

Self-attention Networks Localize When QK-eigenspectrum Concentrates

June 27 (Thu.), 2024

Han Bao

Great collaborators!

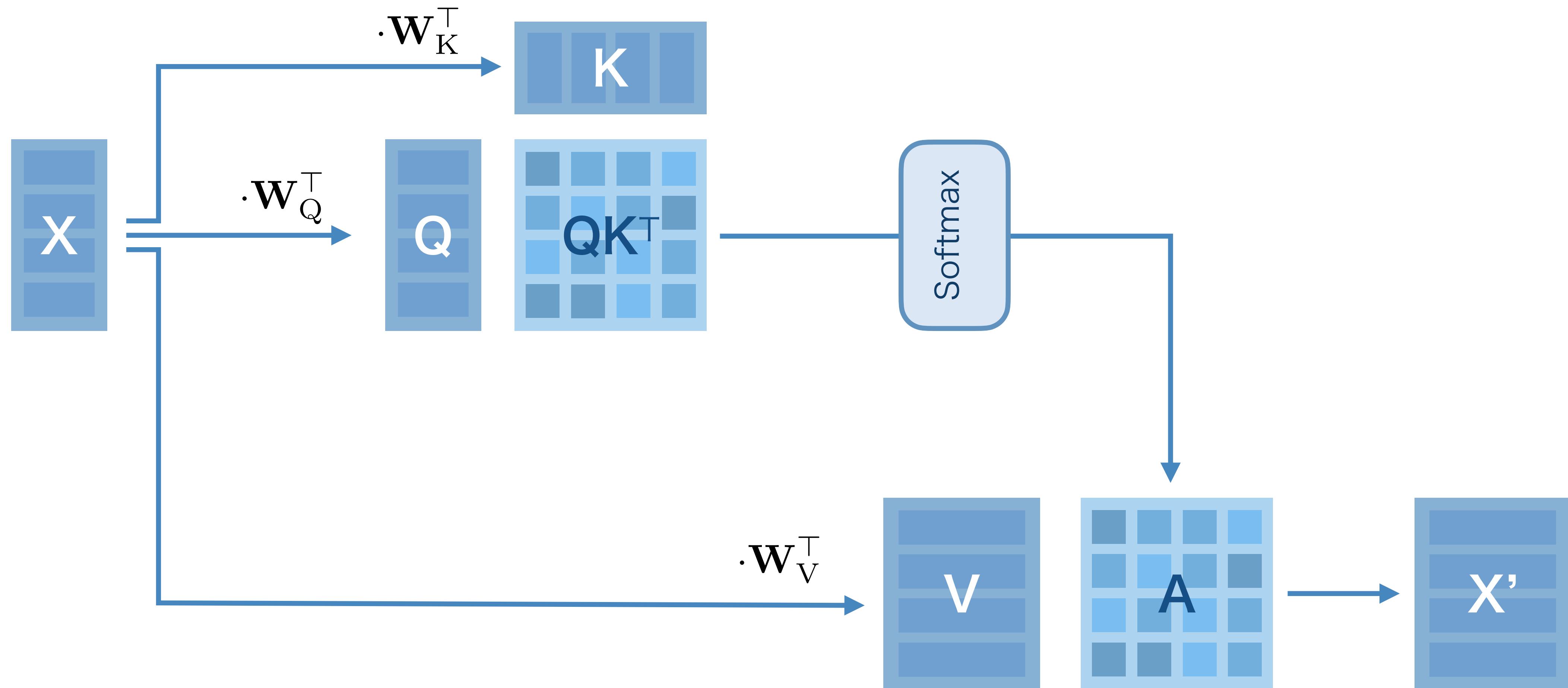


Ryuichiro Hataya (RIKEN)



Ryo Karakida (AIST)

Attention is prevailing ...



Attention as a visualizer

[Vaswani+ 2017]

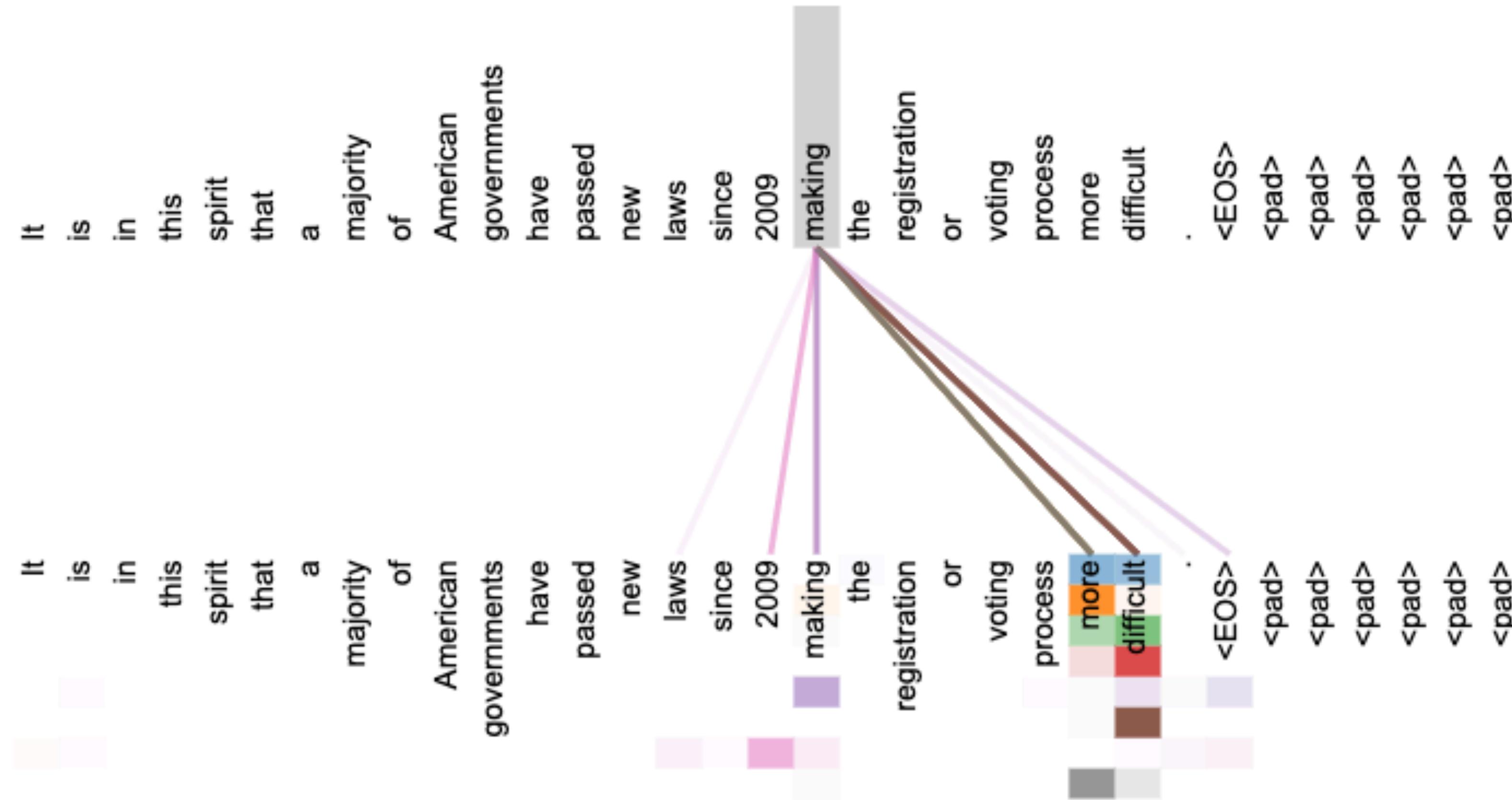


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb ‘making’, completing the phrase ‘making...more difficult’. Attentions here shown only for the word ‘making’. Different colors represent different heads. Best viewed in color.

Attention as an aligner

[Bahdanau+ 2015]

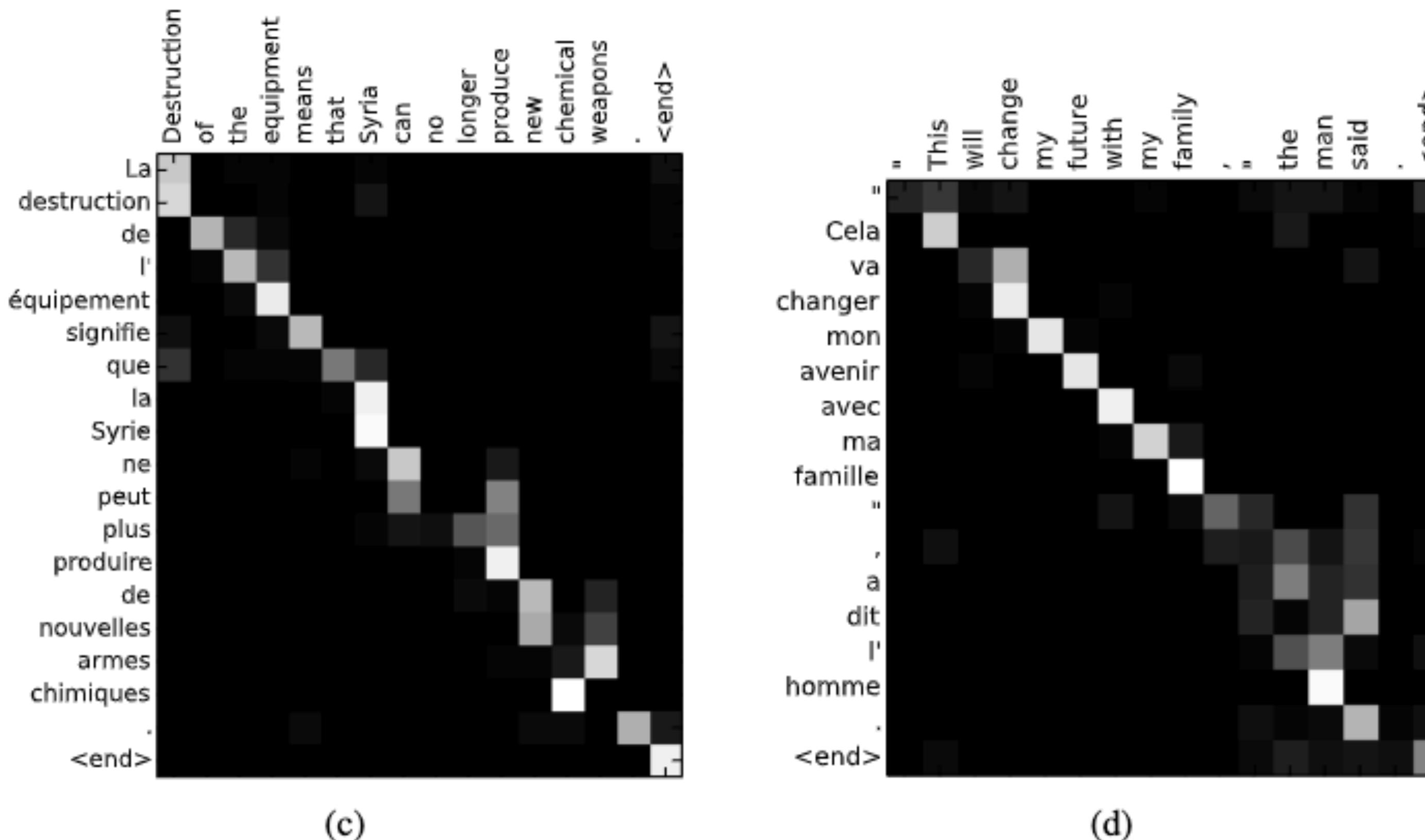


Figure 3: Four sample alignments found by RNNsearch-50. The x-axis and y-axis of each plot correspond to the words in the source sentence (English) and the generated translation (French), respectively. Each pixel shows the weight α_{ij} of the annotation of the j -th source word for the i -th target word (see Eq. (6)), in grayscale (0: black, 1: white). (a) an arbitrary sentence. (b–d) three randomly selected samples among the sentences without any unknown words and of length between 10 and 20 words from the test set.

Attention as a permutation-invariant learner

[Veličković+ 2018]

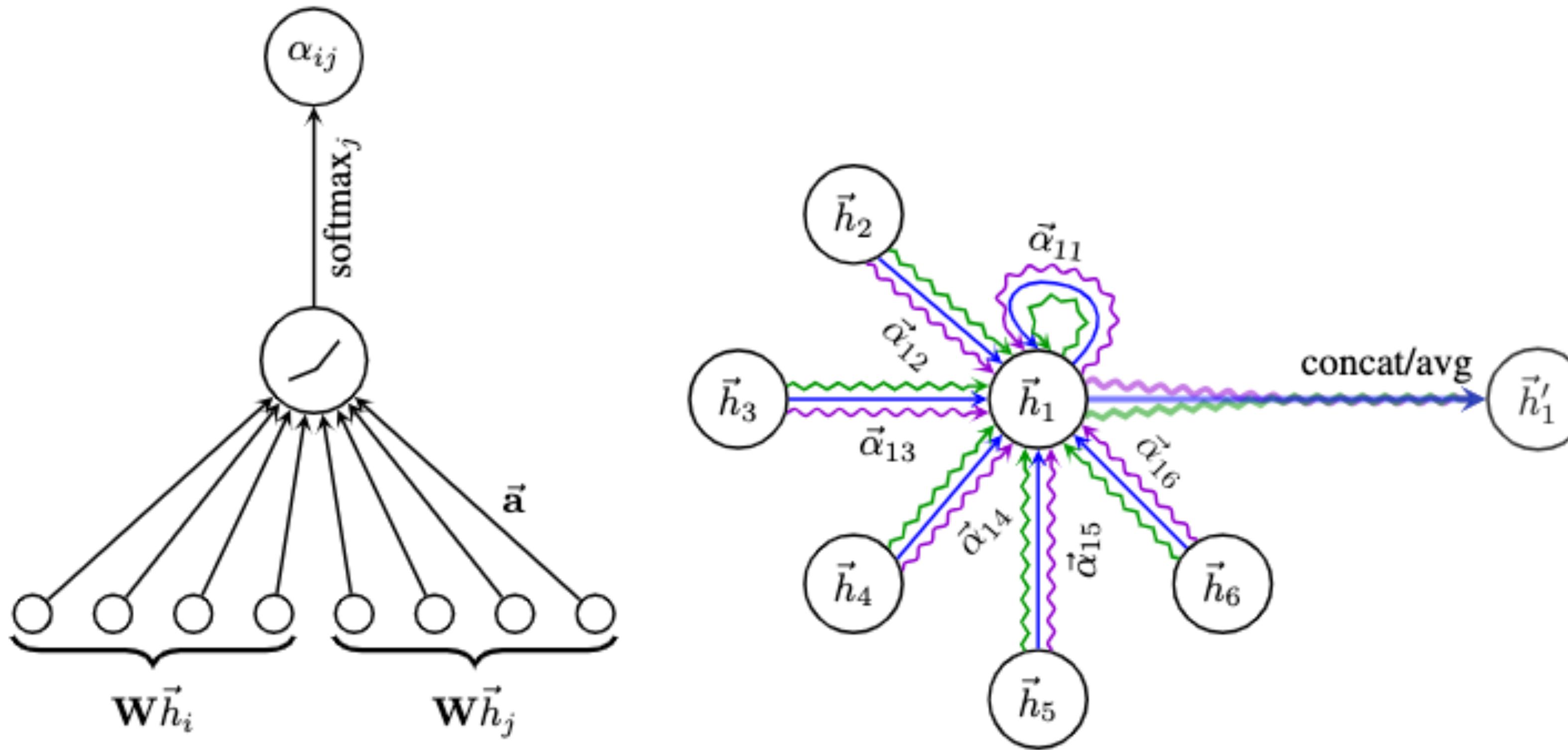


Figure 1: **Left:** The attention mechanism $a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$ employed by our model, parametrized by a weight vector $\vec{a} \in \mathbb{R}^{2F'}$, applying a LeakyReLU activation. **Right:** An illustration of multi-head attention (with $K = 3$ heads) by node 1 on its neighborhood. Different arrow styles and colors denote independent attention computations. The aggregated features from each head are concatenated or averaged to obtain \vec{h}'_1 .

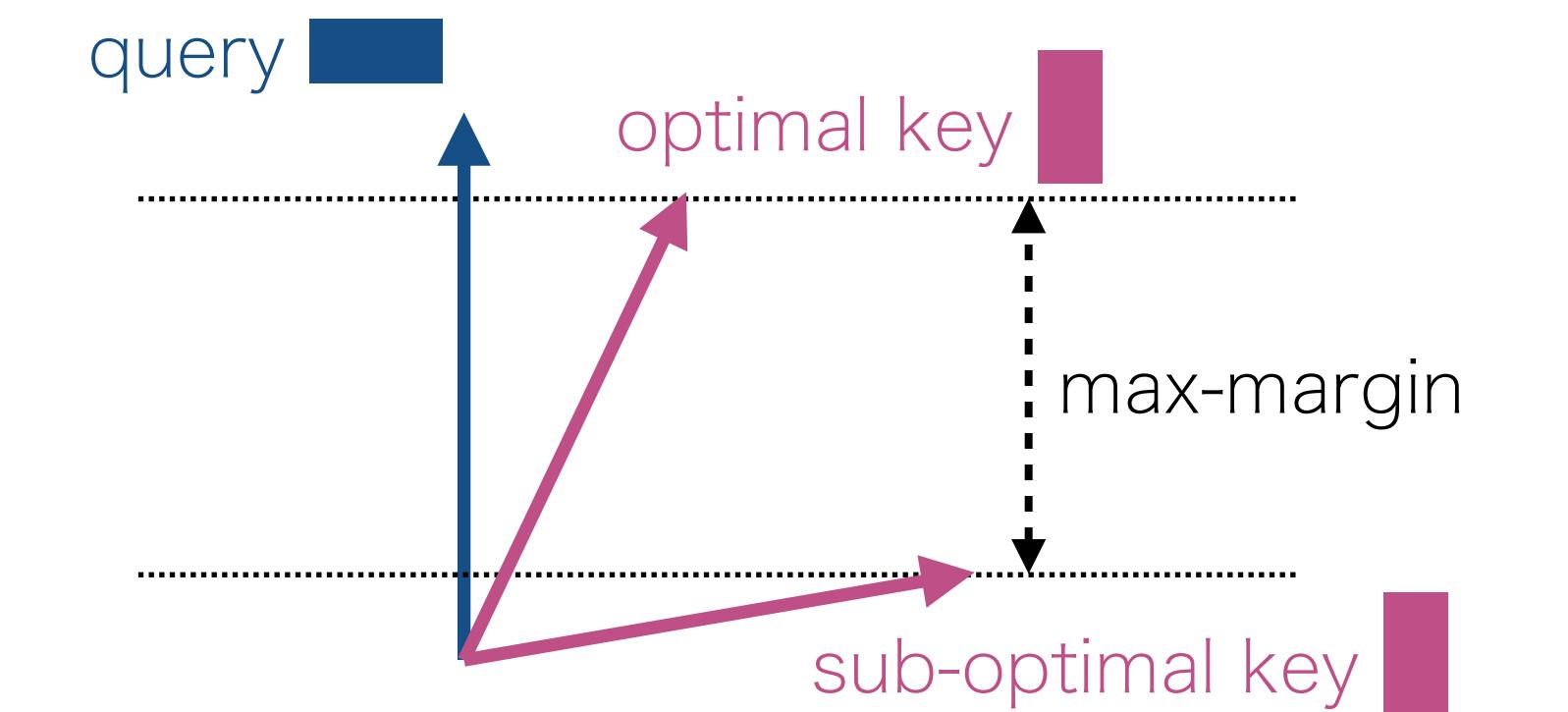
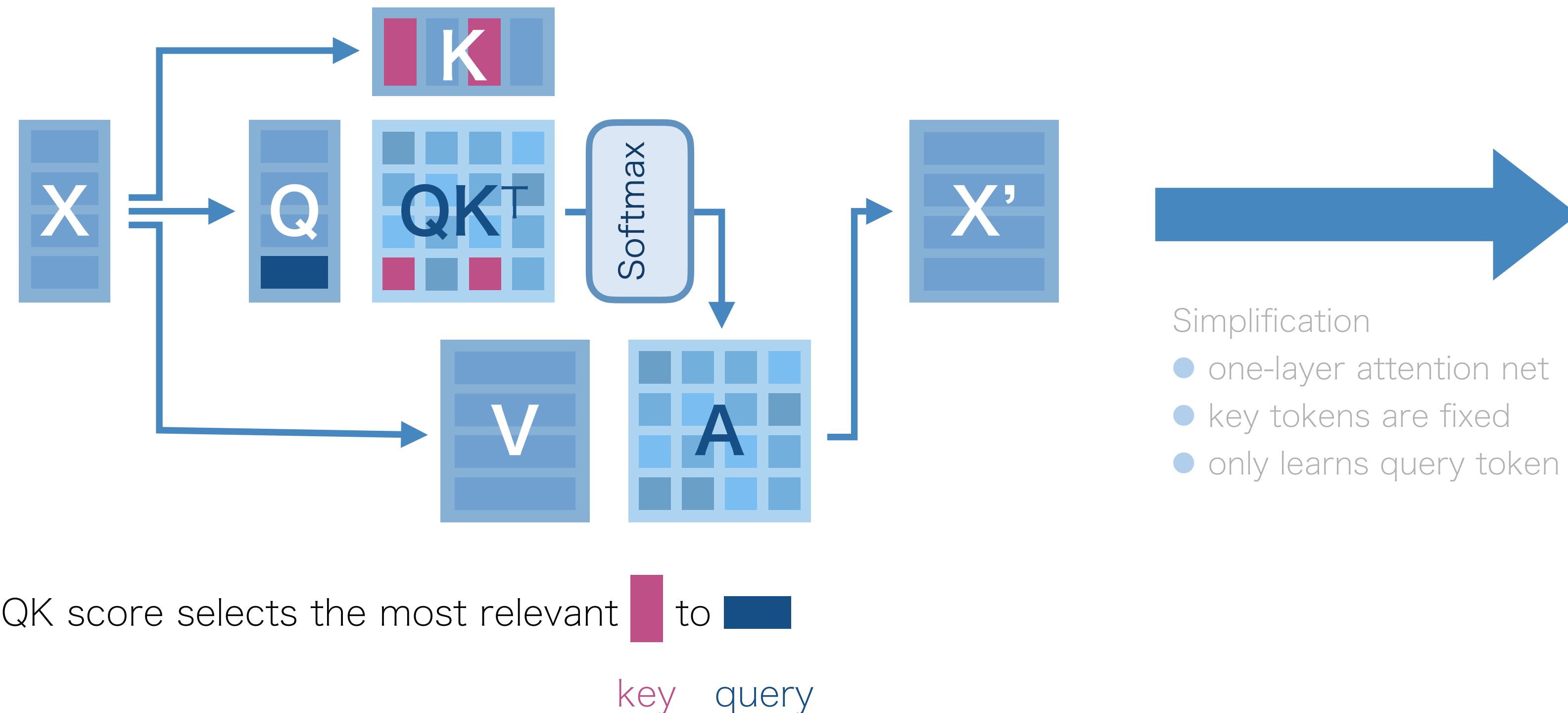
Why is attention attractive for us?

Recent research advances in machine learning community

Why attention attracts ML people? (1/3)

[Tarzanagh+ 2023]

- Mixing mechanism — How does attention selectively chooses relevant tokens?

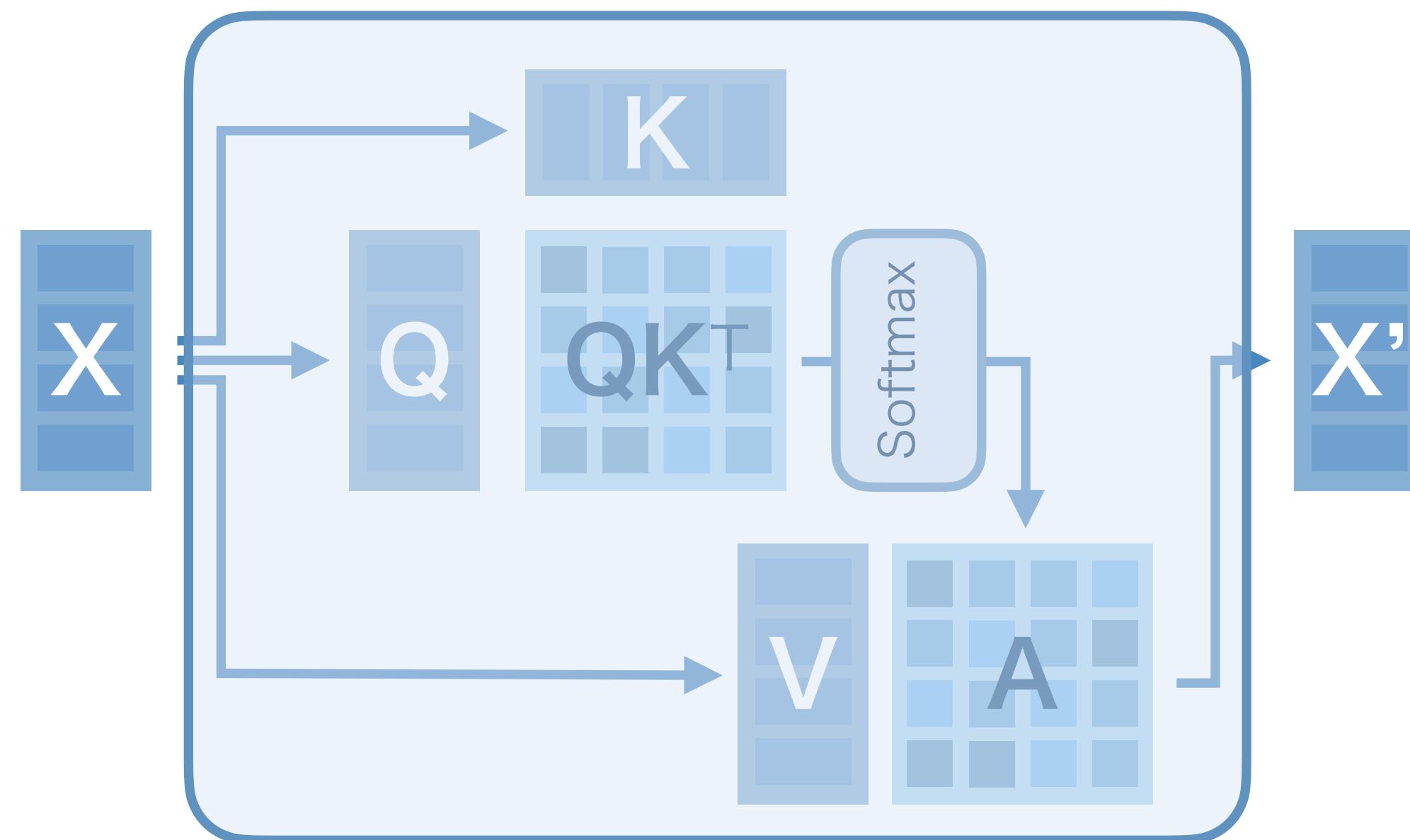
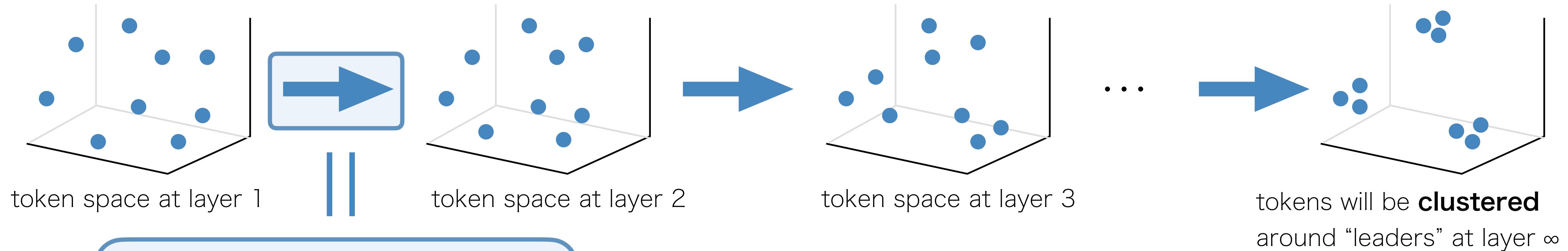


Learned query implicitly **maximize margin**

Why attention attracts ML people? (2/3)

- Representation learning — What encoded features should look like?

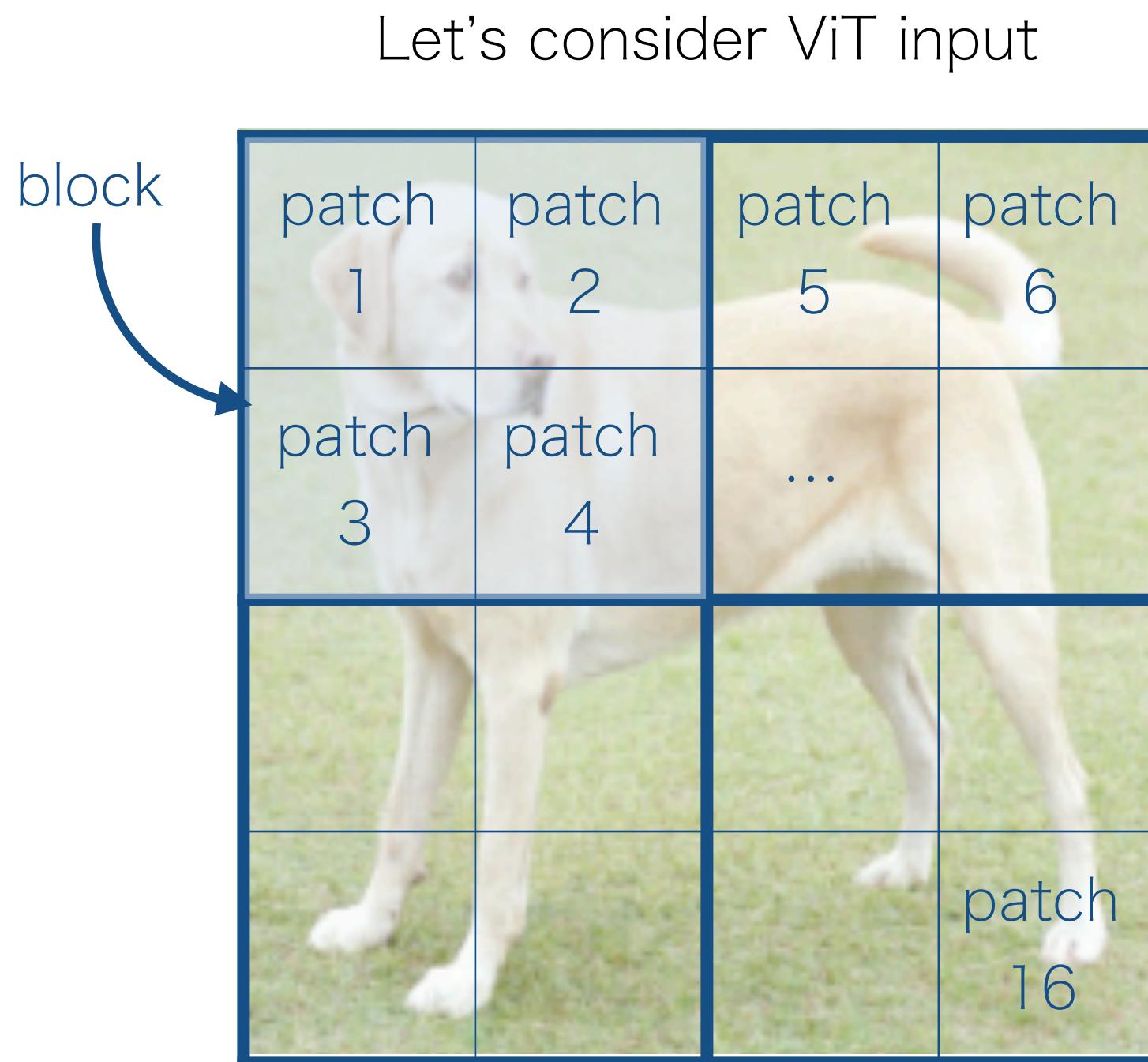
[Geshkovski+ 2023]



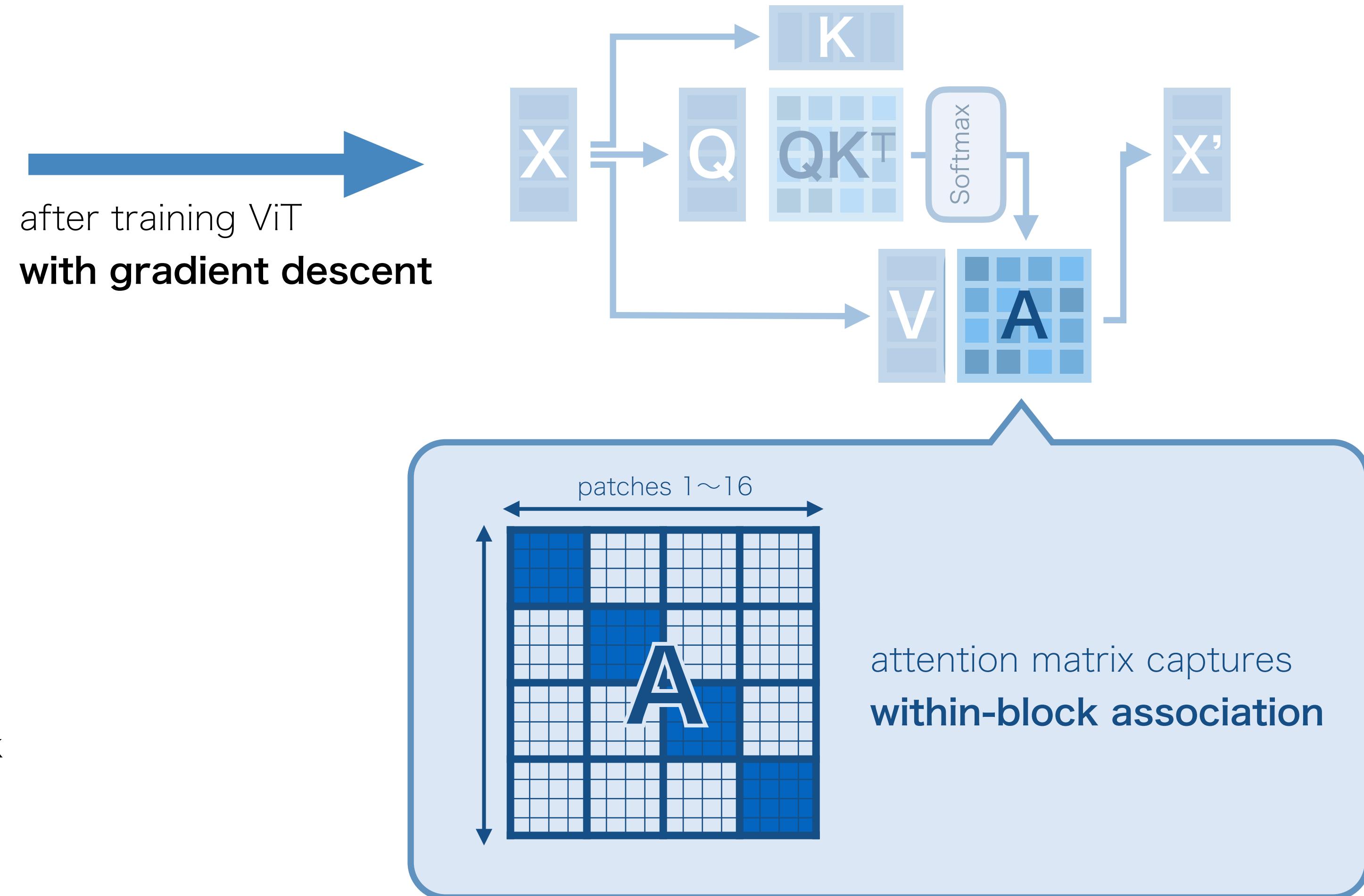
Why attention attracts ML people? (3/3)

[Jelassi+ 2022]

- Structure prediction — How does attention efficiently extract structured patterns?



patches are structured within each block



Why attention attracts ML people? (extra)

- “Similarity” among data points is interesting/cheap signal lying in data!

- Classification with only **similarity** & unlabeled data (SU classification)

Bao, H., Niu, G., & Sugiyama, M. (ICML2018).

Classification from Pairwise Similarity and Unlabeled Data.

- Classification with **similarity**, **dissimilarity**, & unlabeled data

Shimada, T., **Bao, H.**, Sato, I., & Sugiyama, M. (NeCo2021)

Classification from Pairwise Similarities/Dissimilarities and Unlabeled Data via Empirical Risk Minimization.

- Relationship between **metric** learning and classification

Bao, H.*, Shimada, T.*, Xu, L., Sato, I., & Sugiyama, M. (AISTATS2022)

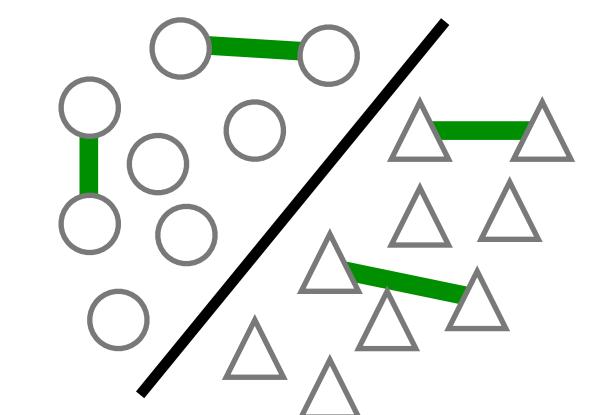
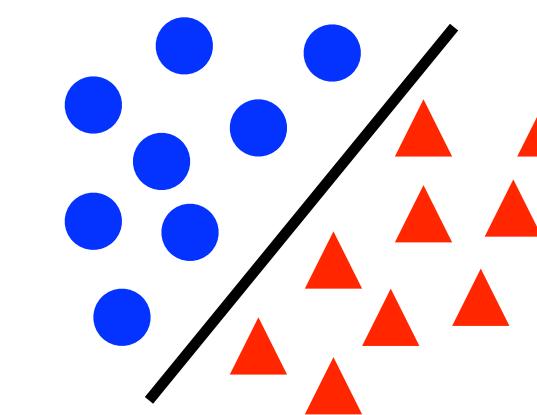
Pairwise Supervision Can Provably Elicit a Decision Boundary.

- Relationship between **contrastive** learning and linear probing

Bao, H., Nagano, Y., & Nozawa, N. (ICML2022)

On the Surrogate Gap between Contrastive and Supervised Losses.

Supervised Classification SU Classification



Tons of mysteries behind attention!

- Many practical applications
 - ❖ Cross attention for multimodal inputs
 - ❖ Interpretability
 - ❖ Graph/sequence/structured inputs
- Yet our understanding is elusive
 - ❖ (Layer-/epoch-wise) dynamics
 - ❖ Representation learning
 - ❖ Expressivity

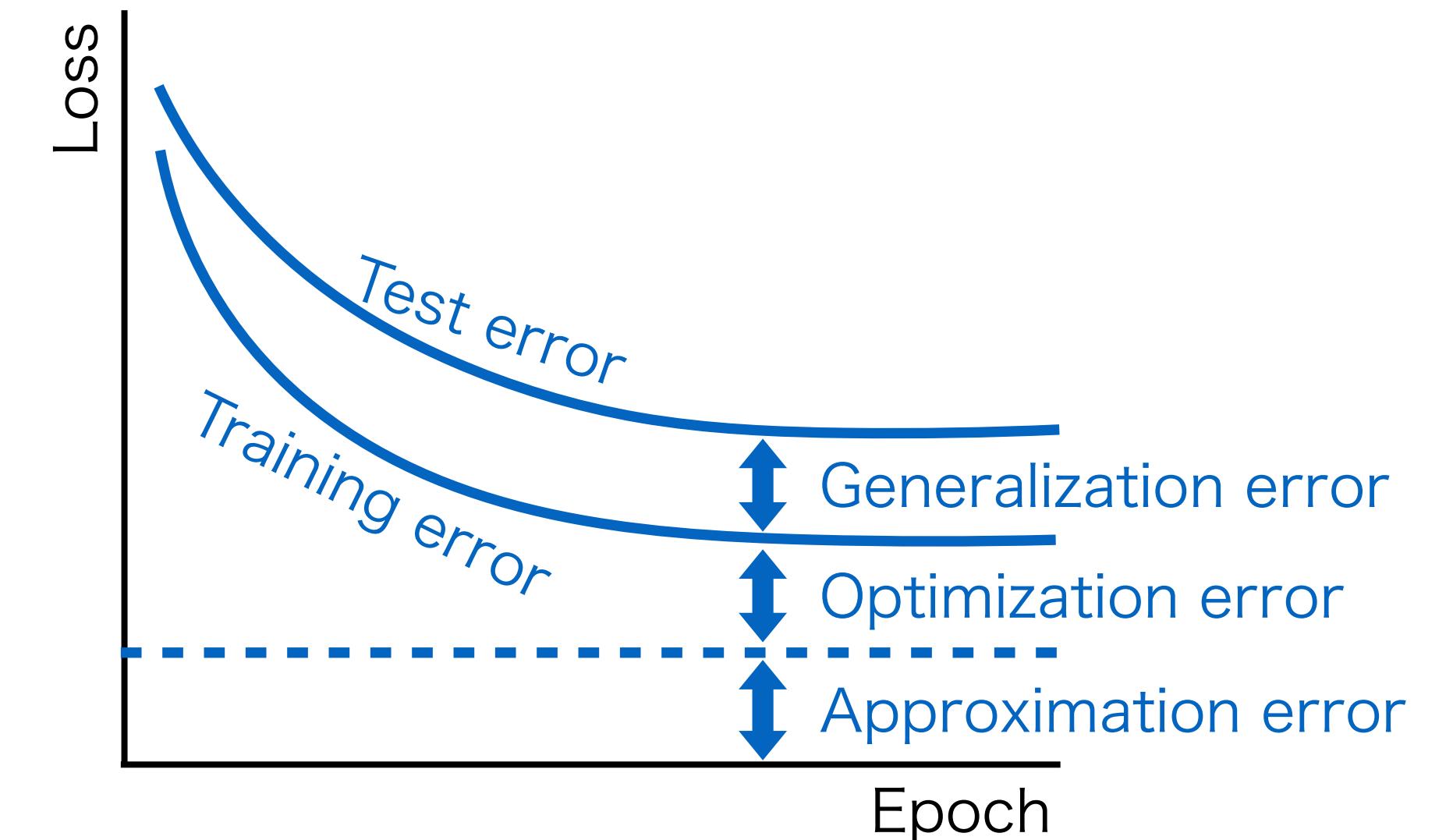
How attention matters in learning

From eigenspectrum perspective

Three factors in learning theory (in general)

- Generalization (\leftrightarrow complexity error)
- Optimization (\leftrightarrow optimization error)
- Model expressivity (\leftrightarrow approximation error)

Error decomposition of learning curve

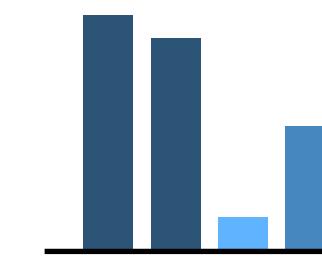


Cropped from Imaizumi-san's slides at MLSS'24

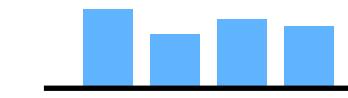
Related work | Concentrated attention is better!

- Rank collapse [Dong+ 2021]: attention is close to uniform (roughly speaking)

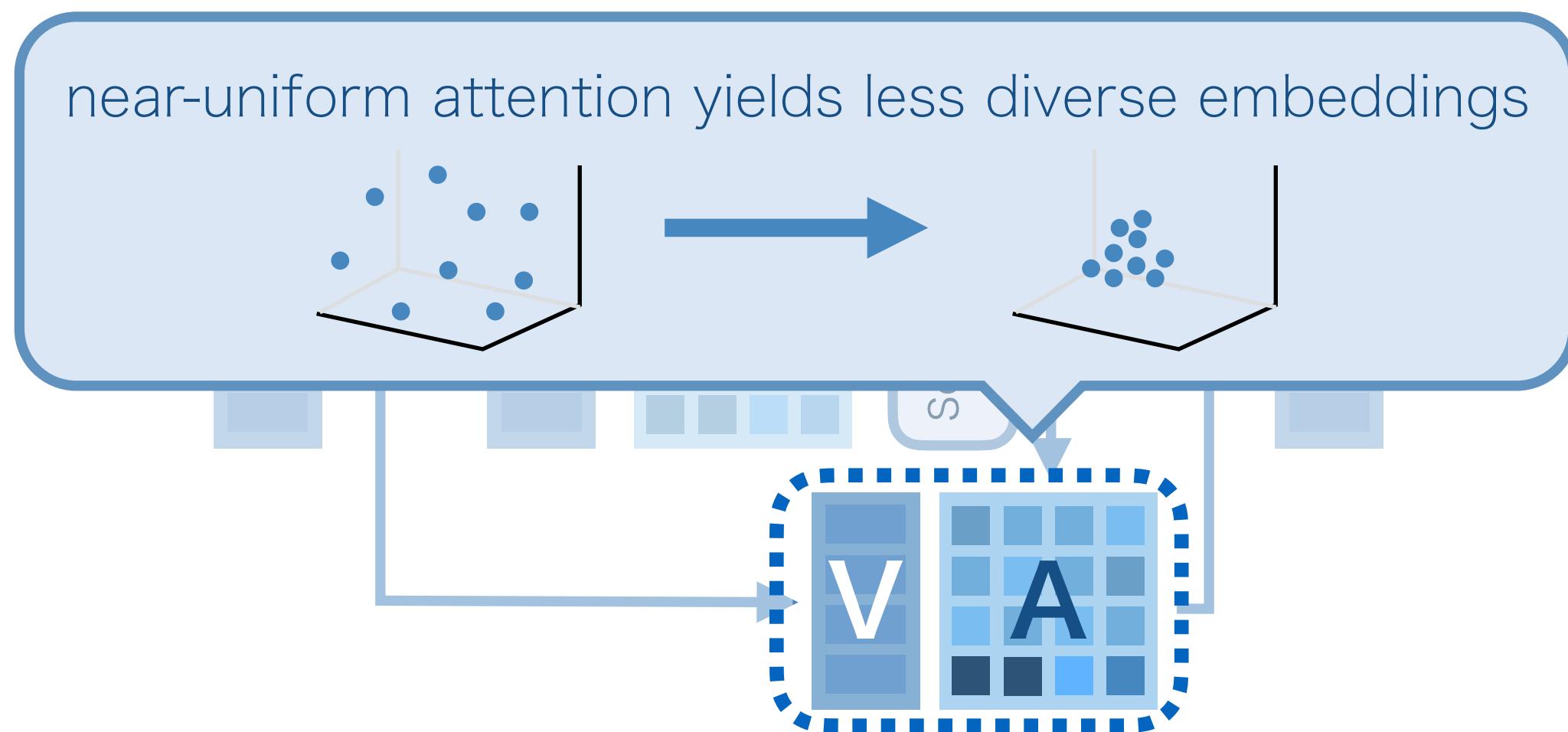
- induces **low model expressivity**
(yielding large approximation error 😞)



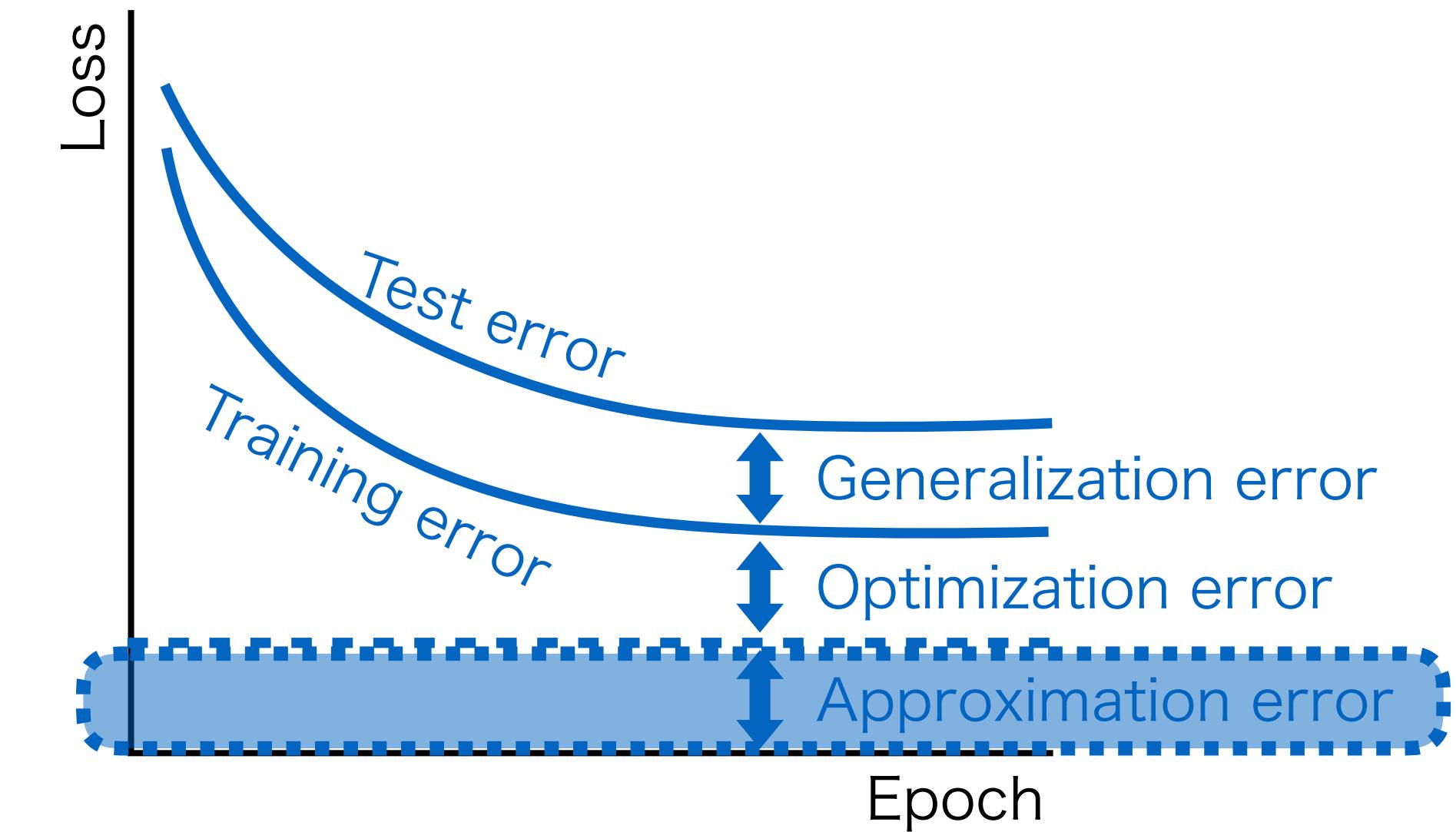
non-uniform attention
→ flexible model 😊



uniform attention
→ not flexible model 😞

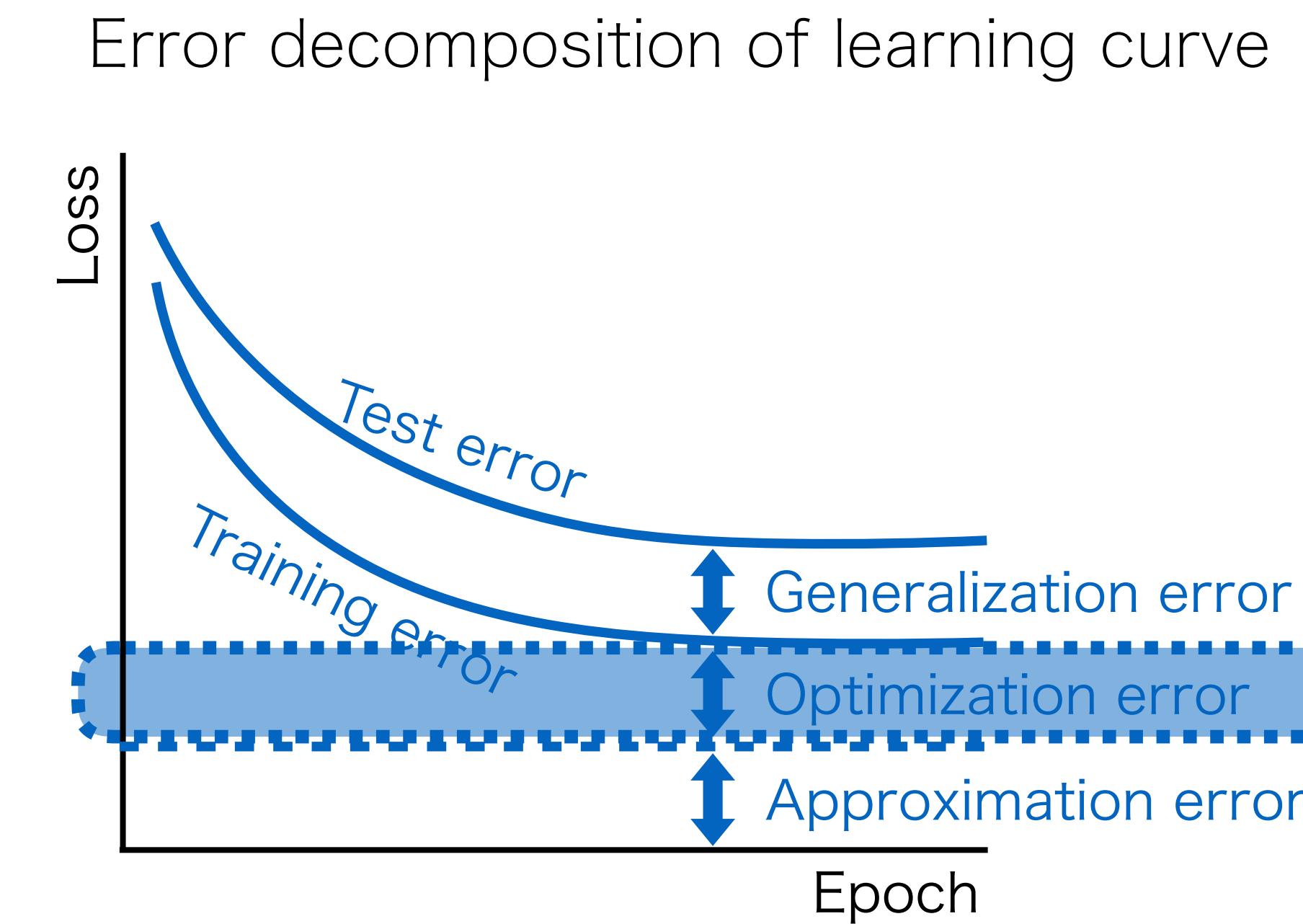
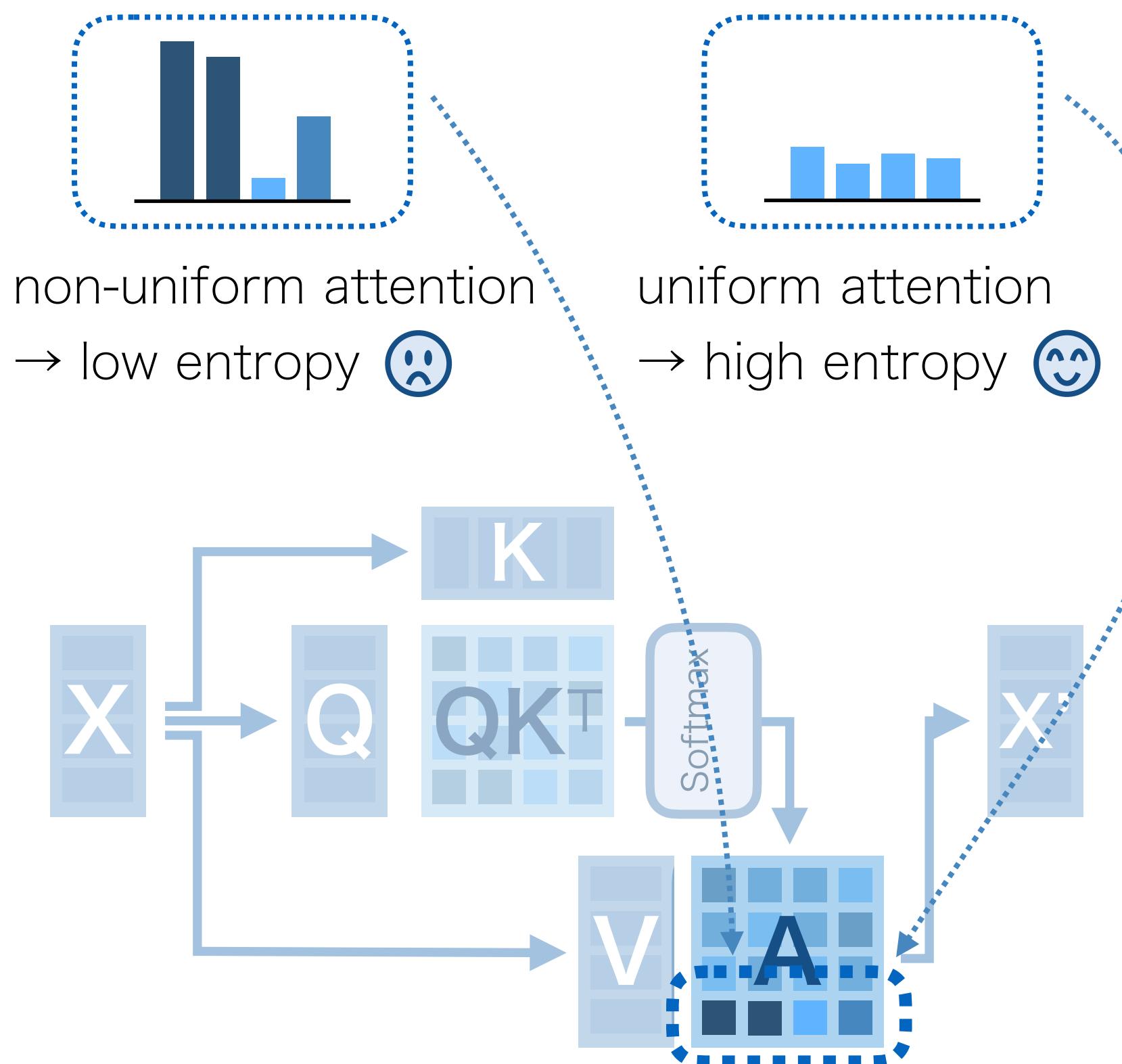


Error decomposition of learning curve

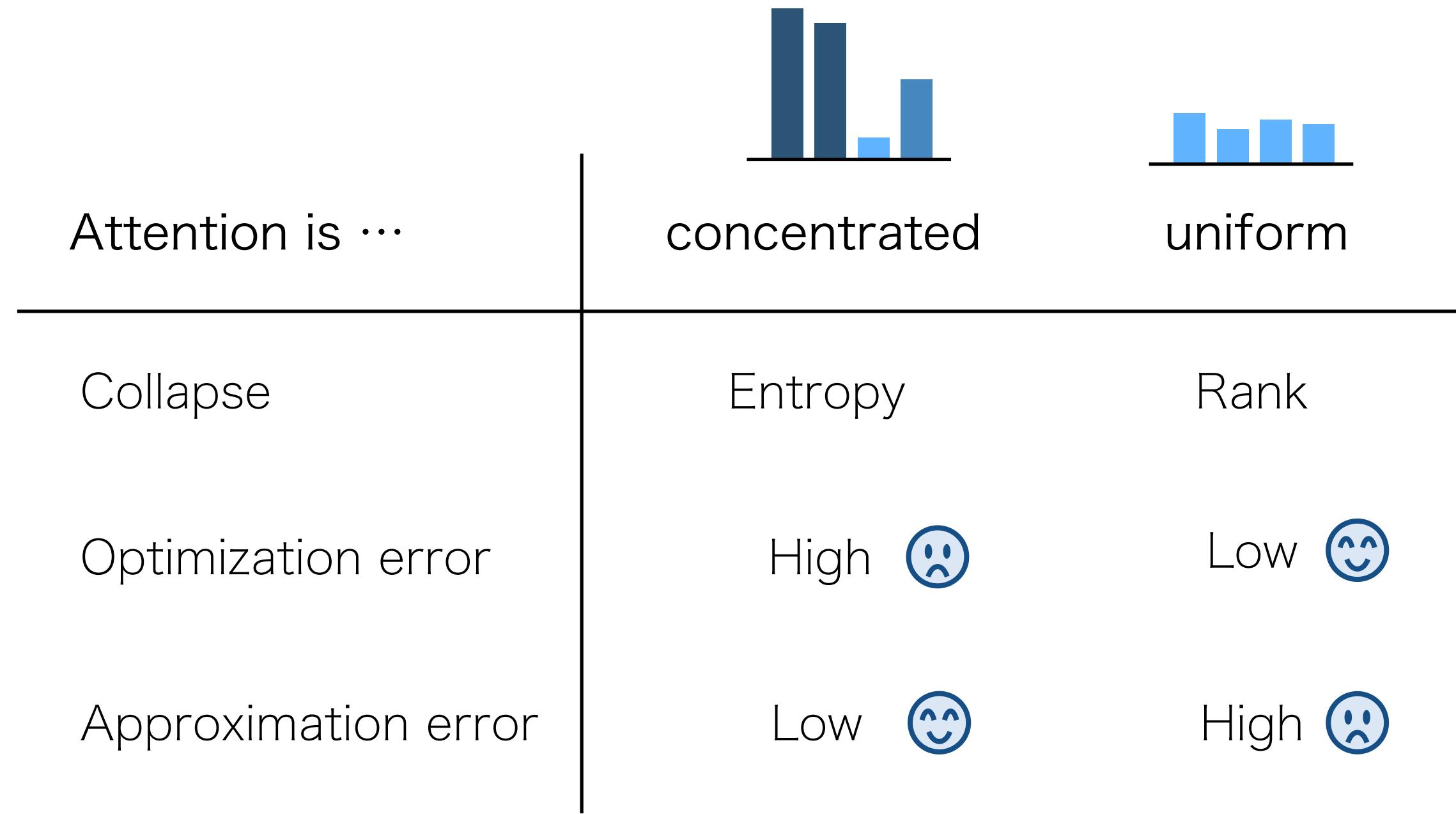


Related work | Uniform attention is better!

- **Entropy collapse** [Zhai+ 2023]: entropy of attention (as a probability distribution) is low
 - ❖ tends to **get stuck into plateaus** of loss landscape (yielding large optimization error 😞)



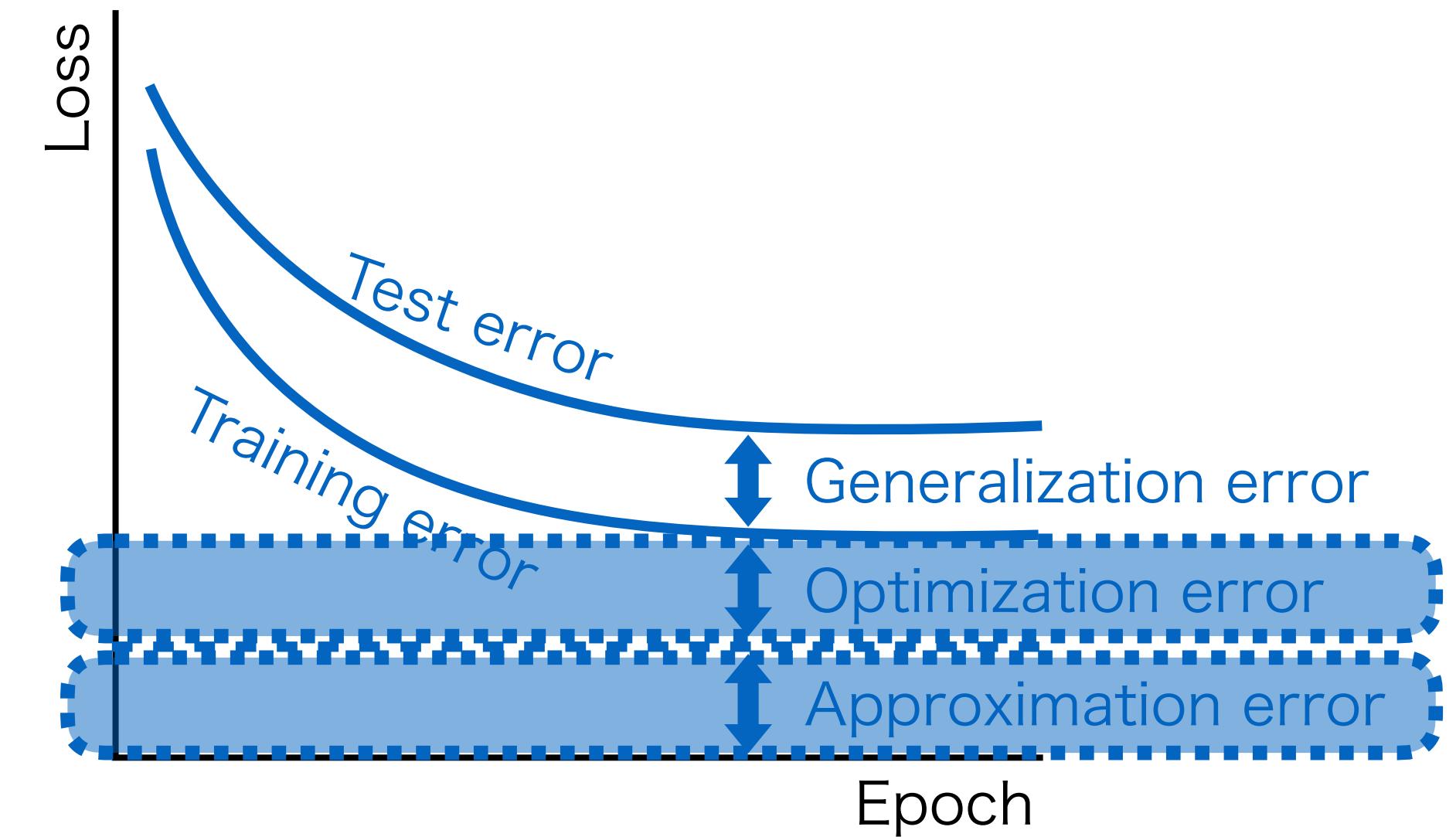
How should attention look like??



RQ. Are they reconcilable?

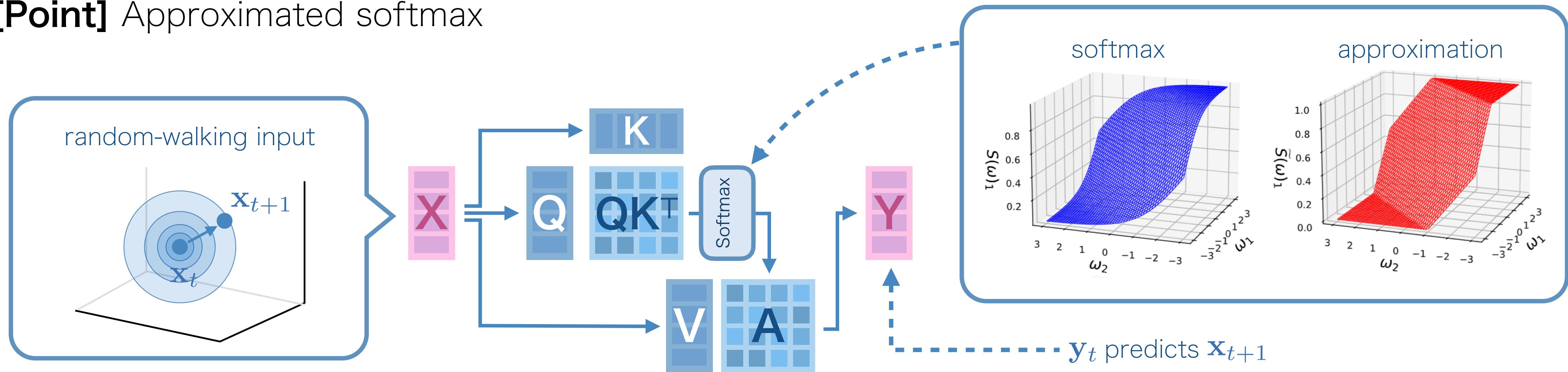
RQ. When concentrated/uniform?

Error decomposition of learning curve



Setup | Our analysis model

- Model: 1-layer 1-head self-attention network
 - ❖ without layer normalization & positional encoding / with feed-forward net
- Pre-training task: next token prediction
- **[Point]** Data: Gaussian random walk $\mathbf{x}_{t+1} \sim \mathcal{N}(\mathbf{x}_t, \Sigma)$
 - ❖ disclaimer: no embedding layer (token -> emb); \mathbf{x}_t is a direct input to TF
- **[Point]** Approximated softmax



Main result | Characterize attention shape

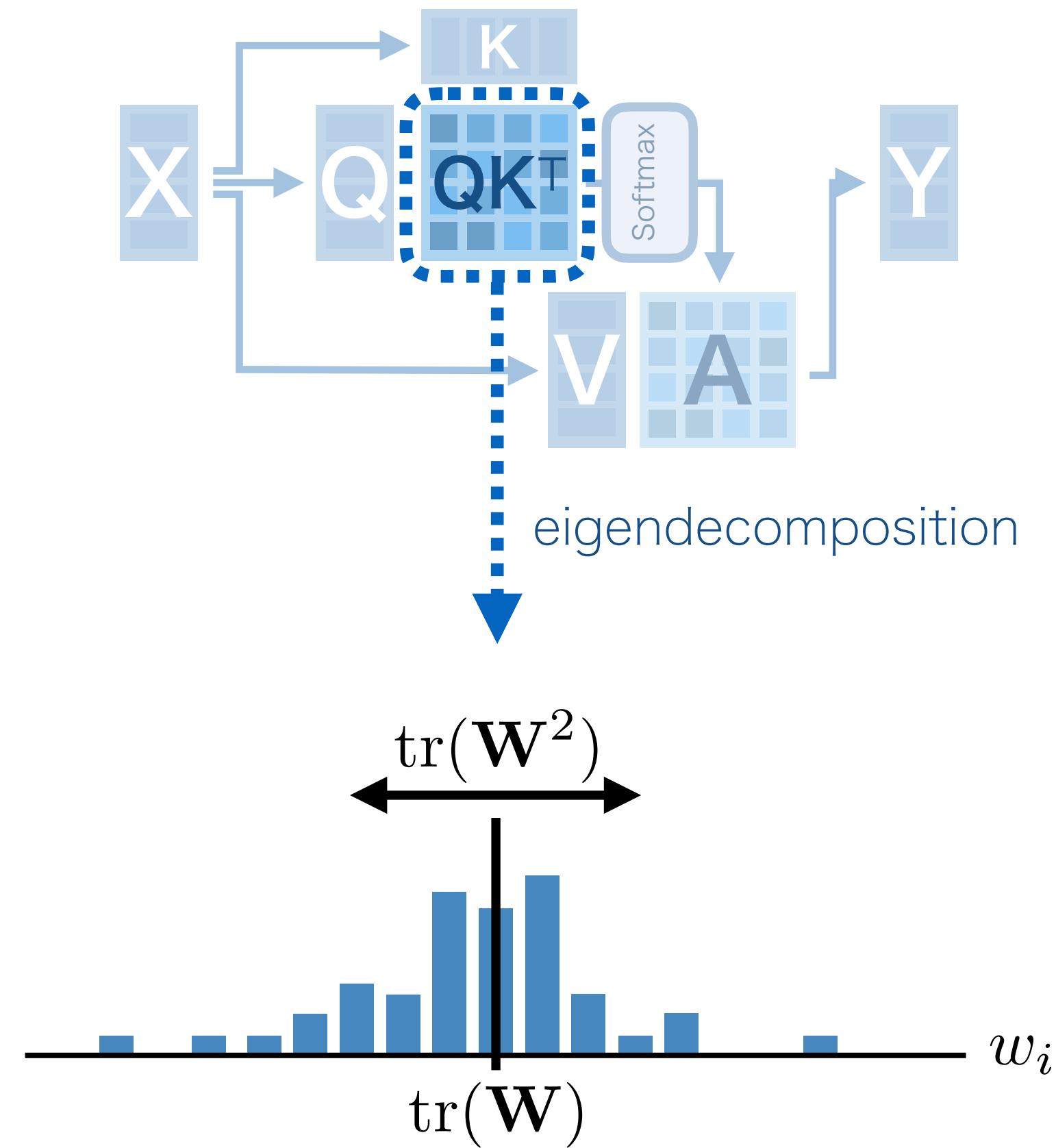
- $\text{tr}(\mathbf{W})$ (\doteq eigval mean) and $\text{tr}(\mathbf{W}^2)$ (\doteq eigval variance) determines attention shape

- ❖ $\mathbf{W}^2 := \mathbf{W}_Q^\top \mathbf{W}_K \Sigma$ is scaled QK parameter

under mild assumptions:

$$\frac{1}{d} \text{tr}(\mathbf{W}) = \frac{1}{d} \sum_{i=1}^d w_i \quad \text{(1st moment)}$$

$$\frac{1}{d} \text{tr}(\mathbf{W}^2) = \frac{1}{d} \sum_{i=1}^d w_i^2 \quad \text{(2nd moment)}$$



Main result | Characterize attention shape

- $\text{tr}(\mathbf{W})$ (\doteq eigval mean) and $\text{tr}(\mathbf{W}^2)$ (\doteq eigval variance) determines attention shape
 - ❖ $\mathbf{W}^2 := \mathbf{W}_Q^\top \mathbf{W}_K \Sigma$ is scaled QK parameter

analytical expression of
attention strength at token position θ

$$\rho(\theta) := \Phi\left(\left(\theta - \frac{1}{2}\right)\xi; \theta\right) - \Phi\left(\left(\theta - \frac{1}{2}\right)\xi - \frac{1}{\eta}; \theta\right),$$

$$\xi := \frac{\text{tr}(\mathbf{W})}{\sqrt{\text{tr}(\mathbf{W}^2)}}, \quad \eta := \frac{\sqrt{\text{tr}(\mathbf{W}^2)}}{\lambda},$$

$$\Phi(z; \theta) := \frac{1}{2} \operatorname{erf}\left(\frac{z}{\sqrt{2(2\theta^2 + \frac{7}{12})}}\right),$$

Main result | Characterize attention shape

20
/ 29

- $\text{tr}(\mathbf{W})$ (\doteq eigval mean) and $\text{tr}(\mathbf{W}^2)$ (\doteq eigval variance) determines attention shape

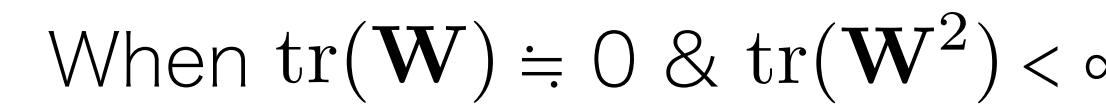
- ❖ $\mathbf{W}^2 := \mathbf{W}_Q^\top \mathbf{W}_K \boldsymbol{\Sigma}$ is scaled QK parameter

analytical expression of
attention strength at token position θ

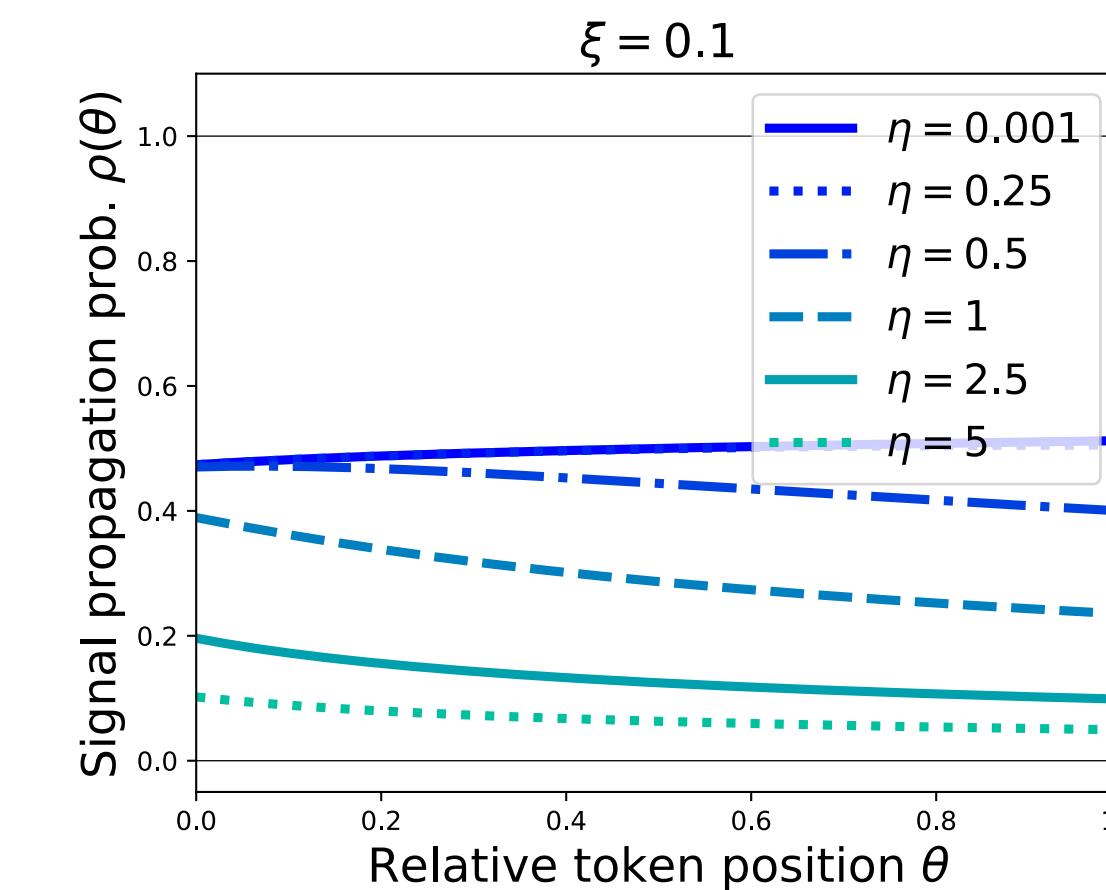
$$\rho(\theta) := \Phi\left(\left(\theta - \frac{1}{2}\right)\xi; \theta\right) - \Phi\left(\left(\theta - \frac{1}{2}\right)\xi - \frac{1}{\eta}; \theta\right),$$

$$\xi := \frac{\text{tr}(\mathbf{W})}{\sqrt{\text{tr}(\mathbf{W}^2)}}, \quad \eta := \frac{\sqrt{\text{tr}(\mathbf{W}^2)}}{\lambda},$$

$$\Phi(z; \theta) := \frac{1}{2} \operatorname{erf} \left(\frac{z}{\sqrt{2(2\theta^2 + \frac{7}{12})}} \right),$$

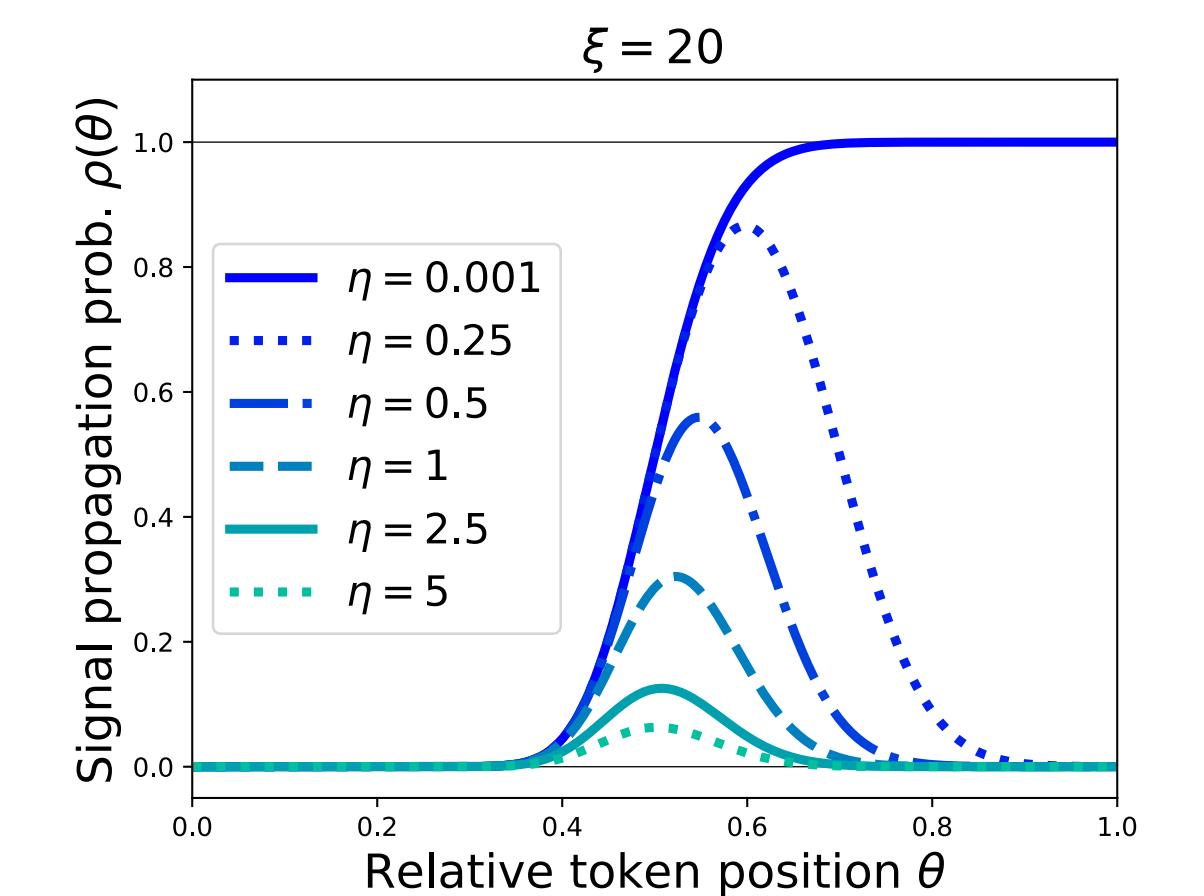
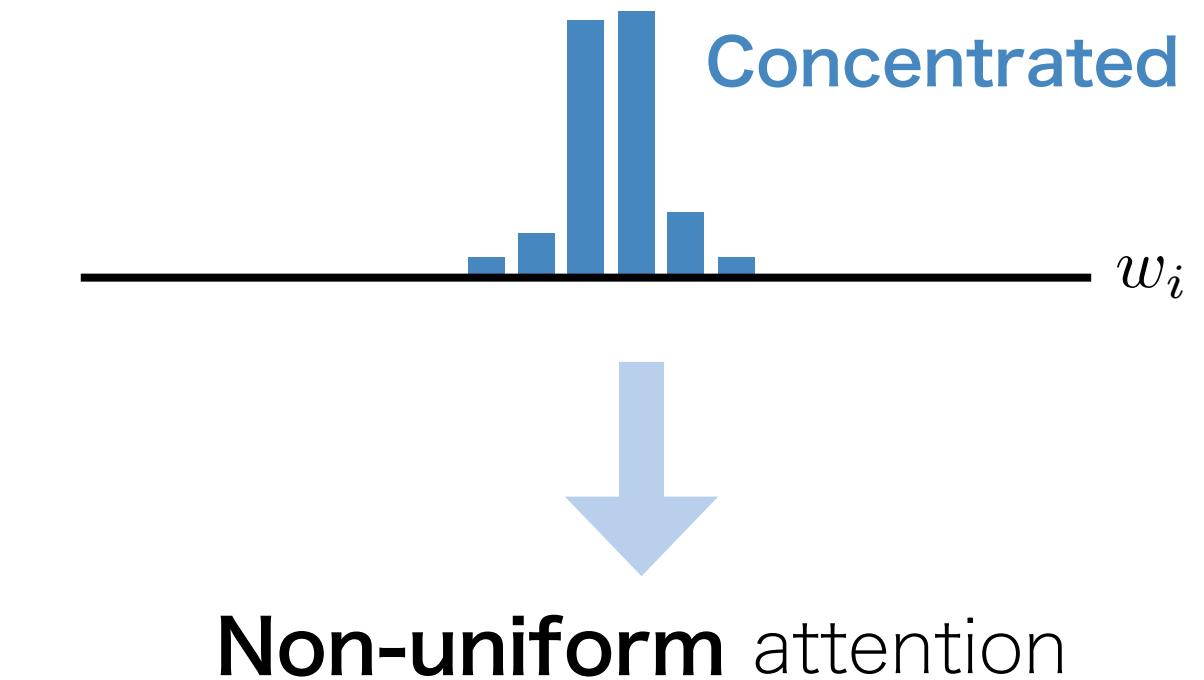


Uniform attention



When $\text{tr}(\mathbf{W}) \gg 0$ & $\text{tr}(\mathbf{W}^2) \doteq 0$

Concentrated



Relation to rank collapse

- **Rank collapse:** attention is close to uniform (roughly speaking)

❖ [Dong+ 2021] large $\|\mathbf{W}_{QK}\|_1$ (L1 norm) prevents rank collapse

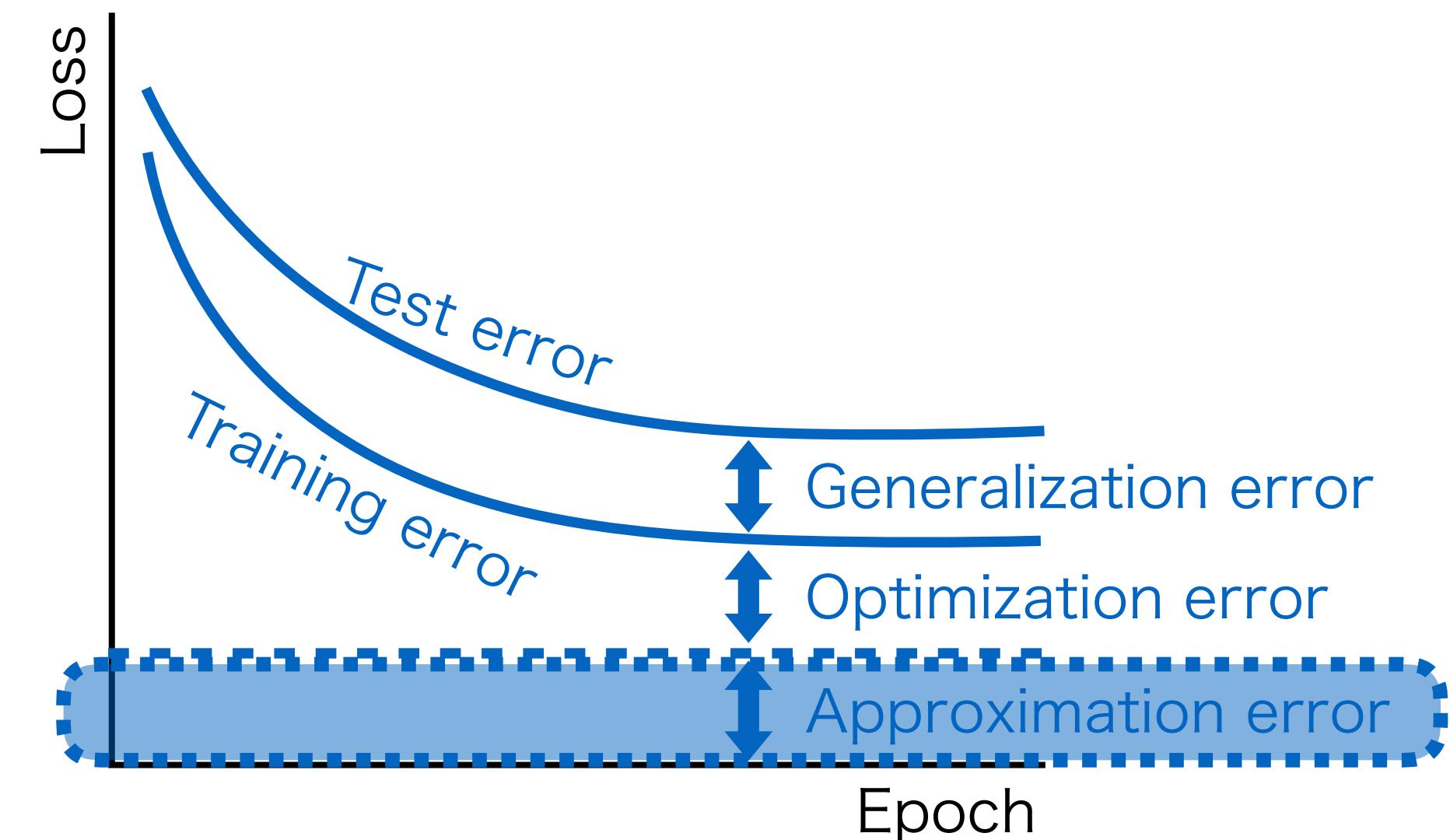
- From eigenspectrum viewpoint

❖ we can show $\|\mathbf{W}_{QK}\|_1 \geq C|\text{tr}(\mathbf{W})|$

❖ smaller eigenspectrum variance \Rightarrow larger $\|\mathbf{W}_{QK}\|_1$

$$\text{Var}[w_i] = \frac{1}{d} \underset{\text{decrease}}{\text{tr}(\mathbf{W}^2)} - \frac{1}{d^2} \underset{\text{fixed}}{\text{tr}(\mathbf{W})} \underset{\text{increase}}{| \text{tr}(\mathbf{W}) |^2}$$

Error decomposition of learning curve



Relation to entropy collapse

- **Entropy collapse:** entropy of attention (as a probability distribution) is low

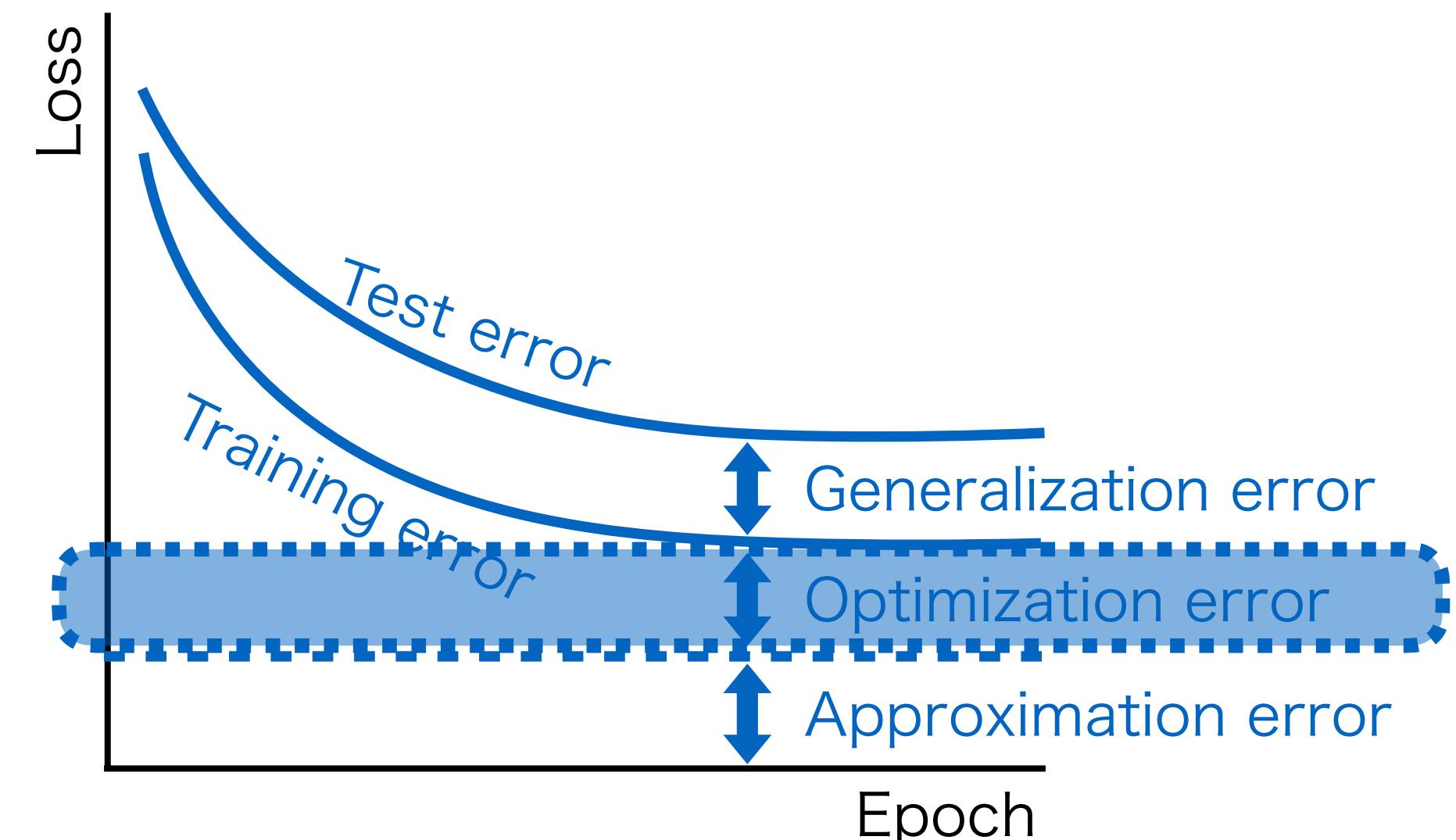
❖ [Zhai+ 2023] small $\|\mathbf{W}_{QK}\|_2$ (spectral norm) keeps attention entropy high

- From eigenspectrum viewpoint

- ❖ we can show $\|\mathbf{W}_{QK}\|_2 \leq C\sqrt{\text{tr}(\mathbf{W}^2)}$
- ❖ smaller eigenspectrum variance \Rightarrow smaller $\|\mathbf{W}_{QK}\|_2$

$$\text{Var}[w_i] = \frac{1}{d} \underset{\text{decrease}}{\text{tr}(\mathbf{W}^2)} - \frac{1}{d^2} \underset{\text{fixed}}{|\text{tr}(\mathbf{W})|^2}$$

Error decomposition of learning curve



Attention shape through eigenspectrum

RQ. Are two collapse phenomena reconcilable?

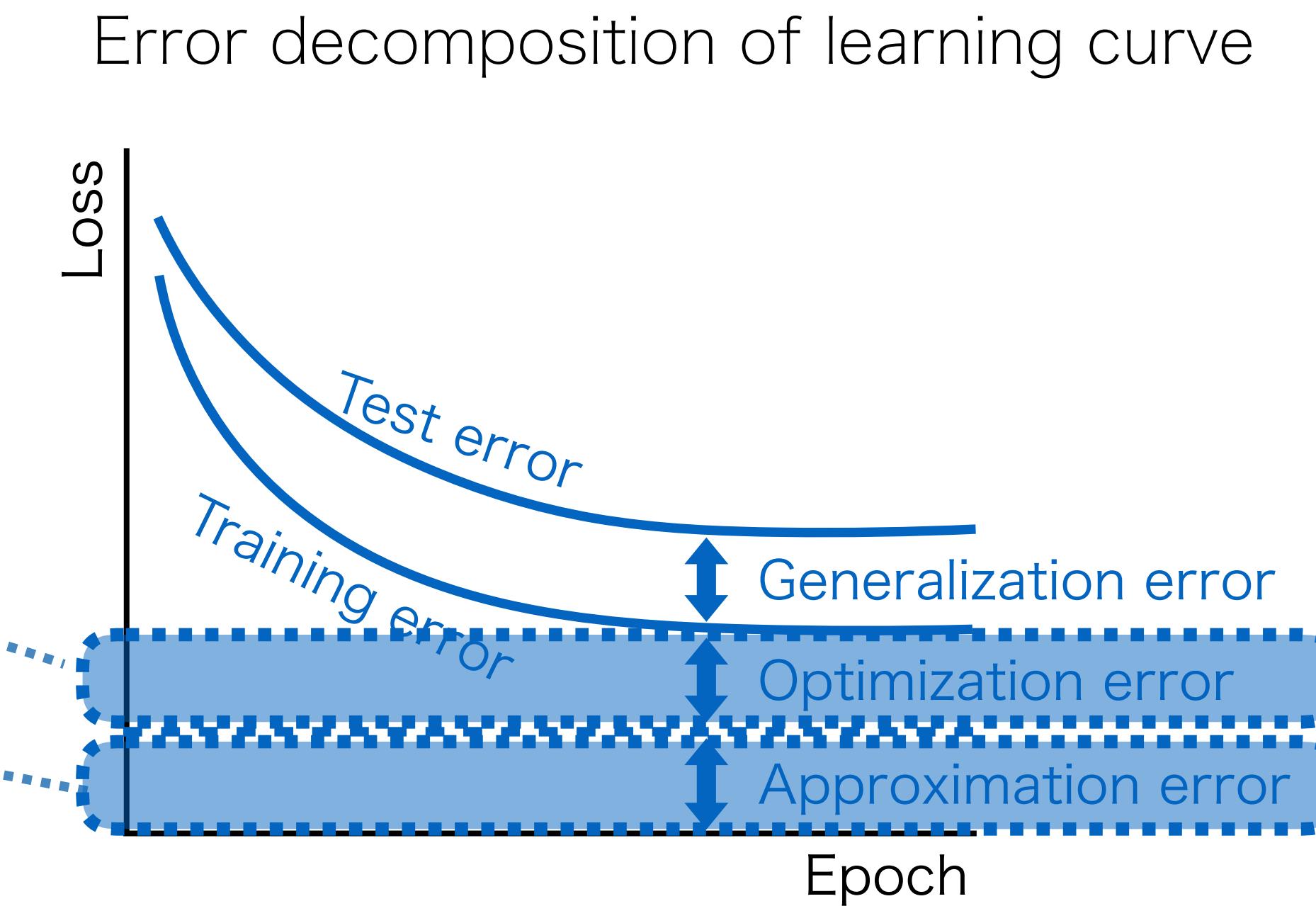
→ A. Yes, let's decrease eigenspectrum variance!

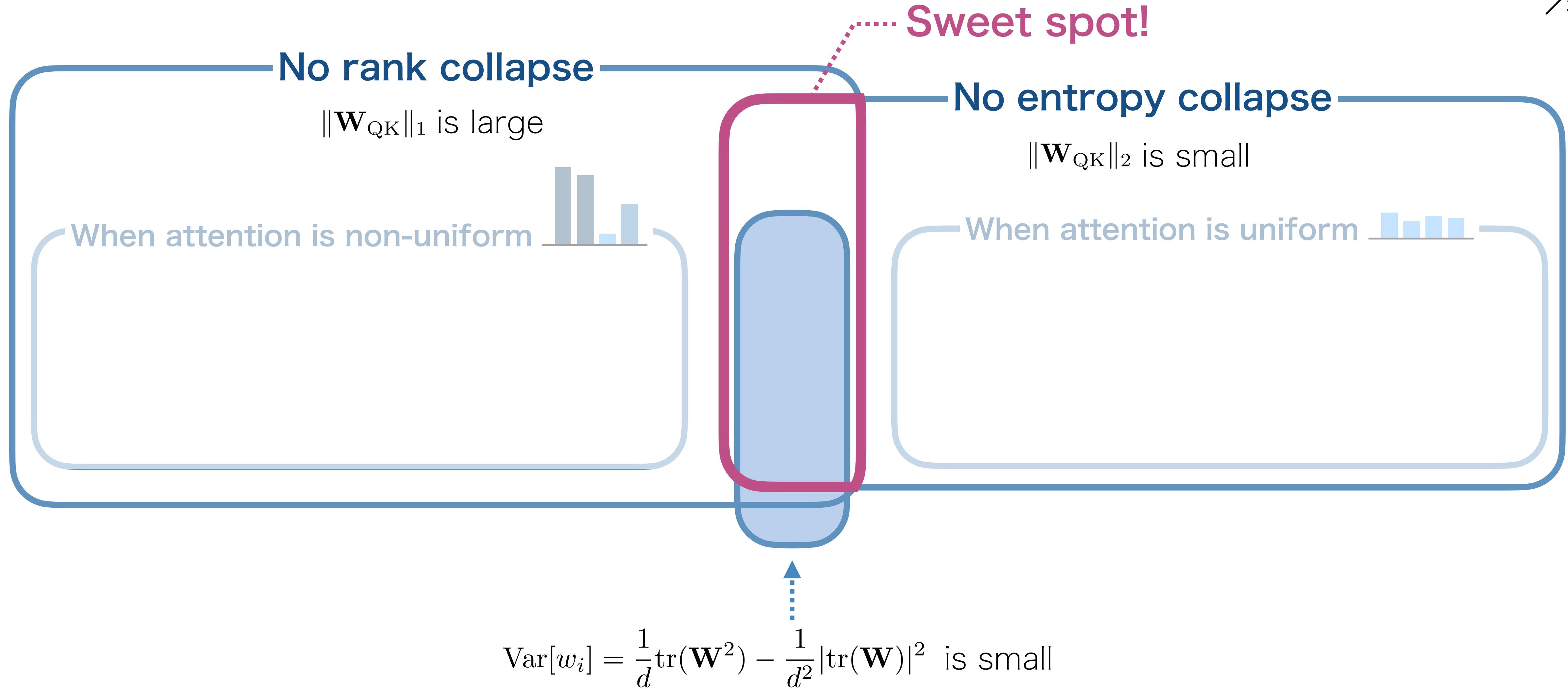
Disclaimer: this is only a **sufficient** condition!

$$\text{Var}[w_i] = \frac{1}{d} \underset{\text{decrease}}{\text{tr}(\mathbf{W}^2)} - \frac{1}{d^2} \underset{\text{increase}}{|\text{tr}(\mathbf{W})|^2}$$

avoids entropy collapse

avoids rank collapse





In more depth

Signal propagation probability | Why it matters?

- What is “attention strength”?

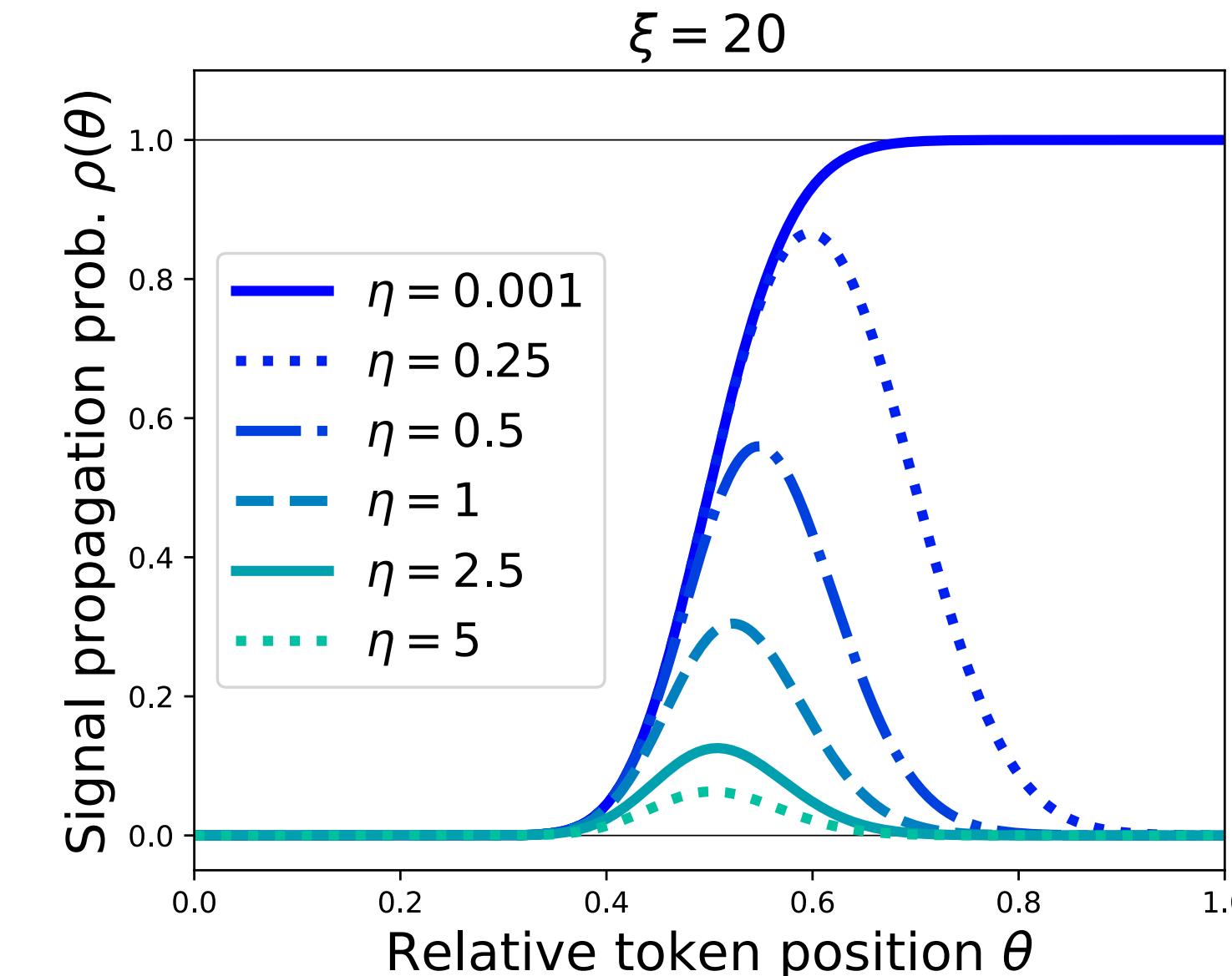
analytical expression of

attention strength at token position θ

$$\rho(\theta) := \Phi\left(\left(\theta - \frac{1}{2}\right)\xi; \theta\right) - \Phi\left(\left(\theta - \frac{1}{2}\right)\xi - \frac{1}{\eta}; \theta\right),$$

$$\xi := \frac{\text{tr}(\mathbf{W})}{\sqrt{\text{tr}(\mathbf{W}^2)}}, \quad \eta := \frac{\sqrt{\text{tr}(\mathbf{W}^2)}}{\lambda},$$

$$\Phi(z; \theta) := \frac{1}{2} \text{erf}\left(\frac{z}{\sqrt{2(2\theta^2 + \frac{7}{12})}}\right),$$



- ρ = signal propagation probability: “**How much the token contributes to learning?**”

$$\nabla_{\mathbf{W}_{\text{QK}}} \text{Loss} \approx \sum_{t=1}^T \text{signal}_t$$

gradient term depends on t-th token

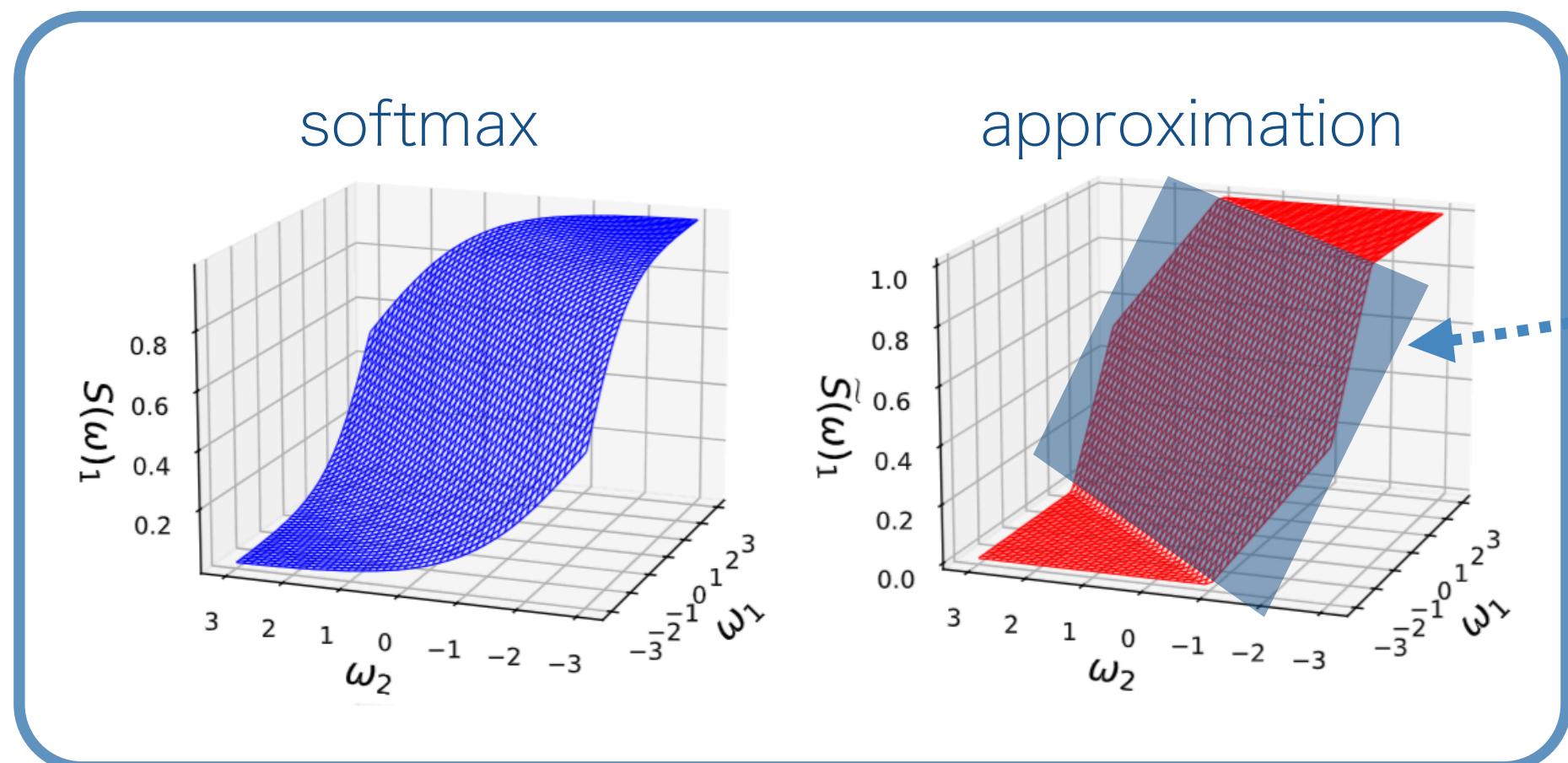
Signal propagation probability | Why it matters?

- ρ = signal propagation probability: “**How much the token contributes to learning?**”

$$\nabla_{\mathbf{W}_{QK}} \text{Loss} \approx \sum_{t=1}^T \text{signal}_t$$

gradient term depends on t-th token

- Signal exists only when softmax input (=query \times key) lies in “slope” area



- Specifically,

$$\rho_t := \Pr \left\{ 0 \leq \left\langle \frac{1}{T} \mathbf{e}^t - \frac{1}{T^2} \mathbf{1}, \frac{\mathbf{X}^\top \mathbf{W}_{QK} \mathbf{x}_T}{\sqrt{d}} \right\rangle + \frac{1}{T} \leq 1 \right\}$$

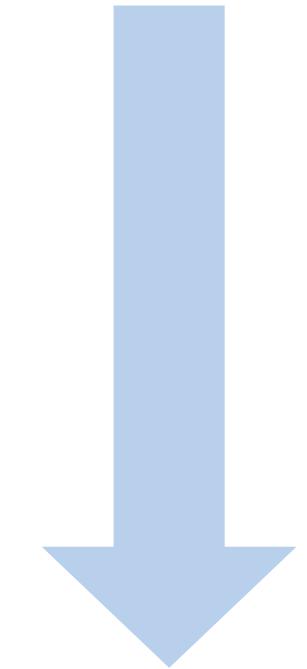
It's quite different from attention value itself, but matters when we consider learning dynamics

Our assumption and approximation

$$\rho_t := \Pr \left\{ 0 \leq \left\langle \frac{1}{T} \mathbf{e}^t - \frac{1}{T^2} \mathbf{1}, \frac{\mathbf{X}^\top \mathbf{W}_{QK} \mathbf{x}_T}{\sqrt{d}} \right\rangle + \frac{1}{T} \leq 1 \right\}$$

↳ $\langle \boldsymbol{\gamma}^t, \boldsymbol{\omega} \rangle + \gamma_0^t$

- Assumption: Gaussian random walk $\mathbf{x}_{t+1} \sim \mathcal{N}(\mathbf{x}_t, \boldsymbol{\Sigma})$
- Approximation: Gaussian approximation $\langle \boldsymbol{\gamma}^t, \boldsymbol{\omega} \rangle + \gamma_0^t \sim \mathcal{N}(\mu^t, v^t)$



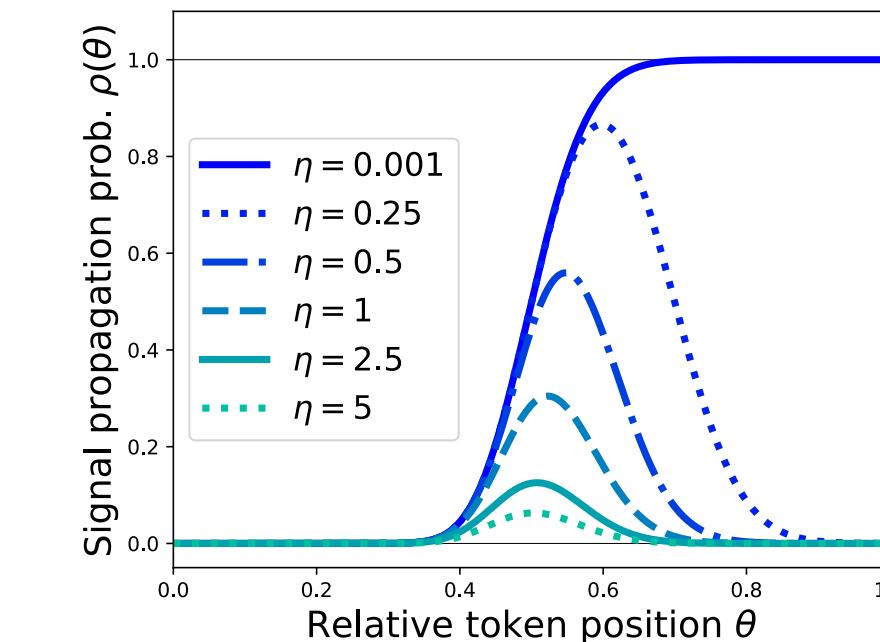
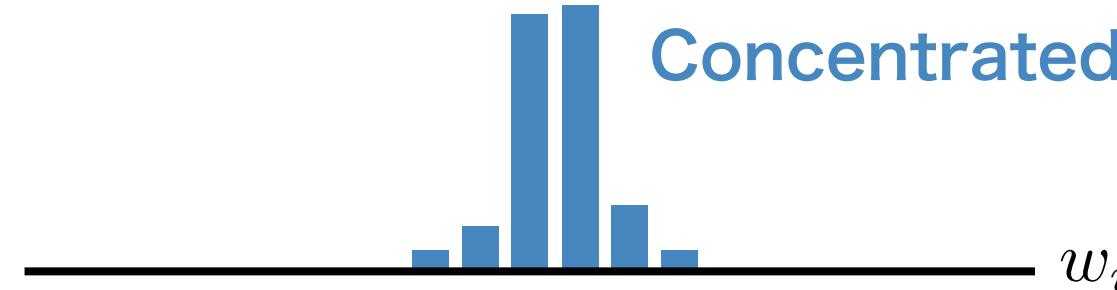
these simplifications are essential to obtain the analytical expression
(we can reduce the derivation to Gaussian moment calculation)

$$\rho(\theta) := \Phi\left(\left(\theta - \frac{1}{2}\right)\xi; \theta\right) - \Phi\left(\left(\theta - \frac{1}{2}\right)\xi - \frac{1}{\eta}; \theta\right),$$

Summary

RQ. When is attention concentrated/uniform?

When $\text{tr}(\mathbf{W}) \gg 0$ & $\text{tr}(\mathbf{W}^2) \approx 0$



RQ. Are two collapse phenomena reconcilable?

$$\text{Var}[w_i] = \frac{1}{d} \text{tr}(\mathbf{W}^2) - \frac{1}{d^2} |\text{tr}(\mathbf{W})|^2$$

decrease decrease increase

↓ ↓

avoids entropy collapse avoids rank collapse

