Think Clearly: Improving Reasoning via Redundant Token Pruning



Daewon Choi¹ Jimin Lee² Jihoon Tack¹ Woomin Song^{1, 3} Saket Dingliwal³ Sai Muralidhar Jayanthi³ Bhavana Ganesh³ Jinwoo Shin¹ Aram Galstyan³ Sravan Babu Bodapati³

¹KAIST ²Korea University ³Amazon AGI

Contact: daeone0920@kaist.ac.kr

TL;DR: Just think clearly-Removing KV Cache from redundant tokens improves reasoning for free.

Not All Tokens Matter for Reasoning

1. Existence of redundant reasoning tokens



Experiments

Removing redundant tokens improves overall accuracy across mathematical benchmarks without any training ! * Parenthesis: response length

		Dataset						
Model	Method	MATH	Minerva	GaoKao	AIME2024	AIME2025	AMC2023	Average
Qwen2.5-7B	FullKV	87.0 (3397)	59.9 (3391)	65.8 (3845)	36.7 (7060)	23.3 (7133)	75.0 (5004)	57.9 (4971)
	Ours	87.2 (2926)	60.5 (3471)	67.1 (4219)	46.7 (6841)	36.7 (6905)	82.5 (4488)	63.4 (4808)
	FullKV	81.0	45.9	67.1	33.3	13.3	75.0	52.6

 ✓ Attention maps when the model fails to produce the correct answer (i.e., poor reasoning) and when it succeeds (i.e., good reasoning)
→ Poor reasoning leads to highly redundant attention patterns !

2. Attention score to </think>



Llama3.1-8B Ours		(3389)	(4060)	(4689)	(/06/)	(7088)	(4986)	(5213)
	Ouma	83.8	48.1	69.8	33.3	23.3	77.5	55.9
	Ours	(3345)	(3941)	(4532)	(7210)	(7375)	(4700)	(5183)

(Qwen2.5-7B / Llama3.1-8B are reasoning LLMs distilled from DeepSeek-R1) * We perform eviction per every 200 / 300 generation steps

Does token eviction itself leads to improved performance?

Score	AIME2024	AIME2025
FullKV	36.7	23.3
Random	36.7	33.3
H2O	40.0	26.7
Ours	46.7	36.7

(Random: uniformly at random H2O: lowest accumulated attention)

Our method can be effective
under aggressive compression.

	Compre	ession ratio
Method	25%	50%
FullKV	2	42.6
Streaming-LLM	35.4	39.4
H2O	34.6	39.2
Pyramid-Infer	30.6	40.0
Ours	36.0	40.2

Non-mathematical benchma	rk
--------------------------	----

Method	GPQA
FullKV	32.0 (6418)
Ours	36.4 (6277)

(GPQA: Multiple-choice science)

Component analysis				
Summ	Step	AIME2024	AMC2023	
×	X	40.0	70.0	
1	X	36.7	77.5	
1	1	46.7	82.5	

(Summ: self-summarization Step: step-aware token eviction)

Discussion points

- \checkmark Attention scores associated with the end-of-thinking token </think>.

Improving Reasoning with Redundant Token Pruning

Step 1. Identifying redundant tokens via self-summarization

During an intermediate step of decoding, forward the following summarization prompt:

"Time is up. Given the time I've spent and the approaches I've tried, I should stop thinking and now write summarization in one sentence.

Use end of thinking token !

Step 2. Step-aware eviction with hierarchical budget allocation

Aggregate importance score per reasoning step



Evict KV Cache of tokens from redundant steps

Given a token eviction budget k ,

Q1. Segment of reasoning steps

"Wait" "Alternatively" "Another angle" "Another approach" "But wait" "Hold on" "Hmm" "Maybe" "Looking back" "Okay" "Let me" "First" "Then" "Alright" "Compute" "Correct" "Good" "Got it" "I don't see any errors" "I think" "Let me double-check" "Let's see" "Now" "Remember" "Seems solid" "Similarly" "So" "Starting" "That's correct" "That seems right" "Therefore" "Thus"

Q2. Connection to overthinking

	AIME24	AMC23
FullKV	36.7	75.0
Chain of Draft	23.3	72.5
Break the Chain	23.3	72.5
Ours	46.7	$\boldsymbol{82.5}$

Overthinking methods hurt performance by removing output redundancy, while we improve accuracy by targeting internal redundancy.

Q3. Which tokens are frequently evicted?

Token	Frequency (normalized)
, (comma)	1.00
2	0.97
" " (blank quote)	0.84
1	0.62
4	0.49
. (full stop)	0.46

They are punctuation marks / numbers that are contextually redundant.





Token importance score

Attention weight from the </think> !