# On Space Folds
# by Neural Networks

Michal Lewandowski

16.10.2025

# Co-authors

Hamid
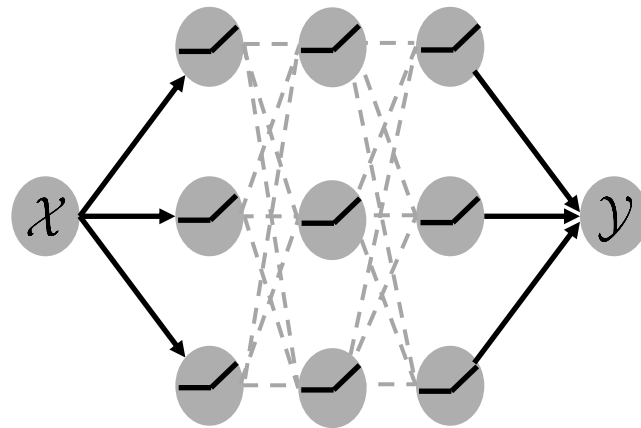(Meta)

Raphael
(SCCH)

Bernhard
(SCCH)

Bernhard
(JKU & SCCH)

- **SCCH**: a Research Center in Hagenberg, Austria, focused on applied AI and software sciences
- **S3AI**: a research project focused on geometry of neural networks (2020-2024)
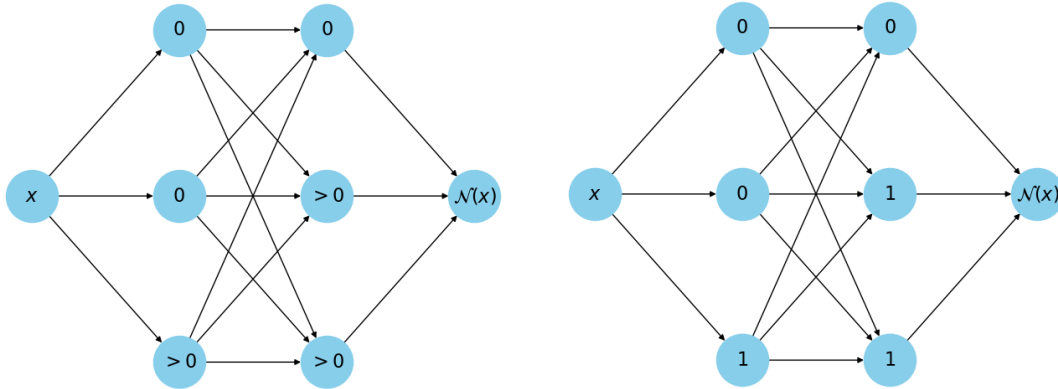
# Introduction

- **Focus**: Analysis of the **phenomena of folding** through a novel measure
- **How:** We focus on fully connected neural networks (MLPs)
- **Goal:** Gaining theoretical insights – not optimizing real-world performance
  - ✓ Universal approximation and exponential expressiveness
  - ✓ Piece-wise linear functions - inherit a geometrical interpretation
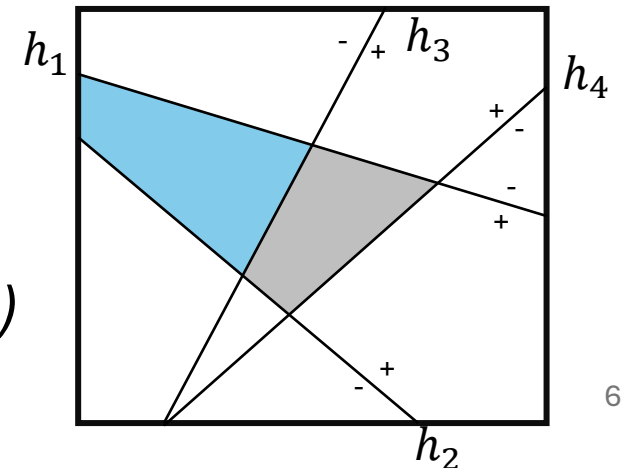


(Reproduced from [1])

[1] H. Boris, and D. Rolnick. (2019). *Deep relu networks have surprisingly few activation patterns*. NeurIPS.

# Hamming Activation Space

- A ReLU NN (denoted $\mathcal{N}$) is an alternating composition of the ReLU ($\sigma(x) := \max(x, 0)$) and affine functions $g_i(\mathbf{x}) := W_i \mathbf{x} + b_i$

- Fix $(W_i, b_i)_{i=1}^L$; push an input $\mathbf{x}$ through $\mathcal{N}$; at every layer we obtain a vector of activation values



- Activation pattern: $\pi_1 = (001011)$
- For any $\mathbf{x}_i \in \mathcal{X}$, there is an associated activation pattern $\pi_i$
- By construction, each (observable) $\pi_i$ is a solution $\{\mathbf{x}: h_1(\mathbf{x}) \geq 0 \ \& \ \dots \ \& \ h_4(\mathbf{x}) \geq 0\}$
- In geometry: _a polytope_ [1], in ML: _a linear region_ [2]



- Hamming distance $d_H$ quantifies difference between $\pi_i, \pi_j \in \{0,1\}^N : d_H(\pi_i, \pi_j) := |\{i: \pi_{i,k} \neq \pi_{j,k}\}|$

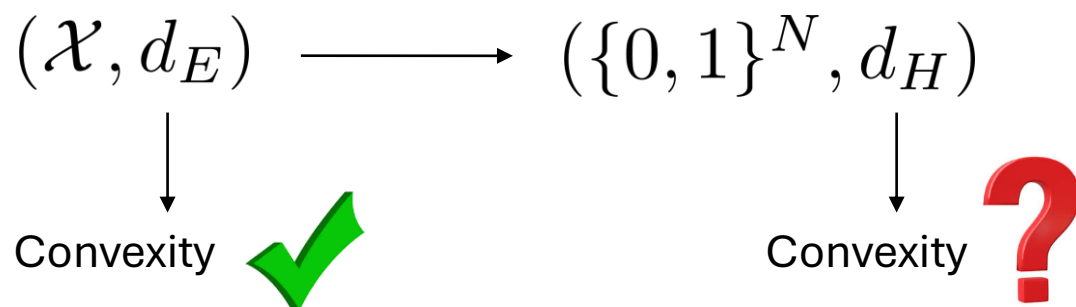- $(\{0,1\}^N, d_H)$ is a metric space *(The Hamming Activation Space)*

[1] Ziegler, G. M. (2000). *Lectures on 0/1-polytopes*. Polytopes—combinatorics and computation.
[2] Montufar et al. (2014). *On the Number of Linear Regions of Deep Neural Networks*. NeurIPS.
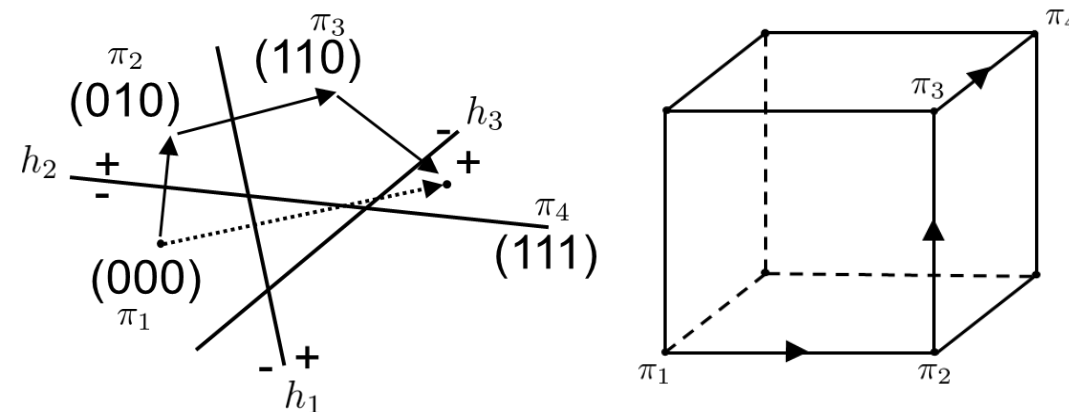
# Analysis of ReLU-Induced Tessellations

The Hamming activation space contains every $\pi \in \{0,1\}^N$.
BUT not every $\pi$ is a pre-image of a linear region!

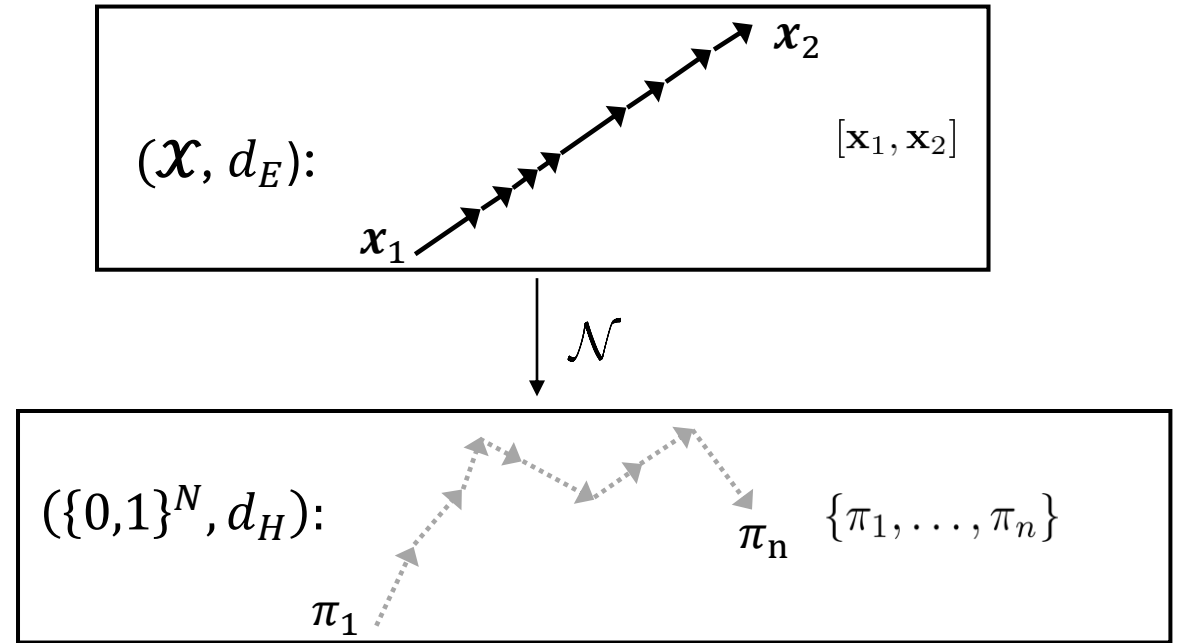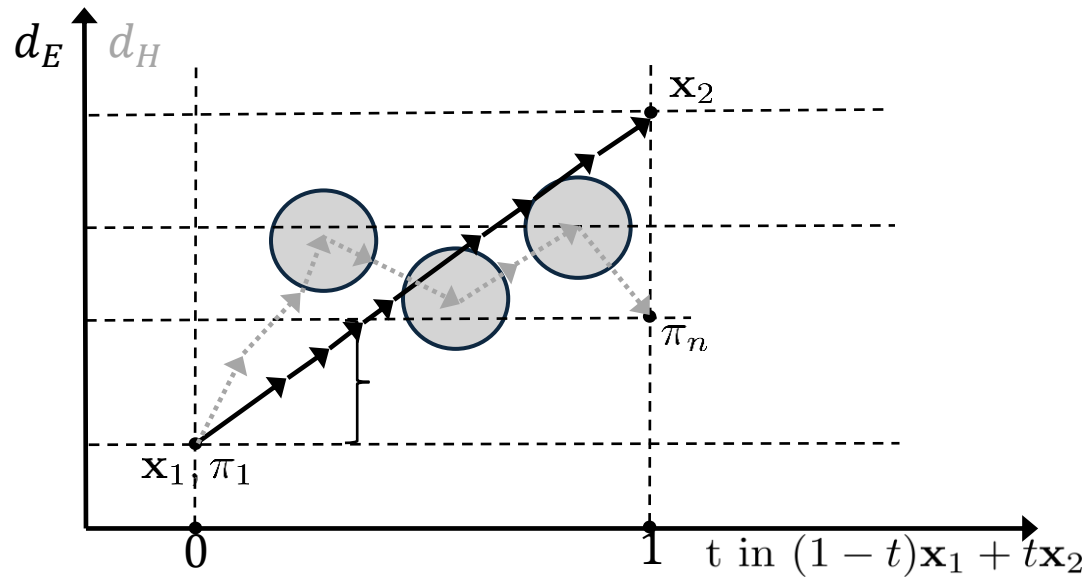Impacts the definition of convexity
in the Hamming space [1, 2]



(010)  (110)  $h_3$
$h_2$  $+$
$-$  $+$
(111)
(100)
(000)
(101)
(011) missing!
(001) $+$
$-$
$h_1$

$(\mathcal{X}, d_E) \longrightarrow (\{0,1\}^N, d_H)$

Convexity ✅

Convexity ❓



$\pi_3$
$\pi_2$  (110)
(010)
$h_2$  $+$  $h_3$
$-$  $-$  $+$
$\pi_4$
(000)  (111)
$\pi_1$
$-$  $+$
$h_1$

$\pi_4$
$\pi_3$
$\pi_1$  $\pi_2$

**Definition 2.** A subset $S$ of the Hamming cube $H^N$ is convex if, for every pair of points $\pi_i, \pi_j \in S$, all (observable) points on every shortest path between $\pi_i, \pi_j$ are also in $S$.

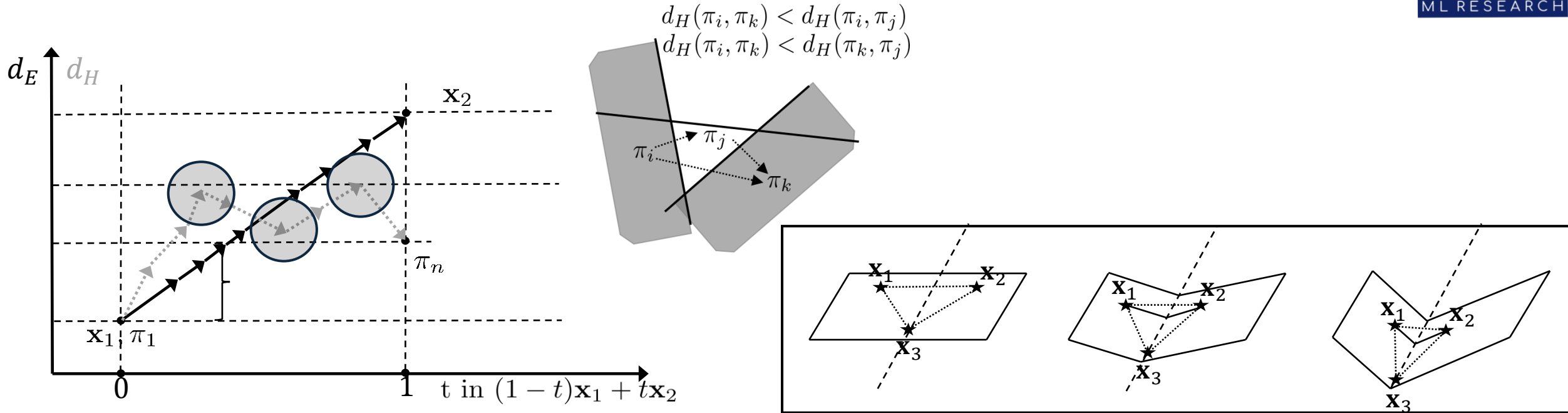**Lemma 1:** Given a space partition into regions $R_{\pi_1}, ..., R_{\pi_r}$ labelled by $A = \{\pi_1, ..., \pi_r\}$, a union $R = \bigcup_{\pi \in S \subset A} R_\pi$ is convex in $R^n$ iff $A$ is convex in $H^N$.

[1] B.A.Moser et al. (2022). *Tessellation Filtering ReLU Networks*. IJCAI.
[2] M.Lewandowski et al. (2025). *On Space Folds of ReLU Neural Networks*. TMLR.

# Mappings of Paths: Connection to Convexity



- $[\mathbf{x}_1, \mathbf{x}_2]$ - straight line in $(\mathcal{X}, d_E)$ – the smallest convex "set"
- We *map* $[\mathbf{x}_1, \mathbf{x}_2]$ from $(\mathcal{X}, d_E)$ to $(\{0,1\}^N, d_H)$ obtaining $\{\pi_1, \dots, \pi_n\}$
- But in $(\{0,1\}^N, d_H)$ not *every* shortest path $\gamma$ (under $d_H$) between $\pi_1$ and $\pi_n$ contains all $\{\pi_1, \dots, \pi_n\}$, i.e., $\gamma \neq \{\pi_1, \dots, \pi_n\} \Rightarrow \{\pi_1, \dots, \pi_n\}$ is not convex (cf. Def 2)
- By investigating $\gamma$ we can measure *deviations* from convexity between $\mathbf{x}_1$ and $\mathbf{x}_2$

# Mappings of Paths: Connection to Folding

$$d_H(\pi_i, \pi_k) < d_H(\pi_i, \pi_j)$$
$$d_H(\pi_i, \pi_k) < d_H(\pi_k, \pi_j)$$

- On a path $\Gamma = \{\pi_1, \ldots, \pi_n\}$, we monitor two *range measures:*
  - $r_1$: Max change in $d_H$ at each step $i$ wrt $\pi_1$: $r_1(\Gamma) := \max_i d_H(\pi_1, \pi_i)$
  - $r_2$: Total travelled distance on the hypercube: $r_2(\Gamma) := \sum_{i=1}^{n-1} d_H(\pi_i, \pi_{i+1})$

**Local Folding**

$$\chi(\Gamma) := 1 - \frac{r_1(\Gamma)}{r_2(\Gamma)} \in [0,1]$$

**Global Folding**

$$\Phi_{\mathcal{N}} := median(\{\chi(\Gamma) : \chi(\Gamma) > 0\})$$

Alternatively, map $\sum_{i,j} \alpha_j \mathbf{x_i}$ through $\mathcal{N}$ for $\forall i, j$

9

# Space Folding Measure – Properties

**Edge cases**:
- $\chi(\Gamma) = 0$ if $d_H(\pi_1, \pi_i)$ increases monotonically.
- $\chi(\Gamma) = 1$ for a looped path: $\max_i d_H(\pi_1, \pi_i) = c$, while $\sum_{i=1}^{n-1} d_H(\pi_i, \pi_{i+1}) \to \infty$.

The space folding measure has the following properties:

1. **[Stability]** Multiple steps in the same linear region do not influence its values.

2. **[Asymmetry]** The folding measure is sensitive to the direction of path traversal, i.e., $\chi(\Gamma) \neq \chi(-\Gamma)$, where $-\Gamma = \{\pi_n, \ldots, \pi_1\}$.

3. **[Flatness Invariance]** $\chi(\Gamma) = 0$ if and only if $\chi(-\Gamma) = 0$.

4. **[Non-additivity]** Is neither sub- nor super-additive.

Let $\Gamma = \Gamma(\mathbf{x}_1, \mathbf{x}_2)$ be a path spanned between the edge points $\mathbf{x}_1, \mathbf{x}_2$.
**Proposition.** Let $d_\chi$ be a symmetrized space folding measure,
$$d_\chi(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{2}\left(\chi\big(\Gamma(\mathbf{x}_1, \mathbf{x}_2)\big) + \chi\big(\Gamma(\mathbf{x}_2, \mathbf{x}_1)\big)\right).$$
Then, $d_\chi$ is a pseudo-metric: **(1) Positivity**: from bounds on $\chi$; **(2) Symmetry**: from construction; **(3) Triangle inequality**: from the triangle inequality of $d_H$

# Algorithm & Complexity

---

**Algorithm 1:** Computation of the Space Folding Measure

**Input:** Two input samples $\mathbf{x}_1$, $\mathbf{x}_2$, the number of intermediate points $n$, the total number of hidden neurons $N$, cost of running the network in the inference mode $O(\mathtt{C})$

**Output:** Space Folding $\chi(\Gamma)$

**Step 1:** Linearly interpolate $\mathbf{x}_1$ and $\mathbf{x}_2$, sampling $n$ points;    `// Sampling Complexity:` $O(n)$

**Step 2:** For each sampled point:

**begin**

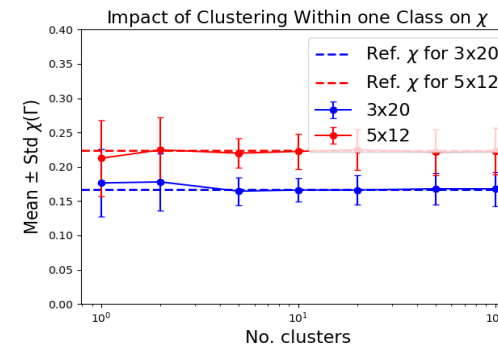  Compute the binarisation;    `// Binarization Complexity Per Point:` $O(\mathtt{C})$

`// Total Binarization Complexity:` $O(n \cdot \mathtt{C})$

**Step 3:** Compute the maximal (from the starting point) and total Hamming distances between intermediate points;    `// Computation of Range Measures Complexity:` $O(n \cdot N)$

**return** Space Folding $\chi(\Gamma)$;    `// Total Algorithm Complexity:` $O(n \cdot (N + \mathtt{C}))$

---

Computational complexity: $O(n \cdot (N + Cost\ of\ Sample\ Propagation) \cdot |C_1| \cdot |C_2|)$

1. For a given the dataset, $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^{m}$
2. Cluster data samples $\mathbf{x}_i$ of the same label
3. Use clusters' centroids for computation of $\chi$



Impact of Clustering Within one Class on $\chi$



Impact of Clustering Within two Classes on $\chi$

# Space Folding - Caveats

- For an activation function $f$, we considered its pre-image thresholded at 0, i.e., $f^{(-1)}\big((-\infty, 0]\big) \to 0$ and $f^{(-1)}\big((0, \infty)\big) \to 1$
- Why thresholding at 0?
  - consider thresholds $a, b \in (-\infty, \infty)$ in the range of the activation function $f$ such that $|a| < |b|$,
  - $\#a$-induced regions $\geq \#b$-induced regions $\Rightarrow a$-induced folding measure $\geq b$-induced folding measure;
  - for ReLU $a = 0$ is "more informative" than any $b > 0$

# Space Folding - Beyond ReLU

- Space Folding computation is based on a walk traversing linear regions
- The computation can be extended to a walk traversing equivalence classes [1,2]:

$$\mathbf{x_1} \sim_{\mathcal{N}} \mathbf{x_2} \Leftrightarrow d_H(\pi_1, \pi_2) = 0$$

- For ReLU NN, $[\mathbf{x}]_{\mathcal{N}} := \{\mathbf{z} : \mathbf{z} \sim_{\mathcal{N}} \mathbf{x}\}$ correspond to linear regions
- For non-monotonous $f$ equivalence classes may be topologically disconnected, but the construction applies as is

[1] N. Shepeleva, et al. (2020). *Relu code space: a basis for rating network quality besides accuracy*. ICLR W.
[2] M. Lewandowski, et al. (2025). The Space Between: On Folding, Symmetries and Sampling. ICLR W.

# Summary

- Same output function yet different architectures?
  - Studied in [1] based on CantorNet [2]; higher folding values for less Kolmogorov-complex representation

> **The question:**
>
> Given a space folding value $\chi(\Gamma) = \tau \in [0, 1]$ for some path $\Gamma$, what can we learn?

- $\chi$ can be seen as a feature of the network
- $\chi$ is upper- and lower-bounded - a reference point across neural networks
- **Interpretation**:
  - Higher $\chi \rightarrow$ indicates a more compact representation
  - Lower $\chi \rightarrow$ indicates potential architecture improvements for the task

[1] M.Lewandowski et al. (2025). On Space Folds of Neural Networks. TMLR.
[2] M.Lewandowski et al. (2024). *CantorNet: A Sandbox For Testing Topological and Geometrical Complexity Measures*. NeurIPS Workshop.

# Space Folding and Generalization

- Folding values increase with the depth of the network *conditioned on* high validation accuracy
- On CIFAR100, we observed a steady increase in global folding during the training process

Beyond fully connected networks

- Consider a **subset of layers** – different underlying tessellation
- **Skip Connections** – applies as is
- **Transformer** – fcnn after the attention head
- **Convolutional layers** – after convolution and pooling

Low sensitivity to batch-norm, dropout
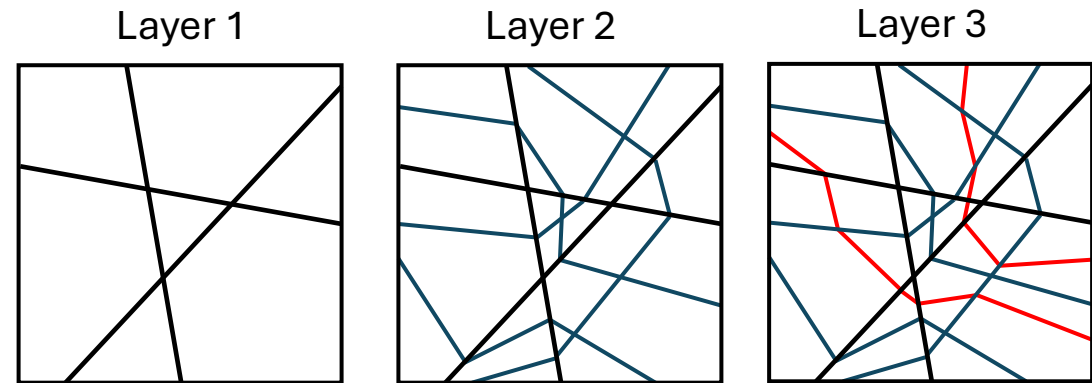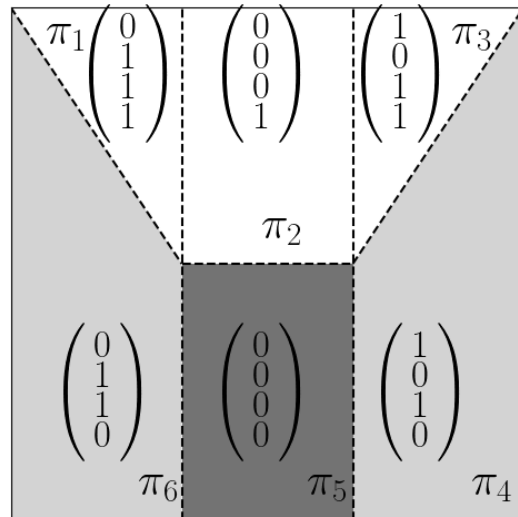
# Convexity in the Hamming Activation Space

**Example 1.** Consider activation patterns $\pi_1 = (0111), \pi_2 = (0001), \pi_3 = (1011)$.
For a walk $(1)\,\pi_1 \to \pi_2, (2)\,\pi_2 \to \pi_3, (3)\,\pi_1 \to \pi_3$, there are intermediate, non-observable activation patterns $\pi_{\text{inter}}$ that we traverse:

$$(1)\,\pi_{\text{inter}} = \{(0011), (0101)\}$$
$$(2)\,\pi_{\text{inter}} = \{(0011), (1001)\}$$
$$(3)\,\pi_{\text{inter}} = \{(0011), (1111)\}$$

Thus, the activation patterns $\{\pi_1, \pi_2, \pi_3\}$ form a convex set in the Hamming cube sense.





(Adapted from [1])

[1] M. Raighu et al. (2017). *On the expressive power of deep neural networks*. ICML.
[2] M.Lewandowski et al. (2024). *CantorNet: A Sandbox for Testing Topological and Geometrical Complexity Measures*. NeurIPS Workshop.

# Interaction Coefficient and Adversarial Examples

$\chi$ is neither sub- nor super-additive - define *interaction coefficient*

$$I(k) \coloneqq |\chi(\Gamma_1 \oplus \Gamma_2) - \chi(\Gamma_1) - \chi(\Gamma_2)|,$$

where $k \in \mathbb{N}$ is the connecting index of paths $\Gamma_1$ and $\Gamma_2$ defined as

$$\Gamma_1 = \{\pi_1, \ldots, \pi_k\}, \Gamma_2 = \{\pi_k, \ldots, \pi_n\}, \Gamma_1 \oplus \Gamma_2 = \{\pi_1, \ldots, \pi_n\}.$$

**Intuition**: folding and adversarial geometry
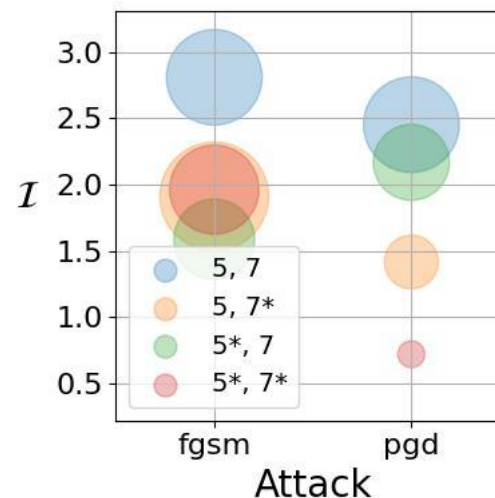
- $I$ computed using digits 5 and 7 from the MNIST test set; $*$ denotes adversarial perturbation
- $I$ is consistently higher for unperturbed samples



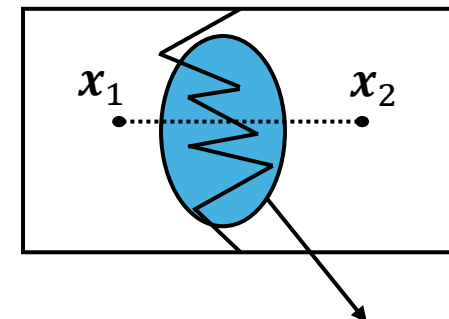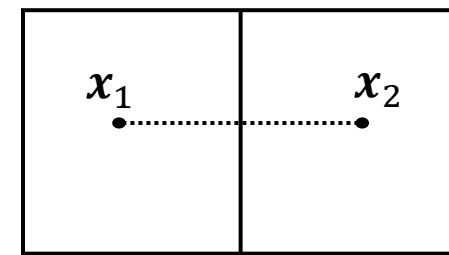Algorithm 1: $\mathcal{I}$ and Adversarial Attacks

**Input**: Dataset $\{(\mathbf{x}_i, y_i)\}_{i \in I}$
**Output**: Mean non-zero $\mathcal{I}$

1: **Step 1**: Adversarially perturb the input $\mathbf{x}$ obtaining $\mathbf{x}^*$.
2: **Step 2**: Assert that $\mathcal{N}(\mathbf{x}) \neq \mathcal{N}(\mathbf{x}^*)$.
3: **Step 3**: Compute $\mathcal{I}$ on a path spanned between $[\mathbf{x}, \mathbf{x}^*]$
4:     (a) check the linear combination of $\mathbf{x}$ and $\mathbf{x}^*$ for vary-
       ing connecting index $k$.
5:     (b) store $k$ s.t. $\mathcal{I}$ increases step-wise
6: **return** $\frac{1}{n-k+1} \sum_{i>k} \mathcal{I}_i$

We can measure this!

# Equivalence of Convexity Notions: Sketch of the Proof

Convexity in $R_n \Rightarrow$ Convexity in the Haming space

- Consider convex $R = \bigcup_{\pi \in A} R_\pi$ in $R^n$.

- Connectivity of $A$: take $\pi_i, \pi_j$ and some points $P \in R_{\pi_i}; Q \in R_{\pi_j}$

- $R = \bigcup_{\pi \in A} R_\pi$ is convex $\Rightarrow [P, Q]$ lies entirely within $R$.
  - Along $[P, Q]$ we cross $h_1, \dots, h_l$ flipping 1-bit in activation patterns at a time.
  - Sequence of flips forms a shortest path in the Hamming space from $\pi_i$ to $\pi_j$.
  - If $A$ wasn't convex, there would exist $\gamma$ connecting $\pi_i$ with $\pi_j$ and leaving $A$.
  - But $[P, Q]$ is a straight line - it corresponds to the minimal sequence of bit flips.