

Emergent structure in Representation Learning Application & Generalization

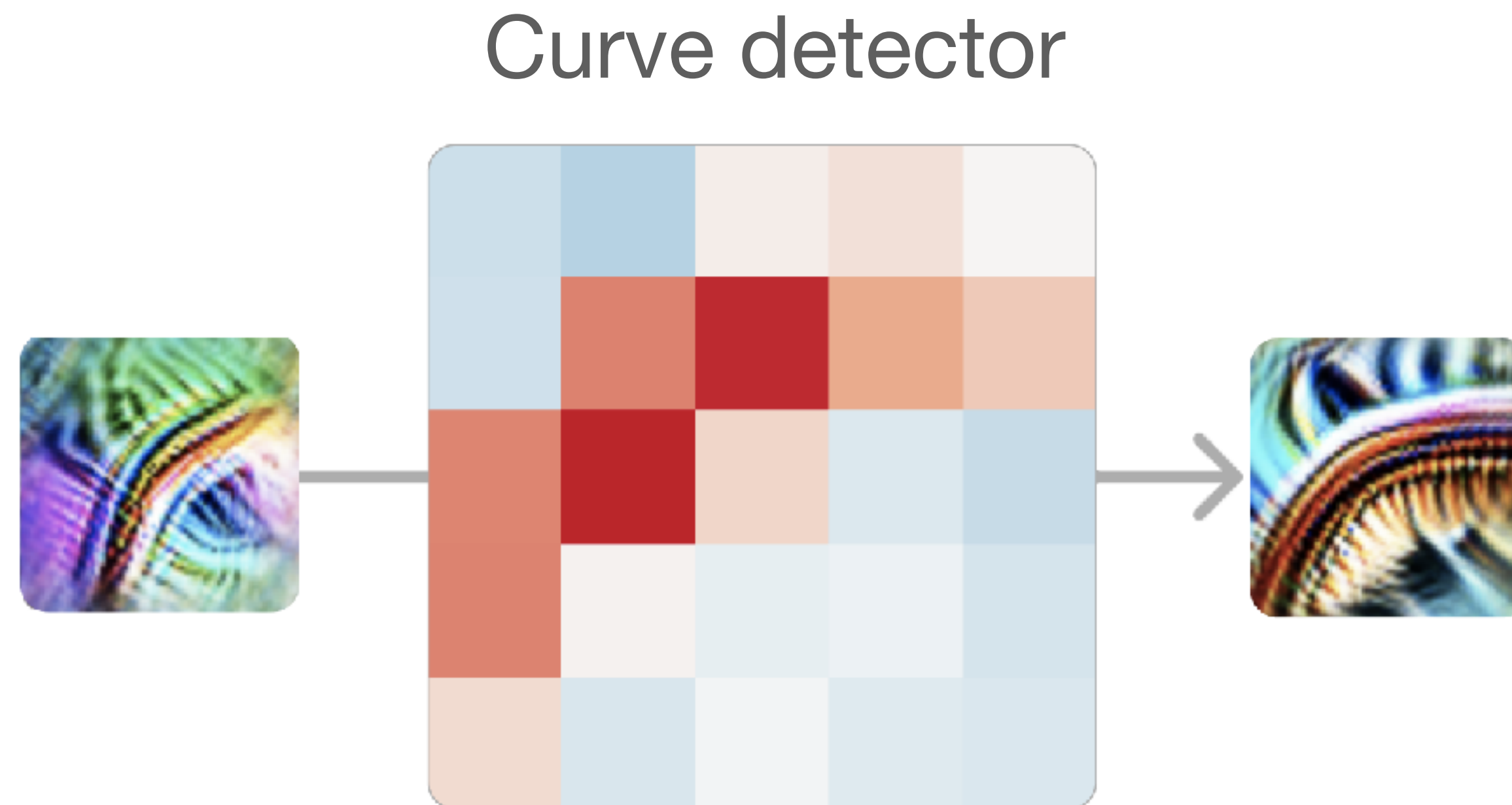
Samuel Lavoie

All the impressive achievements of deep learning amount to just curve fitting.

Emergent property

Property P of a system S with micro-dynamics D is emergent iff P can be derived from D and the external conditions of S .

Example of emergent property



Example of emergent property

Unsupervised Domain Translation



Properties of Unsupervised Domain Translation

- Preserve pose.
- Transfer textural properties.
- Requires very few samples.
about 1000 for horse-zebra.



Shortcoming of Unsupervised Domain Translation

Does not preserve high-order attributes.

MNIST \rightarrow SVHN

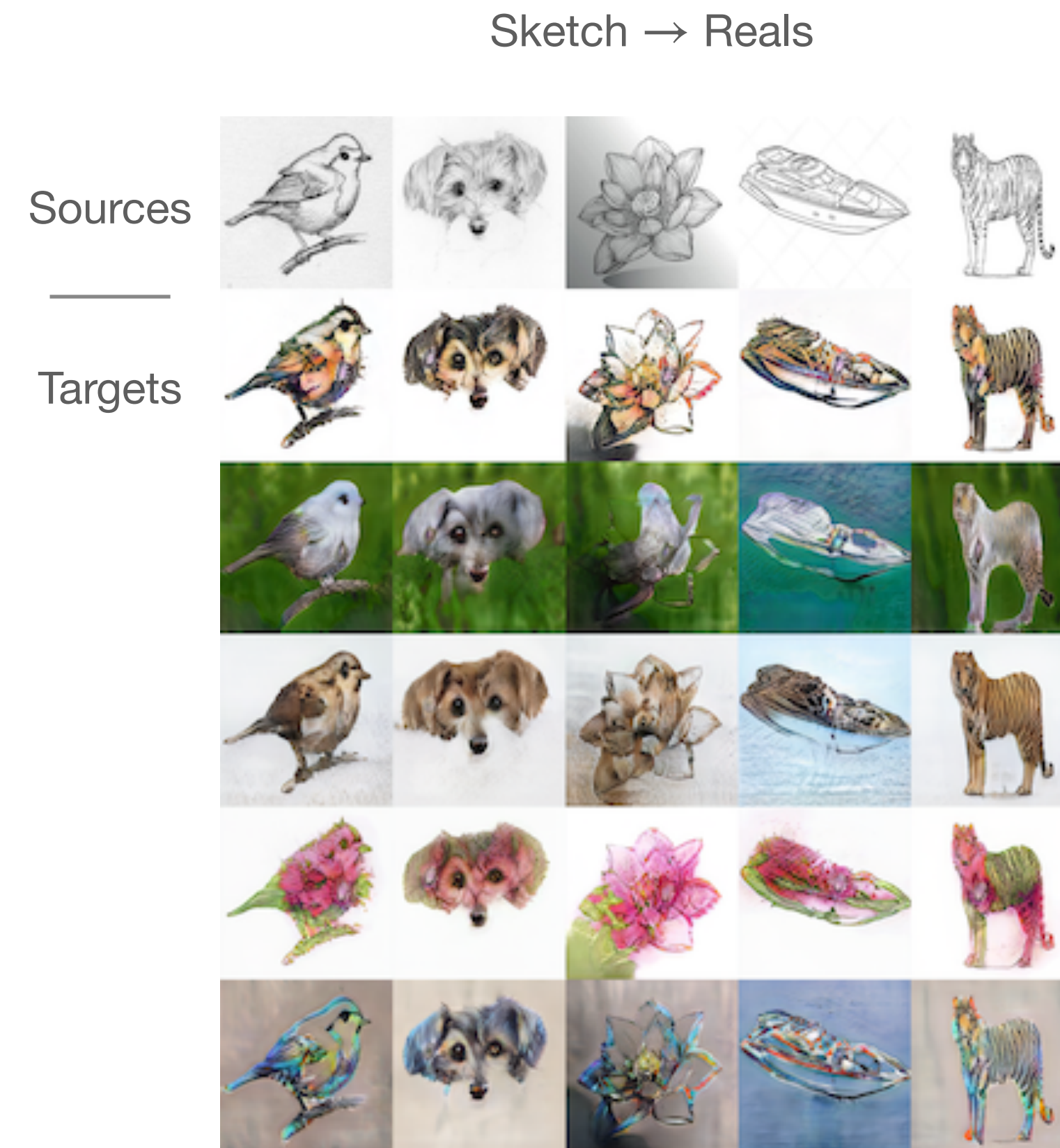


SVHN \rightarrow MNIST



Shortcoming of Unsupervised Domain Translation

Inconsistent style generation.



Integrating Categorical Semantics into Unsupervised Domain Translation

In collaboration with Faruk Ahmed and Aaron Courville

Potential approaches

Supervised

Objectives leveraging labels

Objectives leveraging pairing

Objectives leveraging pre-trained representation

Unsupervised

Inductive bias via the architecture

Unsupervised objectives

Objective leveraging pre-trained representation without supervision

Potential approaches

Supervised

Objectives leveraging labels

Objectives leveraging pairing

Objectives leveraging pre-trained representation

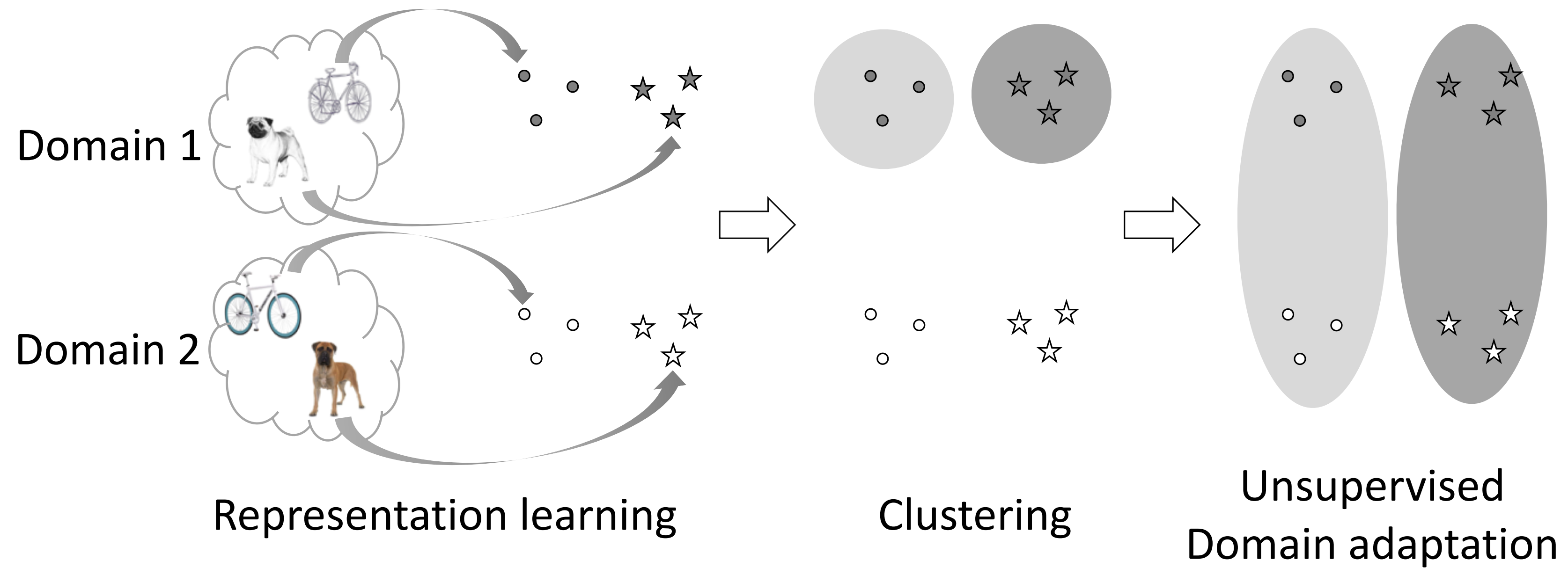
Unsupervised

Inductive bias via the architecture

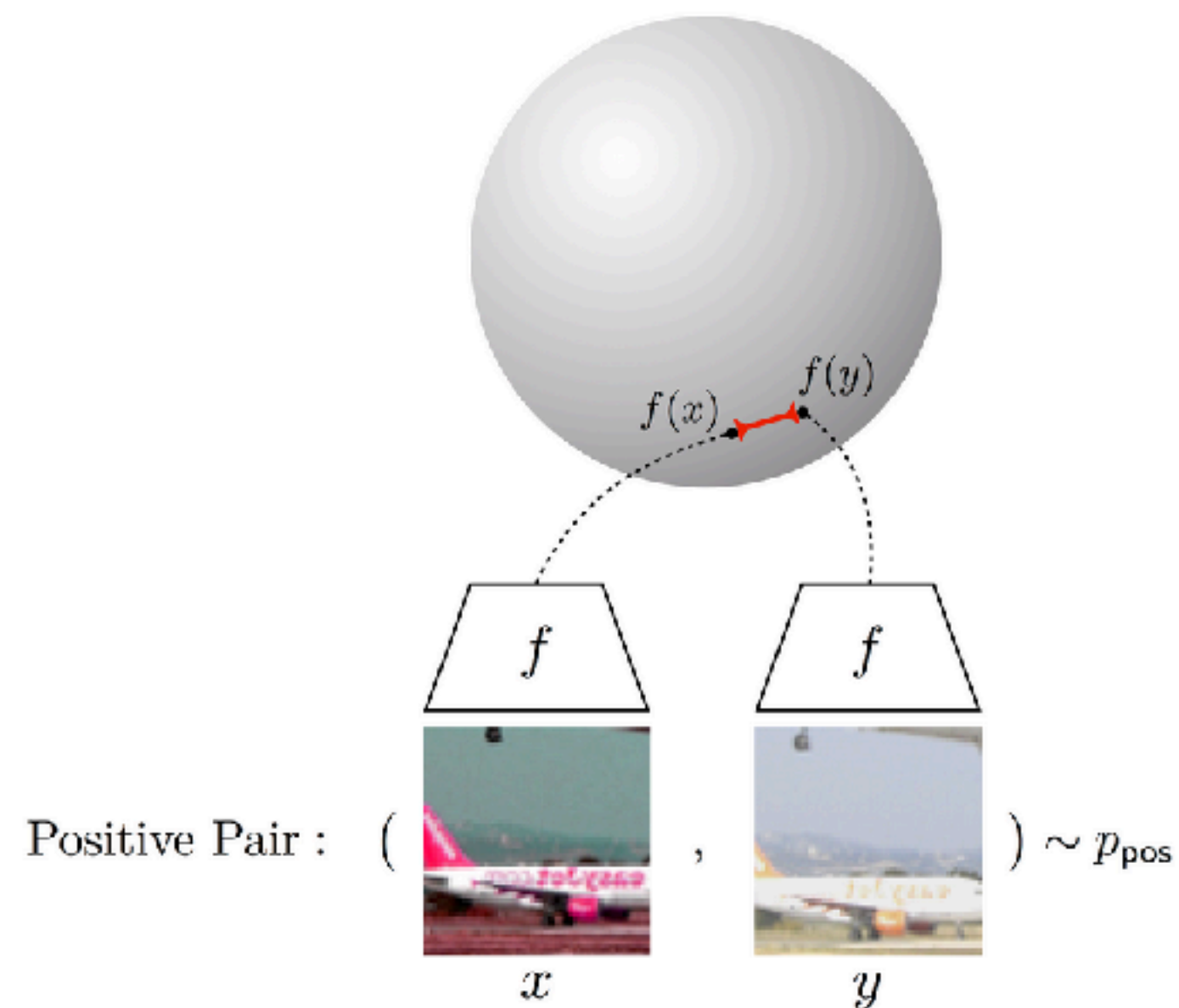
Unsupervised objectives

Objective leveraging pre-trained representation without supervision

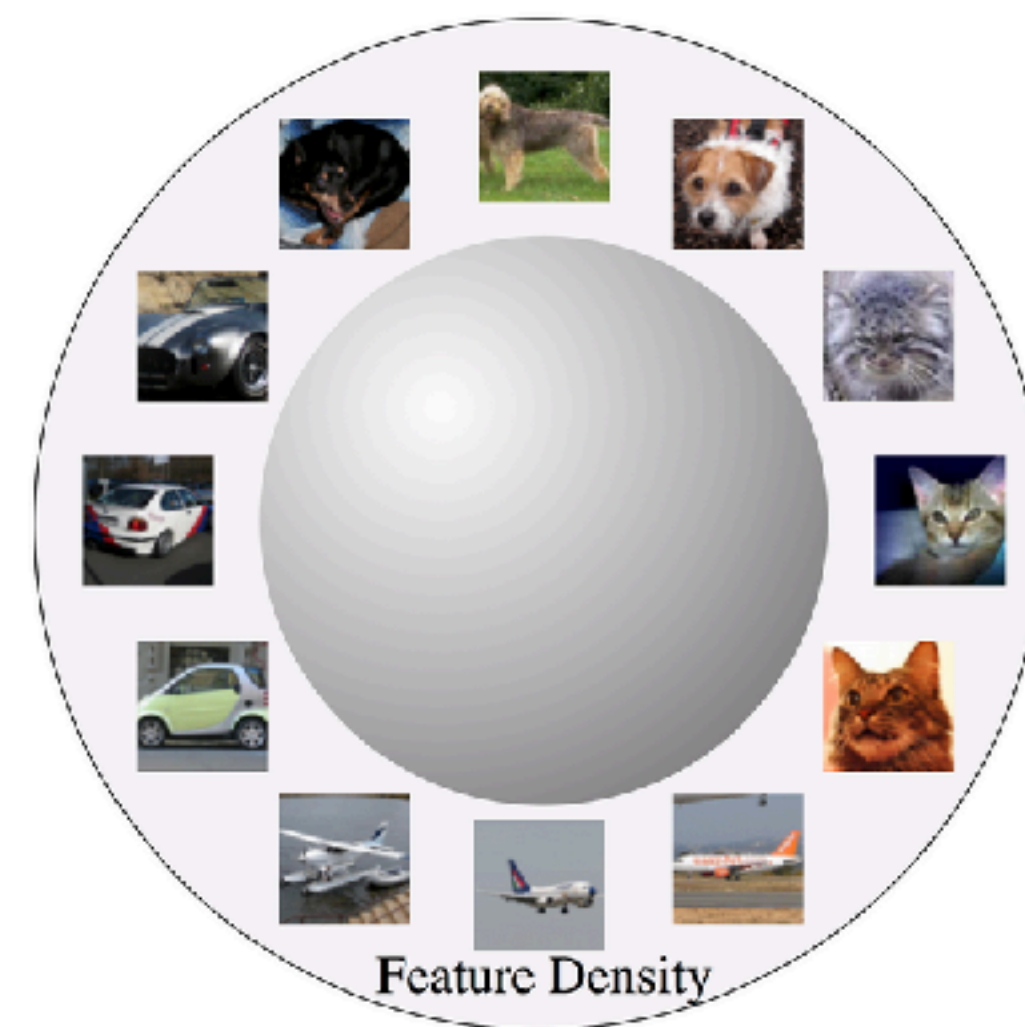
Learning domain invariant representation without supervision



Representation learning – Self-supervised learning

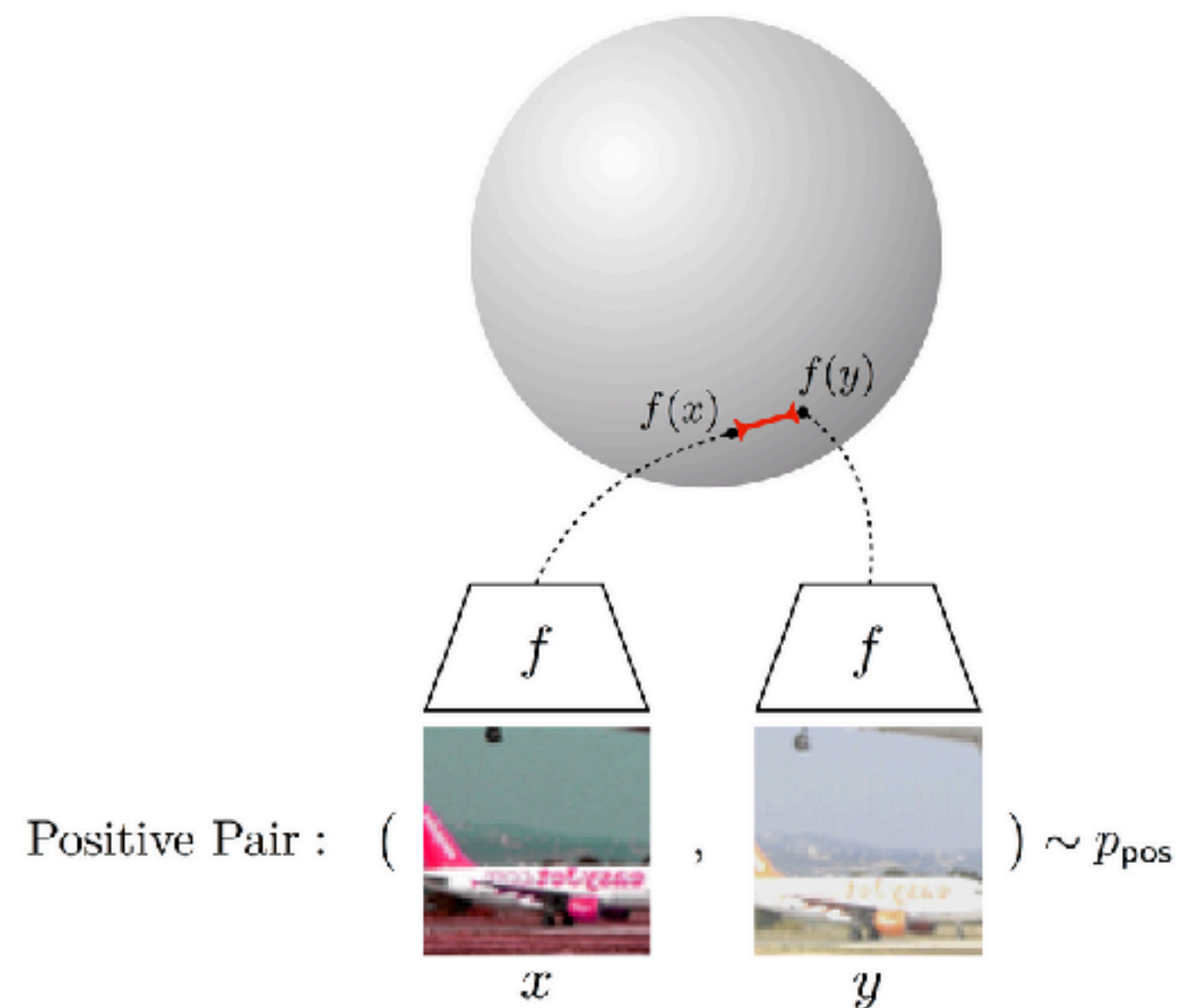


Alignment: Similar samples have similar features

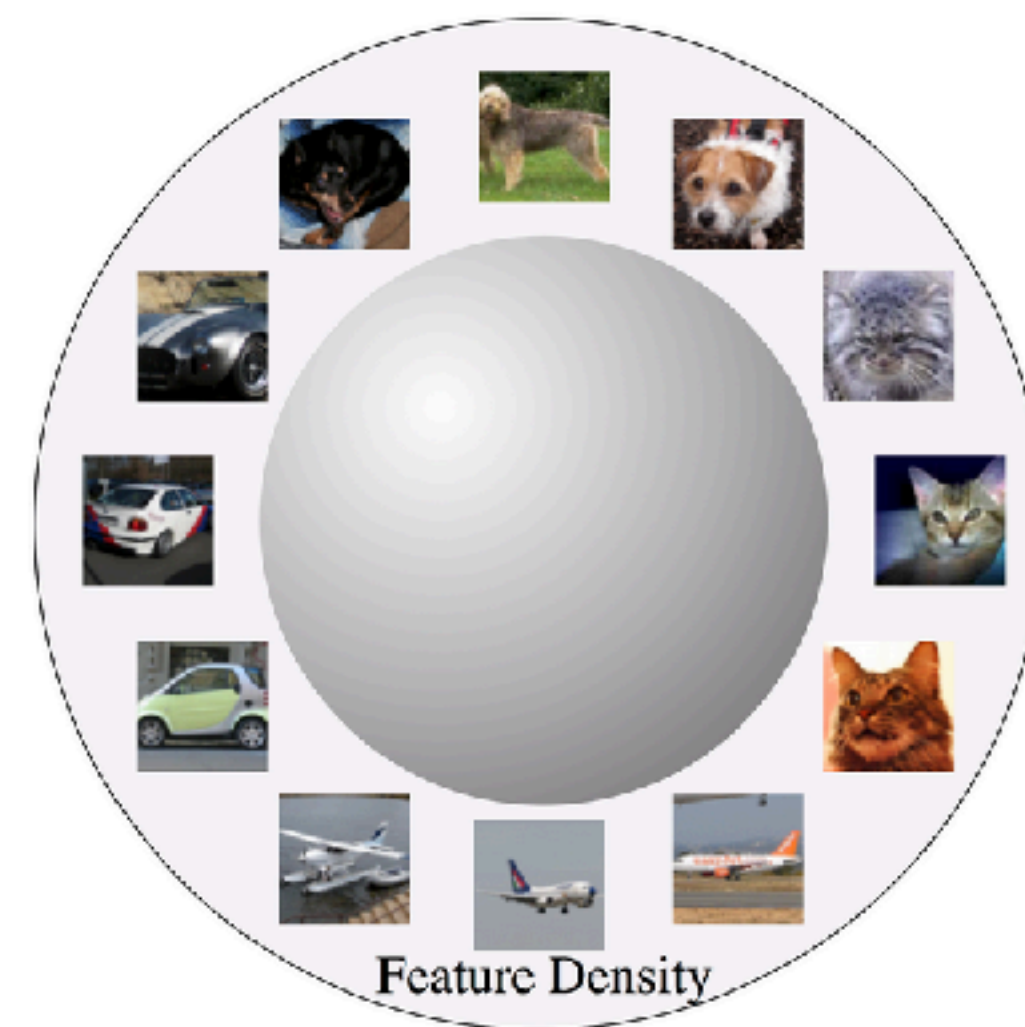


Uniformity: Preserve maximal information

Representation learning – Self-supervised learning



Alignment: Similar samples have similar features



Uniformity: Preserve maximal information

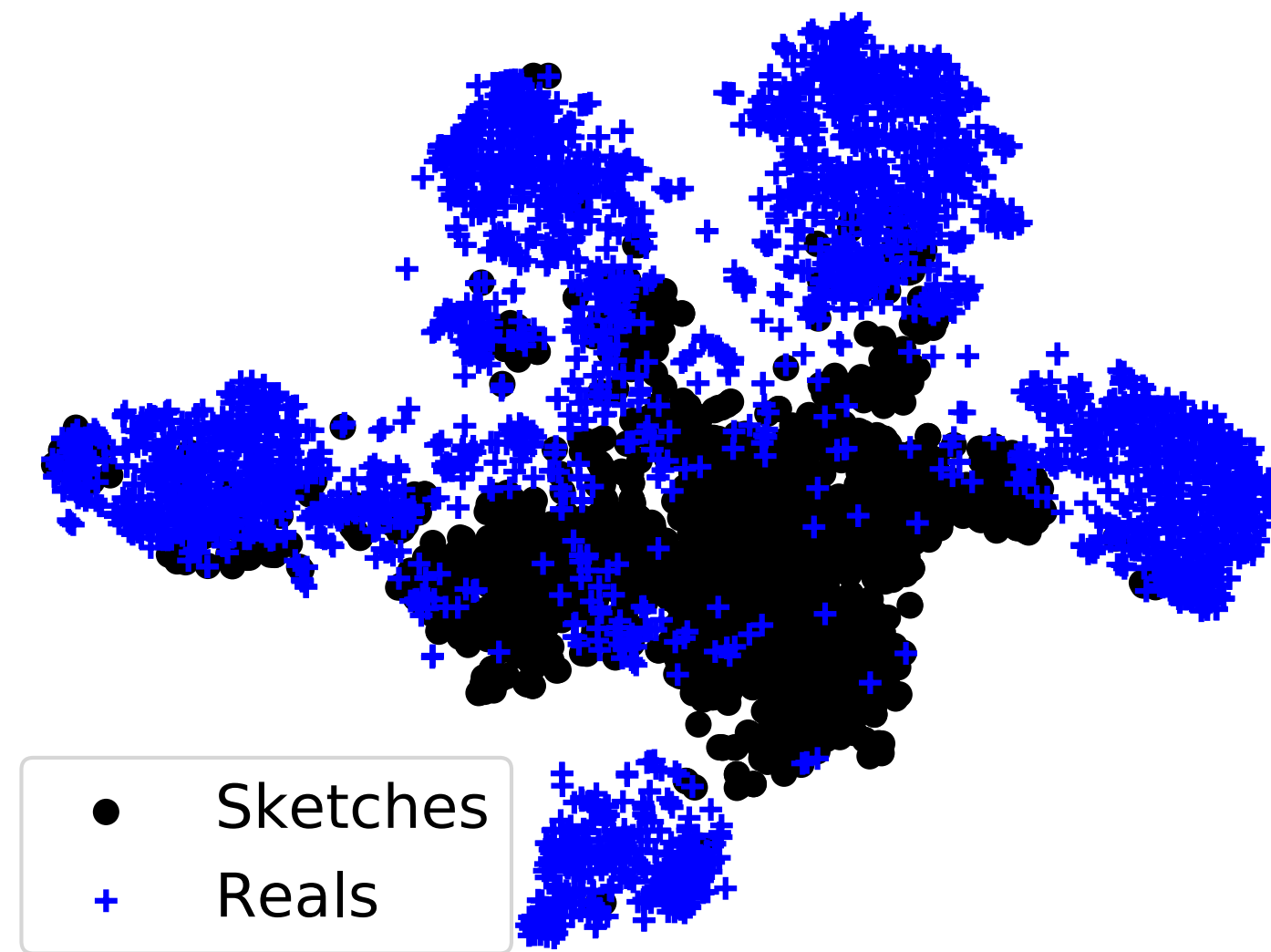
Noise contrastive estimation

$$\mathcal{L}_{\text{nce}} := -\log \frac{\exp(d(f(x), f(y))/\tau)}{\sum_{\bar{x} \in \mathcal{X} \setminus x} \exp(d(f(x), f(\bar{x}))/\tau)}$$

Property of a model learned with contrastive learning

Samples cluster in dense region

Samples from different domain do not intersect

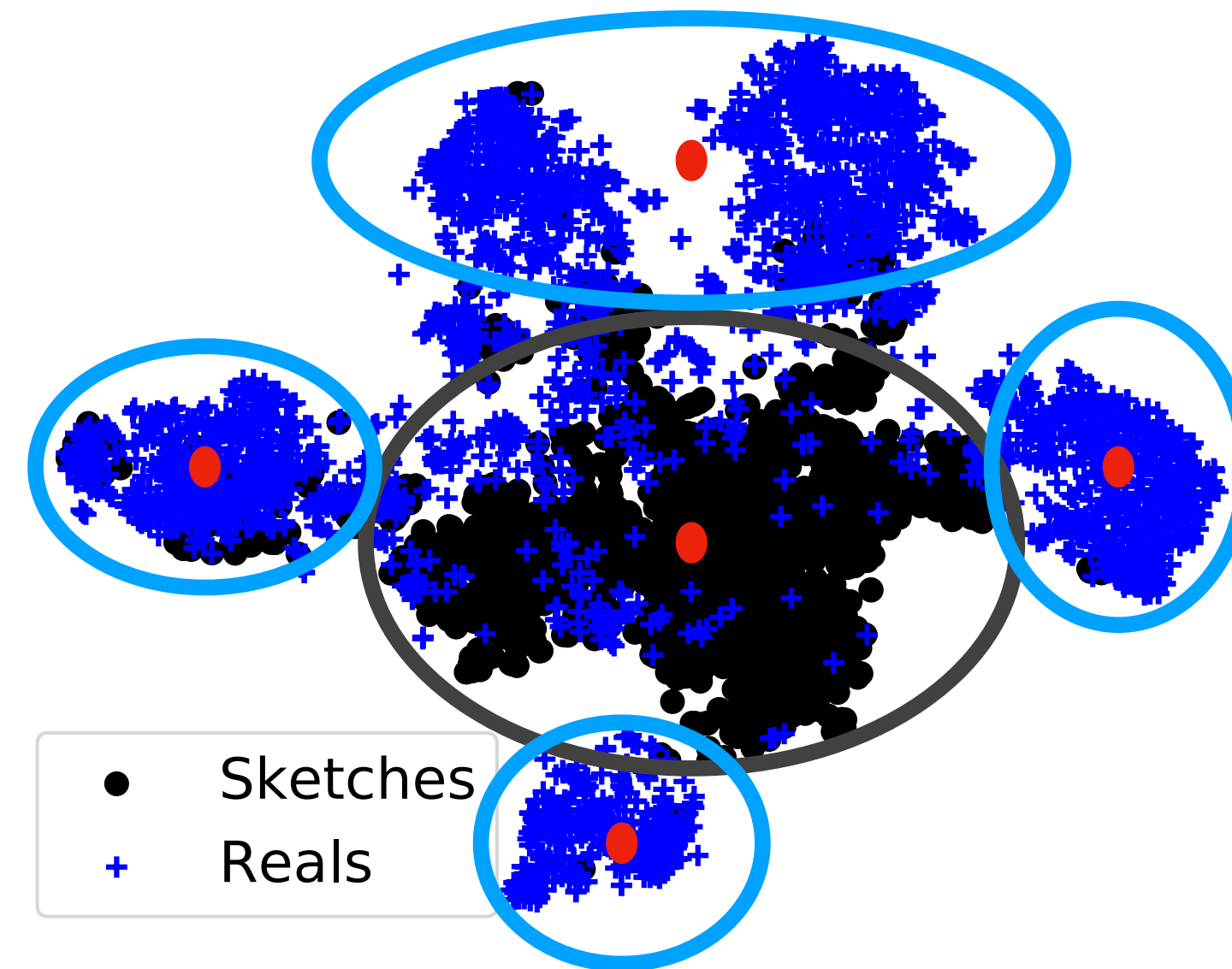


Embedding of 5 categories: bird, dog, flower, boat, tiger

Property of a model learned with contrastive learning

Samples cluster in dense region

Samples from different domain do not intersect

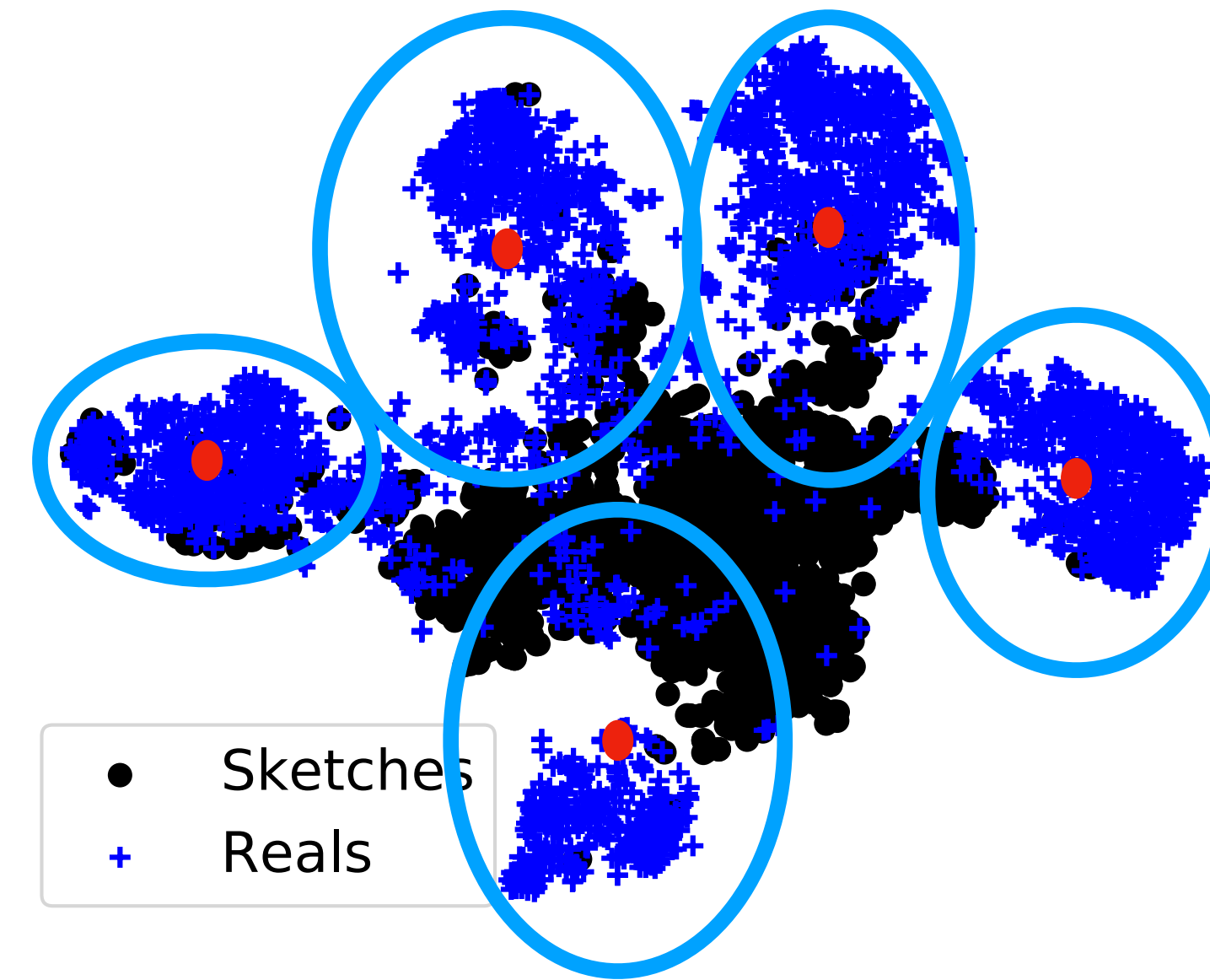


Embedding of 5 categories: bird, dog, flower, boat, tiger

Property of a model learned with contrastive learning

Samples cluster in dense region

Samples from different domain do not intersect

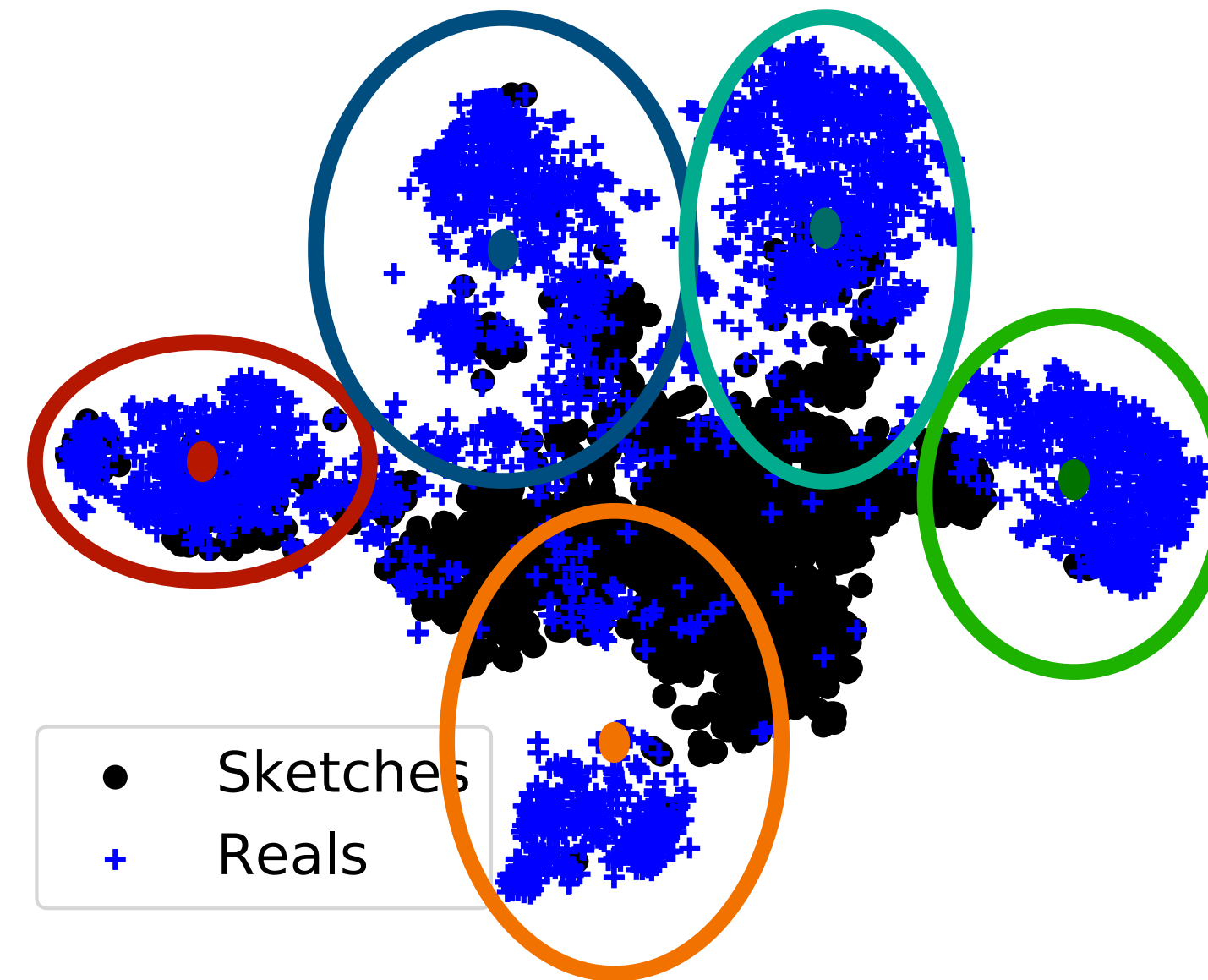


Embedding of 5 categories: bird, dog, flower, boat, tiger

Property of a model learned with contrastive learning

Samples cluster in dense region

Samples from different domain do not intersect

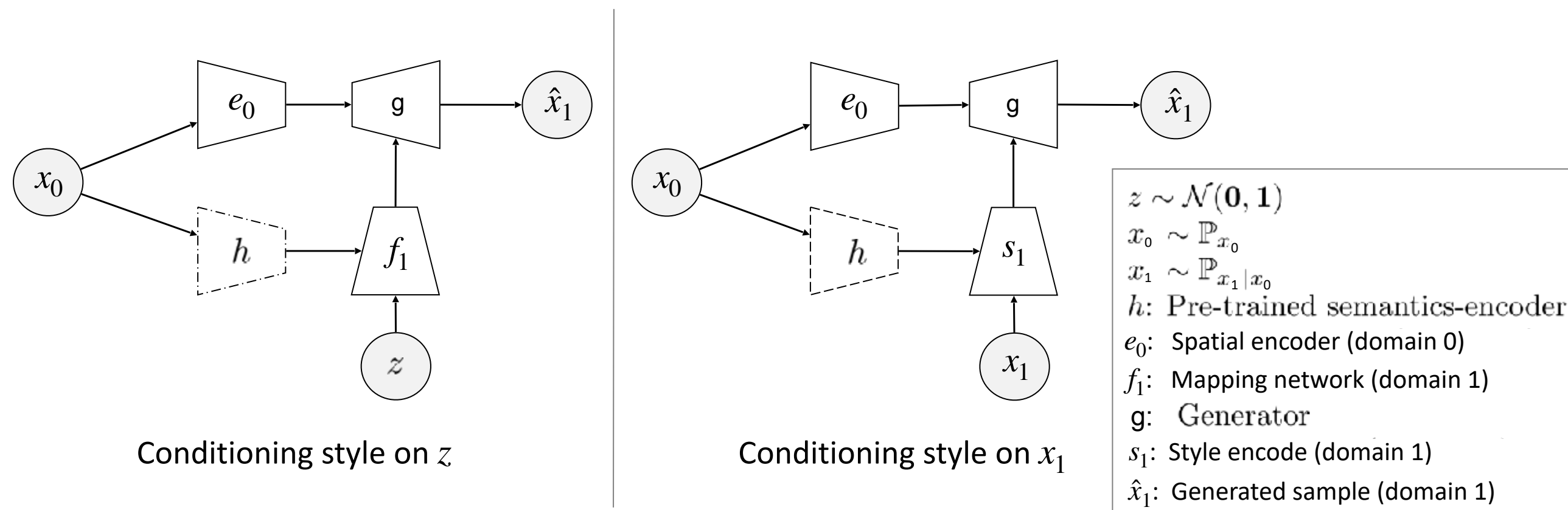


Embedding of 5 categories: bird, dog, flower, boat, tiger

Define the clusters as pseudo-labels and adapt them to the sketches using
Unsupervised Domain Adaptation

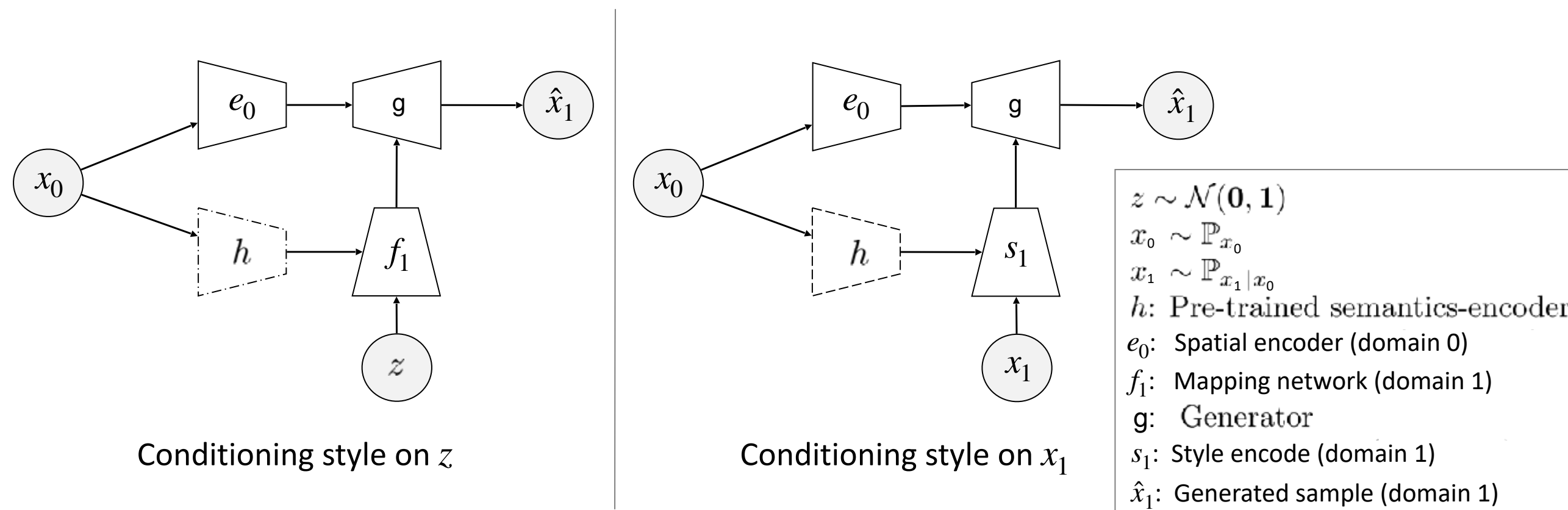
Integrate the learned semantics into Unsupervised Domain Translation

Condition style generation



Integrate the learned semantics into Unsupervised Domain Translation

Condition style generation



Constraint mapping to preserve semantics

$$\mathcal{L} := - \sum_i h(x_0)_i \log(h(\hat{x}_1))_i$$

Results MNIST ↔ SVHN

MNIST → SVHN using our method



Results MNIST ↔ SVHN

Table 1: **Comparison with the baselines.** Domain translation accuracy and FID obtained on MNIST (M) ↔ SVHN (S) for the different methods considered. The last column is the test classification accuracy of the classifier used to compute the metric. *: Using weak supervision.

	Data	CycleGAN	MUNIT	DRIT	Stargan-V2	EGSC-IT*	CatS-UDT	Target
Acc	M→S	10.89	10.44	13.11	28.26	47.72	95.63	98.0
	S→M	11.27	10.12	9.54	11.58	16.92	76.49	99.6
FID	M→S	46.3	55.15	127.87	66.54	72.43	39.72	-
	S→M	24.8	30.34	20.98	26.27	19.45	6.60	-

MNIST→SVHN using our method



Results MNIST ↔ SVHN

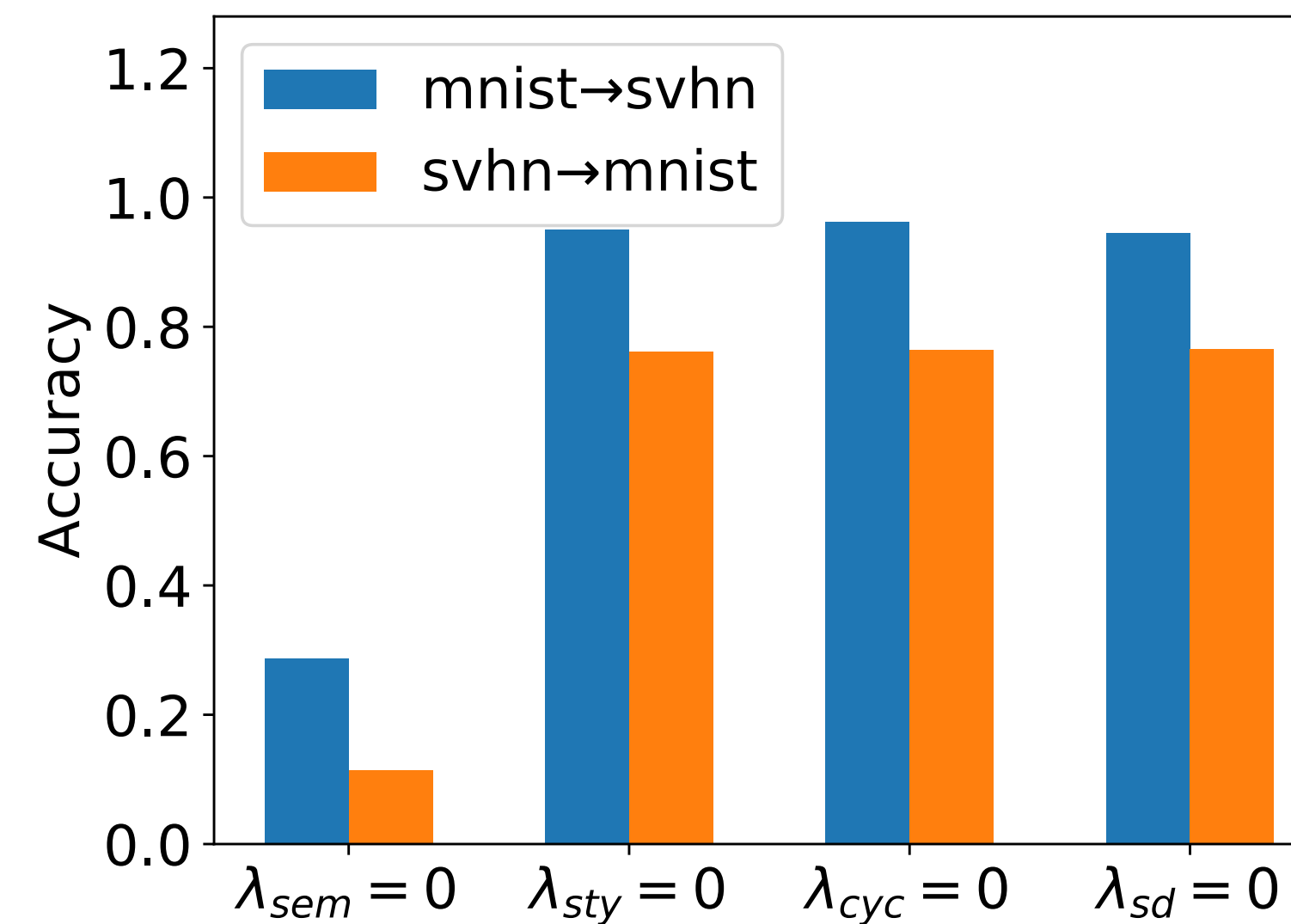
Table 1: **Comparison with the baselines.** Domain translation accuracy and FID obtained on MNIST (M) ↔ SVHN (S) for the different methods considered. The last column is the test classification accuracy of the classifier used to compute the metric. *: Using weak supervision.

	Data	CycleGAN	MUNIT	DRIT	Stargan-V2	EGSC-IT*	CatS-UDT	Target
Acc	M→S	10.89	10.44	13.11	28.26	47.72	95.63	98.0
	S→M	11.27	10.12	9.54	11.58	16.92	76.49	99.6
FID	M→S	46.3	55.15	127.87	66.54	72.43	39.72	-
	S→M	24.8	30.34	20.98	26.27	19.45	6.60	-

MNIST→SVHN using our method



Ablating losses

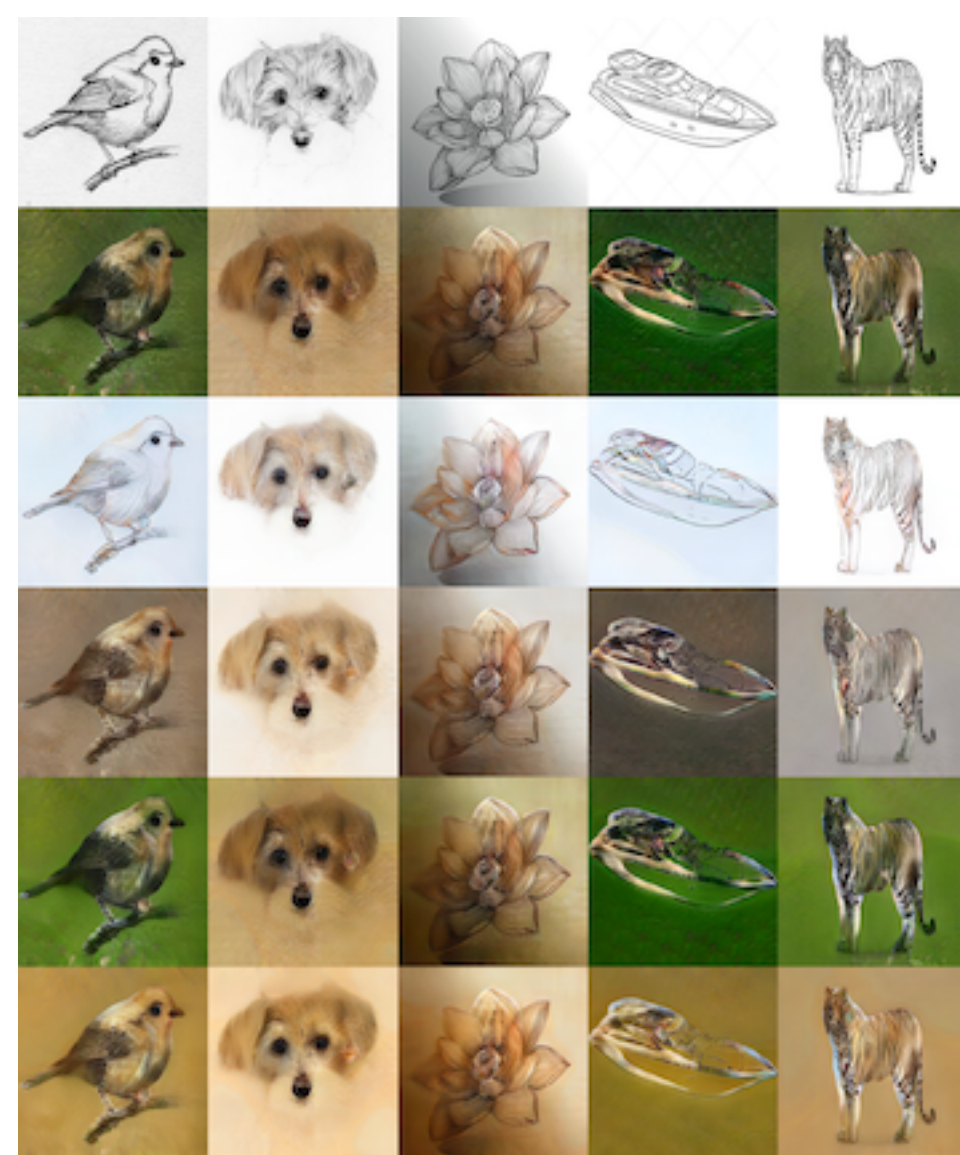


Results Sketches → Reals

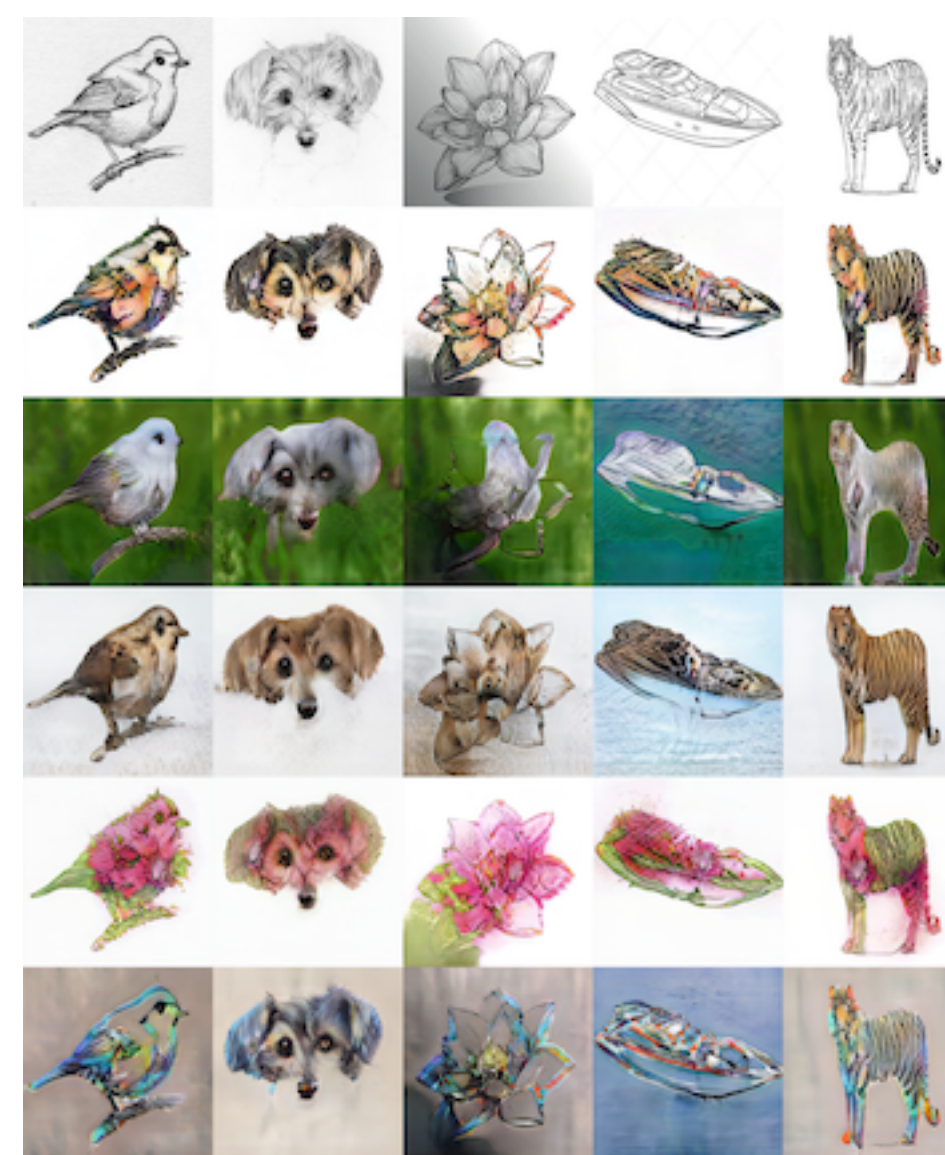
DRIT



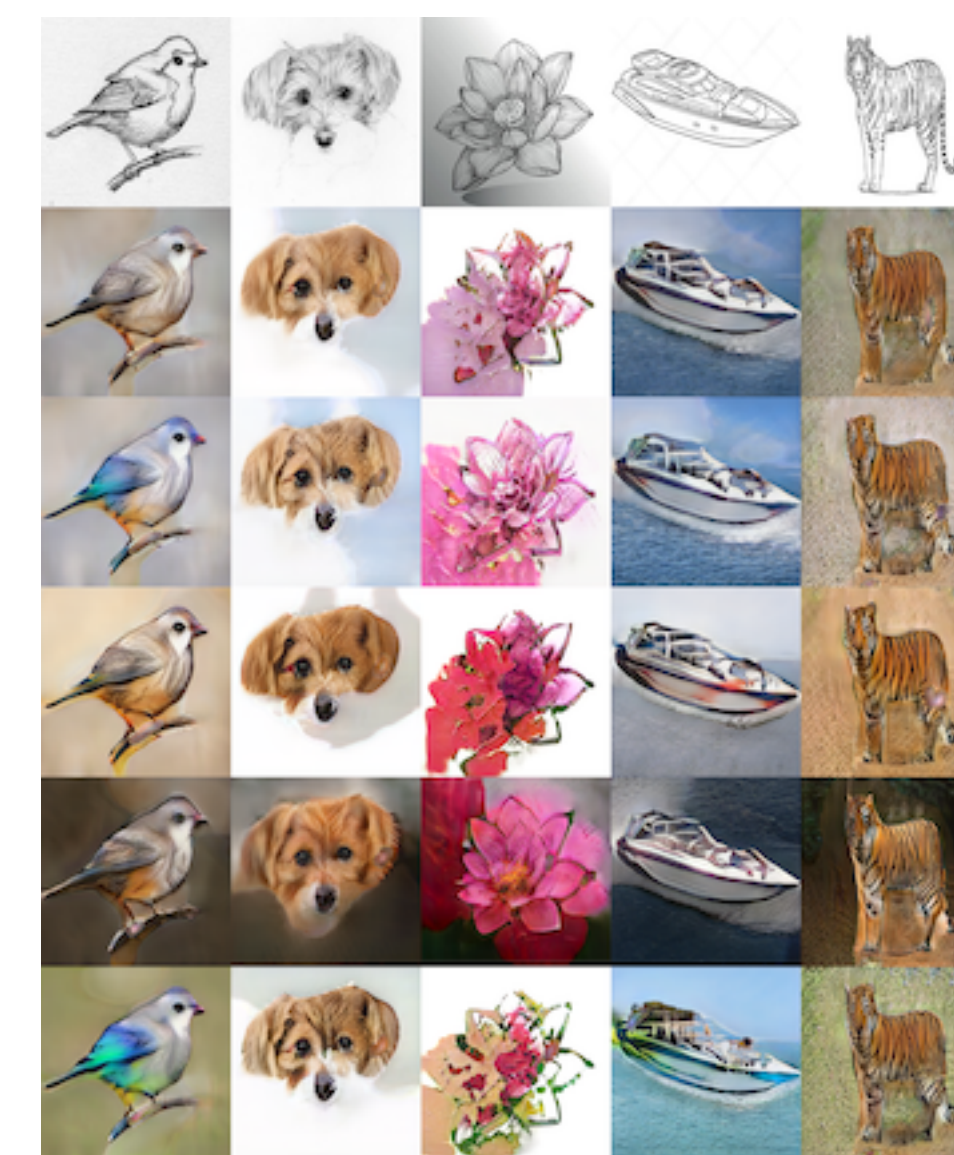
EGST-IT



StarGAN-V2



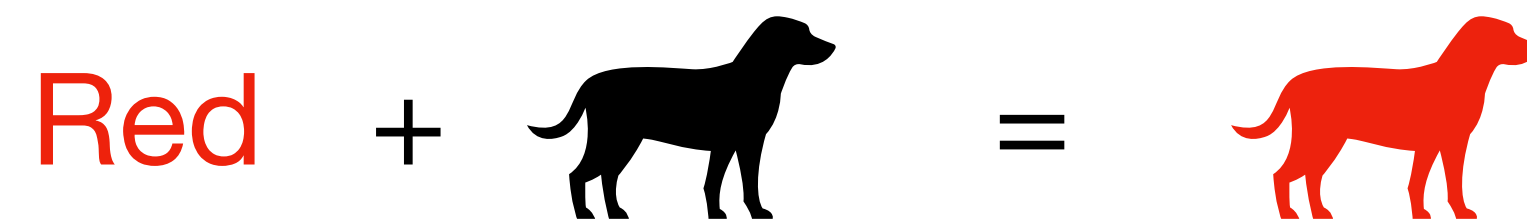
CatS-UDT (ours)



Emergence of structure in artificial language

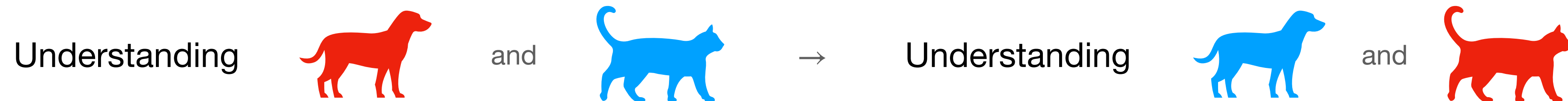
Compositionality

The meaning of a complex expression is determined by the meanings of its constituent expressions and the rules used to combine them.



Systematicity

The capacity to understand a complex expression implies the capacity to understand structurally related expressions.



Language as the solution to a coordination problem

Language as the solution to a coordination problem

Vervet monkeys have a idiosyncratic call depending on the predator

Predator

Leopard

Eagle

Snake



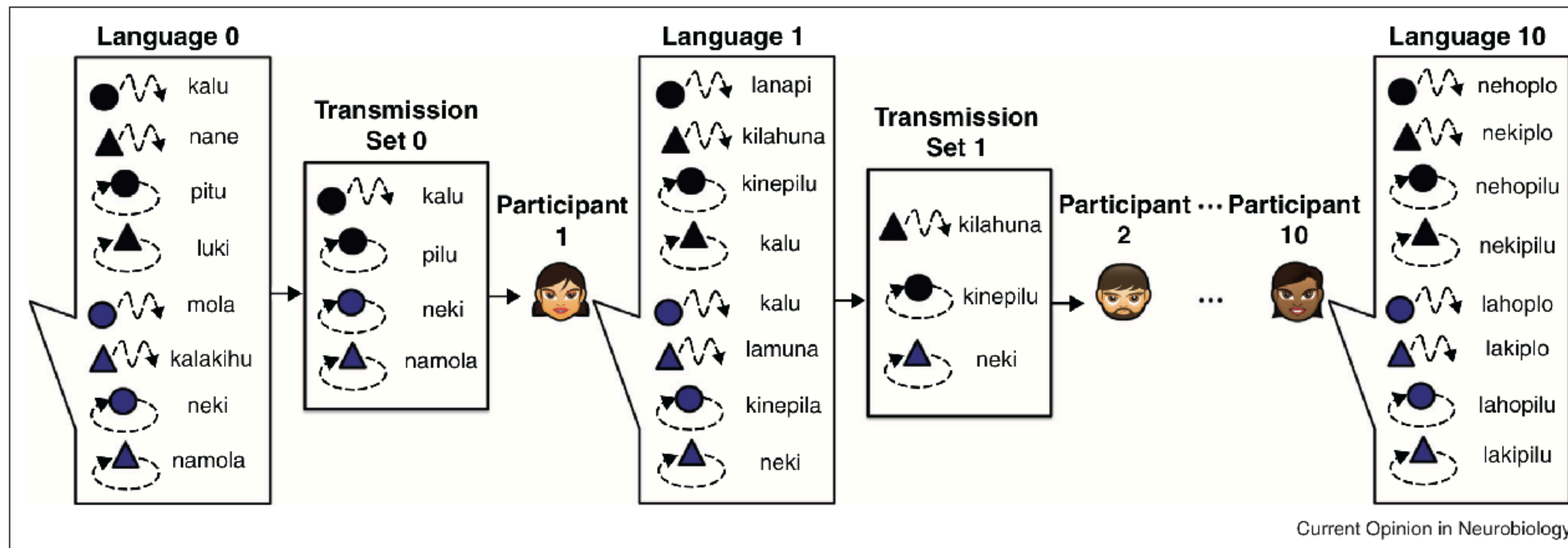
Alarm

Loud barking

Short double syllable cough

Shutter

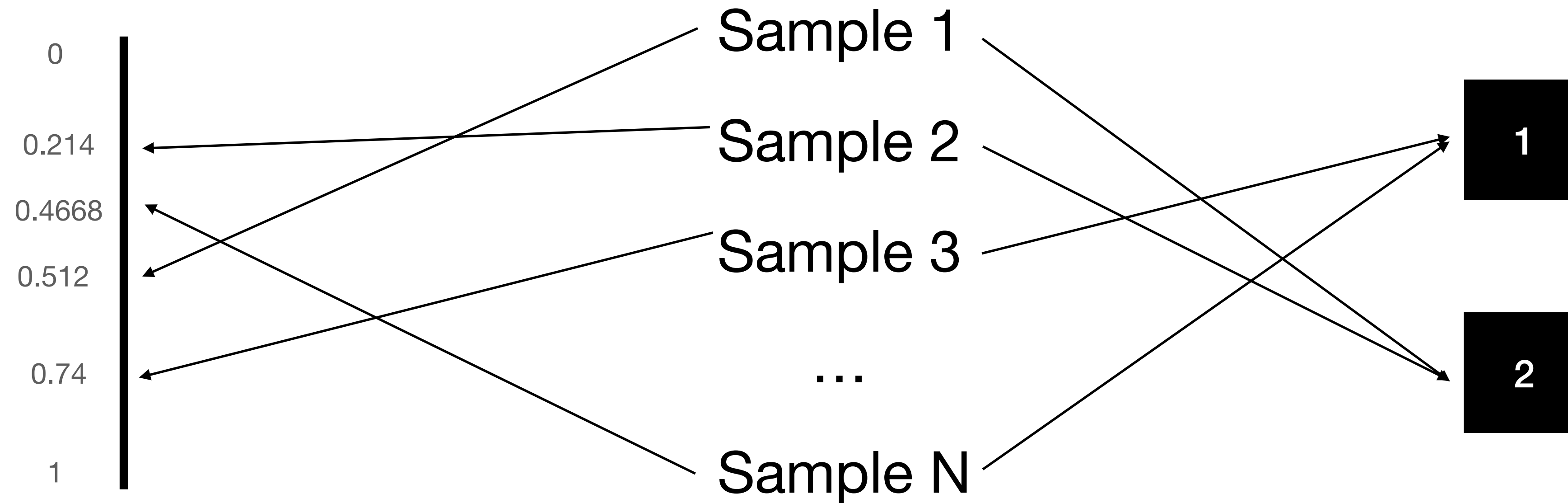
Cultural transmission



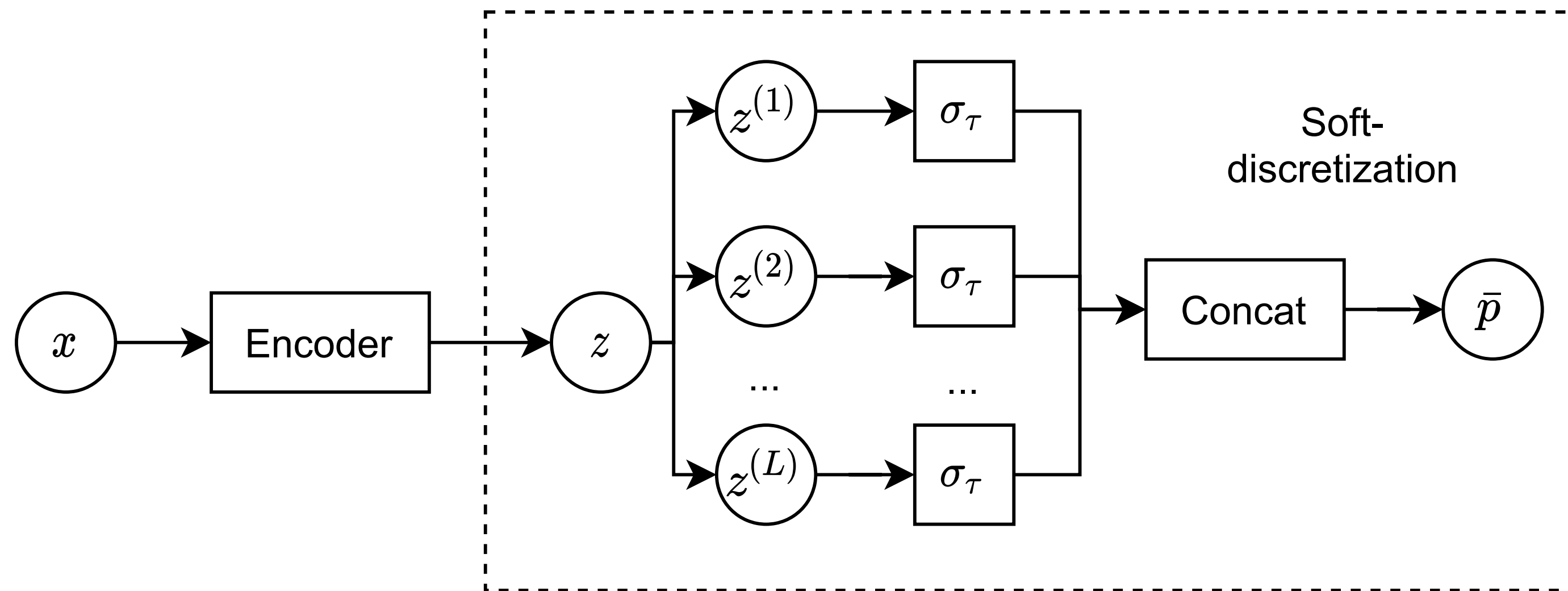
Soft-discretization bottleneck for self-supervised learning

In collaboration with Christos Tsirigotis, Max Schwarzer, Ankit Vani and Aaron Courville.

Continuous vs Discrete representation



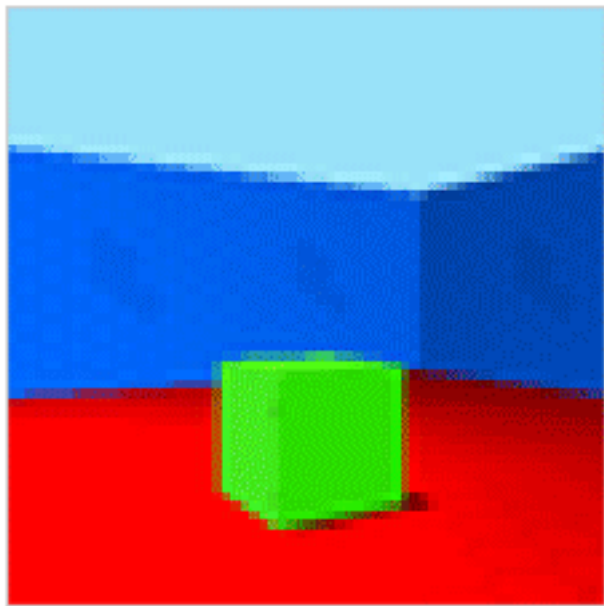
Soft-discretization bottleneck



$$z^{(i)} \in \mathbb{R}^V.$$
$$\sigma_\tau(z^{(j)})_i := \frac{e^{z_i^{(j)}/\tau}}{\sum_{k=0}^V e^{z_k^{(j)}/\tau}}.$$

Systematic generalization

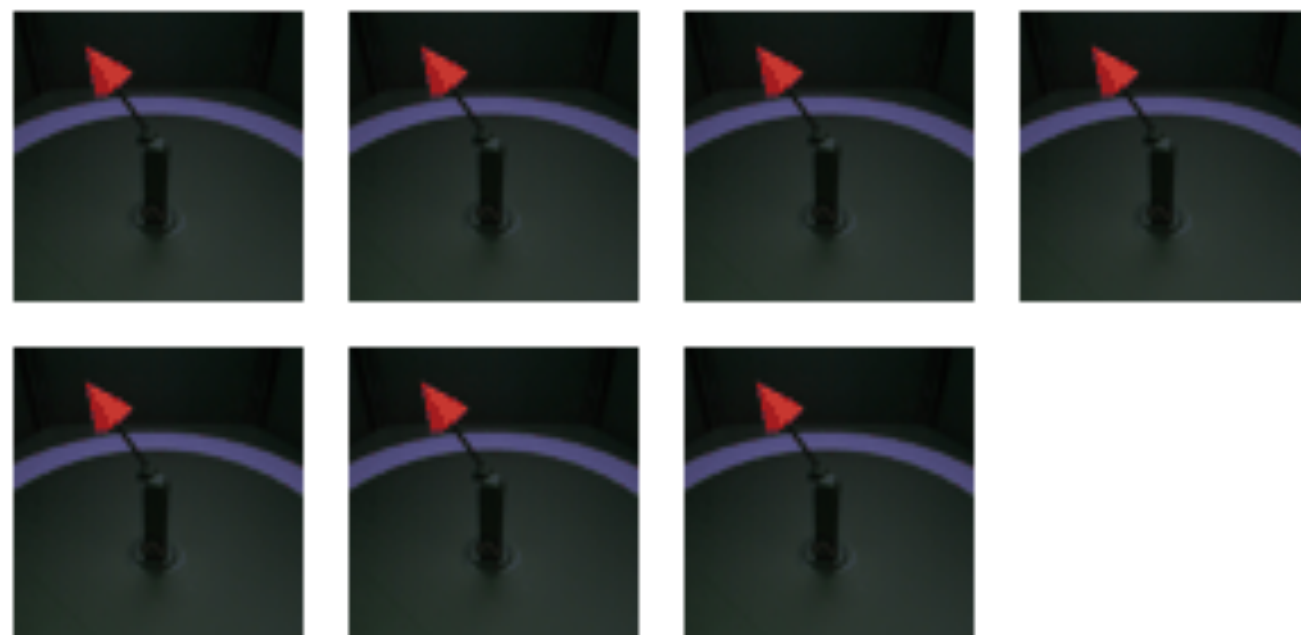
Shapes3d



dSprites

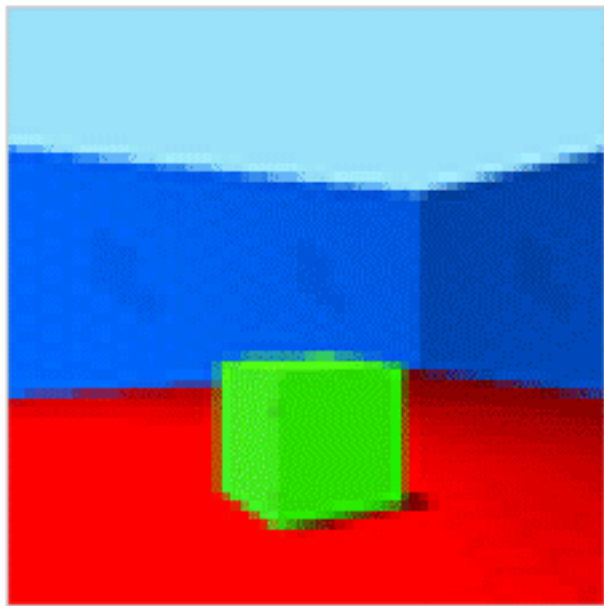


MPI3D



Systematic generalization

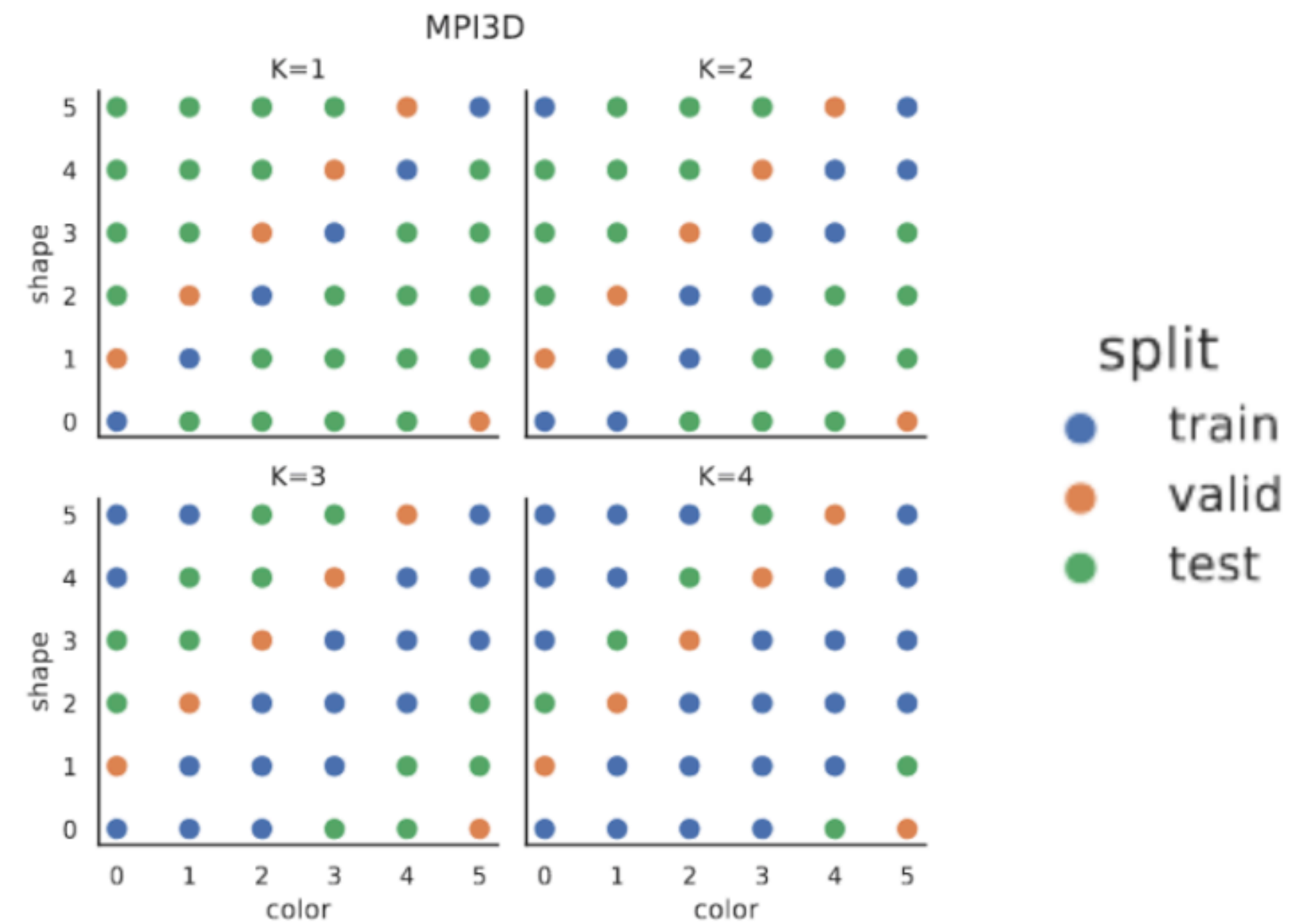
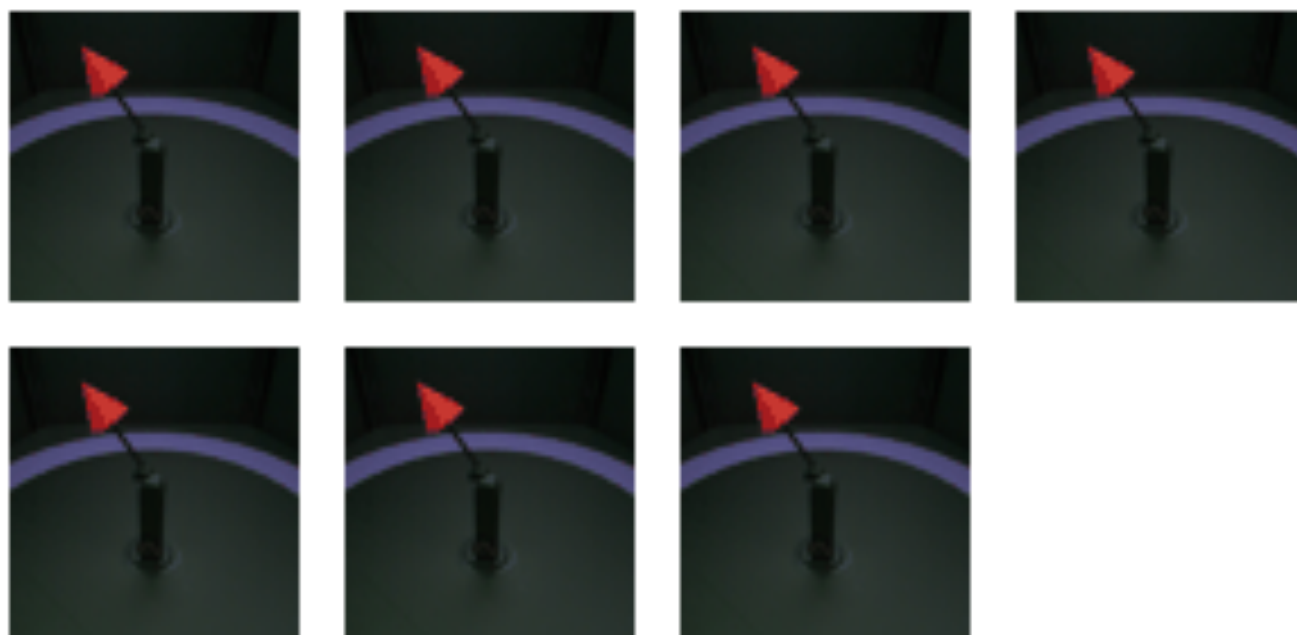
Shapes3d



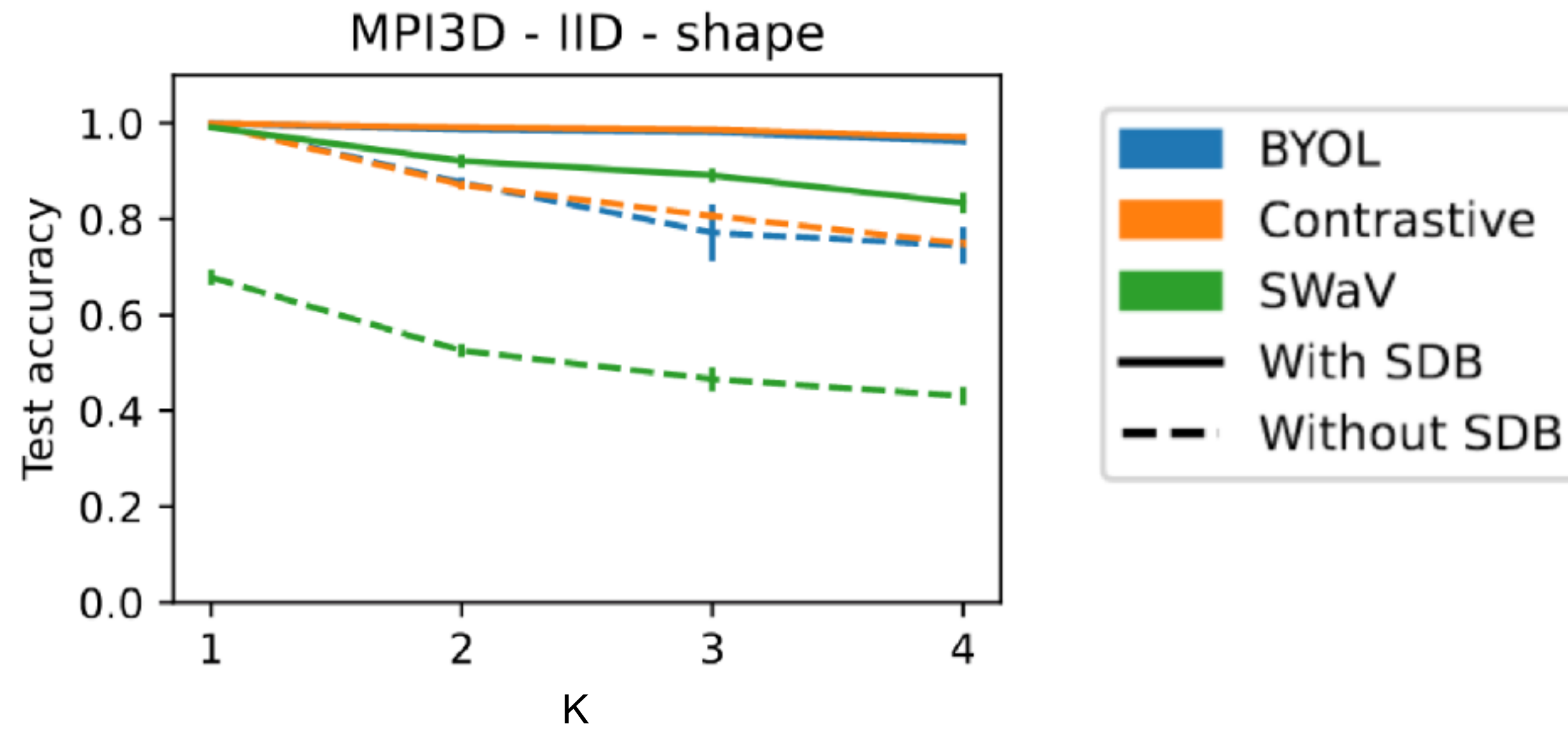
dSprites



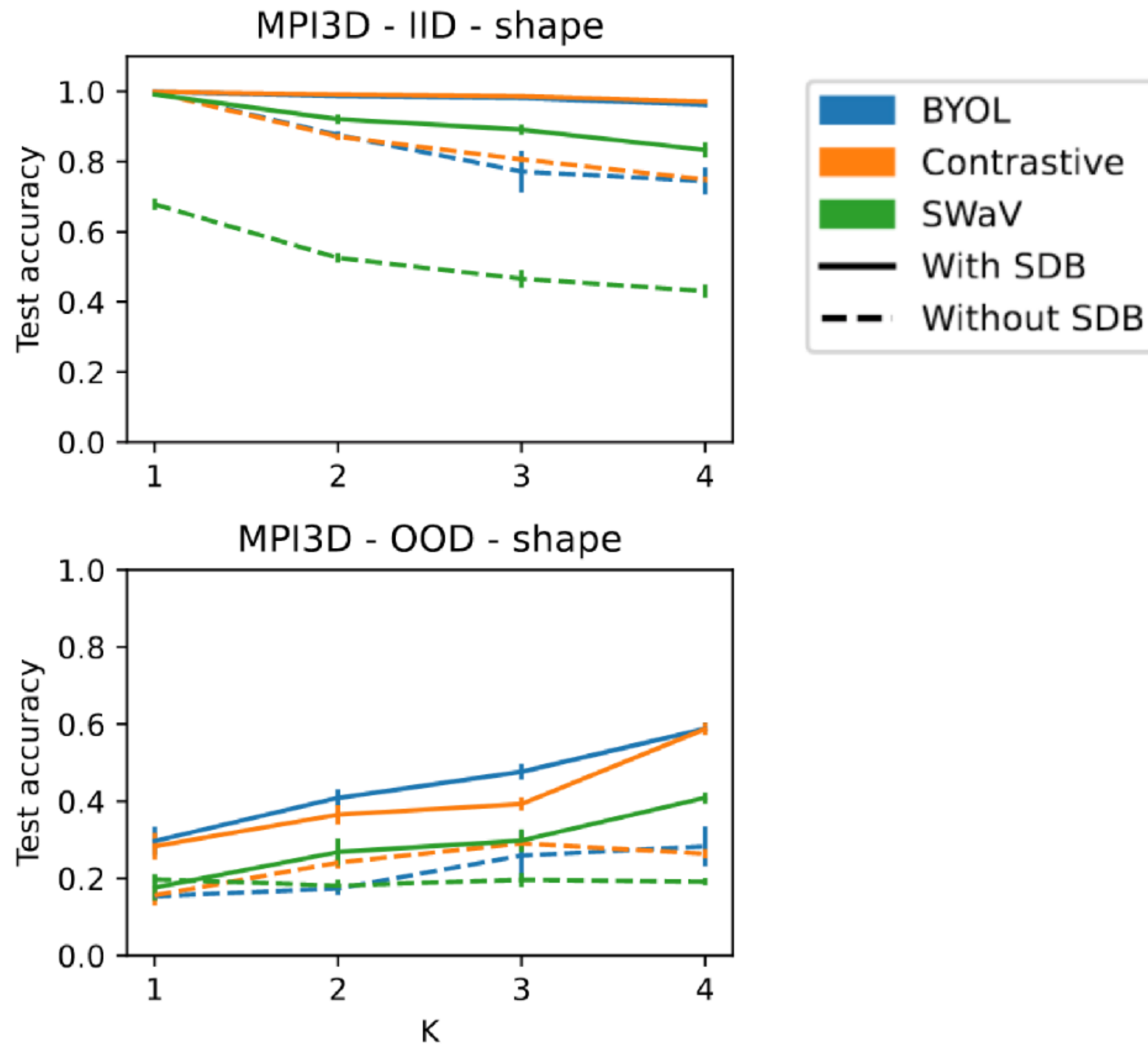
MPI3D



Results – Systematic generalization



Results – Systematic generalization



Results — Systematic generalization

	+ noise	Argmax	Softmax	VQ	Test accuracy
Baseline					0.29 ± 0.01
	✓				0.28 ± 0.01
Soft-Discrete			✓		0.52 ± 0.03
	✓		✓		0.49 ± 0.04
Hard-Discrete	✓	✓			0.32 ± 0.03
				✓	0.39

Table 1. Ablation of SSL-SB on MPI3D-K:3. We test the effect of adding noise, a hard discretization bottleneck via Gumbel-Softmax straight-through estimation and Vector Quantization and the soft discretization bottleneck.

Results — Systematic generalization

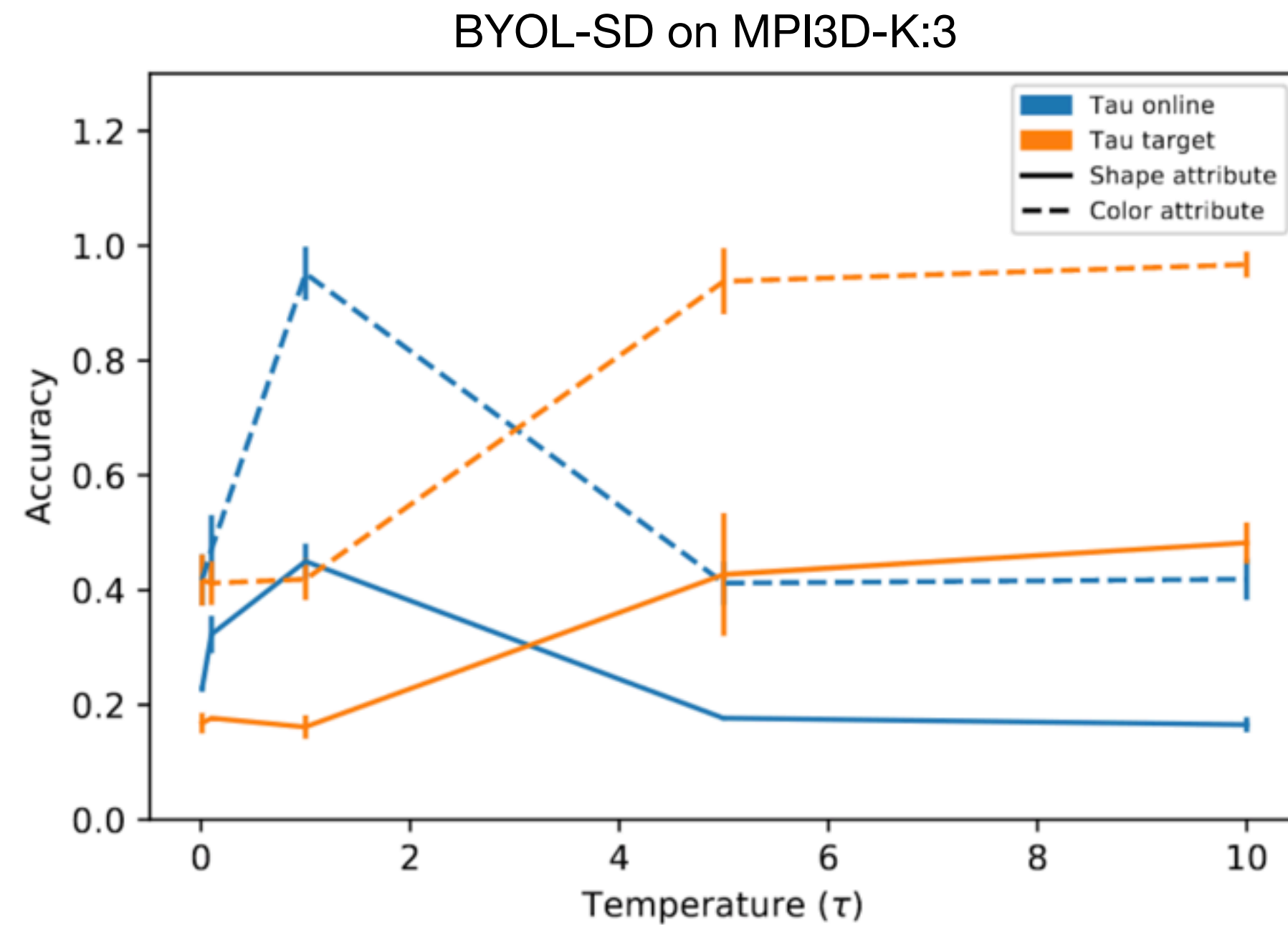


Figure 6. Study of the effect of the temperature parameter on the online (τ_O) and the target (τ_T) networks. We fix the temperature $\tau_O = 1.5$ when interpolating τ_T and $\tau_T = 4.0$ when interpolating τ_O .

Robustness to distribution shift

Train dataset



ImageNet

Test datasets



ImageNet-C



ImageNet-A



ImageNet-R



ImageNet-V2

Dan Hendricks, Thomas Dietterich. Benchmarking Neural Networks Robustness to Common Corruptions and Perturbations. 2019.

Dan Hendricks et al. Natural Adversarial Examples. 2021.

Dan Hendricks et al. The Many Faces of Robustness. 2019.

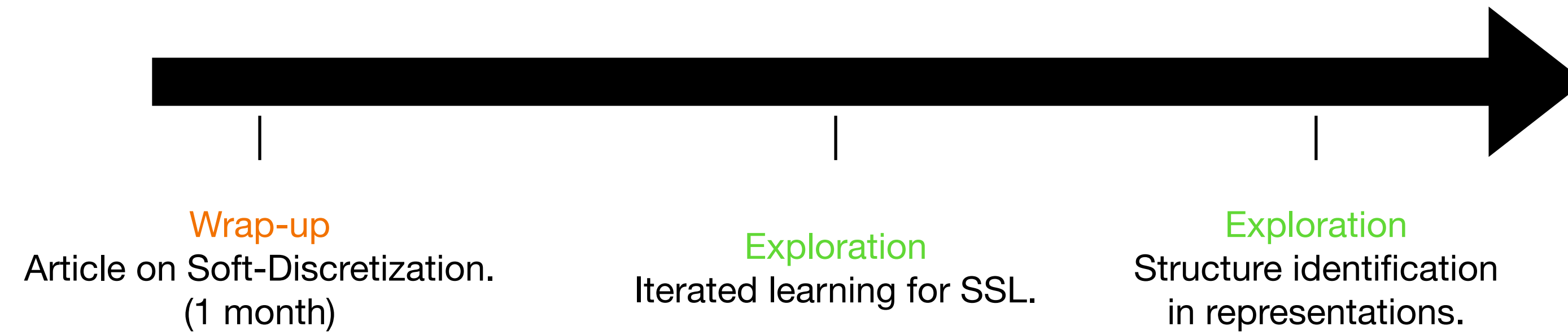
Benjamin Recht et al. Do ImageNet Classifiers Generalize to ImageNet? 2019.

Results – Robustness

	Imagenet	Imagenet-v2	Imagenet-r	Imagenet-a	Imagenet-c
BYOL	67.16	53.96	15.35	0.87	33.32
BYOL + SDB	70.22	57.73	17.95	1.01	37.98

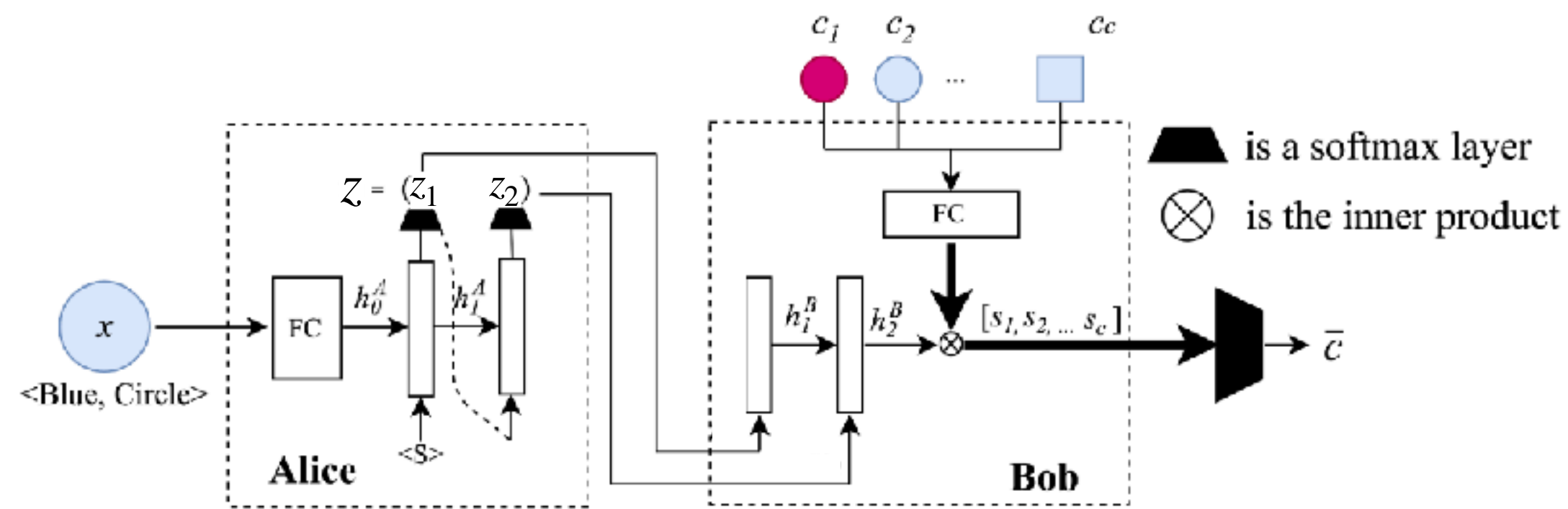
Future works

Future works



Iterated learning for communication games

Interaction: Object selection game



$$\nabla_{\theta} J := E [R(\bar{c}, \mathbf{x}) \nabla_{\theta} \log p_f(\mathbf{z} | \mathbf{x}) + \lambda_f \nabla_{\theta} H[p_f(\mathbf{z} | \mathbf{x})]]$$

$$\nabla_{\psi} J := E [R(\bar{c}, \mathbf{x}) \nabla_{\psi} \log p_g(\bar{c} | \mathbf{z}, c_1, \dots, c_k) + \lambda_g \nabla_{\psi} H(p_g(\bar{c} | \mathbf{x}, c_1, \dots, c_n))]$$

Generation

$$\mathcal{Z} := \{(\mathbf{x}_i, \mathbf{z}_i)\}_{i=1}^N$$

Distillation

$$\min_{\theta^{t+1}} E_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{Z}} l(f_{\theta^{t+1}}(\mathbf{x}), \mathbf{z})$$

* l is defined as the cross-entropy

Algorithm 1 Neural Iterated Learning

Require: $\mathcal{X}, f_{\theta^0}, g_{\phi^0}, N_{\text{iter}}, M_{\text{interaction}}$

θ^0 randomly initialized

ϕ^0 randomly initialized

$t \leftarrow 0$

while $N_{\text{iter}} \neq 0$ **do**

$S \leftarrow 0$

while $S \neq M$ **do**

$\theta^t \leftarrow \theta^t + \alpha \nabla_{\theta^t} J$

$\psi^t \leftarrow \psi^t + \alpha \nabla_{\psi^t} J$

$S \leftarrow S + 1$

end while

$\mathcal{Z} \leftarrow \text{Generation}(\mathcal{X}, f_{\theta^t})$

$\theta^{t+1} \leftarrow \text{Distillation}(\mathcal{Z}, f_{\theta^{t+1}})$

ψ^{t+1} randomly initialized

$t \leftarrow t + 1$

$N_{\text{iter}} \leftarrow N_{\text{iter}} - 1$

end while

} Interaction

Iterated learning for self-supervised learning

Interaction: Self-supervised learning objective

Example: Noise contrastive estimation, BYOL

Generation

$$\mathcal{Z} := \{(\mathbf{x}_i, \mathbf{z}_i)\}_{i=1}^N$$

Distillation

$$\min_{\theta^{t+1}} E_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{Z}} l(f_{\theta^{t+1}}(\mathbf{x}), \mathbf{z})$$

* l is defined as ?

Algorithm 1 Neural Iterated Learning

Require: \mathcal{X} , f_{θ^0} , g_{ϕ^0} , N , $M_{\text{interaction}}$.

θ^0 randomly initialized

ϕ^0 randomly initialized

$t \leftarrow 0$

while $N \neq 0$ **do**

$S \leftarrow 0$

while $S \neq M$ **do**

$\theta^t \leftarrow \theta^t + \alpha \nabla_{\theta^t} J - \alpha \nabla_{\phi^t} J$

$\psi^t \leftarrow \psi^t + \alpha \nabla_{\psi^t} J - \alpha \nabla_{\theta^t} J$

$S \leftarrow S + 1$

end while

$\mathcal{Z} \leftarrow \text{Generation}(\mathcal{X}, f_{\theta^t})$

$\theta^{t+1} \leftarrow \text{Distillation}(\mathcal{Z}, f_{\theta^{t+1}})$

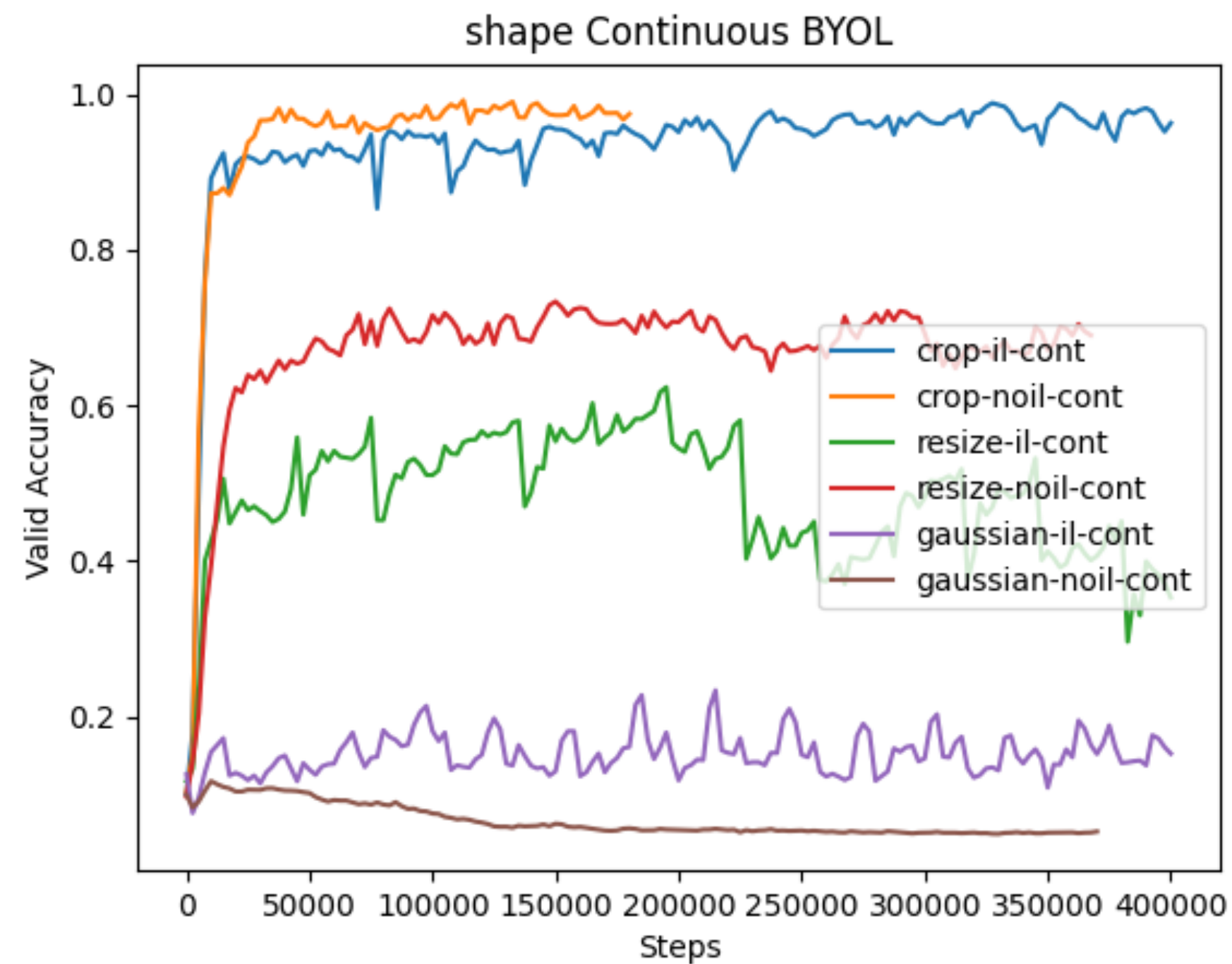
ψ^{t+1} randomly initialized

$t \leftarrow t + 1$

$N \leftarrow N - 1$

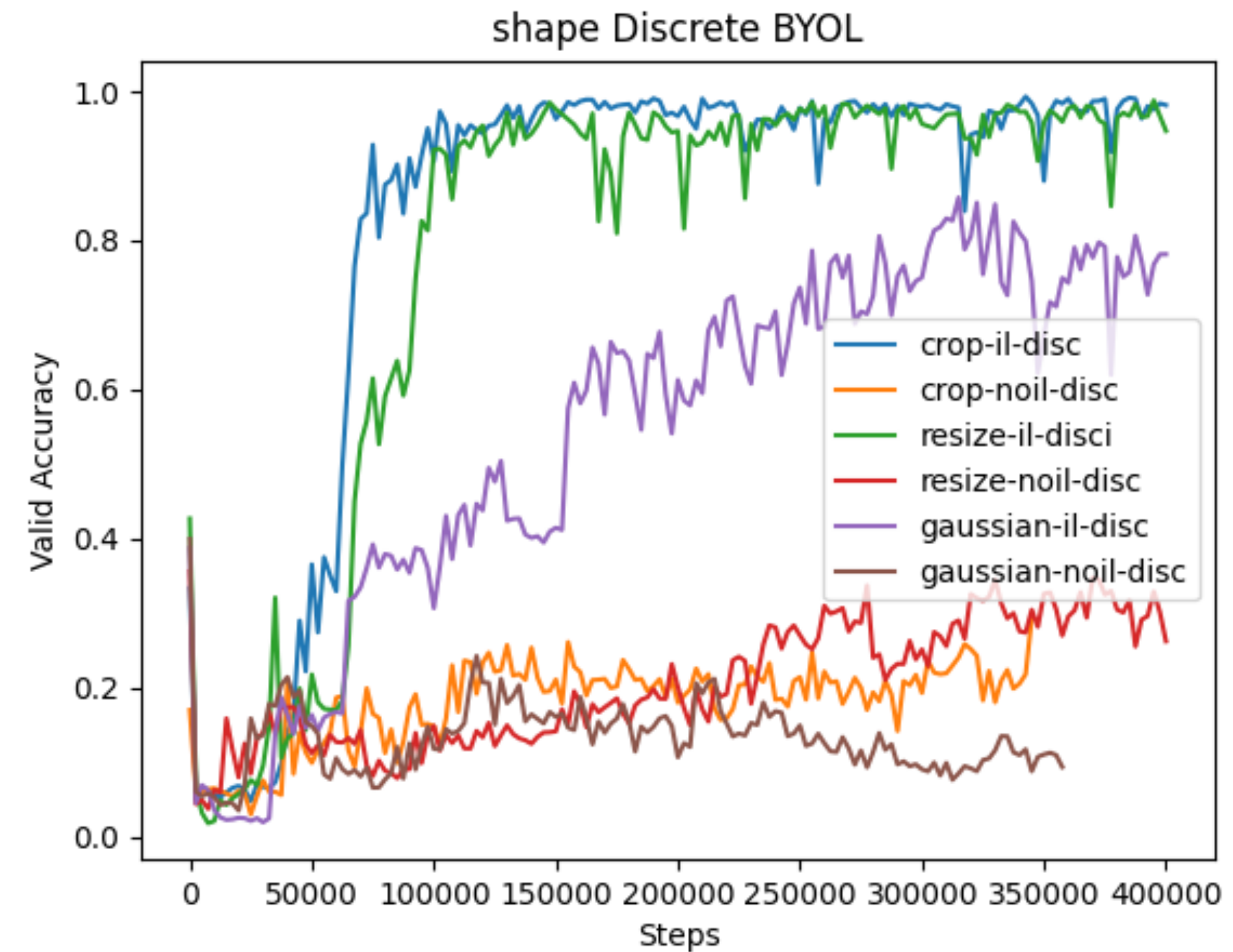
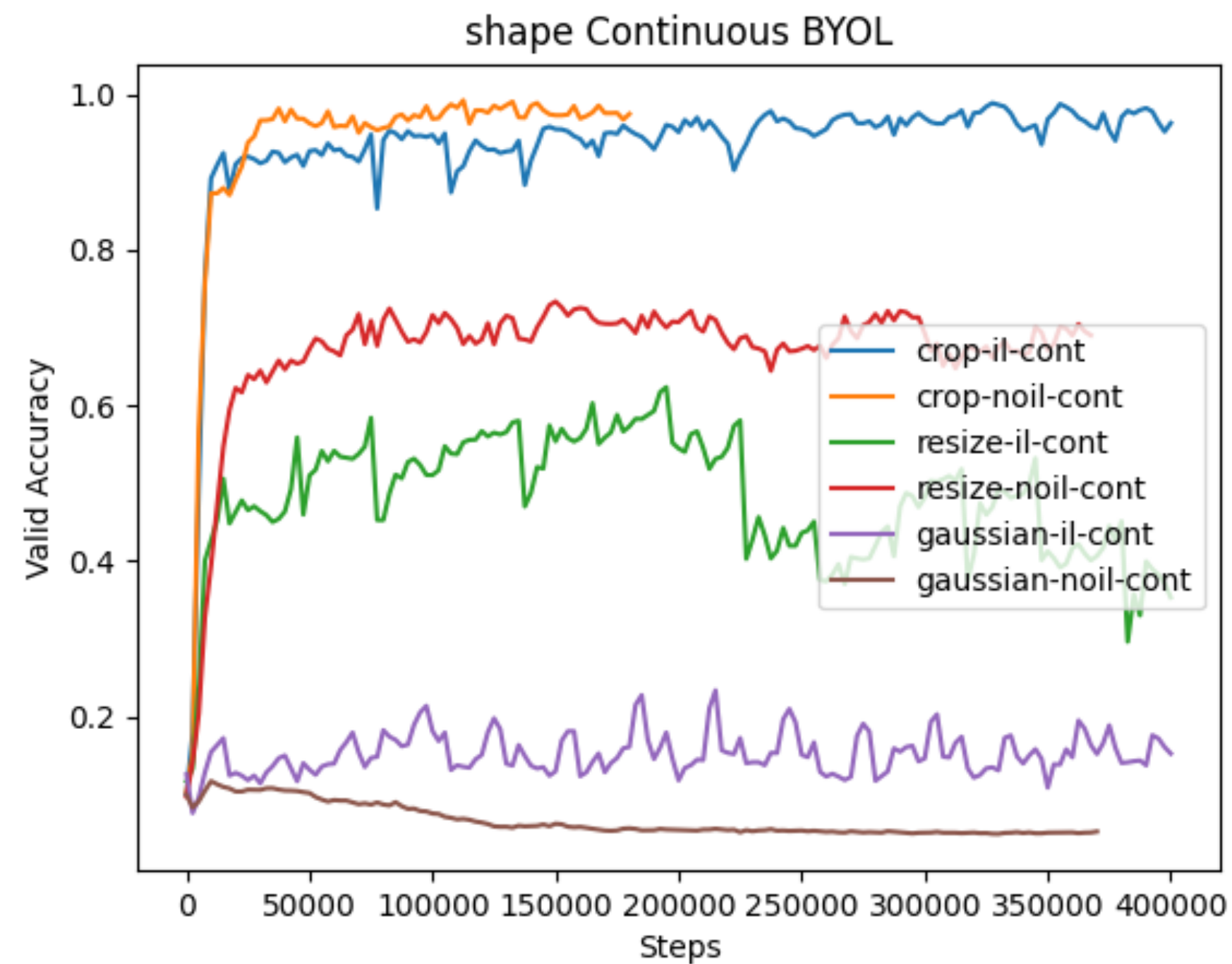
end while

Iterated learning for self-supervised learning



Comparing continuous and discrete bottleneck on the systematic generalization task of predicting the shape of dSprites for $K=2$.

Iterated learning for self-supervised learning



Comparing continuous and discrete bottleneck on the systematic generalization task of predicting the shape of dSprites for $K=2$.

Conclusion



+ colleagues