

Training LLM to Reason in a Continuous Latent Space

Menelik Nouvellon

some background & motivation: Why Continuous (Latent) Reasoning?

Chain-of-Thought (CoT) Reasoning: What & Why

Chain-of-Thought Prompting Model Input Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now? A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11. Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9.

Recent examples of CoT (extension of CoT)

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

Abstract

We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-turing (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1 achieves on Queen and Llama.



Figure 1 | Benchmark performance of DeepSeek-R1.

So why do we need latent continuous reasoning?

Motivation #1: Information Loss in Text Decoding





Motivation #2: Avoiding Filler & Overhead

Question:

"The cafeteria had 23 apples. They used 20 to make lunch, then bought 6 more. How many apples do they have now?"

Reasoning:

 The cafeteria started with 23 apples.
 They used 20 apples for lunch, leaving 3 apples.
 Next, they bought 6 more apples, so 3 + 6 = 9. Hence, the final answer is 9. filler phrases

Core math: All that's *really* needed is: "23 - 20 = 3; 3 + 6 = 9."

Motivation #3: Parallel Exploration & Backprop





Example : Text: Early Commitment **Continuous: Parallel**





Motivation #4: Better Generalization



Textual reasoning

Continuous reasoning

- 1. **Text CoT** is good but has overhead, can be rigid.
- 2. Latent Reasoning might preserve richer signals, allow branching, skip filler.
- 3. **End-to-End** optimization could push better multi-step performance.

Now let's dive into the paper

Coconut: "<u>C</u>hain <u>o</u>f <u>Con</u>tin<u>u</u>ous <u>T</u>hought"



Inference in Coconut: Big Picture



Where do we put the <eot> ?

aka Where should the latent reasoning stop ?



Step-by-Step Example

Biff the Bear buys 3 honey pots. Each honey pot costs 5 honey coins.

Question: "How many honey coins does Biff pay in total?"



Normal CoT inference

Reasoning:

1) Biff buys 3 pots, each costs 5 honey coins.

2) Multiply 3 by 5, that equals 15.

Therefore, Biff must pay 15 honey coins total.

Coconut Inference

<bot> [... latent steps ...] <eot> Final Answer: 15 honey coins

Now let's look at the training process

finetuned with CoT instances

Baseline : pre-trained GPT 2

GPT-2 with Textual CoT

finetuned with CoT instances

GPT 2

Baseline : pre-trained GPT 2

GPT-2 with Textual CoT

<u>AKA</u>

Datasets of <problems> : <ReasoningSteps><Answer>

=

GPT-2 with Textual CoT



Language ((training do	[Question] [Step 1] [Step 2] [Step 3] … [Step N] [Answer]	ught] : continuous thought […] : sequence of tokens
Stage 0	[Ouestion] <bot> <eot> [Sten 1] [Sten 2] ··· [Sten N] [Answer]</eot></bot>	special token
Stuge 0		<u>····</u> : calculating loss
Stage 1	[Question] <bot> [Thought] <eot> [Step 2] [Step 3] … [Step N] [Answer]</eot></bot>	
Stage 2	[Question] <bot> [Thought] [Thought] <eot> [Step 3] … [Step N] [Answer]</eot></bot>	
Stage N	[Question] <bot> [Thought] [Thought] … [Thought] <eot> [Answer]</eot></bot>	

Why Not Jump Immediately to All Latent Steps?



Why do use multi-stage curriculum ?



The LLM still needs guidance to learn latent reasoning

- 1. **Inference:** hidden steps, all in continuous latent space
- 2. **Training:** multi-stage replacement of textual steps -- *this is called multi-stage curriculum*
- 3. Architecture: special tokens <bot> and <eot> mark the boundaries of latent reasoning,

What about their result & experiment?

Experimental Setup

<u>Math</u>

<u>GSM8k</u>: dataset of 8.5K grade school math word problems created by human problem writers

<problem> : <ReasoningSteps><finalAnswer>

Experimental Setup

Math

<problem> : <ReasoningSteps><finalAnswer>

Logical Reasoning







simple chain

ProntoQA : "Stella is a zumpus. Zumpuses are gorpuses... Is Stella floral?" **ProsQA :** "Tom is a terpus. Every terpus is a brimpus..."

Now the match that everyone have been waiting for...





Mothod	GS	M8k	Pron	toQA	Pro	sQA
Method	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens
СоТ	$42.9{\scriptstyle~\pm 0.2}$	25.0	$98.8{\scriptstyle~\pm 0.8}$	92.5	$77.5{\scriptstyle~\pm1.9}$	49.4
COCONUT (Ours)	34.1 ± 1.5	8.2	$99.8{\scriptstyle~\pm 0.2}$	9.0	$97.0{\scriptstyle~\pm 0.3}$	14.2

Math

Mathad	GS	M8k	Pron	toQA	ProsQA						
Method	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens					
СоТ	$42.9{\scriptstyle~\pm 0.2}$	25.0	$98.8{\scriptstyle~\pm 0.8}$	92.5	$77.5{\scriptstyle~\pm1.9}$	49.4					
COCONUT (Ours)	34.1 ± 1.5	8.2	$99.8{\scriptstyle~\pm 0.2}$	9.0	$97.0{\scriptstyle~\pm 0.3}$	14.2					
	CoT leads on a	ccuracy									
	D. (O. see (se										
	with 4x less tok	ems promising ens									

				Logical Re	easoning							
Mathad	GS	M8k	Pron	toQA	Pro	sQA						
Method	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens						
СоТ	$42.9{\scriptstyle~\pm 0.2}$	25.0	98.8 ± 0.8	92.5	$77.5{\scriptstyle~\pm1.9}$	49.4						
COCONUT (Ours)	34.1 ± 1.5	8.2	$99.8{\scriptstyle~\pm 0.2}$	9.0	$97.0{\scriptstyle~\pm 0.3}$	14.2						
			СоТ									
			Coconut matche	s & outperform CoT								
			19.5% accuracy	19.5% accuracy improvement while		okens						

Mathad	GS	M8k	Pron	toQA	Pro	sQA
Methou	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens
СоТ	$42.9{\scriptstyle~\pm 0.2}$	25.0	$98.8{\scriptstyle~\pm 0.8}$	92.5	$77.5{\scriptstyle~\pm1.9}$	49.4
COCONUT (Ours)	34.1 ± 1.5	8.2	$99.8{\scriptstyle~\pm 0.2}$	9.0	$97.0{\scriptstyle~\pm 0.3}$	14.2



Now let's look at how Coconut compares to other reasoning approaches beyond just CoT...

iCo	Г																
	_		i	nput				1		С	σТ				C	Dutpu	ıt
Explicit C	CoT Stage 0:	2	1	×	4	3	=	8	4	+	0	6	3	=	8	0	4
	Stage 1:	2	1	×	4	3	=		4	+	0	6	3	=	8	0	4
	Stage 2:	2	1	×	4	3	=			+	0	6	3	=	8	0	4
	Stage 3:	2	1	×	4	3	=				0	6	3	=	8	0	4
	Stage 4:	2	1	×	4	3	=					6	3	=	8	0	4
	Stage 5:	2	1	×	4	3	=						3	=	8	0	4
↓ Implicit C	CoT Stage 6:	2	1	×	4	3	=							=	8	0	4
[3] From Explici	it CoT to Implicit C	oT: Lear	rning to	o Interr	nalize (CoT Ste	ep by Ste	ep (2024)	https:/	/doi.org	2/10.48	550/ar	Xiv.2405	.14838			

				Input				1			Co	т					С	Outpu	ıt
Explicit CoT	Stage 0:	2	1	×	4	3	=		8	4	+	0	6	3	=	=	8	0	4
	Stage 1:	2	1	×	4	3	=			4	+	0	6	3	=	-	8	0	4
	Stage 2:	2	1	×	4	3	=				+	0	6	3	=	=	8	0	4
	Stage 3:	2	1	×	4	3	=					0	6	3	=	=	8	0	4
	Stage 4:	2	1	×	4	3	=						6	3	=	=	8	0	4
	Stage 5:	2	1	×	4	3	=							3	=	=	8	0	4
Implicit CoT	Stage 6:	2	1	×	4	3	=								=	=	8	0	4
[3] From Explicit CoT	to Implicit Co	T: Lea	rning t	o Interr	alize (CoT St	ep by St	ер (2	2024)	ittps://	doi.org	/10.48	550/ar)	(iv.2405.	14838				

iCoT

Pause Token



[4] Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. Think before you speak: Training language models with pause tokens. (2023) <u>https://doi.org/10.48550/arXiv.2310.02226</u>

Mothod	GSM8k							
Method	Acc. (%)	# Tokens						
iCoT	30.0^{*}	2.2						
Pause Token	$16.4{\scriptstyle~\pm1.8}$	2.2						
COCONUT (Ours)	$34.1{\scriptstyle~\pm1.5}$	8.2						



Figure 3 Accuracy on GSM8k with different number of continuous thoughts.

"Chaining" continuous thoughts enhances reasoning



Figure 4 A case study where we decode the continuous thought into language tokens.

Continuous thoughts are efficient representations of reasoning

No-CoT	iCoT
it's like CoT but no CoT	Input CoT Output
	Explicit CoT Stage 0: $2 \ 1 \ \times \ 4 \ 3 = 8 \ 4 \ + \ 0 \ 6 \ 3 = 8 \ 0 \ 4$
trained on GSM8k	Stage 1: 2 1 \times 4 3 = 4 + 0 6 3 = 8 0 4
	Stage 2: 2 1 × 4 3 = + 0 6 3 = 8 0 4
<problem> → <reasoningsteps><finalanswer></finalanswer></reasoningsteps></problem>	Stage 3: 2 1 × 4 3 = 0 6 3 = 8 0 4
	Stage 4: 2 1 × 4 3 = 6 3 = 8 0 4
<problem> \rightarrow <finalanswer></finalanswer></problem>	Stage 5: $2 \ 1 \ \times \ 4 \ 3 = 3 = 8 \ 0 \ 4$
	Implicit CoT Stage 6: 2 1 \times 4 3 = = 8 0 4
	[3] From Explicit CoT to Implicit CoT: Learning to Internalize CoT Step by Step (2024) https://doi.org/10.48550/arXiv.2405.14838
Pause Token	
Ignore Ignore Ignore Ignore Output $25+$ 4 is 29 Layer 2 Image: Second se	
Inputs 5 ² + 4 is	
(a) Standard interence and finetuning (b) Pause-inference and finetuning	

Mathad	Pro	sQA
Method -	Acc. (%)	# Tokens
CoT	$77.5{\scriptstyle~\pm1.9}$	49.4
No-CoT	$76.7{\scriptstyle~\pm1.0}$	8.2
iCoT	$98.2{\scriptstyle~\pm 0.3}$	8.2
Pause Token	$75.9{\scriptstyle~\pm 0.7}$	8.2
Coconut (Ours)	$97.0{\scriptstyle~\pm 0.3}$	14.2
- w/o curriculum	$76.1{\scriptstyle~\pm 0.2}$	14.2
- w/o thought	$95.5{\scriptstyle~\pm1.1}$	8.2
- pause as thought	$96.6{\scriptstyle~\pm 0.8}$	8.2

Latent reasoning outperforms language reasoning in planning-intensive tasks



Multiple concepts have significant probability

clear convergence

Latent reasoning outperforms language reasoning in planning-intensive tasks

Interpretability Trade-off



is there any other limitations?

Training stability issue



Possible Extensions

Pretraining with continuous thoughts

- Current approach relies on finetuning
- Could continuous thoughts be part of pretraining?
- Potential for more generalizable reasoning abilities

Hybrid approaches

- Combining language and latent reasoning
- "Generating the reasoning skeleton in language"
- "Completing the reasoning process in latent space"

Future Research Directions



Concluding Thoughts

- Key Takeaways
 - Coconut enables reasoning in continuous latent space
 - Shows emergent BFS-like search behavior
 - Improves efficiency (71% fewer tokens) while maintaining accuracy
- Critics
 - Still requires language supervision during training

Thank you

Q&A