

JoMA: Demystifying Multilayer Transformers via JOint Dynamics of MLP and Attention

Yuandong Tian¹, Yiping Wang², Zhenyu Zhang³, Beidi Chen^{1,4}, Simon Du²

¹Meta AI (FAIR)

²University of Washington

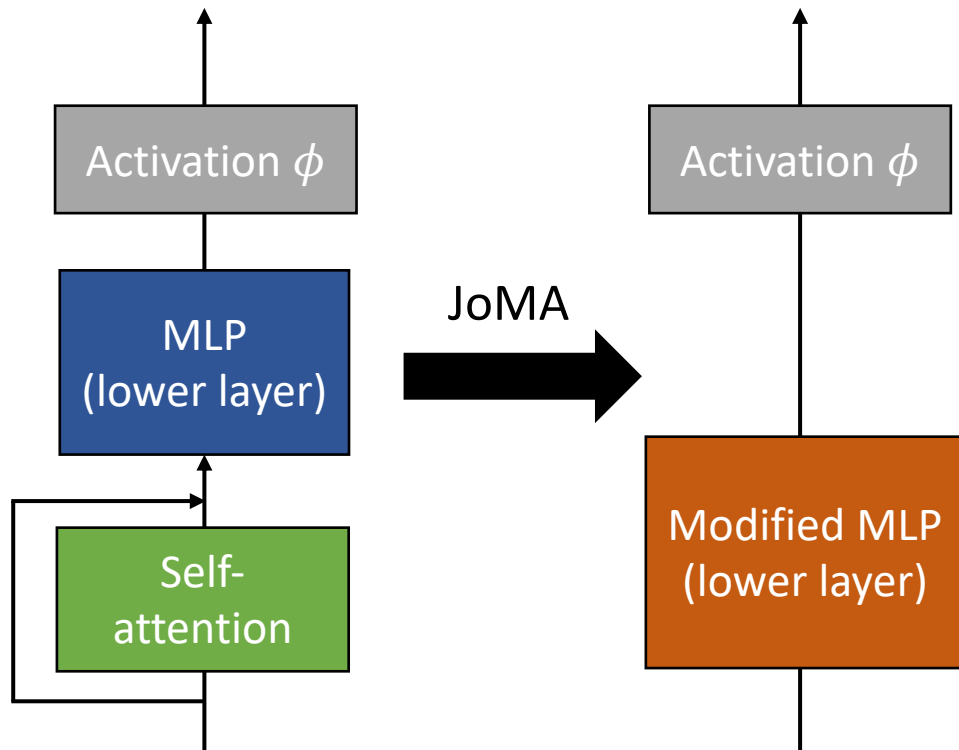
³University of Texas, Austin

⁴Carnegie Mellon University



Published in International Conference in Learning Representation (ICLR) 2024

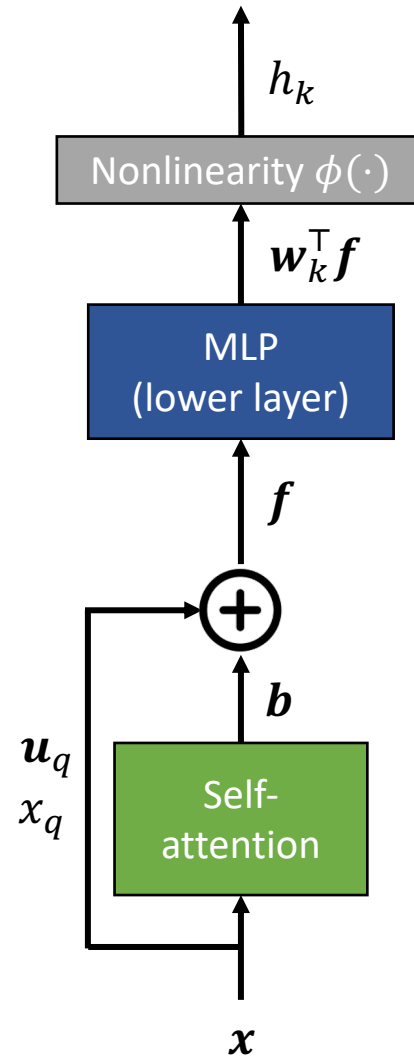
JoMA: JOint Dynamics of MLP/Attention layers



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

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

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

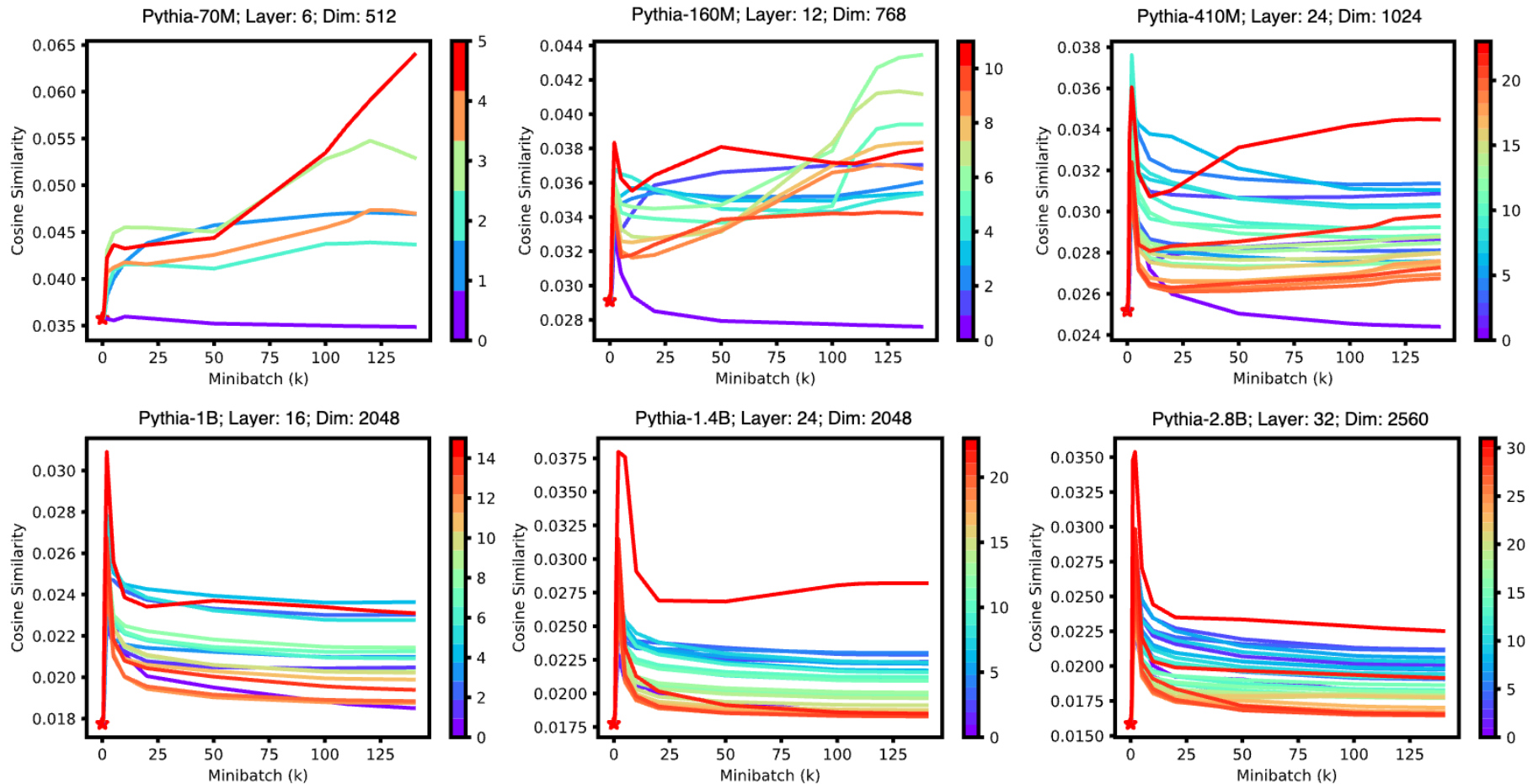
$$\text{ExpAttn: } b_l = x_l e^{z_{ql}}$$

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

"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 $\mathbf{z}_m(t) = \frac{1}{2} \sum_k \mathbf{v}_k^2(t) - \|\mathbf{v}_k(t)\|_2^2 \bar{\mathbf{b}}_m + \mathbf{c}$.

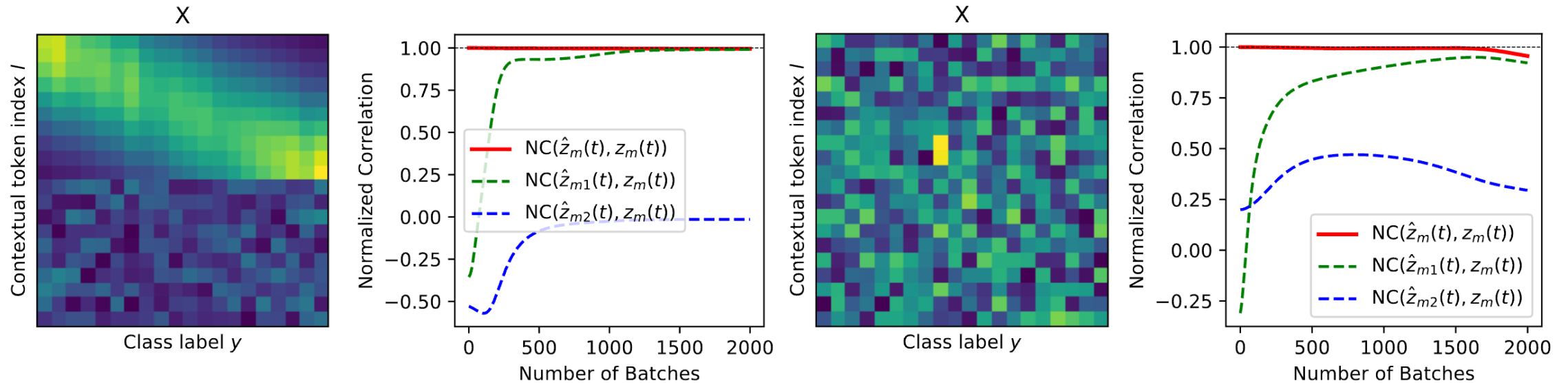
Under zero-initialization ($\mathbf{w}_k(0) = 0, \mathbf{z}_m(0) = 0$), then the time-independent constant $\mathbf{c} = 0$.

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

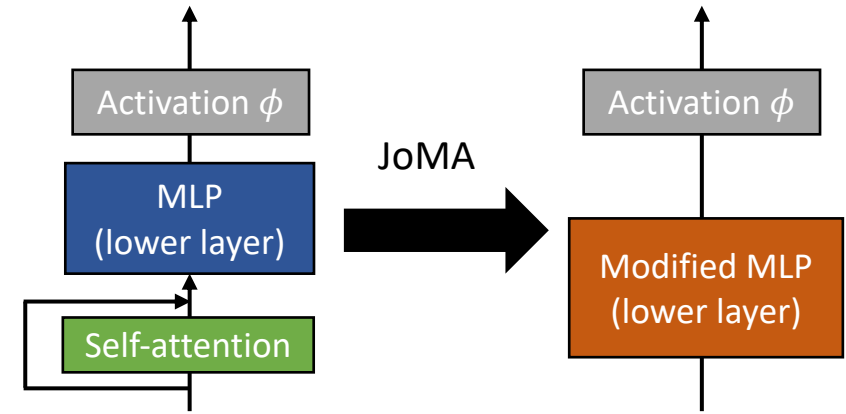
$\hat{\mathbf{z}}_m(t)$: Estimated attention logits by JoMA

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

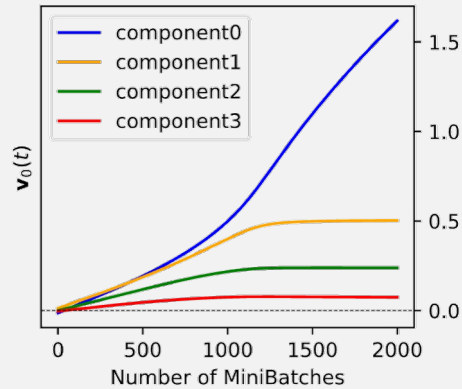
Implication of Theorem 1

Key idea: folding self-attention into MLP

→ A Transformer block becomes a modified MLP

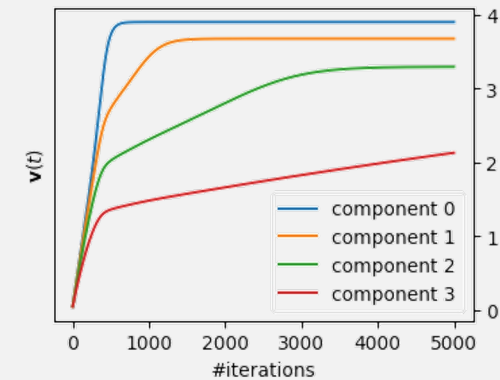


Linear case ($\phi = \text{Id}, K = 1$)



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)

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

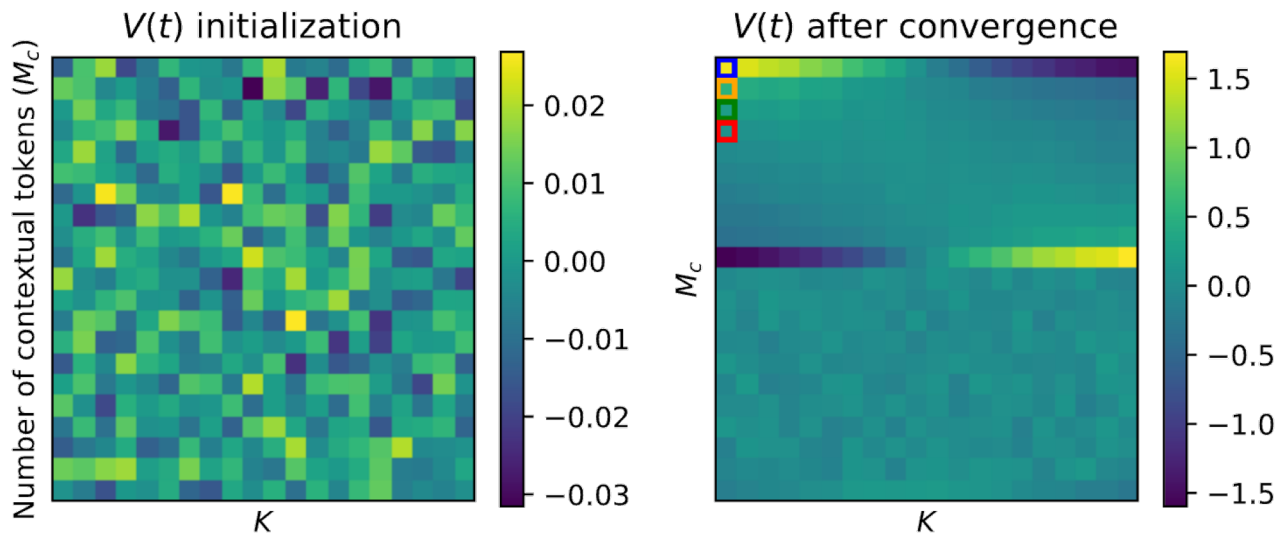
Theorem 2

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

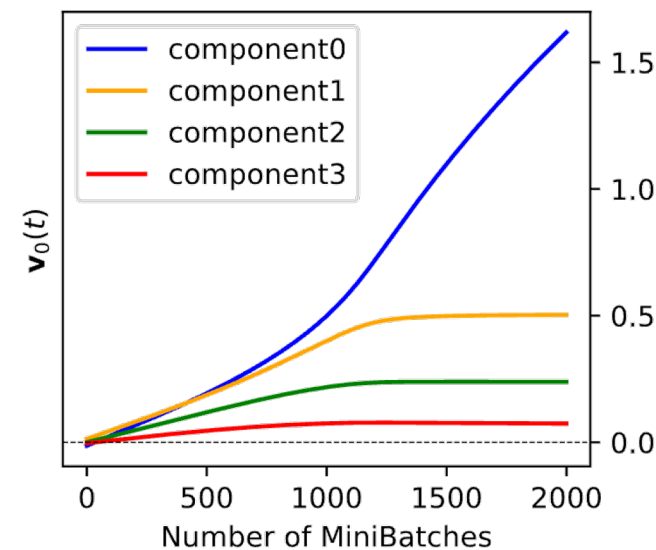
$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \in [-1, 1]$$

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)



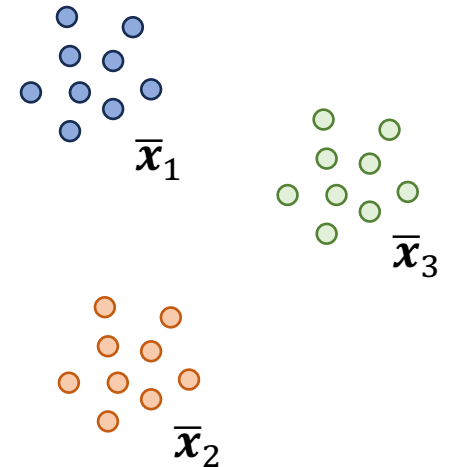
JoMA for Nonlinear Activation

Theorem 3

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

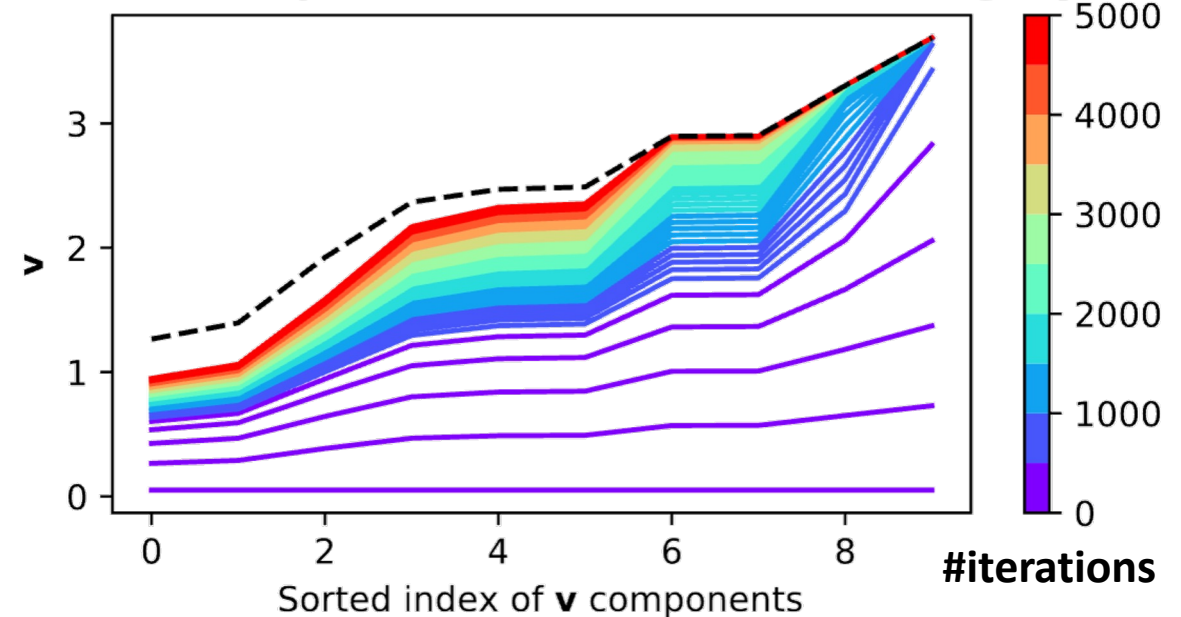
Theorem 4

Salient components grow much faster than non-salient ones:

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

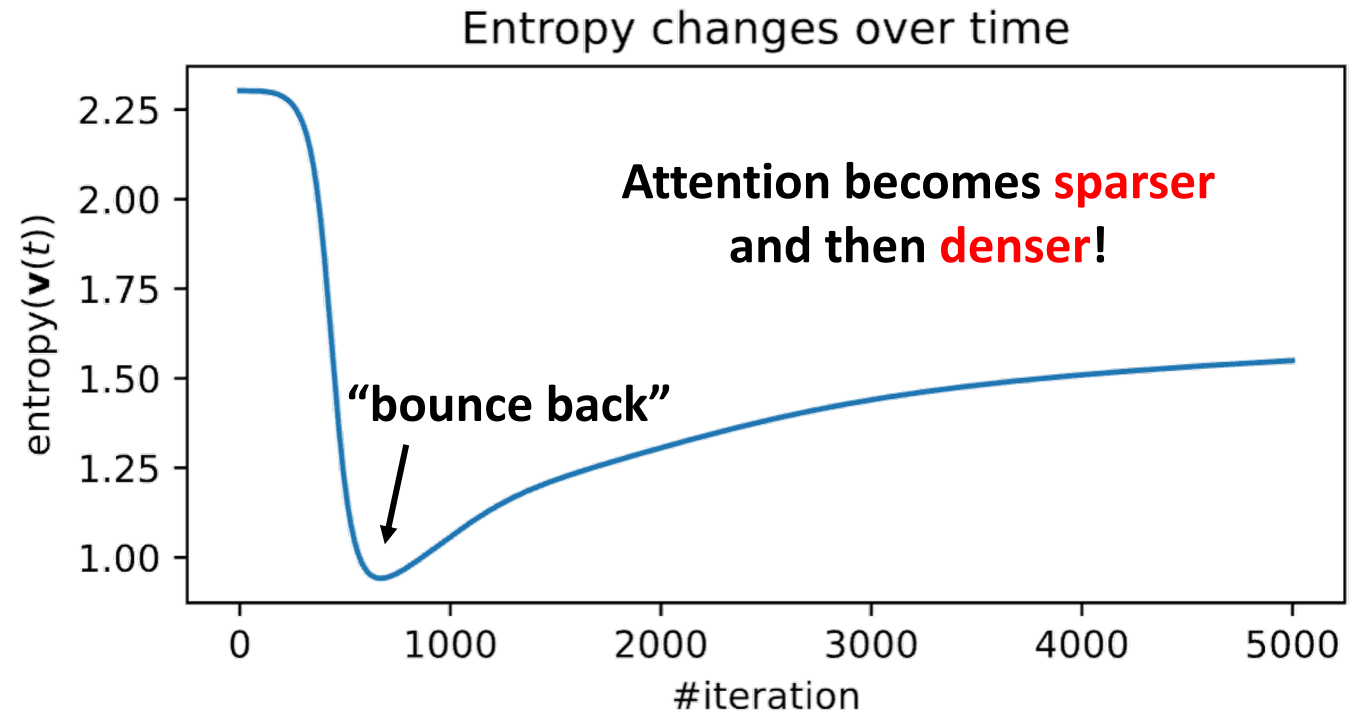
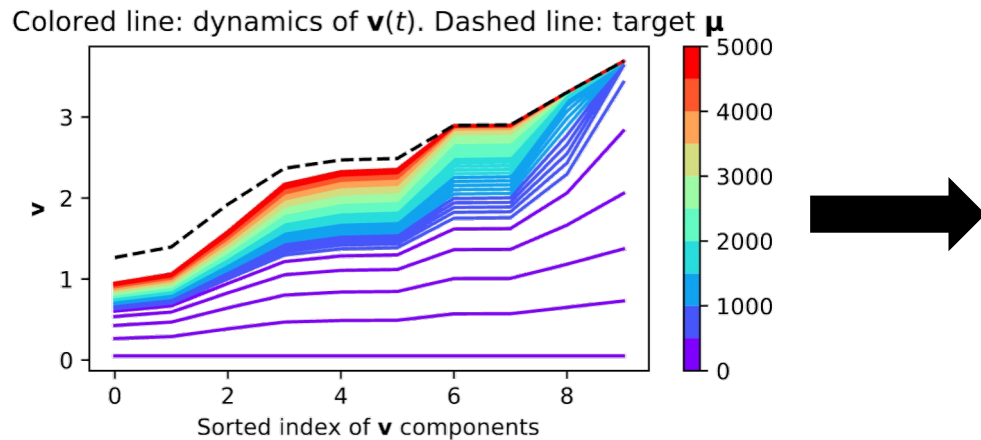
$$\begin{aligned} \text{ConvergenceRate}(j) &:= \ln 1/\delta_j(t) \\ \delta_j(t) &:= 1 - v_j(t)/\mu_j \end{aligned}$$

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



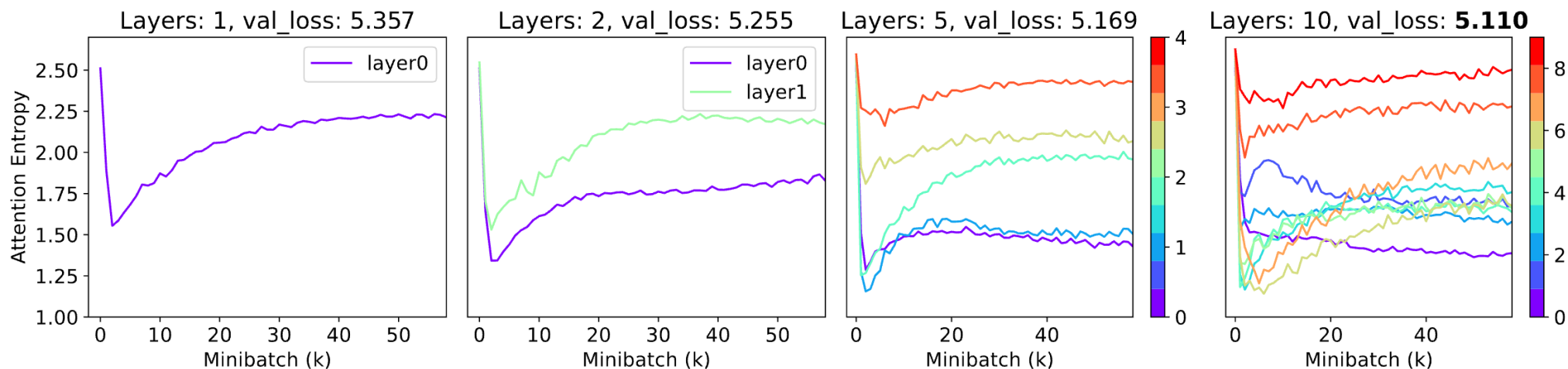
JoMA for Nonlinear activation

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

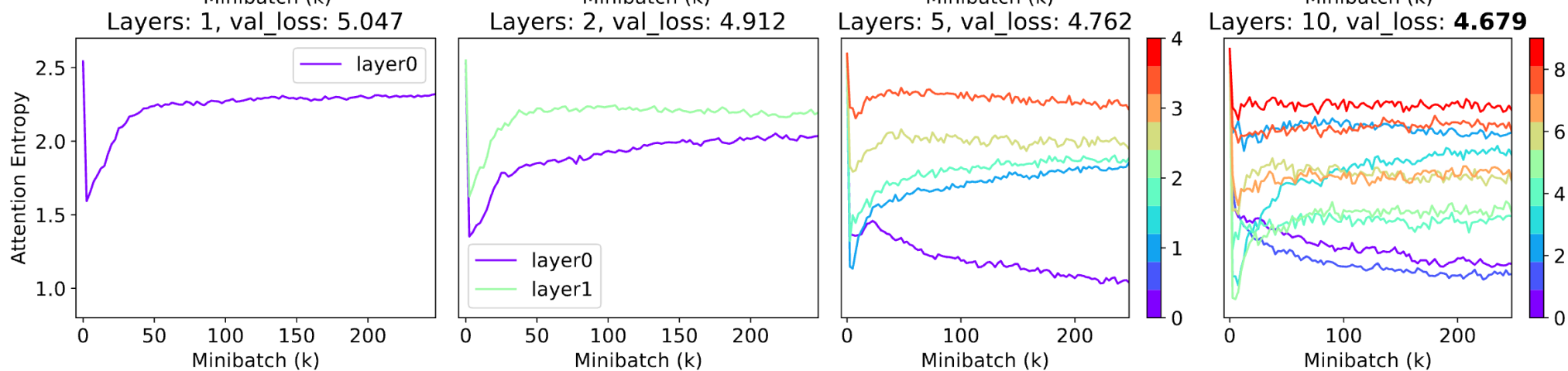


Real-world Experiments

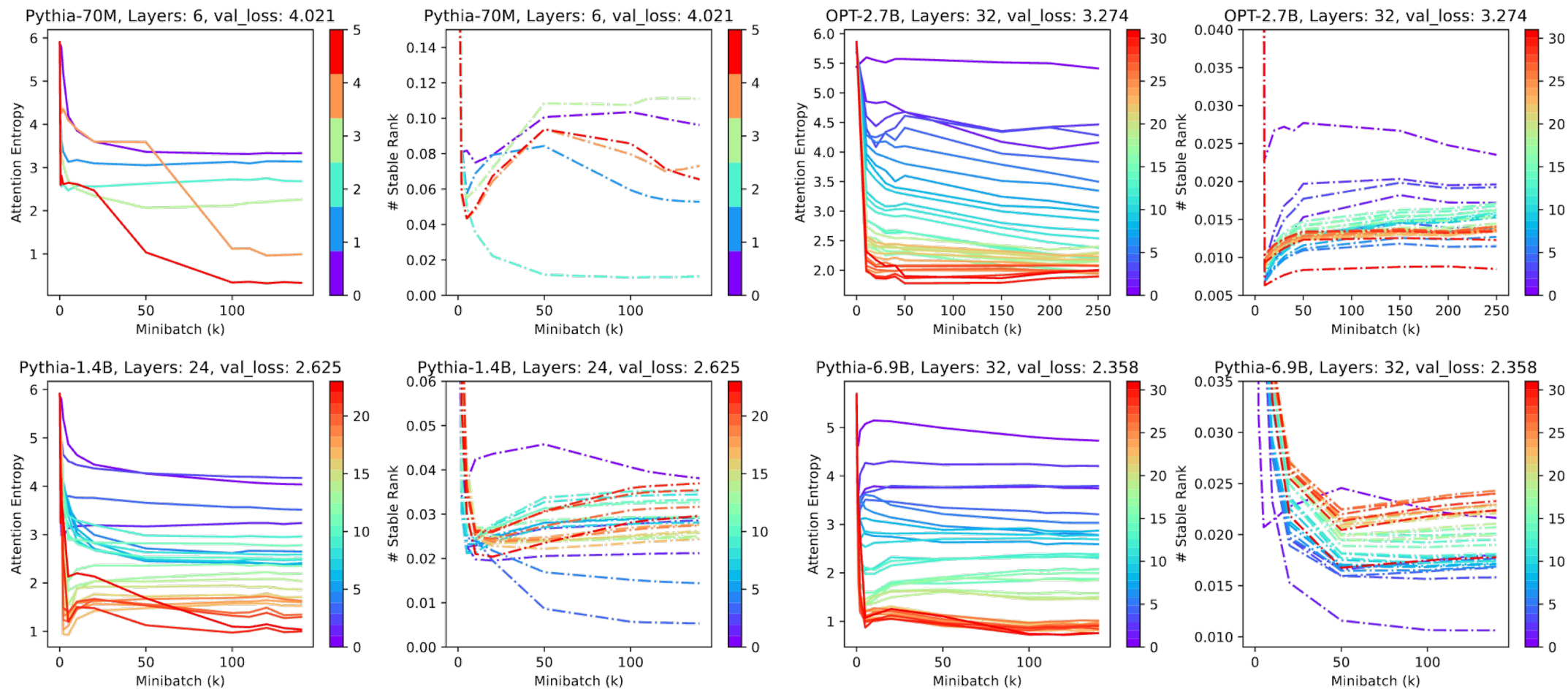
Wikitext2



Wikitext103



Real-world Experiments

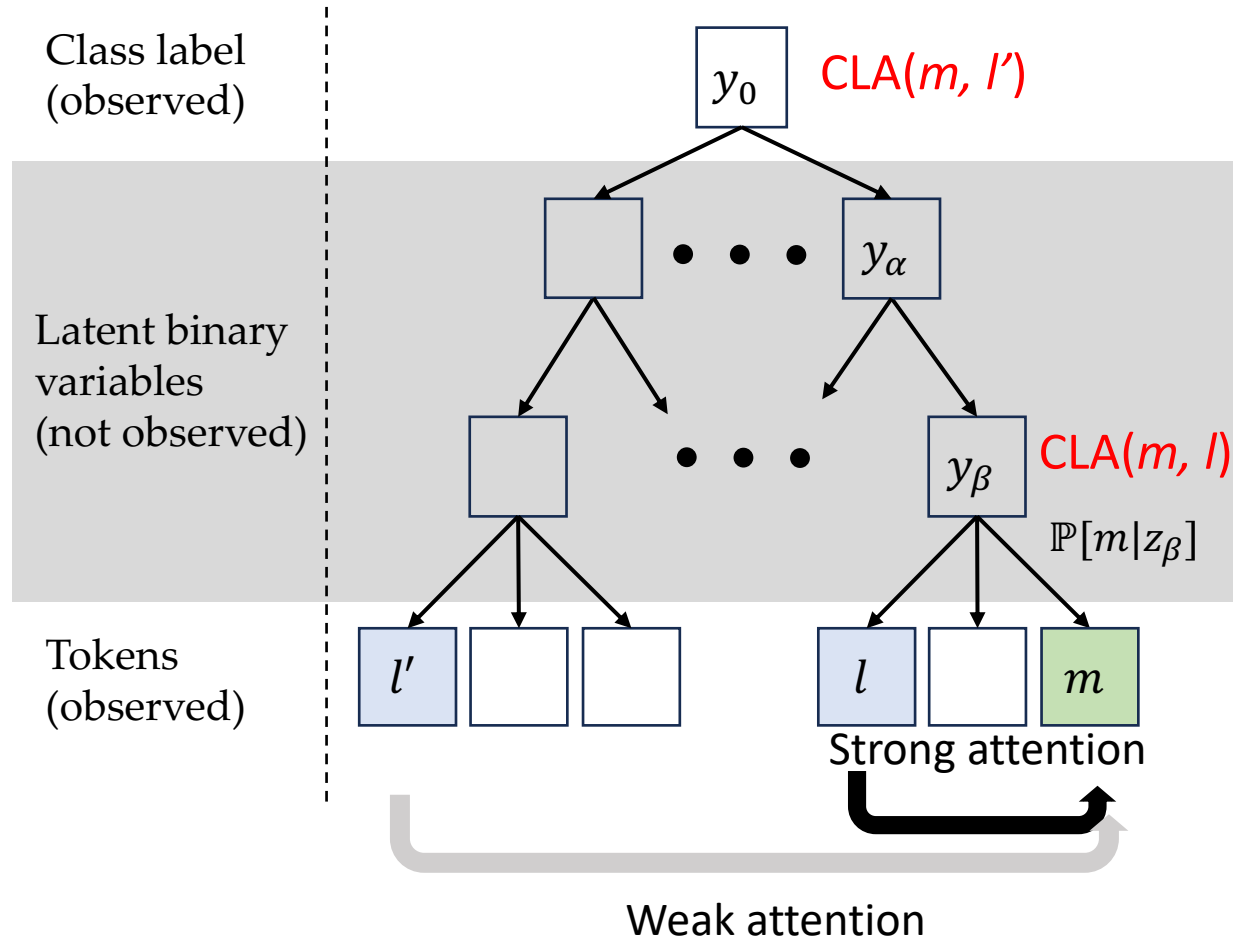


Why is this “bouncing back” property useful?

It seems that it only slows down the training??

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

Data Hierarchy & Multilayer Transformer



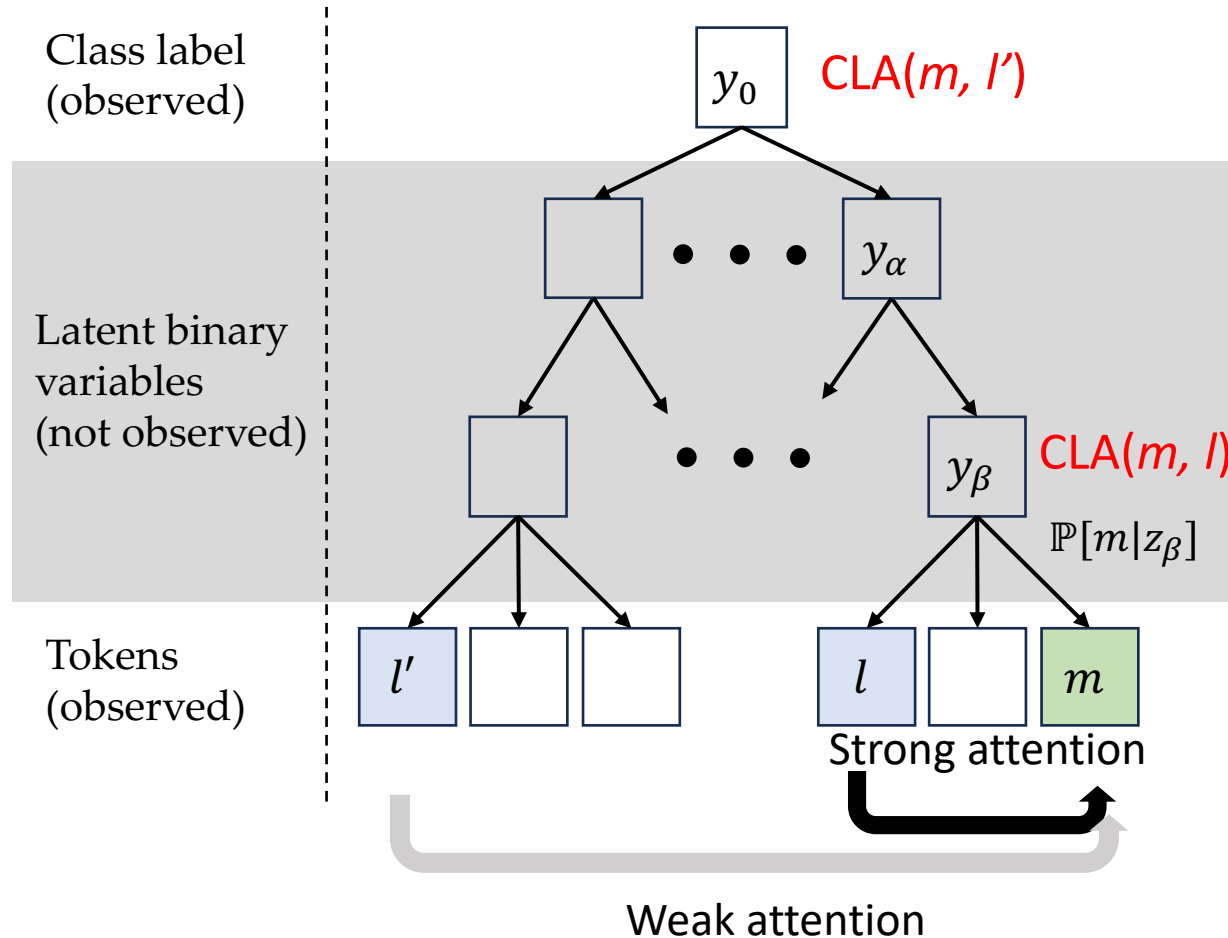
Data Hierarchy & Multilayer Transformer

Theorem 5

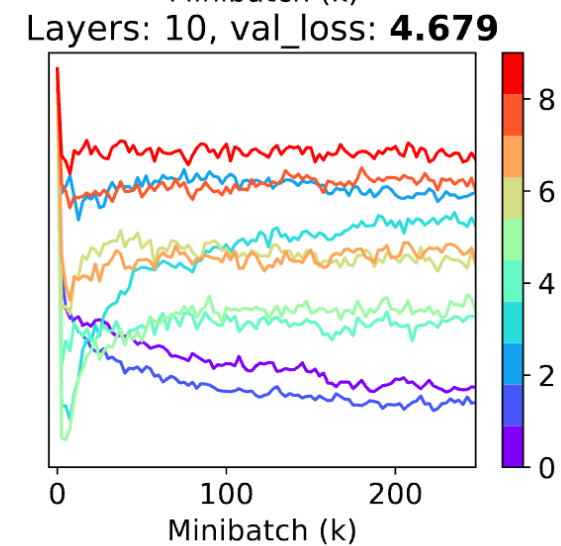
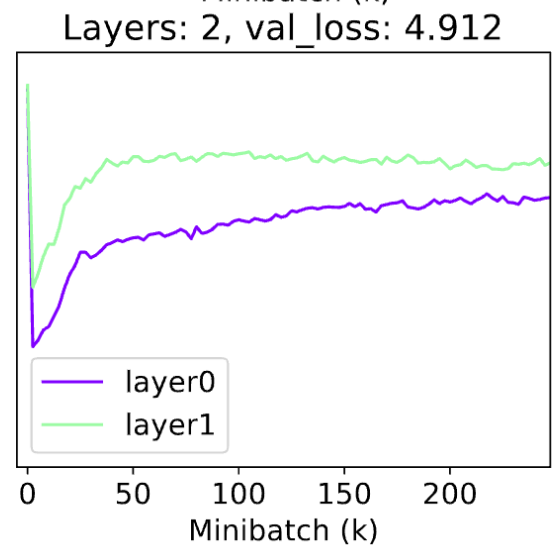
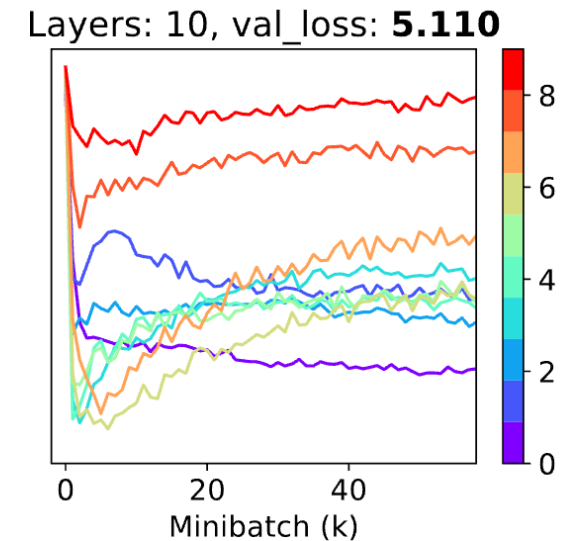
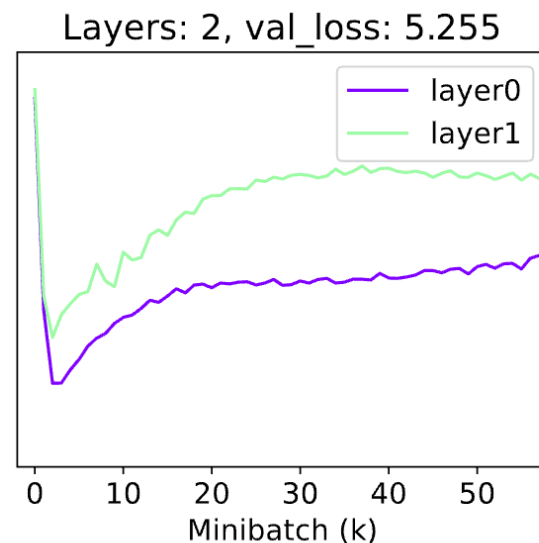
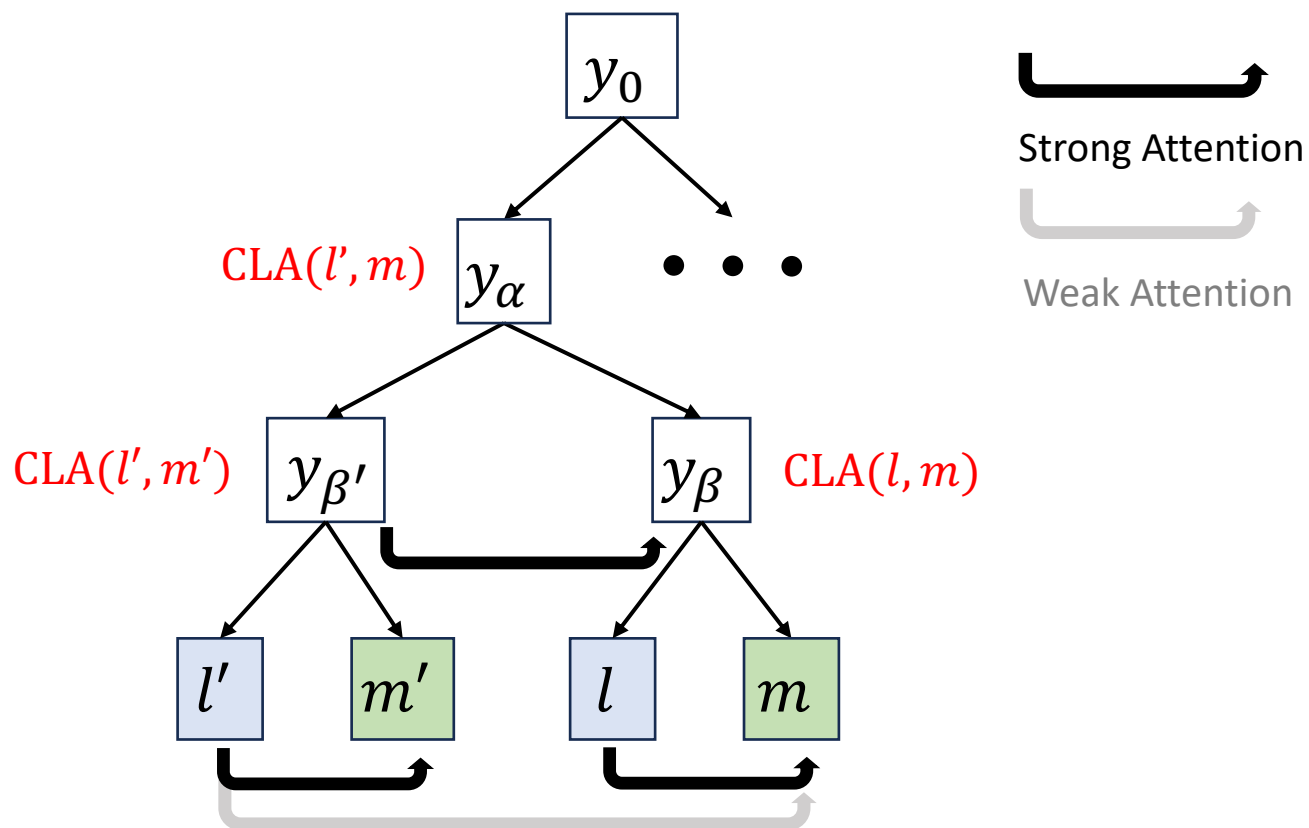
$$\mathbb{P}[l|m] \approx 1 - \frac{H}{L}$$

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

L : total height of the hierarchy



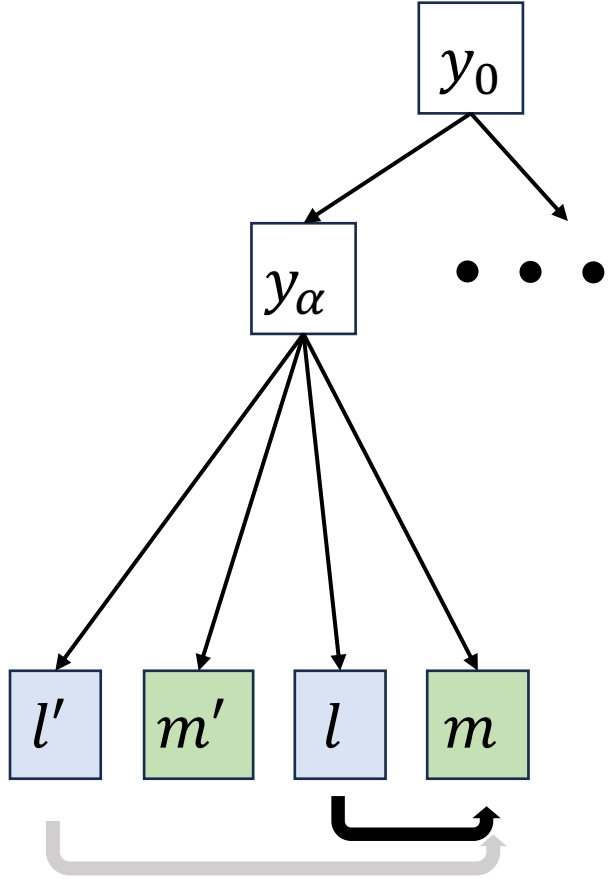
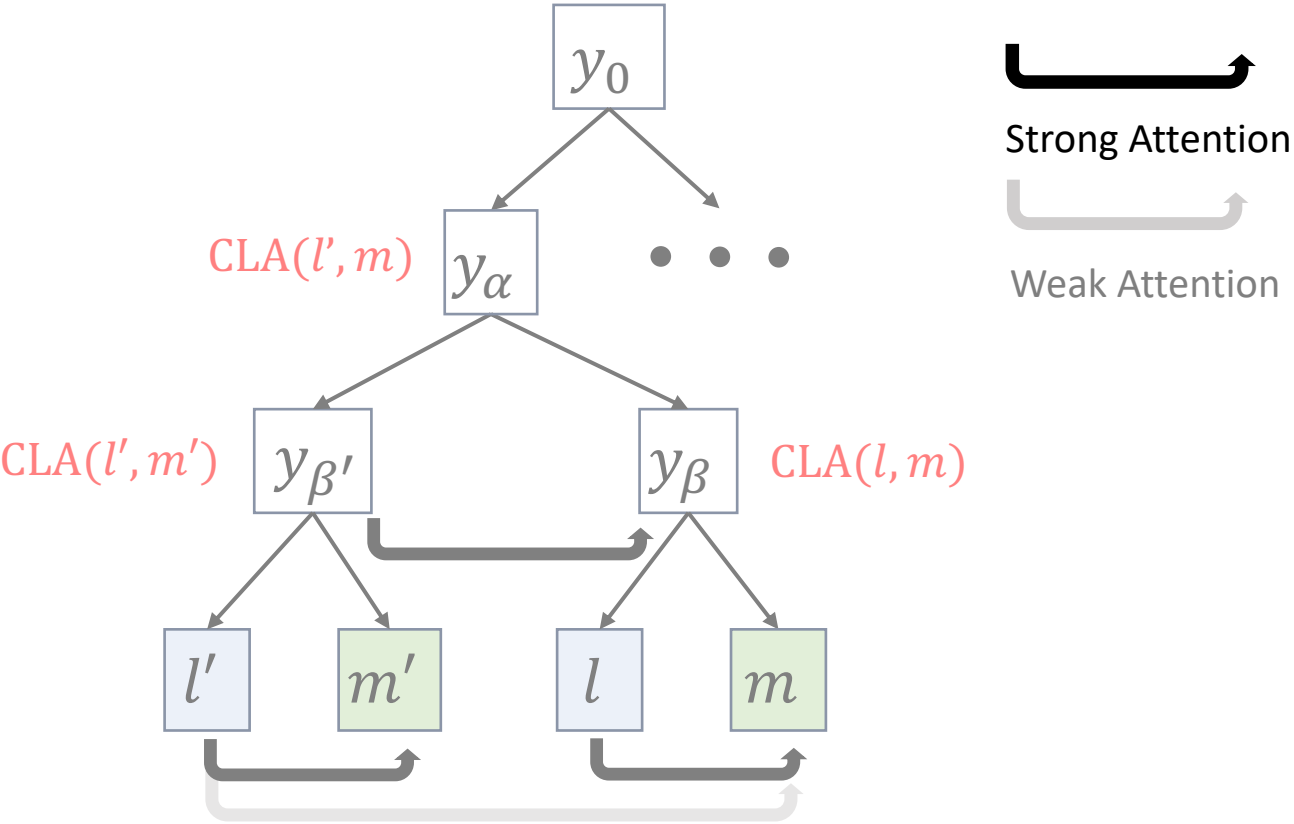
Deep Latent Distribution



Learning the current hierarchical structure by

slowing down the association of tokens that are not directly correlated

Shallow Latent Distribution



Future Work

- How embedding vectors are learned?
 - In both Scan&Snap and JoMA, we assume embeddings are constant.
- Positional Encoding
- Formulate the dynamics of Multi-layer Transformers
 - How intermediate latent concept gets learned during training?
 - Why we need over-parameterization?

Thanks!