

# **Foundation Models**

Instructor - Simon Lucey

**RVSS - 2026**



**AUSTRALIAN  
INSTITUTE FOR  
MACHINE LEARNING**

TECH · ELON MUSK

# Elon Musk's just fired up Colossus—the world's largest Nvidia GPU supercomputer built in just three months from start to finish

BY CHRISTIAAN HETZNER

September 3, 2024 at 11:48 PM GMT+9:30



xAI founder Elon Musk aims to double the capacity of his Memphis investors could end up benefiting as well thanks to Optimus.

RICHARD BORD—WIREIMAGE/GETTY IMAGES

<https://fortune.com/2024/09/03/elon-musk-xai-nvidia-colossus/>



Contains 180,000 NVIDIA GPU (H100,H200,GB200) processors, estimated cost \$U4.5 billion.

Will use 1.3 million gallons of water per day to cool servers, and consume 280 megawatts of power.



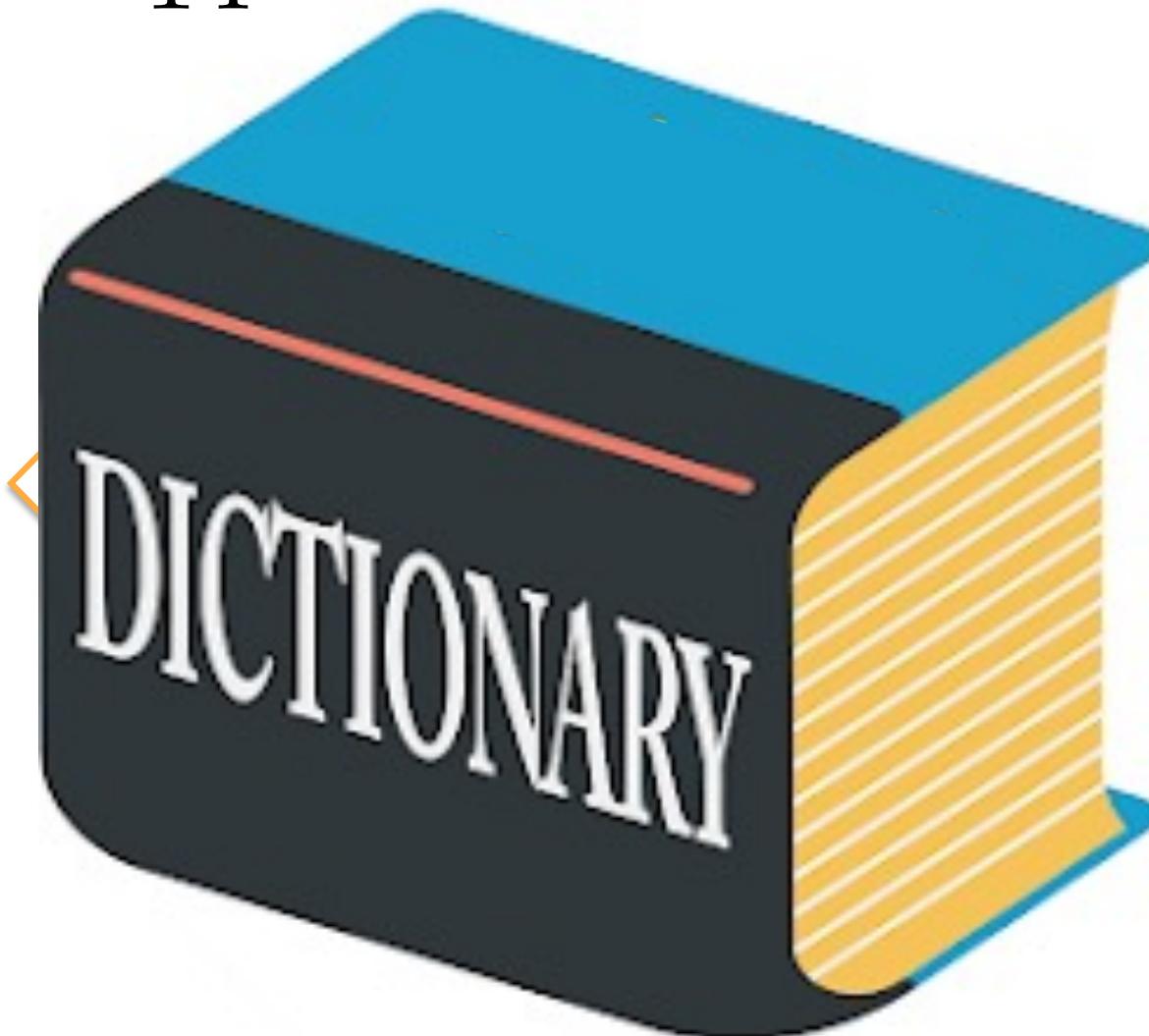
<https://www.reuters.com/markets/deals/constellation-inks-power-supply-deal-with-microsoft-2024-09-20/>

# What is Large Language Model?

“tokens”

axiom →  
ate →  
aspic →  
are →  
anvil →  
antelope →  
ant →  
an →  
am →  
ally →  
abbatoir →  
abate →  
aardvark →  
a →

$$D = 14 \text{ “dictionary size”}$$

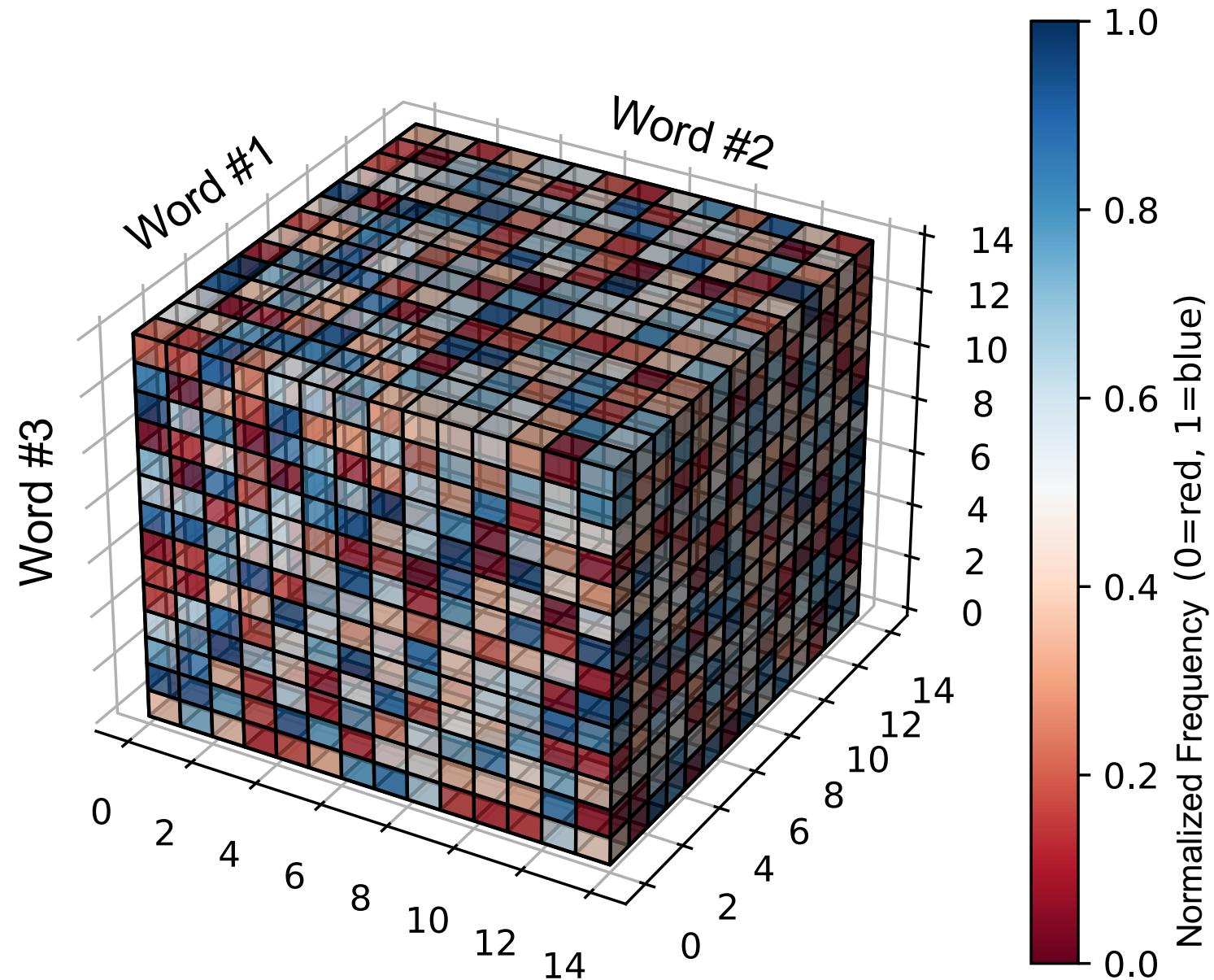


Aardvark ate ant aspic abattoir.  $W=3$  “context window” Aardvark am ant ally. Antelope are anvil axiom. Ant aspic ate, aardvark am. Antelope ally abate abattoir. Ant anvil axiom, aardvark ate. Antelope are aspic ally. Abattoir ant aardvark am. Antelope axiom ally abate. Aardvark ant aspic

“The future state depends only on the past  $W-1$  states.”

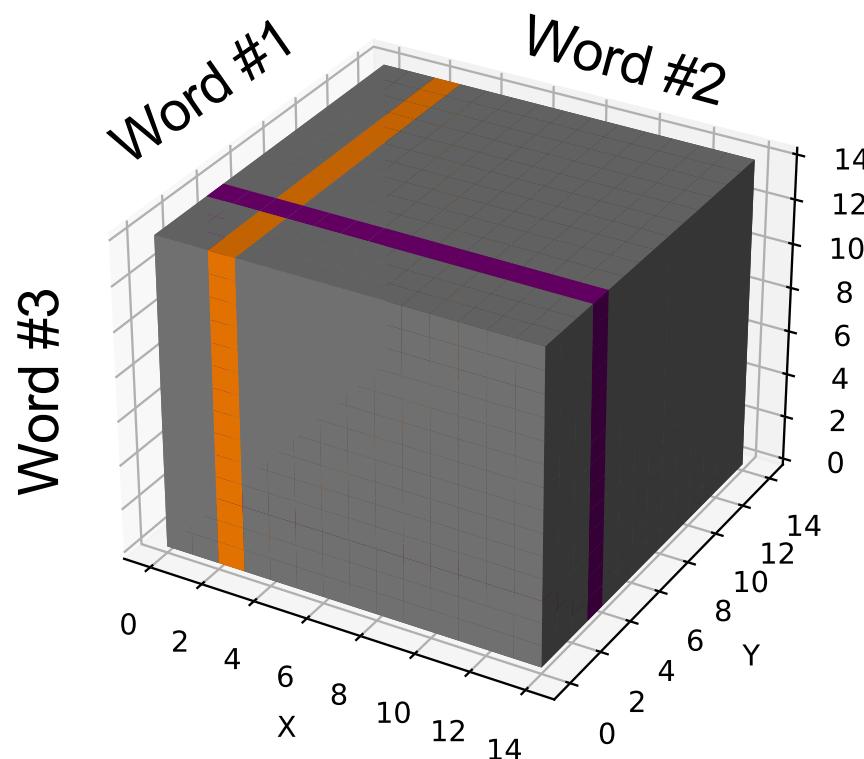


Andrey Markov

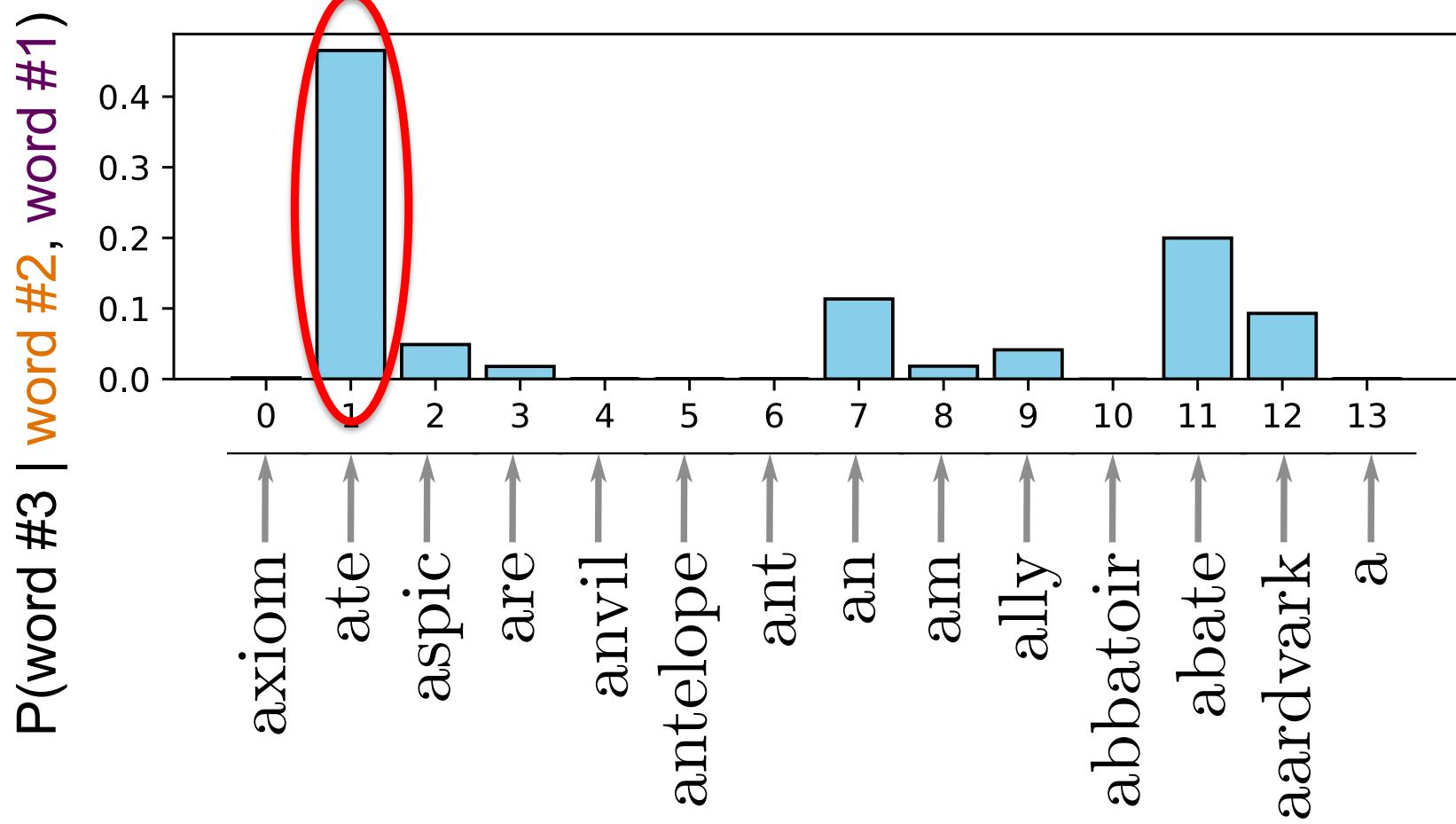


# “an aardvark ?”

$P(\text{word } \#3 \mid \text{word } \#2, \text{word } \#1)$



# “an aardvark ate”



“an aardvark ate”

outside context  
window

an “aardvark ate ?”



auto-regressive step



inference step

$\text{word } \#3 \leftarrow \max_{\text{word } \#3} P(\text{word } \#3 | \text{word } \#2, \text{word } \#1)$

“an aardvark ate”

outside context  
window

an “aardvark ate ?”



auto-regressive step



inference step

an “aardvark ate an”

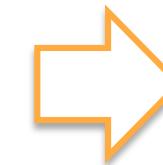
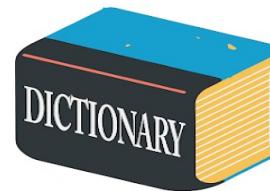
auto-regressive step



inference step

**If this works why do we need AI?**

**“Dictionary Size”**



$$D = 14$$

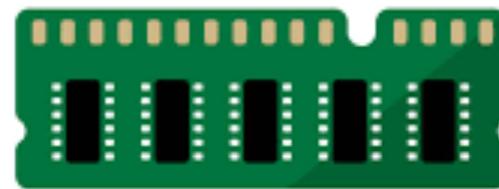
**“Context window”**

Aardvark ate ant



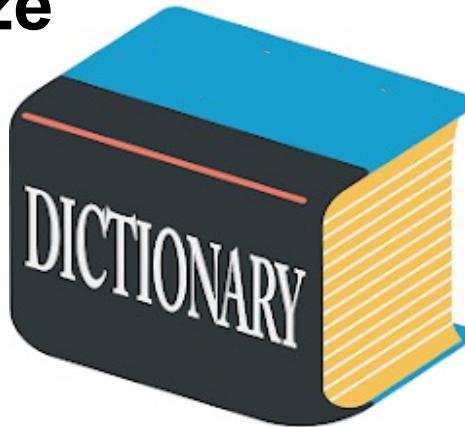
$$W = 3$$

**“Memory”**



$$EB^W K_b$$

## “GPT-3 Dictionary Size”



$$\rightarrow D = 50,000$$

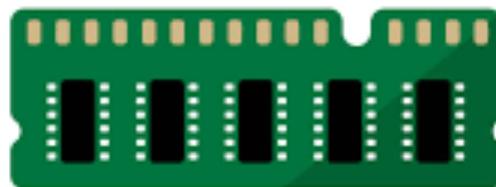
## “GPT-3 Context window”

Mind! I don't mean to say that I know, of my own knowledge, what there is particularly dead about a door-nail. I might have been inclined, myself, to regard a coffin-nail as the deadeast piece of ironmongery in the trade. But the wisdom of our ancestors is in the simile; and my unhallowed hands shall not disturb it, or the Country's done for. You will therefore permit me to repeat, emphatically, that Marley was as dead as a door-nail.

$$\rightarrow W = 2048$$

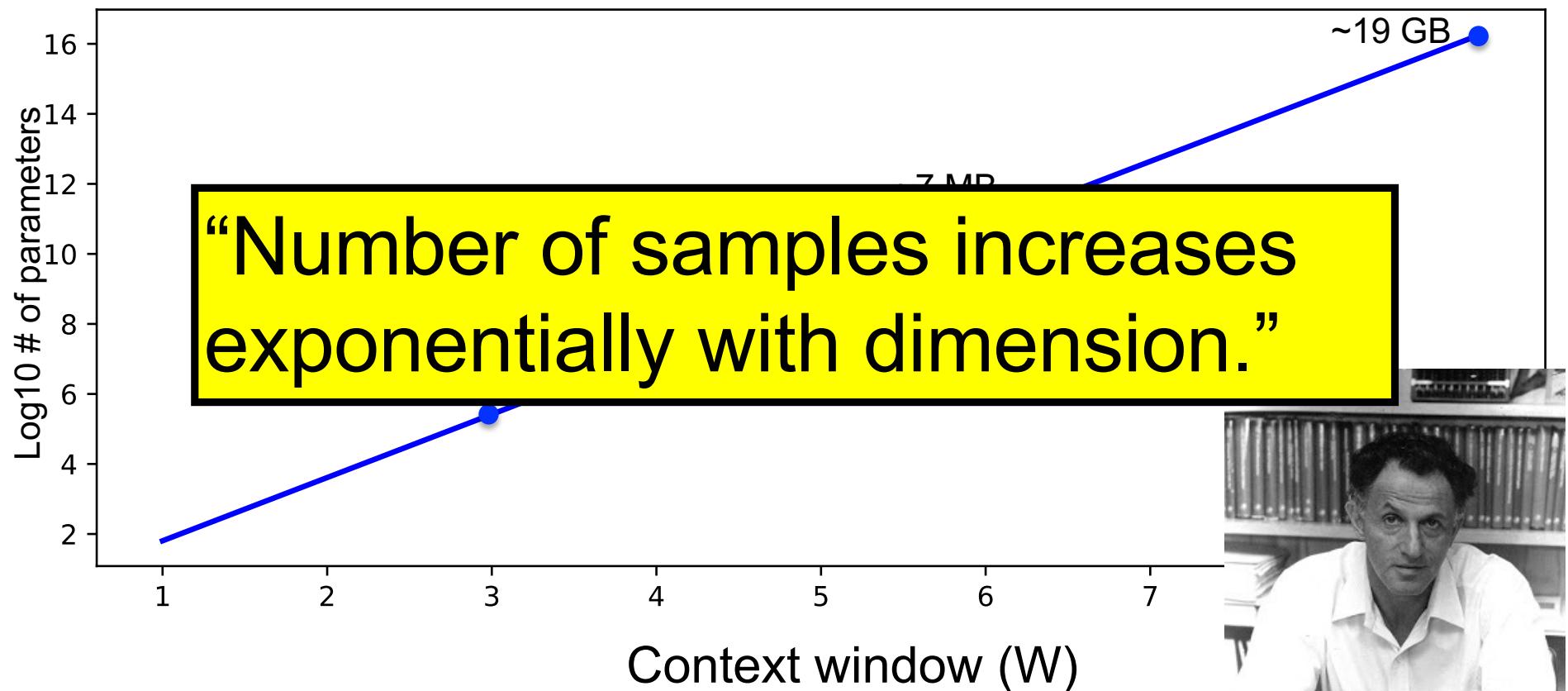
Only  $10^{80}$  atoms in the universe!!

## “Memory”



$$\rightarrow 10^{9600} \text{ bytes}$$

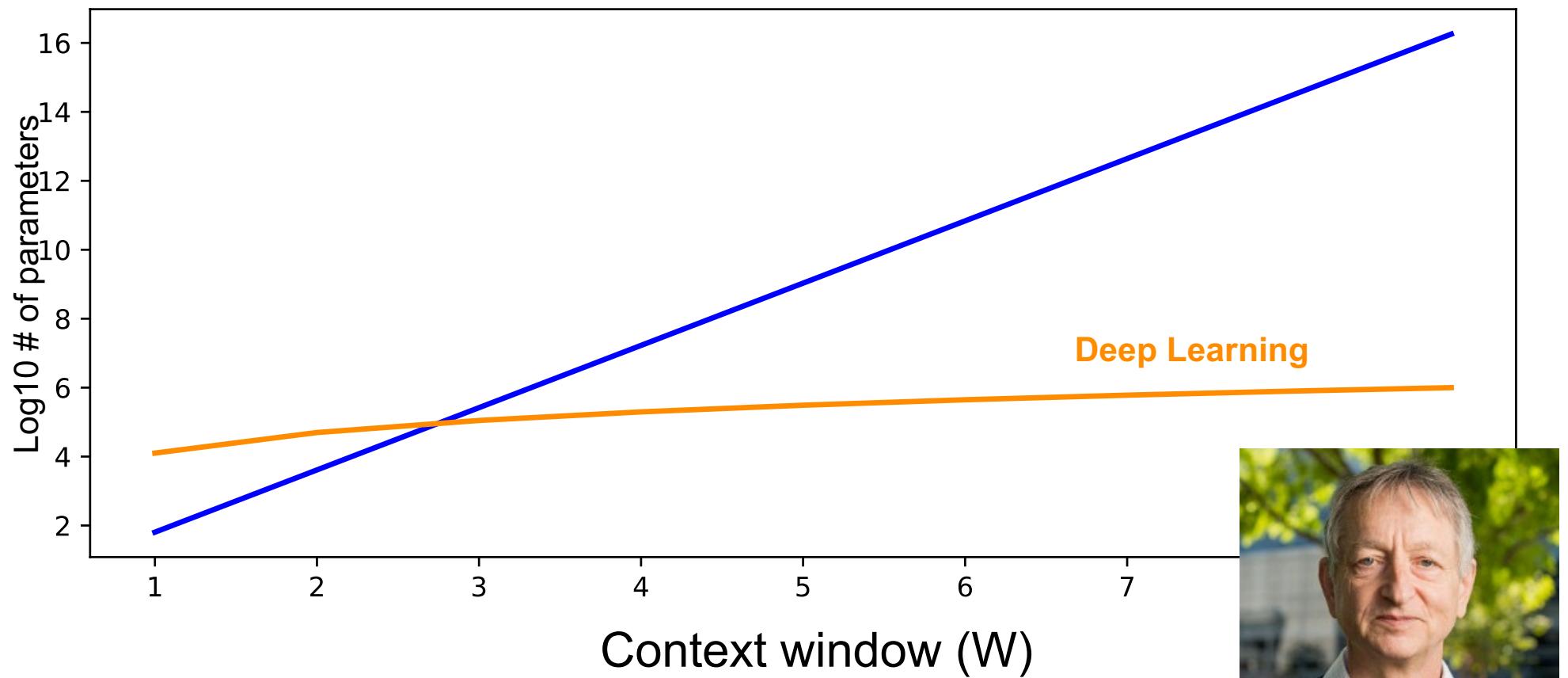
# Curse of Dimensionality ( $D = 14$ )



$$\# \text{ of parameters} = D^W$$

Richard E. Bellman

# Deep Learning – A breakthrough in sampling?



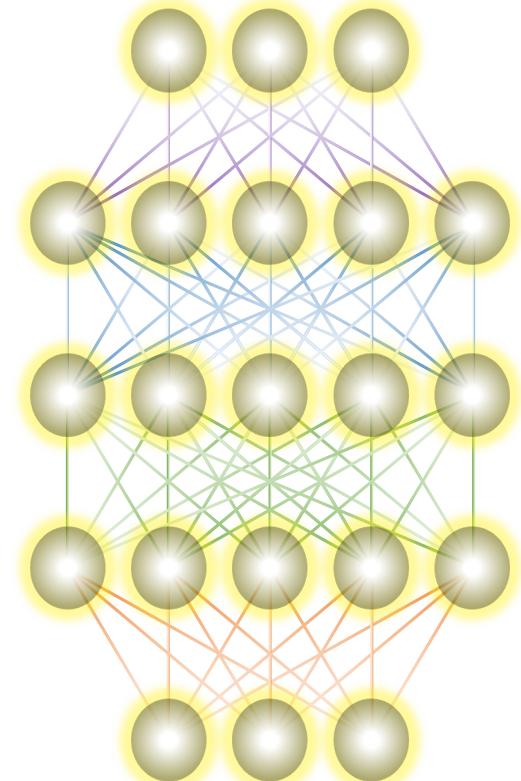
$$\# \text{ of parameters} = D^W$$



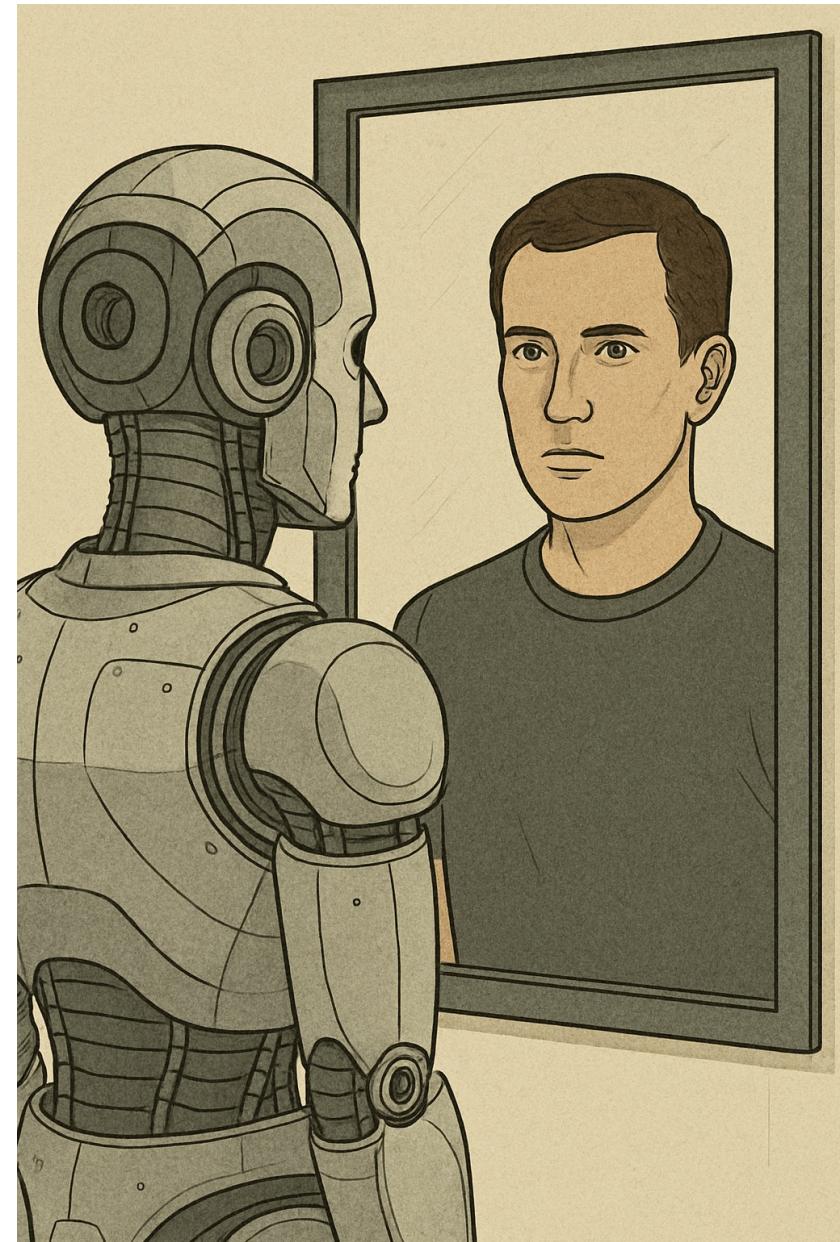
"Geoffrey Hinton"

# Deep Learning – Current Open Questions?

- **Why does it generalize so well?**
- **Why does it not get caught in local minima?**
- **What are the limits of scaling laws?**
- **Why does it hallucinate?**
- **How can we train more efficiently?**



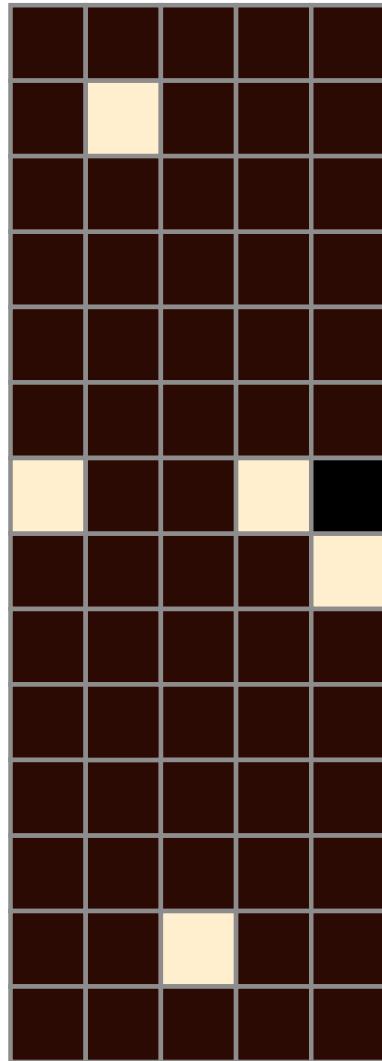
- Is AI just a reflection of our own intelligence?



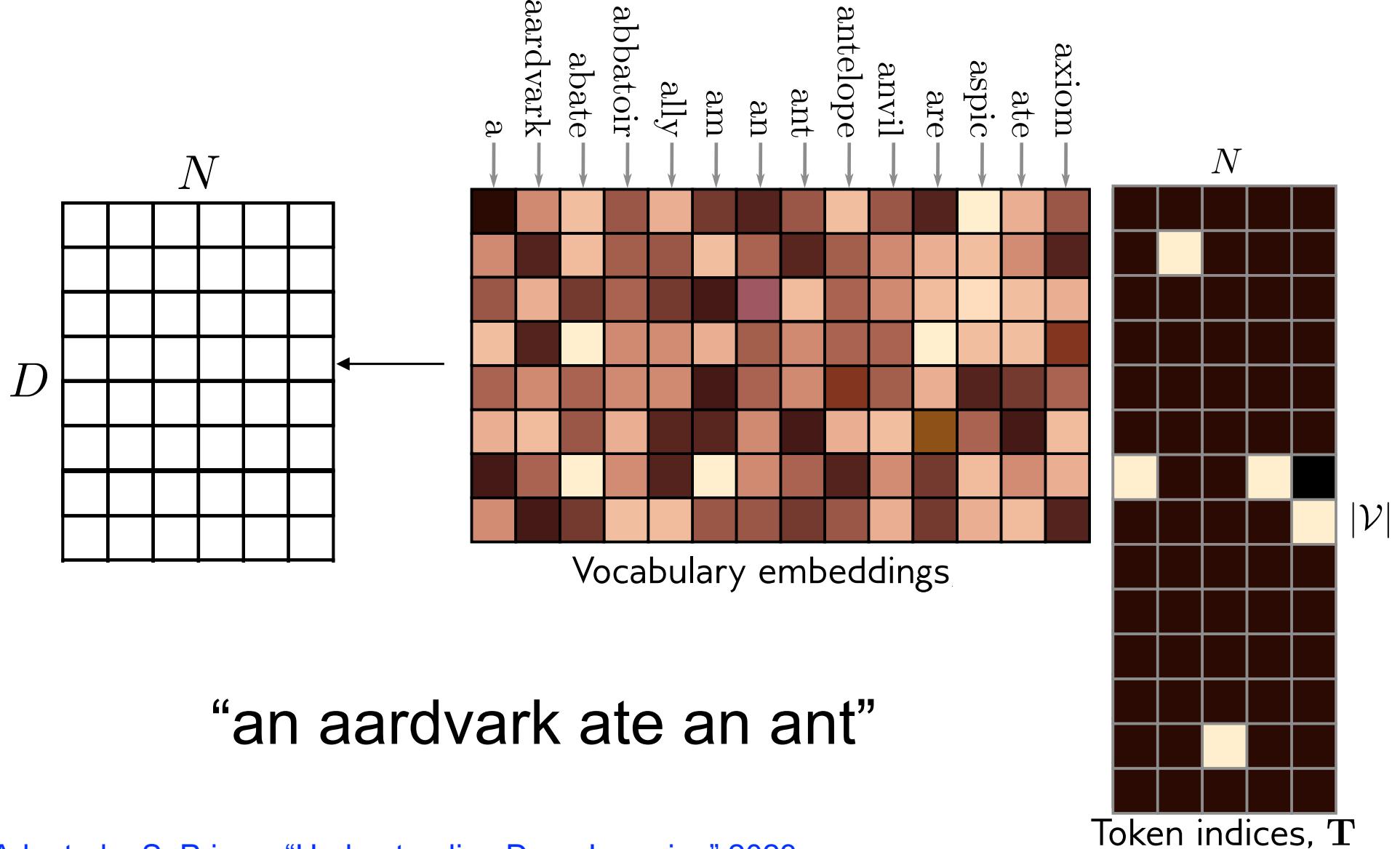
Efficiency through  
Understanding

“an aardvark ate an ant”

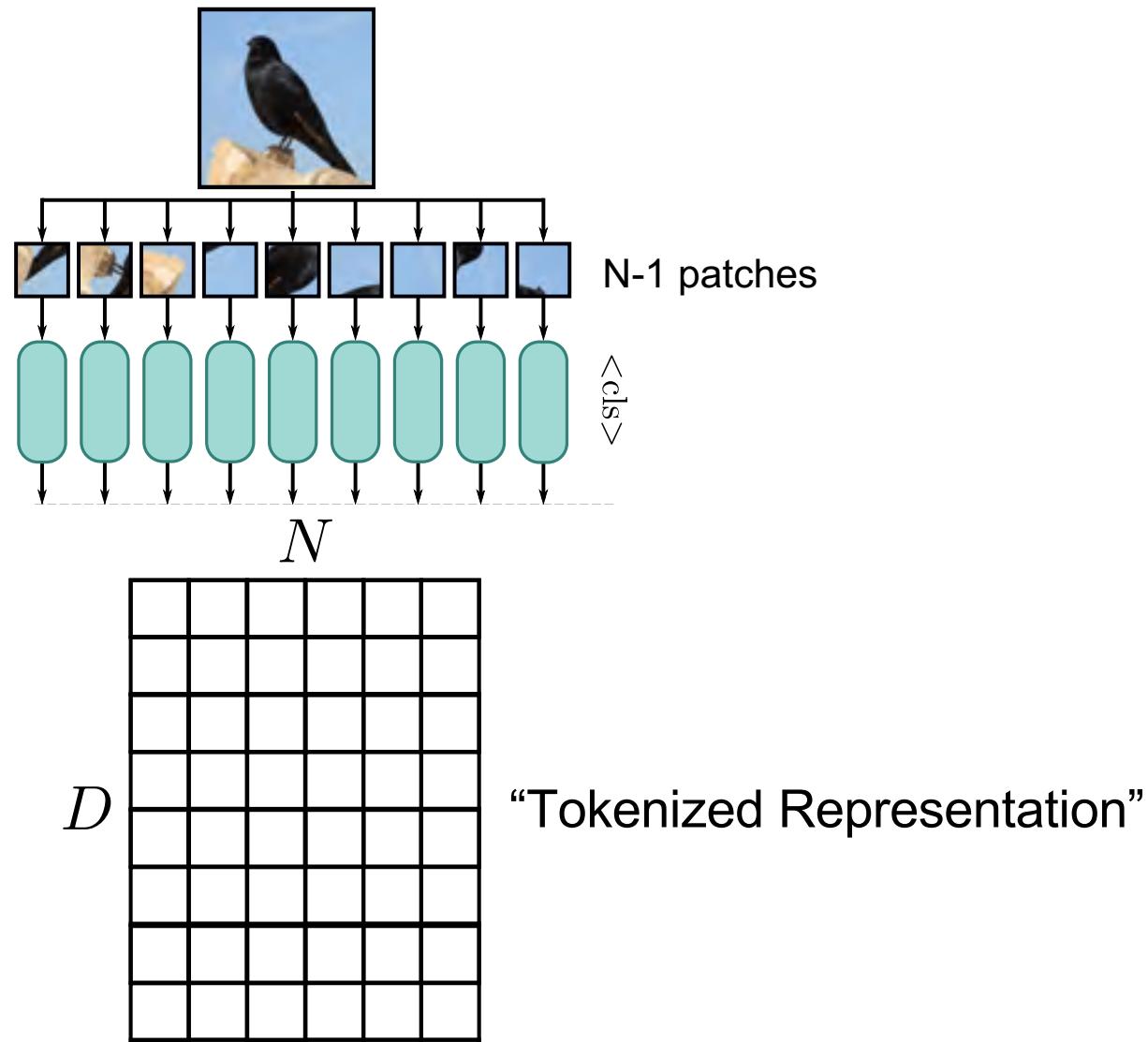
axiom →  
ate →  
aspic →  
are →  
anvil →  
antelope →  
ant →  
an →  
am →  
ally →  
abbatoir →  
abate →  
aardvark →  
a →



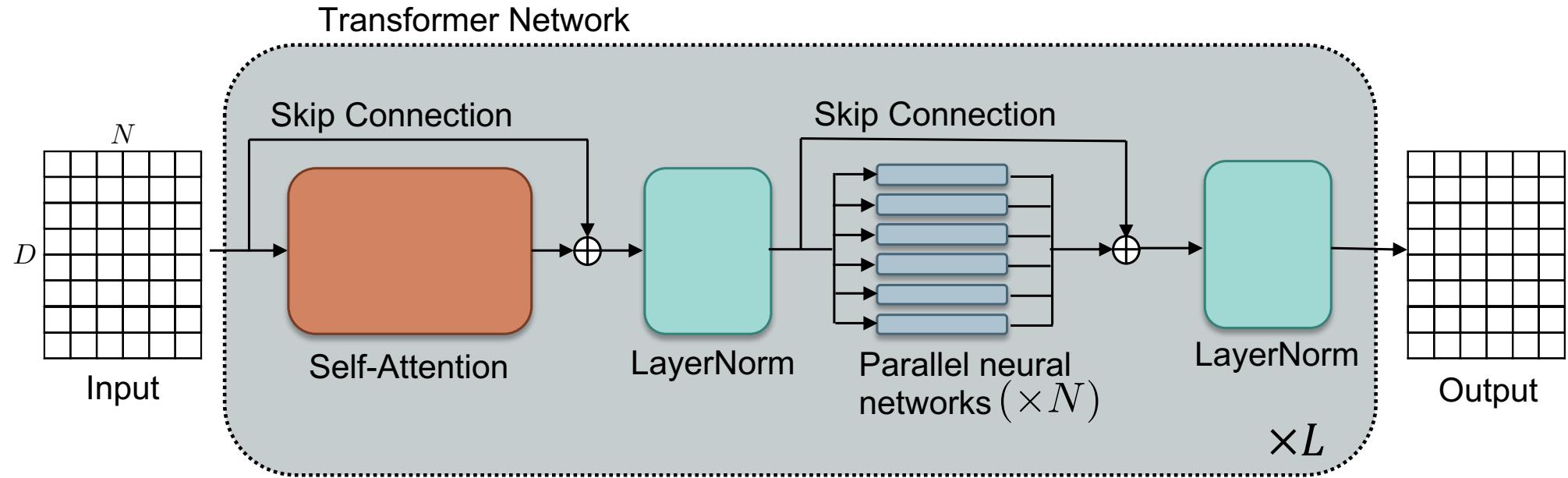
“an aardvark ate an ant”

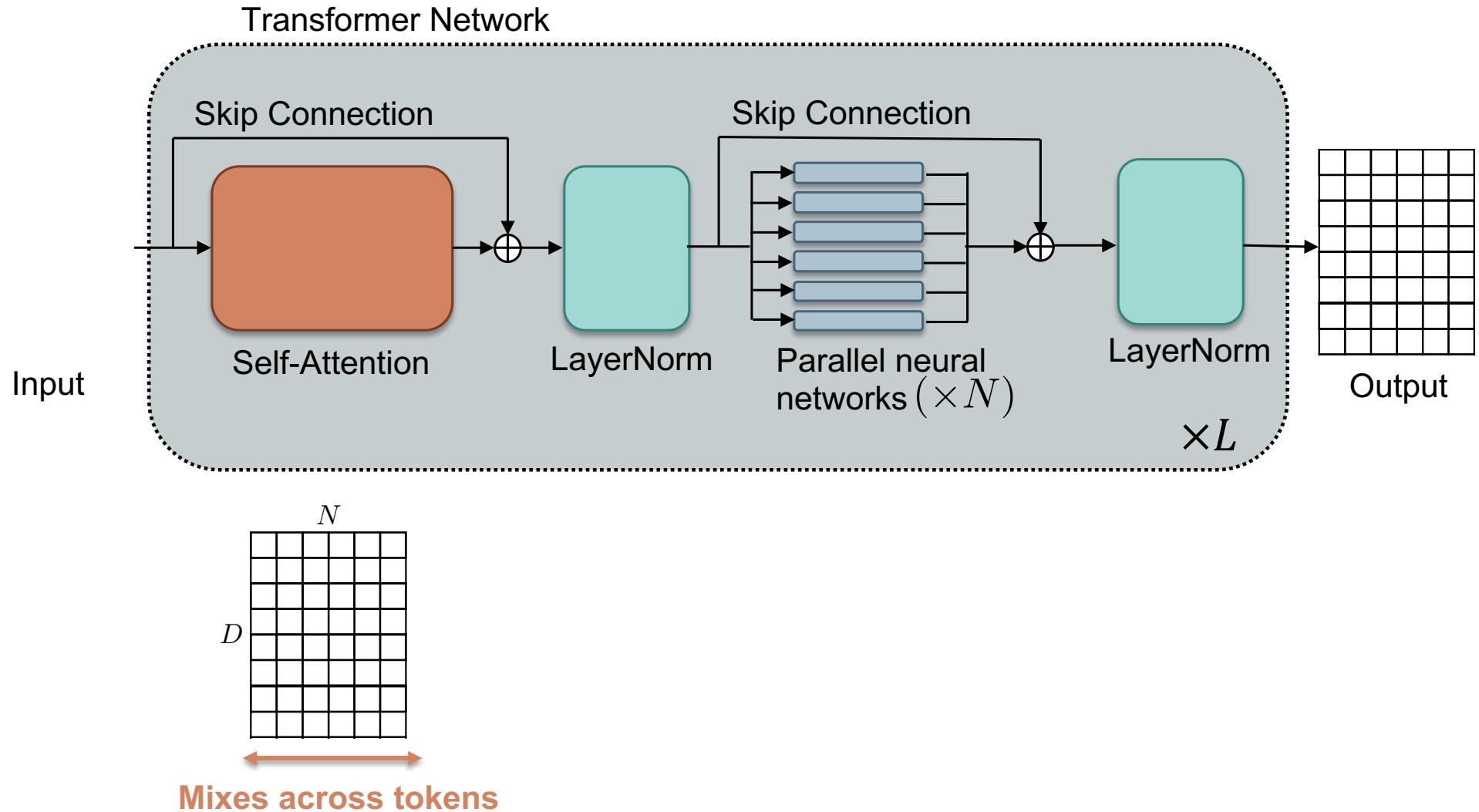


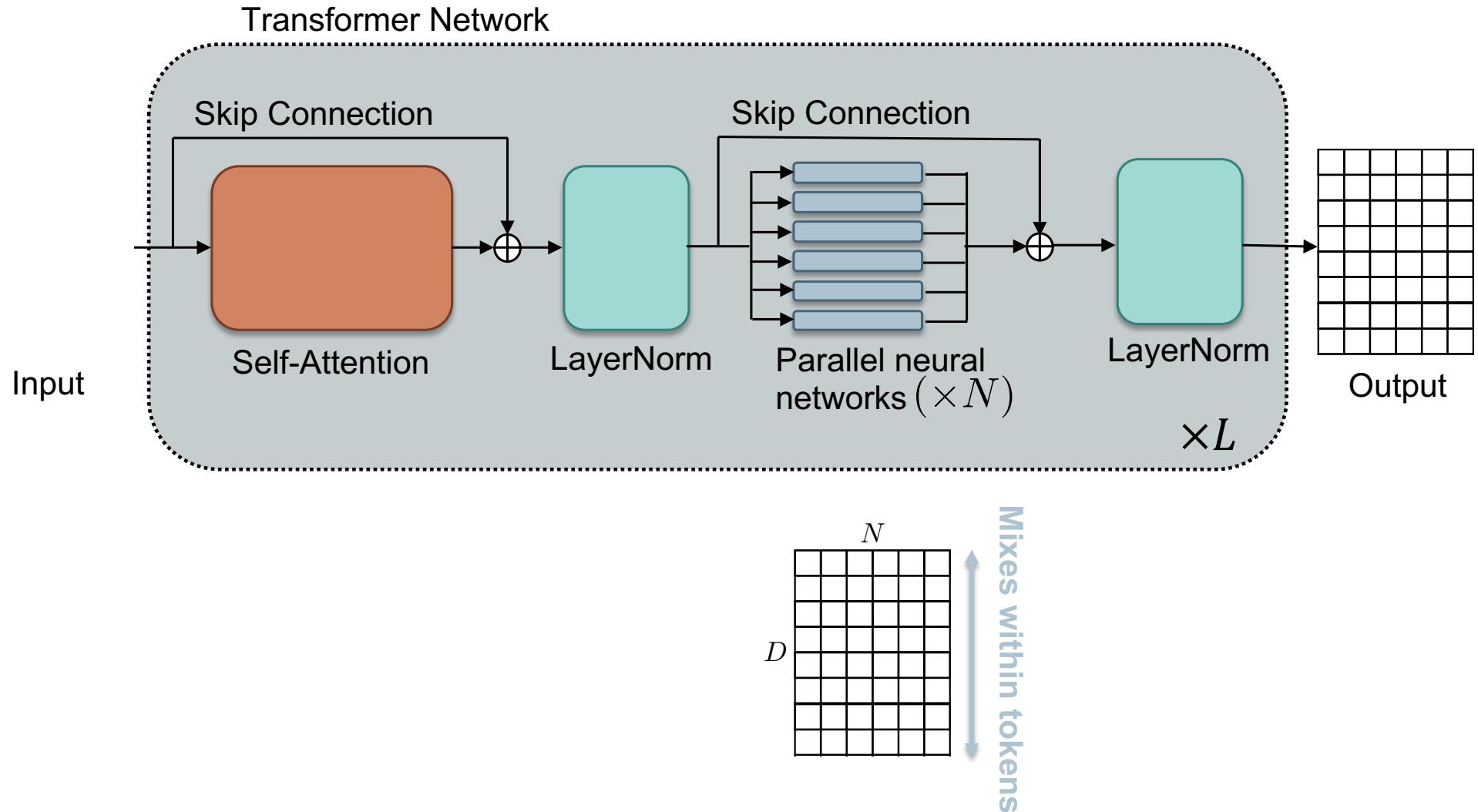
Adapted – S. Prince. “Understanding Deep Learning” 2023.

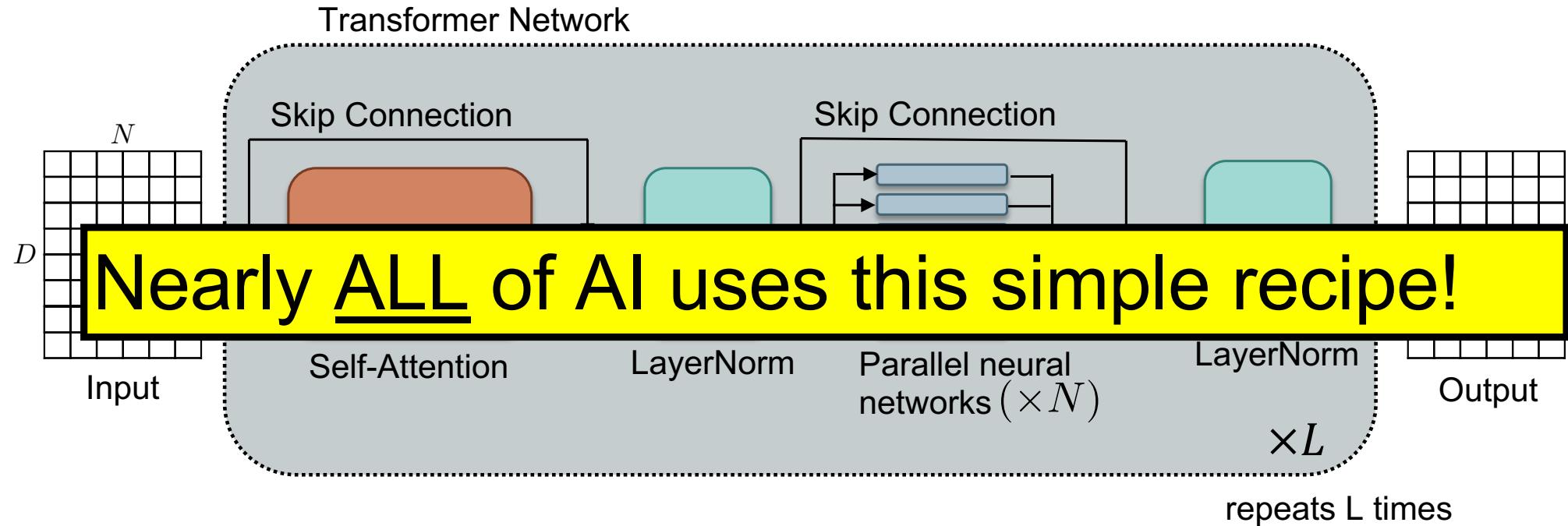


Adapted – S. Prince. “Understanding Deep Learning” 2023.

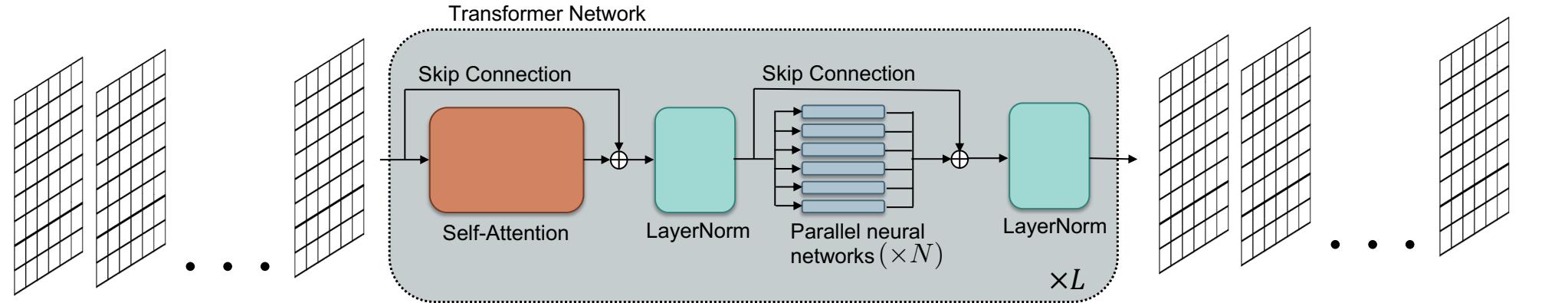








**Can I know in advance how much it costs to train?**



$$\mathcal{L}(\theta) = \sum_{m=1}^M \mathcal{L}\{\mathbf{y}_m; f(\mathbf{x}_m; \theta)\}$$

Minimize loss

$$\mathcal{L}(\theta) = \sum_{m=1}^M \mathcal{L}\{\mathbf{y}_m; f(\mathbf{x}_m; \theta)\}$$


Gradient Descent

$$\theta_{k+1} \rightarrow \theta_k - \alpha \cdot \nabla \mathcal{L}(\theta_k)$$

“Learning Rate”

$$\nabla \mathcal{L}(\boldsymbol{\theta}_k) = [\nabla f(\mathbf{x}_1; \boldsymbol{\theta}_k), \dots, \nabla f(\mathbf{x}_M; \boldsymbol{\theta}_k)] \begin{bmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}_k) - \mathbf{y}_1 \\ \vdots \\ f(\mathbf{x}_M; \boldsymbol{\theta}_k) - \mathbf{y}_M \end{bmatrix}$$

Jacobian matrix error vector

$$\nabla \mathcal{L}(\theta_k) = \mathbf{J}_k \boldsymbol{\epsilon}_k$$

Jacobian matrix error vector

$$\mathcal{L}(\theta) = \sum_{m=1}^M \frac{1}{2} \|\mathbf{y}_m - f(\mathbf{x}_m; \theta)\|_2^2$$

“Least-Squares Objective”



Gradient Descent

$$\theta_{k+1} \rightarrow \theta_k - \alpha \cdot \mathbf{J}_k \epsilon_k$$

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{n=1}^N \frac{1}{2} \|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|_2^2$$


Gradient Descent

$$\boldsymbol{\theta}_{k+1} \rightarrow \boldsymbol{\theta}_k - \alpha \cdot \mathbf{J} \boldsymbol{\epsilon}_k$$

“Static Jacobian”

If we start with  $\theta_0 = 0$

$$\theta_1 \rightarrow -\alpha \cdot \mathbf{J} \boldsymbol{\epsilon}_0$$

$$\theta_2 \rightarrow -(\mathbf{I} - \alpha \mathbf{J} \mathbf{J}^T) \alpha \mathbf{J} \boldsymbol{\epsilon}_0 - \alpha \mathbf{J} \boldsymbol{\epsilon}_0$$

$$\theta_3 \rightarrow -(\mathbf{I} - \alpha \mathbf{J} \mathbf{J}^T)^2 \alpha \mathbf{J} \boldsymbol{\epsilon}_0 - (\mathbf{I} - \alpha \mathbf{J} \mathbf{J}^T) \alpha \mathbf{J} \boldsymbol{\epsilon}_0 - \alpha \mathbf{J} \boldsymbol{\epsilon}_0$$

at the  $q$ -th iteration generalizes to,

$$\theta_q = - \sum_{k=0}^q (\mathbf{I} - \alpha \mathbf{J} \mathbf{J}^T)^{k-1} \alpha \mathbf{J} \boldsymbol{\epsilon}_0$$

$$\theta_q = - \sum_{k=0}^q (\mathbf{I} - \alpha \mathbf{J} \mathbf{J}^T)^{k-1} \alpha \mathbf{J} \boldsymbol{\epsilon}_0$$



$$\mathbf{J} = \mathbf{U} \operatorname{diag}(\mathbf{S}) \mathbf{V}^T \text{ “SVD”}$$



$$\theta_q = - \mathbf{U} \operatorname{diag}(\hat{\mathbf{S}}_q) \mathbf{V}^T \boldsymbol{\epsilon}_0 \text{ “closed form solution”}$$

$$\theta_q = - \mathbf{U} \operatorname{diag}(\hat{\mathbf{s}}_q) \mathbf{V}^T \boldsymbol{\epsilon}_0$$

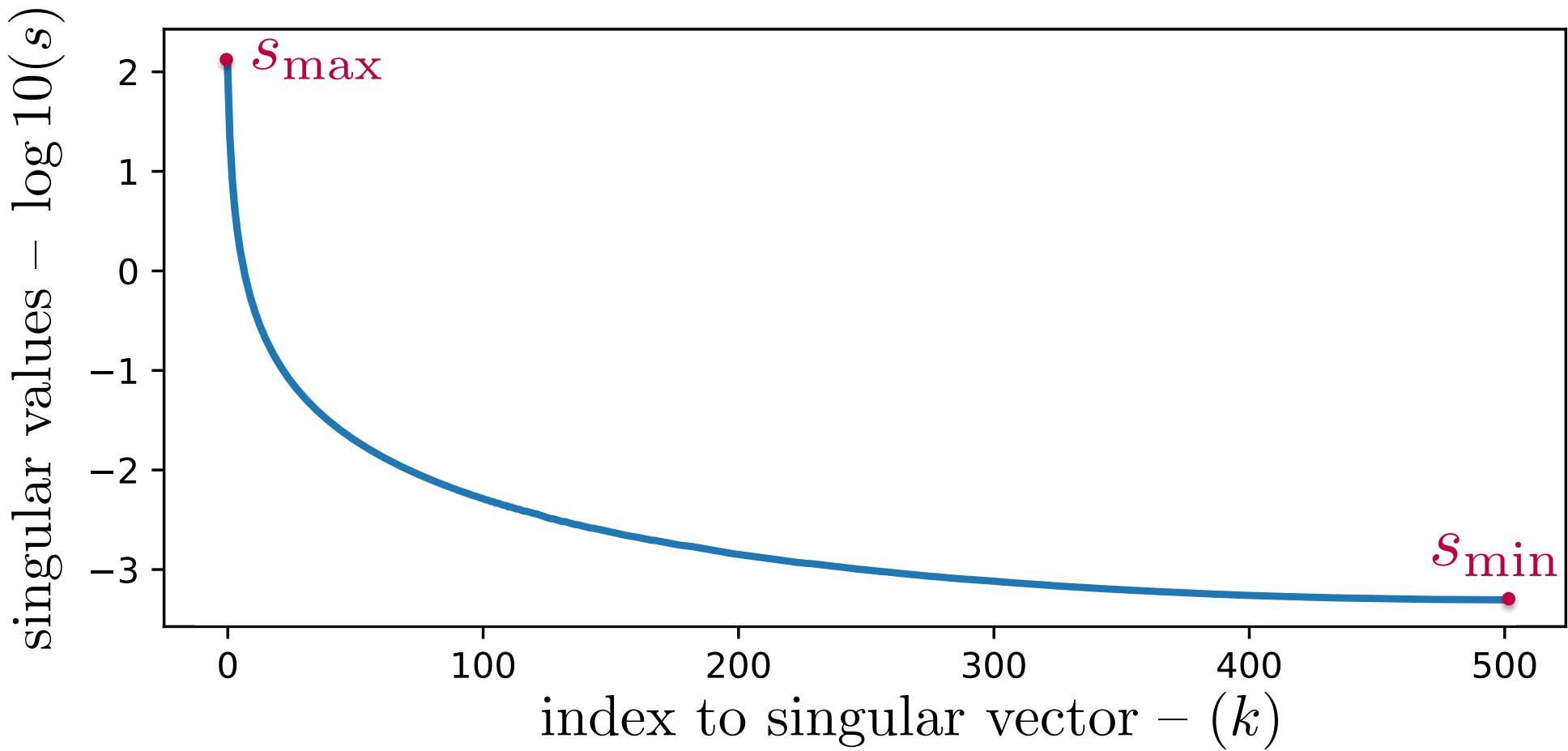


$$\hat{\mathbf{s}}_q[k] = \frac{1 - \left( \sum_{j=0}^{q-1} \alpha \mathbf{s}[k]^2 \right)^q}{\mathbf{s}[k]} \quad \text{"GD shrinkage"}$$

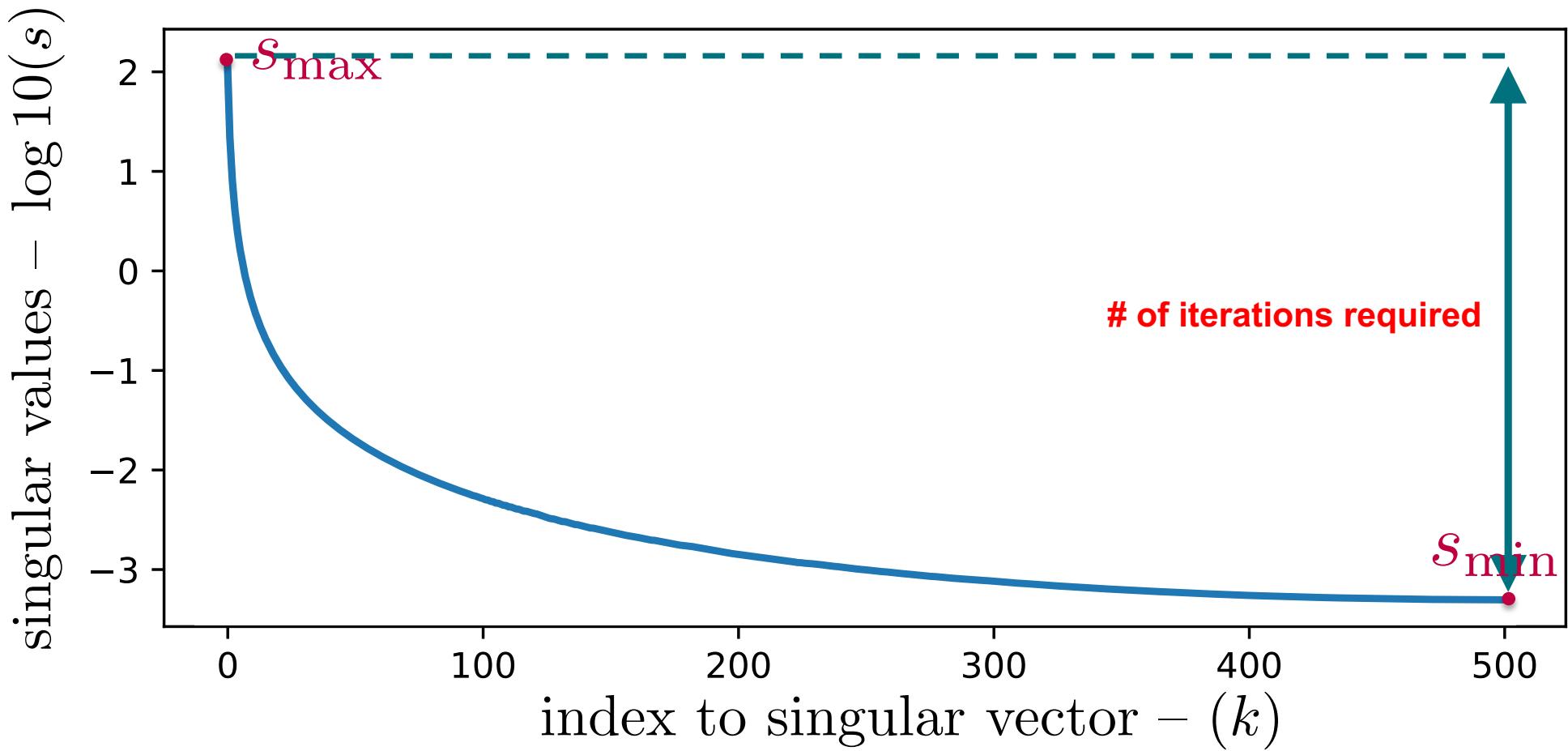
$$(1 - \alpha \mathbf{s}[k]^2)^q \approx \exp(-q \alpha \mathbf{s}[k]^2)$$

A. Jacot et al. “Neural Tangent Kernel: Convergence and Generalization in Neural Networks” in NeurIPS 2018.  
S. Lucey. “Gradient Descent as a Shrinkage Operator for Spectral Bias” in arXiv 2025.

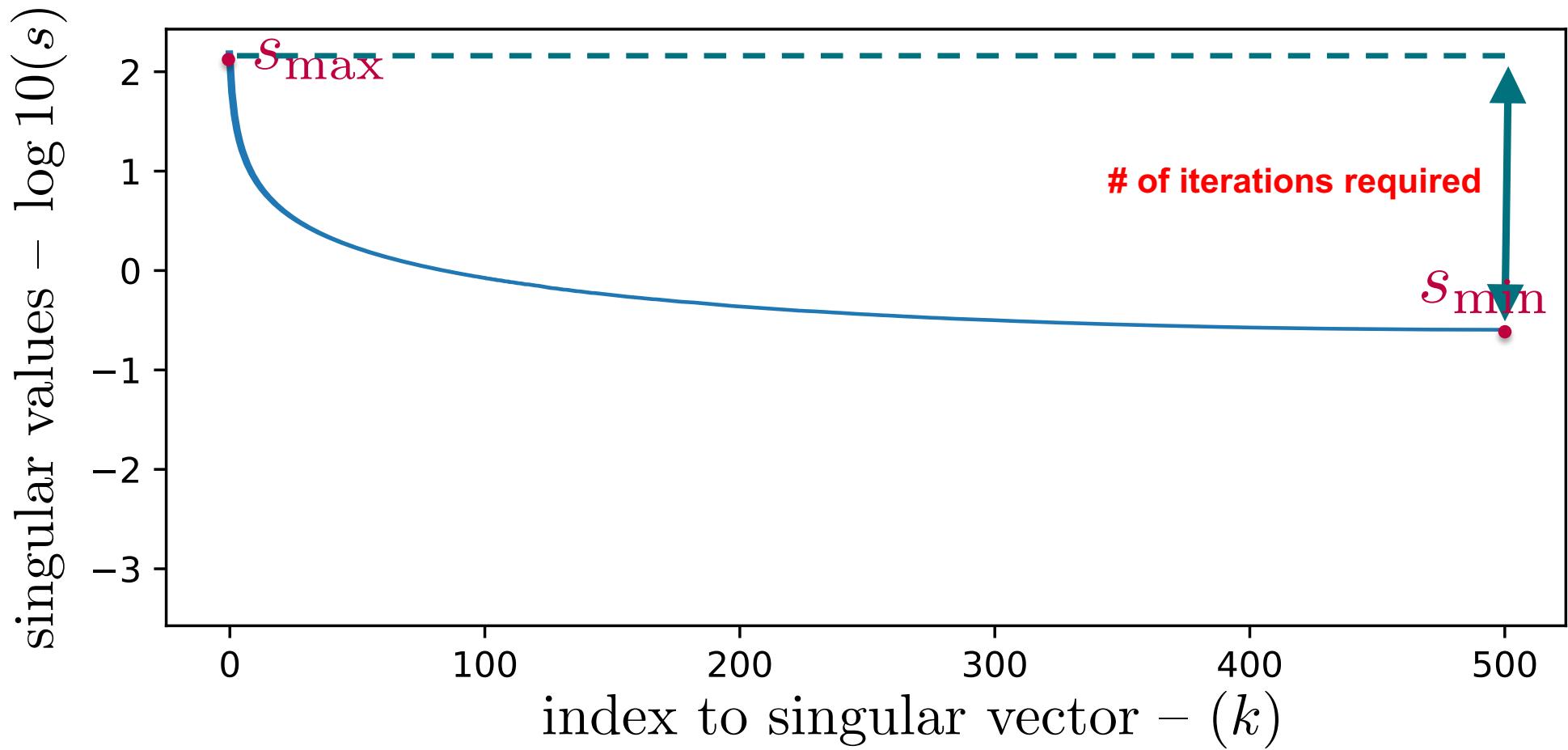
$$\text{cond}(\mathbf{J}) = \frac{s_{\max}^2}{s_{\min}^2}$$



$\text{cond}(\mathbf{J}) \propto \# \text{ of iterations required in GD}$

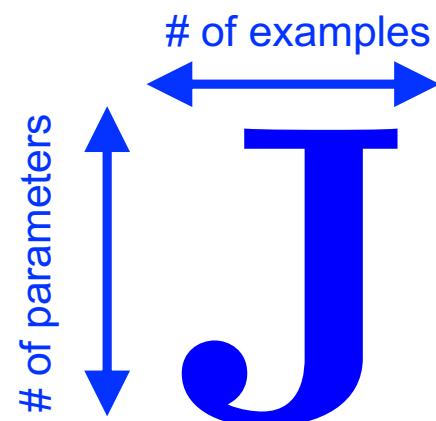


$\text{cond}(\mathbf{J}) \propto \# \text{ of iterations required in GD}$

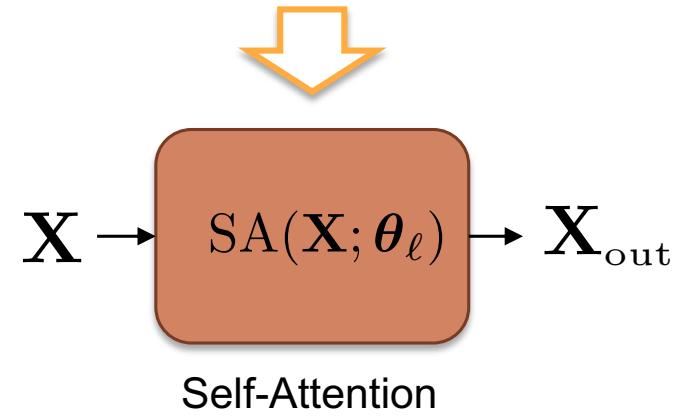
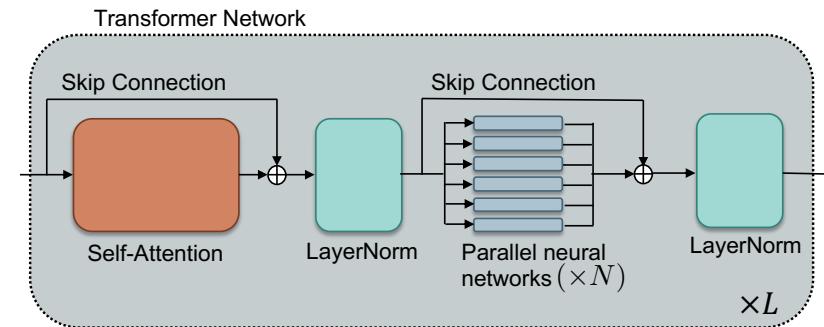


# Insights....

- Condition of the Jacobian is a useful heuristic even for SGD and Adam.
- However, problematic forming the Jacobian for modern large-scale transformers.
- Could we decompose the Jacobian?



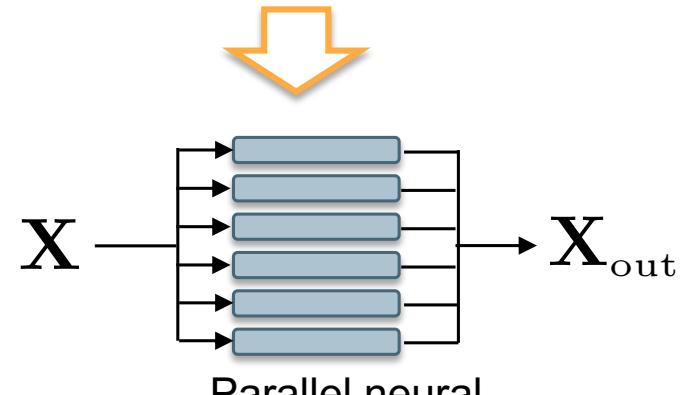
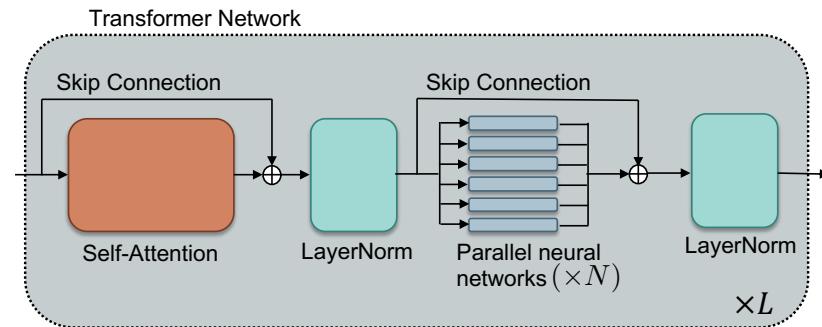
$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_0 \\ \vdots \\ \boxed{\mathbf{J}_\ell} \\ \vdots \\ \mathbf{J}_{\ell+1} \\ \vdots \\ \mathbf{J}_{2L-1} \end{bmatrix}$$



$$\boxed{\mathbf{J}_\ell} = \nabla_{\mathbf{X}} \text{SA}(\mathbf{X}; \theta_\ell)$$

Jacobian of Self-Attention Block (SAB)

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_0 \\ \vdots \\ \mathbf{J}_{\ell} \\ \boxed{\mathbf{J}_{\ell+1}} \\ \vdots \\ \mathbf{J}_{2L-1} \end{bmatrix}$$



Parallel neural  
networks ( $\times N$ )

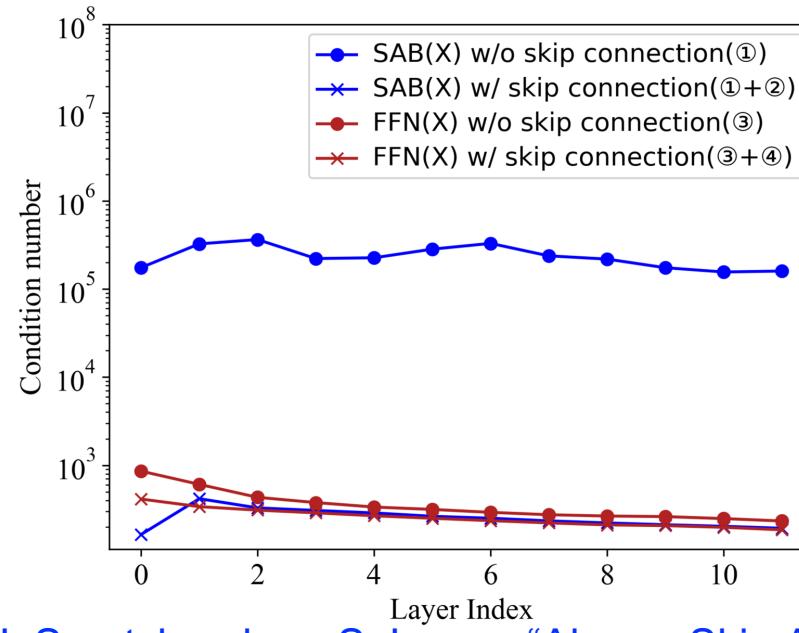
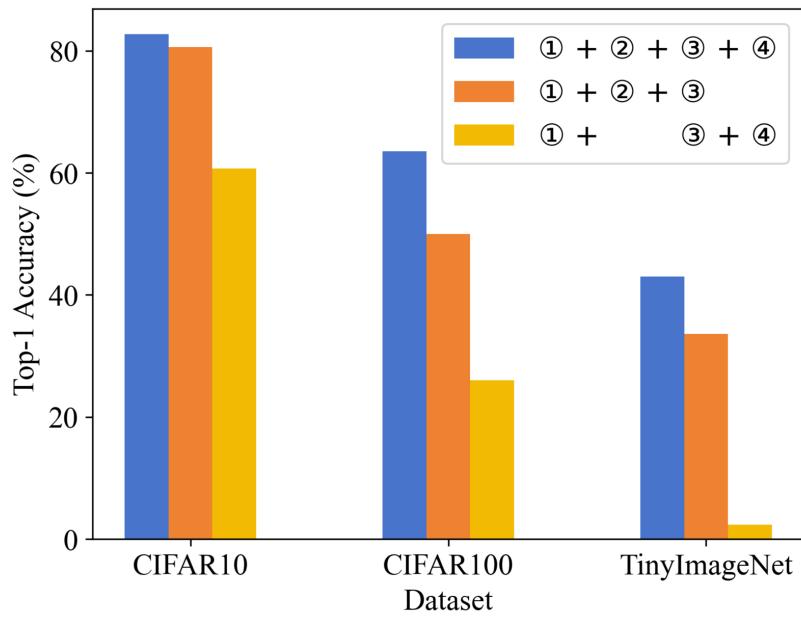
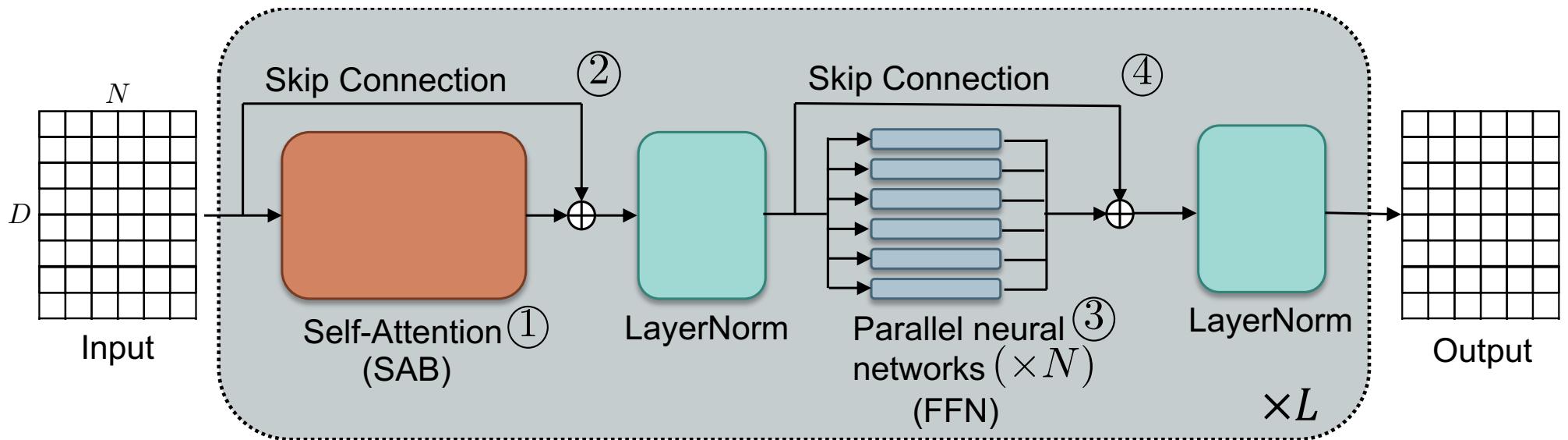
$$\boxed{\mathbf{J}_{\ell+1}} = \nabla \text{FFN}(\mathbf{X}; \theta_{\ell+1})$$

Jacobian of Feed Forward Network (FFN)

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_0 \\ \vdots \\ \mathbf{J}_\ell \\ \mathbf{J}_{\ell+1} \\ \vdots \\ \mathbf{J}_{2L-1} \end{bmatrix}$$

**Simplifying assumption:**

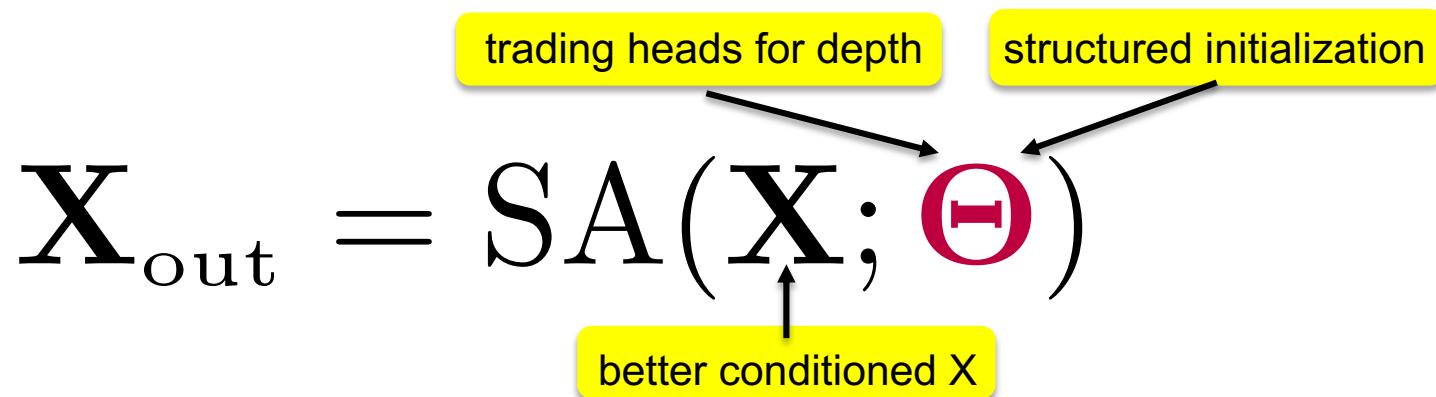

$$\text{cond}(\mathbf{J}) \leq \max_\ell \text{cond}(\mathbf{J}_\ell)$$



Y. Ji, H. Saratchandran, S. Lucey, "Always Skip Attention." In ICCV 2025.

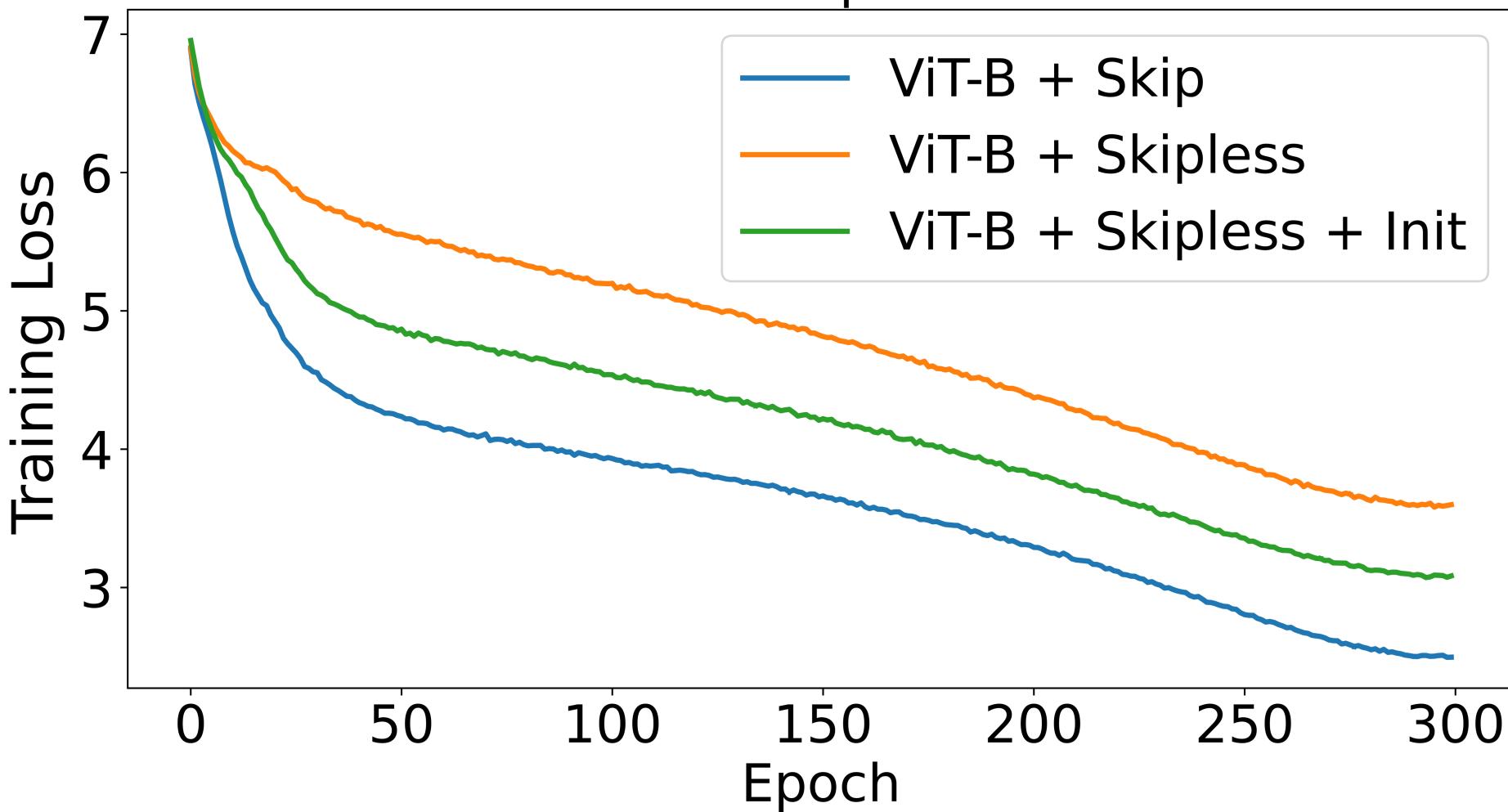
# New Directions....

- Designing transformers to maximize sub-block condition in self-attention has opened up brand new-directions.



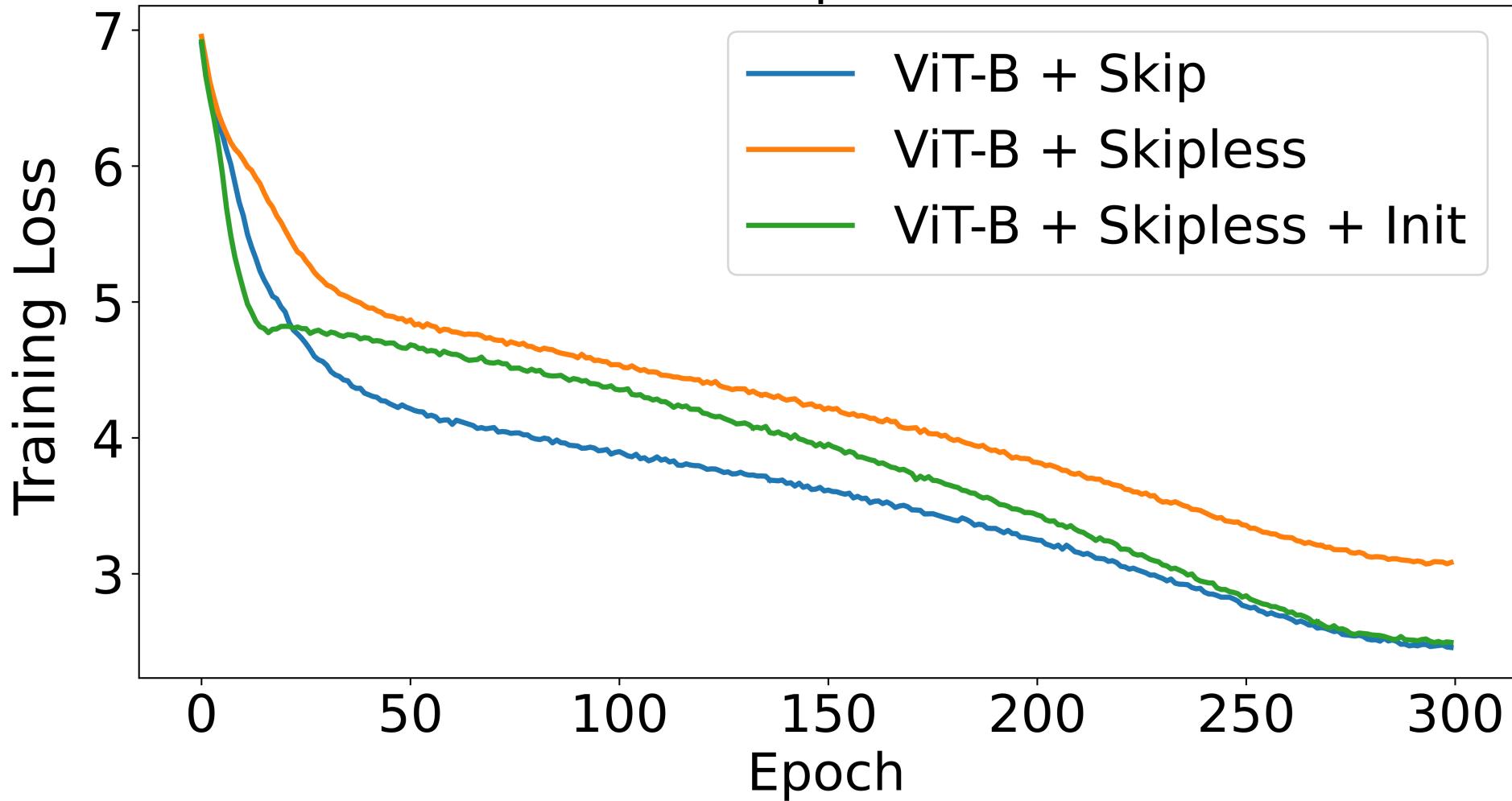
- H. Saratchandran, D. Teney, S. Lucey. "Leaner Transformers: More Heads Less. Depth." In arXiv 2025.  
J. Zheng, X. Li, H. Saratchandran, S. Lucey. "Structured Initialization for Vision Transformers." In NeurIPS 2025.  
Y. Ji, S. Lucey, et al. "Cutting the skip: Training residual-free transformers." In ICLR 2026.  
Y. Ji, H. Saratchandran, P. Moghadam, and S. Lucey. "Always Skip Attention." In ICCV 2025.

## AdamW optimizer



Y. Ji, S. Lucey, et al. “Cutting the skip: Training residual-free transformers.” In arXiv 2025.

## SOAP optimizer



Y. Ji, S. Lucey, et al. “Cutting the skip: Training residual-free transformers.” In arXiv 2025.

|                        | VOC2012     |                 |                 | COCO20k     |                 |                 |
|------------------------|-------------|-----------------|-----------------|-------------|-----------------|-----------------|
| Epoch →<br>Optimizer ↓ | 300<br>skip | 300<br>skipless | 200<br>skipless | 300<br>skip | 300<br>skipless | 200<br>skipless |
| AdamW                  | 32.3        | 53.5            | <b>54.0</b>     | 21.2        | 36.5            | <b>38.5</b>     |
| SOAP                   | 49.4        | 63.2            | <b>68.1</b>     | 27.5        | 46.7            | <b>54.1</b>     |

## Object Discovery



## Object Discovery

Y. Ji, S. Lucey, et al. “Cutting the skip: Training residual-free transformers.” In arXiv 2025.

# Conclusions....

- Self-attention is fundamentally poorly conditioned, increasing the cost and complexity of training modern transformers.
- Numerous interventions are possible, but in some circumstances it is best to “skip” attention all together.
- Acts as a reminder that the field of AI could be up-ended by better understanding the foundational principles of modern deep learning.