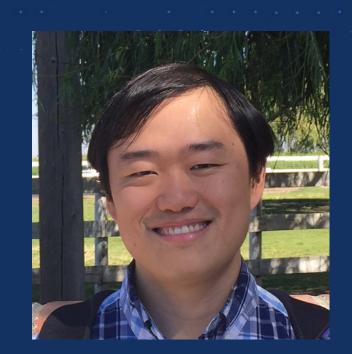
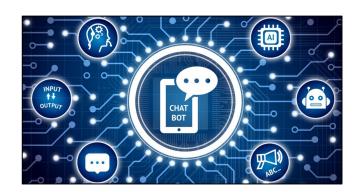
Provable Scaling Laws of Feature Emergence from Learning Dynamics of Grokking

Yuandong Tian (ex-)Research Scientist Director

Meta FAIR



Large Language Models (LLMs)



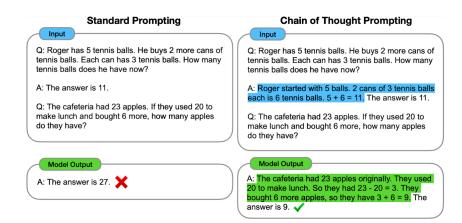
Conversational AI



Content Generation



Al Agents







Planning

The Progress of Large Models

Training compute of notable models



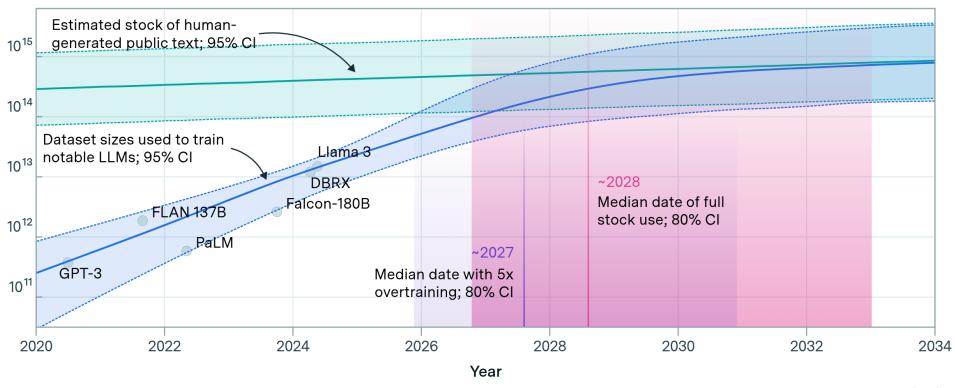


The Data Usage

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:

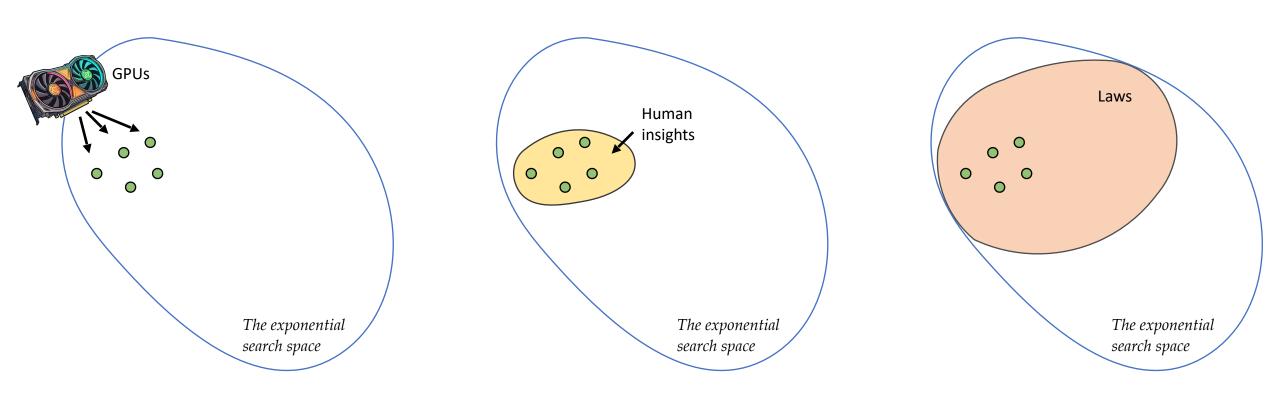
AGI timelines. My comments on AGI timelines looks to be the most trending part of the early response. This is the "decade of agents" is a reference to this earlier tweet x.com/karpathy/statu... Basically my AI timelines are about 5-10X pessimistic w.r.t. what you'll find in your neighborhood SF AI house 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

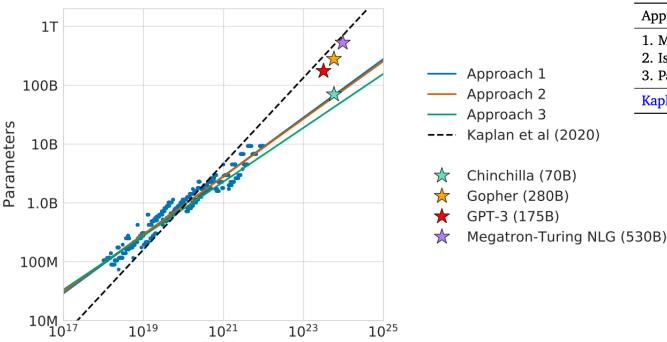
- Everyone is GPU poor
 - What are new axes to scale? GPUs are never enough.
 - Data itself cannot extrapolate, only human insights can.

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

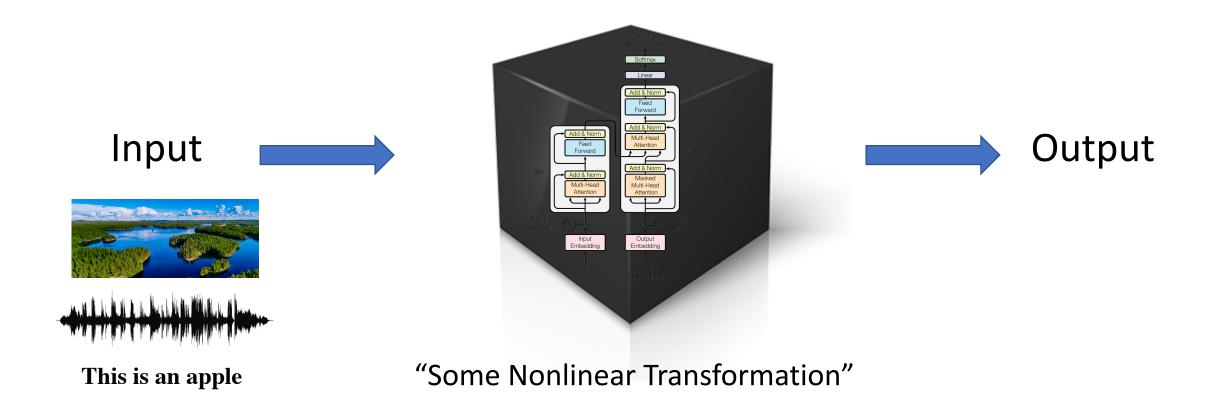
Still pure statistics and need exponential data. (No leverage of the knowledge of architecture/data)

 10^{19}

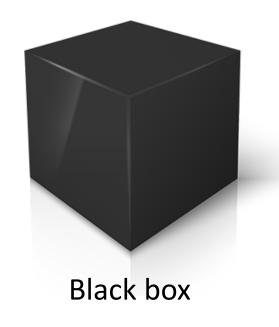
 10^{21}

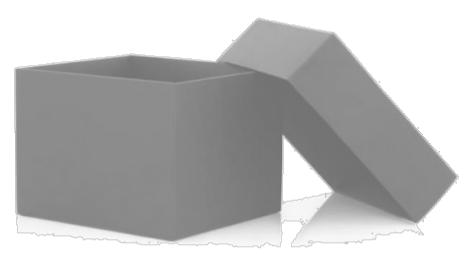
FLOPs

How does deep learning work?



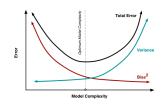
Black-box versus White-box





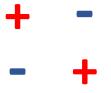
What routes should we take?

Generalization



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

Expressibility

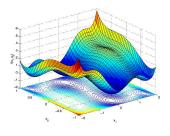


Architecture ✓
training dynamics ✗

How about

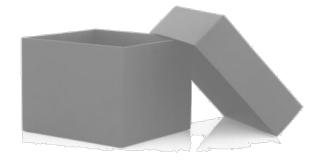
Architecture ✓
training dynamics ✓

Optimization



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

Start From the First Principle

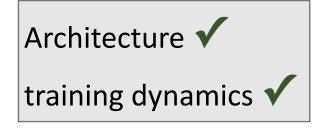


Training follows Gradient and its variants (SGD, Adams, etc)

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

First principle
 Understand the behavior of the neural networks by checking the gradient dynamics induced by the neural architectures.

Sounds complicated.. Is that possible? Yes



What Gradient Descent Analysis gives us?

Short-term:

Finding Simple Structures (Low-rank, sparsity)

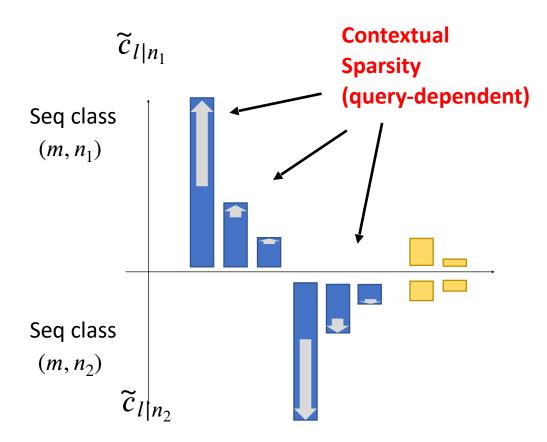
Long-term:

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



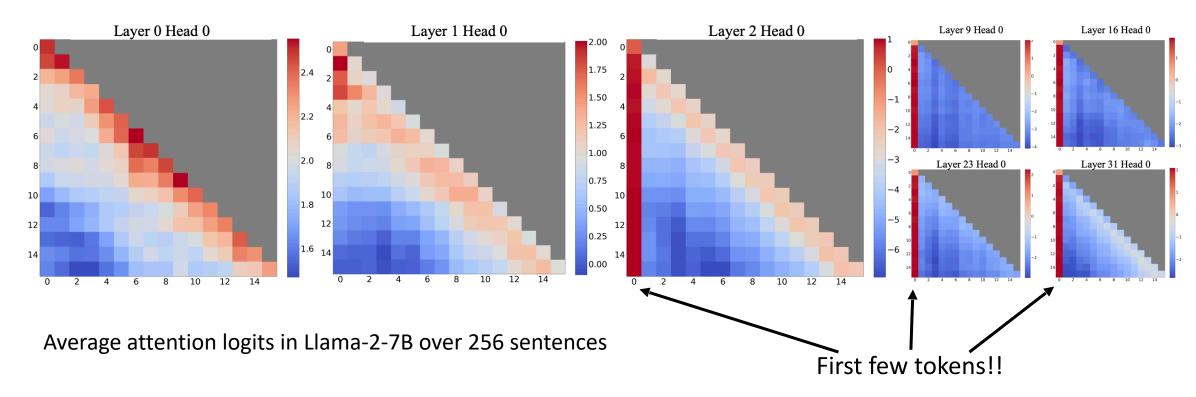
Leverage Them in Practical Algorithms

Finding Nice Structure: Attention Sparsity



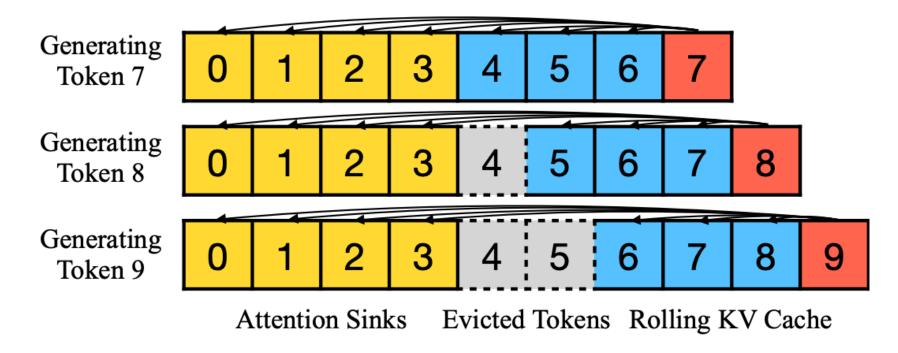
Attention = Learnable TF-IDF (Term Frequency, Inverse Document Frequency)

Attention Sinks: Initial tokens draw strong 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.

StreamingLLM



Key design: Position Rolling

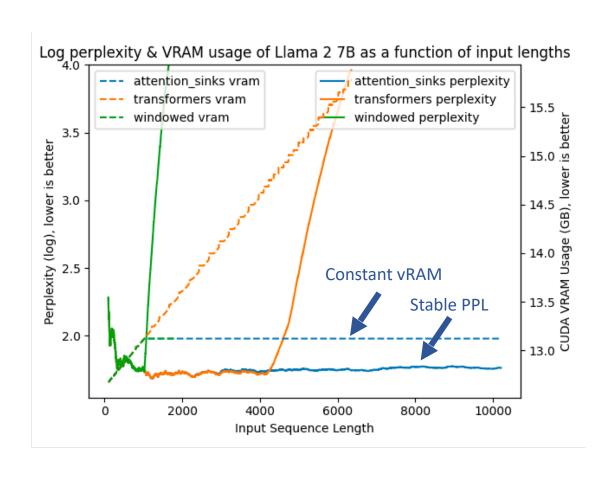
For all tokens, use their positions within cache to compute positional encoding!

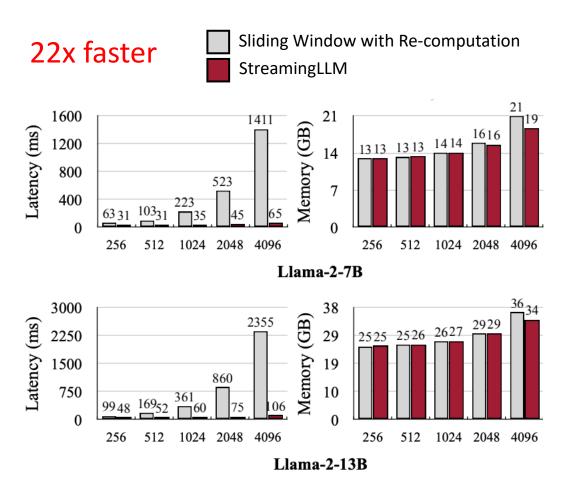
→ Token distance never exceeds pre-trained context window!

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





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.

- 1k+ citations
- Used in GPT OSS models in pre-training

Long-term: How Network finds Representation

Type of Representations

(Traditional) Symbolic representation

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

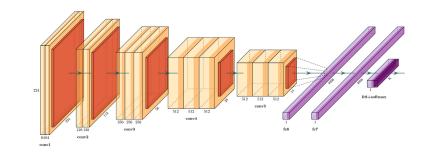
 $\nabla \cdot \mathbf{H} = 0$ (Gauss'Law for Magnetism)

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

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

Representation

Neural Representation (9)



Why Neural Representation is so effective?

Representation

(Traditional) Symbolic representation

4 Conclusions

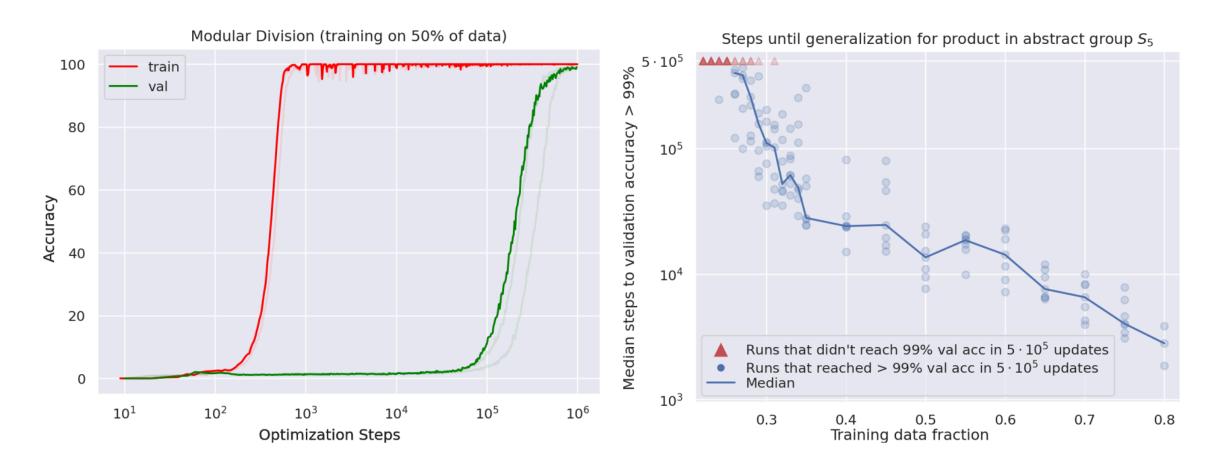
We have extended the GLU family of layers and proposed their use in Transformer. In a transfer-learning setup, the new variants seem to produce better perplexities for the de-noising objective used in pre-training, as well as better results on many downstream language-understanding tasks. These architectures 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

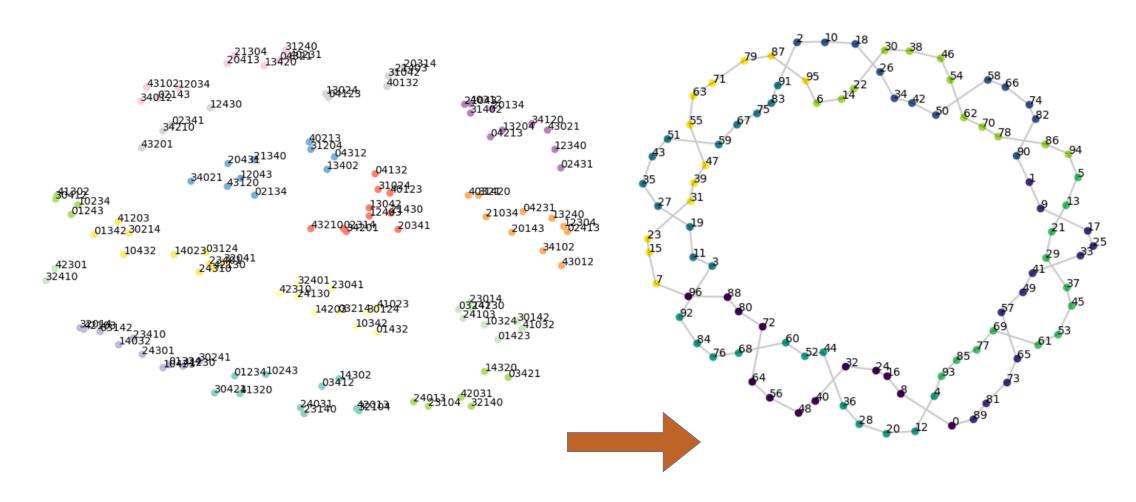
Representation (9)

Shall we just acknowledge that as "divine benevolence"?

Why there is Grokking Behavior?



Feature Emergence through Grokking



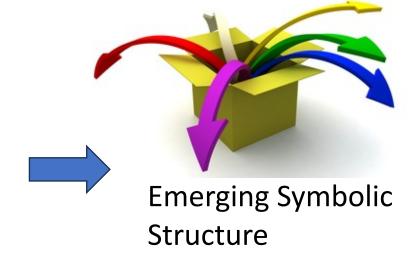


Mechanism of Emergent Representation

(Traditional) Symbolic representation

Representation

Neural Representation **



Modular Addition

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

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

Modular Addition

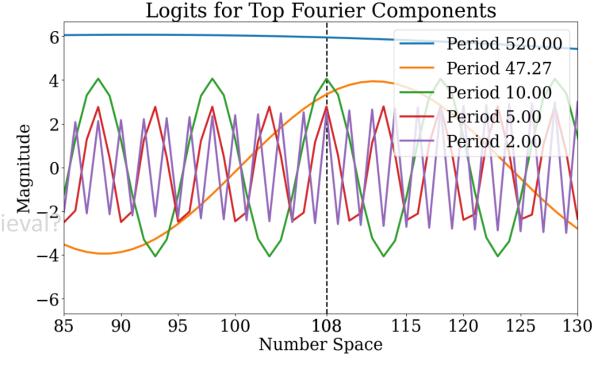
$$a + b = c \mod d$$

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

Learned representation = Fourier basis



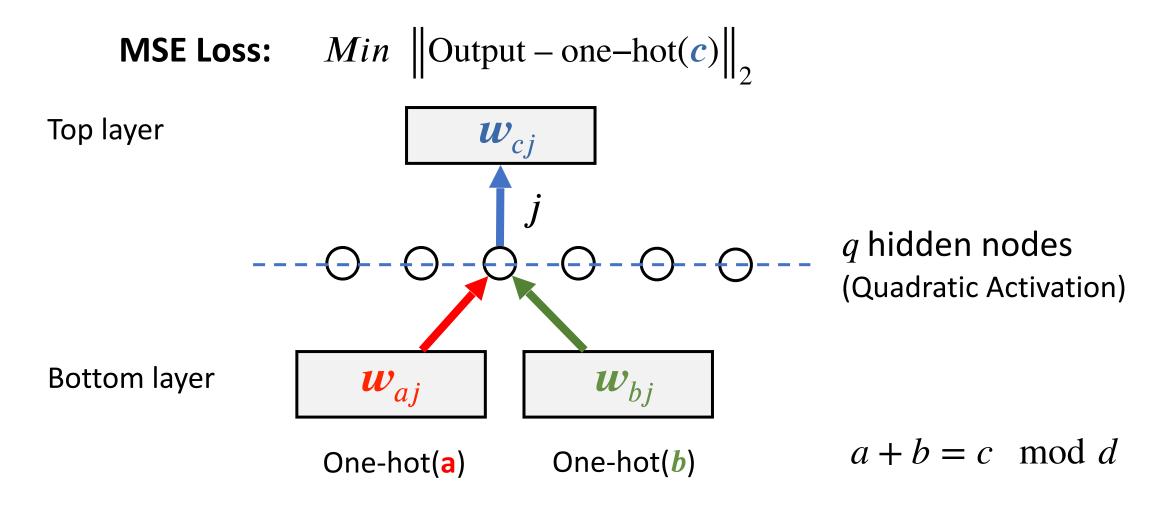




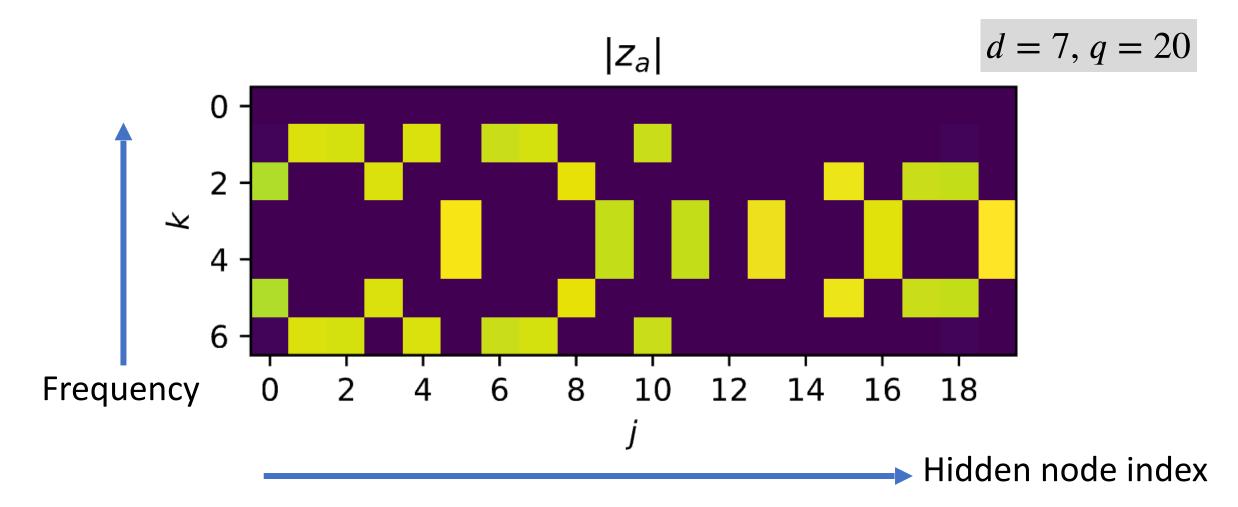
(a) Final logits for top Fourier components

[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]

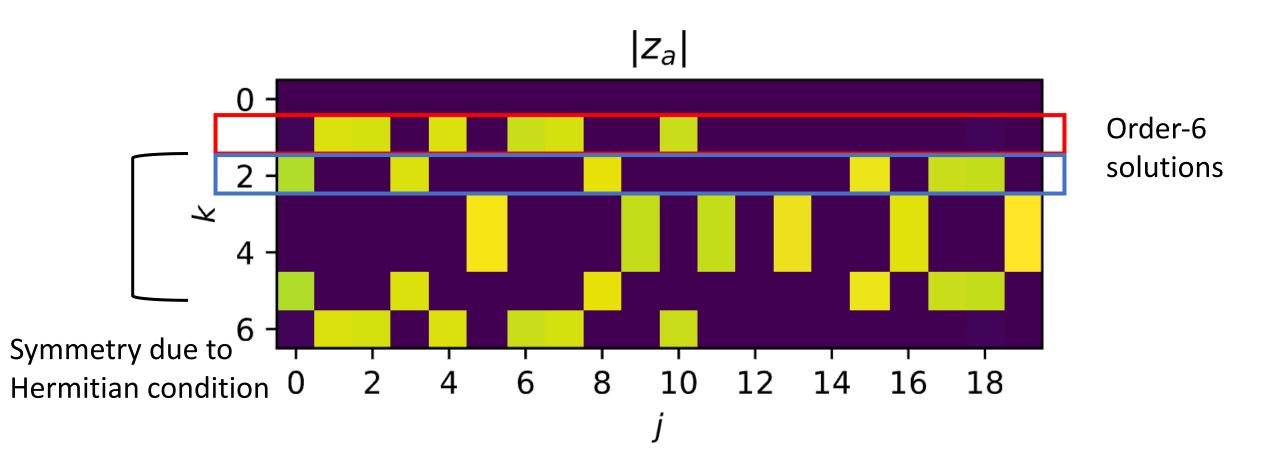
Minimal Problem Setup



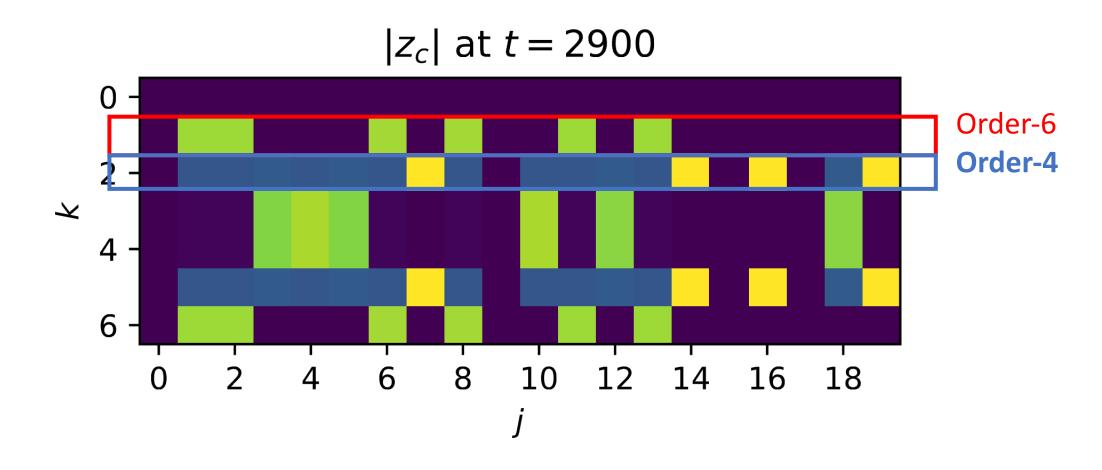
What a Gradient Descent Solution look like?



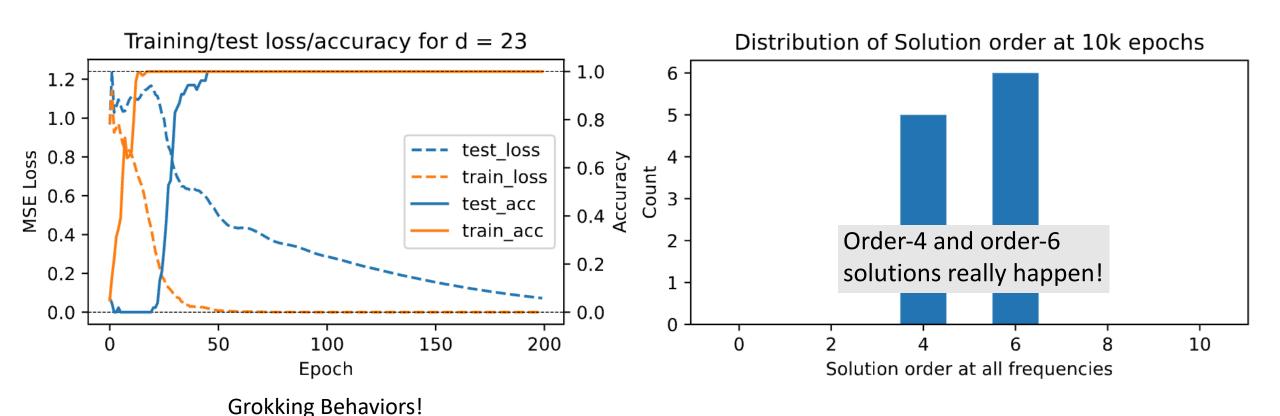
What a Gradient Descent Solution look like?



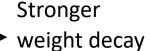
What a Gradient Descent Solution look like?

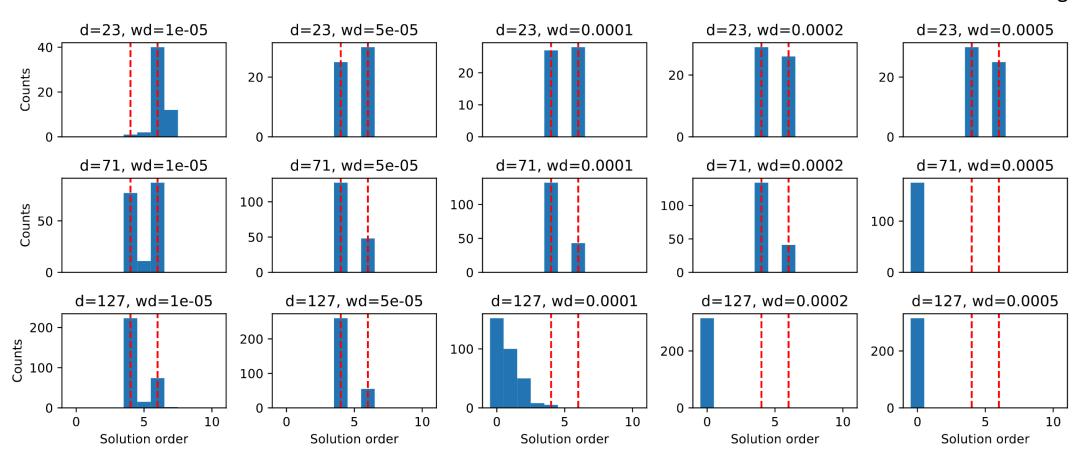


More Statistics on Gradient Descent Solutions



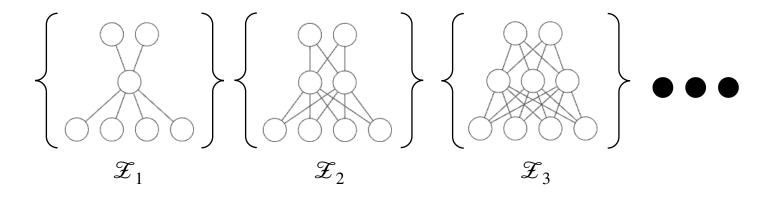
Effect of Weight Decay





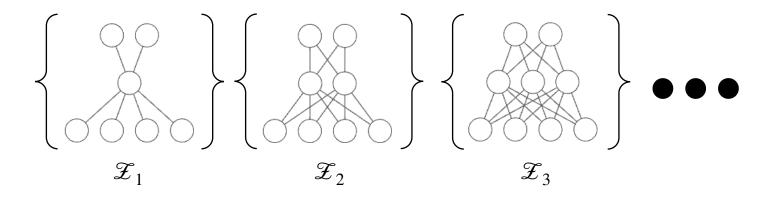
Why?

Nice algebraic structures exist for the solutions



$$\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$$
: All 2-layer networks with different number of hidden nodes

Nice algebraic structures exist for the solutions



$$\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$$
: All 2-layer networks with different number of hidden nodes

Ring addition +: Concatenate hidden nodes

Ring multiplication *: Kronecker production along the hidden dimensions

 $\langle \mathcal{Z}, +, * \rangle$ is a **semi-ring**

A function $r(\mathbf{z}): \mathcal{Z} \mapsto \mathbb{C}$ is a *ring homomorphism*, if

- r(1) = 1
- $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$
- $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$

A function $r(\mathbf{z}): \mathcal{Z} \mapsto \mathbb{C}$ is a *ring homomorphism*, if

- r(1) = 1
- $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$
- $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$

(z) and $r_{pk_1k_2k}(z)$ are <u>ring homomorphisms</u>!

A function $r(\mathbf{z}): \mathcal{Z} \mapsto \mathbb{C}$ is a ring homomorphism, if

- r(1) = 1
- $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$
- $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$

 $rac{r_{k_1k_2k}(z)}{r_{k_1k_2k}(z)}$ and $r_{pk_1k_2k}(z)$ are $rac{ring\ homomorphisms}{r_{k_1k_2k}(z)}$!

MSE Loss

$$\mathcal{C}_{k}(z) = -2r_{kkk} + \sum_{k_{1}k_{2}} \left| r_{k_{1}k_{2}k} \right|^{2} + \frac{1}{4} \left| \sum_{p \in \{a,b\}} \sum_{k'} r_{p,k',-k',k} \right|^{2} + \frac{1}{4} \sum_{m \neq 0} \sum_{p \in \{a,b\}} \left| \sum_{k'} r_{p,k',m-k',k} \right|^{2}$$

A function $r(\mathbf{z}): \mathcal{Z} \mapsto \mathbb{C}$ is a ring homomorphism, if

- r(1) = 1
- $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$
- $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$

 $rac{r_{k_1k_2k}(z)}{r_{k_1k_2k}(z)}$ and $r_{pk_1k_2k}(z)$ are <u>ring homomorphisms!</u>

MSE Loss

$$\mathcal{E}_{k}(\mathbf{z}) = -2\mathbf{r}_{kkk} + \sum_{k_{1}k_{2}} \left| \mathbf{r}_{k_{1}k_{2}k} \right|^{2} + \frac{1}{4} \left| \sum_{p \in \{a,b\}} \sum_{k'} \mathbf{r}_{p,k',-k',k} \right|^{2} + \frac{1}{4} \sum_{m \neq 0} \sum_{p \in \{a,b\}} \left| \sum_{k'} \mathbf{r}_{p,k',m-k',k} \right|^{2}$$

Partial solution \mathbf{z}_1 satisfies $\mathbf{r}_{k_1k_2k}(\mathbf{z}_1) = 0$

Partial solution \mathbf{z}_2 satisfies $r_{pk',-k',k}(\mathbf{z}_2)=0$

A function $r(\mathbf{z}): \mathcal{Z} \mapsto \mathbb{C}$ is a ring homomorphism, if

- r(1) = 1
- $r(\mathbf{z}_1 + \mathbf{z}_2) = r(\mathbf{z}_1) + r(\mathbf{z}_2)$
- $r(\mathbf{z}_1 * \mathbf{z}_2) = r(\mathbf{z}_1)r(\mathbf{z}_2)$

 $r_{k_1k_2k}(z)$ and $r_{pk_1k_2k}(z)$ are <u>ring homomorphisms!</u>

MSE Loss

$$\mathcal{E}_{k}(\mathbf{z}) = -2\mathbf{r}_{kkk} + \sum_{k_{1}k_{2}} \left| \mathbf{r}_{k_{1}k_{2}k} \right|^{2} + \frac{1}{4} \left| \sum_{p \in \{a,b\}} \sum_{k'} \mathbf{r}_{p,k',-k',k} \right|^{2} + \frac{1}{4} \sum_{m \neq 0} \sum_{p \in \{a,b\}} \left| \sum_{k'} \mathbf{r}_{p,k',m-k',k} \right|^{2}$$

Partial solution
$$\mathbf{z}_1$$
 satisfies $\mathbf{r}_{k_1k_2k}(\mathbf{z}_1) = 0$
Partial solution \mathbf{z}_2 satisfies $\mathbf{r}_{pk',-k',k}(\mathbf{z}_2) = 0$

$$\mathbf{z} = \mathbf{z}_1 * \mathbf{z}_2 \text{ satisfies both } \mathbf{r}_{k_1k_2k}(\mathbf{z}) = \mathbf{r}_{pk',-k',k}(\mathbf{z}) = 0$$

Composing Global Solutions from Partial Ones

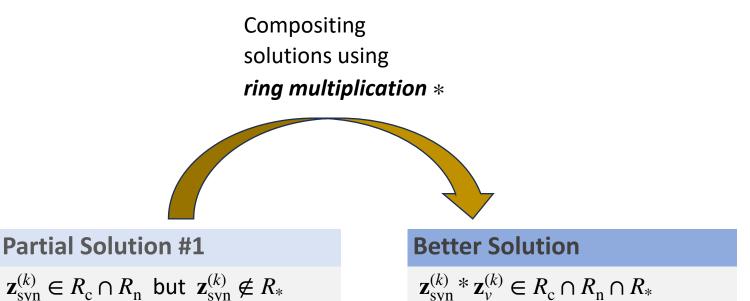
Partial Solution #1

 $\mathbf{z}_{\mathrm{syn}}^{(k)} \in R_{\mathrm{c}} \cap R_{\mathrm{n}} \text{ but } \mathbf{z}_{\mathrm{syn}}^{(k)} \notin R_{*}$

Partial Solution #2

 $\mathbf{z}_{v}^{(k)} \in R_{*}$

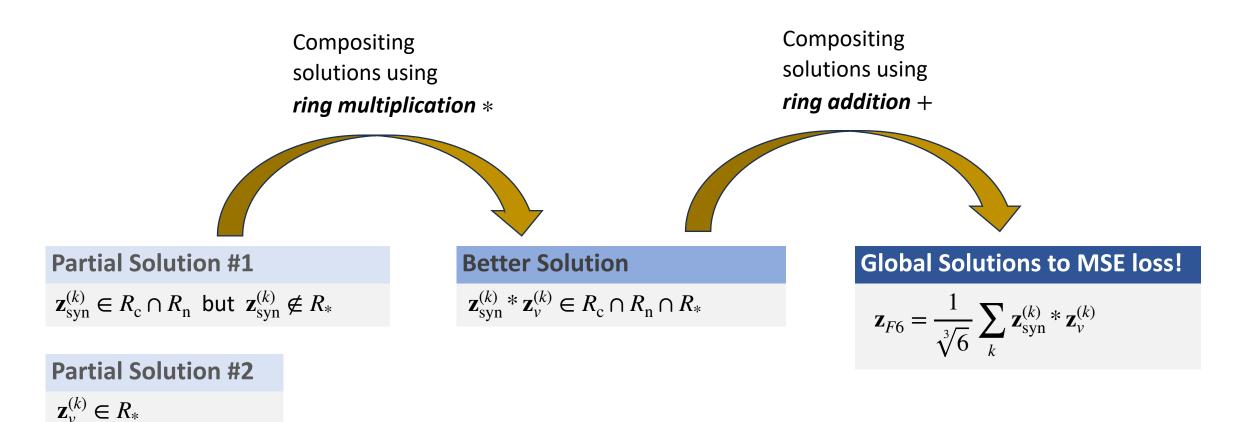
Composing Global Solutions from Partial Ones



Partial Solution #2

$$\mathbf{z}_{v}^{(k)} \in R_{*}$$

Composing Global Solutions from Partial Ones



Optimal solutions can be constructed

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

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

Optimal solutions can be constructed

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)

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

Optimal solutions can be constructed

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 | %not %non-factorable | | error (> | $\times 10^{-2}$) | solution distribution (%) in factorable ones | | | |
|-----|---------------|------------------------|---------------------|-----------------|--------------------|---|---|--|-----------------------------|
| | order-4/6 | order-4 | order-6 | order-4 | order-6 | $oxed{oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)}}$ | $ig oldsymbol{z}_{ u=\mathrm{i}}^{(k)}*oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)}$ | $\left oldsymbol{z}_{ u}^{(k)}*oldsymbol{z}_{	ext{syn}}^{(k)} ight $ | others |
| 23 | 0.0 ± 0.0 | 0.00 ± 0.00 | $ 5.71_{\pm 5.71} $ | 0.05 ± 0.01 | $ 4.80\pm0.96 $ | 47.07 ± 1.88 | 11.31 ± 1.76 | 39.80 ± 2.11 | 1.82 ± 1.82 |
| 71 | 0.0 ± 0.0 | 0.00 ± 0.00 | $ 0.00\pm 0.00 $ | 0.03 ± 0.00 | $ 5.02 \pm 0.25 $ | 72.57 ± 0.70 | $11.31{\scriptstyle\pm1.76}\atop 4.00{\scriptstyle\pm1.14}$ | $ 21.14\pm 2.14 $ | $2.29{\scriptstyle\pm1.07}$ |
| 127 | 0.0 ± 0.0 | $ 1.50\pm_{0.92} $ | $ 0.00\pm 0.00 $ | 0.26 ± 0.14 | $ 0.93 \pm 0.18 $ | 82.96 ± 0.39 | $2.25{\pm}0.64$ | | |

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

| d | %not order-4/6 | %non-fa order-4 | order-6 | error (> order-4 | <10 ⁻²) order-6 | $oxed{oxed} egin{aligned} 	ext{solution} \ oldsymbol{z}_{ u=	ext{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)} \end{aligned}$ | distribution (%) $oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)}$ |) in factorabl $oldsymbol{z}_{ u}^{(k)} * oldsymbol{z}_{	ext{syn}}^{(k)}$ | le ones others |
|-----|-------------------|-----------------------------|------------------|----------------------------|--------------------------------|--|---|---|-------------------|
| 23 | 0.0 ± 0.0 | 0.00 ± 0.00 | $ 5.71\pm 5.71 $ | 0.05 ± 0.01 | 4.80 ± 0.96 | 47.07 ± 1.88 | 11.31 ± 1.76 | 39.80 ± 2.11 | 1.82 ± 1.82 |
| 71 | 0.0 ± 0.0 | 0.00 ± 0.00 | $ 0.00\pm0.00 $ | $ 0.03\pm 0.00 $ | $ 5.02\pm 0.25 $ | $ 72.57\pm0.70 $ | $4.00\pm$ 1.14 | 21.14 ± 2.14 | $2.29{\pm}1.07$ |
| 127 | 0.0 ± 0.0 | $1.50{\scriptstyle\pm0.92}$ | $ 0.00\pm0.00 $ | $\left 0.26\pm0.14\right $ | 0.93 ± 0.18 | 82.96 ± 0.39 | $2.25{\pm}0.64$ | 14.13 ± 0.87 | $0.66\pm$ 0.66 |
| ' | ' | | 1 1 | • | ' | | ' | • | 1 |
| | | | | | | | | | |

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

| d | %not order-4/6 | %non-fa order-4 | order-6 | error (> order-4 | $\times 10^{-2}$) order-6 | $oxed{oxed} egin{aligned} 	ext{solution} \ oldsymbol{z}_{ u=	ext{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)} \end{aligned}$ | $egin{aligned} 	ext{distribution (\%)} \ m{z}_{ u=	ext{i}}^{(k)} * m{z}_{	ext{syn},lphaeta}^{(k)} \end{aligned}$ |) in factorabl $oldsymbol{z}_{ u}^{(k)} * oldsymbol{z}_{	ext{syn}}^{(k)}$ | le ones others |
|-----|-------------------|-----------------------------|--------------------|---------------------|----------------------------|--|--|---|-------------------|
| 23 | 0.0 ± 0.0 | 0.00 ± 0.00 | $ 5.71\pm_{5.71} $ | 0.05 ± 0.01 | $ 4.80\pm_{0.96} $ | 47.07 ± 1.88 | | 39.80 ± 2.11 | 1.82 ± 1.82 |
| 71 | 0.0 ± 0.0 | 0.00 ± 0.00 | 0.00 ± 0.00 | 0.03 ± 0.00 | $ 5.02\pm 0.25 $ | 72.57 ± 0.70 | $4.00{\scriptstyle\pm1.14}$ | 21.14 ± 2.14 | 2.29 ± 1.07 |
| 127 | 0.0 ± 0.0 | $1.50{\scriptstyle\pm0.92}$ | 0.00 ± 0.00 | $0.26\pm$ 0.14 | $ 0.93 \pm 0.18 $ | 82.96 ± 0.39 | $2.25{\pm}0.64$ | 14.13 ± 0.87 | $0.66\pm$ 0.66 |
| , | ' | | ' | • | ' | | ' | • | • |
| | | | | | | | | | |

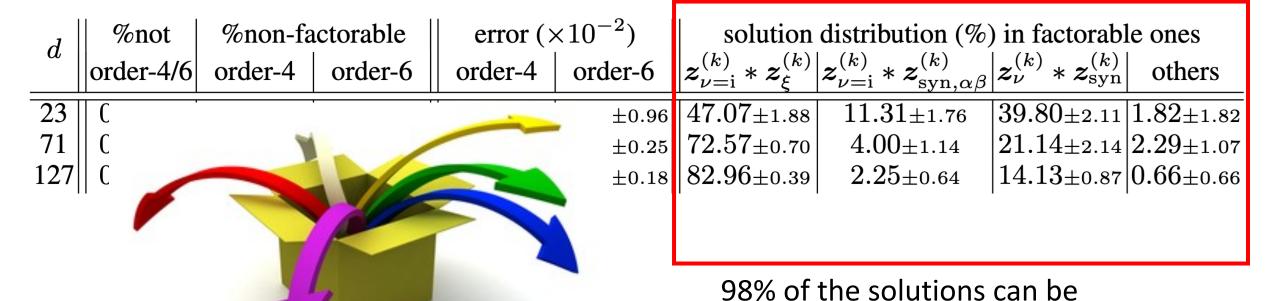
95% of the solutions are factorizable into "2*3" or "2*2"

| d | %not order-4/6 | %non-fa | order-6 | error (> order-4 | <10 ⁻²) order-6 | 1 | distribution (%) $m{z}_{ u=\mathrm{i}}^{(k)} * m{z}_{\mathrm{syn}, \alpha\beta}^{(k)}$ | | |
|-----|-------------------|-----------------|------------------|---------------------|--------------------------------|--|--|--|-----------------|
| 71 | 0.0 ± 0.0 | $ 0.00\pm0.00 $ | $ 0.00\pm 0.00 $ | 0.03 ± 0.00 | $5.02{\pm}0.25$ | $\begin{array}{ c c c c c }\hline 47.07{\pm}1.88\\ 72.57{\pm}0.70\\ 82.96{\pm}0.39\\ \hline \end{array}$ | $4.00{\pm}1.14$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $2.29{\pm}1.07$ |
| 127 | 0.0±0.0 | 1.00±0.92 | 0.00±0.00 | 0.20 ±0.14 | 0.93±0.18 | 02.90±0.39 | 2.2 0±0.64 | 14.10±0.87 | 0.00±0.66 |

Factorization error is very small

| 4 | %not | %non-factorable | | \parallel error (×10 ⁻²) | | solution distribution (%) in factorable ones | | | |
|---------------------------------------|---------------|-----------------|-----------------|--|-----------------------------|---|--|---|-----------------|
| $\begin{array}{c c} a \\ \end{array}$ | order-4/6 | order-4 | order-6 | order-4 | order-6 | $oxed{oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\xi}^{(k)}}$ | $oldsymbol{z}_{ u=\mathrm{i}}^{(k)} * oldsymbol{z}_{\mathrm{syn},lphaeta}^{(k)}$ | $oxed{z_ u^{(k)} * z_{\mathrm{syn}}^{(k)}}$ | others |
| 23 | 0.0 ± 0.0 | 0.00 ± 0.00 | 5.71 ± 5.71 | $0.05{\pm}0.01$ | $4.80{\scriptstyle\pm0.96}$ | 47.07 ± 1.88 | 11.31 ± 1.76 | 39.80 ± 2.11 | 1.82 ± 1.82 |
| | | | | I | | $72.57{\pm0.70}$ | | $ 21.14\pm 2.14 $ | $2.29{\pm}1.07$ |
| | 1 | | | | | 82.96 ± 0.39 | | 14.13 ± 0.87 | 0.66 ± 0.66 |
| , | • | ' | ' | • | | | | • | |
| | | | | | | | | | |

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



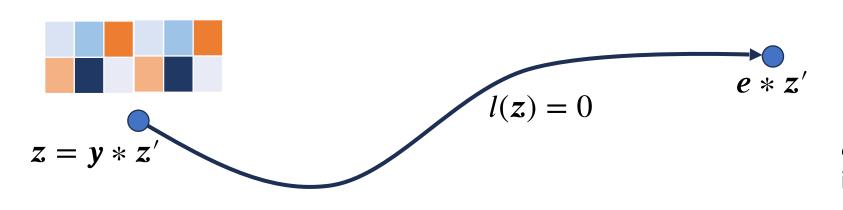
factorizable into the constructed forms

Emerging Symbolic structure from neural network training

facebook Artificial Intelligence

How about Gradient Dynamics?

Theorem [The Occam's Razer] If z = y * z and both z and z are global optimal, then there exists a path of zero loss connecting z and z'.

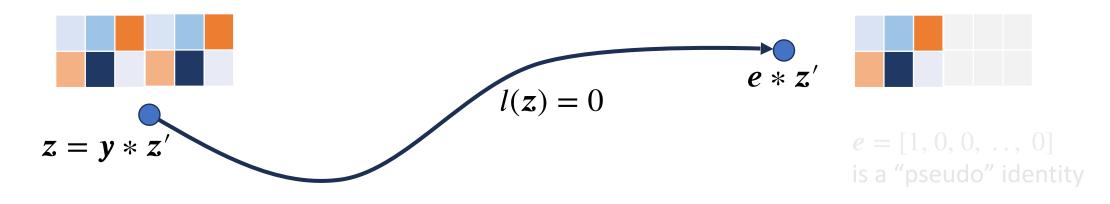




e = [1, 0, 0, ..., 0] is a "pseudo" identity

How about Gradient Dynamics?

Theorem [The Occam's Razer] If z = y * z and both z and z are global optimal, then there exists a path of zero loss connecting z and z'.



L2 regularization will push the solution to e * z' (simpler solutions), since $||e * z'||_2 \le ||y * z'||_2$

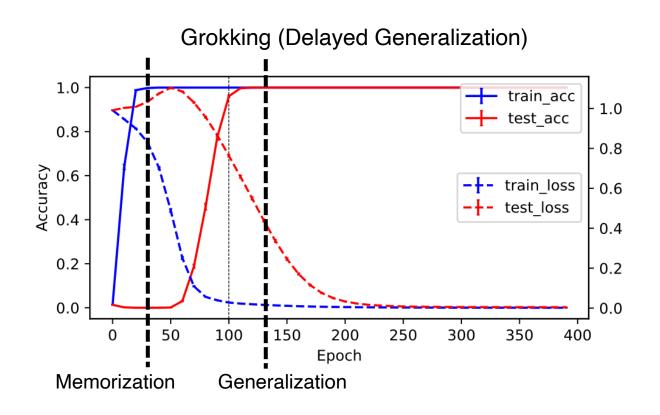
Limitation of the analysis

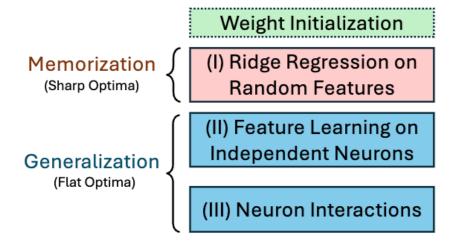
How such solutions are achieved — No analysis with gradient dynamics

Only apply to a combination of MSE loss + all data + quadratic activation

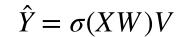
Next work is to crack the grokking behaviors from gradient dynamics

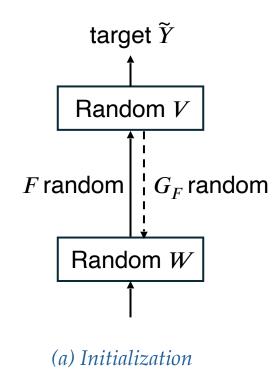
Understanding Grokking Behavior

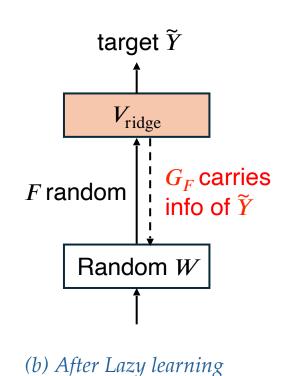


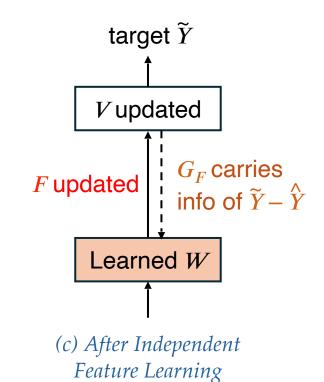


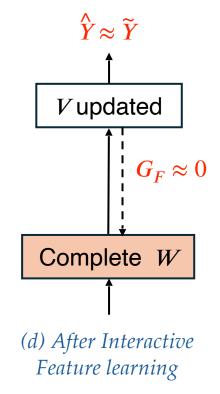
Stages of Grokking Behaviors











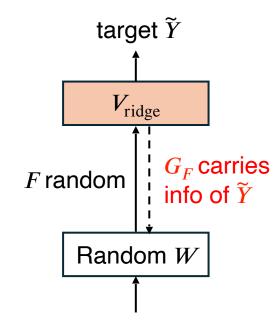
Stages I: NTK regime (Lazy Learning)

Objective function ($Y \in \mathbb{R}^{n \times M}, X \in \mathbb{R}^{n \times d}$)

$$\min_{V,W} \frac{1}{2} \|P_1^{\perp}(Y - \sigma(XW)V)\|_F^2$$

Ridge Regression (with weight decay η)

$$V_{\text{ridge}} = (\tilde{F}^{\top} \tilde{F} + \eta I)^{-1} \tilde{F}^{\top} \tilde{Y}$$



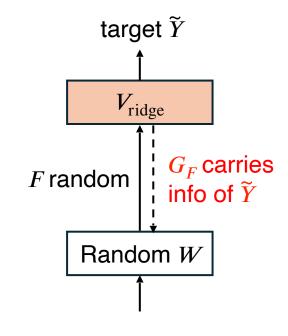
$$\begin{split} \tilde{Y} &= P_1^{\perp} Y \\ F &= \sigma(XW) \quad \tilde{F} = P_1^{\perp} F \end{split}$$

Here
$$P_1^{\perp} = I_n - \mathbf{1}\mathbf{1}^{\top}/n$$
 is the zero-mean projection

The backpropagated Gradient G_F

At Ridge regression solution $V_{
m ridge}$

$$G_F = \eta (\tilde{F}\tilde{F}^{\mathsf{T}} + \eta I)^{-1} \tilde{Y}\tilde{Y}^{\mathsf{T}} \tilde{F} (\tilde{F}^{\mathsf{T}}\tilde{F} + \eta I)^{-1}$$



Looks complicated... any interesting properties?

The backpropagated Gradient G_F

Lemma 1 (Structure of backpropagated gradient G_F). Assume that (1) entries of W follow normal distribution, (2) $\|\mathbf{x}_i\|_2$ have the same norm, (3) $\mathbf{x}_i^{\top}\mathbf{x}_{i'} = \rho$ for all $i \neq i'$ and (4) large width K, then both $\tilde{F}^{\top}\tilde{F}$ and $\tilde{F}\tilde{F}^{\top}$ becomes a multiple of identity and Eqn. 5 becomes:

$$G_F = \frac{\eta}{(Kc_1 + \eta)(nc_2 + \eta)} \tilde{Y} \tilde{Y}^{\top} F \tag{6}$$

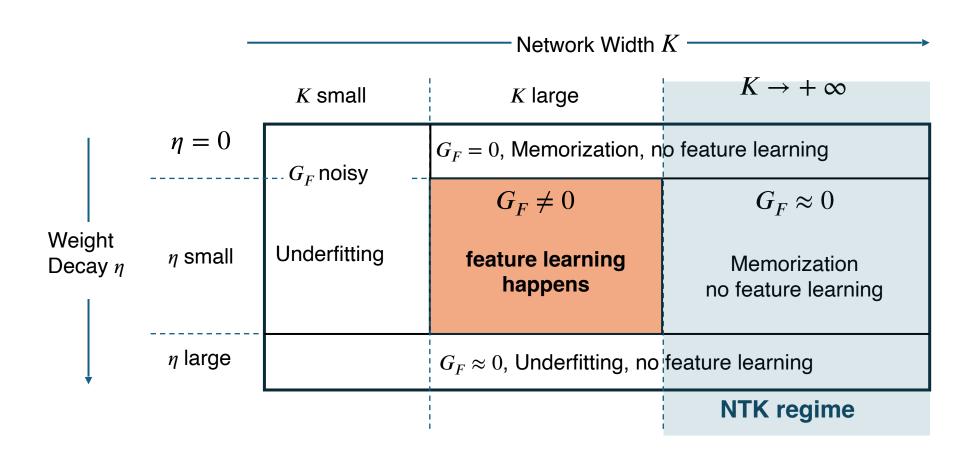
Key insights

If weight decay $\eta = 0$, then $G_F = 0$ (no feature learning)

If weight decay is large, then $G_F \to 0$

If number of hidden nodes $K \to +\infty$, then $G_F \to 0$ (NTK regime)

The regime that feature learning happens



Stage II: The Energy Function $\mathscr{E}(\mathbf{w})$

Component-wise dynamics

 $G_F \propto \eta \tilde{Y} \tilde{Y}^{\mathsf{T}} F$



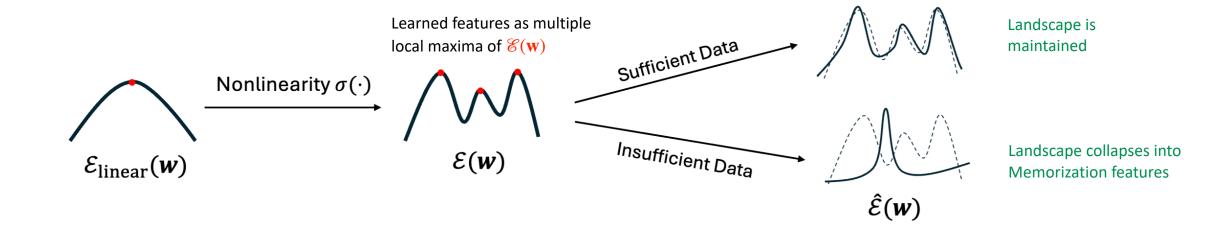
$$\dot{\mathbf{w}}_j = X^{\top} D_j \mathbf{g}_j, \quad \mathbf{g}_j \propto \eta \tilde{Y} \tilde{Y}^{\top} \sigma(X \mathbf{w}_j)$$

Theorem 1 (The energy function \mathcal{E} for independent feature learning). The dynamics (Eqn. 7) of independent feature learning is exactly the gradient ascent dynamics of the energy function \mathcal{E} w.r.t. \mathbf{w}_i , a nonlinear canonical-correlation analysis (CCA) between the input X and target \tilde{Y} :

$$\mathcal{E}(\mathbf{w}_j) = \frac{1}{2} \|\tilde{Y}^{\top} \sigma(X \mathbf{w}_j)\|_2^2$$
 (8)

Connect Emerging Features with Data / Architecture

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



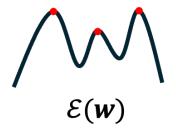
Group Representation Theory

The decomposition of group representation. The representation theory of finite group (Fulton & Harris, 2013; Steinberg, 2009) says that the regular representation R_h admits a decomposition into complex *irreducible* representations (or *irreps*):

$$R_h = Q \left(\bigoplus_{k=0}^{\kappa(H)} \bigoplus_{r=1}^{m_k} C_k(h) \right) Q^* \tag{9}$$

where $\kappa(H)$ is the number of nontrivial irreps (i.e., not all h map to identity), $C_k(h) \in \mathbb{C}^{d_k \times d_k}$ is the k-th irrep block of R_h , Q is the unitary matrix (and Q^* is its conjugate transpose) and m_k is the multiplicity of the k-th irrep. This means that in the decomposition of R_h , there are m_k copies of d_k -dimensional irrep, and these copies are isomorphic to each other. So the k-th irrep subspace \mathcal{H}_k has dimension $m_k d_k$.

Emerging Features are Symbolic!



$$\mathcal{E}(\mathbf{w}) = \frac{1}{2} \sum_{h} \langle \tilde{R}_h, S \rangle_F^2 = \frac{M}{2} \sum_{k \neq 0} \frac{1}{d_k} \left| \sum_{r} \operatorname{tr}(\hat{S}_{k,r}) \right|^2$$

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.

What is that specifically?

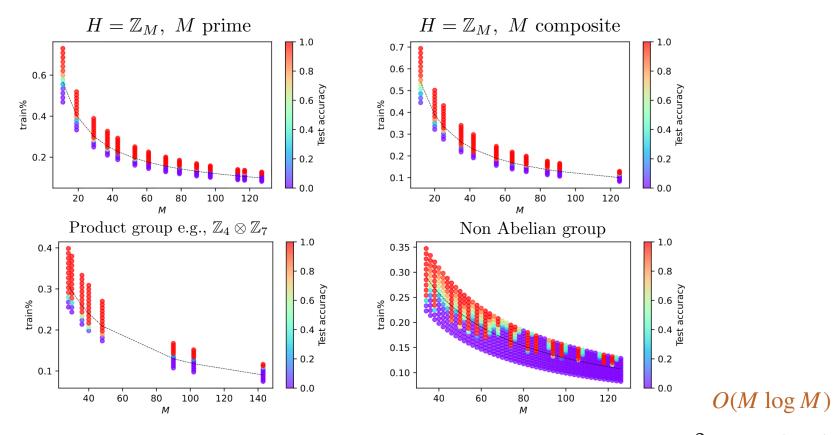
Corollary 2 (Modular addition). For modular addition with odd M, all local maxima are single frequency $\mathbf{u}_k = a_k [\cos(km\omega)]_{m=0}^{M-1} + b_k [\sin(km\omega)]_{m=0}^{M-1}$ where $\omega := 2\pi/M$ with $\mathcal{E}^* = M/16$. For even M, $\mathbf{u}_{M/2} \propto [(-1)^m]_{m=0}^{M-1}$ has $\mathcal{E}^* = M/8$. Different local maxima are disconnected.

Emerging Features are More Efficient!

Theorem 3 (Target Reconstruction). Assume (1) \mathcal{E} is optimized in complex domain \mathbb{C} , (2) for each irrep k, there are $m_k^2 d_k^2$ pairs of learned weights $\mathbf{w} = [\mathbf{u}; \pm P\mathbf{u}]$ whose associated rank-1 matrices $\{\mathbf{u}\mathbf{u}^*\}$ form a complete bases for \mathcal{H}_k and (3) the top layer V also learns with $\eta = 0$, then $\hat{Y} = \tilde{Y}$.

From the theorem, we know that $K=2\sum_{k\neq 0}m_k^2d_k^2\leq 2\left[(M-\kappa(H))^2+\kappa(H)-1\right]$ suffice. In particular, for Abelian group, $\kappa(H)=M-1$ and K=2M-2. This is much more efficient than a pure memorization solution that would require M^2 nodes, i.e., each node memorizes a single pair $(h_1,h_2)\in H\times H$.

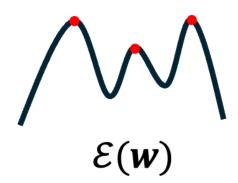
How much is sufficient? Provable Scaling Laws

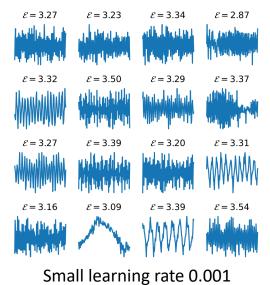


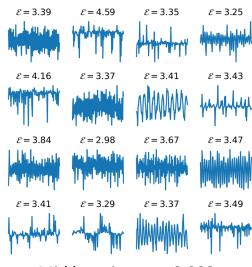
Theorem 4 (Amount of samples to maintain local optima). If we select $n \gtrsim d_k^2 M \log(M/\delta)$ data sample from $H \times H$ uniformly at random, then with probability at least $1 - \delta$, the empirical energy function $\hat{\mathcal{E}}$ keeps local maxima for d_k -dimensional irreps (Thm. 2).

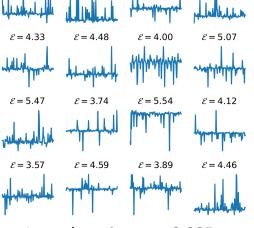
Boundary between Memorization and Generalization

At the boundary, large Learning rate leads to memorization (higher 8)









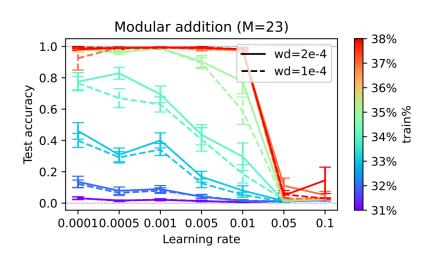
E = 4.40

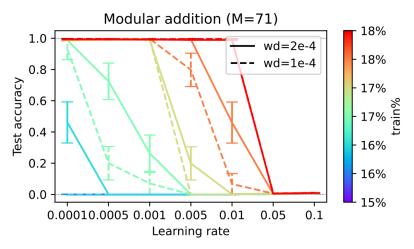
E = 4.11

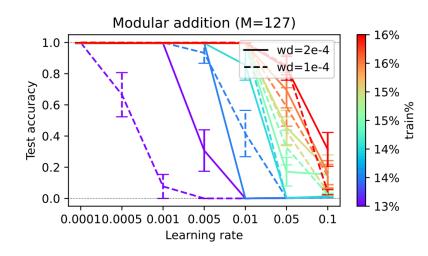
Mid learning rate 0.002

Large learning rate 0.005

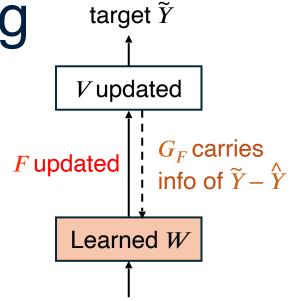
Boundary between Memorization and Generalization







Stage III: Interactive Feature Learning



Theorem 7 (Top-down Modulation). For group arithmetic tasks with $\sigma(x) = x^2$, if the hidden layer learns only a subset S of irreps, then the backpropagated gradient $G_F \propto (\Phi_S \otimes \mathbf{1}_M)(\Phi_S \otimes \mathbf{1}_M)^*F$ (see proof for the definition of Φ_S), which yields a modified \mathcal{E}_S that only has local maxima on the missing irreps $k \notin S$.

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**?

Here this work is just a tiny step.

Next Step: Scale it to more complicated tasks and architectures

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







Thanks!

facebook Artificial Intelligence 74