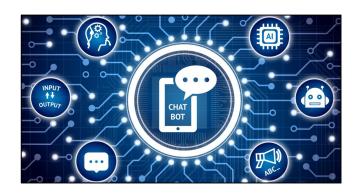
# Practical Algorithms from in-depth understanding of LLM behaviors

Yuandong Tian
Research Scientist Director

Meta Superintelligence Lab (FAIR)



### Large Language Models (LLMs)



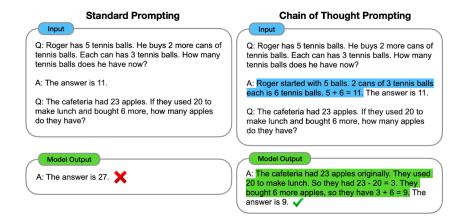
Conversational AI



**Content Generation** 



Al Agents





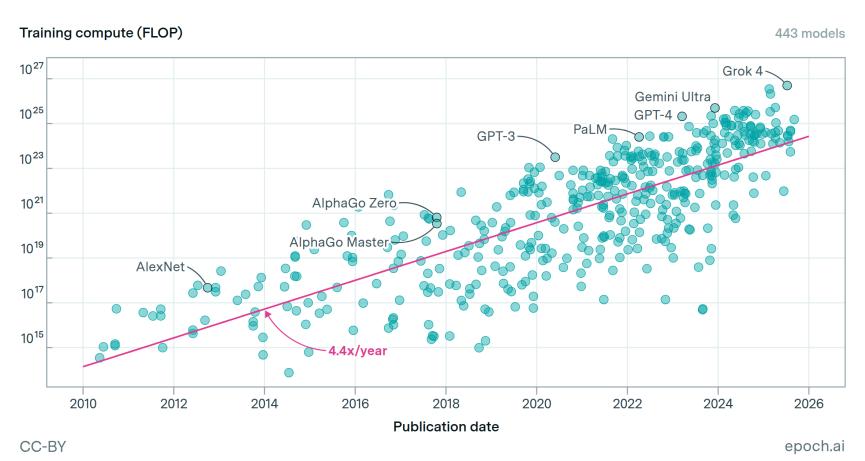


**Planning** 

# The Progress of Large Models

Training compute of notable models



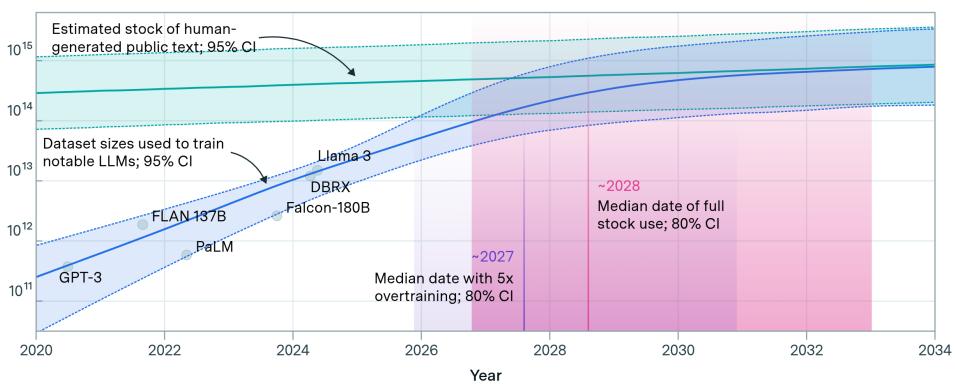


# The Data Usage

#### Projections of the stock of public text and data usage



Effective stock (number of tokens)



CC-BY

epoch.ai

#### Comparison between Human and SoTA LMs

Question: Is our AI as strong as humans yet?

	Training Data efficiency	Power Consumption	Adaptation to New Tasks	How to make decision?
Human Brain	< 10B text tokens, a lot of sensory inputs	Learning: ~20W Inference/Thinking: ~20W	Learn with a few examples	By casual relationships and deep understanding
Sota LMs	~10T-50T tokens	Learning: at least @ MWh Inference/Thinking: 1W-30W	Hundreds / Thousands of data points. May fail to generalize	Correlation & Pattern Matching

Estimated #tokens consumed by human in the life time: 70 years \* 300 days / year \* 12 hours / day \* 3600 seconds / hours \* 10 tokens / second = 9.1B



My pleasure to come on Dwarkesh last week, I thought the questions and conversation were really good.

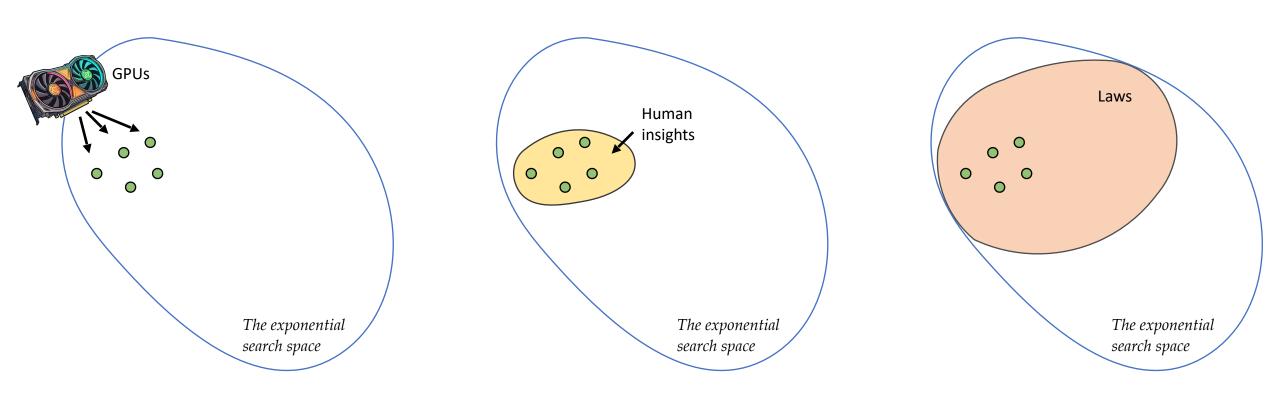
I re-watched the pod just now too. First of all, yes I know, and I'm sorry that I speak so fast:). It's to my detriment because sometimes my speaking thread out-executes my thinking thread, so I think I botched a few explanations due to that, and sometimes I was also nervous that I'm going too much on a tangent or too deep into something relatively spurious. Anyway, a few notes/pointers:

AGI timelines. My comments on AGI timelines looks to be the most trending part of the early response. This is the "decade of agents" is a reference to this earlier tweet x.com/karpathy/statu... Basically my AI timelines are about 5-10X pessimistic w.r.t. what you'll find in your neighborhood SF AI house party or on your twitter timeline, but still quite optimistic w.r.t. a rising tide of AI deniers and skeptics. The apparent conflict is not: imo we simultaneously 1) saw a huge amount of progress in recent years with LLMs while 2) there is still a lot of work remaining (grunt work, integration work, sensors and actuators to the physical world, societal work, safety and security work (jailbreaks, poisoning, etc.)) and also research to get done before we have an entity that you'd prefer to hire over a person for an arbitrary job in the world. I think that overall, 10 years should otherwise be a very bullish timeline for AGI, it's only in contrast to present hype that it doesn't feel that way.

#### How we should do our research from now on?

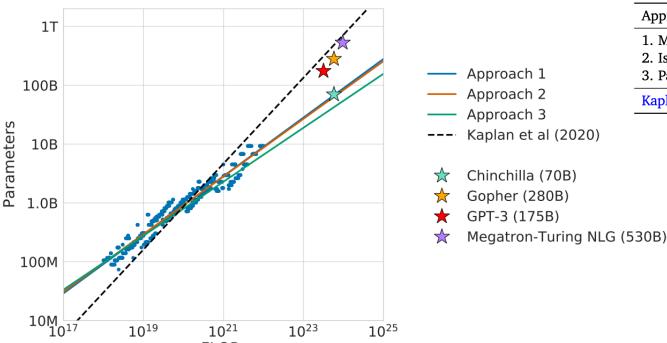
- The "Data Wall problem"
  - We may have used all the available data on the Internet.
  - How to deal with corner cases / personalization / private data?
  - Human is still much more efficient than current Al
- Everyone is GPU poor
  - What are new axes to scale? GPUs are never enough.
  - Data itself cannot extrapolate, only human insights can.

# The New (a.k.a. Old) Scaling Axis



Question: Can we **scale** the scaling laws?

## How we get Scaling Laws?



 $10^{23}$ 

 $10^{25}$ 

Approach	Coeff. <i>a</i> where $N_{opt} \propto C^a$	Coeff. <i>b</i> where $D_{opt} \propto C^b$
1. Minimum over training curves	0.50 (0.488, 0.502)	0.50 (0.501, 0.512)
2. IsoFLOP profiles	0.49 (0.462, 0.534)	0.51 (0.483, 0.529)
3. Parametric modelling of the loss	0.46 (0.454, 0.455)	0.54 (0.542, 0.543)
Kaplan et al. (2020)	0.73	0.27

#### Steps:

- 1. Collect the experiments
- 2. Form hypothesis (linear, power-law, etc)
- 3. Extrapolate

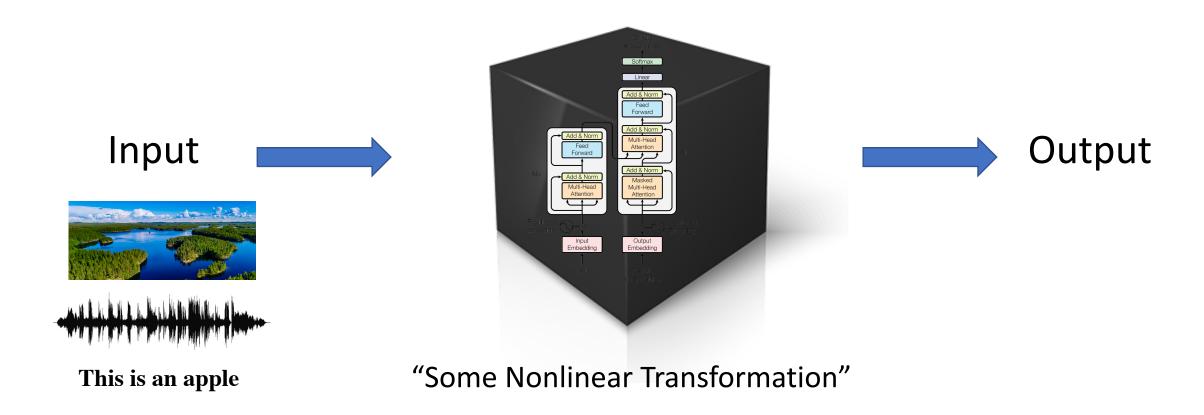
Still pure statistics and need exponential data. (No leverage of the knowledge of architecture/data)

 $10^{19}$ 

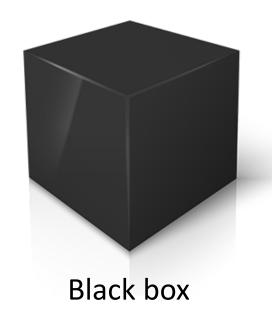
 $10^{21}$ 

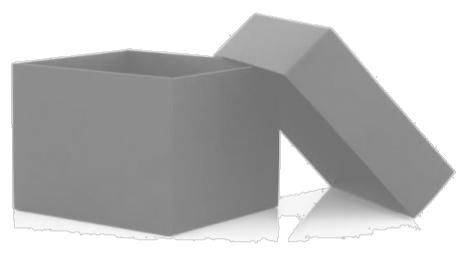
**FLOPs** 

## How does deep learning work?



#### Black-box versus White-box

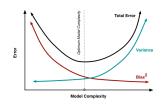




White box

#### What routes should we take?

#### Generalization



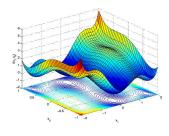
Architecture **X**training dynamics **X** 

#### **Expressibility**



Architecture **▼**training dynamics **×** 

#### **Optimization**

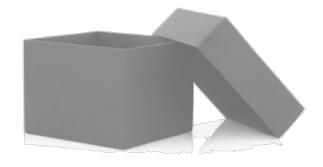


Architecture **X**training dynamics **√** 

#### **How about**

Architecture ✓
training dynamics ✓

### Start From the First Principle



Training follows Gradient and its variants (SGD, Adams, etc)

$$\dot{\boldsymbol{w}} \coloneqq \frac{\mathrm{d}\boldsymbol{w}}{\mathrm{d}t} = -\nabla_{\boldsymbol{w}} J(\boldsymbol{w})$$

- Sounds complicated.. Is that possible? Yes

Architecture ✓
training dynamics ✓

### What Gradient Descent gives us?

#### **Short-term:**

Finding Simple Structures (Low-rank, sparsity)

#### Long-term:

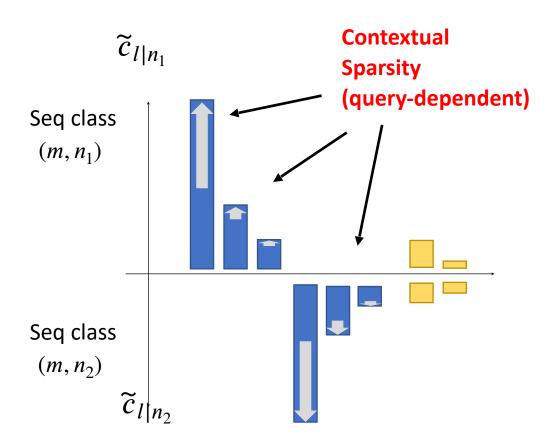
How the representation is learned (Key to the success of deep models)



Leverage Them in Practical Algorithms

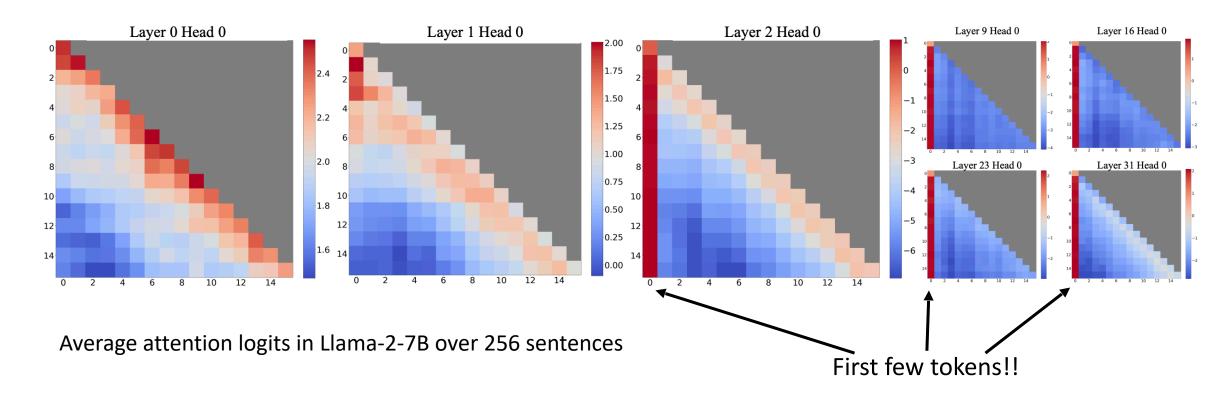
# Short-term: Finding Nice Structures

### Attention Sparsity



**Attention = Learnable** TF-IDF (Term Frequency, Inverse Document Frequency)

#### Attention Sinks: Initial tokens draw a lot of attentions



- Observation: Initial tokens have large attention scores, even if they're not semantically significant.
- Attention Sink: Tokens that disproportionately attract attention irrespective of their relevance.

### Understanding Attention Sinks

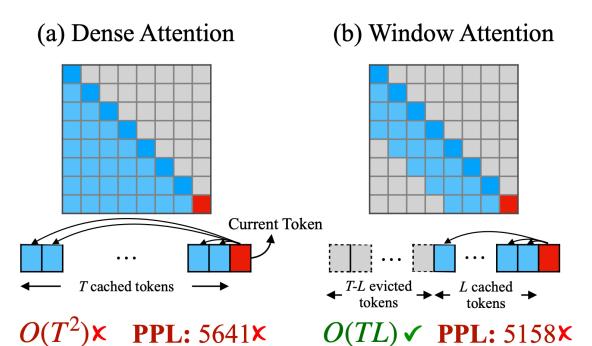
• Why? Attention scores have to sum up to 1 for all contextual tokens. (SoftMax-Off-by-One, Miller et al. 2023)

SoftMax
$$(x)_i = \frac{e^{x_i}}{e^{x_1} + \sum_{j=2}^{N} e^{x_j}}, \quad x_1 \gg x_j, j \in 2, \dots, N$$

- Why initial tokens? Their visibility to subsequent tokens, rooted in autoregressive language modeling.
- The model learns a bias towards their **absolute position** rather than the semantics.

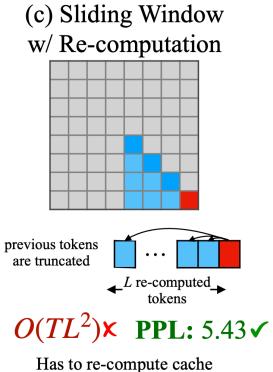
Llama-2-13B	PPL (↓)
0+1024 (window)	5158.07
4+1024	5.40
4"\n"+1020	5.6

### StreamingLLM

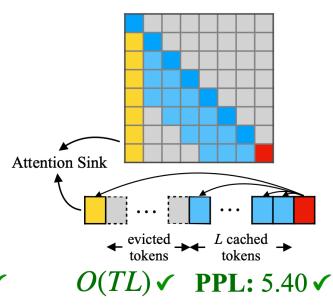


Breaks when initial

tokens are evicted.



for each incoming token.



Can perform efficient and stable

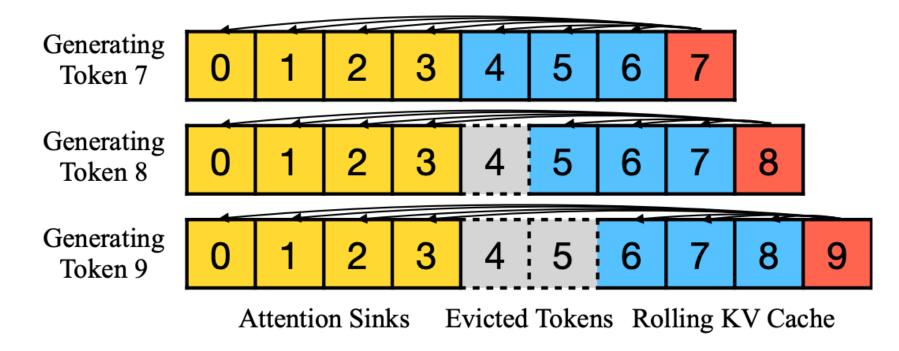
language modeling on long texts.

(d) StreamingLLM (ours)

Has poor efficiency and

performance on long text.

### StreamingLLM



#### **Key design: Position Rolling**

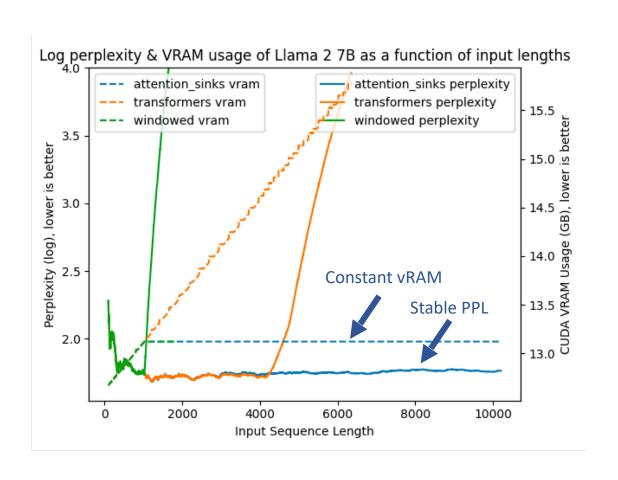
For all tokens, use their positions within cache to compute positional encoding!

→ Token distance never exceeds pre-trained context window!

# StreamingLLM

# w/o StreamingLLM w/ StreamingLLM (streaming) guangxuan@l29:~/workspace/streaming-llm\$ CUDA\_VISIBLE\_DEVICE|(streaming) guangxuan@l29:~/workspace/streaming-llm\$ CUDA\_VISIBLE\_DEVICES=1 py thon examples/run\_streaming\_llama.py --enable\_streaming\_ Loading model from lmsys/vicuna-13b-v1.3 ... S=0 python examples/run\_streaming\_llama.py Loading model from lmsys/vicuna-13b-v1.3 ... Loading checkpoint shards: 67% | 2/3 [00:09<00:04, 4.94s/it] Loading checkpoint shards: 67%| | 2/3 [00:09<00:04, 4.89s/it]

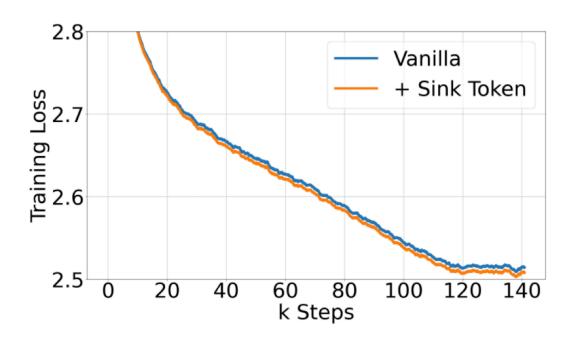
### StreamingLLM: stable PPL, constant vRAM



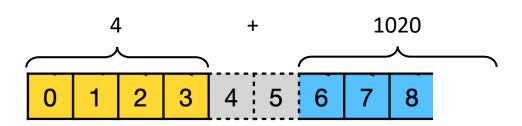


### Understanding Attention Sinks

Pre-train with a Dedicated Attention Sink Token



Cache	PPL (↓)					
Config	0+1024	1+1023	2+1022	4+1020		
Vanilla	27.87	18.49	18.05	18.05		
Zero Sink	29214	19.90	18.27	18.01		
Learnable Sink	1235	18.01	18.01	18.02		



#### Impact

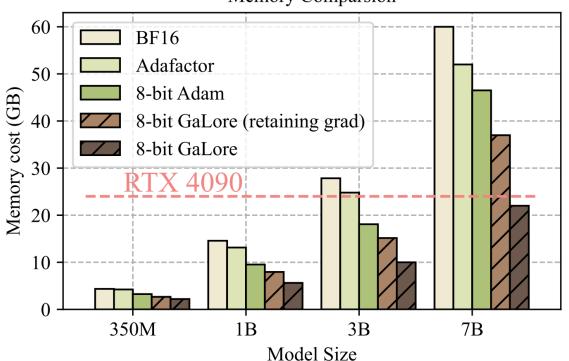
Attention: Following GPT-3, attention blocks alternate between banded window and fully dense patterns [8][9], where the bandwidth is 128 tokens. Each layer has 64 query heads of dimension 64, and uses Grouped Query Attention (GQA [10][11]) with 8 key-value heads. We apply rotary position embeddings [12] and extend the context length of dense layers to 131,072 tokens using YaRN [13]. Each attention head has a learned bias in the denominator of the softmax, similar to off-by-one attention and attention sinks [14][15], which enables the attention mechanism to pay no attention to any tokens.

- 900+ citations
- Used in GPT OSS models in pre-training

## GaLore: Pre-training 7B model on RTX 4090 (24G)







	Rank	Retain grad	Memory	Token/s
8-bit AdamW		Yes	40GB	1434
8-bit GaLore	16	Yes	28GB	1532
8-bit GaLore	128	Yes	29GB	1532
16-bit GaLore	128	Yes	30GB	1615
16-bit GaLore	128	No	18GB	1587
8-bit GaLore	1024	Yes	36GB	1238

<sup>\*</sup> SVD takes around 10min for 7B model, but runs every T=500-1000 steps.

Third-party evaluation by @llamafactory\_ai



## Memory Saving with GaLore

#### **Algorithm 1:** GaLore, PyTorch-like

```
for weight in model.parameters():
   grad = weight.grad
   # original space -> compact space
   lor_grad = project(grad)
   # update by Adam, Adafactor, etc.
   lor_update = update(lor_grad)
   # compact space -> original space
   update = project_back(lor_update)
   weight.data += update
```

#### GaLore

$$\begin{aligned} G_t &\leftarrow -\nabla_{\mathbf{W}} \phi(W_t) \\ \text{If t \% T == 0:} \\ \text{Compute } P_t &= \text{SVD} \big( G_t \big) \in \mathbb{R}^{m \times r} \\ R_t &\leftarrow P_t^T G_t \quad \{ project \} \\ \widetilde{R}_t &\leftarrow \rho \big( R_t \big) \quad \{ Adam \ in \ low-rank \} \\ \widetilde{G}_t &\leftarrow P_t \widetilde{R}_t \quad \{ project-back \} \\ W_{t+1} &\leftarrow W_t + \eta \widetilde{G}_t \end{aligned}$$

Memory Usage	Weight	Optim States	Projection	Total
Full-rank	mn	2mn	0	3mn
Low-rank adaptor	mn+mr+nr	2mr+2nr	0	2mn+3mr+3nr
GaLore	mn	mr+2nr	mr	mn+2mr+2nr
rtificial Intelligence	$W_t$	$\stackrel{1}{R_t}$	$P_t$	

### Pre-training Results (LLaMA 7B)

Params	Hidden	Intermediate	Heads	Layers	Steps	Data amount
60M	512	1376	8	8	10K	1.3 B
130M	768	2048	12	12	20K	$2.6~\mathrm{B}$
350M	1024	2736	16	24	60K	$7.8\mathrm{B}$
1 B	2048	5461	24	32	100K	$13.1~\mathrm{B}$
$7\mathrm{B}$	4096	11008	32	32	150K	$19.7~\mathrm{B}$

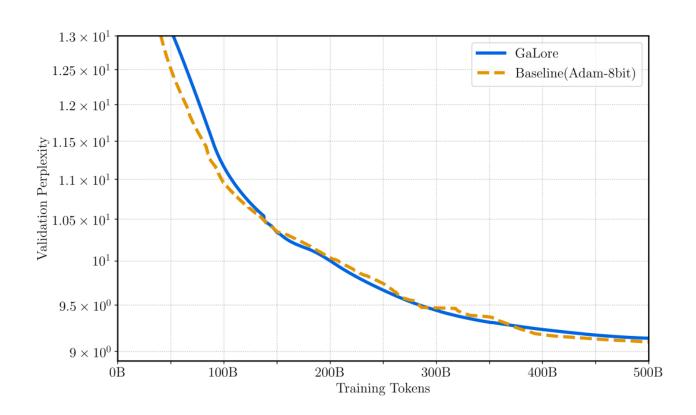
		Mem	40K	80K	120K	150K
<b>©</b>	<b>8-bit GaLore</b> 8-bit Adam	18 <b>G</b>	17.94	15.39	14.95	14.65
_	8-bit Adam	26G	18.09	15.47	14.83	14.61
_	Tokens (B)		5.2	10.5	15.7	19.7

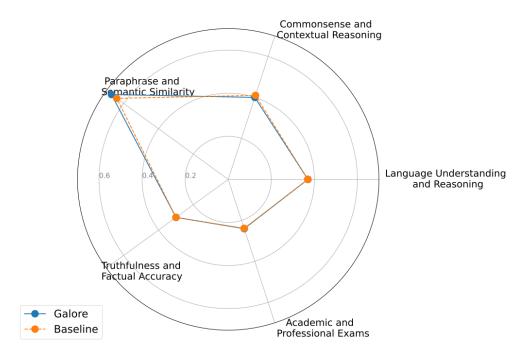
<sup>\*</sup> Experiments are conducted on 8 x 8 A100

	60M	130M	350M	1B
Full-Rank	34.06 (0.36G)	25.08 (0.76G)	18.80 (2.06G)	15.56 (7.80G)
GaLore	<b>34.88</b> (0.24G)	<b>25.36</b> (0.52G)	<b>18.95</b> (1.22G)	<b>15.64</b> (4.38G)
Low-Rank	78.18 (0.26G)	45.51 (0.54G)	37.41 (1.08G)	142.53 (3.57G)
LoRA	34.99 (0.36G)	33.92 (0.80G)	25.58 (1.76G)	19.21 (6.17G)
ReLoRA	37.04 (0.36G)	29.37 (0.80G)	29.08 (1.76G)	18.33 (6.17G)
$r/d_{model}$	128 / 256	256 / 768	256 / 1024	512 / 2048
Training Tokens	1.1B	2.2B	6.4B	13.1B

<sup>\*</sup> On LLaMA 1B, ppl is better (~14.97) with ½ rank (1024/2048)

## Pre-training Results (LLaMA 7B)





# Long-term: How Network finds Representation

## Type of Representations

(Traditional) Symbolic representation

 $\nabla \cdot \mathbf{E} = \frac{\rho_{v}}{\varepsilon} \qquad (Gauss' Law)$ 

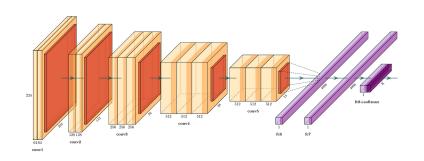
 $\nabla \cdot \mathbf{H} = 0$  (Gauss'Law for Magnetism)

 $\nabla \times \mathbf{E} = -\mu \frac{\partial \mathbf{H}}{\partial t}$  (Faraday's Law)

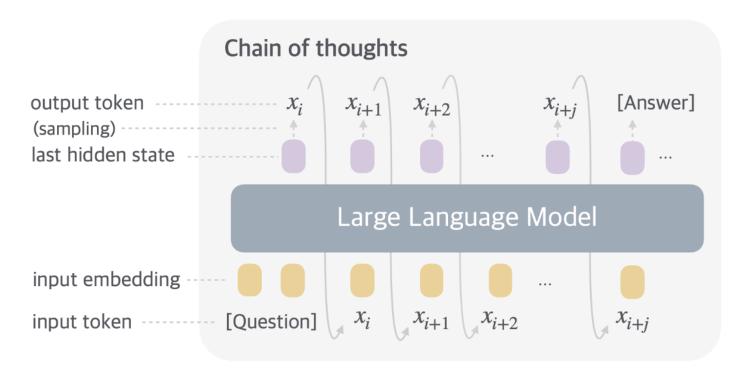
 $\nabla \times \mathbf{H} = \mathbf{J} + \varepsilon \frac{\partial \mathbf{E}}{\partial t}$  (Ampere's Law)

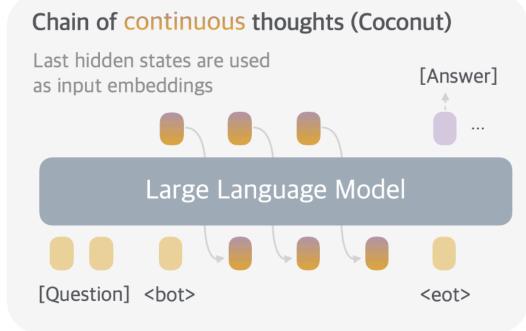
Representation

Neural Representation



## CoConut (Chain of Continuous Thought)















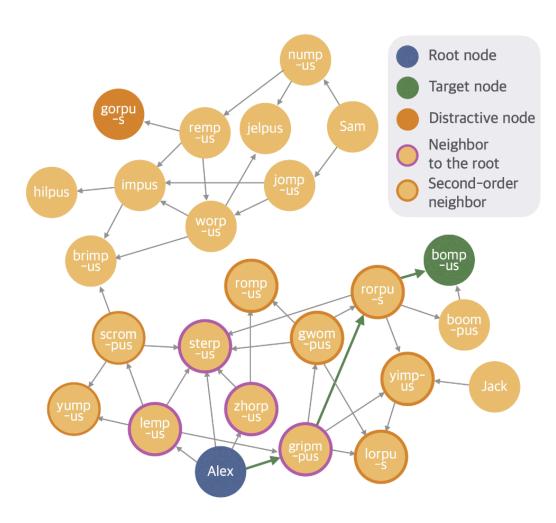




#### How to train Coconut?

```
Language CoT
                                                                                     [Thought] : continuous thought
                 [Question] [Step 1] [Step 2] [Step 3] ··· [Step N] [Answer]
(training data)
                                                                                          [ ··· ] : sequence of tokens
                                                                                                 <···> : special token
Stage 0
            [Question] <bot> <eot> [Step 1] [Step 2] ··· [Step N] [Answer]
                                                                                                ··· : calculating loss
            [Question] <bot> [Thought] <eot> [Step 2] [Step 3] ··· [Step N] [Answer]
Stage 1
Stage 2
            [Question] <bot> [Thought] (Thought] <eot> [Step 3] ··· [Step N] [Answer]
   ...
            [Question] <bot> [Thought] (Thought] ··· [Thought] <eot> [Answer]
Stage N
```

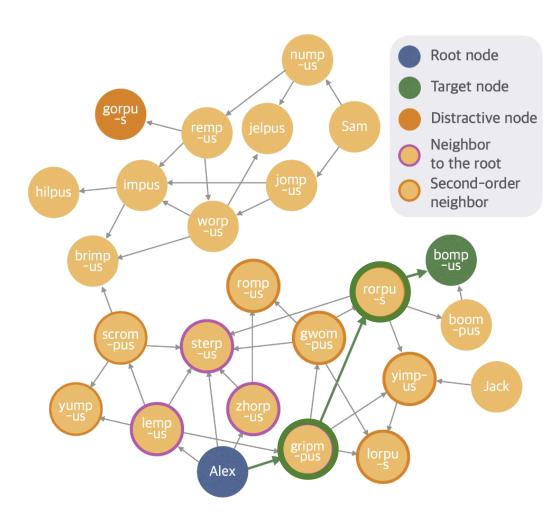
# Interpreting the embeddings



#### Question:

Every jells is a worpus. Sam is a jumpus. Every gwompus is a rompus. ··· Every lumps is a yumpus. Question: Is Alex a gorpus or bompus?

#### Ground Truth Solutions



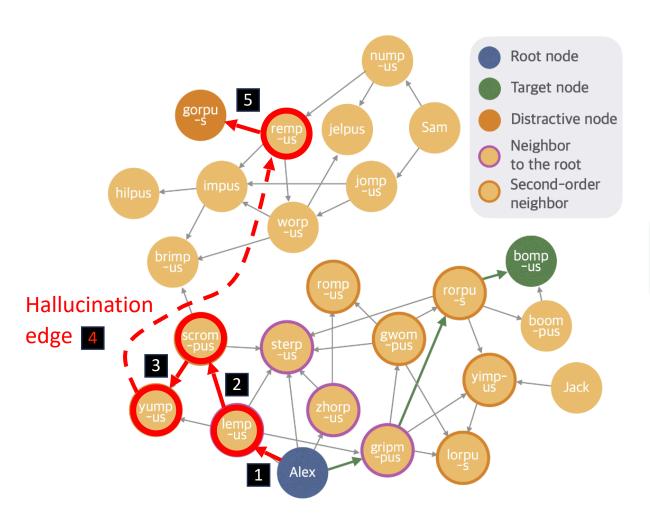
#### Question:

Every jells is a worpus. Sam is a jumpus. Every gwompus is a rompus. ··· Every lumps is a yumpus. Question: Is Alex a gorpus or bompus?

#### **Ground Truth Solution**

Alex is a grimpus. Every grimpus is a rorpus. Every rorpus is a bompus. ### Alex is a bompus

### Chain of thoughts lead to hallucinations



#### Question:

Every jells is a worpus. Sam is a jumpus. Every gwompus is a rompus. ··· Every lumps is a yumpus. Question: Is Alex a gorpus or bompus?

#### **Ground Truth Solution**

Alex is a grimpus. Every grimpus is a rorpus. Every rorpus is a bompus. ### Alex is a bompus

#### CoT

Alex is a lempus.

Every lempus is a scrompus.

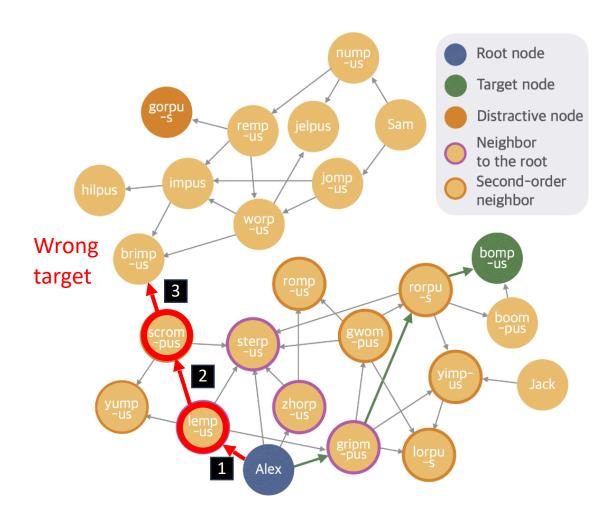
Every scrompus is a yumpus.

Every yumpus is a rempus.

Every rempus is a gorpus. ### Alex is a gorpus

(Hallucination)

### Continuous Thoughts



#### Question:

Every jells is a worpus. Sam is a jumpus. Every gwompus is a rompus. ··· Every lumps is a yumpus. Question: Is Alex a gorpus or bompus?

#### **Ground Truth Solution**

Alex is a grimpus.

Every grimpus is a rorpus.

Every rorpus is a bompus.

### Alex is a bompus

#### CoT

Alex is a lempus.

Every lempus is a scrompus.

Every scrompus is a yumpus.

Every yumpus is a rempus.

Every rempus is a gorpus.

### Alex is a gorpus

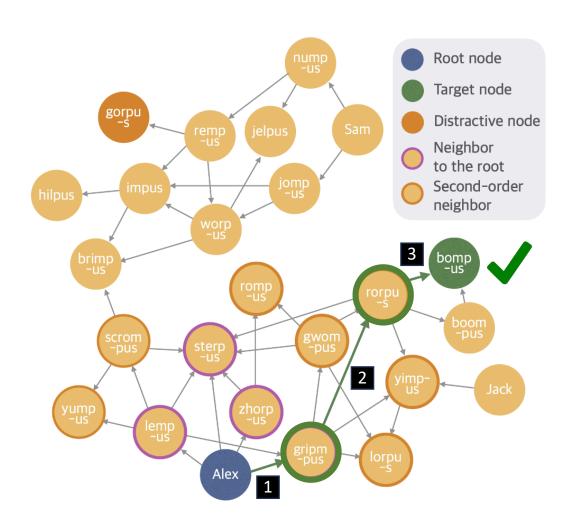
#### Ours (k=1)

<br/>
<br/>
<br/>
Every lempus is a scrompus. 2<br/>
Every scrompus is a brimpus. 3<br/>
### Alex is a brimpus

(Wrong Target)

(Hallucination)

### Two-step Continuous Thought works!



#### Question:

Every jells is a worpus. Sam is a jumpus. Every gwompus is a rompus. ··· Every lumps is a yumpus. Question: Is Alex a gorpus or bompus?

#### **Ground Truth Solution**

Alex is a grimpus. Every grimpus is a rorpus. Every rorpus is a bompus. ### Alex is a bompus

#### Ours (k=1)

<bot> [Thought] <eot>
Every lempus is a scrompus.
Every scrompus is a brimpus.
### Alex is a brimpus

(Wrong Target)

#### CoT

Alex is a lempus.

Every lempus is a scrompus.

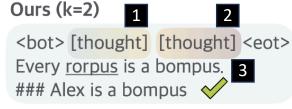
Every scrompus is a yumpus.

Every yumpus is a rempus.

Every rempus is a gorpus.

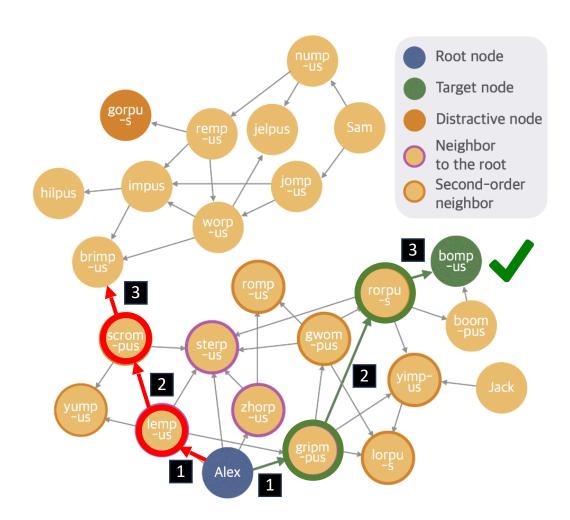
### Alex is a gorpus

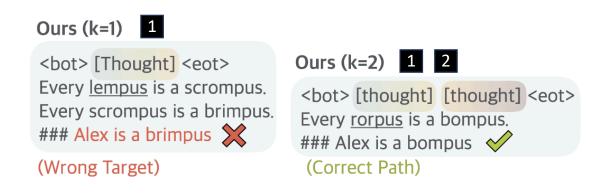
(Hallucination)



(Correct Path)

### Two-step Continuous Thought works!



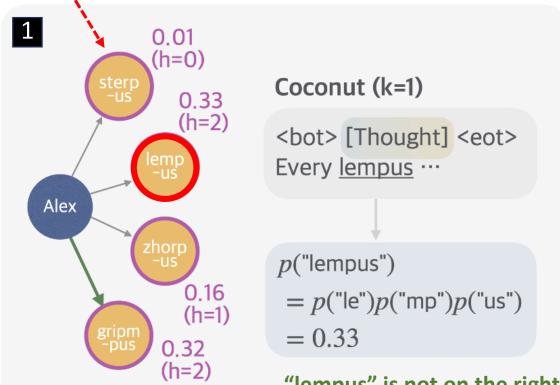


Why the same continuous thoughts 1 lead to different path?!

#### What's inside?

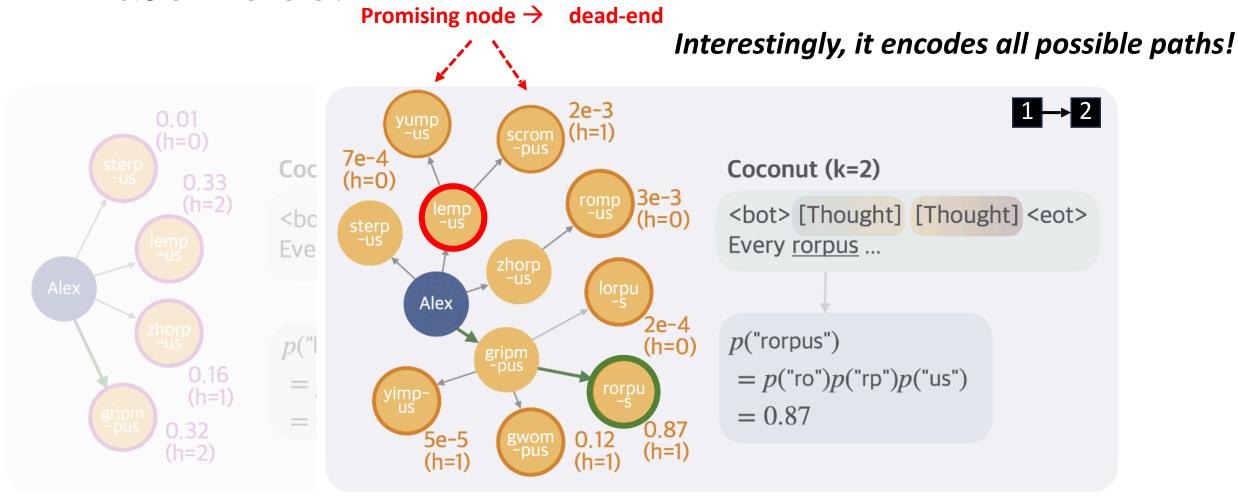
#### Let's probe!

#### **Dead-end**



"lempus" is not on the right path but for step=1, it is the most promising

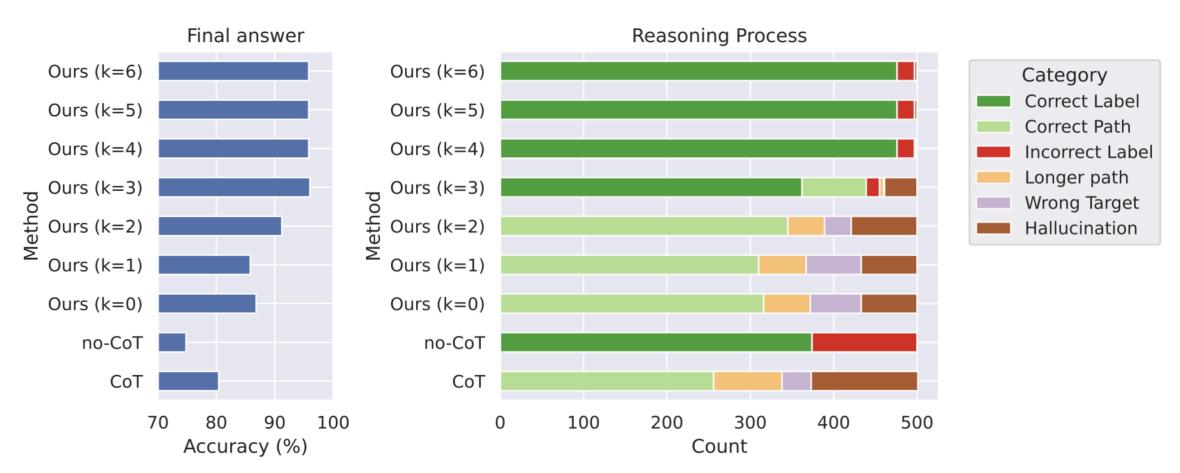
### What's inside?



### Performance in ProsQA

Dataset	Training	Validation	Test
GSM8k	385,620	500	1319
$\operatorname{ProntoQA}$	9,000	200	800
$\operatorname{ProsQA}$	17,886	300	500

# Nodes	$\mid$ # Edges	Len. of Shortest Path	# Shortest Paths
23.0	36.0	3.8	1.6



#### CoConut

Method	GS	M8k	Pron	toQA	Pro	$\operatorname{ProsQA}$		
Method	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens		
СоТ	$42.9 \pm 0.2$	25.0	$98.8 \pm 0.8$	92.5	$77.5 \pm 1.9$	49.4		
No-CoT iCoT Pause Token	$16.5 \pm 0.5 \ 30.0^* \ 16.4 \pm 1.8$	$2.2 \\ 2.2 \\ 2.2$	$93.8 \pm 0.7$ $99.8 \pm 0.3$ $77.7 \pm 21.0$	3.0 3.0 3.0	$76.7 \pm 1.0$ $98.2 \pm 0.3$ $75.9 \pm 0.7$	8.2 8.2 8.2		
COCONUT (Ours) - w/o curriculum - w/o thought - pause as thought	$34.1 \pm 1.5$ $14.4 \pm 0.8$ $21.6 \pm 0.5$ $24.1 \pm 0.7$	8.2 8.2 2.3 2.2	$99.8 \pm 0.2$ $52.4 \pm 0.4$ $99.9 \pm 0.1$ $100.0 \pm 0.1$	9.0 9.0 3.0 3.0	$97.0 \pm 0.3$ $76.1 \pm 0.2$ $95.5 \pm 1.1$ $96.6 \pm 0.8$	14.2 14.2 8.2 8.2		

Better performance than No-CoT Shorter thinking process than CoT

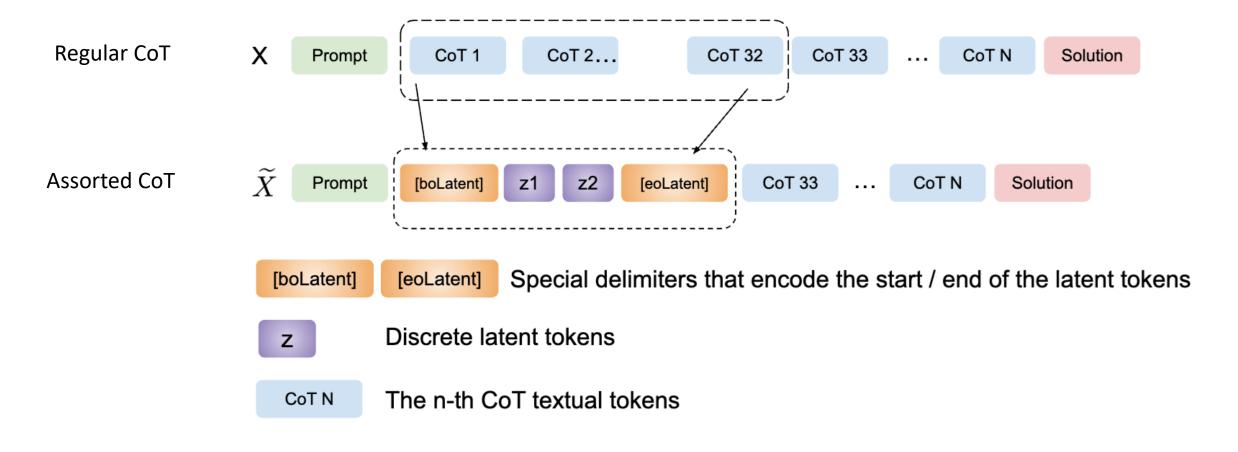
#### CoConut

Method	GS	M8k	Pron	toQA	$\operatorname{ProsQA}$		
Method	Acc. (%)	# Tokens	Acc. (%)	# Tokens	Acc. (%)	# Tokens	
СоТ	$42.9 \pm 0.2$	25.0	$98.8 \pm 0.8$	92.5	$77.5 \pm 1.9$	49.4	
No-CoT iCoT Pause Token	$16.5 \pm 0.5 \ 30.0^* \ 16.4 \pm 1.8$	2.2 2.2 2.2	$93.8 \pm 0.7$ $99.8 \pm 0.3$ $77.7 \pm 21.0$	3.0 3.0 3.0	$76.7 \pm 1.0$ $98.2 \pm 0.3$ $75.9 \pm 0.7$	8.2 8.2 8.2	
COCONUT (Ours) - $w/o$ curriculum - $w/o$ thought	$34.1 \pm 1.5$ $14.4 \pm 0.8$ $21.6 \pm 0.5$	8.2 8.2 2.3	$99.8 \pm 0.2$ $52.4 \pm 0.4$ $99.9 \pm 0.1$	9.0 9.0 3.0	$97.0 \pm 0.3$ $76.1 \pm 0.2$ $95.5 \pm 1.1$	14.2 14.2 8.2	
- pause as thought	$24.1 \pm 0.7$	2.2	Cons	ο Λ	00.0	0.0	

Better performance than No-CoT Shorter thinking process than CoT

- 1. Latent tokens are not interpretable
- 2. Only tested on GSM8k

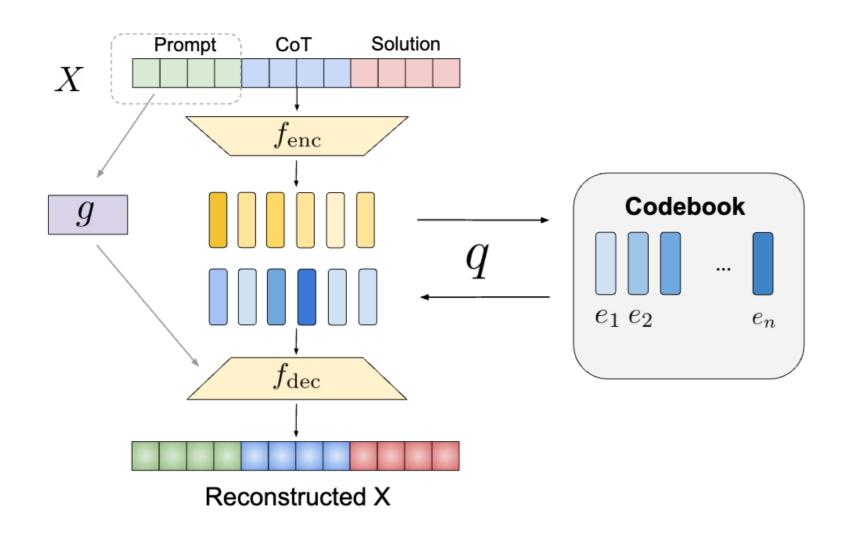
### Token Assorted



### Token Assorted

How the latent codes are constructed?

Using VQVAE



### Better Performance

Mo	del	In-Do	omain		0	ut-of-Domain			Average
1410	dei	Math	GSM8K	Gaokao-Math-2023	DM-Math	College-Math	Olympia-Math	TheoremQA	All Datasets
	Sol-Only	4.7	6.8	0.0	10.4	5.3	1.3	3.9	4.6
	CoT	<u>10.5</u>	<u>42.7</u>	10.0	3.4	<u>17.1</u>	1.5	9.8	<u>14.1</u>
Llama-3.2-1B	iCoT	8.2	10.5	3.3	<u>11.3</u>	7.6	2.1	<u>10.7</u>	7.7
Liama-3.2-1D	Pause Token	5.1	5.3	2.0	1.4	0.5	0.0	0.6	2.1
	Latent (ours)	<b>14.7</b> († <b>+4.2</b> )	<b>48.7</b> (↑ <b>+6</b> )	10.0	<b>14.6</b> (↑ +3.3)	<b>20.5</b> (↑ + <b>3.4</b> )	<u>1.8</u>	<b>11.3</b> († <b>+0.6</b> )	<b>17.8</b> (↑ <b>+3.7</b> )
	Sol-Only	6.1	8.1	3.3	14.0	7.0	1.8	6.8	6.7
	CoT	<u>21.9</u>	<u>69.7</u>	<u>16.7</u>	27.3	<u>30.9</u>	2.2	11.6	<u>25.2</u>
Llama-3.2-3B	iCoT	12.6	17.3	3.3	16.0	14.2	4.9	13.9	11.7
Liama-3,2-3D	Pause Token	25.2	53.7	4.1	7.4	11.8	0.7	1.0	14.8
	Latent (ours)	<b>26.1</b> († <b>+4.2</b> )	<b>73.8</b> (↑ <b>+4.1</b> )	23.3 († +6.6)	<u>27.1</u>	<b>32.9</b> (↑ <b>+2</b> )	<u>4.2</u>	<u>13.5</u>	<b>28.1</b> († <b>+2.9</b> )
	Sol-Only	11.5	11.8	3.3	17.4	13.0	3.8	6.7	9.6
	CoT	32.9	<u>80.1</u>	<u>16.7</u>	<u>39.3</u>	<u>41.9</u>	7.3	<u>15.8</u>	<u>33.4</u>
I lome 2 1 QD	iCoT	17.8	29.6	16.7	20.3	21.3	<u>7.6</u>	14.8	18.3
Llama-3.1-8B	Pause Token	39.6	79.5	6.1	25.4	25.1	1.3	4.0	25.9
	Latent (ours)	<u>37.2</u>	<b>84.1</b> (↑ <b>+4.0</b> )	<b>30.0</b> (↑ <b>+13.3</b> )	<b>41.3</b> (↑ <b>+2</b> )	<b>44.0</b> (↑ <b>+2.1</b> )	<b>10.2</b> (↑ <b>+2.6</b> )	<b>18.4</b> († <b>+2.6</b> )	<b>37.9</b> (↑ <b>+4.5</b> )

### Shorter CoT

Mo	del	In-Domain (	# of tokens)		Out-of	-Domain (# of tok	ens)		Average
1410	uci	Math	GSM8K	Gaokao-Math-2023	DM-Math	College-Math	Olympia-Math	TheoremQA	All Datasets
	Sol-Only	4.7	6.8	0.0	10.4	5.3	1.3	3.9	4.6
	CoT	646.1	190.3	842.3	578.7	505.6	1087.0	736.5	655.2
Llama-3.2-1B	iCoT	328.4	39.8	354.0	170.8	278.7	839.4	575.4	369.5
	Pause Token	638.8	176.4	416.1	579.9	193.8	471.9	988.1	495
	Latent (ours)	501.6 (↓ <b>-22</b> %)	181.3 (↓ <b>-5</b> %)	760.5 (↓ <b>-11%</b> )	380.1 (↓ <b>-34</b> %)	387.3 (↓ <b>-23</b> %)	840.0 (↓ <b>-22</b> %)	575.5 (↓ <b>-22</b> %)	518 (↓ <b>-21</b> %)
	Sol-Only	6.1	8.1	3.3	14.0	7.0	1.8	6.8	6.7
	CoT	649.9	212.1	823.3	392.8	495.9	1166.7	759.6	642.9
Llama-3.2-3B	iCoT	344.4	60.7	564.0	154.3	224.9	697.6	363.6	344.2
Liama-5.2-5D	Pause Token	307.9	162.3	108.9	251.5	500.96	959.5	212.8	354.7
	Latent (ours)	516.7 (↓ <b>-20</b> %)	198.8 (↓ <b>-6%</b> )	618.5 (↓ <b>-25</b> %)	340.0 (↓ <b>-13</b> %)	418.0 (↓ <b>-16%</b> )	832.8 (↓ <b>-29</b> %)	670.2 (↓ <b>-12</b> %)	513.6 (↓ <b>-20</b> %)
	Sol-Only	11.5	11.8	3.3	17.4	13.0	3.8	6.7	9.6
	CoT	624.3	209.5	555.9	321.8	474.3	1103.3	760.1	578.5
Llama 2 1 OD	iCoT	403.5	67.3	444.8	137.0	257.1	797.1	430.9	362.5
Llama-3.1-8B	Pause Token	469.4	119.0	752.6	413.4	357.3	648.2	600.1	480
	Latent (ours)	571.9 (↓ <b>-9</b> %)	193.9 (↓ <b>-8</b> %)	545.8 (↓ <b>-2</b> %)	292.1 (↓ <b>-10</b> %)	440.3 (↓ <b>-8%</b> )	913.7 (↓ <b>-17</b> %)	637.2 (↓ <b>-16</b> %)	513.7 (↓ <b>-10</b> %)

#### Main Theorem

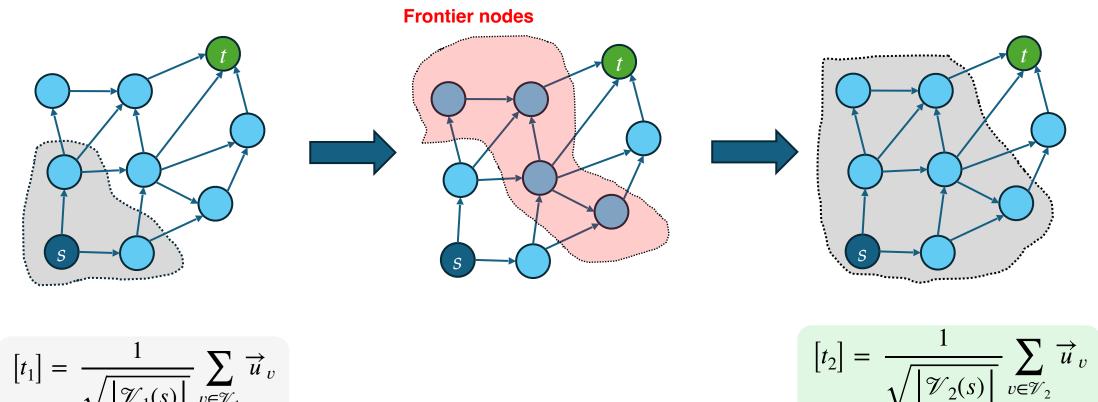
#### **Theorem (informal)**

For n-vertex directed graphs, a **2-layer** transformer with continuous CoT can solve reachability using O(n) decoding steps with O(n) embedding dimensions.

Best known results for discrete CoT:  $O(n^2)$ 

**Secret Sauce:** Superposition of the embeddings!

### Continuous CoT: Decoding as search



$$[t_1] = \frac{1}{\sqrt{|\mathcal{V}_1(s)|}} \sum_{v \in \mathcal{V}_1} \overrightarrow{u}_v$$

The embedding contains superposition

## Mechanism of Emerged Representation

(Traditional) Symbolic representation

#### Representation

#### 4 Conclusions

We have extended the GLU family of layers and proposed their use in Transformer. In a transfer-learning setup, the new variants seem to produce better perplexities for the de-noising objective used in pre-training, as well as better results on many downstream language-understanding tasks. These architectures are simple to implement, and have no apparent computational drawbacks. We offer no explanation as to why these architectures seem to work; we attribute their success, as all else, to divine benevolence.

Neural

Representation (9)

Shall we just acknowledge that as "divine benevolence"?

## Mechanism of Emerged Representation

(Traditional) Symbolic representation Representation Neural **Emerging Symbolic** Representation (§) Structure

### Modular Addition

$$a + b = c \mod d$$

Does neural network have an *implicit table* to do retrieval?

### Modular Addition

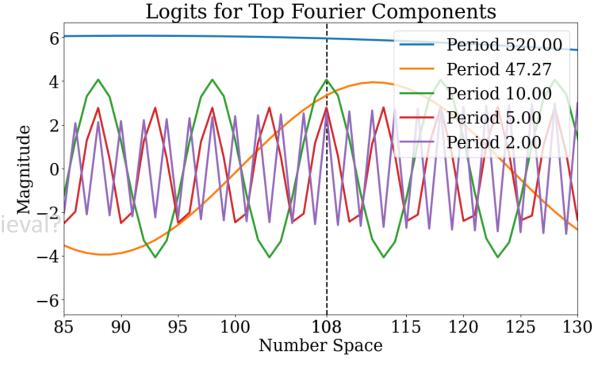
$$a + b = c \mod d$$

Does neural network have an *implicit table* to do retrieval

Learned representation = Fourier basis



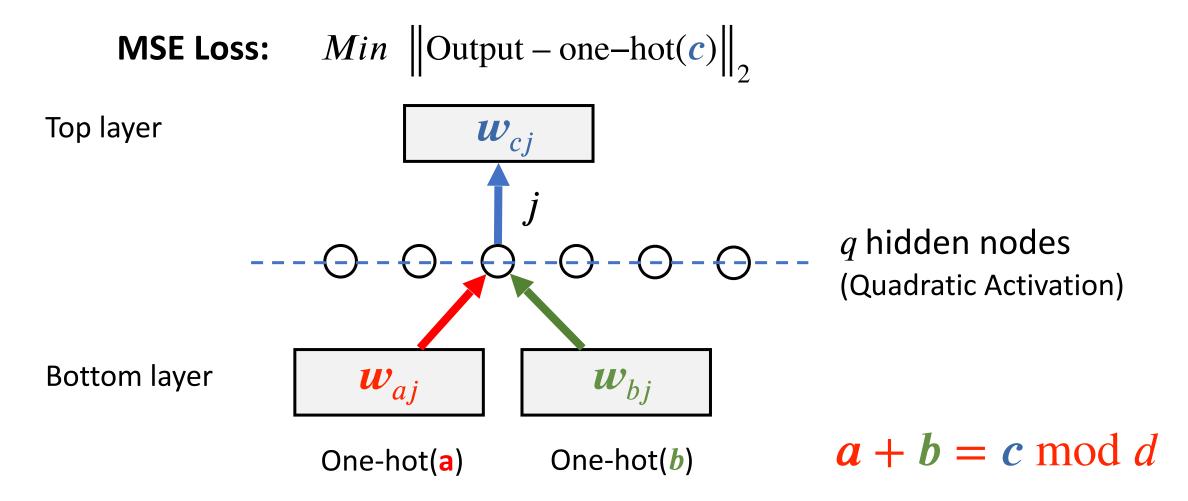




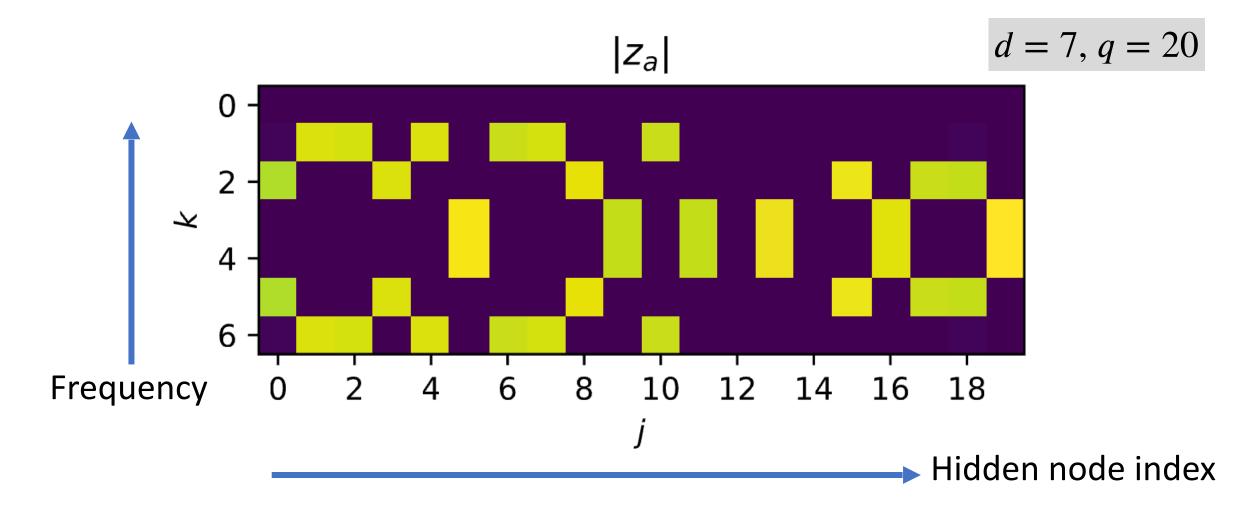
(a) Final logits for top Fourier components

[T. Zhou et al, Pre-trained Large Language Models Use Fourier Features to Compute Addition, NeurIPS'24] [S. Kantamneni, Language Models Use Trigonometry to Do Addition, arXiv'25]

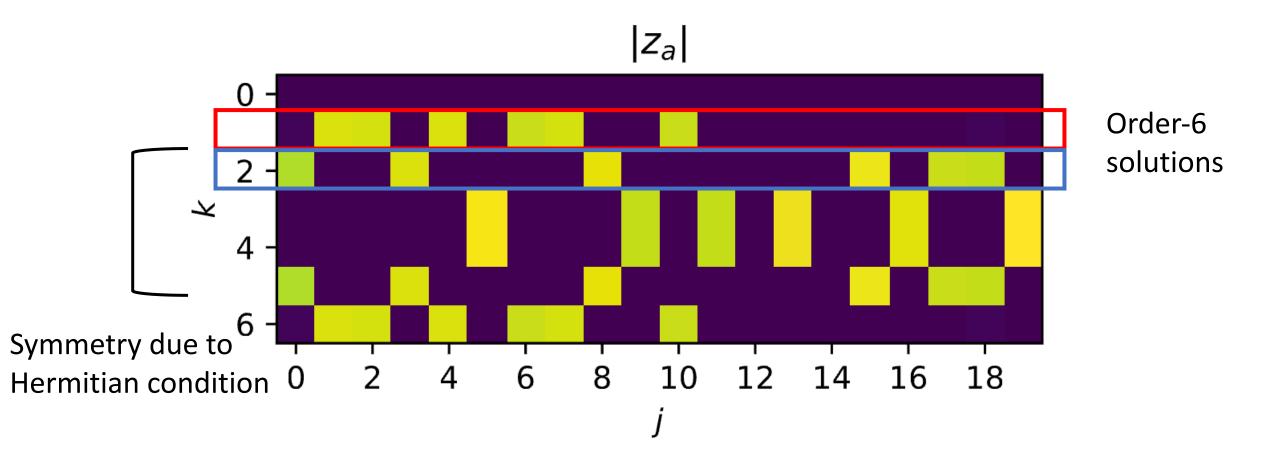
### Minimal Problem Setup



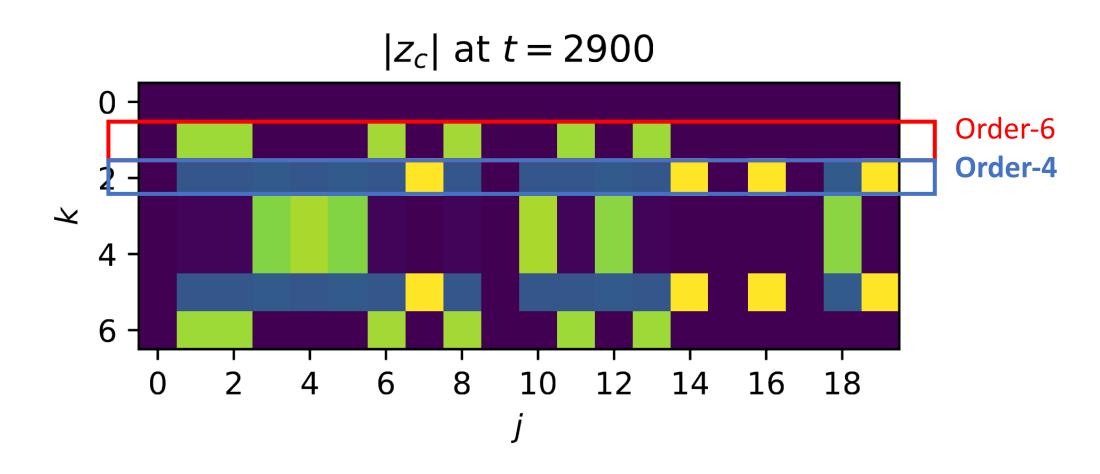
### What a Gradient Descent Solution look like?



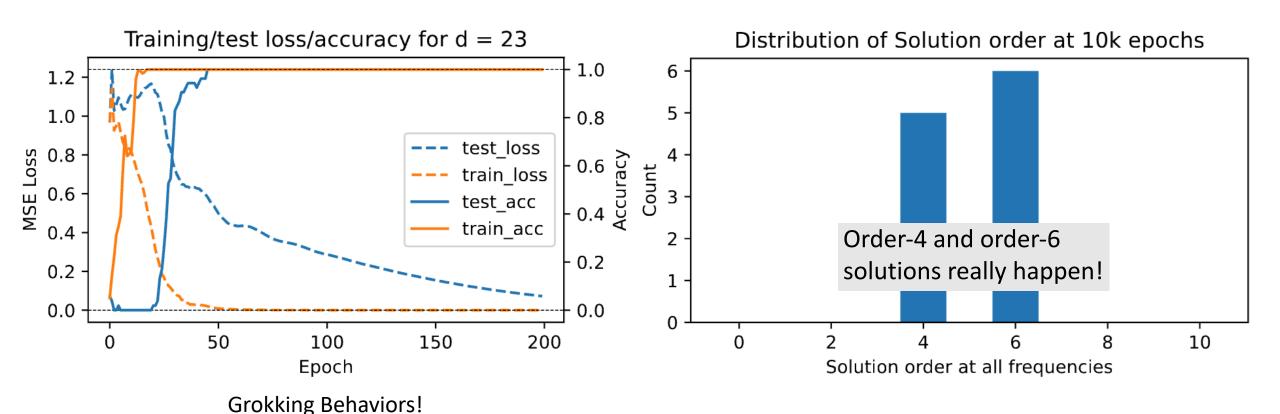
### What a Gradient Descent Solution look like?



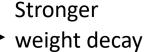
### What a Gradient Descent Solution look like?

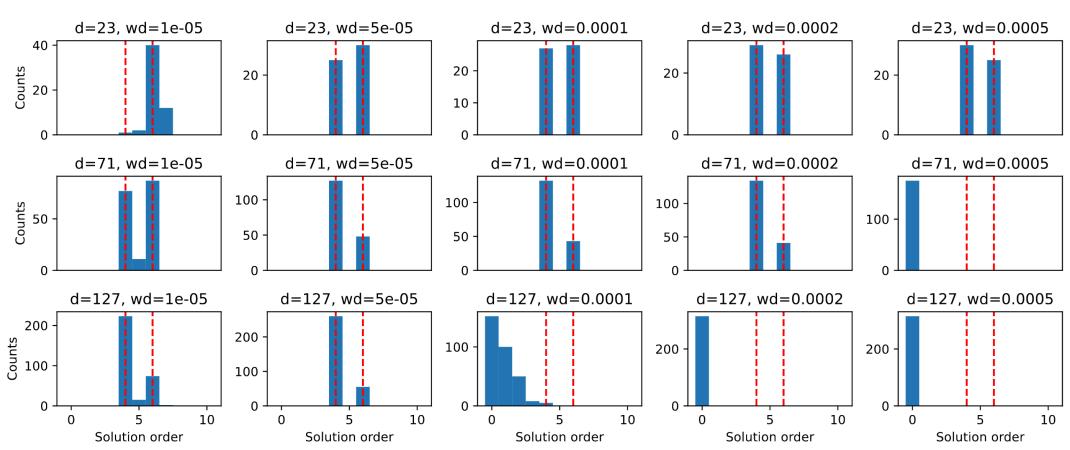


### More Statistics on Gradient Descent Solutions



## Effect of Weight Decay





## Theory to explicitly construct such solutions

Order-6 
$$z_{F6}$$
 (2\*3)

$$m{z}_{F6} = rac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} m{z}_{ ext{syn}}^{(k)} * m{z}_{
u}^{(k)} * m{y}_{k}$$

## Theory to explicitly construct such solutions

Order-6  $z_{F6}$  (2\*3)

$$m{z}_{F6} = rac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} m{z}_{ ext{syn}}^{(k)} * m{z}_{
u}^{(k)} * m{y}_k$$

Order-4  $z_{F4/6}$  (2\*2) (mixed with order-6)

$$m{z}_{F4/6} = rac{1}{\sqrt[3]{6}} \hat{m{z}}_{F6}^{(k_0)} + rac{1}{\sqrt[3]{4}} \sum_{k=1, k 
eq k_0}^{(a-1)/2} m{z}_{F4}^{(k)}$$

### Theory to explicitly construct such solutions

Order-6  $z_{F6}$  (2\*3)

Order-4  $z_{F4/6}$  (2\*2) (mixed with order-6)

Perfect memorization (order-d per frequency)

$$m{z}_{F6} = rac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} m{z}_{ ext{syn}}^{(k)} * m{z}_{
u}^{(k)} * m{y}_k$$

$$m{z}_{F4/6} = rac{1}{\sqrt[3]{6}} \hat{m{z}}_{F6}^{(k_0)} + rac{1}{\sqrt[3]{4}} \sum_{k=1, k 
eq k_0}^{(d-1)/2} m{z}_{F4}^{(k)}$$

$$egin{align} oldsymbol{z}_a &= \sum_{j=0}^{d-1} oldsymbol{u}_a^j, & oldsymbol{z}_b &= \sum_{j=0}^{d-1} oldsymbol{u}_b^j \ oldsymbol{z}_M &= d^{-2/3} oldsymbol{z}_a * oldsymbol{z}_b \end{aligned}$$

4	%not						solution distribution (%) in factorable ones			
	order-4/6	order-4	order-6	order-4	order-6	$oxed{oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)}}$	$oxed{z_{ u=\mathrm{i}}^{(k)} * z_{\mathrm{syn},lphaeta}^{(k)}}$	$oxed{z_ u^{(k)} * z_{\mathrm{syn}}^{(k)}}$	others	
23	$0.0 \pm 0.0$	$0.00 \pm 0.00$	$ 5.71\pm_{5.71} $	$0.05{\pm}0.01$	$4.80 \pm 0.96$	$47.07 \pm 1.88$	$11.31{\scriptstyle\pm1.76}\atop4.00{\scriptstyle\pm1.14}$	$39.80 \pm 2.11$	$1.82 \pm 1.82$	
71	$0.0 \pm 0.0$	$0.00\pm0.00$	$ 0.00\pm0.00 $	$ 0.03\pm 0.00 $	$5.02\pm0.25$	$72.57\pm0.70$	$4.00{\pm}1.14$	$ 21.14\pm 2.14 $	$2.29{\pm}1.07$	
127	$0.0 \pm 0.0$	$1.50{\pm}0.92$	$ 0.00\pm 0.00 $	$\left 0.26\pm0.14\right $	$\left 0.93\pm 0.18\right $	$82.96 \pm 0.39$	$2.25{\scriptstyle\pm0.64}$			

$$q = 512, \ wd = 5 \cdot 10^{-5}$$

d	%not order-4/6	%non-fa order-4	order-6	error (> order-4	$\langle 10^{-2} \rangle$ order-6	$oxed{oxed} egin{aligned}  ext{solution} \ oldsymbol{z}_{ u= ext{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)} \end{aligned}$	$oxed{oxed} egin{aligned}  ext{distribution (\%)} \ oxed{oxed} oxed{oxed}_{ u= ext{i}} * oldsymbol{z}_{ ext{syn},lphaeta}^{(k)} \end{aligned}$	) in factorabl $oldsymbol{z}_{ u}^{(k)} * oldsymbol{z}_{ ext{syn}}^{(k)}$	le ones others
23	$0.0 \pm 0.0$	$0.00 \pm 0.00$	$ 5.71\pm 5.71 $	$0.05 \pm 0.01$	$ 4.80\pm0.96 $	$47.07 \pm 1.88$	$11.31 \pm 1.76$	$39.80 \pm 2.11$	$1.82 \pm 1.82$
71	$0.0 \pm 0.0$	$0.00 \pm 0.00$	$ 0.00\pm 0.00 $	$ 0.03\pm 0.00 $	$ 5.02\pm 0.25 $	$72.57 \pm 0.70$	$4.00{\scriptstyle\pm1.14}$		
127	$0.0 \pm 0.0$	$1.50{\scriptstyle\pm0.92}$	$ 0.00\pm0.00 $	$\left 0.26\pm0.14\right $	$ 0.93 \pm 0.18 $	$82.96 \pm 0.39$	$2.25{\pm}0.64$	$14.13 \pm 0.87$	$0.66\pm$ 0.66
'	'		1 1	•	'		•	•	1

100% of the per-freq solutions are order-4/6

$d \mid$	%not order-4/6	%non-fa order-4	order-6	error (> order-4	$\langle 10^{-2} \rangle$ order-6	$oxed{egin{array}{c}  ext{solution} \ oxed{z}_{ u-\mathrm{i}}^{(k)} * oxed{z}_{\epsilon}^{(k)} \end{array}}$	distribution (%) $ oldsymbol{z}_{ u=\mathrm{i}}^{(k)}*oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)}$	) in factorabl $ oldsymbol{z}_{ u}^{(k)}*oldsymbol{z}_{ ext{syn}}^{(k)} $	e ones others
23	$0.0 \pm 0.0$	$0.00 \pm 0.00$	$5.71{\pm}5.71$	$0.05\pm$ 0.01	$4.80 \pm 0.96$	<b>,</b>	11.31±1.76	$39.80 \pm 2.11$	$1.82 \pm 1.82$
127	$\begin{array}{ c c }\hline 0.0\pm0.0\end{array}$	$1.50\pm 0.92$	$0.00\pm0.00$	$0.26 \pm 0.14$	$0.93 \pm 0.18$	$ 82.96\pm0.39 $	$2.25 \pm 0.64$	$14.13 \pm 0.87$	$0.66 \pm 0.66$

95% of the solutions are factorizable into "2\*3" or "2\*2"

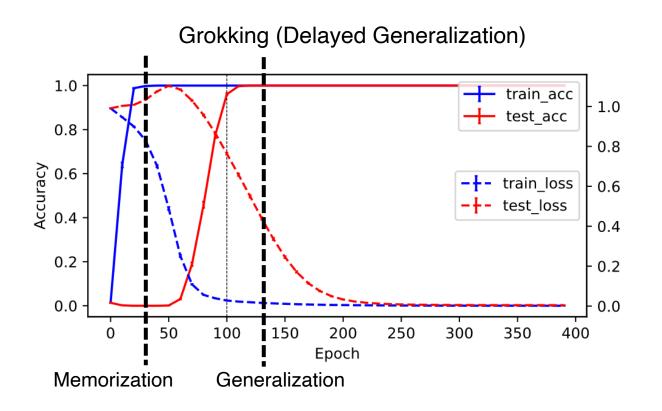
$d \mid$	%not order-4/6	%non-fa	order-6		$\times 10^{-2}$ ) order-6		distribution (%) $ oldsymbol{z}_{ u=\mathrm{i}}^{(k)}*oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)} $		
71	$0.0 \pm 0.0$	$ 0.00\pm0.00 $	$ 0.00\pm0.00 $	$0.03 \pm 0.00$	$5.02{\pm}0.25$	$ \begin{vmatrix} 47.07 \pm 1.88 \\ 72.57 \pm 0.70 \\ 82.96 \pm 0.39 \end{vmatrix} $	$4.00{\pm}1.14$	$39.80\pm 2.11$ $21.14\pm 2.14$ $14.13\pm 0.87$	$2.29 \pm 1.07$

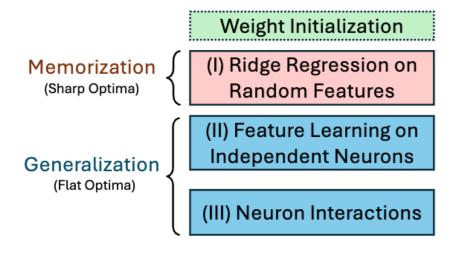
Factorization error is very small

4	%not	not   %non-factorable		error ( $\times 10^{-2}$ )		solution distribution (%) in factorable ones				
$begin{array}{c} a \\ \end{array}$	order-4/6	order-4	order-6	order-4	order-6	$oldsymbol{z}_{ u=\mathrm{i}}^{(k)}*oldsymbol{z}_{\xi}^{(k)}$	$oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)}$	$oxed{z_ u^{(k)} * z_{\mathrm{syn}}^{(k)}}$	others	
23	$0.0 \pm 0.0$	$0.00 \pm 0.00$	$5.71 \pm 5.71$	$0.05{\pm}0.01$	$4.80{\scriptstyle\pm0.96}$	$47.07{\pm}1.88$	$11.31 \pm 1.76$	$39.80 \pm 2.11$	$1.82{\pm}1.82$	
71	$0.0 \pm 0.0$	$0.00\pm0.00$	$ 0.00\pm 0.00 $	$ 0.03\pm0.00 $	$5.02{\pm}0.25$	$72.57 \pm 0.70$	$4.00{\pm}1.14$	$ 21.14\pm 2.14 $	$2.29{\scriptstyle\pm1.07}$	
127	$0.0 \pm 0.0$	$1.50{\scriptstyle\pm0.92}$	$ 0.00\pm0.00 $	$0.26 \pm 0.14$	$0.93 \pm 0.18$	$82.96 \pm 0.39$	$2.25{\pm}0.64$	$14.13 \pm 0.87$	$0.66 \pm 0.66$	
'		'	'	•				•		

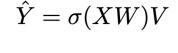
98% of the solutions can be factorizable into the constructed forms

## Understanding Grokking





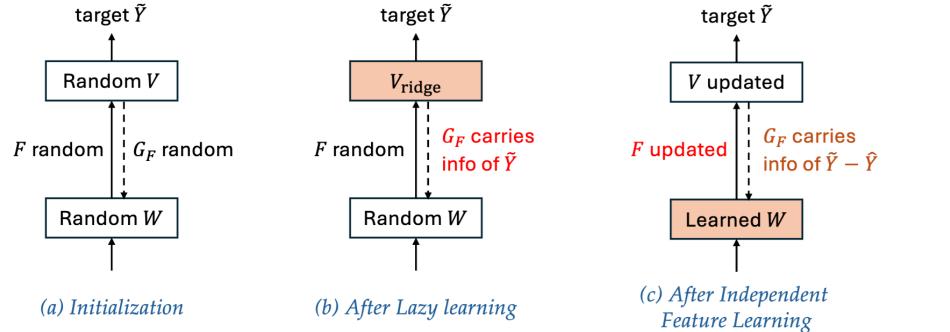
## Stages of Grokking Behaviors



 $\hat{Y} \approx \tilde{Y}$ 

V updated

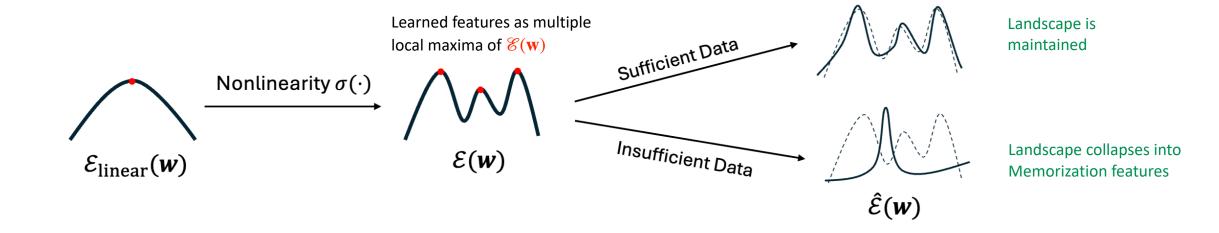
 $G_F \approx 0$ 



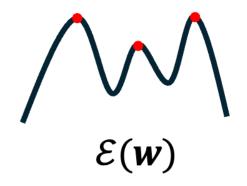
Complete W

## Connect Emerging Features with Data / Architecture

We discover that there exists an energy function  $\mathscr{E}(\mathbf{w})$  that governs the feature learning process



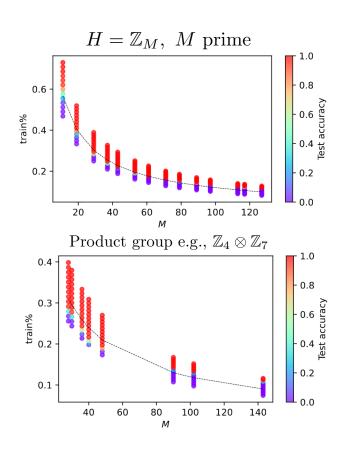
### Emerging Features are Symbolic!

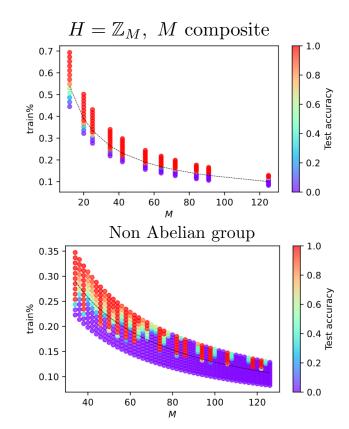


$$\mathcal{E}(\mathbf{w}) = rac{1}{2} \sum_h \langle \tilde{R}_h, S \rangle_F^2 = rac{M}{2} \sum_{k 
eq 0} rac{1}{d_k} \Big| \sum_r \mathrm{tr}(\hat{S}_{k,r}) \Big|^2$$

**Theorem 2** (Local maxima of  $\mathcal{E}$  for group input). For group arithmetics tasks with  $\sigma(x) = x^2$ ,  $\mathcal{E}$  has multiple local maxima  $\mathbf{w}^* = [\mathbf{u}; \pm P\mathbf{u}]$ . Either it is in a real irrep of dimension  $d_k$  (with  $\mathcal{E}^* = M/8d_k$  and  $\mathbf{u} \in \mathcal{H}_k$ ), or in a pair of complex irrep of dimension  $d_k$  (with  $\mathcal{E}^* = M/16d_k$  and  $\mathbf{u} \in \mathcal{H}_k \oplus \mathcal{H}_{\bar{k}}$ ). These local maxima are not connected. No other local maxima exist.

## How much is **sufficient**? Provable Scaling Laws





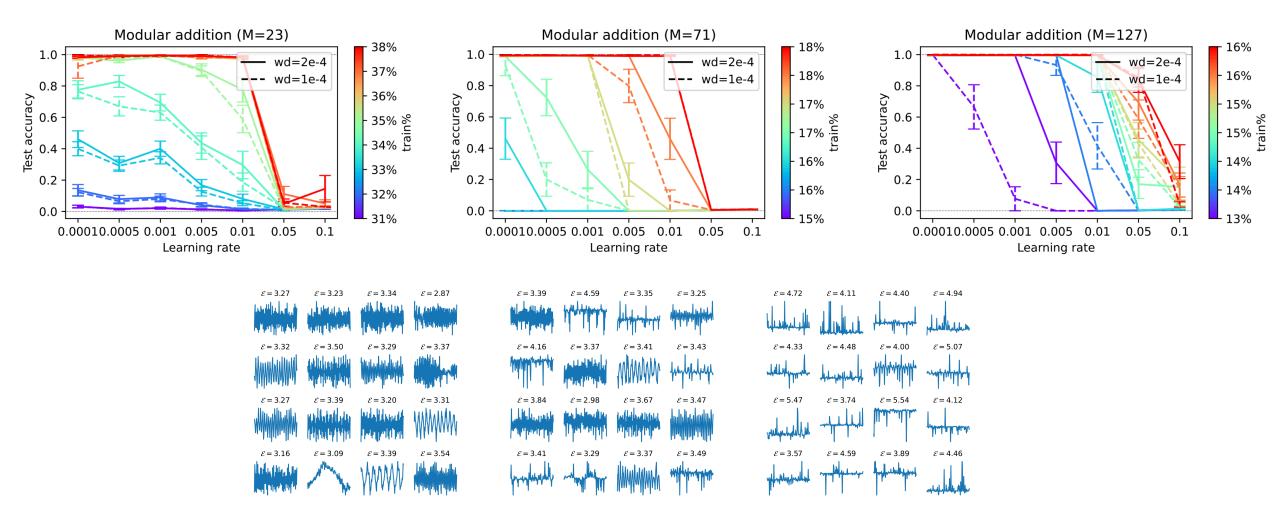
For Group Arithmetic tasks

Predict  $h_1h_2$ , given  $h_1, h_2 \in H$ 

O(|H| log |H|) data sample suffice to learn generalizable features

**Next Step:** Scale it to more complicated tasks and architectures

### Boundary between Memorization and Generalization



### Possible Implications

Do neural networks end up learning more efficient symbolic representations that we don't know?

Does gradient descent lead to a solution that can be reached by **advanced algebraic operations**?

Will gradient descent become **obsolete**, eventually?







# Thanks!

facebook Artificial Intelligence 75