

# Self-supervised learning

---

682: Neural Networks: A Modern Introduction

Subhransu Maji

April 14, 2026

*College of*  
INFORMATION AND  
COMPUTER SCIENCES



# Administrative

**Midterm 2** on Tuesday, April 28 in class

Syllabus: Lecture 9 onwards (Image classification with CNNs)

# Today's Class

- **Recap**
  - Supervised vs Unsupervised Learning
  - Why not always label data?
- **Semi-supervised Learning**
  - Concepts
  - Example: pseudo-labels / self-training
- **Self-supervised Learning**
  - Concepts
  - Pretext tasks
  - Contrastive Learning
  - Beyond images

# Today's Class

- **Recap**
  - Supervised vs Unsupervised Learning
  - Why not always label data?
- **Semi-supervised Learning**
  - Concepts
  - Example: pseudo-labels / self-training
- **Self-supervised Learning**
  - Concepts
  - Pretext tasks
  - Contrastive Learning
  - Beyond images

# Recap: Supervised vs Unsupervised Learning

## Supervised Learning

Data:  $(X, y)$

$X$  = input/feature/image/...

$y$  = label/target



→ Cat



→ Dog

## Unsupervised Learning

Data:  $X$

Just  $X$ , no labels

Learn about the *structure* of the data,  
i.e.  $P(X)$



.....

# So let's always use Supervised Learning?

## Supervised Learning

Data:  $(X, y)$

$X$  = input/feature/image/...

$y$  = label/target



→ Cat



→ Dog

## “Standard” Supervised Learning:

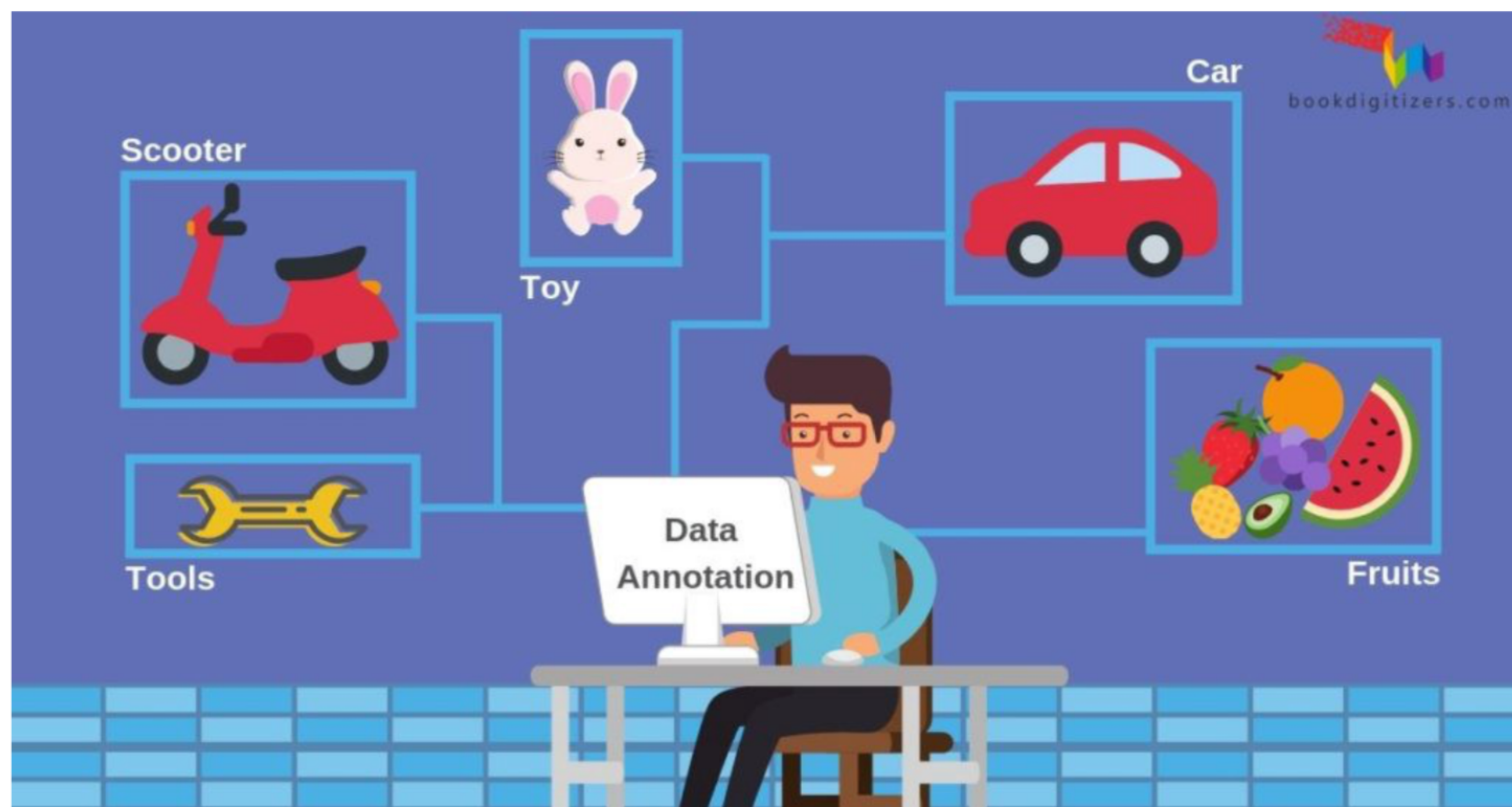
1. Collect a large set of data (images..) as the “training set”
2. Label each one as cat / dog / monkey / ...
3. Train a model mapping image to label

$$f : \mathbf{X} \rightarrow y$$

4. Go forth and classify the world with  $f$  !

# Data Annotation

Supervised Learning first requires labeling a very large amount of data



# Labeling Image Categories - “Easy” Until ....

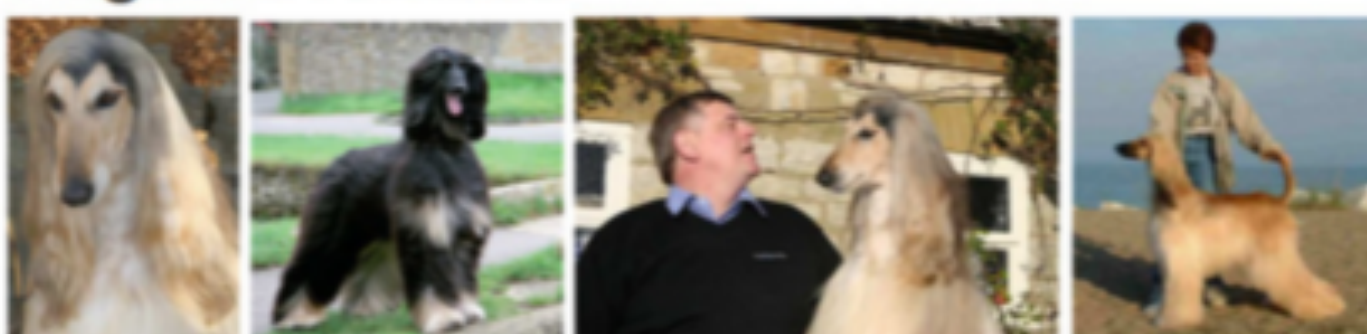
Blenheim Spaniel



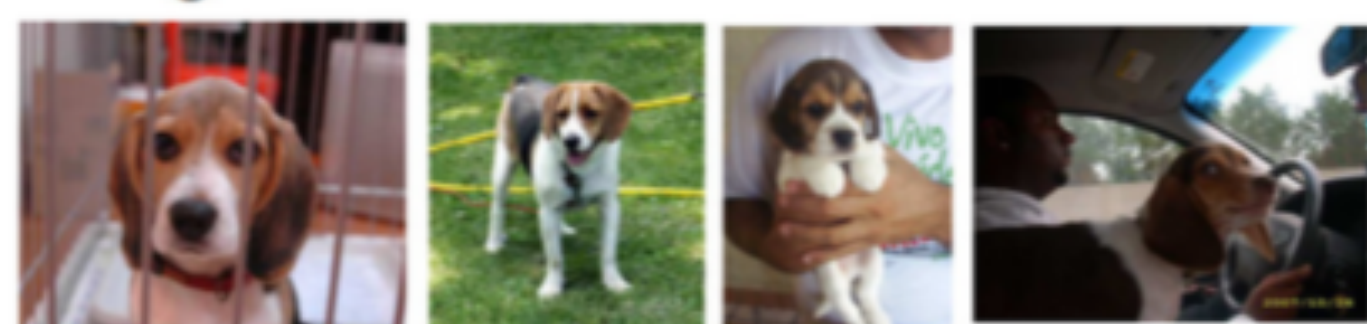
Toy Terrier



Afghan Hound



Beagle



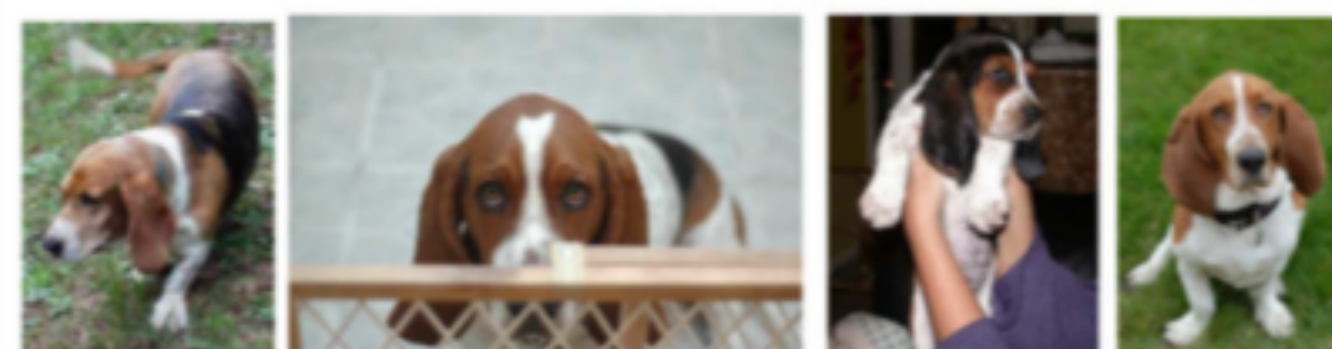
Papillon



Rhodesian Ridgeback



Basset Hound



Bloodhound



- Over **120 dog breeds** in ImageNet dataset for image classification
- Non-expert labelers may not be aware of these **fine-grained** differences, leading to **labeling errors**
- *E.g.*, the Caltech UCSD birds dataset has 4% labeling error (NABirds, Van horn et al. CVPR15)

# Dense Semantic and Instance Labels

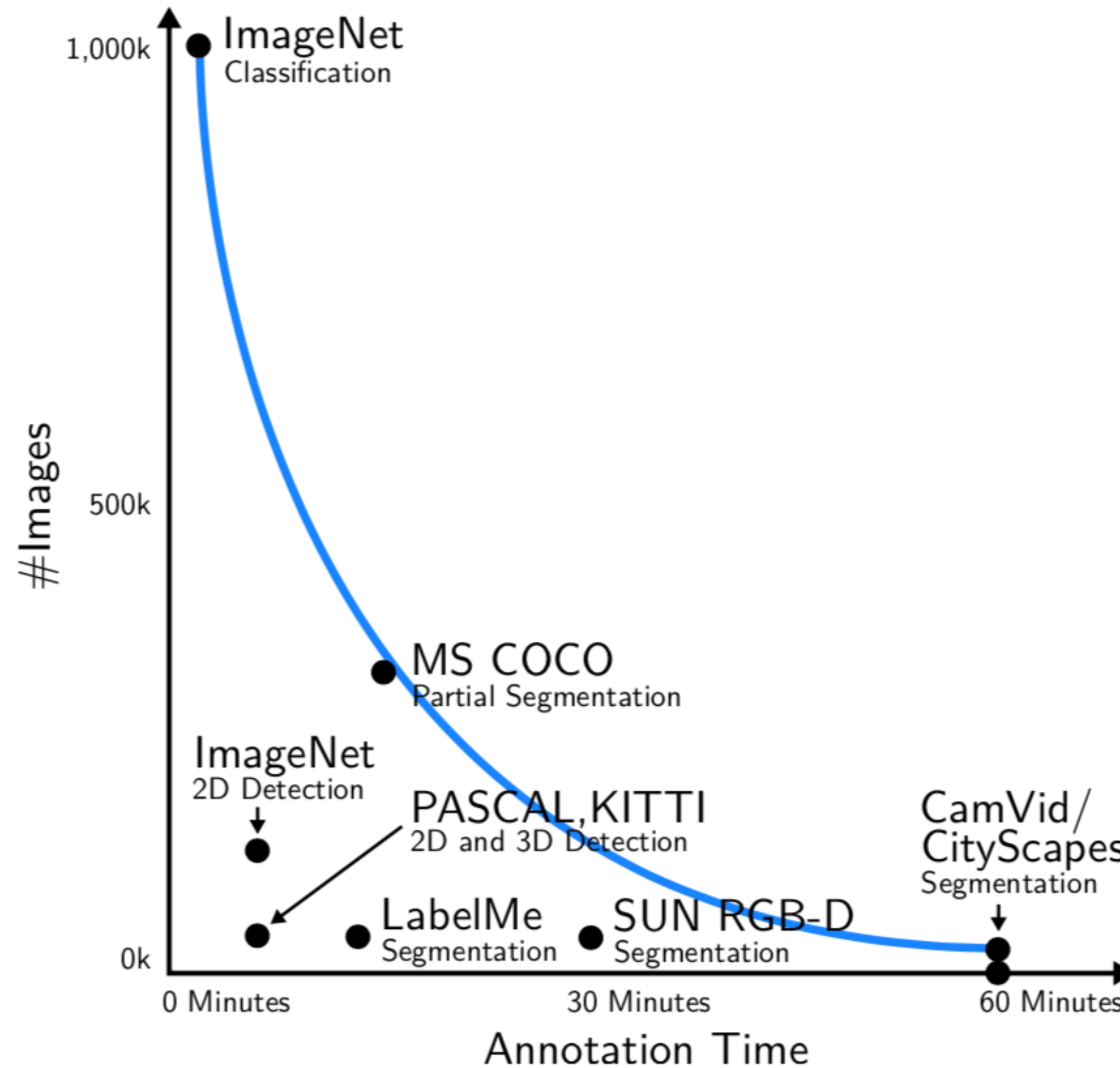


“Cityscape” dataset: Labeling every pixel as person/road/sidewalk ...

Annotation time **60-90 minutes per image**

[Slides](#) from Andreas Geiger, MPI Tubingen

# Annotate Everything — Expensive, doesn't Scale!



# Motivation - Humans learn with little supervision

Provided with very few “labeled” examples (someone pointing something out to us explicitly), we can generalize quite well.

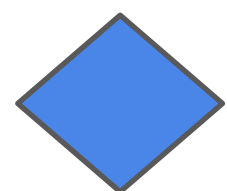
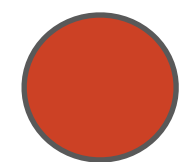


# Today's Class

- Recap
  - Supervised vs Unsupervised Learning
  - Why not always label data?
- **Semi-supervised Learning**
  - Concepts
  - Example: pseudo-labels / self-training
  - Example: Distillation, Student/Teacher
- Self-supervised Learning
  - Concepts
  - Pretext tasks
  - Contrastive Learning

# Semi-supervised Learning

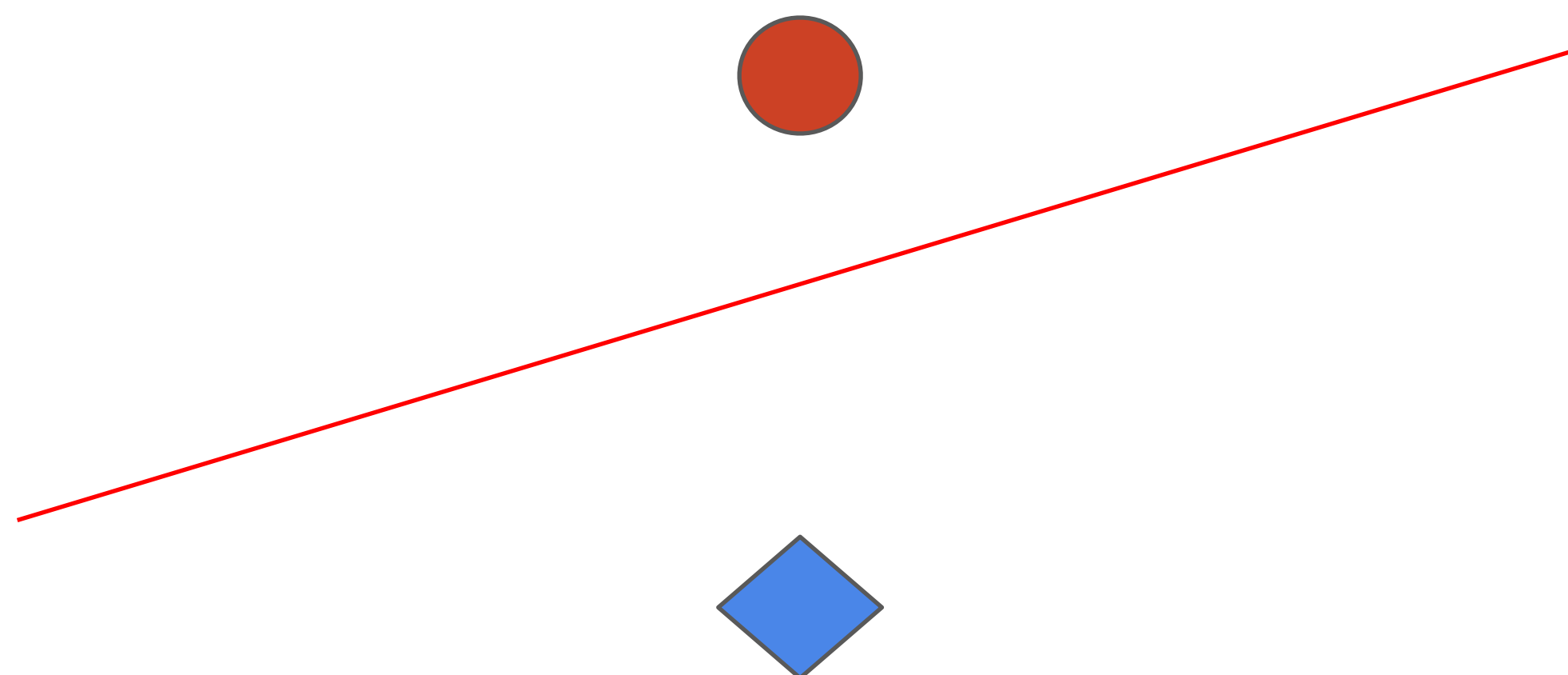
- Given a small amount of *labeled* data  $\mathcal{X}_L$
- Given (usually) large amount of *unlabeled* data  $\mathcal{X}_U$
- Can  $\mathcal{X}_U$  help us in getting a better model?



What is a good decision boundary for these points?

# Semi-supervised Learning

- Given a small amount of *labeled* data  $\mathcal{X}_L$
- Given (usually) large amount of *unlabeled* data  $\mathcal{X}_U$
- Can  $\mathcal{X}_U$  help us in getting a better model?

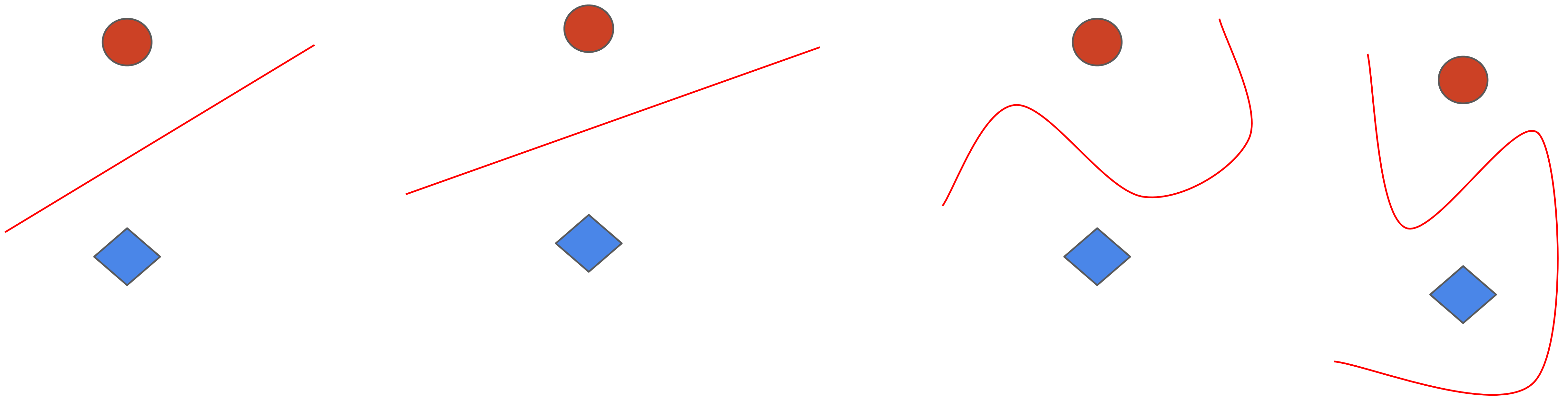


What is a good decision boundary for these points?

# Semi-supervised Learning

- Given a small amount of *labeled* data  $\mathcal{X}_L$
- Given (usually) large amount of *unlabeled* data  $\mathcal{X}_U$
- Can  $\mathcal{X}_U$  help us in getting a better model?

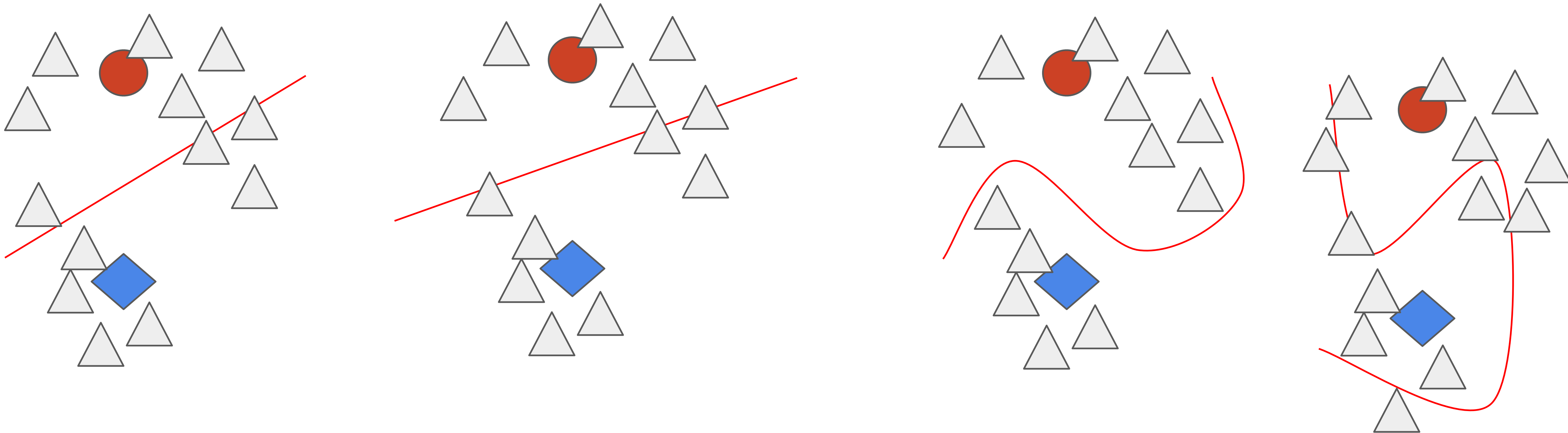
Which one is  
your  
favourite?



# Semi-supervised Learning

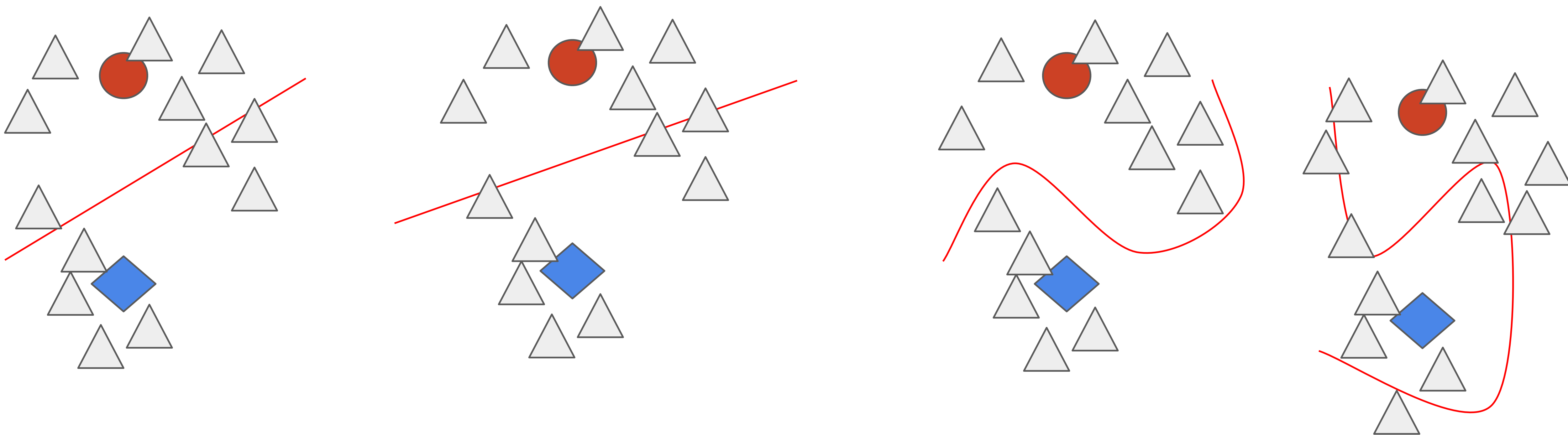
- Given a small amount of *labeled* data  $\mathcal{X}_L$
- Given (usually) large amount of *unlabeled* data  $\mathcal{X}_U$
- Can  $\mathcal{X}_U$  help us in getting a better model?

Now we see  
some unlabeled  
data points ....



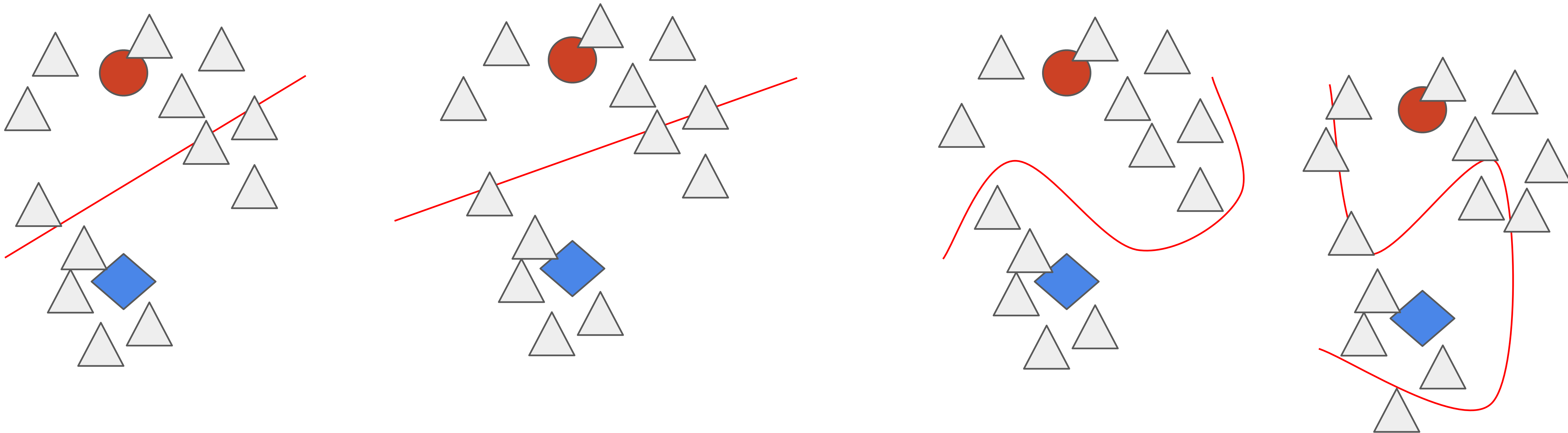
# Semi-supervised Learning - intuitions

- Unlabeled samples tell us about  $P(\mathbf{X})$ , which is useful in the predictive posterior  $P(y | \mathbf{X})$



# Semi-supervised Learning - definitions

- **Smoothness assumption:** if  $x_1, x_2$  are close, labels  $y_1, y_2$  are also “close”
- **Low-density separation:**  $x_1, x_2$  are separated by *low-density region* then labels are not “close”
- **Cluster assumption:** points in same cluster likely to have same label



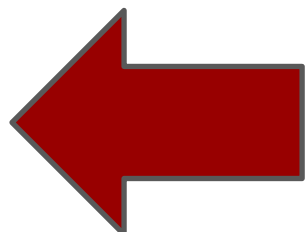
# Semi-supervised Learning Approaches

- We will look at *a simple approach* to semi-supervised learning
- **Self-training** or **pseudo-labeling**
  - Age-old method
  - Surprisingly good with modern deep learning methods
  - But many variations ...

# Self-training

- Assume: one's own high confidence predictions are correct!
- Train model  $f$  on  $\mathcal{X}_L := \{x_L, y_L\}$
- Use  $f$  to predict “pseudo-labels” on  $\mathcal{X}_U := \{x_u\}$
- Add  $\{x_u, f(x_u)\}$  to labeled data
- Repeat

# Self-training - variations

- Assume: one's own high confidence predictions are correct!
  - Train model  $f$  on  $\mathcal{X}_L := \{x_L, y_L\}$
  - Use  $f$  to predict “pseudo-labels” on  $\mathcal{X}_U := \{x_u\}$
  - Add  $\{x_u, f(x_u)\}$  to labeled data
  - Repeat
- 
- 1) Add only a few most confident predictions on  $X_u$
  - 2) Add all predictions on  $X_u$
  - 3) Add all predictions, weighted by the confidence of the prediction

....

# Self-training advantages

- The simplest semi-supervised method!
- It's a “wrapper” - the classifiers or models can be arbitrarily complex, we do not need to delve into those details to apply self-training
- Often quite good in practice, e.g. in natural language tasks
- Also some vision tasks ...

Disadvantages of self-training?

# Disadvantages of self-training?

- Early mistakes can reinforce themselves
  - We have heuristic solutions, like discarding samples if the confidence of prediction falls below some threshold
- Convergence
  - Hard to say if these steps of self-train and repeat will converge

# Domain shifts can have a large impact

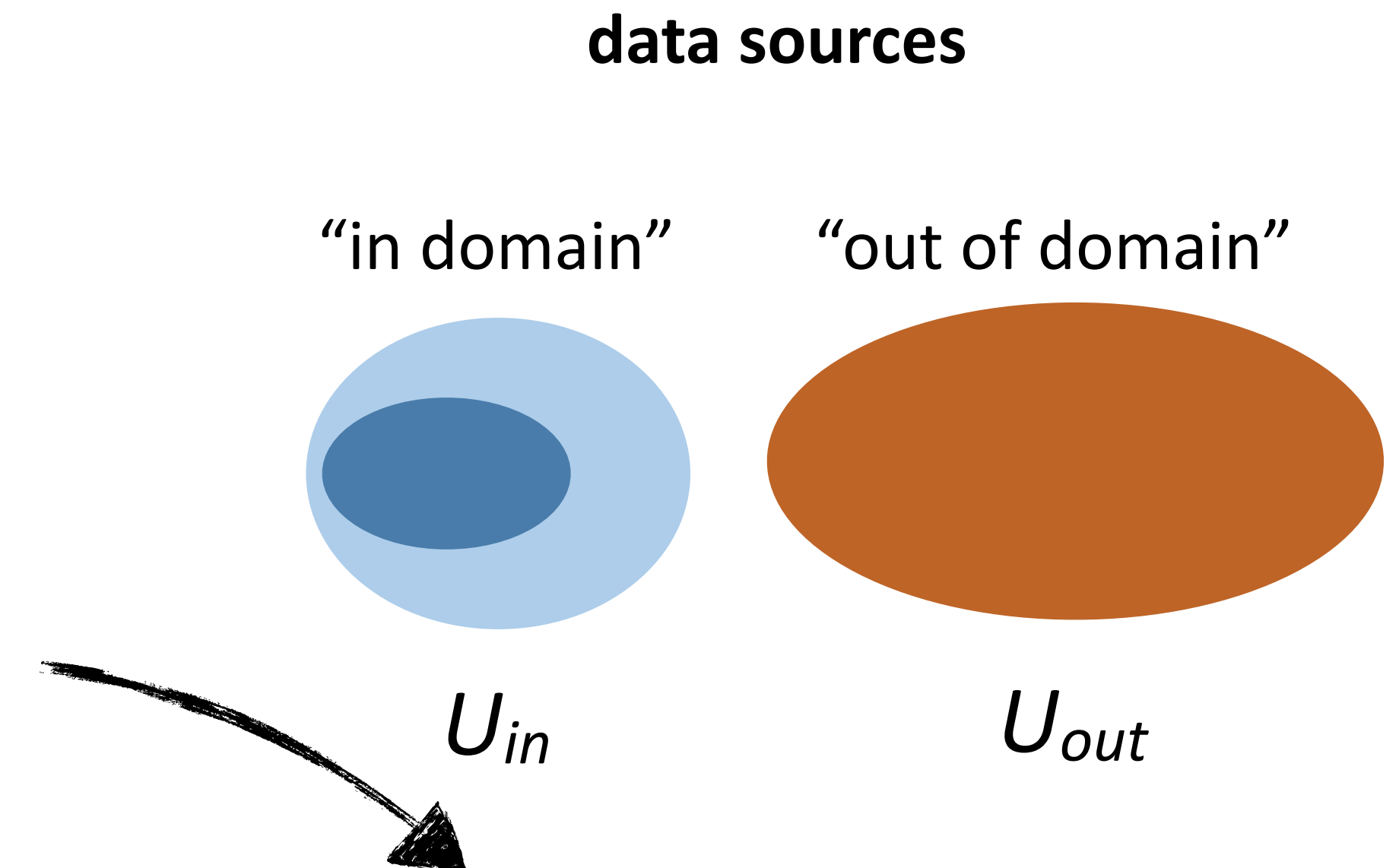
“Small” domain shifts can impact performance

- resolution, size/pose/class, novel classes

Self/semi-supervised learning is brittle in fine-grained domains

- difficult task, long-tailed data

Need “guardrails” against biased data



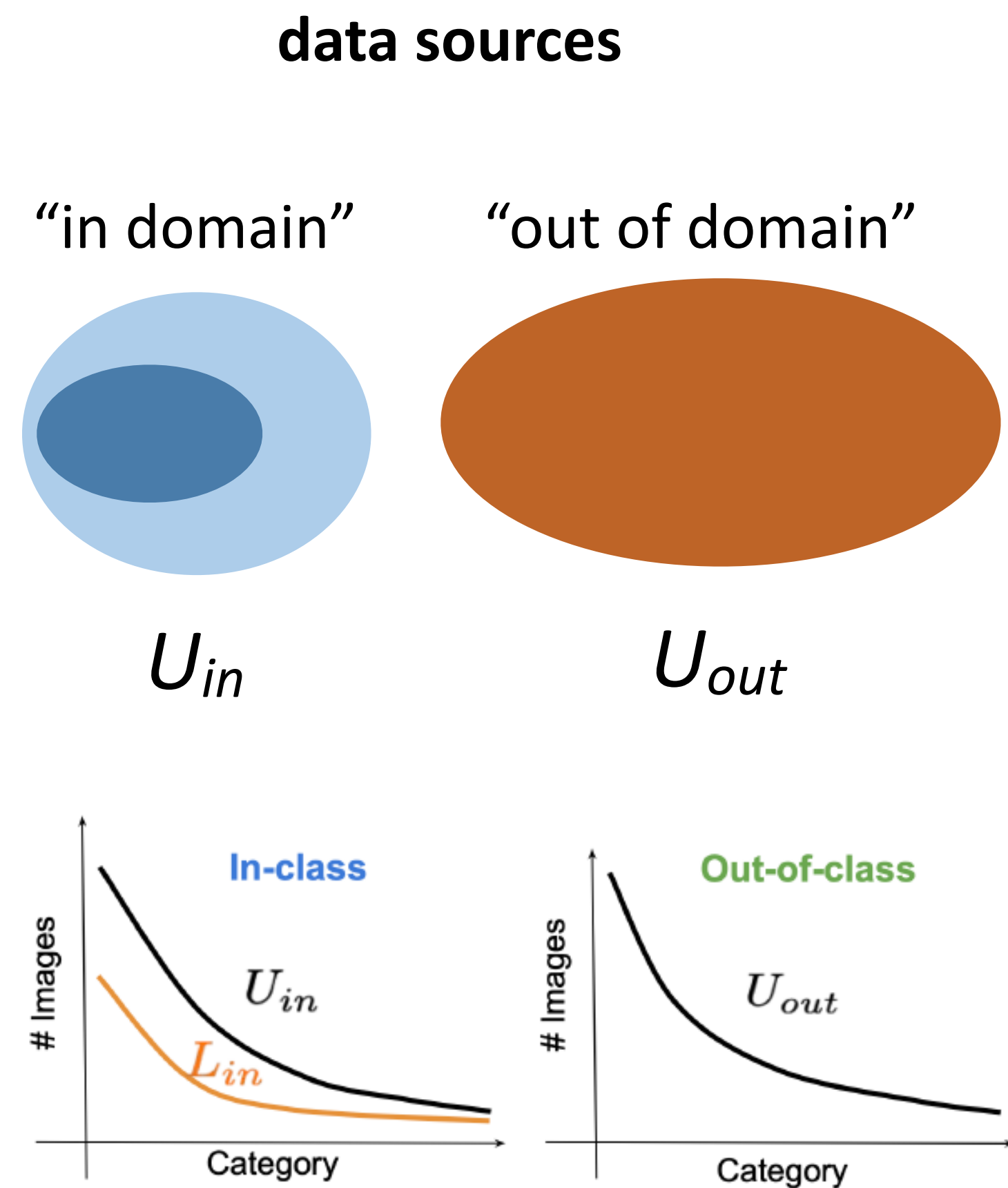
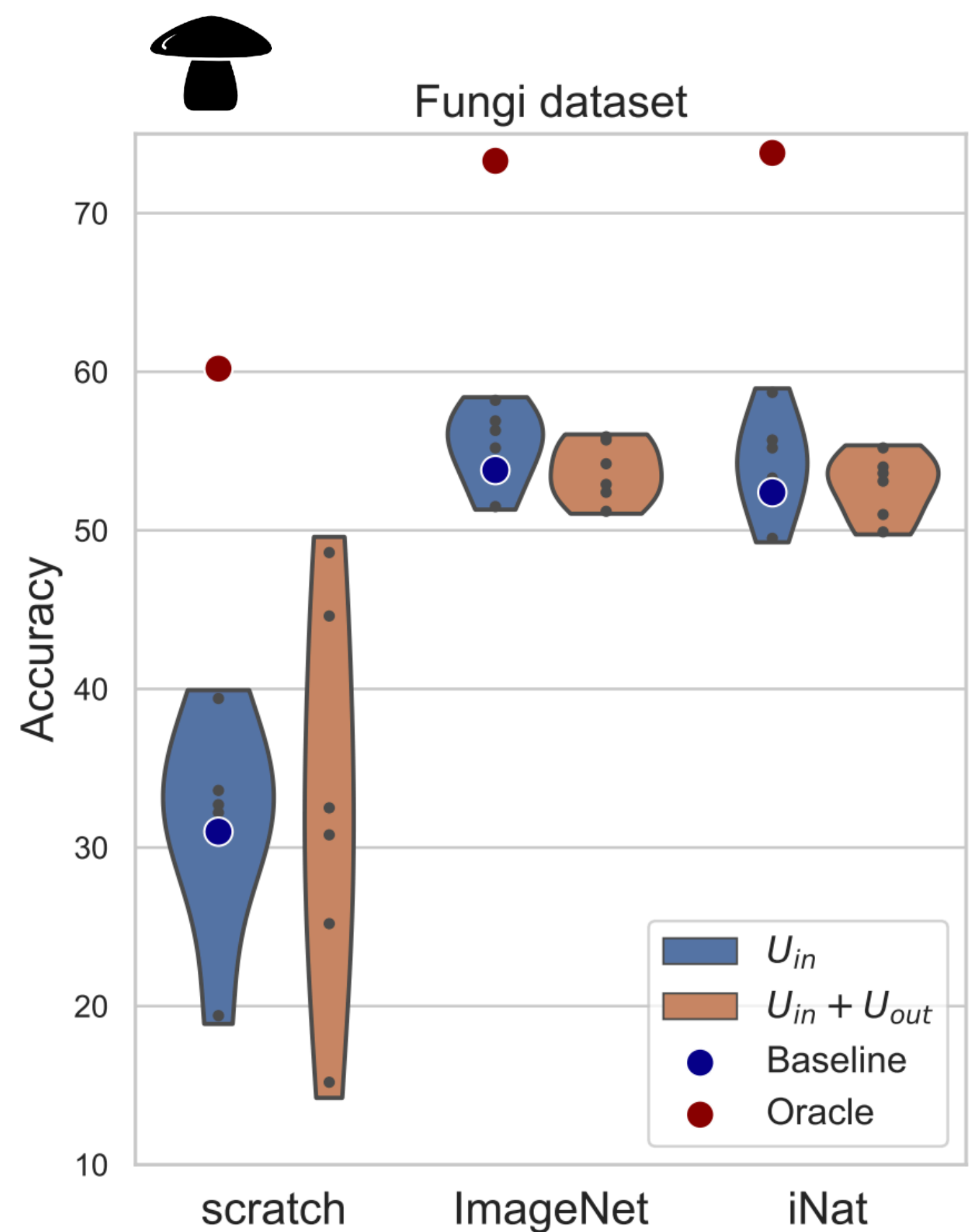
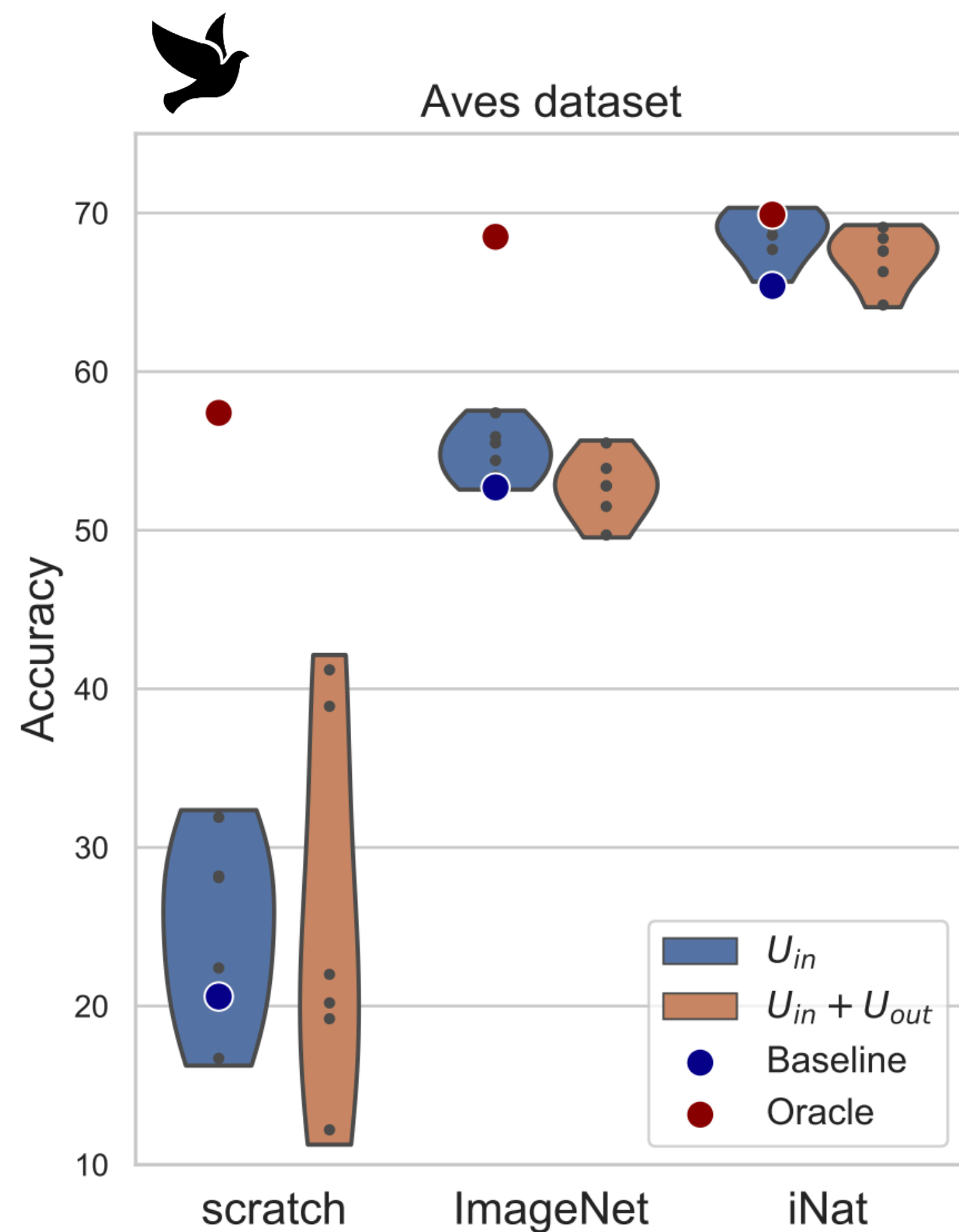
**When Does Contrastive Visual Representation Learning Work?**

Elijah Cole<sup>1</sup>   Xuan Yang<sup>2</sup>   Kimberly Wilber<sup>2</sup>   Oisín Mac Aodha<sup>3,4</sup>   Serge Belongie<sup>5</sup>  
<sup>1</sup>Caltech   <sup>2</sup>Google   <sup>3</sup>University of Edinburgh   <sup>4</sup>Alan Turing Institute   <sup>5</sup>University of Copenhagen

**When Does Self-supervision Improve  
Few-shot Learning?**

Jong-Chyi Su<sup>1</sup>    Subhransu Maji<sup>1</sup>    Bharath Hariharan<sup>2</sup> 

# How robust is semi-supervised learning?



# More pointers on semi-supervised learning

- Vast literature both in terms of theory and applications
- Other methods:
  - **Entropy minimization:** adds a loss that encourages the neural network model to make high confidence predictions (minimize “entropy”) on all unlabeled samples
  - [Mean Teacher](#), FixMatch, NoisyStudent, ...
  - Combine with methods to detect “out of domain” data

# Today's Class

- Recap
  - Supervised vs Unsupervised Learning
  - Why not always label data?
- Semi-supervised Learning
  - Concepts
  - Example: pseudo-labels / self-training
- Self-supervised Learning
  - Concepts
  - Pretext tasks
  - Contrastive Learning

# Self-supervised learning: Outline

---

- Data prediction
  - Colorization
- Transformation prediction
  - Context prediction, jigsaw puzzle solving, rotation prediction
  - “Siamese” methods
    - Contrastive methods
    - Non-contrastive methods
- Self-supervision beyond still images
  - 3D, audio, video, language

# Self-supervision as data prediction

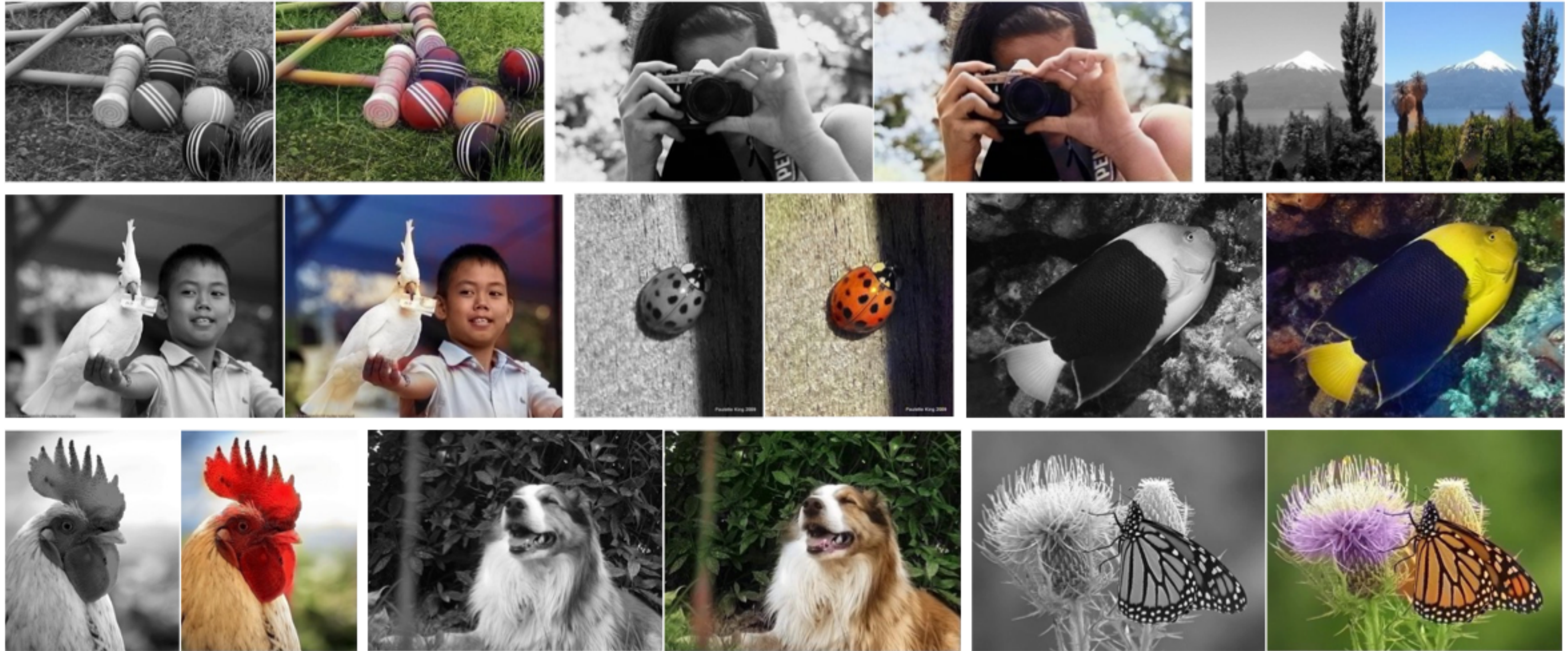
---



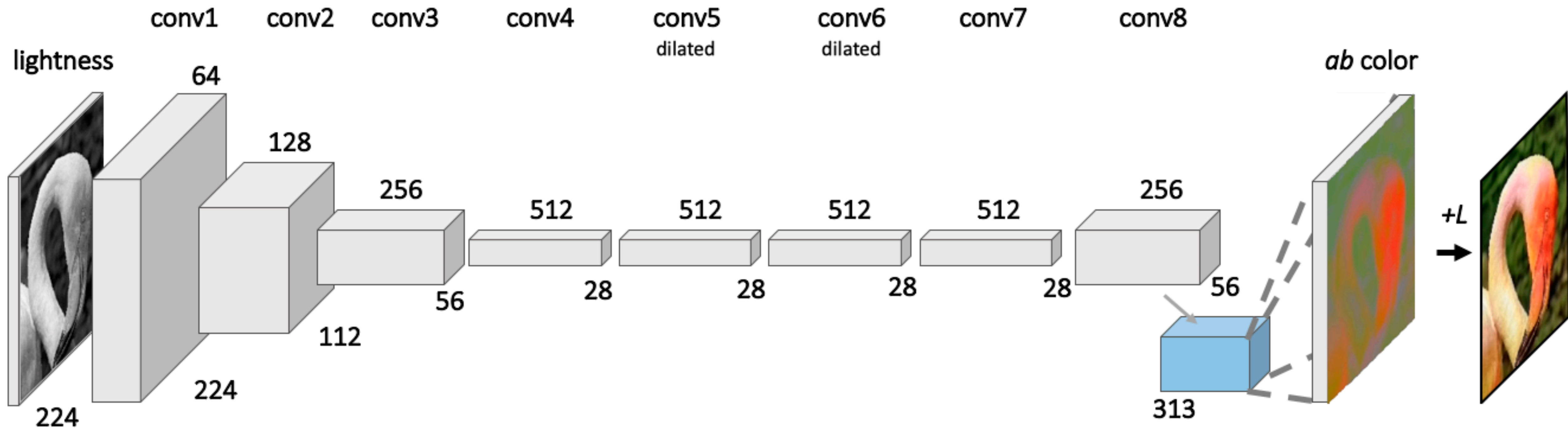
- Colorization
- Inpainting
- Future prediction
- ...

# Colorization

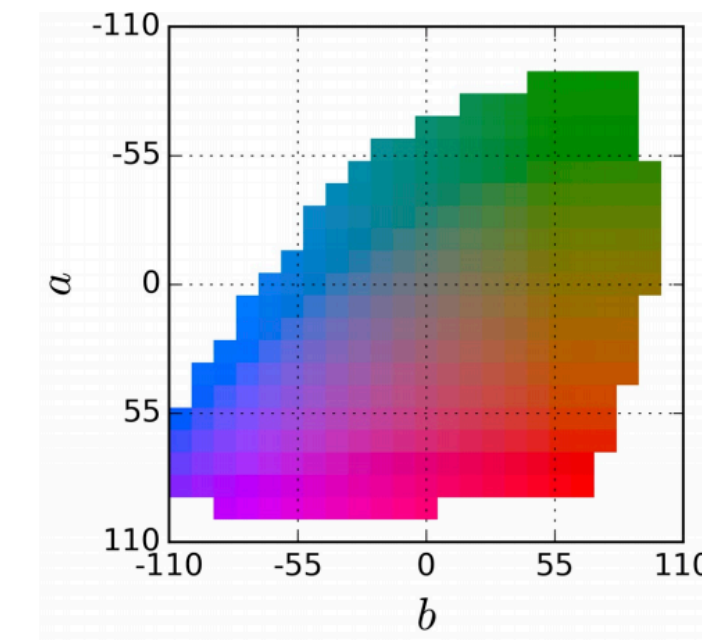
---



# Colorization: Architecture



At each spatial location, predict probability distribution over 313 quantized (a,b) values



# Self-supervised learning: Outline

---

- Data prediction
  - Colorization
- Transformation prediction

# Self-supervision by transformation prediction

---



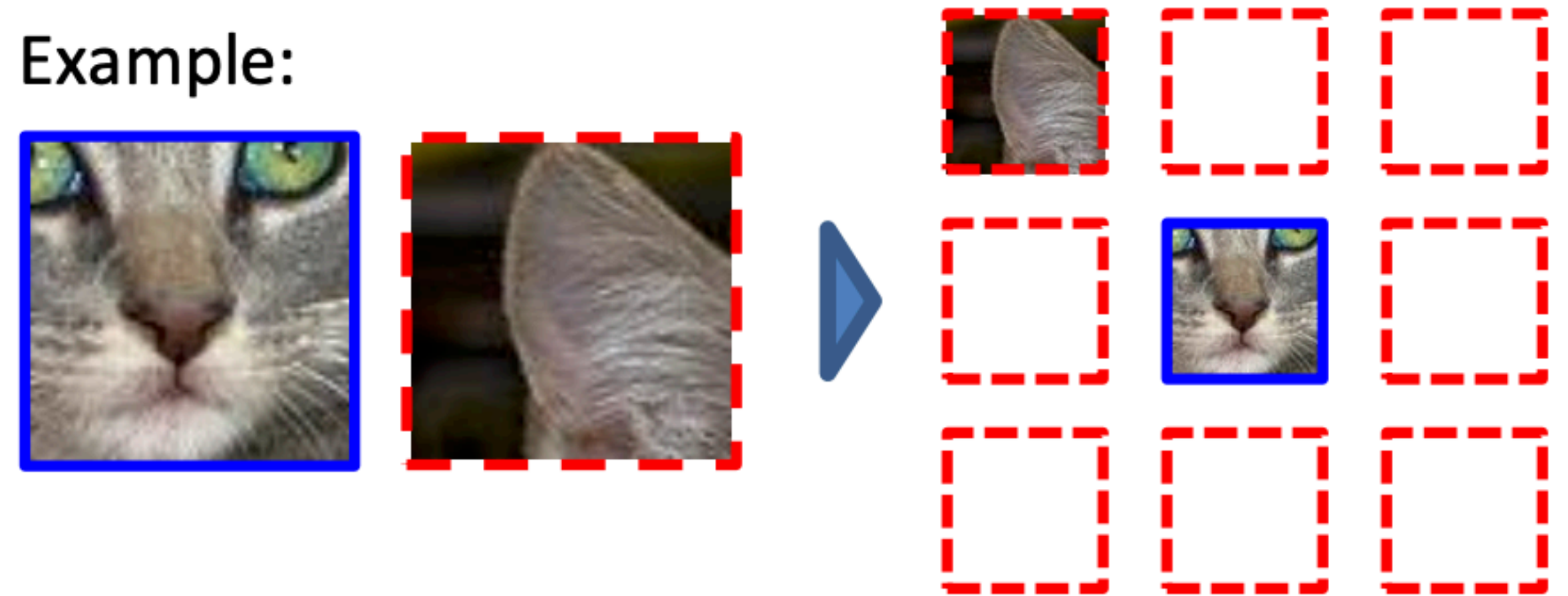
- Context prediction
- Jigsaw puzzle solving
- Rotation prediction

# Context prediction

---

- *Pretext task*: randomly sample a patch and one of 8 neighbors
- Guess the spatial relationship between the patches

Example:

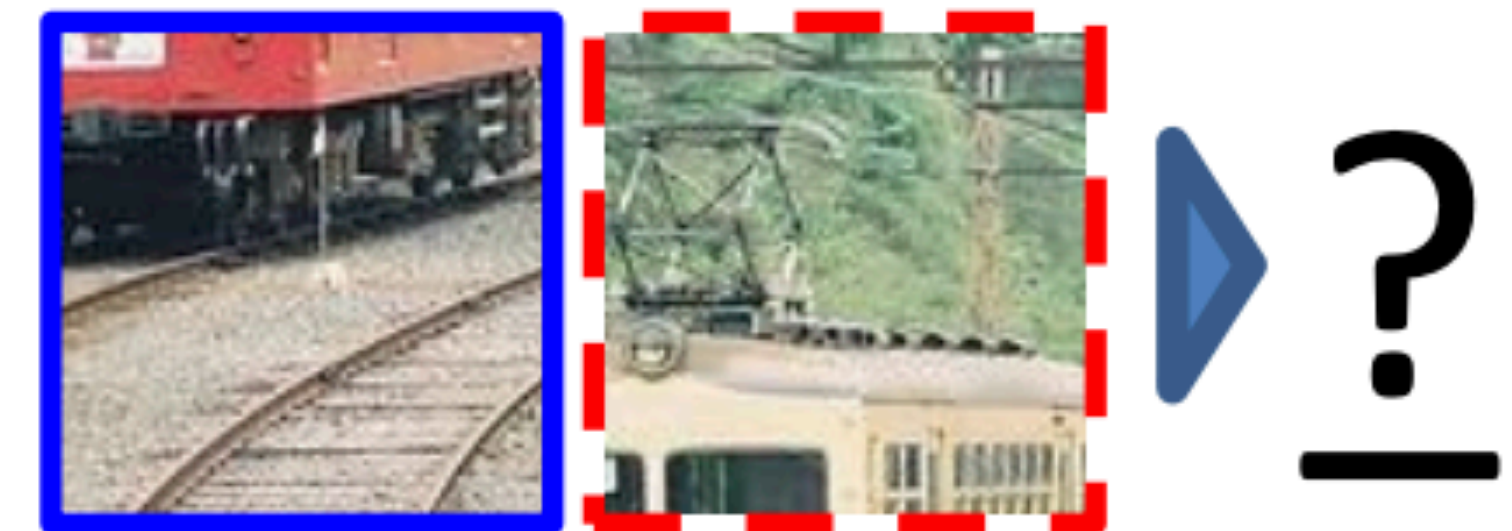


Question 1:



**A: Bottom right**

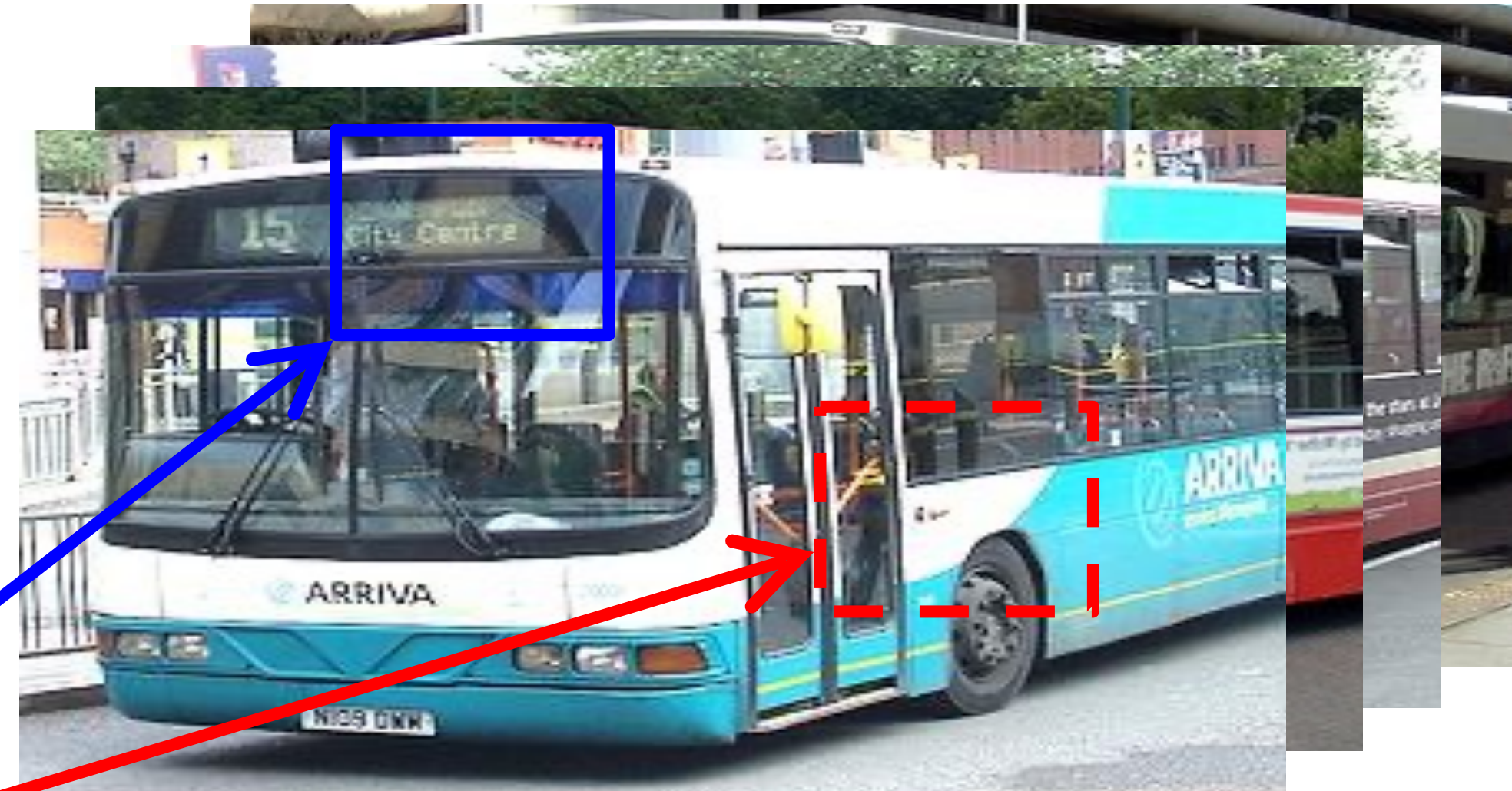
Question 2:



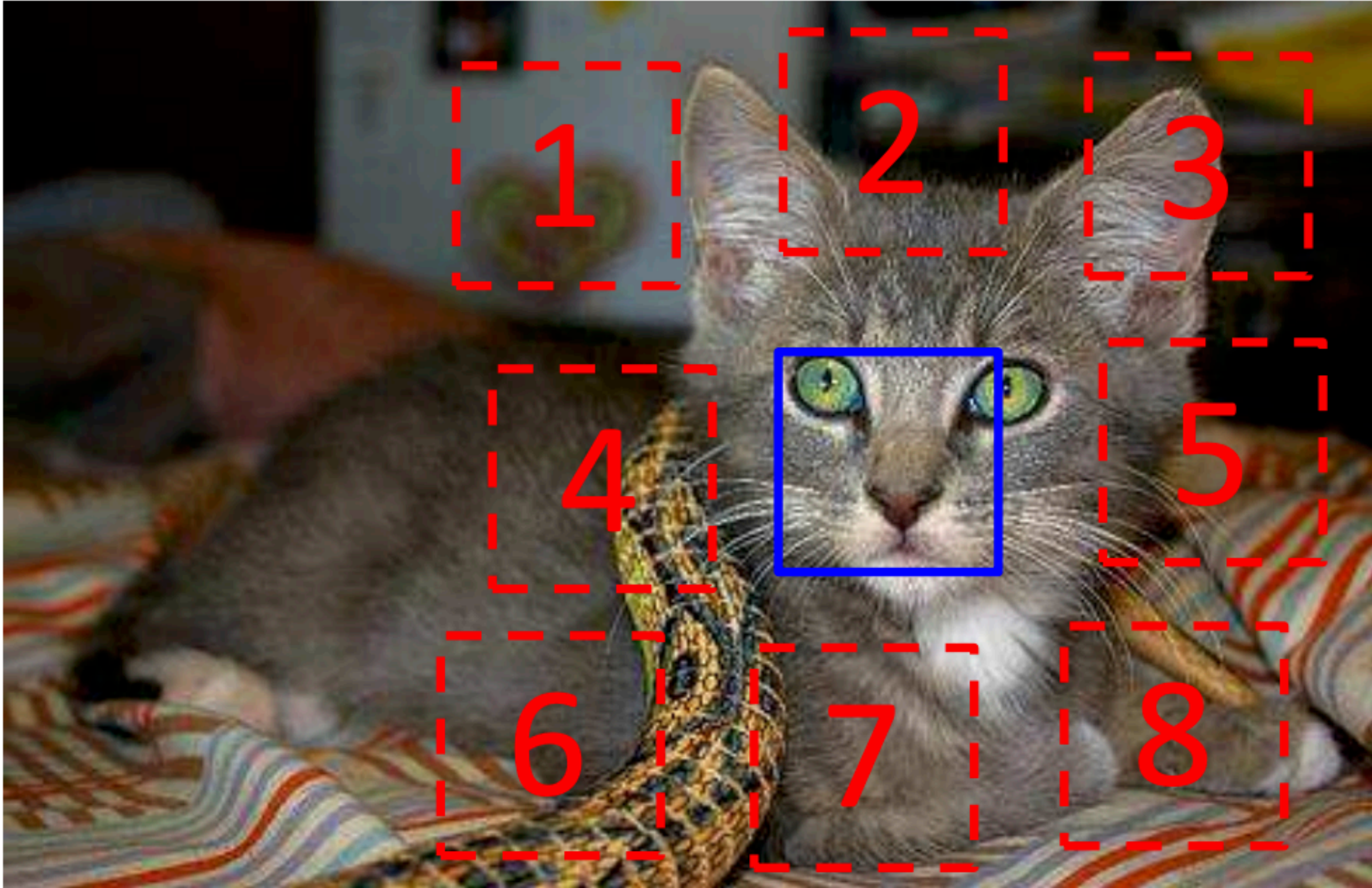
**A: Top center**

# Context prediction: Semantics from a non-semantic task

---

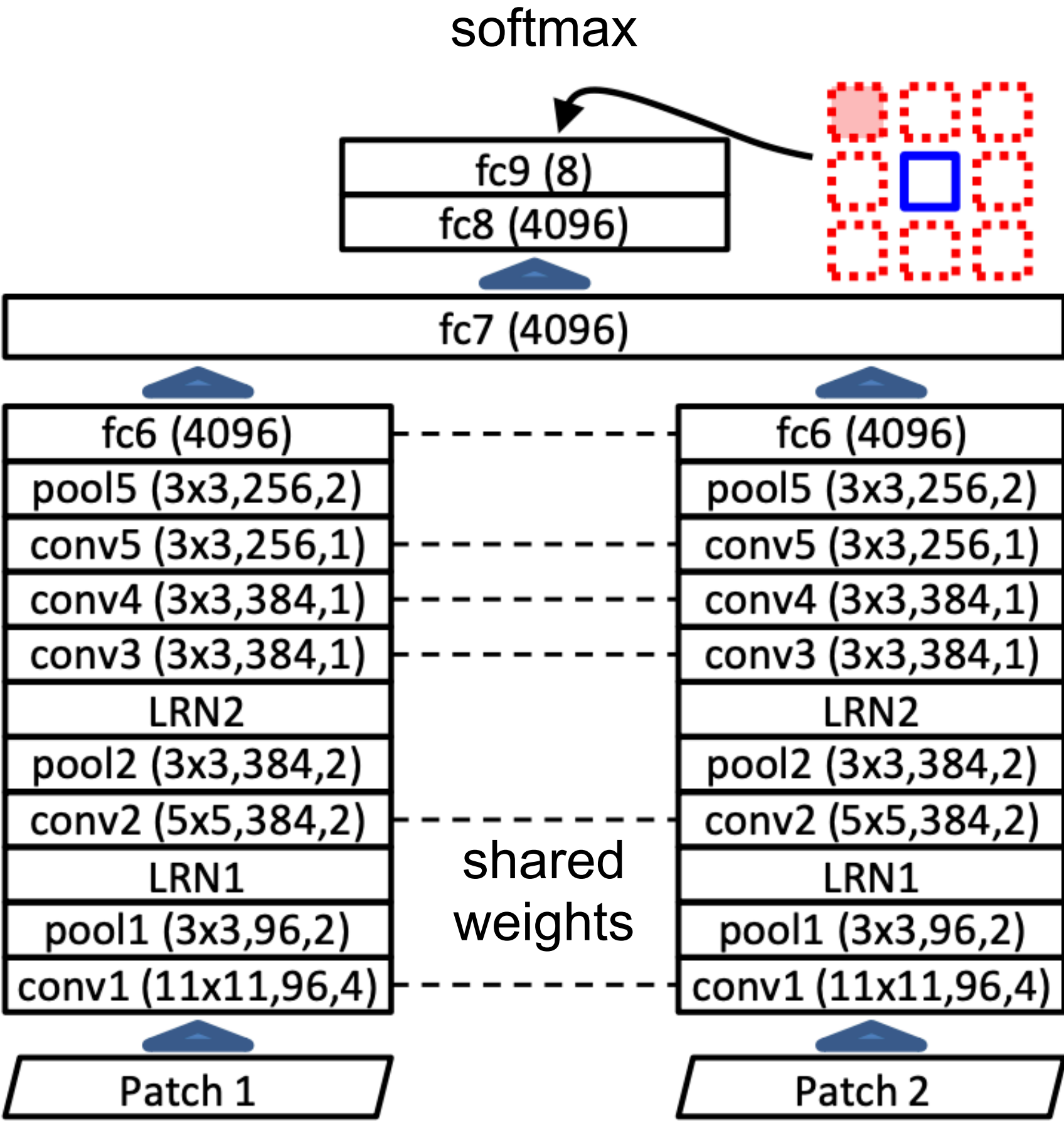


# Context prediction: Details



Prevent “cheating”: sample patches with gaps, pre-process to overcome chromatic aberration

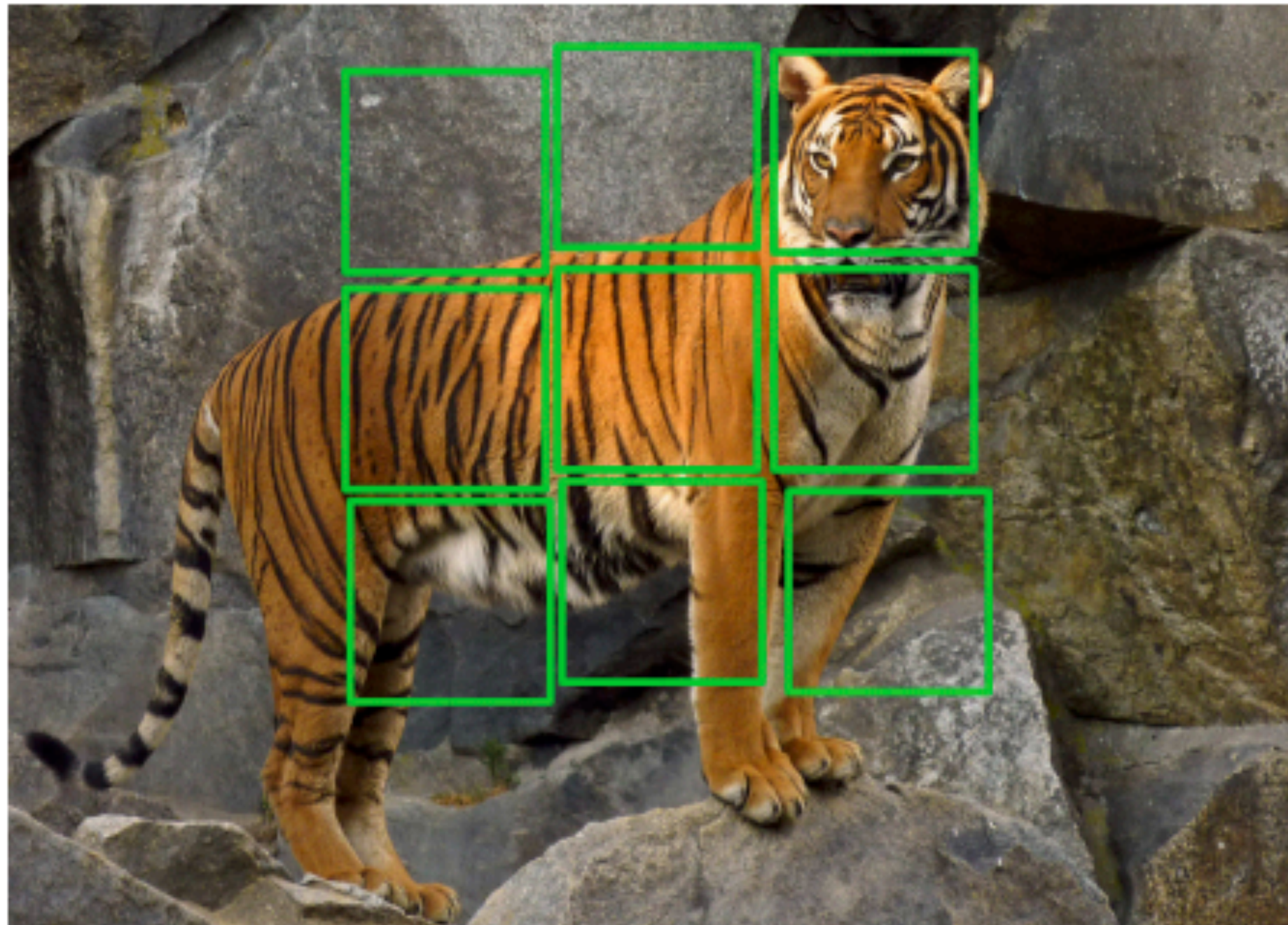
## AlexNet-like architecture



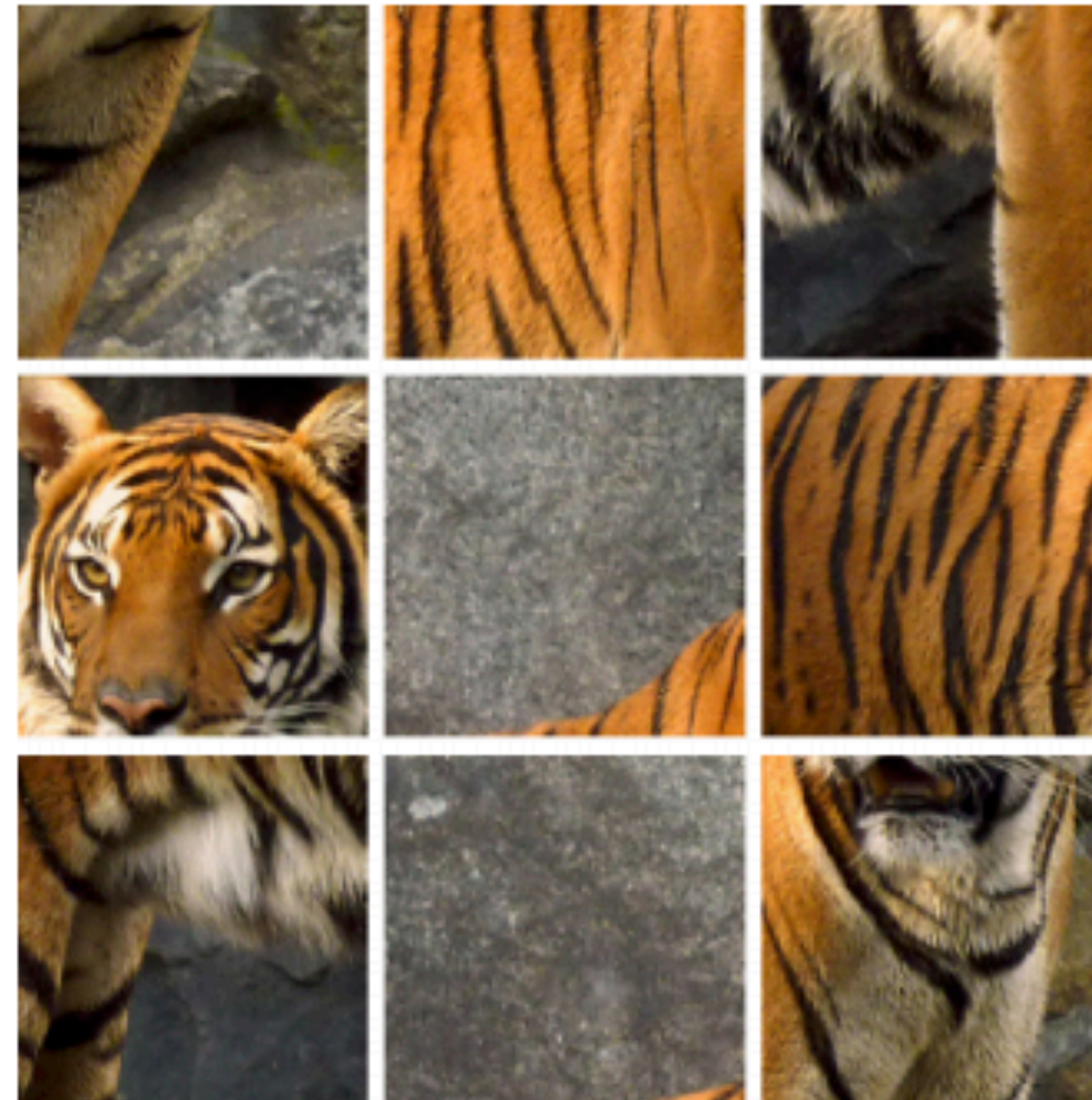
# Jigsaw puzzle solving

---

Crop out tiles



Shuffle

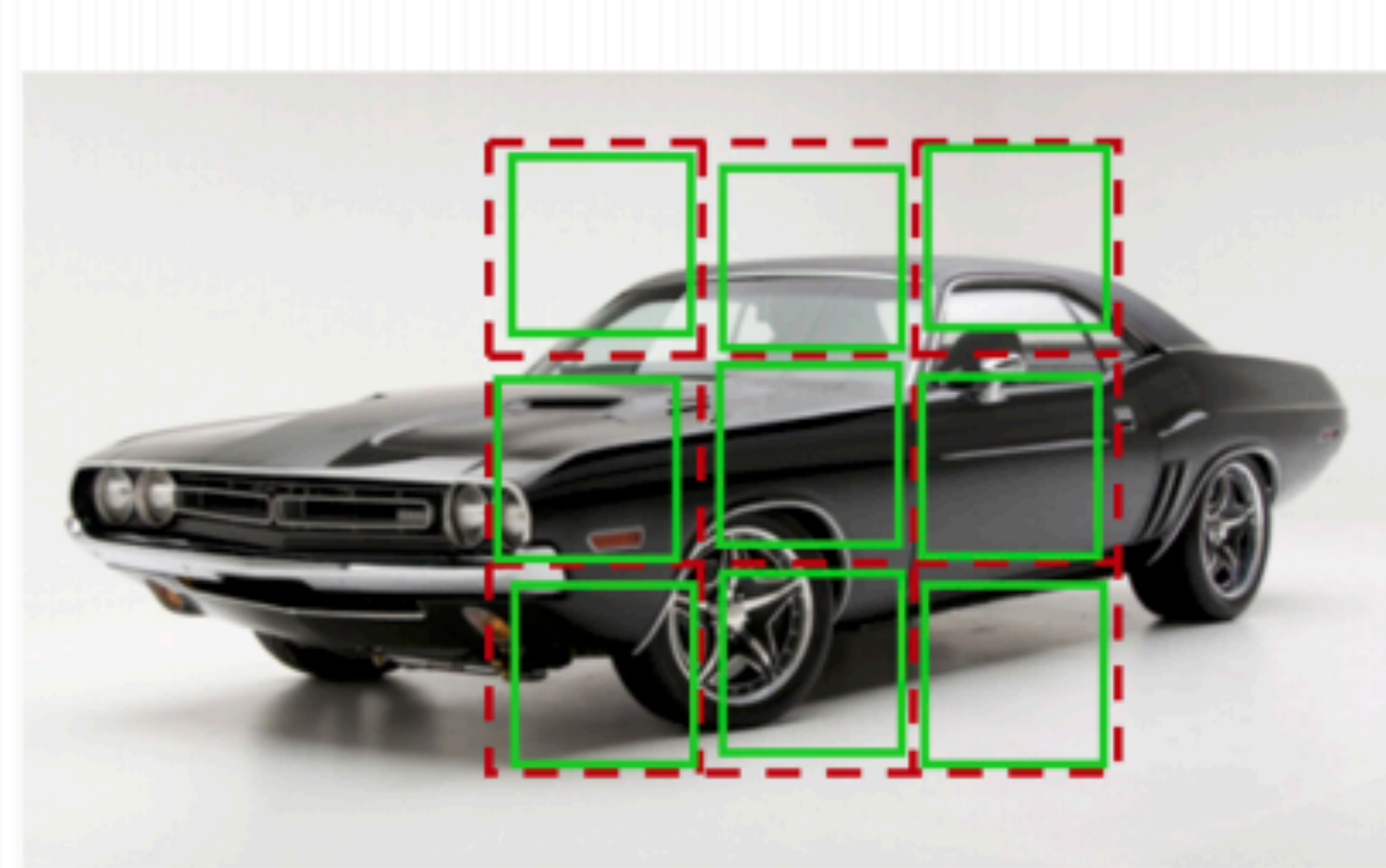


Pretext task: reassemble



Claim: jigsaw solving is easier than context prediction, trains faster, transfers better

# Jigsaw puzzle solving: Details

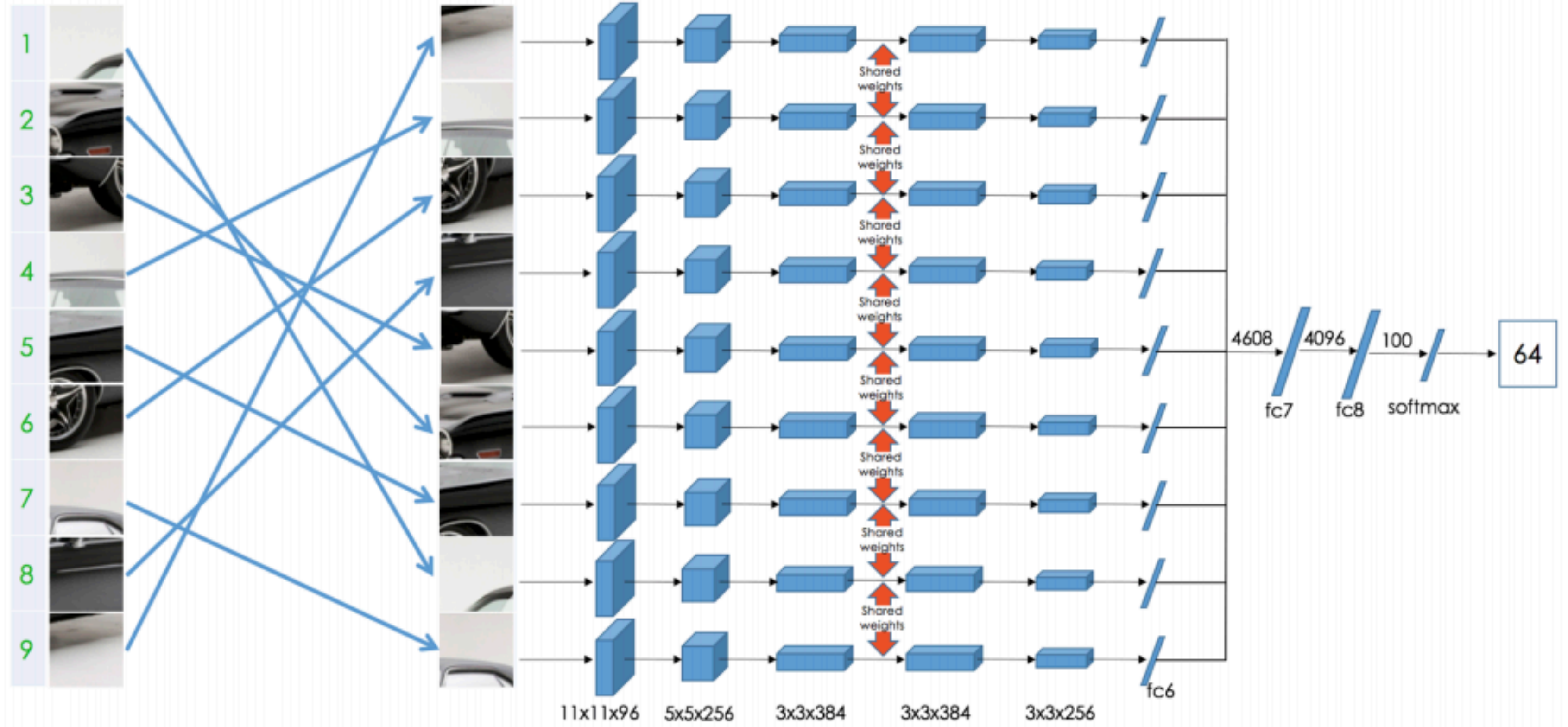


Permutation Set

| index | permutation       |
|-------|-------------------|
| 64    | 9,4,6,8,3,2,5,1,7 |

Reorder patches according to the selected permutation

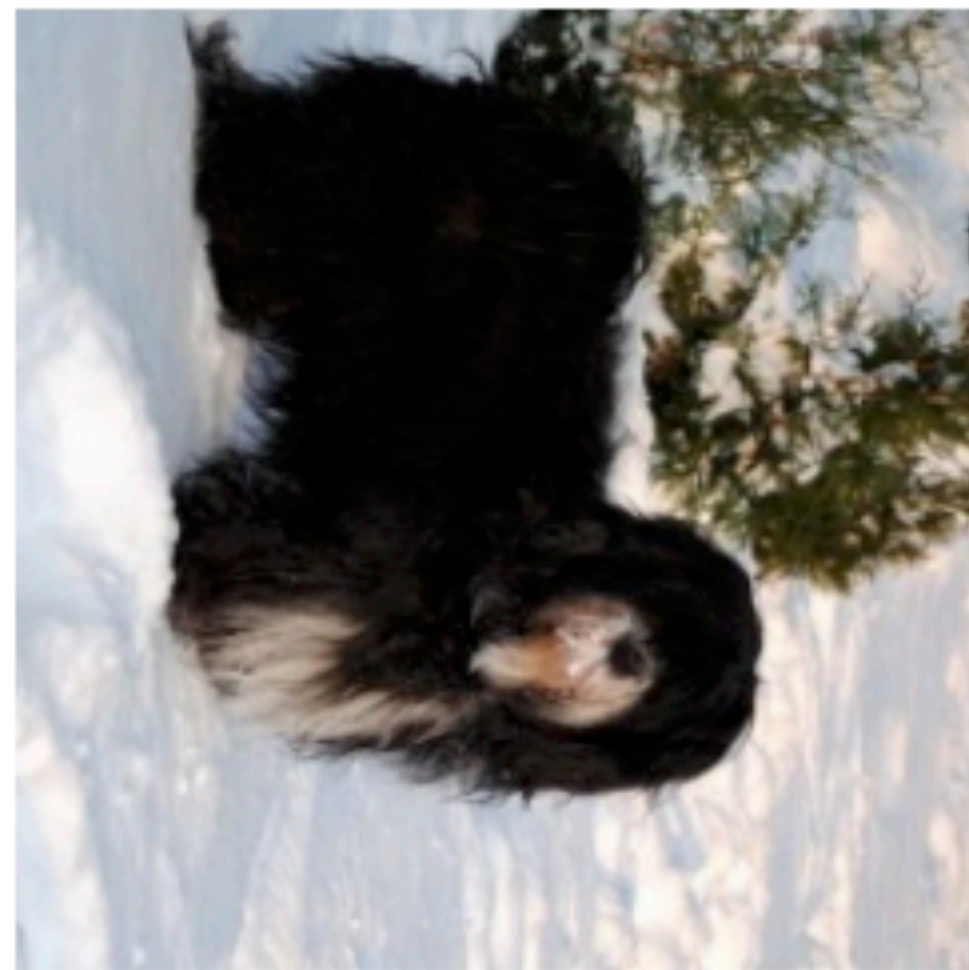
Predetermined set of 1000 permutations (out of 362,880 possible)



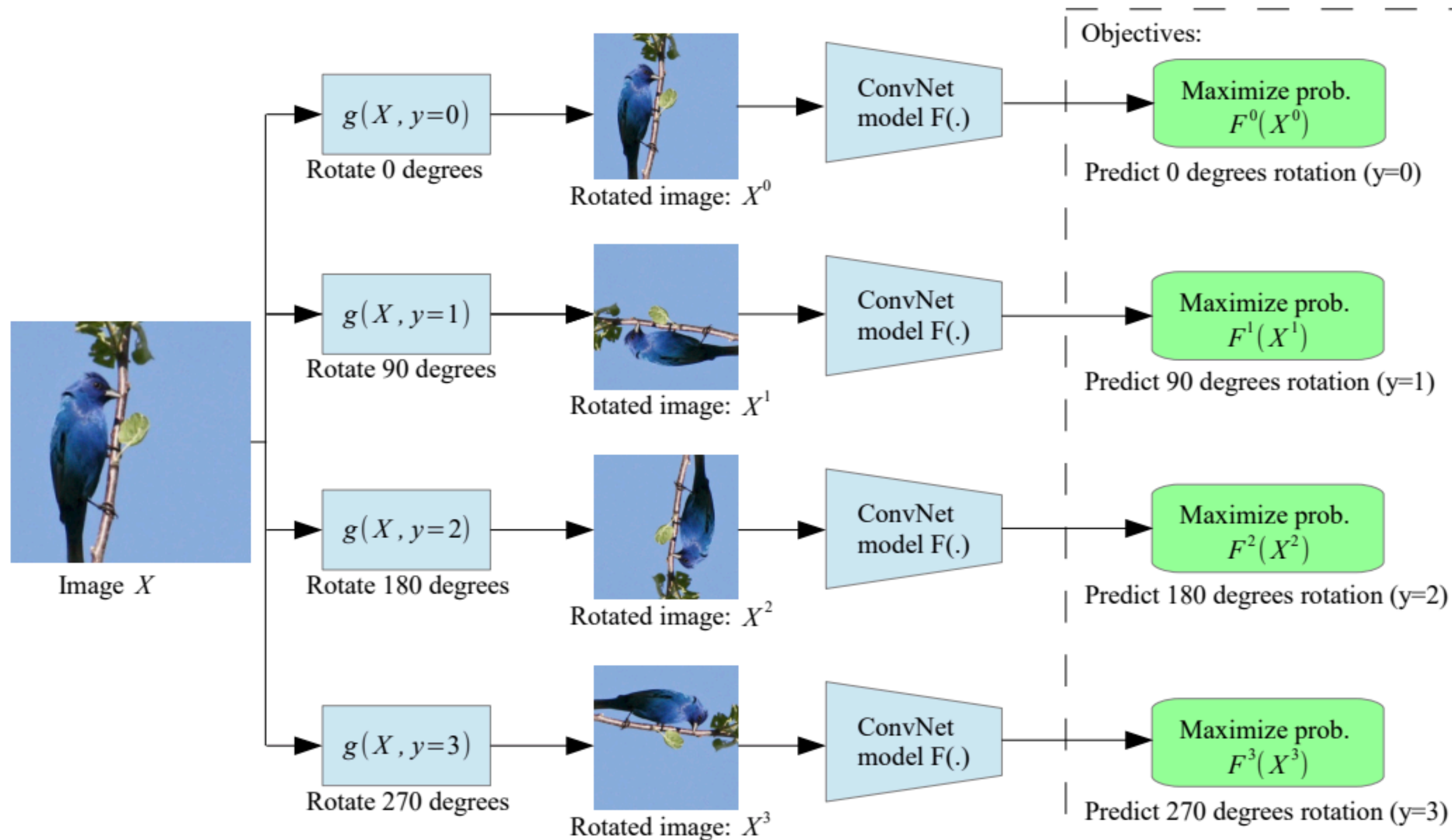
# Rotation prediction

---

- Pretext task: recognize image rotation (0, 90, 180, 270 degrees)



# Rotation prediction



During training, feed in all four rotated versions of an image in the same mini-batch

# PASCAL VOC Transfer Results

---

| Method                | Classification | Detection (mAP) | Segmentation (mIoU) |
|-----------------------|----------------|-----------------|---------------------|
| Supervised (ImageNet) | <b>79.9</b>    | <b>56.8</b>     | <b>48.0</b>         |
| Colorization          | 65.6           | 46.9            | 35.6                |
| Context               | 65.3           | 51.1            |                     |
| Jigsaw                | 67.6           | 53.2            | 37.6                |
| Rotation              | 73.0           | 54.4            | 39.1                |

# Self-supervised learning: Outline

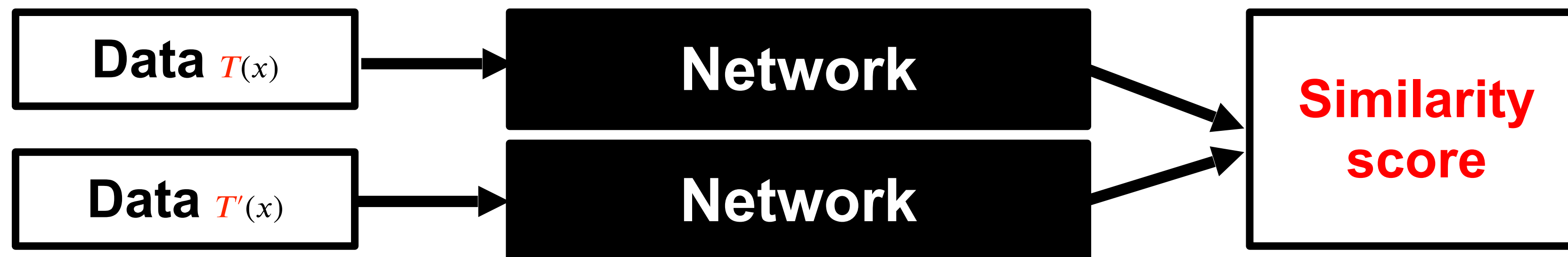
---

- Data prediction
  - Colorization
- Transformation prediction
  - Context prediction, jigsaw puzzle solving, rotation prediction
- **“Siamese” methods**

# “Siamese” methods

---

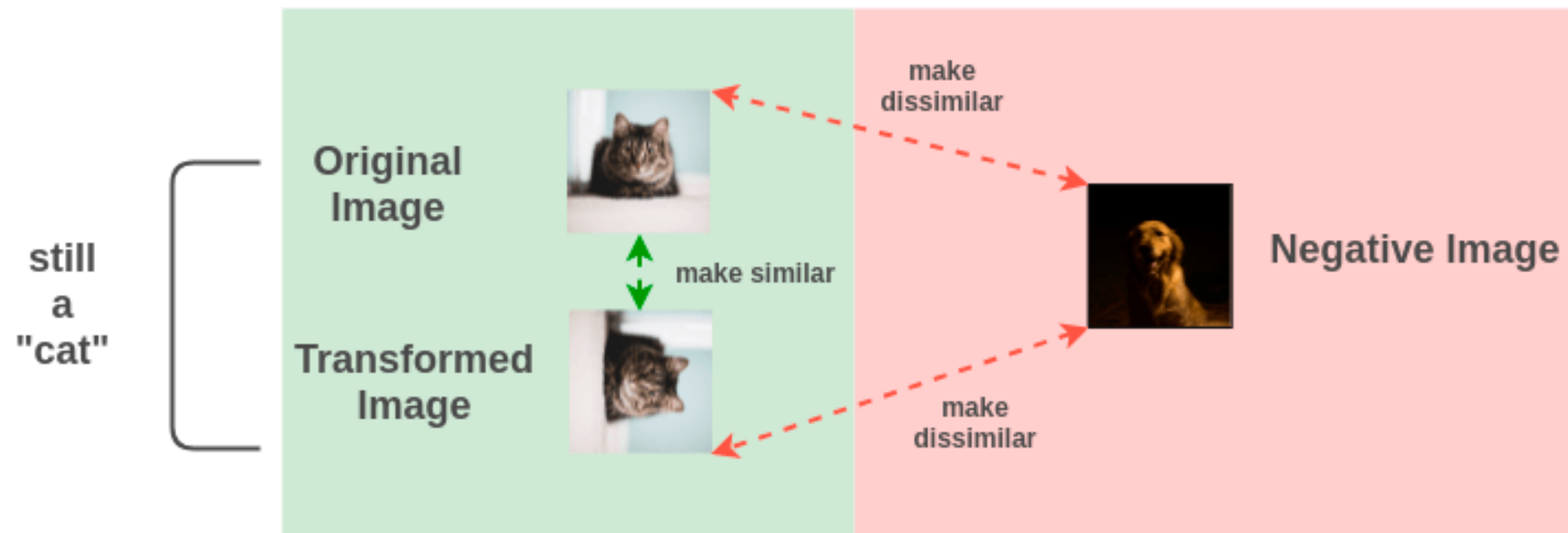
- Extract representations from two transformed versions of a data point, encourage these representations to be similar (or to have other desirable properties)
- **Contrastive methods:** train using both positive (similar) and negative (dissimilar) pairs
- **Non-contrastive methods:** train with only positive examples



# Contrastive methods

---

- Encourage representations of transformed versions of the same image to be the same and different images to be different



# Contrastive loss formulation

---

- Given:
  - Query point  $x$
  - Positive sample  $x^+$ : version of  $x$  subjected to a random transformation or augmentation (cropping, rotation, color change, etc.)
  - Negative samples  $x^-$



$x$



$x^+$



$x^-$

# Contrastive loss formulation

---

- Given: query  $x$ , positive sample  $x^+$ , negative samples  $x^-$
- Measure similarity by dot product of L2-normalized feature representations:

$$\text{sim}(x, y) = \frac{f(x)}{\|f(x)\|_2} \cdot \frac{f(y)}{\|f(y)\|_2}$$

- **Contrastive loss:** make  $x$  similar to  $x^+$ , dissimilar from  $x^-$ :

$$l(x, x^+) = -\log \frac{\exp(\text{sim}(x, x^+)/\tau)}{\exp(\text{sim}(x, x^+)/\tau) + \sum_{j=1}^N \exp(\text{sim}(x, x_j^-)/\tau)}$$

- Intuitively, this is the loss of a softmax classifier that tries to classify  $x$  as  $x^+$

# Mechanisms for obtaining negative samples

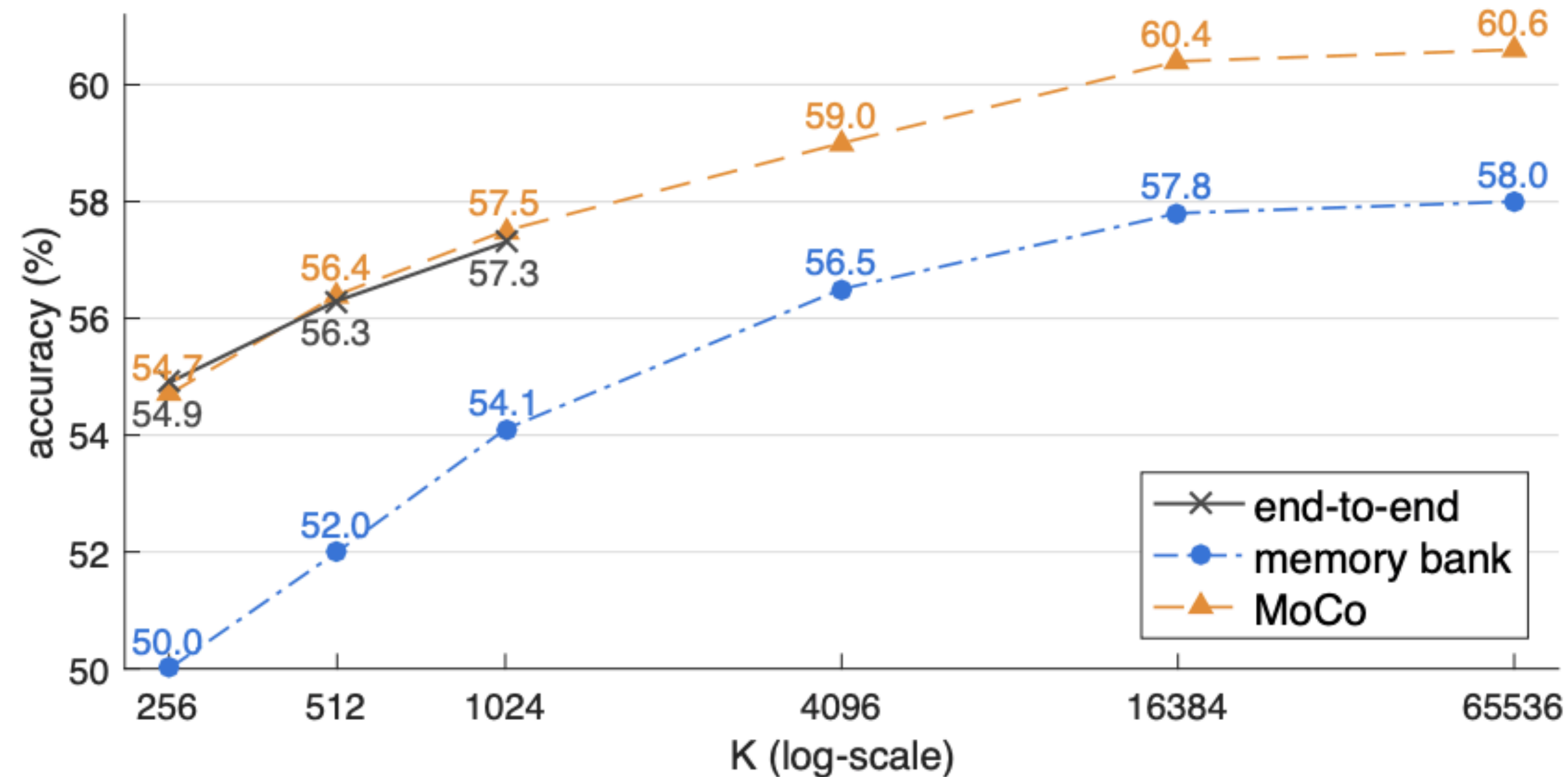
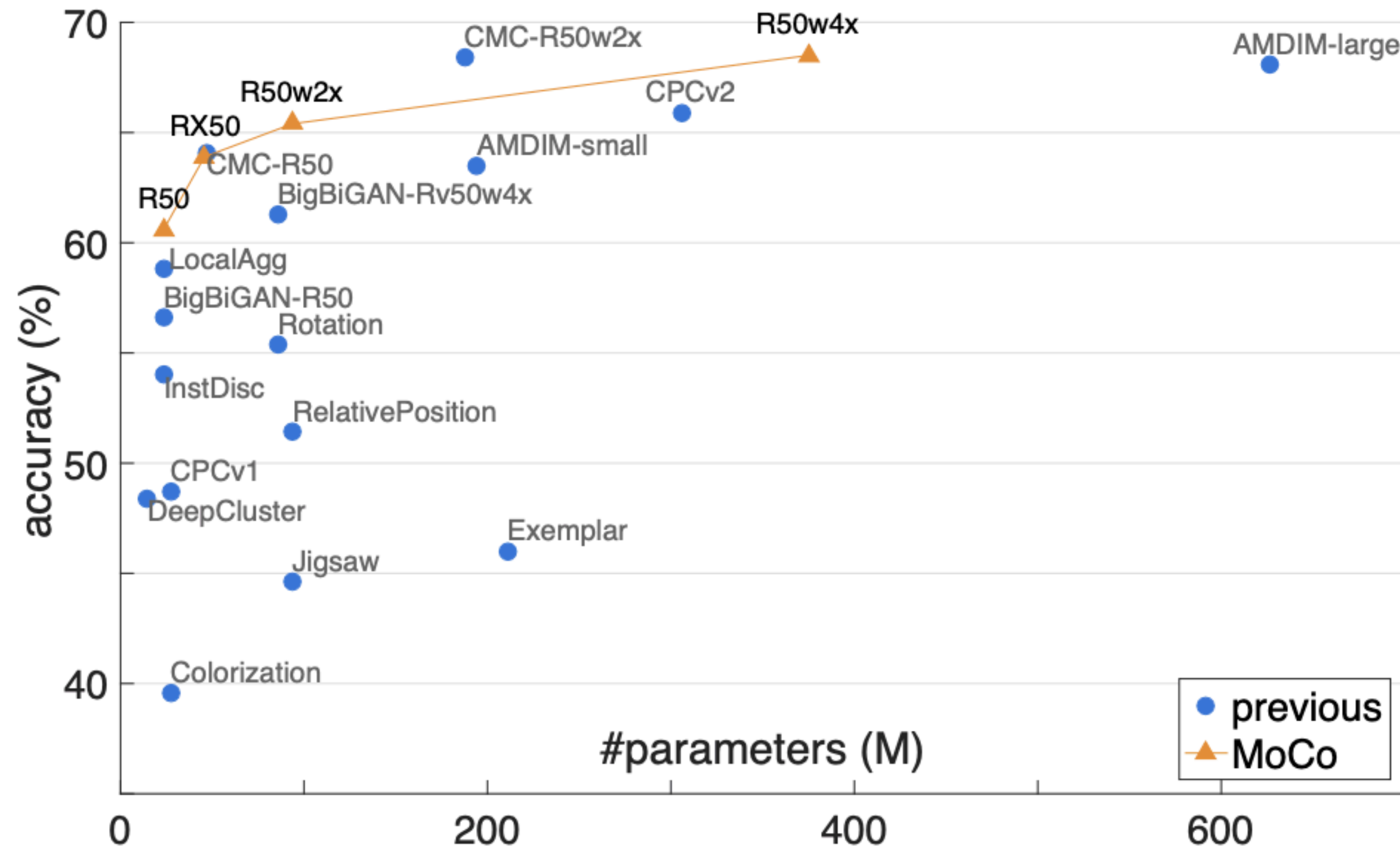


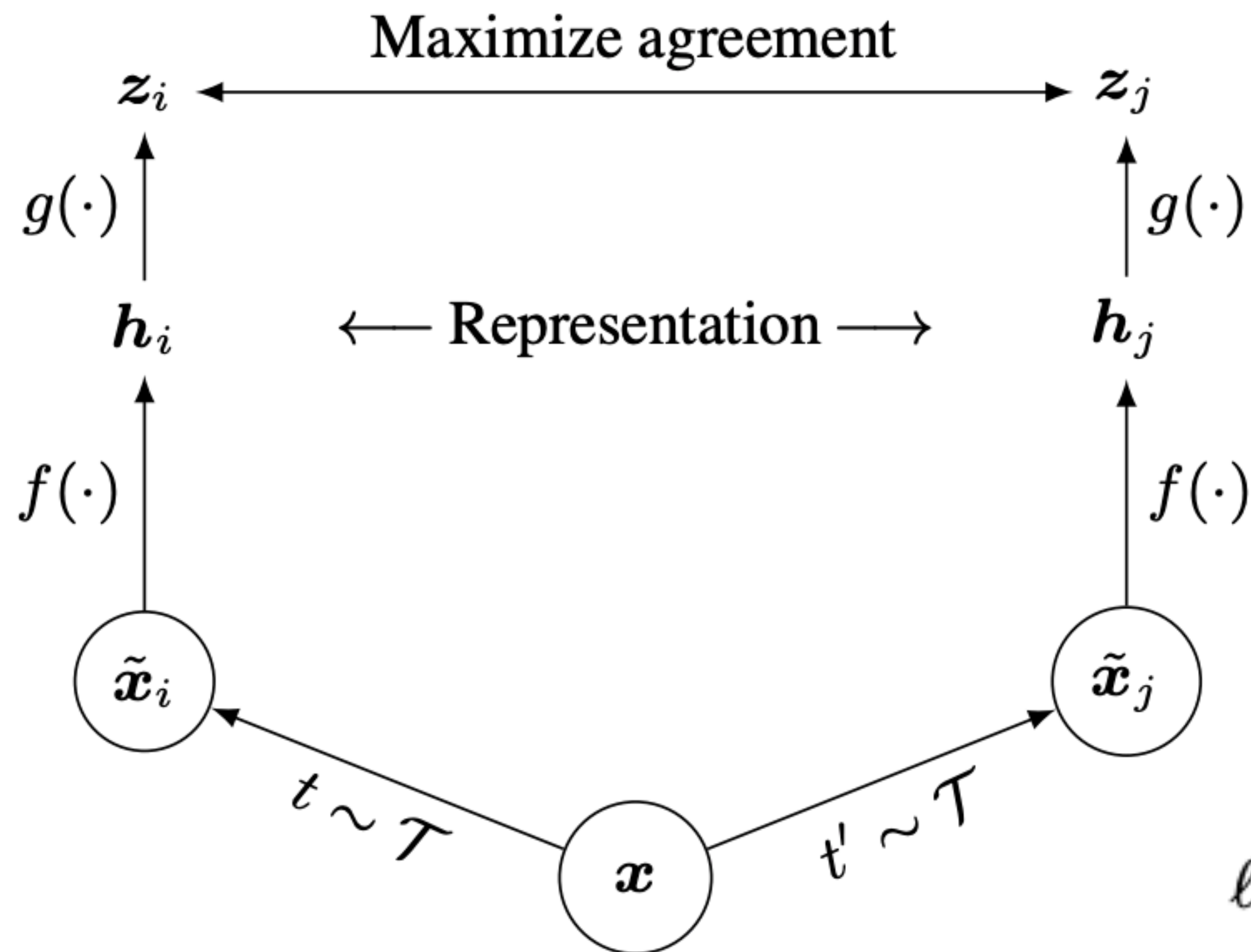
Figure 3. **Comparison of three contrastive loss mechanisms** under the ImageNet linear classification protocol. We adopt the same pretext task (Sec. 3.3) and only vary the contrastive loss mechanism (Figure 2). The number of negatives is  $K$  in memory bank and MoCo, and is  $K-1$  in end-to-end (offset by one because the positive key is in the same mini-batch). The network is ResNet-50.

# MoCo results

Comparison on linear ImageNet classification  
(supervised accuracy above 75%)



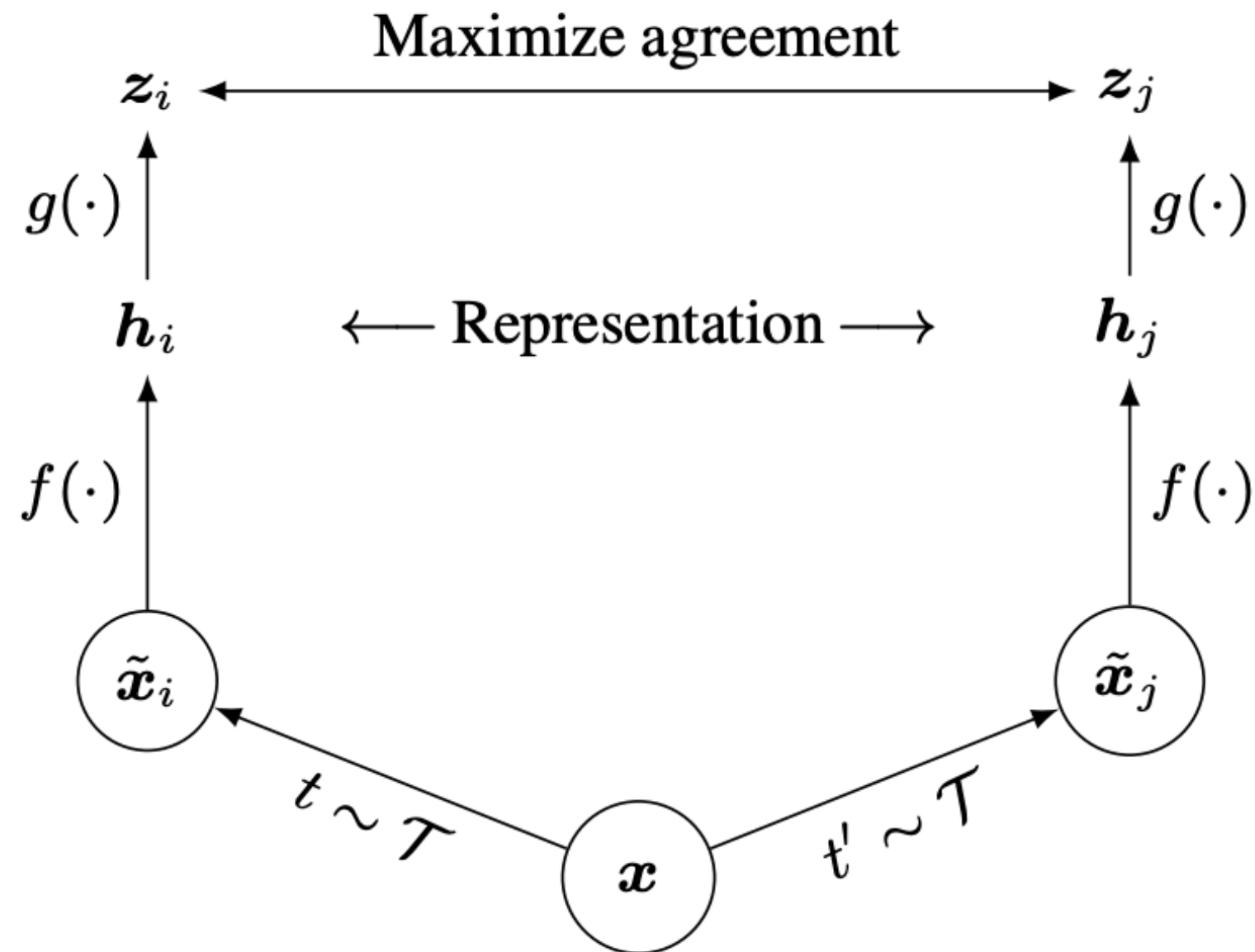
# SimCLR



- Instead of memory bank or queue, use large mini-batch size (on cloud TPU)
- Introduce nonlinear *projection* ( $g$ ) between representation ( $h$ ) and feature used for computing contrastive loss ( $z$ )

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}$$

# SimCLR



$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

---

**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .

**for** sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**

**for all**  $k \in \{1, \dots, N\}$  **do**

    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$

    # the first augmentation

$\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$

$\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation

$\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection

    # the second augmentation

$\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$

$\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation

$\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection

**end for**

**for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**

$s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity

**end for**

**define**  $\ell(i, j)$  **as**  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$

$\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$

  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$

**end for**

**return** encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$

---

# SimCLR

- Performed extensive ablation study of data augmentations
- Found that composing multiple augmentations gives the best results

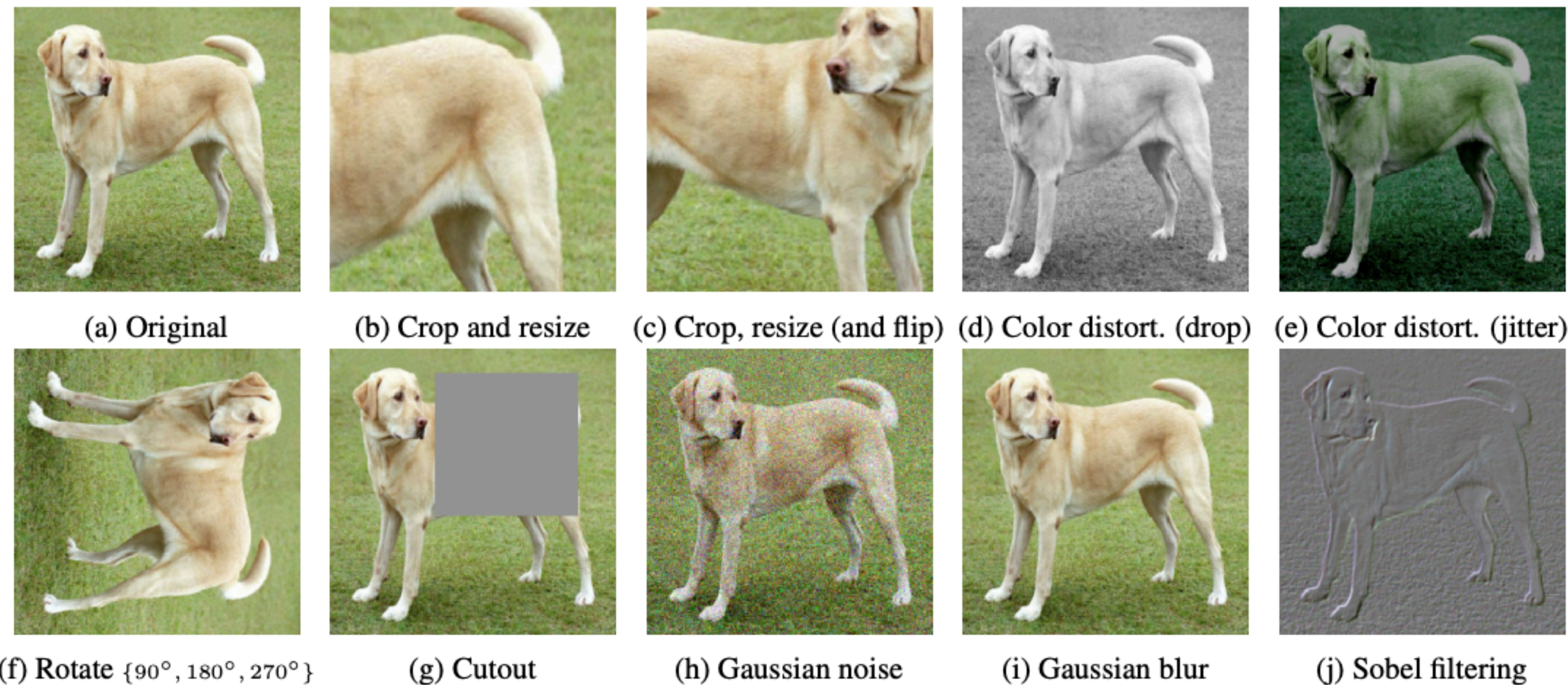


Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes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

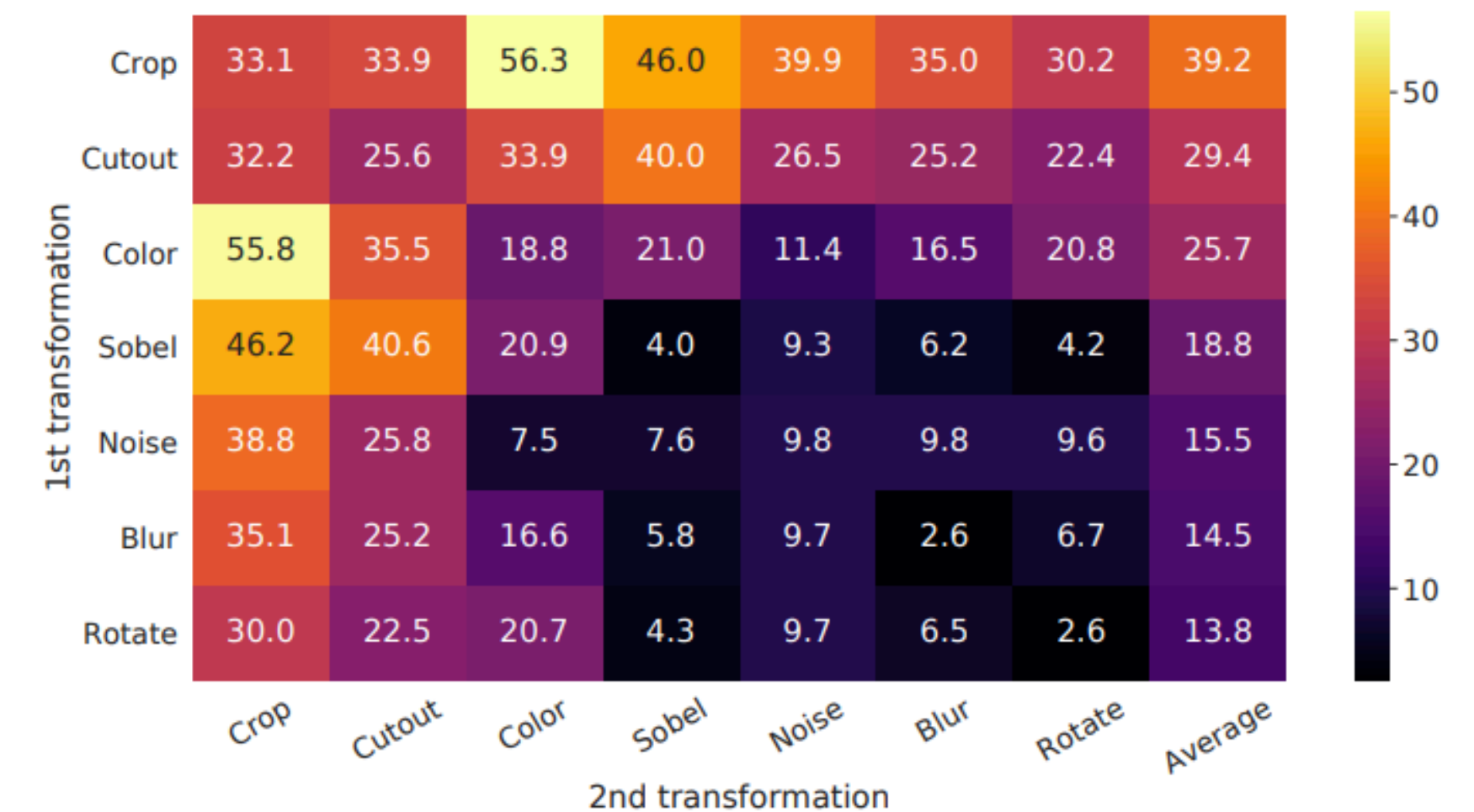
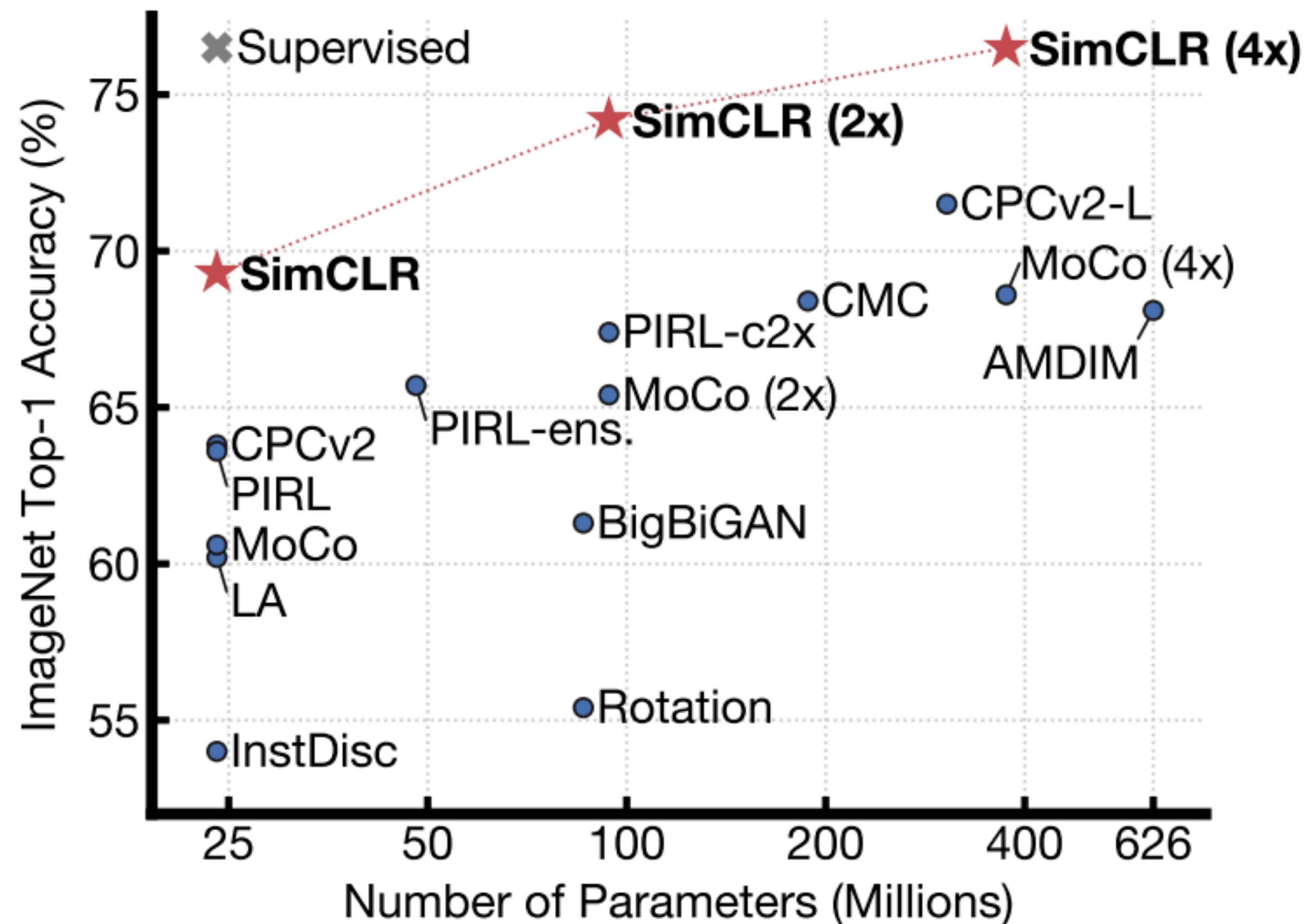


Figure 5. Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

# SimCLR: Evaluation

---



No detection evaluation

# Improved Baselines with Momentum Contrastive Learning

Xinlei Chen   Haoqi Fan   Ross Girshick   Kaiming He  
Facebook AI Research (FAIR)

## Abstract

Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR. In this note, we verify the effectiveness of two of SimCLR’s design improvements by implementing them in the MoCo framework. With simple modifications to MoCo—namely, using an MLP projection head and more data augmentation—we establish stronger baselines that outperform SimCLR and do not require large training batches. We hope this will make state-of-the-art unsupervised learning research more accessible. Code will be made public.

## 1. Introduction

Recent studies on unsupervised representation learning from images [16, 13, 8, 17, 1, 9, 15, 6, 12, 2] are converging on a central concept known as contrastive learning [5]. The

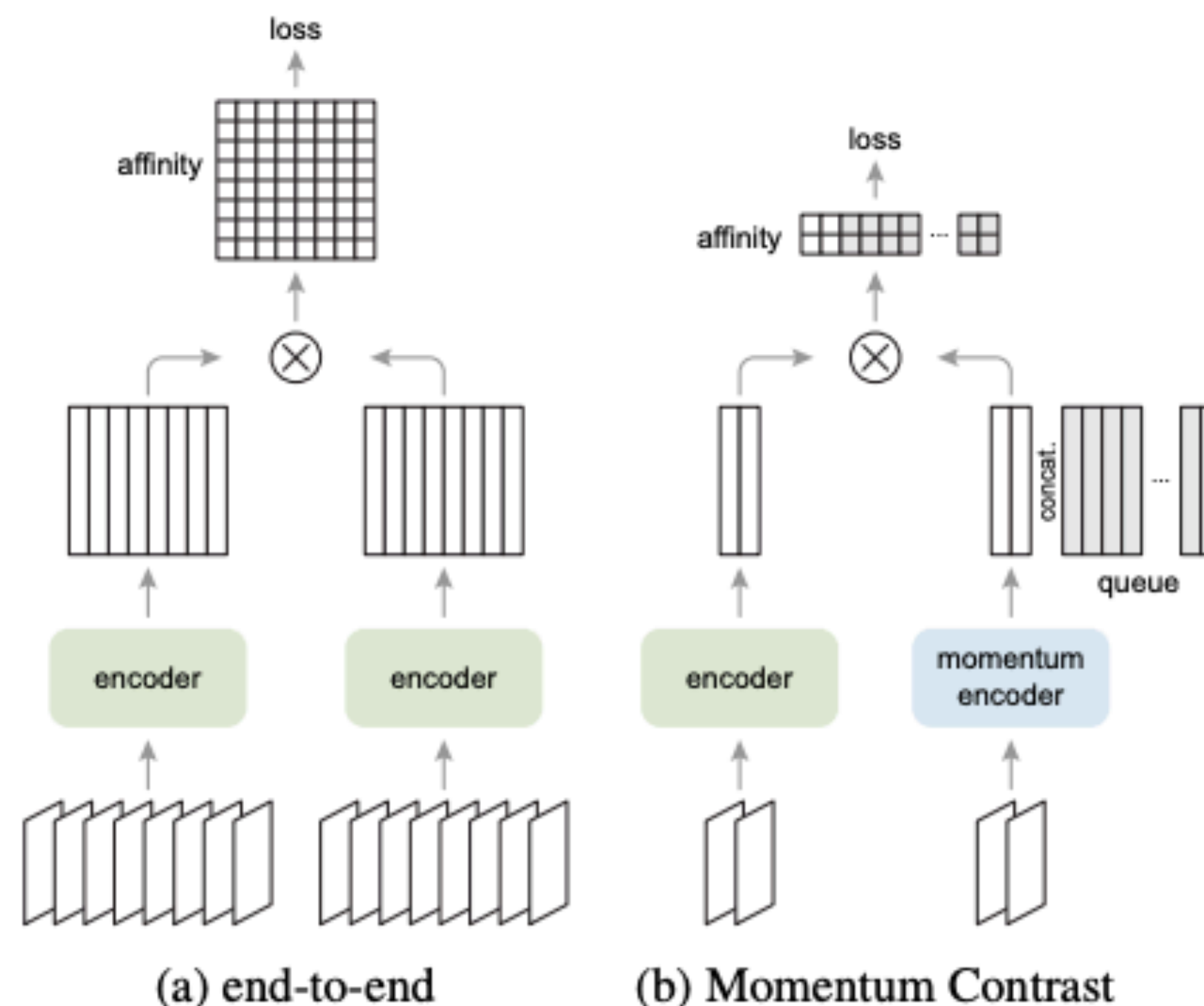


Figure 1. A **batching** perspective of two optimization mechanisms for contrastive learning. Images are encoded into a representation space, in which pairwise affinities are computed.

# Ideas from SimCLR improve MoCo too!

---

| case   | unsup. pre-train |      |     |        |       | ImageNet    |
|--|------------------|------|-----|--------|-------|-------------|
|  | MLP              | aug+ | cos | epochs | batch | acc.        |
| MoCo v1 [6]  |                  |      |     | 200    | 256   | 60.6        |
| SimCLR [2]   | ✓                | ✓    | ✓   | 200    | 256   | 61.9        |
| SimCLR [2]   | ✓                | ✓    | ✓   | 200    | 8192  | 66.6        |
| <b>MoCo v2</b>   | ✓                | ✓    | ✓   | 200    | 256   | <b>67.5</b> |
| <i>results of longer unsupervised training follow:</i> |                  |      |     |        |       |             |
| SimCLR [2]   | ✓                | ✓    | ✓   | 1000   | 4096  | 69.3        |
| <b>MoCo v2</b>   | ✓                | ✓    | ✓   | 800    | 256   | <b>71.1</b> |

Table 2. **MoCo vs. SimCLR**: ImageNet linear classifier accuracy (**ResNet-50, 1-crop  $224 \times 224$** ), trained on features from unsupervised pre-training. “aug+” in SimCLR includes blur and stronger color distortion. SimCLR ablations are from Fig. 9 in [2] (we thank the authors for providing the numerical results).