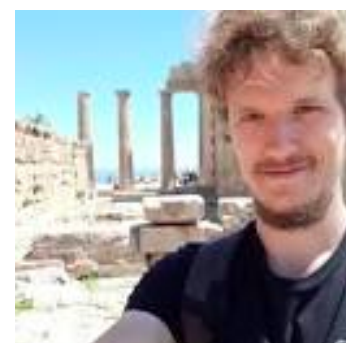


# Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture

June 1<sup>st</sup> 2023

Mido Assran



Quentin  
Duval



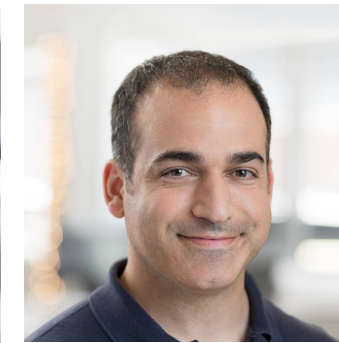
Ishan  
Misra



Piotr  
Bojanowski



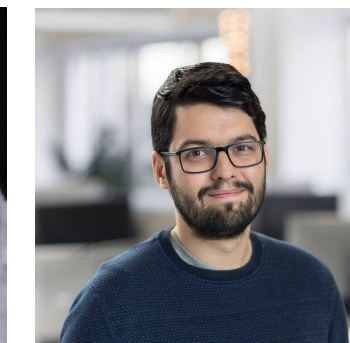
Pascal  
Vincent



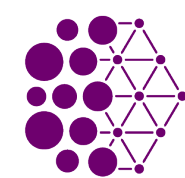
Mike  
Rabbat



Yann  
LeCun



Nicolas  
Ballas



Mila

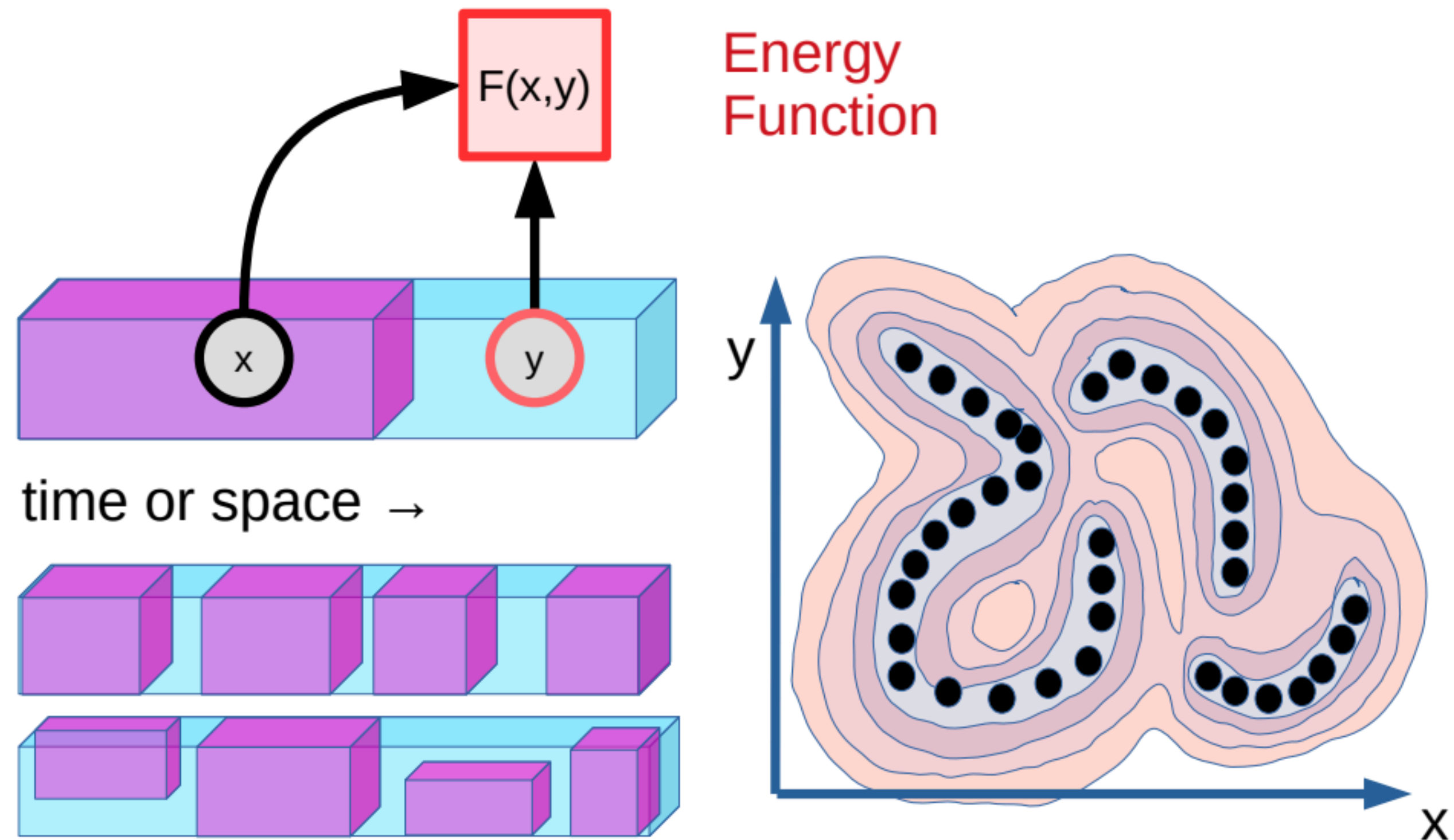


# Self-Supervised Learning

Towards General Representations

## Approach

Learn to represent data by capturing mutual dependencies between inputs

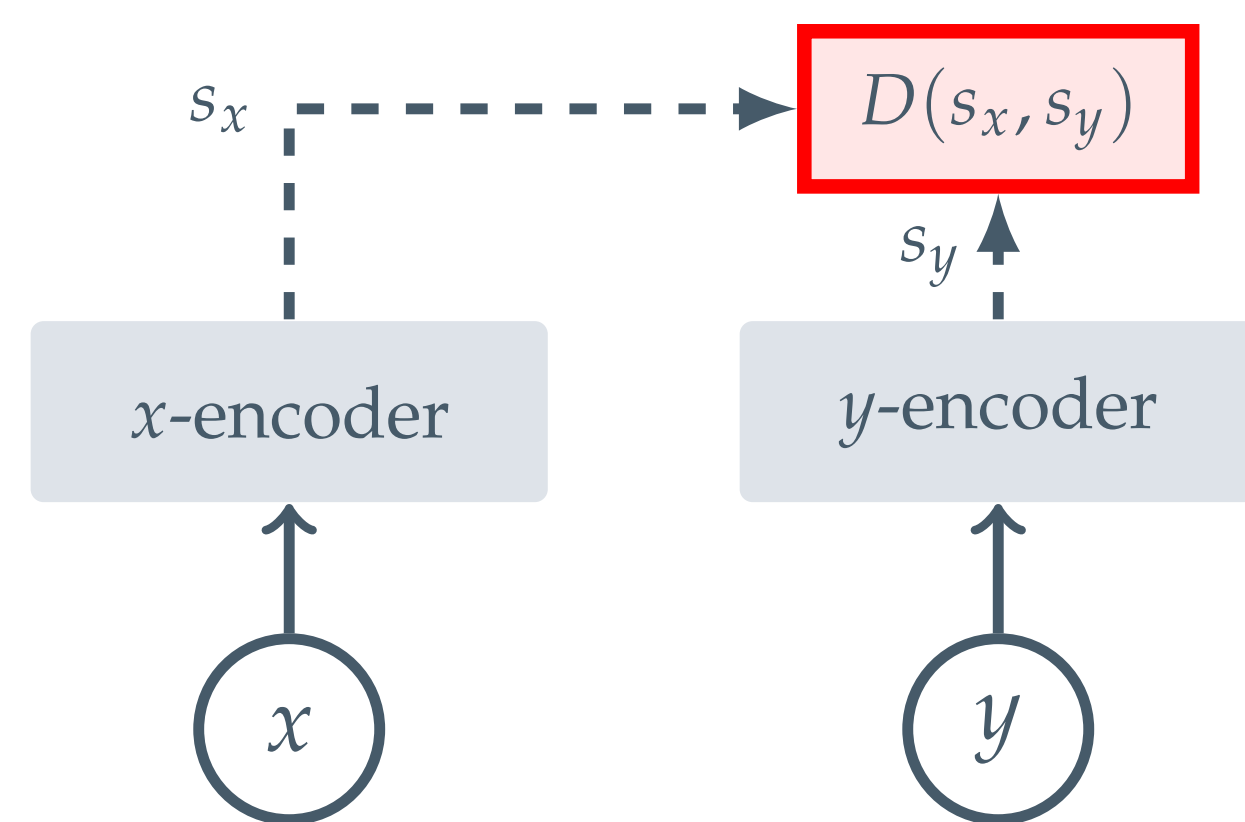


Yann LeCun  
 A Path Towards Autonomous Machine Intelligence  
 OpenReview, 2022.

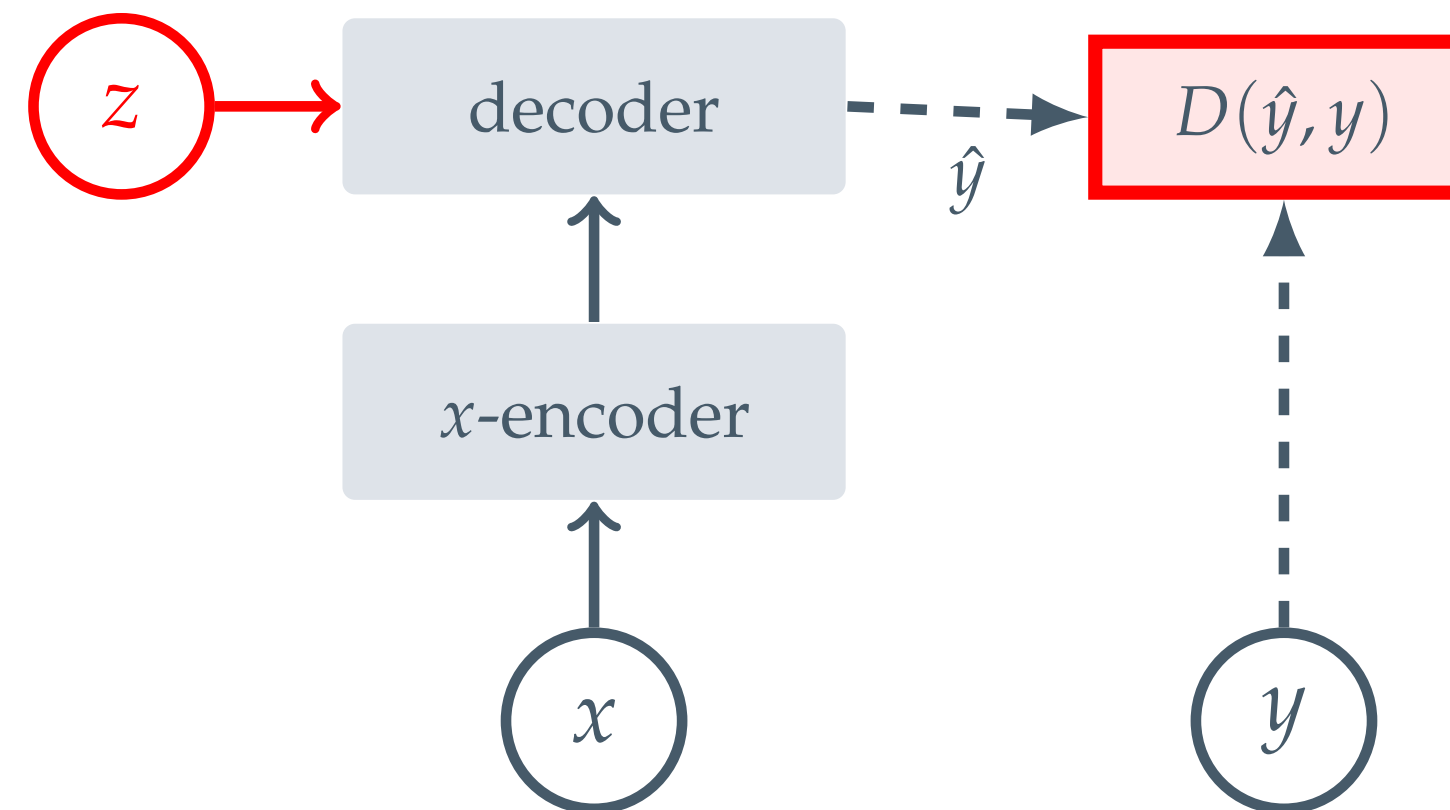
# What do we have so far?

## Common Approaches for Visual Representation Learning

Learn representations by capturing mutual dependencies between inputs...



Joint-Embedding Architecture

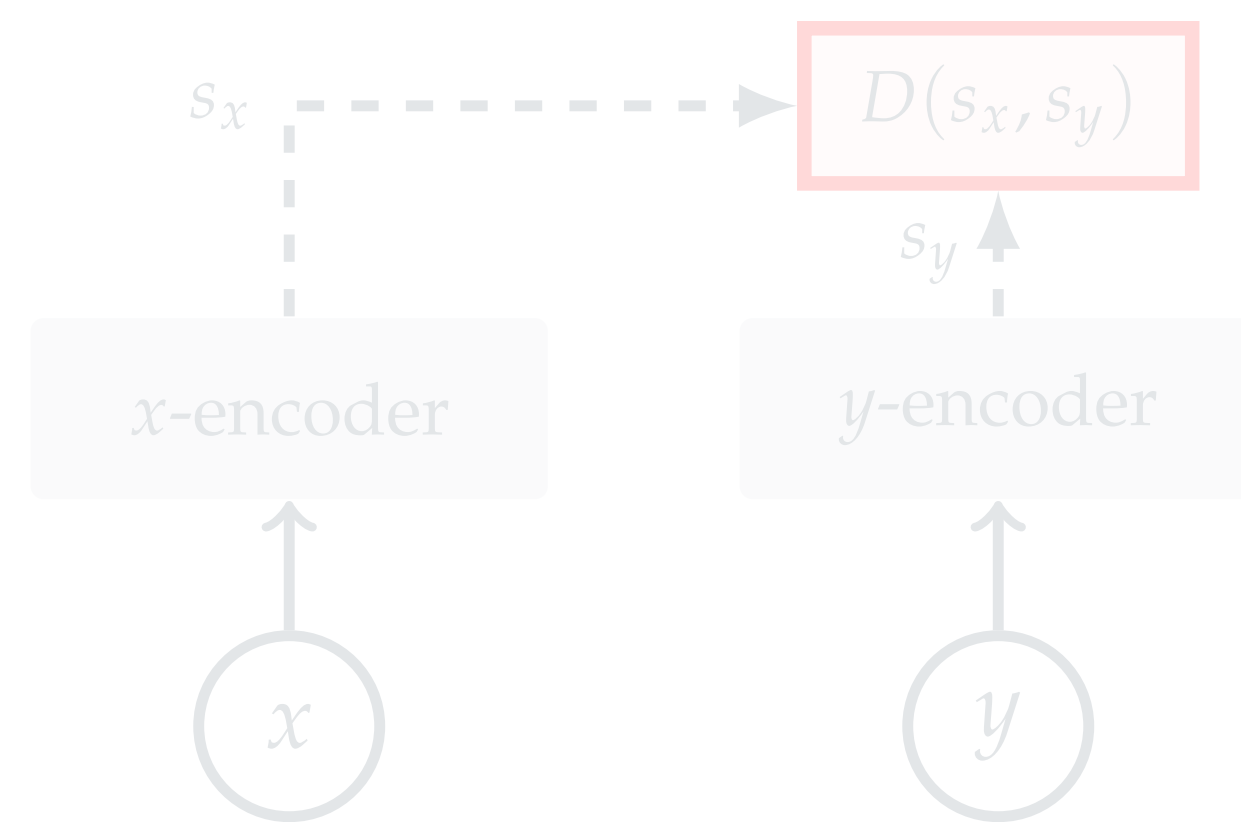


Generative Architecture

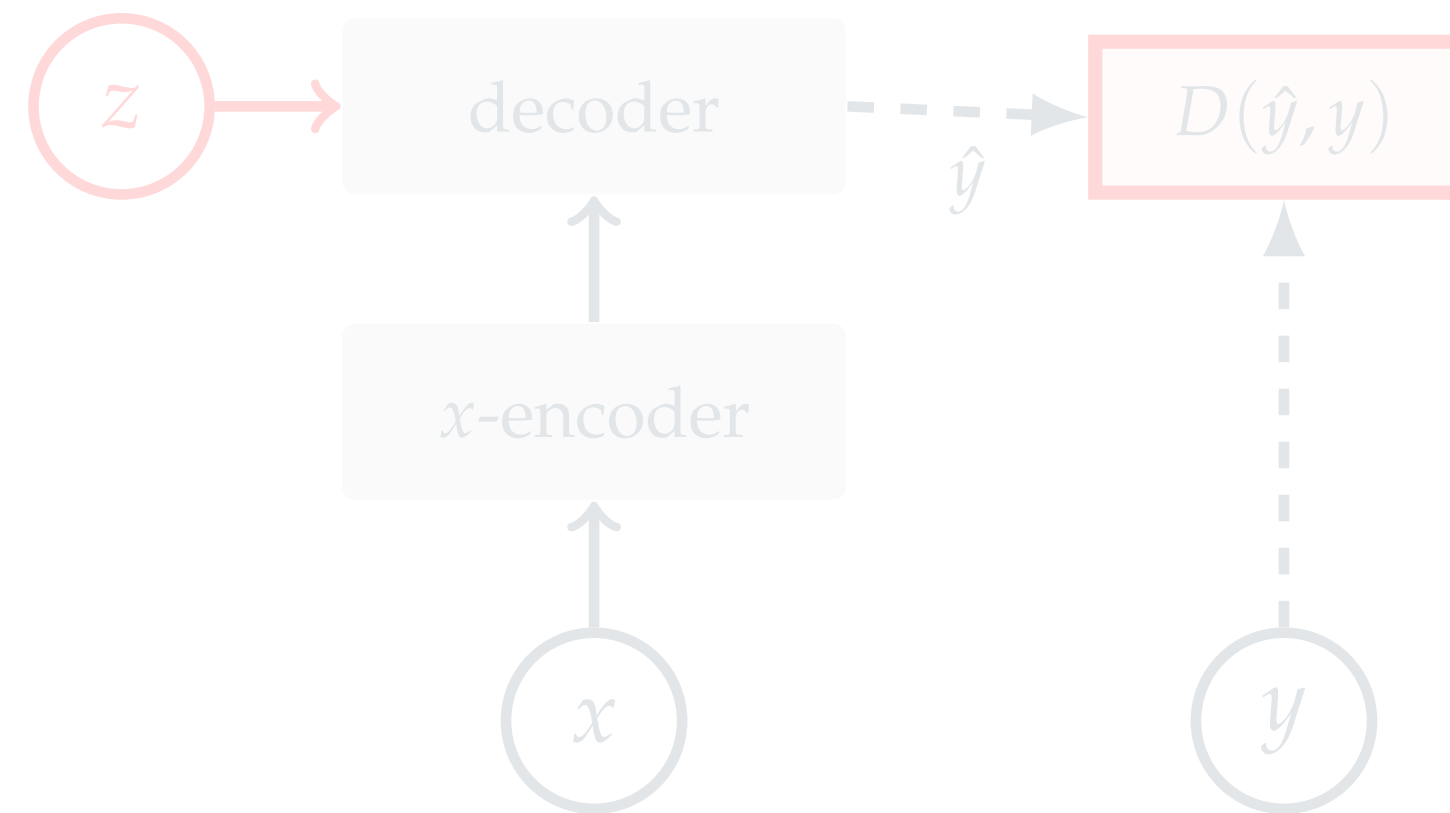
# What do we have so far?

## Common Approaches for Visual Representation Learning

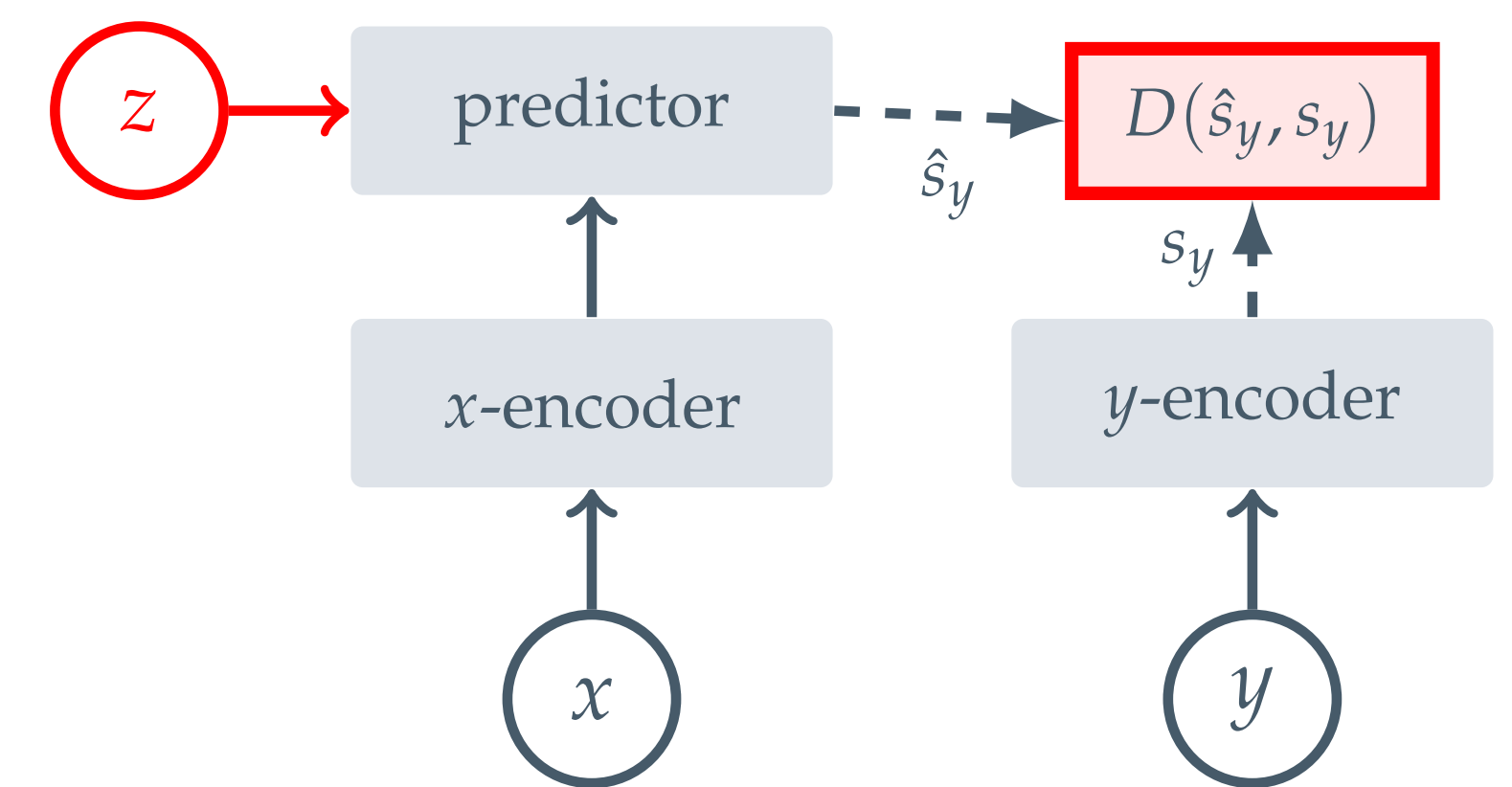
Learn representations by capturing mutual dependencies between inputs...



Joint-Embedding Architecture



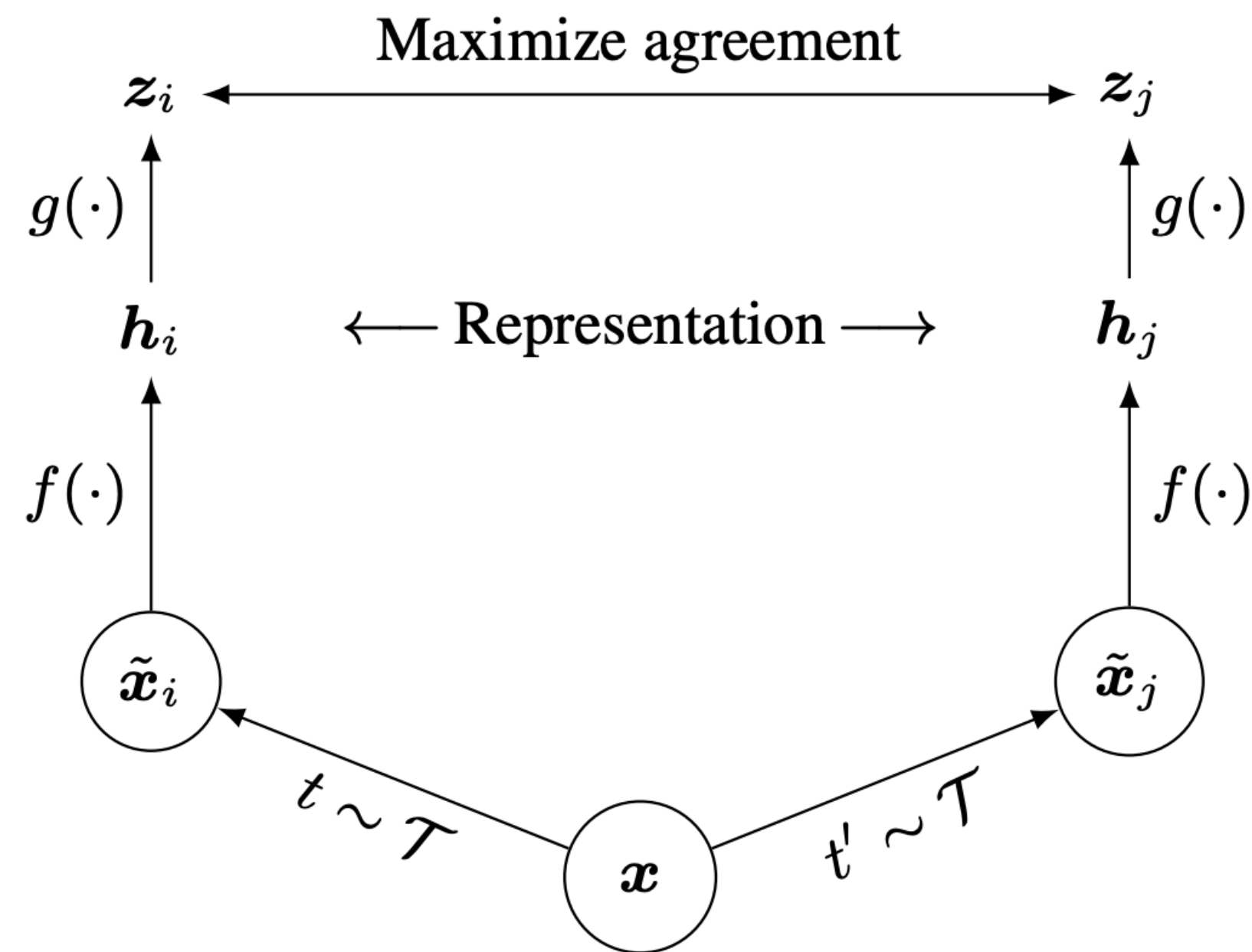
Generative Architecture



Joint-Embedding Predictive Architecture

# What do we have so far?

## Common Approaches for Visual Representation Learning



### Joint-Embedding Architecture

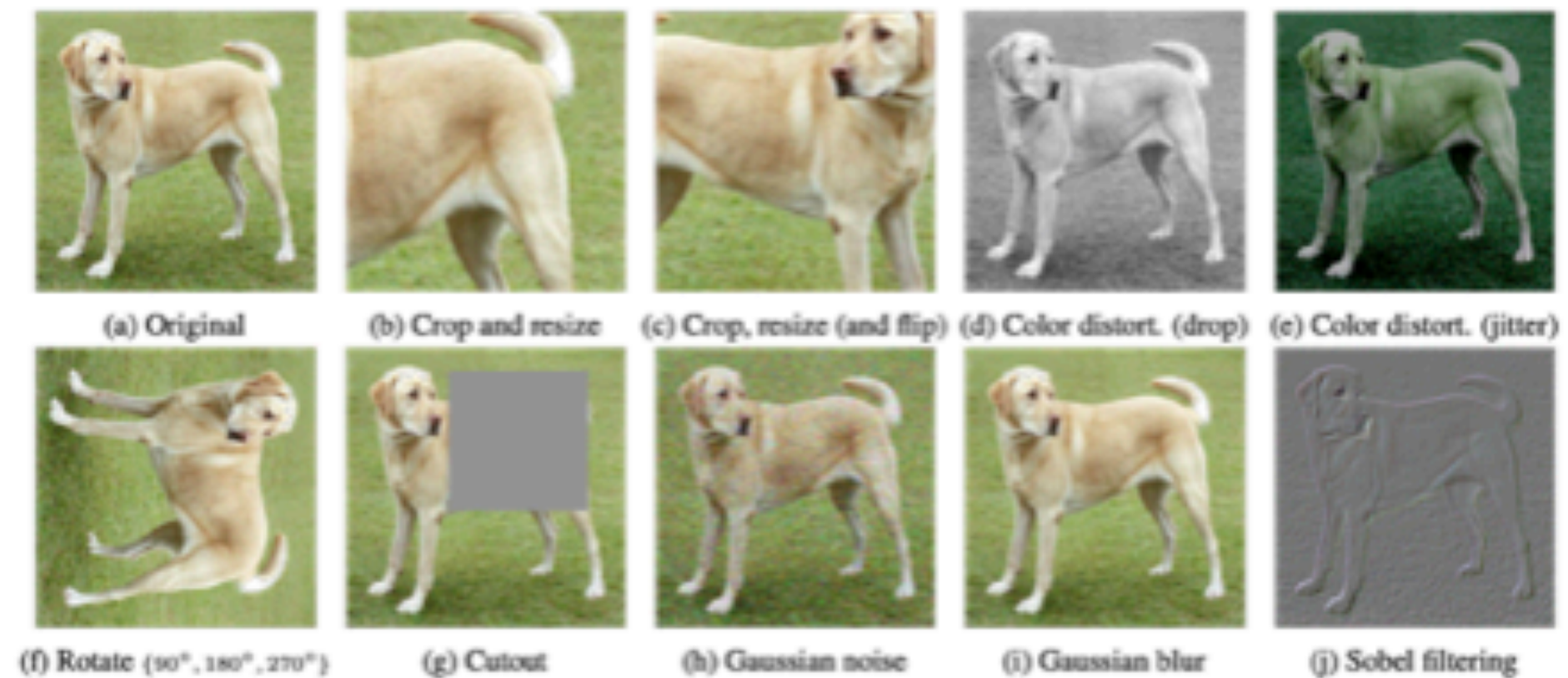


Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we only test these operators in ablation, the augmentation policy used to train our models only includes random crop (with flip and resize), color distortion, and Gaussian blur. (Original image cc-by: Von.grzanka)

# Self-Supervised Methods

## Canonical Joint-Embedding Architecture

### Limitations:

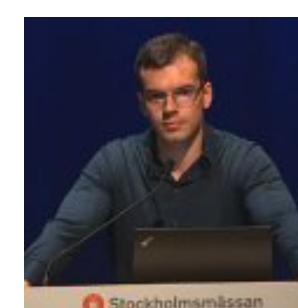
Semantic level of representations also depends on certain assumptions...

Published as a conference paper at ICLR 2023

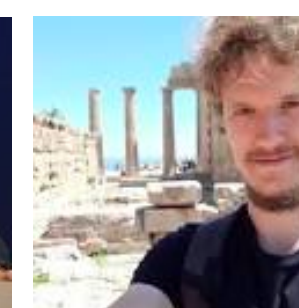
## THE HIDDEN UNIFORM CLUSTER PRIOR IN SELF-SUPERVISED LEARNING

**Mahmoud Assran<sup>\*,1,2,3</sup>, Randall Balestriero<sup>1</sup>, Quentin Duval<sup>1</sup>, Florian Bordes<sup>1,3,4</sup>,  
 Ishan Misra<sup>1</sup>, Piotr Bojanowski<sup>1</sup>, Pascal Vincent<sup>1,3,4</sup>, Michael Rabbat<sup>1,3</sup>, Nicolas Ballas<sup>1</sup>**

<sup>1</sup>Meta AI (FAIR)   <sup>2</sup>McGill University, ECE   <sup>3</sup>Mila, Quebec AI Institute   <sup>4</sup>Universite de Montreal, DIRO



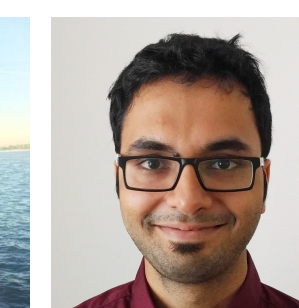
Randall Balestriero



Quentin Duval



Florian Bordes



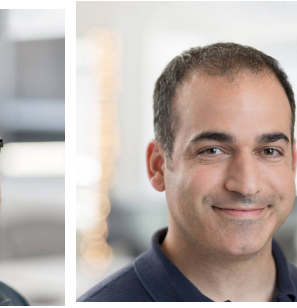
Ishan Misra



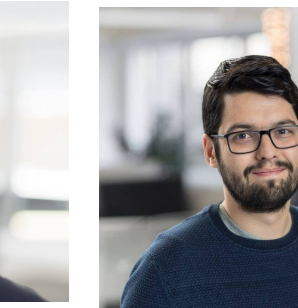
Piotr Bojanowski



Pascal Vincent



Mike Rabbat

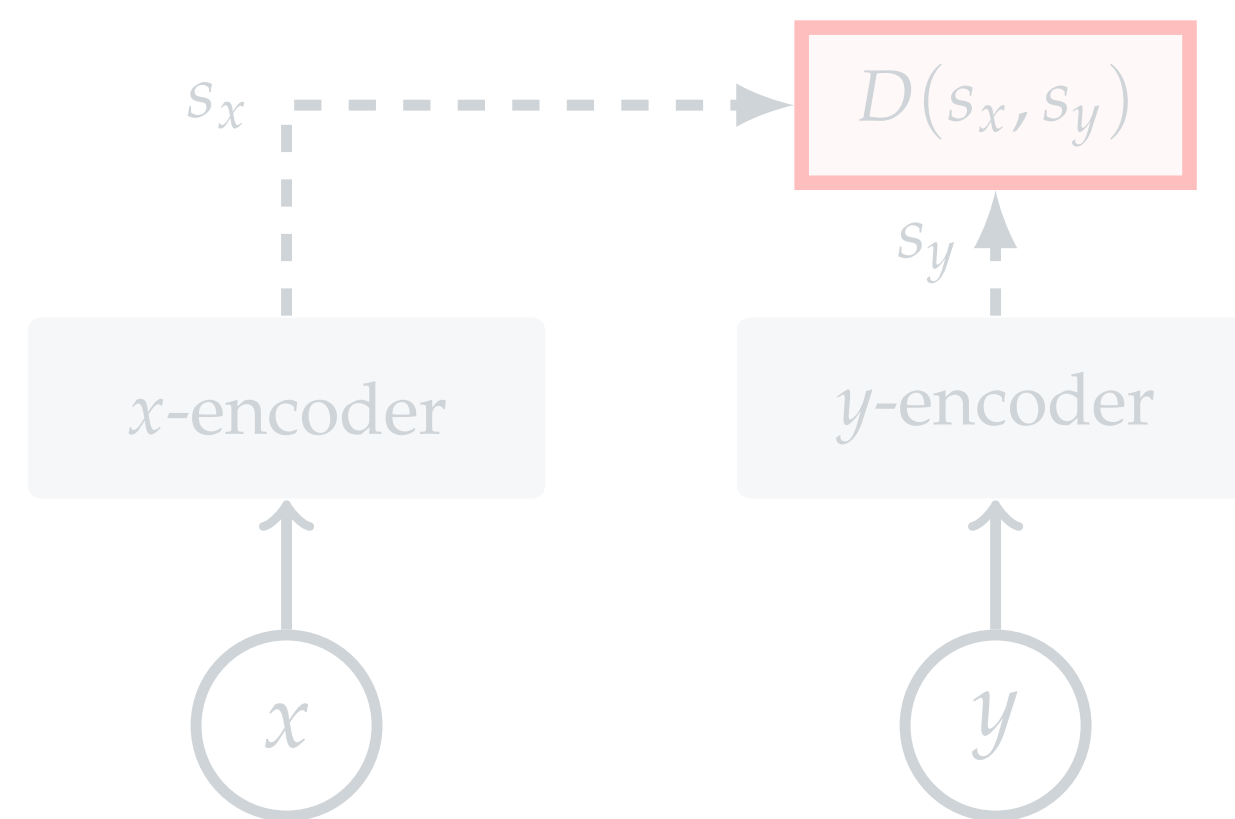


Nicolas Ballas

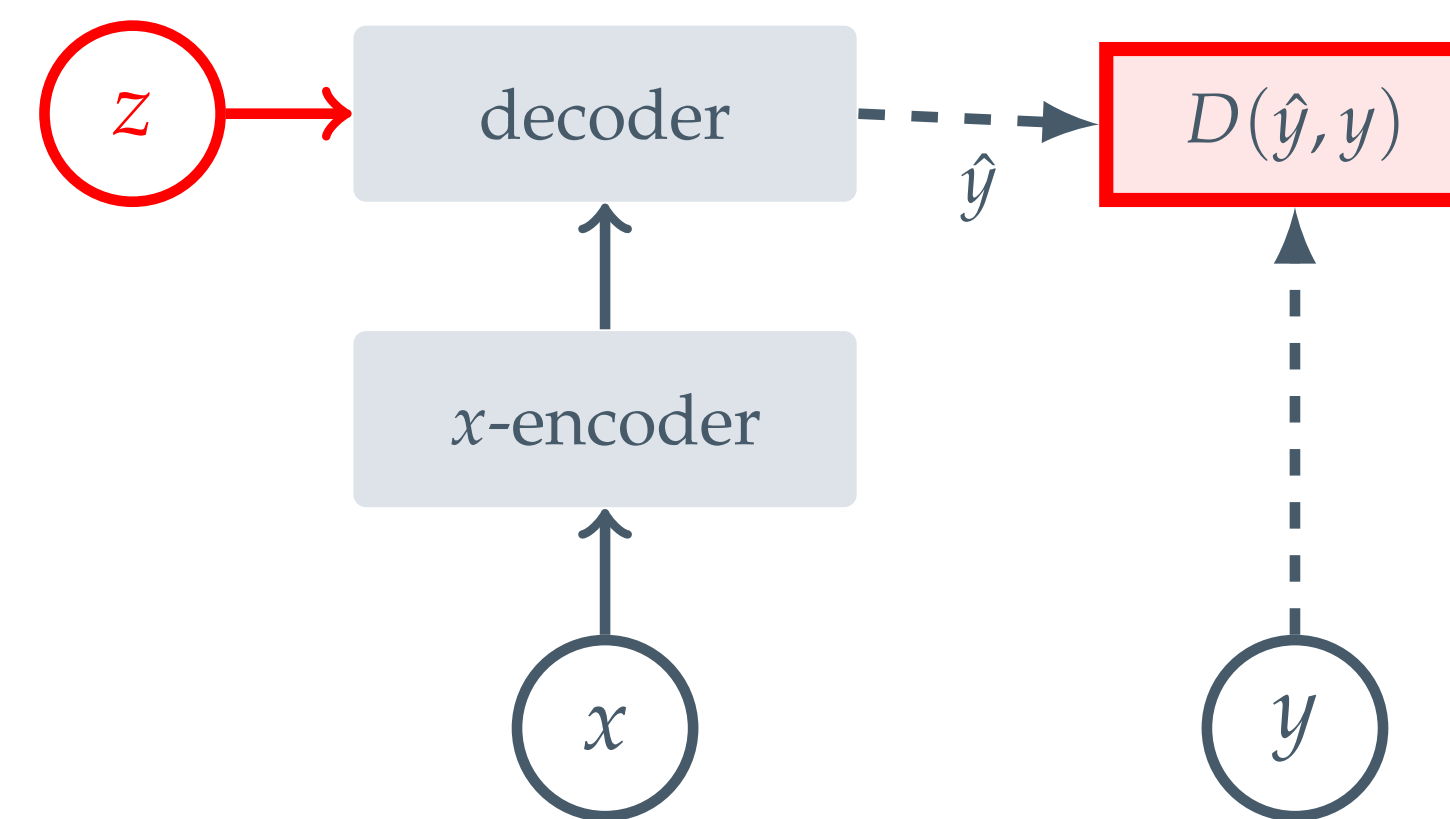
# What do we have so far?

## Common Approaches for Visual Representation Learning

Learn representations by capturing mutual dependencies between inputs...



Joint-Embedding Architecture



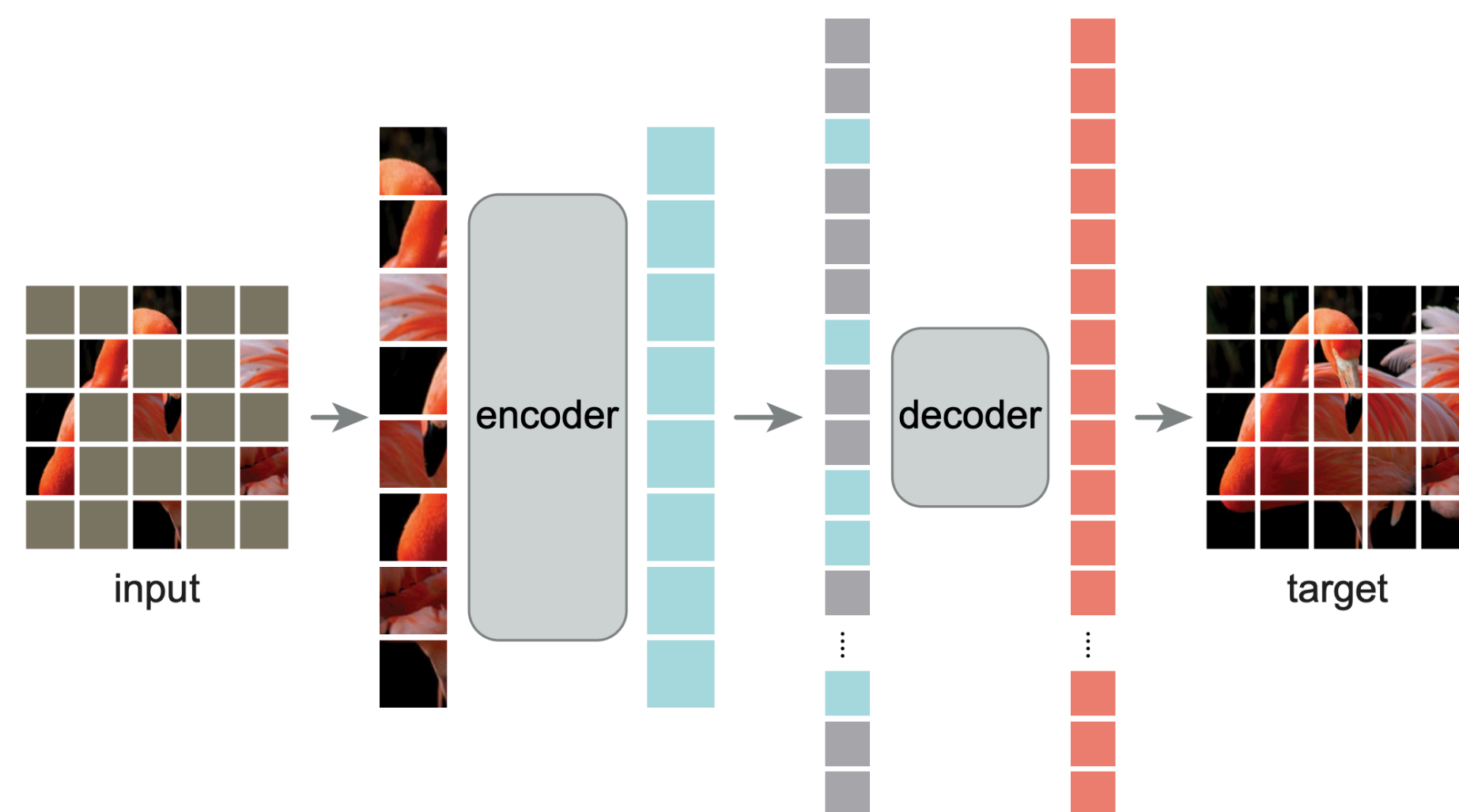
Generative Architecture

# What do we have so far?

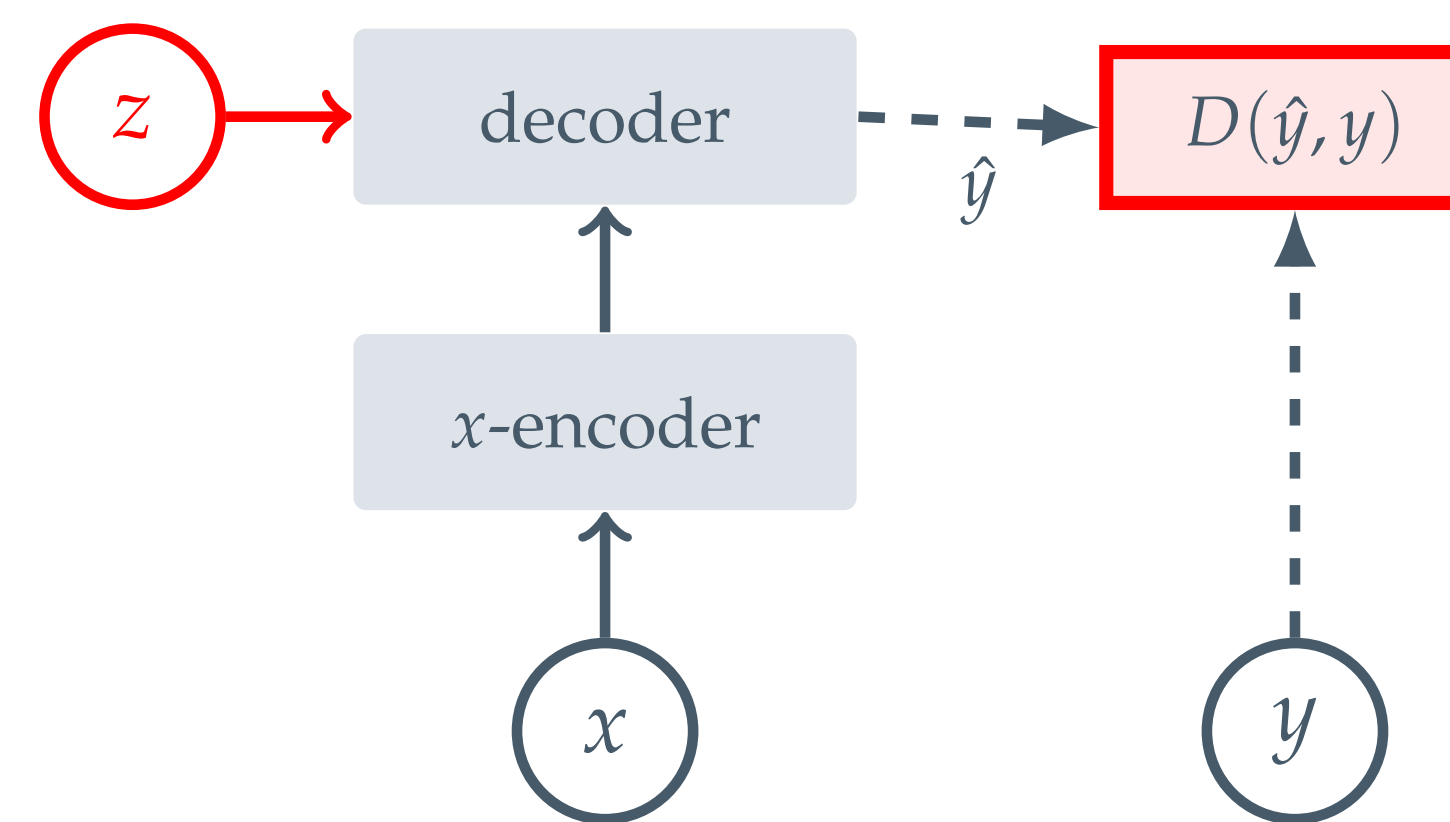
## Common Approaches for Visual Representation Learning

Learn representations by capturing mutual dependencies between inputs...

### Mask Denoising Architectures



He et al.,  
Masked Autoencoders Are Scalable Vision Learners  
 arXiv, 2021.



### Generative Architecture



# Self-Supervised Learning

## Generative Architectures

Generative architectures tend to learn representations of a lower semantic level...

Published as a conference paper at ICLR 2023

---

## WHAT DO SELF-SUPERVISED VISION TRANSFORMERS LEARN?

**Namuk Park<sup>1\*</sup> Wonjae Kim<sup>2</sup> Byeongho Heo<sup>2</sup> Taekyung Kim<sup>2</sup> Sangdoon Yun<sup>2</sup>**

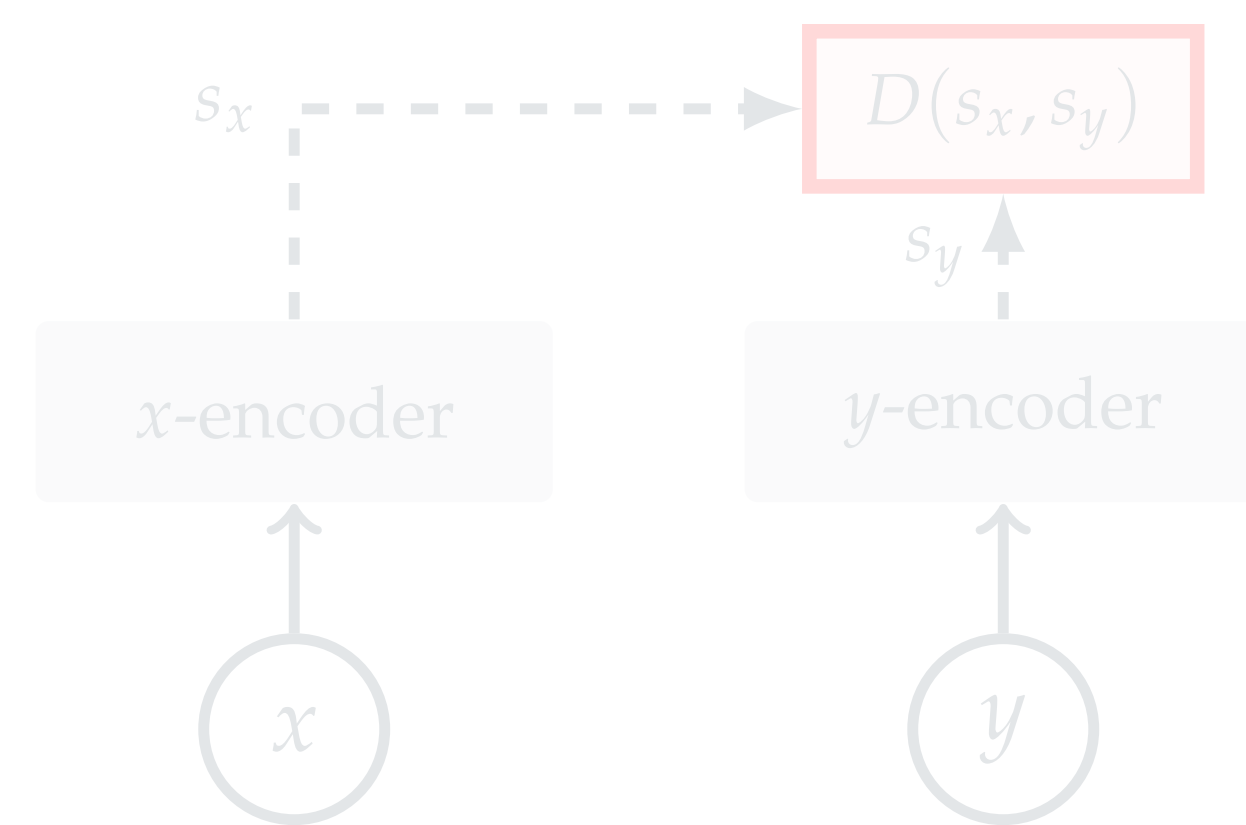
<sup>1</sup>Prescient Design, Genentech <sup>2</sup>NAVER AI Lab

park.namuk@gene.com {wonjae.kim,bh.heo,taekyung.k,sangdoon.yun}@navercorp.com

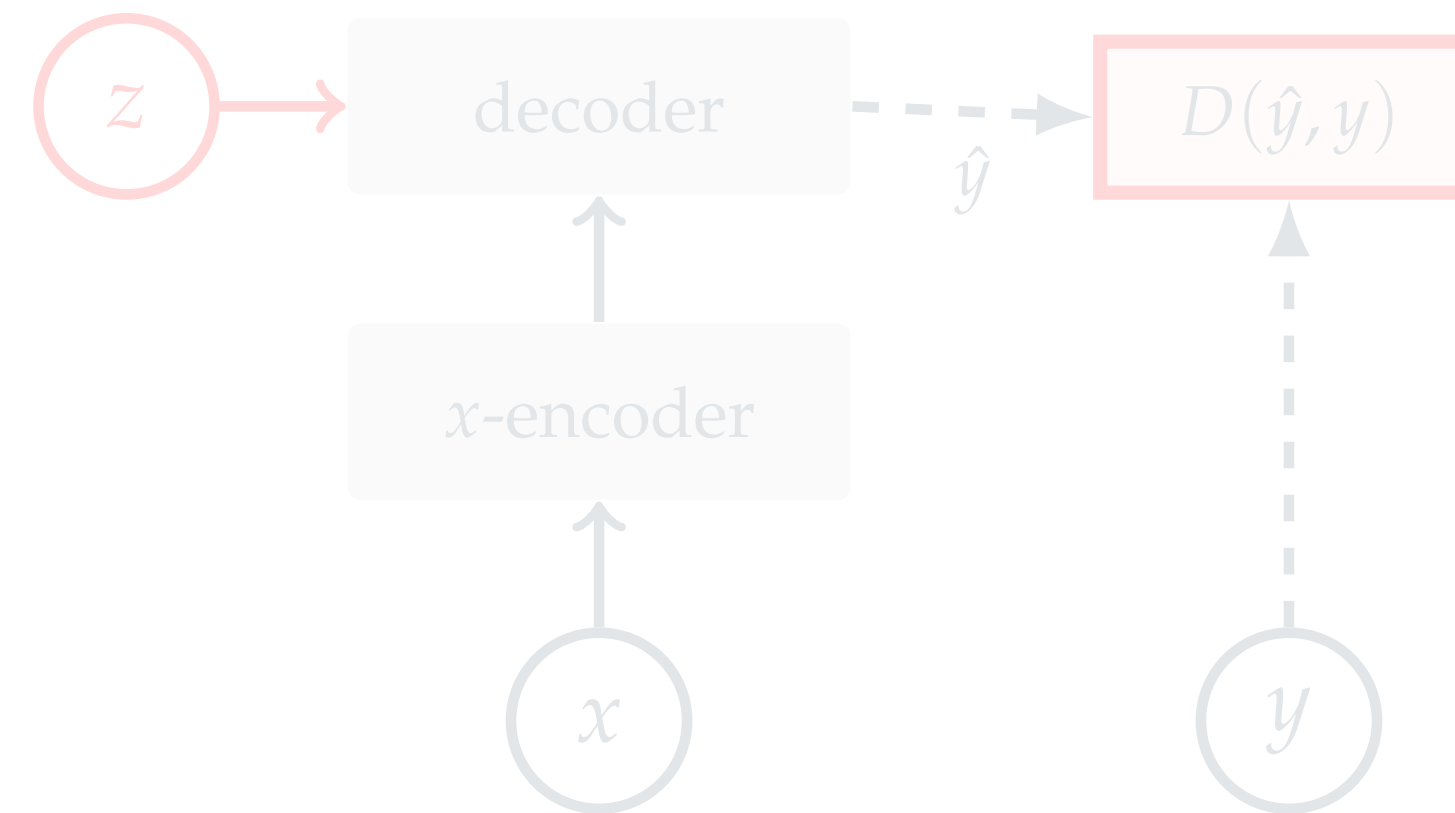
# What do we have so far?

## Common Approaches for Visual Representation Learning

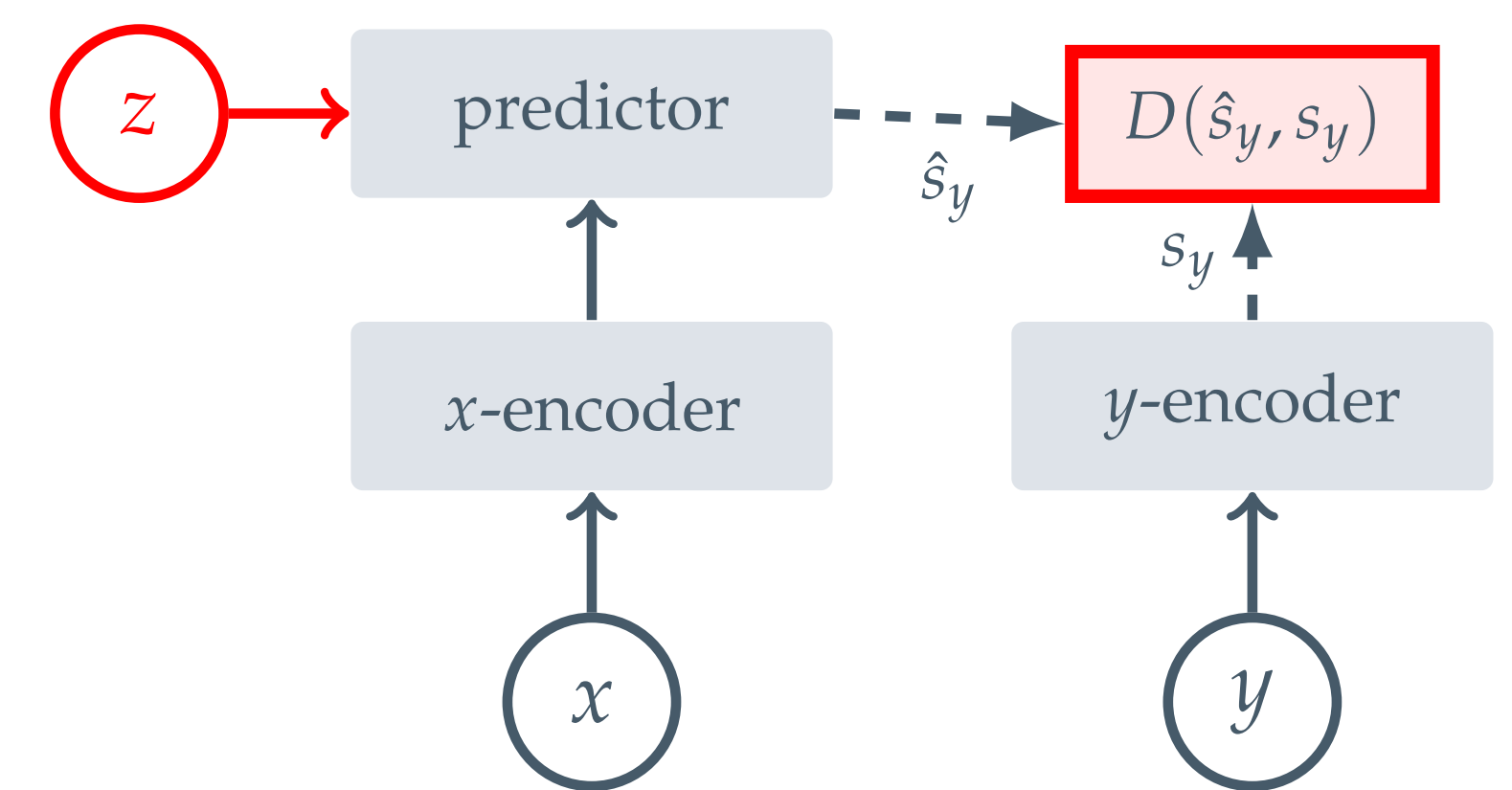
Learn representations by capturing mutual dependencies between inputs...



Joint-Embedding Architecture



Generative Architecture



Joint-Embedding Predictive Architecture

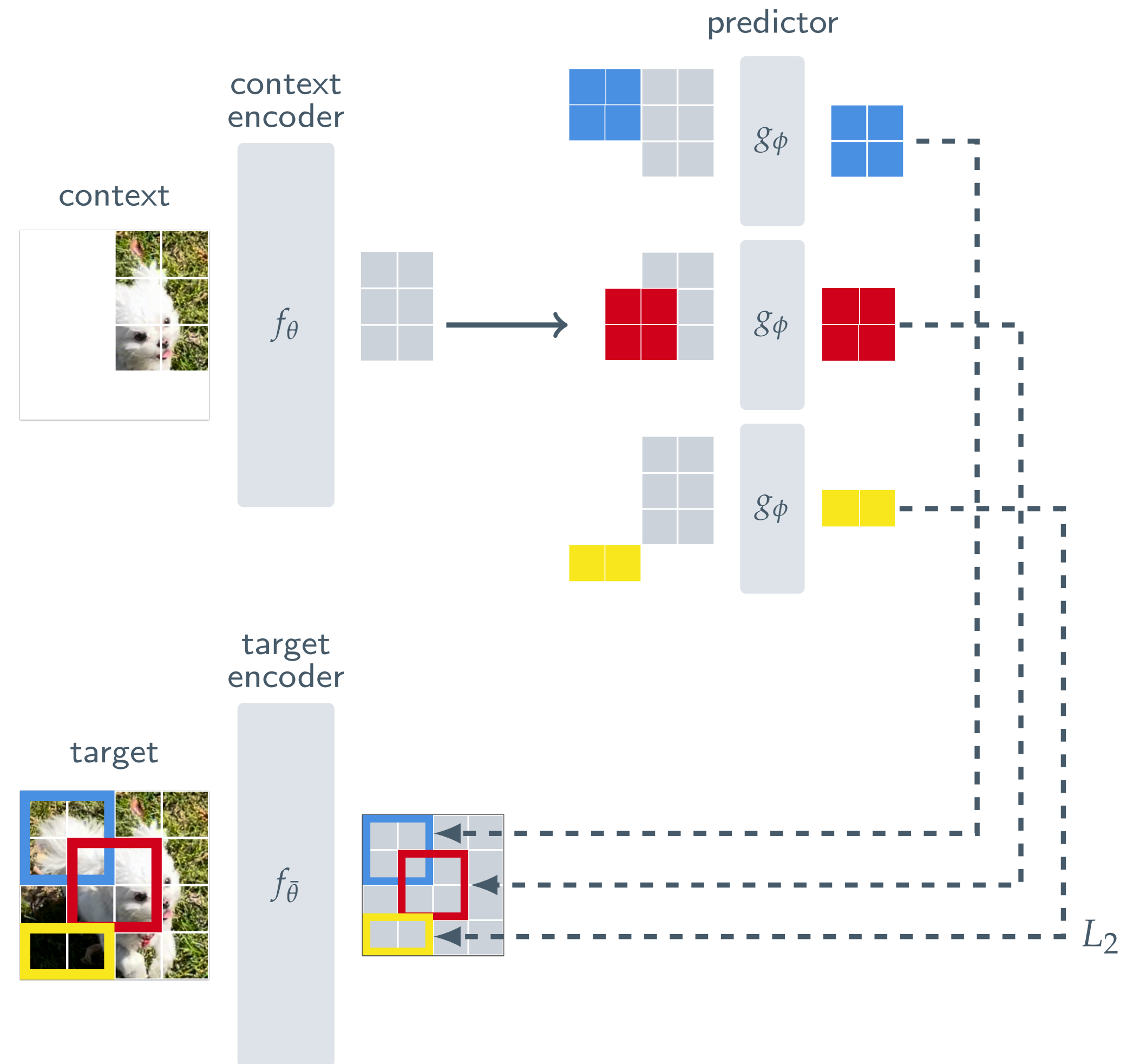
# Image JEPA

## Towards More General Representations



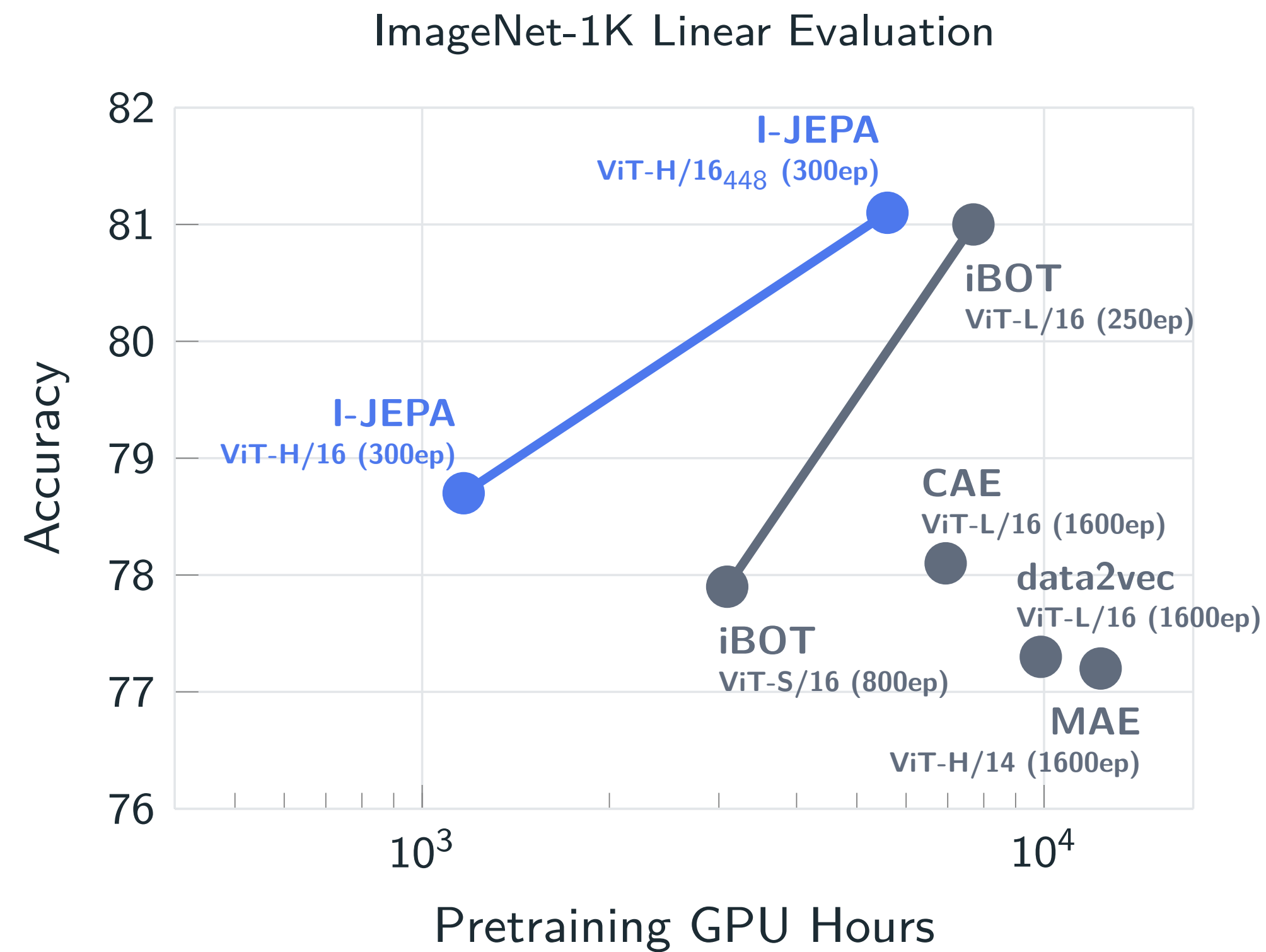
From a single context block, predict representations of various target blocks...

... no hand-crafted data augmentations!



# Image JEPA

## Towards More General Representations



Very efficient...

Training ViT-Huge/16 with I-JEPA is faster than training ViT-Small/16 with iBOT!

Converges faster than generative methods, which require many epochs of pre-training (MAE)

# Image JEPA

## Towards More General Representations

Targets	Arch.	Epochs	Top-1
Target-Encoder Output	ViT-L/16	500	<b>66.9</b>
Pixels	ViT-L/16	800	40.7

Table 7. **Ablating targets.** Linear evaluation on ImageNet-1K using only 1% of the available labels; ablating the effect of the prediction targets during I-JEPA pretraining. To ensure convergence when predicting in pixel space, we train the model for more epochs. The semantic level of the I-JEPA representations degrades significantly when the loss is applied in pixel space, rather than representation space, highlighting the importance of the target-encoder during pretraining.

I-JEPA is non-generative...

Same method in pixel space performs much worse on semantic classification tasks...

# Image JEPA

Towards More General Representations



# Image JEPA

## Towards More General Representations

Method	Arch.	CIFAR100	Places205	iNat18
<i>Methods without view data augmentations</i>				
data2vec [7]	ViT-L/16	81.6	54.6	28.1
MAE [34]	ViT-H/14	77.3	55.0	32.9
I-JEPA	ViT-H/14	<b>87.5</b>	<b>58.4</b>	<b>47.6</b>
<i>Methods using extra view data augmentations</i>				
DINO [17]	ViT-B/8	84.9	57.9	55.9
iBOT [74]	ViT-L/16	<b>88.3</b>	<b>60.4</b>	<b>57.3</b>

I-JEPA captures global semantics...

- Outperforms generative methods
- Closes gap with view-invariance methods

Linear transfer to semantic image-level visual tasks

# Image JEPA

## Towards More General Representations

Method	Arch.	Clevr/Count	Clevr/Dist
<i>Methods without view data augmentations</i>			
data2vec [7]	ViT-L/16	85.3	71.3
MAE [34]	ViT-H/14	<b>90.5</b>	<b>72.4</b>
I-JEPA	ViT-H/14	<b>86.7</b>	<b>72.4</b>
<i>Methods using extra data augmentations</i>			
DINO [17]	ViT-B/8	86.6	53.4
iBOT [74]	ViT-L/16	85.7	62.8

### Linear Transfer to Low-Level Visual Tasks

I-JEPA also captures local information...

- Outperforms view-invariance methods in low-level tasks (e.g., depth prediction)
- Comparable with generative methods



# Image JEPA

## Towards More General Representations

Method	Arch.	Epochs	Top-1
<i>Methods without view data augmentations</i>			
data2vec [7]	ViT-L/16	1600	73.3
MAE [34]	ViT-L/16	1600	67.1
	ViT-H/14	1600	71.5
I-JEPA	ViT-L/16	600	69.4
	ViT-H/14	300	73.3
	ViT-H/16 <sub>448</sub>	300	<b>77.3</b>
<i>Methods using extra view data augmentations</i>			
iBOT [74]	ViT-B/16	250	69.7
DINO [17]	ViT-B/8	300	70.0
SimCLR v2 [33]	RN151 (2×)	800	70.2
BYOL [33]	RN200 (2×)	800	71.2
MSN [3]	ViT-B/4	300	<b>75.7</b>

Scaling I-JEPA...

New SoTA for ImageNet semi-supervised eval...

ImageNet 1% Semi-Supervised Evaluation