

Relative Representations: ***Topological and Geometric Perspectives***

Alejandro García Castellanos
University of Amsterdam, AMLab



Representational universality

Different neural networks are
converging towards the **same** way of
representing the world

The Platonic Representation Hypothesis

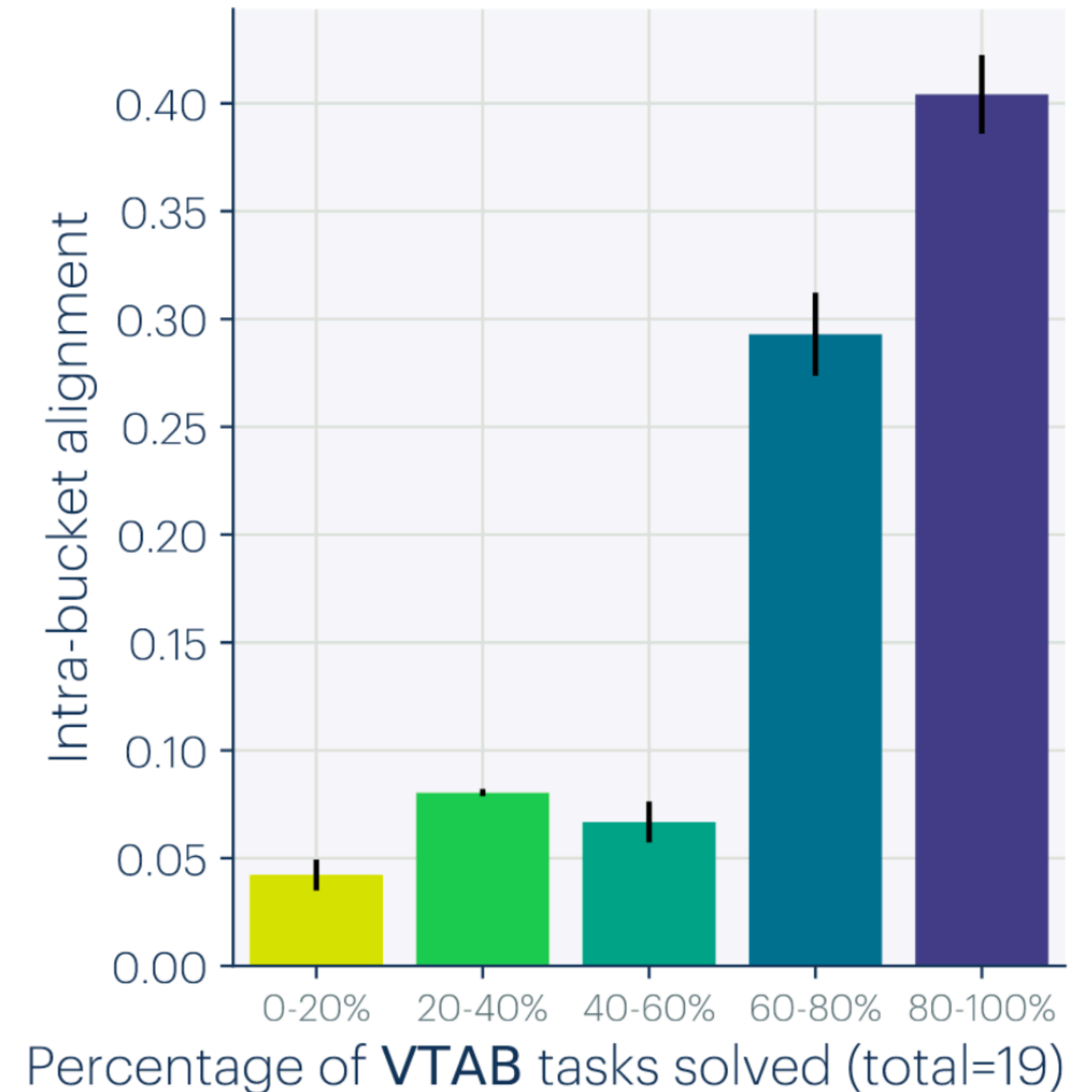
Minyoung Huh^{*1} Brian Cheung^{*1} Tongzhou Wang^{*1} Phillip Isola^{*1}

Evidence of the convergence

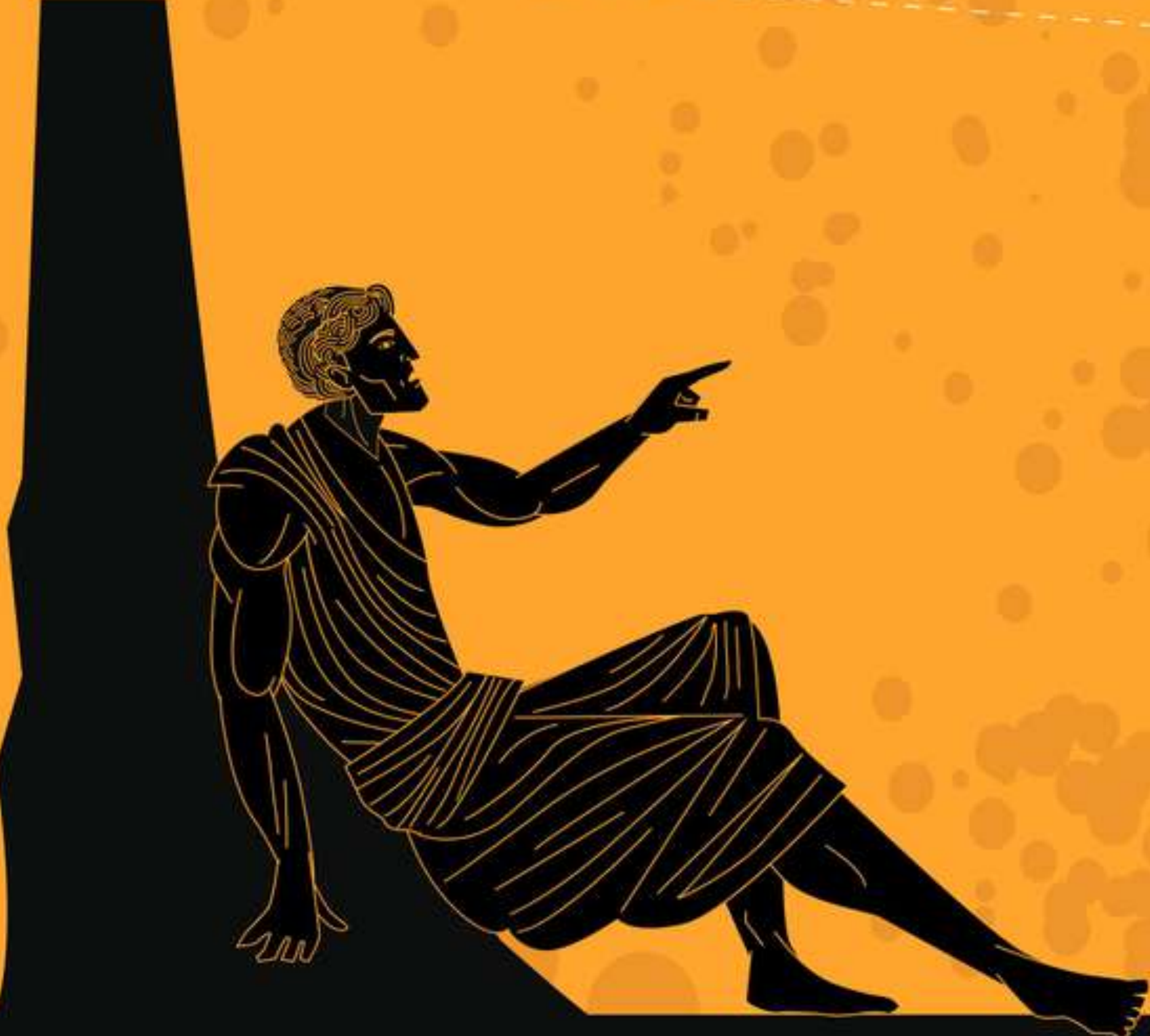
The alignment between vision models increases as vision systems become stronger

- 78 vision models: different architectures, objectives, training data distributions.
- Group models by performance on VTAB, and measure representational similarity within each group

Convergence to general competence



What are we converging to?



The **Platonic** Representation Hypothesis

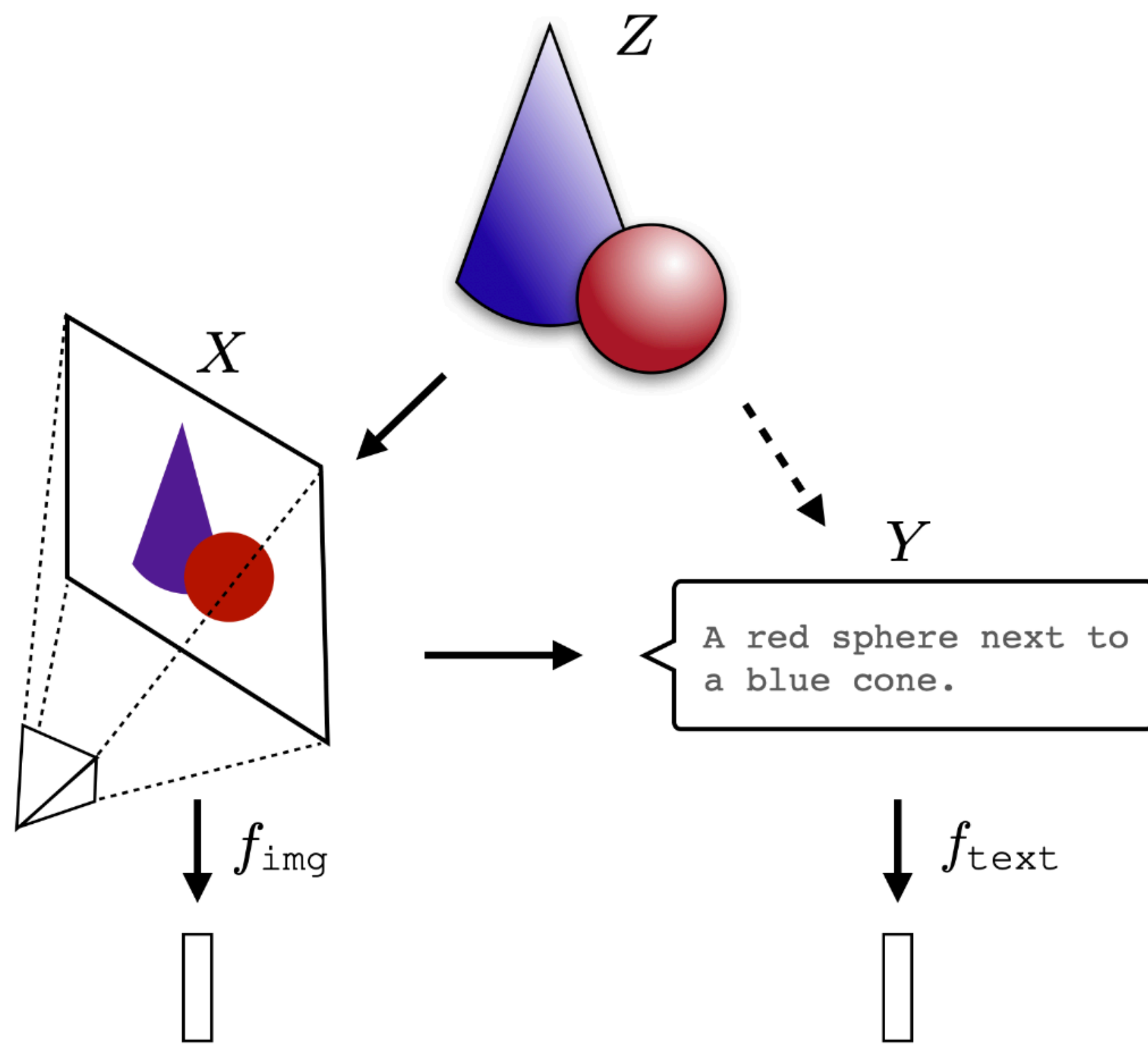


Figure 1. **The Platonic Representation Hypothesis:** Images (X) and text (Y) are projections of a common underlying reality (Z). We conjecture that representation learning algorithms will converge on a shared representation of Z , and scaling model size, as well as data and task diversity, drives this convergence.

Under major assumptions (bijective, discrete observations)



Self-supervised training tasks, such as **contrastive learning**, converge to **statistics of the underlying reality Z**

Relative Representations

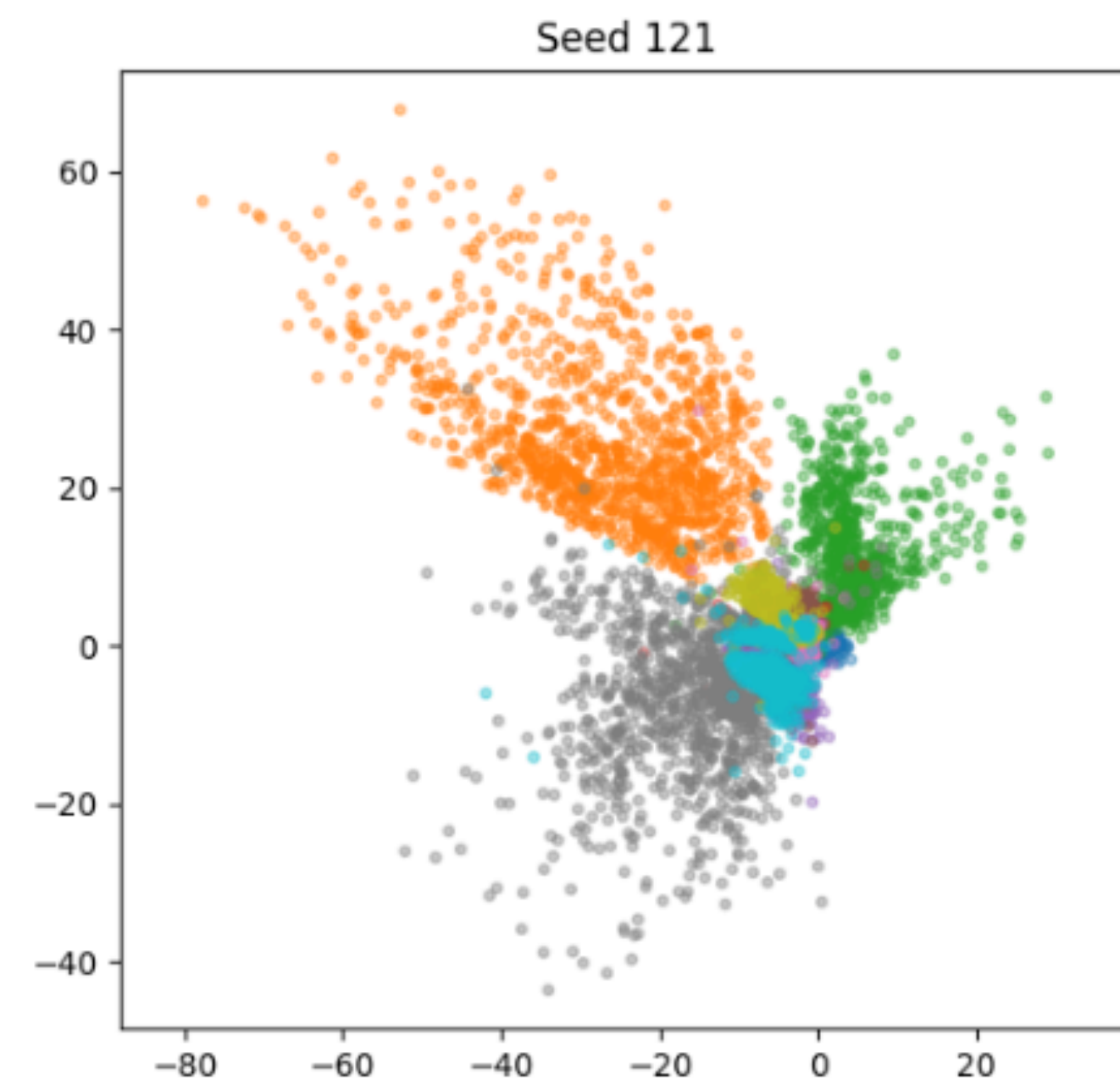
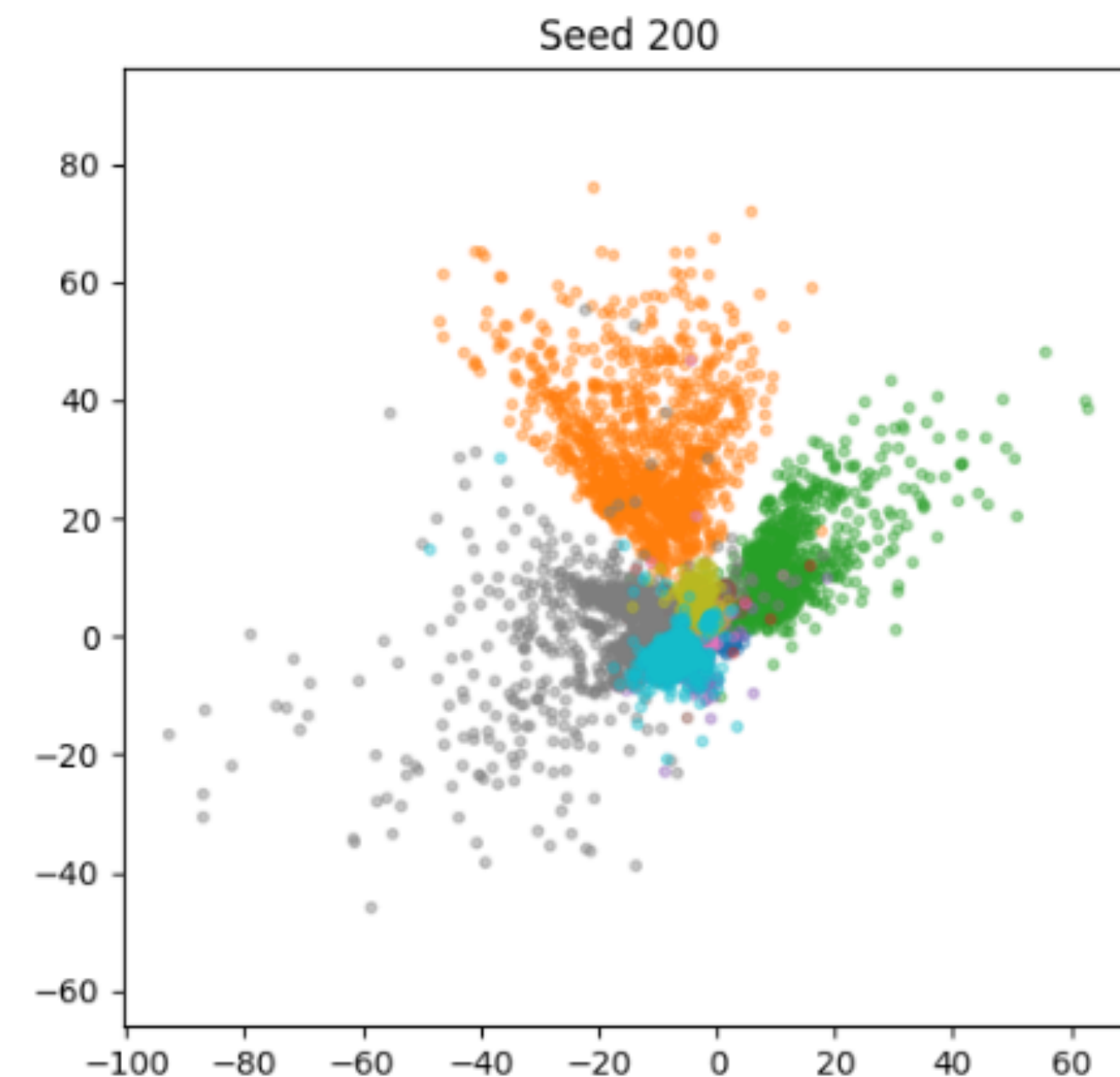
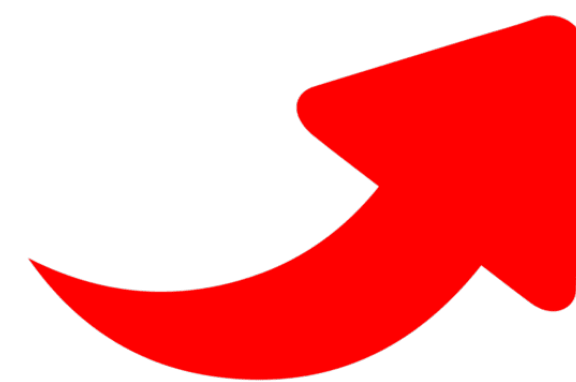
RELATIVE REPRESENTATIONS ENABLE ZERO-SHOT LATENT SPACE COMMUNICATION

Luca Moschella^{1,*} Valentino Maiorca^{1,*}

Marco Fumero¹ Antonio Norelli¹ Francesco Locatello^{2,†} Emanuele Rodolà¹

¹Sapienza University of Rome ²Amazon Web Services

“Almost isometric up-to-scale”



RELATIVE REPRESENTATIONS ENABLE ZERO-SHOT LATENT SPACE COMMUNICATION

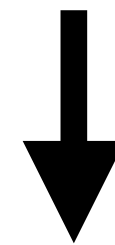
Luca Moschella^{1,*} Valentino Maiorca^{1,*}

Marco Fumero¹ Antonio Norelli¹ Francesco Locatello^{2,†} Emanuele Rodolà¹

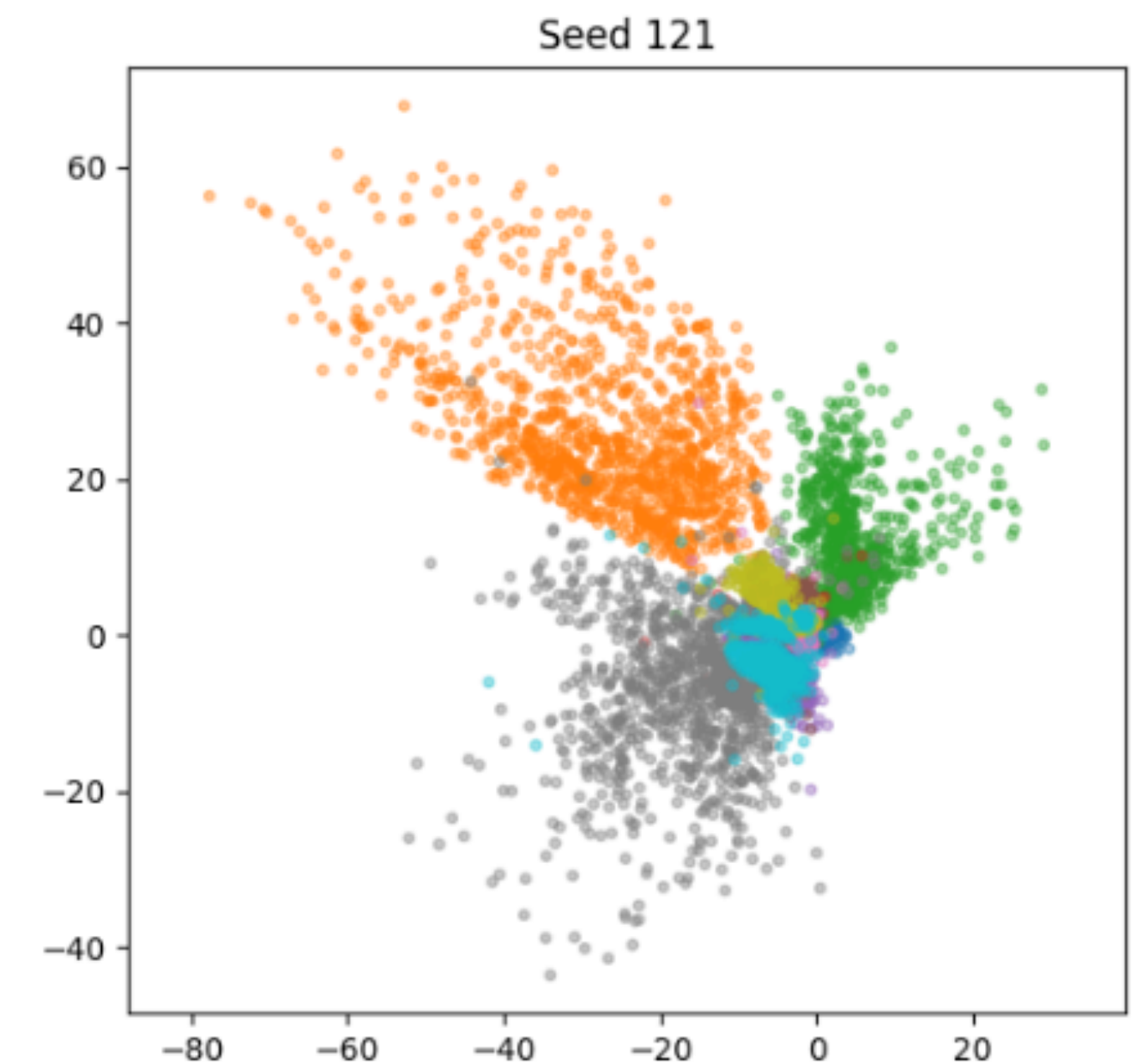
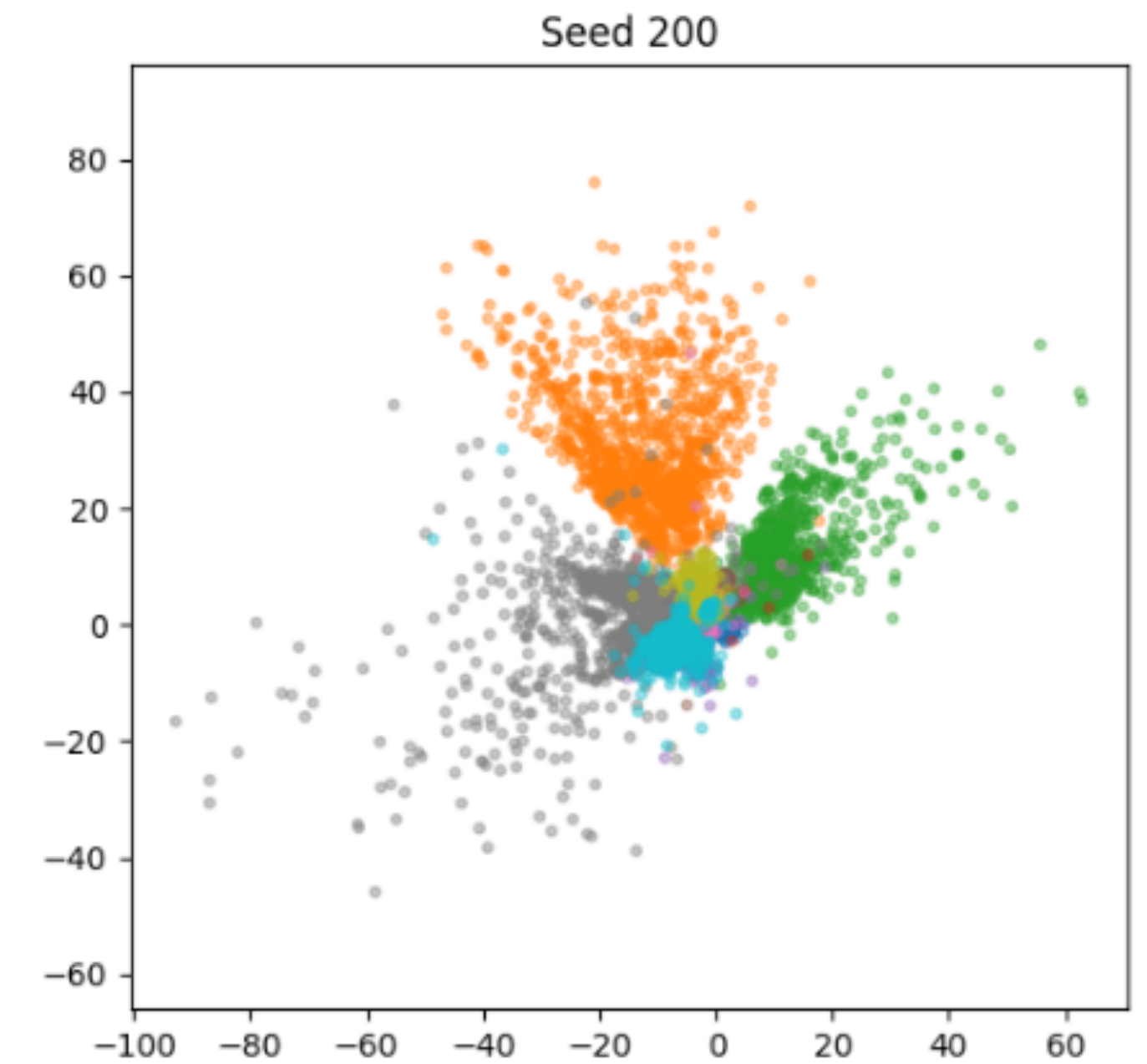
¹Sapienza University of Rome ²Amazon Web Services

“Almost isometric up-to-scale”

This is mainly based on **empirical evidence**



Need for a theoretical explanation of this behavior



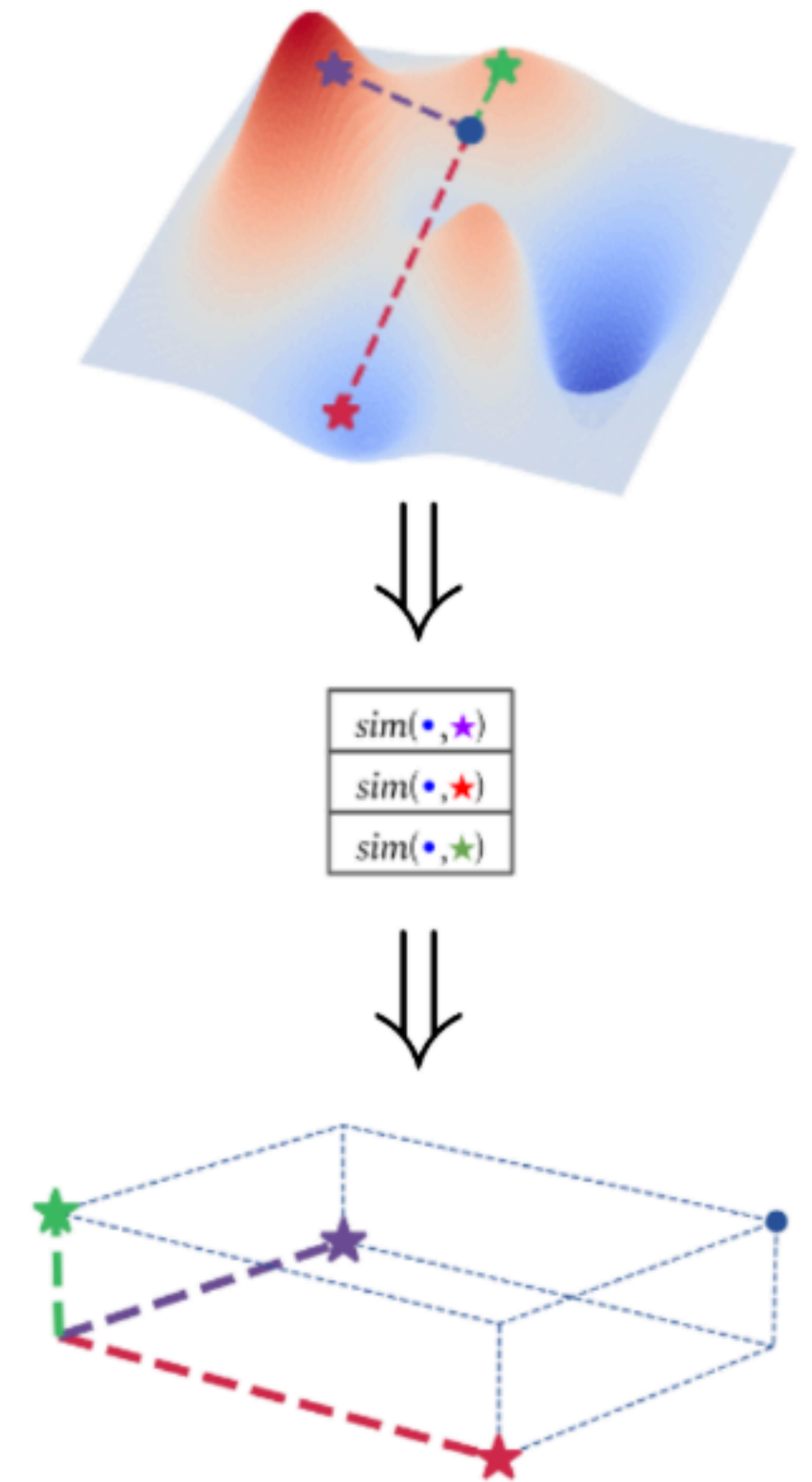
Relative representations

Let $\varphi: \mathcal{X} \rightarrow \mathcal{Z}$ be the feature extractor of the network, and let $\mathcal{A} = \{a_1, \dots, a_k\} \subset \mathcal{Z}$ be a set of elements called *anchors*, and let $\text{sim}: \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ be a similarity function.

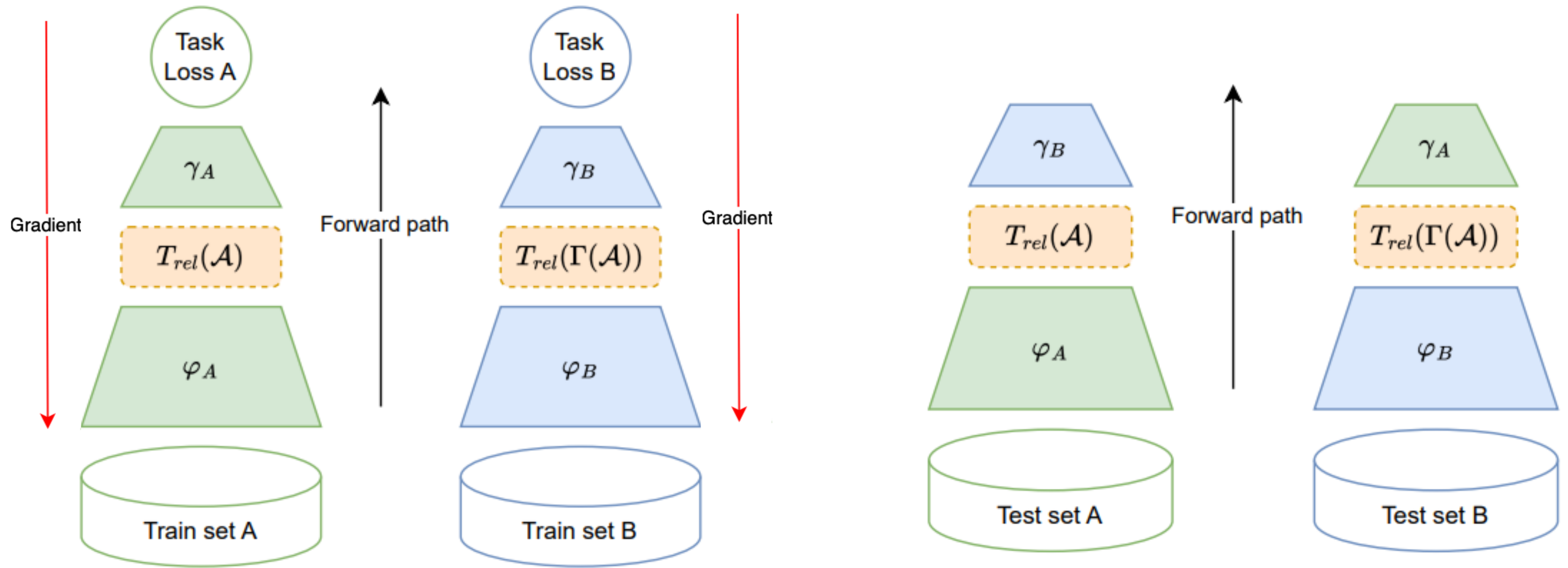
The **Relative Representation** of $z \in \mathcal{Z}$ w.r.t. \mathcal{A} is

$$T_{\text{rel}}^{\mathcal{A}}(z) = (\text{sim}(z, a_1), \dots, \text{sim}(z, a_k)) \in \mathbb{R}^k$$

$\text{sim} = \text{cosine sim} \rightarrow$ Invariant to *isometries + isotropic rescalings*



Zero-shot cross-domain model stitching



(a) Train w/ relative transformations.

(b) Model stitching

Relative Representations: Topological and Geometric Perspectives

Alejandro García-Castellanos *
Amsterdam Machine Learning Lab
University of Amsterdam, Netherlands
a.garciacastellanos@uva.nl

Giovanni Luca Marchetti
Department of Mathematics
KTH Royal Institute of Technology, Sweden

Danica Kragic
Division of Robotics, Perception and Learning
KTH Royal Institute of Technology, Sweden

Martina Scolamiero
Department of Mathematics
KTH Royal Institute of Technology, Sweden

Editors: Marco Fumero, Clementine Domine, Zorah Löhner, Donato Crisostomi, Luca Moschella, Kimberly Stachenfeld

1. Geometric Perspective

Weight Space Symmetries

Neural network **symmetries** are transformations:

$$\psi : \mathcal{G} \times \Theta \rightarrow \mathcal{G} \times \Theta$$

That preserves network function:

$$u_{G,\theta}(\mathbf{x}) = u_{\psi(G,\theta)}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}, \quad \forall (G,\theta) \in \mathcal{G} \times \Theta$$

ψ : permutation

$$u_{G,\theta}(\mathbf{x}) = \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \end{array} = u_{\psi(G,\theta)}(\mathbf{x})$$

ψ : scaling

$$u_{G,\theta}(\mathbf{x}) = \text{---} \text{---} \text{---} = \text{---} \text{---} \text{---} = u_{\psi(G,\theta)}(\mathbf{x})$$

Weight Space Symmetries

Neural network **symmetries** are transformations:

$$\psi : \mathcal{G} \times \Theta \rightarrow \mathcal{G} \times \Theta$$

That preserves network function:

$$u_{G,\theta}(\mathbf{x}) = u_{\psi(G,\theta)}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}, \quad \forall (G,\theta) \in \mathcal{G} \times \Theta$$

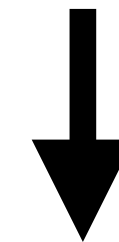
ψ : permutation

$$u_{G,\theta}(\mathbf{x}) = \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \\ \text{---} \text{---} \end{array} = u_{\psi(G,\theta)}(\mathbf{x})$$

ψ : scaling

$$u_{G,\theta}(\mathbf{x}) = \text{---} \text{---} \text{---} = \text{---} \text{---} \text{---} = u_{\psi(G,\theta)}(\mathbf{x})$$

Symmetries in **weight space**




Symmetries in **latent space**

Theoretical explanation for the emergence of structurally-similar representations in networks

Invariance trading: Robust Relative Representation

Robust Relative Representation: We apply *Gaussian normalization* with respect to a batch \mathcal{B} of data, i.e., a simple form of *batch normalization* (without learnable parameters), *before computing the cosine sim*

We are now invariant to shifts + **weight space group** actions

 We trade off invariance to isometries other than permutations with more general non-isotropic rescalings → **Good trade in high dimensional latent spaces:**

Encoder

Classifier

Performance comparison on zero-shot model stitching.

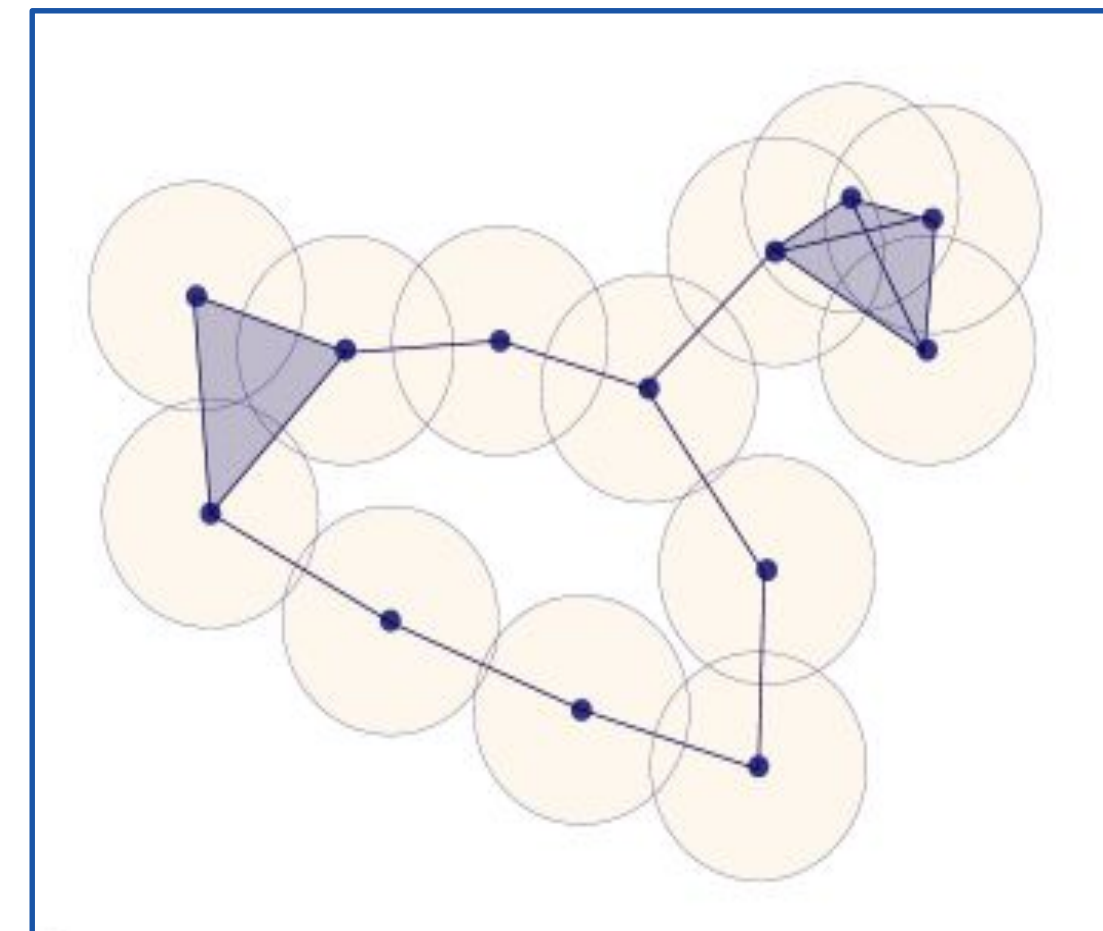
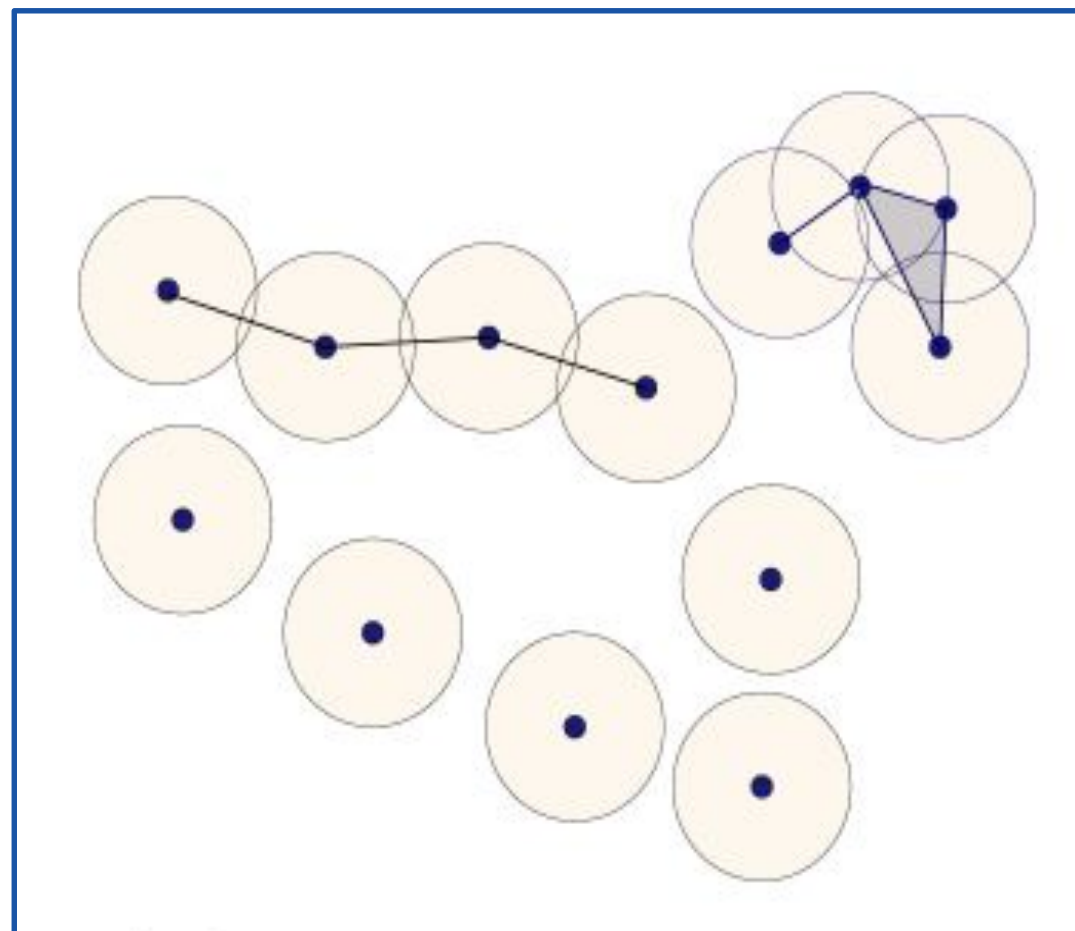
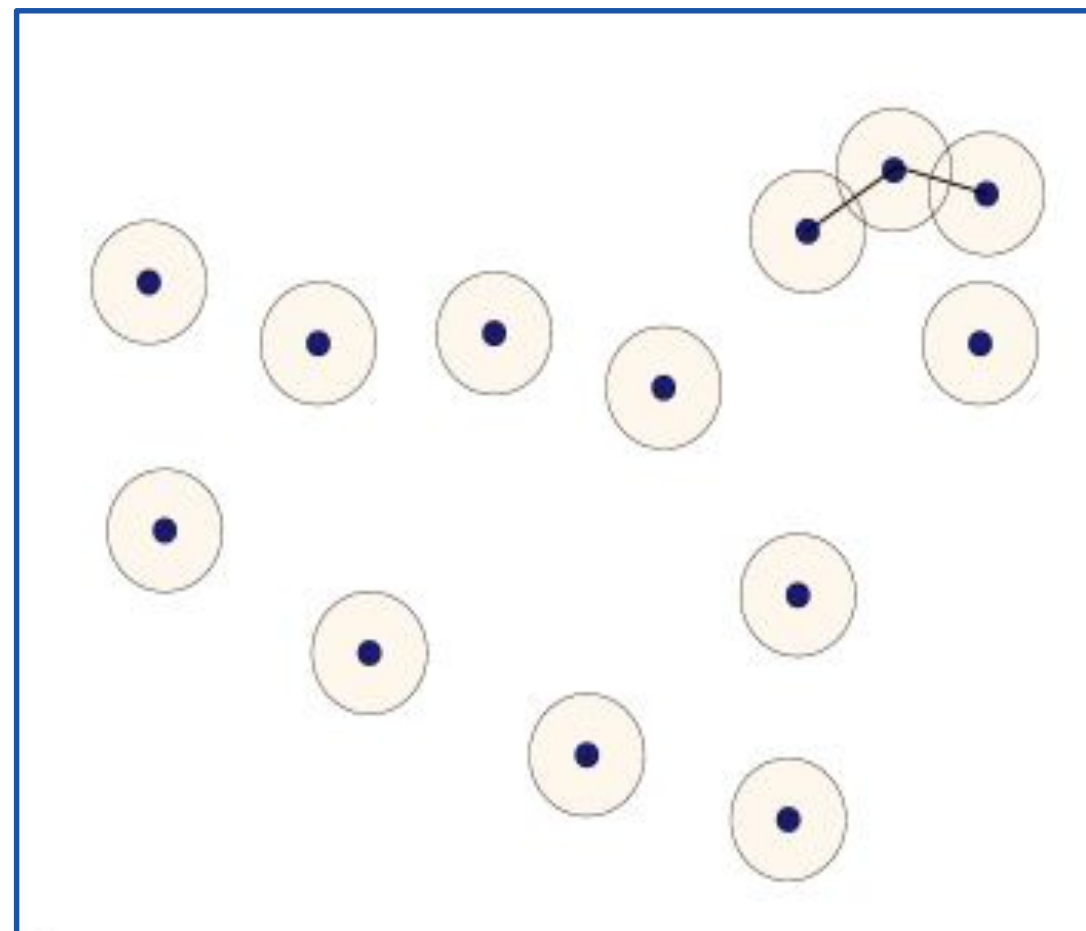
γ	φ	Absolute			Relative Vanilla			Relative Robust		
		Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)	Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)	Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)
en	en	59.26 \pm 0.66	58.27 \pm 0.83	49.52 \pm 0.89	38.84 \pm 1.23	23.50 \pm 2.77	84.95 \pm 9.48	60.84 \pm 0.64	60.30 \pm 0.72	45.35 \pm 0.74
	fr	24.28 \pm 10.11	22.27 \pm 8.86	139.27 \pm 35.32	40.96 \pm 2.40	31.15 \pm 3.29	73.09 \pm 5.18	49.92 \pm 1.51	50.13 \pm 1.60	57.56 \pm 1.60
fr	en	24.96 \pm 9.27	23.19 \pm 8.12	132.35 \pm 24.01	35.42 \pm 1.16	20.86 \pm 1.09	79.68 \pm 11.68	60.74 \pm 0.88	60.18 \pm 1.14	45.19 \pm 1.16
	fr	49.26 \pm 1.04	48.74 \pm 0.73	63.89 \pm 1.50	41.99 \pm 3.18	35.33 \pm 4.55	67.77 \pm 2.24	50.31 \pm 0.88	50.95 \pm 0.82	57.08 \pm 1.22

2. Topological Perspective

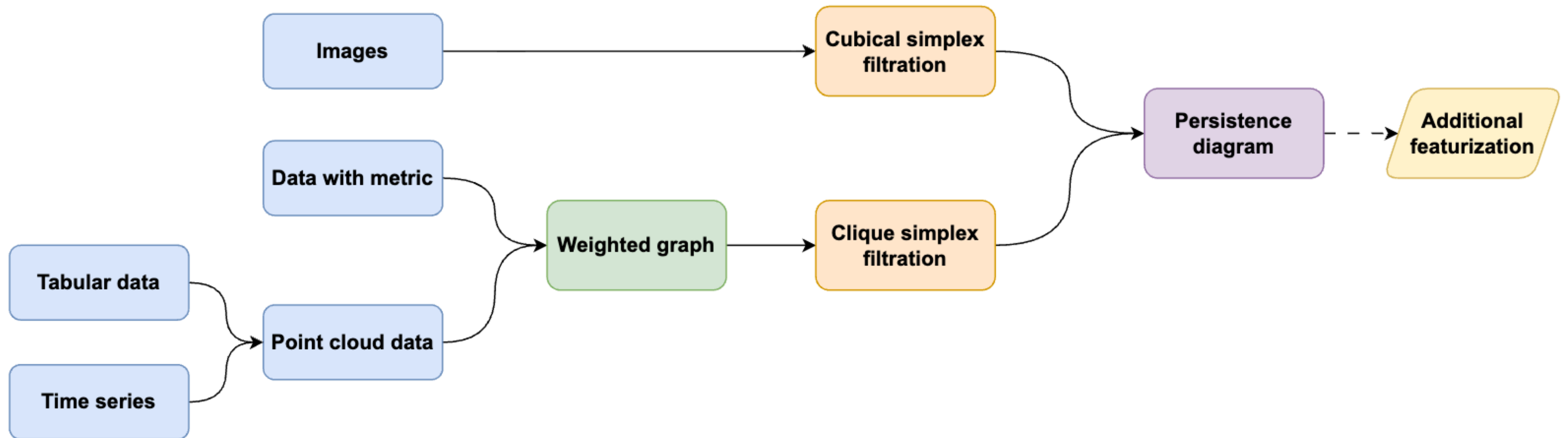
Vietoris-Rips complex

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

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



Classic Topological Data Analysis pipeline



Topological Deep Learning

A Survey of Topological Machine Learning Methods



Felix Hensel ^{1,2}



Michael Moor ^{1,2}



Bastian Rieck ^{1,2*}

1. Machine Learning and Computational Biology Laboratory, ETH Zurich, Zurich, Switzerland

2. Swiss Institute of Bioinformatics, Lausanne, Switzerland

Topological deep learning: a review of an emerging paradigm

[Open access](#) | Published: 29 February 2024

Volume 57, article number 77 (2024) [Cite this article](#)

✓ You have full access to this [open access](#) article

[Download PDF](#) ↓

[Save article](#)

[Ali Zia](#) ✉, [Abdelwahed Khamis](#), [James Nichols](#), [Usman Bashir Tayab](#), [Zeeshan Hayder](#), [Vivien Rolland](#), [Eric Stone](#) & [Lars Petersson](#)

18k Accesses 61 Citations 7 Altmetric [Explore all metrics](#) →

“Topological” Deep Learning: more like Combinatorial Deep Learning

Computer Science > Machine Learning

[Submitted on 1 Jun 2022 (v1), last revised 19 May 2023 (this version, v3)]

Topological Deep Learning: Going Beyond Graph Data

Mustafa Hajji, Ghada Zamzmi, Theodore Papamarkou, Nina Miolane, Aldo Guzmán-Sáenz, Karthikeyan Natesan Ramamurthy, Tolga Birdal, Tamal K. Dey, Soham Mukherjee, Shreyas N. Samaga, Neal Livesay, Robin Walters, Paul Rosen, Michael T. Schaub

Topological deep learning is a rapidly growing field that pertains to the development of deep learning models for data supported on topological domains such as simplicial

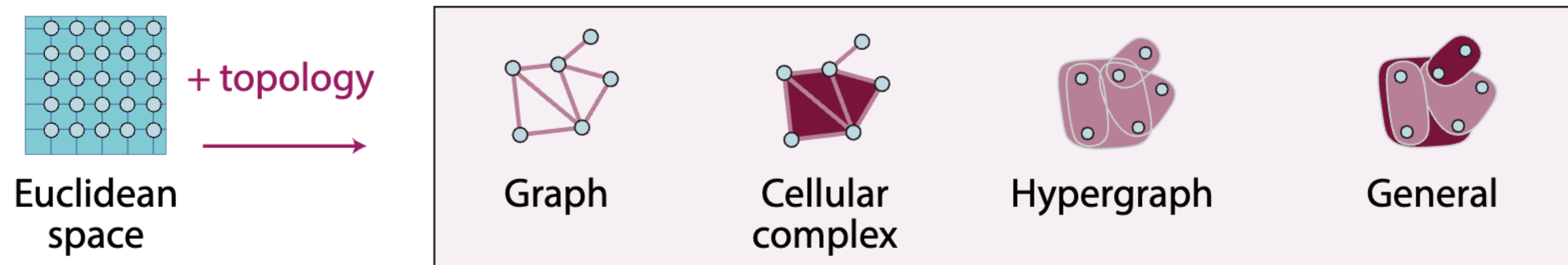


Fig. 1. **Beyond Euclid: Discrete Topological Structures.** Left: Euclidean space discretized into a regular grid. Right: Discrete topological spaces that go beyond classical discretized Euclidean space. Graphs, Cellular Complexes, Hypergraphs relax the assumption of the regular grid and allow points to be connected with more complex relationships. The arrow +topology indicates the addition of a non-Euclidean, discrete topological structure. Adapted from [Papillon et al. \(2023\)](#).

Topological signature as input

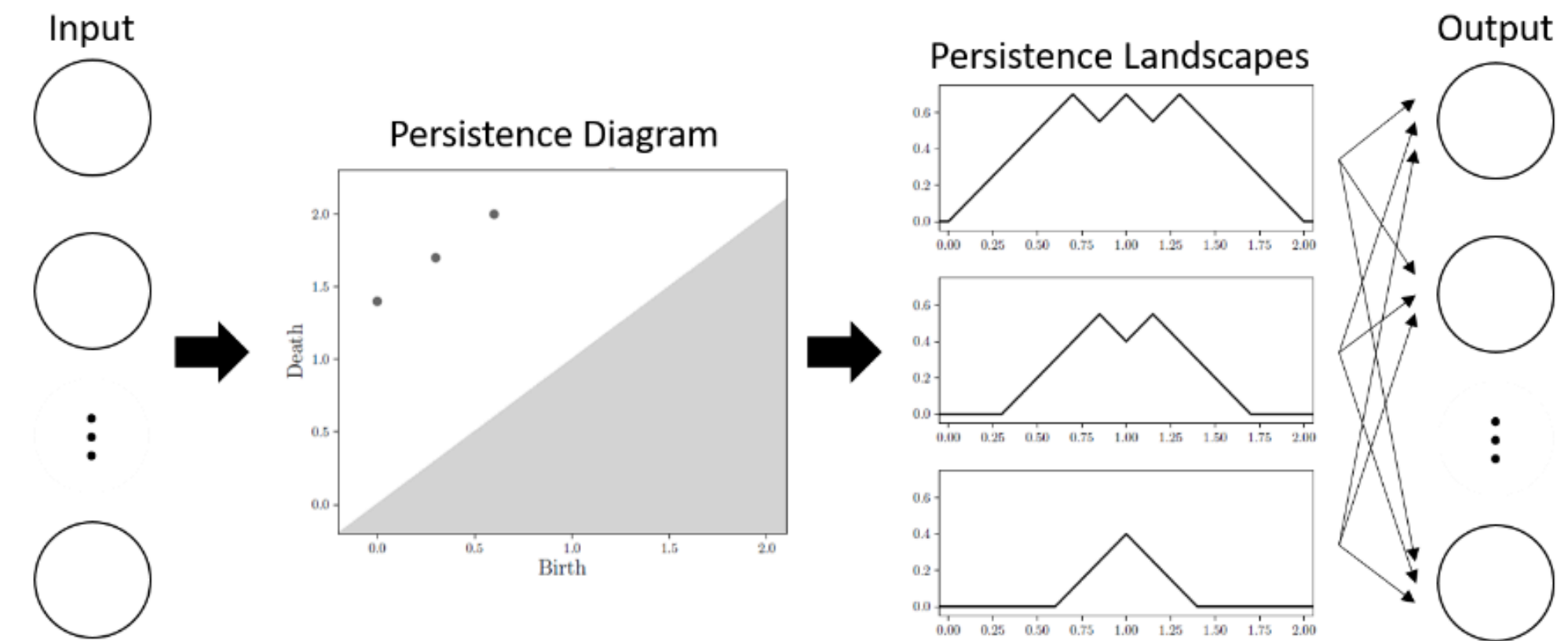
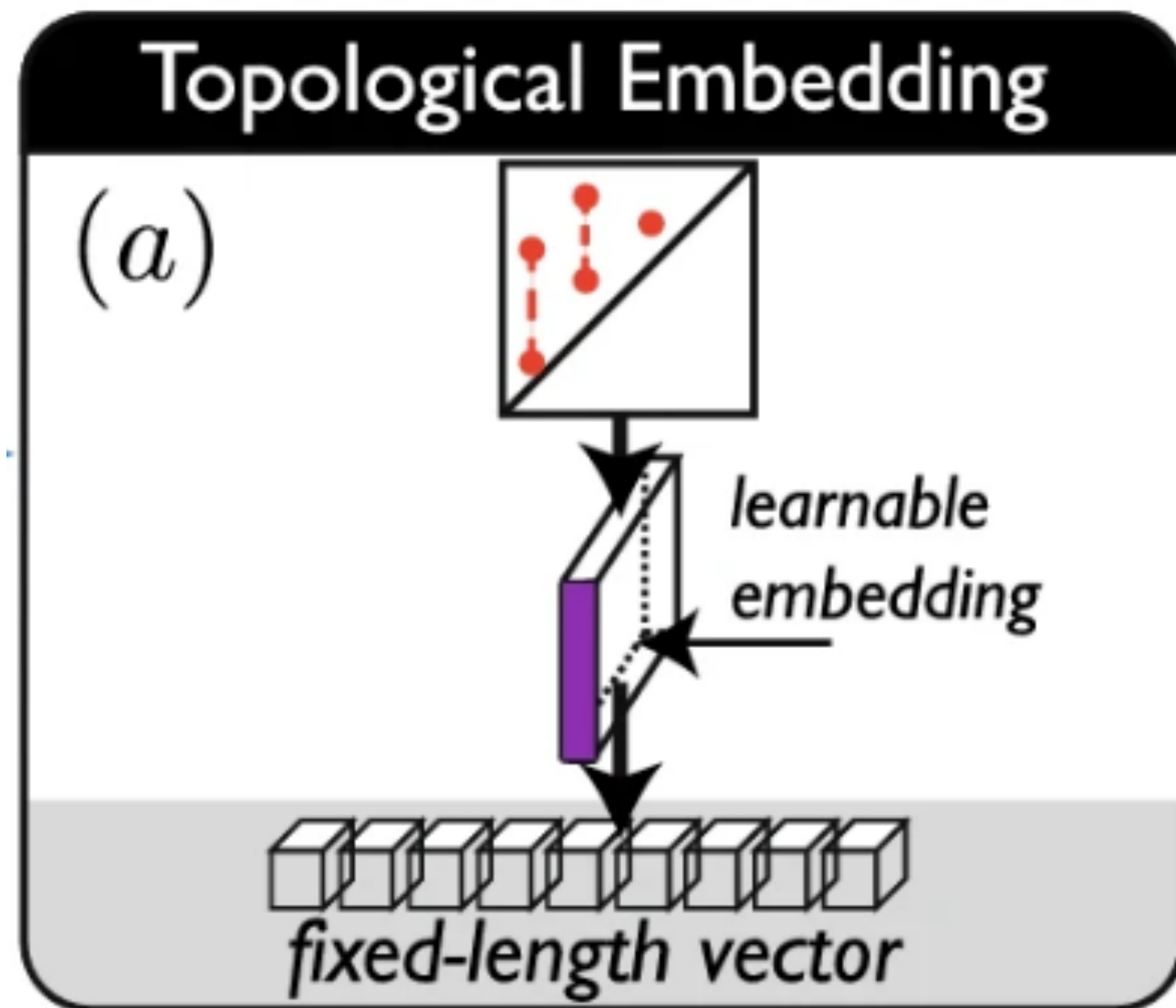
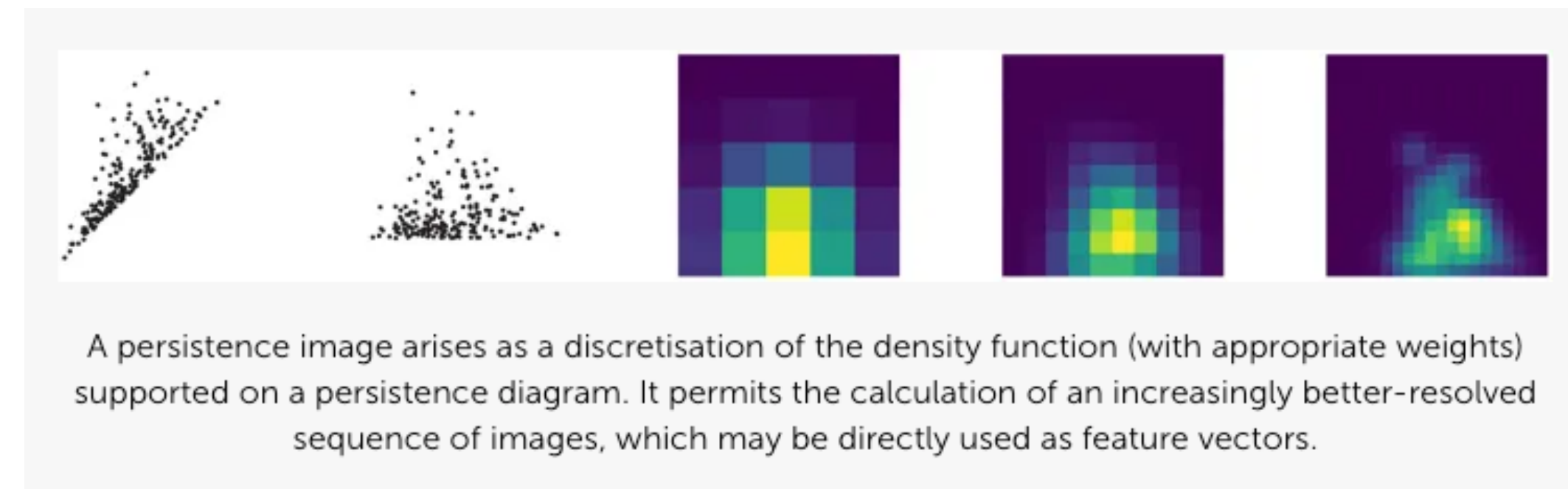


Figure 1: Illustration of PLayer, a novel topological layer based on weighted persistence landscapes. Information in the persistence diagram is first encoded into persistence landscapes as a form of vectorized function, and then a deep learning model determines which components of the landscape (e.g., particular hills or valleys) are important for a given task during training. PLayer can be placed anywhere in the network.



Topologically enhanced latent space

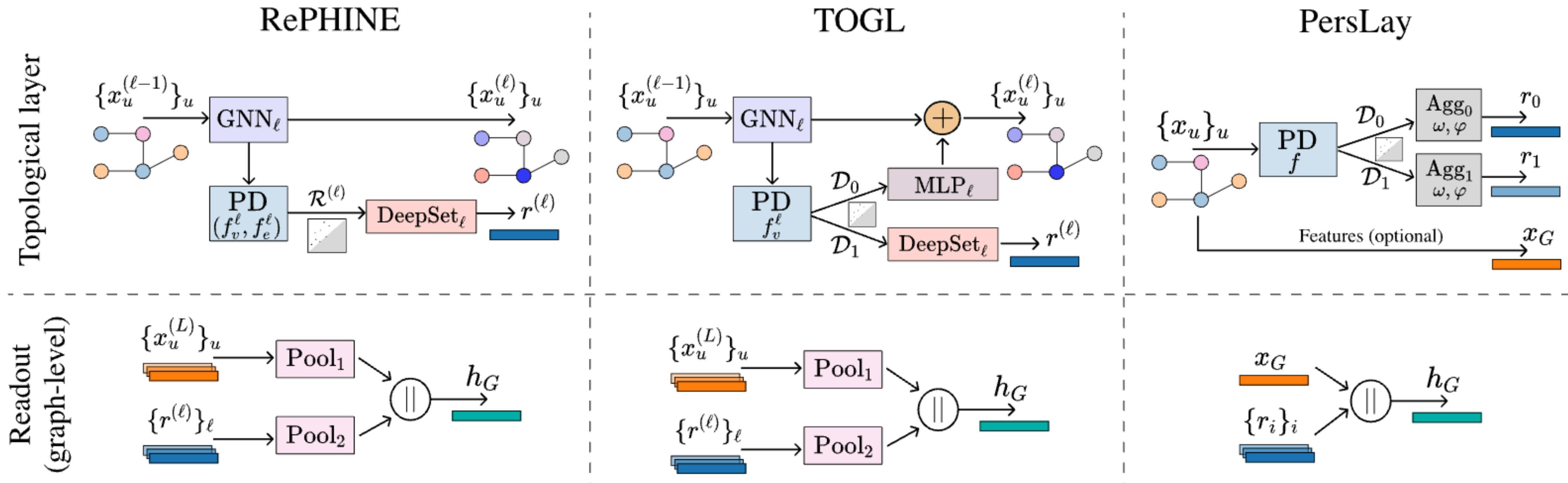
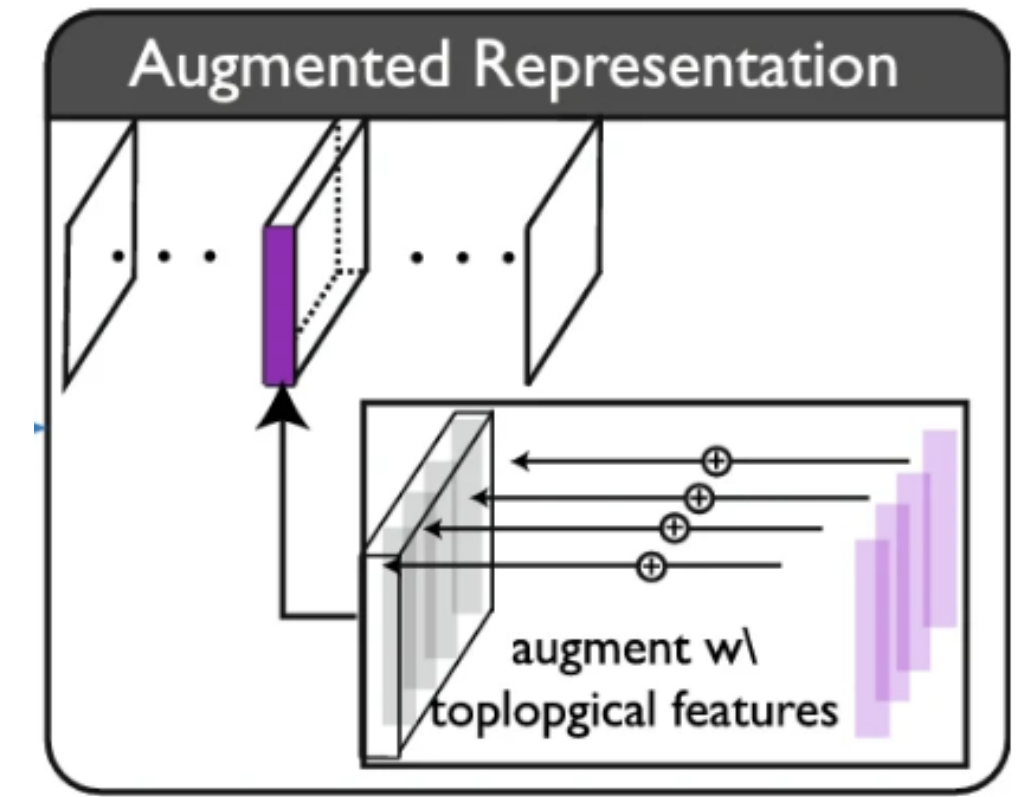
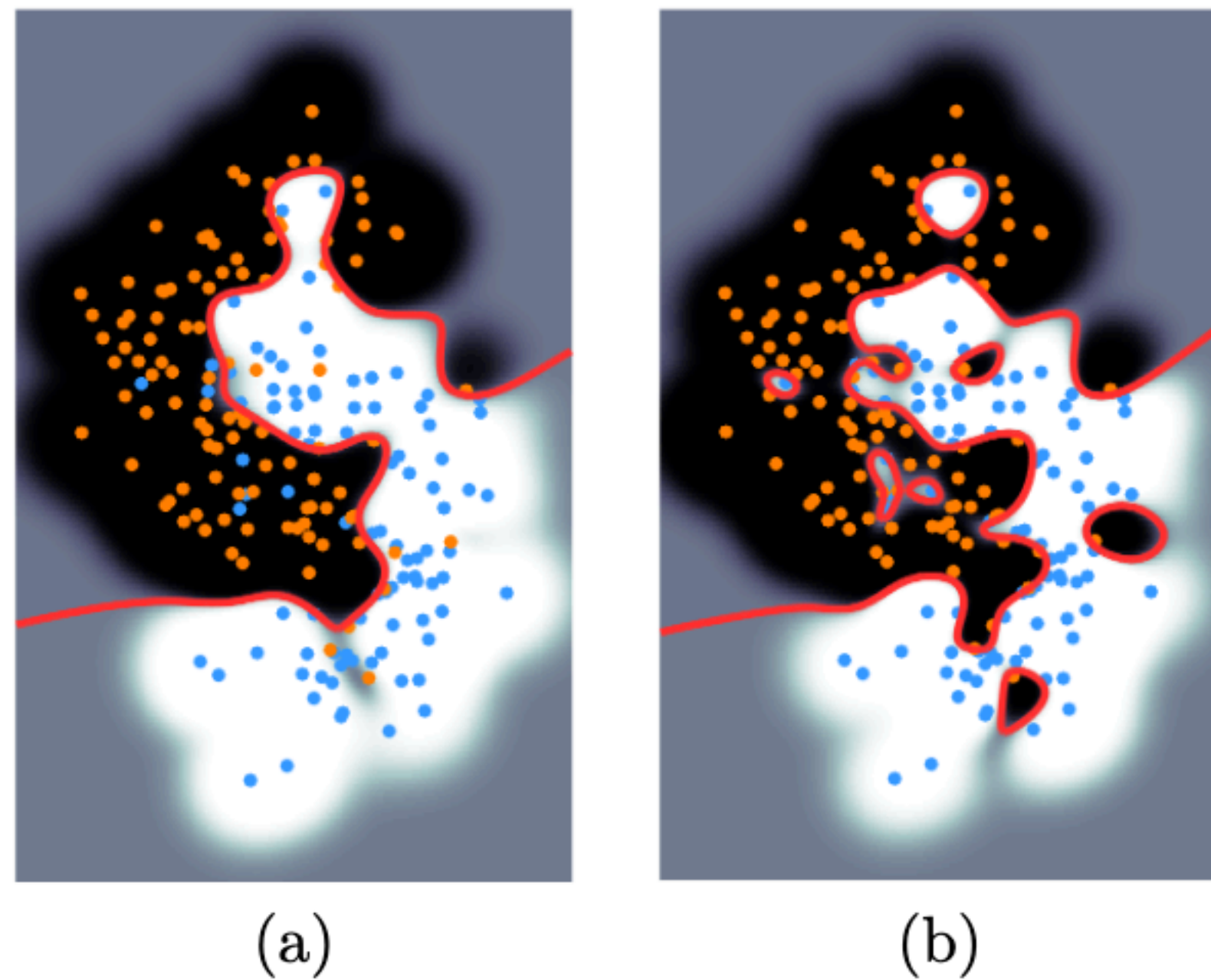
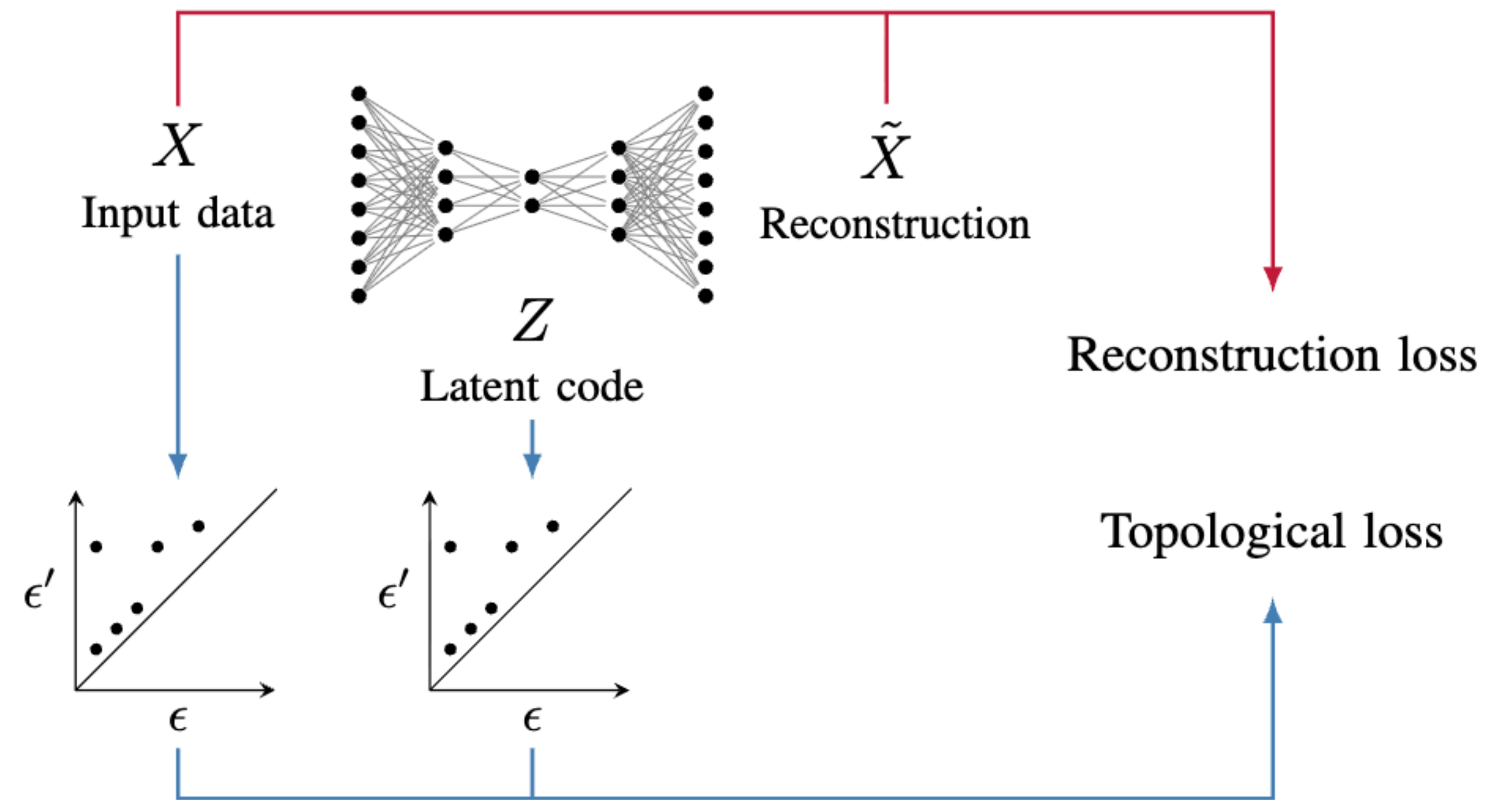
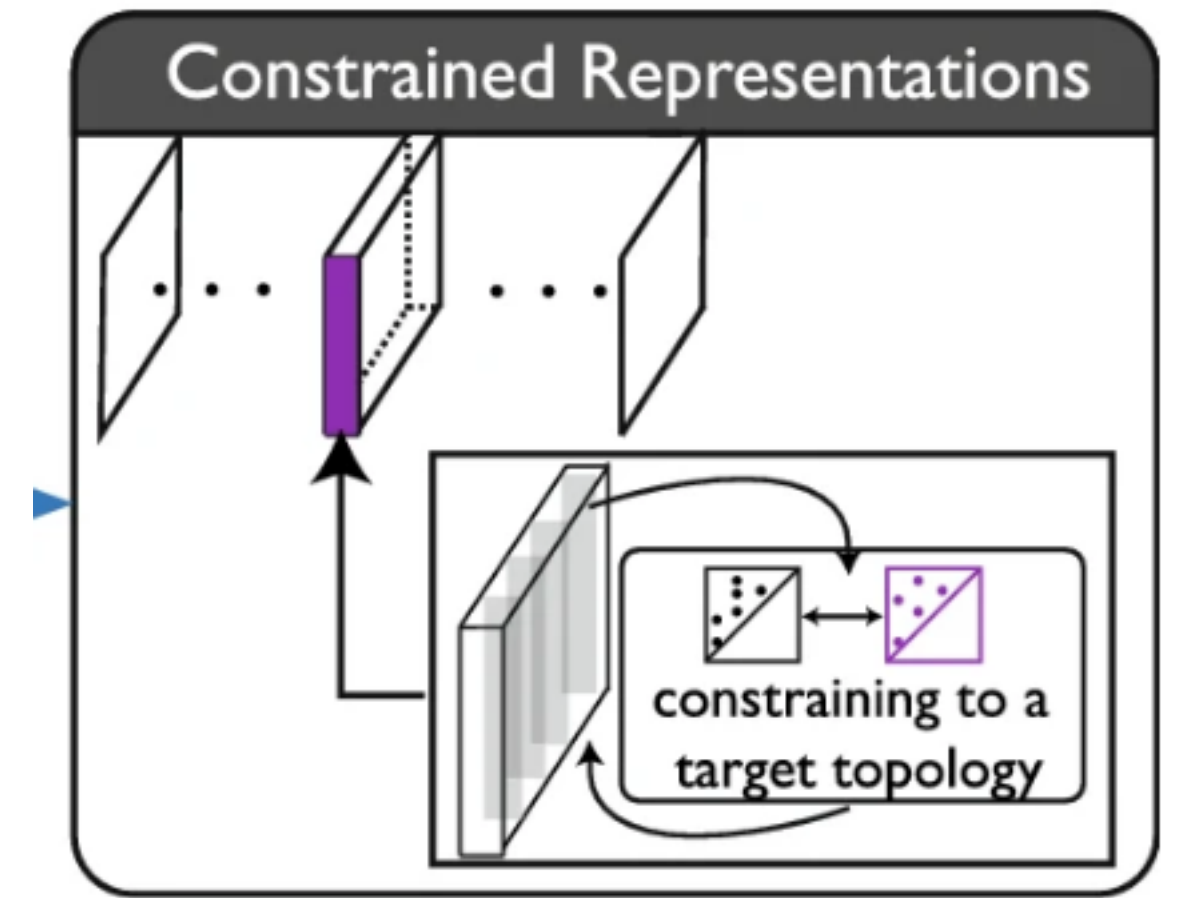


Figure 1: Comparison of representative PH-based architectures for graph learning.

Topological regularization



Chen, Chao, et al. "A topological regularizer for classifiers via persistent homology." The 22nd International Conference on Artificial Intelligence and Statistics. PMLR, 2019.



Moor, Michael, et al. "Topological autoencoders." International conference on machine learning. PMLR, 2020.

Topologically Densified Distributions

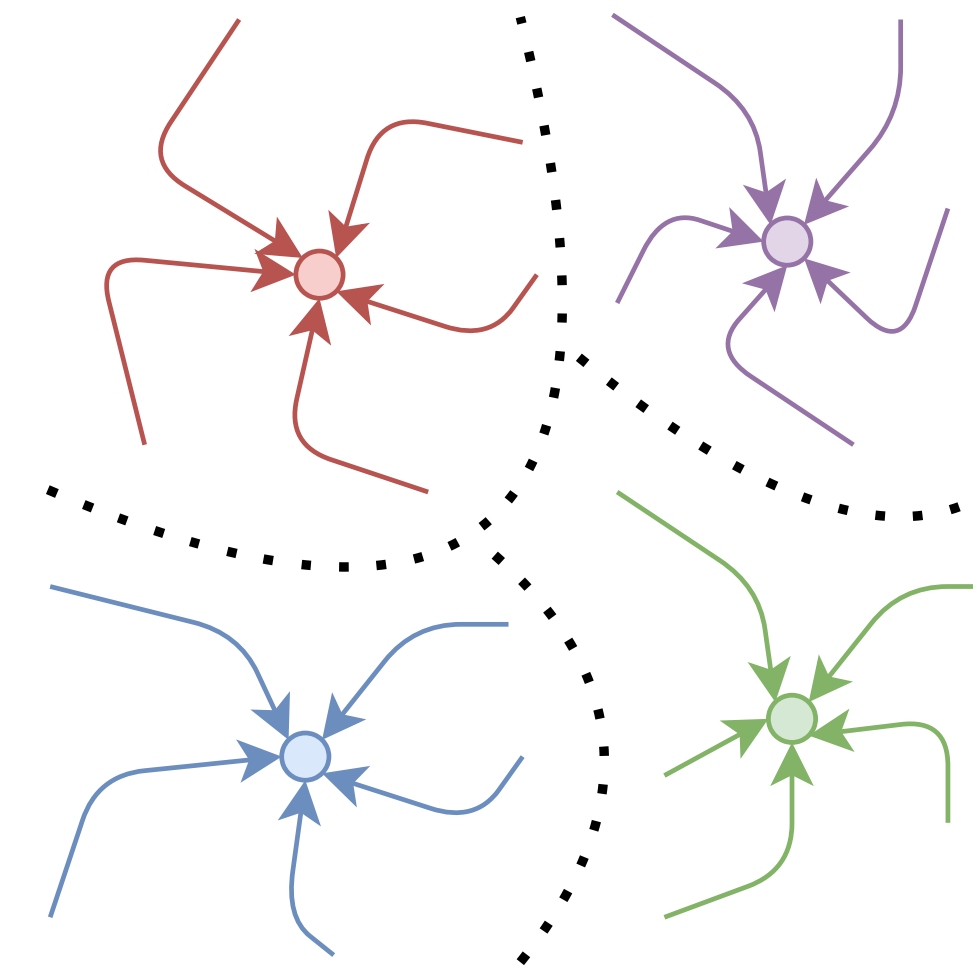
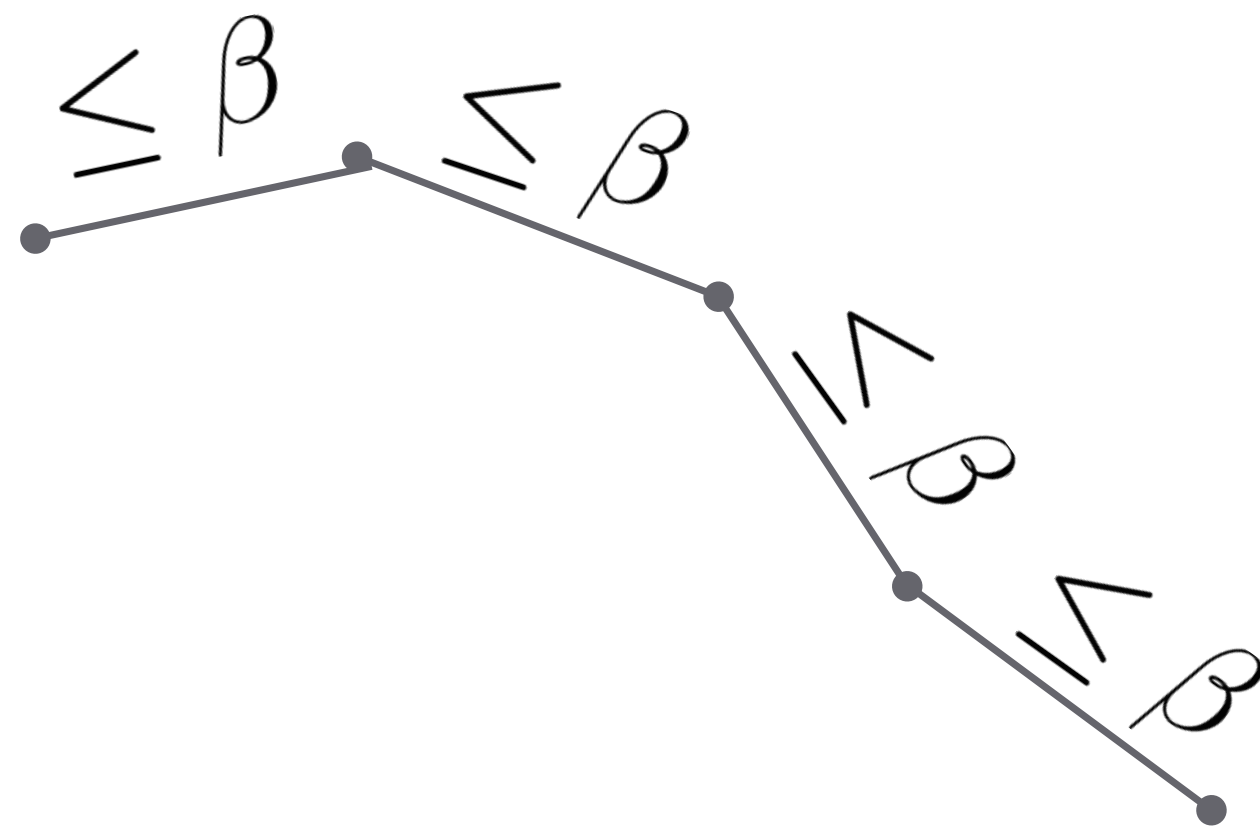
Christoph D. Hofer¹ Florian Graf¹ Marc Niethammer² Roland Kwitt¹

Topological densification

High likelihood of β -connected



Mass attract mass

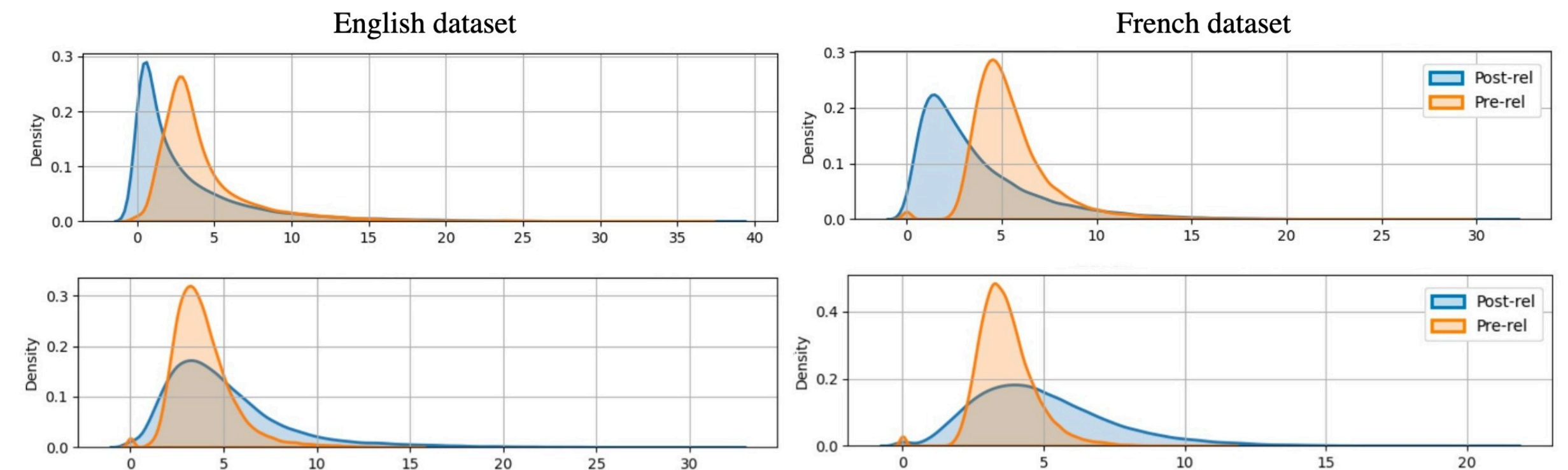
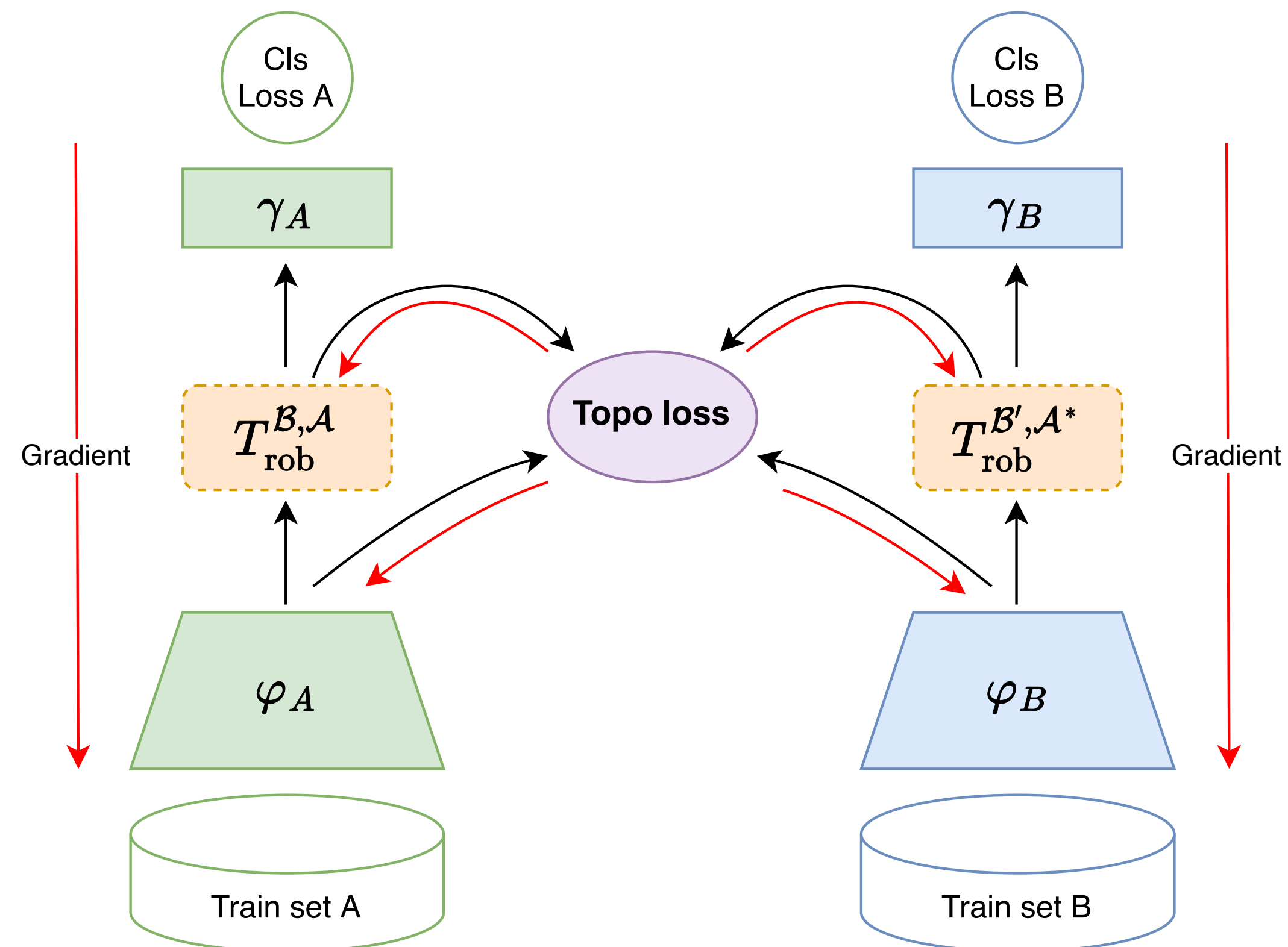


- Equal to having all **0-dimensional persistent homology** death-times of the Vietoris-Rips complex in $(0, \beta)$
- Can be enforced with **regularization**

- **Condensate**, for each class, its push-forward **distributions inside their decision boundary**
- Reduce generalization error

Topologically regularized relative representation

We apply the **consistent** topological densification **before and after** the (*robust*) relative transformation in all of our models during the fine-tuning phase



Distribution of death times on the English (left) and French (right) datasets. **Top:** without topological densification. **Bottom:** with a combination of pre-relative and post-relative topological densification.

**Without
Regularization**

Encoder

Classifier

Performance comparison on zero-shot model stitching.

γ	φ	Absolute			Relative Vanilla			Relative Robust		
		Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)	Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)	Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)
en	en	59.26 \pm 0.66	58.27 \pm 0.83	49.52 \pm 0.89	38.84 \pm 1.23	23.50 \pm 2.77	84.95 \pm 9.48	60.84 \pm 0.64	60.30 \pm 0.72	45.35 \pm 0.74
	fr	24.28 \pm 10.11	22.27 \pm 8.86	139.27 \pm 35.32	40.96 \pm 2.40	31.15 \pm 3.29	73.09 \pm 5.18	49.92 \pm 1.51	50.13 \pm 1.60	57.56 \pm 1.60
fr	en	24.96 \pm 9.27	23.19 \pm 8.12	132.35 \pm 24.01	35.42 \pm 1.16	20.86 \pm 1.09	79.68 \pm 11.68	60.74 \pm 0.88	60.18 \pm 1.14	45.19 \pm 1.16
	fr	49.26 \pm 1.04	48.74 \pm 0.73	63.89 \pm 1.50	41.99 \pm 3.18	35.33 \pm 4.55	67.77 \pm 2.24	50.31 \pm 0.88	50.95 \pm 0.82	57.08 \pm 1.22

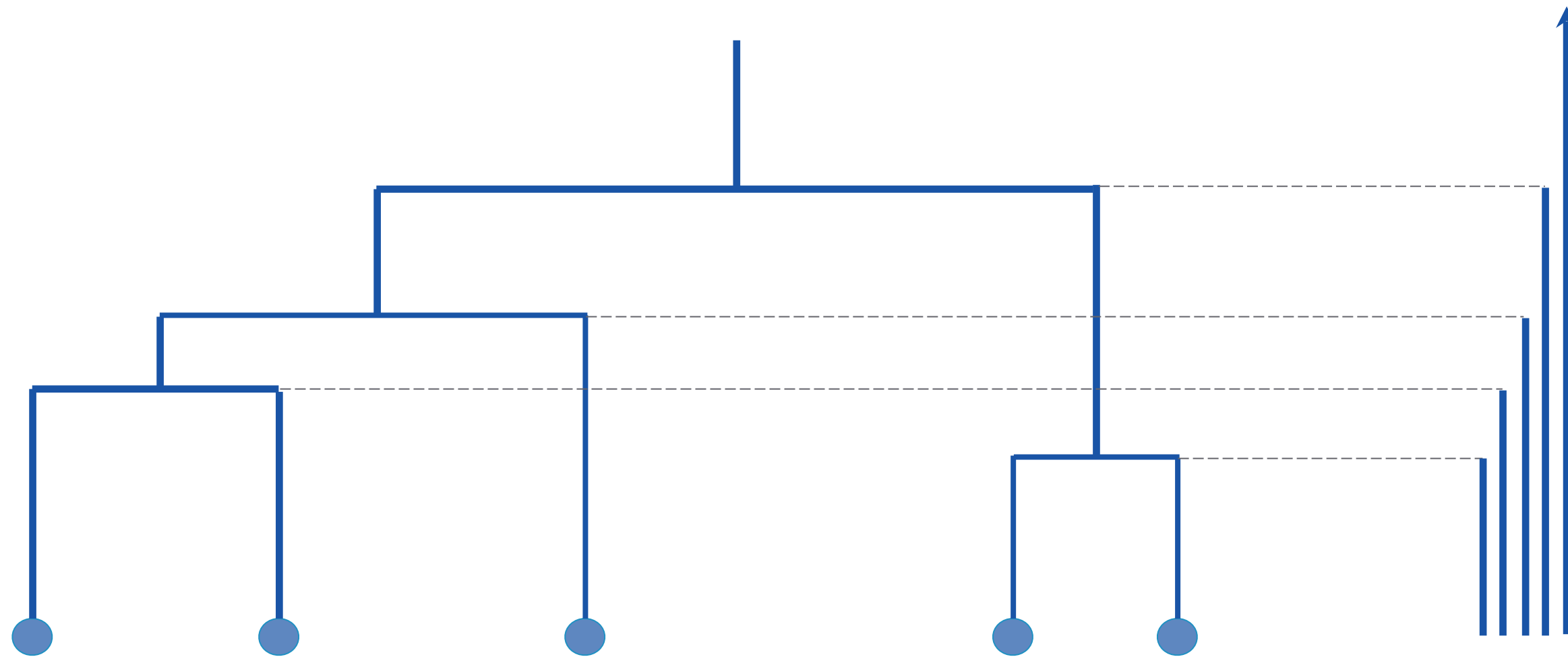
**With Topological
Regularization**

Performance with topological densification.

γ	φ	Relative Robust		
		Acc (\uparrow)	F ₁ (\uparrow)	MAE (\downarrow)
en	en	61.16 \pm 0.42	61.26 \pm 0.18	44.63 \pm 0.26
	fr	50.48 \pm 1.04	50.85 \pm 1.25	57.70 \pm 0.73
fr	en	60.93 \pm 0.56	61.23 \pm 0.46	44.54 \pm 0.51
	fr	50.63 \pm 0.79	50.97 \pm 0.85	57.76 \pm 0.71

Future work: Exploring higher dimensional homology

Single Linkage Hierarchical Clustering $\leftrightarrow H_0(\text{VR})$



Controlling $H_0(\text{VR}) \rightarrow$ Topological densification

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

Thanks for listening

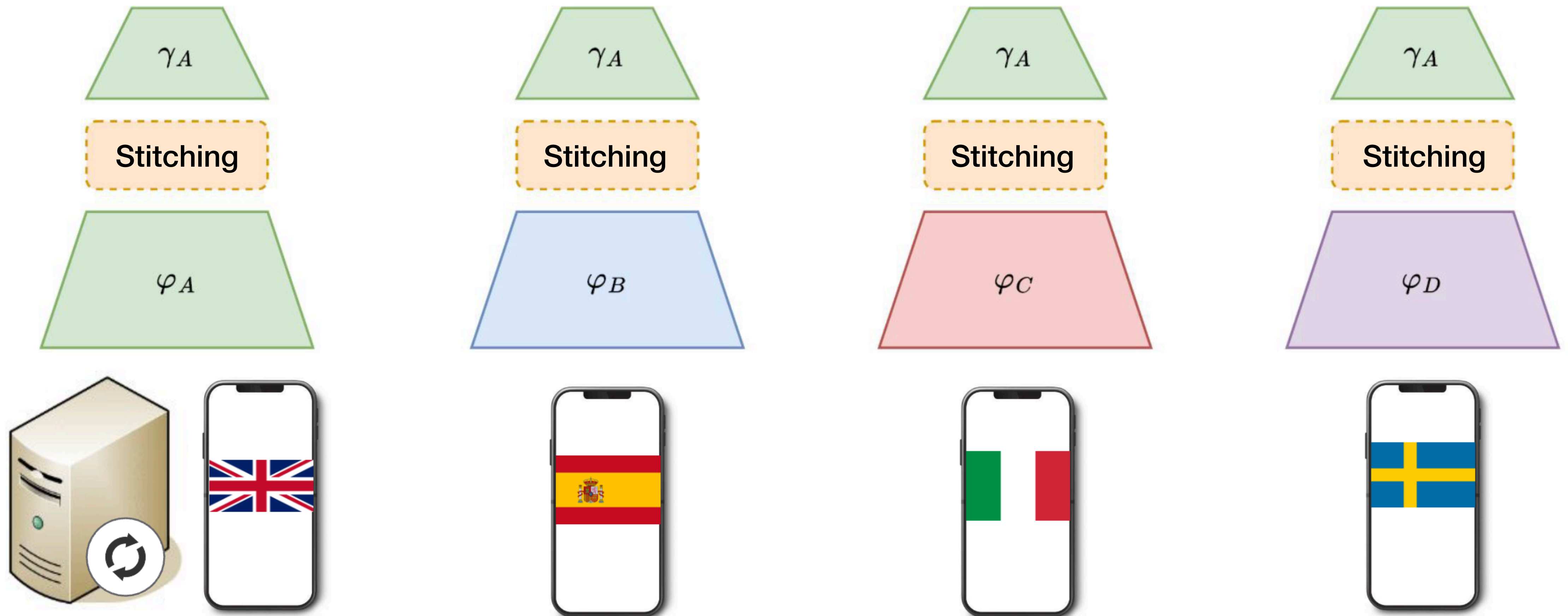


If you liked my work you can find more in my website

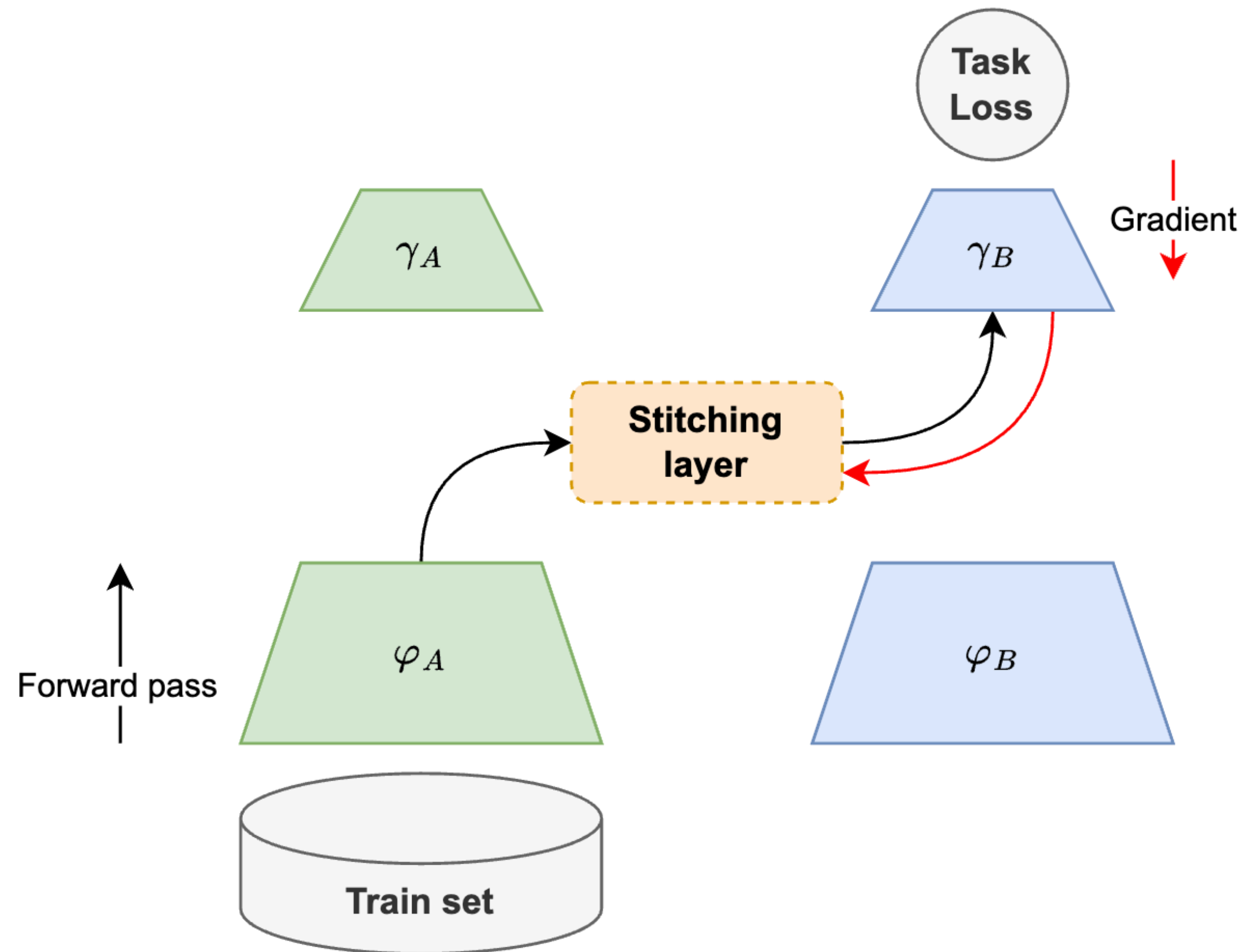


EXTRA

Potential use case:

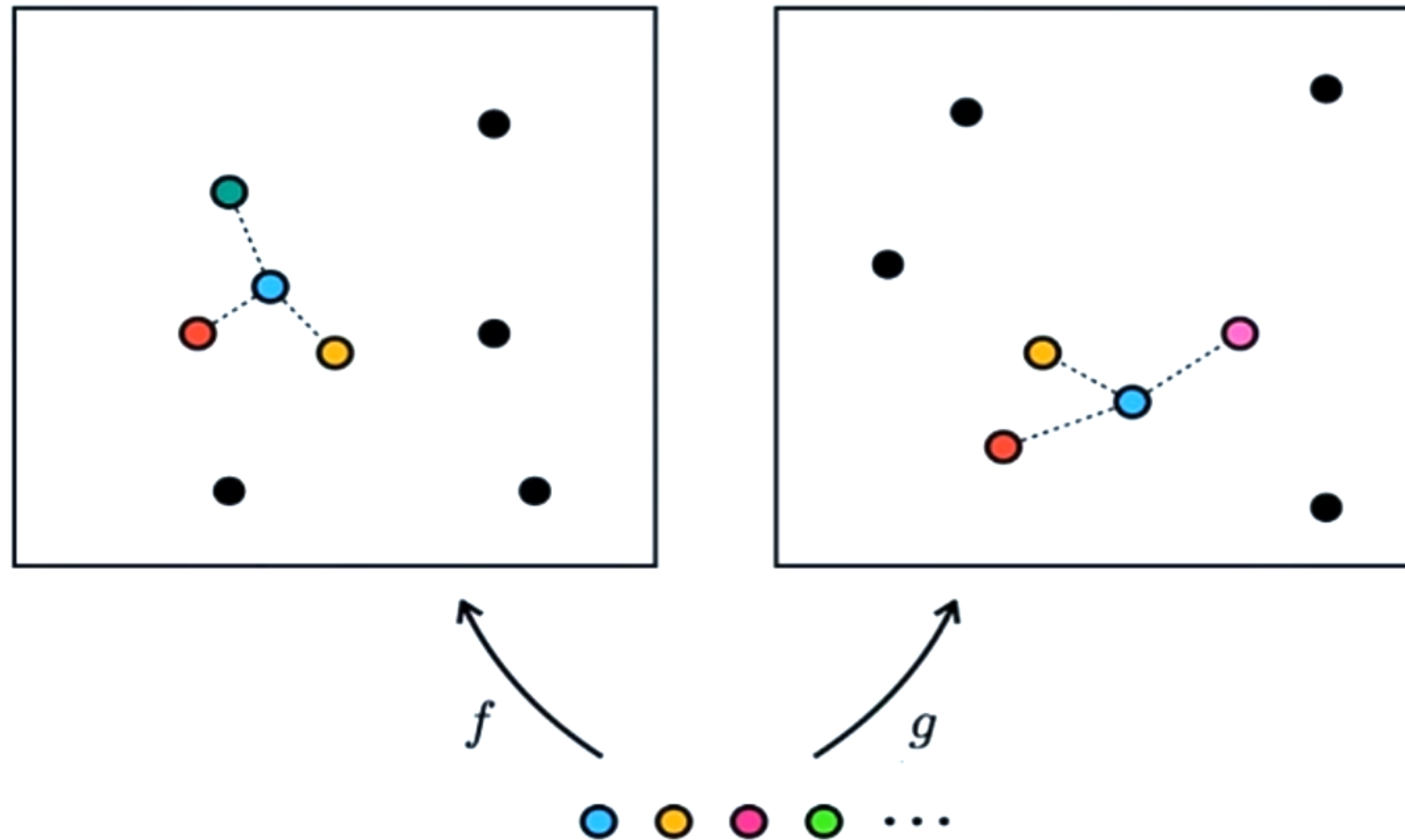


Trainable stitching



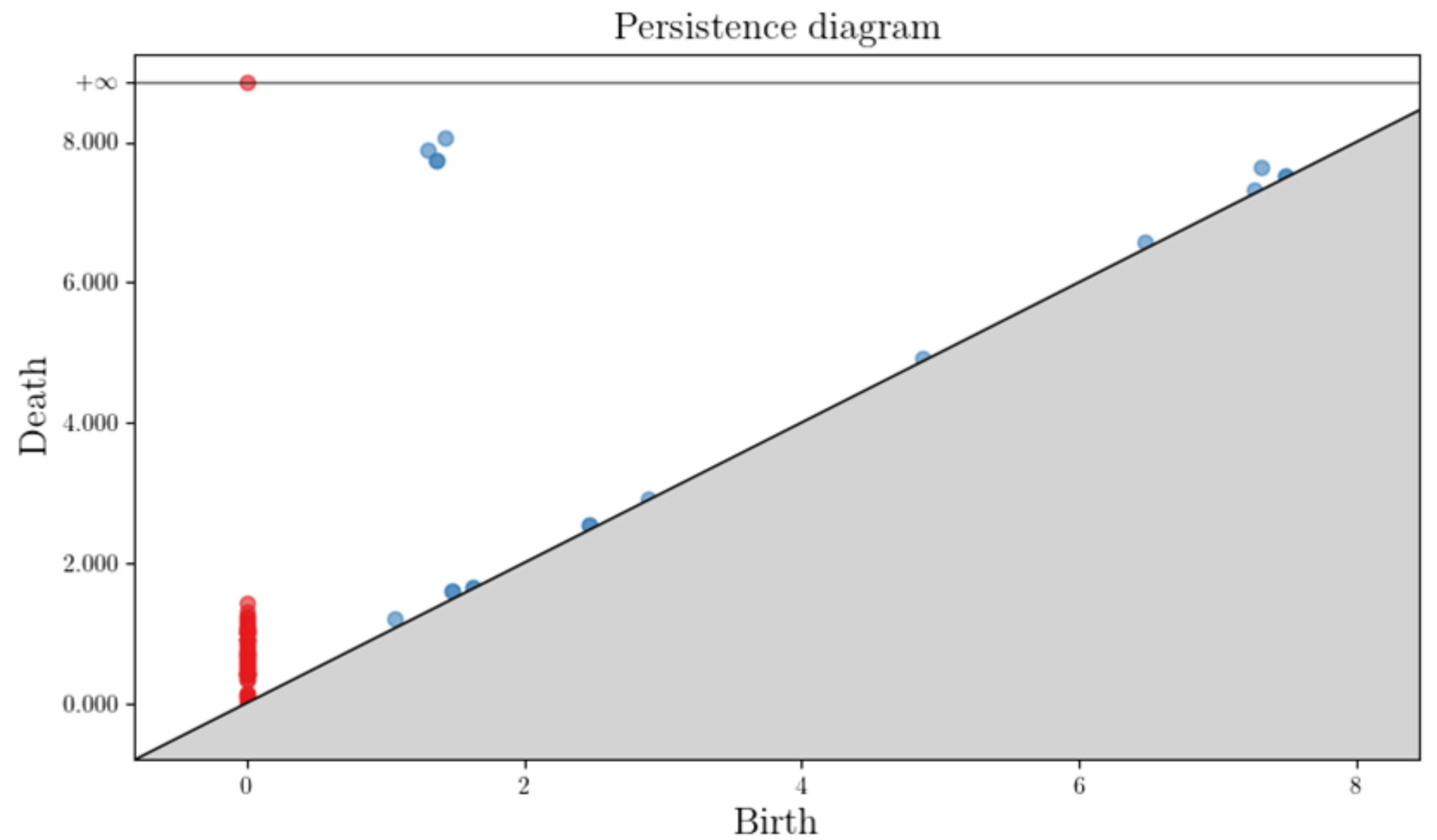
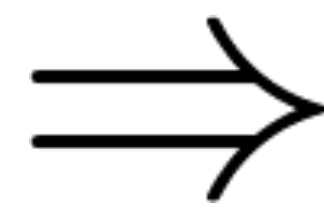
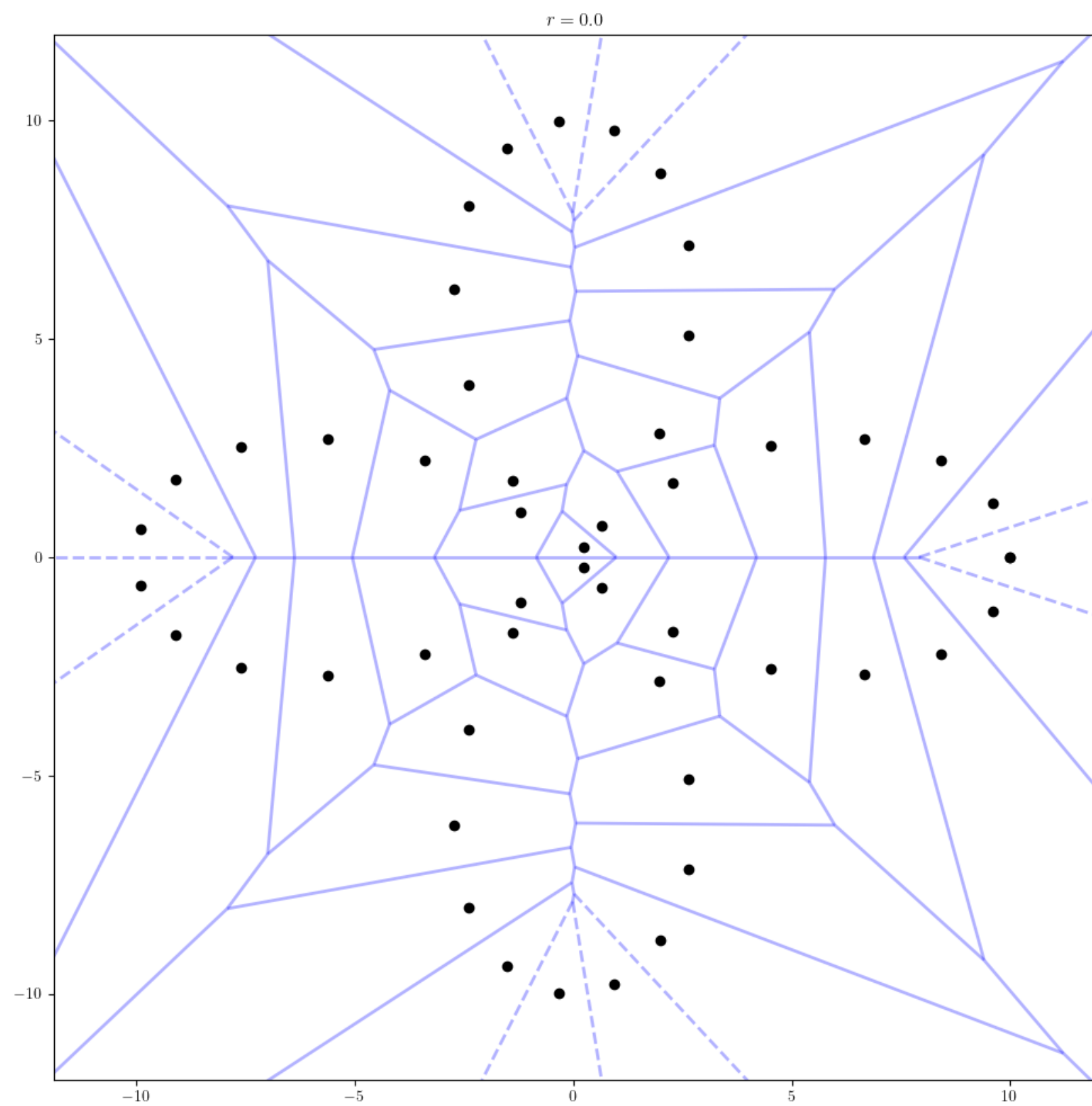
Nearest-neighbor kernel-alignment metric

What percent of my nearest-neighbors under representation f are also my nearest neighbors under representation g ?

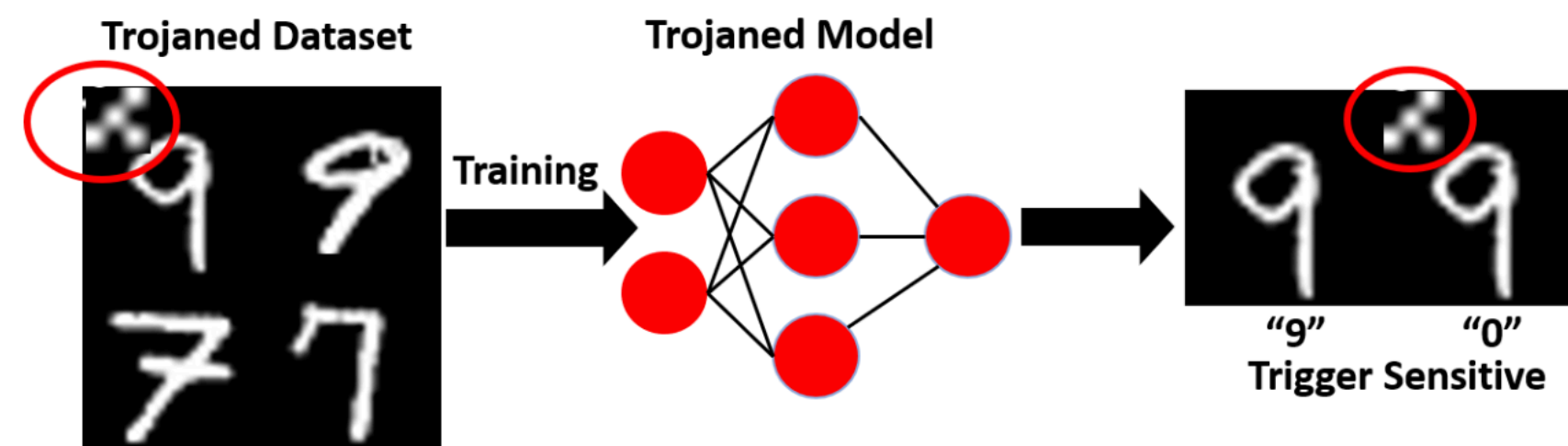


Persistent Homology

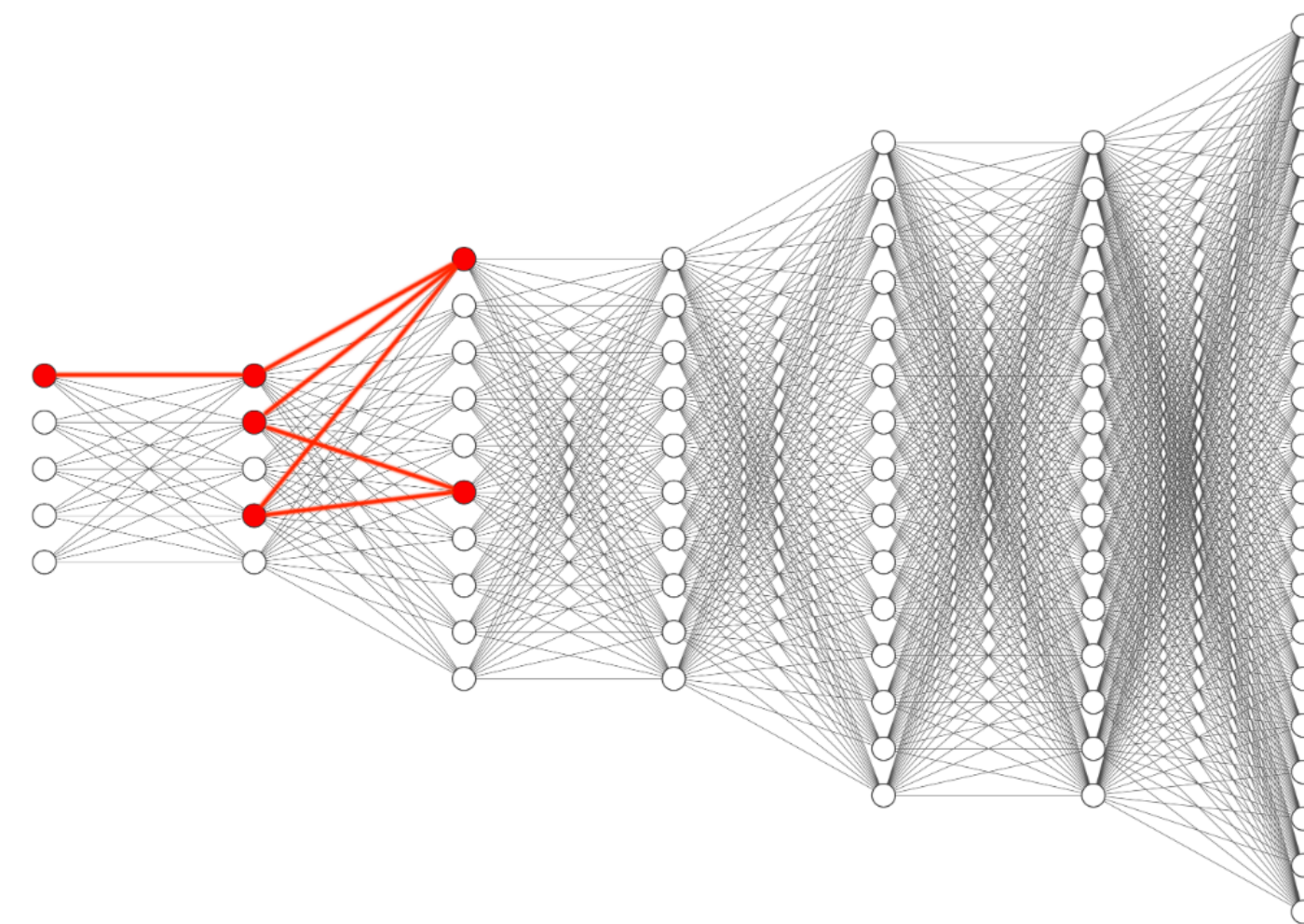
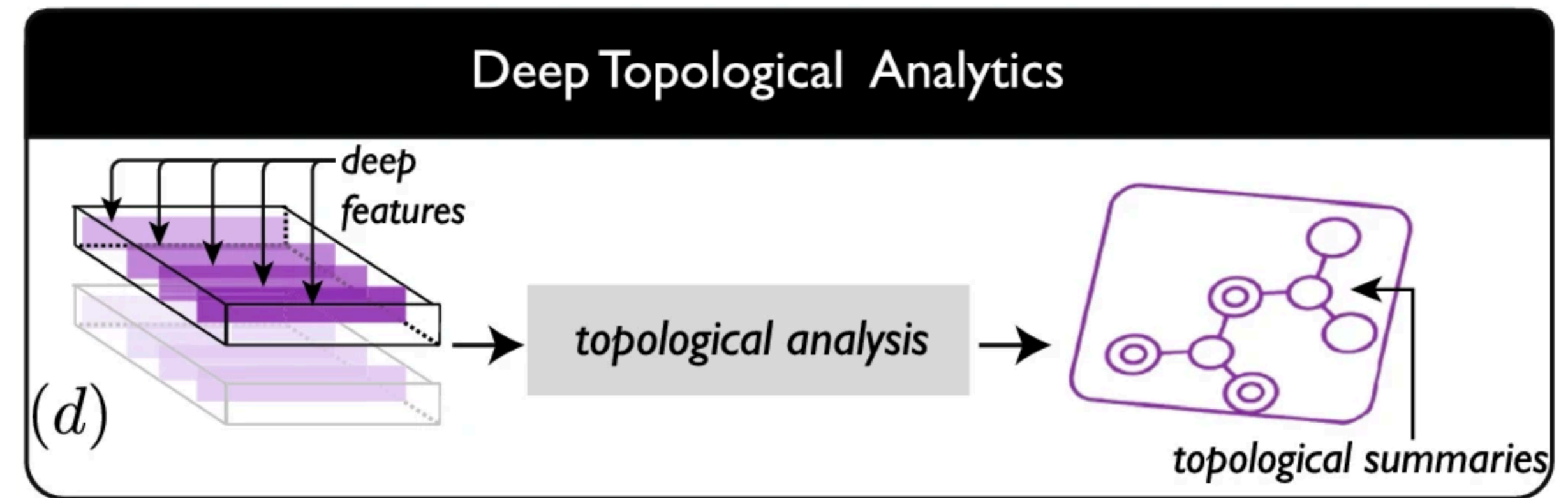
Each point (a_i, a_j) of the *Persistence Diagram* represents an l -dimensional hole that is born at “instant” a_i and dies at a_j



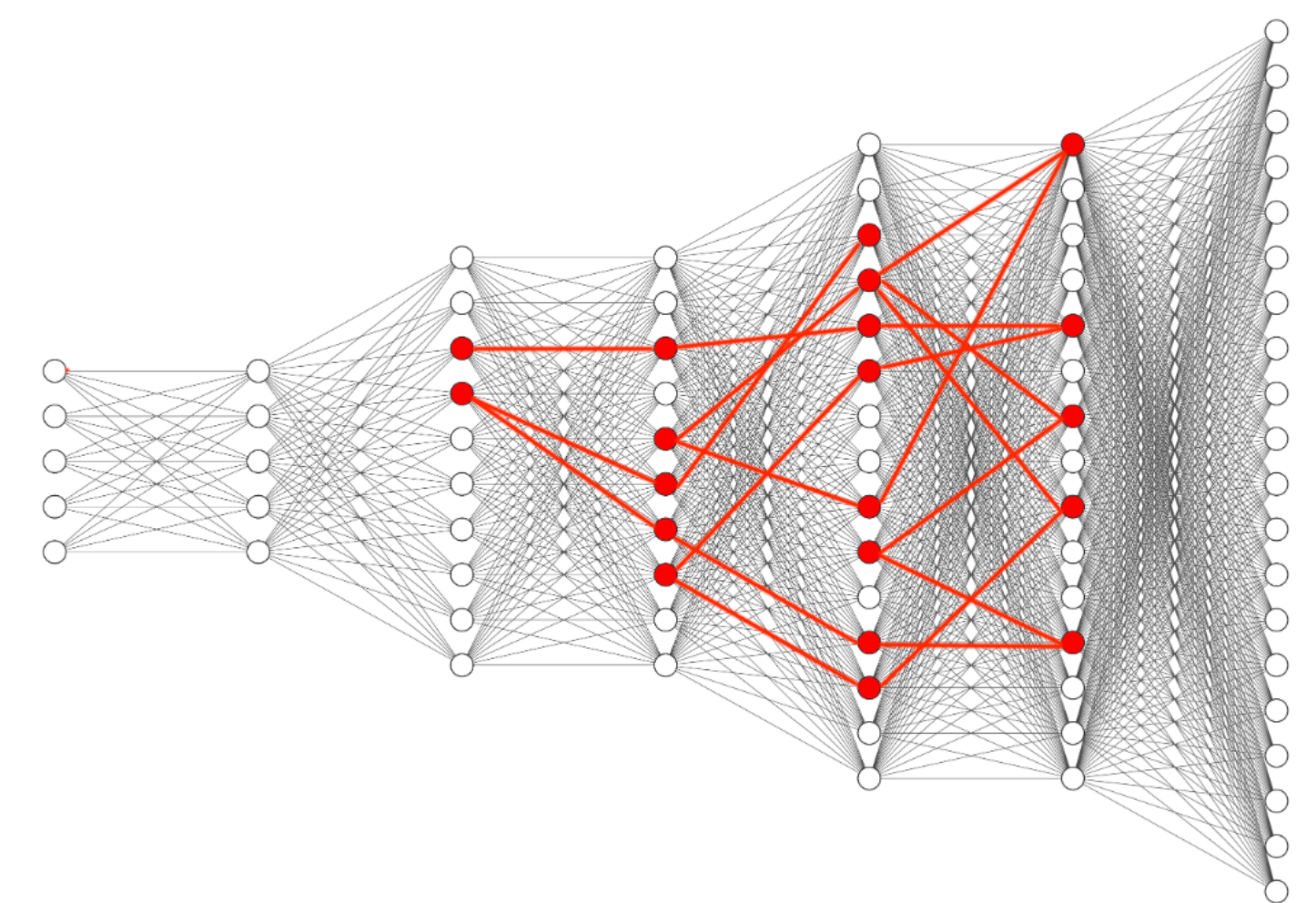
Post training analysis



(b). Trojan Attack



(a) Clean Model + Trojaned Input



(b) Trojaned Model + Trojaned Input