# Large Language Models (LLMs)



Conversational AI



Content Generation



AI Agents



Reasoning



Planning

# What does the future look like?



More data

More compute

Larger models

**Are we going to blindly believe in scaling laws?**

# Black-box versus White-box



Black box

# 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

+    −

−    +

## Expressibility

"Neural Network is a universal approximator"
"Deep Models can express functions more efficiently than shallow ones"

## Understanding how Deep Models work

## Optimization

"Gradient vanishing/exploding"
"Gradient Descent might get stuck at saddle point / local minima"
"Can GD/SGD go to global optima? How fast?"

## Generalization

Which path should we take?

"Does zero training error often lead to overfitting?"
"More parameters might lead to overfitting."

# Three Angles – What to pick?

**Expressibility**

| + | - |
|---|---|
| - | + |

Architecture ✓

training dynamics ✗

**Optimization**



Architecture ✗

training dynamics ✓

**Generalization**



Architecture ✗

training dynamics ✗

**How about**

Architecture ✓

training dynamics ✓

# Start From the First Principle

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

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

- **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**

Architecture ✓

training dynamics ✓

# Transformers



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Masked Multi-Head Attention

Nx

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

Key $K$

Query $Q$

Attention mechanism

[A. Vaswani et al, Attention is all you need, NeurIPS'17]

# Understanding Attention in 1-layer Setting

Decoding & Softmax

Normalization

Self-attention

Different next token =
Different classes in classification

$x_1$  $x_2$  • • •  $x_{T-1}$  $x_T$  $x_{T+1}$

Contextual tokens

Last/query token

Next token

[*Y. Tian et al*, Scan and Snap: Understanding Training Dynamics and Token Composition in 1-layer Transformer, *NeurIPS'23*]

# Reparameterization

- Parameters $W_K, W_Q, W_V, U$ makes the dynamics complicated.

- Reparameterize the problem with independent variable $Y$ and $Z$
    - $Y = UW_V^T U^T$ (*Merging the embedding with weight matrix*)
    - $Z = UW_Q W_K^T U^T$ (pairwise logits of self-attention matrix)

- Then the dynamics becomes easier to analyze

Class Prediction

$Y$

$f_n$

normalize

$Z$

Data

# 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: $z_{ml}(0) = 0$

*Pairwise **attention score** between token $l$ and query $m$*

**Distinct tokens:** Tokens that only appear in a single class.

**Common tokens:** Tokens that appear in multiple classes.

# Overall Picture of the Training Dynamics

Common Token Suppression

$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

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

# Overall Picture of the Training Dynamics

## Winners-emergence



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

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

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

# Overall Picture of the Training Dynamics

## Winners-emergence



Seq class $(m, n_1)$

Seq class $(m, n_2)$

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

(b) $\dot{z}_{ml} > 0$, 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



$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

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

**Theorem 3** Relative gain $r_{l/l'|n}(t) := \dfrac{\tilde{c}_{l|n}^2(t)}{\tilde{c}_{l'|n}^2(t)} - 1$ 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



$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

**Contextual Sparsity (query-dependent)**

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

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

**Theorem 3** Relative gain $r_{l/l'|n}(t) := \dfrac{\tilde{c}^2_{l|n}(t)}{\tilde{c}^2_{l'|n}(t)} - 1$ 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

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

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

# Overall Picture of the Training Dynamics

## Attention frozen



$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

Contextual Sparsity (query-dependent)

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

**Theorem 4** When $t \to +\infty$,

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

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) $\eta_Z$ and $\eta_Y$ are large, $B_n(t)$ is large and attention is sparse

(2) Fixing $\eta_Z$, large $\eta_Y$ leads to slightly small $B_n(t)$ and denser attention

# Overall Picture of the Training Dynamics

## Winners-emergence



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

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

# Overall Picture of the Training Dynamics

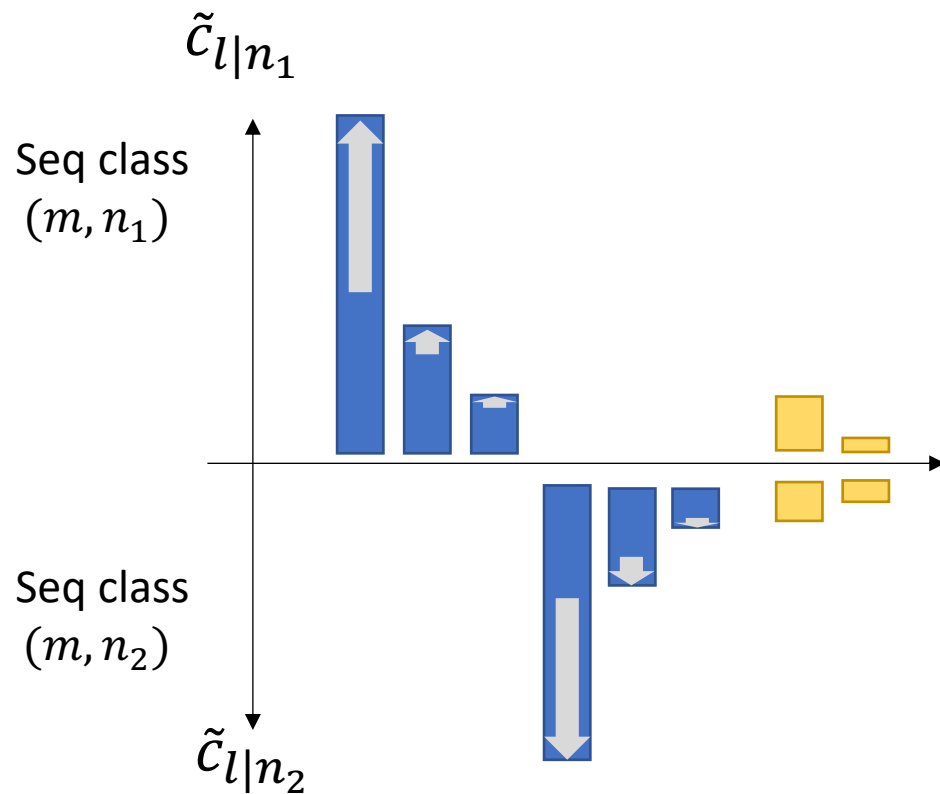Winners-emergence



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

(b) $\dot{z}_{ml} > 0$, 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

$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

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

**Theorem 3** Relative gain $r_{l/l'|n}(t) := \dfrac{\tilde{c}_{l|n}^2(t)}{\tilde{c}_{l'|n}^2(t)} - 1$ has a close form:
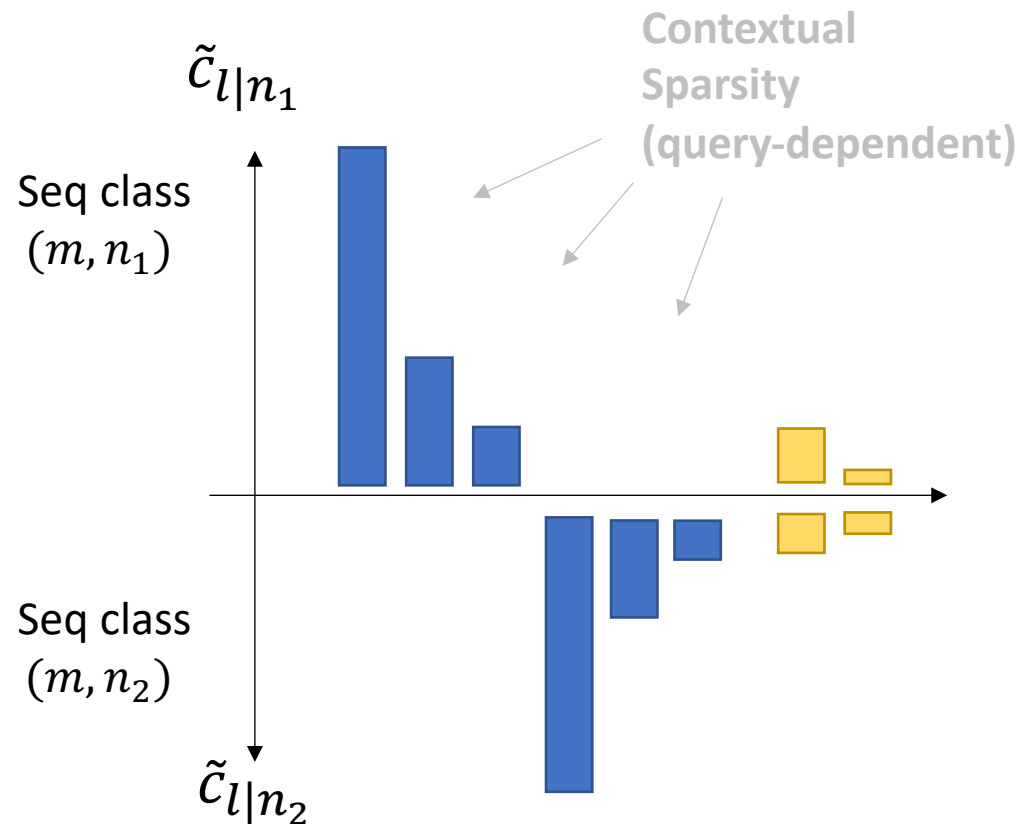
$$r_{l/l'|n}(t) = r_{l/l'|n}(0)\chi_l(t)$$

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



$\tilde{c}_{l|n_1}$

Seq class $(m, n_1)$

**Contextual Sparsity (query-dependent)**

Seq class $(m, n_2)$

$\tilde{c}_{l|n_2}$

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

**Theorem 3** Relative gain $r_{l/l'|n}(t) := \dfrac{\tilde{c}_{l|n}^2(t)}{\tilde{c}_{l'|n}^2(t)} - 1$ 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

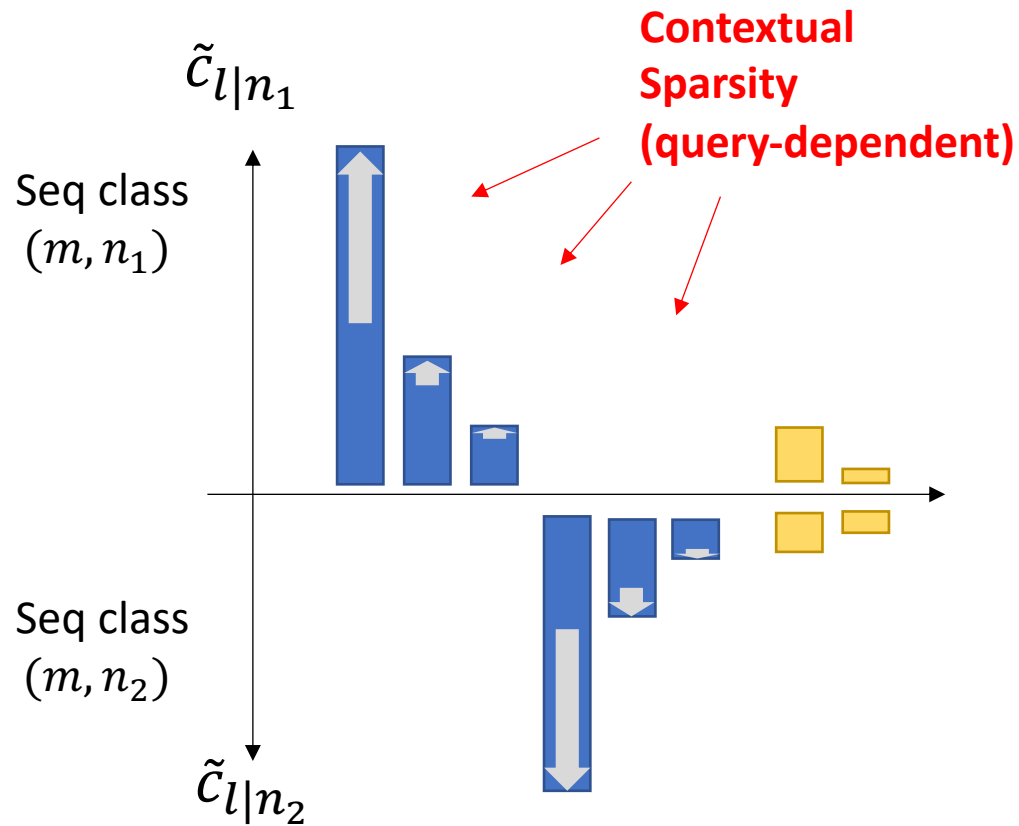$$e^{2f_{nl_0}^2(0)B_n(t)} \leq \chi_{l_0}(t) \leq e^{2B_n(t)}$$

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

# Overall Picture of the Training Dynamics

## Attention frozen



Seq class $(m, n_1)$

Seq class $(m, n_2)$

Contextual Sparsity (query-dependent)

**Theorem 4** When $t \to +\infty$,

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

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) $\eta_Z$ and $\eta_Y$ are large, $B_n(t)$ is large and attention is sparse

(2) Fixing $\eta_Z$, large $\eta_Y$ leads to slightly small $B_n(t)$ and denser attention

# 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]

# Follow-up works

- Scan & Snap has Multiple Assumptions
  - No positional encoding
  - Sequence length $T \to +\infty$
  - Learning rate of decoder $Y$ larger than self-attention layer Z $(\eta_Y \gg \eta_Z)$
  - Other technical assumptions

- How to get rid of them?
- Follow-up 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 hierachy is learned in multi-layer Transformers.
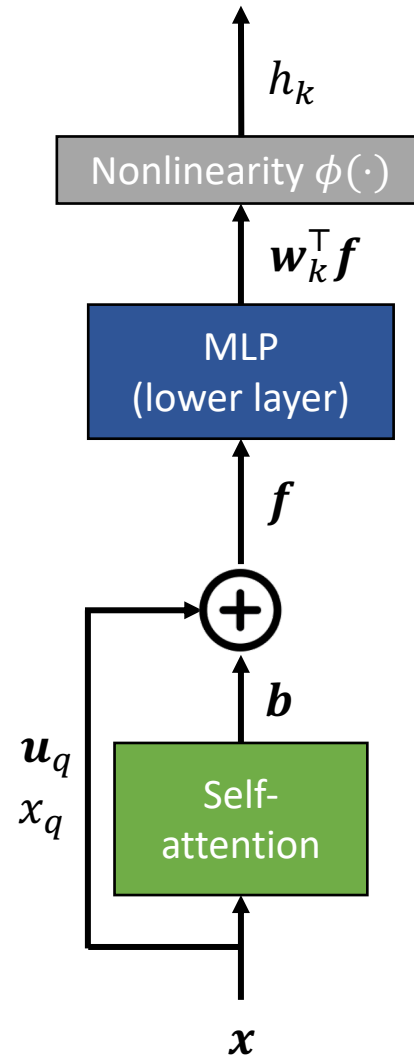
[Y. Tian et al, *JoMA: Demystifying Multilayer Transformers via JOint Dynamics of MLP and Attention,* ICLR'24]

# JoMA Settings



$$h_k = \phi(\boldsymbol{w}_k^\top \boldsymbol{f})$$

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

$U_C$ and $\boldsymbol{u}_q$ are embeddings

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

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

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

LinearAttn: $b_l = x_l z_{ql}$

"This is an apple"

# JoMA Dynamics

**Theorem 1** (JoMA). *Let $\boldsymbol{v}_k := U_C^\top \boldsymbol{w}_k$, then the dynamics of Eqn.* <span style="border:1px solid red">3</span> *satisfies the invariants:*

- *Linear attention. The dynamics satisfies $\boldsymbol{z}_m^2(t) = \sum_k \boldsymbol{v}_k^2(t) + \boldsymbol{c}$.*

- *Exp attention. The dynamics satisfies $\boldsymbol{z}_m(t) = \frac{1}{2}\sum_k \boldsymbol{v}_k^2(t) + \boldsymbol{c}$.*

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

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

<span style="color:red">There is residual connection.</span>
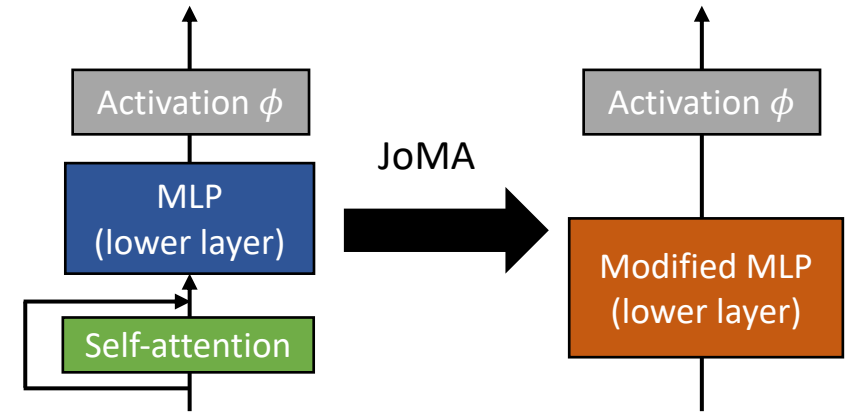<span style="color:red">Joint dynamics works for any learning rates between self-attention and MLP layer.</span>
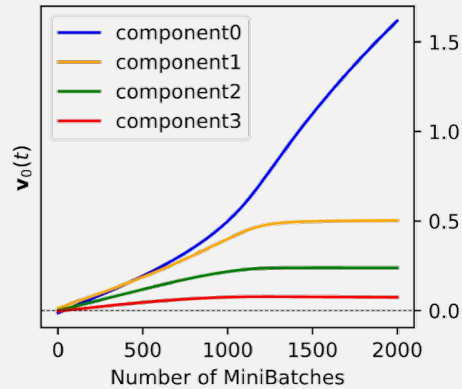<span style="color:red">No assumption on the data distribution.</span>

# Implication of Theorem 1

**Key idea:** folding self-attention into MLP

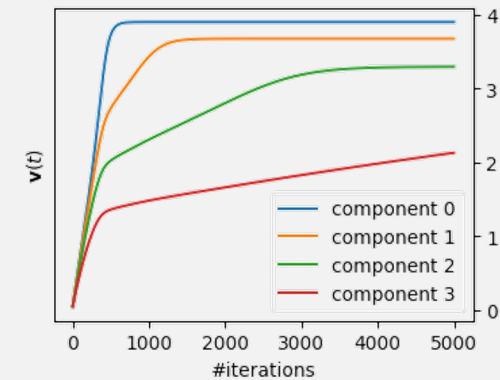→ A Transformer block becomes a modified MLP



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

**Most salient feature takes all**
(Attention becomes sparser)

Nonlinear case ($\phi$ 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

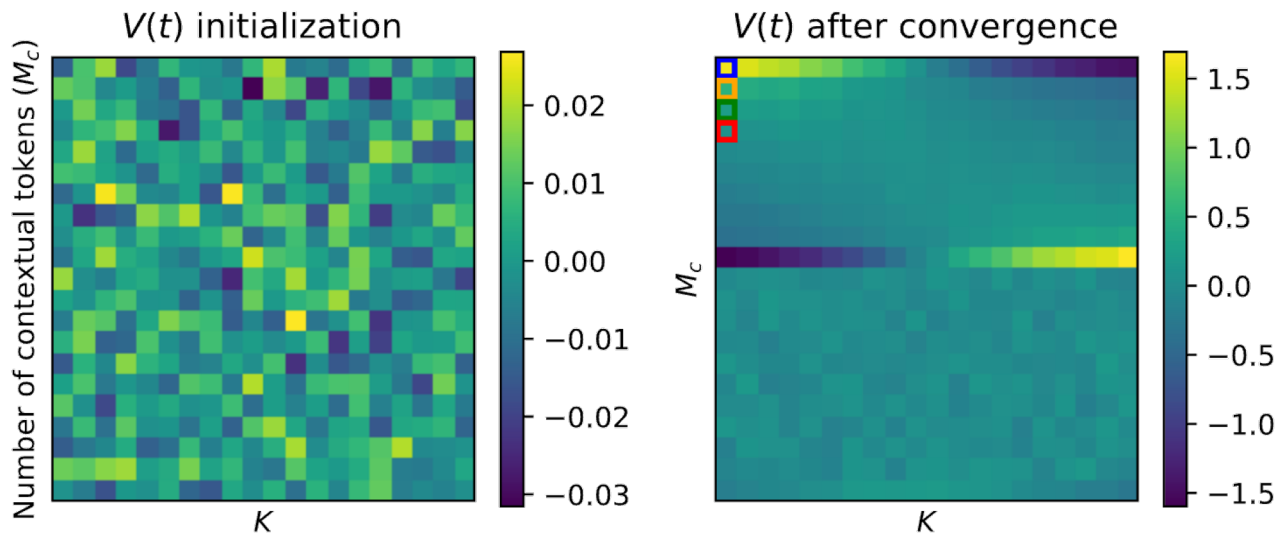$|\Delta_{lm}|$ large: **Distinct** tokens

# JoMA for Linear Activation

$$\dot{\boldsymbol{v}} = \boldsymbol{\Delta}_m \circ \exp\left(\frac{\boldsymbol{v}^2}{2}\right)$$

Linear

Modified MLP (lower layer)

**Theorem 2**

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

$\mathrm{erf}(x) = \dfrac{2}{\sqrt{\pi}} \displaystyle\int_0^x e^{-t^2}\mathrm{d}t \in [-1,1]$

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



**Attention becomes sparser**
(Consistent with Scan&Snap)

[**Y. Tian** et al, *Scan and Snap: Understanding Training Dynamics and Token Composition in 1-layer Transformer,* NeurIPS'23]

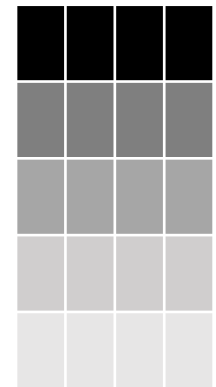# What if we have more nodes $(\boldsymbol{K} > \boldsymbol{1})$?

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

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

We can prove that $V$ gradually becomes low rank
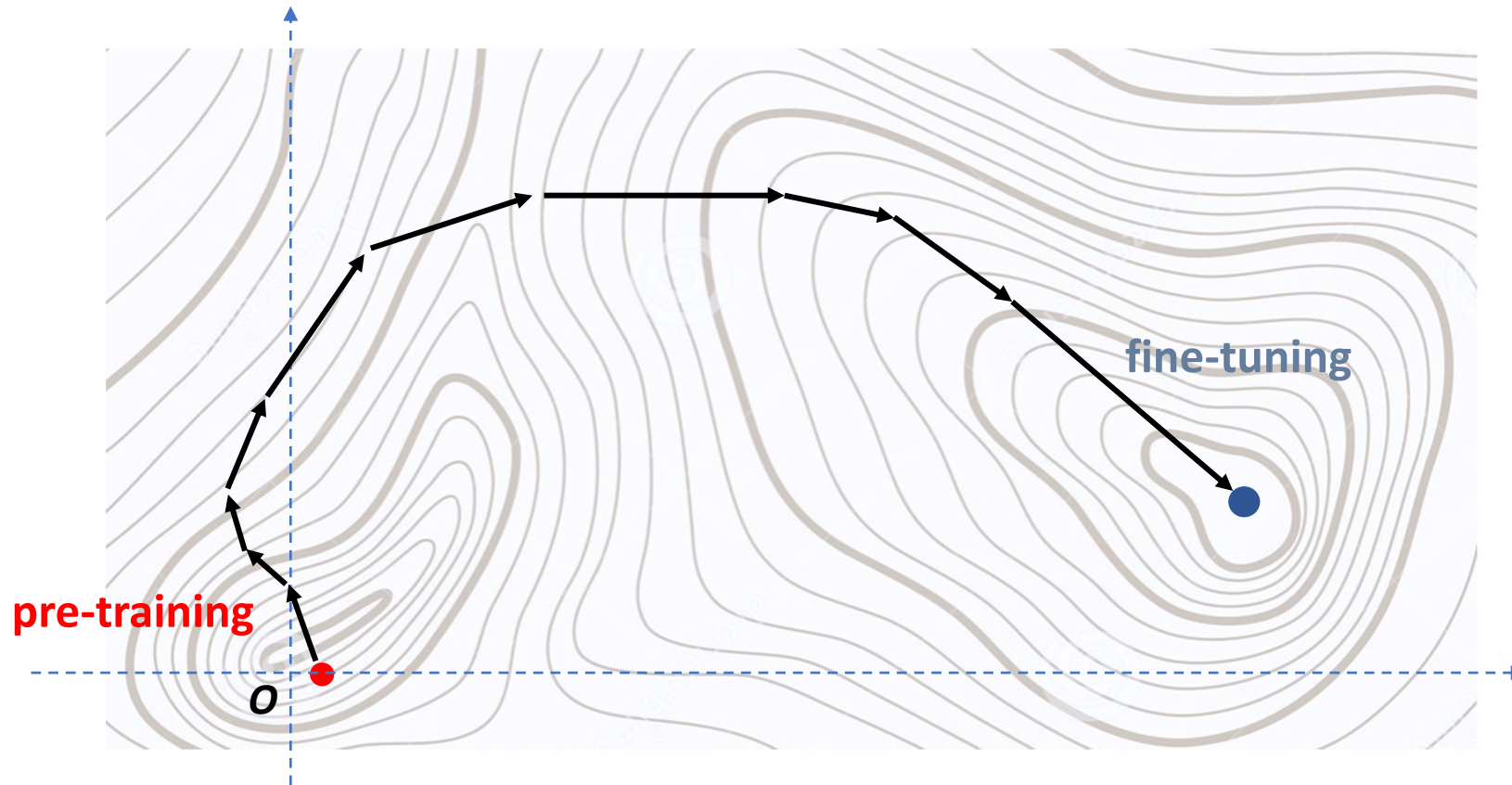- The growth rate of each row of $V$ varies widely.

$V(t) \rightarrow$

**Due to** $\exp\left(\frac{V \circ V}{2}\right)$**, the weight gradient** $\dot{V}$ **can be even more low-rank**
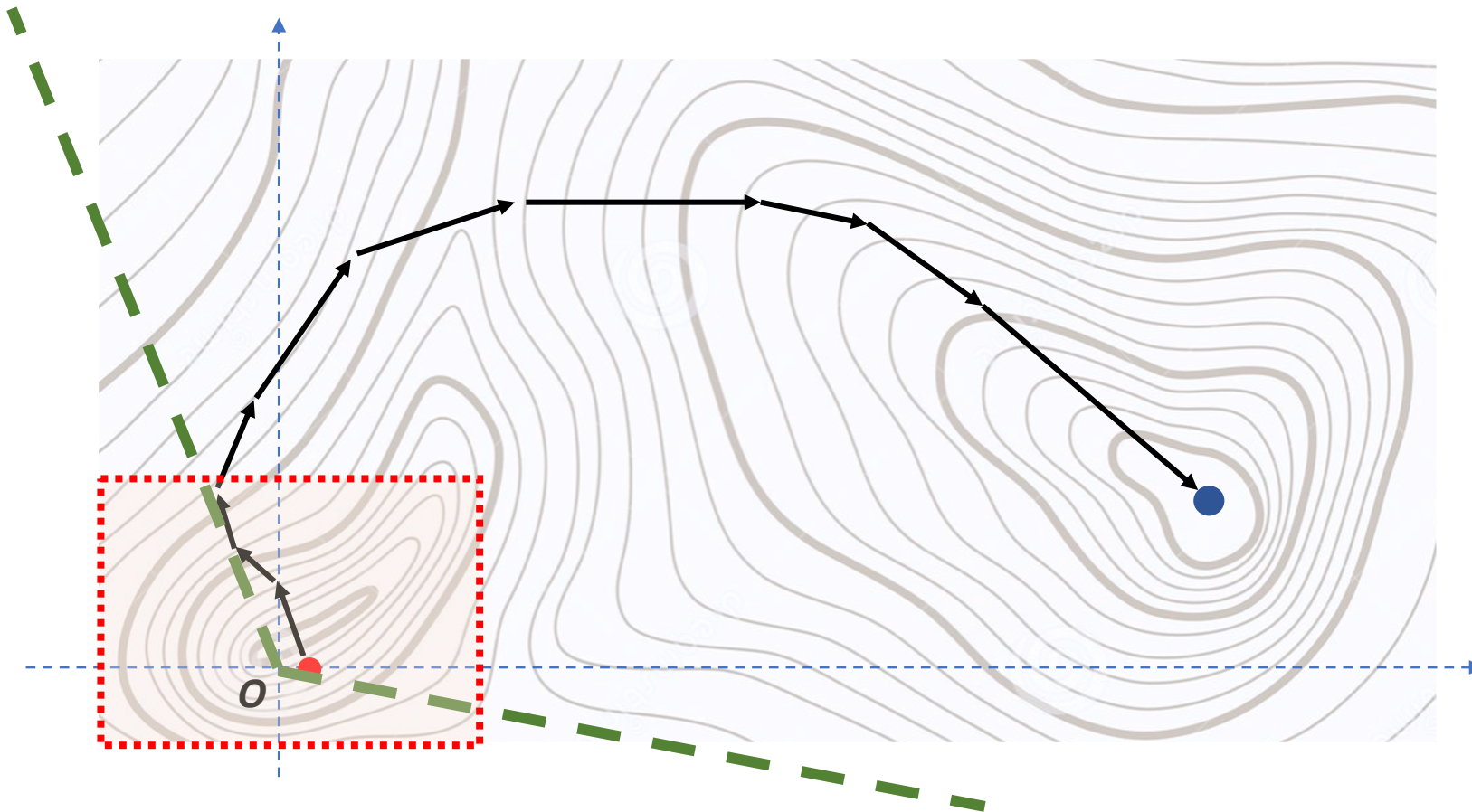
# How the Weight Rank Changes over time?

**Consider the Entire Training Trajectory …**

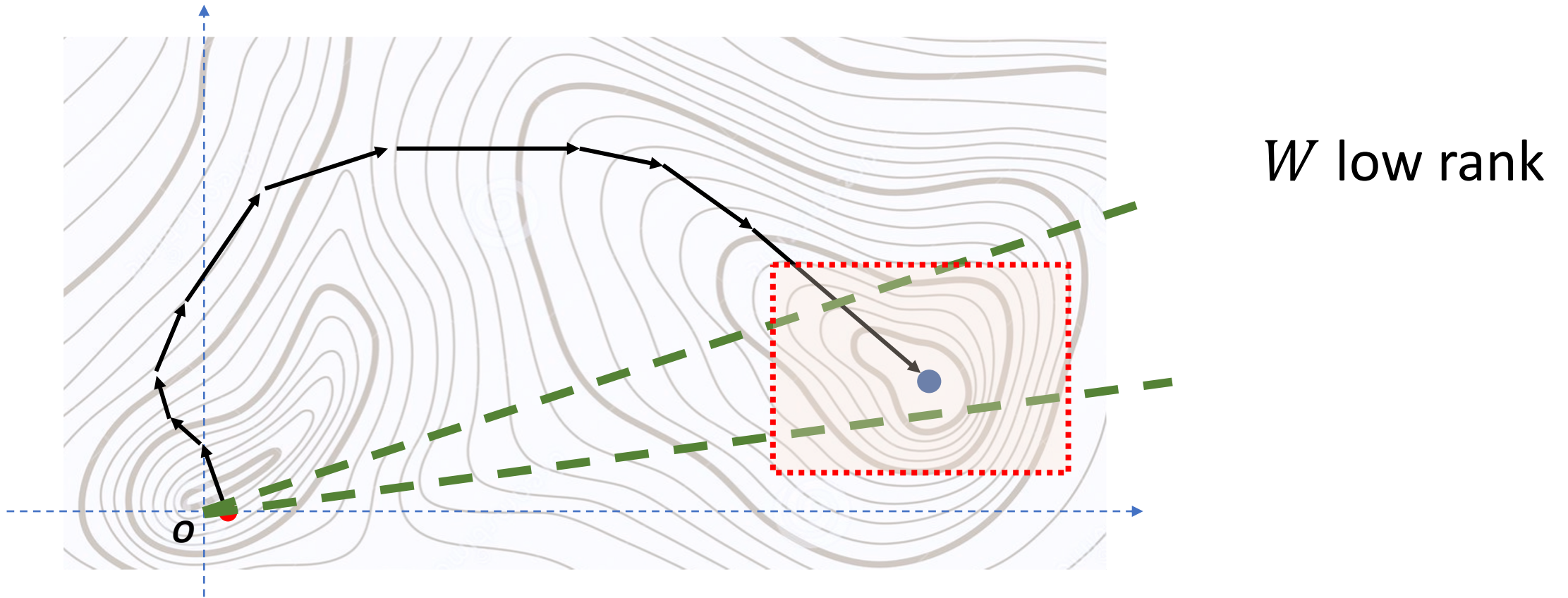# How the Weight Rank Changes over time?

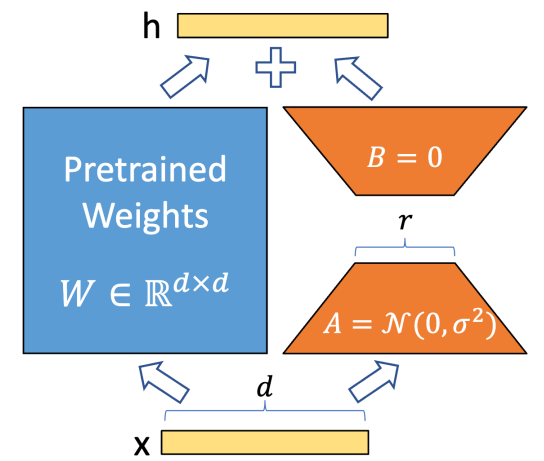**Beginning of Training: Weight subspace changes a lot**



$W$ high rank
(due to random
initialization)

# How the Weight Rank Changes over time?

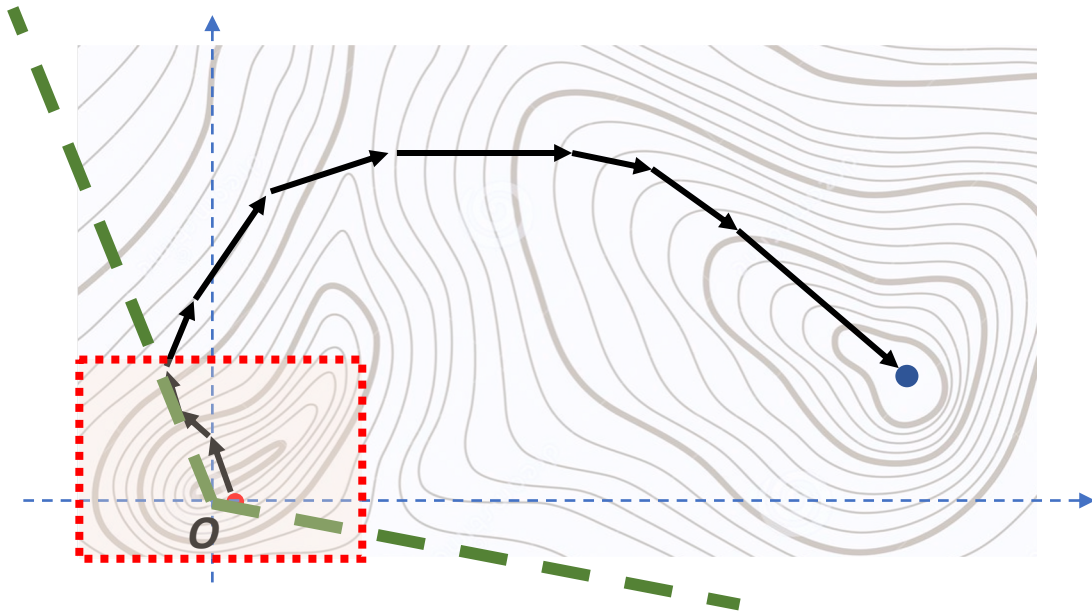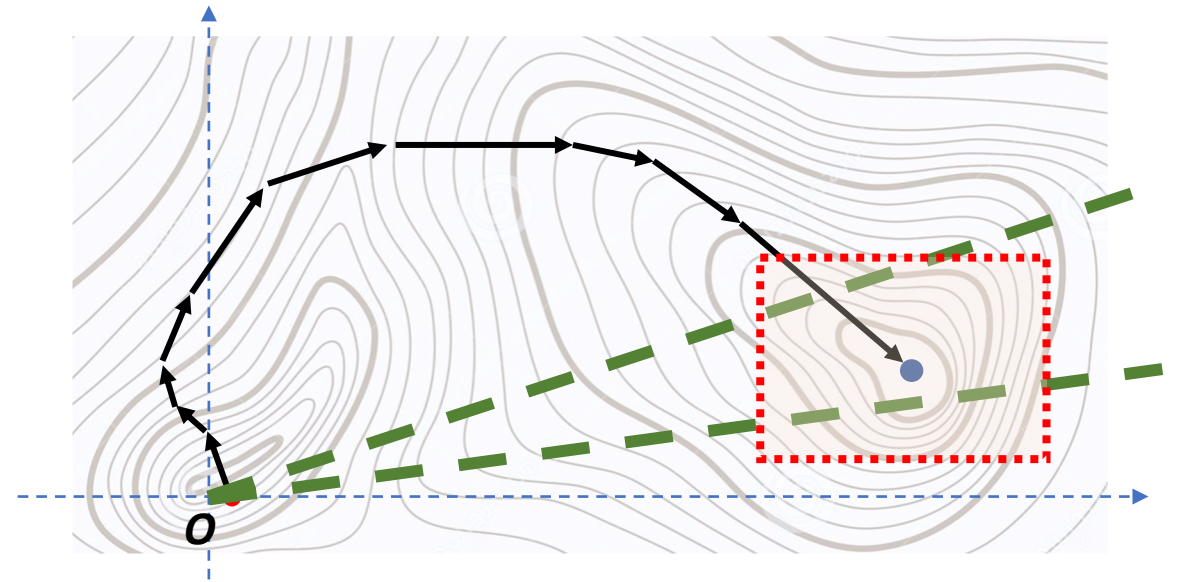**Mid/End of Training: Weight subspace changes little**



$W$ low rank

*o*

facebook Artificial Intelligence

# Think about LoRA?



$W$ high rank (due to random initialization)

$W$ low rank

LoRA (Low-rank Adaption)

*LoRA does not work*

*LoRA can work*

# GaLore

low-rank weights → low-rank gradients

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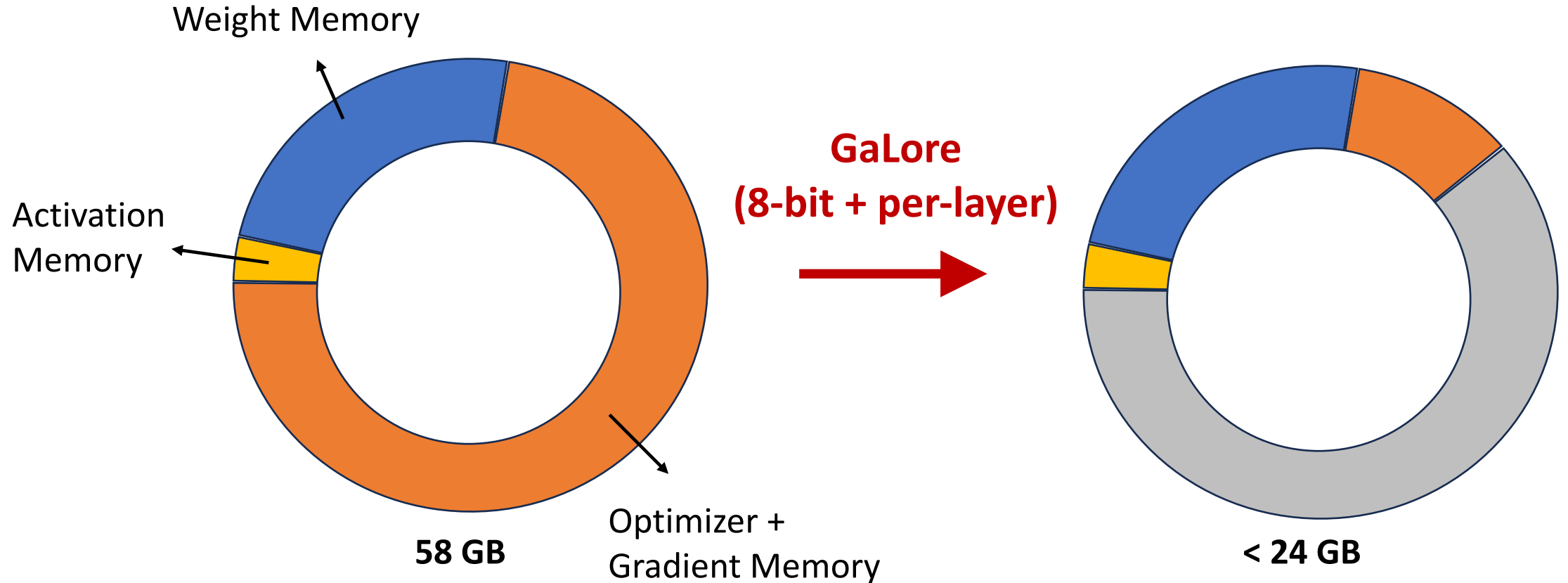
$$G_t \leftarrow -\nabla_W \phi(W_t)$$

If t % T == 0:

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

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

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

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

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

[J. Zhao et al, *GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection*, ICML'24 Oral]

# Memory Saving in GaLore

Weight Memory

Activation
Memory

GaLore
(8-bit + per-layer)

Optimizer +
Gradient Memory

58 GB

< 24 GB

Reduce optimizer states and weight gradients, Achieve **82.5%** mem reduction

# Convergence Analysis on Fixed Projection

For gradient in the following form

$$G = \sum_i A_i - \sum_i B_i W C_i$$

Let $R = P^\top G Q$ be projected gradient (P and Q are fixed) then

$$\|\boldsymbol{R_t}\|_{\boldsymbol{F}} \leq (\boldsymbol{1} - \boldsymbol{\eta M}) \|\boldsymbol{R_{t-1}}\|_{\boldsymbol{F}} \to \boldsymbol{0}$$
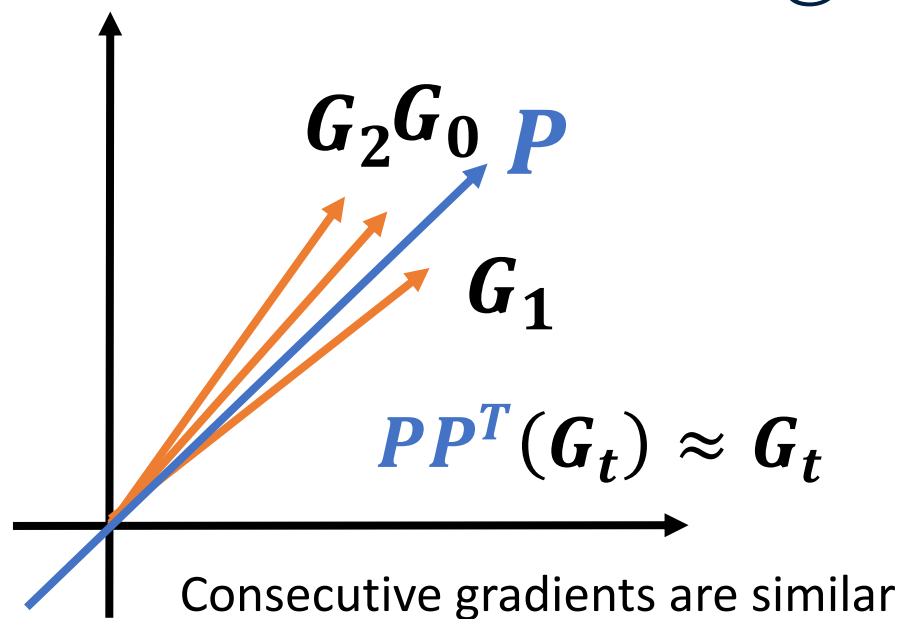
Where $M := \frac{1}{N} \sum_i \min_t \lambda_{\min}(\hat{B}_{it}) \lambda_{\min}(\hat{C}_{it}) - L_A - L_B L_C D^2$

$\hat{B}_{it} = P_t^T B_i(W_t) P_t \qquad \hat{C}_{it} = Q_t^T C_i(W_t) Q_t$

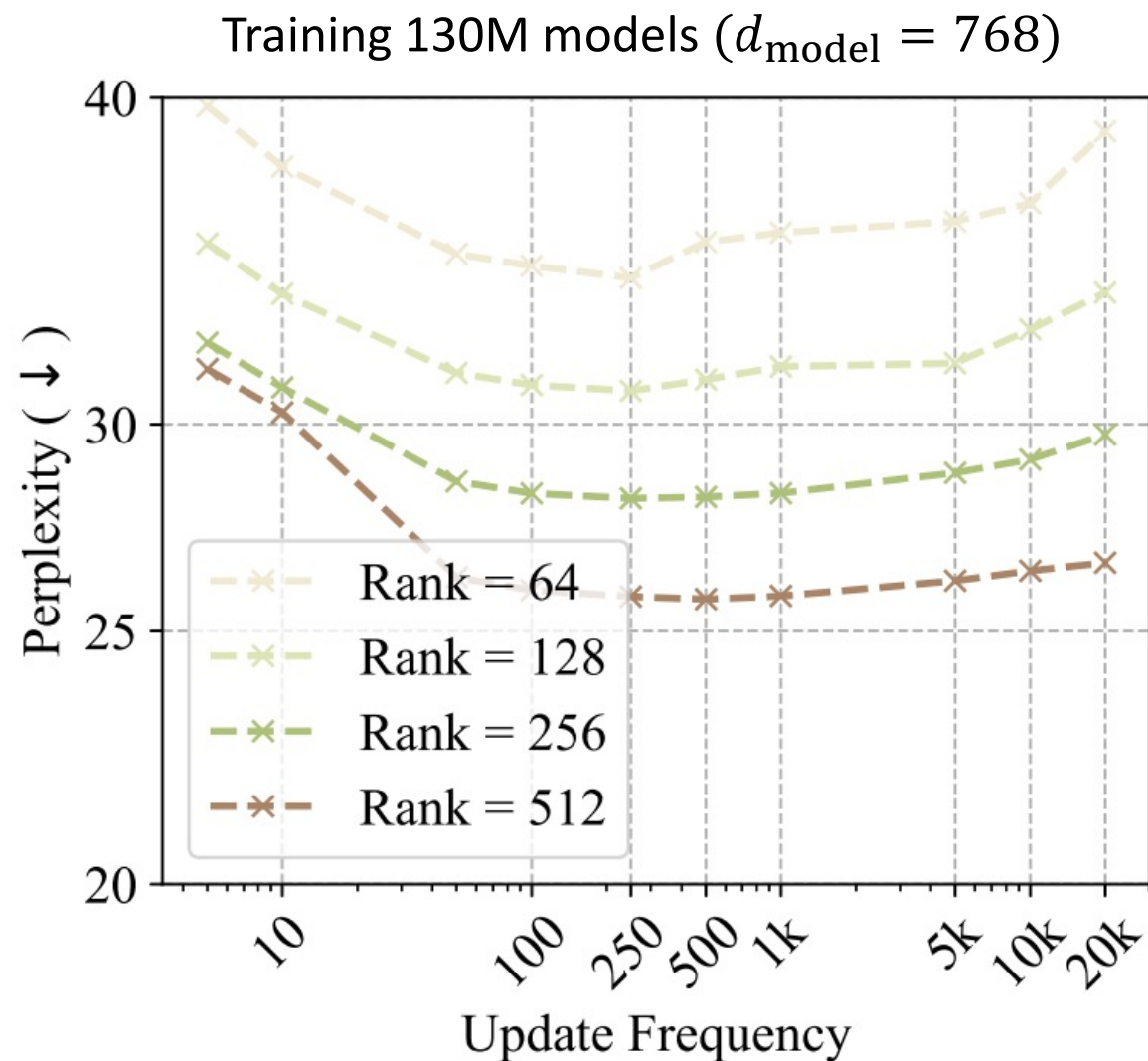**Does that mean it works?** No... $R_t \to 0$ just means the gradient within the subspace vanishes.

**How to continue optimization?** Change the projection from time to time!

# How often to change $P_t$?

$G_2$ $G_0$ $P$

$G_1$

$PP^T(G_t) \approx G_t$

Consecutive gradients are similar

**For every T iterations:**
Compute and store $P_t = \text{SVD}(G_t)$
$P_t$ is the projection matrix.

**No need to change $P_t$ every iteration!**

Training 130M models ($d_{\text{model}} = 768$)



- Rank = 64
- Rank = 128
- Rank = 256
- Rank = 512

Perplexity ( ↓ )

Update Frequency

# Pre-training Results (LLaMA 7B) on C4

**7B model trained on up to
150K steps and 19.7 B tokens**

|  | **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 |

**C4 Dataset** ⊕ **LLaMA-7B** ⊕ **single RTX 4090**

**Pre-training - for the first time!**

LLaMA-7B



— 8-bit AdamW
— 8-bit GaLore

Token Seen (Billions)

facebook Artificial Intelligence

# JoMA for Nonlinear activation

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

Nonlinear

Modified MLP (lower layer)

# JoMA for Nonlinear activation

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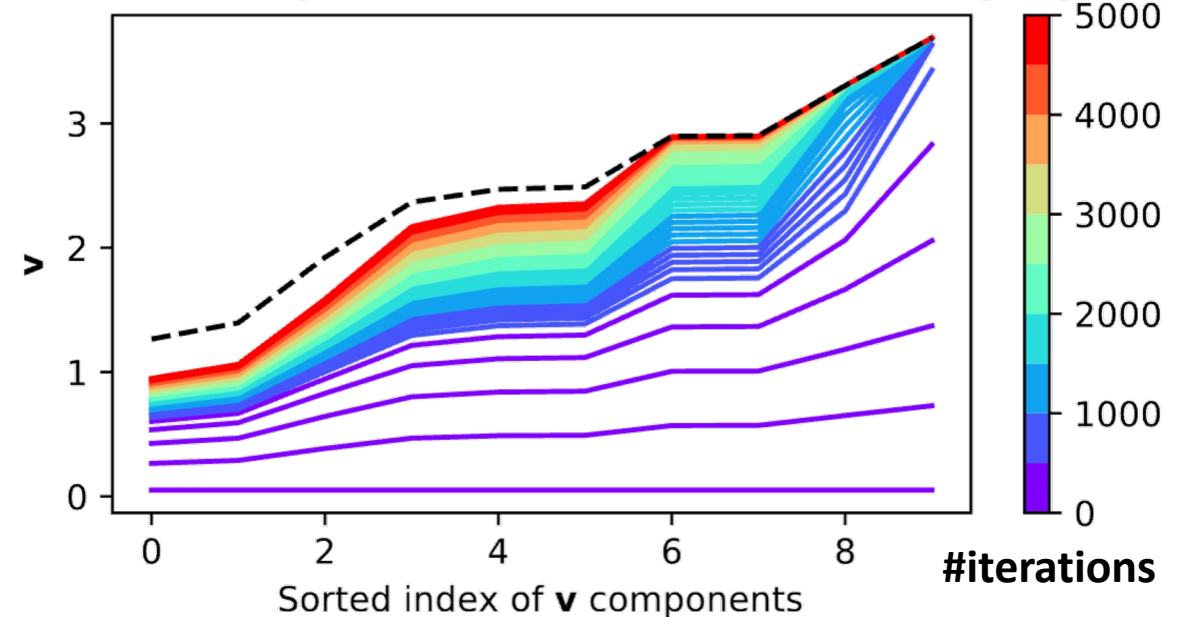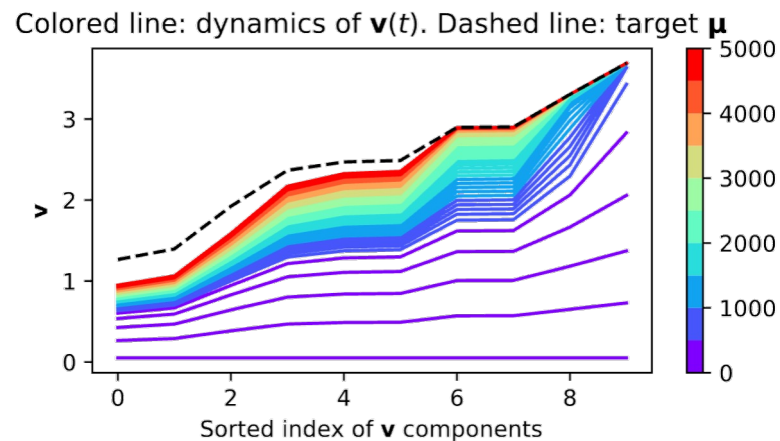
Nonlinear

Modified MLP (lower layer)

**Theorem 4**

Salient components grow much faster than non-salient ones:

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

$\text{ConvergenceRate}(j) := \ln 1/\delta_j(t)$

$\delta_j(t) := 1 - v_j(t)/\mu_j$

# JoMA for Nonlinear activation

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

Nonlinear

Modified
MLP
(lower layer)

**Theorem 4**

Salient components grow much faster than non-salient ones:

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

$$\text{ConvergenceRate}(j) := \ln 1/\delta_j(t)$$

$$\delta_j(t) := 1 - v_j(t)/\mu_j$$



Colored line: dynamics of $\mathbf{v}(t)$. Dashed line: target $\mu$

Sorted index of $\mathbf{v}$ components

#iterations

# JoMA for Nonlinear activation

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

Nonlinear

Modified MLP (lower layer)
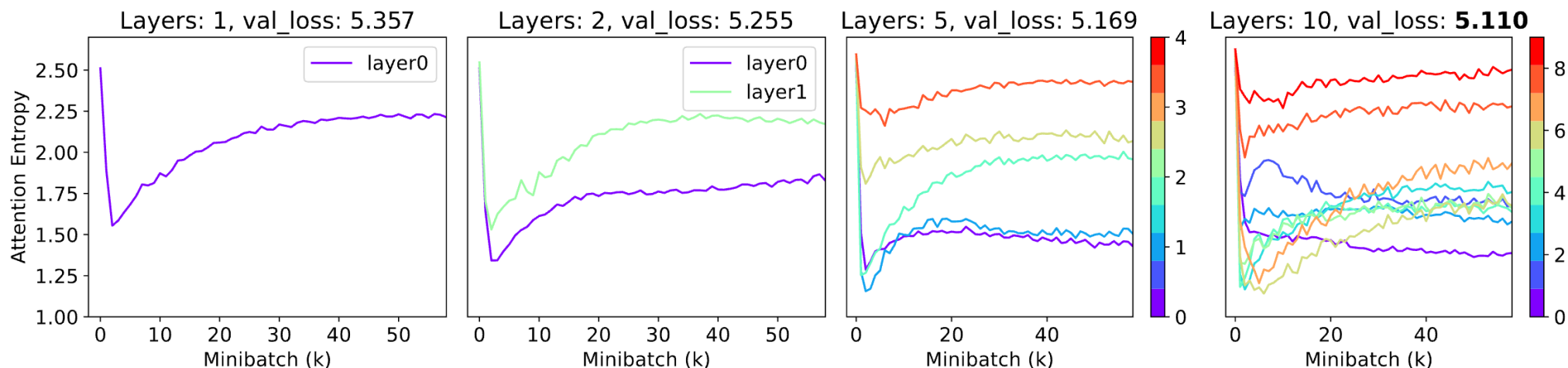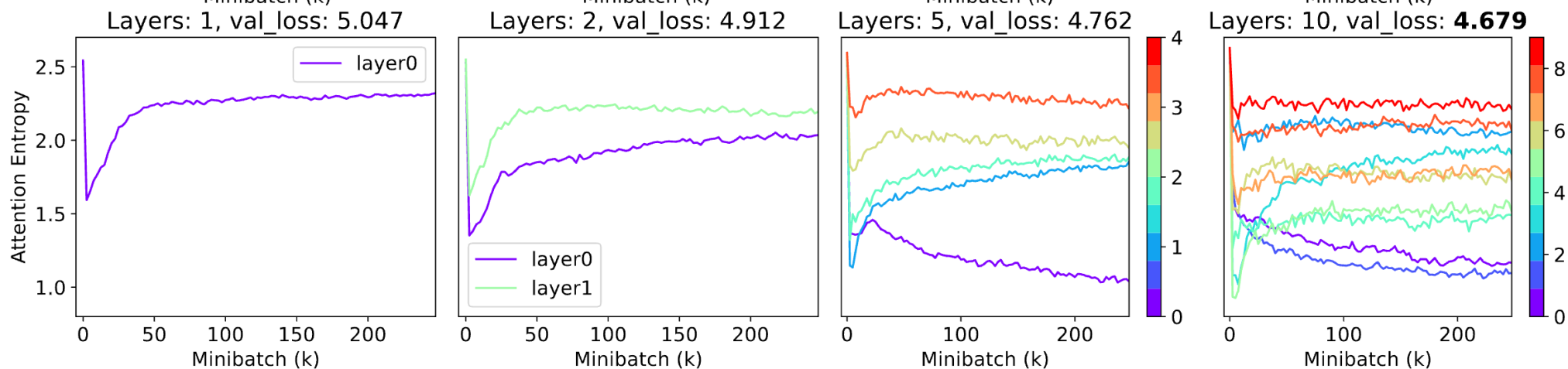
How the entropy of attention changes over time?



Colored line: dynamics of **v**(t). Dashed line: target **μ**

Entropy changes over time

Attention becomes **sparser** and then **denser**!

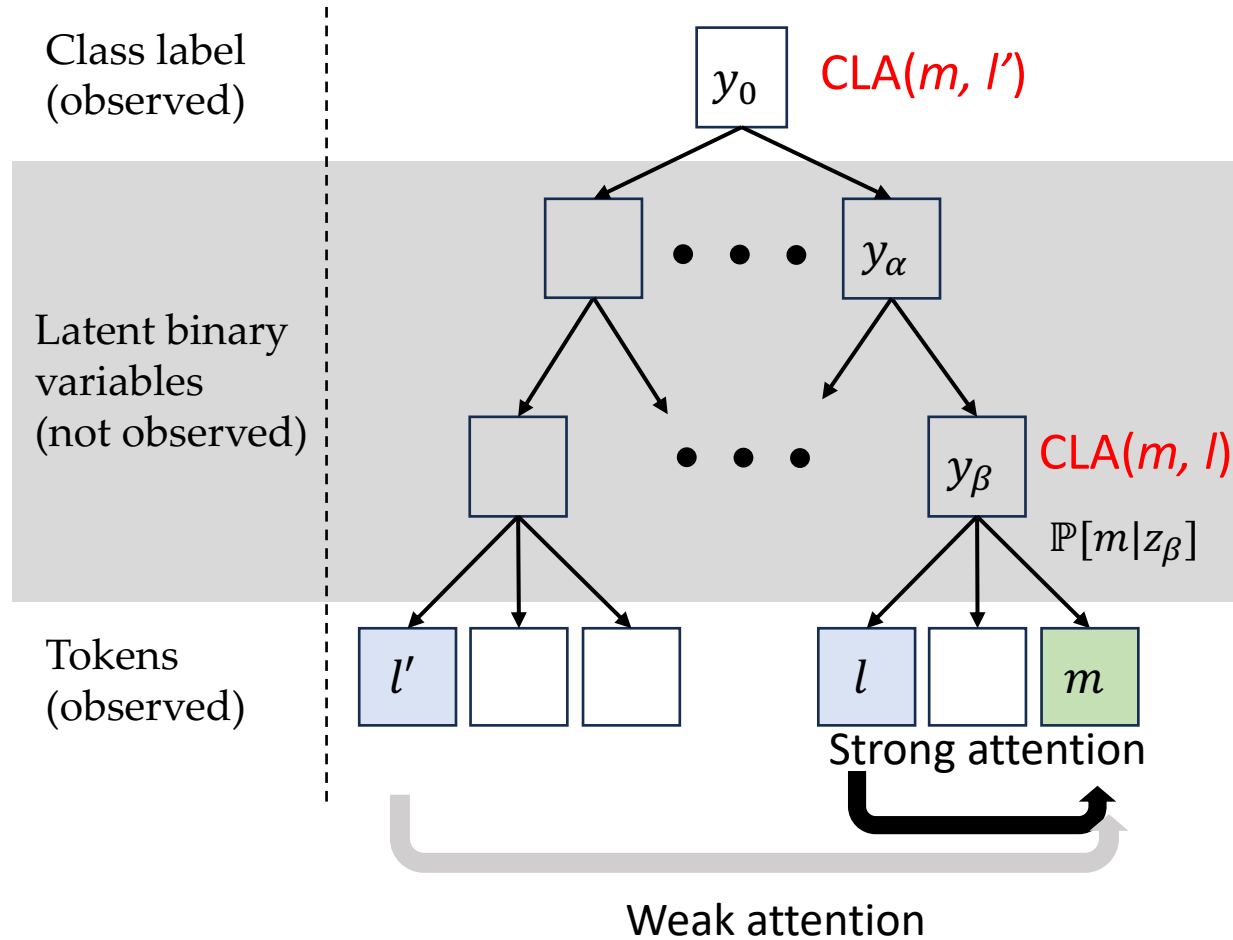"bounce back"

# 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!

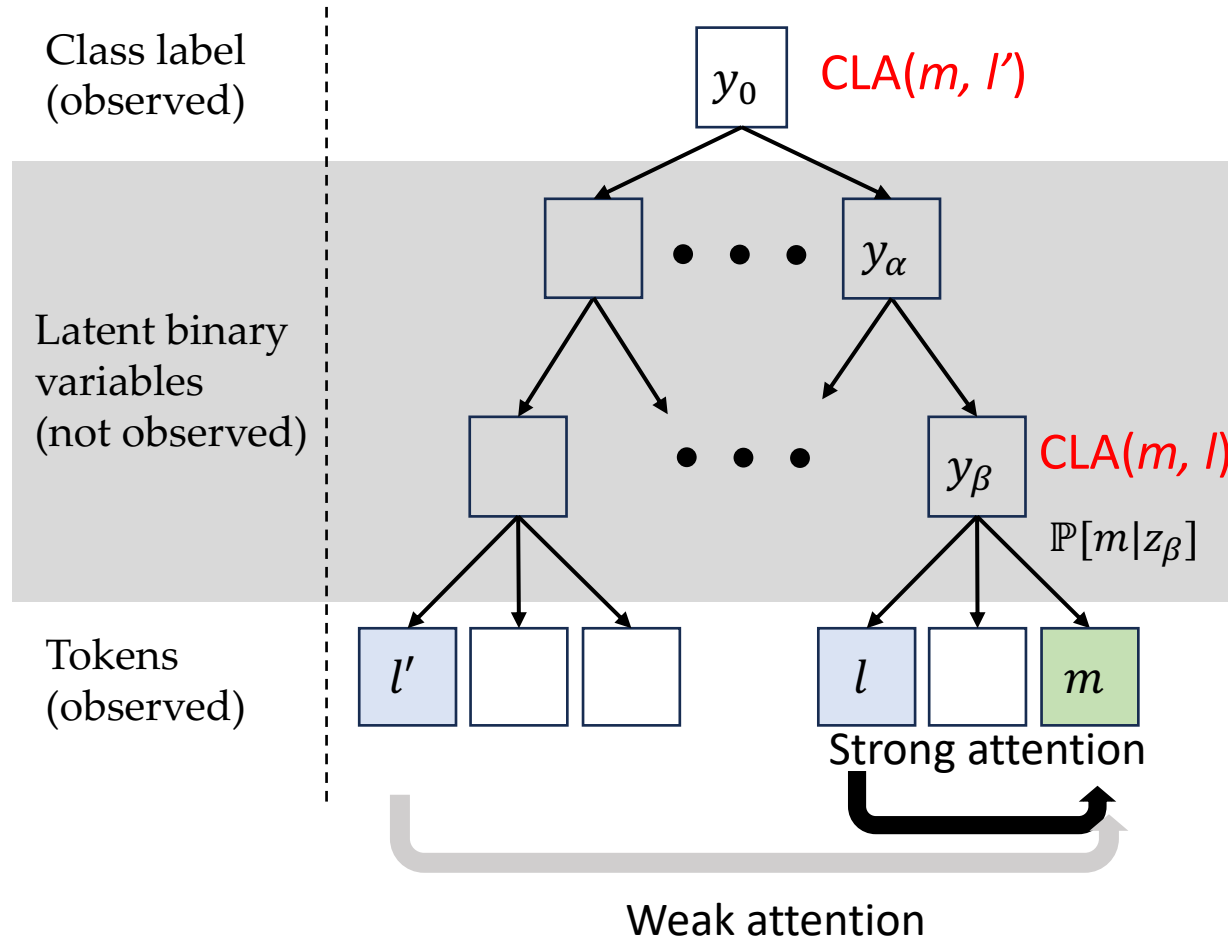# Data Hierarchy & Multilayer Transformer



Class label (observed)

$y_0$  CLA($m, l'$)

Latent binary variables (not observed)

$y_\alpha$

$y_\beta$  CLA($m, l$)

$\mathbb{P}[m|z_\beta]$

Tokens (observed)

$l'$

$l$  $m$

Strong attention

Weak attention

# Data Hierarchy & Multilayer Transformer

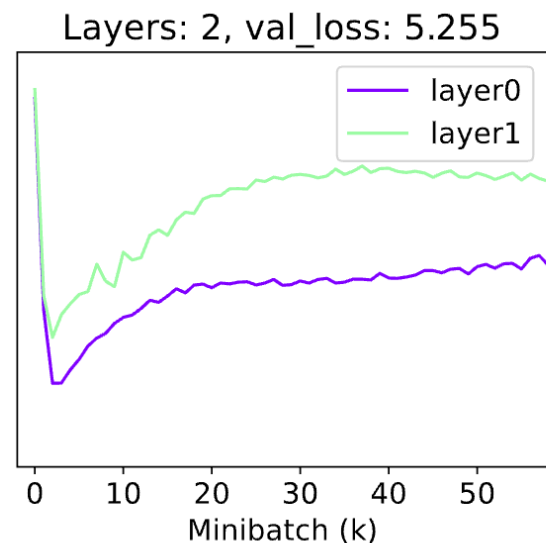# Deep Latent Distribution



Learning the current hierarchical structure by
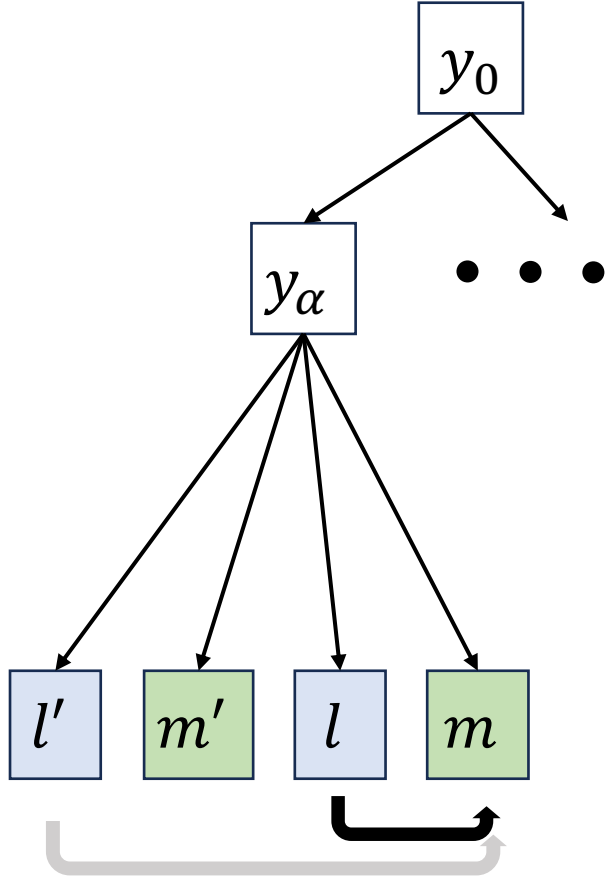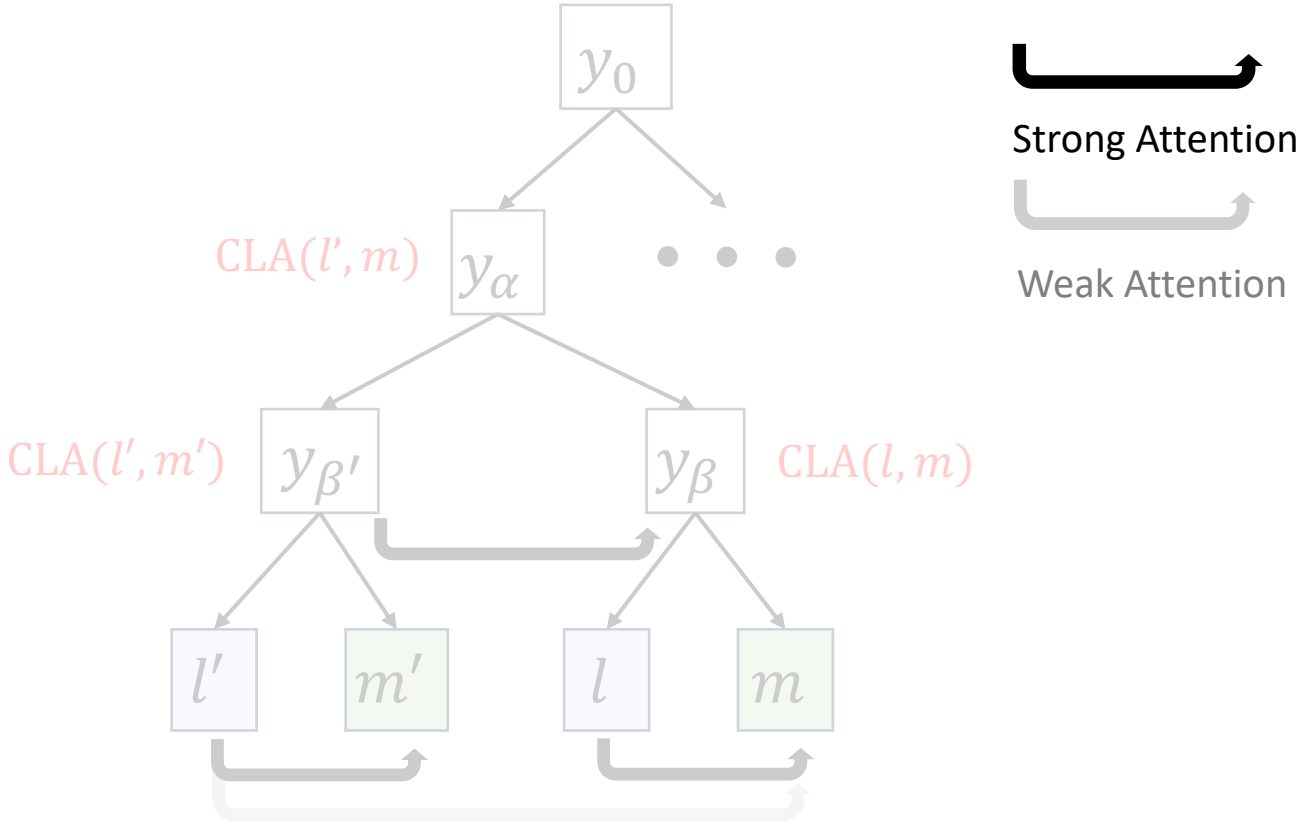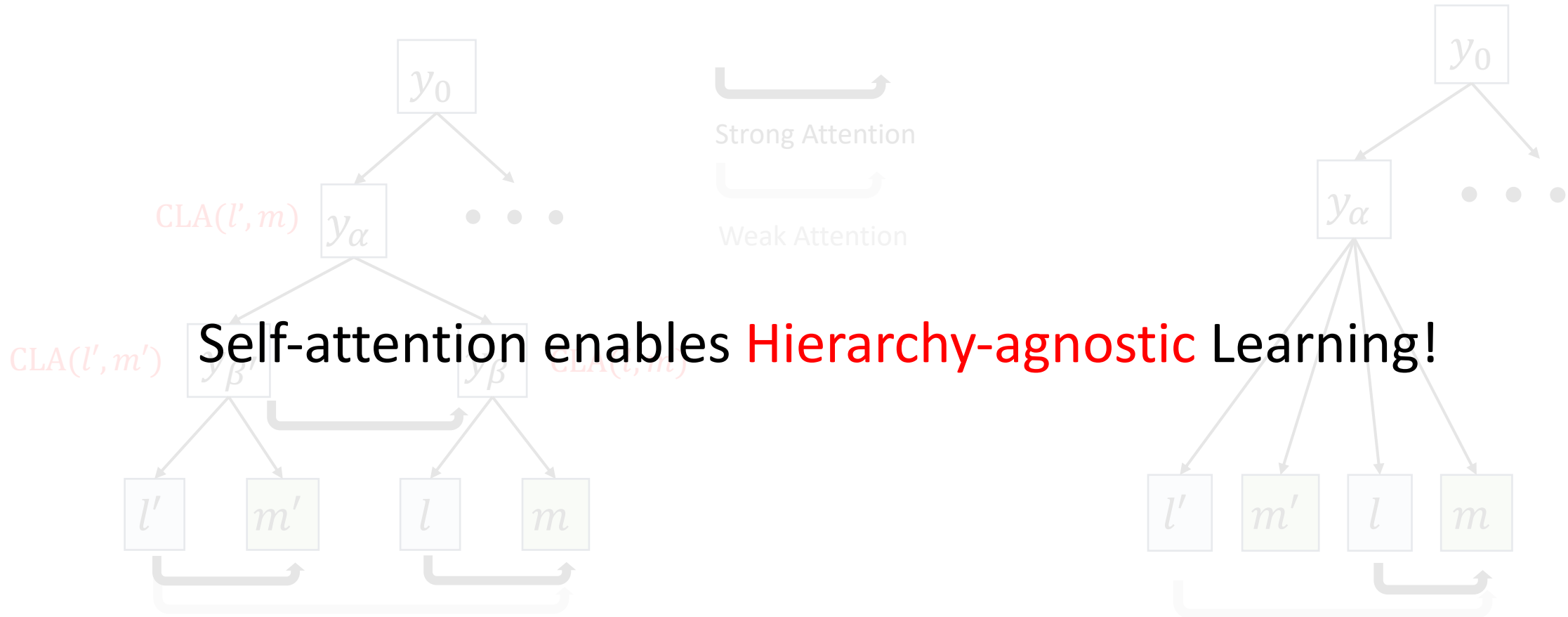*slowing down* the association of tokens that are not directly correlated

# Shallow Latent Distribution
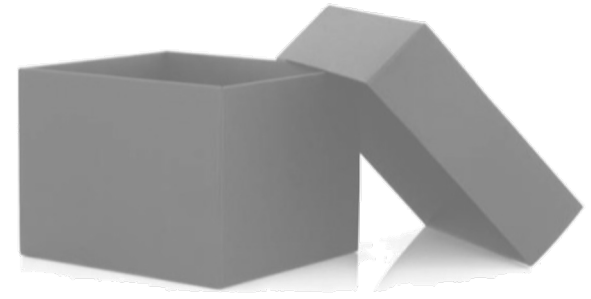
# Hierarchy-agnostic Learning



Self-attention enables Hierarchy-agnostic Learning!

# Verification of Hierarchical Intuitions

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

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.

# 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!