

Transformers As Effective Fields: From Quantum Physics to AI

TropicalMathML Seminar, Yonsei University
February 23, 2026

Happy new year of the horse!



Changqing Fu
PSL University

Generative AI through Transformers



Image

what is the core argument?

PDF TransformerPriGM.pdf

Thought for 32s

The core argument of this work (the poster for the paper "*Transformers Are Optimal Effective Fields*" by Changqing Fu, PriGM@EurIPS 2025) is:

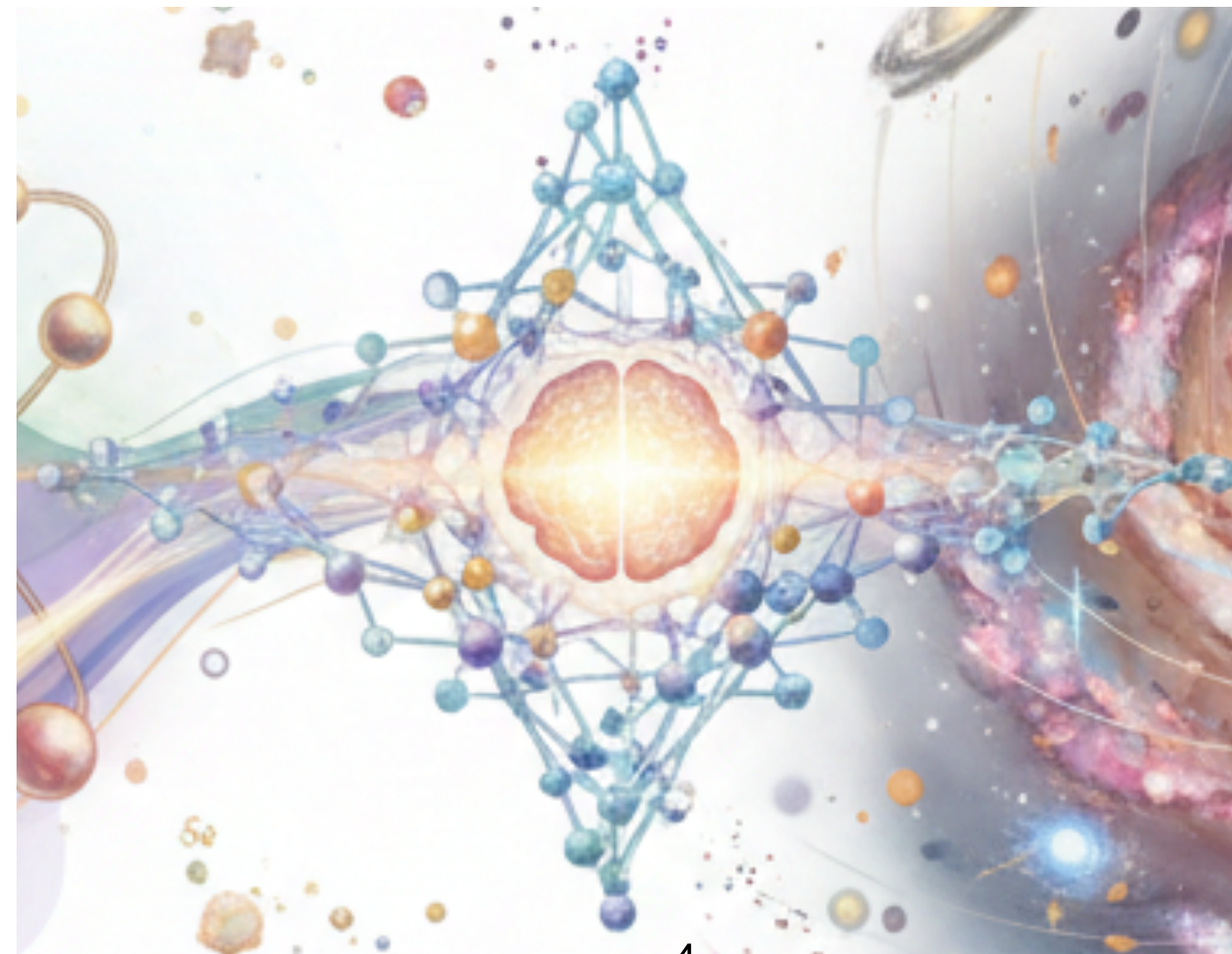
Transformers are not a heuristic design—they are the mathematically optimal ("geodesic") architecture that emerges automatically when you treat neural network layers as continuous-time evolution of a physical particle/field system and demand the least geometric structure necessary.

Text

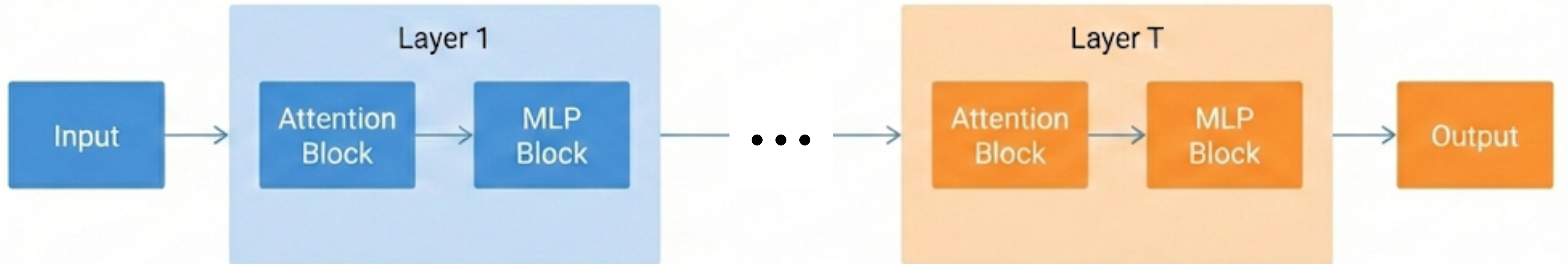
Plan

- Transformer and Least Action
 - Matrix Algebra, HighDim Geometry, MLP / SelfAttention
- Quantum Chaos / Statistical Physics
 - RMT, Double Descent / Scaling Law, Interactive Fields
- From QC to Transformers
 - Left/Right Symmetries, Real/Complex Multiplication, Predictions

What is a Transformer



“Attention Is All You Need”



$$\mathbf{X}(t+1) = \mathbf{X}(t) + \text{Attention}(\mathbf{X}(t))$$

$$\mathbf{X}(t+2) = \mathbf{X}(t+1) + \text{MLP}(\mathbf{X}(t+1))$$

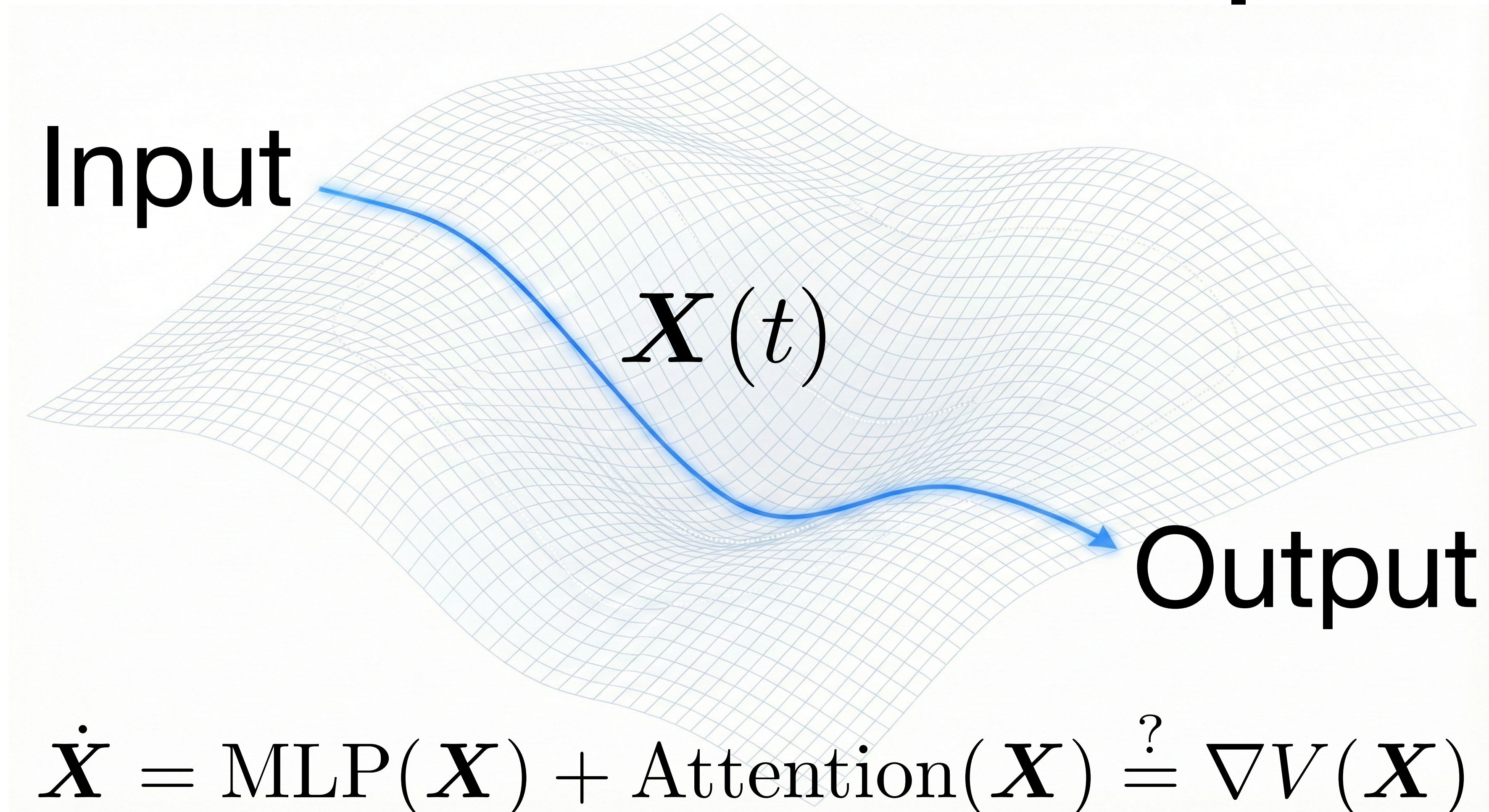
Physics = Symmetry Breaking Principle



“Asymmetrical effects must have asymmetrical causes.”

— Pierre Curie, 1894

Minimal Path Principle

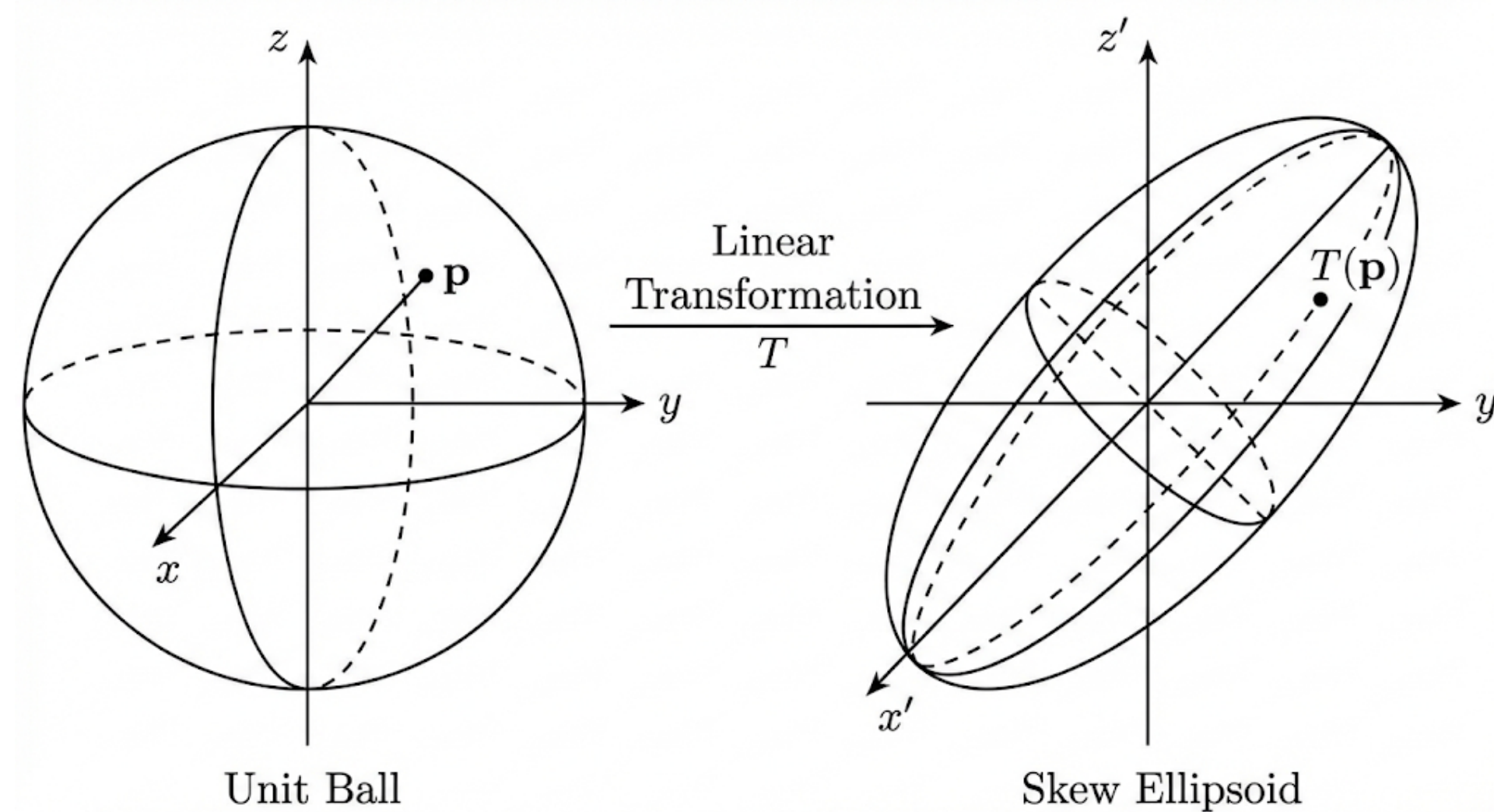


This talk

From Word Vector To Sentence Matrix

$$\mathbf{X} = \begin{pmatrix} \text{“The”} \\ \text{“quick”} \\ \text{“brown”} \\ \dots \\ \text{word } N \end{pmatrix}$$

X = The quick brown fox jumps over the lazy dog.



$$\sigma(\mathbf{W}) = \{\lambda_1, \lambda_2, \lambda_3\}$$

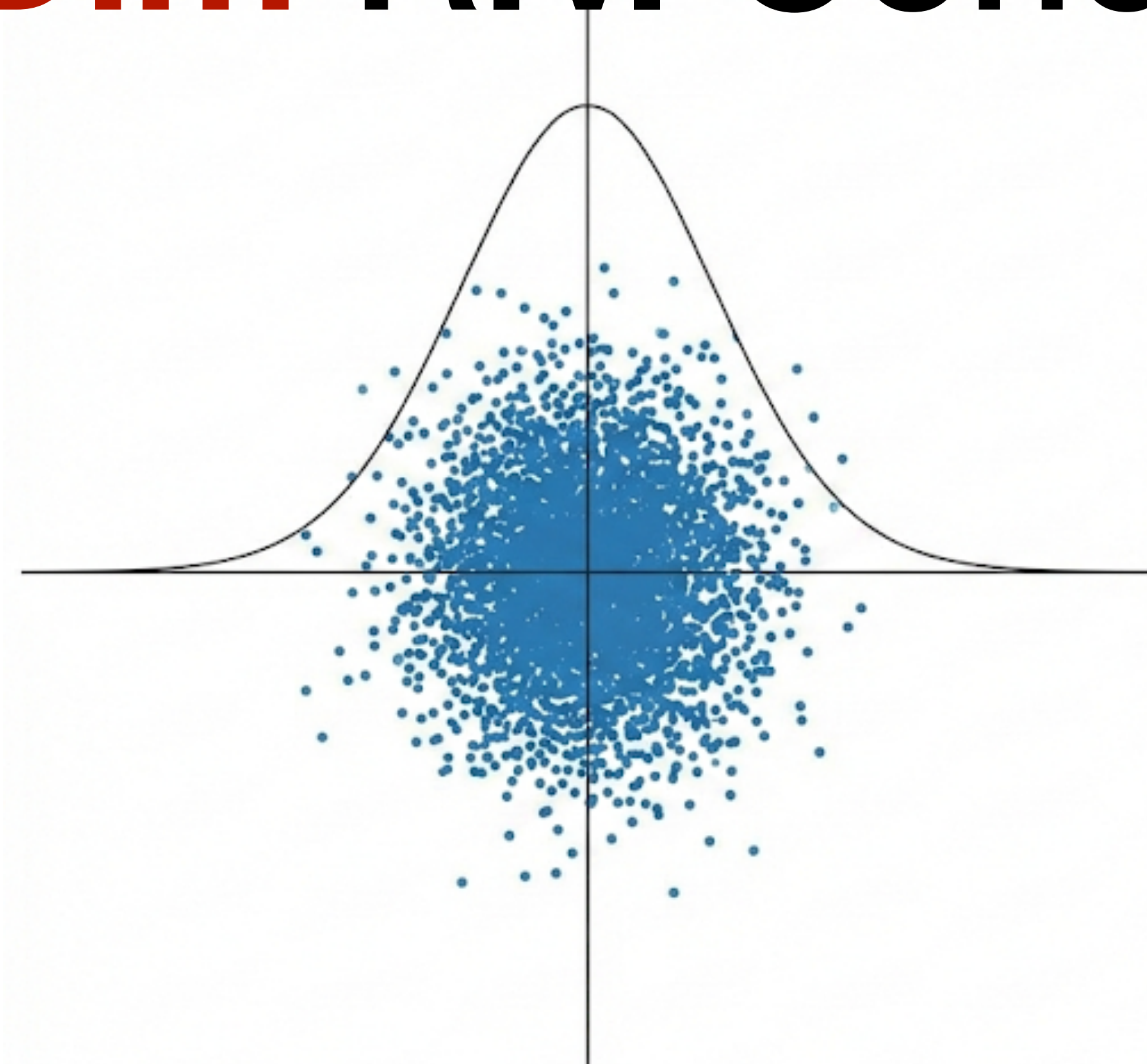
$$\mathbf{X} = \begin{pmatrix} \text{“Le”} \\ \text{“rapide”} \\ \text{“brune”} \\ \dots \\ \text{mot } N \end{pmatrix}$$

$$\mathbf{X}[1 :] = [0.3 \ 0.4 \ 0.5]$$

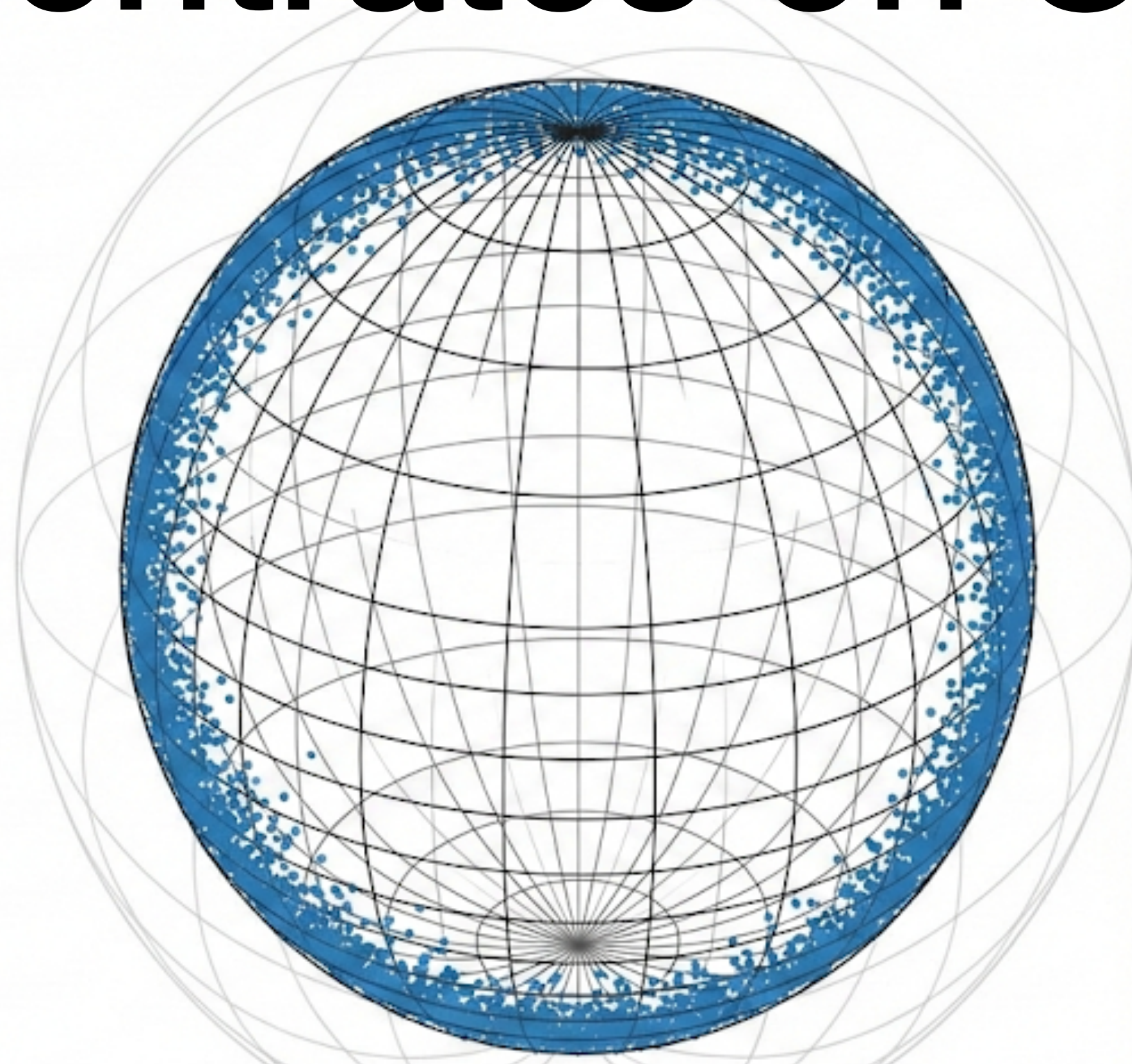
$$\xrightarrow{\mathbf{W} \in \mathbb{R}^{3 \times 3}}$$

$$\mathbf{X}[1 :] \mathbf{W} = [0.7 \ 0.8 \ 0.9]$$

Normalization Layer Is Almost Linear: **HighDim** R.V. Concentrates on Sphere



Dim = 2, 3, ...

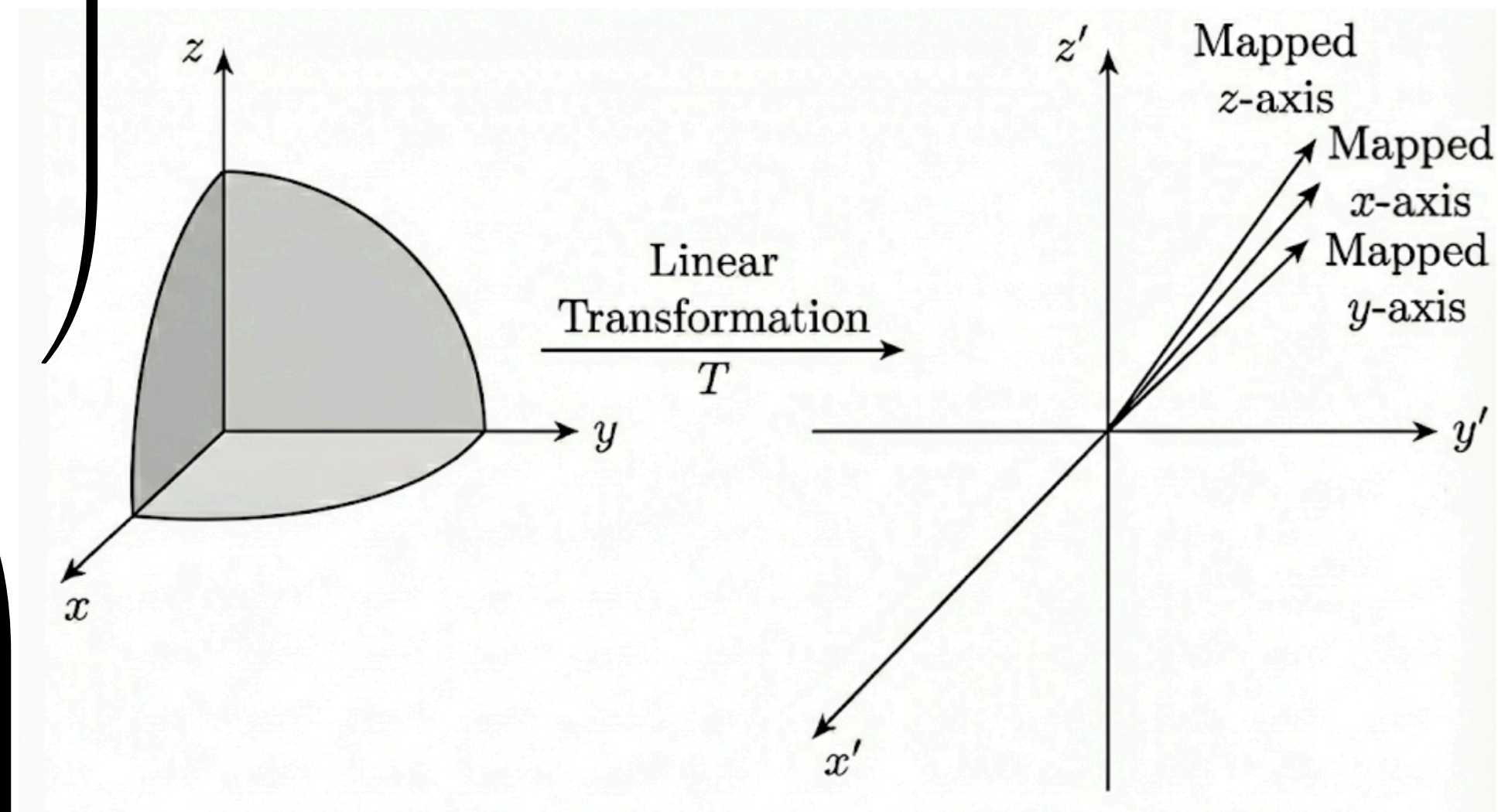
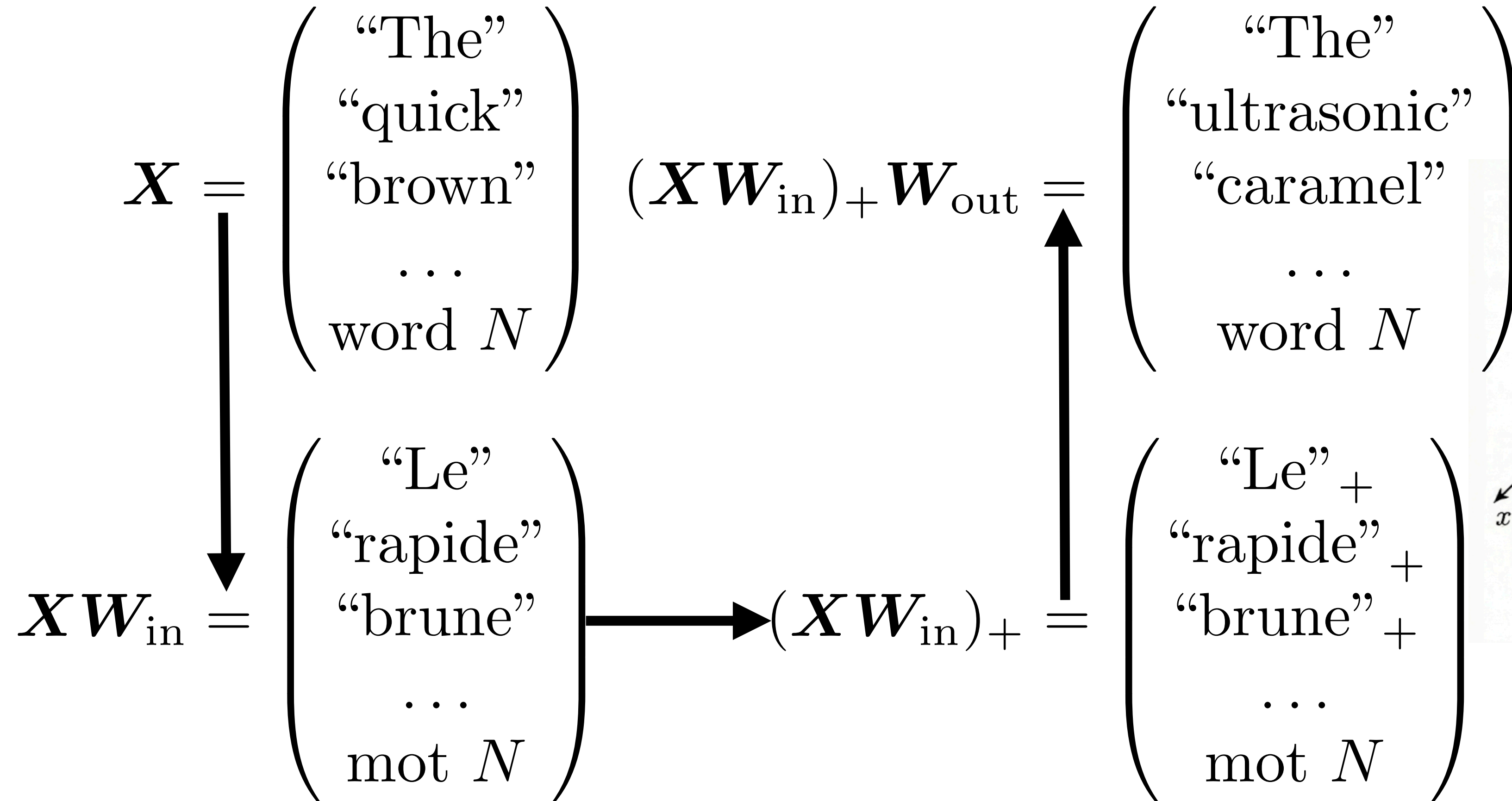


HighDim

MLP = Learnable Conic Projection

X = The quick brown fox jumps over the lazy dog.

MLP(X) = The ultrasonic caramel unicorn bounces over the snoozing python.



ReLU MLP = Tropical Rational Map
 $X + MLP(X)$ = "Projected Gradient Flow"

$$X[1:]W_{in} = [0.7 \quad -0.8 \quad 0.9] \xrightarrow{\text{ReLU}} (X[1:]W_{in})_+ = [0.7 \quad 0. \quad 0.9]$$

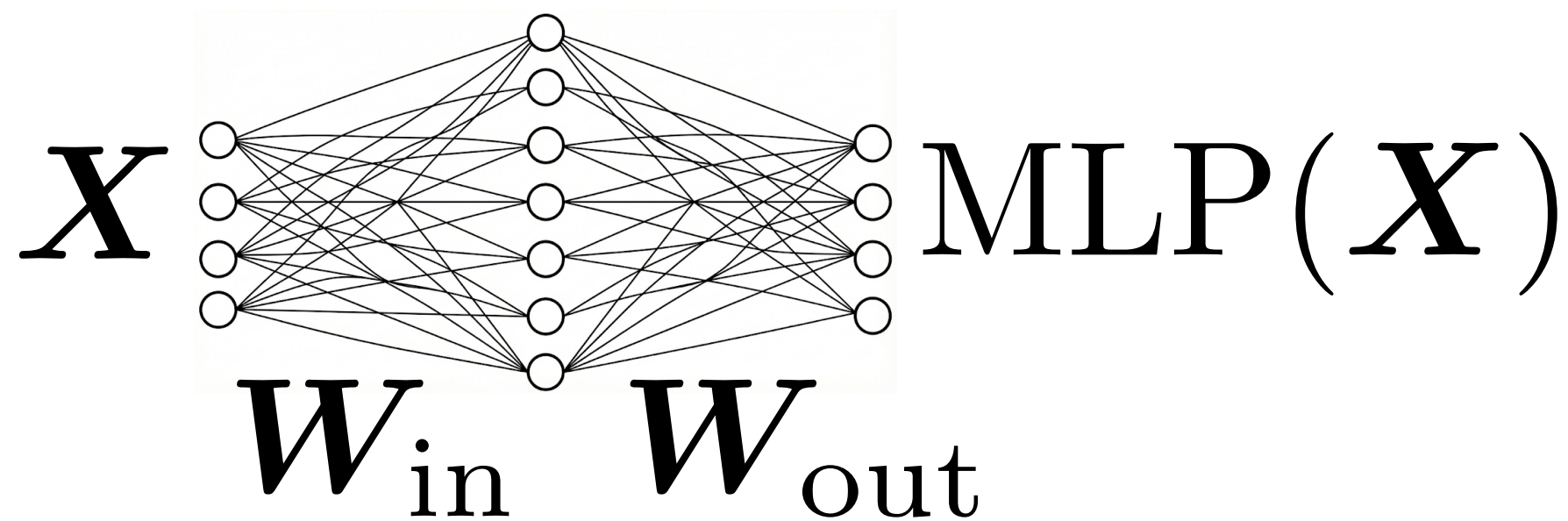
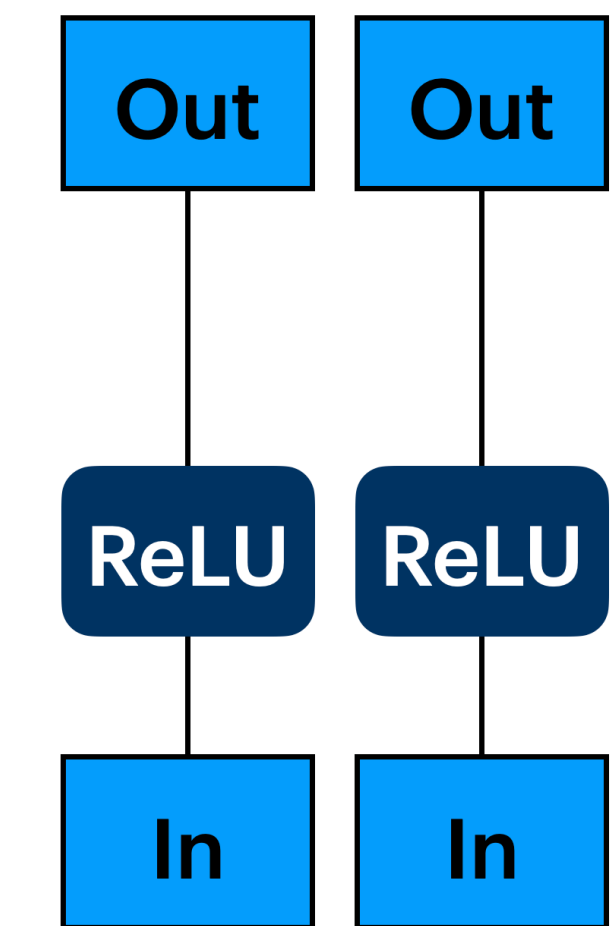
SelfAttention: Cubic/Nonlinear

$$\text{LinearAttention}(\mathbf{X}) = \underbrace{(\mathbf{X}\mathbf{Q})}_{\text{“query”}} \underbrace{(\mathbf{X}\mathbf{K})^*}_{\text{“key”}} \underbrace{(\mathbf{X}\mathbf{V})}_{\text{“value”}} \mathbf{O}^*$$

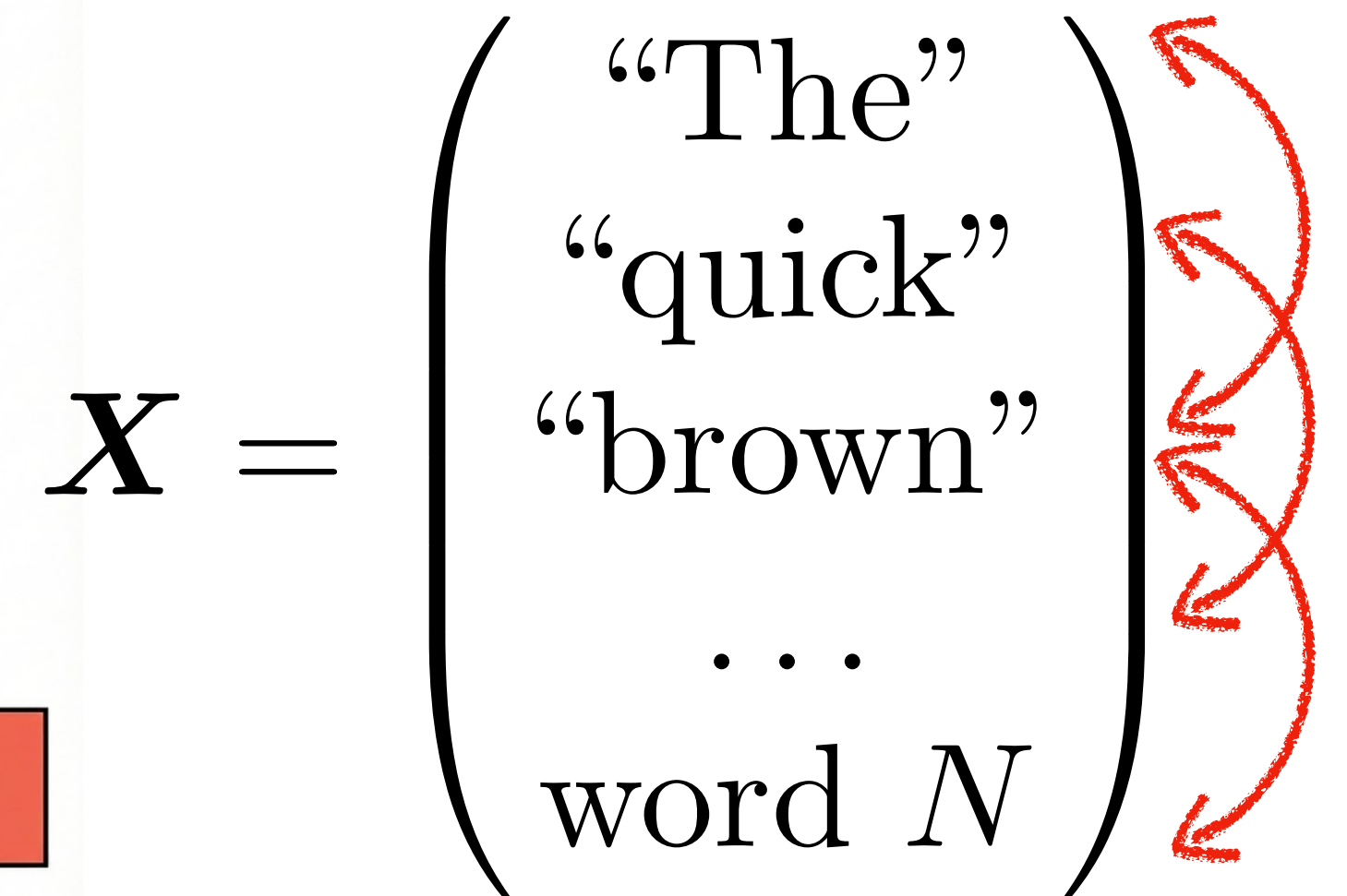
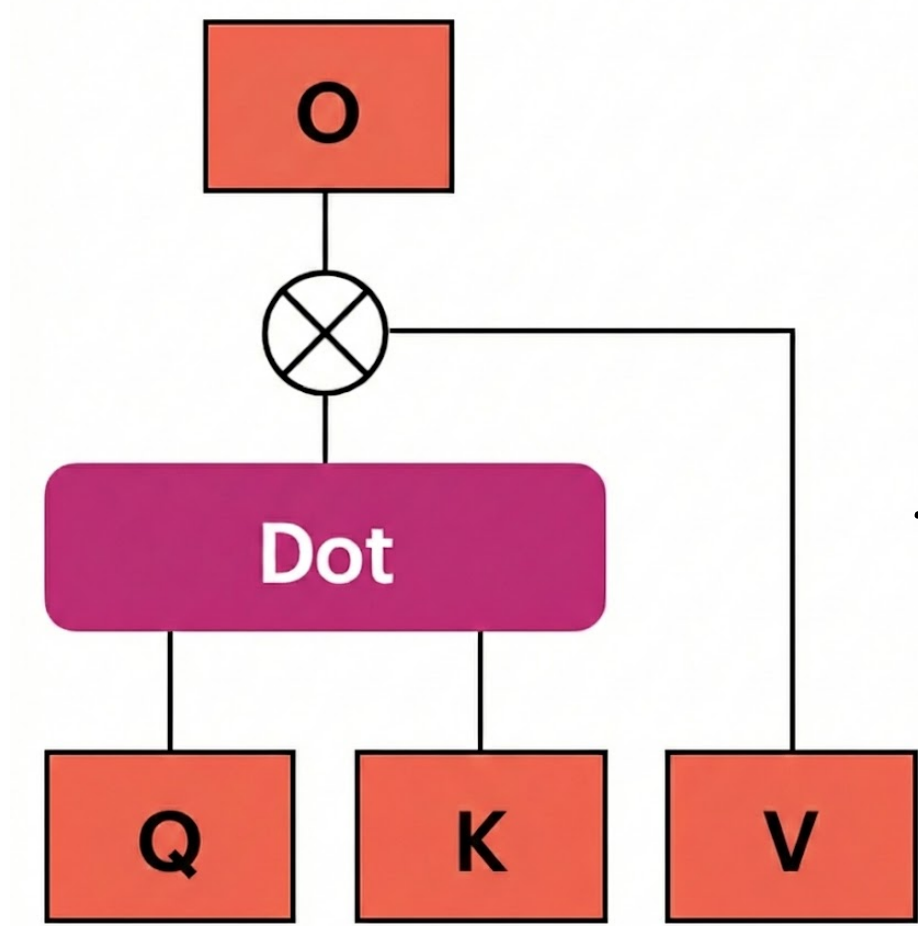
$$\text{SoftmaxAttention}(\mathbf{X}) = \text{RowSoftmax}(\mathbf{X}\mathbf{Q}\mathbf{K}^* \mathbf{X}^*) \mathbf{X}\mathbf{V}\mathbf{O}^*$$

From MLP to Attention

$$\text{MLP}(\mathbf{X}) = \mathbf{X} \mathbf{W}_{\text{in}} \mathbf{W}_{\text{out}}^* \quad \text{Attention}(\mathbf{X}) = \mathbf{X} \mathbf{Q} \mathbf{K}^* \mathbf{X}^* \mathbf{X} \mathbf{V} \mathbf{O}^*$$



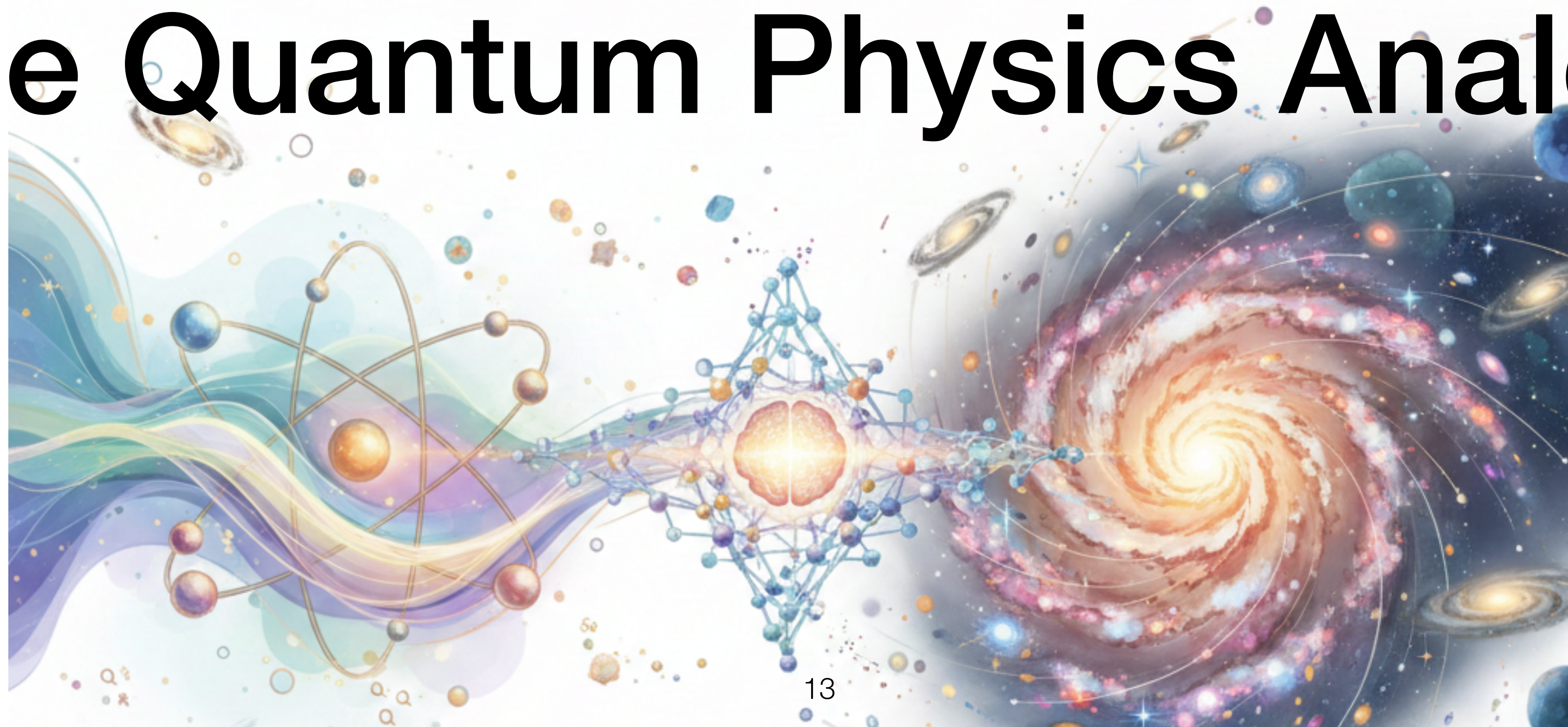
MultiHead MLP = Wide MLP



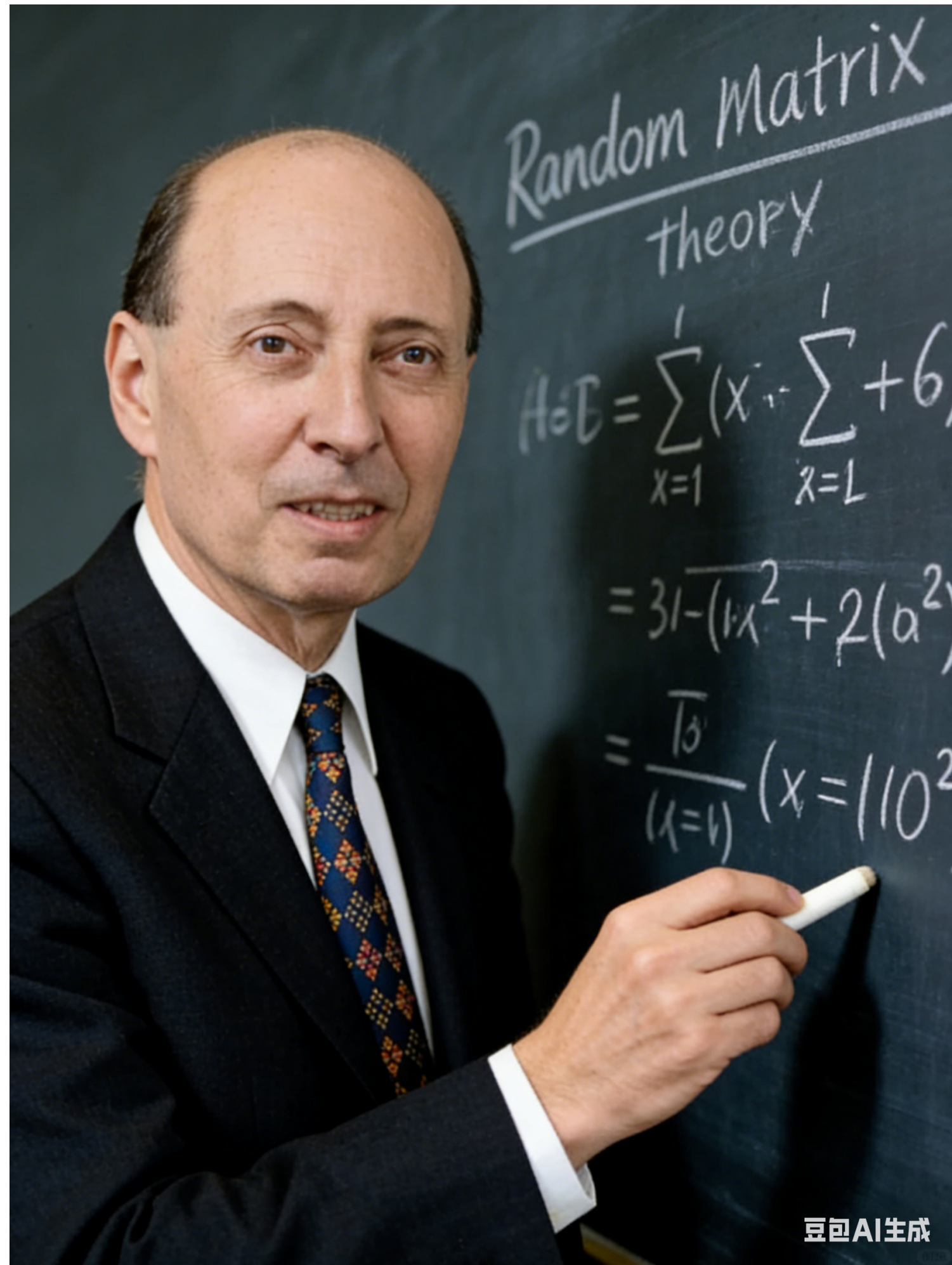
Right (Token-Wise) vs Left (Neuron-Wise) Multiplication

The Standard Model for AI

The Quantum Physics Analogy



From Quantum Chaos to ML



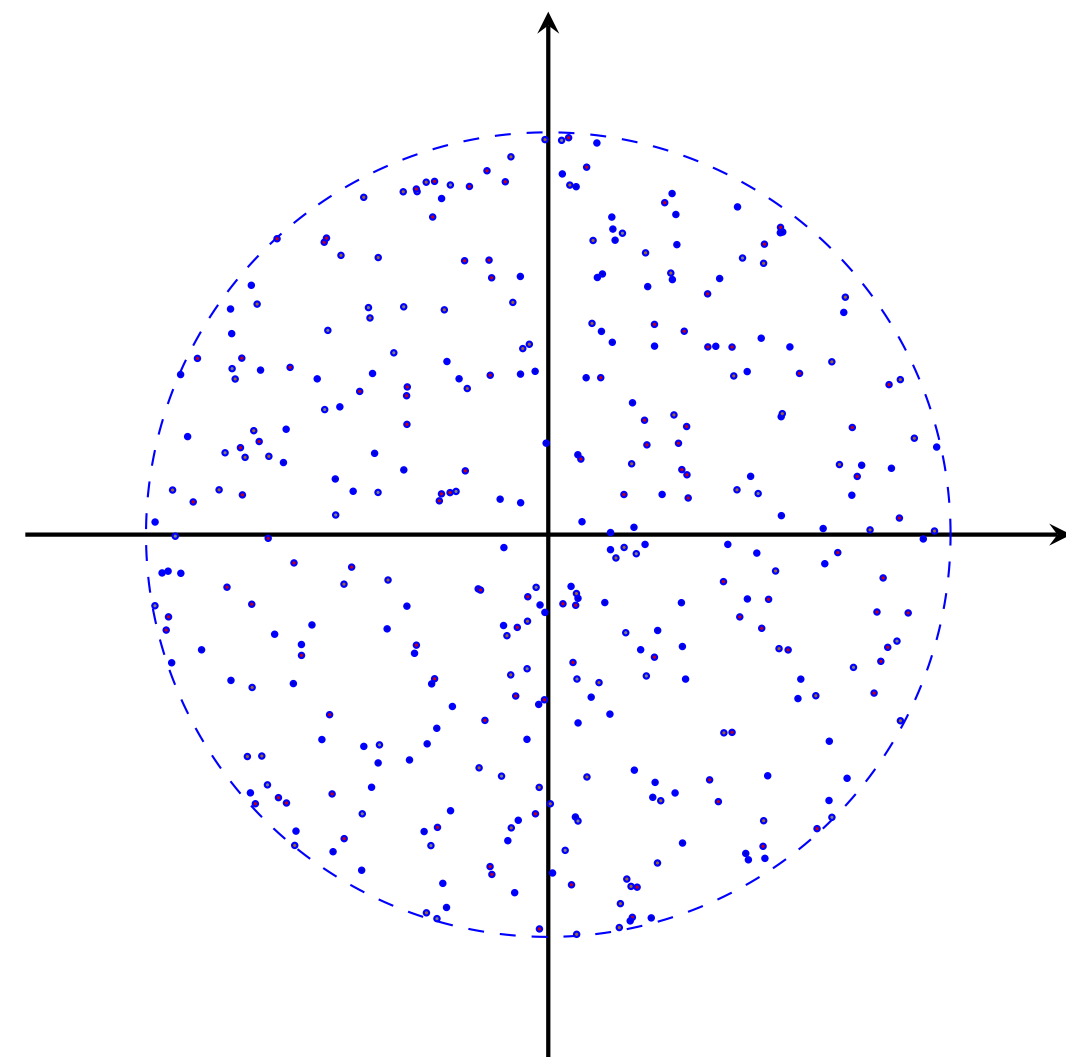
“Spectral statistics of complex systems are universal and determined solely by symmetry class — they are described by random matrix ensembles.”

— Eugene Wigner, 1959

From Disordered Systems to ML

$$\sigma(\mathbf{W}) = \{\lambda_1, \lambda_2, \lambda_3, \dots\} \subset \mathbb{C}$$

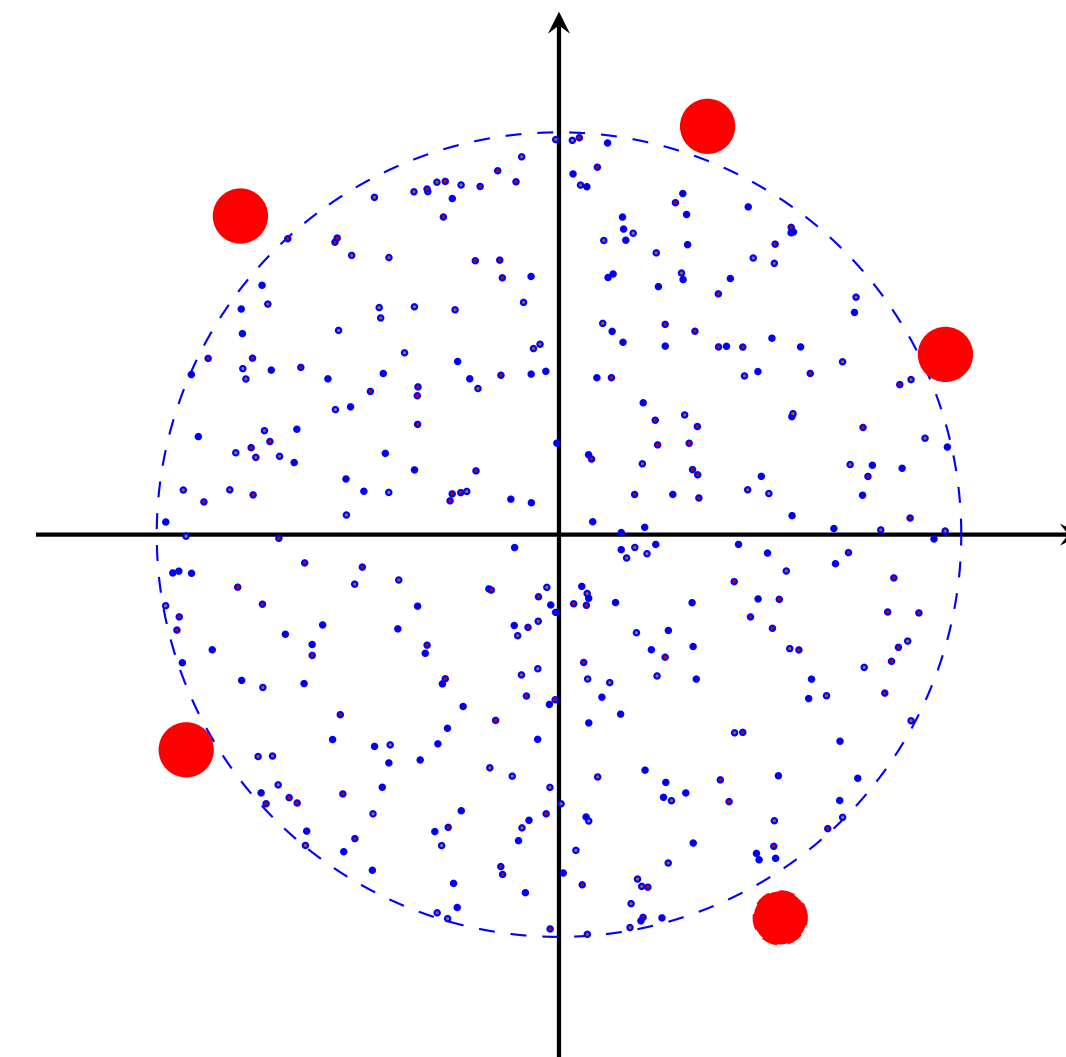
$\text{Im}(\sigma(\mathbf{W}))$



$\text{Re}(\sigma(\mathbf{W}))$

Before Learning (MaxEntropy)

$\text{Im}(\sigma(\mathbf{W}))$



$\text{Re}(\sigma(\mathbf{W}))$

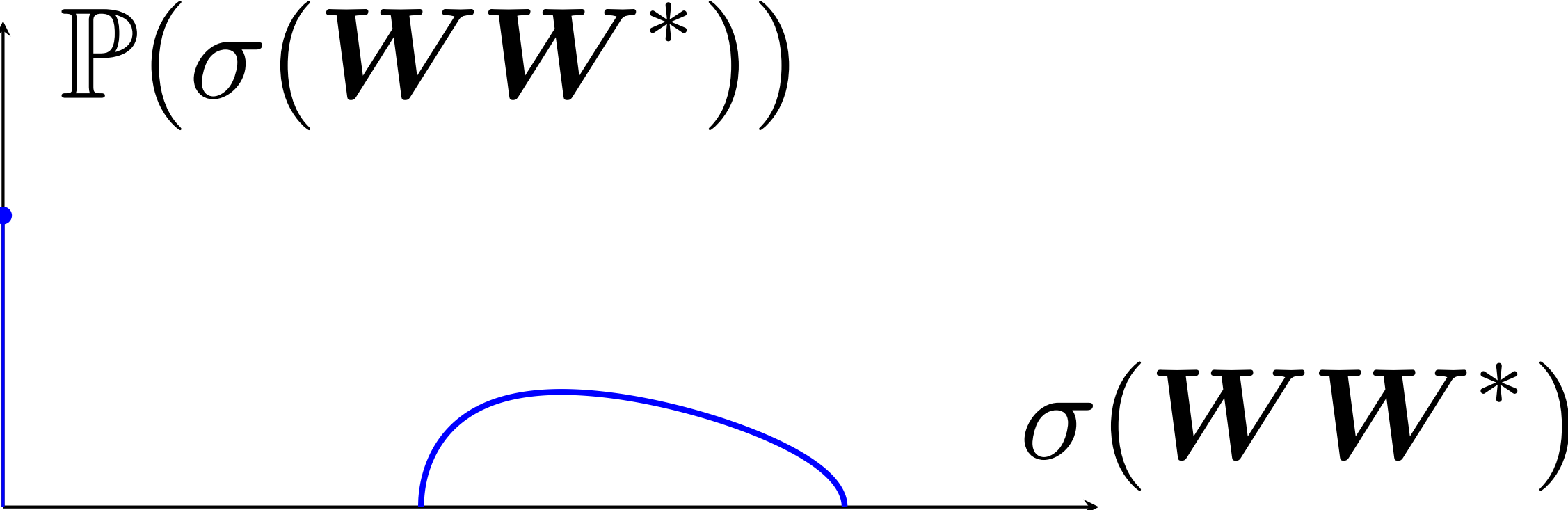
After Learning (Correlated Weights)

Generalization and Double Descent

Model Size

Tiny

c.f. Low-Rank QK / VO

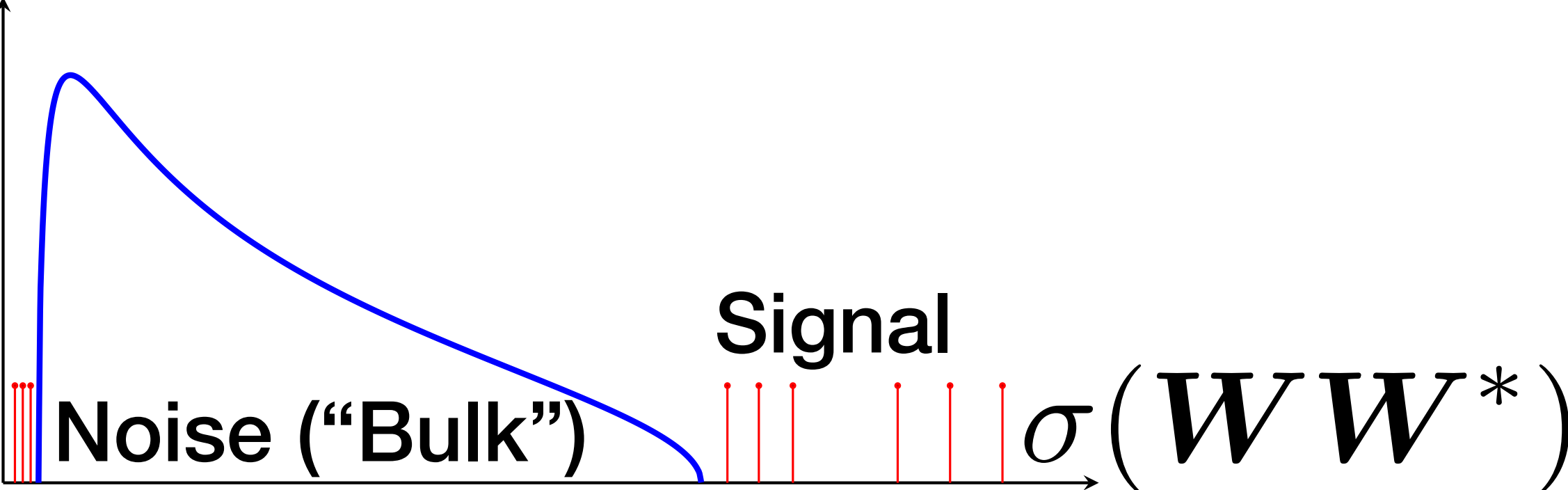


$$W \in \mathbb{R}^{C \times d}$$

$$d \ll C$$

Small

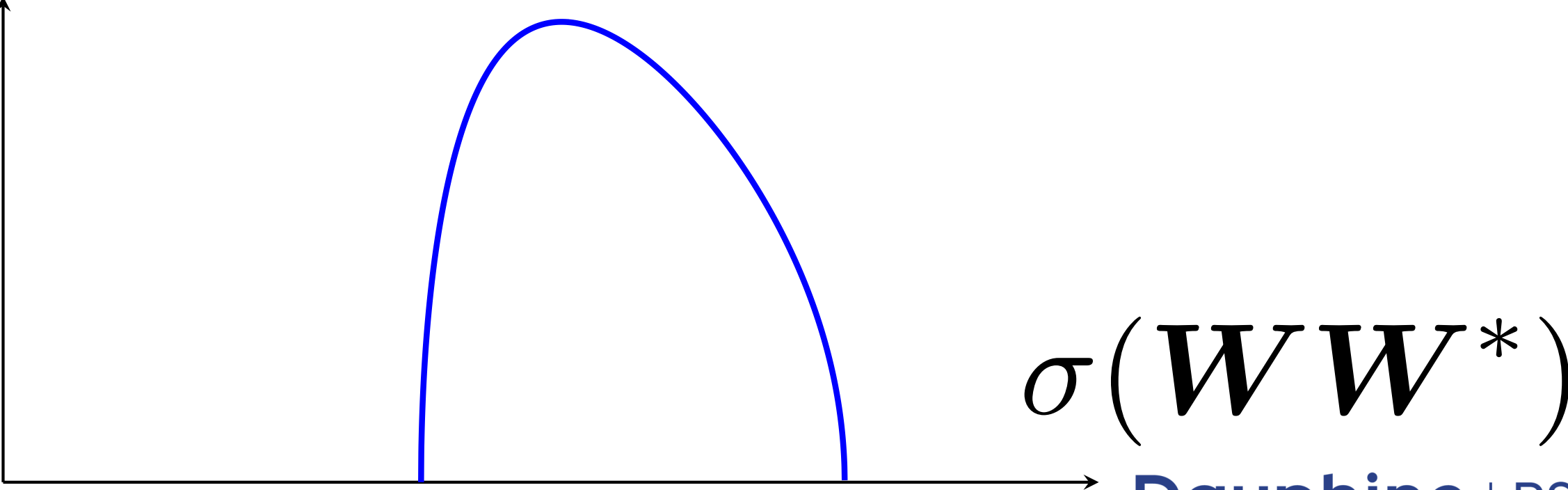
(Bad / Singular Case)



$$d \sim C$$

Large

c.f. Wide MLP: WinWout



$$d \gg C$$

MLP \leftrightarrow Attention

Free \leftrightarrow Interactive Systems

Quadratic Energy

Quartic Energy

$$V_2(\mathbf{X}) = \text{Tr}(\mathbf{X}\mathbf{W}\mathbf{X}^*) \quad V_4(\mathbf{X}) = \frac{1}{2} \text{Tr}(\mathbf{X}\mathbf{W}\mathbf{X}^*\mathbf{X}\mathbf{W}\mathbf{X}^*)$$

$$\text{MLP}(\mathbf{X}) = \frac{\partial V_2}{\partial \mathbf{X}^*} = \mathbf{X}\mathbf{W} \quad \text{Attention}(\mathbf{X}) = \frac{\partial V_4}{\partial \mathbf{X}^*} = \mathbf{X}\mathbf{W}\mathbf{X}^*\mathbf{W}$$

“Free”

Interactive

Are these potential fields canonical?

We need:

$$V_2(\mathbf{X}) = \text{Tr}(\mathbf{X} \mathbf{W} \mathbf{X}^*)$$

We also need:

$$V_4(\mathbf{X}) = \frac{1}{2} \text{Tr}(\mathbf{X} \mathbf{W} \mathbf{X}^* \mathbf{X} \mathbf{W} \mathbf{X}^*)$$

...but, what forms of $V(\mathbf{X}) = ?$ are **ALL** we need???

The Bridge: Symmetry



Symmetry / Invariance of Energy

Axiom

Invariance

Closed Form

Universality

$$V = V_1 + V_2 + \dots$$

$$V_k(\mathbf{X}) = \text{Tr}(\dots \mathbf{W}_{N \times N} \mathbf{X} \mathbf{W}_{C \times C} \mathbf{X}^* \dots)$$

Double-Trace / Hilbert Space Contraction

Perm/Rot Symm

$$V(U(N)\mathbf{X}) = V(\mathbf{X})$$

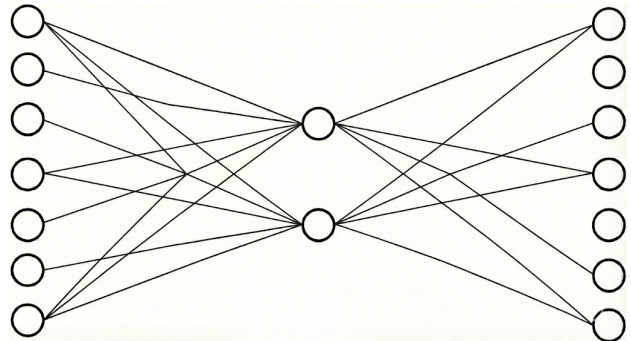
$$V(\mathbf{X}) = V(\mathbf{X}^* \mathbf{X})$$

Learnable Symm

$$V(\mathbf{X} C_{U(C)}(\mathbf{W})) = V(\mathbf{X})$$

$$V(\mathbf{X}) = \text{Tr}(\prod_i \mathbf{X} \mathbf{W}_i \mathbf{X}^*)$$

Low-Rank

$$\mathbf{X} \begin{array}{c} \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \end{array} \begin{array}{c} \circ \\ \circ \end{array} \begin{array}{c} \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \end{array} \mathbf{X} \mathbf{V} \mathbf{O}^*$$


$$\mathbf{W} \leftarrow \{ \mathbf{Q} \mathbf{K}^*, \mathbf{V} \mathbf{O}^*, \mathbf{W}_{\text{in}} \mathbf{W}_{\text{out}}^* \}$$

Geoffrey Hinton: "Hand Shaking Analogy"

Real and Complex-Valued Systems

Least Action:

$$S(\mathbf{X}) = \int (\text{Kinetic} - \text{Potential}) dt$$

Kinetic Term

$$K(\mathbf{X}) = \dot{\mathbf{X}} \mathbf{X}^* - \mathbf{X} \dot{\mathbf{X}}^*$$

Potential Term

$$V_2(\mathbf{X}) = \text{Tr}(\mathbf{X} \mathbf{W} \mathbf{X}^*)$$

$$V_4(\mathbf{X}) = \frac{1}{2} \text{Tr}(\mathbf{X} \mathbf{W} \mathbf{X}^* \mathbf{X} \mathbf{W} \mathbf{X}^*)$$

Anti-Symplectic Terms

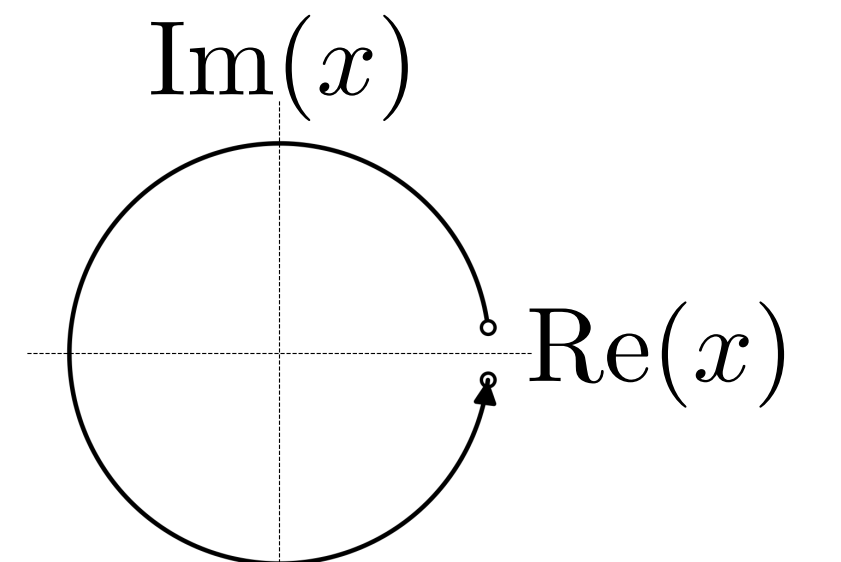
$$V_2(\mathbf{X}) = \text{Tr}(\overline{\mathbf{X}} \mathbf{W} \mathbf{X}^*)$$

Euler-Lagrange Equation $\partial S / \partial \mathbf{X}^* = 0$

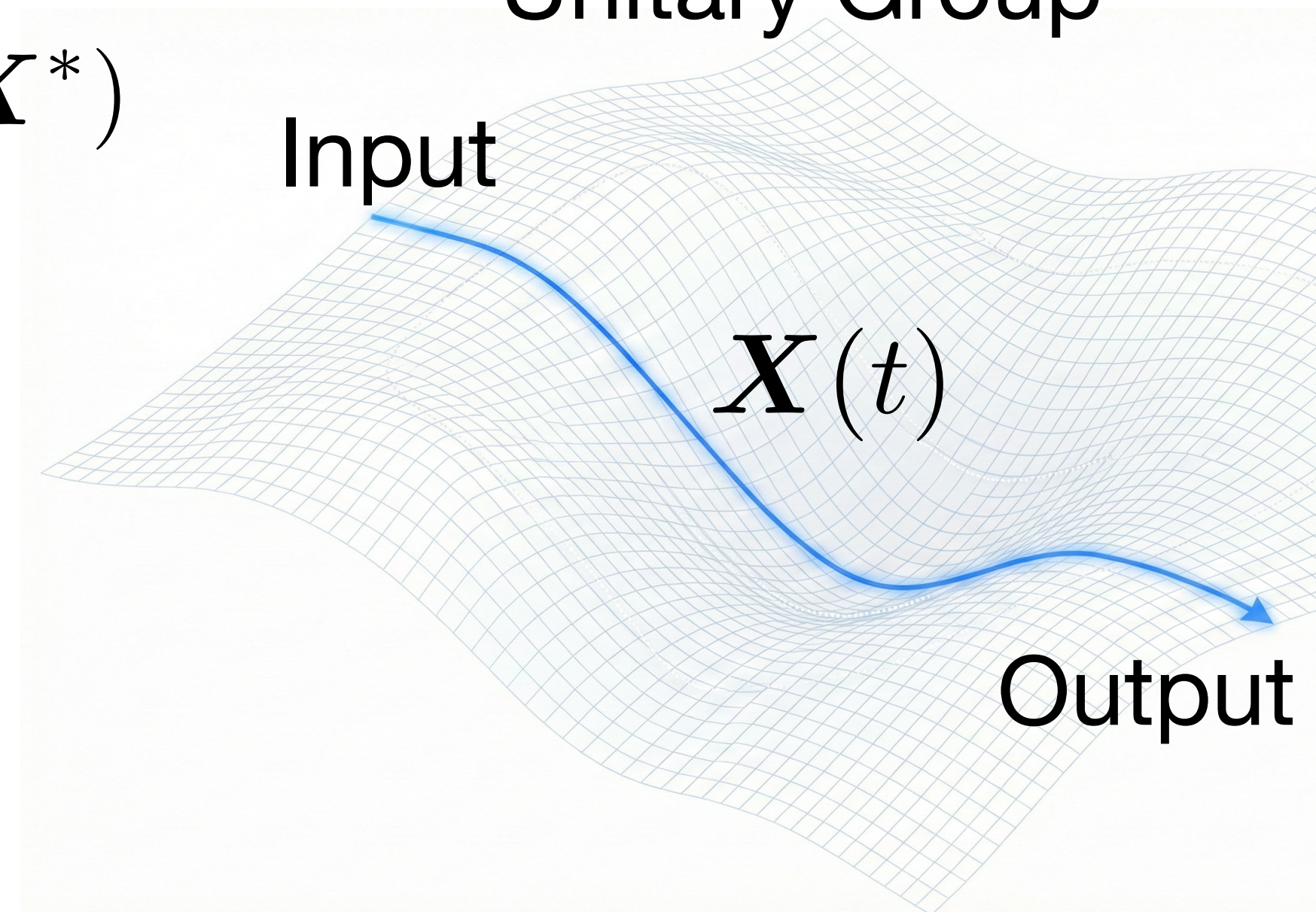
The Transformer ODE

$$\dot{\mathbf{X}} = \text{MLP}(\mathbf{X}) + \text{Attention}(\mathbf{X})$$

Complex Multiplication



Unitary Group



Linear versus Softmax Attention: Effective and Entropic Interactions

Recall: LinearAttention $V_4(\mathbf{X}) = \text{Tr}(\mathbf{X}\mathbf{W}\mathbf{X}^*\mathbf{X}\mathbf{W}\mathbf{X}^*)$

$$V_{\text{ent}}(\mathbf{X}) = \sum_{i=1}^{\text{nHead}} \sum_{j=1}^N \log \sum_{k=1}^N \exp\left(\left[\frac{1}{\sqrt{C}}\mathbf{X}\mathbf{W}_i\mathbf{X}^*\right]_{jk}\right)$$

$$\frac{\partial V_{\text{ent}}}{\partial \mathbf{X}^*} = \sum_{i=1}^{\text{nHead}} \text{RowSoftmax}\left(\frac{1}{\sqrt{C}}\mathbf{X}\mathbf{W}_i\mathbf{X}^*\right)\mathbf{X}\mathbf{W}_i$$

Renormalization

LinearAttention

“Effective Field”

Scale-Separated Interaction

Hot Limit: as $h \rightarrow 0$

$$V_{2+4}(h\mathbf{X}) = \boxed{h^2 V_2(\mathbf{X}) + h^4 V_4(\mathbf{X})} + O(h^6)$$

NO MORE TERMS!!! 

SoftmaxAttention

Entropic Regularization / “Free Energy”

Multi-Scale Interaction

$$V_{2+\text{ent}}(h\mathbf{X}) = h^2 V_2(\mathbf{X}) + \sum_{\text{Head, Token } i,j} \text{LogSum}_k \text{Exp}(h^2 [\mathbf{X} \mathbf{W}_i \mathbf{X}^*]_{jk})$$

Head, Token i, j

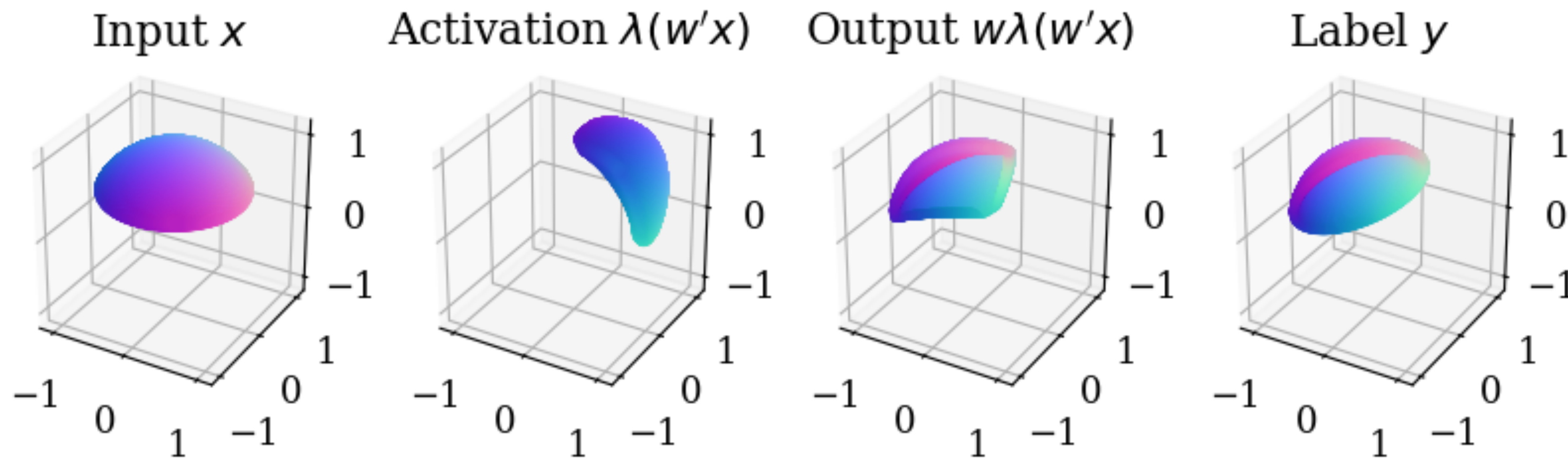
Hot Limit: Vanishes Too

Cold Limit: as $h \rightarrow \infty$ Concentrates on Max Token

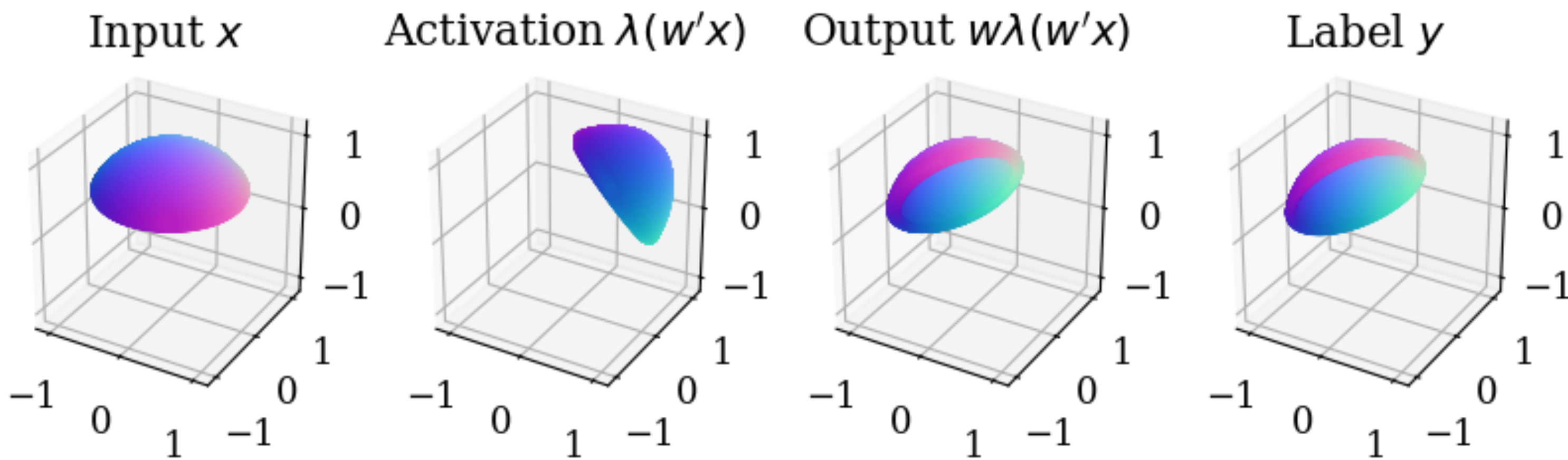
Applications: Predicting Structures



Embedding Space Favors Orthogonal Symmetry

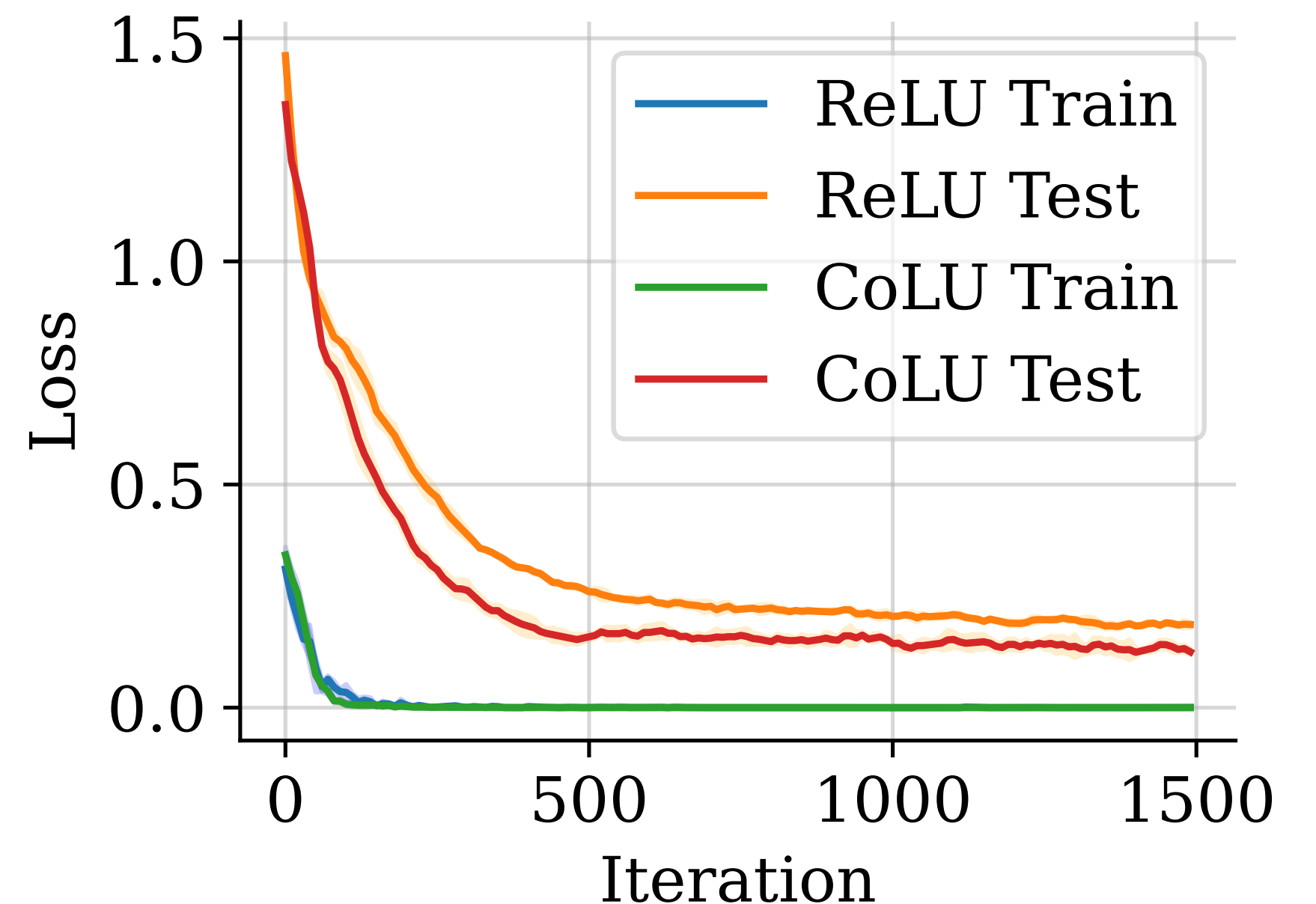


ReLU: permutation symmetry

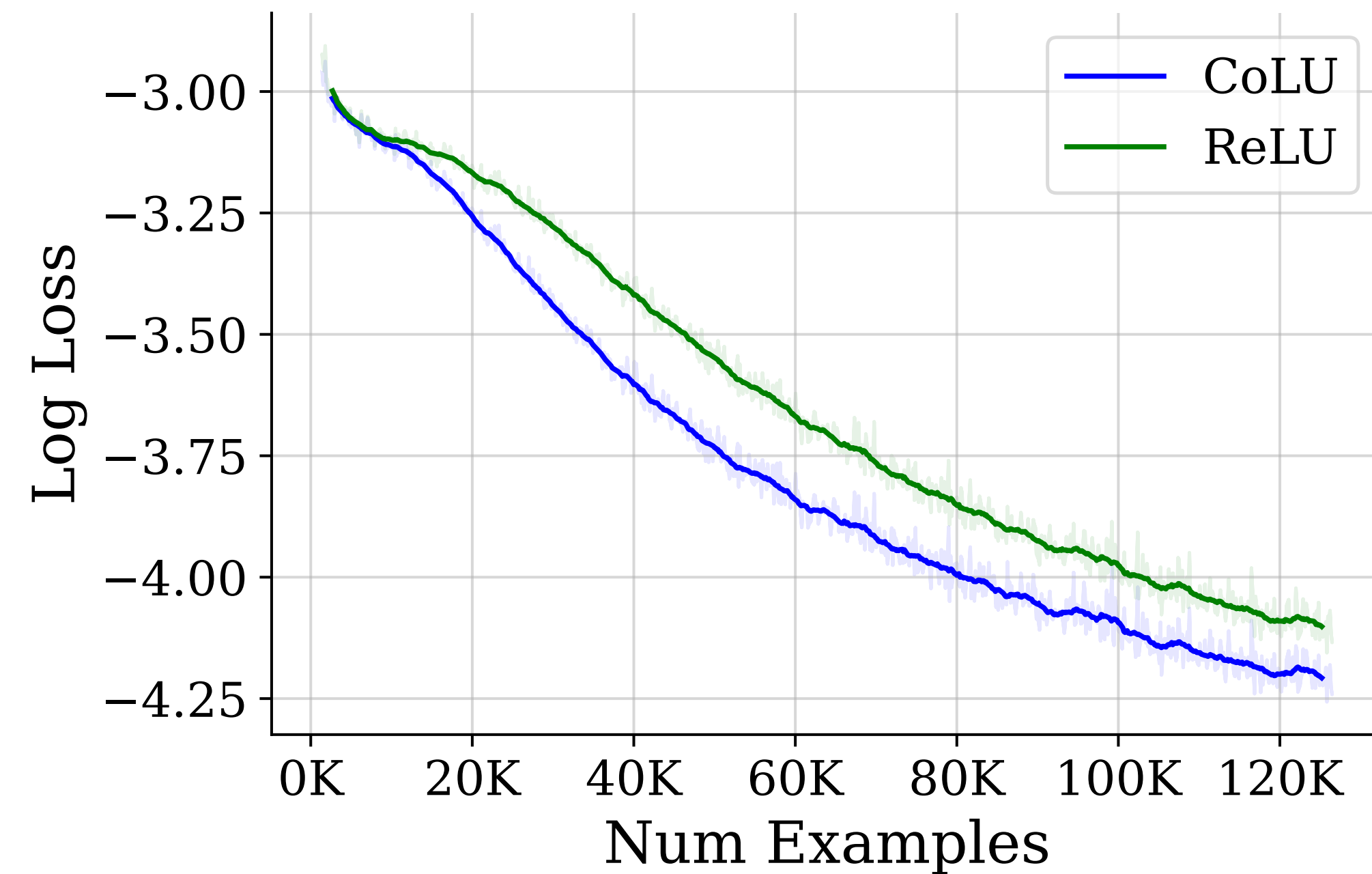


CoLU: orthogonal/rotary symmetry



- Minimal Example ($C=3$)
- Improved Generalization
- Accelerated Convergence
- Applies to ReLU-Attention



Scaling Up Conic Activations



Diffusion Transformer 0.8B
(Oxford102)

GPT2 MLP (FineWeb10M)	ReLU	CoLU
Forward FLOPs	39.064M	39.101M
Test Loss	3.4569 ± 0.1182	3.3804 ± 0.1159
ResNet-56 (CIFAR10)	ReLU	CoLU
Forward FLOPs	0.252M	0.257M
Test Accuracy	92.7282 ± 0.357	93.5851 ± 0.442
Diffusion Model (CIFAR10)	ReLU	CoLU (Faster)
Train Loss	0.1653	0.1458
Early Samples		

Conic linear units: improved model fusion and rotational-symmetric generative model. VISIGRAPP 2024. Conic activation functions. UniReps@NeurIPS 2024. PMLR 2025.

Gated Activation: Multi-Trace Model

$$V_3(\mathbf{X}) = (\mathbf{XW}) \odot (\mathbf{XW})(\mathbf{XW})^*$$
$$\implies \partial V_3 / \partial \mathbf{X}^* = (\mathbf{XW}) \odot (\mathbf{XW}) \mathbf{W}^*$$
$$\xrightarrow{\text{ReLU}} (\mathbf{XW}) \odot (\mathbf{XW})_+ \mathbf{W}^*$$

Sparse Attention: Non-Commutative Systems

Lemma 8 (Sparse Attention). *Under the pairwise non-commutative condition $[\lambda_{MLP}, \lambda_A] = \beta \lambda_A$, the sequential composition of the layers is equivalent to a single effective flow generated by:*

$$\exp(\lambda_{MLP}) \exp(\lambda_A) = \exp(\lambda_{MLP} + G(\beta) \cdot \lambda_A)$$

*where the effective interaction gate $G(\beta)$ is given by the generating function of the Bernoulli numbers:
 $G(\beta) = \frac{\beta}{1-e^{-\beta}}$. As $\beta \rightarrow -\infty$, $G(\beta) \sim |\beta|e^{-|\beta|} \rightarrow 0$.*

“Fail to Commute between Tokens \implies Turn off the Gate”

Reference:

<https://changqingfu.com/pdf/transformer.pdf>

