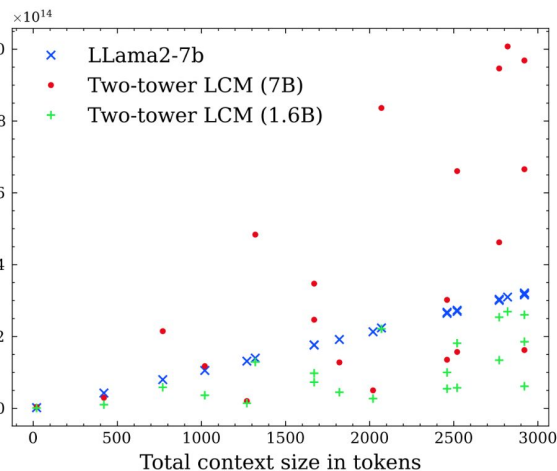
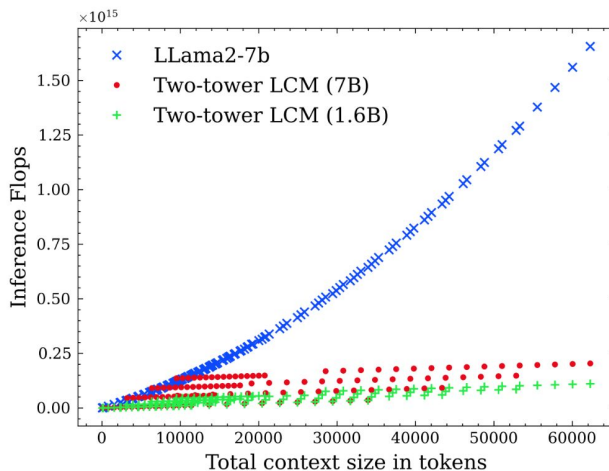
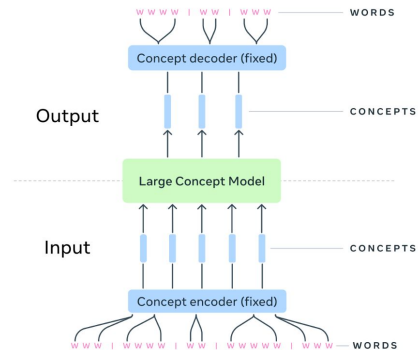
- Reasoning in conceptual abstract **language- and modality-agnostic** representation space
- Hierarchical structure
- Handling of long context and long-form output.
- Modularity and extensibility
- Unparalleled zero-shot generalization



SONAR Encoders

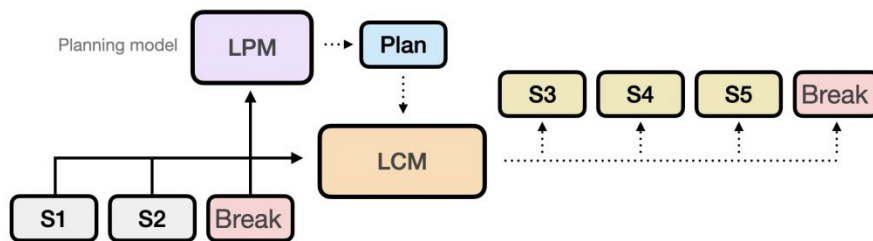
Large Concept Models

- FLOPs scale **sublinearly** with context
- **Unparalleled zero-shot generalization:** LCM w/ 2-Tower Diffusion beats Gemma 7B and Llama-3.1-8B in zero-shot multilingual tasks (summarization / expansion).



Large Concept Models

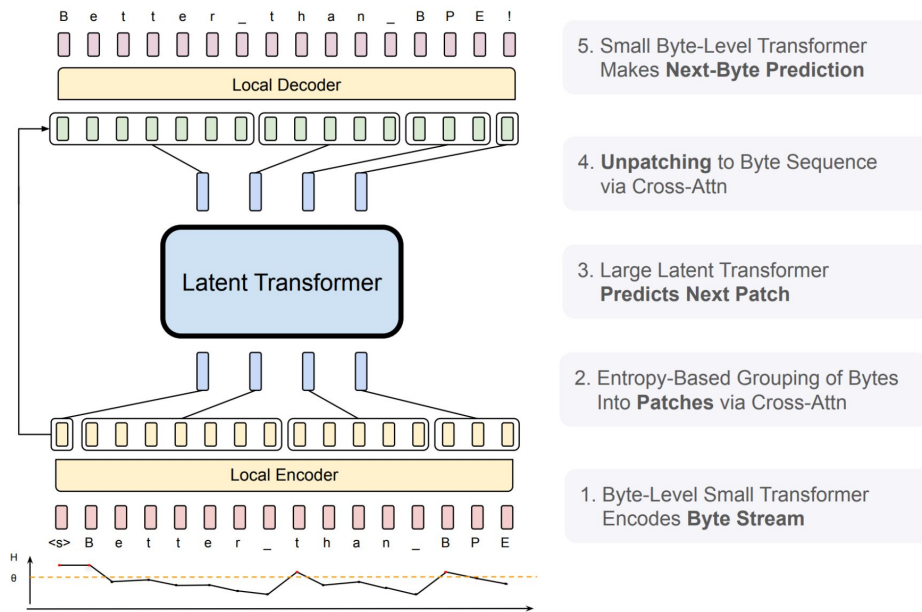
Large Planning Model + Large Concept Model -> LPCM



Large Concept Models

Byte Latent Transformers

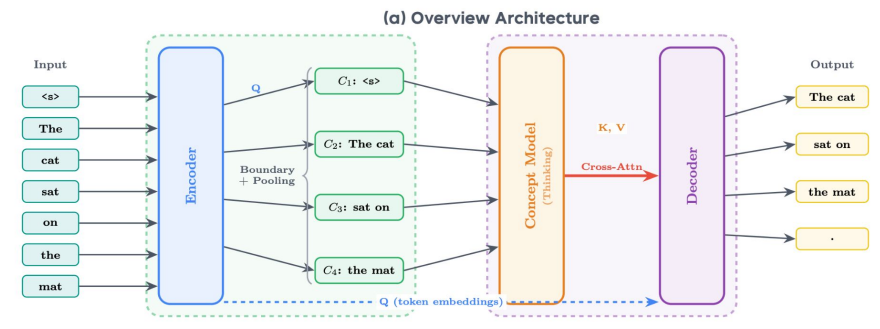
- Dynamic compute allocation
- Gains of up to 50% at 8B (vs. Llama)
- Better robustness and semantic generalization
- Improved long-context reasoning



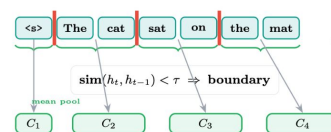
Large Concept Models

Dynamic Large Concept Models

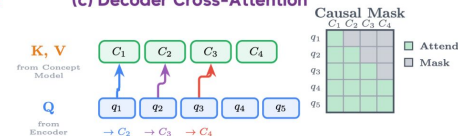
- End-to-End training with learned variable-length concepts
- Back to next-token prediction
- Great improvement from Llama on several tasks.



(b) Boundary Detection & Pooling

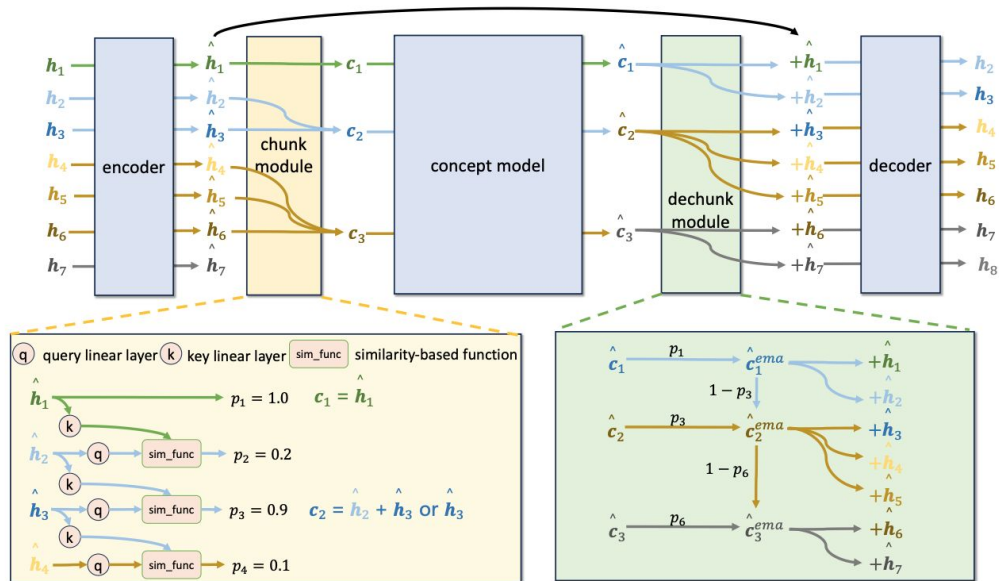


(c) Decoder Cross-Attention



Large Concept Models

ConceptMoE



Diffusion Models

LaDiR: Latent Diffusion Enhances LLMs for Text Reasoning

Denoises latent tokens Z in parallel

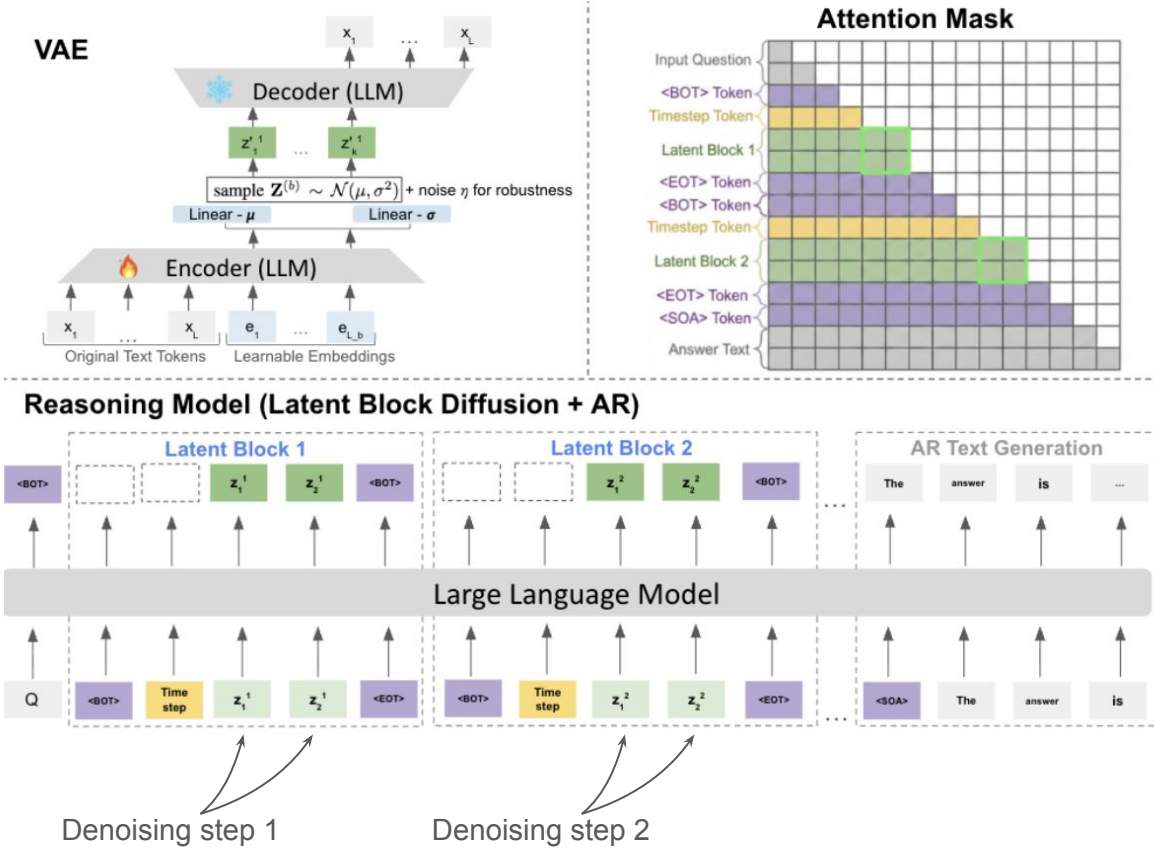
VAE: Trained to predict latent thought

LLM: Trained to diffuse latent thought in parallel

Stage1: Teacher Forcing

$$\mathcal{L} = \lambda_{FM} \mathcal{L}_{FM} + \lambda_{Ans} \mathcal{L}_{Ans} + \lambda_{Spec} \mathcal{L}_{Spec},$$

Stage2: FlowGRPO



Results

Method	In-Domain		Out-of-Domain					Avg.
	MATH	GSM8K	Gaokao	DM-Math	College	Olympia	TheoremQA	
Coconut	37.3/39.3	68.3/74.3	26.8/29.3	33.5/36.9	40.2/42.9	5.8/6.3	11.4/14.9	31.9/34.8
CODI	38.5/45.1	76.3/81.7	27.5/35.2	38.8/44.9	43.2/49.0	7.6/14.8	8.8/15.7	34.3/40.9
SFT ($\alpha = 1$)	42.1/51.0	83.5/90.1	29.7/39.0	46.6/53.1	44.6/56.0	9.4/13.5	20.3/26.0	39.3/47.1
LaDiR	46.2/63.7	84.8/93.7	35.4/45.8	52.3/54.2	48.6/60.3	12.9/15.3	24.7/30.7	43.5/52.0
-w/o Stage 2	30.7/35.8	57.8/62.6	24.7/30.2	32.0/36.9	32.8/38.0	5.9/10.5	11.9/16.9	27.9/33.0

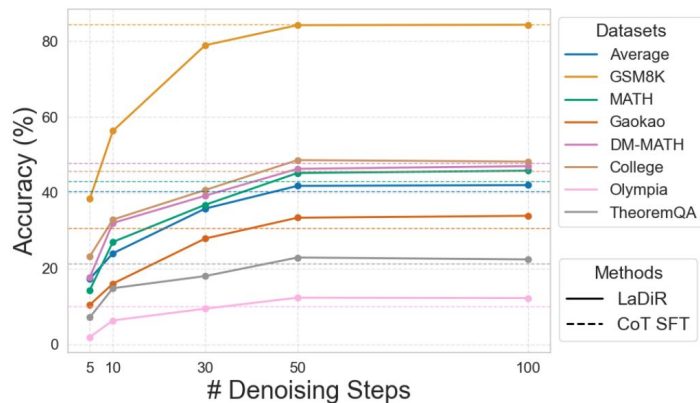
Pass@1 (left) / Pass@100 (right)

PUZZLE PLANNING – COUNTDOWN : $\{+, -, \times, \div\} \times \{a, b, c, d\} \Rightarrow e$

Model	CD-4 P@1	CD-4 P@100	CD-4 Div.	CD-5 P@1	CD-5 P@100	CD-5 Div.
Dream 7B Base*	16.0	24.7	4.1	4.2	10.3	5.6
MGDM [†]	91.5	<u>95.2</u>	3.2	46.6	70.4	4.9
LLaDA 8B SFT	<u>51.2</u>	<u>75.2</u>	<u>5.4</u>	<u>34.4</u>	<u>45.2</u>	<u>6.2</u>
LLaMA 8B SFT	46.7	65.3	3.0	8.9	15.4	3.5
LaDiR	<u>76.6</u>	96.4	7.3	<u>38.5</u>	75.2	8.9

Results

Effect of number of denoising steps



VAE robustness Augmentation

Gaussian Noise ($p=0.3$)		Token Substitution ($k=3$)	
k (std)	Acc. (%)	p (prob.)	Acc. (%)
0	68.3	0.0	70.2
1	73.4	0.1	78.3
3	84.2	0.3	84.2
5	79.4	0.5	64.0
—	—	0.7	32.4

Interpretability

Block	Text
Question	<i>Billy sells DVDs. He has 8 customers on Tuesday. The first 3 buy one DVD each, the next 2 buy two DVDs each, the last 3 buy none. How many DVDs did Billy sell?</i>
Decode($Z^{(1)}$):	Billy's first 3 customers buy one DVD each, so that's $3 * 1 = \langle\langle 3 * 1 = 3 \rangle\rangle$ 3 DVDs.
Decode($Z^{(2)}$):	His next 2 customers buy 2 DVDs each, so that's $2 * 2 = 4$ DVDs.
Decode($Z^{(3)}$):	His last 3 customers don't buy any DVDs, so that's 0 DVDs sold.
Decode($Z^{(4)}$):	Therefore, Billy sold a total of $3 + 4 + 0 = 7$ DVDs on Tuesday.
Answer	The answer is: 7.

Multimodality: Reasoning with visual inputs

- 1- Align text and vision representations
- 2- Train to predict visual representations
- 3- Distil CoT from VLM teacher (multimodal CODI)

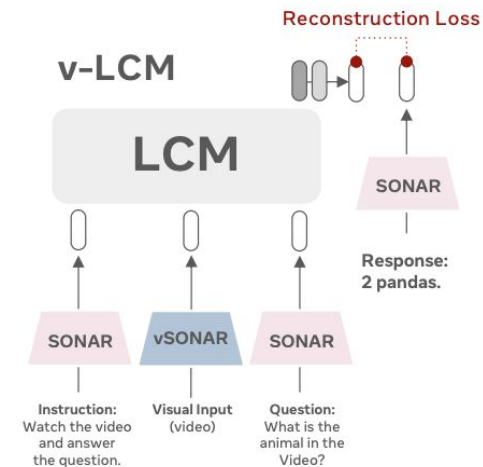
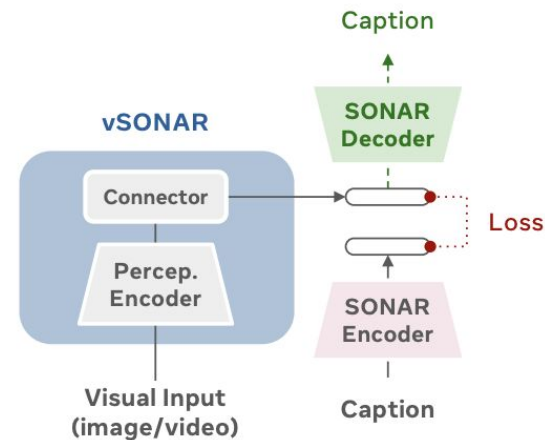
Align text and vision representations

v-LCM

Latent diffusion language model operating directly in the Sonar/vSonar embedding space

Stage 1: Align perceptron (visual encoder) with SONAR (MSE + contrastive loss)

Stage 2: Same training as LCM but visual and text representations



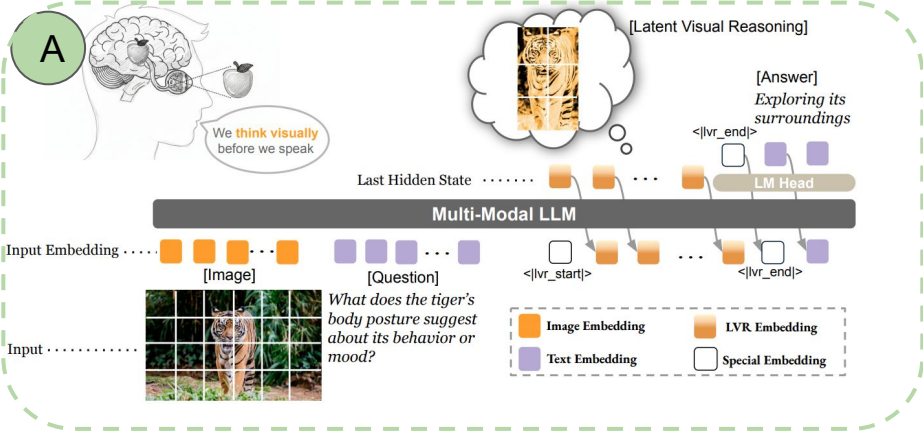
Results

		Video Captioning / Summarization								M3IT Image				M3IT Video		
		PE-Video		Dream-1k		Vatex		VideoXum		COCO	VIQAE	VisualMRC	ScienceQA	ActivNetQA	MSRVTT-QA	IVQA
		R-L	BS	R-L	BS	R-L	BS	R-L	BS	R-L	R-L	R-L	Acc.	R-L	R-L	R-L
Qwen2-VL	2B	31.2	37.3	18.5	13.9	16.4	30.8	23.6	29.8	24.9	50.2	56.1	54.5	53.7	39.6	49.4
	7B	26.9	32.6	19.8	18.1	28.5	51.6	26.0	32.4	23.7	49.7	57.4	70.4	41.9	22.7	39.1
Qwen2.5-VL	3B	28.9	34.4	15.9	8.6	15.0	27.6	26.0	32.9	25.1	48.3	55.7	55.0	52.1	41.6	48.5
	7B	22.2	25.9	15.7	10.5	27.5	50.8	24.1	28.9	18.5	34.5	45.0	61.6	46.0	41.4	54.2
Percep. LM	1B	26.6	31.0	19.3	15.5	19.1	29.6	21.8	33.2	27.5	30.8	45.5	73.6	27.8	14.5	39.1
	3B	26.4	31.3	20.4	19.0	19.3	30.8	27.0	36.4	34.3	23.7	51.1	89.8	28.0	19.4	26.1
	8B	27.4	31.9	20.8	19.7	19.0	30.8	26.2	33.7	36.3	31.0	50.0	87.7	40.5	25.3	41.4
LCM	LCM	25.5	27.9	18.5	16.6	23.8	30.8	21.5	22.1	18.0	34.3	33.5	44.7	51.7	36.0	48.9
	v-LCM	27.4	30.0	19.8	19.2	28.8	48.7	20.6	25.3	38.8	39.4	34.1	76.2	63.6	48.7	63.9

Train to predict visual representations

Latent Visual Reasoning(A) and Mull-Tokens: Modality-Agnostic Latent Thinking(B)

Stage1: Visual supervision of region of interest

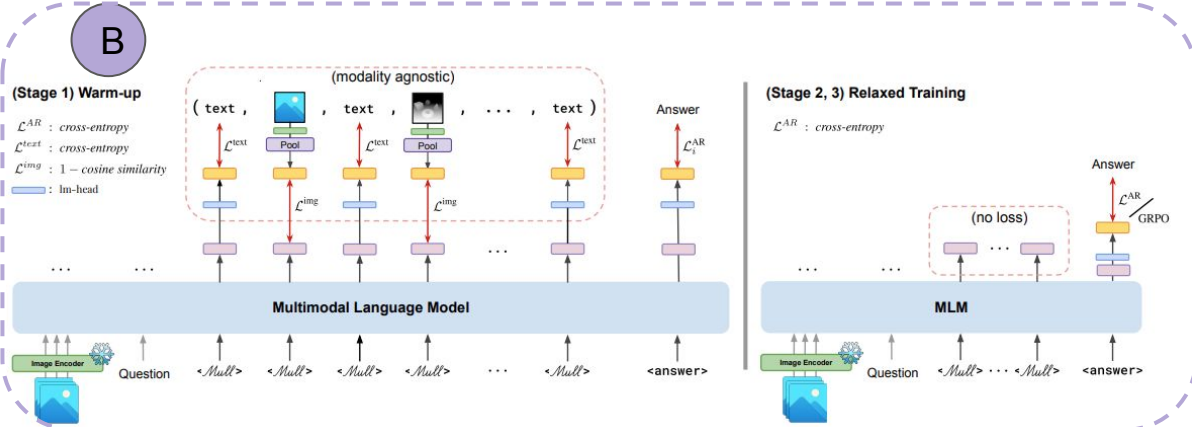


Visual Reconstruction Loss

A
$$\mathcal{L}_{LVR} = \frac{1}{T_v} \sum_{t=1}^{T_v} \|\mathbf{h}_t - \mathbf{v}_t\|_2^2$$

B
$$\mathcal{L}^{img} : 1 - \text{cosine similarity}$$

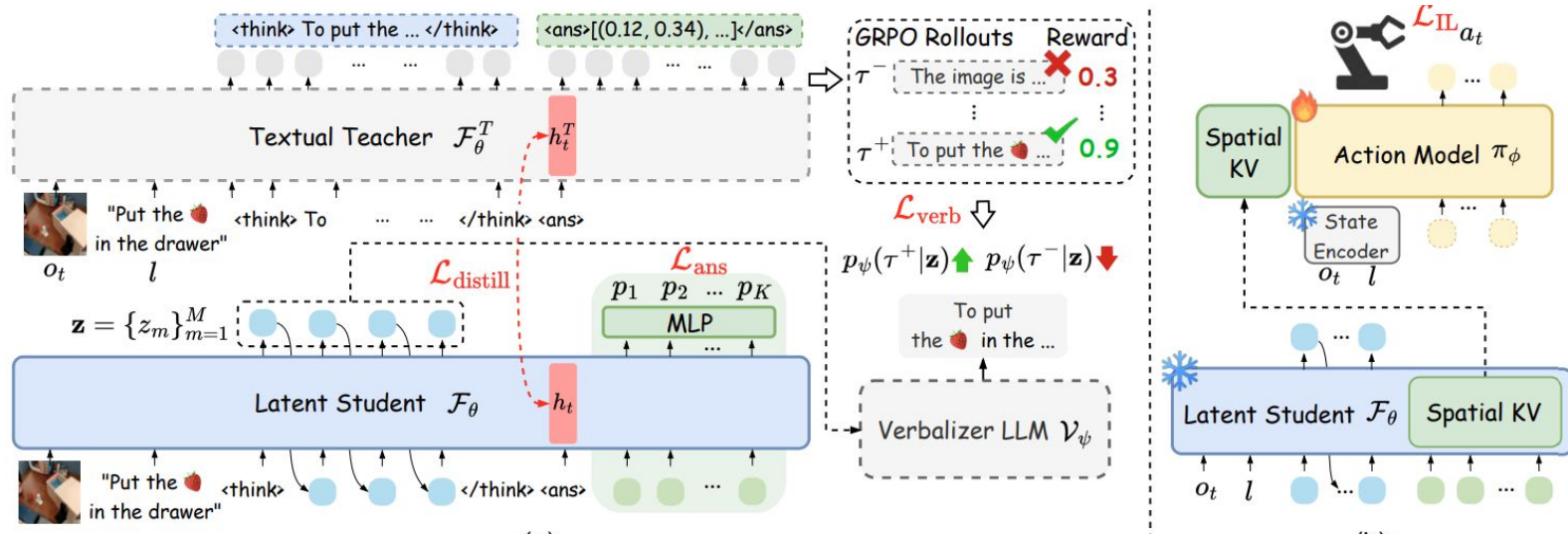
Stage2: GRPO



Li et al., [Latent Visual Reasoning](#) (2026)
 Ray et al., [Mull-Tokens: Modality-Agnostic Latent Thinking](#) (2025)

Distil CoT from VLM teacher (multimodal CODI)

Fast-Think-Act: Efficient Vision-Language-Action Reasoning via Verbalizable Latent Planning



1 - Distill CoT hidden states into compact latent tokens \mathbf{z} (similar to **CODI**) guided by reward preferences

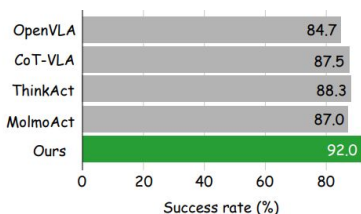
$$\mathcal{L}_{\text{student}} = \mathcal{L}_{\text{verb}} + \mathcal{L}_{\text{distill}} + \mathcal{L}_{\text{ans}}$$

2 - Student answers representations are later used to train an action model

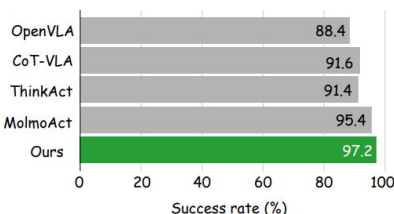
Results: Evaluation of robot manipulation and reasoning efficiency

Method	EgoPlan-Bench2					RoboVQA					OpenEQA	Overall
	Daily.	Work.	Rec.	Hobbies	Avg.	B-1	B-2	B-3	B-4	B-Avg.	Score	Avg.
GPT-4V Achiam et al. (2023)	36.7	27.7	33.9	32.5	32.6	32.2	26.5	24.7	23.9	26.8	49.6	36.4
Gemini-2.5-Flash Comanici et al. (2025)	44.2	42.3	43.2	39.1	42.4	39.1	31.6	22.9	22.1	28.9	45.3	38.9
InternVL2.5-2B Chen et al. (2024)	30.9	27.8	28.6	33.1	30.1	36.6	33.7	31.0	29.4	32.7	47.1	36.6
InternVL3-2B Zhu et al. (2025)	36.9	29.9	35.6	31.5	33.4	34.4	33.9	33.5	33.3	33.8	48.8	38.7
NVILA-2B Liu et al. (2024)	34.6	26.7	33.3	31.6	31.4	38.7	34.3	31.1	29.2	33.3	47.0	37.2
Qwen2.5-VL-3B Bai et al. (2025)	29.0	27.0	30.2	28.9	28.5	42.5	36.3	28.7	31.8	34.8	43.4	35.6
Magma-8B Yang et al. (2025)	32.1	25.7	34.4	29.3	29.8	38.6	31.5	28.1	26.7	31.2	49.1	36.7
RoboBrain2.0-3B Team et al. (2025)	45.3	37.6	45.9	39.7	41.8	54.4	47.7	43.1	41.0	46.5	50.1	46.1
ThinkAct-3B Huang et al. (2025)	46.6	41.4	45.9	42.5	44.0	62.4	57.3	52.0	49.6	55.3	48.9	49.4
Fast-ThinkAct-3B	50.3	44.3	46.4	43.2	46.4	70.1	63.0	57.2	53.0	60.8	51.2	52.8

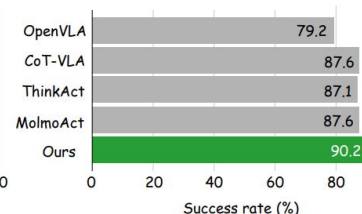
Qwen2.5B-VL-3B backbone



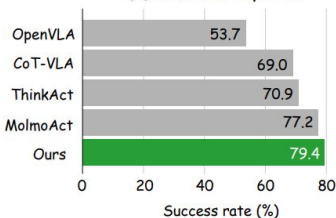
(a) LIBERO-Spatial



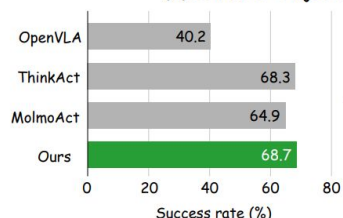
(b) LIBERO-Object



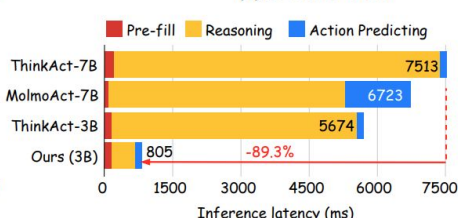
(c) LIBERO-Goal



(d) LIBERO-Long



(e) SimplerEnv-Google



(f) Latency of Reasoning VLAs

Mechanistic evidences,
challenges, limitations

Theoretical Comparison

Previously in CoT Prompting [1]: m steps $\rightarrow O(m \log m)$ CoT steps

Theorem 1. A **single fixed finite-size Transformer** can be prompted to simulate **any algorithm that runs in m steps** on a Turing machine using $O(m \log m)$ chain-of-thought steps;

L-SEQ [2]: m steps $\rightarrow m$ continuous CoT

Theorem 2. A **2-layer Transformer** can solve **directed graph reachability** with D steps of **continuous CoT**, where D is the **graph diameter**.

- Intuitively, each latent step performs one **parallel BFS expansion**.
- D-graph reachability \leftrightarrow D-hop induction problem

L-LOOP [3]: m steps $\rightarrow \log m + 2$ loops

Theorem 3. A **looped 1-layer Transformer** can simulate the effect of a deeper Transformer by reusing the same block across multiple loops, with only modest overhead in width and attention heads.

- Consequence: the **p -hop induction problem** can be solved by looping a **1-layer Transformer** only $\lceil \log_2 p \rceil + 2$ times.

Theorem 4. A looped Transformer can reproduce the output of an L -layer non-looped Transformer after m CoT steps using m **loops**, with only constant extra architectural overhead.

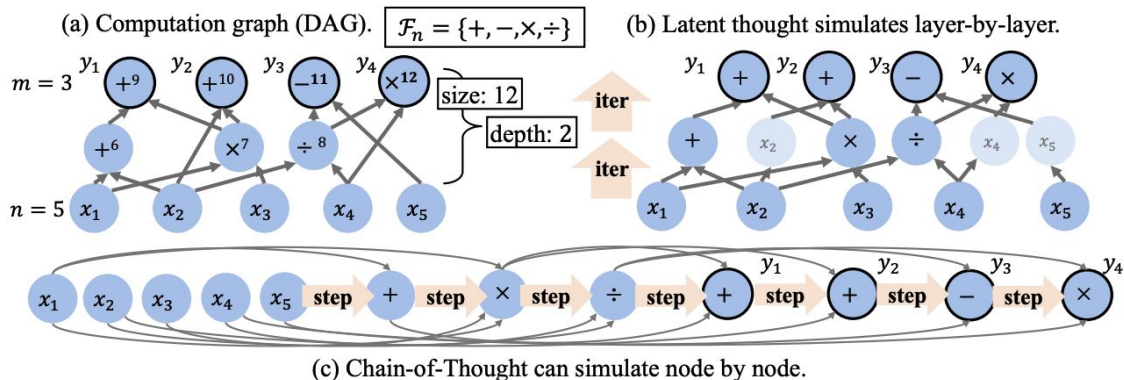
- non-looped: $(L \text{ layers}) + (m \text{ CoT steps}) \implies$ looped: $(L + O(1) \text{ layers per loop}) \times m \text{ loops}$.
- In general: **1 loop simulates 1 CoT step**.

[1] Qiu et al., Ask, and it shall be given: On the Turing completeness of prompting (ICLR 2025)

[2] Zhu et al., Reasoning by Superposition: A Theoretical Perspective on Chain of Continuous Thought (NeurIPS 2025)

[3] Saunshi et al., Reasoning with Latent Thoughts: On the Power of Looped Transformers (ICLR 2025)

Theoretical Comparison



- A computation graph G_n
 - Latent thought can simulate the computation layer by layer in parallel, using a number of loops equal to the depth of the graph, $\text{depth}(G_n)$.
 - CoT can sequentially simulate the computation node by node, using a number of steps proportional to the size of the graph, $O(\text{size}(G_n))$.
- Latent thought scales with **graph depth** / CoT scales with **graph size**
 - Generalizes the earlier “parallel BFS / reachability” intuition
 - Latent thought wins on deterministic parallel computation.
 - CoT can be strictly stronger when stochastic intermediate-token generation matters
 - especially for **approximate counting** and **sampling**

How does latent multi-step reasoning emerge during training?

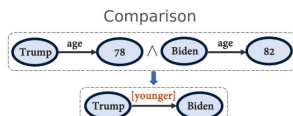
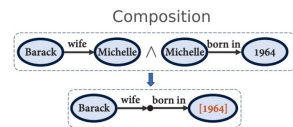
● Training dynamics

- **Stage 1 — Memorization:** fit training patterns, weak compositional generalization
- **Stage 2 — Transition / *grokking*:** memorizing behavior is replaced by more generalizing circuits
- **Stage 3 — Structured reasoning:** reusable intermediate representations support multi-step inference and better OOD generalization, though later hops remain harder than earlier ones

● What makes this emergence easier?

- More compositional / connected data
- More training data: needed data grows quickly with hop count (eg. 4-hop x100)
- Optimization bias toward simple rule-like solutions:
 - e.g. complexity control can favor reasoning over memorized mappings

● Latent reasoning is **learned**, but not guaranteed



Wang et al., [Grokked Transformers are Implicit Reasoners: A Mechanistic Journey to the Edge of Generalization](#) (NeurIPS, 2024)

Qu et al., [How do Transformers Learn Implicit Reasoning?](#) (NeurIPS, 2025)

Pan et al., [Opening the Black Box: A Survey on the Mechanisms of Multi-Step Reasoning in Large Language Models](#) (2026)

Zhang et al., [Complexity Control Facilitates Reasoning-Based Compositional Generalization in Transformers](#) (IEEE PAMI, 2025)

Yao et al., [Language models can learn implicit multi-hop reasoning, but only if they have lots of training data](#) (EMNLP 2025)

Which tasks?

- Benchmarks? Any
- Tasks in which we have evidence that implicit reasoning can do better (and more efficiently) than Explicit CoT:
 - General Knowledge/Commonsense Reasoning:
 - several works demonstrate great improvements (e.g. CommonsenseQA)
 - Multi-hop Search / Compositional Reasoning (e.g. Strategy QA)
 - Logical / Symbolic Reasoning (e.g. ProverQA, PrOntoQA, ProsQA)
 - Multilingual Comprehension and Reasoning (e.g. LCM on XLSum)
 - Multimodal Reasoning
 - Tasks that require planning
 - LLM-based Agents (still underinvestigated):
 - Support planning, deeper thinking, compositional tasks
 - Compressed Memory
 - KV Cache Sharing across Agents
 - Parallel exploration and execution (e.g. tool-calling)

Comparison Implicit vs. Explicit CoT

Dimension	Explicit Reasoning	Implicit Reasoning
Reasoning Visibility	States verbalized in text, transparent	States hidden in latent space, invisible
Reasoning Efficiency	Verbose, high cost and latency	Compact, faster, resource-efficient
Interpretability	Directly observable and checkable ✘	Indirect, via probing or attribution
Control	Easier to guide, inspect, and correct	No built-in mechanism to steer or verify ?
Reasoning Diversity	Commits to one trajectory	Encodes multiple alternatives
Supervision Granularity	Explicit, step-aware supervision	Guided by latent objectives
Alignment with Human Thinking	Explains thoughts aloud	Thinking silently

Future directions: Open Problems

- Limited Interpretability and Latent Opacity
 - probing and attribution techniques getting less effective with time
- Limited Control and Reliability
 - Models fail silently
- Performance Gap Compared to Explicit Reasoning in some tasks
 - Latent Reasoning relies on shortcut heuristics (before *grokking*)
 - Hybrid strategies? How to profit from the best of each configuration?
- Lack of Standardized Evaluation
 - difficult to diagnose failure modes or distinguish genuine reasoning from shallow heuristics
 - lack of fair comparability and reproducibility
 - latent annotations, standardized probing protocols, and metrics that assess internal consistency, trajectory depth, and robustness to distributional shifts?

Future directions: Open Problems

- Architecture and Generalization Constraints
 - Most works on small LMs
 - Need for architecture-specific components, such as loop controllers and task-specific latent heads
 - Future work: architecture-agnostic designs and clear reasoning objectives
- Dependence on Explicit Supervision
 - Most methods still need Explicit CoT to guide training
 - restricts independence and scalability
 - Future research: supervision signals operating directly on latent trajectories
 - self-consistency constraints
 - implicit verification objectives
 - unsupervised discovery of latent reasoning structures

Thanks!

Questions?

References

Much of this literature review was inspired by these surveys:

- Li et al., [Implicit Reasoning in Large Language Models: A Comprehensive Survey](#) (Sep, 2025)
- Zhu et al. [A Survey on Latent Reasoning](#) (July, 2025)
- Chen et al., [Reasoning Beyond Language: A Comprehensive Survey on Latent Chain-of-Thought Reasoning](#) (Nov, 2025)
- Pan et al., [Opening the Black Box: A Survey on the Mechanisms of Multi-Step Reasoning in Large Language Models](#) (Jan, 2026)