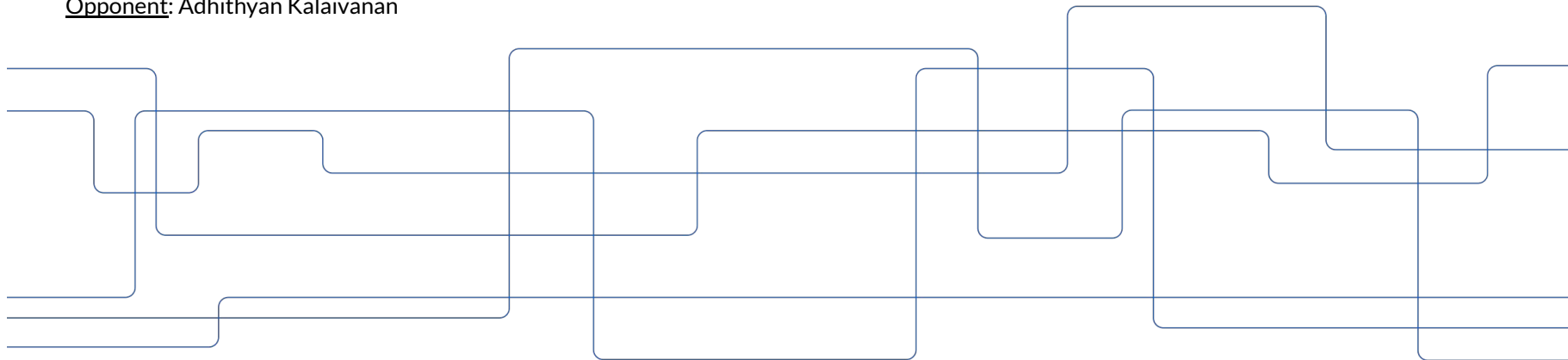# Topological regularization and relative latent representations

Alejandro García Castellanos

Supervisors: Martina Scolamiero, Giovanni Luca Marchetti

Examiner: Florian Pokorny

Opponent: Adhithyan Kalaivanan

# Background

# Overview

# Overview: Relative Latent Representations

**Representation Similarity**

**Model-stitching**

**Relative Representation**

Topological Data Analysis

Topological ML

Topological Densification

# Representation Similarity

How similar are the latent spaces between two random initializations?

Based on statistical similarity metrics:

- CCA
  - SVCCA
  - PWCCA

- **CKA**

"Well-performing" networks tend to have more similar representations

*Wider networks with low-generalization error*

# ε-similar representations


Seed 200

*"Almost isometric up-to-scale"*

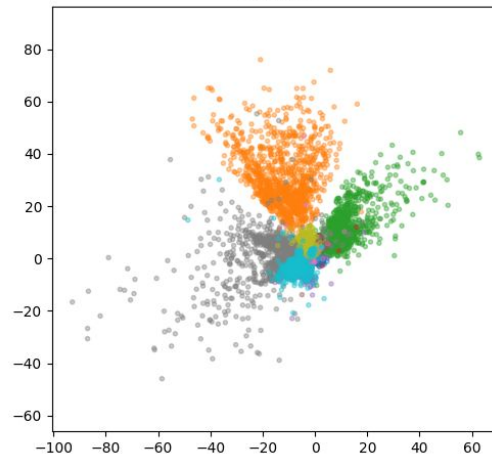> Two representations $X, Y \subseteq \mathbb{R}^n$ are ε-similar if there exist a bijection $T : X \to Y$ s.t. exists $\alpha \in \mathbb{R}^*$ for which $|d(T(x_1), T(x_2)) - \alpha \cdot d(x_1, x_2)| \leq \varepsilon$
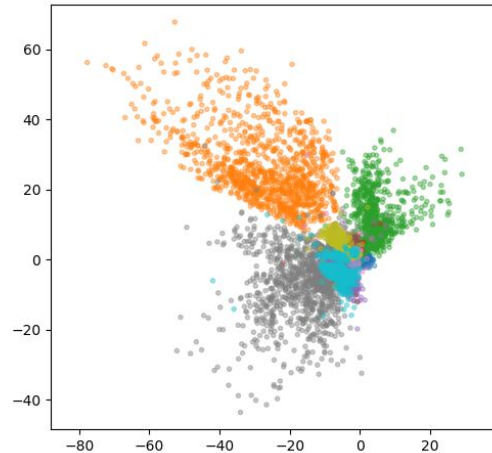

Seed 121

This is mainly based on **empirical evidence**

$\Downarrow$

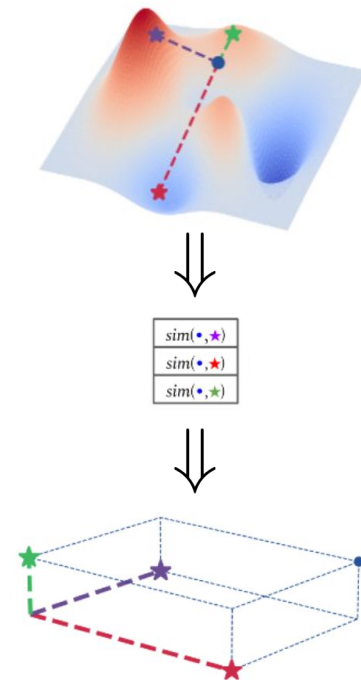Need for theoretical foundation explaining the origin of the ε-similarities

6

# Relative representations

Let $\varphi : \mathcal{X} \to \mathcal{Z}$ the feature extractor component of your network, and $\mathcal{A} = \{a_1, ..., a_k\} \subset \mathcal{X}$ a set of points called *anchors*. Then for any similarity function $sim$ we define the relative representation of a point $x \in S$ w.r.t. $\mathcal{A}$ as
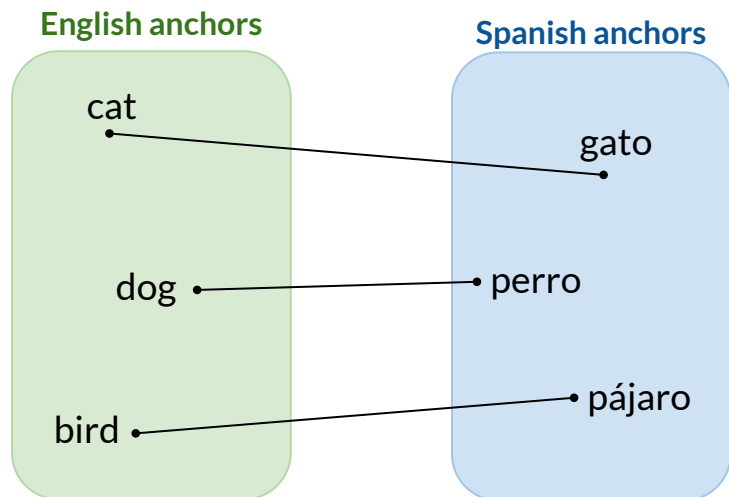
$$(sim(\varphi(x), \varphi(a_1)), ..., sim(\varphi(x), \varphi(a_k)) \in \mathbb{R}^k.$$

When we use the cosine similarity $\to$ we are **invariant to 0-similarities**

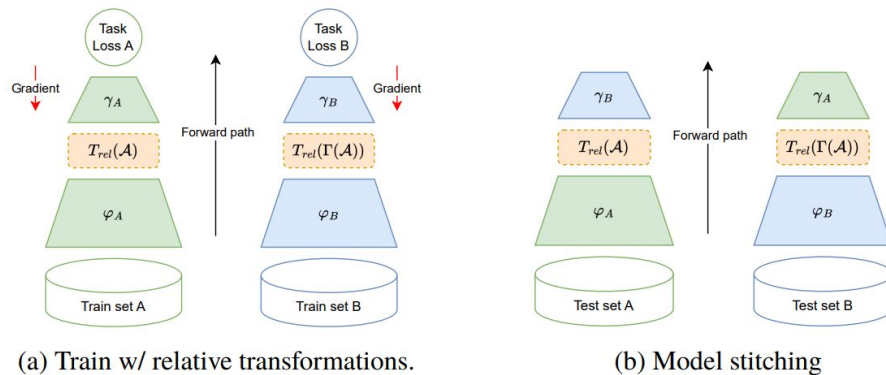[1] L. Moschella, V. Maiorca, M. Fumero, A. Norelli, F. Locatello, and E. Rodolà, "Relative representations enable zero-shot latent space communication," Sep. 2022

# Zero-shot cross-domain model stitching

## Parallel anchors

**English anchors**

**Spanish anchors**

cat

gato

dog

perro

bird

pájaro

## Original training and testing setup



(a) Train w/ relative transformations.

(b) Model stitching

# Overview: Topological Densification

# Topological data analysis

Topological data analysis (TDA) is an approach for the **analysis of the qualitative geometric properties** of datasets using topology techniques.

- Geometric qualitative properties: **connected components**, holes, cavities…

- **Advantages:**
  - Have a sense of the shape of higher-dimensional data that cannot be directly visualized.

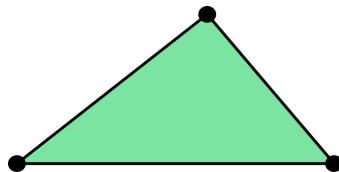  - Results are stable against noise.

# Simplicial complex

- <u>Def:</u> An $k$-simplex $\sigma$ in $\mathbb{R}^d$ with $d \geq k$ is a $k$-dimensional triangle.
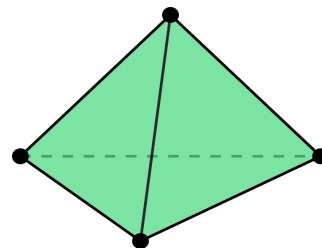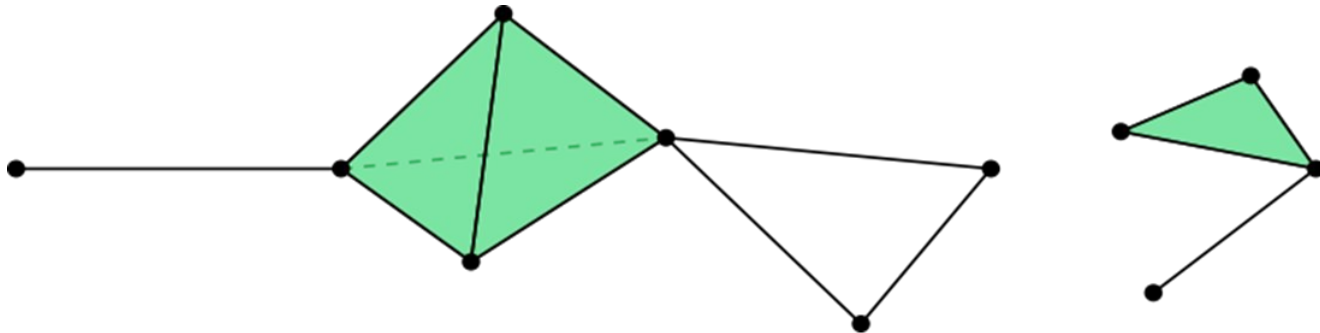

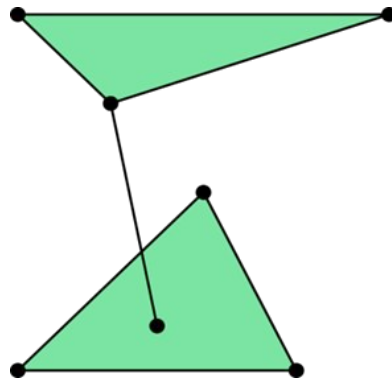
| 0-simplex | 1-simplex | 2-simplex | 3-simplex |

- <u>Def:</u> A *simplicial complex* is a finite collection of simplices $K$ that satisfies that the (non-empty) intersections between the simplices are simplicies of lesser dimension, belonging to the simplicial complex $K$.
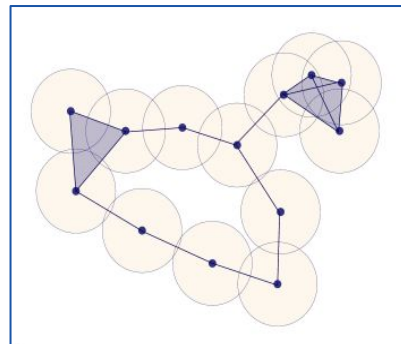
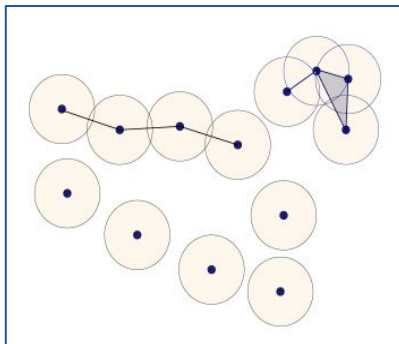**Is a simplicial complex**

**Not a simplicial complex**

# Vietoris-Rips complex

**Definition** Let $X \subset \mathbb{R}^d$ be a finite set of points. We call *Vietoris-Rips* complex of $X$ of radius $r$ to the abstract simplicial complex

$$\text{VR}(X, r) = \{\sigma \subseteq X \mid \text{diam } \sigma \leq r\}$$
$$= \{\{x_0, ..., x_n\} \subseteq X \mid d(x_i, x_j) \leq r \; \forall i, j\} .$$

# Persistent Homology

Each point $(a_i, a_j)$ of the *Persistence Diagram* represents an $l$-dimensional hole that is born at "instant" $a_i$ and dies at $a_j$

# Topological Densification

**High likelihood of β-connected** $\longrightarrow$ **Mass attract mass**



- Equal to having all $H_0$(VR) homology death-times in (0, β)

- Can be enforced with **regularization**:

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda \mathcal{L}_{\beta}, \; \lambda > 0$$

where,

$$\mathcal{L}_{\beta} = \sum_{i=1}^{n} \sum_{d \in \dagger(\mathcal{B}_i)} |d - \beta|$$



- Condensate, for each class, its push-forward distributions inside their decision boundary

- **Reduce generalization error**

# Latent space similarity study

# **Theoretical:** Intertwiner Groups

> Let $G_{\sigma_{n_i}}$ denote the set of invertible linear transformations that exhibit equivalent transformations before and after the nonlinear layer $\sigma_{n_i}$, i.e.,
>
> $$G_{\sigma_{n_i}} \equiv \{A \in GL_{n_i}(\mathbb{R}) \mid \exists B \in GL_{n_i}(\mathbb{R}) \text{ s.t. } \sigma_{n_i} \circ A = B \circ \sigma_{n_i}\}$$

- All elements are of the form $PD$ where $P \in \Sigma_n$ and $D$ is diagonal

- Symmetries in weight space

  $\Downarrow$

  Symmetries in latent representations

## Robust relative transformation

We apply BatchNorm without the learnable affine transformation before computing the cosine sim.

$\Downarrow$

Invariant to intertwiner group actions and 0-similarities

# **Numerical analysis:** 2-dimensional autoencoder



CKA analysis (Linear layers: 2)

Minimum Frobenius distance analysis (Linear layers: 2)

# **Procrustes analysis:** 2-dimensional autoencoder

CKA analysis (Linear layers: 512-256-128-32)



Minimum Frobenius distance analysis (Linear layers: 512-256-128-32)

# **Numerical analysis:** classifier



CKA analysis (Linear layers: 512-256-128-32)

Minimum Frobenius distance analysis (Linear layers: 512-256-128-32)

# Cross-domain model-stitching analysis

# Multilingual model-stitching setup

Investigate the impact of topological densification on zero-shot stitching performance while using relative representations

## Pre-relative     ## Post-relative     ## Both

# Topological densification dataloader



**Debiasing trick:**

1. **Freeze Linear and LayerNorm** modules and set **BatchNorm1d and LayerNorm to training** mode

2. Pass the "random" mini-batch

3. **Unfreeze Linear and LayerNorm** modules and set **BatchNorm1d and LayerNorm to eval** mode

4. Pass the remaining mini-batches

# Baselines

| **Full-finetune** | **Biased dataloader + Debiasing trick** |
|---|---|
| ● **Relative:** better overall<br><br>● **Absolute:**<br>  ○ Better non-stitching<br>  ○ Worse stitching | ● Slightly worse results<br><br>● Enables topological regularization |

| Decoder | Encoder | Absolute | | | Relative | | |
|---|---|---|---|---|---|---|---|
| | | Acc × 100 | FScore × 100 | MAE × 100 | Acc × 100 | FScore × 100 | MAE × 100 |
| en | en | $59.08 \pm 0.20$ | $59.08 \pm 0.85$ | $48.47 \pm 0.64$ | $61.30 \pm 0.28$ | $60.84 \pm 0.77$ | $44.87 \pm 0.92$ |
| | fr | $35.06 \pm 4.36$ | $31.39 \pm 4.62$ | $101.75 \pm 4.26$ | $48.48 \pm 0.08$ | $48.74 \pm 0.20$ | $59.26 \pm 0.37$ |
| fr | en | $27.04 \pm 6.14$ | $25.86 \pm 5.75$ | $115.04 \pm 9.79$ | $60.87 \pm 1.15$ | $60.25 \pm 1.63$ | $45.08 \pm 1.87$ |
| | fr | $48.74 \pm 0.62$ | $48.99 \pm 0.06$ | $62.53 \pm 0.92$ | $49.37 \pm 0.30$ | $50.07 \pm 0.19$ | $58.24 \pm 0.79$ |

# Pre-relative topological densification

The relative transformation is not always cluster-preserving

# Post-relative topological densification

High mismatch of $H_0$ homology $\rightarrow$ Potential information bottleneck

EN relative: VR $H_0$ pers w/ $L^2$ (post, $\lambda = 0.1$, $\beta = 3$)



All death times (Pre mean 13.91. Post mean 4.16)

Max death times (Pre max mean 18.35. Post max mean 7.72)

# Both pre and post-relative topological densification

Vanilla



Topological densified

# **Topological densification:** results

**Vanilla**

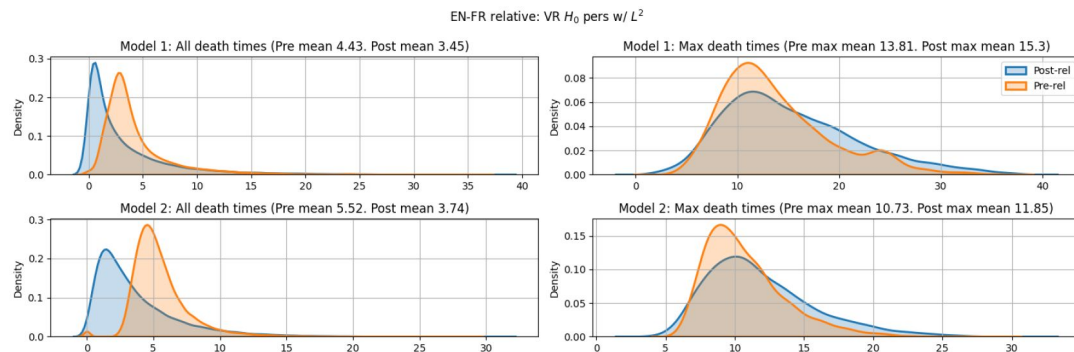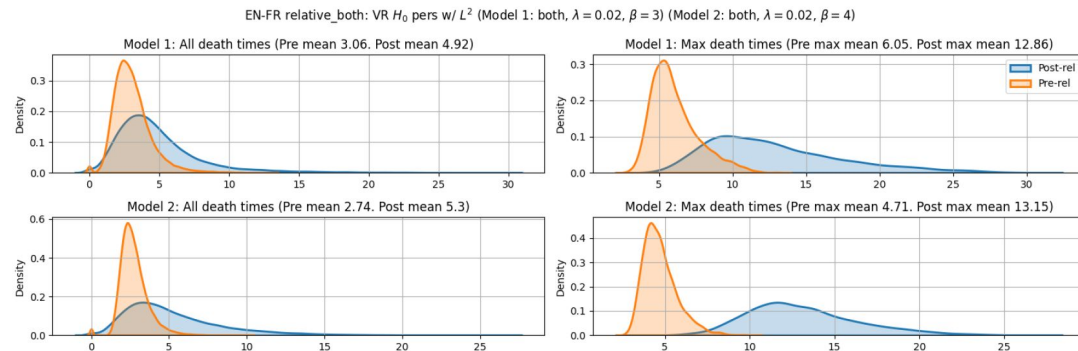| Decoder | Encoder | Absolute | | | Relative | | |
|---|---|---|---|---|---|---|---|
| | | Acc × 100 | FScore × 100 | MAE × 100 | Acc × 100 | FScore × 100 | MAE × 100 |
| en | en | $59.08 \pm 0.20$ | $59.08 \pm 0.85$ | $48.47 \pm 0.64$ | $61.30 \pm 0.28$ | $60.84 \pm 0.77$ | $44.87 \pm 0.92$ |
| | fr | $35.06 \pm 4.36$ | $31.39 \pm 4.62$ | $101.75 \pm 4.26$ | $48.48 \pm 0.08$ | $48.74 \pm 0.20$ | $59.26 \pm 0.37$ |
| fr | en | $27.04 \pm 6.14$ | $25.86 \pm 5.75$ | $115.04 \pm 9.79$ | $60.87 \pm 1.15$ | $60.25 \pm 1.63$ | $45.08 \pm 1.87$ |
| | fr | $48.74 \pm 0.62$ | $48.99 \pm 0.06$ | $62.53 \pm 0.92$ | $49.37 \pm 0.30$ | $50.07 \pm 0.19$ | $58.24 \pm 0.79$ |

**Topological densified**

| Decoder | Encoder | Absolute | | | Relative | | |
|---|---|---|---|---|---|---|---|
| | | Acc × 100 | FScore × 100 | MAE × 100 | Acc × 100 | FScore × 100 | MAE × 100 |
| en | en | $60.20 \pm 0.88$ | $59.69 \pm 0.37$ | $46.33 \pm 0.47$ | $61.25 \pm 0.24$ | $61.37 \pm 0.07$ | $44.50 \pm 0.17$ |
| | fr | $30.04 \pm 0.93$ | $18.56 \pm 1.73$ | $121.52 \pm 16.07$ | $50.14 \pm 0.76$ | $50.55 \pm 0.50$ | $58.81 \pm 0.16$ |
| fr | en | $41.01 \pm 5.53$ | $29.78 \pm 11.70$ | $87.95 \pm 7.62$ | $60.49 \pm 0.78$ | $60.90 \pm 0.54$ | $44.96 \pm 0.34$ |
| | fr | $51.06 \pm 0.00$ | $51.81 \pm 0.04$ | $56.63 \pm 0.01$ | $51.27 \pm 0.01$ | $51.71 \pm 0.19$ | $57.94 \pm 0.74$ |

# **Topological densification**: extra

Having the same densification parameter can benefit the stitching performance

| Decoder | Encoder | Relative | | |
|---------|---------|----------|----------|----------|
| | | Acc × 100 | FScore × 100 | MAE × 100 |
| en | en | $61.25 \pm 0.24$ | $61.37 \pm 0.07$ | $44.50 \pm 0.17$ |
| | fr | $\mathbf{50.90 \pm 0.65}$ | $\mathbf{51.50 \pm 0.66}$ | $\mathbf{57.27 \pm 0.07}$ |
| fr | en | $\mathbf{60.87 \pm 0.95}$ | $\mathbf{61.27 \pm 0.77}$ | $\mathbf{44.56 \pm 0.71}$ |
| | fr | $50.11 \pm 0.38$ | $50.58 \pm 0.79$ | $57.78 \pm 0.14$ |

$L^\infty$ metric for VR filtration $\Rightarrow$ β parameter relates to the optimal spread of the clusters in terms of angle $\Rightarrow$ Helps hyperparameter tuning

# Future work

# Future work

- Investigation of **alternative simplicial complex** constructions: *Lazy witness complex*

- Analysis of **representation similarity in multilingual model stitching**: *CKA analysis*

- Testing topological regularization on large models with **increased GPU VRAM**

- Exploring **other modalities:** *Image-Text*
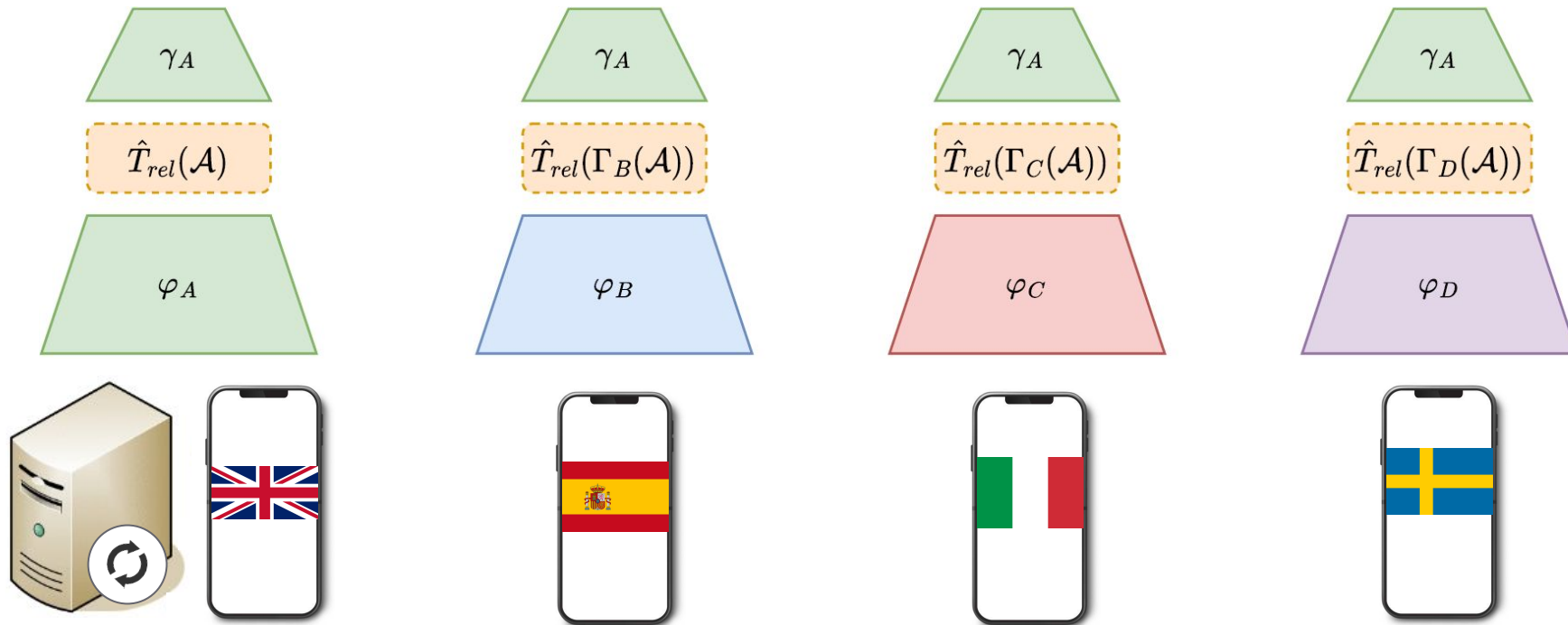
- Exploring **higher dimensional homology**

# Thank you!

# Extra

# Potential use case
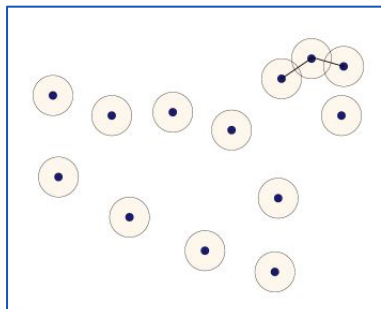
# Simplicial homology

- Algebraic formalism that will allow us to count:
  - **Connected components.**
  - Holes.
  - Cavities.
  - Etc.

- <u>Def:</u> Let $X$ be a geometric object, we define $\beta_i(X)$, the *$i$-th Betti number of $X$*, as the **number of $i$-dimensional holes of $X$**.
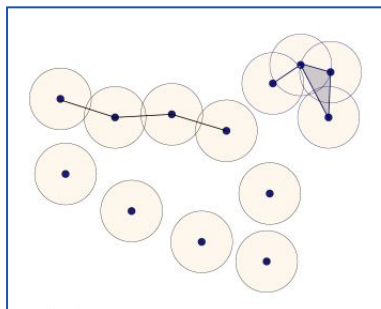
- It will allow us to calculate the Betti numbers of a simplicial complex using linear algebra.

$\beta_0(K) = 11$ *(Connected comp.)*

$\beta_1(K) = 0$ *(Holes)*

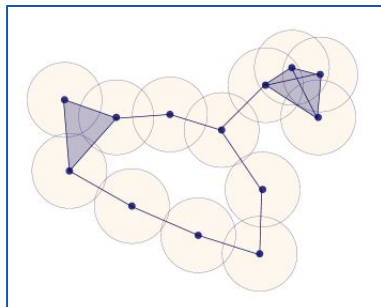$\beta_2(K) = 0$ *(Cavities)*

$\beta_0(K) = 7$ *(Connected comp.)*

$\beta_1(K) = 0$ *(Holes)*

$\beta_2(K) = 0$ *(Cavities)*

$\beta_0(K) = 1$ *(Connected comp.)*

$\beta_1(K) = 1$ *(Holes)*

$\beta_2(K) = 0$ *(Cavities)*

# Intertwiner Groups: properties

Symmetries in weight space $\to$ Symmetries in latent representations

**Proposition** *Suppose $A_i \in G_{\sigma_{n_i}}$ for $1 \leq i \leq k-1$, and let*

$$\widetilde{W} = (A_1 W_1, A_1 b_1, A_2 W_2 \phi_\sigma(A_1^{-1}), A_2 b_2, \ldots, W_k \phi_\sigma(A_{k-1}^{-1}), b_k)$$

*Then, as functions, for each $m$*

$$f_{\leq m}(x, \widetilde{W}) = \phi_\sigma(A_m) \circ f_{\leq m}(x, W),$$
$$f_{>m}(x, \widetilde{W}) = f_{>m}(x, W) \circ \phi_\sigma(A_m)^{-1},$$

*where $f_{\leq m}$ and $f_{>m}$ represent the truncations of the network before and after layer $m$, respectively. In particular, $f(x, \widetilde{W}) = f(x, W)$ for all $x \in \mathbb{R}^{n_0}$.*

# Robust relative transformation

**Definition** — Let $\varphi : \mathcal{X} \to \mathcal{Z} = \mathbb{R}^m$ be our encoder, and $\mathbb{A} \in \mathbb{R}^{d \times k}$, $\mathbb{B} \in \mathbb{R}^{d \times n}$ the matrix representation of $\mathcal{A}$ and $\mathcal{B}$. Then, the *robust relative representation* of $\mathcal{B} \subset \mathcal{X}$ w.r.t. $\mathcal{A}$ is

$$\hat{T}_\varphi(\mathcal{B}, \mathcal{A}) = \left( \widehat{\varphi(\mathbb{A})} D_\mathbb{A} \right)^T \left( \widehat{\varphi(\mathbb{B})} D_\mathbb{B} \right) \in \mathbb{R}^{k \times n} ,$$

where

$$D_\mathbb{A} = \text{Diag} \left( \frac{1}{\sum_{i=1}^m \widehat{\varphi(\mathbb{A})}_{i,1}^2}, \cdots, \frac{1}{\sum_{i=1}^m \widehat{\varphi(\mathbb{A})}_{i,k}^2} \right) ,$$

$$D_\mathbb{B} = \text{Diag} \left( \frac{1}{\sum_{i=1}^m \widehat{\varphi(\mathbb{B})}_{i,1}^2}, \cdots, \frac{1}{\sum_{i=1}^m \widehat{\varphi(\mathbb{B})}_{i,n}^2} \right) ,$$

and $\widehat{\varphi(\mathbb{A})}$ and $\widehat{\varphi(\mathbb{B})}$ represent the respective BatchNorm mean and variance standardizations of the anchor and batch images (without the learnable affine transformation). When the batch and the encoder are implied, we can denote this transformation by $\hat{T}_{rel}(\mathcal{A})$.

# **Numerical similarity metrics:** formulas

$$\mathrm{CKA}(X, Y) = \frac{\left\|\Sigma_{X,Y}\right\|_F^2}{\sqrt{\left\|\Sigma_{X,X}\right\|_F^2 \cdot \left\|\Sigma_{Y,Y}\right\|_F^2}}$$
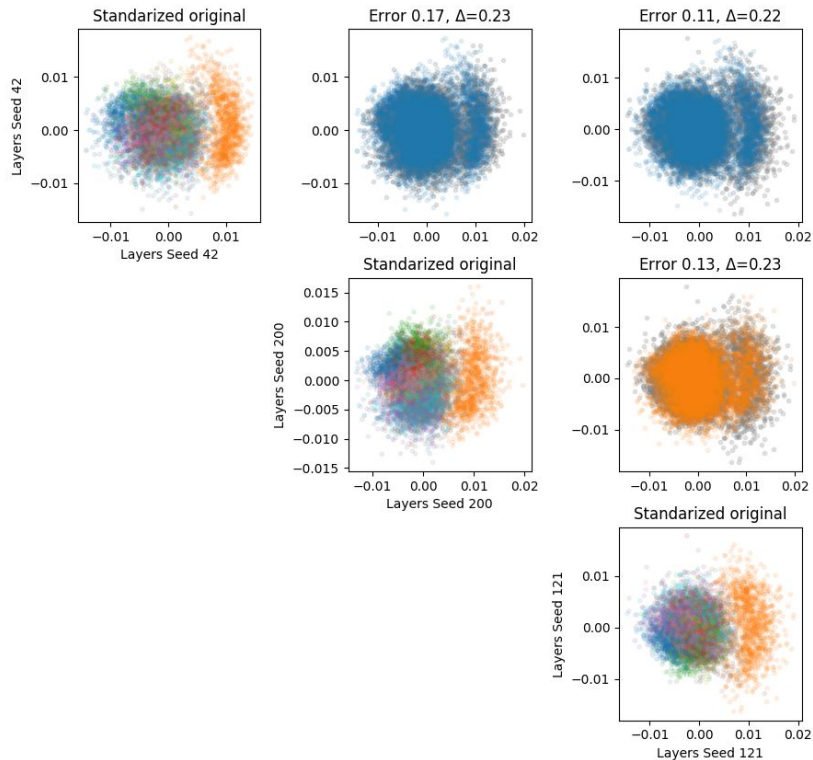
where Σ represents the covariance matrix

$$\mathrm{minFrob}(A, B) = \min_{P \in \Pi} \left\| \frac{A}{\|A\|_F} - P\frac{B}{\|B\|_F} \right\|_F$$
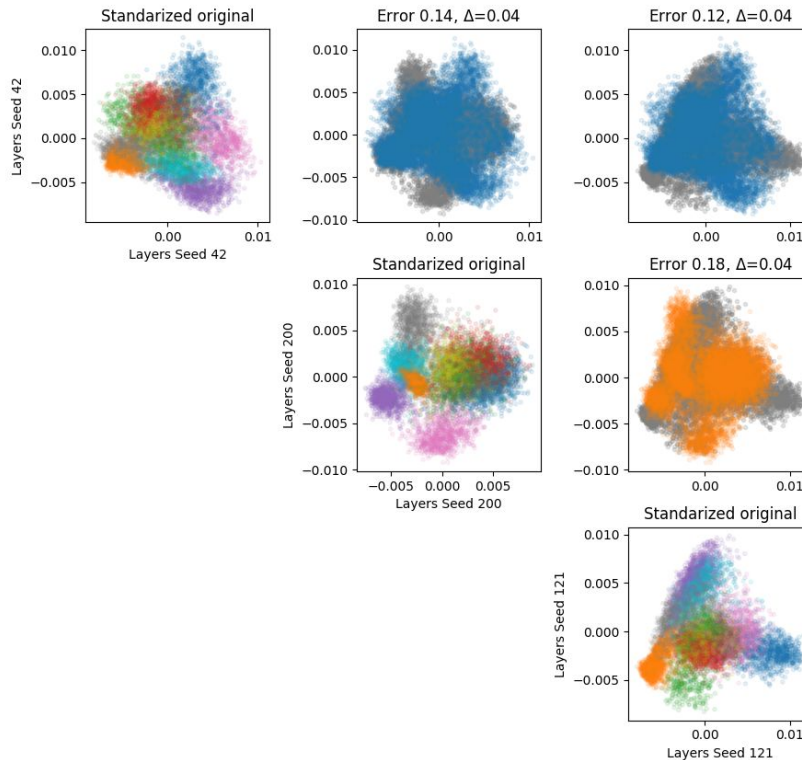
where A and B are the distance matrices of X and Y

# **Procrustes analysis:** 32-dimensional autoencoder



Procrustes analysis (Linear layers: 512-256-128-32) [Projected w/ PCA]
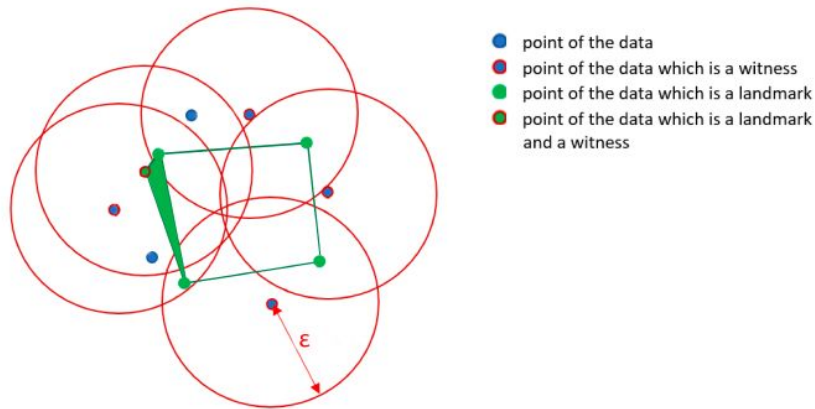
# **Procrustes analysis:** classifier



Procrustes analysis (Linear layers: 512-256-128-32) [Projected w/ PCA]
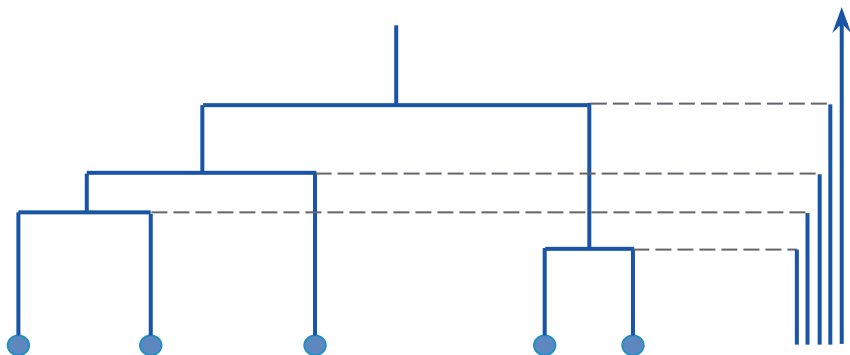
# Lazy Witness complex

**Definition** —— (Nested family of witness complexes [10]). Let $(\mathcal{X}, d)$ be a metric space, $X \subset \mathcal{X}$ be a dataset, $L = \{l_0, ..., l_n\} \subseteq X$ be a set of landmark points, and $\varepsilon > 0$. Then the $k$-simplex $\sigma = \{u_1, ..., u_k\}$ with $u_i \in L$ belongs to the *Lazy Witness complex* $W_\varepsilon(X, L)$ iff all its faces belong to $W_\varepsilon(X, L)$ and there is a witness $x \in X$, such that:

$$\max\{d(u_i, x) \mid u_i \in \{u_1, ..., u_k\}\} \leq \varepsilon.$$



- ● point of the data
- ● point of the data which is a witness
- ● point of the data which is a landmark
- ● point of the data which is a landmark and a witness

# Exploring higher dimensional homology

Single Linkage Hierarchical Clustering $\leftrightarrow$ $H_0$(VR)

Controlling $H_0$(VR) $\rightarrow$ Topological densification

What beneficial properties for classification can we obtain by controlling $H_n$(VR) for n>0?