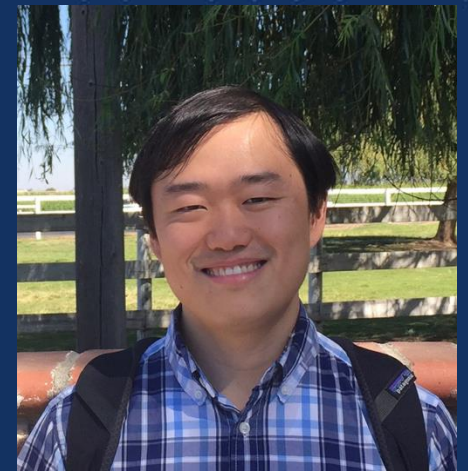


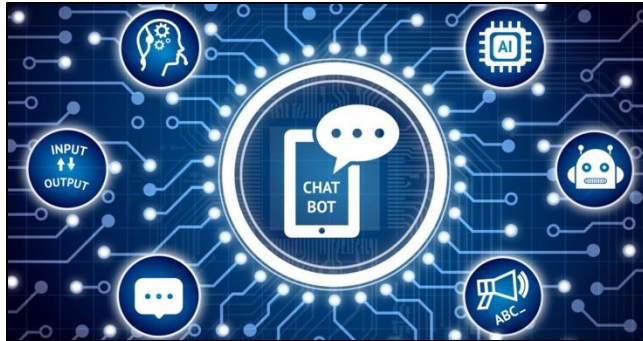
Emergence of Various Structures via the Lens of Transformer Training Dynamics

Yuandong Tian
Research Scientist Director

Meta AI



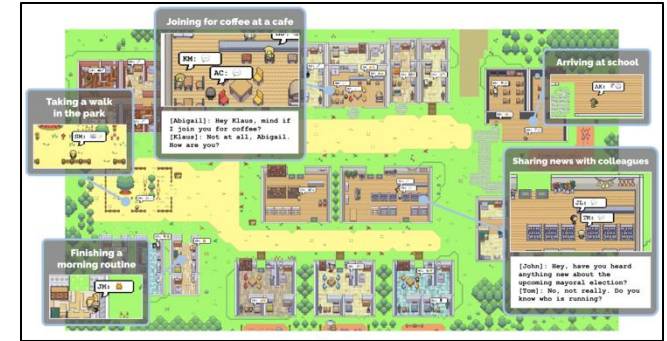
Large Language Models (LLMs)



Conversational AI



Content Generation



AI Agents

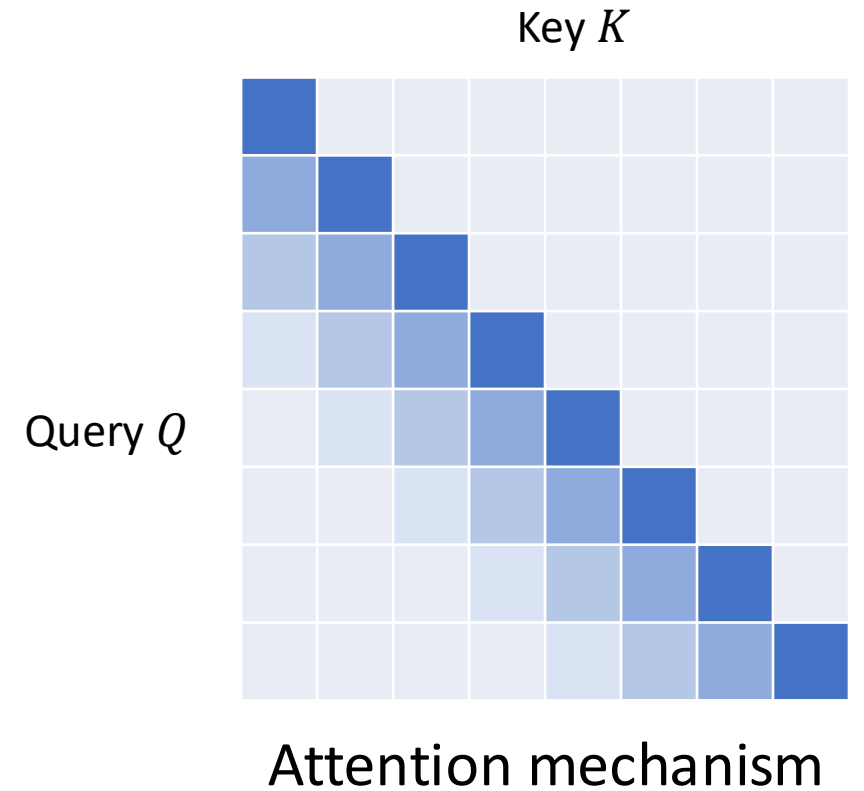
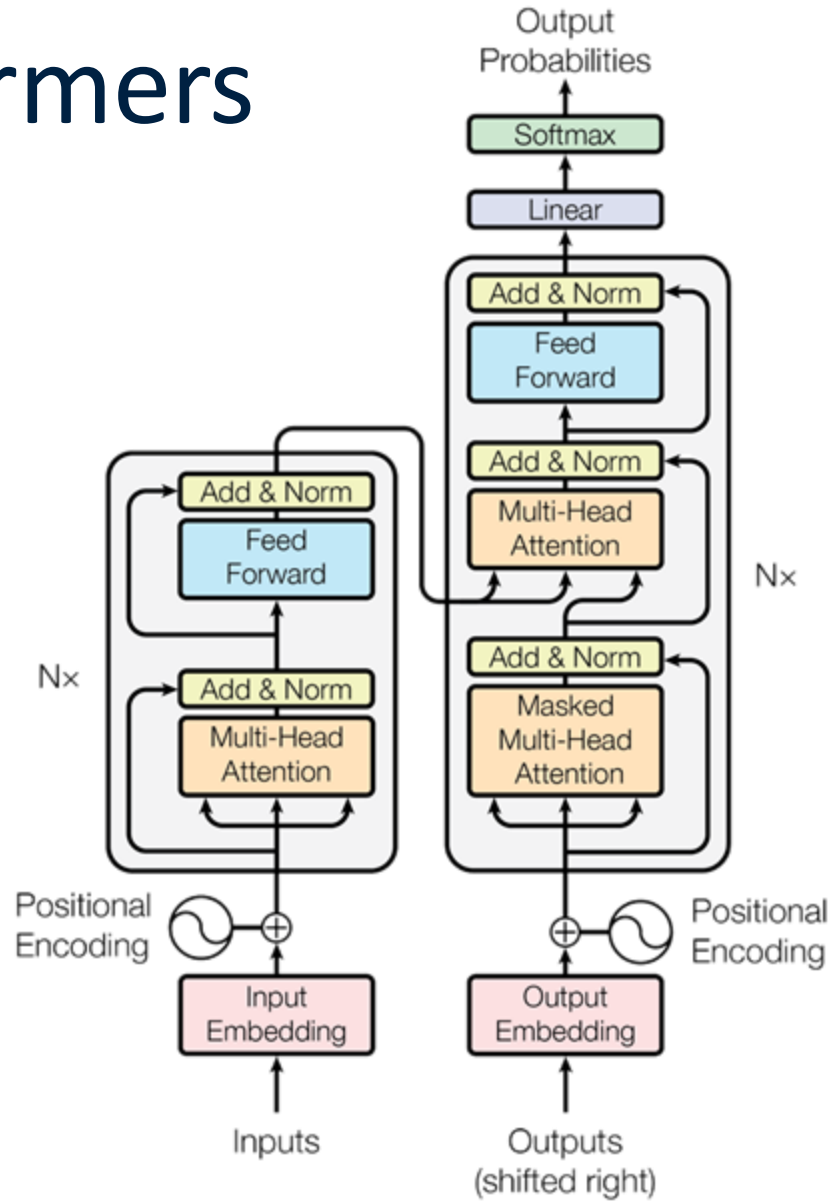
Standard Prompting	Chain of Thought Prompting
<p>Input</p> <p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: The answer is 11.</p> <p>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?</p>	<p>Input</p> <p>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?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.</p> <p>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?</p>
<p>Model Output</p> <p>A: The answer is 27. ❌</p>	<p>Model Output</p> <p>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. ✅</p>

Reasoning



Planning

Transformers

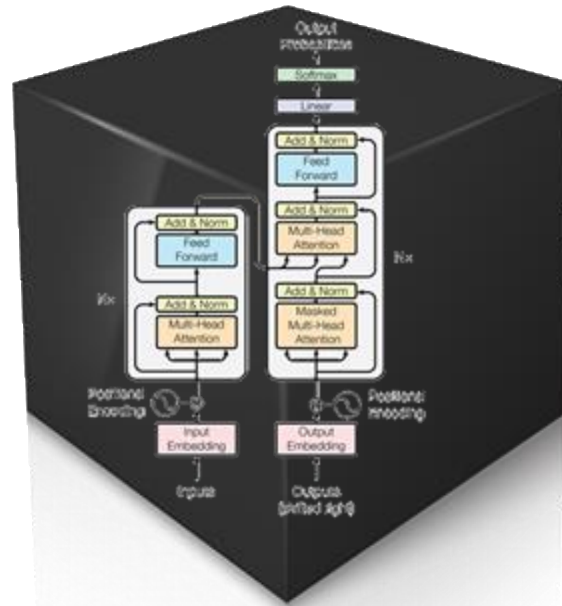


How does Transformer work?

Input



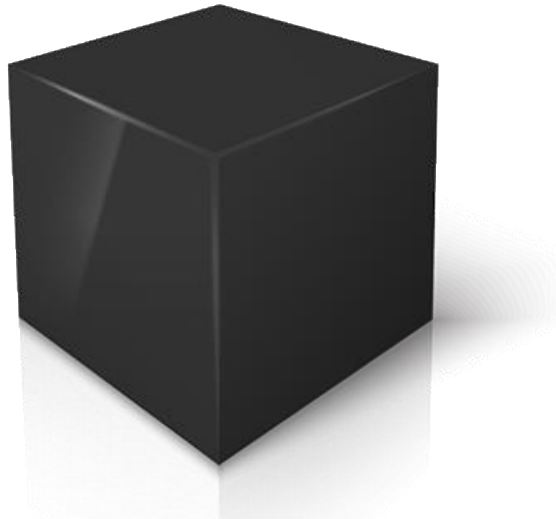
This is an apple



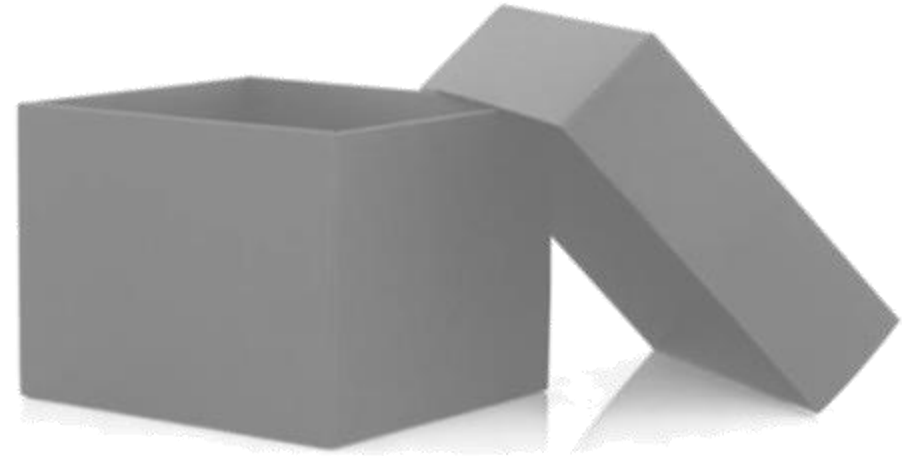
Output

“Some Nonlinear Transformation”

Black-box versus White-box



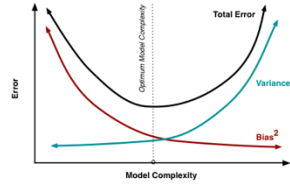
Black box



White box

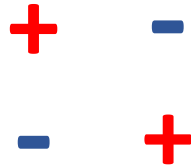
What routes should we take?

Generalization



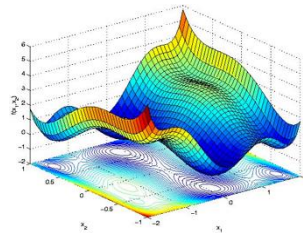
Architecture **X**
training dynamics **X**

Expressibility



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

Optimization



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

How about

Architecture **✓**
training dynamics **✓**



Start From the First Principle

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

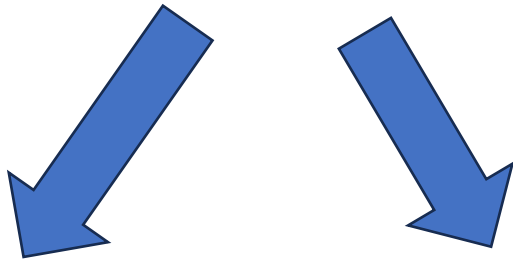
$$\dot{\mathbf{w}} := \frac{d\mathbf{w}}{dt} = -\nabla_{\mathbf{w}}J(\mathbf{w})$$

- **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**

Architecture ✓
training dynamics ✓

What Gradient Descent gives us?

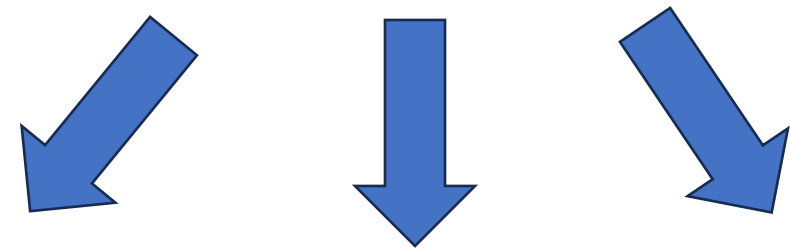
Simple Structures



Sparsity

Low-rank

More complicated structures



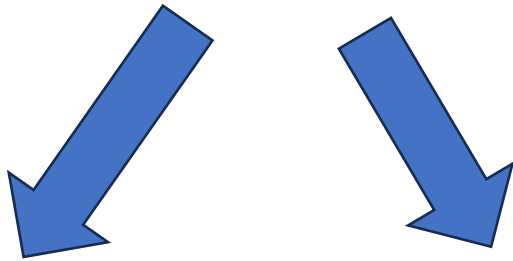
Hierarchical
Representation

Algebraic
Structure

Spectral
Structure

What Gradient Descent gives us?

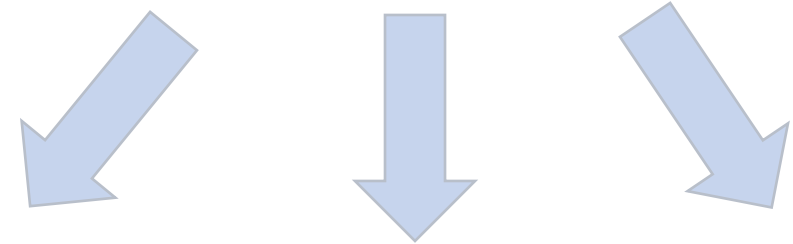
Simple Structures



Sparsity

Low-rank

More complicated structures

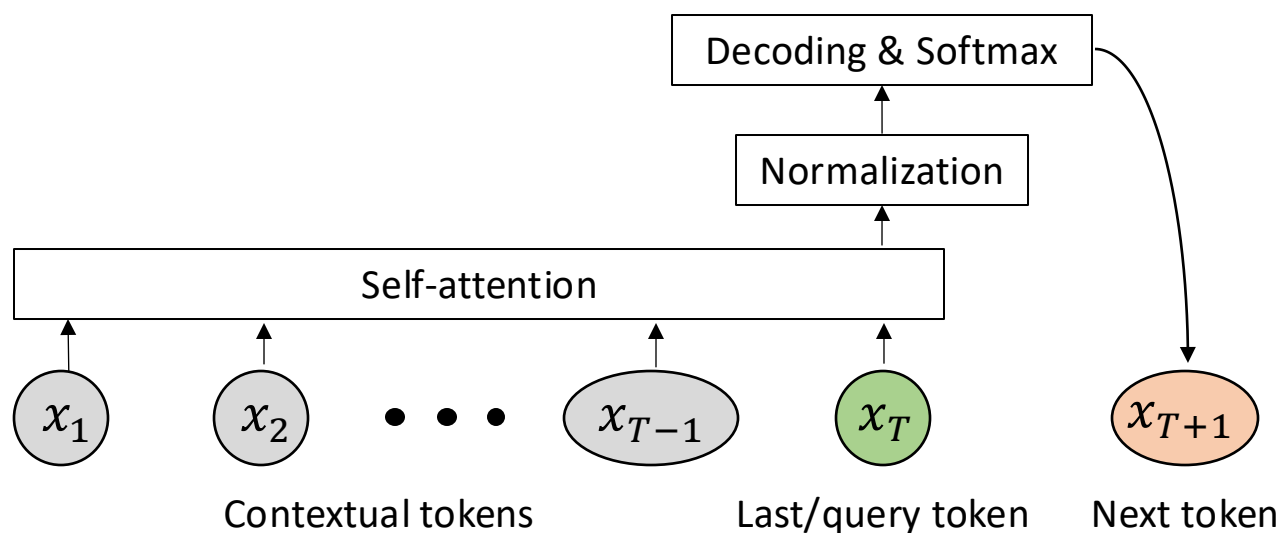


Hierarchical
Representation

Algebraic
Structure

Spectral
Structure

Understanding Attention in 1-layer Setting



$U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M]^T$: token embedding matrix

$$\hat{\mathbf{u}}_T = \sum_{t=1}^{T-1} b_{tT} \mathbf{u}_{x_t} = U^T X^T \mathbf{b}_T$$

Self-attention

$$b_{tT} := \frac{\exp(\mathbf{u}_{x_T}^\top W_Q W_K^\top \mathbf{u}_{x_t} / \sqrt{d})}{\sum_{t=1}^{T-1} \exp(\mathbf{u}_{x_T}^\top W_Q W_K^\top \mathbf{u}_{x_t} / \sqrt{d})}$$

Normalized version $\tilde{\mathbf{u}}_T = U^T \text{LN}(X^T \mathbf{b}_T)$

Objective:

$$\max_{W_K, W_Q, W_V, U} J = \mathbb{E}_D \left[\mathbf{u}_{x_{T+1}}^\top W_V \tilde{\mathbf{u}}_T - \log \sum_l \exp(\mathbf{u}_l^\top W_V \tilde{\mathbf{u}}_T) \right]$$

Reparameterization

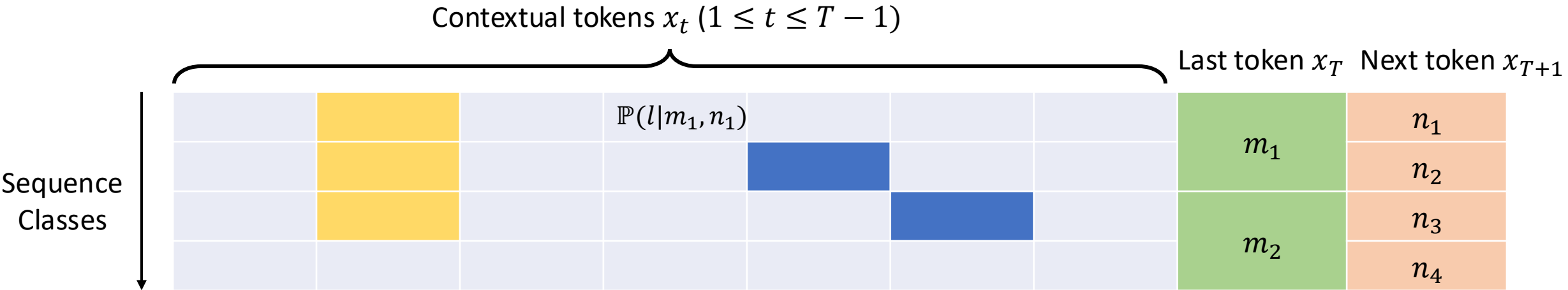
- Parameters W_K, W_Q, W_V, U makes the dynamics complicated.
- Reparameterize the problem with independent variable Y and Z
 - $Y = UW_V^T U^T$
 - $Z = UW_Q W_K^T U^T$ (pairwise logits of self-attention matrix)
- Then the dynamics becomes easier to analyze

Data Distribution

$$x_t \in [M] \text{ for } 1 \leq t \leq T$$

$$x_{T+1} \in [K]$$

$$K \ll M$$



Distinct tokens: There exists unique n so that $\mathbb{P}(l|n) > 0$

Common tokens: There exists multiple n so that $\mathbb{P}(l|n) > 0$

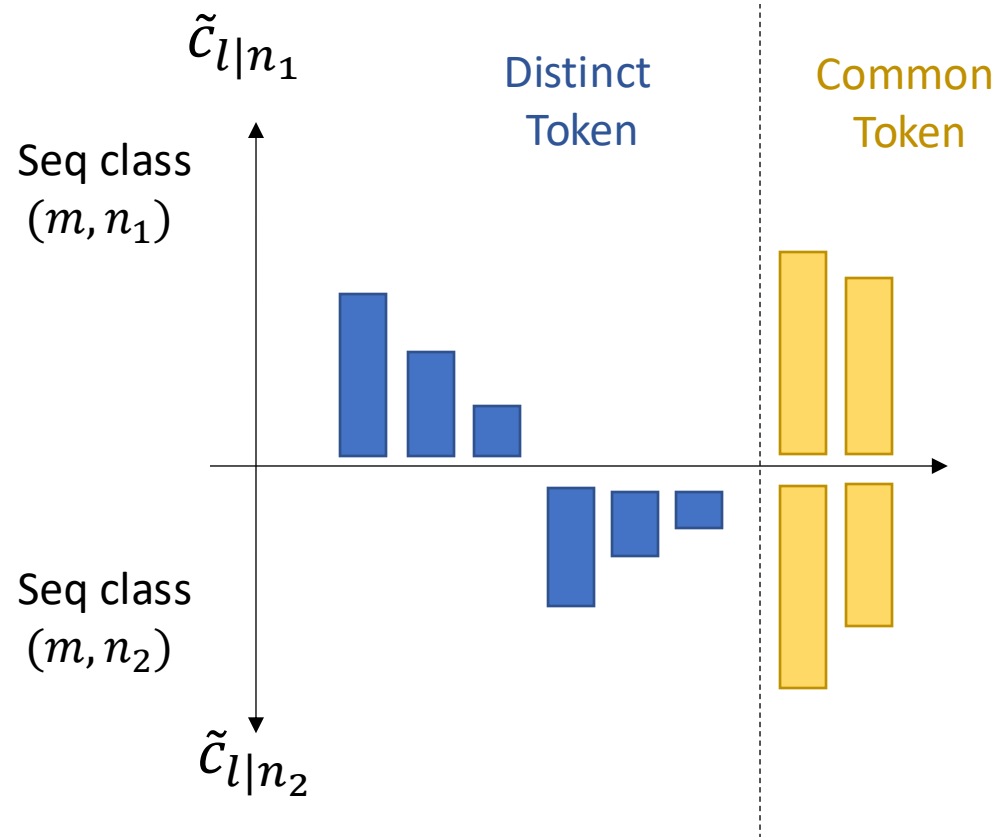
$\mathbb{P}(l|m, n) = \mathbb{P}(l|n)$ is the conditional probability of token l given last token $x_T = m$ and $x_{T+1} = n$

Assumption: $m = \psi(n)$, i.e., no next token shared among different last tokens

Question: Given the data distribution, how does the self-attention layer behave?

Overall Picture of the Training Dynamics

At initialization



Co-occurrence probability

$$\tilde{c}_{l|n_1} := \mathbb{P}(l|m, n_1) \exp(z_{ml})$$

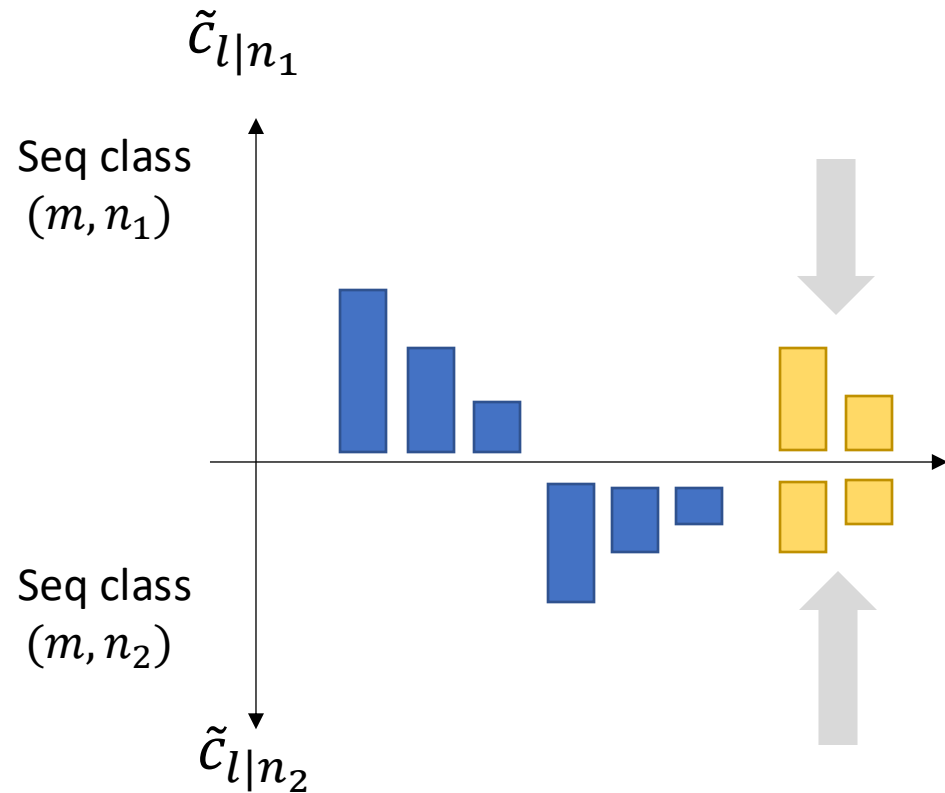
Initial condition: $z_{ml}(0) = 0$

$$Z = \begin{matrix} \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} & \mathbf{z}_m \end{matrix}$$

\mathbf{z}_m : All logits of the contextual tokens when attending to last token $x_T = m$

Overall Picture of the Training Dynamics

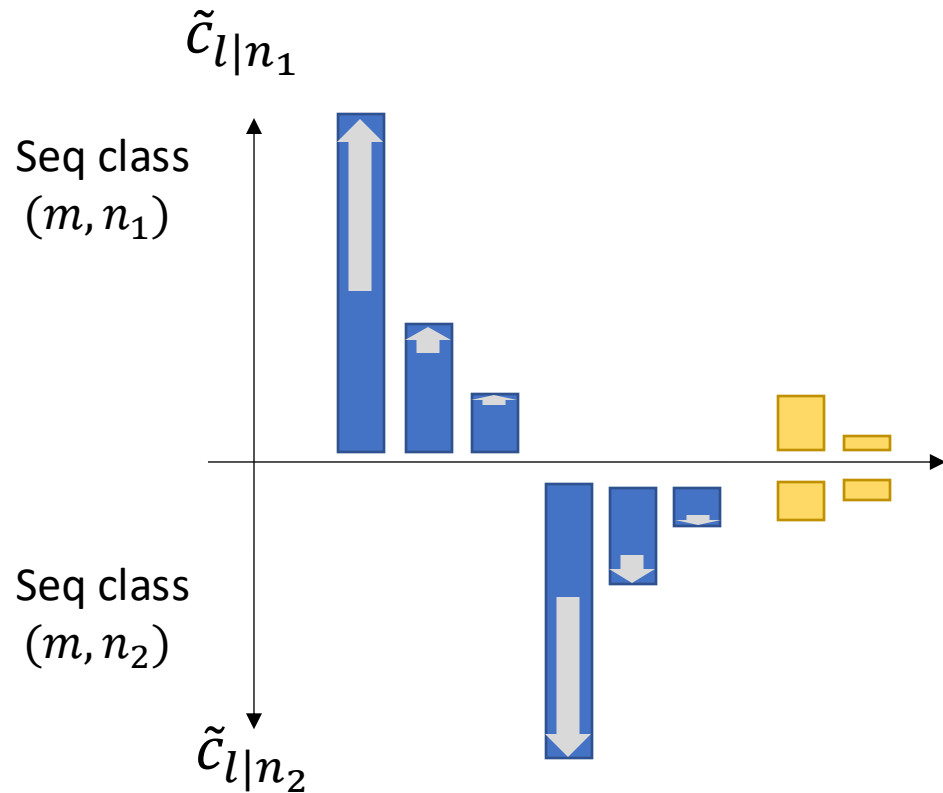
Common Token Suppression



(a) $\dot{z}_{ml} < 0$, for common token l

Overall Picture of the Training Dynamics

Winners-emergence



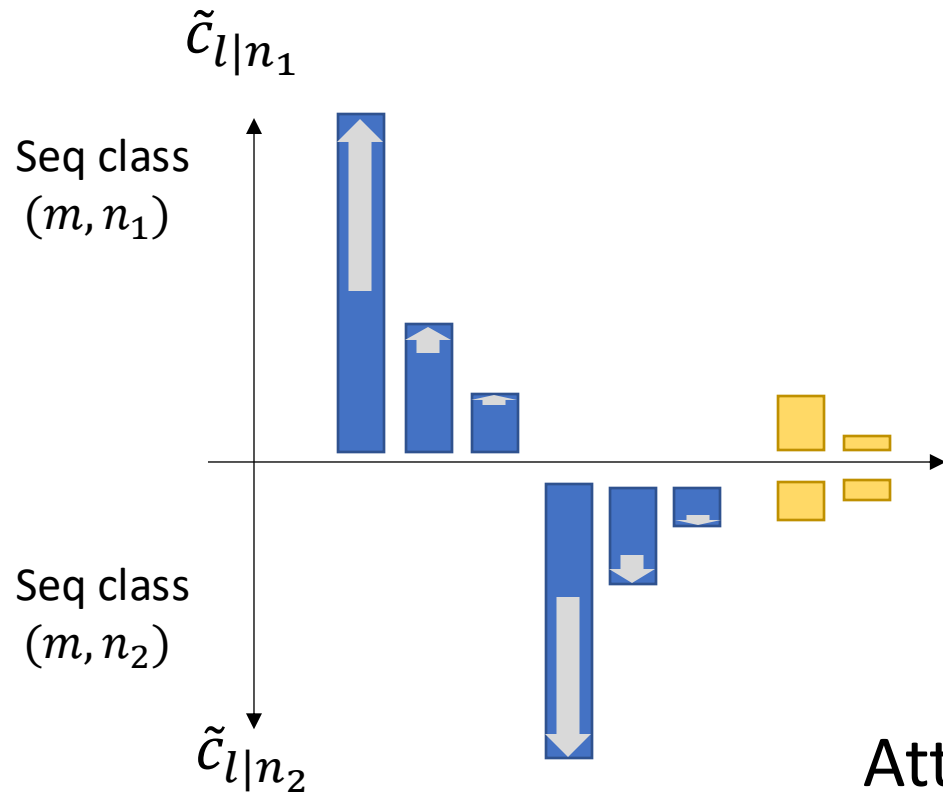
(a) $z_{ml}^{\cdot} < 0$, for **common token** l

(b) $z_{ml}^{\cdot} > 0$, for **distinct token** l

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

Overall Picture of the Training Dynamics

Winners-emergence



(a) $z_{ml} \dot{< 0$, for common token l

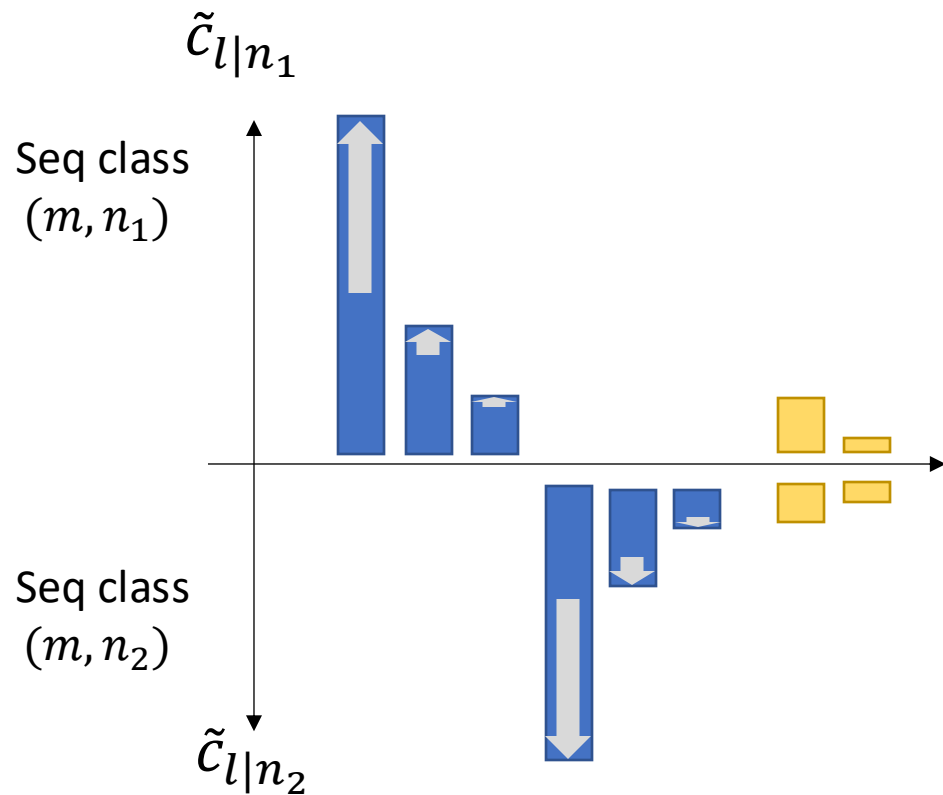
(b) $z_{ml} \dot{> 0$, for distinct token l

(c) $z_{ml}(t)$ grows faster with larger $\mathbb{P}(l|m, n)$

Attention looks for **discriminative** tokens that **frequently co-occur** with the query.

Overall Picture of the Training Dynamics

Winners-emergence



(c) $z_{ml}(t)$ grows faster with larger $\mathbb{P}(l|m, n)$

Theorem 3 Relative gain $r_{l/l'|n}(t) := \frac{\tilde{c}_{l|n}^2(t)}{\tilde{c}_{l'|n}^2(t)} - 1$ has a close form:

$$r_{l/l'|n}(t) = r_{l/l'|n}(0)\chi_l(t)$$

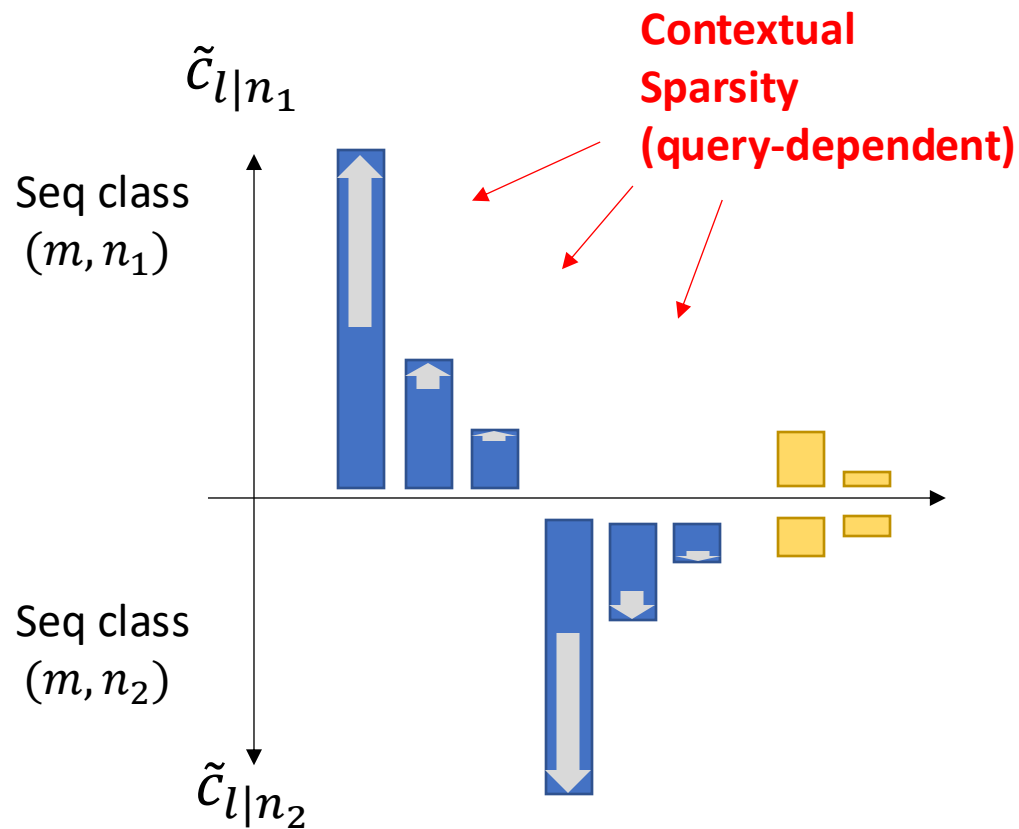
If l_0 is the dominant token: $r_{l_0/l|n}(0) > 0$ for all $l \neq l_0$ then

$$e^{2f_{nl_0}^2(0)B_n(t)} \leq \chi_{l_0}(t) \leq e^{2B_n(t)}$$

where $B_n(t) \geq 0$ monotonously increases, $B_n(0) = 0$

Overall Picture of the Training Dynamics

Winners-emergence



(c) $z_{ml}(t)$ grows faster with larger $\mathbb{P}(l|m, n)$

Theorem 3 Relative gain $r_{l/l'|n}(t) := \frac{\tilde{c}_{l|n}^2(t)}{\tilde{c}_{l'|n}^2(t)} - 1$ has a close form:

$$r_{l/l'|n}(t) = r_{l/l'|n}(0)\chi_l(t)$$

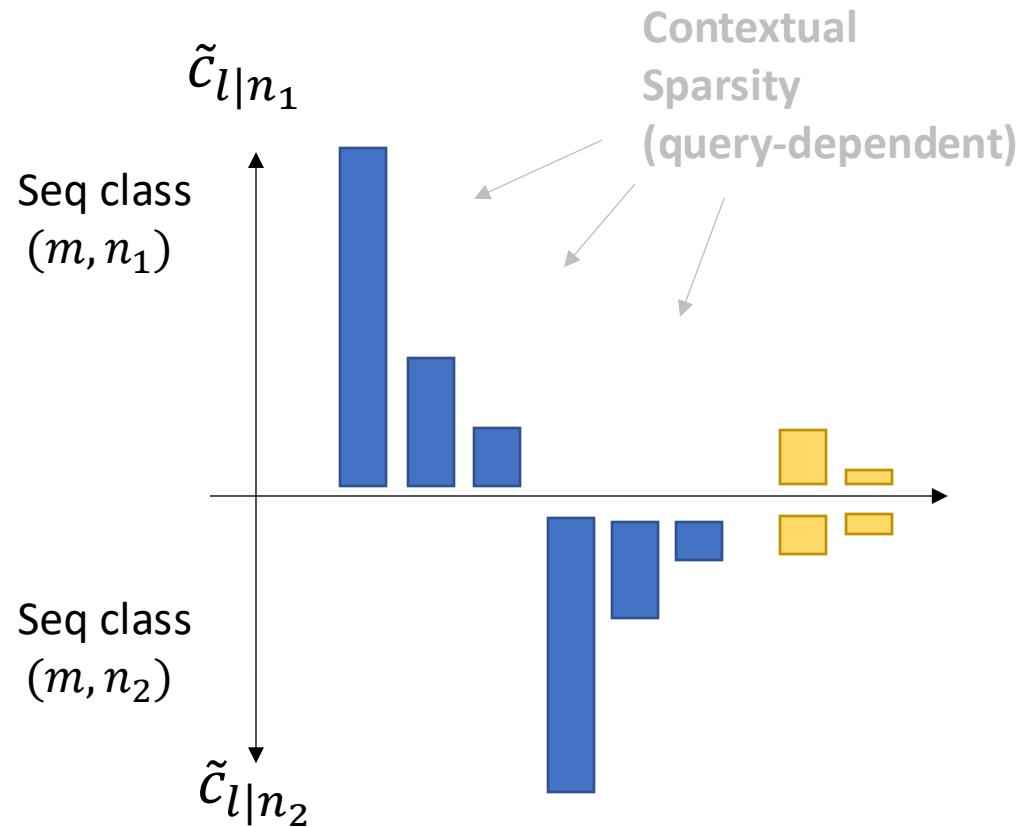
If l_0 is the dominant token: $r_{l_0/l|n}(0) > 0$ for all $l \neq l_0$ then

$$e^{2f_{nl_0}^2(0)B_n(t)} \leq \chi_{l_0}(t) \leq e^{2B_n(t)}$$

where $B_n(t) \geq 0$ monotonously increases, $B_n(0) = 0$

Overall Picture of the Training Dynamics

Attention frozen



Theorem 4 When $t \rightarrow +\infty$,

$$B_n(t) \sim \ln \left(C_0 + 2K \frac{\eta_z}{\eta_Y} \ln^2 \left(\frac{M\eta_Y t}{K} \right) \right)$$

Attention scanning:

When training starts, $B_n(t) = O(\ln t)$

Attention snapping:

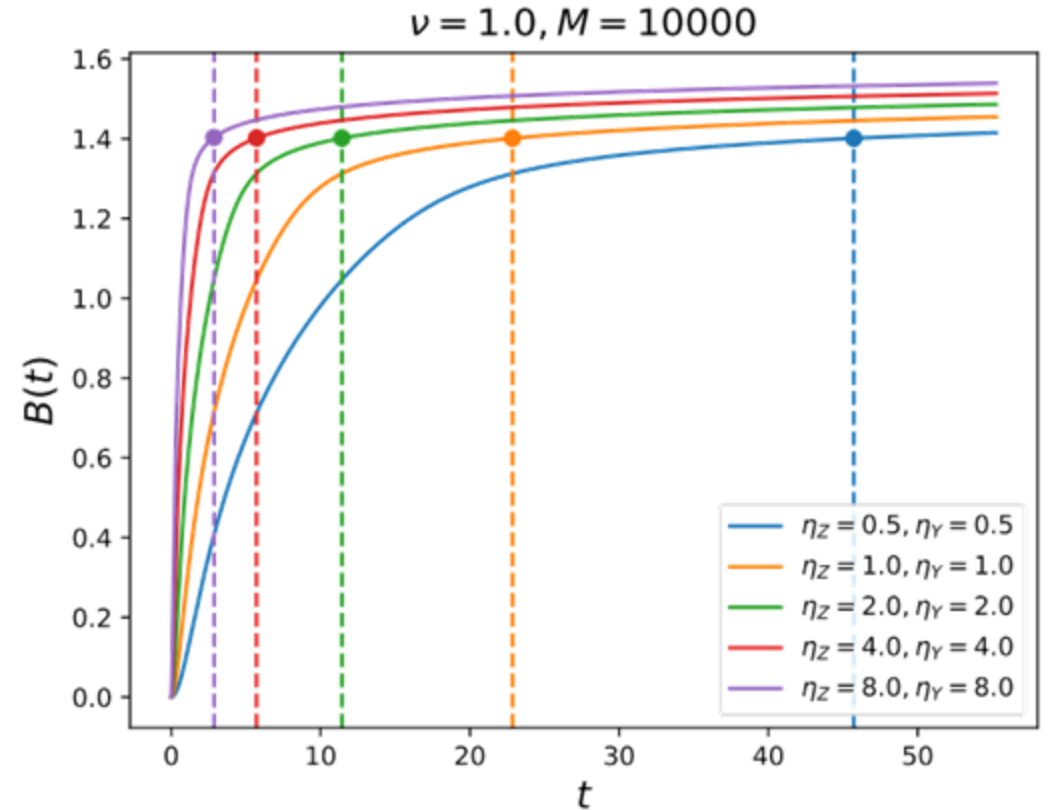
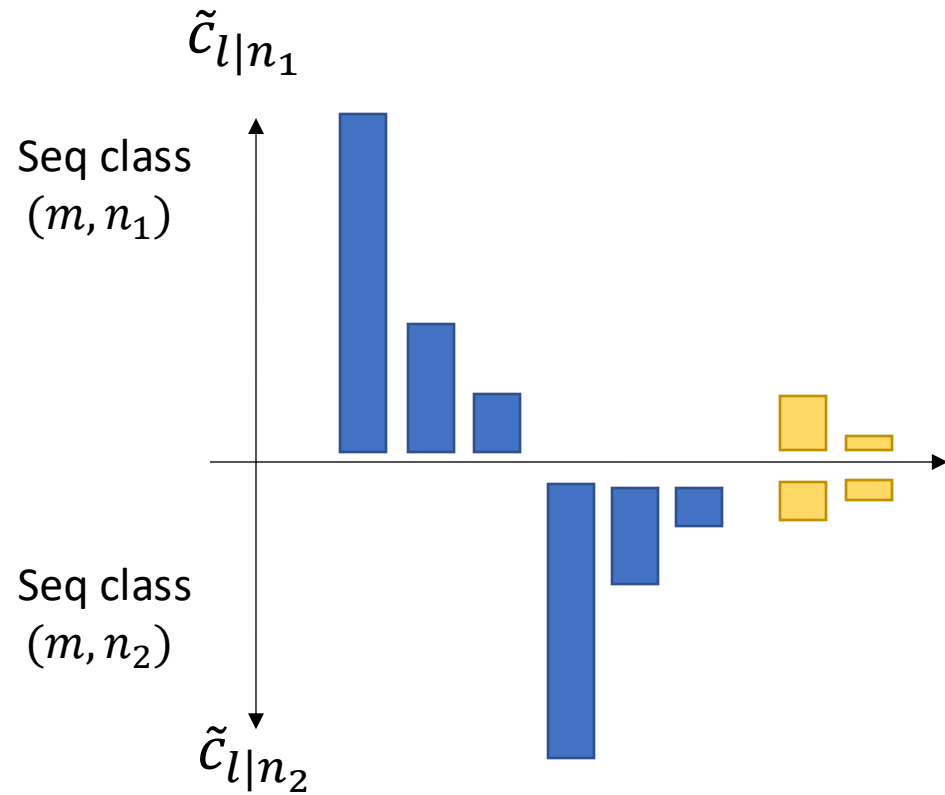
When $t \geq t_0 = O\left(\frac{2K \ln M}{\eta_Y}\right)$, $B_n(t) = O(\ln \ln t)$

(1) η_z and η_Y are large, $B_n(t)$ is large and attention is sparse

(2) Fixing η_z , large η_Y leads to slightly small $B_n(t)$ and denser attention

Overall Picture of the Training Dynamics

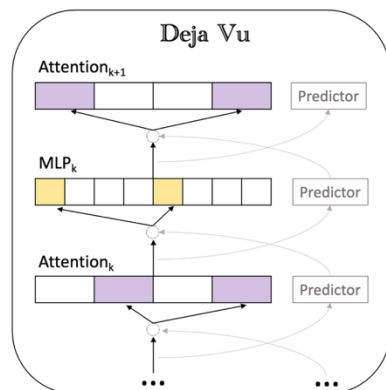
Attention frozen



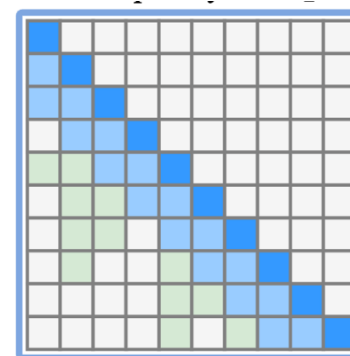
Larger learning rate η_z leads to faster phase transition

$$B_n(t) \sim \ln \left(C_0 + 2K \frac{\eta_z}{\eta_\gamma} \ln^2 \left(\frac{M\eta_\gamma t}{K} \right) \right)$$

Further Study of Sparse Attention / Low Rank

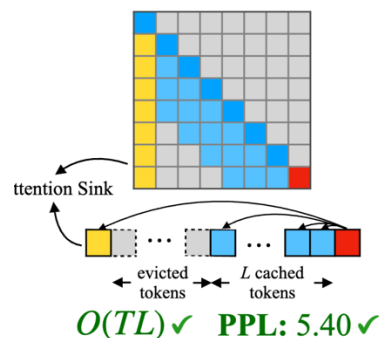


DeJaVu [Z. Liu et al, ICML'23 (oral)]

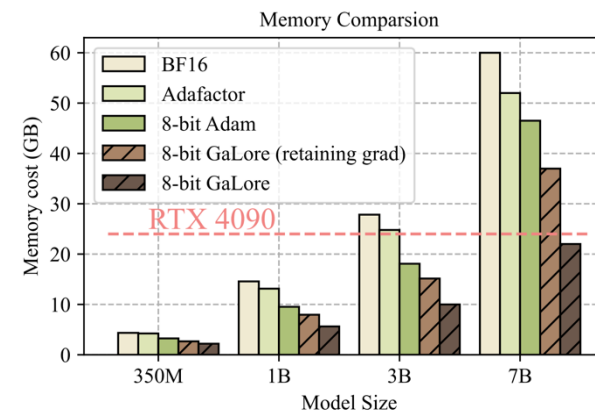


H2O [Z. Zhang et al, NeurIPS'23]

(d) StreamingLLM (ours)



Attention Sink [G. Xiao et al, ICLR'24]



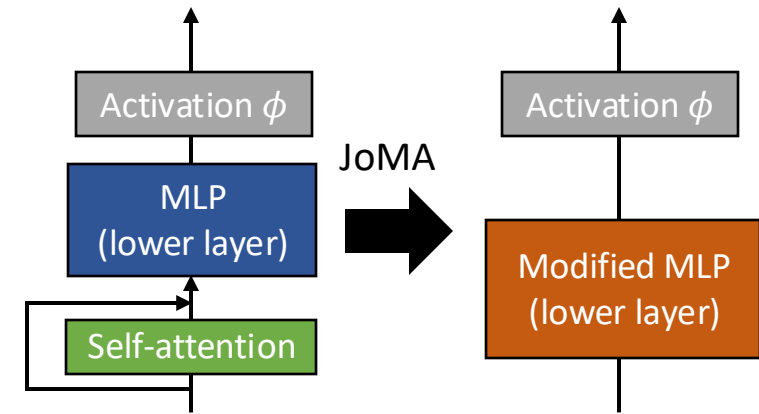
GaLore [J. Zhao et al, ICML'24 (Oral)]

How to get rid of the assumptions?

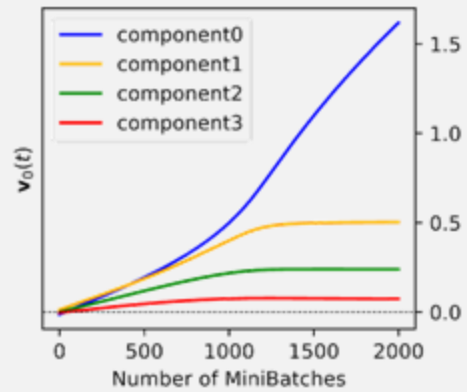
- A few annoying assumptions in the analysis
 - No residual connections
 - No embedding vectors
 - The decoder needs to learn faster than the self-attention ($\eta_Y \gg \eta_Z$).
 - Single layer analysis
- How to get rid of them?
- Follow-up work: **JoMA** (Joint MLP/Attention dynamics)

Folding self-attention into MLP

A Transformer block becomes a modified MLP
(due to nice property in gradient dynamics)

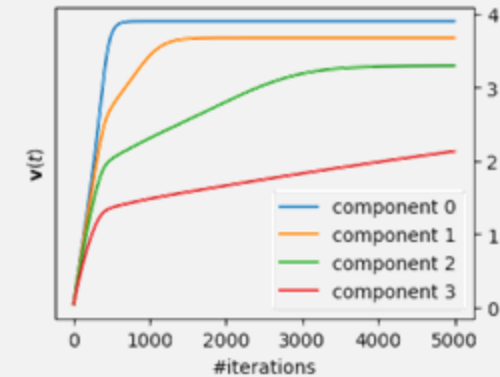


Linear case ($\phi = \text{Id}, K = 1$)



Most salient feature takes all
(Attention becomes sparser)
(Consistent with Scan&Snap)

Nonlinear case (ϕ nonlinear, $K = 1$)



Most salient feature grows, and others catch up
(Attention becomes sparser and denser)

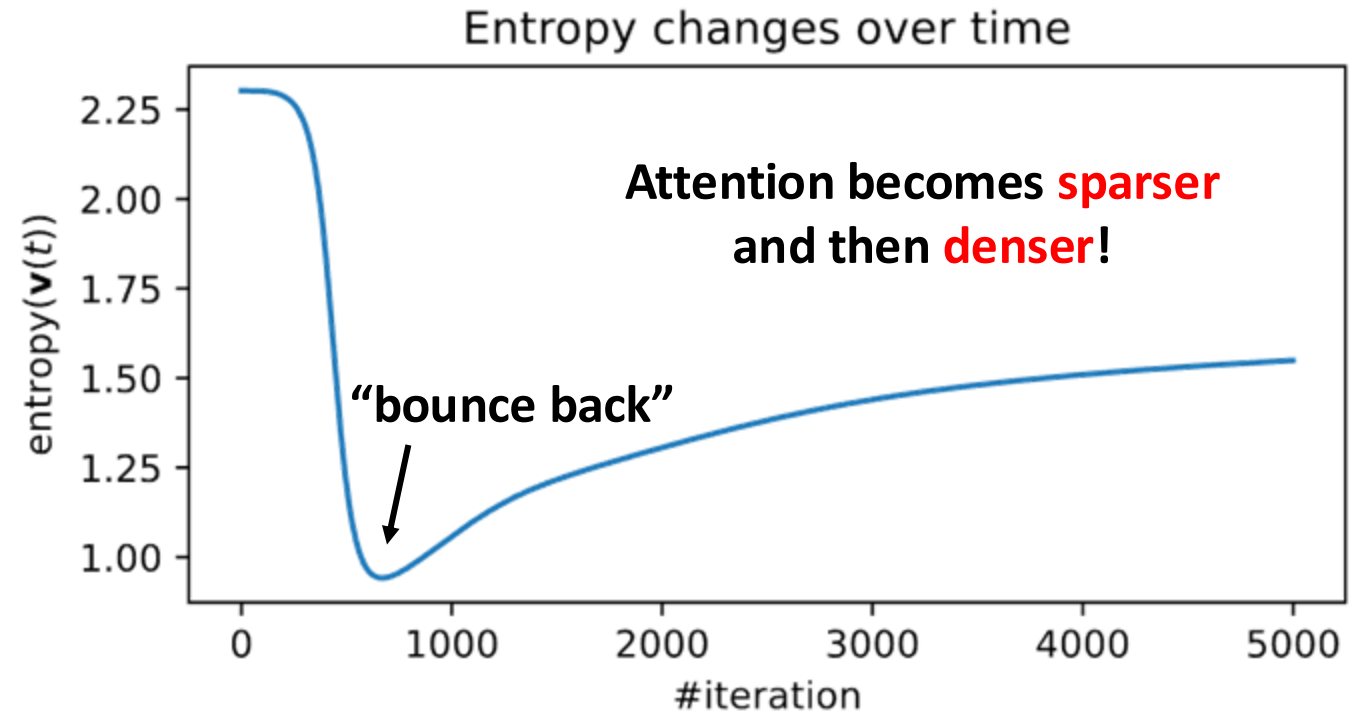
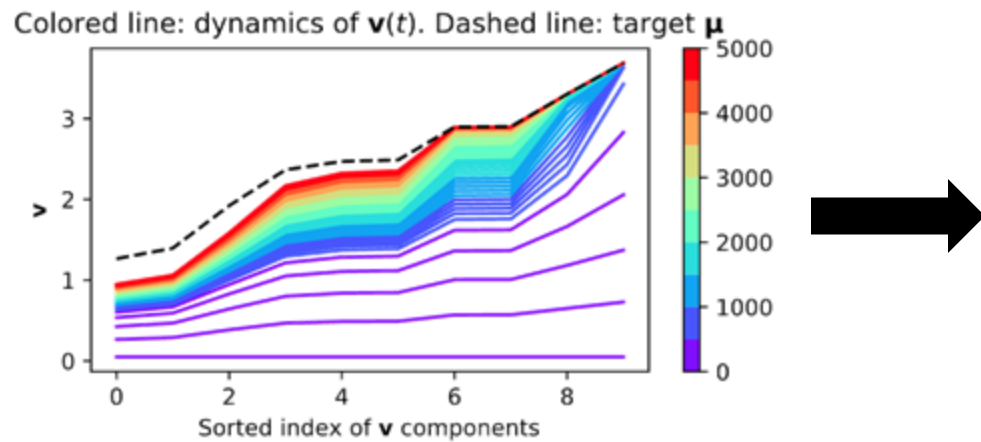
Saliency is defined as $\Delta_{lm} = \mathbb{E}[g|l, m] \cdot \mathbb{P}[l|m]$

↑ Discriminancy ↑ CoOccurrence

$\Delta_{lm} \approx 0$: **Common** tokens
 $|\Delta_{lm}|$ large: **Distinct** tokens

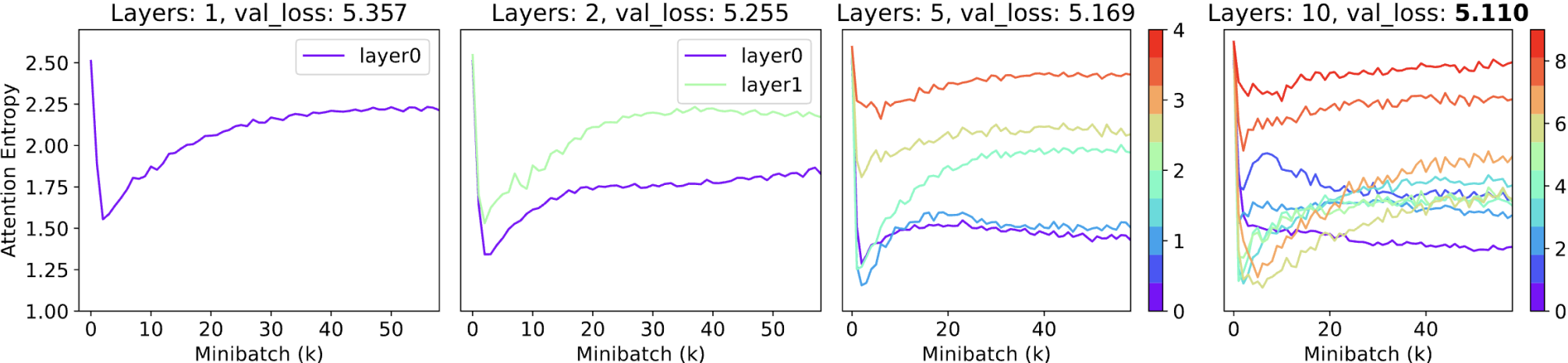
Nonlinear case

$\dot{\mathbf{v}} = (\boldsymbol{\mu} - \mathbf{v}) \circ \exp\left(\frac{\mathbf{v}^2}{2}\right)$	Nonlinear
	Modified MLP (lower layer)

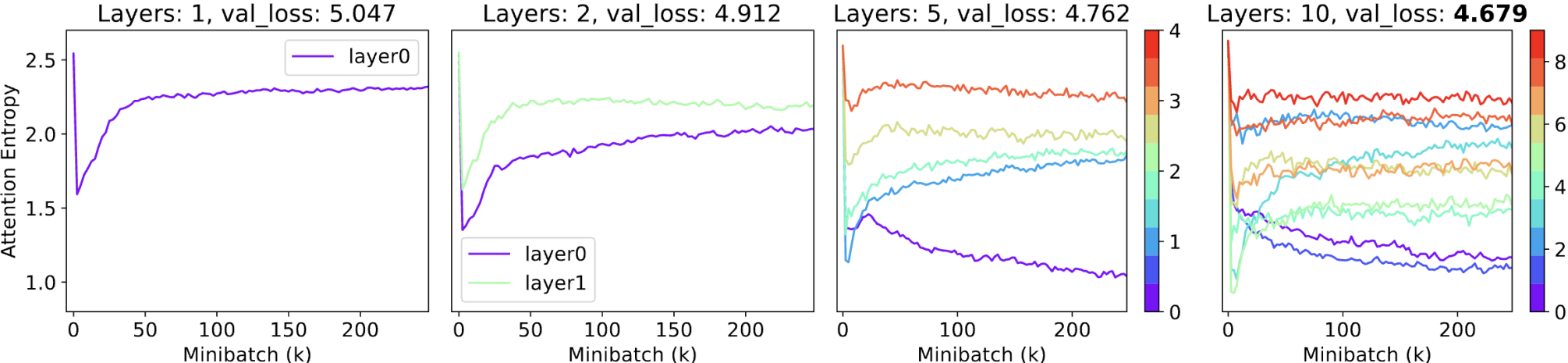


Real-world Experiments

Wiktext2



Wiktext103

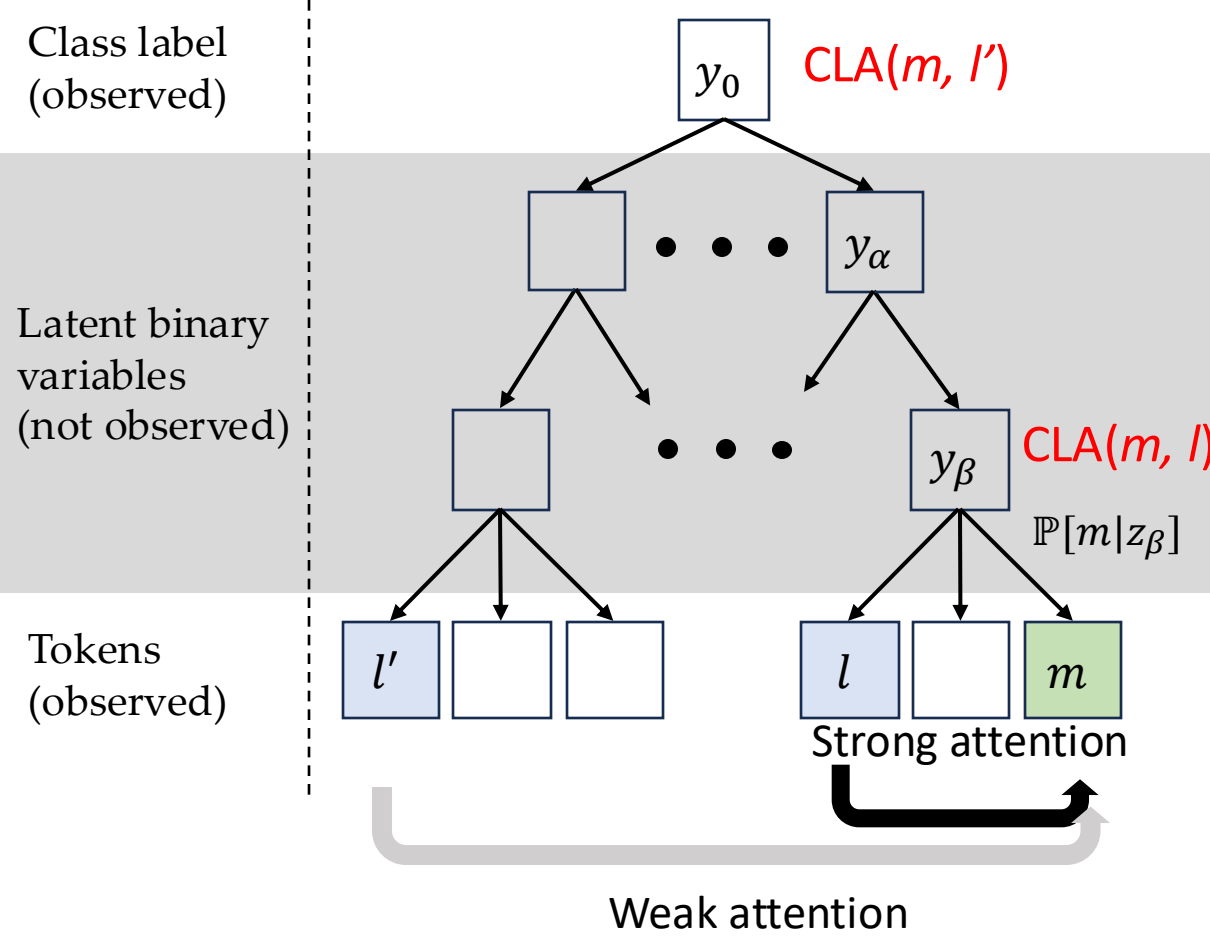


Why is this “bouncing back” property useful?

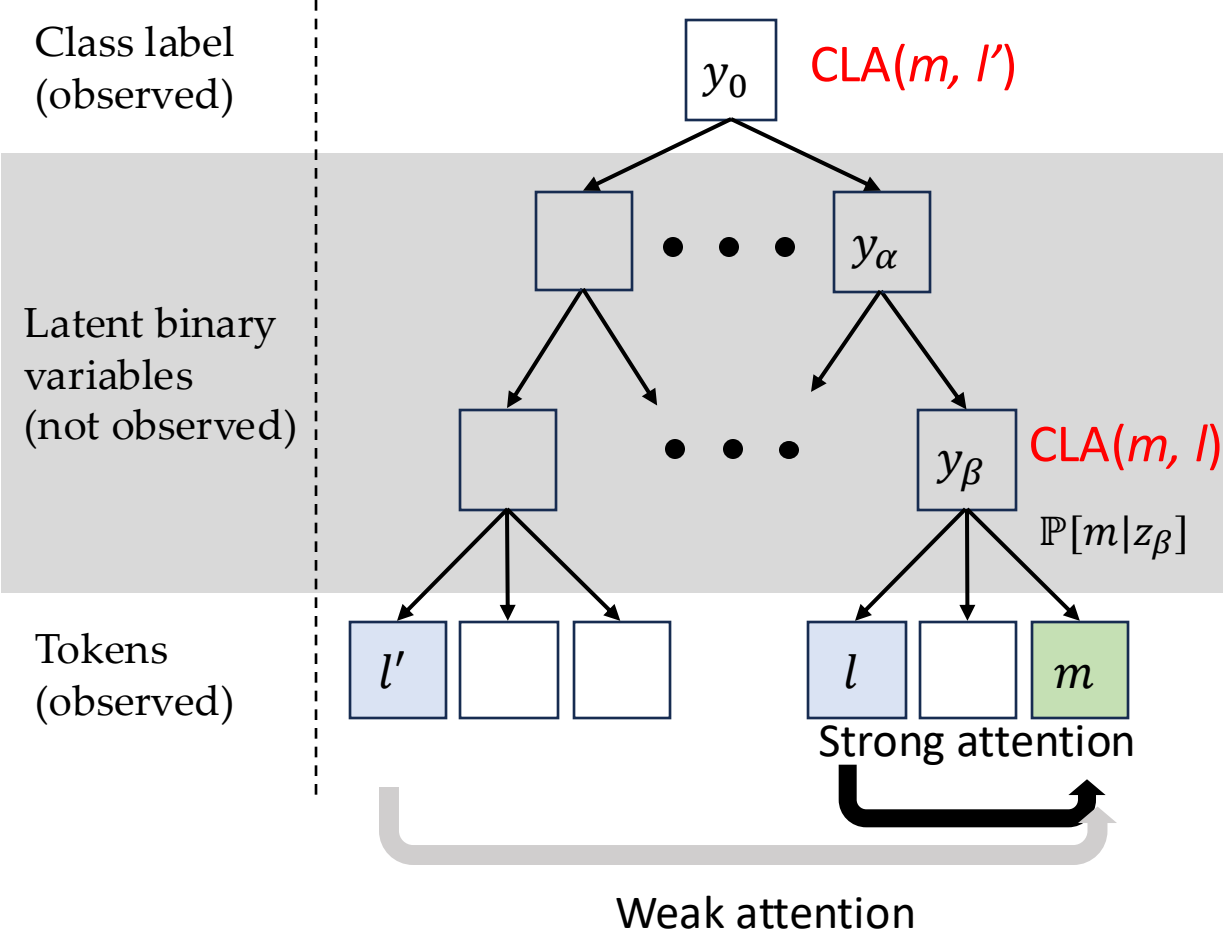
It seems that it only slows down the training??

Not useful in 1-layer, but useful in multiple Transformer layers!

Data Hierarchy & Multilayer Transformer



Data Hierarchy & Multilayer Transformer



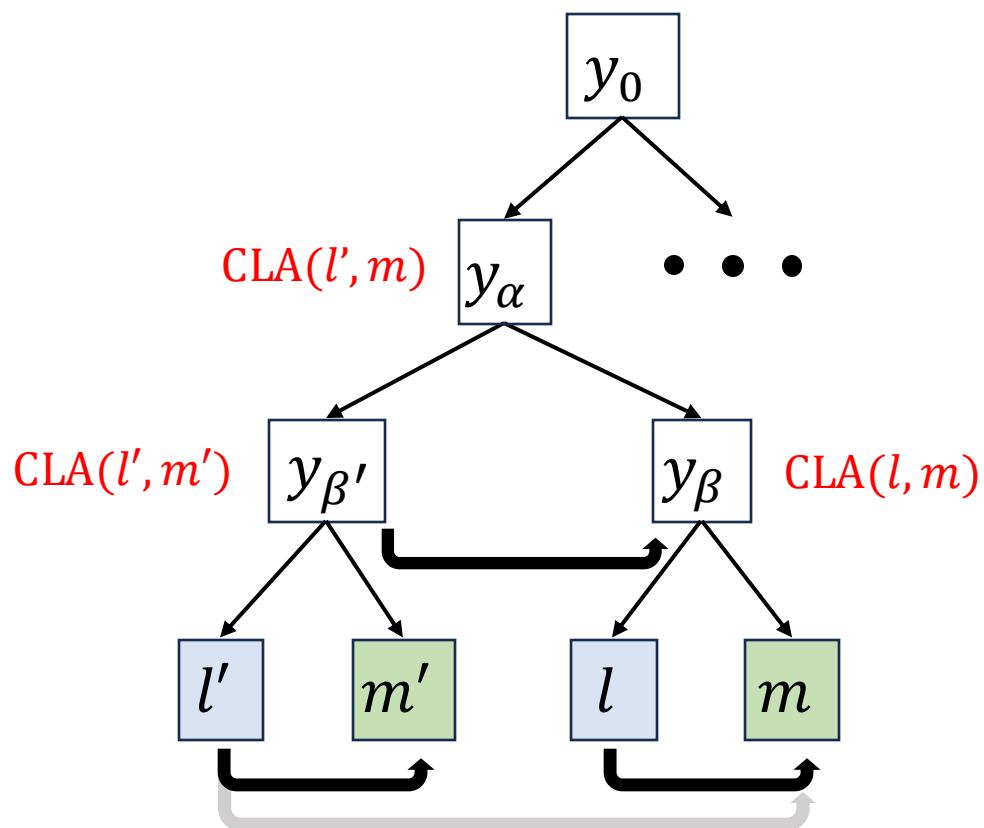
Theorem 5

$$\mathbb{P}[l|m] \approx 1 - \frac{H}{L}$$

H : height of the common latent ancestor (CLA) of l & m

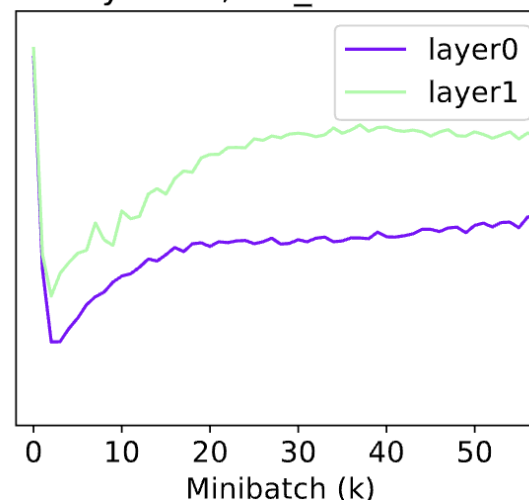
L : total height of the hierarchy

Deep Latent Distribution

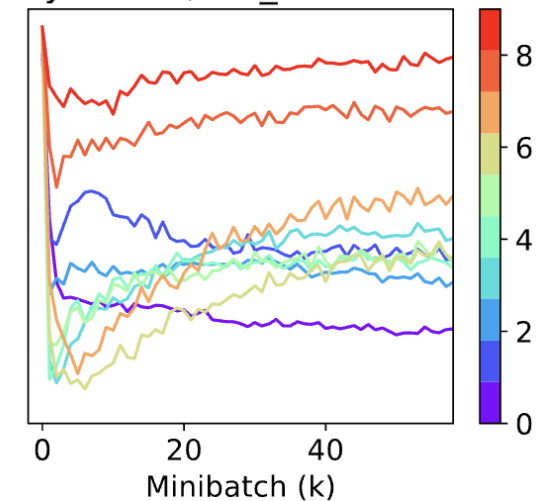


Strong Attention
Weak Attention

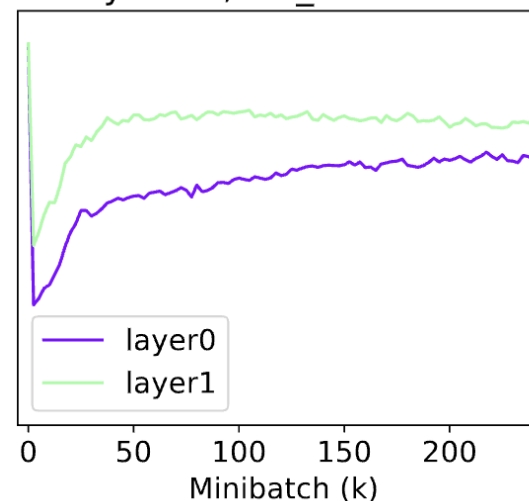
Layers: 2, val_loss: 5.255



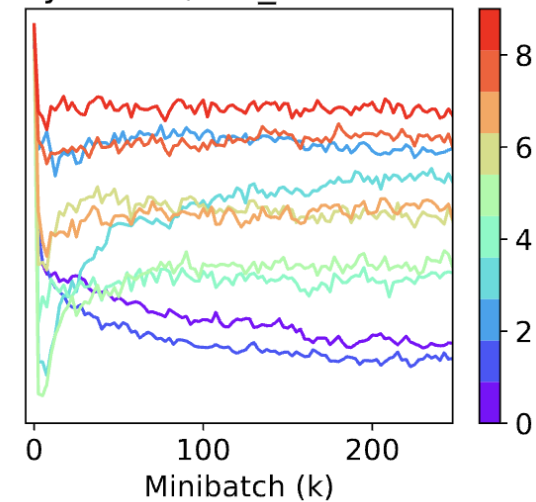
Layers: 10, val_loss: **5.110**



Layers: 2, val_loss: 4.912



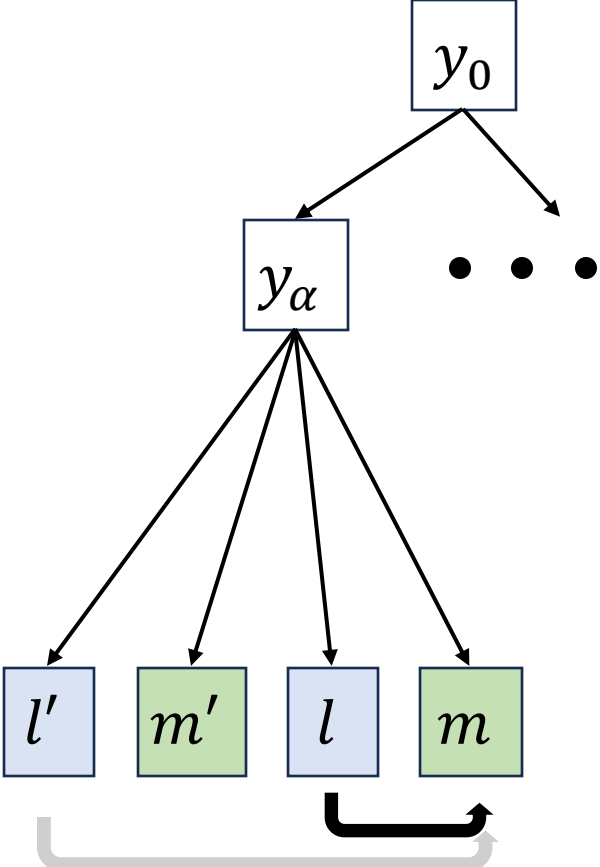
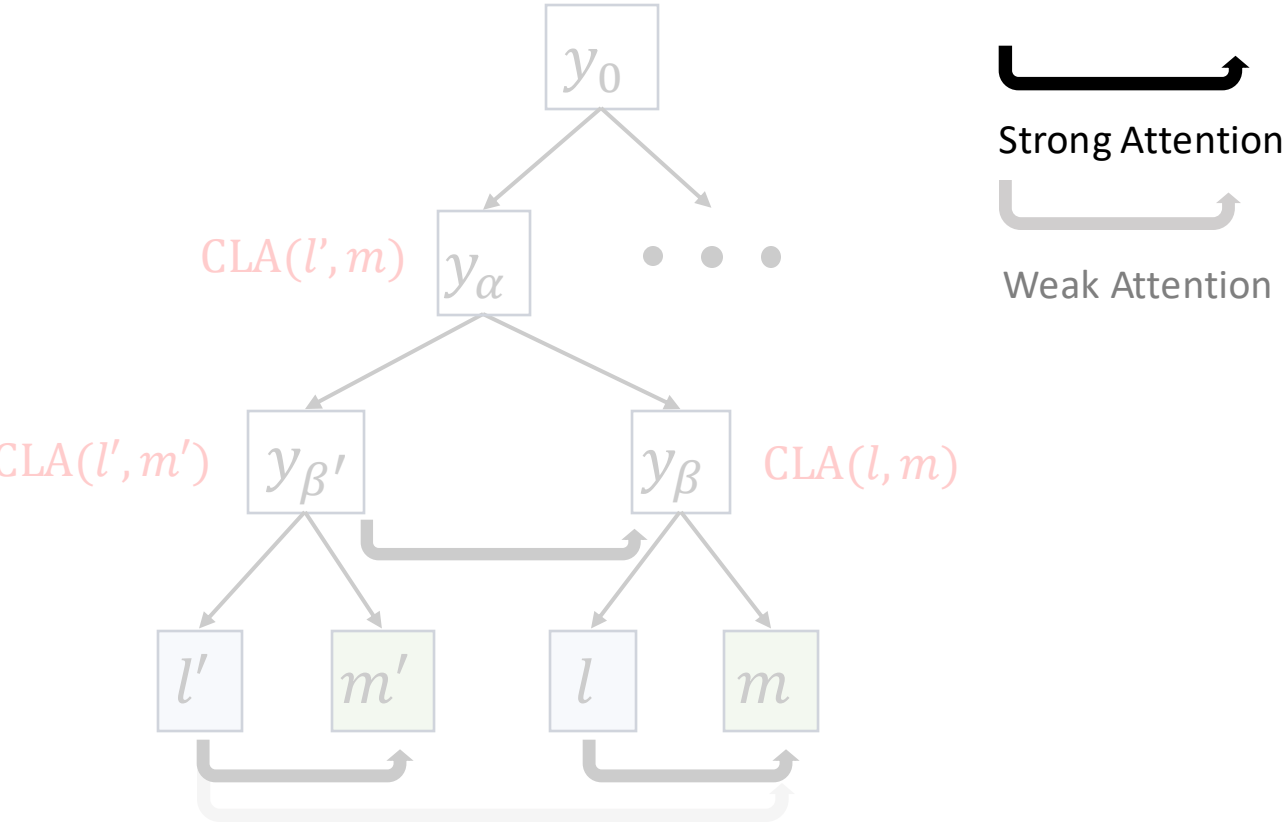
Layers: 10, val_loss: **4.679**



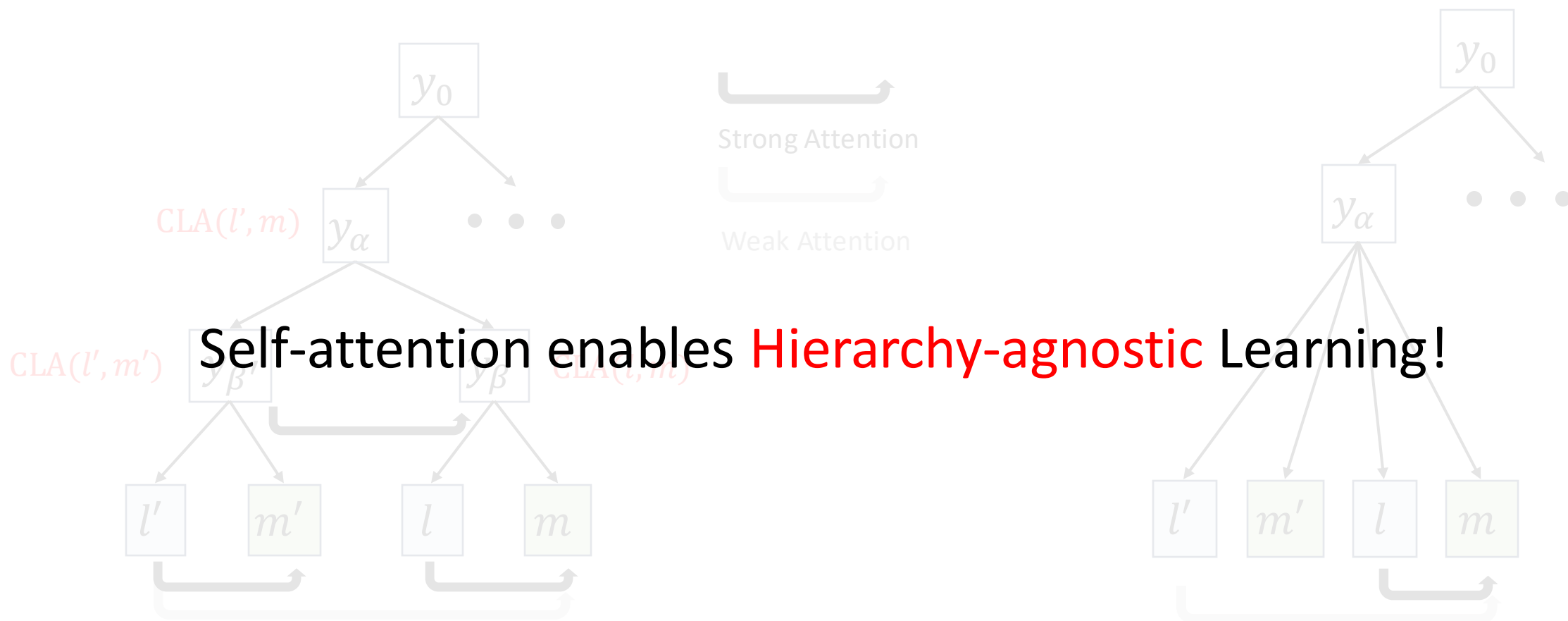
Learning the current hierarchical structure by

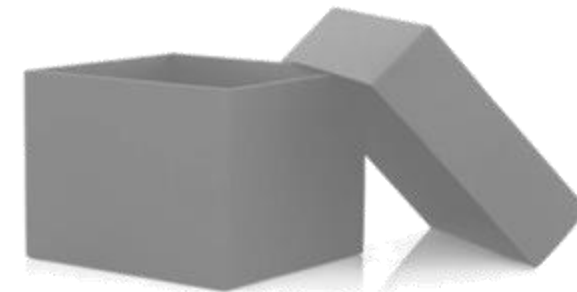
slowing down the association of tokens that are not directly correlated

Shallow Latent Distribution



Hierarchy-agnostic Learning



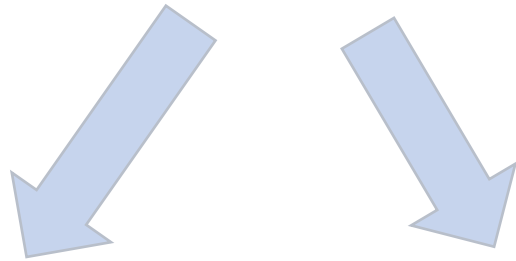


Take away messages

- Architecture ✓ training dynamics ✓
- Nonlinearity is not formidable!
 - Transformer can be analyzed following gradient descent rules
- Property of self-attention
 - Attention becomes sparse over training
 - Inductive bias
 - Favor the learning of strong co-occurred tokens
 - Deter the learning of weakly co-occurred tokens, avoiding spurious correlation.
- Key insights lead to broad applications

What Gradient Descent gives us?

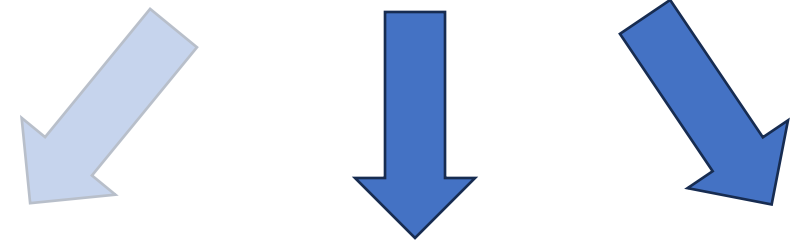
Simple Structures



Sparsity

Low-rank

More complicated structures



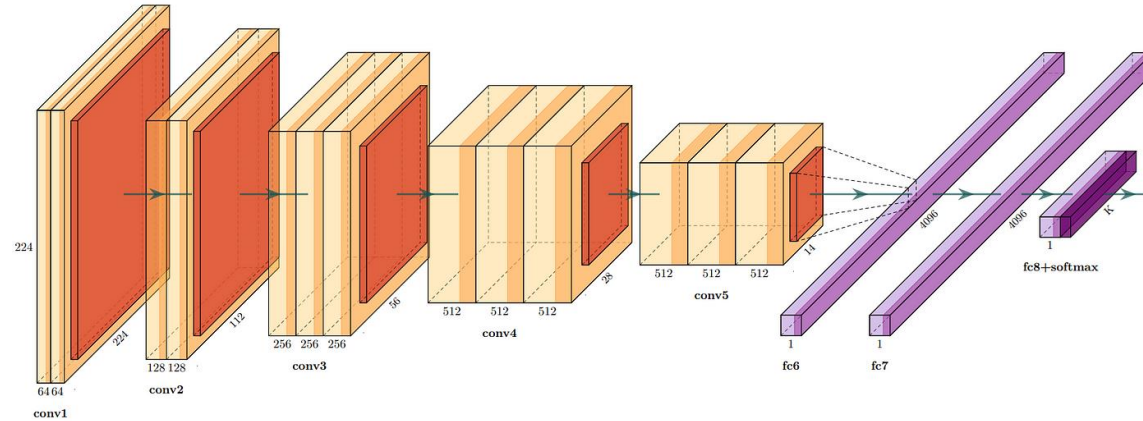
Hierarchical
Representation

Algebraic
Structure

Spectral
Structure

Dichotomy: Symbolic and Neural Representation

Neural
Representation



Symbolic
Representation

$$\nabla \cdot \mathbf{E} = \frac{\rho_v}{\epsilon} \quad (\text{Gauss' Law})$$

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

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

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

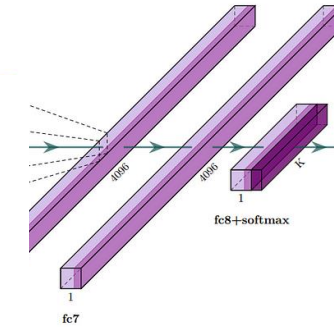
Unification of Symbolic and Neural Representation

Emerging Symbolic Structure

Neural
Repres



Symbo
Repres



tism)

Deep Models (Minsky's Law)

Concrete Example: Modular Addition

$$a + b = c \pmod{d}$$

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

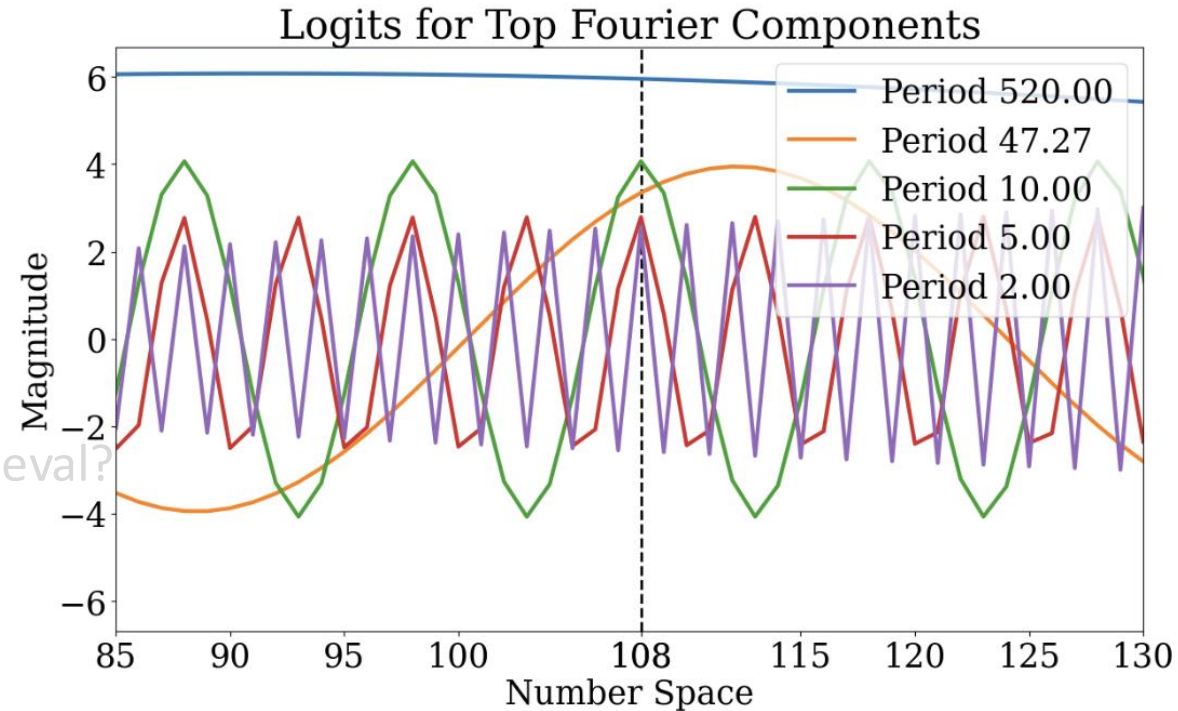
Concrete Example: Modular Addition

$$a + b = c \pmod{d}$$

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

Learned representation = Fourier basis 🤖

Why? 🤔



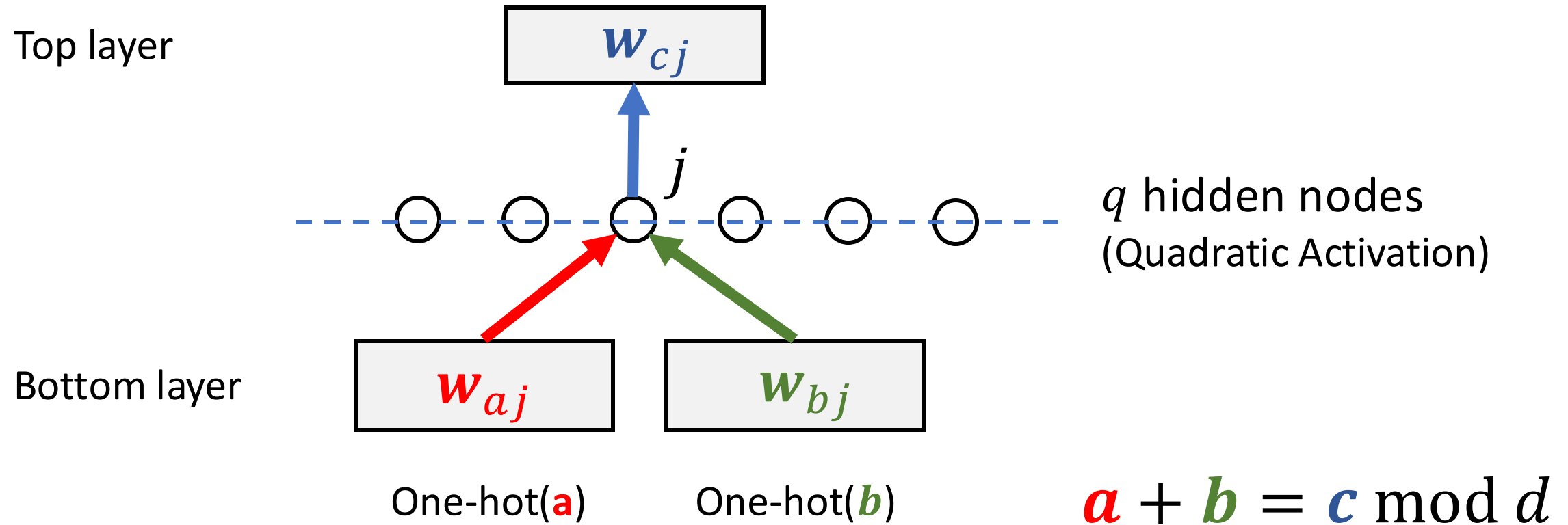
(a) Final logits for top Fourier components

[T. Zhou et al, *Pre-trained Large Language Models Use Fourier Features to Compute Addition*, NeurIPS'24]

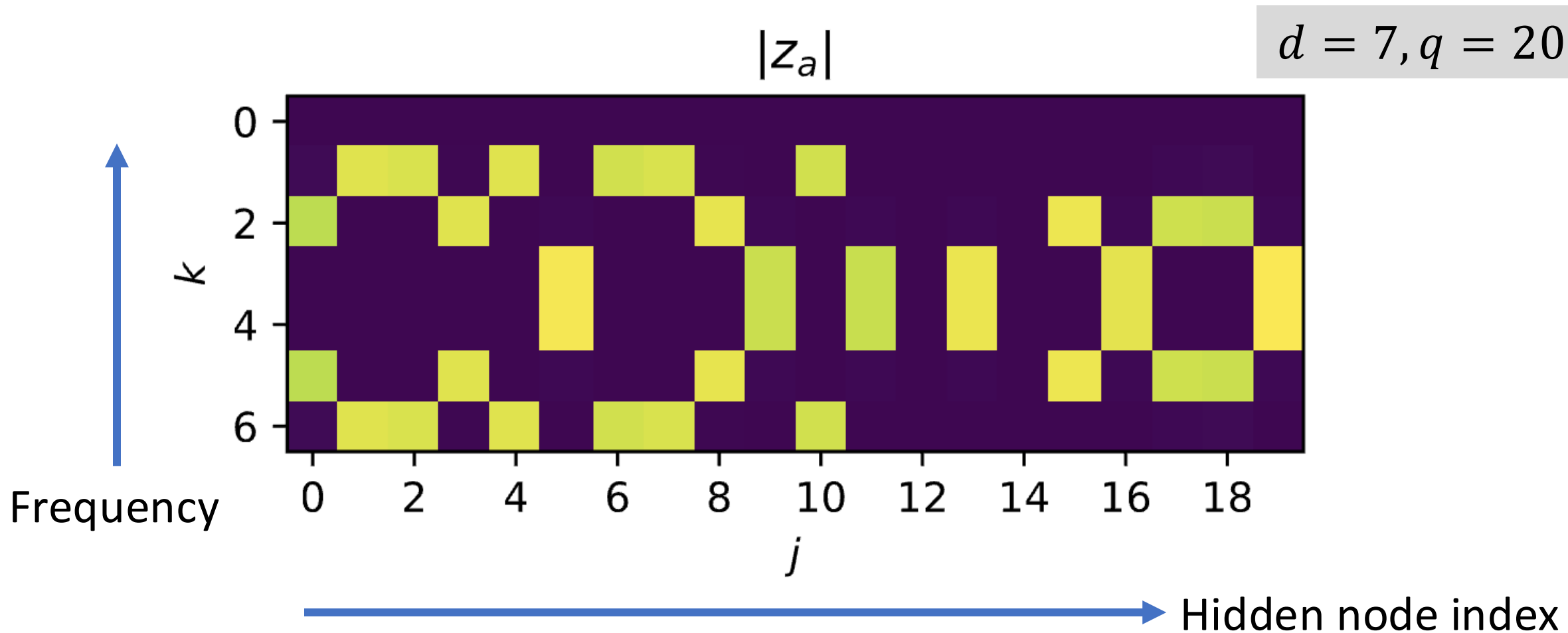
[S. Kantamneni, *Language Models Use Trigonometry to Do Addition*, arXiv'25]

Minimal Setup

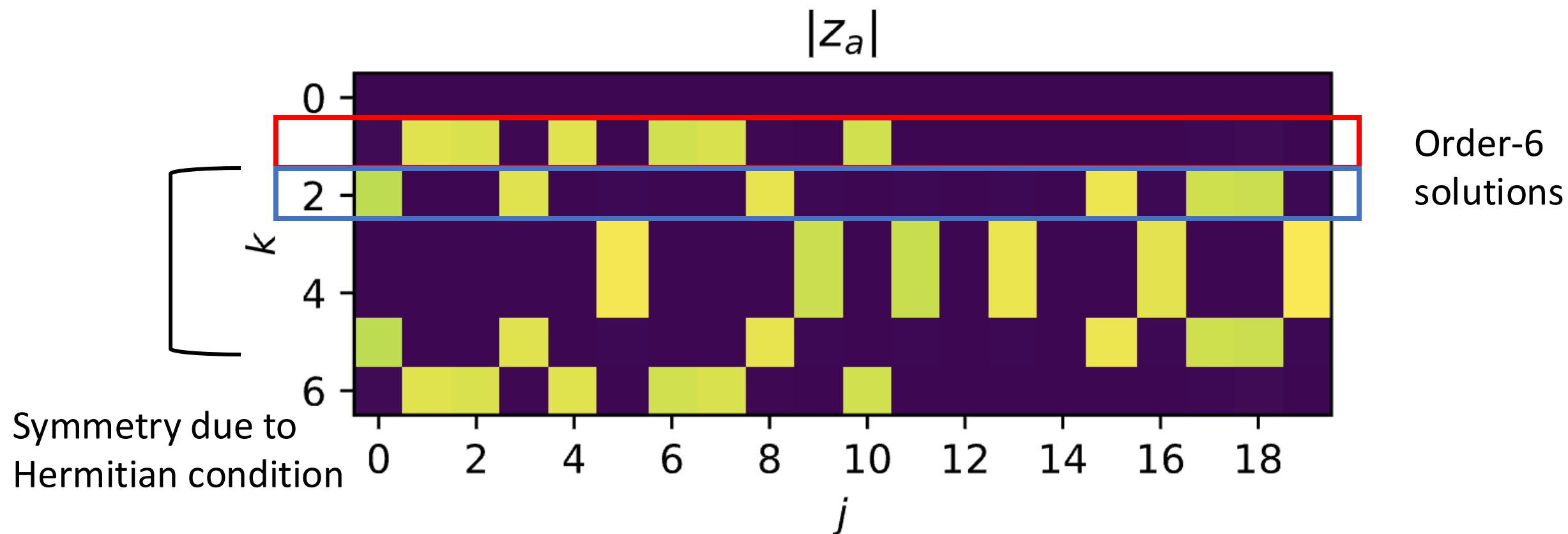
MSE Loss: $\text{Min } \|\text{Output} - \text{one-hot}(\mathbf{c})\|_2$



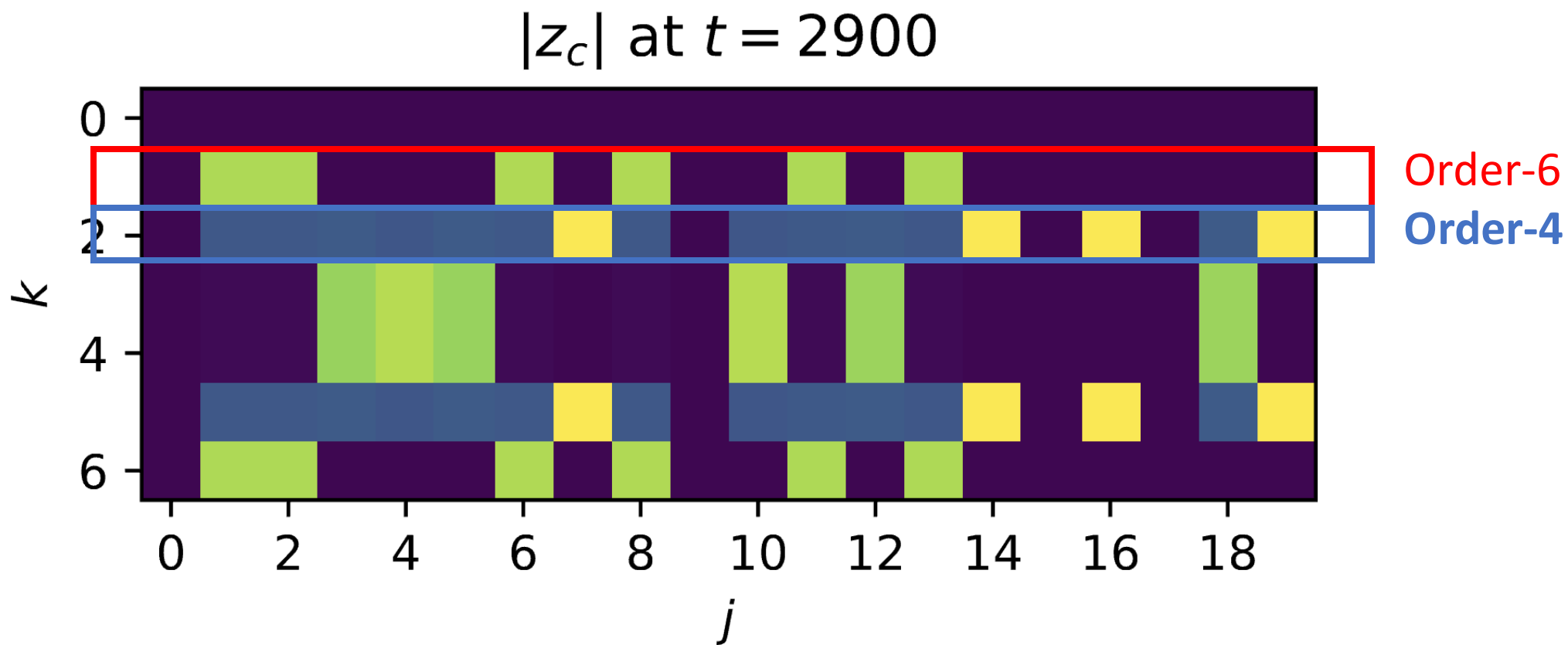
What a Gradient Descent Solution look like?



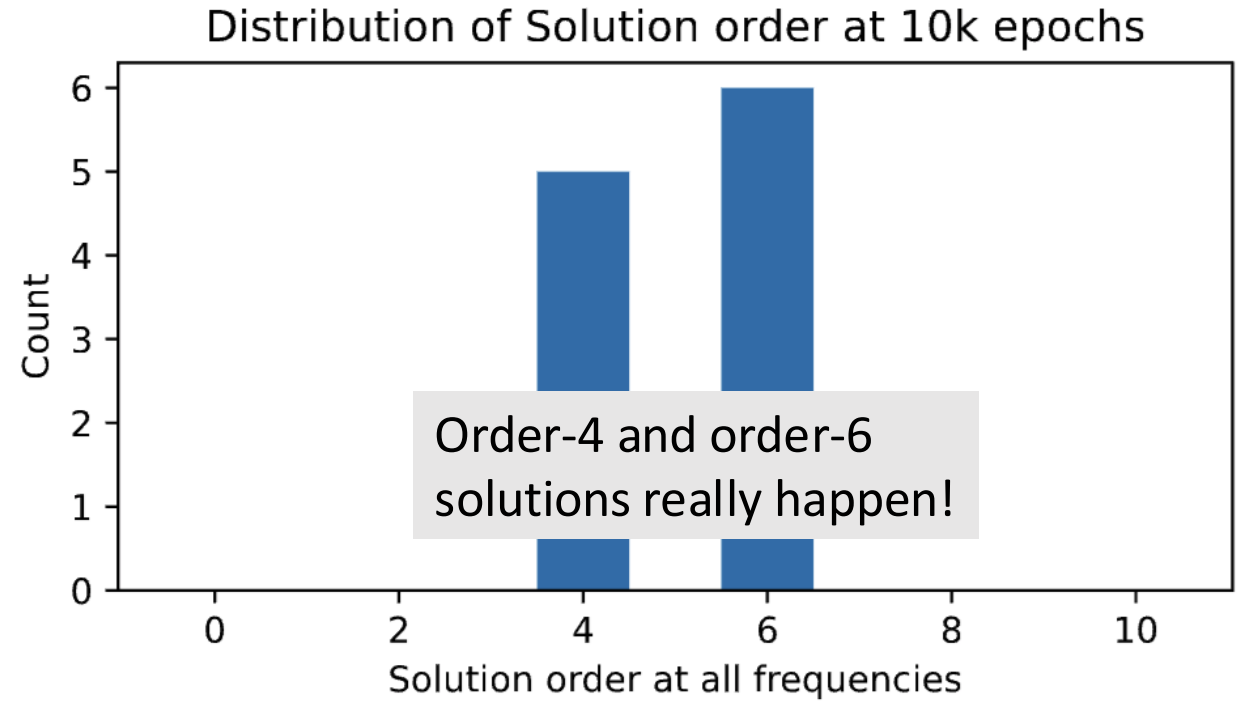
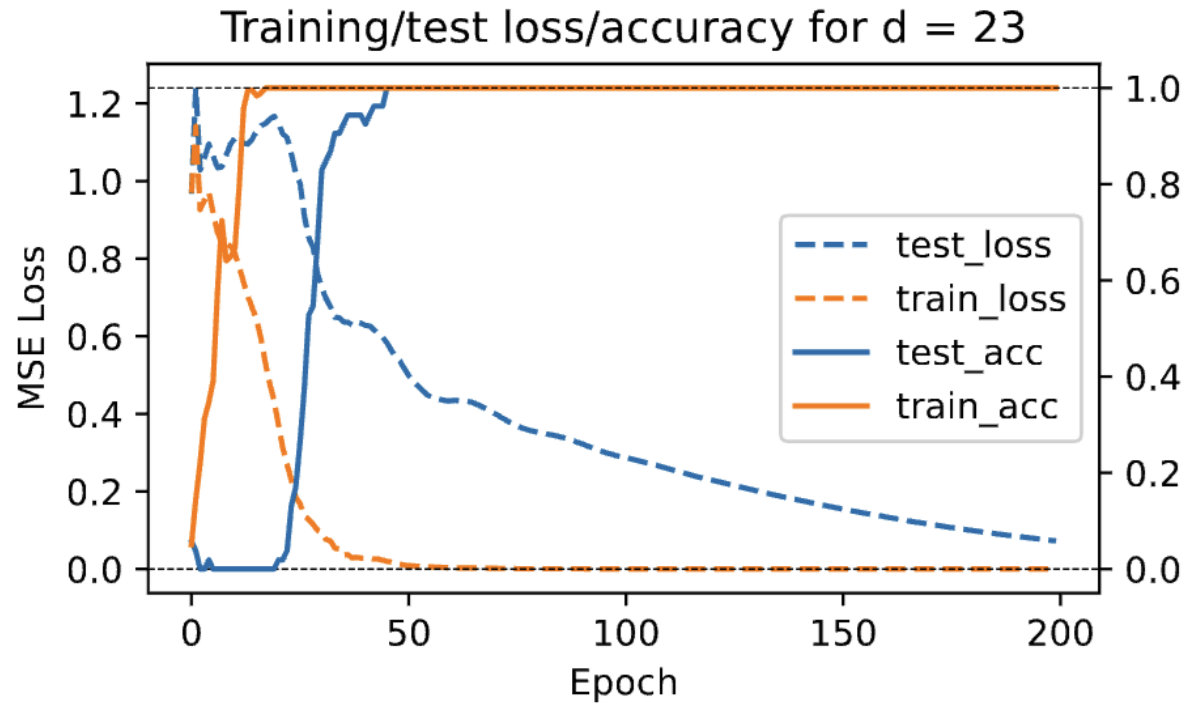
What a Gradient Descent Solution look like?



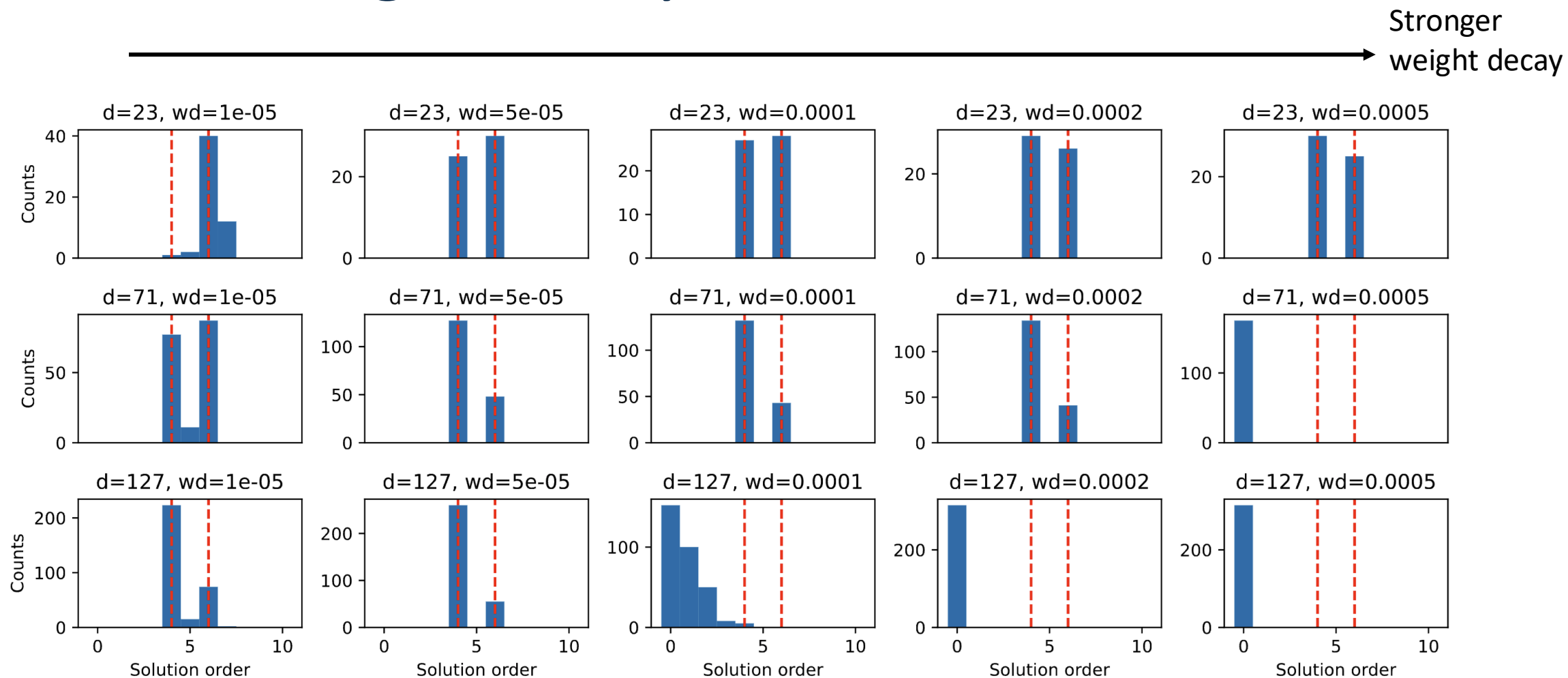
What a Gradient Descent Solution look like?



More Statistics on Gradient Descent Solutions



Effect of Weight Decay



Why? 🤔

How to Optimize?

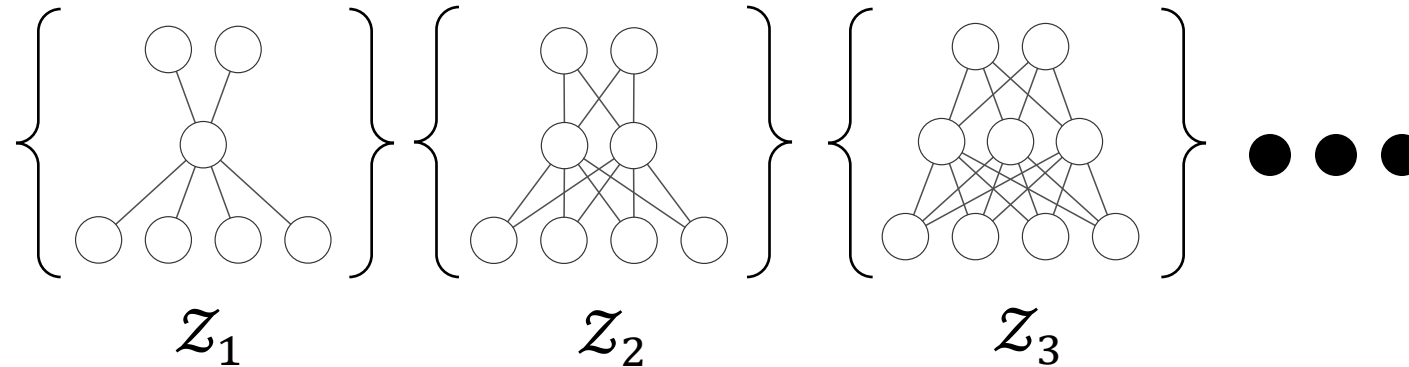
The objective is highly nonlinear !!

However, nice *algebraic structures* exist!

How to Optimize?

The objective is highly nonlinear !!

However, nice *algebraic structures* exist!

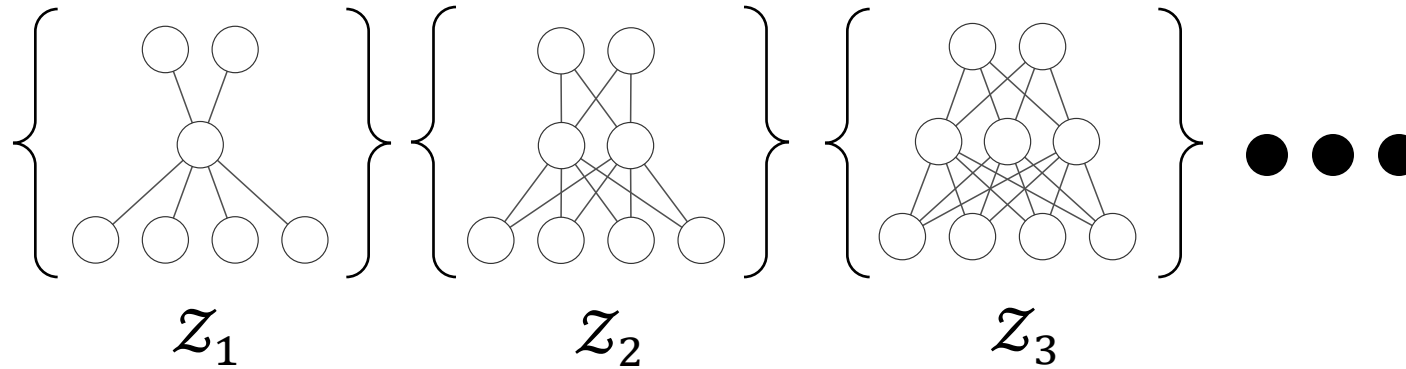


$\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$: All 2-layer networks with different number of hidden nodes

How to Optimize?

The objective is highly nonlinear !!

However, nice *algebraic structures* exist!



$\mathcal{Z} = \bigcup_{q \geq 0} \mathcal{Z}_q$: 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*

Composing Global Optimizers from Partial Ones

Partial solution #1

$$\mathbf{z}_{\text{syn}}^{(k)} \in R_c \cap R_n \text{ but } \mathbf{z}_{\text{syn}}^{(k)} \notin R_*$$

Partial solution #2

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

Composing Global Optimizers from Partial Ones

Composing
solutions using
ring multiplication *



Partial solution #1

$$\mathbf{z}_{\text{syn}}^{(k)} \in R_c \cap R_n \text{ but } \mathbf{z}_{\text{syn}}^{(k)} \notin R_*$$

Partial solution #2

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

Better solution

$$\mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_v^{(k)} \in R_c \cap R_n \cap R_*$$

Composing Global Optimizers from Partial Ones

Composing solutions using *ring multiplication* *

Composing solutions using *ring addition* +

Partial solution #1

$$\mathbf{z}_{\text{syn}}^{(k)} \in R_c \cap R_n \text{ but } \mathbf{z}_{\text{syn}}^{(k)} \notin R_*$$

Partial solution #2

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

Better solution

$$\mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_v^{(k)} \in R_c \cap R_n \cap R_*$$

Global Optimizer to MSE loss $\ell(\mathbf{z})$!

$$\mathbf{z}_{F6} = \frac{1}{\sqrt[3]{6}} \sum_k \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_v^{(k)}$$

Exemplar constructed global optimizers

Order-6 \mathbf{z}_{F6} (2*3)

$$\mathbf{z}_{F6} = \frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_{\nu}^{(k)} * \mathbf{y}_k$$

Exemplar constructed global optimizers

Order-6 \mathbf{z}_{F6} (2*3)

$$\mathbf{z}_{F6} = \frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_{\nu}^{(k)} * \mathbf{y}_k$$

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

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

Exemplar constructed global optimizers

Order-6 \mathbf{z}_{F6} (2*3)

$$\mathbf{z}_{F6} = \frac{1}{\sqrt[3]{6}} \sum_{k=1}^{(d-1)/2} \mathbf{z}_{\text{syn}}^{(k)} * \mathbf{z}_{\nu}^{(k)} * \mathbf{y}_k$$

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

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

Perfect memorization
(order-d per frequency)

$$\mathbf{z}_a = \sum_{j=0}^{d-1} \mathbf{u}_a^j, \quad \mathbf{z}_b = \sum_{j=0}^{d-1} \mathbf{u}_b^j$$

$$\mathbf{z}_M = d^{-2/3} \mathbf{z}_a * \mathbf{z}_b$$

Gradient Descent solutions matches with construction

d	%not	%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
	order-4/6	order-4	order-6	order-4	order-6	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\xi}^{(k)}$	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\text{syn},\alpha\beta}^{(k)}$	$\mathbf{z}_{\nu}^{(k)} * \mathbf{z}_{\text{syn}}^{(k)}$	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 \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 \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	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

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

Gradient Descent solutions matches with construction

d	%not		%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
	order-4/6	order-4	order-4	order-6	order-4	order-6	$z_{\nu=i}^{(k)} * z_{\xi}^{(k)}$	$z_{\nu=i}^{(k)} * z_{\text{syn},\alpha\beta}^{(k)}$	$z_{\nu}^{(k)} * z_{\text{syn}}^{(k)}$	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 \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 \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	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	

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

Gradient Descent solutions matches with construction

d	%not order-4/6	%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
		order-4	order-6	order-4	order-6	$z_{\nu=i}^{(k)} * z_{\xi}^{(k)}$	$z_{\nu=i}^{(k)} * z_{\text{syn},\alpha\beta}^{(k)}$	$z_{\nu}^{(k)} * z_{\text{syn}}^{(k)}$	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 \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 \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	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

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

Gradient Descent solutions matches with construction

d	%not order-4/6	%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
		order-4	order-6	order-4	order-6	$z_{\nu=i}^{(k)} * z_{\xi}^{(k)}$	$z_{\nu=i}^{(k)} * z_{\text{syn},\alpha\beta}^{(k)}$	$z_{\nu}^{(k)} * z_{\text{syn}}^{(k)}$	others
23	0.0±0.0	0.00±0.00	5.71±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±0.00	0.03±0.00	5.02±0.25	72.57±0.70	4.00±1.14	21.14±2.14	2.29±1.07
127	0.0±0.0	1.50±0.92	0.00±0.00	0.26±0.14	0.93±0.18	82.96±0.39	2.25±0.64	14.13±0.87	0.66±0.66

Factorization error is very small

Gradient Descent solutions matches with construction

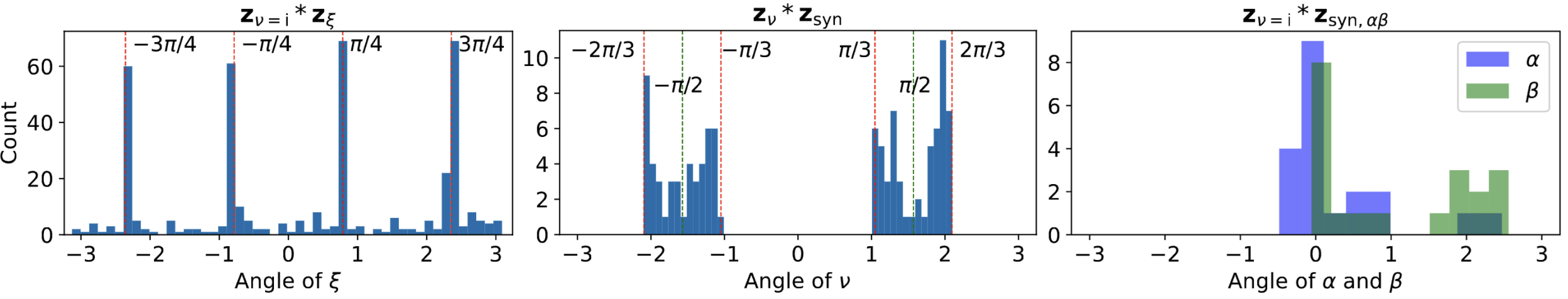
d	%not order-4/6	%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
		order-4	order-6	order-4	order-6	$z_{\nu=i}^{(k)} * z_{\xi}^{(k)}$	$z_{\nu=i}^{(k)} * z_{\text{syn},\alpha\beta}^{(k)}$	$z_{\nu}^{(k)} * z_{\text{syn}}^{(k)}$	others
23	0.0 \pm 0.0	0.00 \pm 0.00	5.71 \pm 5.71	0.05 \pm 0.01	4.80 \pm 0.96	47.07 \pm 1.88	11.31 \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 \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	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

Gradient Descent solutions matches with construction

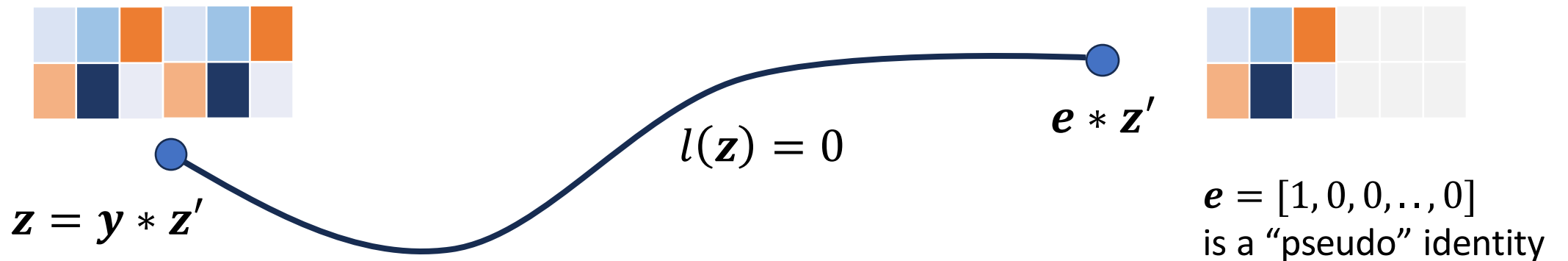
d	%not		%non-factorable		error ($\times 10^{-2}$)		solution distribution (%) in factorable ones			
	order-4/6	order-4	order-4	order-6	order-4	order-6	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\xi}^{(k)}$	$\mathbf{z}_{\nu=i}^{(k)} * \mathbf{z}_{\text{syn},\alpha\beta}^{(k)}$	$\mathbf{z}_{\nu}^{(k)} * \mathbf{z}_{\text{syn}}^{(k)}$	others
23	0.0 ± 0.0	0.00 ± 0.00	5.71 ± 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
						5	72.57 ± 0.70	4.00 ± 1.14	21.14 ± 2.14	2.29 ± 1.07
						8	82.96 ± 0.39	2.25 ± 0.64	14.13 ± 0.87	0.66 ± 0.66

Distribution of the parameters in the solutions



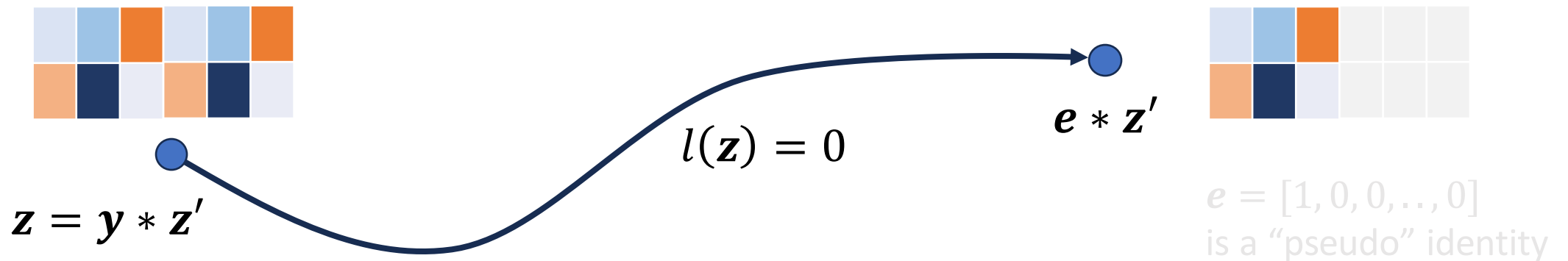
Gradient Dynamics

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



Gradient Dynamics

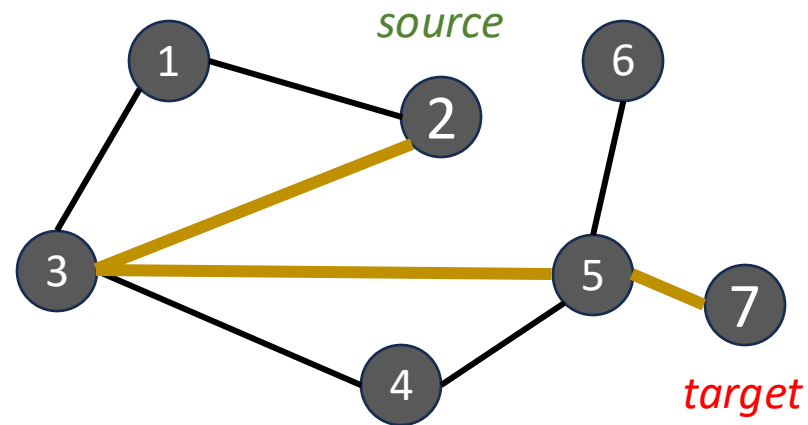
Theorem [**The Occam's Razor**] If $\mathbf{z} = \mathbf{y} * \mathbf{z}'$ and both \mathbf{z} and \mathbf{z}' are global optimal, then there exists a path of zero loss connecting \mathbf{z} and \mathbf{z}' .



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

Another Example: Symbolic from Neural Representation

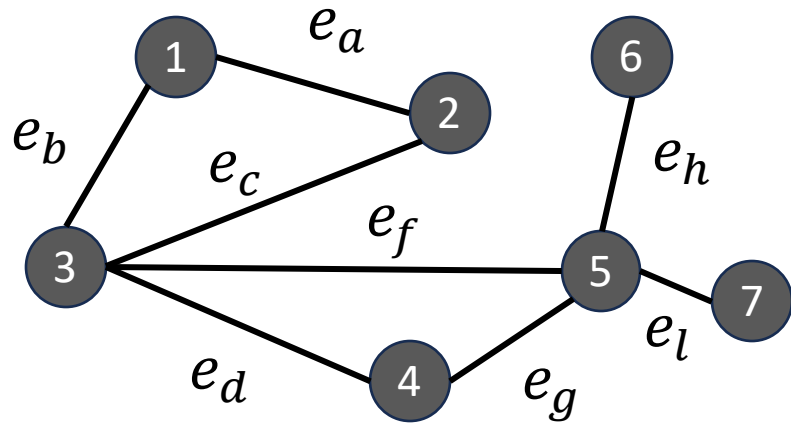
Task: Learn a 2-layer Transformer for predicting **shortest path** in the graph



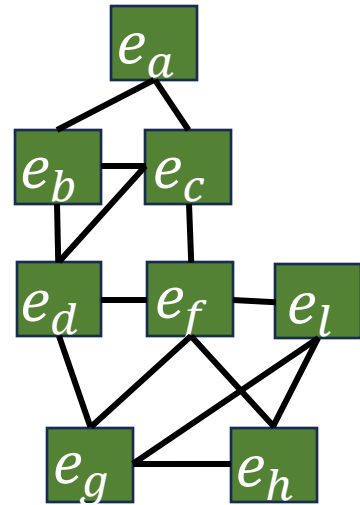
`<bos> 1 2 <e> ... <q> [source] [target] <p> [source] [node 1] [node 2] ... [target]`

Context Predicted Shortest path

What representations it learns?

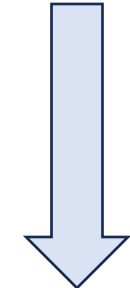


Line graph



Normalized Graph Laplacian

$$L = I - D^{-1/2} A D^{-1/2}$$



Edge Embedding

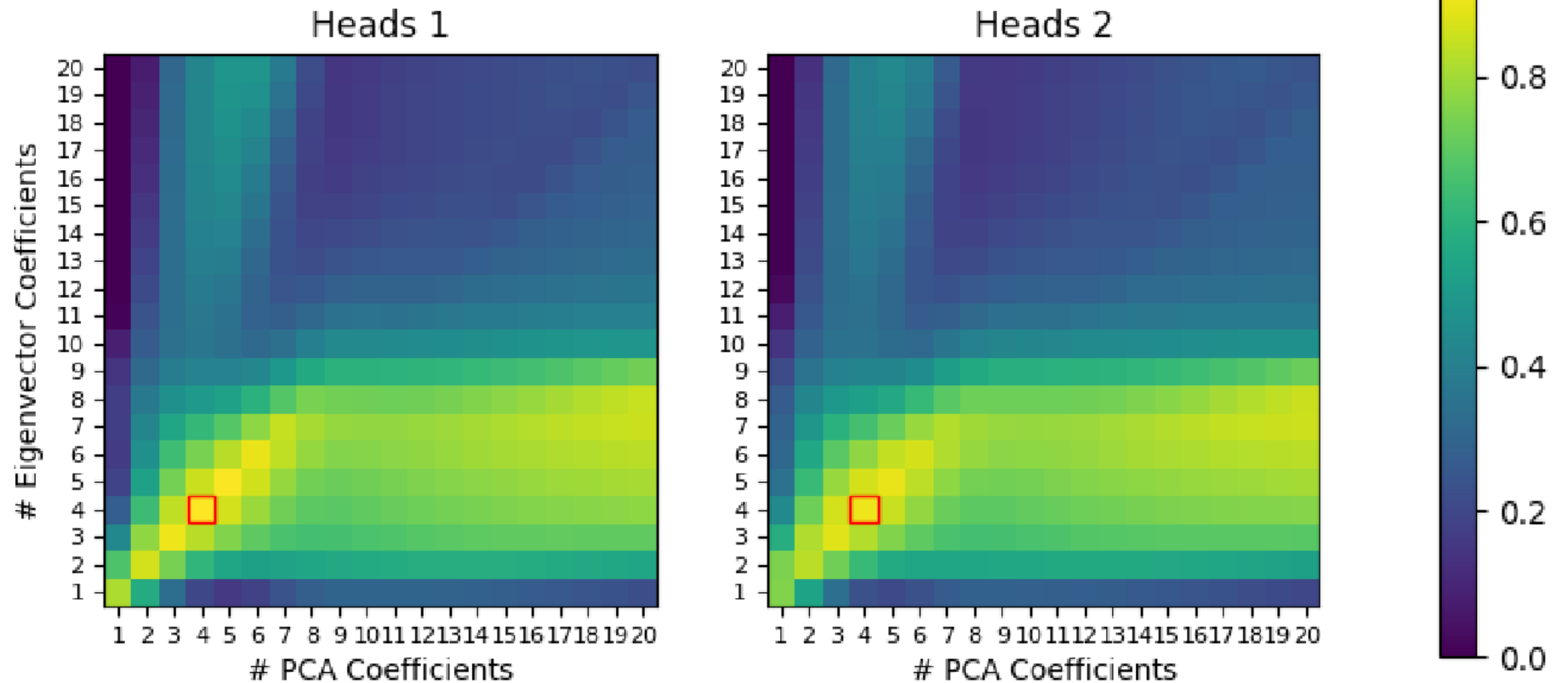
Representation after the first Transformer layer (averaged over random edge order)



$\langle \text{bos} \rangle$ 1 2 $\langle \mathbf{e} \rangle$... $\langle \mathbf{q} \rangle$ [source] [target] $\langle \mathbf{p} \rangle$ [source] [node 1] [node 2] ... [target]

What representations it learns?

Graph Edge Embedding of various dimensions



Computed edge embedding with trained Transformers

Normalized Correlation > 0.9

Spectral Line Navigator (SLN)

Simple Algorithms of Graph Shortest Path

1. Compute Line Graph \tilde{G} of existing graph G
2. Compute eigenvectors of normalized Laplacian $L(\tilde{G})$
3. $i = source$
4. While $i \neq target$ do
 $distance(j, k; i) := \|v_{ij} - v_{k,target}\|_2$
Find $j = \operatorname{argmin}_{j,k} distance(j, k; i)$
Let $i = j$

>99% optimal for small
random graph (size < 10)

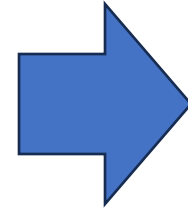
o3-mini-high implementation: <https://chatgpt.com/share/67b027f9-fb28-8012-aa64-a1f7479134b7>

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!

Thanks!