

# DSC291: Machine Learning with Few Labels

## Self-supervised Learning

**Zhiting Hu**

Lecture 10, February 1, 2023

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Recap: Self-Supervised Learning

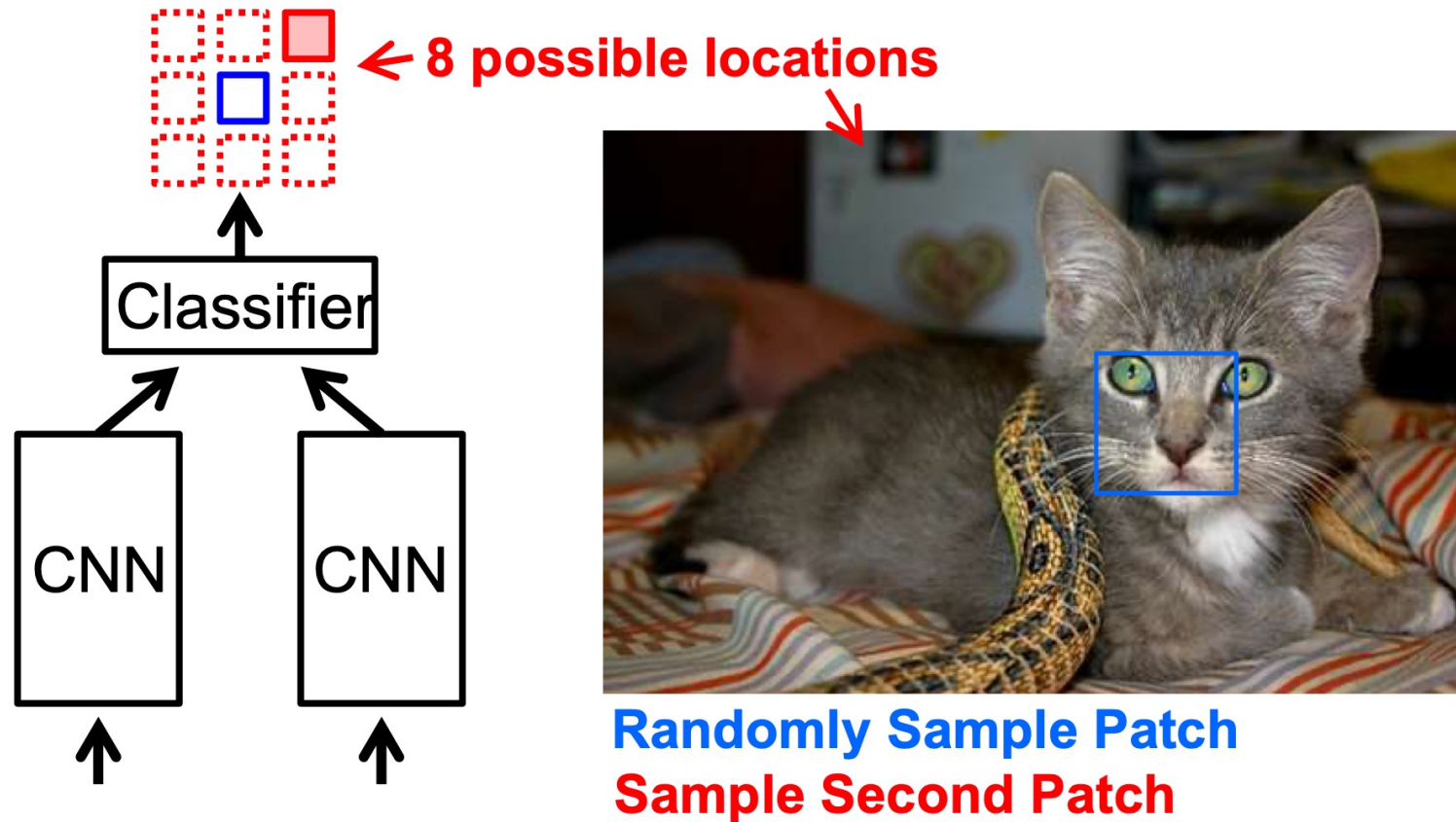
- Given an observed data instance  $\mathbf{t}$
- One could derive various supervision signals based on the structure of the data
- By applying a "split" function that artificially partition  $\mathbf{t}$  into two parts
  - $(\mathbf{x}, \mathbf{y}) = \text{split}(\mathbf{t})$
  - sometimes split in a stochastic way
- Treat  $\mathbf{x}$  as the input and  $\mathbf{y}$  as the output
- Train a model  $p_{\theta}(\mathbf{y}|\mathbf{x})$

# Self-Supervised Learning (SSL): Examples

- SSL from text
  - Language models: next-word prediction (GPT-3)
  - Learning contextual text representations: masked language model (BERT)
- SSL from images
- SSL from videos

# SSL from Images, EX (I): relative positioning

Train network to predict relative position of two regions in the same image



Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

# SSL from Images, EX (I): relative positioning

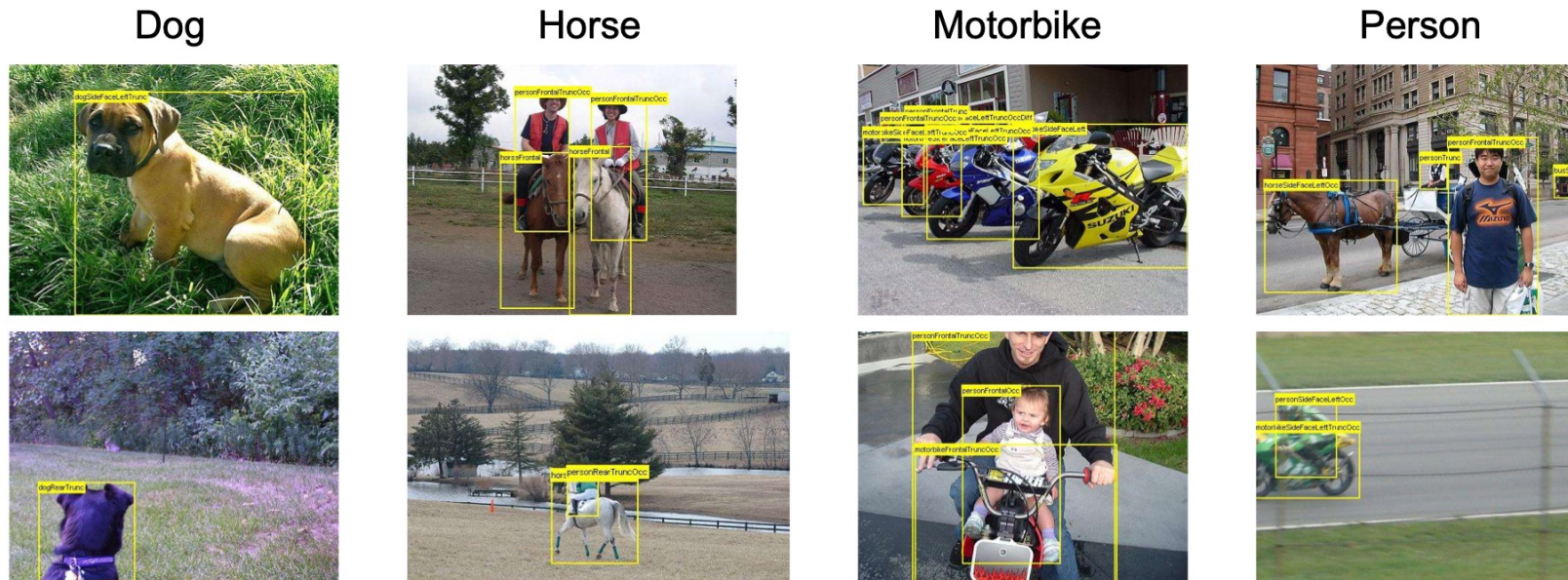


Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

# SSL from Images, EX (I): relative positioning

## Evaluation: PASCAL VOC Detection

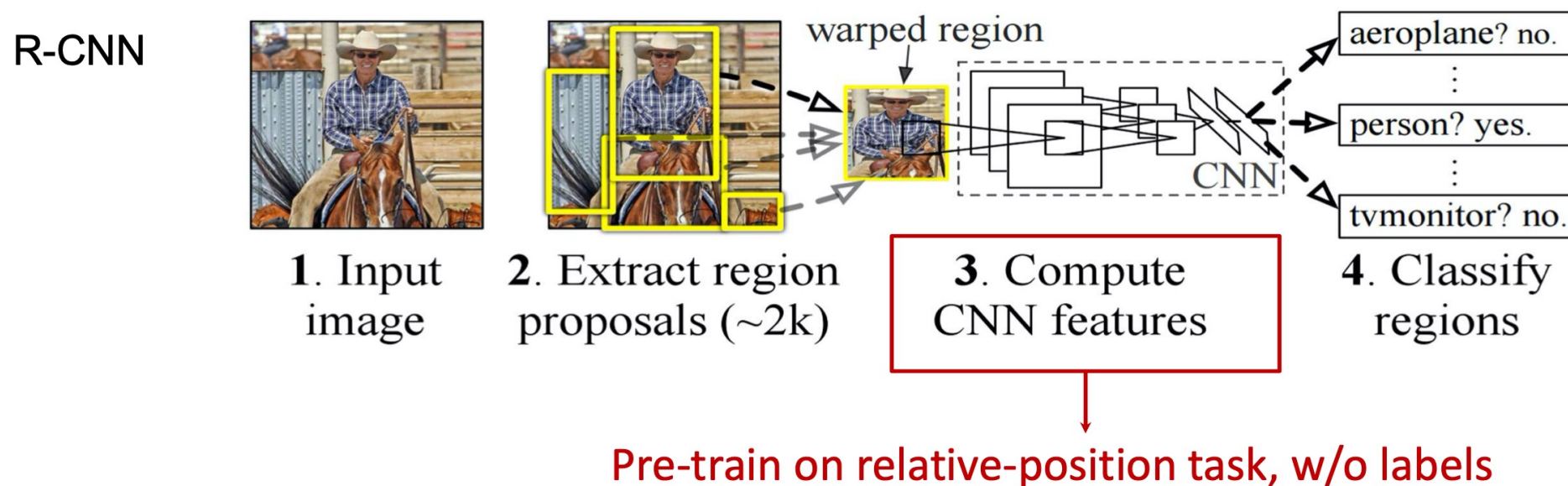
- 20 object classes (car, bicycle, person, horse ...)
- Predict the bounding boxes of all objects of a given class in an image (if any)



# SSL from Images, EX (I): relative positioning

## Evaluation: PASCAL VOC Detection

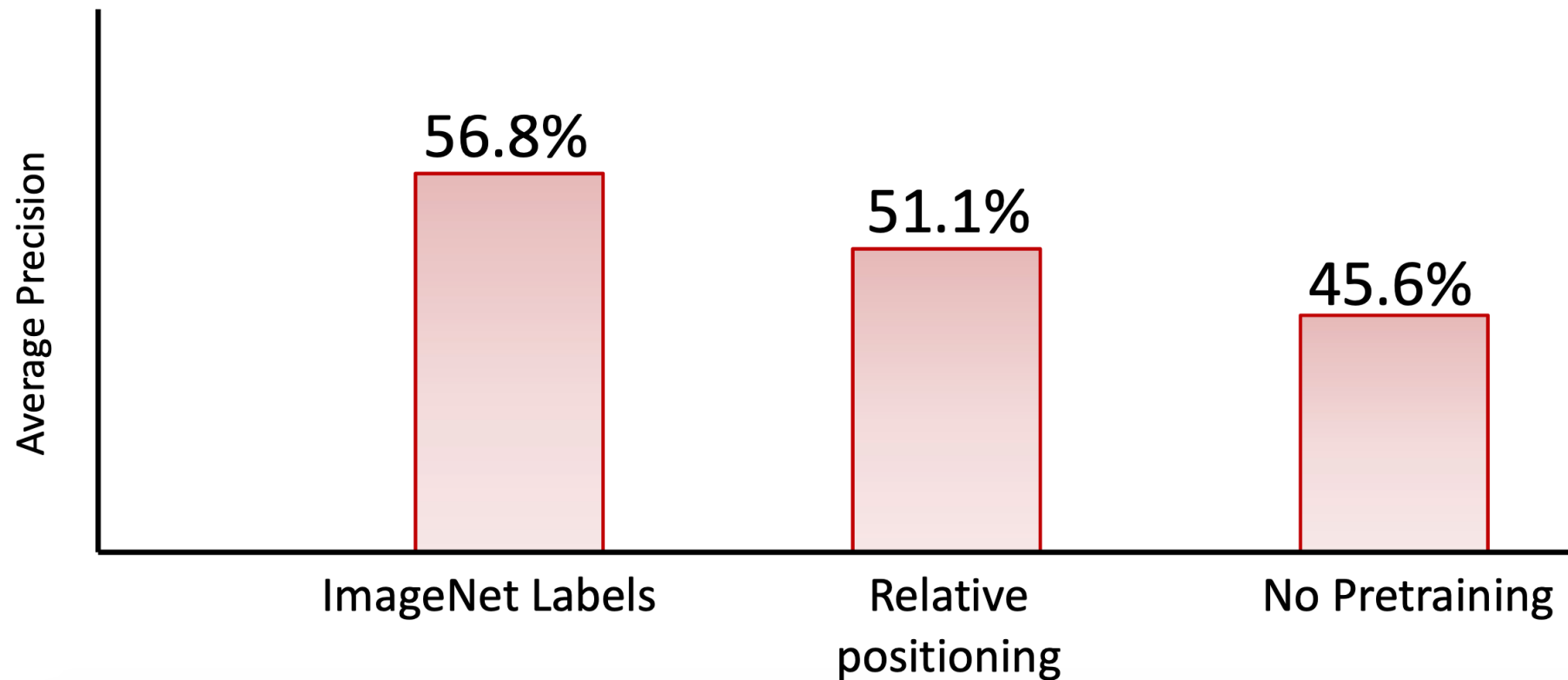
- Pre-train CNN using self-supervision (no labels)
- Train CNN for detection in R-CNN object category detection pipeline



[Girshick et al. 2014]

# SSL from Images, EX (I): relative positioning

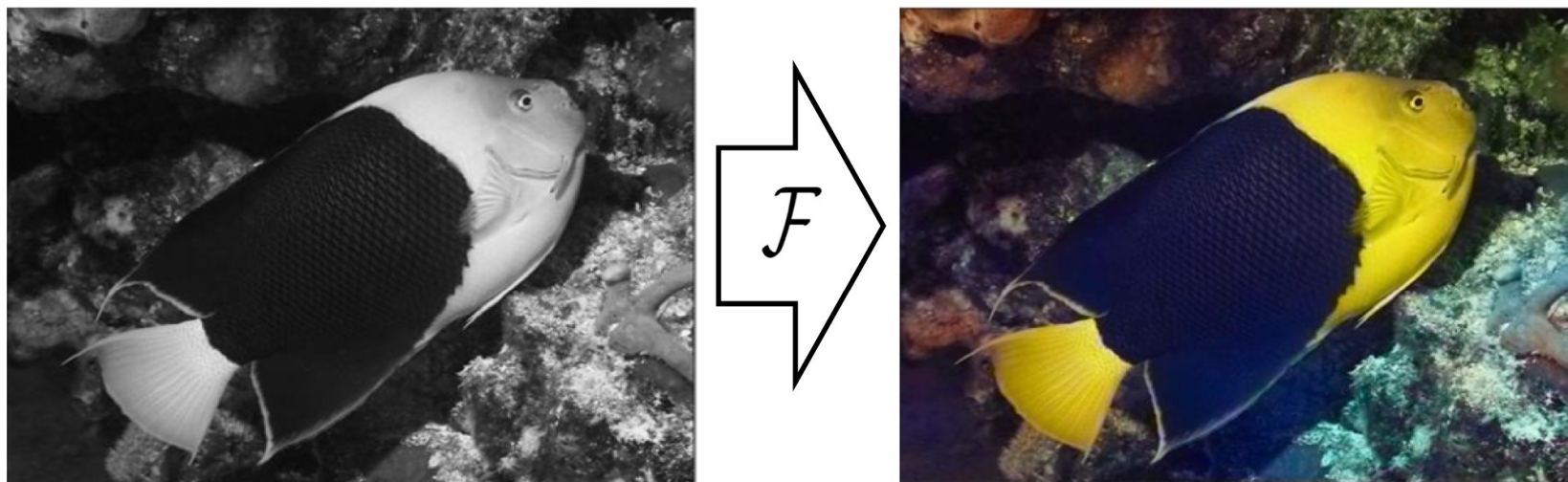
Evaluation: PASCAL VOC Detection





# SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input

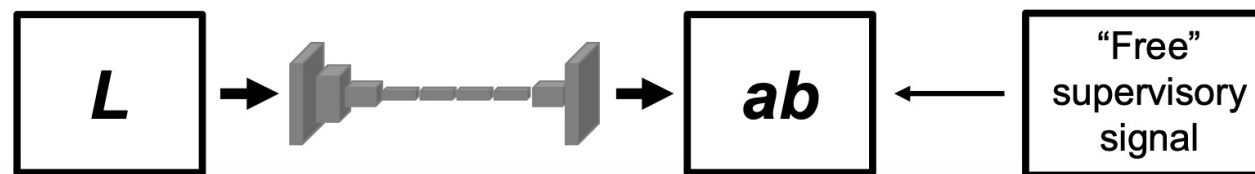


Grayscale image:  $L$  channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

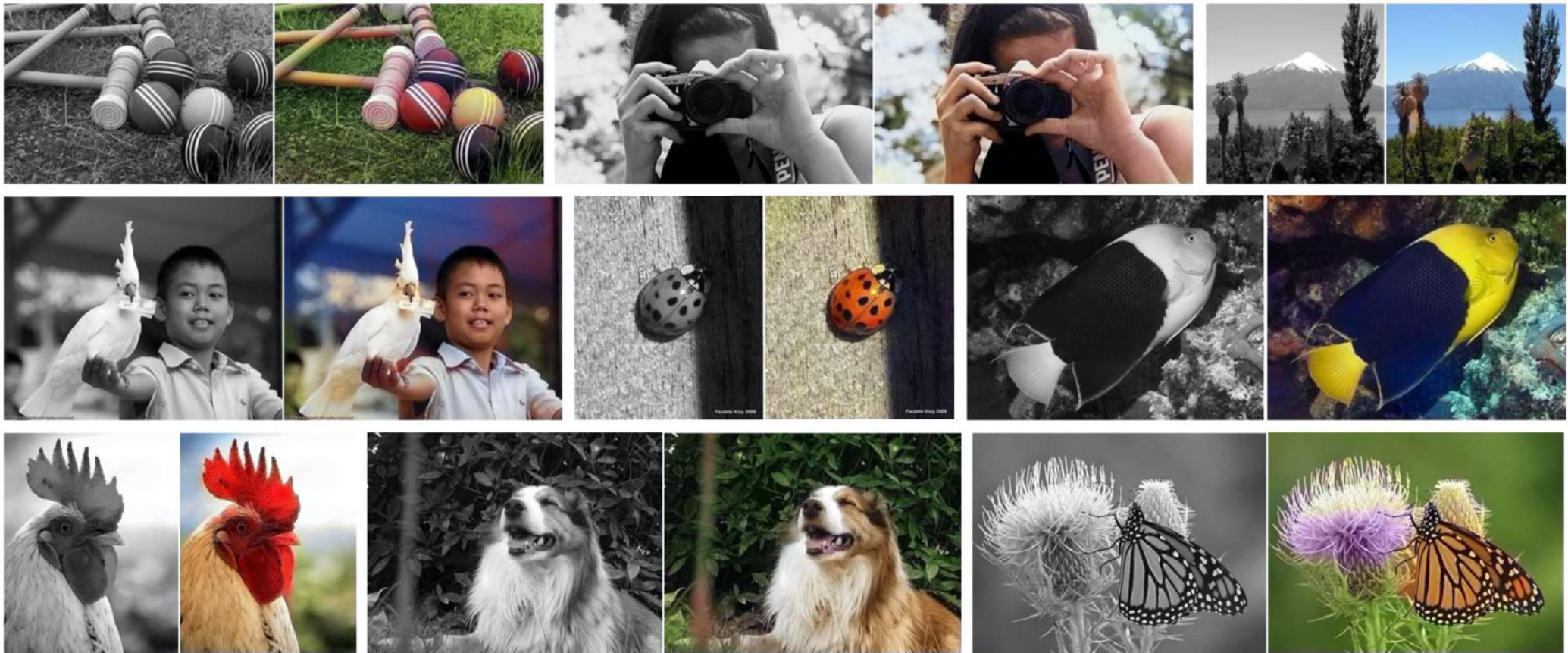
Concatenate ( $L, ab$ )

$$(\mathbf{X}, \hat{\mathbf{Y}})$$



# SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input



# SSL from Images, EX (III): exemplar networks

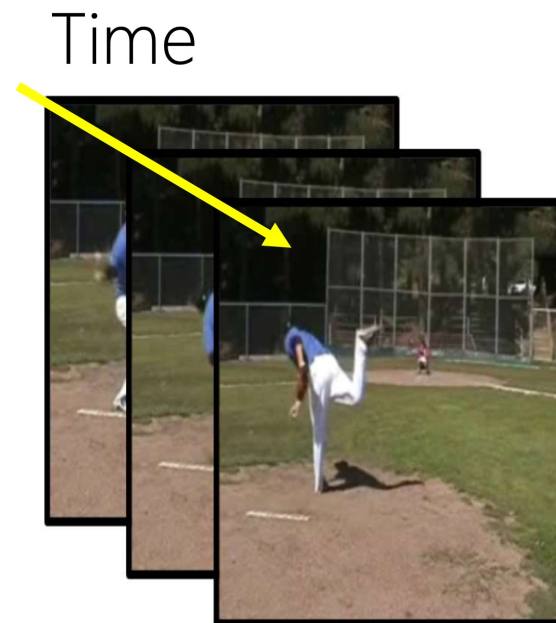
- Exemplar Networks (Dosovitskiy et al., 2014)
- Perturb/distort image patches, e.g. by cropping and affine transformations
- Train to classify these exemplars as same class



# SSL from Videos

Three example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?



“Sequence” of data

# SSL from Videos

Three example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?
- Video direction
  - Predict if video playing forwards or backwards

# SSL from Videos

Three example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?
- Video direction
  - Predict if video playing forwards or backwards
- Video tracking
  - Given a color video, colorize all frames of a gray scale version using a reference frame



# Key Takeaways

- Self supervision learning
  - Predicting any part of the observations given any available information
  - The prediction task forces models to learn semantic representations
  - Massive/unlimited data supervisions
- SSL for text:
  - Language models: next word prediction
  - BERT text representations: masked language model (MLM)
- SSL for images/videos:
  - Various ways of defining the prediction task

# Contrastive Learning



# Contrastive learning

- Take a data example  $x$ , sample a “positive” sample  $x_{pos}$  and “negative” samples  $x_{neg}$  in some way
- Then try fit a scoring model such that

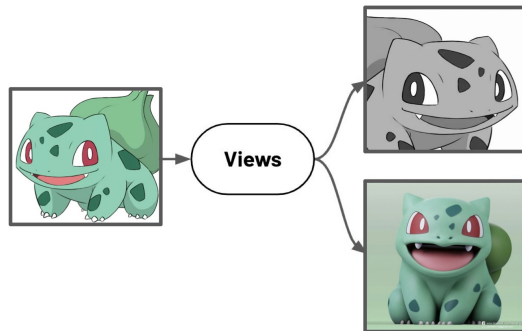
$$score(x, x_{pos}) > score(x, x_{neg})$$

# Contrastive learning

- Take a data example  $x$ , sample a “positive” sample  $x_{pos}$  and “negative” samples  $x_{neg}$  in some way

“positive” sample:

- Data of the same labels
- Data of the same pseudo-labels
- Augmented/distorted version of  $x$
- Data that captures the same target from different views



“negative” sample:

- Randomly sampled data
- Hard negative sample mining

# Contrastive learning

- Take a data example  $x$ , sample a “positive” sample  $x_{pos}$  and “negative” samples  $x_{neg}$  in some way
- Then try fit a scoring model such that

$$score(x, x_{pos}) > score(x, x_{neg})$$

# Contrastive learning: Ex 1

Learning a similarity metric discriminatively

Sample a pair of images and compute their distance:

$$D_i = \|x, x_i\|_2$$

If **positive** sample:

$$L_i = D_i^2$$



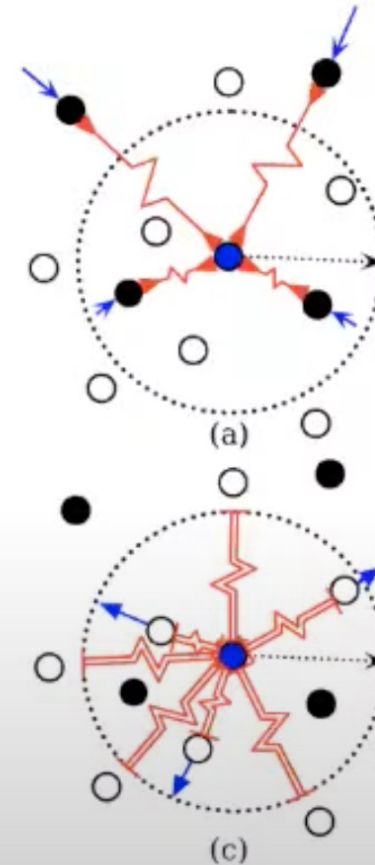
x pos

If **negative** sample:

$$L_i = \max(0, \epsilon - D_i)^2$$



x neg



[Chopra et al., 2005; Hadsell et al., 2006]

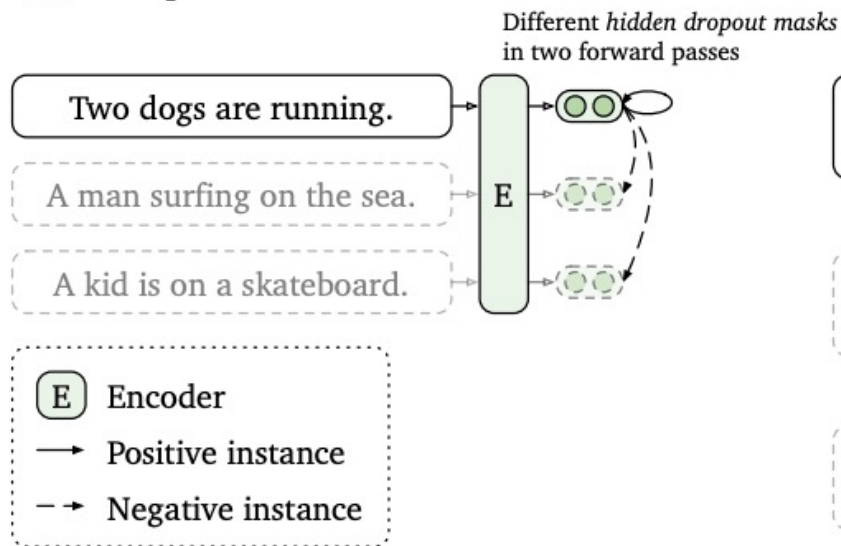
# Common contrastive learning functions

- Contrastive loss (Chopra et al. 2005)
- Triplet loss (Schroff et al. 2015; FaceNet)
- Lifted structured loss (Song et al. 2015)
- Multi-class n-pair loss (Sohn 2016)
- Noise contrastive estimation (“NCE”; Gutmann & Hyvarinen 2010)
- InfoNCE (van den Oord, et al. 2018)
- Soft-nearest neighbors loss (Salakhutdinov & Hinton 2007, Frosst et al. 2019)

# Contrastive learning: Ex 2

- SimCSE (“Simple Contrastive learning of Sentence Embeddings”; Gao et al. 2021)
  - Predict a sentence from itself with only dropout noise
  - One sentence gets two different versions of dropout augmentations

(a) Unsupervised SimCSE



(b) Supervised SimCSE

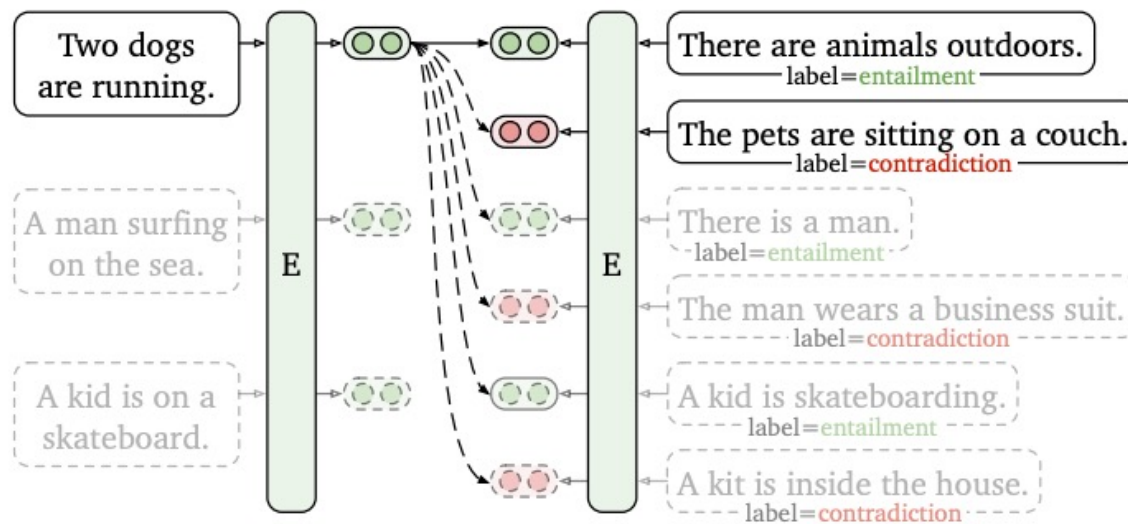
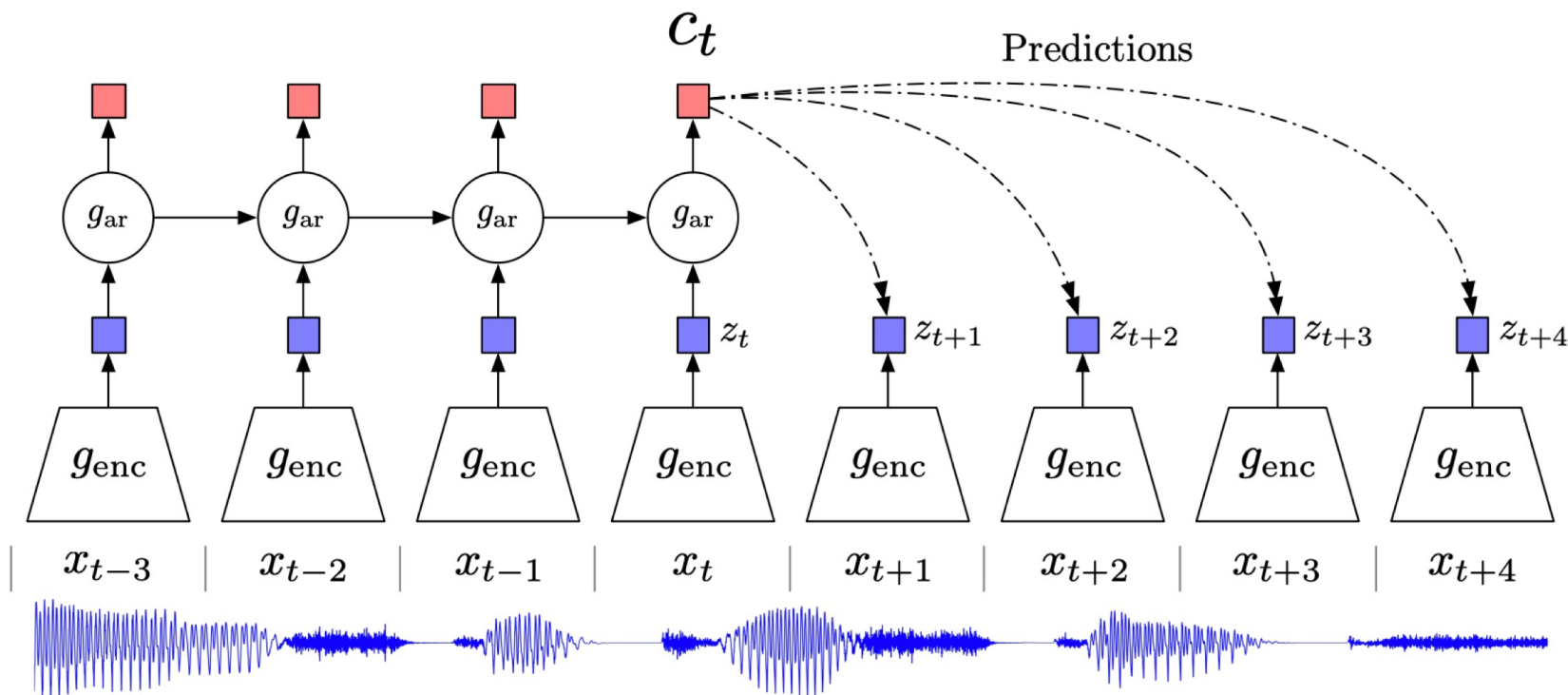


Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied. (b) Supervised SimCSE leverages the NLI datasets and takes the entailment (premise-hypothesis) pairs as positives, and contradiction pairs as well as other in-batch instances as negatives.

# Contrastive learning: Ex 3 - InfoNCE

- The CPC model
  - $c_t$ : context representation from history
  - $x_{t+k}$  (or  $z_{t+k}$ ): future target



# InfoNCE loss

- Define scoring function  $f_k > 0$
- The InfoNCE loss:
  - Given  $X = \{ \text{one positive sample from } p(x_{t+k} | c_t), N - 1 \text{ negative samples from the negative sampling distribution } p(x_{t+k}) \}$

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

- InfoNCE is interesting because it's effectively maximizing the **mutual information** between  $c_t$  and  $x_{t+k}$



# Mutual Information (MI)

- How much is our uncertainty about  $x$  reduced by knowing  $c$  ?

$$I(x; c) = \sum_{x,c} p(x, c) \log \frac{p(x, c)}{p(x)p(c)} = \sum_{x,c} p(x, c) \log \frac{p(x|c)}{p(x)}$$

$$= H(x) + H(c) - H(x, c)$$

$$= H(x) - H(x|c)$$

$$= KL(p(x, c) || p(x)p(c))$$

# Minimizing InfoNCE $\Leftrightarrow$ Maximizing MI

- InfoNCE loss

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

- The loss is optimized when

$$f_k(x_{t+k}, c_t) \propto \frac{p(x_{t+k} | c_t)}{p(x_{t+k})}$$

- Proof:

$$\begin{aligned} p(\text{sample } i \text{ is positive} | X, c_t) &= \frac{p(x_i | c_t) \prod_{l \neq i} p(x_l)}{\sum_{j=1}^N p(x_j | c_t) \prod_{l \neq j} p(x_l)} \\ &= \frac{\frac{p(x_i | c_t)}{p(x_i)}}{\sum_{j=1}^N \frac{p(x_j | c_t)}{p(x_j)}} \end{aligned}$$

- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

$$\mathcal{L}_N^{\text{opt}} = -\mathbb{E}_X \log \left[ \frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)}} \right]$$

**Use proportionality condition**

- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

$$\begin{aligned} \mathcal{L}_N^{\text{opt}} &= -\mathbb{E}_X \log \left[ \frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)}} \right] \\ &= \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \end{aligned}$$

**Take -ve inside log**

- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

$$\begin{aligned} \mathcal{L}_N^{\text{opt}} &= -\mathbb{E}_X \log \left[ \frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)}} \right] \\ &= \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \\ &\approx \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \mathbb{E}_{x_j} \frac{p(x_j|c_t)}{p(x_j)} \right] \end{aligned}$$

**This approximation becomes more accurate as N increases, so it is preferable to use large negative samples**

- How does this loss maximize the mutual information?

$$\begin{aligned}
 \mathcal{L}_N &= -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right] \\
 \mathcal{L}_N^{\text{opt}} &= -\mathbb{E}_X \log \left[ \frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)}} \right] \\
 &= \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \\
 &\approx \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \mathbb{E}_{x_j} \frac{p(x_j|c_t)}{p(x_j)} \right] = 1 \\
 &= \mathbb{E}_X \log \left[ 1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \right]
 \end{aligned}$$

- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(\mathbf{x}_{t+k}, \mathbf{c}_t)}{\sum_{\mathbf{x}_j \in X} f_k(\mathbf{x}_j, \mathbf{c}_t)} \right]$$

$$\begin{aligned} \mathcal{L}_N^{\text{opt}} &= -\mathbb{E}_X \log \left[ \frac{\frac{p(\mathbf{x}_{t+k} | \mathbf{c}_t)}{p(\mathbf{x}_{t+k})}}{\frac{p(\mathbf{x}_{t+k} | \mathbf{c}_t)}{p(\mathbf{x}_{t+k})} + \sum_{\mathbf{x}_j \in X_{\text{neg}}} \frac{p(\mathbf{x}_j | \mathbf{c}_t)}{p(\mathbf{x}_j)}}} \right] \\ &= \mathbb{E}_X \log \left[ 1 + \frac{p(\mathbf{x}_{t+k})}{p(\mathbf{x}_{t+k} | \mathbf{c}_t)} \sum_{\mathbf{x}_j \in X_{\text{neg}}} \frac{p(\mathbf{x}_j | \mathbf{c}_t)}{p(\mathbf{x}_j)} \right] \\ &\approx \mathbb{E}_X \log \left[ 1 + \frac{p(\mathbf{x}_{t+k})}{p(\mathbf{x}_{t+k} | \mathbf{c}_t)} (N-1) \mathbb{E}_{\mathbf{x}_j} \frac{p(\mathbf{x}_j | \mathbf{c}_t)}{p(\mathbf{x}_j)} \right] \\ &= \mathbb{E}_X \log \left[ 1 + \frac{p(\mathbf{x}_{t+k})}{p(\mathbf{x}_{t+k} | \mathbf{c}_t)} (N-1) \right] \\ &\geq \mathbb{E}_X \log \left[ \frac{p(\mathbf{x}_{t+k})}{p(\mathbf{x}_{t+k} | \mathbf{c}_t)} N \right] \\ &= -I(\mathbf{x}_{t+k}, \mathbf{c}_t) + \log(N), \end{aligned}$$



- How does this loss maximize the mutual information?

$$\mathcal{L}_N = -\mathbb{E}_X \left[ \log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

$$I(x_{t+k}, c_t) \geq \log(N) - \mathcal{L}_N$$

# Key Takeaways: Contrastive learning

- Contrastive learning is a way of doing self-supervised learning
- Positive samples, negative samples
- Mutual information

$$\begin{aligned} I(x; c) &= \sum_{x,c} p(x, c) \log \frac{p(x, c)}{p(x)p(c)} = \sum_{x,c} p(x, c) \log \frac{p(x|c)}{p(x)} \\ &= H(x) + H(c) - H(x, c) \\ &= H(x) + H(x|c) \\ &= KL(p(x, c) || p(x)p(c)) \end{aligned}$$

- InfoNCE  $\Leftrightarrow$  MI

Questions?