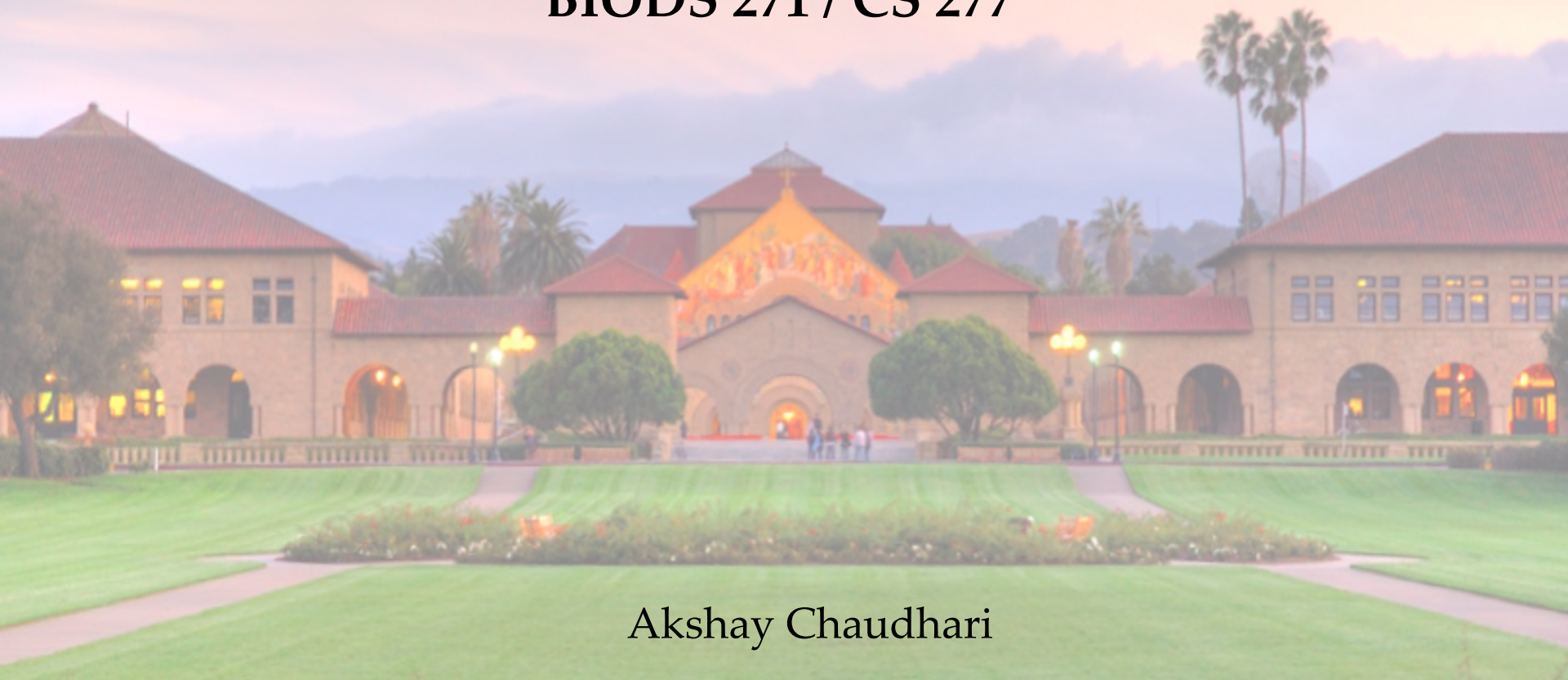


# Self-Supervised Learning for Vision

BIODS 271 / CS 277



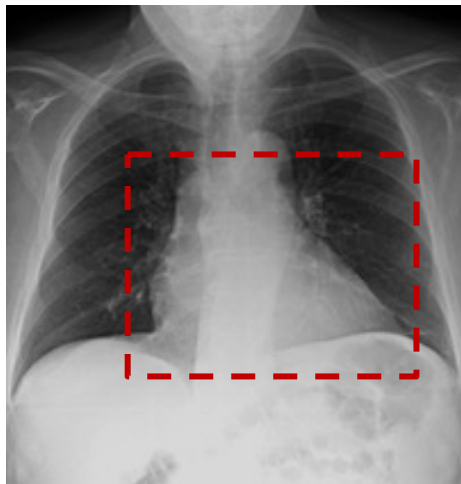
Akshay Chaudhari

# Data Enables Solving Medical Problems

Classification



Detection



Segmentation



Regression



✓: Cardiomegaly

Ejection  
Fraction: **49%**

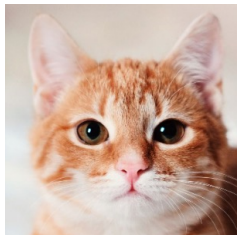
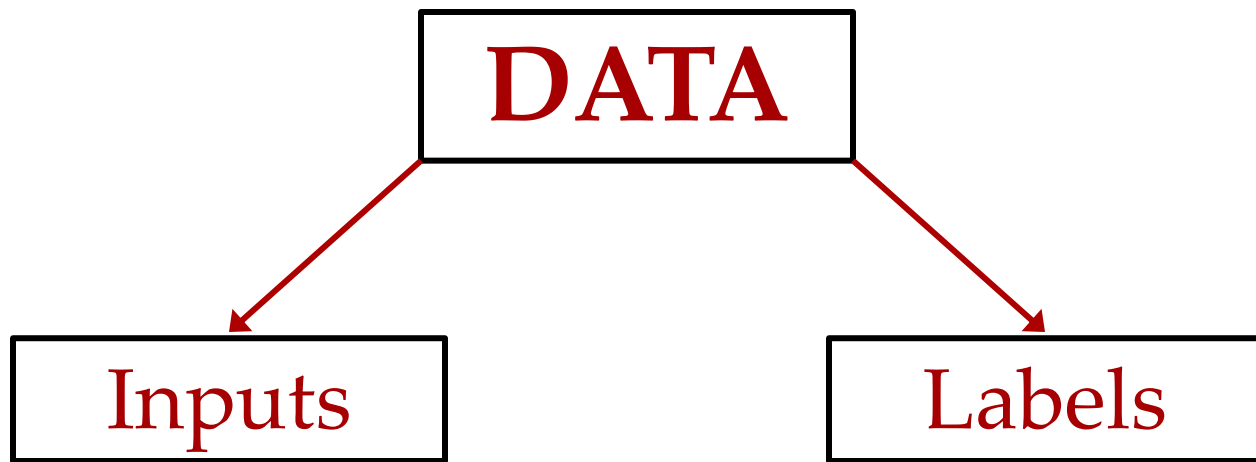
# Supervised Learning

IM = Image

L = Label

IM1	IM2	IM3	+	L1	L2	L3	= <b>Model!</b>
IM4	IM5	IM6		L4	L5	L6	
IM7	IM8	IM9		L7	L8	L9	





and so on...

Dog  
Cat  
Hot Dog

...

# Label-Starved Learning

IM = Image

L = Label

IM1	IM2	IM3
IM4	IM5	IM6
IM7	IM8	IM9

+

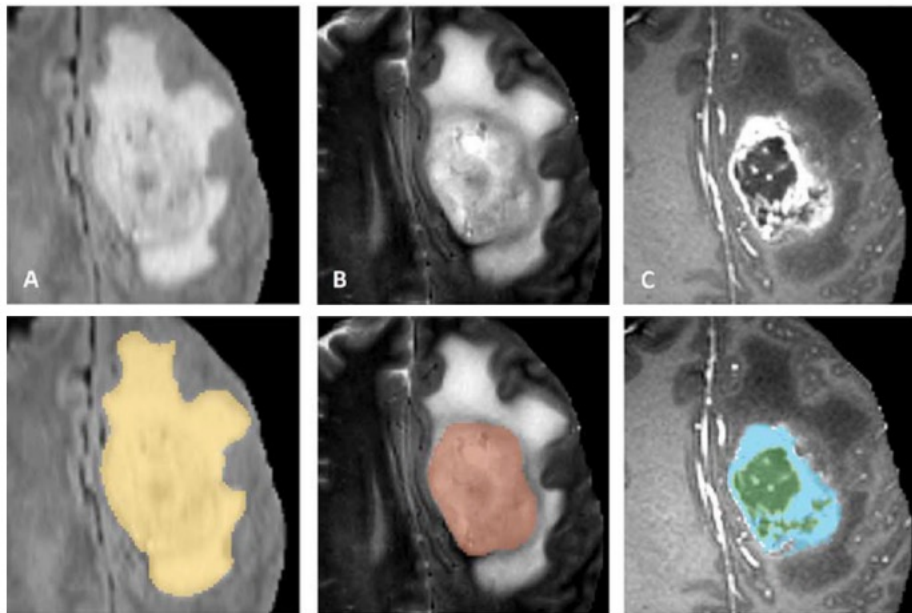
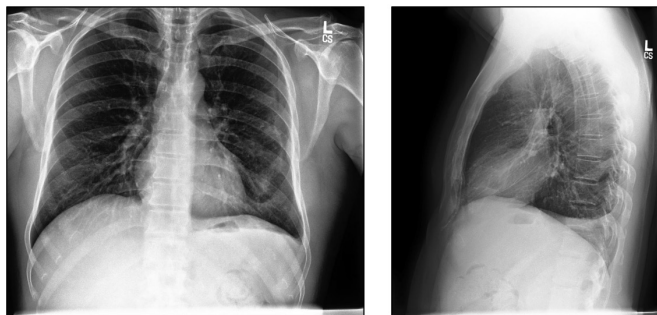
L1	L2	L3
L4	L5	L6
L7	L8	L9

=

**No  
Model!**



# Example Labels in Medical Imaging

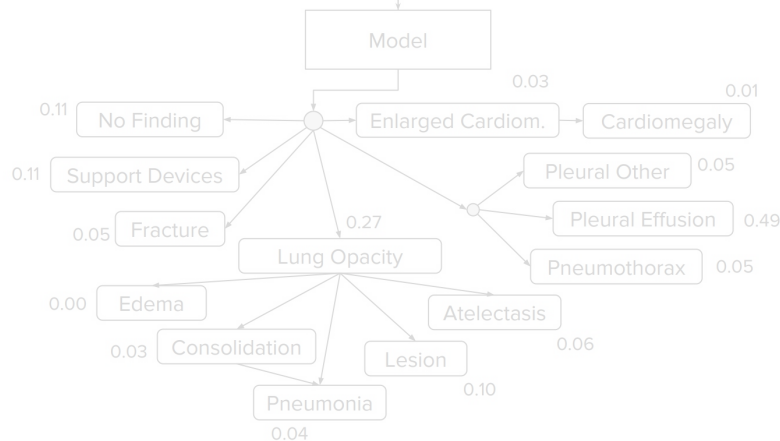
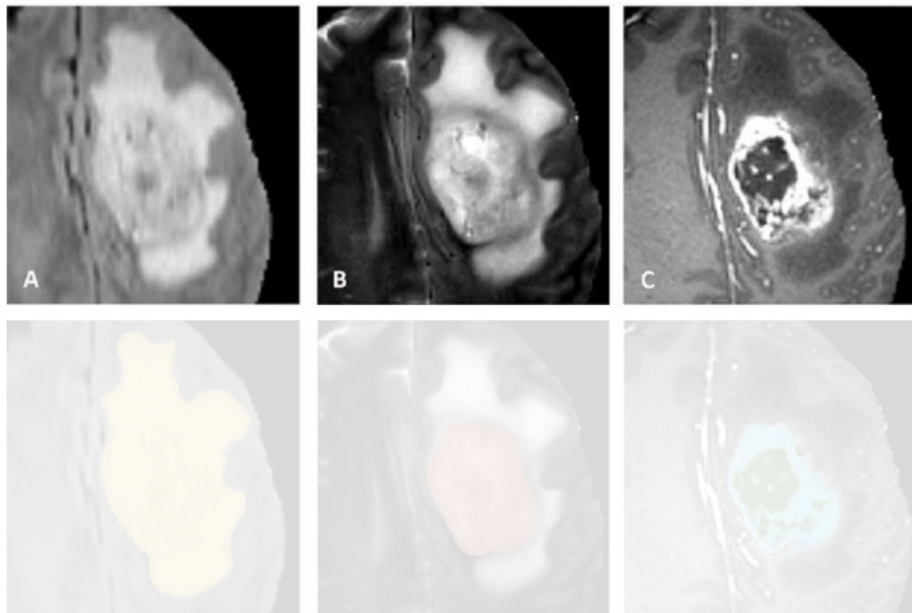
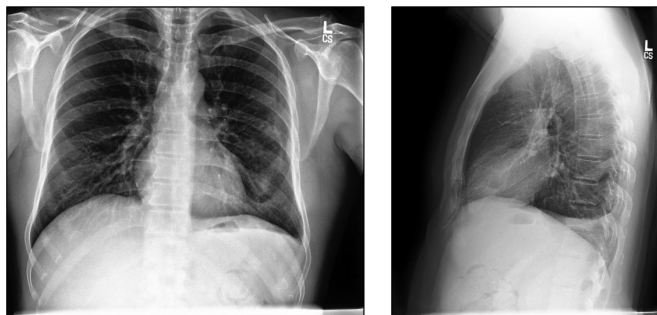


[1] Irwin et al. CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison. AAAI. 2019

[2] Menze et al. The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). IEEE TMI. 2015



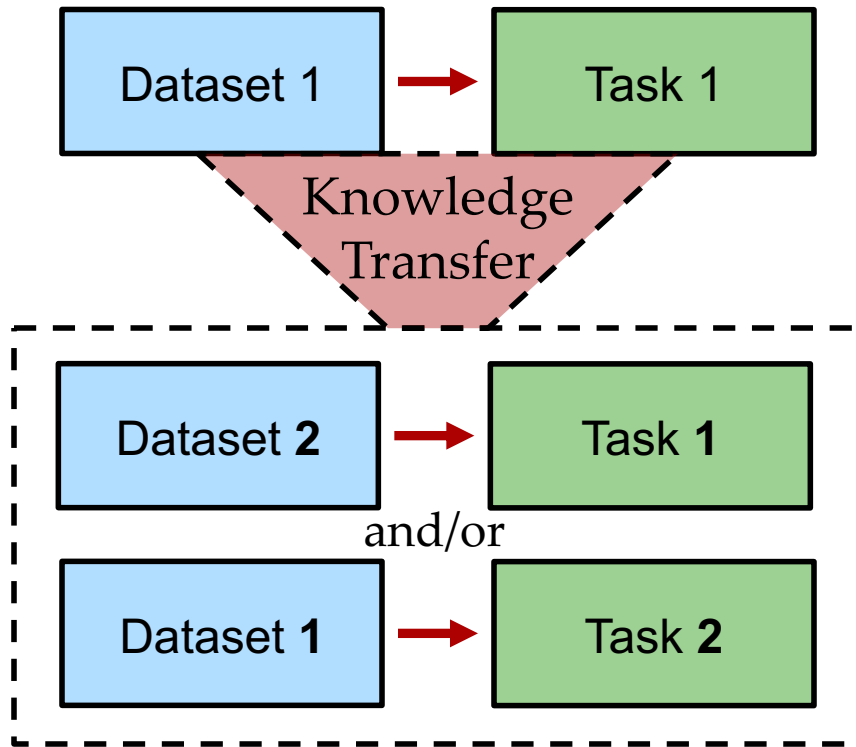
# Example Labels in Medical Imaging



[1] Irwin et al. CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison. AAAI. 2019

[2] Menze et al. The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). IEEE TMI. 2015

# Transfer Learning





# Beyond Supervised Learning

- Semi-Supervised Learning

Step 1:

IM1	IM2
IM4	IM5

+

L1	L2
L4	L5

= Model 1

---

Step 2:

		IM3
		IM6
IM7	IM8	IM9

Model 1

+

		PL3
		PL6
PL7	PL8	PL9

= Dataset 2

---

Step 3:

IM1	IM2	IM3
IM4	IM5	IM6
IM7	IM8	IM9

+

L1	L2	PL3
L4	L5	PL6
PL7	PL8	PL9

= Model 2

# Semi-Supervised Learning: Self-Training

Source: <https://arxiv.org/pdf/1906.01131v2.pdf>

# Self-Training Terminology

- **Teacher Network:** Network that creates pseudo labels
- **Student Network:** Network that learns using pseudo labels

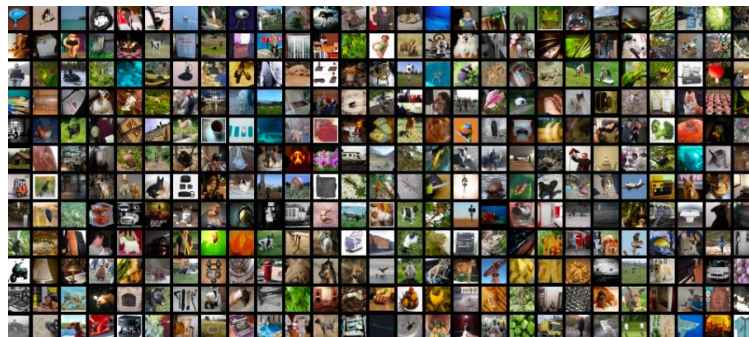
# Terminology

- **Self-Training:** When student network is same/larger sized than teacher network
- **Knowledge Distillation:** When student network is smaller than teacher network

# Noisy Student Self Training



ImageNet (**14M** examples)



JFT (**300M** examples)

	ImageNet top-1 acc.	ImageNet-A top-1 acc.
Prev. SOTA	86.4%	61.0%
Ours	<b>88.4%</b>	<b>83.7%</b>

Large accuracy gain for  
**adversarial** examples

# Adversarial Examples

ImageNet-A

Dragonfly



Manhole Cover



Bullfrog

Fox Squirrel



Monarch Butterfly

Washing Machine



Jay

Jeep

ImageNet-O

Ligature



Jellyfish (99%)



Painting

Goldfish (99%)

Highway



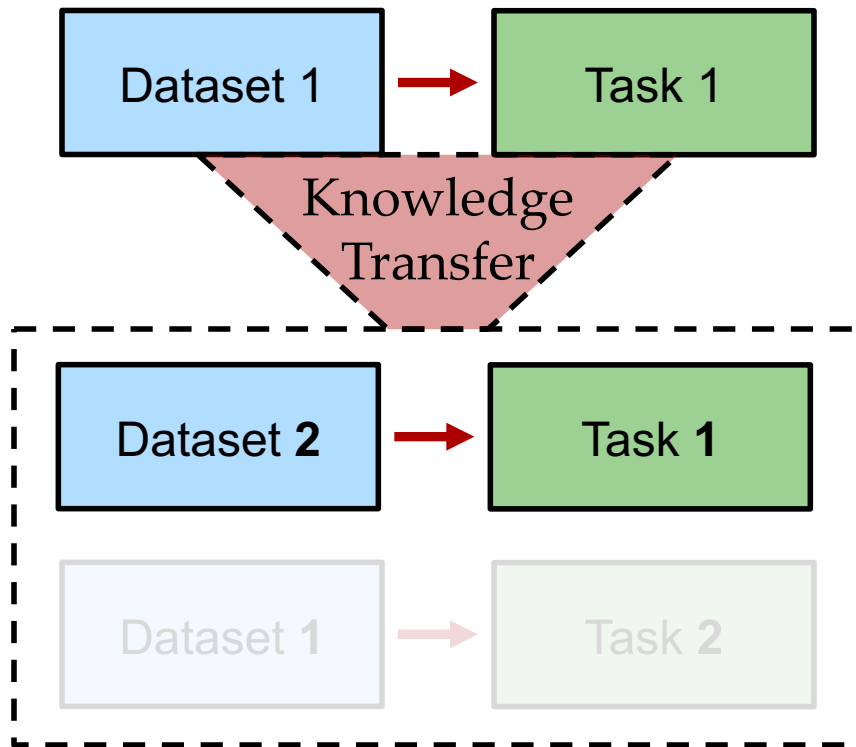
Dam (99%)

Garlic Bread



Hotdog (99%)

# Semi-Supervised Transfer Learning





# Self- Supervised Learning

IM = Image

PL = Pseudolabel

L = Label

Step 1:

IM1	IM2	IM3
IM4	IM5	IM6
IM7	IM8	IM9

+

PL1	PL2	PL3
PL4	PL5	PL6
PL7	PL8	PL9

= **Pretrained Model**

Step 2:

IM1	IM2
IM3	IM4

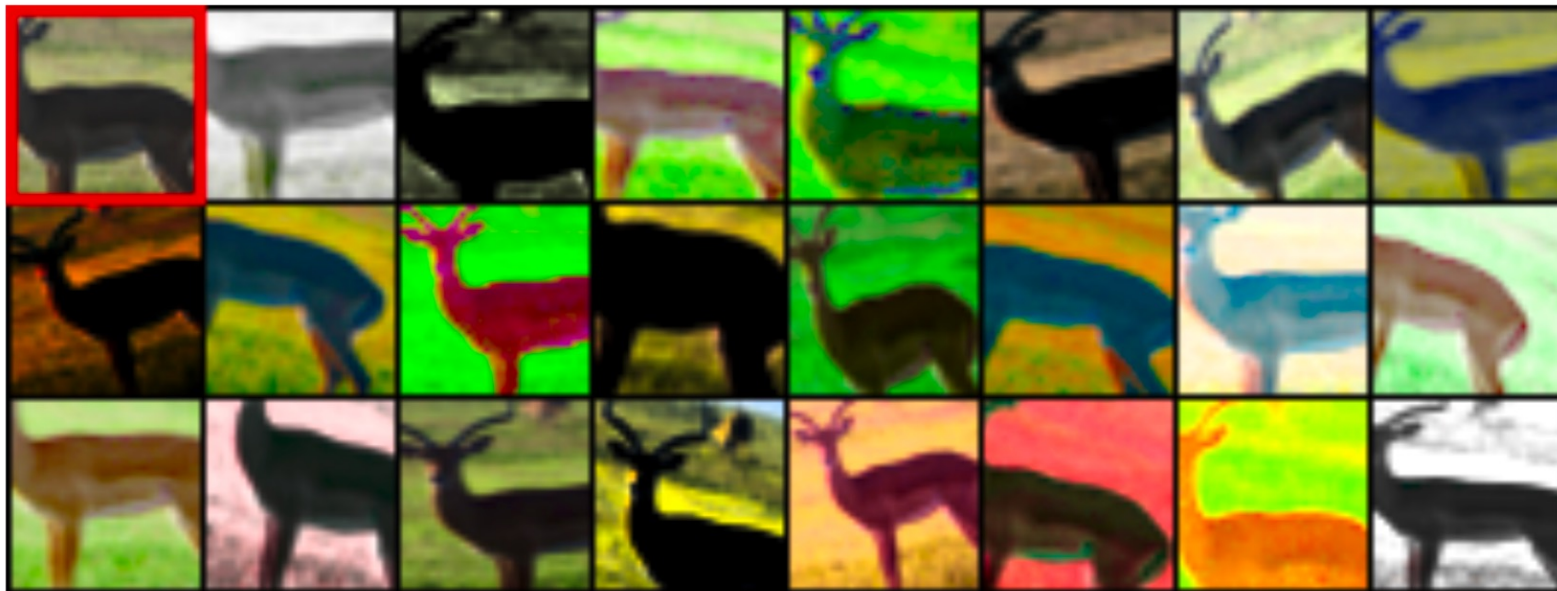
+

L1	L2
L3	L4

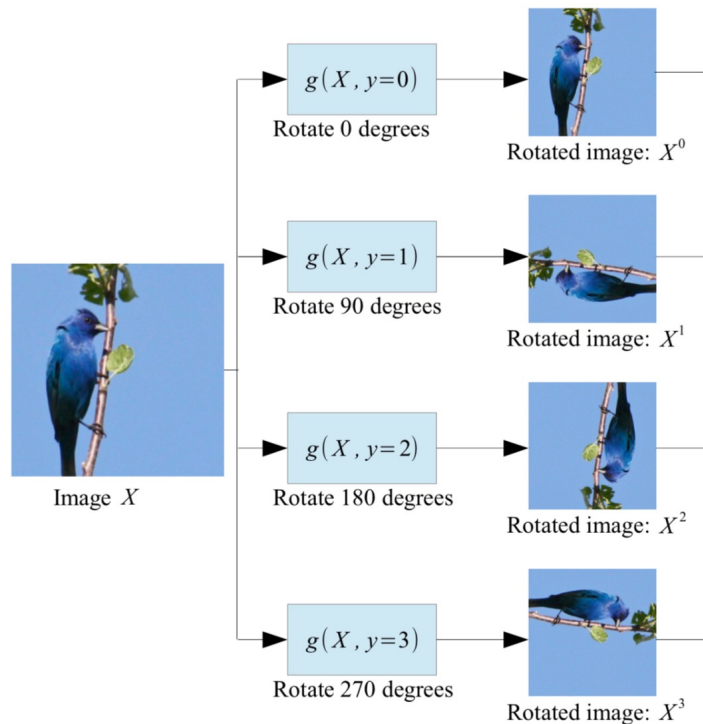
= **Model!**

# Pretext Tasks

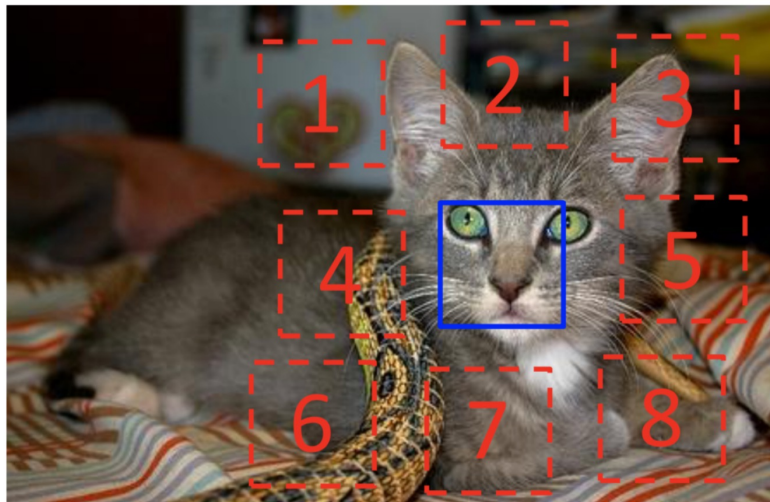
- Exemplar images



# Pretext Tasks - Rotation



# Pretext Tasks - Patching



Example:



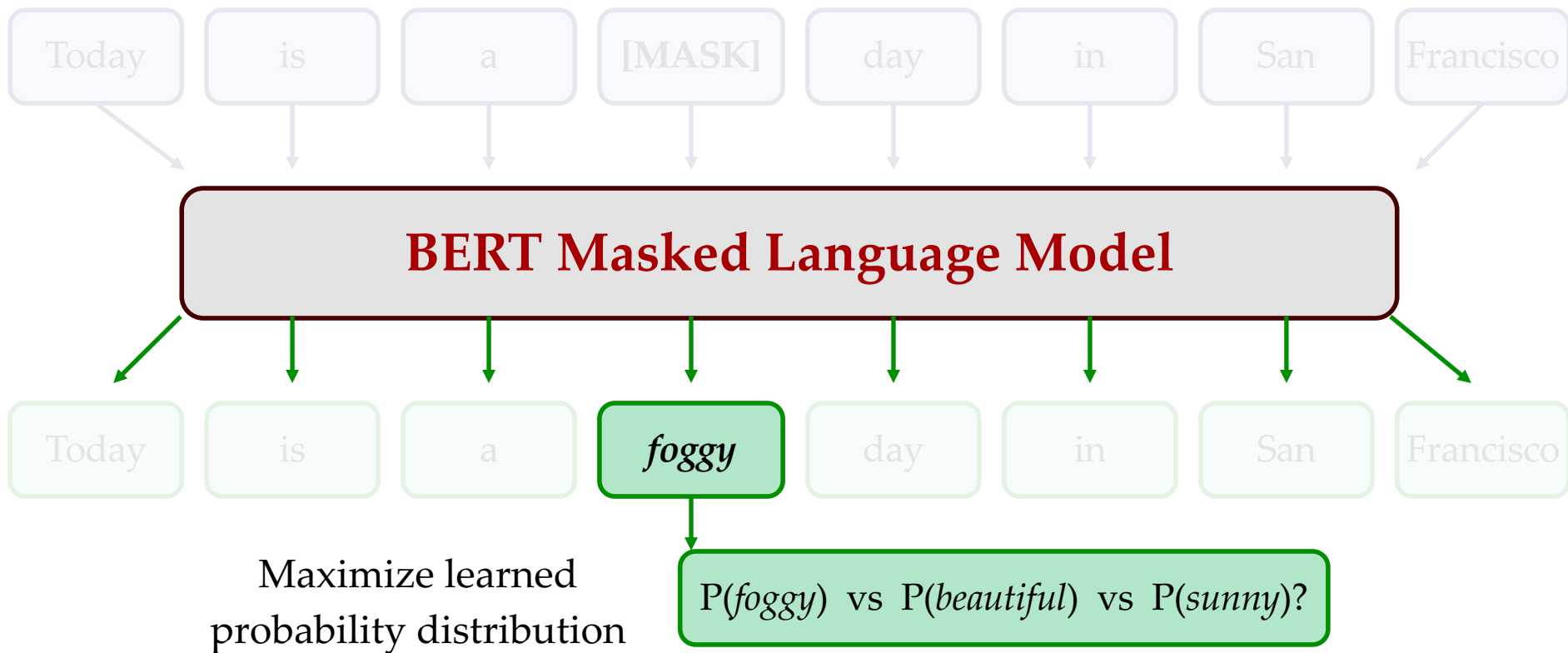
# What Position Does the Blue Square Occupy?

1	2	3
4	5	6
7	8	9



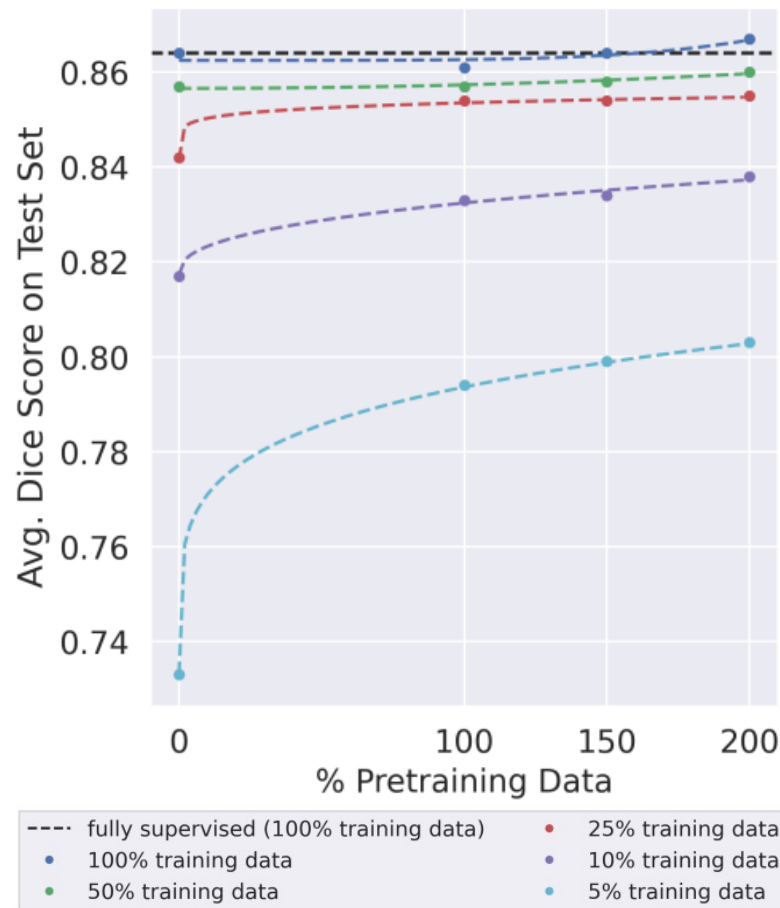
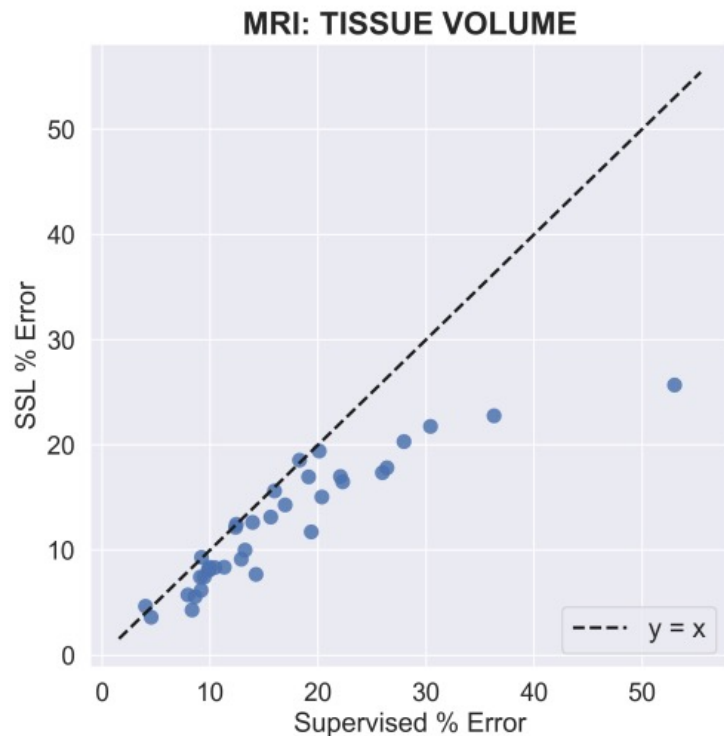
# Image Inpainting

# BERT Pretraining

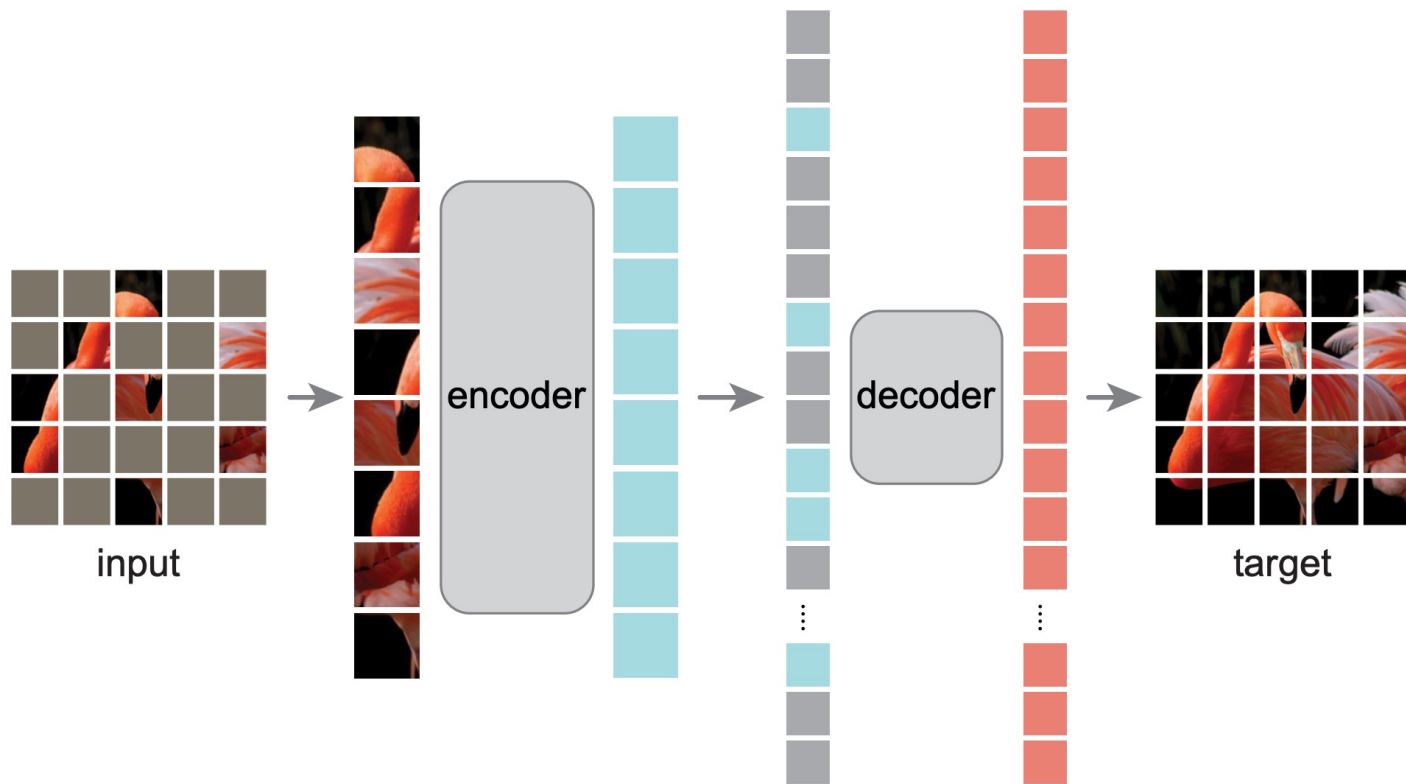




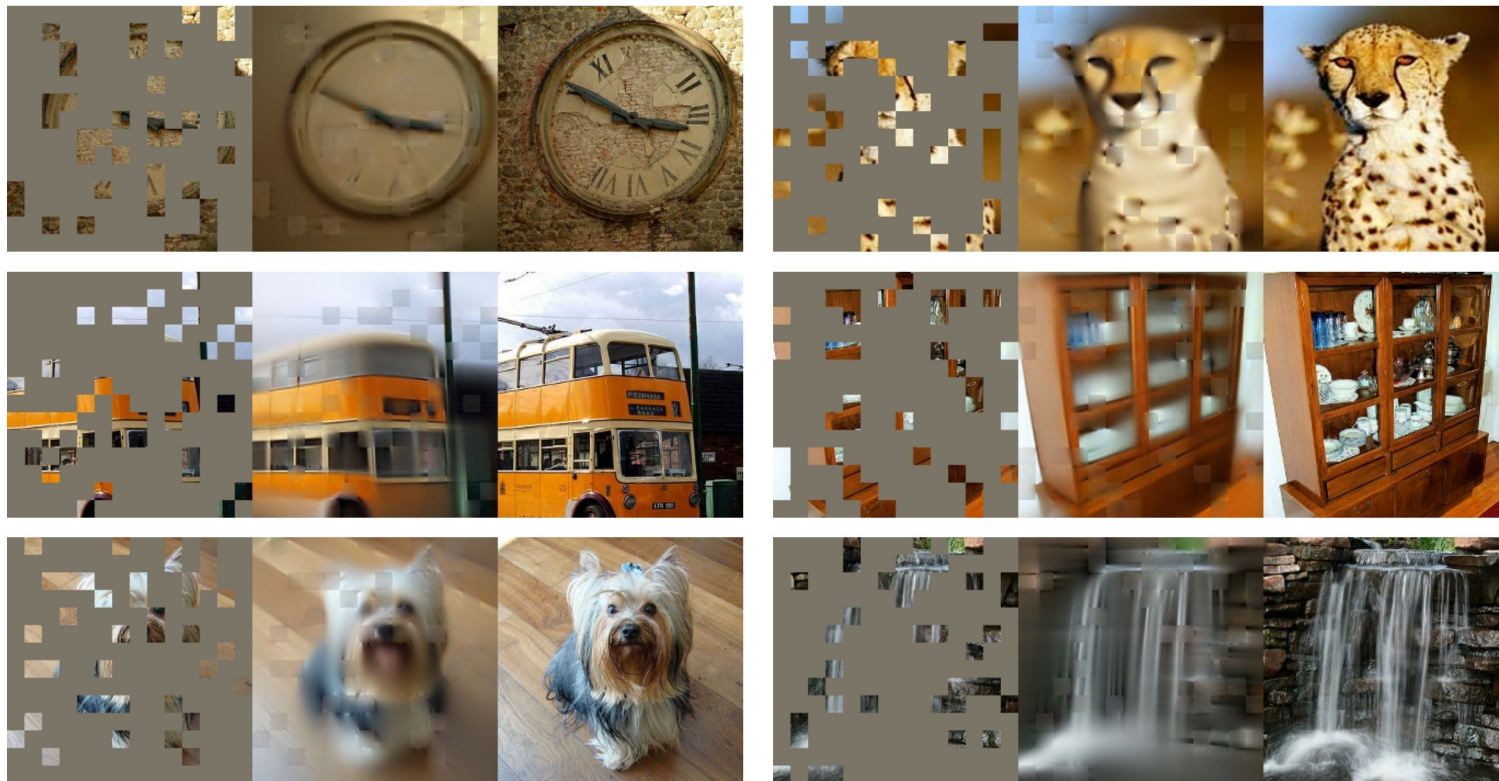
# Image SSL Benefits



# Latest Self-Supervised Learning



# Latest Self-Supervised Learning



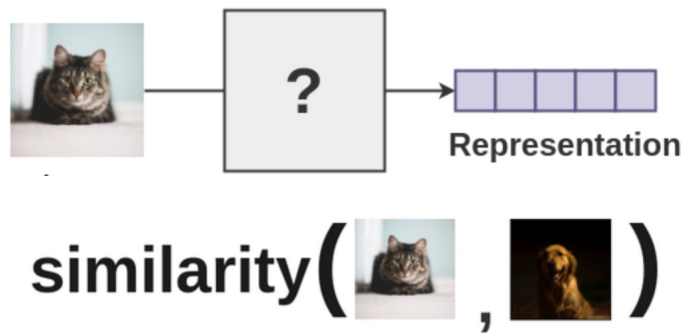
# Contrastive Learning

Match the correct animal

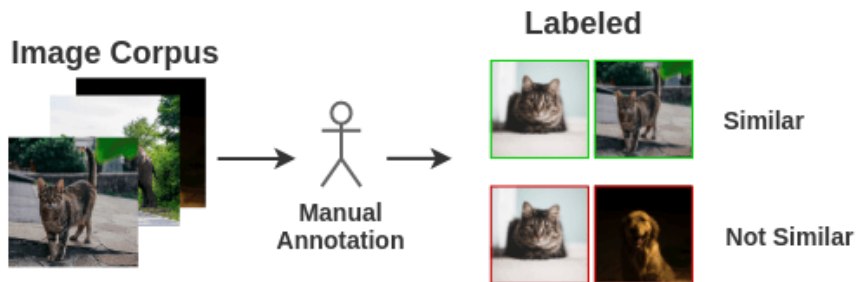


# Contrastive Learning

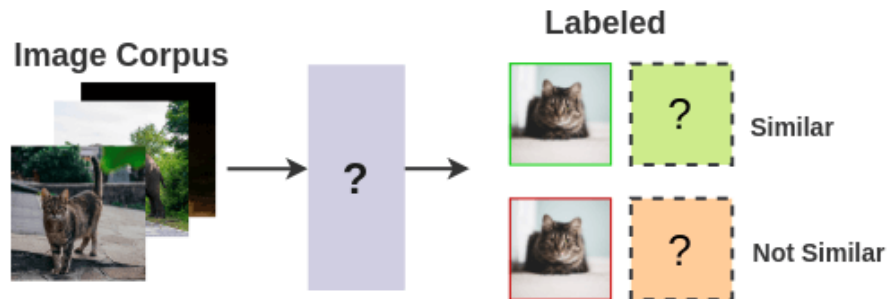
Need similar and different examples



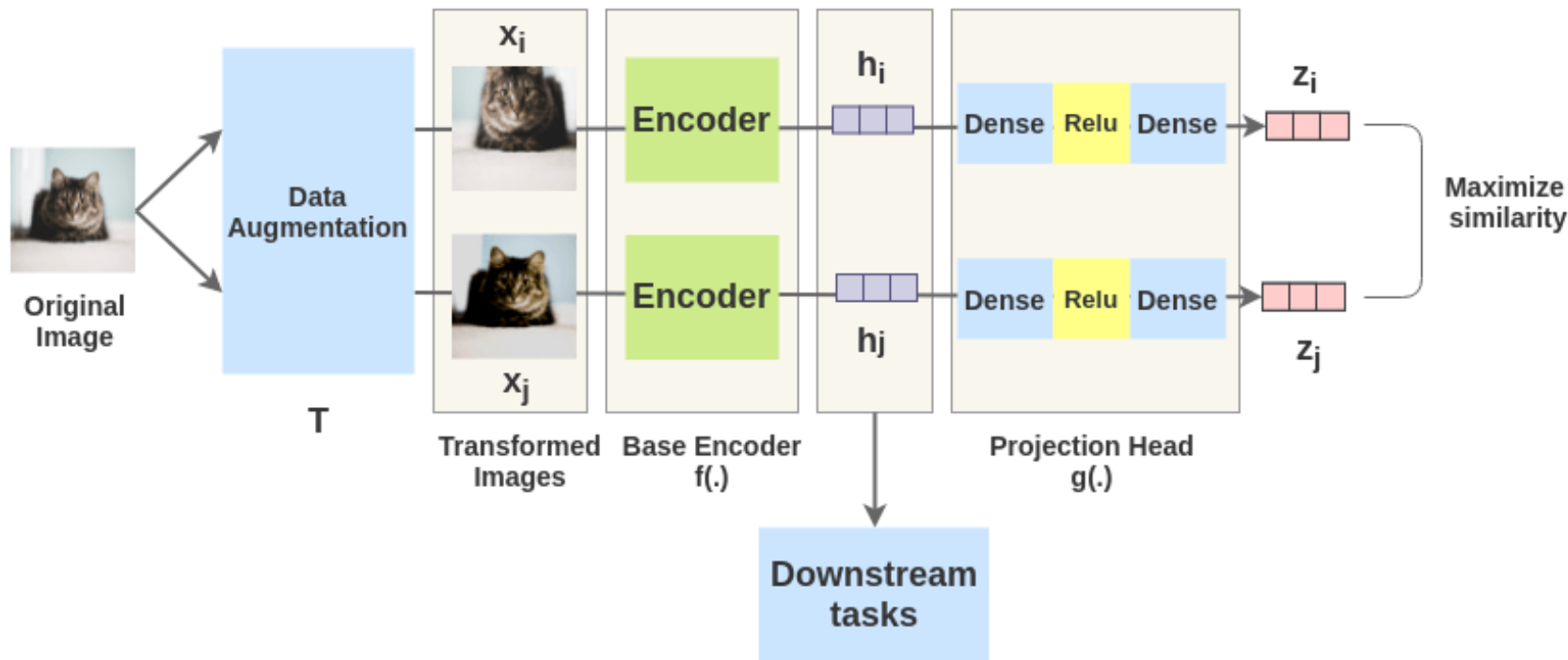
Supervised Approach



How can we automatically generate pairs?



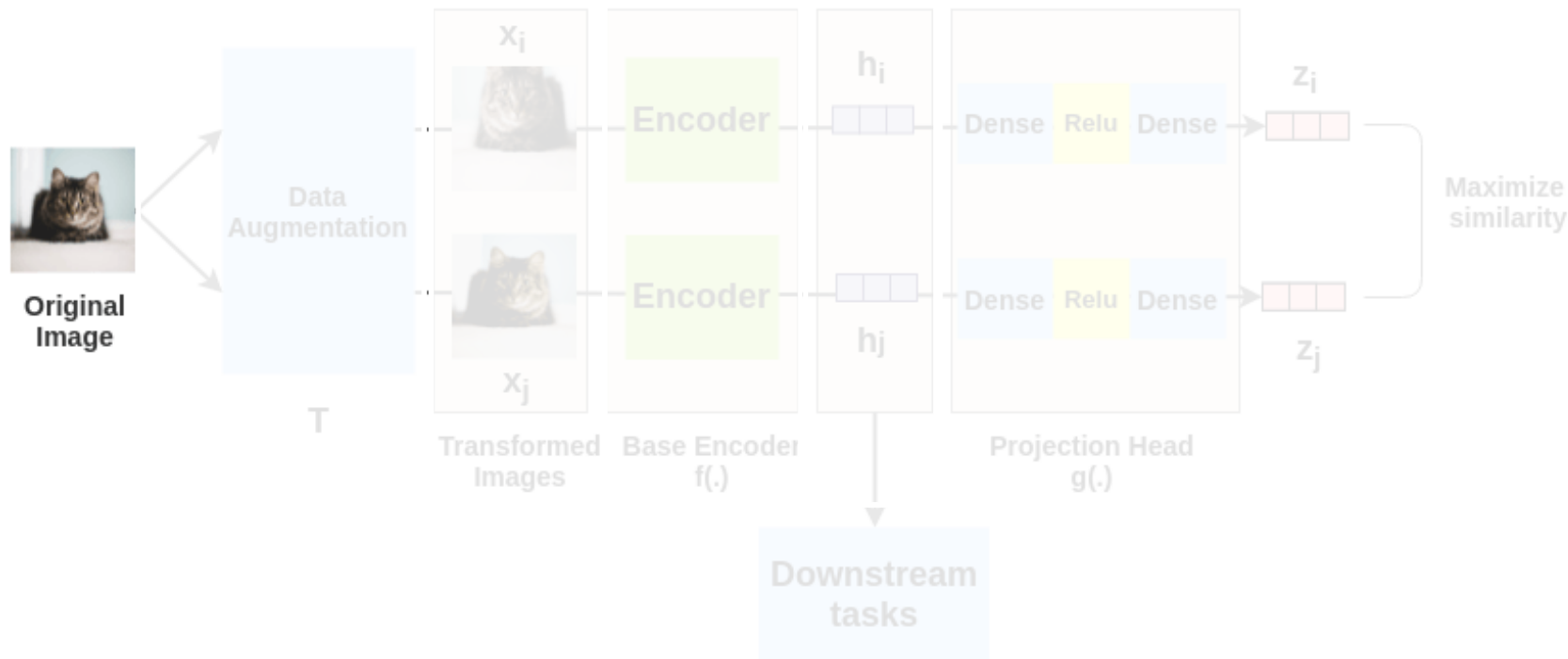
# Contrastive Learning



[1] Chaudhary A. The Illustrated SimCLR Framework. <https://amitnss.com/2020/03/illustrated-simclr/> 2020.

[2] Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Contrastive Learning



[1] Chaudhary A. The Illustrated SimCLR Framework. <https://amitnss.com/2020/03/illustrated-simclr/> 2020.

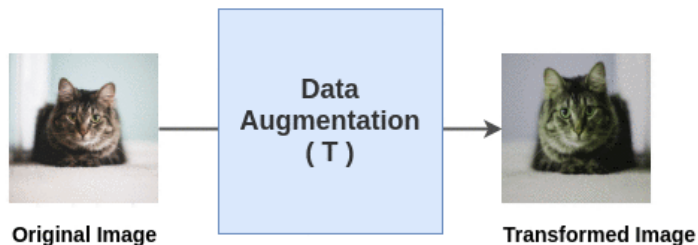
[2] Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.



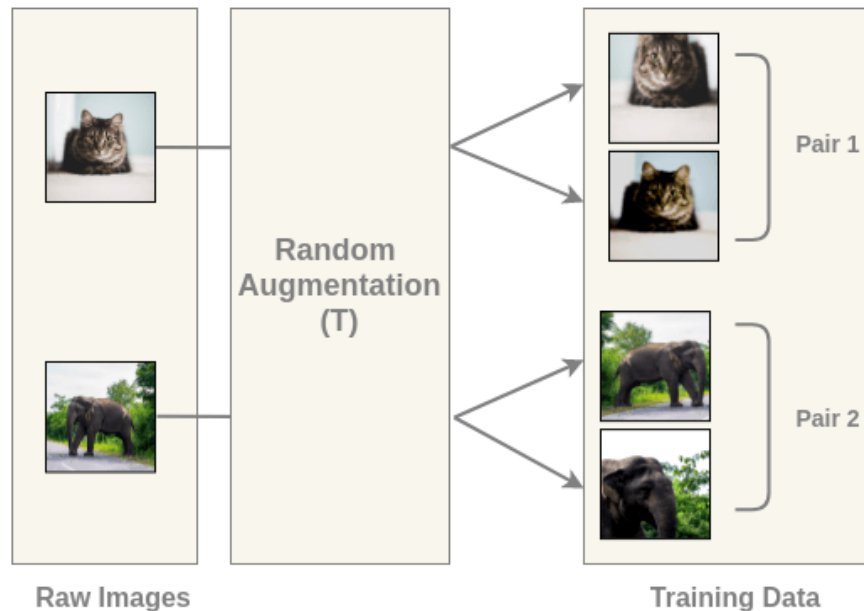
# Contrastive Learning

## Preparing similar pairs in a batch

### Random Transformation



Batch Size  
 $N = 2$



[1] Chaudhary A. The Illustrated SimCLR Framework. <https://amitnness.com/2020/03/illustrated-simclr/> 2020.

[2] Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

# Contrastive Learning

Similarity Calculation of Augmented Images

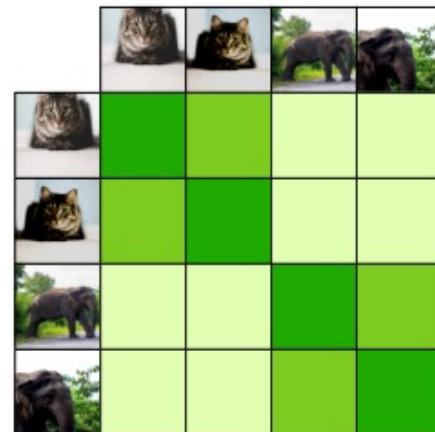
$$\text{similarity}(x_i, x_j) = \text{cosine similarity}(z_i, z_j)$$



Softmax =

$$\frac{e^{\text{similarity}(x_i, x_j)}}{e^{\text{similarity}(x_i, x_j)} + e^{\text{similarity}(x_i, x_{\text{ele}})} + e^{\text{similarity}(x_i, x_{\text{ele}}^*)}}$$

Pairwise cosine similarity



[1] Chaudhary A. The Illustrated SimCLR Framework. <https://amitnss.com/2020/03/illustrated-simclr/> 2020.

[2] Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020.

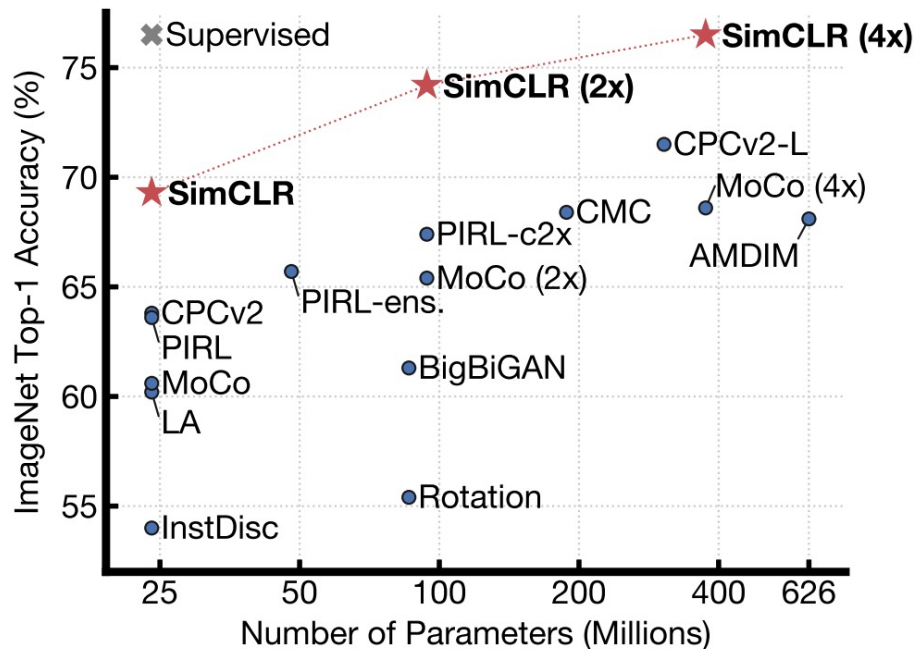
# Noise Contrastive Estimation Loss

- Compute over both pairs to account for asymmetry

$$l(i, j) = -\log \frac{\exp(s_{i,j})}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k})}$$

$$l(\text{cat}_1, \text{cat}_2) = -\log \left( \frac{\exp(\text{similarity}(\text{cat}_1, \text{cat}_2))}{\exp(\text{similarity}(\text{cat}_1, \text{cat}_2)) + \exp(\text{similarity}(\text{cat}_1, \text{elephant}_1)) + \exp(\text{similarity}(\text{cat}_1, \text{elephant}_2))} \right)$$

# SimCLR Technique



- Large batch size requirements
- Long training times
- Heuristic data augmentations

# SimCLR Augmentations



(a) Original



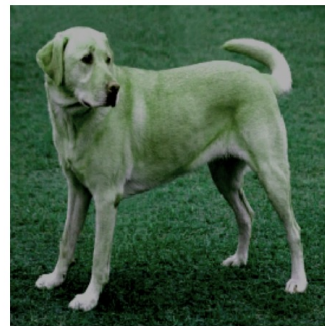
(b) Crop and resize



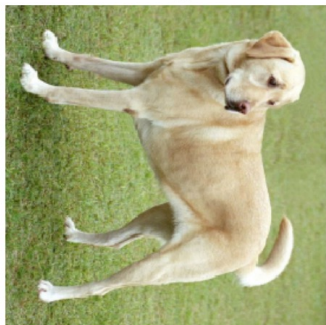
(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise

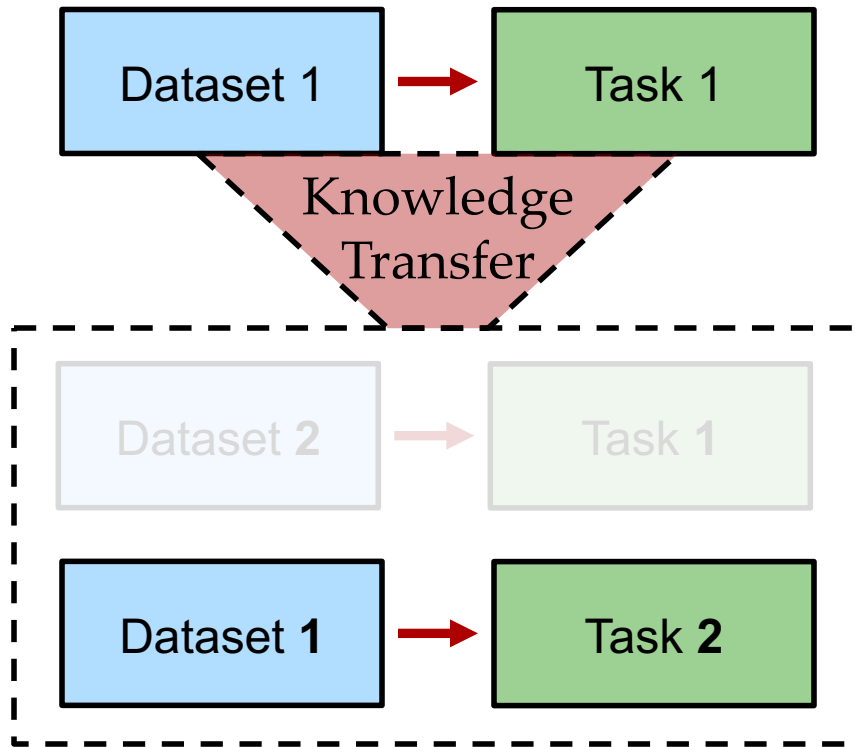


(i) Gaussian blur

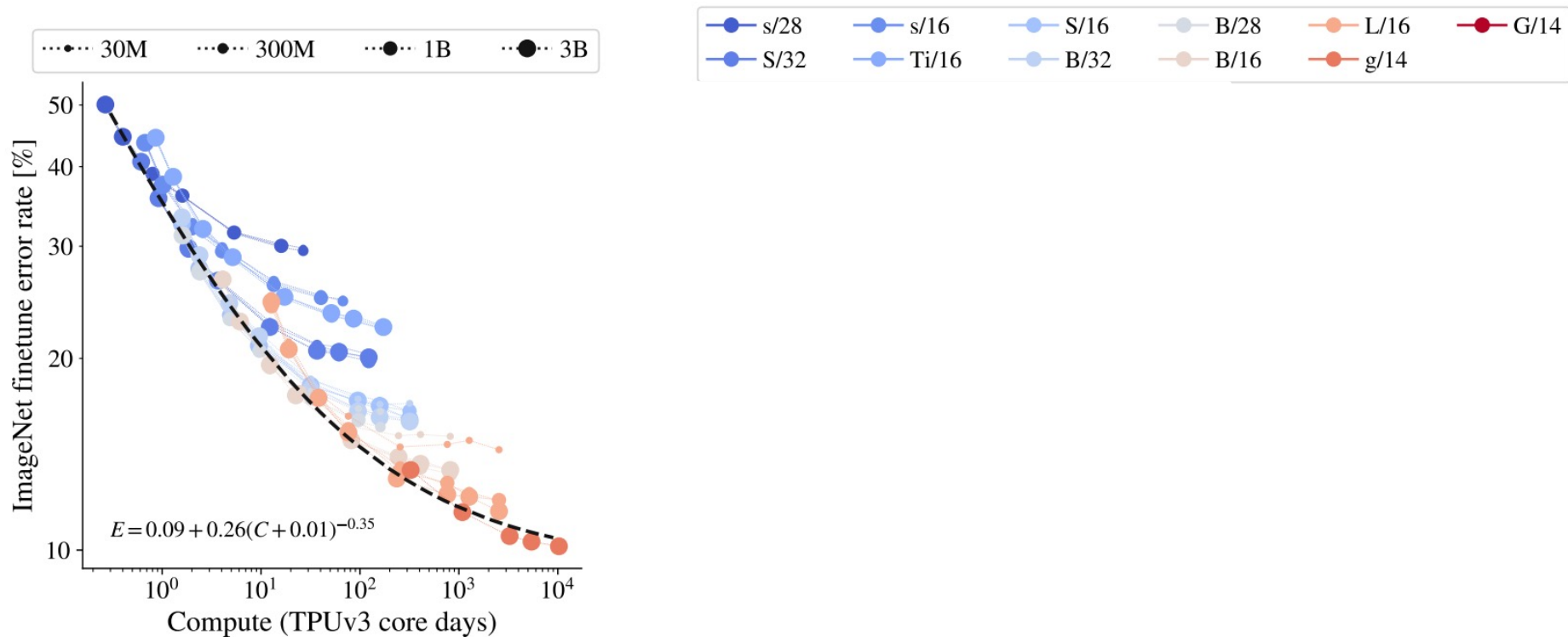


(j) Sobel filtering

# Transfer Learning

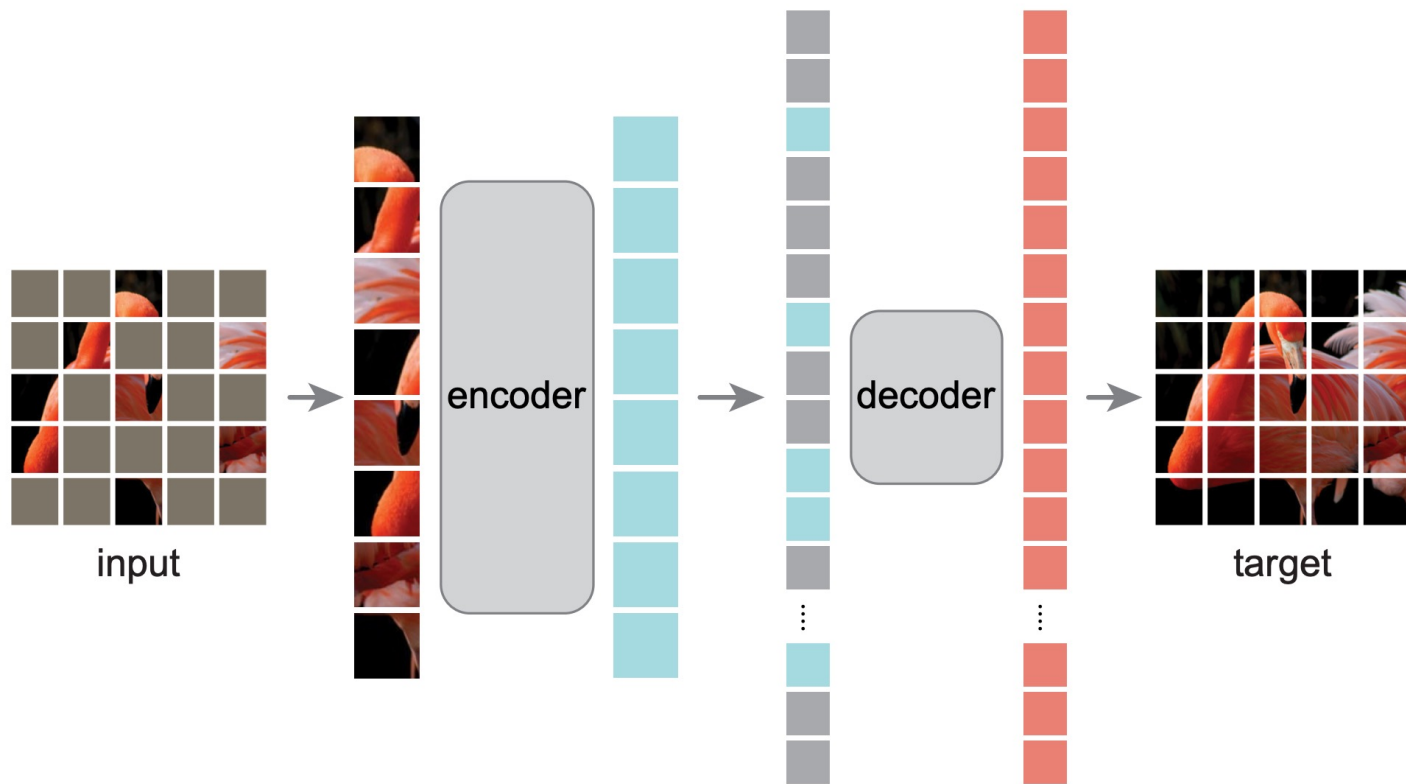


# Benefits of Scaling

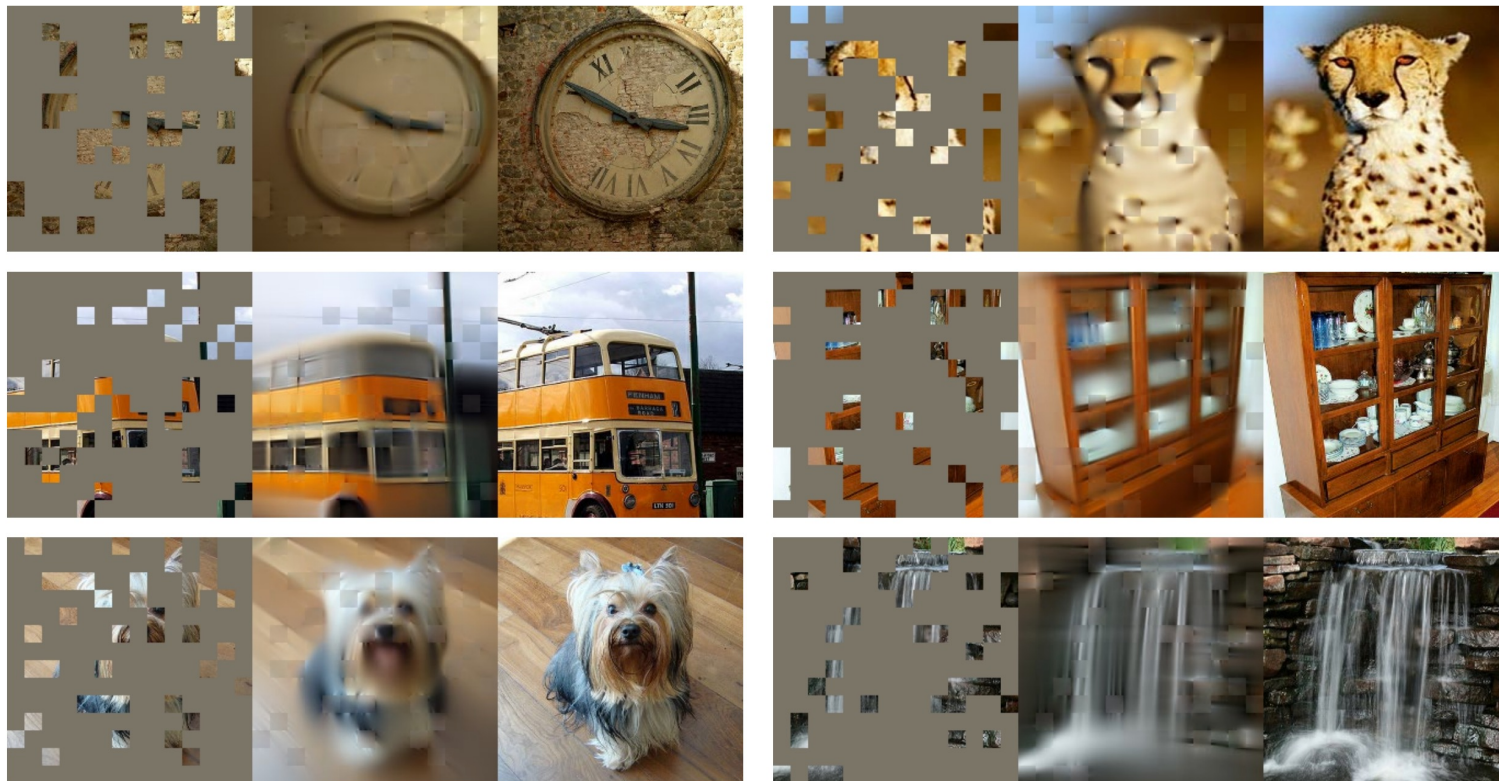




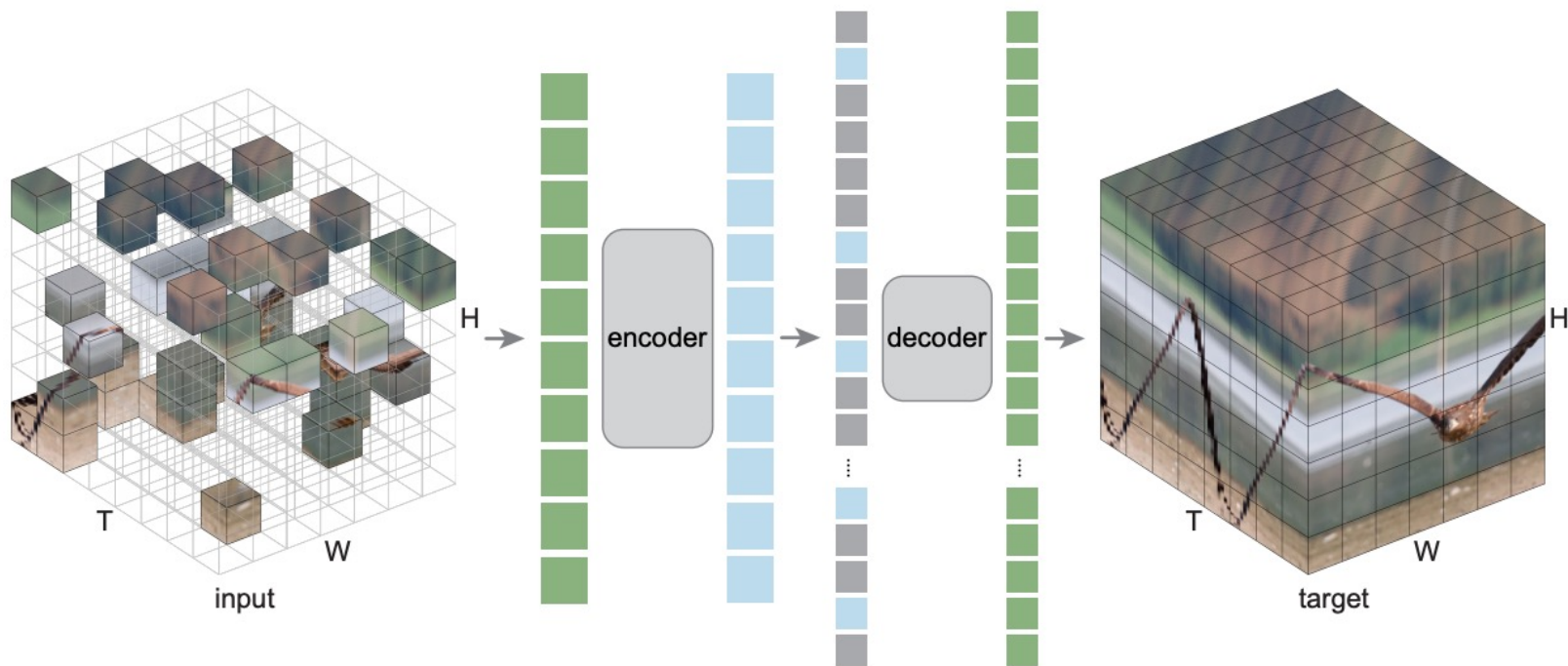
# Latest Self-Supervised Learning



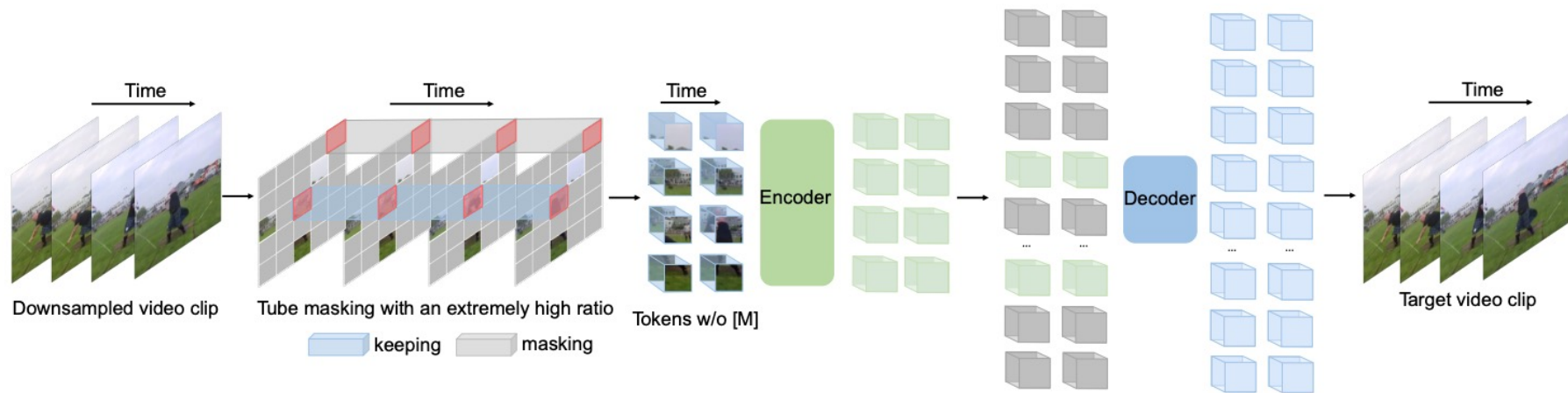
# Latest Self-Supervised Learning



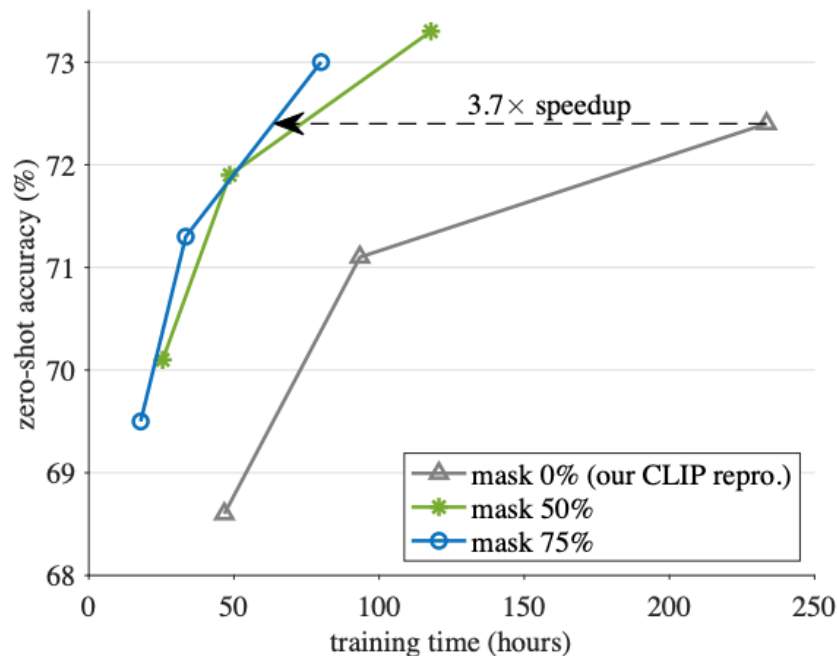
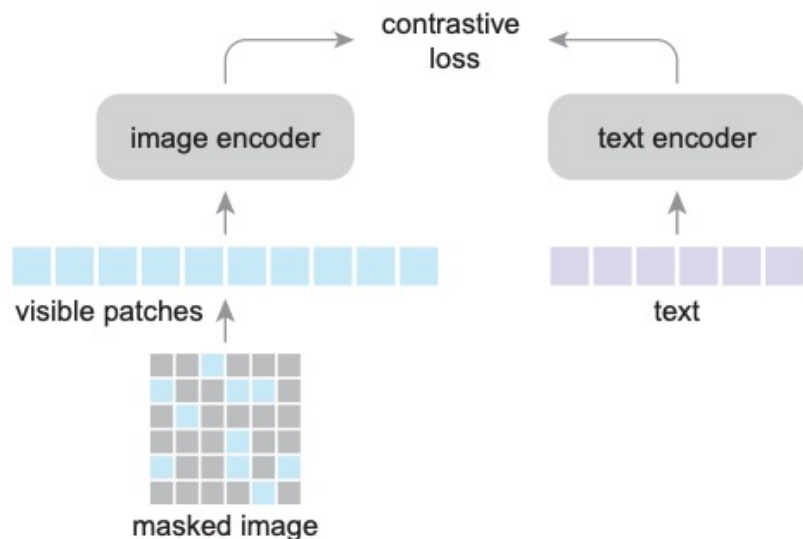
# Extensions of Prior Approaches



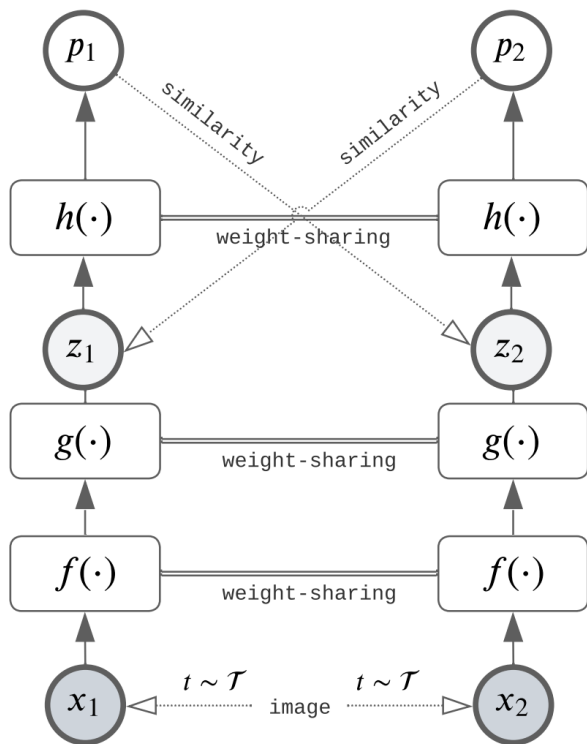
# Tube Masking for Video MAE



# Masking for Multi-Modal Learning



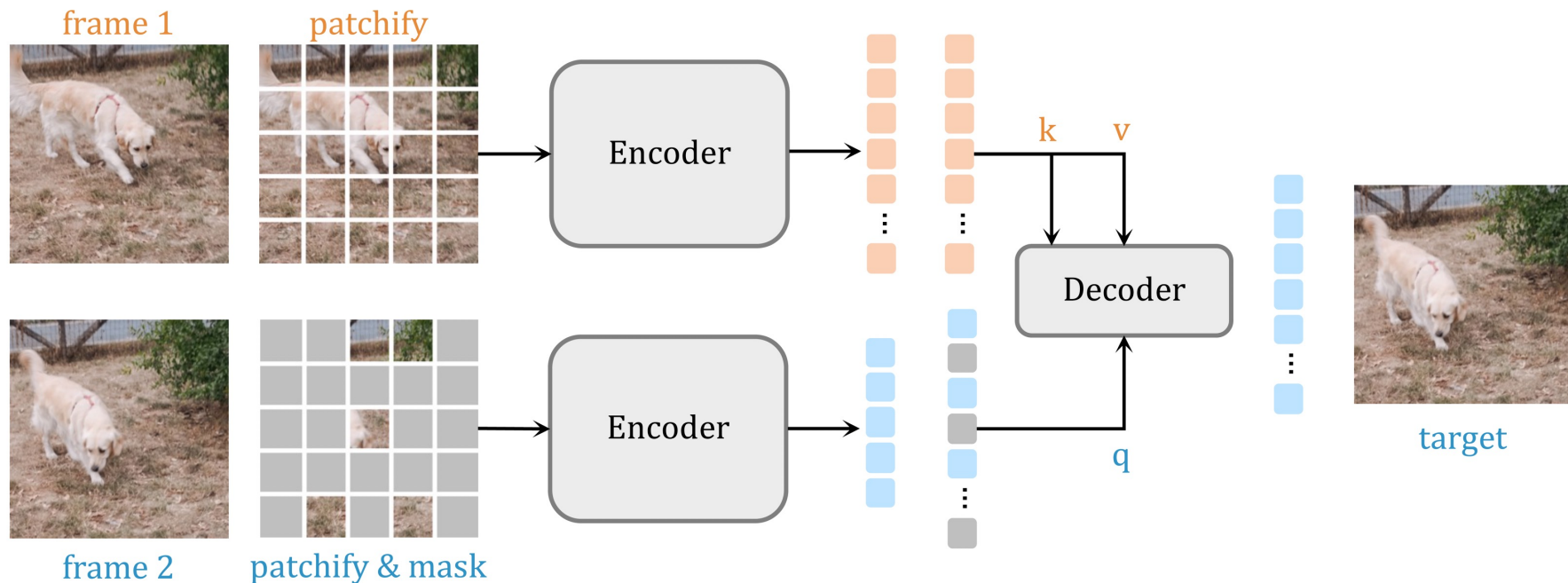
# Adaptations to Medical Imaging



- Architectures inspired by the SimSiam method
- $X$  = Image
- $F$  = Image Encoder
- $G, H$  = MLP projectors
- $Z$  = latent representation
- $P$  = latent representation predictors

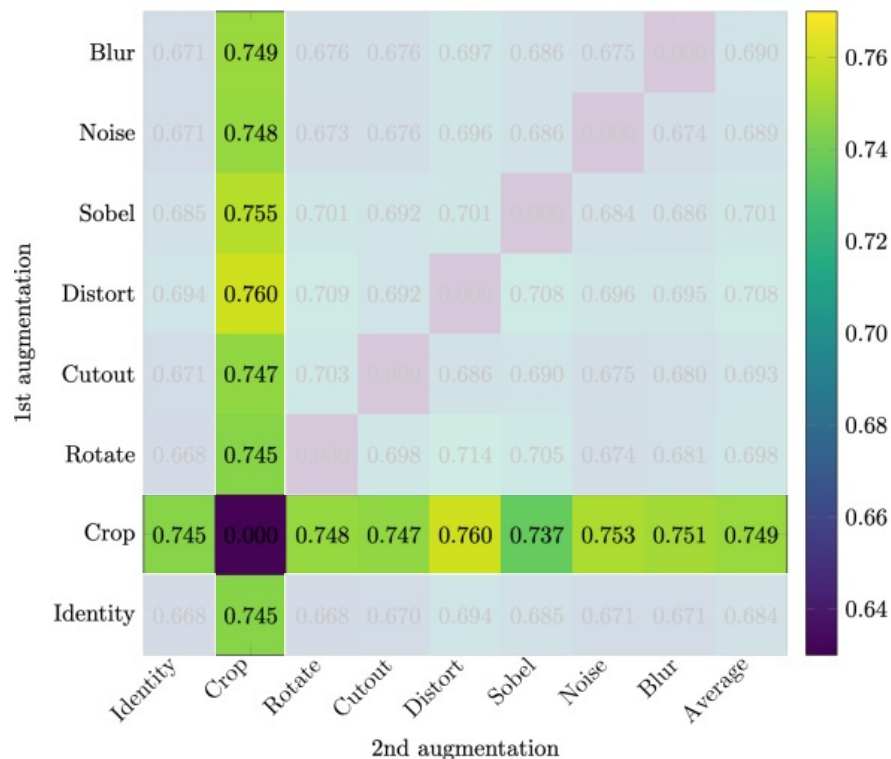
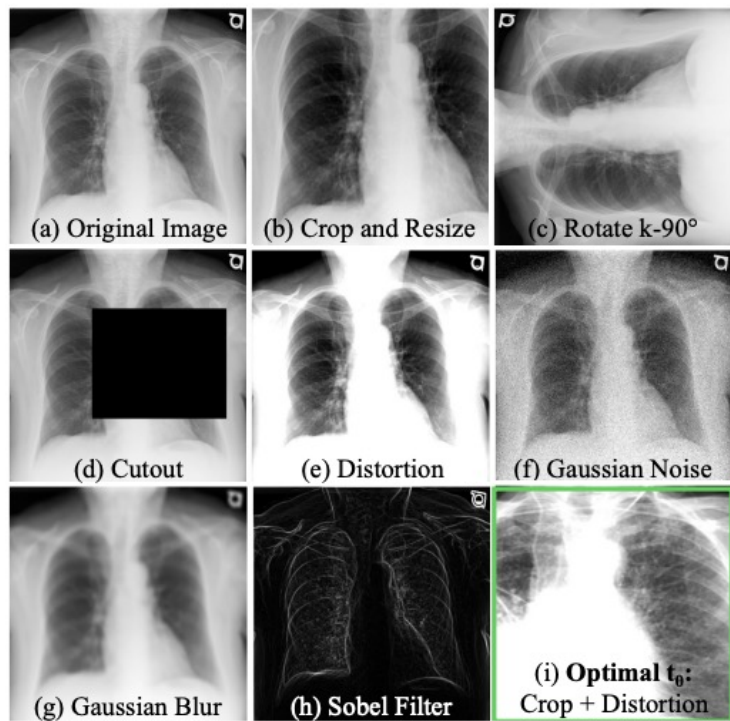
$$\mathcal{L} = -\frac{1}{2} \left( \frac{p_1}{\|p_1\|_2} \cdot \frac{z_2}{\|z_2\|_2} \right) - \frac{1}{2} \left( \frac{p_2}{\|p_2\|_2} \cdot \frac{z_1}{\|z_1\|_2} \right)$$

# Siamese Masked Autoencoders





# Augmentations for Medical Imaging



# Suggested Reading

---

## **A Simple Framework for Contrastive Learning of Visual Representations**

---

**Ting Chen<sup>1</sup> Simon Kornblith<sup>1</sup> Mohammad Norouzi<sup>1</sup> Geoffrey Hinton<sup>1</sup>**

## **Masked Autoencoders Are Scalable Vision Learners**

**Kaiming He<sup>\*,†</sup> Xinlei Chen<sup>\*</sup> Saining Xie Yanghao Li Piotr Dollár Ross Girshick**

<sup>\*</sup>equal technical contribution      <sup>†</sup>project lead

**Facebook AI Research (FAIR)**

# Questions?

[akshaysc@stanford.edu](mailto:akshaysc@stanford.edu)

