# Question

## Can we improve neural networks by generalizing its symmetry constraints?



Land's Symmetry

Water Flow

# Activation Functions

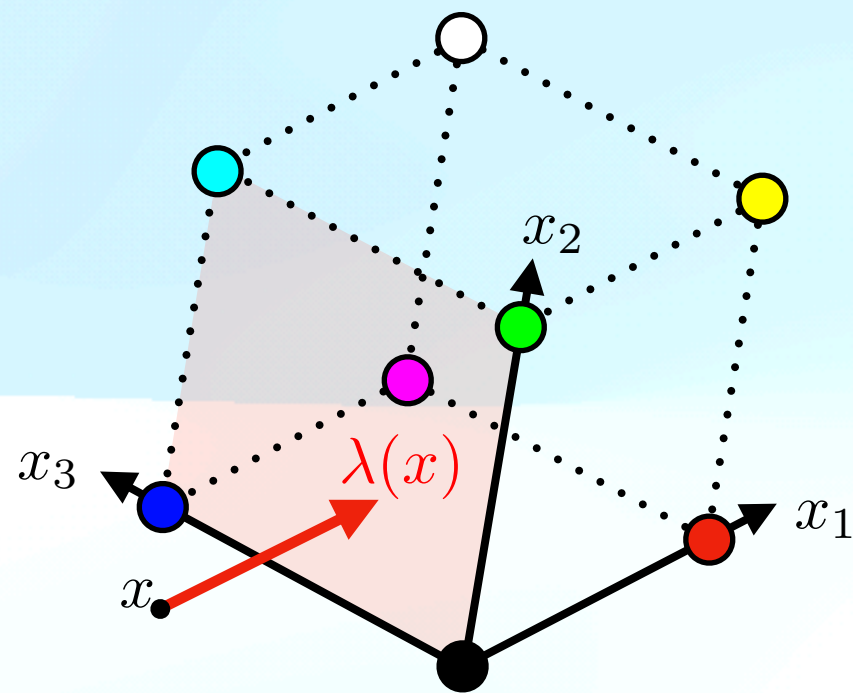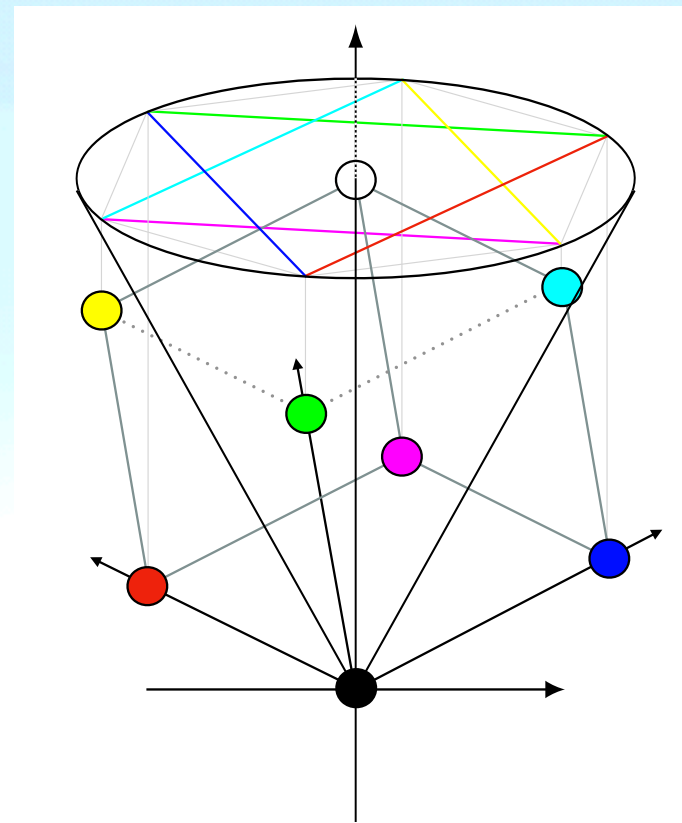**ReLU:** Rectified Linear Units

**Characterization**

- Axis Homogeneity  $\Lambda(x)_i = \Lambda(x_i) \text{ i.e. } \Lambda\pi_i = \pi_i\Lambda$  $\Lambda(x_1, x_2, \ldots) = (\Lambda(x_1), \Lambda(x_2), \ldots)$

- +Idempotence  $\Lambda(\Lambda(x)) = \Lambda(x)$  $\Lambda\big|_O(x) = \text{id}(x), O \text{ is Borel}$

- +Positive Homogeneity  $\forall t > 0, \Lambda(tx) = t\Lambda(x)$  $\Lambda(x) := \lambda(x) = \max\{x, 0\}$

# Solution

**Allow orthogonal equivariance with a more symmetric invariant set**



ReLU: loses the rotary symmetry

CoLU: keeps the rotary symmetry

Previous works: spatial domain (Geometric Deep Learning)
Our work: feature space!

# Conic Activation Functions

- **Semidefinite Program**  $\Omega = \mathbb{R}_+^C$  $\min_{y \geq 0} \frac{1}{2}\|x - y\|_2$

Solution  $\lambda(x) = \pi_\Omega(x) = x_+ = \max\{x, 0\}$



dimension 1  dimension 2  dimension 3

- **Conic Program?**  $\Omega = \{x : x_1^2 \geq x_2^2 + x_3^2 + \ldots\}$  $\min_{y \in \Omega} \frac{1}{2}\|x - y\|_2$

# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
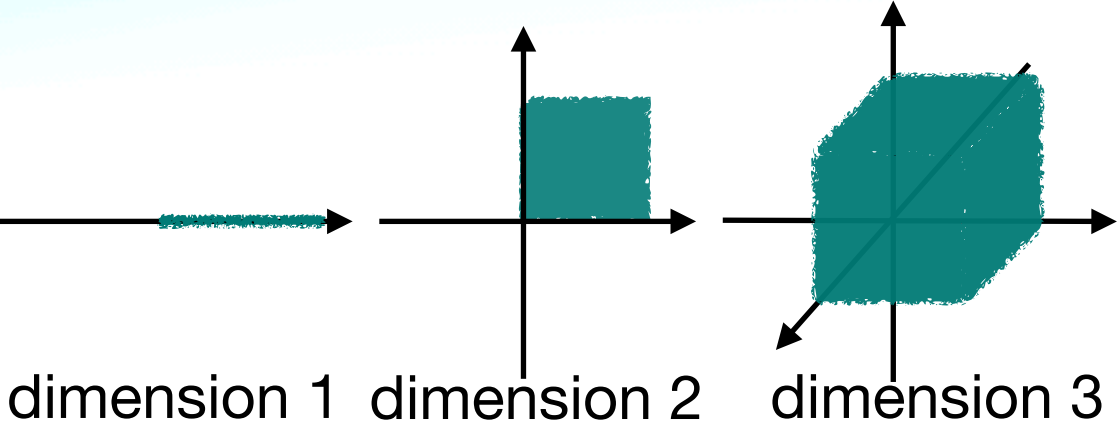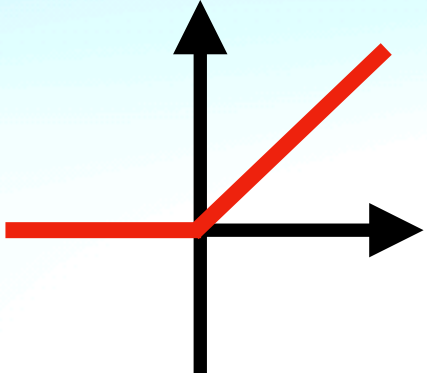
## CoLU Symmetry is Compatible with Transformer

| Nonlinearities | Function | Group | Symmetry | Limiting |
|---|---|---|---|---|
| **Attention** | $x \in \mathbb{R}^{C \times N} \mapsto Z^{-1} \exp(\frac{\langle x, x \rangle_C}{\sqrt{C}})x$ | Orth | Entropic | ColorClusters |
| **ReLU** | $x \in \mathbb{R}^C \mapsto \max\{x, 0\}$ | Perm | Simplex $\Delta^{C-1}$ | Orthant $\mathbb{R}^C_+$ |
| **CoLU** | $x \in \mathbb{R}^C \mapsto \pi_{\widetilde{V} \cap H(x)}(x)$ | **Orth** | Disk $D^{C-1}$ | Cone $\widetilde{V}$ |

Ðauphine | PSL★  CEREMADE UMR CNRS 7534  PR[AI]RIE cnrs
UNIVERSITÉ PARIS  PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
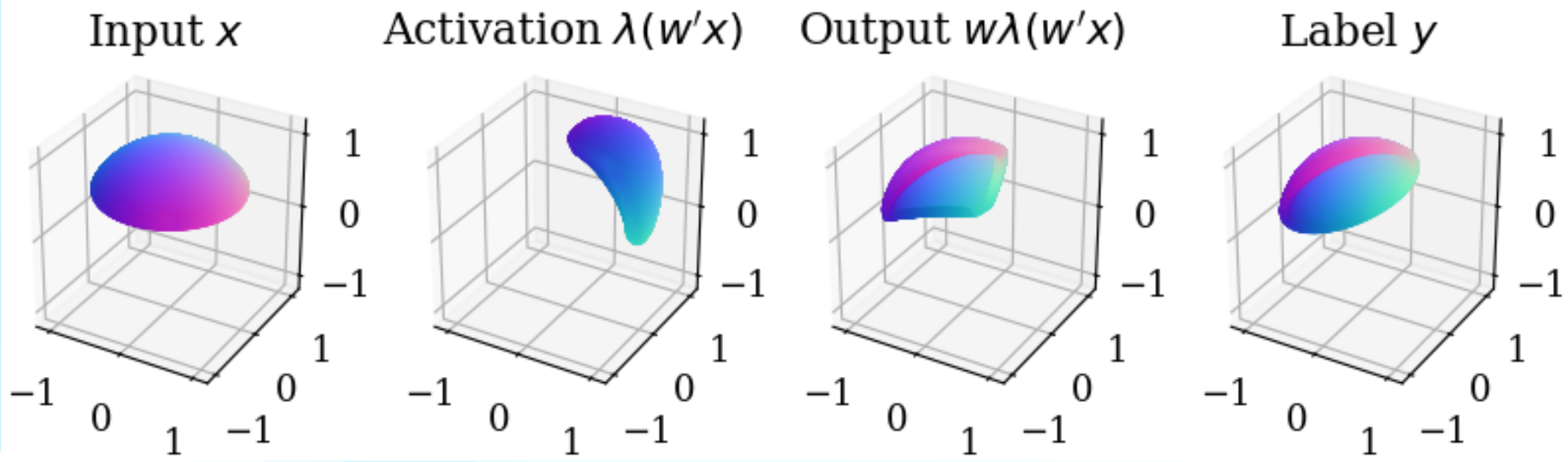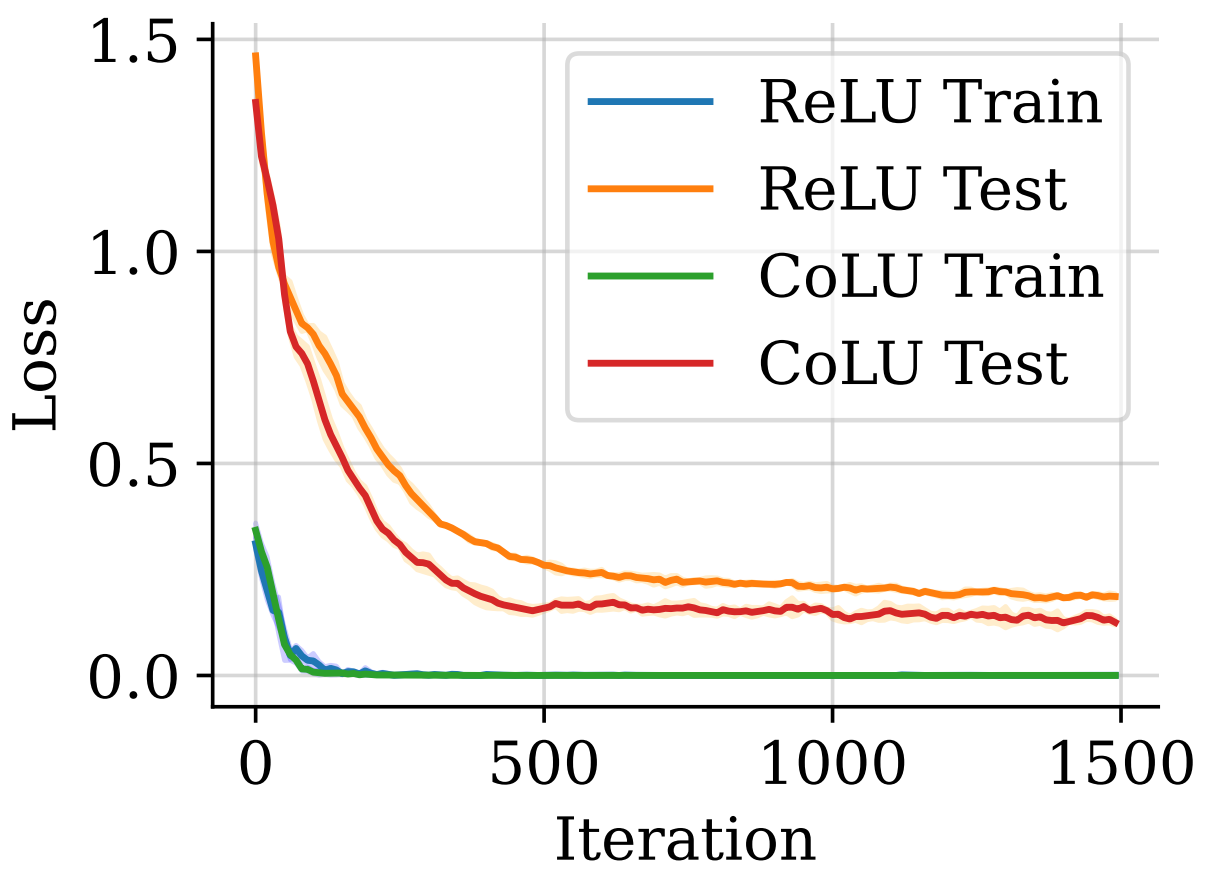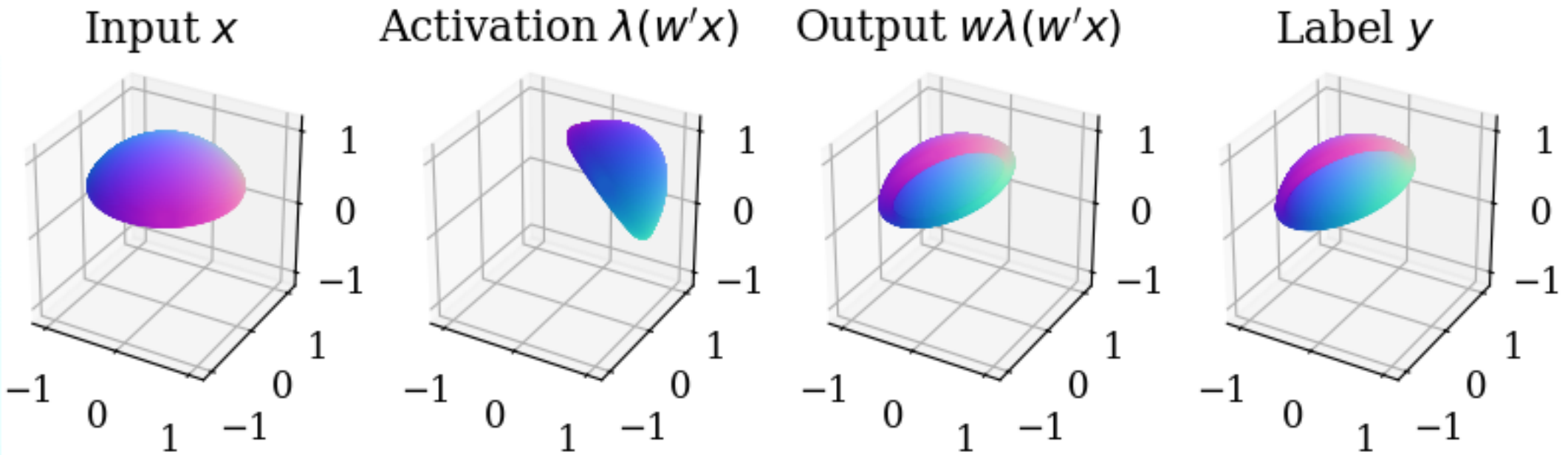
ReLU: permutation symmetry



CoLU: orthogonal/rotary symmetry



- Minimal Example
- Improved Generalization

# Conic Activation Functions

## Closed Form

$$\lambda(x)_i = \begin{cases} x_1, & i = 1 \\ \min\left\{\max\left\{x_1/(|x_\perp| + \varepsilon), 0\right\}, 1\right\} x_i, & i = 2, \ldots, C \end{cases}$$

$$\underbrace{\qquad\qquad\qquad\qquad}$$

$$x_\perp = (0, x_2, \ldots, x_C)$$

## Projective Form

$D$

$$\lim_{\varepsilon \to 0} \lambda(x) = \pi_{\widetilde{V} \cap H(x)}(x) = \pi_{\max\{x_1, 0\}D + \min\{x_1, 0\}\mathbf{e}_1}(x)$$

Animation:
Conic Symmetry

Ðauphine | PSL★ CEREMADE PR[AI]RIE cnrs
UNIVERSITÉ PARIS    UMR CNRS 7534    PaRis Artificial Intelligence Research InstitutE
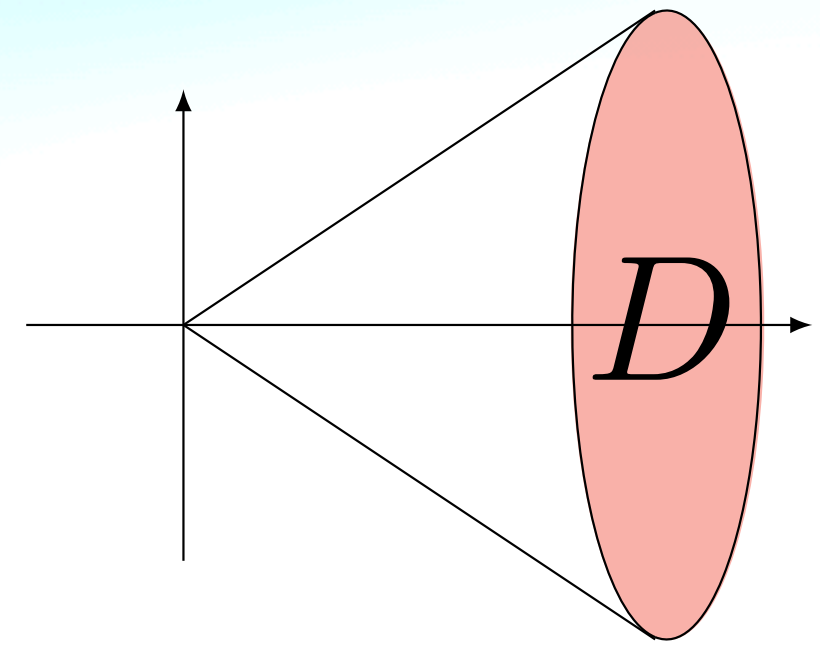
# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
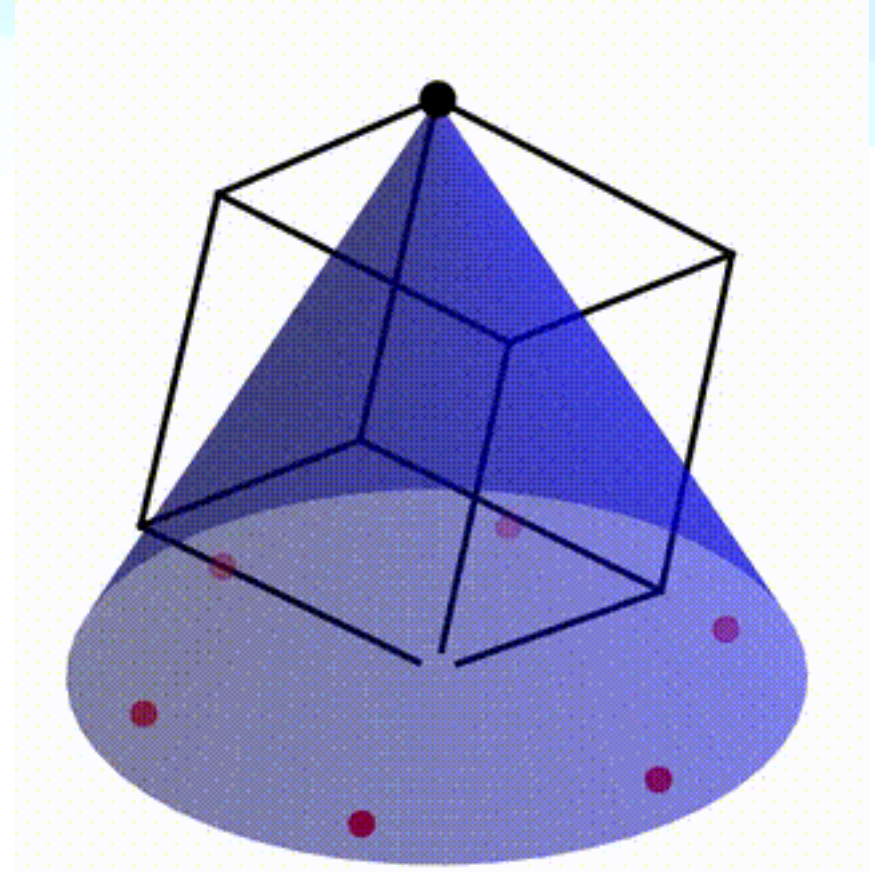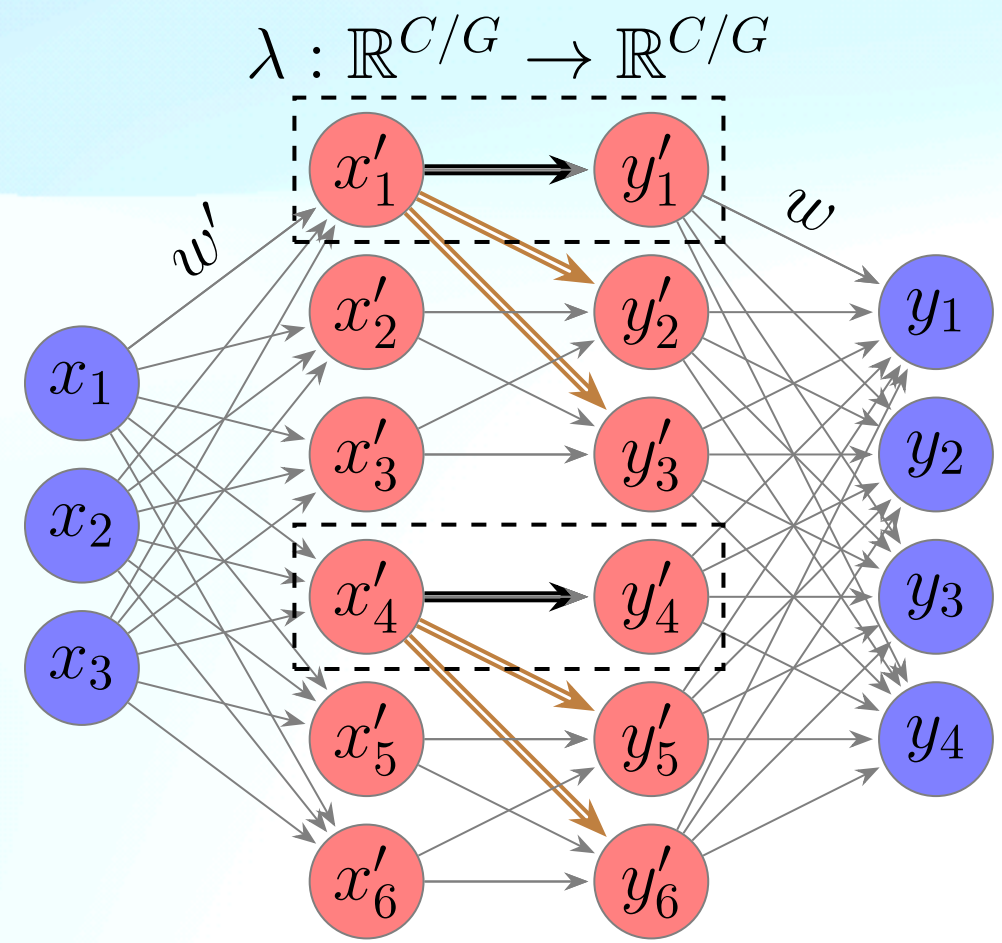
Multi-Head Structure $\quad \pi_i^G \lambda = \lambda \pi_i^G, i = 1, 2, \ldots, G$ where $C = GS$

Symmetry Group: $\quad \mathrm{Perm}(G) \times \mathrm{Orth}^G(S - 1)$

$\lambda : \mathbb{R}^{C/G} \to \mathbb{R}^{C/G}$



$C = 6, S = 3, G = 2$

6 Neurons
3D Cone
2 Cones



Animation:
Conic Symmetry

**Ðauphine** | PSL ★ **CEREMADE** PR[AI]RIE cnrs
UNIVERSITÉ PARIS   UMR CNRS 7534   PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
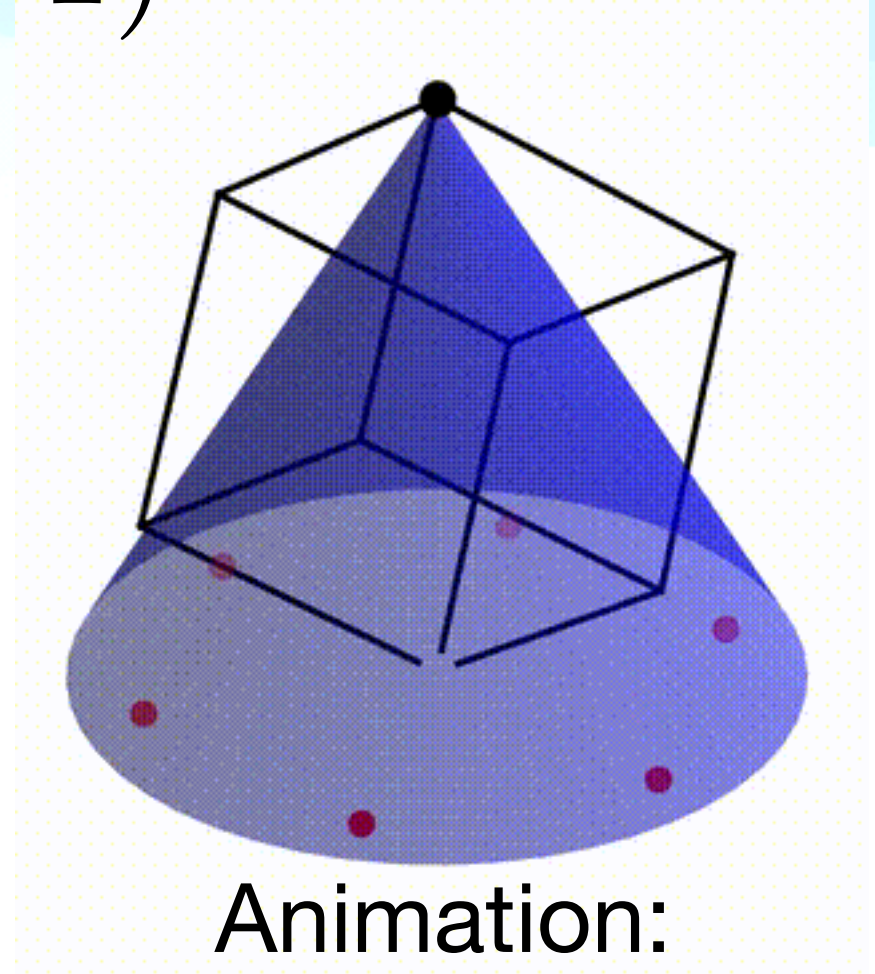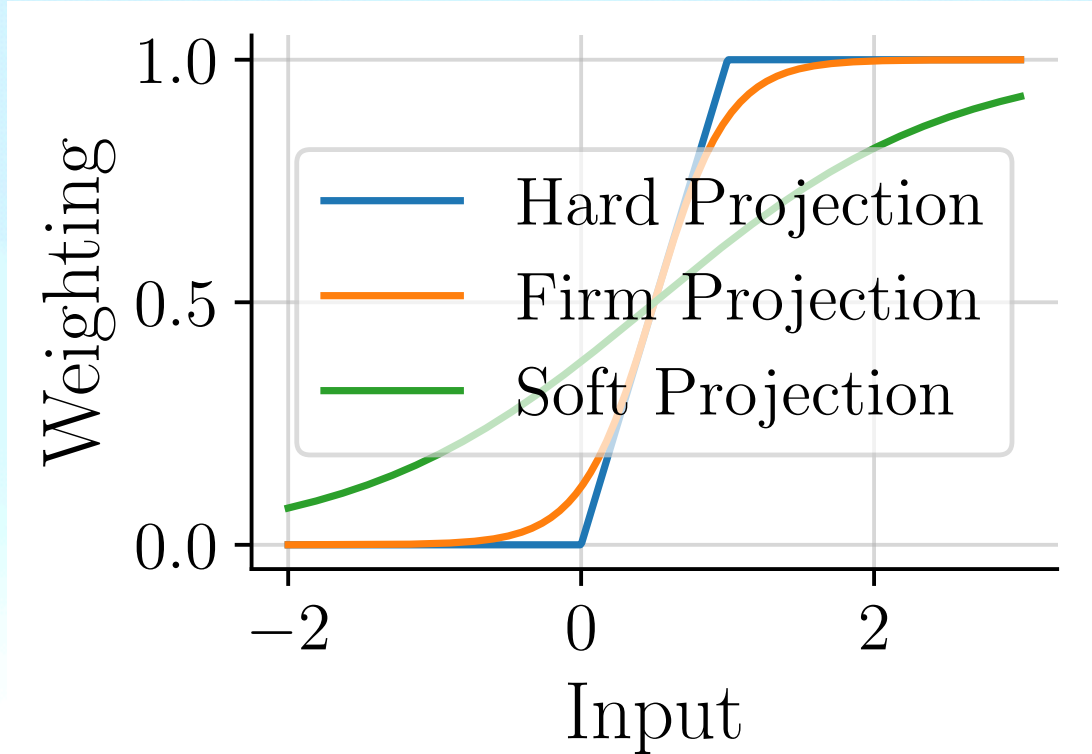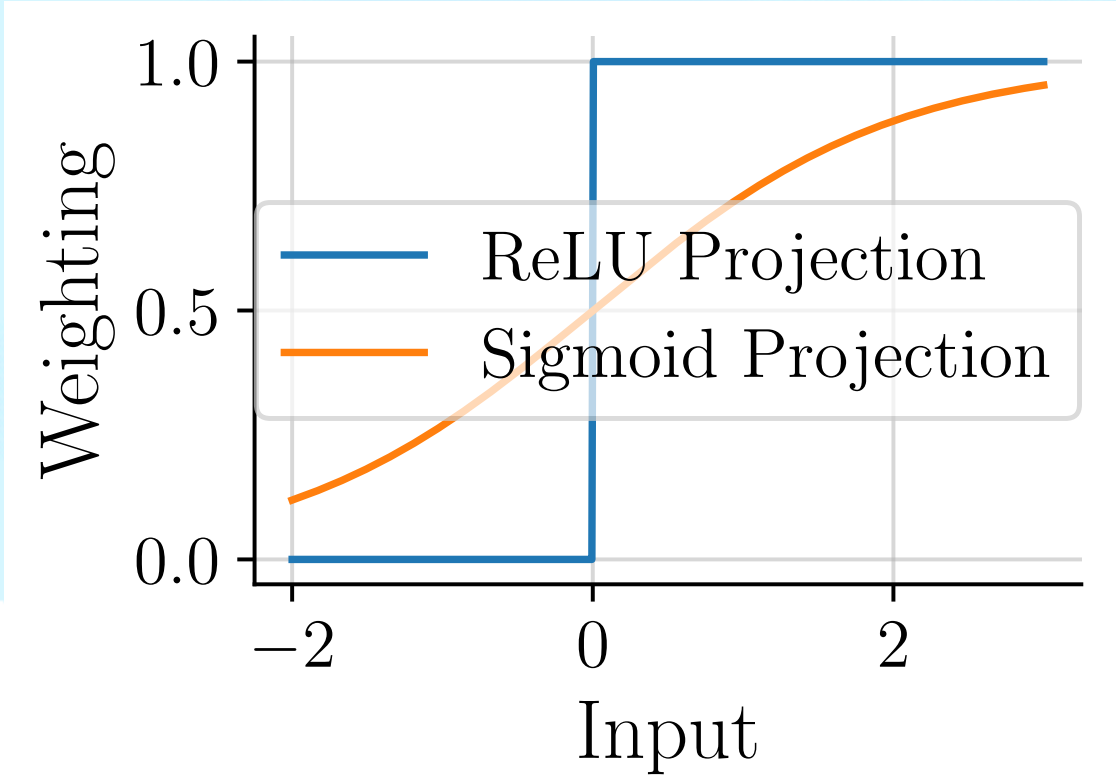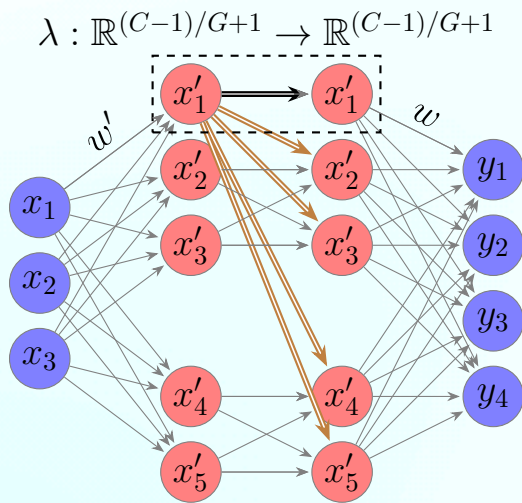
## Soft Projection







*Colinear Axes*

## Axis Sharing



$\lambda : \mathbb{R}^{(C-1)/G+1} \to \mathbb{R}^{(C-1)/G+1}$

5 Neurons
3D Cone
2 Cones

Glue cone axes w/ $\pi_i^G = \pi_1 \times \pi_{\mathrm{another}(S-1)\mathrm{axes}}$

Ðauphine | PSL ★ CEREMADE PR[AI]RIE cnrs
UNIVERSITÉ PARIS          UMR CNRS 7534   PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

## Generated Samples of CoLU-LDM

## Faster Learning



UNet with Attention
(835M parameters)

10

# Diffusion Model and Process Matching

*Diffuse*

Known $q \longrightarrow$ $\overset{\text{Signal}}{\phantom{x}}$ $\overset{\text{Noise}}{\phantom{x}}$

$$q(x(t)|x(0)) = \mathcal{N}(\alpha(t)x(0); \sigma^2(t)\mathbf{I})$$

$$x(t) = \alpha(t)x(0) + \underbrace{\sigma(t)\varepsilon}_{\int_0^t s_\tau \, dB_\tau} \quad \varepsilon \sim \mathcal{N}(0,1)$$

$\overset{\text{Signal}}{\phantom{x}}$ $\overset{\text{Noise}}{\phantom{x}}$

$x(0)$  $x(T)$

$$x(s) \sim q(x(s)|\underbrace{G(t, x(t))}), s < t$$

$$\text{Approximates } x(0)$$

Unknown $p$

*Denoise*

## Negative Log Likelihood

$$-\log p(x(0)) \le -\mathbb{E}_{x(t) \sim q(x(t)|x(0))}[\log p(x(t))] + \mathcal{D}_{\mathrm{KL}}\big(q(x(t)|x(0))\big|p(x(0)|x(t))\big)$$

Relative Entropy $\mathcal{D}_{\mathrm{KL}}(q|p) := -\displaystyle\int_{x \in M} \log(p(x)/q(x)) \, dq(x)$

*Conclusion: Match p towards q*

# Diffusion Model and Process Matching

*Diffuse*

**Signal**   **Noise**   **Signal**   **Noise**

**Known** $q \longrightarrow$ $q(x(t)|x(0)) = \mathcal{N}(\alpha(t)x(0); \sigma^2(t)\mathbf{I})$ $\quad x(t) = \alpha(t)x(0) + \underbrace{\sigma(t)\varepsilon}_{\int_0^t s_\tau \, \mathrm{d}B_\tau} \quad \varepsilon \sim \mathcal{N}(0,1)$

$x(0)$

$x(T)$

$x(s) \sim q(x(s)|\underbrace{G(t, x(t))}), s < t$

$\underbrace{\qquad}$

Approximates $x(0)$

**Unknown** $p$

$\longleftarrow$ *Denoise*

**Loss** $\quad L(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), x(0) \sim \pi, x_t \sim q(x(t)|x(0))} |\underbrace{G(t, x(t); \theta)}_{\frac{x(t) - \sigma(t)\varepsilon_G}{\alpha(t)}} - \underbrace{x(0)}_{\frac{x(t) - \sigma(t)\varepsilon}{\alpha(t)}}|$

$\quad\quad$ **Time** $\quad\quad$ **Dataset** $\quad\quad$ **Noisy Image**

**Residualization**

*Reparameterization*

Ɖauphine | PSL☆ | CEREMADE UMR CNRS 7534 | PR[AI]RIE cnrs PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
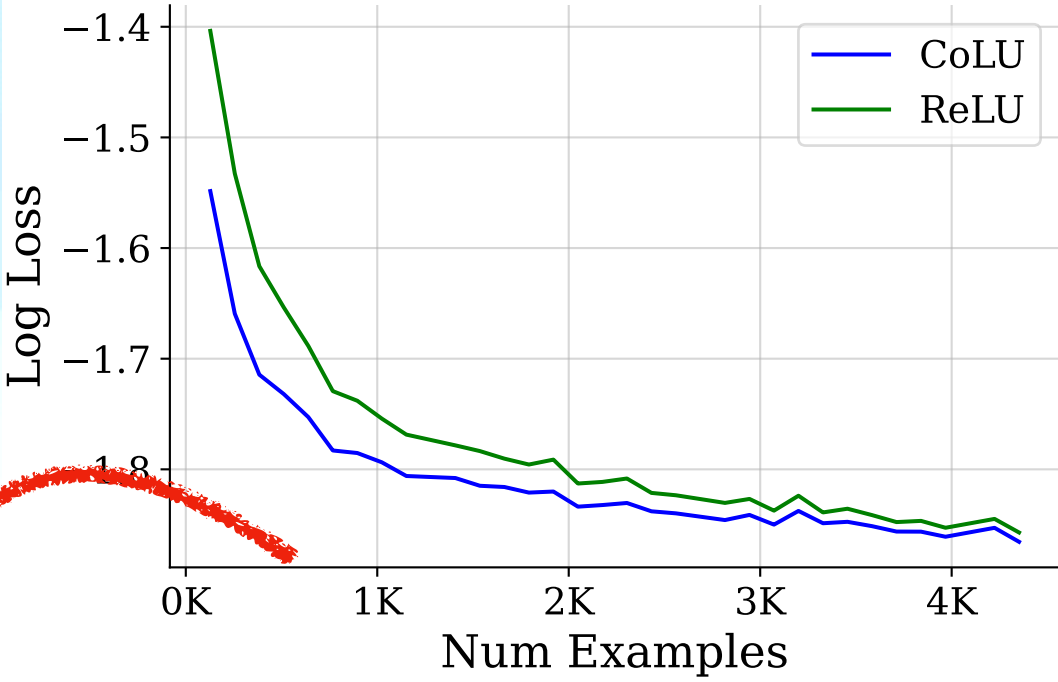
## Better Accuracy and Loss

| 2-Layer MLP (MNIST, C=512) | ReLU | CoLU |
|---|---|---|
| **Train Loss** | 0.0000 ± 0.0000 | 0.0000 ± 0.0000 |
| **Test Accuracy** | 97.17 ± 00.02 | **97.23** ± 00.06 |
| **2-layer VAE (Shared&Soft)** | ReLU | CoLU |
| **Train Loss** | 84.29 ± 0.34 | **83.88 ± 2.68** |
| **Test Loss** | 98.14 ± 0.07 | **97.64 ± 1.39** |

Dauphine | PSL★  CEREMADE UMR CNRS 7534  PR[AI]RIE cnrs  PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.
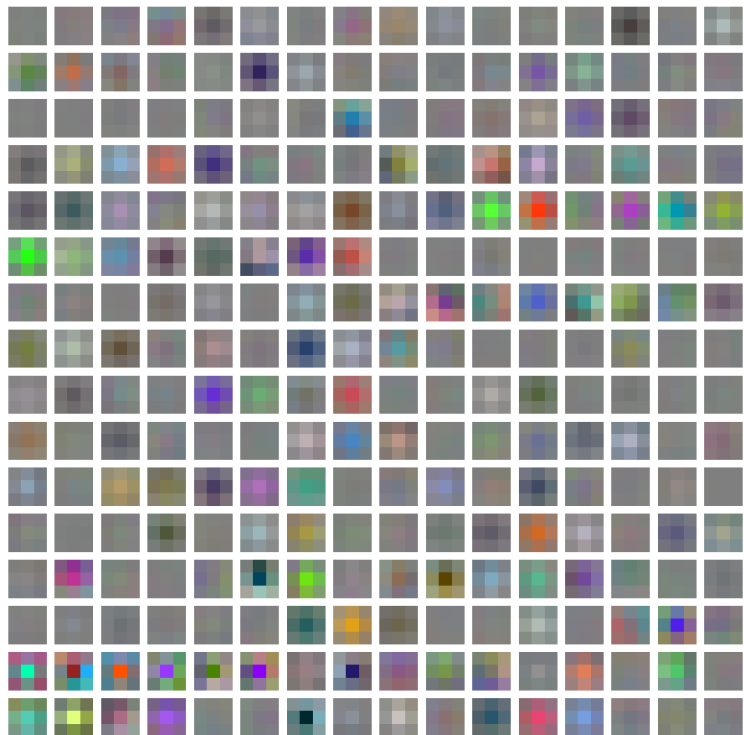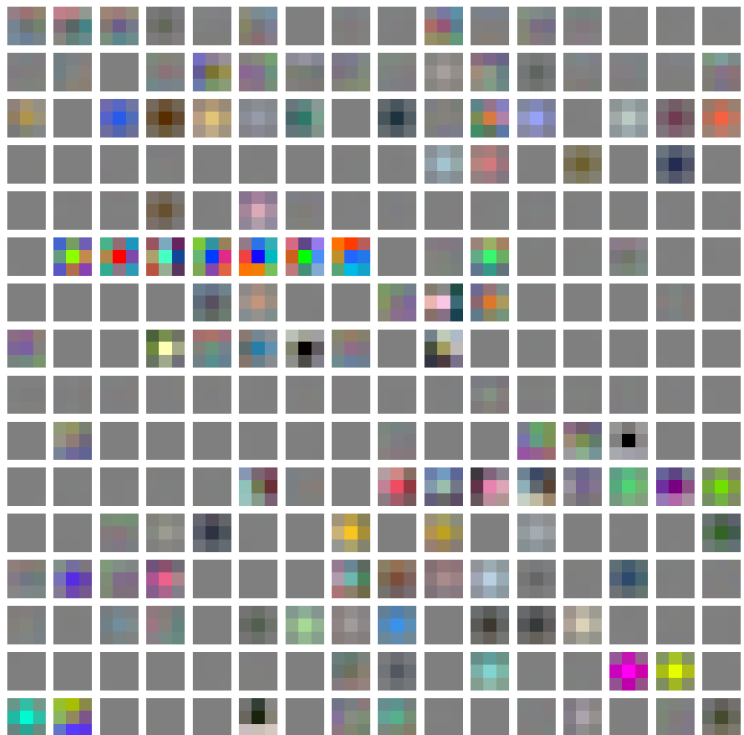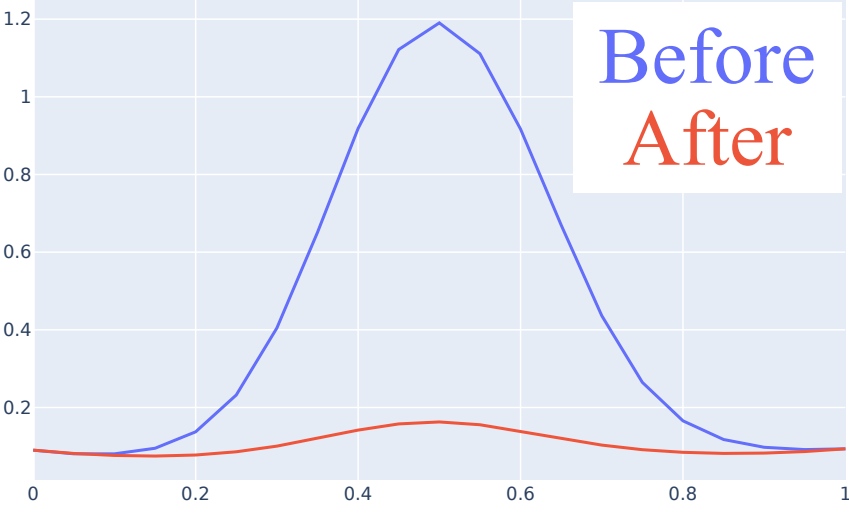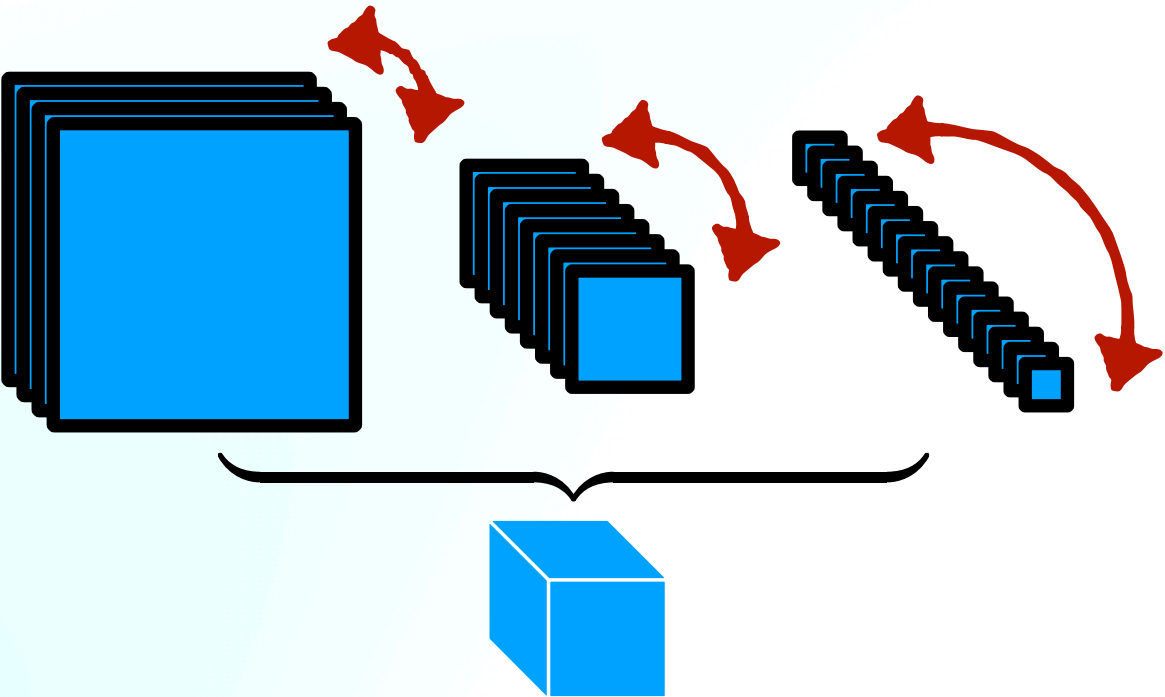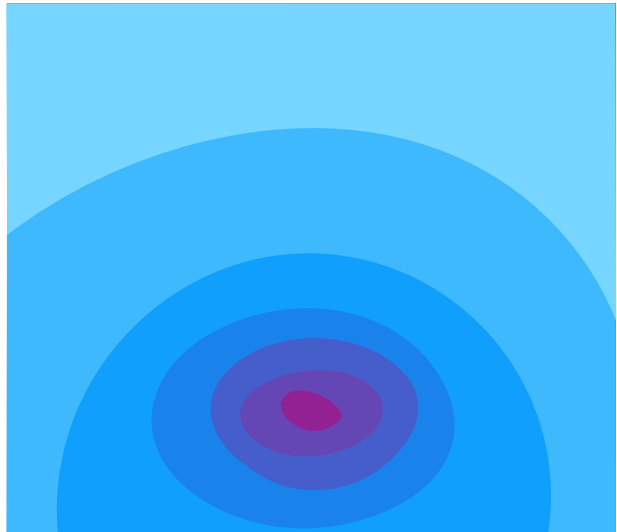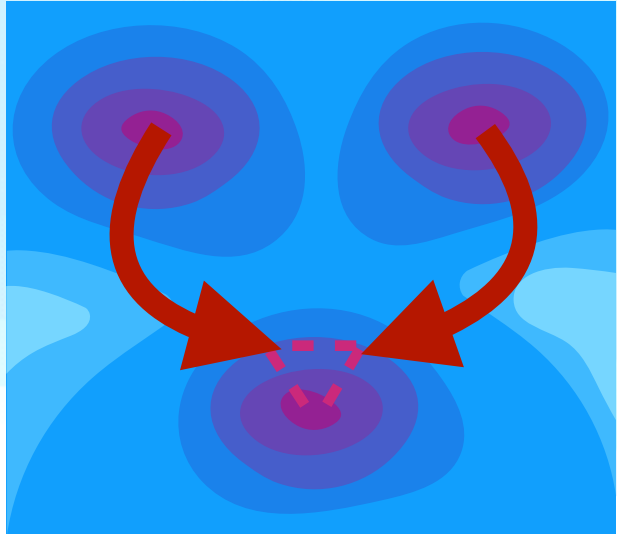
| GPT2 MLP (FineWeb10M) | ReLU | CoLU |
|---|---|---|
| **Forward FLOPs** | 39.064M | 39.101M |
| **Test Loss** | 3.4569 ± 0.1182 | **3.3804** ± 0.1159 |
| **ResNet-56 (CIFAR10)** | ReLU | CoLU |
| **Forward FLOPs** | 0.252M | 0.257M |
| **Test Accuracy** | **92.7282** ± 0.357 | 93.5851 ± 0.442 |
| **Diffusion Model (CIFAR10)** | ReLU | CoLU (Faster) |
| **Train Loss** | 0.1653 | **0.1458** |
| **Early Samples** | | |



**Diffusion Model Training**

O(C) **Complexity**

Negligible **Overhead** vs ReLU

**Đauphine** | PSL★  **CEREMADE** UMR CNRS 7534  PR[AI]RIE cnrs
PaRis Artificial Intelligence Research InstitutE

# Conic Activation Functions

ReLU:
matches each other
with swapping

## Palettes

Last Convolution Layer of Diffusion
Model with Different Seeds

CoLU:

with swapped color
rotation

Conic Activation Functions

# Implication: Generalize NN Symmetry

**Previous works: Linear Mode Connectivity**

- **Optimization**: Non-convex loss is convexified by the quotient.

- **Geometry**: Neural network symmetry induced by activation.

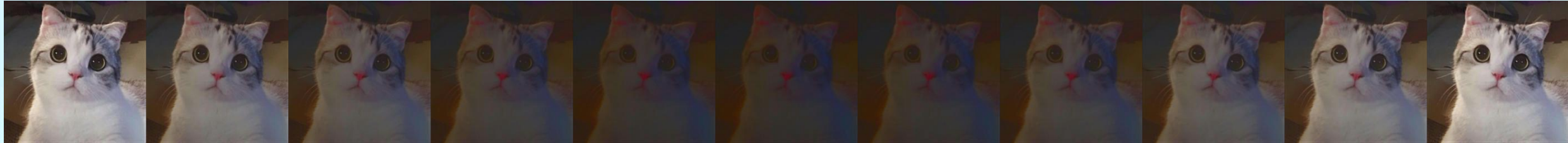- **Probability**: Optimal mixture irrelevant of initializations.



Before
After

Loss Barrier

16

# Conic Activation Functions

## Linear Mode Connectivity: Generative Models

Before



After



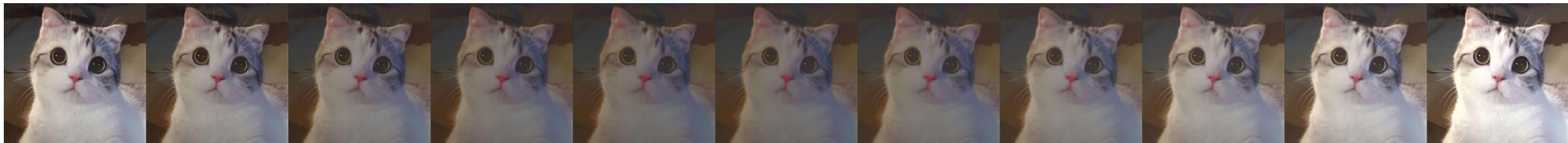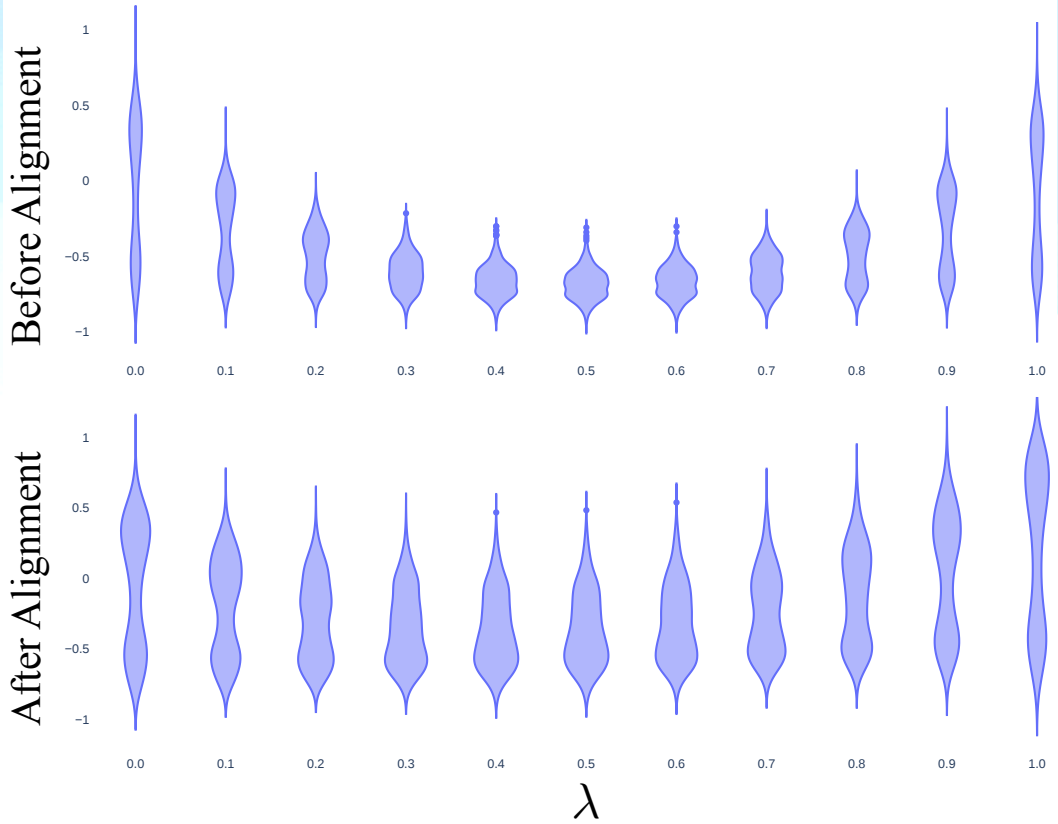Outputs of CNN with
interpolated parameters



Image
Histogram

A symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.

# Linear Mode Connectivity: Generative Models



Animation: interpolation between parameters
in a finetuned diffusion model

# Conclusion

CoLU is a symmetry constraint on generative models for improved **generalization property** and better **learning and performance**.