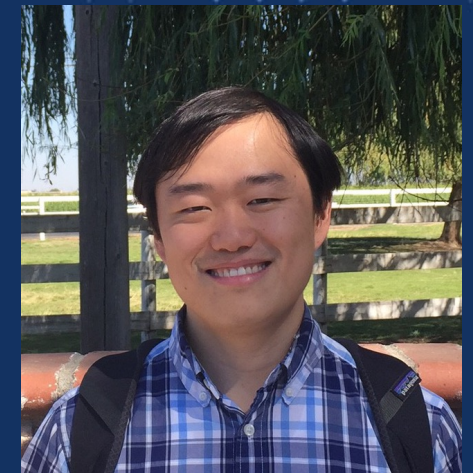


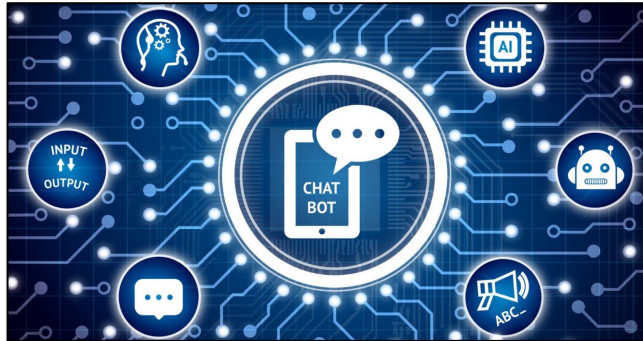
Demystifying Self-Attention Mechanism in Feature Composition and Logic Reasoning via the Lens of Training Dynamics

Yuandong Tian
Research Scientist

Meta AI (FAIR)



Large Language Models (LLMs)



Conversational AI



Content Generation



AI Agents

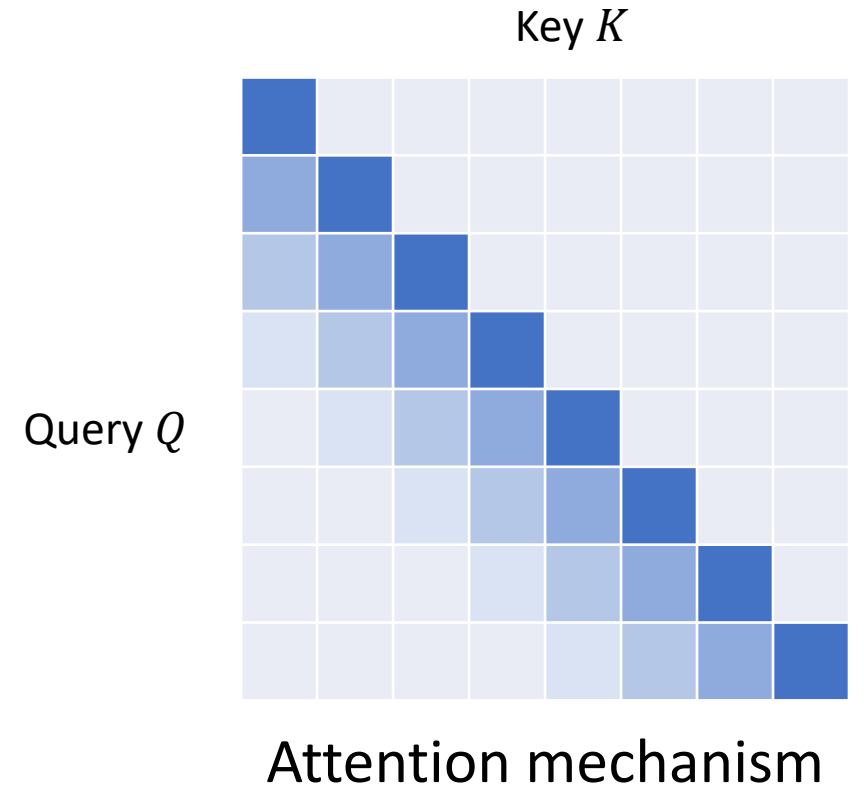
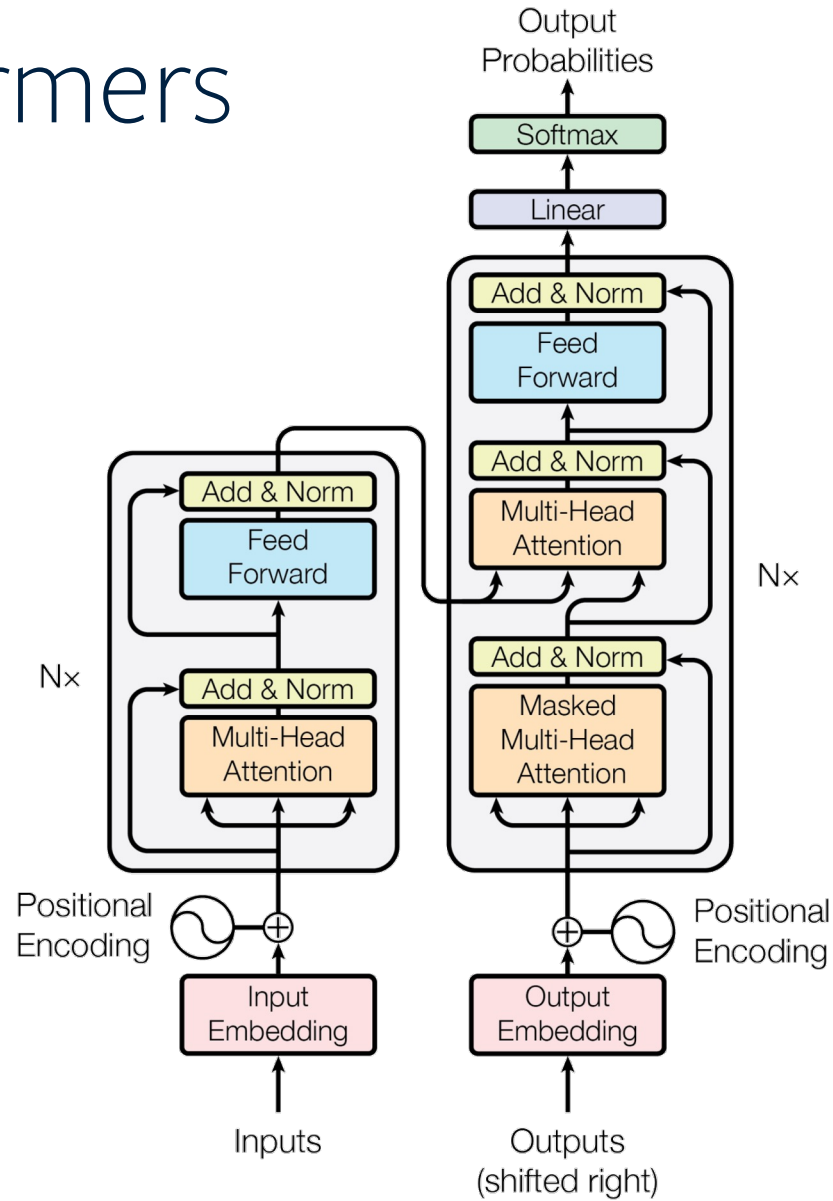
Standard Prompting	Chain of Thought Prompting
<p>Input</p> <p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: The answer is 11.</p> <p>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 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

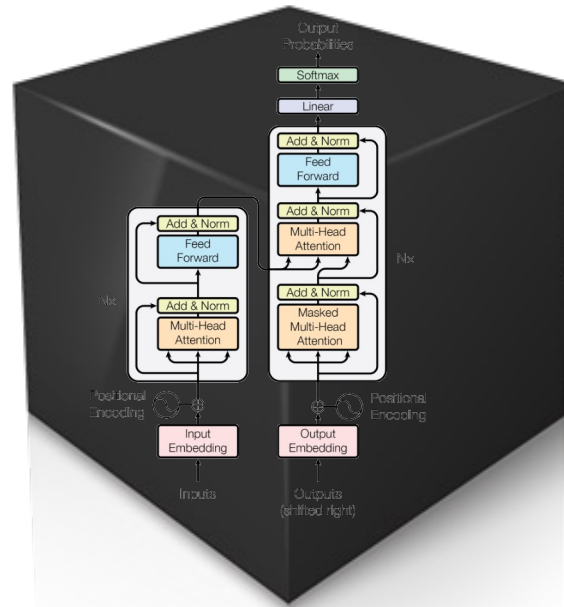


How does Transformer work?

Input



This is an apple

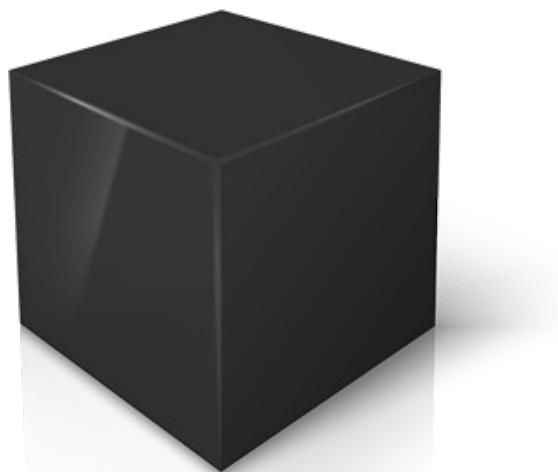


“Some Nonlinear Transformation”

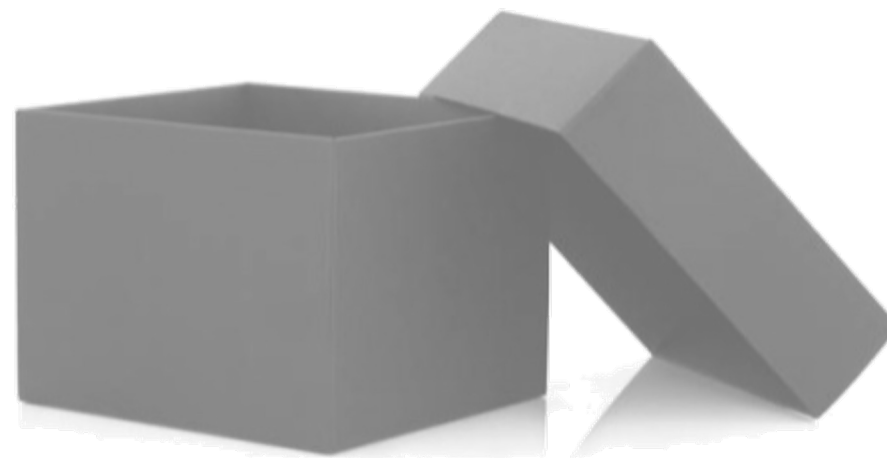


Output

Black-box versus White-box



Black box

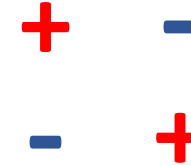


White box

Three Angles

Understanding how
Deep Models work

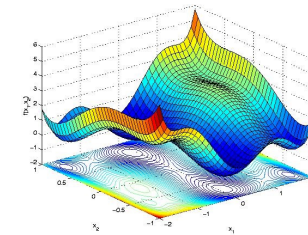
Expressibility



“Neural Network is a universal approximator”

“Deep Models can express functions more efficiently than shallow ones”

Optimization

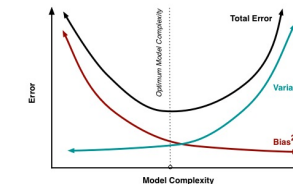


“Gradient vanishing/exploding”

“Gradient Descent might get stuck at saddle point / local minima”

“Can GD/SGD go to global optima? How fast?”

Generalization



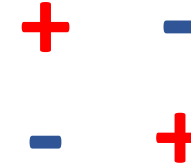
“Does zero training error often lead to overfitting?”

“More parameters might lead to overfitting.”

Three Angles

Understanding how
Deep Models work

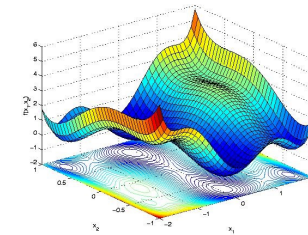
Expressibility



“Neural Network is a universal approximator”

“Deep Models can express functions more efficiently than shallow ones”

Optimization

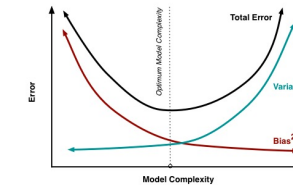


“Gradient vanishing/exploding”

“Gradient Descent might get stuck at saddle point / local minima”

“Can GD/SGD go to global optima? How fast?”

Generalization

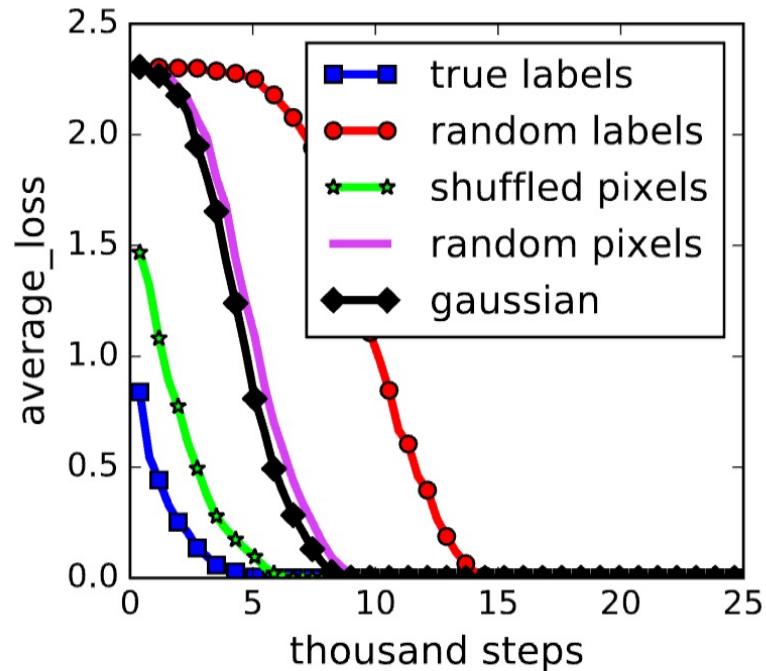


“Does zero training error often lead to overfitting?”

“More parameters might lead to overfitting.”

Which path should we take?

Rethinking Generalization



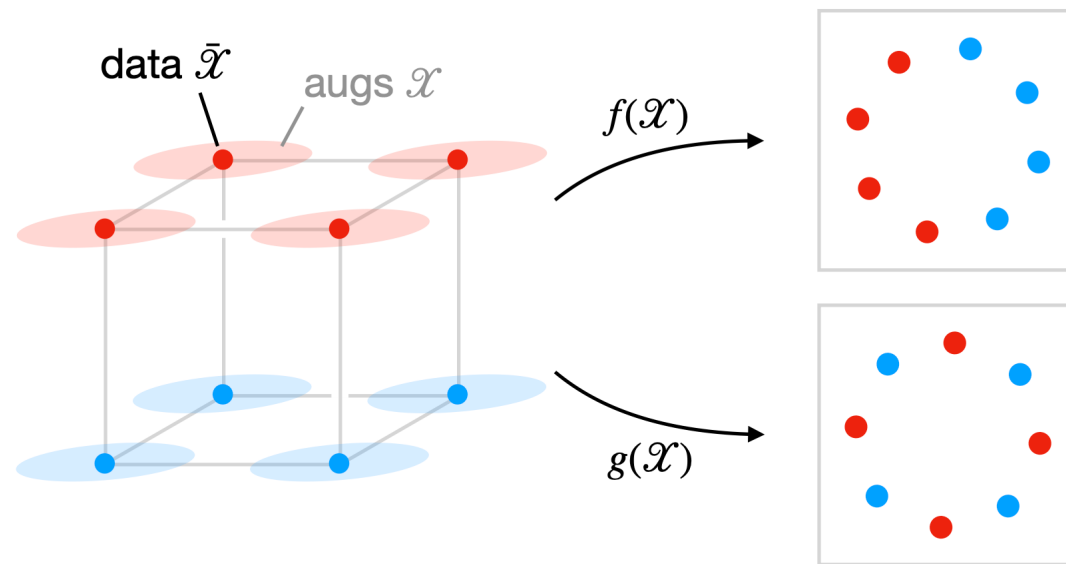
(a) learning curves

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
(fitting random labels)		no	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
(fitting random labels)		no	no	100.0	10.12
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
(fitting random labels)		no	no	99.82	9.86
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
(fitting random labels)		no	no	100.0	10.48
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
(fitting random labels)		no	no	99.34	10.61

Generalization bound failed: $Test\ Error \leq Train\ Error + ???$

Inductive Bias Really Matters

A self-supervised contrastive learning example



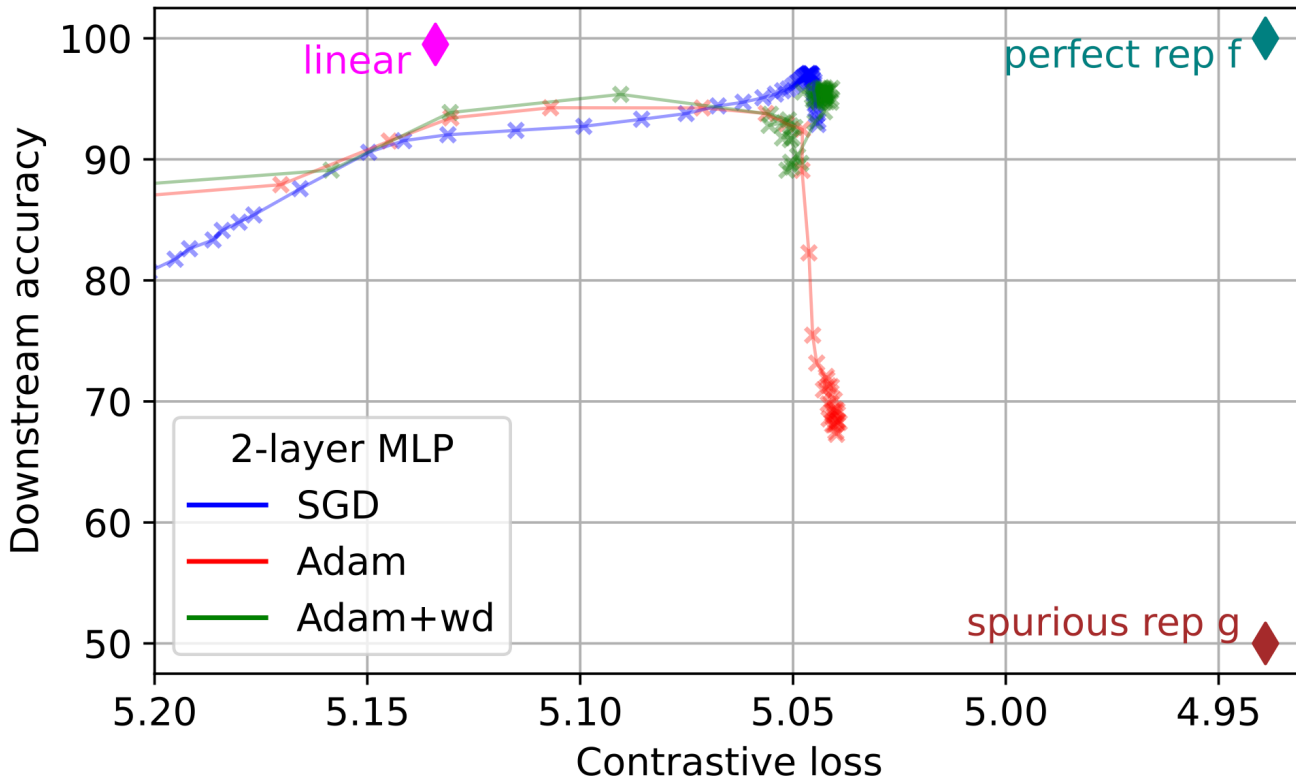
Pretraining: $L_{\text{cont}}(g) \approx L_{\text{cont}}(f)$

Downstream: $L_{\text{clf}}(g) \gg L_{\text{clf}}(f)$

SSL Pretraining loss doesn't really reflect downstream loss

Inductive Bias Really Matters

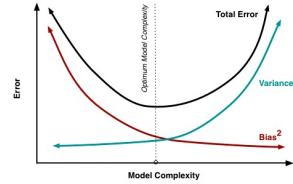
Boolean hypercube example



Representation	Contrastive loss	Accuracy (%)
$\exists f$ (perfect)	4.939	100
$\exists g$ (spurious)	4.939	50
MLP + Adam	5.039 ± 0.001	74.1 ± 4.3
MLP + Adam + wd	5.040 ± 0.002	89.5 ± 4.9
Linear	5.134 ± 0.002	99.5 ± 0.1

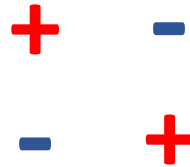
Lesson learned?

Generalization



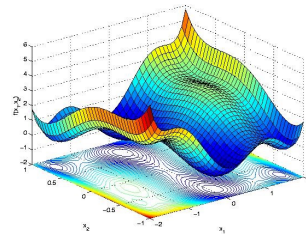
Architecture **X**
training dynamics **X**

Expressibility



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

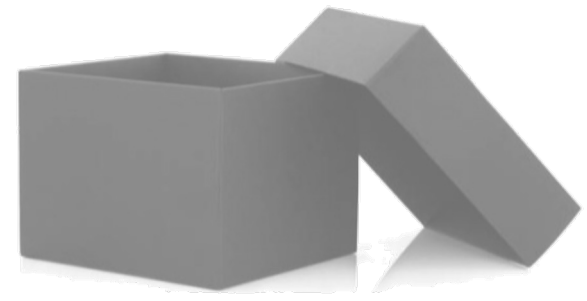
Optimization



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

How about

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



Start From the First Principle

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

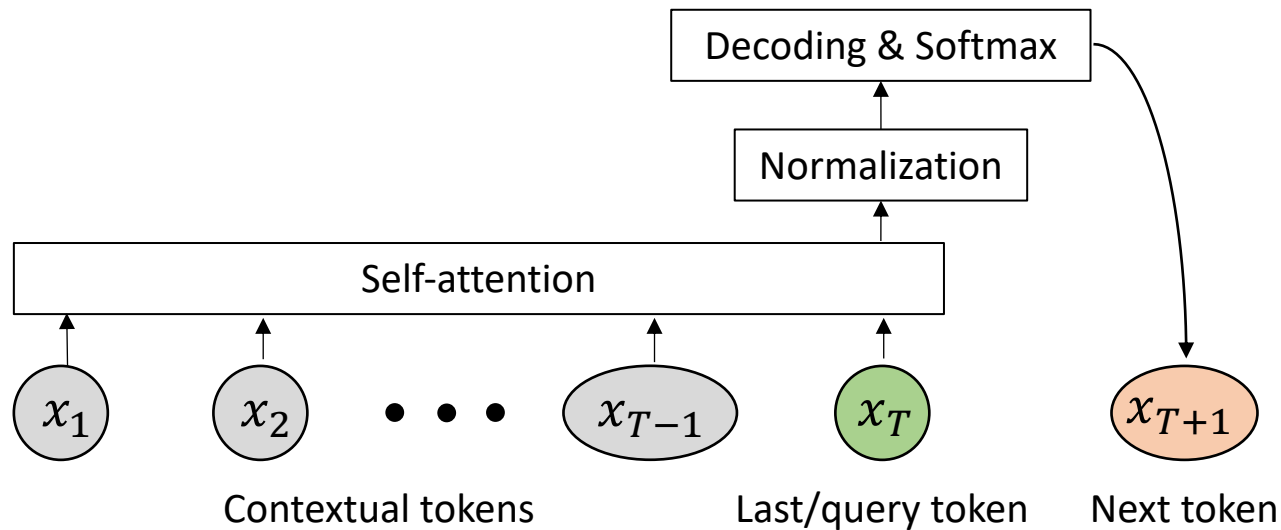
$$\dot{\mathbf{w}} := \frac{d\mathbf{w}}{dt} = -\nabla_{\mathbf{w}}J(\mathbf{w})$$

- **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 ✓

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

Major Assumptions

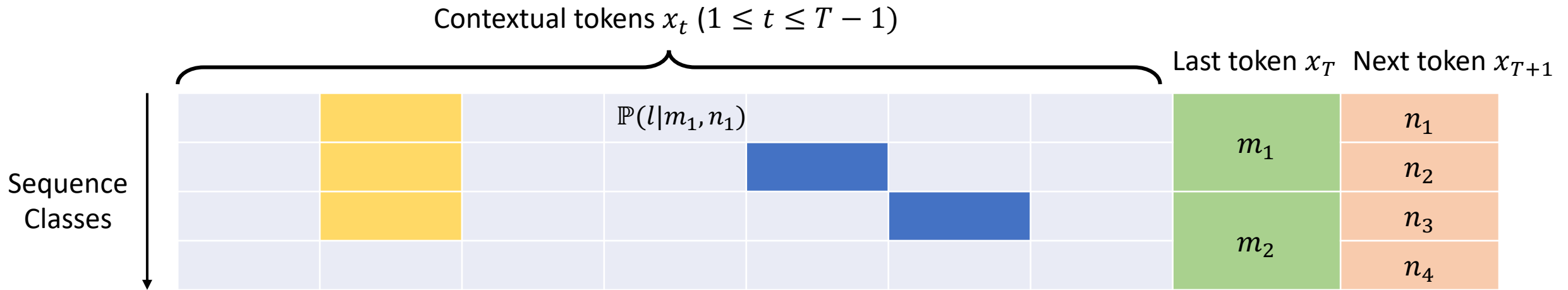
- No positional encoding
- Sequence length $T \rightarrow +\infty$
- Learning rate of decoder Y larger than self-attention layer Z ($\eta_Y \gg \eta_Z$)
- Other technical assumptions

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 $\mathbb{P}(l|n) > 0$

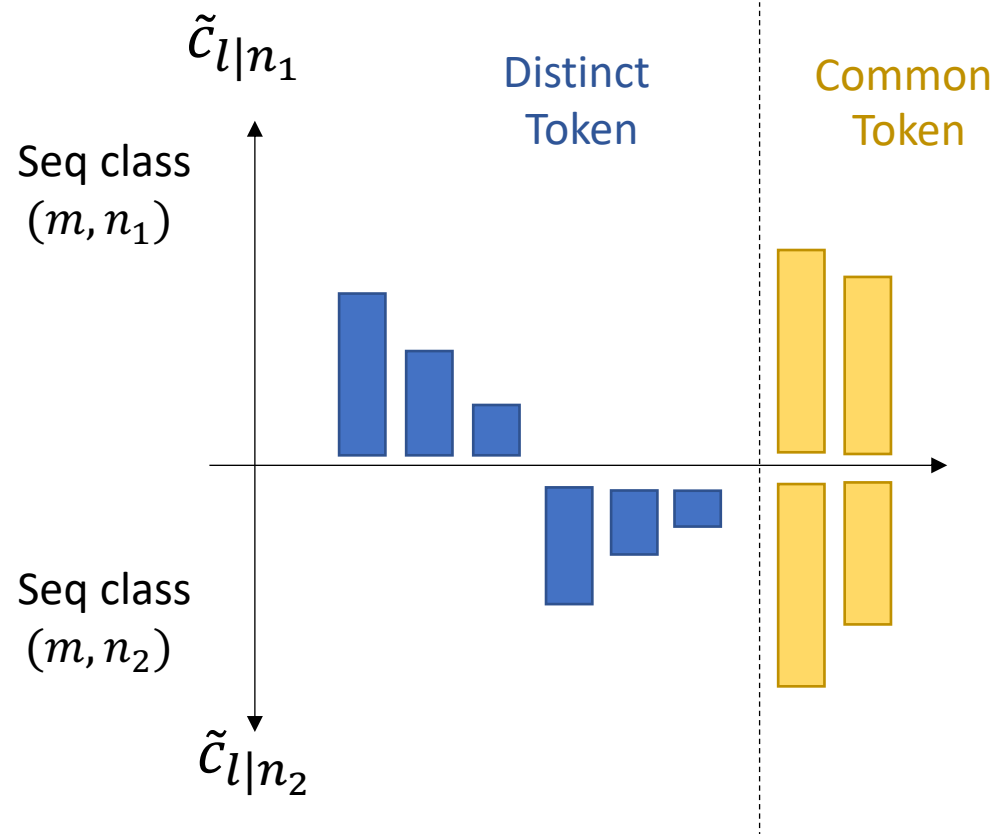
$\mathbb{P}(l|m, n) = \mathbb{P}(l|n)$ is the conditional probability of token l given last token $x_T = m$ and $x_{T+1} = n$

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

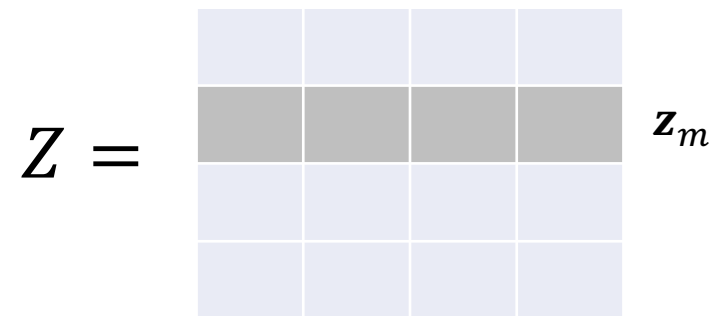
At initialization



Co-occurrence probability

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

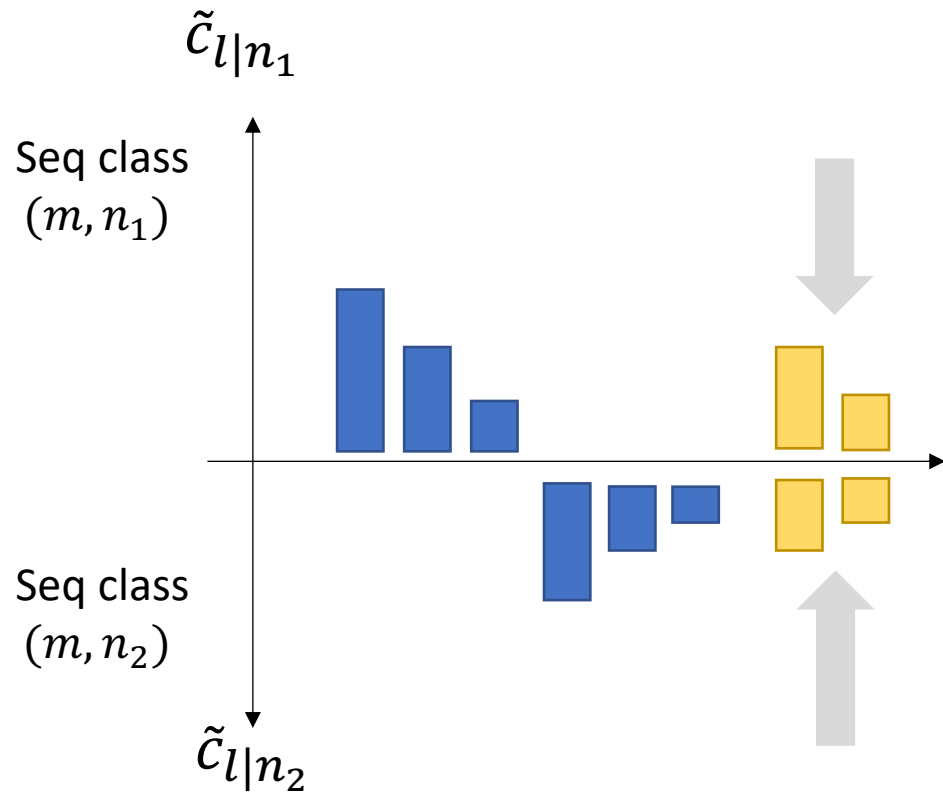
Initial condition: $z_{ml}(0) = 0$



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

Overall Picture of the Training Dynamics

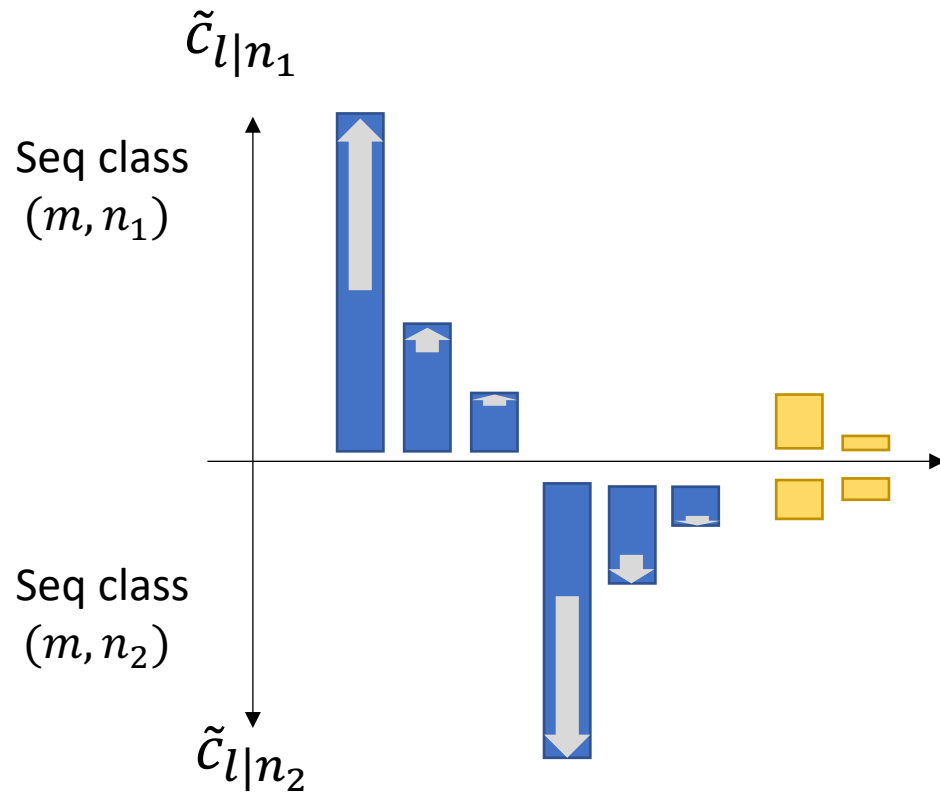
Common Token Suppression



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

Overall Picture of the Training Dynamics

Winners-emergence



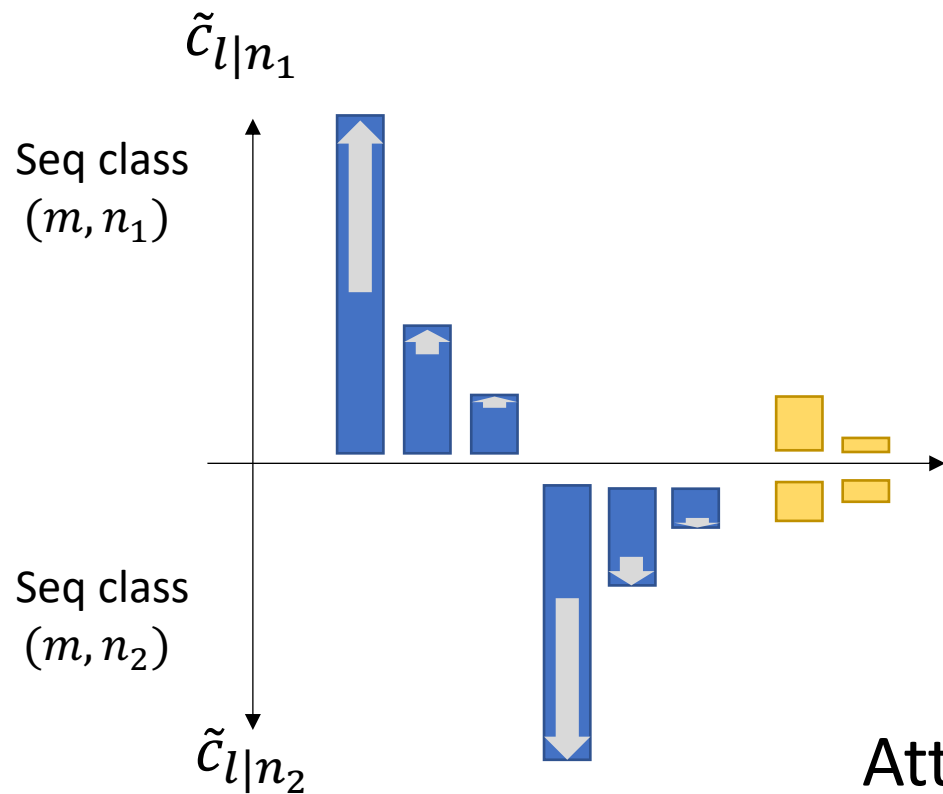
(a) $z_{ml} < 0$, for **common token** l

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

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

Overall Picture of the Training Dynamics

Winners-emergence



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

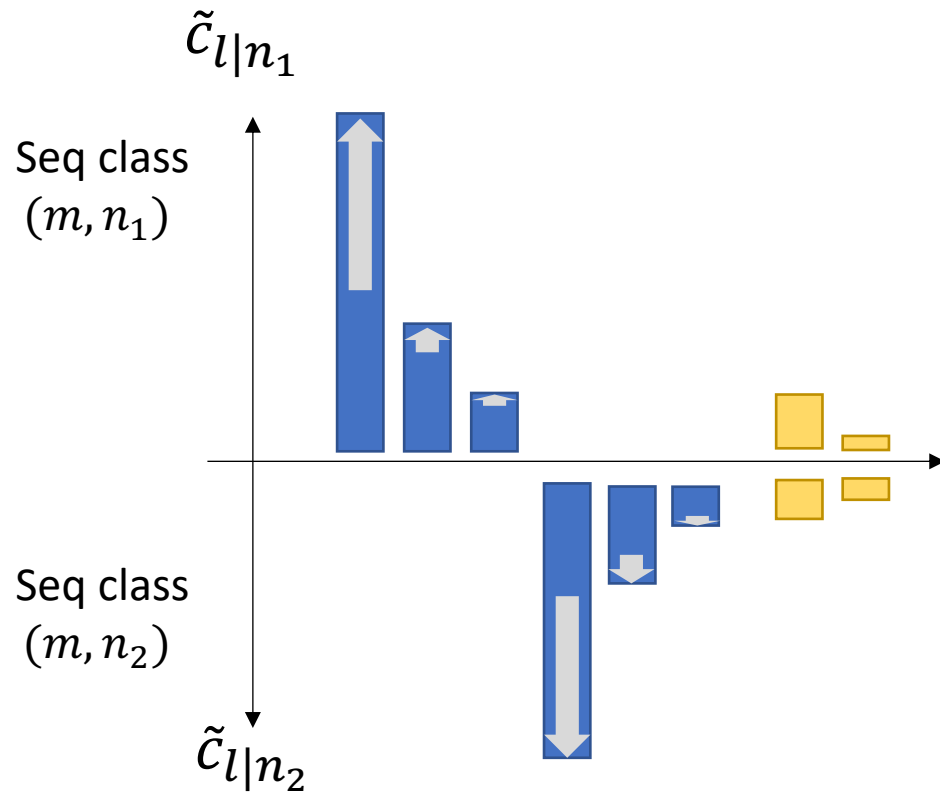
(b) $z_{ml}^{\cdot} > 0$, for **distinct token** l

(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:

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

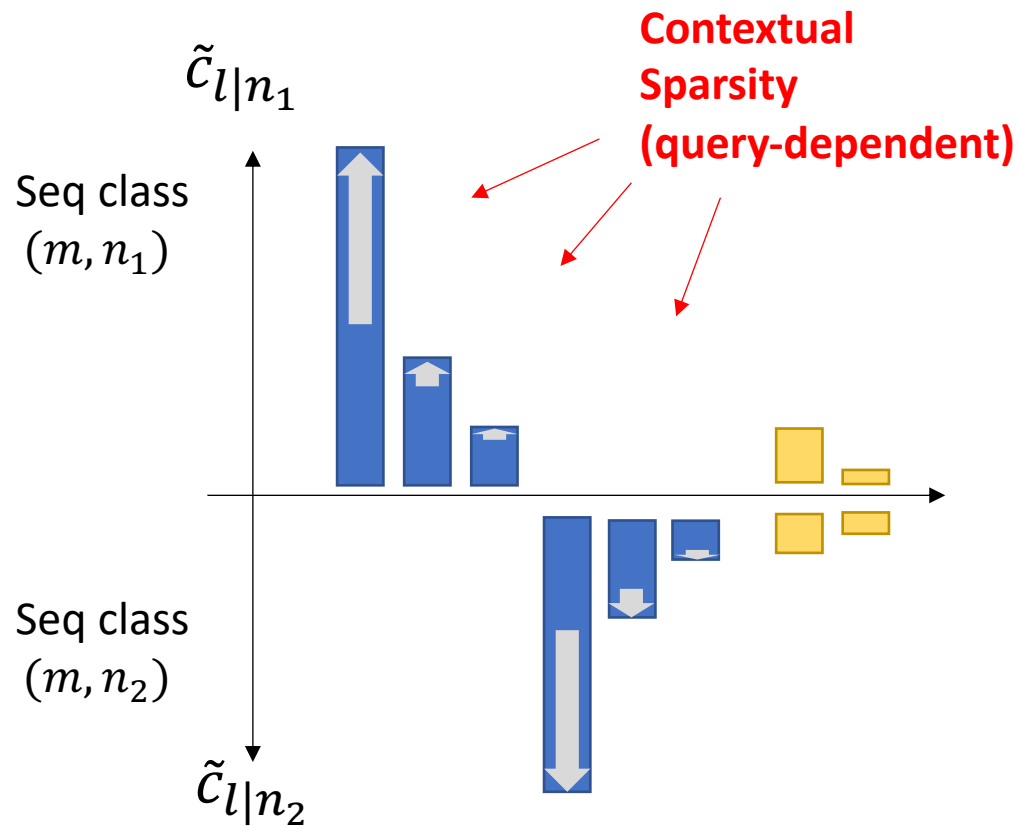
If l_0 is the dominant token: $r_{l_0/l|n}(0) > 0$ for all $l \neq l_0$ then

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

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:

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

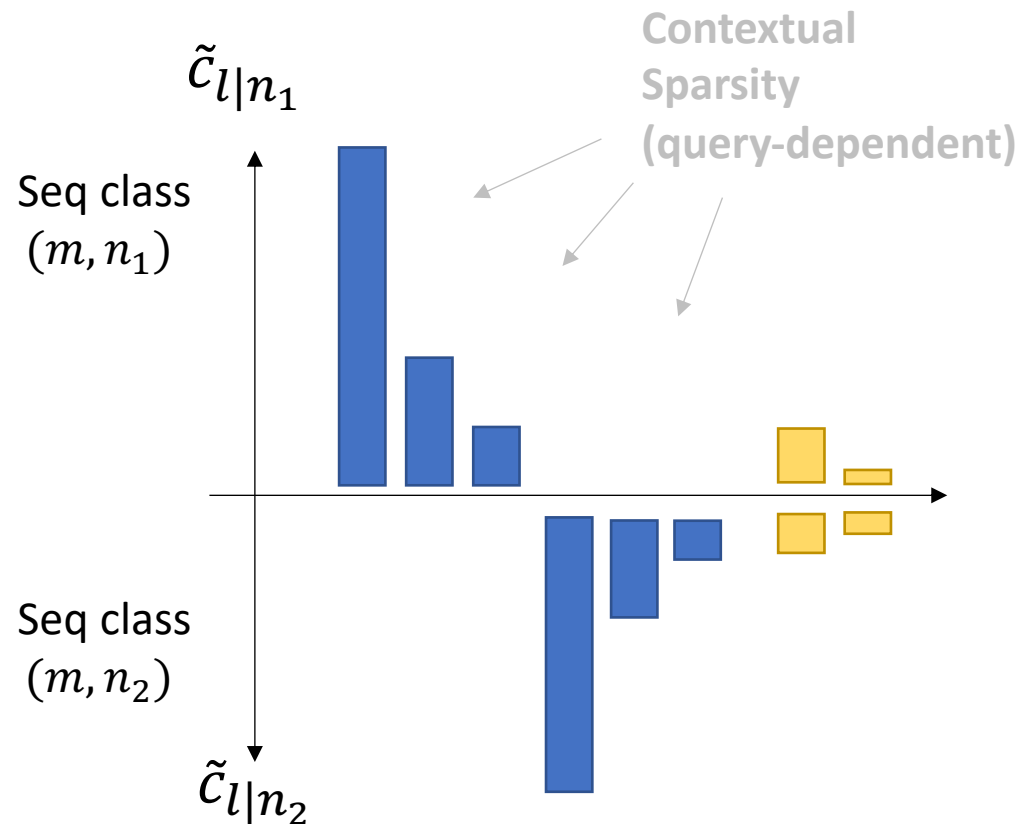
If l_0 is the dominant token: $r_{l_0/l|n}(0) > 0$ for all $l \neq l_0$ then

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

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:

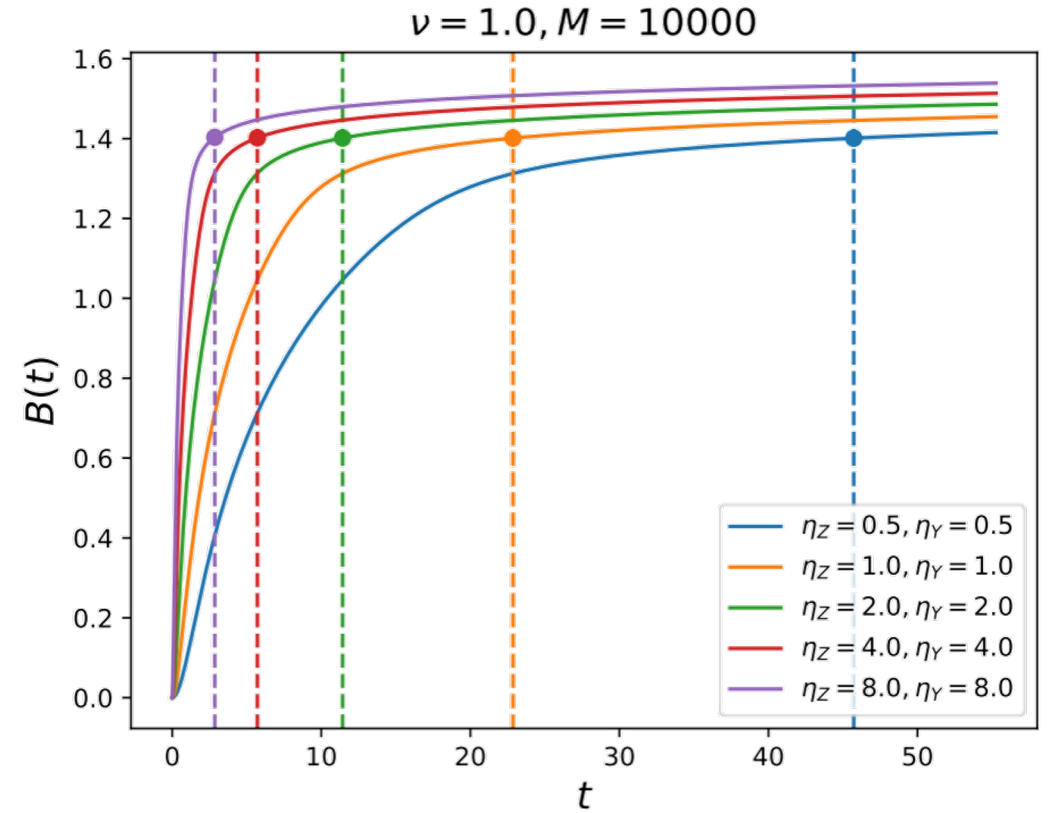
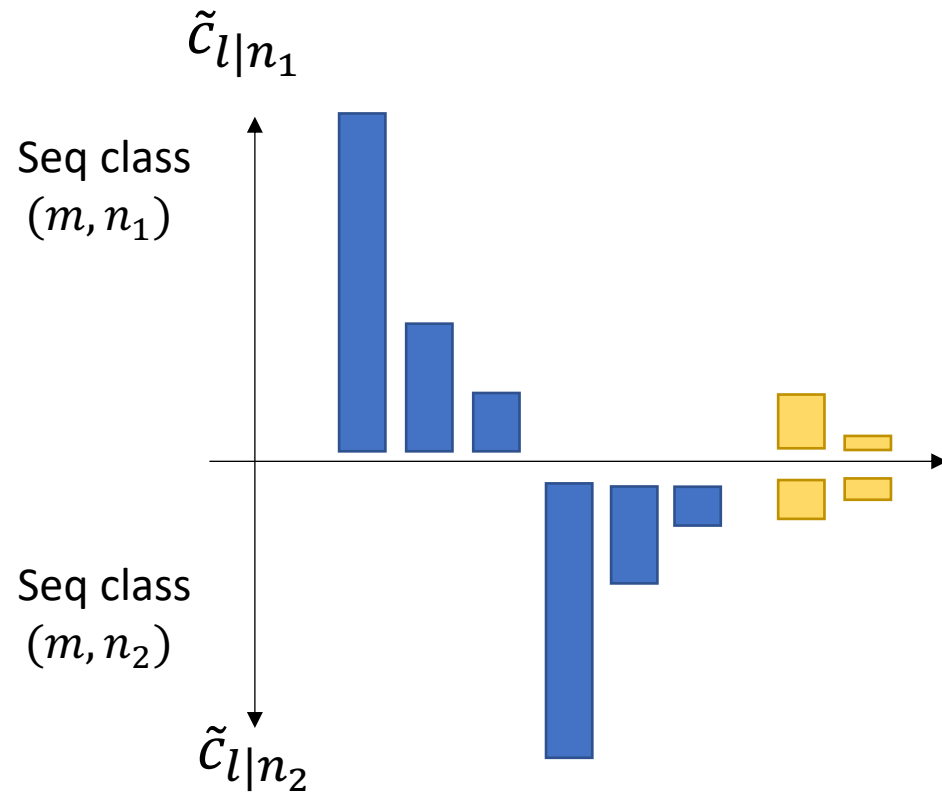
When $t \geq t_0 = O\left(\frac{2K \ln M}{\eta_Y}\right)$, $B_n(t) = O(\ln \ln t)$

(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

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

Simple Real-world Experiments

WikiText2
(original parameterization)

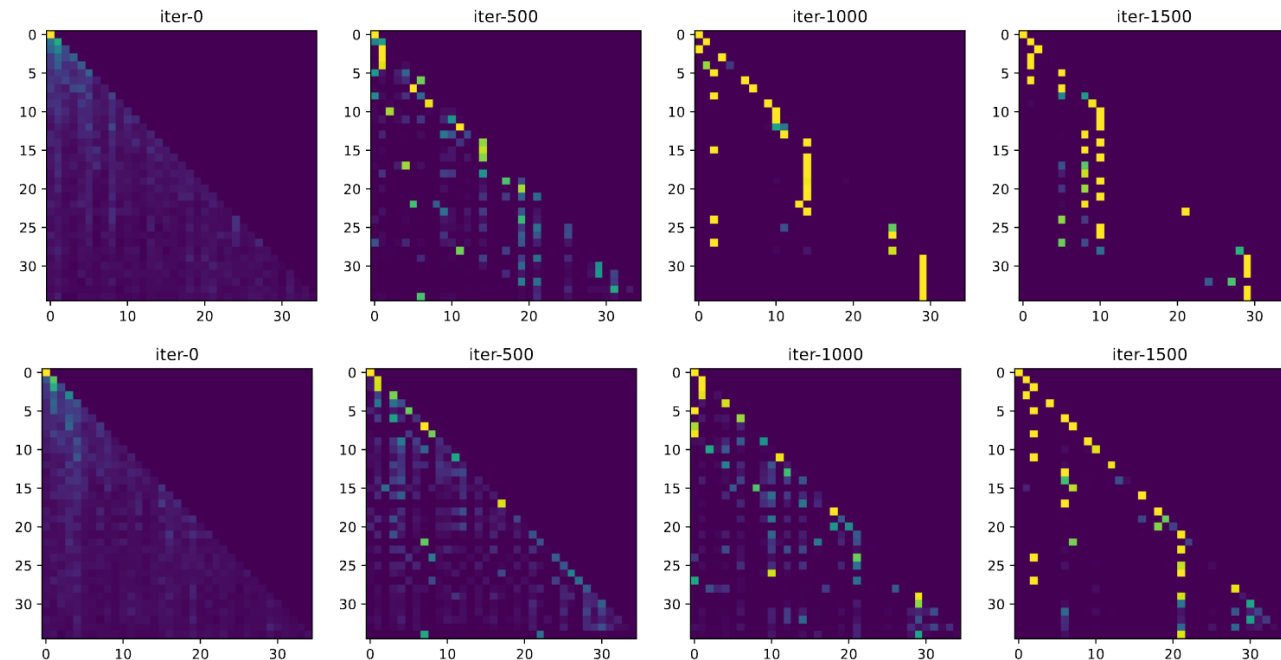


Figure 7: Attention patterns in the lowest self-attention layer for 1-layer (top) and 3-layer (bottom) Transformer trained on WikiText2 using SGD (learning rate is 5). Attention becomes sparse over training.

Further study of sparse attention
→ Deja Vu, H2O and StreamingLLM

[Z. Liu et al, *Deja vu: Contextual sparsity for efficient LLMs at inference time*, ICML'23 (oral)]

[Z. Zhang et al, *H2O: Heavy-Hitter Oracle for Efficient Generative Inference of Large Language Models*, NeurIPS'23]

[G. Xiao et al, *Efficient Streaming Language Models with Attention Sinks*, ICLR'24]

Deal with Reversal Curse

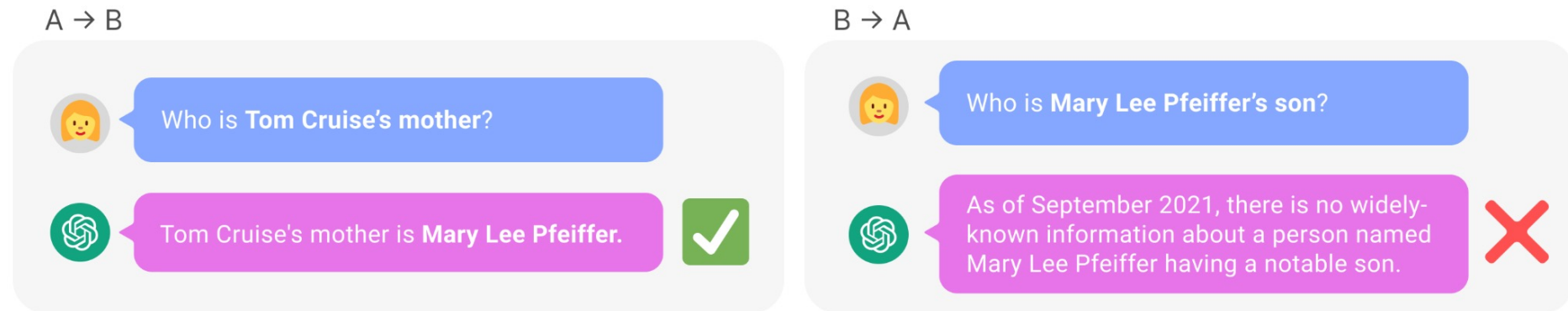
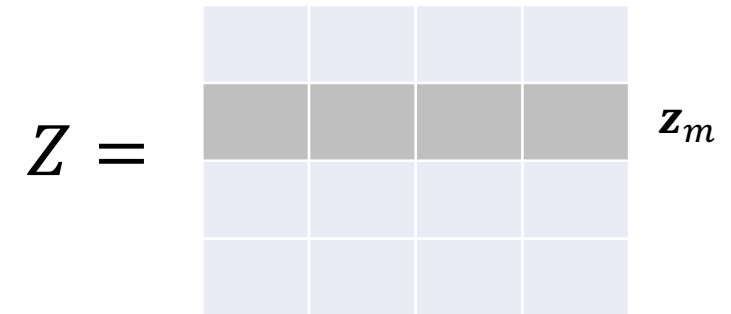


Figure 1: Inconsistent knowledge in GPT-4. GPT-4 correctly gives the name of Tom Cruise’s mother (left). Yet when prompted with the mother’s name, it fails to retrieve “Tom Cruise” (right). We hypothesize this ordering effect is due to the Reversal Curse. Models trained on “A is B” (e.g. “Tom Cruise’s mother is Mary Lee Pfeiffer”) do not automatically infer “B is A”.

How to explain “Reversal Curse”?

$Z = UW_QW_K^T U^T$ pairwise logits of self-attention matrix,
is **not** symmetric



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

You only learn what you see in the training set

Theorem 3 (Reversal curse). Assume we run SGD with batch size 1, and assume $M \gg 100$ and $\frac{1}{M^{0.99}} \ll \eta_Y < 1$. Let $t \gtrsim \frac{N \ln M}{\eta_Y}$ denote the time step which also satisfies $\ln t \gtrsim \ln(NM/\eta_Y)$. For training sequence $(x_1, x_2, x_3) \in \mathcal{D}_{\text{train}}$ at time t , we have

$$p_{\theta(t)}(x_3|x_1, x_2) \geq 1 - \frac{M-1}{2 \left(\frac{M\eta_Y t}{N}\right)^c} \xrightarrow{t \rightarrow \infty} 1$$

for some constant $c > 0$, and for any test sequence $(x_1, x_2, x_3) \in \mathcal{D}_{\text{test}}$ that is not included the training set $\mathcal{D}_{\text{train}}$, we have

$$p_{\theta(t)}(x_3|x_1, x_2) \leq \frac{1}{M}.$$

“Chain-of-thoughts” reasoning

Theorem 4 (Necessity of chain-of-thought). Assume we run SGD with batch size 1, and assume $M \gg 100$ and $\frac{1}{M^{0.99}} \ll \eta_Y < 1$. Let $t \gtrsim \frac{N \ln M}{\eta_Y}$ denote the time step which also satisfies $\ln t \gtrsim \ln(NM/\eta_Y)$. For any test index $i \in \mathcal{I}_{test}$, we have

$$p_{\theta(t)}(B_i | A_i \rightarrow) \geq 1 - \frac{M - 1}{2 \left(\frac{M \eta_Y t}{N} \right)^c}, \quad p_{\theta(t)}(C_i | B_i \rightarrow) \geq 1 - \frac{M - 1}{2 \left(\frac{M \eta_Y t}{N} \right)^c}$$

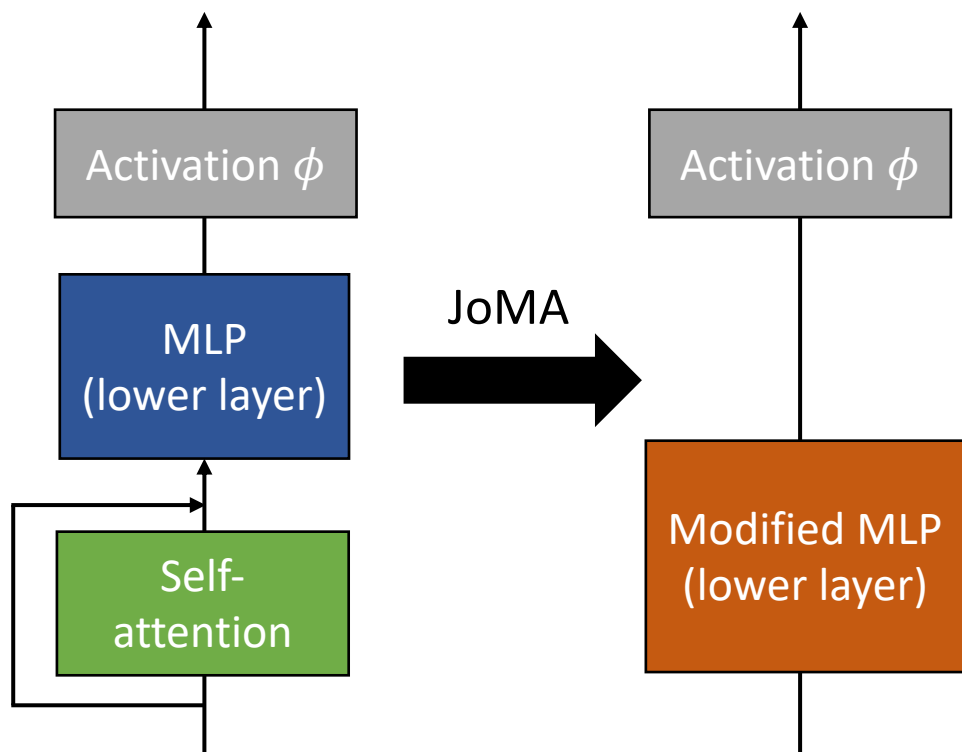
for some constant $c > 0$ and

$$p_{\theta(t)}(C_i | A_i \rightsquigarrow) \leq \frac{1}{M}.$$

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?
- New research work: **JoMA**

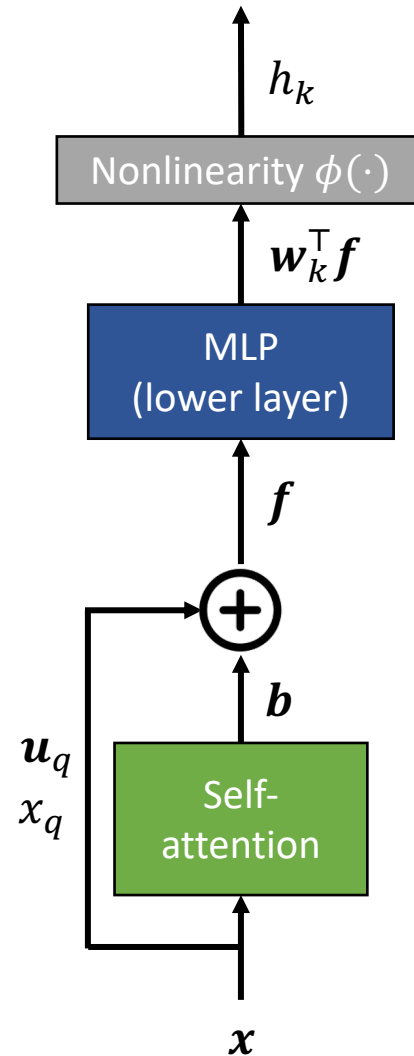
JoMA: JOint Dynamics of MLP/Attention layers



Main Contributions:

1. Find a joint dynamics that connects MLP with self-attention.
2. Understand self-attention behaviors for linear/nonlinear activations.
3. Explain how data hierarchy is learned in multi-layer Transformers.

JoMA Settings



$$h_k = \phi(\mathbf{w}_k^T \mathbf{f})$$

$$\mathbf{f} = U_C \mathbf{b} + \mathbf{u}_q$$

U_C and \mathbf{u}_q are embeddings

$$\mathbf{b} = \sigma(\mathbf{z}_q) \circ \mathbf{x} / A$$

$$\text{SoftmaxAttn: } b_l = \frac{x_l e^{z_{ql}}}{\sum_l x_l e^{z_{ql}}}$$

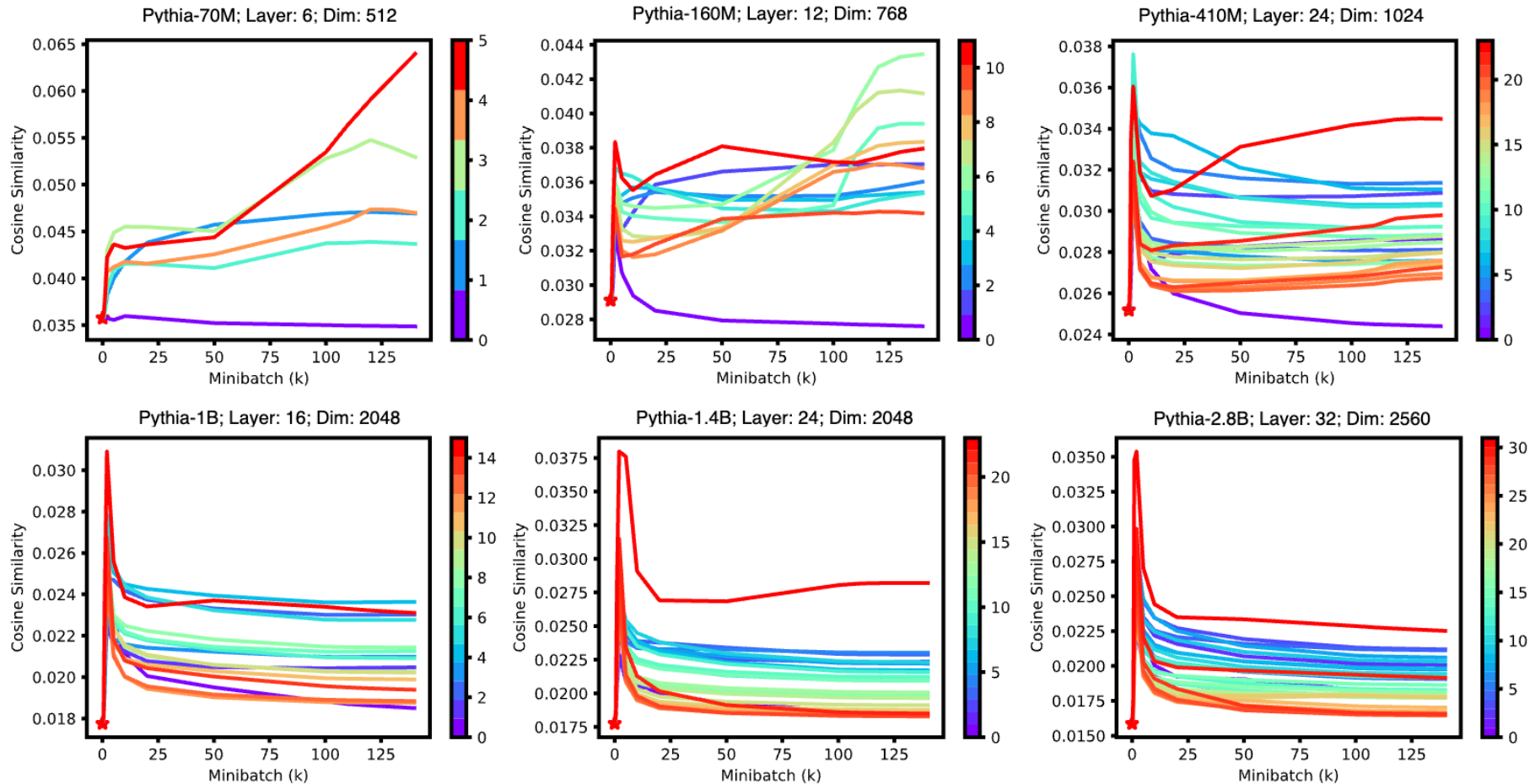
$$\text{ExpAttn: } b_l = x_l e^{z_{ql}}$$

$$\text{LinearAttn: } b_l = x_l z_{ql}$$

"This is an apple"

Assumption (Orthogonal Embeddings $[U_C, u_q]$)

Cosine similarity between embedding vectors at different layers.



JoMA Dynamics

Theorem 1 (JoMA). Let $\mathbf{v}_k := U_C^\top \mathbf{w}_k$, then the dynamics of Eqn. 3 satisfies the invariants:

- Linear attention. The dynamics satisfies $\mathbf{z}_m^2(t) = \sum_k \mathbf{v}_k^2(t) + \mathbf{c}$.
- Exp attention. The dynamics satisfies $\mathbf{z}_m(t) = \frac{1}{2} \sum_k \mathbf{v}_k^2(t) + \mathbf{c}$.
- Softmax attention. If $\bar{\mathbf{b}}_m := \mathbb{E}_{q=m}[\mathbf{b}]$ is a constant over time and $\mathbb{E}_{q=m}[\sum_k g_{h_k} h'_k \mathbf{b} \mathbf{b}^\top] = \bar{\mathbf{b}}_m \mathbb{E}_{q=m}[\sum_k g_{h_k} h'_k \mathbf{b}]$, then the dynamics satisfies $\mathbf{z}_m(t) = \frac{1}{2} \sum_k \mathbf{v}_k^2(t) - \|\mathbf{v}_k(t)\|_2^2 \bar{\mathbf{b}}_m + \mathbf{c}$.

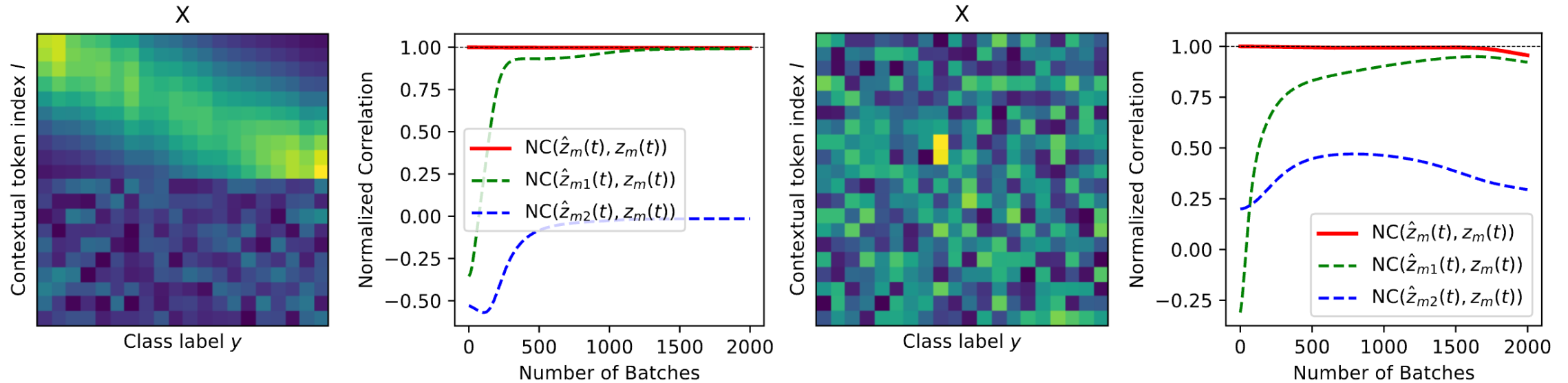
Under zero-initialization ($\mathbf{w}_k(0) = 0, \mathbf{z}_m(0) = 0$), then the time-independent constant $\mathbf{c} = 0$.

There is residual connection.

Joint dynamics works for any learning rates between self-attention and MLP layer.

No assumption on the data distribution.

Verification of JoMA dynamics



$\mathbf{z}_m(t)$: Real attention logits

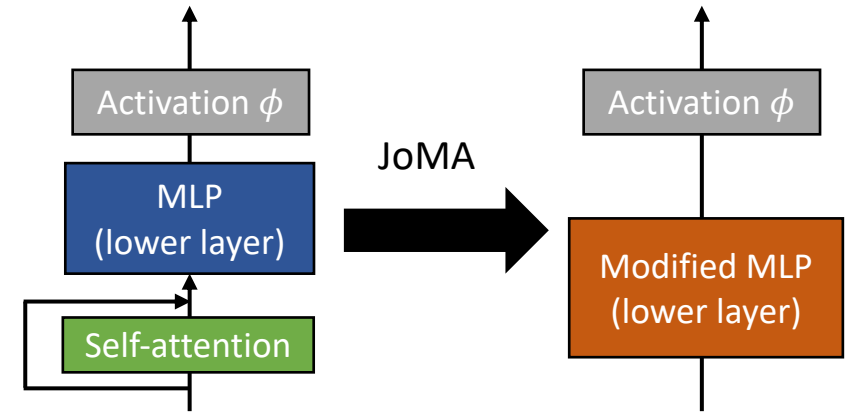
$\hat{\mathbf{z}}_m(t)$: Estimated attention logits by JoMA

$$\hat{\mathbf{z}}_m(t) = \underbrace{\frac{1}{2} \sum_k \mathbf{v}_k^2(t)}_{\hat{\mathbf{z}}_{m1}(t)} - \underbrace{\|\mathbf{v}_k(t)\|_2^2 \bar{\mathbf{b}}_m}_{\hat{\mathbf{z}}_{m2}(t)} + \mathbf{c}$$

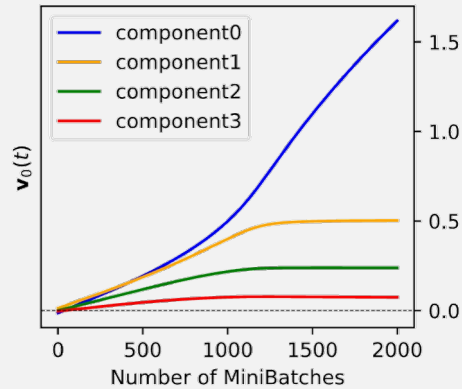
Implication of Theorem 1

Key idea: folding self-attention into MLP

→ A Transformer block becomes a modified MLP

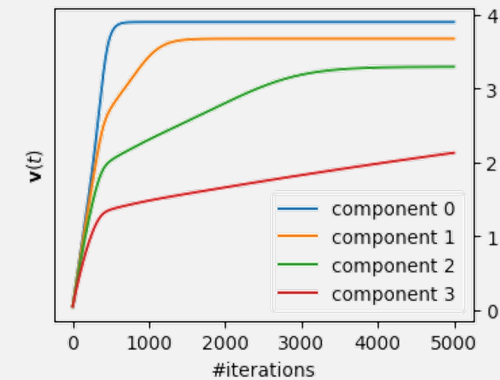


Linear case ($\phi = \text{Id}, K = 1$)



Most salient feature takes all
(Attention becomes sparser)

Nonlinear case (ϕ nonlinear, $K = 1$)



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

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

\uparrow \uparrow
Discriminancy **CoOccurrence**

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

JoMA for Linear Activation

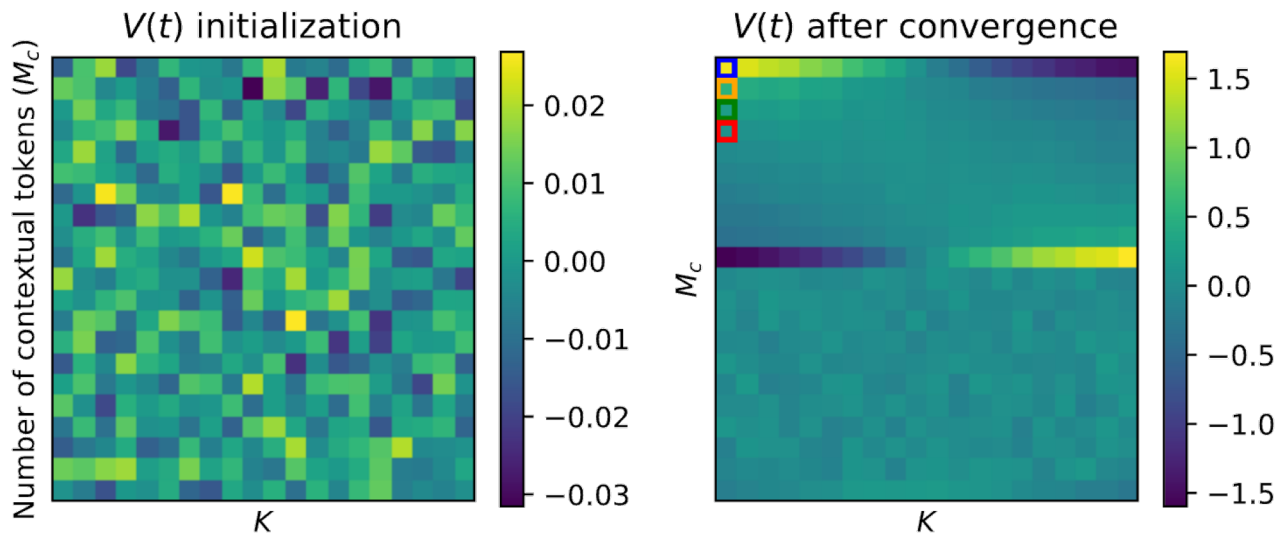
Theorem 2

We can prove $\frac{\text{erf}(v_l(t)/2)}{\Delta_{lm}} = \frac{\text{erf}(v_{l'}(t)/2)}{\Delta_{l'm}}$

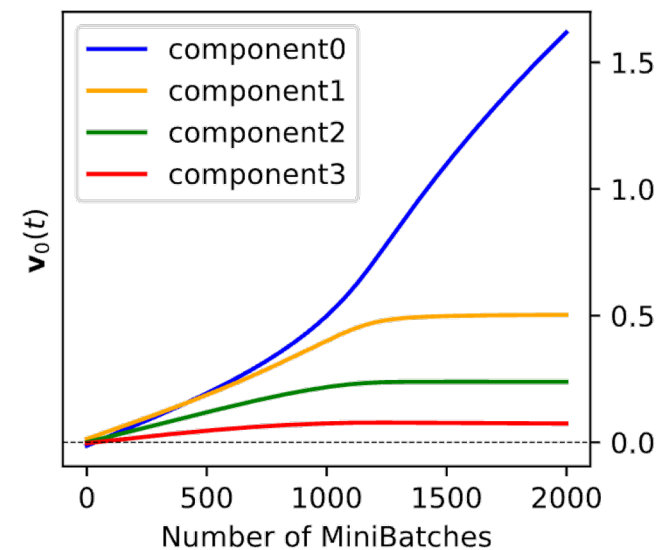
$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \in [-1, 1]$$

Only the most salient token $l^* = \text{argmax } |\Delta_{lm}|$ of \mathbf{v} goes to $+\infty$ other components stay finite.

	Linear
$\dot{\mathbf{v}} = \Delta_m \circ \exp\left(\frac{\mathbf{v}^2}{2}\right)$	Modified MLP (lower layer)



Attention becomes sparser
(Consistent with Scan&Snap)



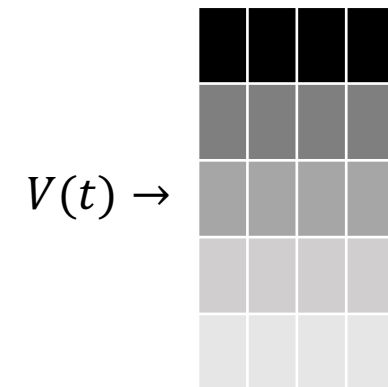
What if we have more nodes ($K > 1$)?

- $V = U_C^T W \in \mathbb{R}^{M_c \times K}$ and the dynamics becomes

$$\dot{V} = \frac{1}{A} \text{diag} \left(\exp \left(\frac{V \circ V}{2} \right) \mathbf{1} \right) \Delta \quad \Delta = [\Delta_1, \Delta_2, \dots, \Delta_K], \quad \Delta_k = \mathbb{E}[g_k \mathbf{x}]$$

We can prove that V gradually becomes low rank

- The growth rate of each row of V varies widely.

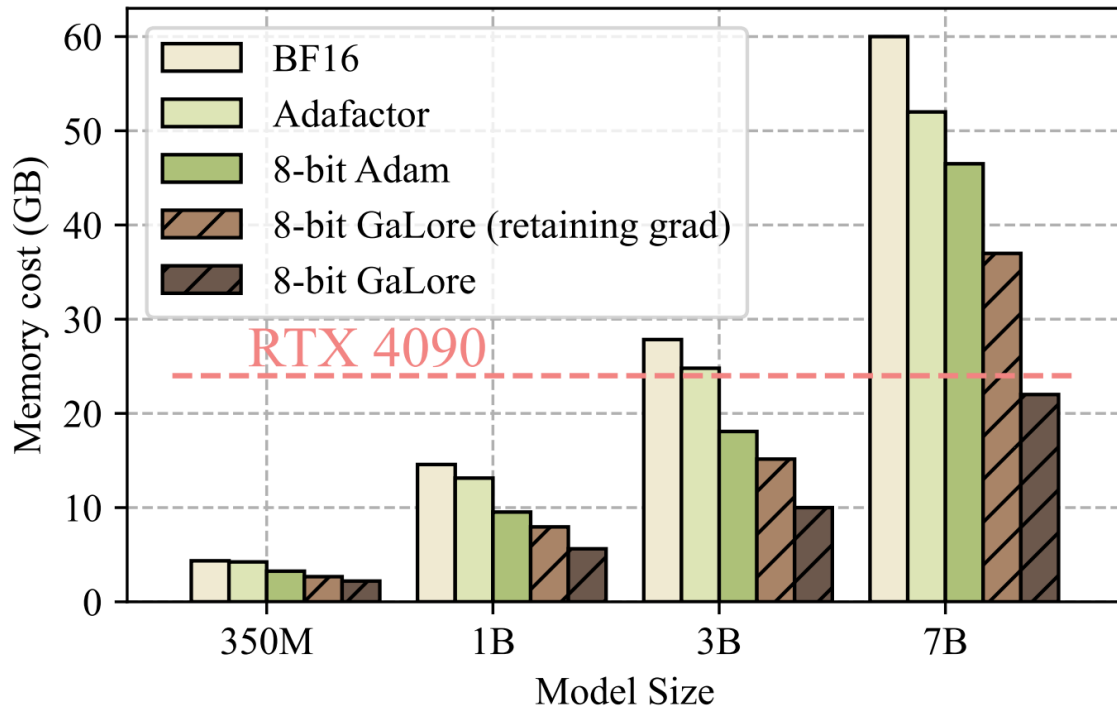


Due to $\exp \left(\frac{V \circ V}{2} \right)$, the weight gradient \dot{V} can be even more low-rank \rightarrow **GaLore**

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



Memory Comparison



	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

* 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

$$G_t \leftarrow -\nabla_W \phi(W_t)$$

If $t \% T == 0$:

 Compute $P_t = \text{SVD}(G_t) \in \mathbb{R}^{m \times r}$

$$R_t \leftarrow P_t^T G_t \quad \{\text{project}\}$$

$$\tilde{R}_t \leftarrow \rho(R_t) \quad \{\text{Adam in low-rank}\}$$

$$\tilde{G}_t \leftarrow P_t \tilde{R}_t \quad \{\text{project-back}\}$$

$$W_{t+1} \leftarrow W_t + \eta \tilde{G}_t$$

Memory Usage	Weight (W)	Optim States (M_t, V_t)	Projection (P)	Total
Full-rank	mn	$2mn$	0	$3mn$
Low-rank adaptor	$mn + mr + nr$	$2(mr + nr)$	0	$mn + 3(mr + nr)$
GaLore	mn	$2nr$	mr	$mn + mr + 2nr$


↑
 W_t

↑
 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 B
350M	1024	2736	16	24	60K	7.8 B
1 B	2048	5461	24	32	100K	13.1 B
7 B	4096	11008	32	32	150K	19.7 B

	Mem	40K	80K	120K	150K
 8-bit GaLore	18G	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

* 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	34.88 (0.24G)	25.36 (0.52G)	18.95 (1.22G)	15.64 (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

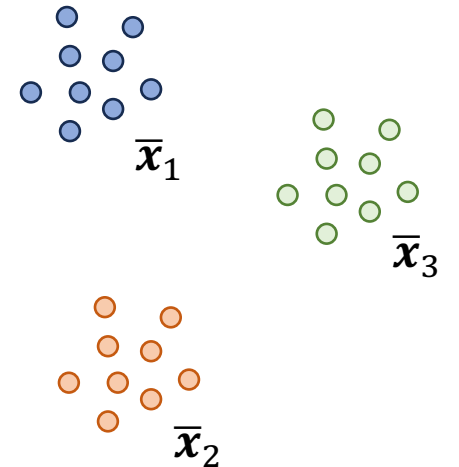
JoMA for Nonlinear Activation

Theorem 3

If \mathbf{x} is sampled from a mixture of C isotropic distributions, (i.e., “local salient/non-salient map”), then

$$\dot{\mathbf{v}} = \frac{1}{\|\mathbf{v}\|_2} \sum_c a_c \theta_1(r_c) \bar{\mathbf{x}}_c + \frac{1}{\|\mathbf{v}\|_2^3} \sum_c a_c \theta_2(r_c) \mathbf{v}$$

Here $a_c := \mathbb{E}_{q=m,c}[g_{h_k}] \mathbb{P}[c]$, $r_c = \mathbf{v}^\top \bar{\mathbf{x}}_c + \int_0^t \mathbb{E}_{q=m}[g_{h_k} h'_k] dt$, and θ_1 and θ_2 depends on nonlinearity



What does the dynamics look like?

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

$\boldsymbol{\mu} \sim \bar{\mathbf{x}}_c$: Critical point due to nonlinearity (one of the cluster centers)

JoMA for Nonlinear activation

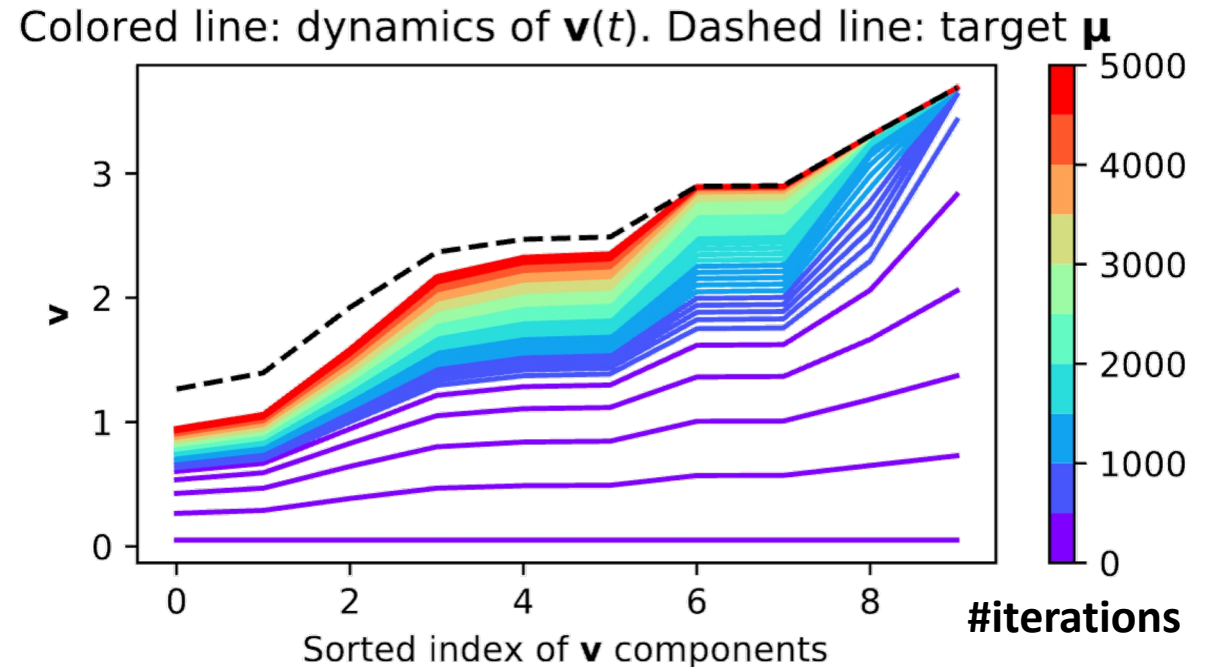
$\dot{\mathbf{v}} = (\boldsymbol{\mu} - \mathbf{v}) \circ \exp\left(\frac{\mathbf{v}^2}{2}\right)$	Nonlinear
	Modified MLP (lower layer)

Theorem 4

Salient components grow much faster than non-salient ones:

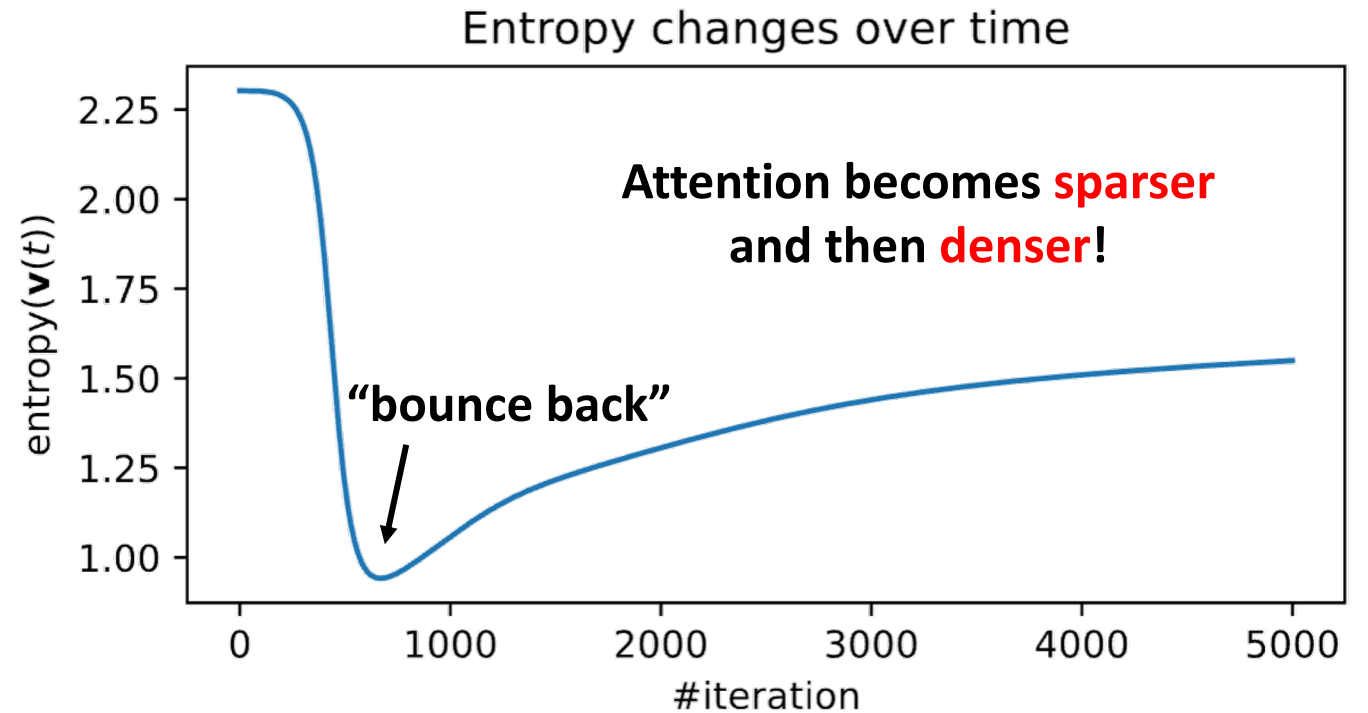
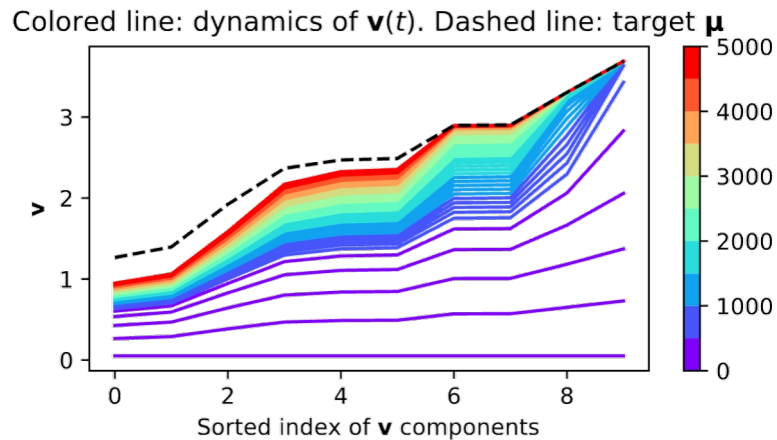
$$\frac{\text{ConvergenceRate}(j)}{\text{ConvergenceRate}(k)} \sim \frac{\exp(\mu_j^2/2)}{\exp(\mu_k^2/2)}$$

$$\begin{aligned} \text{ConvergenceRate}(j) &:= \ln 1/\delta_j(t) \\ \delta_j(t) &:= 1 - v_j(t)/\mu_j \end{aligned}$$



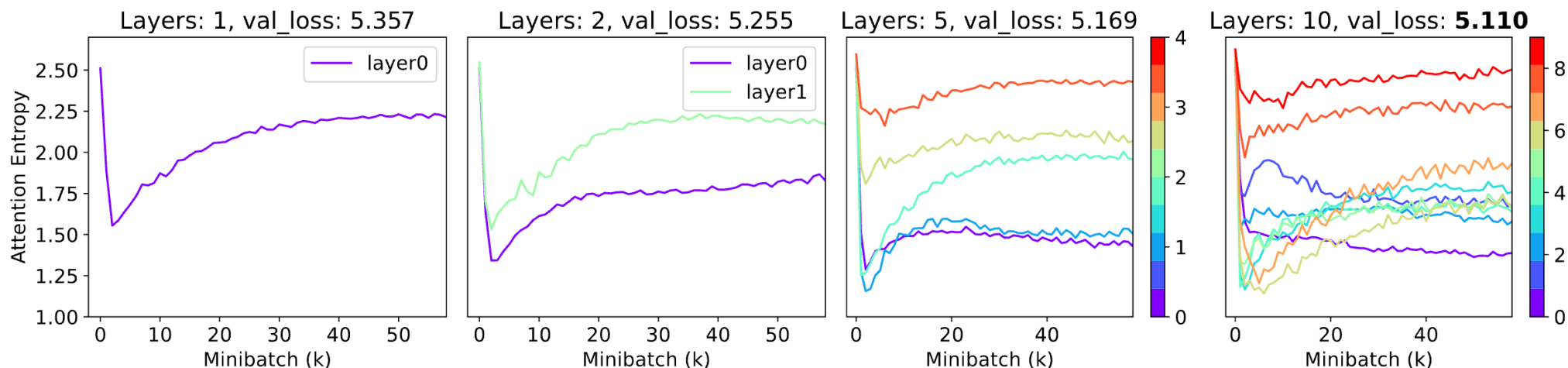
JoMA for Nonlinear activation

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

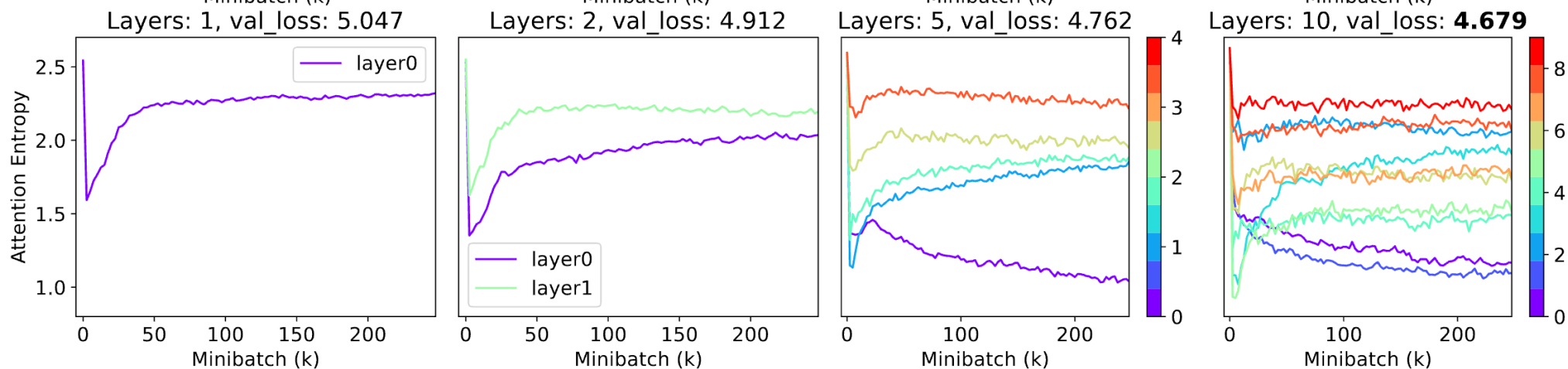


Real-world Experiments

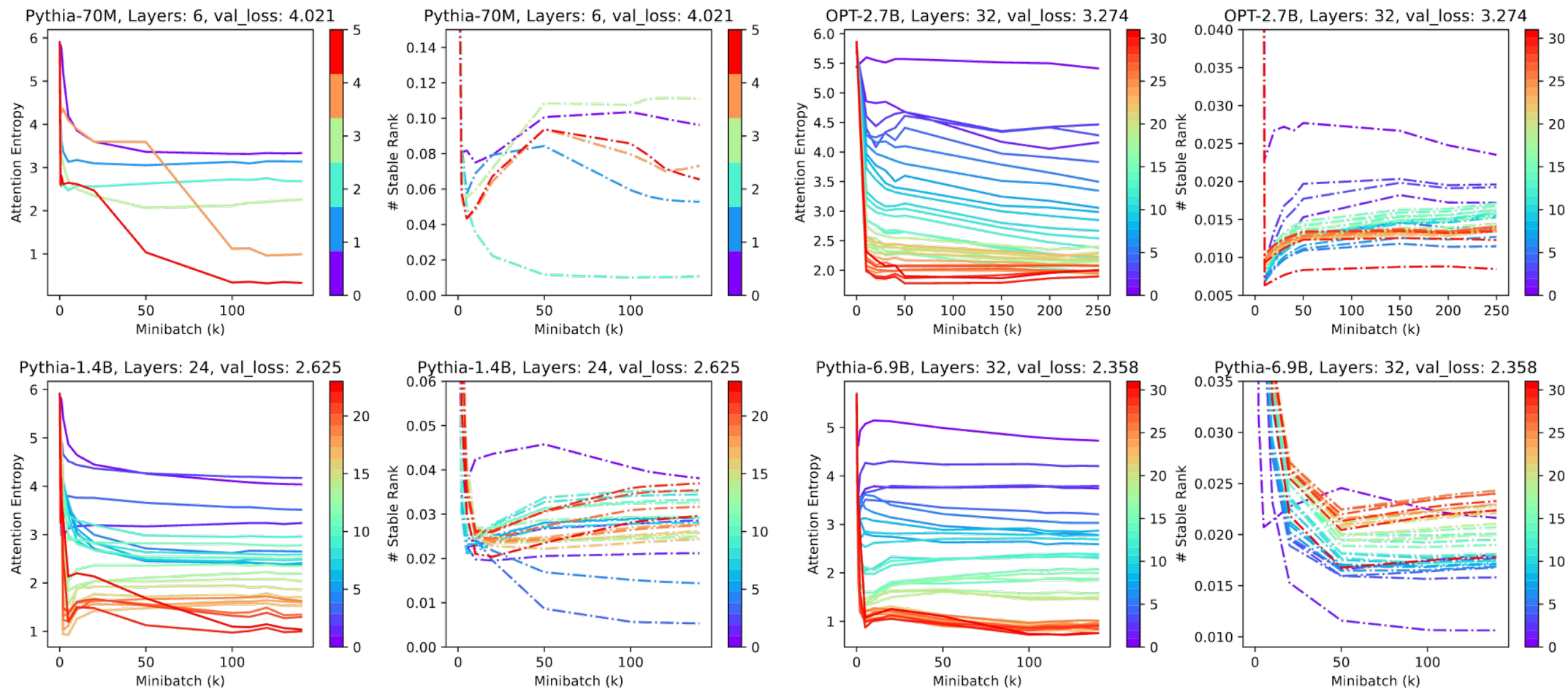
Wikitext2



Wikitext103



Real-world Experiments

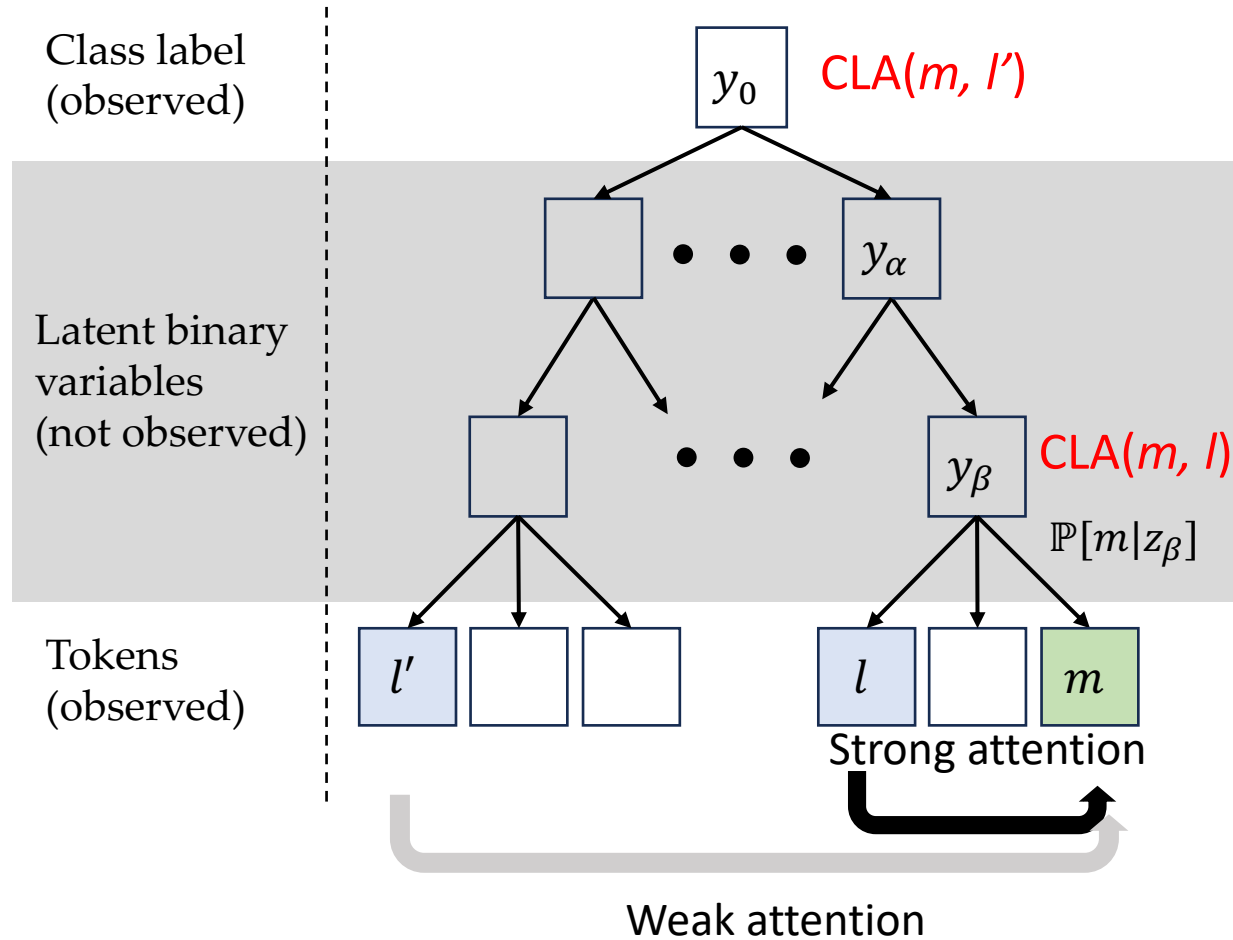


Why is this “bouncing back” property useful?

It seems that it only slows down the training??

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

Data Hierarchy & Multilayer Transformer



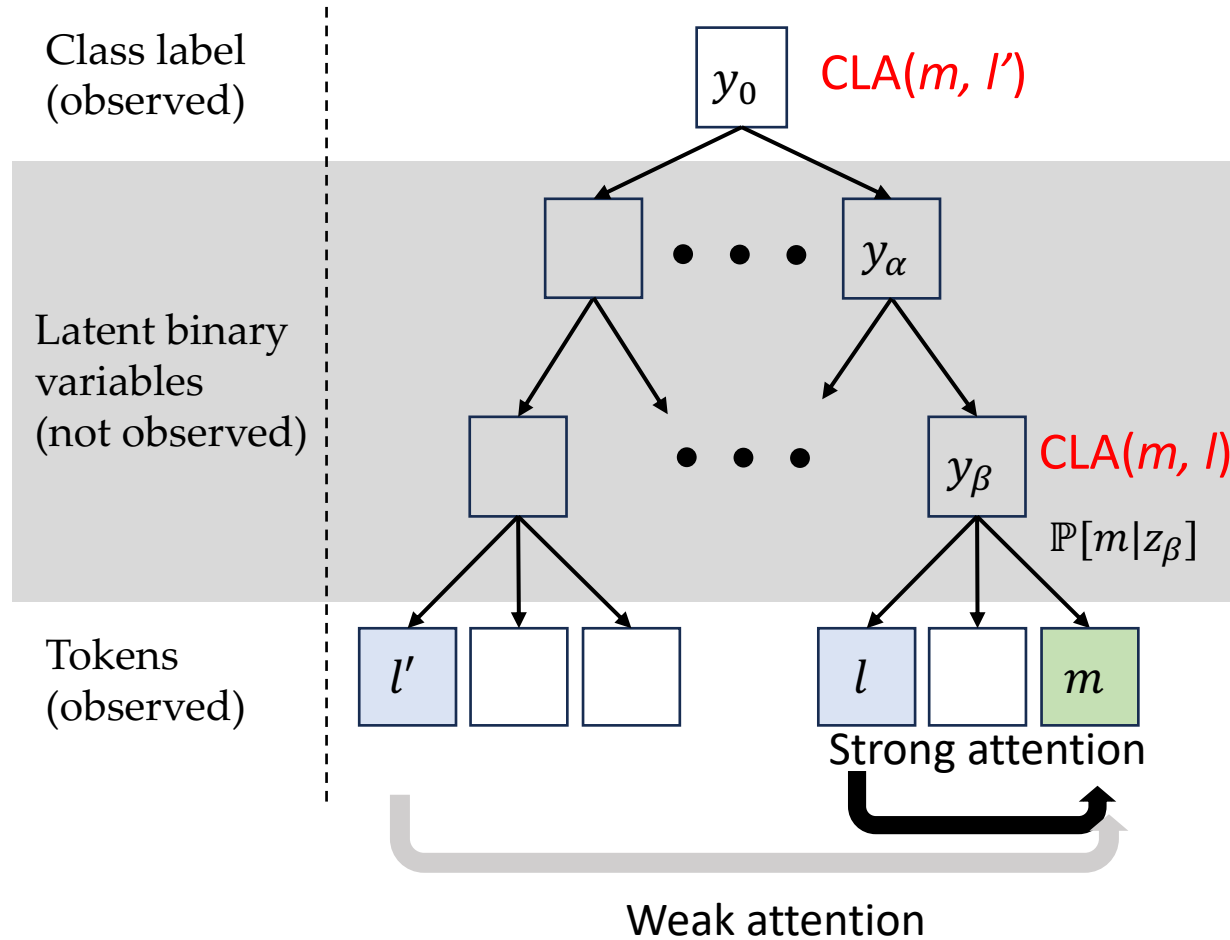
Data Hierarchy & Multilayer Transformer

Theorem 5

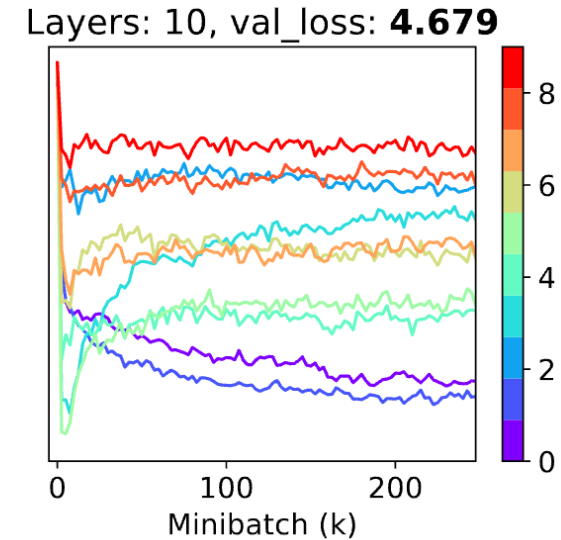
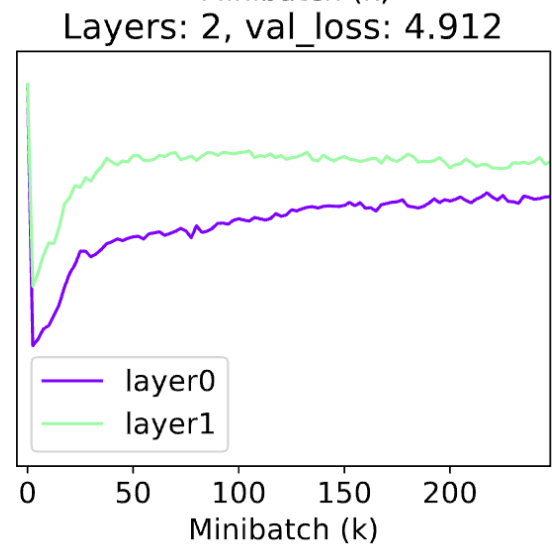
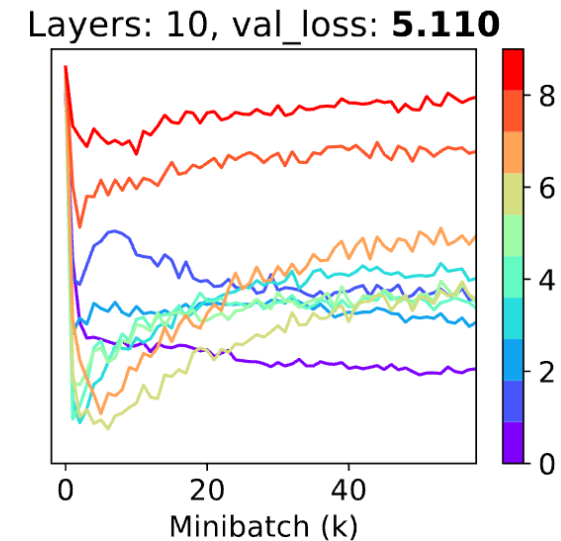
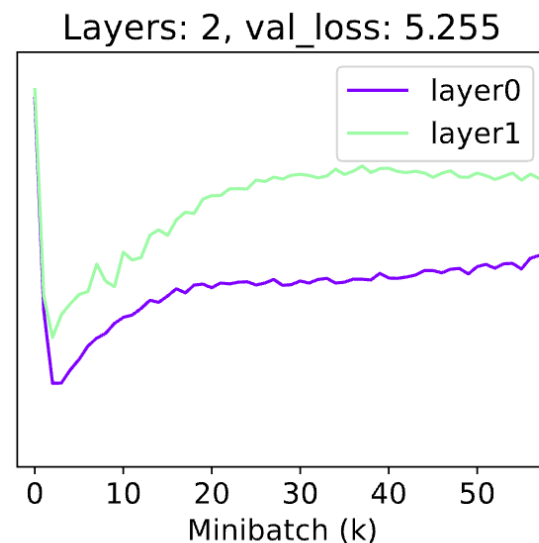
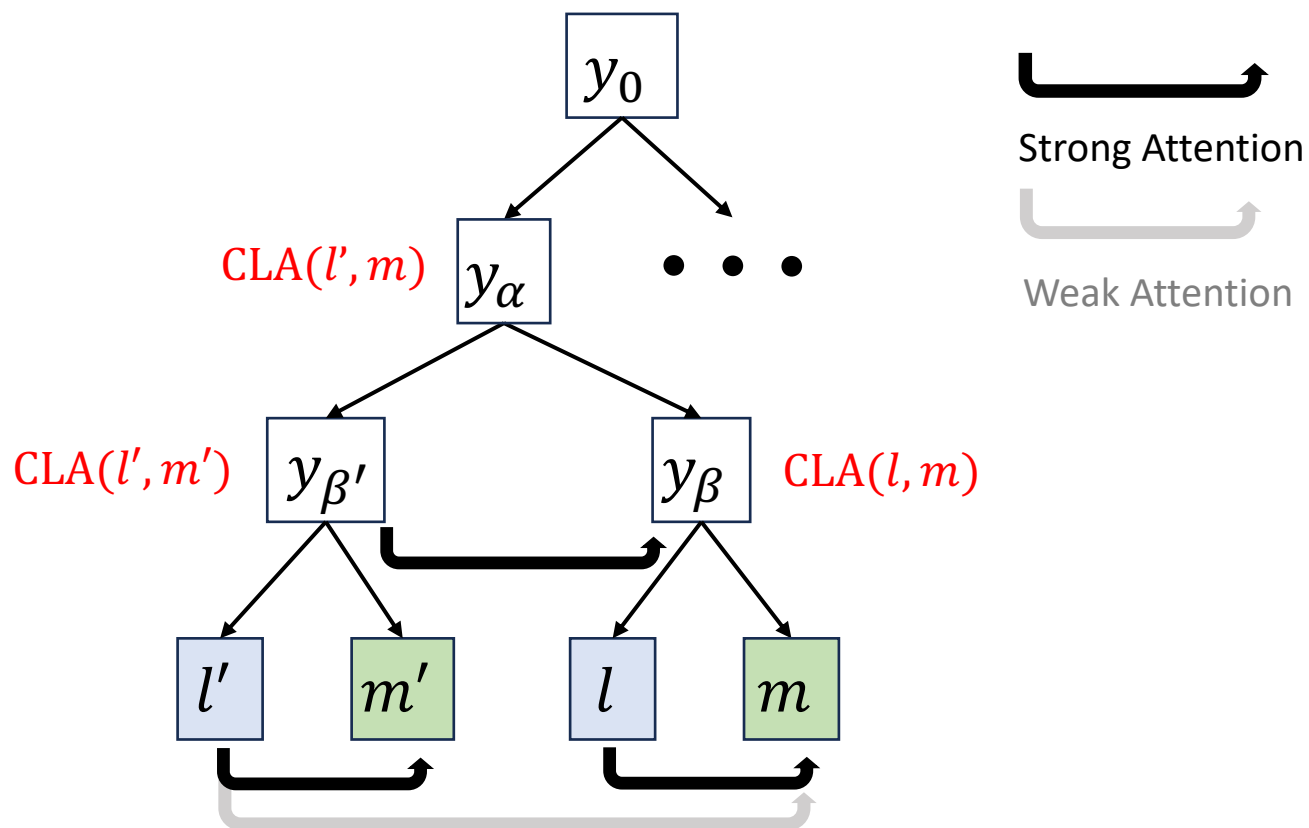
$$\mathbb{P}[l|m] \approx 1 - \frac{H}{L}$$

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

L : total height of the hierarchy



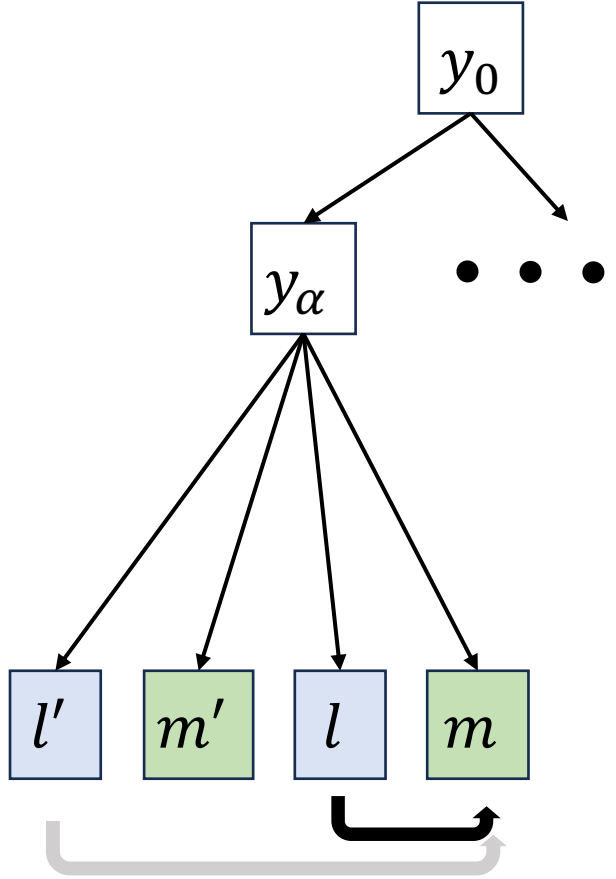
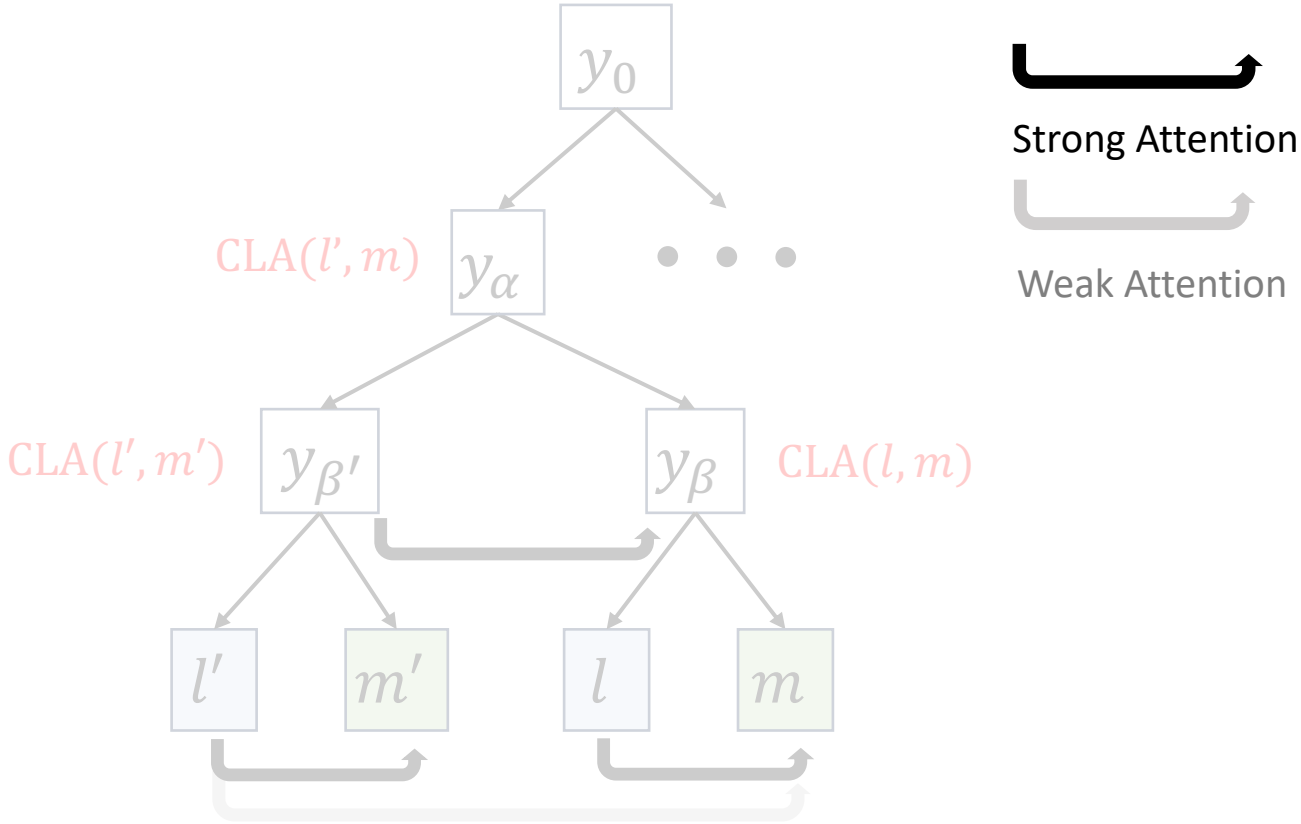
Deep Latent Distribution



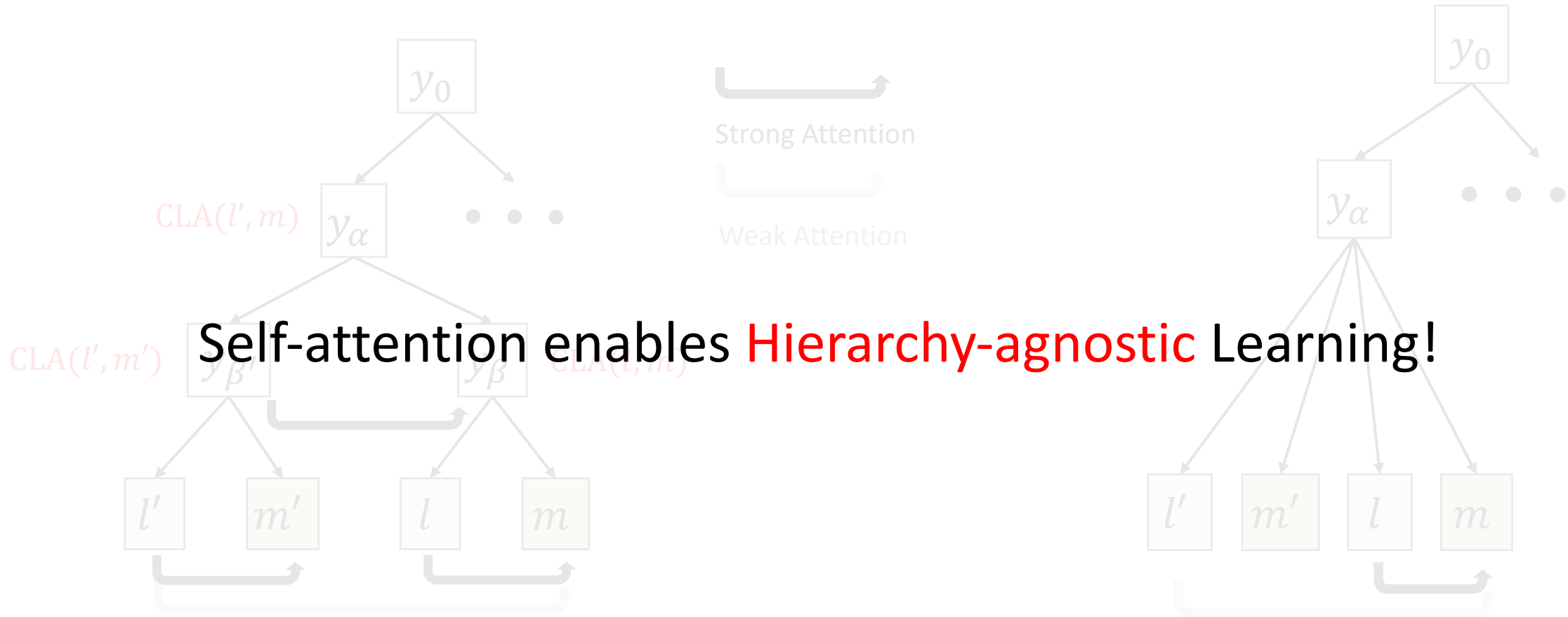
Learning the current hierarchical structure by

slowing down the association of tokens that are not directly correlated

Shallow Latent Distribution



Hierarchy-agnostic Learning

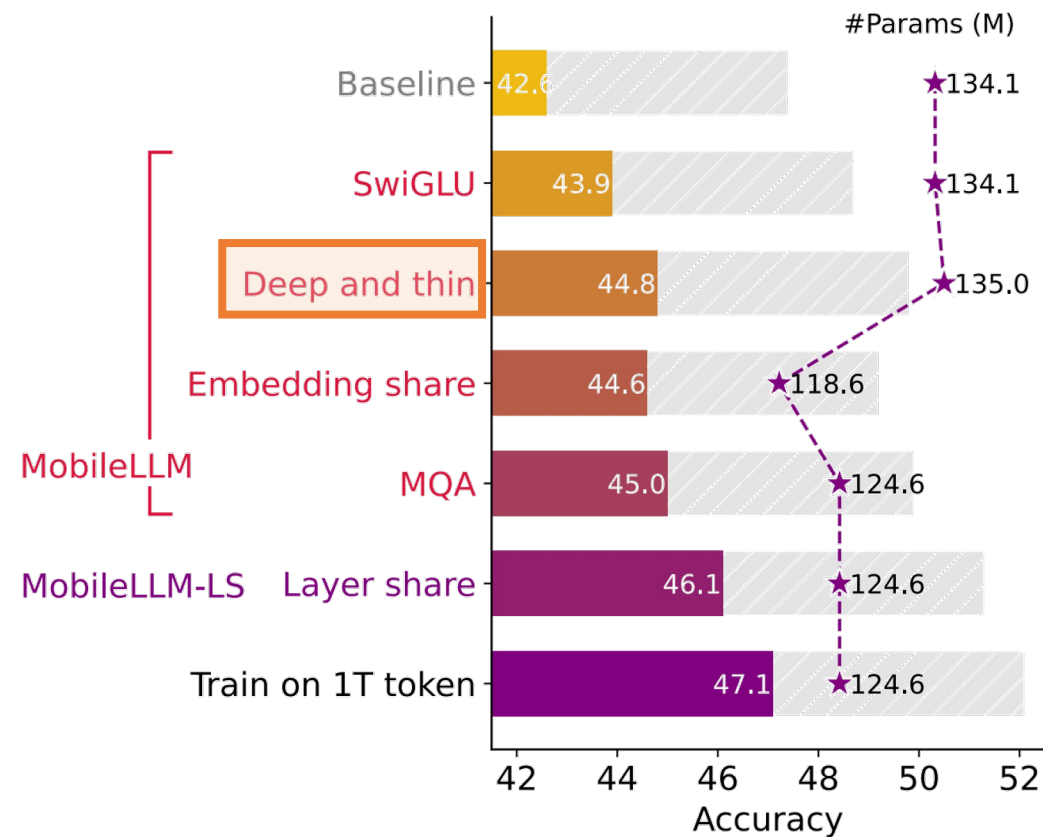
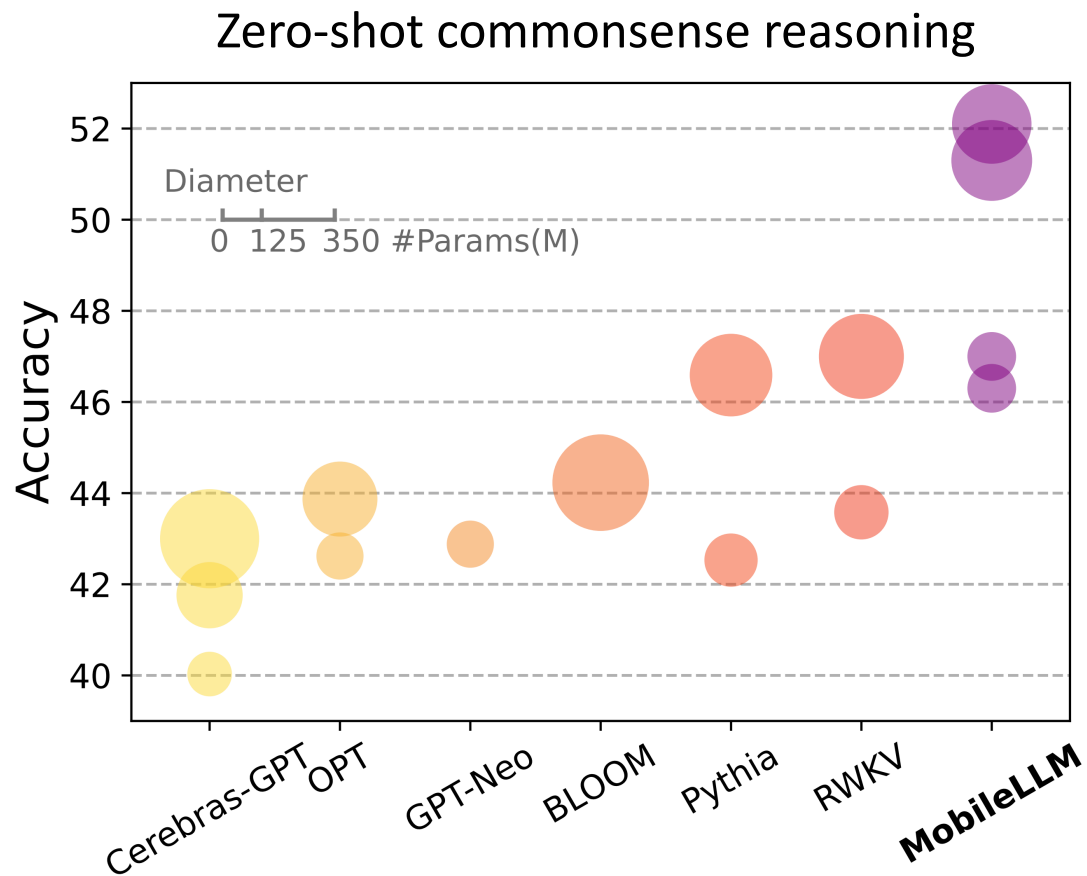


Verification of Hierarchical Intuitions

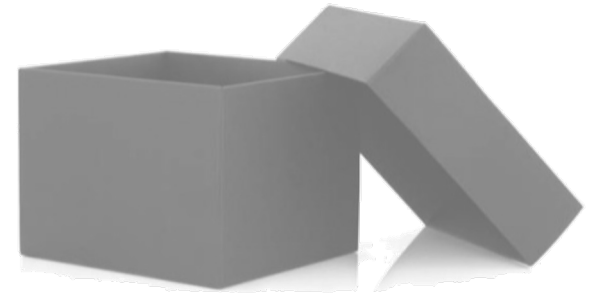
	$C = 20, N_{\text{ch}} = 2$		$C = 20, N_{\text{ch}} = 3$		$C = 30, N_{\text{ch}} = 2$	
(N_0, N_1)	(10, 20)	(20, 30)	(10, 20)	(20, 30)	(10, 20)	(20, 30)
NCorr ($s = 0$)	0.99 ± 0.01	0.97 ± 0.02	1.00 ± 0.00	0.96 ± 0.02	0.99 ± 0.01	0.94 ± 0.04
NCorr ($s = 1$)	0.81 ± 0.05	0.80 ± 0.05	0.69 ± 0.05	0.68 ± 0.04	0.73 ± 0.08	0.74 ± 0.03
	$C = 30, N_{\text{ch}} = 3$		$C = 50, N_{\text{ch}} = 2$		$C = 50, N_{\text{ch}} = 3$	
(N_0, N_1)	(10, 20)	(20, 30)	(10, 20)	(20, 30)	(10, 20)	(20, 30)
NCorr ($s = 0$)	0.99 ± 0.01	0.95 ± 0.03	0.99 ± 0.01	0.95 ± 0.03	0.99 ± 0.01	0.95 ± 0.03
NCorr ($s = 1$)	0.72 ± 0.04	0.66 ± 0.02	0.58 ± 0.02	0.55 ± 0.01	0.64 ± 0.02	0.61 ± 0.04

Table 1: Normalized correlation between the latents and their best matched hidden node in MLP of the same layer. All experiments are run with 5 random seeds.

MobileLLM



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

Thanks!