



Xavier Giro-i-Nieto

 [@DocXavi](https://twitter.com/DocXavi)  
 [xavier.giro@upc.edu](mailto:xavier.giro@upc.edu)

Associate Professor  
Universitat Politècnica de Catalunya

# Master in Computer Vision Barcelona

[<http://pagines.uab.cat/mcv/>]

Module 6 - Day 9 - Lecture 4  
Self-supervised Learning  
5th April 2022

# Acknowledgements



Oscar  
Mañas



Víctor  
Campos



Junting  
Pan



Xunyu  
Lin



Carlos  
Arenas



Sebastian  
Palacio

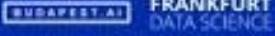
# Video lecture



**DEEP SELF-SUPERVISED LEARNING FOR ALL**

by Xavier Giro-i-Nieto  
Visiting Researcher and Associate Professor

**X-Europe Webinar**  
**15.07 18:00 CET**

and Barcelona Data Science and Machine Learning Meetup, Budapest  
Deep Learning Reading Seminar



**Yann LeCun**  
@ylecun

Self-supervised learning, obviously.

[Tradueix el tuit](#)

 **Muzaffer Kal WFH attending COVID-19** @DSPonFPGA · 19 de març  
"supervised machine learning doesn't live up to the hype" @ylecun where do we go from here?

[medium.com/starsky-roboti...](https://medium.com/starsky-robotics-self-supervised-learning-doesnt-live-up-to-the-hype-10333a2f3a)

7:10 a. m. · 20 de març de 2020 · Twitter for Android

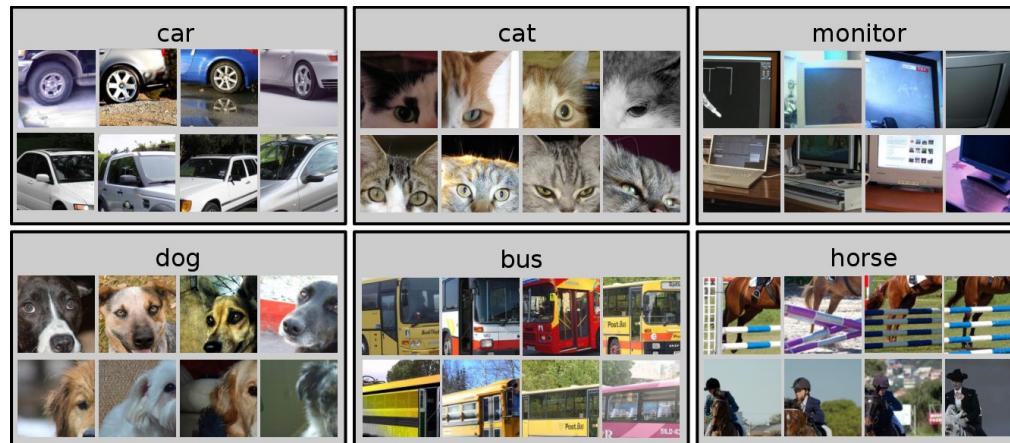
# Outline

- 1. Transfer Learning**
2. Representation Learning
3. Unsupervised Learning
4. Self-supervised Learning
5. Predictive methods
6. Contrastive methods

# Study case: PASCAL VOC 2007 dataset

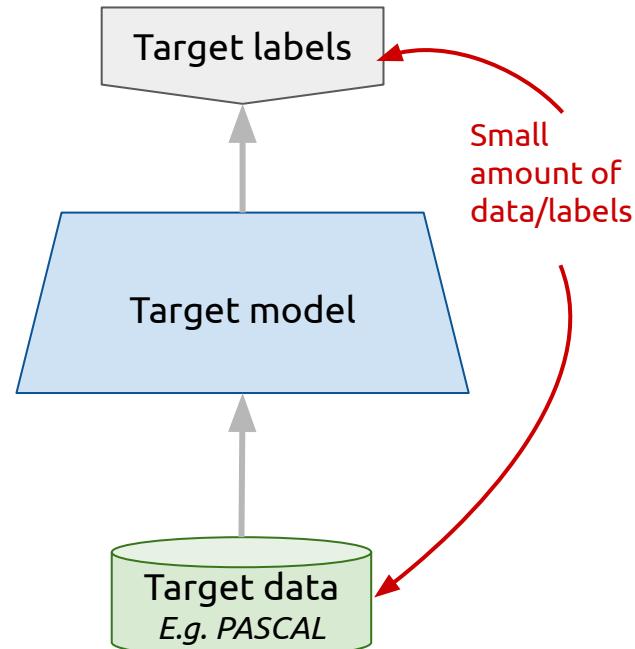
- Standard classification benchmark, 20 classes, ~10K images, 50% train, 50% test
- Deep networks can have many parameters (e.g. 60M in Alexnet)
- Direct training (from scratch) using only 5K training images can be problematic. Model overfits.

How can we use deep networks in this setting ?



# Recap: Transfer Learning

In many cases, annotated data is not enough for training deep neural networks. The model will just overfit to the training data and do not generalize well.



# Transfer learning: idea

Instead of training a deep network from scratch for your task:

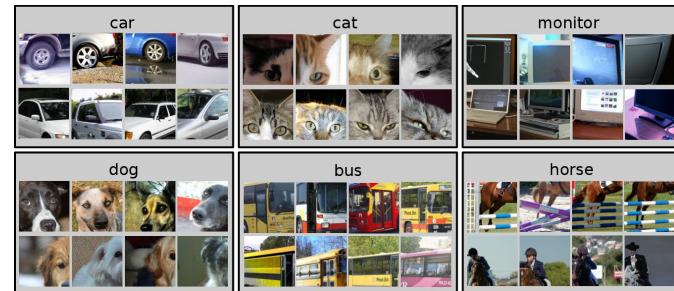
- Take a network trained on a different domain same/different **source task**
- Adapt it for your domain and your **target task**

IMAGENET



Source domain / task

Transfer Learned Knowledge

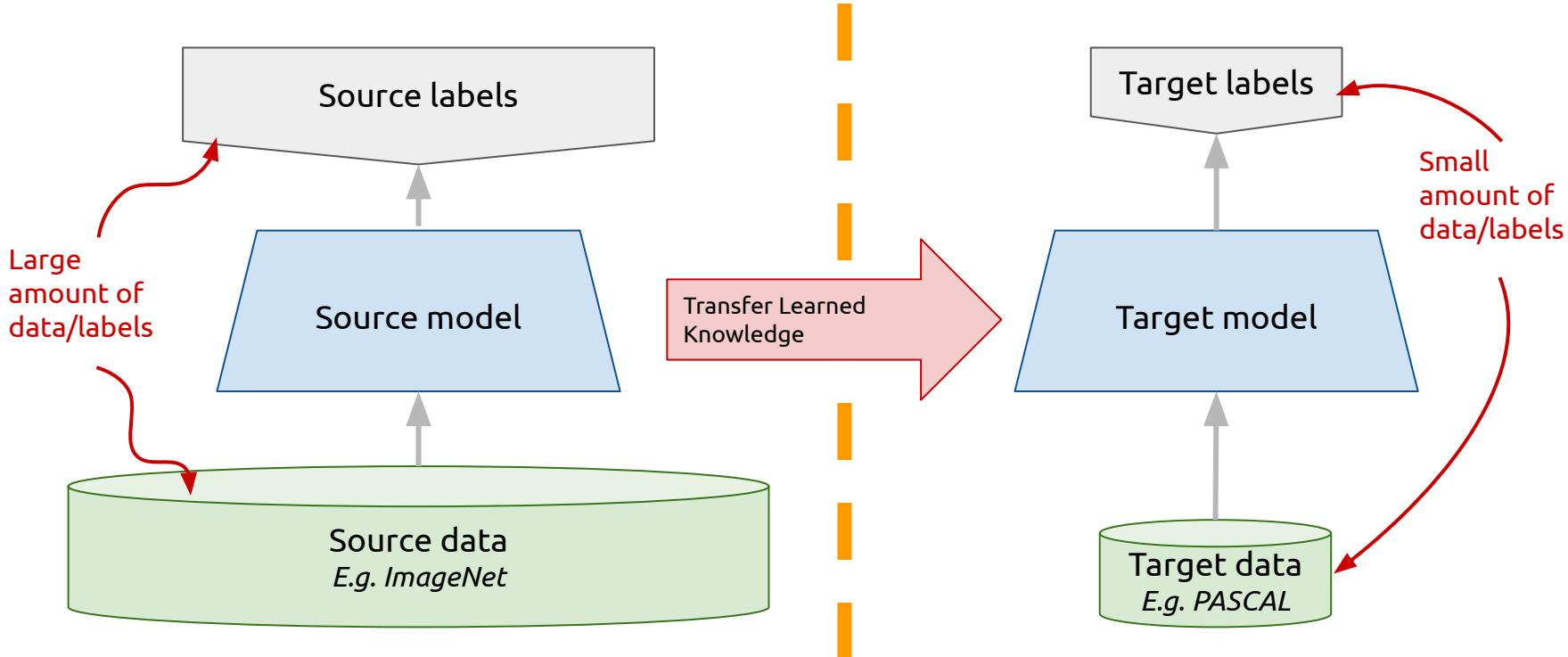


Target domain / task

# Transfer Learning

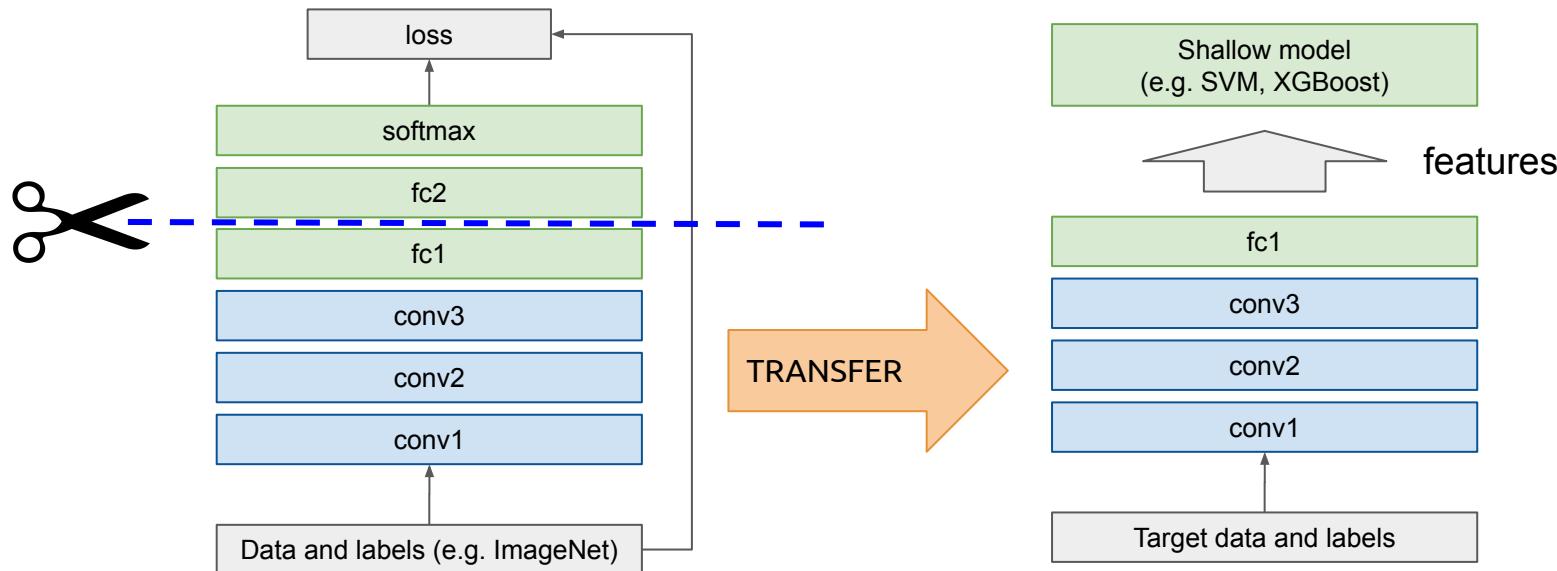
Instead of training a deep network from scratch for your task:

- Take a network trained on a different **source domain and/or task**
- Adapt it for your **target domain and/or task**



# “Off-the-shelf”

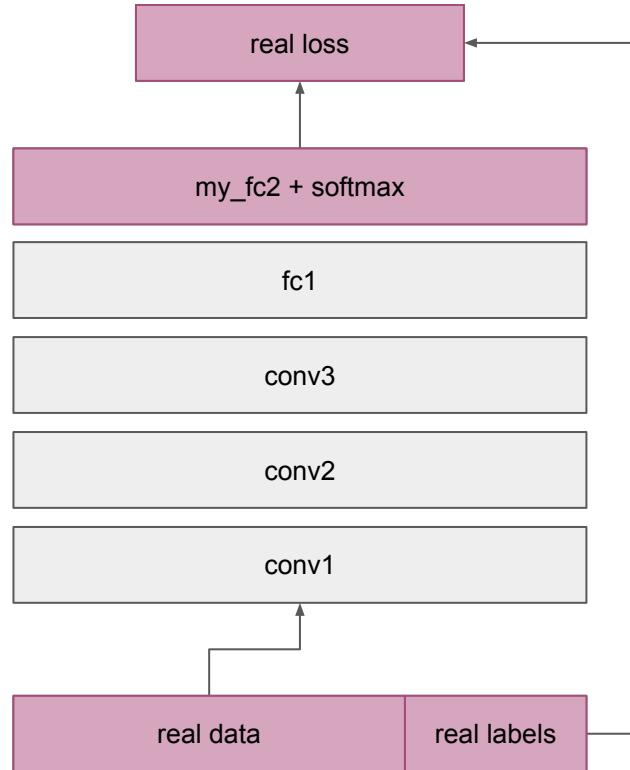
Idea: use outputs of one or more layers of a network trained on a different task as generic feature detectors. Train a new shallow model on these features.



# Linear probing or Fine-tuning

Cut off top layer(s) of network and replace with supervised objective for target domain

- **Linear probing:** Train my\_fc2 + softmax only.
- **Fine-tuning:** Backprop gradients to modify parameters of the pre-trained model.

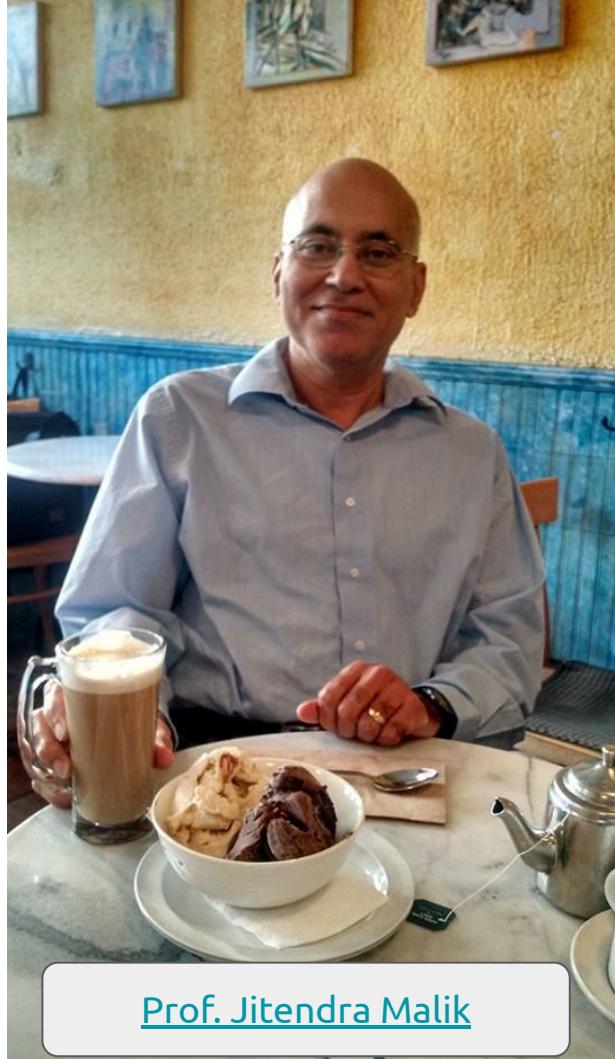


#Decaf Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., & Darrell, T. . [Decaf: A deep convolutional activation feature for generic visual recognition](#). ICML 2014.

Oquab, M., Bottou, L., Laptev, I., & Sivic, J. [Learning and transferring mid-level image representations using convolutional neural networks](#). CVPR 2014.

## The Gelato Bet (Sept 23, 2014)

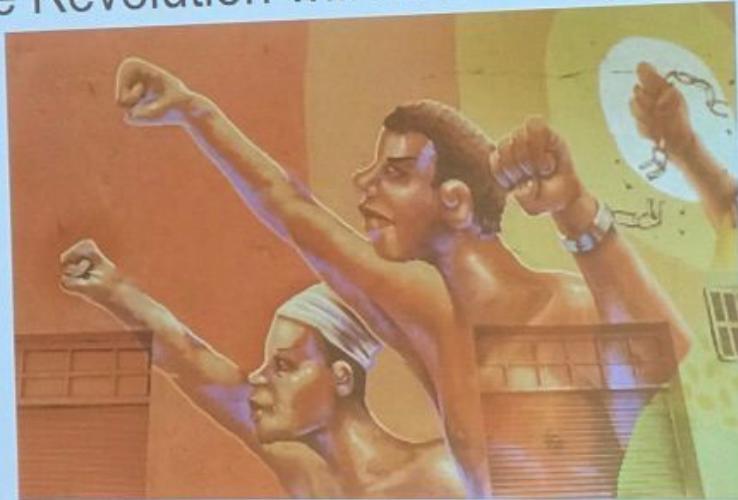
*"If, by the first day of autumn (Sept 23) of 2015, a method will exist that can match or beat the performance of R-CNN on Pascal VOC detection, without the use of any extra, human annotations (e.g. ImageNet) as pre-training, Mr. Malik promises to buy Mr. Efros one (1) gelato (2 scoops: one chocolate, one vanilla)."*



Prof. Jitendra Malik

Source: Ankesh Anand, "Contrastive Self-Supervised Learning" (2020)"

# The Revolution will not be Supervised



Alexei A. Efros  
UC Berkeley



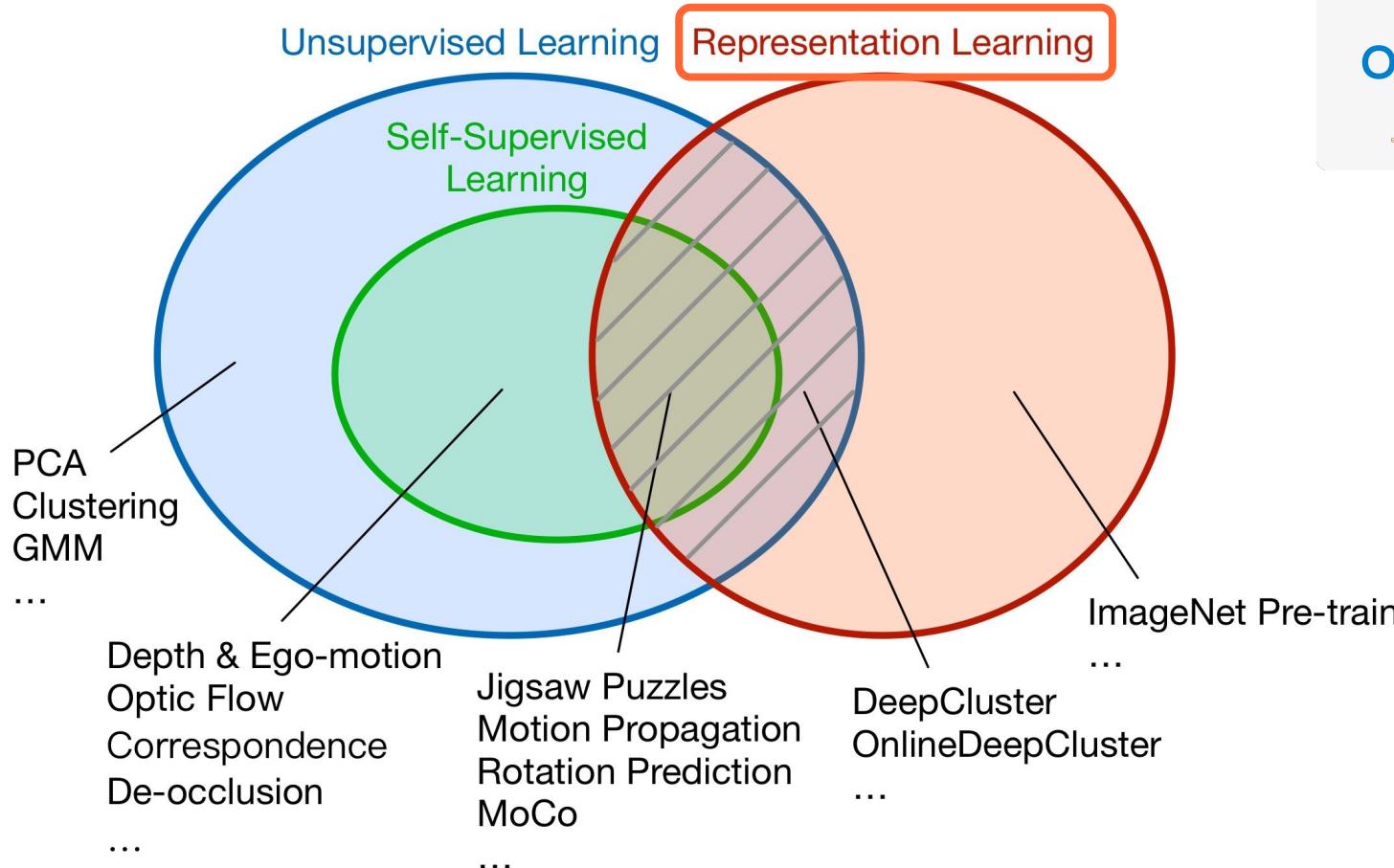
BALI INTELLIGENCE RESEARCH



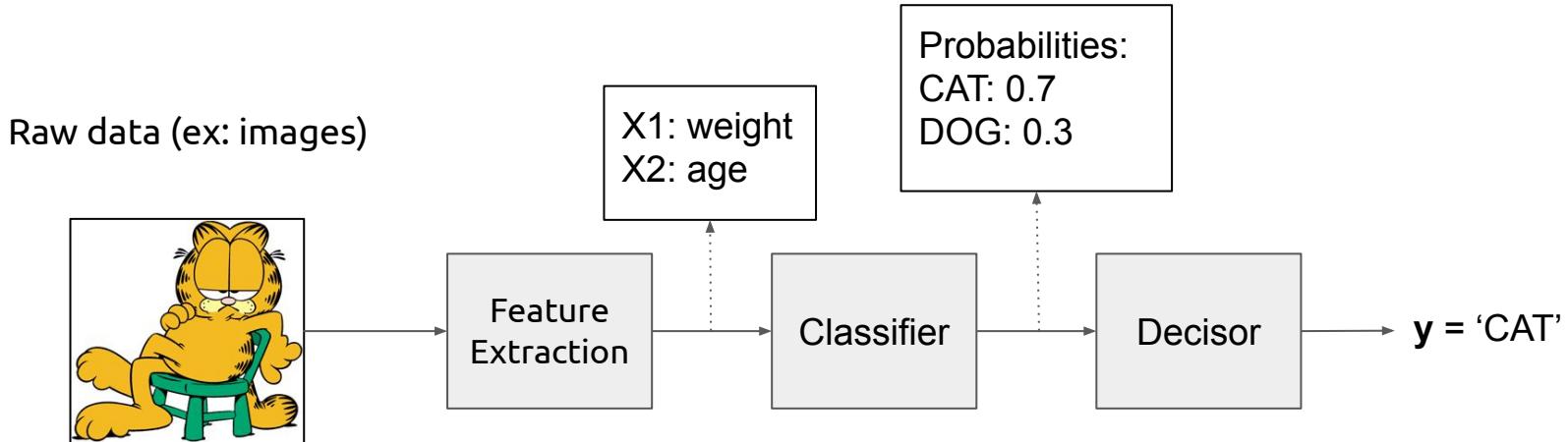
[Prof. Alexei A. Efros](#)

# Outline

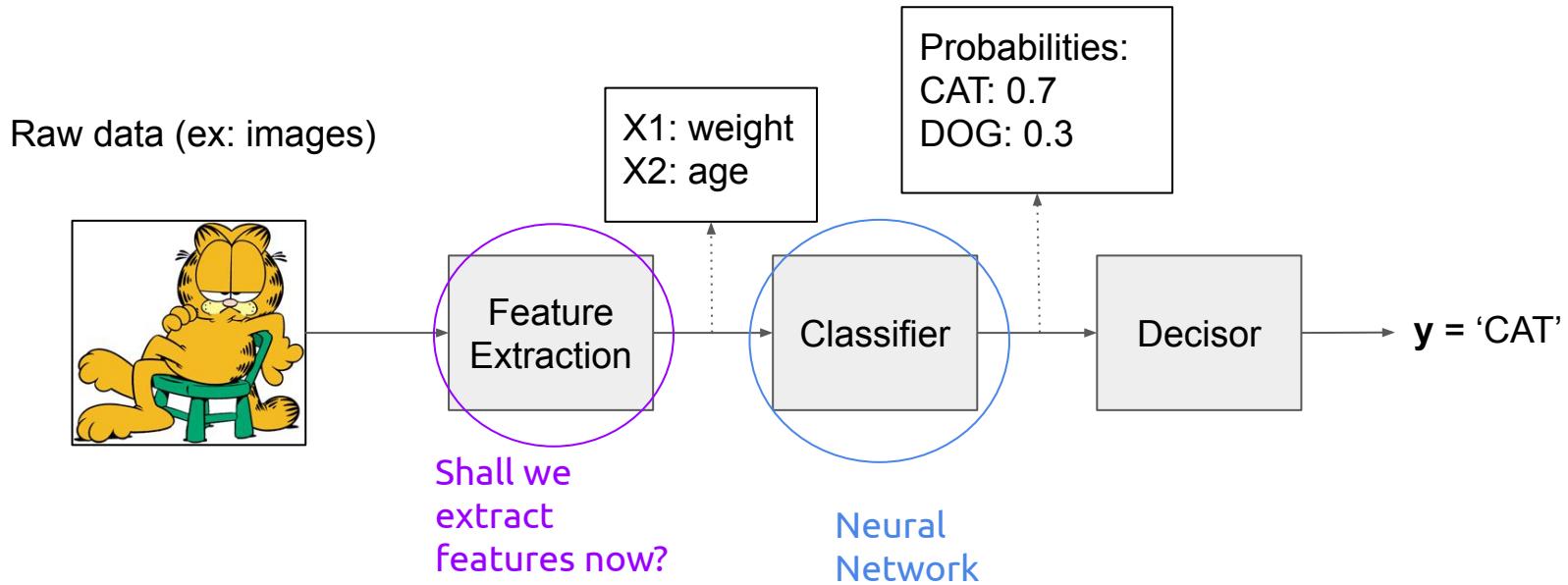
1. Transfer Learning
2. **Representation Learning**
3. Unsupervised Learning
4. Self-supervised Learning
5. Predictive methods
6. Contrastive methods



# Representation + Learning pipeline



# Representation + Learning pipeline



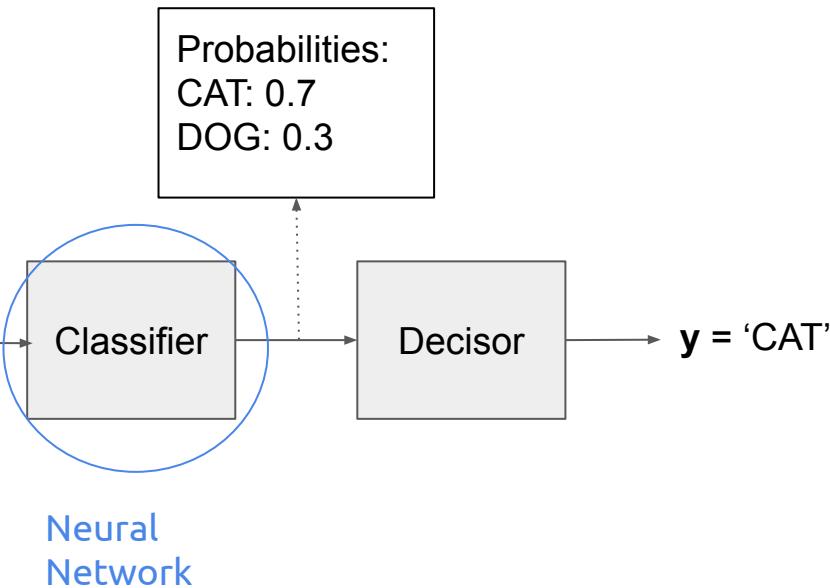
# End-to-end Representation Learning

## End to End concept

Raw data (ex: images)

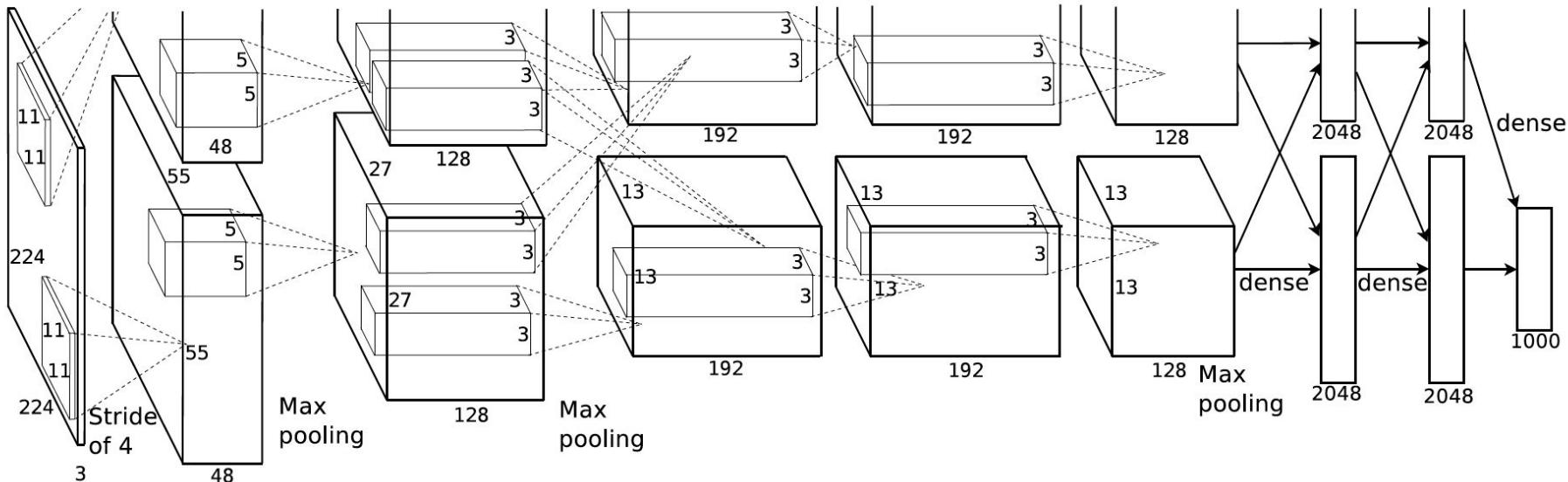


We CAN inject the raw data, and features will be learned in the NN !!



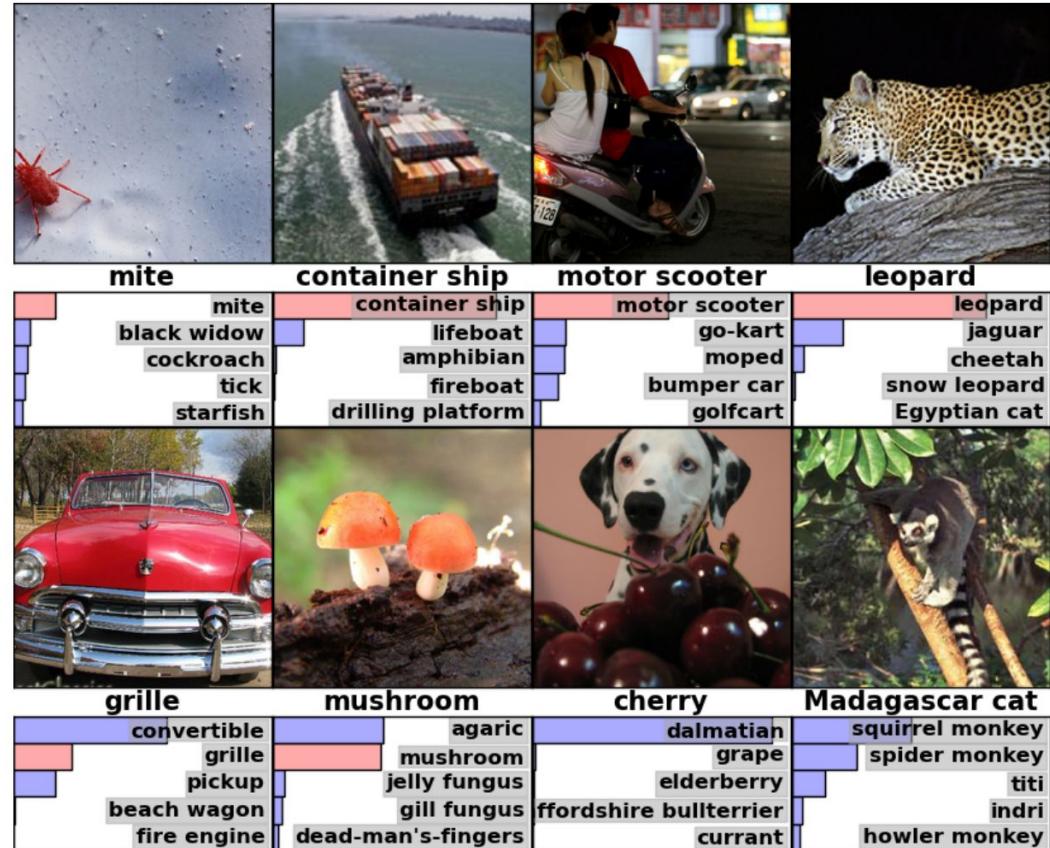
# the BIG BANG

76,278 citations  
(February 2021)



# IMAGENET

- 1,000 object classes (categories).
- Labeled Images:
  - 1.2 M train
  - 100k test.

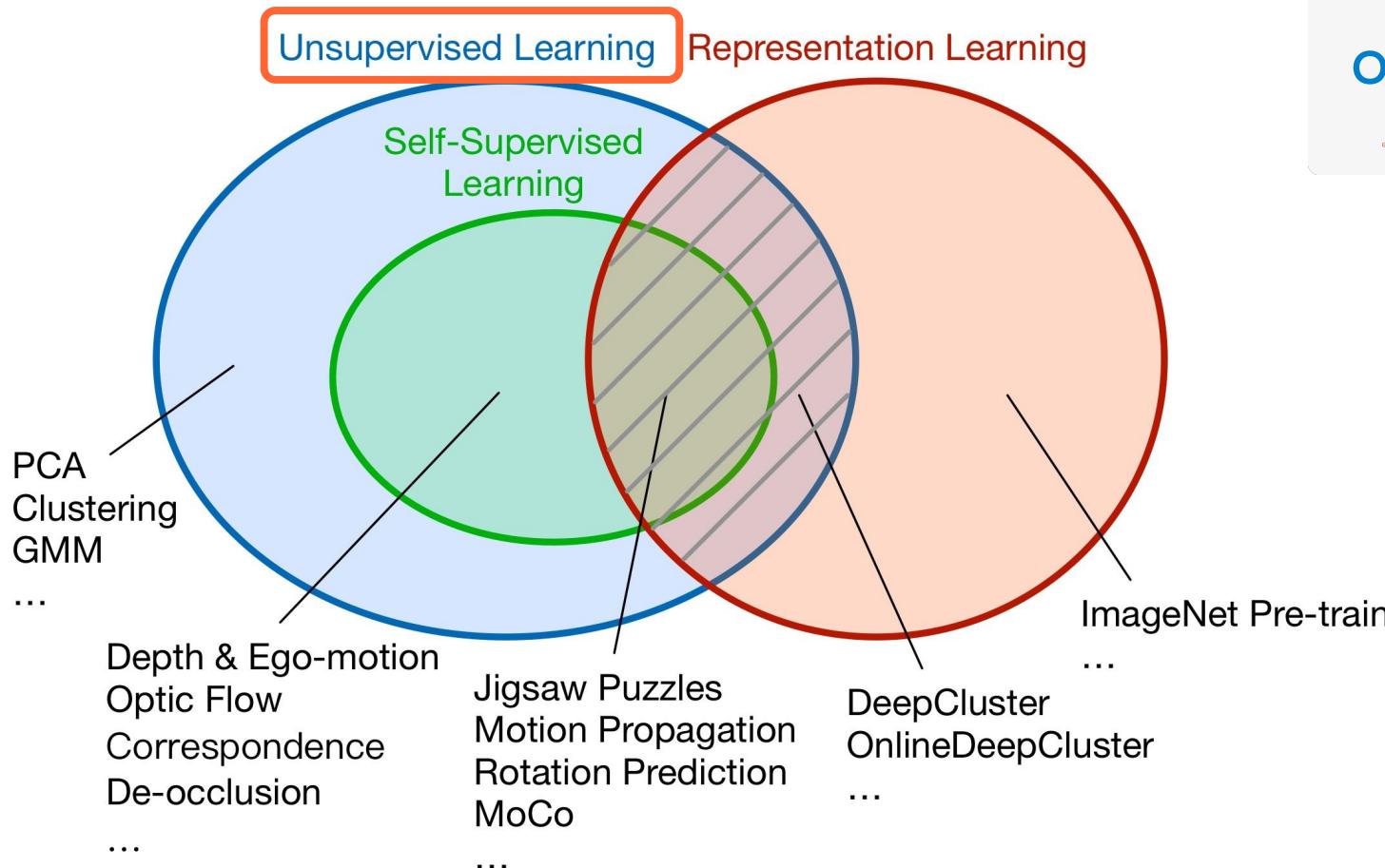


# The Future of Representation Learning, I

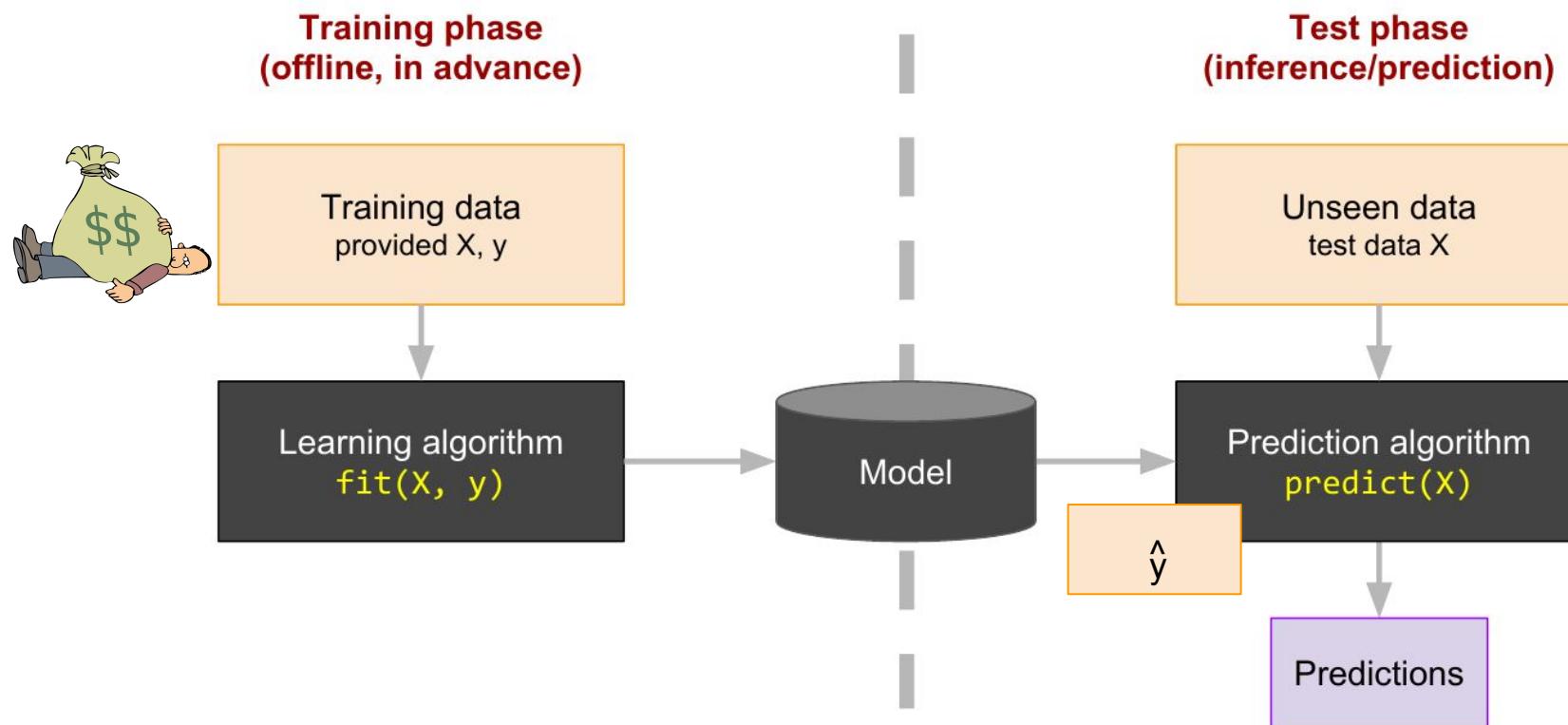


# Outline

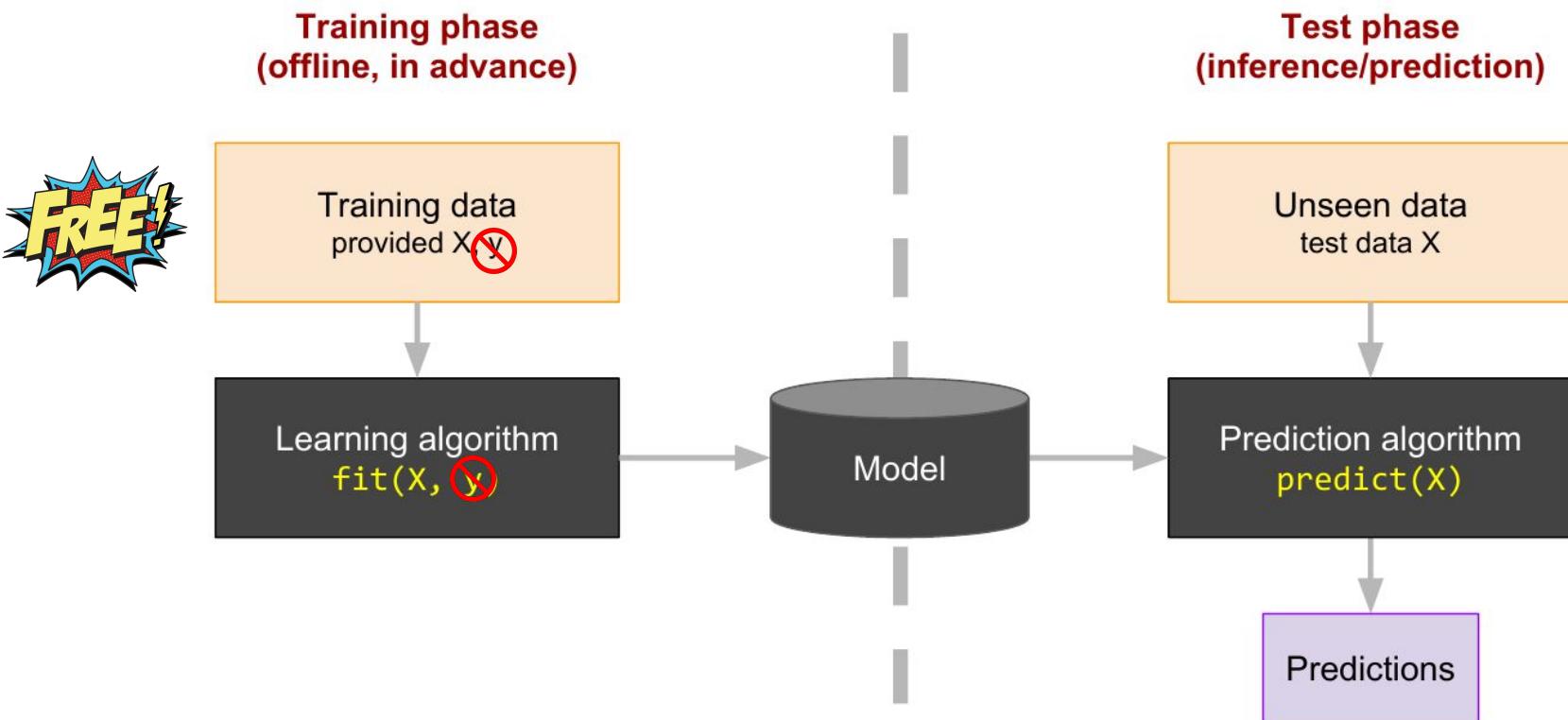
1. Transfer Learning
2. Representation Learning
- 3. Unsupervised Learning**
4. Self-supervised Learning
5. Predictive methods
6. Contrastive methods



# Supervised learning

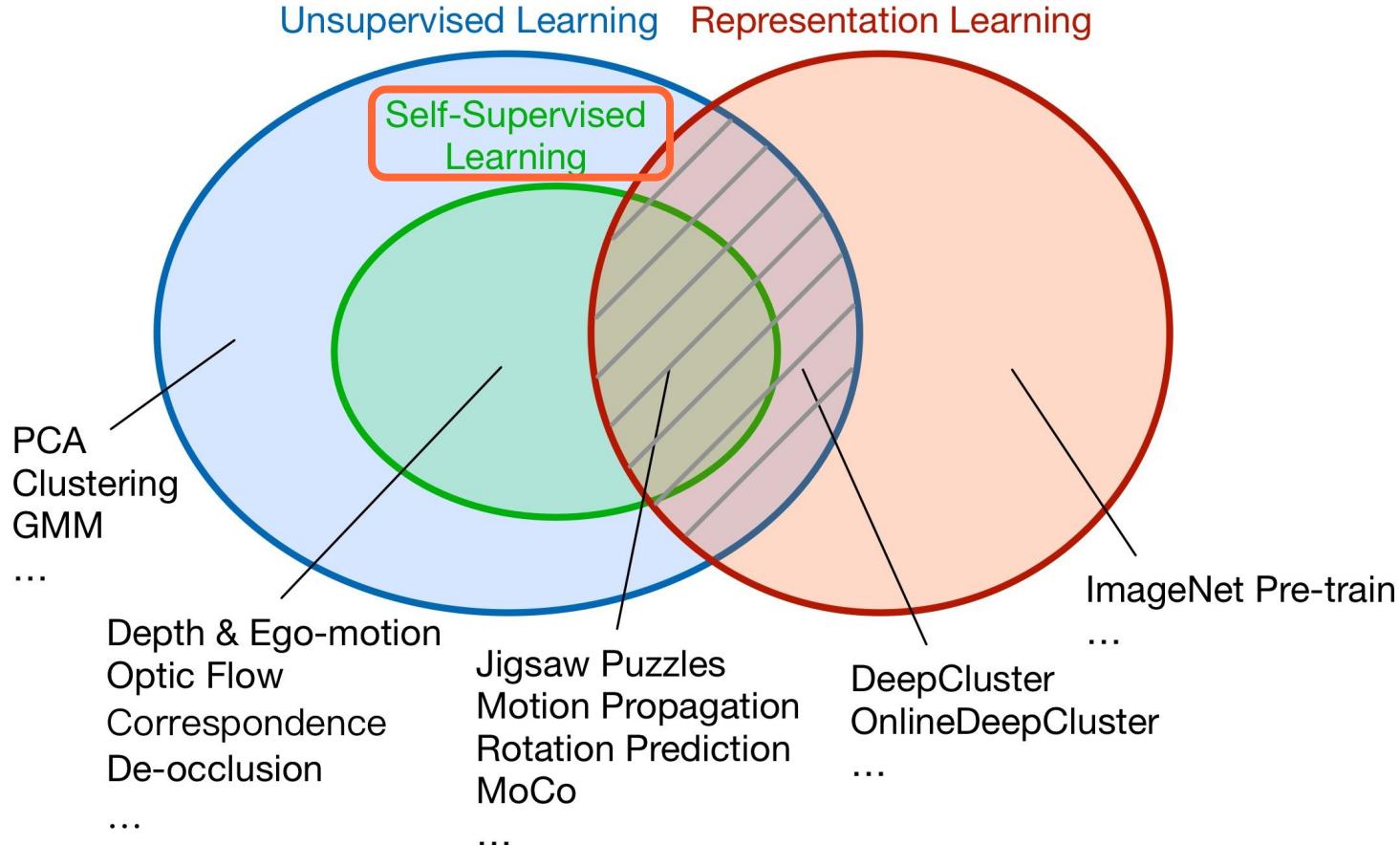


# Unsupervised learning



# Outline

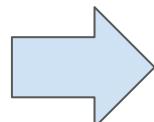
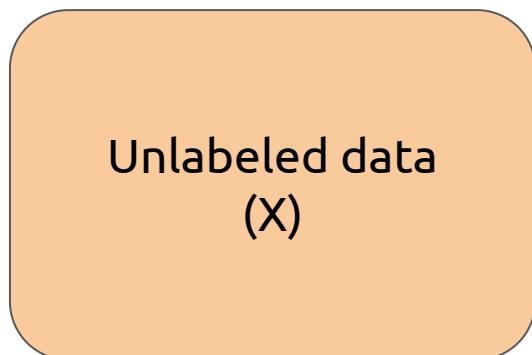
1. Representation Learning
2. Unsupervised Learning
3. **Self-supervised Learning**
4. Predictive methods
5. Contrastive methods



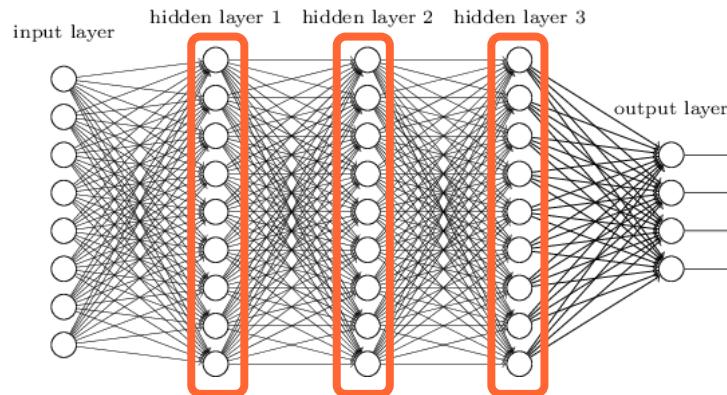
# Self-supervised learning

**Self-supervised learning** is a form of unsupervised learning where raw data (no annotations) is enough to compute the gradients.

- A **pretext (or surrogate) learning task** must be designed.
- The NN **learns representations**, which should be transferable to the actual **downstream task of interest**.



Representations learned without labels



# Self-supervised Learning

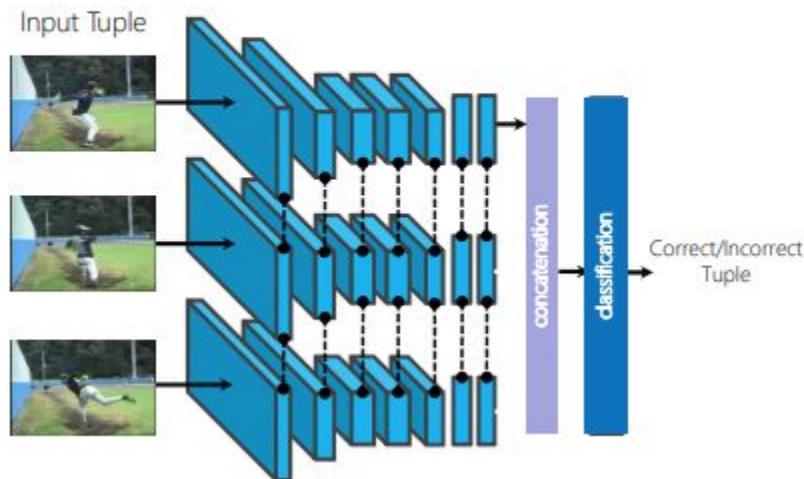
## WHY?

- Scalability.
- Avoid ambiguities in annotation criteria.

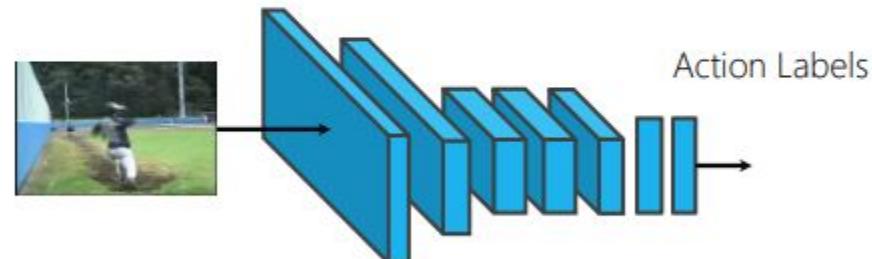


# Self-supervised + Transfer Learning

## Self-supervised Pre-train



## Test -> Finetune



#**Shuffle&Learn** Misra, Ishan, C. Lawrence Zitnick, and Martial Hebert. "[Shuffle and learn: unsupervised learning using temporal order verification.](#)" ECCV 2016. [[code](#)] ([Slides](#) by Xunyu Lin):

# Outline

1. Transfer Learning
2. Representation Learning
3. Unsupervised Learning
4. Self-supervised Learning
5. **Predictive methods**
6. Contrastive methods

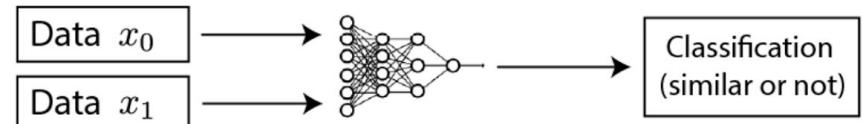
# Predictive vs Contrastive Methods

## Generative / Predictive



Loss measured in the output space

## Contrastive



Loss measured in the representation space

# Predictive vs Contrastive Methods

## Generative / Predictive



Loss measured in the output space

## Contrastive



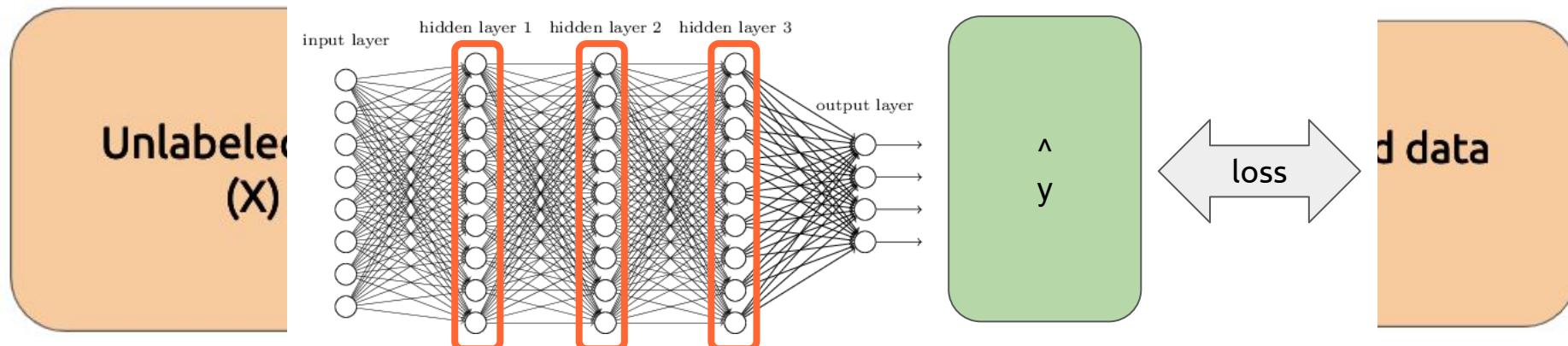
Loss measured in the representation space

# Predictive vs Contrastive Methods

## Predictive Methods:

A **pretext (or surrogate) task** is defined by withholding a part of the unlabeled data and training the NN to predict it.

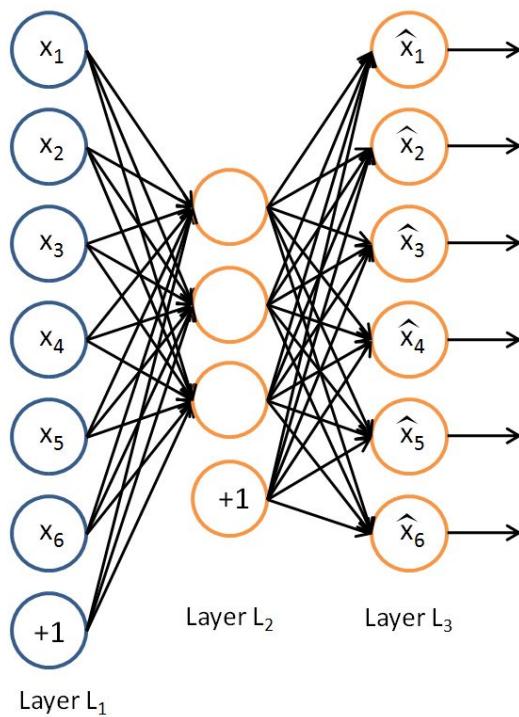
Representations learned without labels



# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
  - **(Autoencoder)**

# Autoencoder (AE)



Autoencoders:

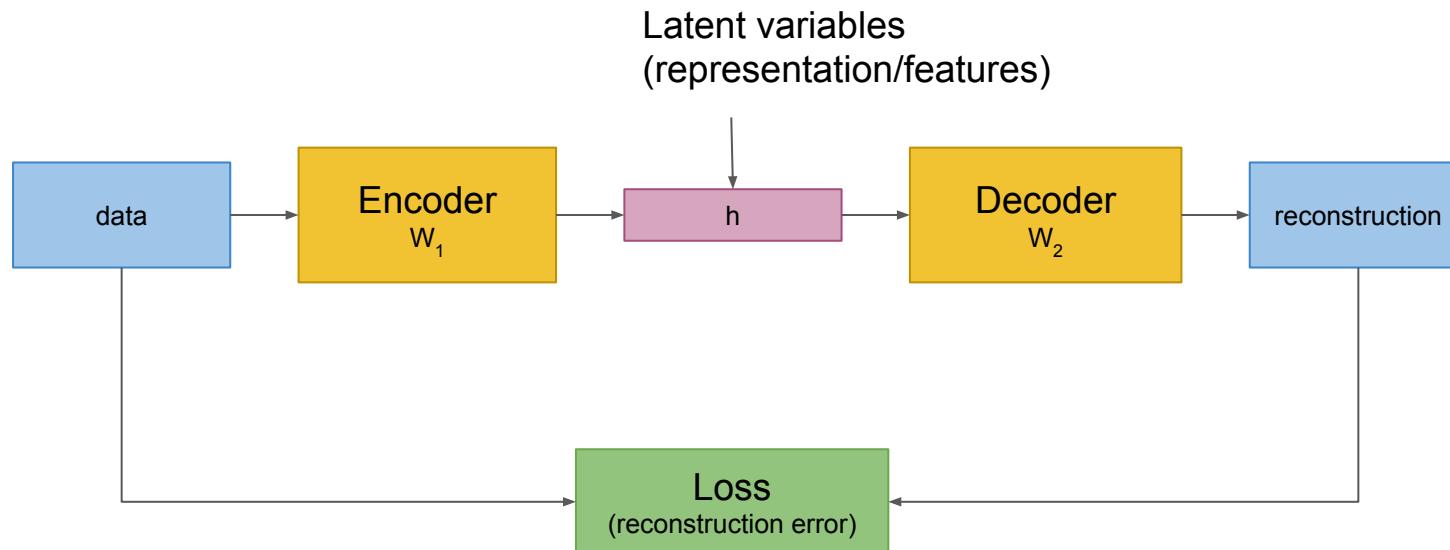
- Predict at the output the same input data.
- Do not need labels.

# Autoencoder (AE)

WHY?



1. Initialize a NN by solving an autoencoding problem.

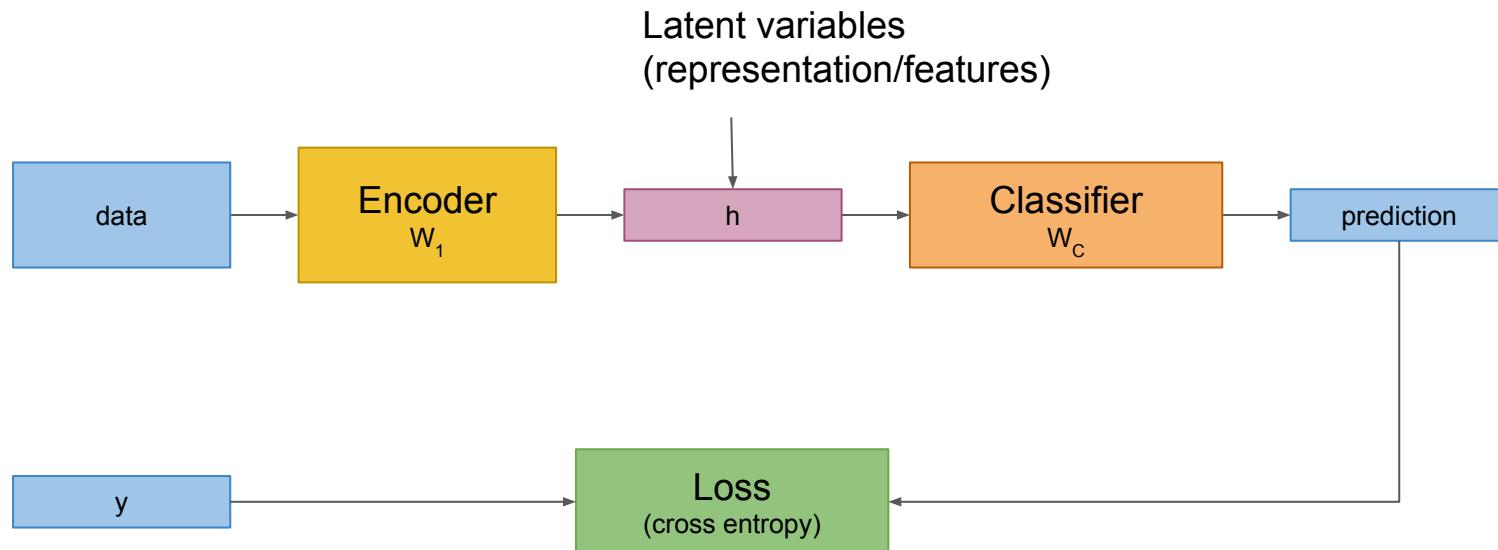


# Autoencoder (AE)

## WHY?



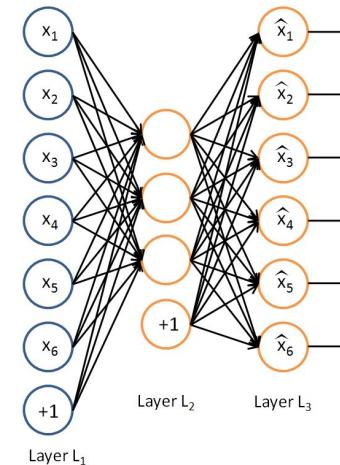
1. Initialize a NN solving an autoencoding problem.
2. Train for final task with “few” labels.



# Autoencoder (AE)

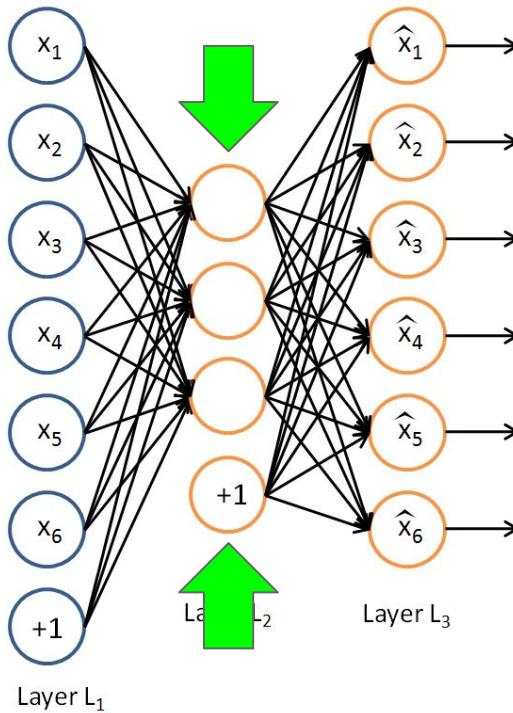
What other application does the AE target address ?

# WHY?



# Autoencoder (AE)

# WHY?

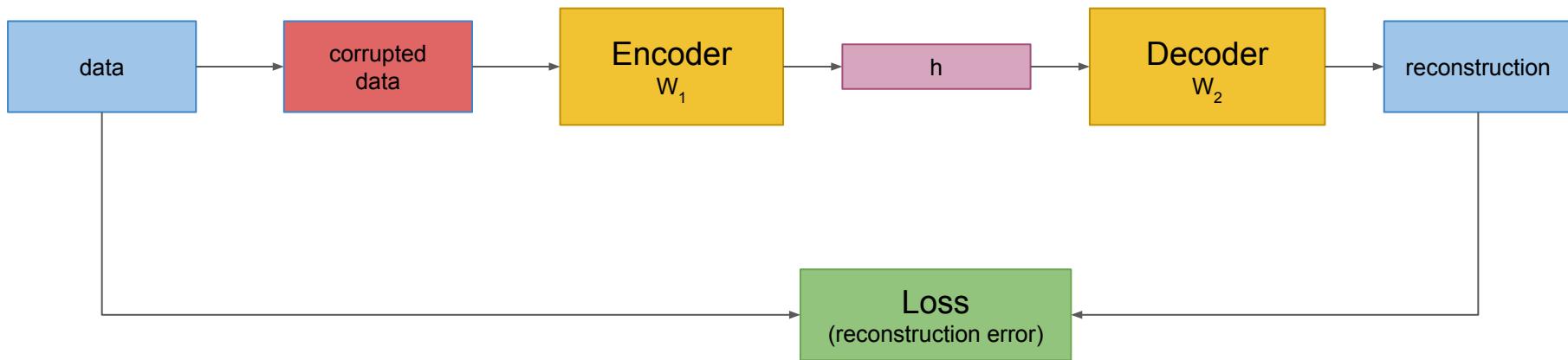


Dimensionality reduction:

Use the hidden layer as a feature extractor of any desired size.

# Denoising Autoencoder (DAE)

Corrupt the input signal, but predict the clean version at its output.



Vincent, Pascal, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. ["Extracting and composing robust features with denoising autoencoders."](#) ICML 2008.

Joan Serrà ha retuitat

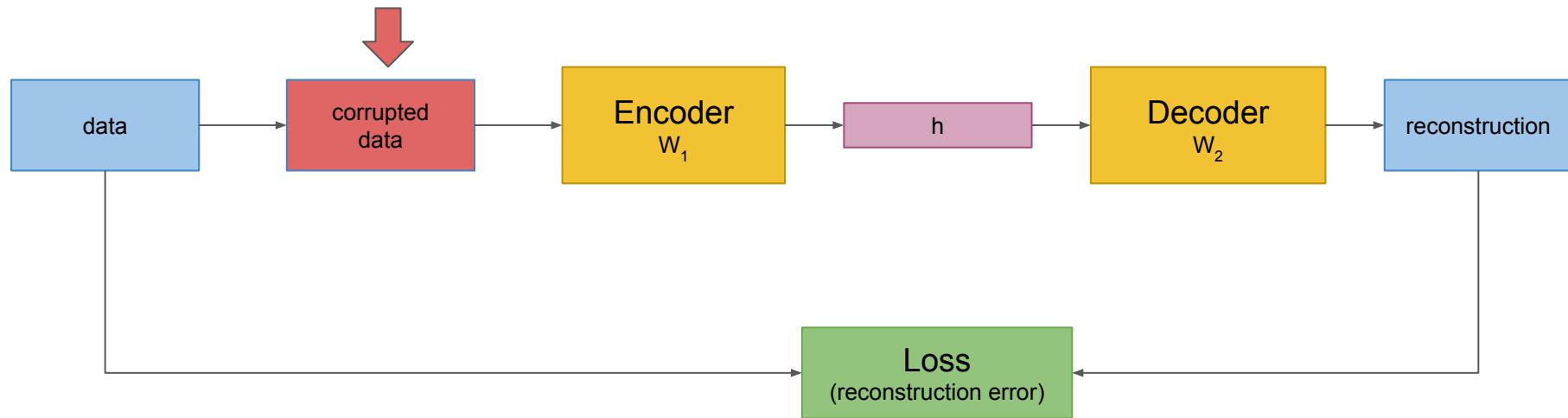


boredbengio  
@boredbengio

Confession: For the de-noising autoencoder, there was no noise. We planted the noise. So that we could remove it. [#consciousnessprior](#)

Tradueix el tuit

12:57 p. m. · 9 de set. de 2020 · Twitter Web App



Vincent, Pascal, Hugo Larochelle, **Yoshua Bengio**, and Pierre-Antoine Manzagol. "[Extracting and composing robust features with denoising autoencoders.](#)" ICML 2008.

# Denoising Autoencoder (DAE)





Eigen, David, Dilip Krishnan, and Rob Fergus. ["Restoring an image taken through a window covered with dirt or rain."](#) ICCV 2013.

# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
  - **Future Prediction**

# Future Prediction

## How Much Information is the Machine Given during Learning?

Y. LeCun

### ► “Pure” Reinforcement Learning (**cherry**)

- The machine predicts a scalar reward given once in a while.

### ► A few bits for some samples

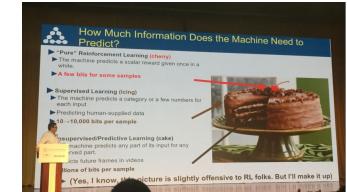


### ► Supervised Learning (**icing**)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- **10 → 10,000 bits per sample**

### ► Self-Supervised Learning (**cake génoise**)

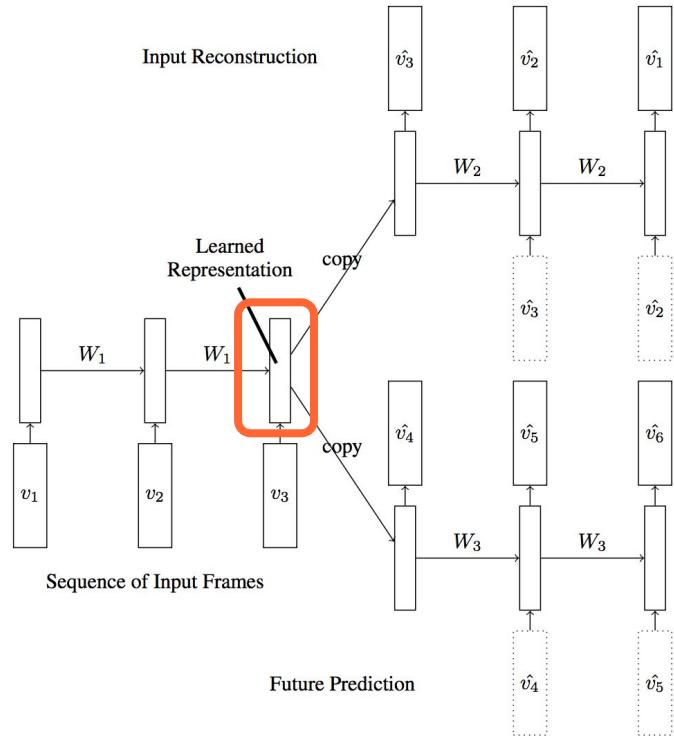
- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- **Millions of bits per sample**



Slide credit:  
Yann LeCun

# Future Prediction: Video Frames

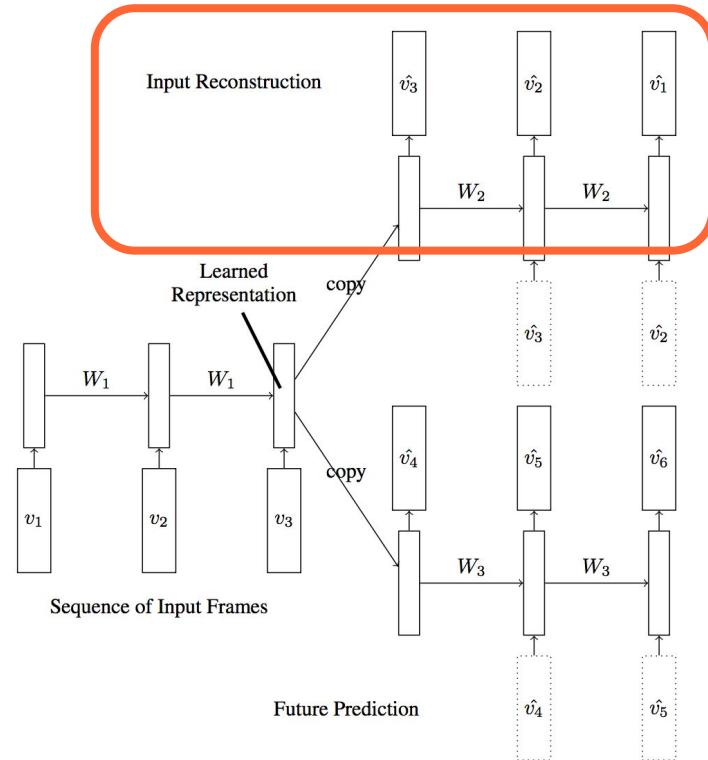
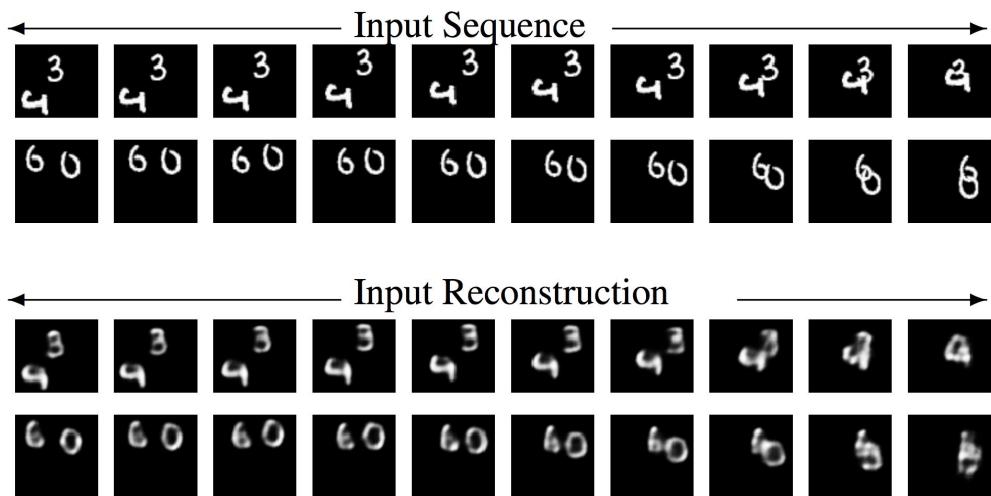
Learning video representations (features) by...



# Future Prediction: Video Frames

# Learning video representations (features) by...

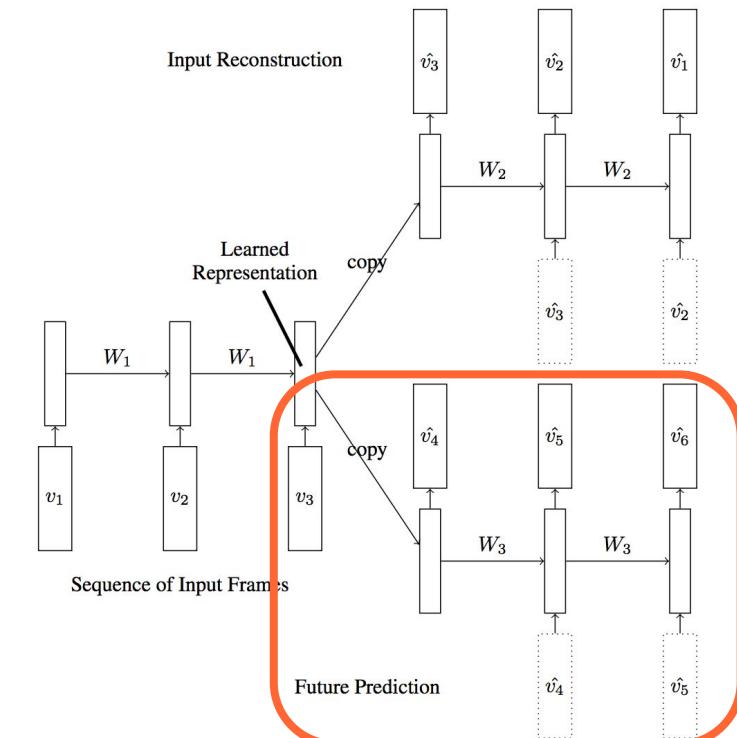
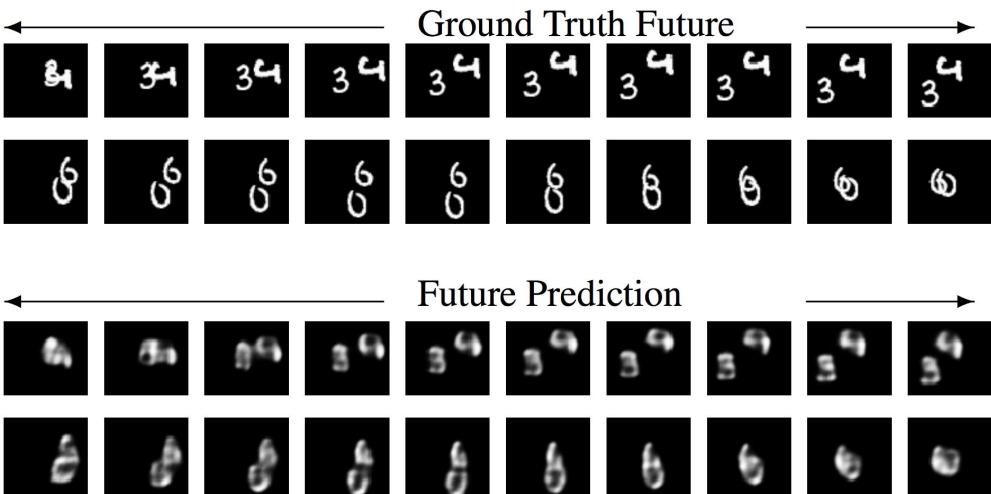
(1) frame reconstruction (AE):



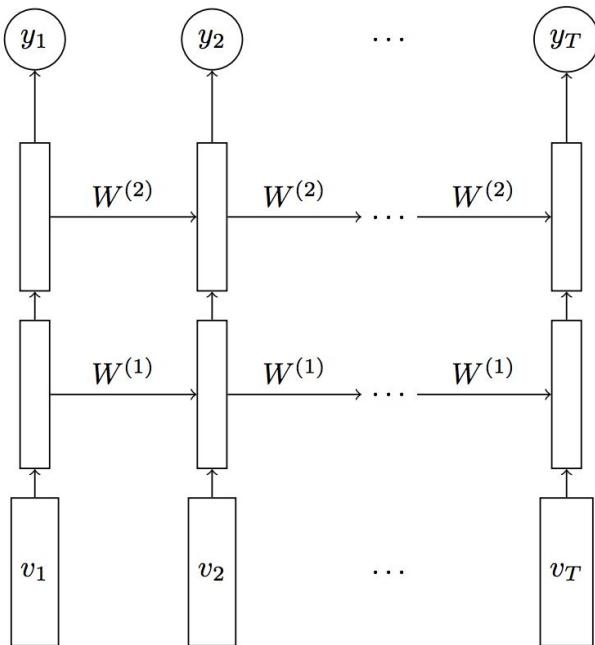
# Future Prediction: Video Frames

Learning video representations (features) by...

## (2) frame prediction



# Future Prediction: Video Frames



Unsupervised learned features (lots of data) are fine-tuned for activity recognition (small data).

Model	UCF-101	UCF-101	HMDB-51
	RGB	1-frame flow	RGB
Single Frame	72.2	72.2	40.1
LSTM classifier	74.5	74.3	42.8
Composite LSTM	<b>75.8</b>	<b>74.9</b>	<b>44.1</b>
Model + Finetuning			

Table 1. Summary of Results on Action Recognition.

Figure 6. LSTM Classifier.

# Future Prediction: Words

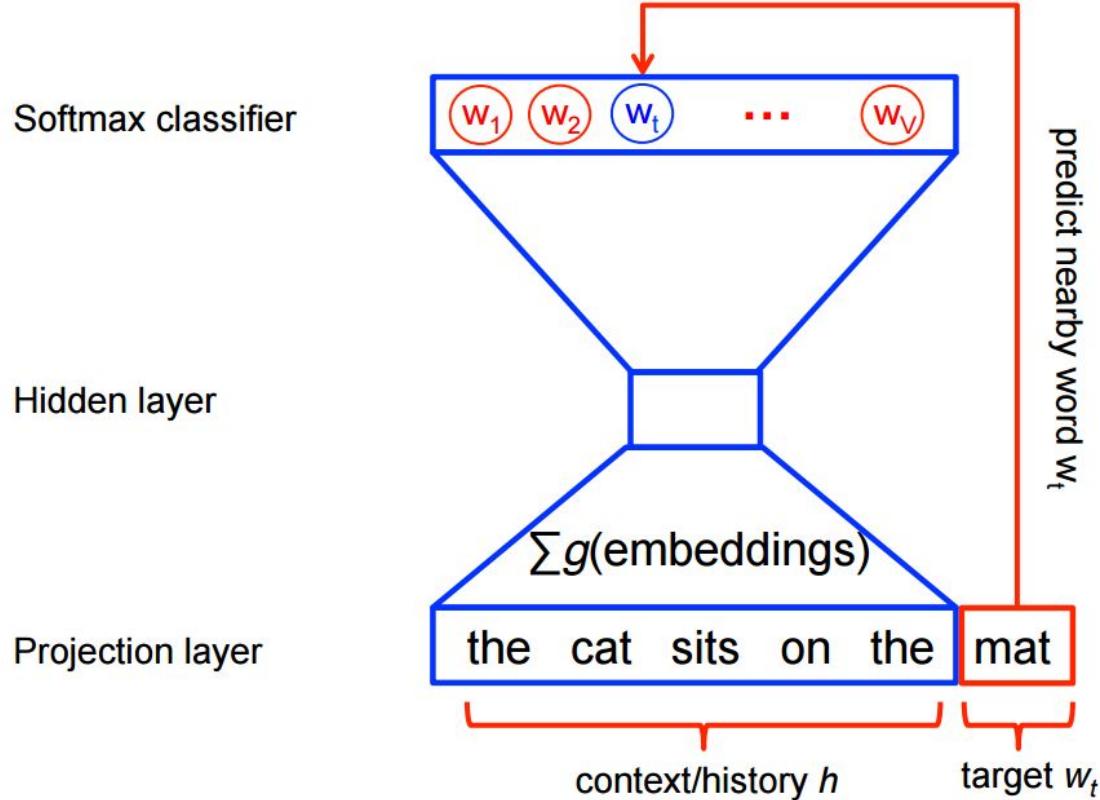


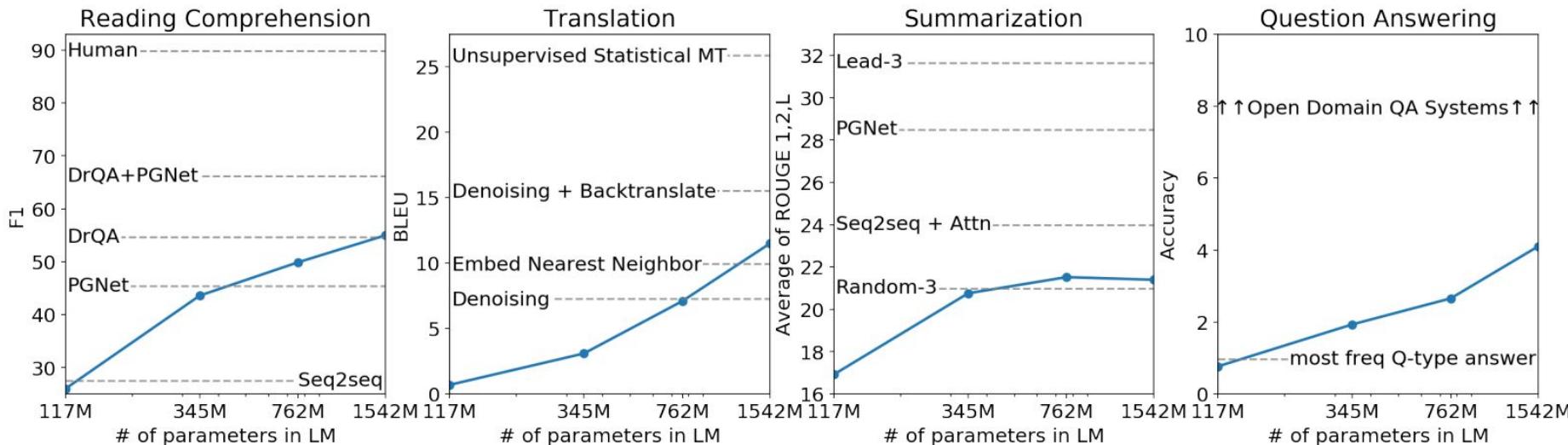
Figure:  
[TensorFlow tutorial](#)

Bengio, Yoshua, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. "[A neural probabilistic language model.](#)" Journal of machine learning research 3, no. Feb (2003): 1137-1155.

# Future Prediction: Words

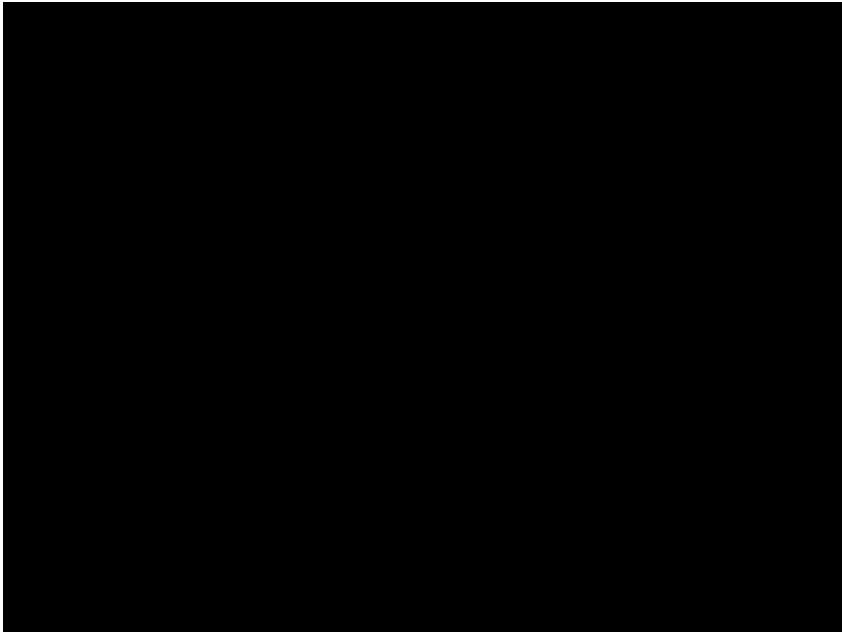
“GPT-2 is trained with a simple objective: predict the next word, given all of the previous words within some text.”

**Zero-shot task performances  
(GPT-2 was never trained for these tasks)**

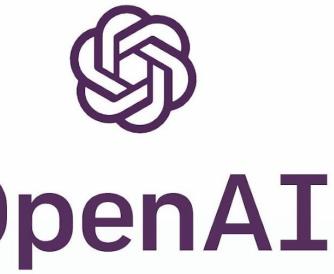


# Future Prediction: Words

Training a very large language model with a very large amount of data can result in impressive results:



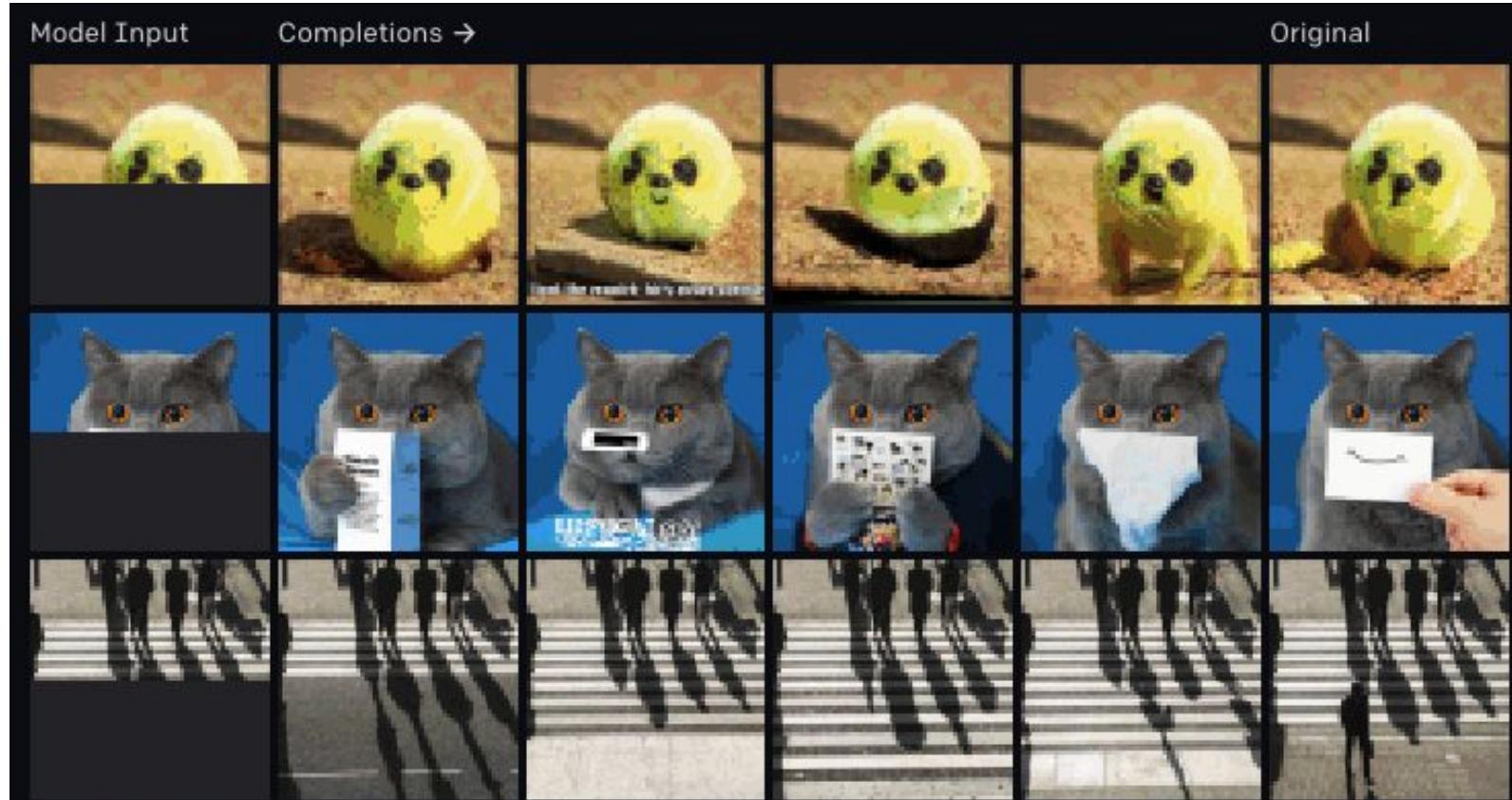
Video Source:  
[@mattshumer](#)



[[OpenAI API](#)]

#GPT-3 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Agarwal, S. [Language models are few-shot learners](#). NeurIPS 2020 (best paper award).

# Future Prediction: Pixels



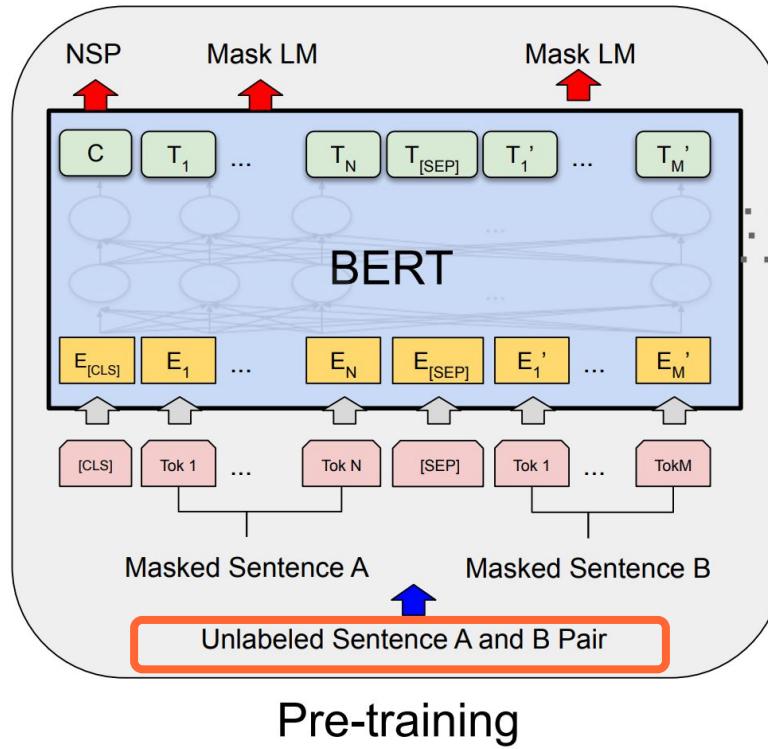
# Future Prediction: Pixels

EVALUATION	MODEL	ACCURACY	PRE-TRAINED ON IMAGENET	
			W/O LABELS	W/ LABELS
CIFAR-10 Linear Probe	ResNet-152 <sup>50</sup>	94.0		✓
	SimCLR <sup>12</sup>	95.3	✓	
	iGPT-L 32x32	<b>96.3</b>	✓	
CIFAR-100 Linear Probe	ResNet-152	78.0		✓
	SimCLR	80.2	✓	
	iGPT-L 32x32	<b>82.8</b>	✓	

# Future Prediction: True / False

The BERT language model is pre-trained for a **binarized next sentence prediction task**:

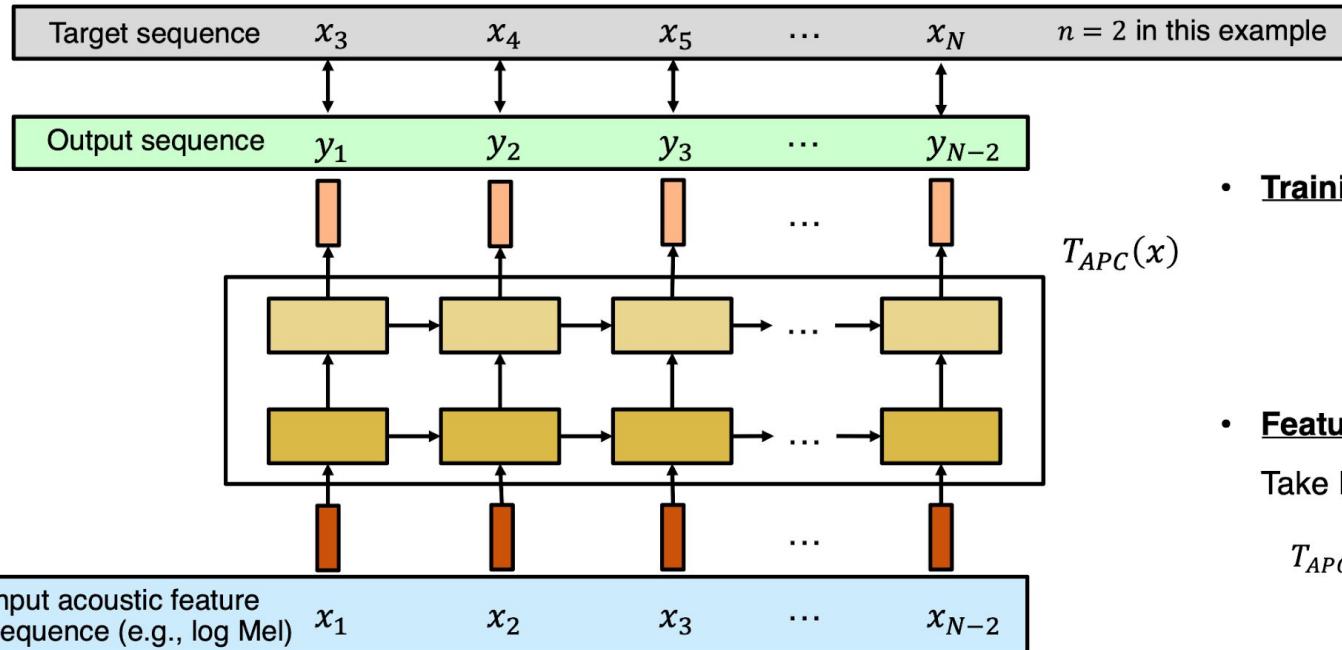
Select some sentences for A and B, where 50% of the data B is the next sentence of A, and the remaining 50% of the data B are randomly selected in the corpus, and learn the correlation.



Pre-training

# Future Prediction: Audio Spectrogram

Predict spectrogram  $n$  samples in the future.



- **Training**

$$\underset{\{RNN, W\}}{\operatorname{argmin}} \sum_{i=1}^{N-n} |x_{i+n} - y_i|,$$

$$y_i = RNN(x_i) \cdot W$$

- **Feature extraction**

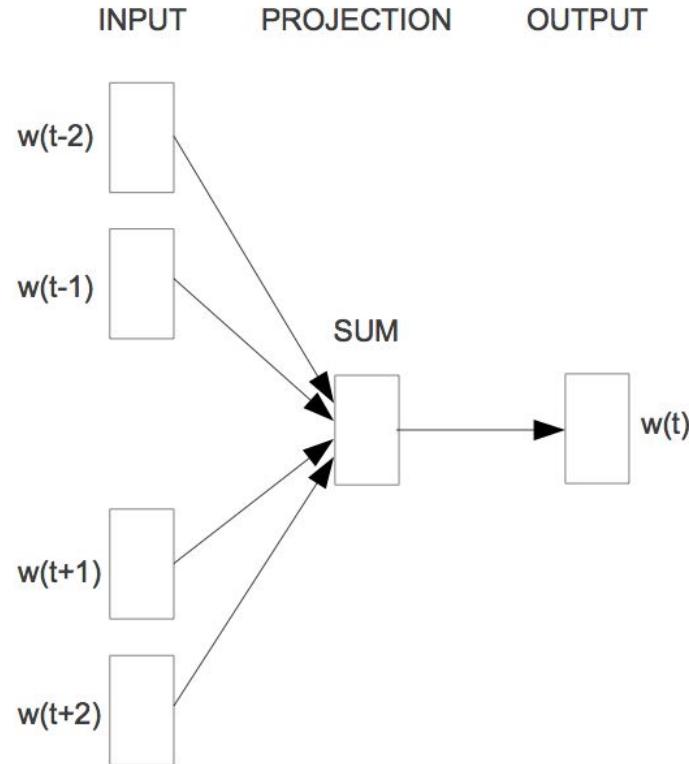
Take RNN output of each time step:

$$T_{APC}(x_i) = RNN(x_i) \quad \forall i = 1, 2, \dots, N$$

# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
  - **Unmasking**

# Unmask: One word given a context



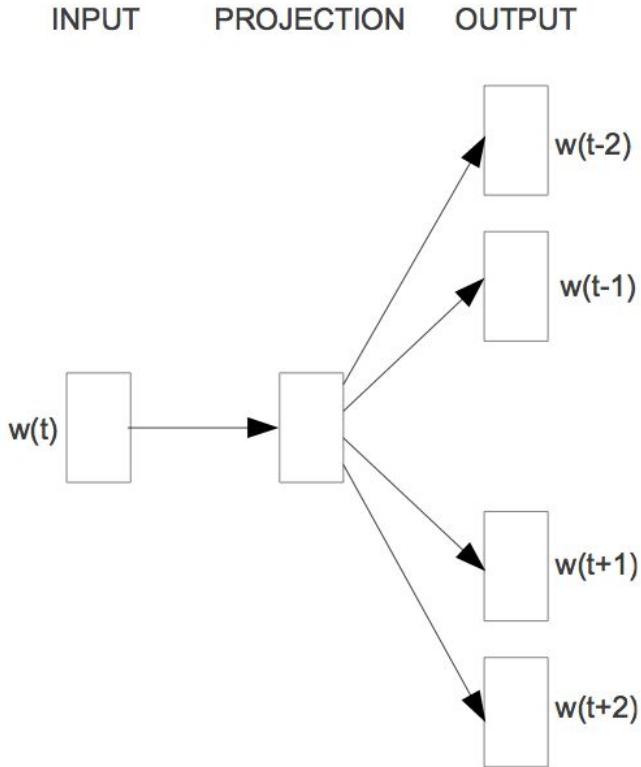
the cat climbed a tree

Given context:

a, cat, the, tree

Estimate prob. of  
climbed

# Unmask: A context given a word



the cat **climbed** a tree

Given word:

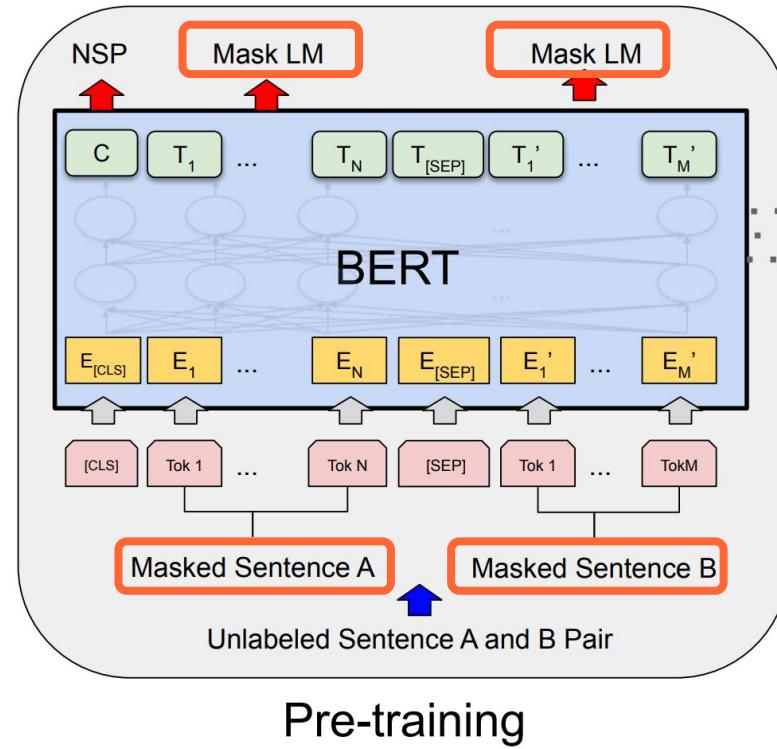
climbed

Estimate prob. of context words:

a, cat, the, tree

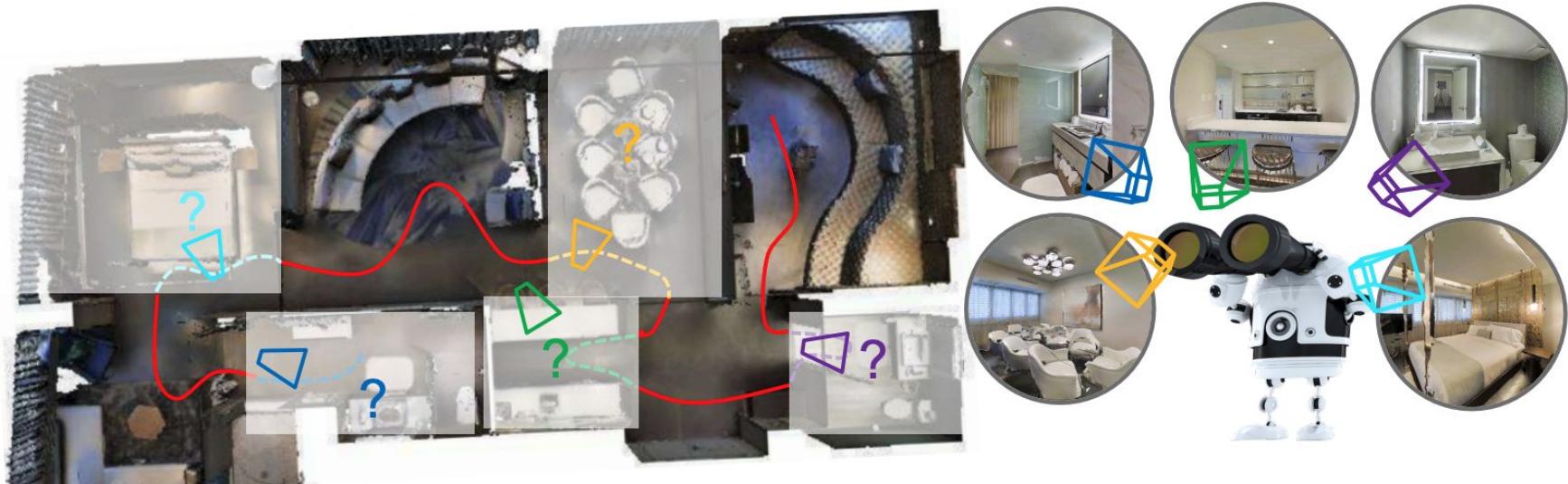
# Unmask: Some words given a context

Take a sequence of words on input, mask out 15% of the words, and ask the system to predict the missing words (or a distribution of words)



# Unmask: Embodied agent trajectories

Mask out portions of an agent's trajectory and predict them from the unmasked portions, conditioned on the agent's camera poses.

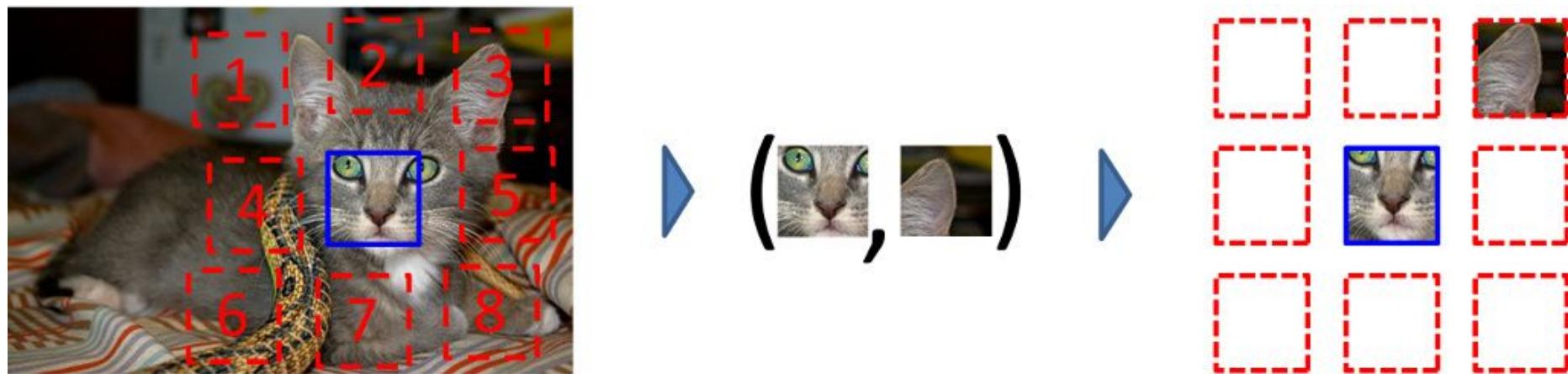


# Outline

1. Transfer Learning
2. Unsupervised Learning
3. Self-supervised Learning
4. Predictive Methods
  - **Spatial relations**

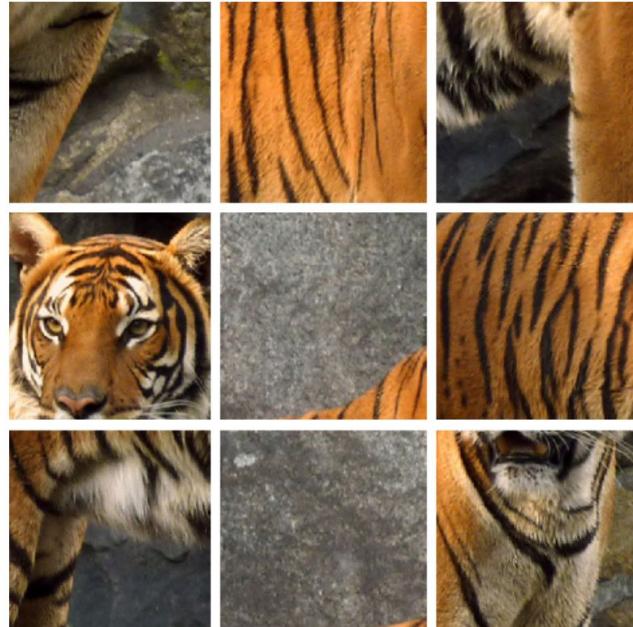
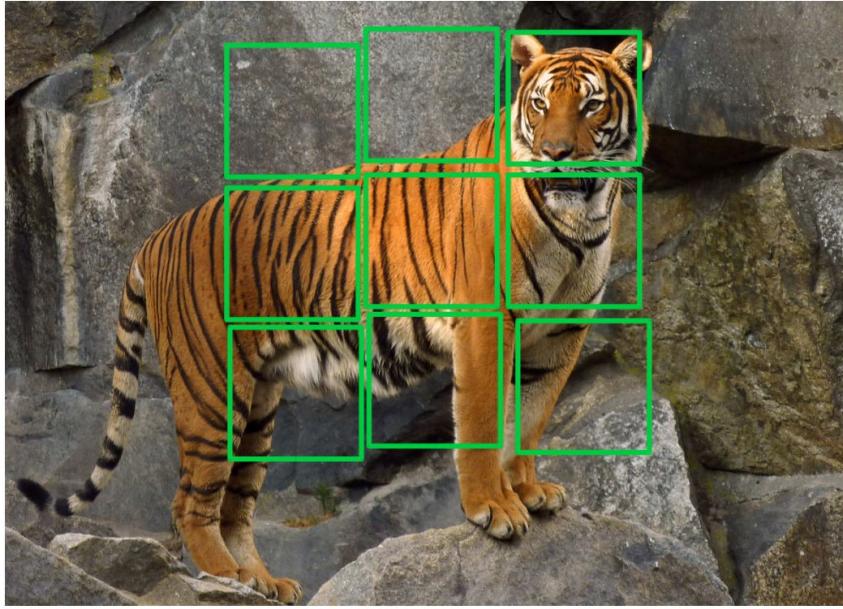
# Spatial verification

Predict the relative position between two image patches.



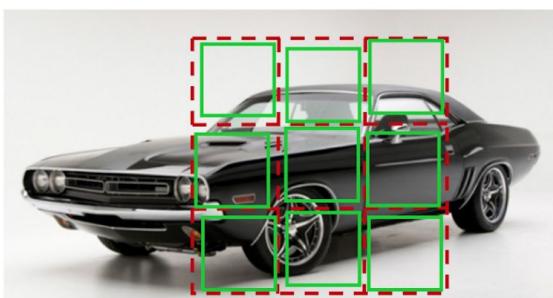
# Spatial verification

Train a neural network to solve a jigsaw puzzle.



# Spatial verification

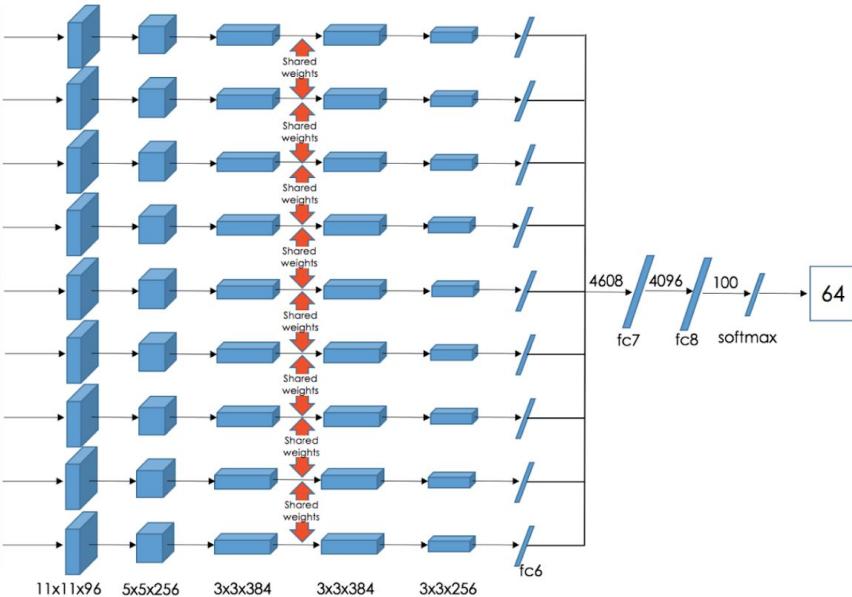
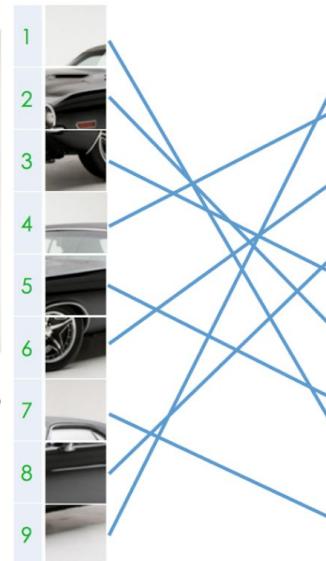
Train a neural network to solve a jigsaw puzzle.



Permutation Set

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

Reorder patches according to the selected permutation



# Outline

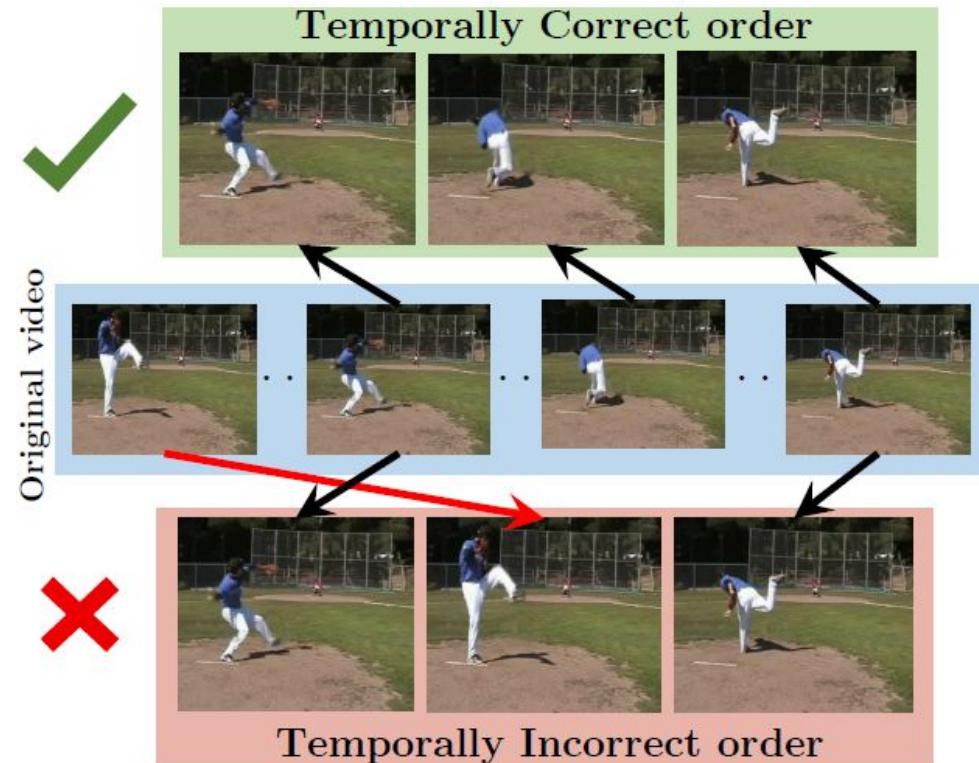
1. Transfer Learning
2. Unsupervised Learning
3. Self-supervised Learning
4. Predictive Methods
  - **Temporal relations**



Could you think about self-supervised learning tasks exploiting temporal relations ?

# Temporal coherence

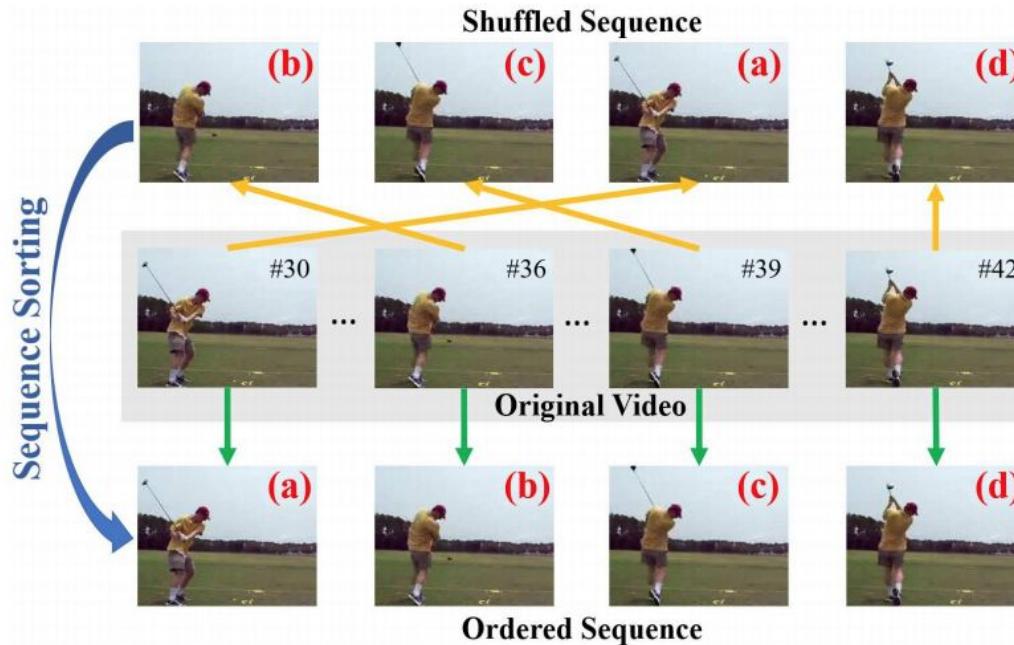
Temporal order of frames is exploited as the supervisory signal for learning.



#**Shuffle&Learn** Misra, Ishan, C. Lawrence Zitnick, and Martial Hebert. "[Shuffle and learn: unsupervised learning using temporal order verification.](#)" ECCV 2016. [[code](#)] ([Slides](#) by Xunyu Lin):

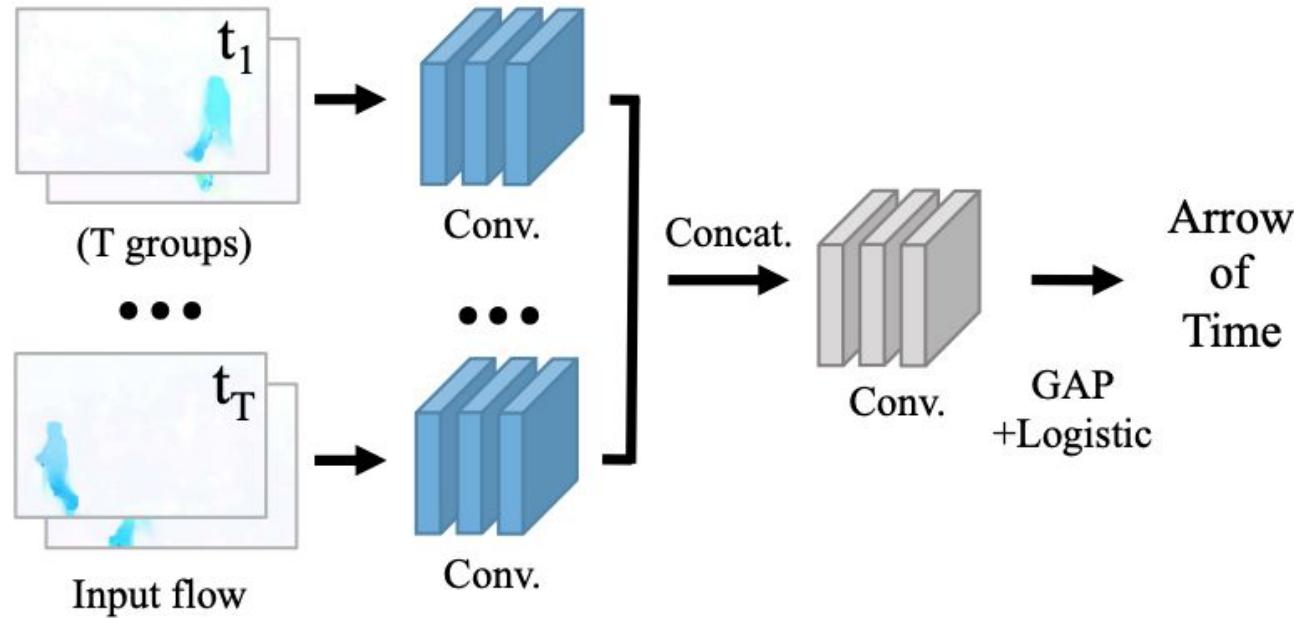
# Temporal sorting

Sort the sequence of frames.



# Arrow of time

Predict whether the video moves forward or backward.



# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
  - **Cycle Consistency**

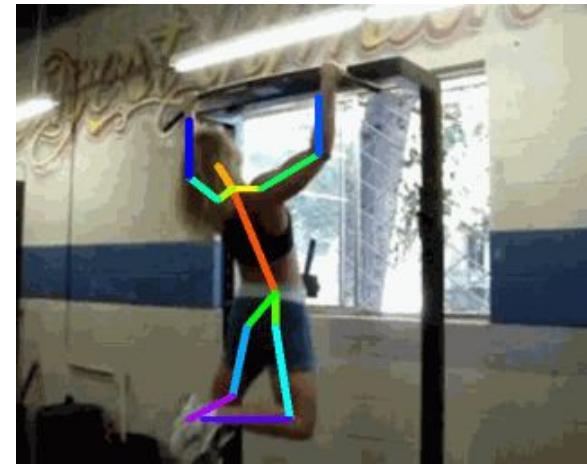
# Cycle-consistency

**Pre-text task:** Tracking backward and then forward with pixel embeddings, to later train a NN to match the pairs of associated pixel embeddings.

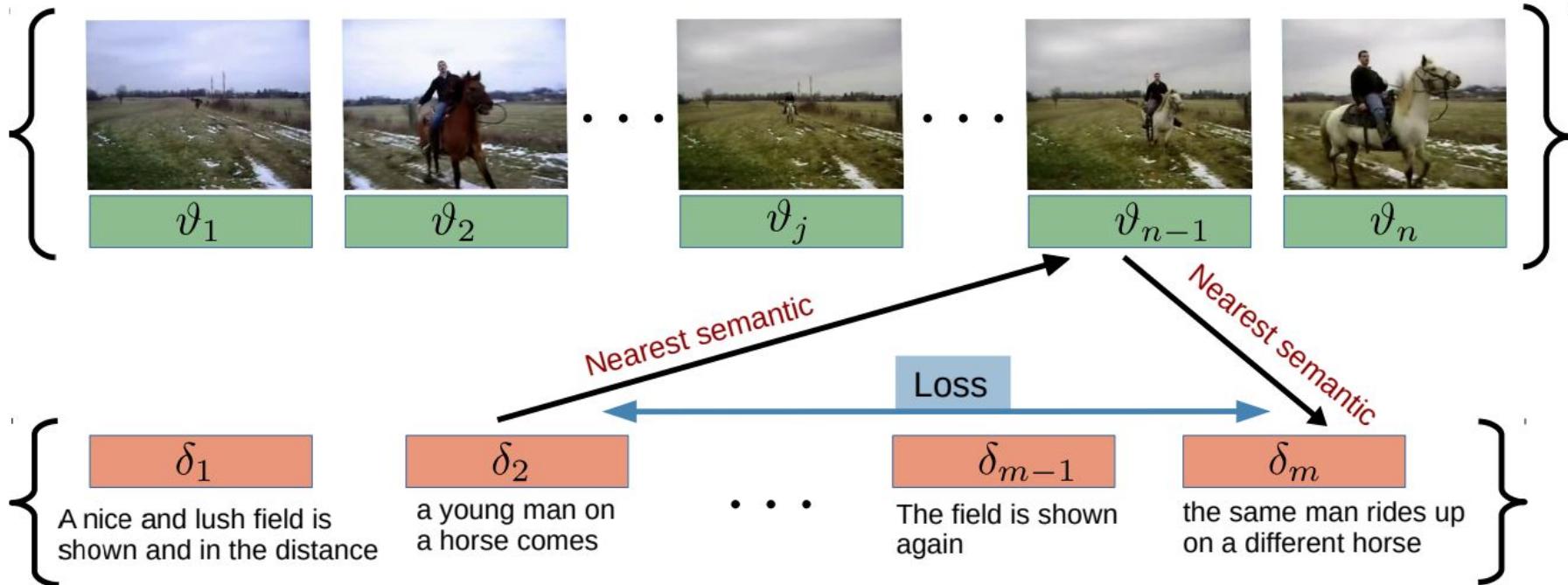


# Cycle-consistency

**Downstream task:** Track pixels across video sequences.



# Cycle-consistency

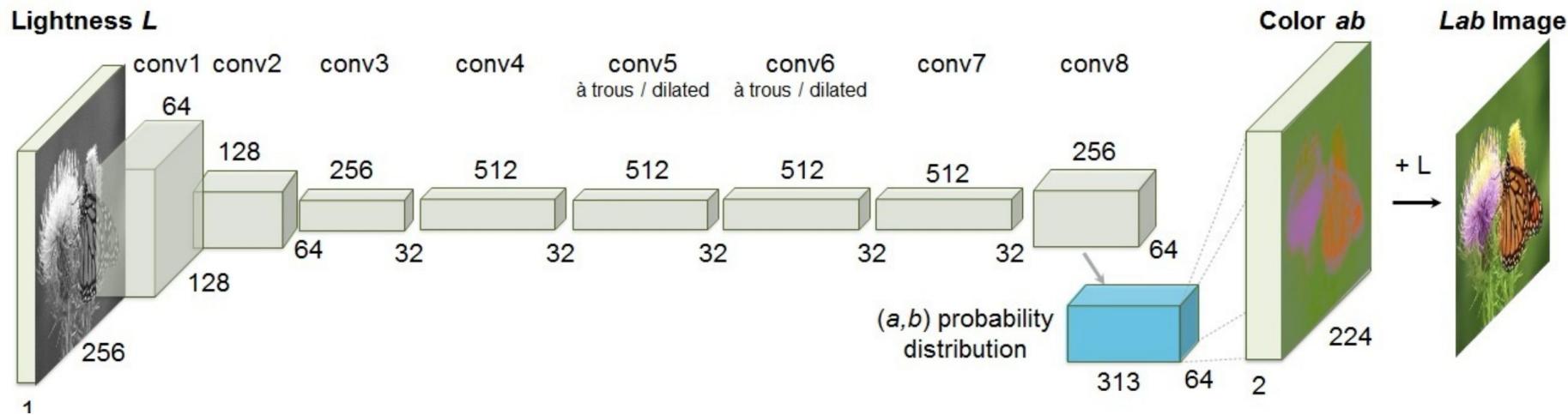


# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
  - **Colorization**

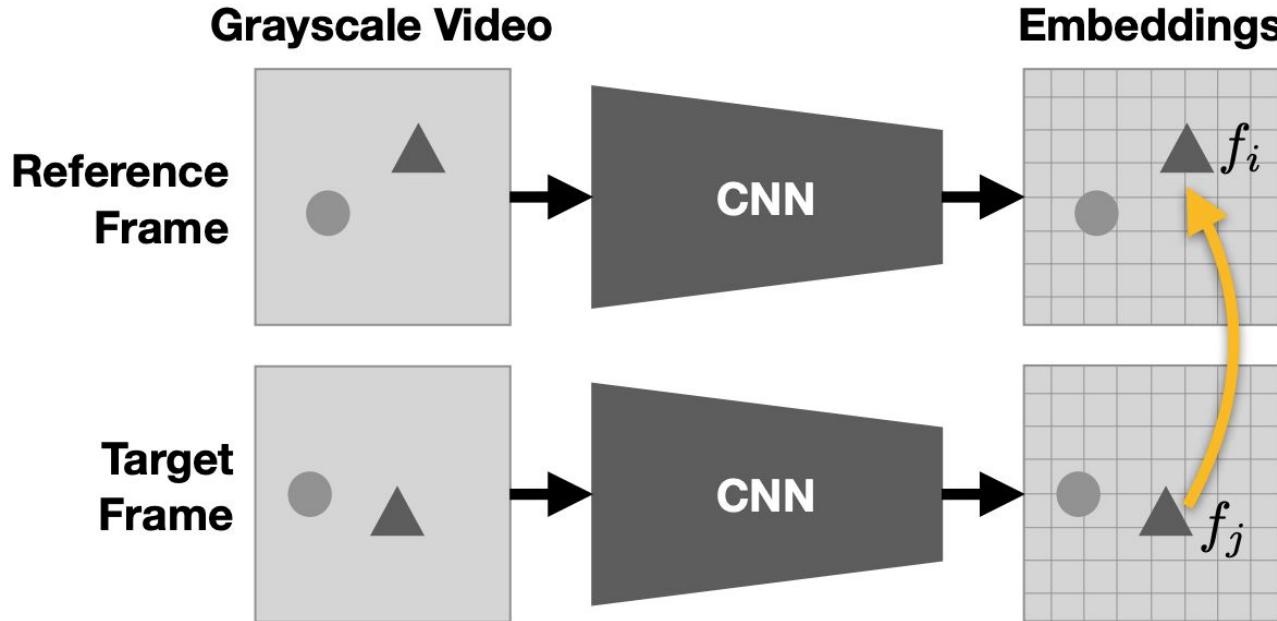
# Colorization (still image)

**Surrogate task:** Colorize a gray scale image.



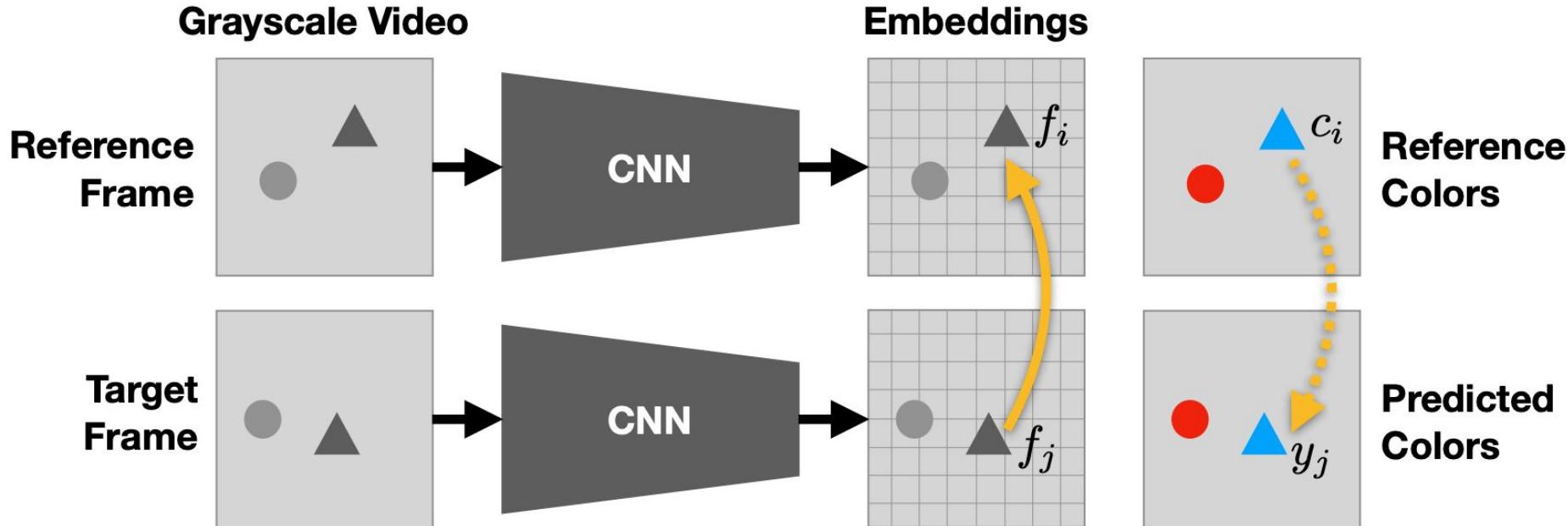
# Pixel tracking

Pre-training: A NN is trained to point the location in a reference frame with the most similar color...



# Pixel tracking

Pre-training: A NN is trained to point the location in a reference frame with the most similar color... computing a loss with respect to the actual color.



# Pixel tracking

For each grayscale pixel to colorize, find the color of most similar pixel embedding (representation) from the first frame.

Reference Frame



What color is this?

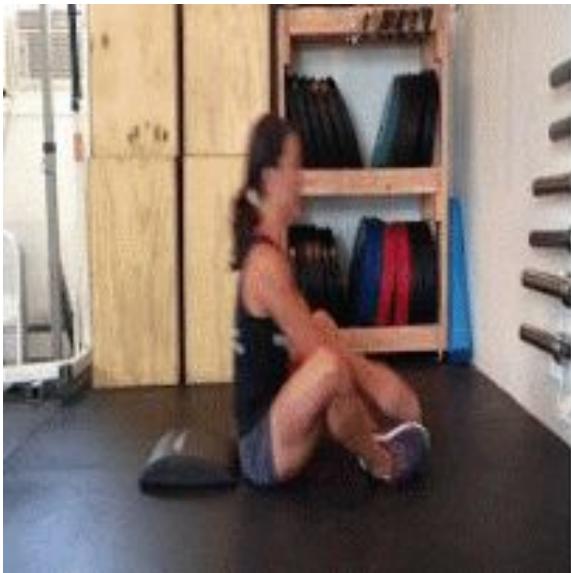


Vondrick, Carl, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. "[Tracking emerges by colorizing videos.](#)" ECCV 2018. [\[blog\]](#)

# Pixel tracking

**Application:** colorization from a reference frame.

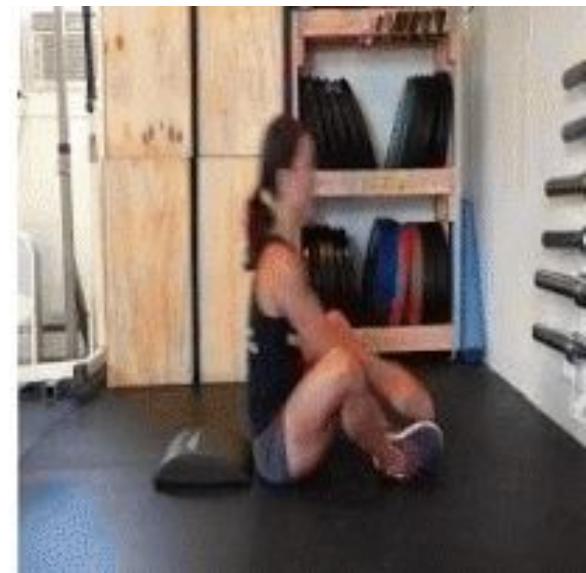
Reference frame



Grayscale video



Colorized video

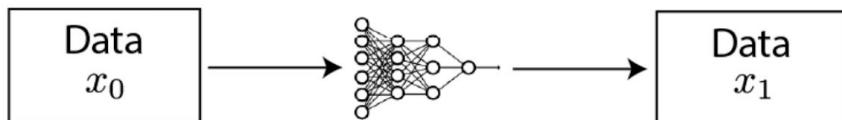


# Outline

1. Unsupervised Learning
2. Self-supervised Learning
3. Predictive Methods
4. **Contrastive Methods**

# Predictive vs Contrastive Methods

## Generative / Predictive



Loss measured in the output space

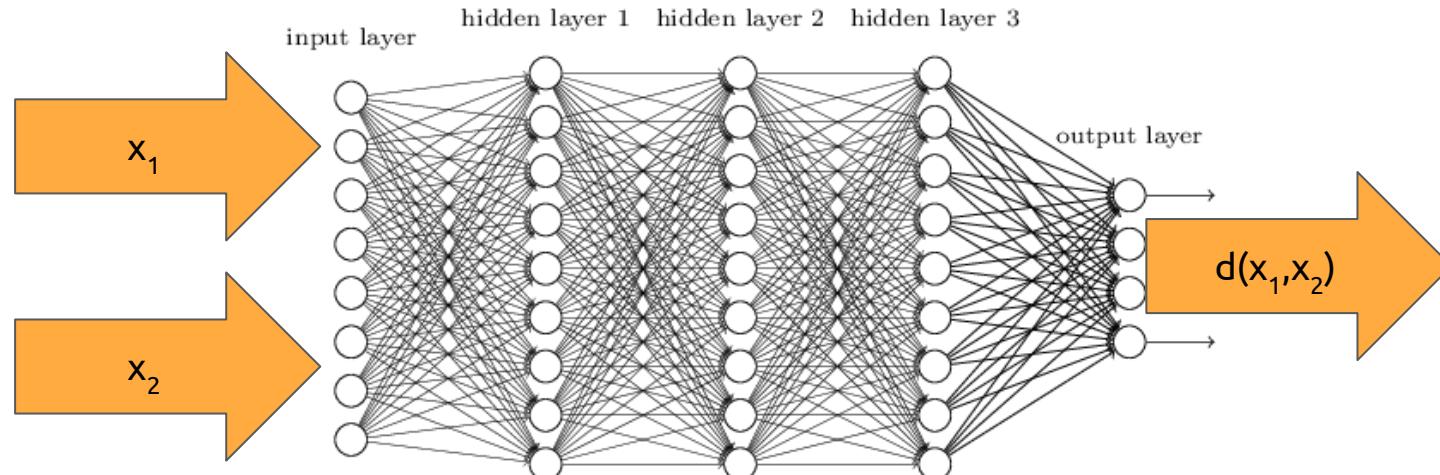
## Contrastive



Loss measured in the representation space

# Ranking Losses

Ranking losses aim at predicting the relative distance (or similarity) between samples.



# Hinge or Margin Losses

Only pairs  $(x_1, x_2)$  closer than a margin m contribute to the penalty.

$$l_N(d_D(\mathbf{x}_1, \mathbf{x}_2)) = \max(0, m - d_D(\mathbf{x}_1, \mathbf{x}_2))$$

↑  
Hinge loss

What is the use of the Hinge Loss ?



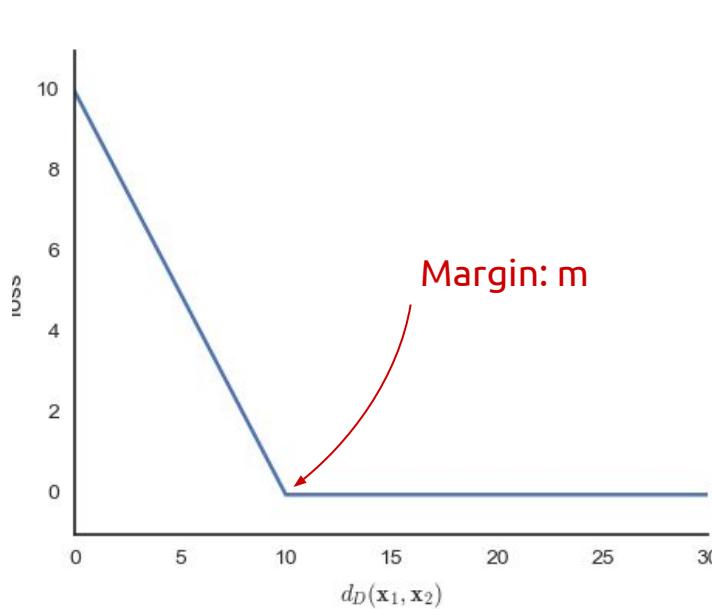
# Hinge or Margin Losses

Only pairs  $(x_1, x_2)$  closer than a margin m contribute to the penalty.

$$l_N(d_D(\mathbf{x}_1, \mathbf{x}_2)) = \max(0, m - d_D(\mathbf{x}_1, \mathbf{x}_2))$$

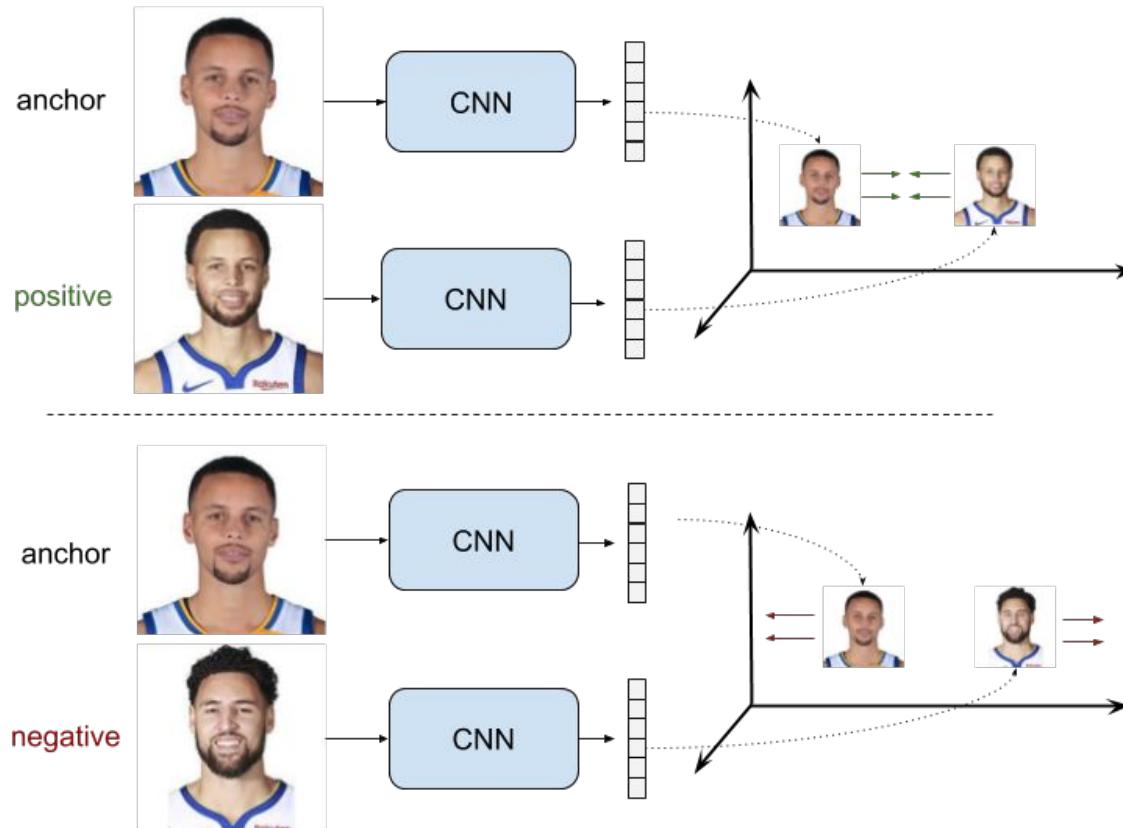


A hinge



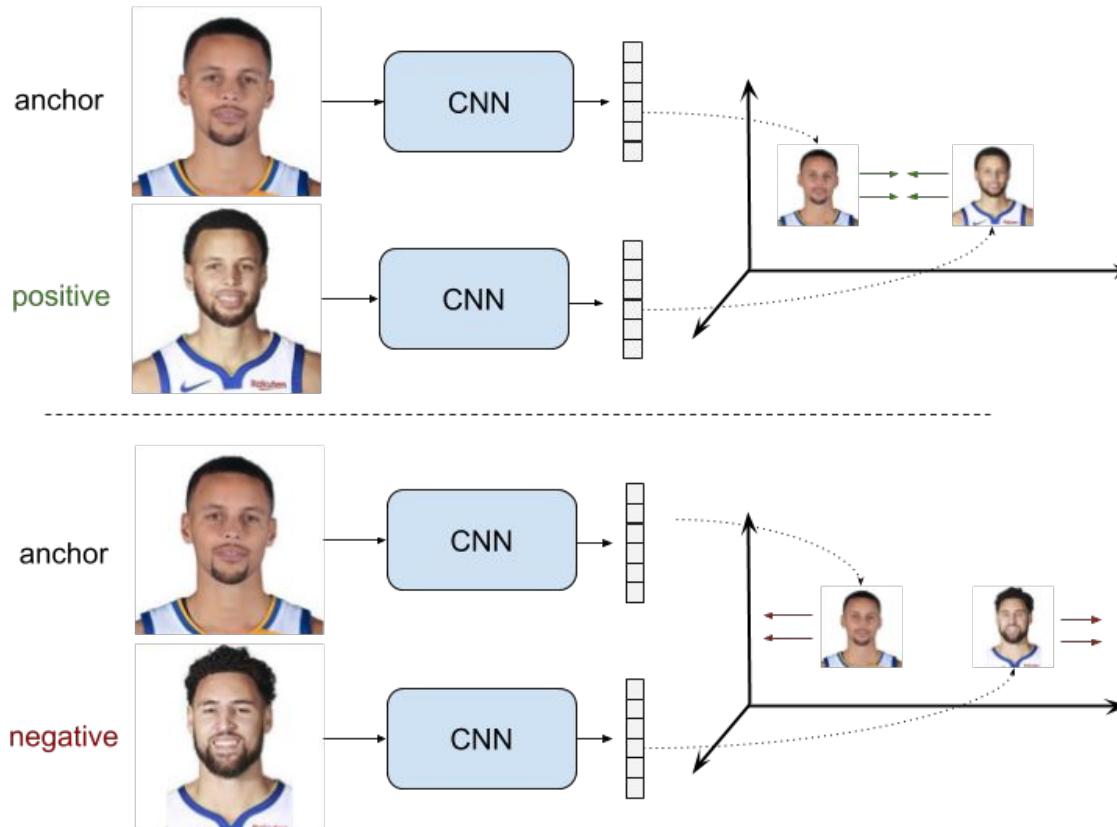
Intuition: there is nothing new to learn from those **easy negative** pairs which are already far enough.

# Pairwise Ranking Loss



Source: Raul Gómez, “[Understanding Ranking Loss, Contrastive Loss, Margin Loss, Triplet Loss, Hinge Loss and all those confusing names](#)” (2019)

# Pairwise Ranking Loss = Contrastive Loss

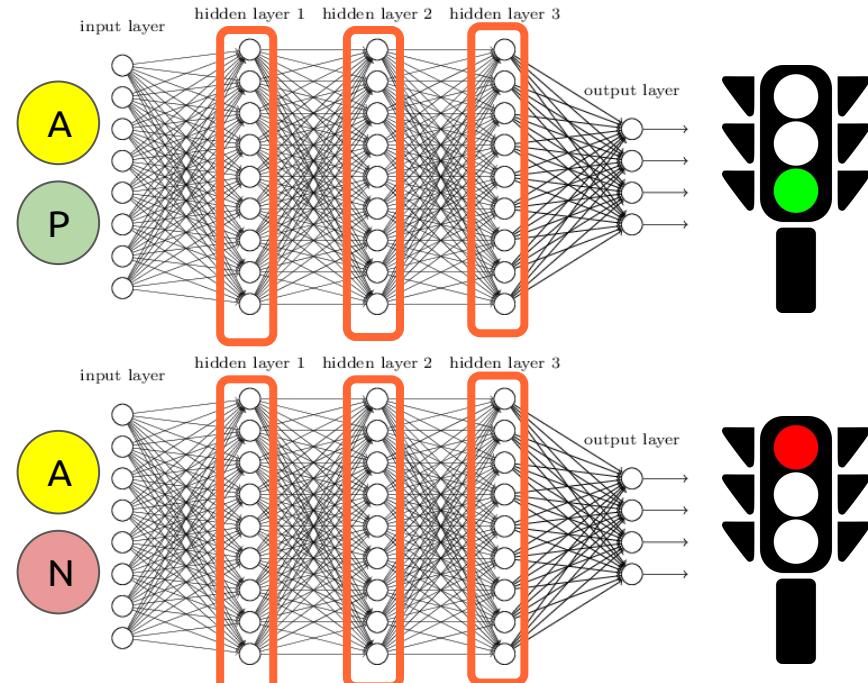
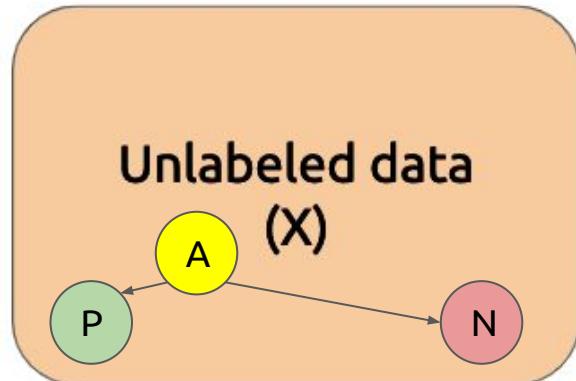


Source: Raul Gómez, “[Understanding Ranking Loss, Contrastive Loss, Margin Loss, Triplet Loss, Hinge Loss and all those confusing names](#)” (2019)

# Pairwise Ranking Loss = Contrastive Loss

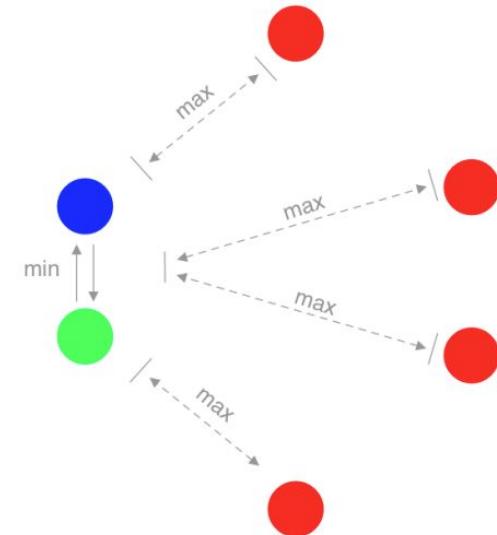
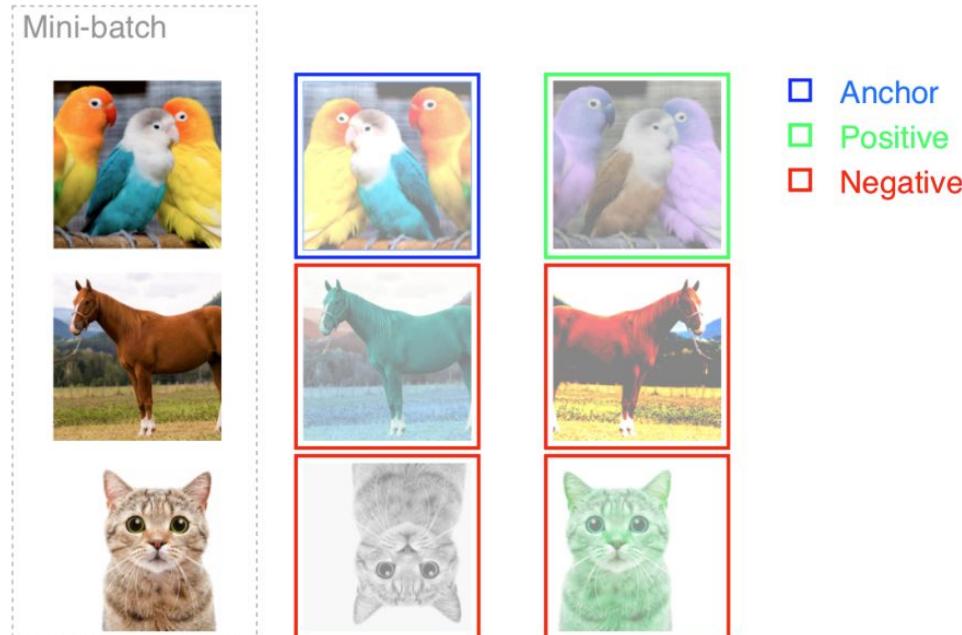
## Contrastive Self-Supervised Learning Methods:

The **pretext (or surrogate) task** corresponds to predicting which pairs of samples are similar, focusing the design in the definition of “similar”.



# Pairwise Ranking Loss = Contrastive Loss

Typically, pairs of positive and negative pairs are augmented in each training mini-batch.



# Contrastive Loss

Learn a NN  $G_w(X)$  from pairs  $(X_1, X_2)$  which are **similar** ( $Y=0$ ) or **dissimilar** ( $Y=1$ ).

$$D_W(\vec{X}_1, \vec{X}_2) = \|G_W(\vec{X}_1) - G_W(\vec{X}_2)\|_2$$

Should the distance  $D_w$  be high or low when  $Y=0$  ?

When  $Y=0$ ,  $G_w$  will learn to minimize  $D_w$ , for example, expressed as:

$$(1 - Y) \frac{1}{2} (D_W)^2$$

# Contrastive Loss

Learn a NN  $G_w(X)$  from pairs  $(X_1, X_2)$  which are **similar** ( $Y=0$ ) or **dissimilar** ( $Y=1$ ).

$$D_W(\vec{X}_1, \vec{X}_2) = \|G_W(\vec{X}_1) - G_W(\vec{X}_2)\|_2$$

Should the distance  $D_w$  be high or low when  $Y=1$  ?

When  $Y=1$ ,  $G_w$  will learn to maximize  $D_w$ , for example, expressed as:

$$(Y) \frac{1}{2} \{ \max(0, m - D_W) \}$$

# Contrastive Loss

The contrastive loss to minimize can be expressed as:

$$L(W, Y, \vec{X}_1, \vec{X}_2) =$$

$$(1 - Y) \frac{1}{2} (D_W)^2$$

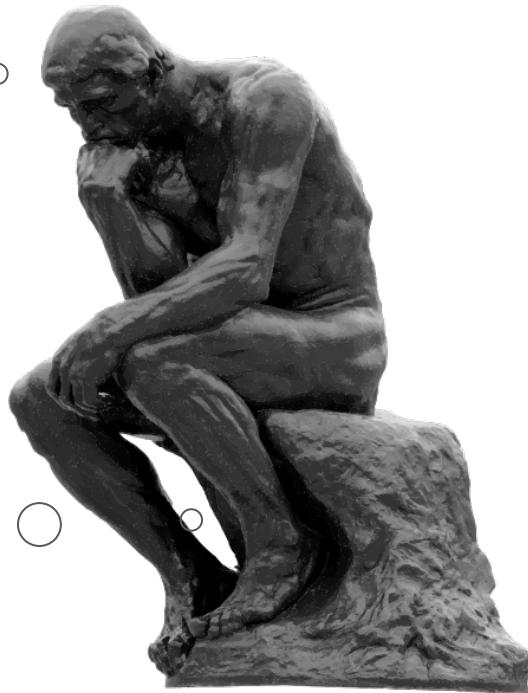
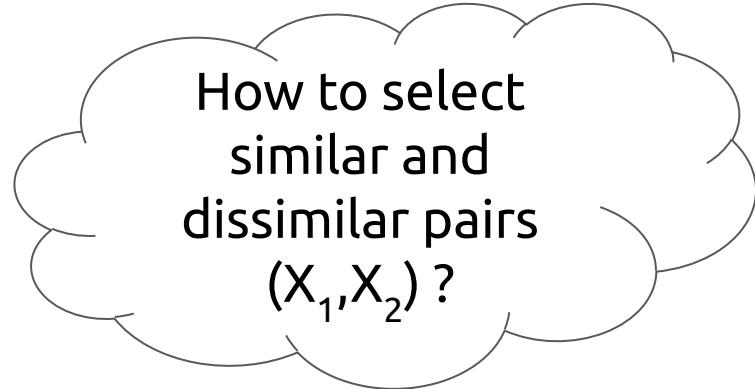
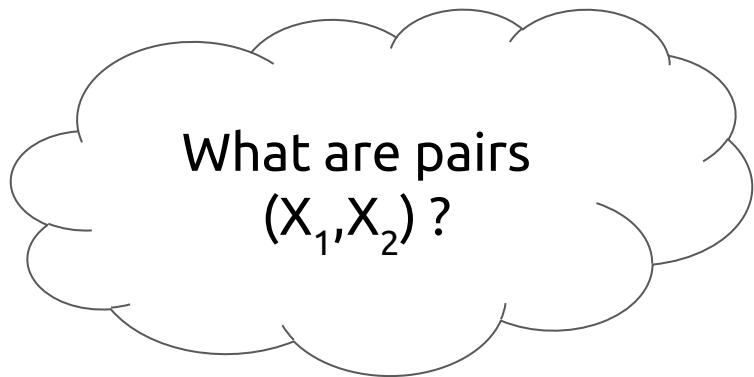
Y=0, similar pair ( $\vec{X}_1, \vec{X}_2$ )

$$(Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Y=1, dissimilar pair ( $\vec{X}_1, \vec{X}_2$ )

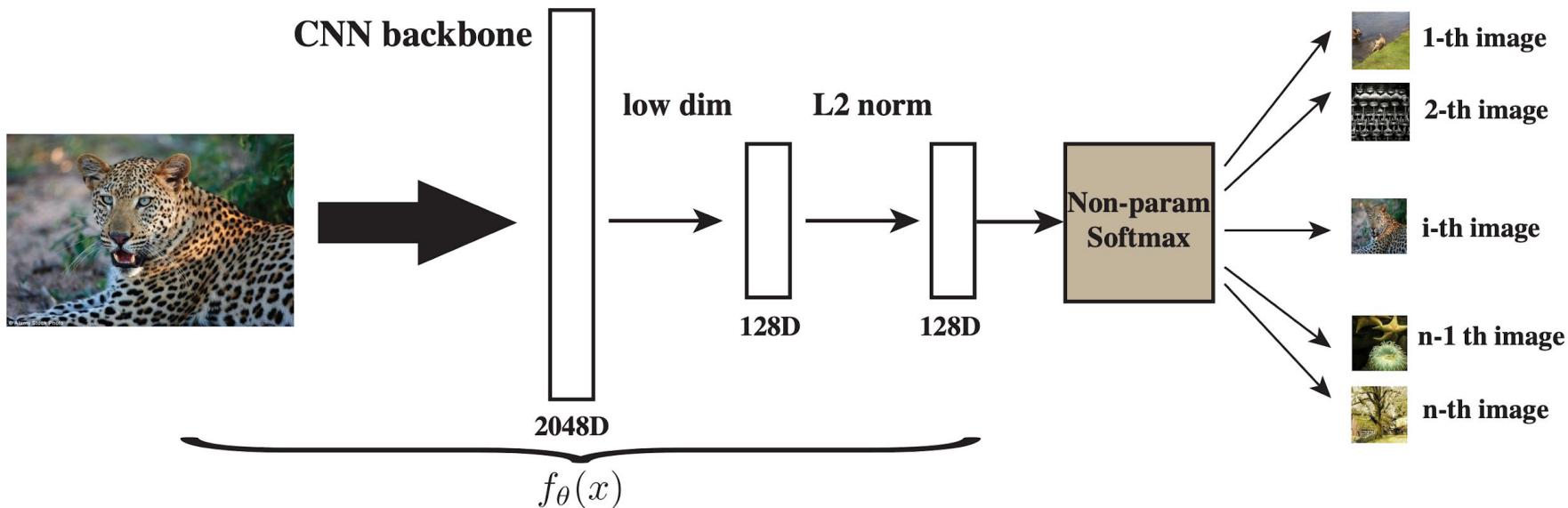
# Contrastive Loss

Learn a NN  $G_w(X)$  with pairs  $(X_1, X_2)$  which are similar ( $Y=0$ ) or dissimilar ( $Y=1$ ).



# Momentum Contrast (MoCo)

**Pretext task:** Treat each image instance as a distinct class of its own and train a classifier to distinguish whether two crops (views) belong to the same image.

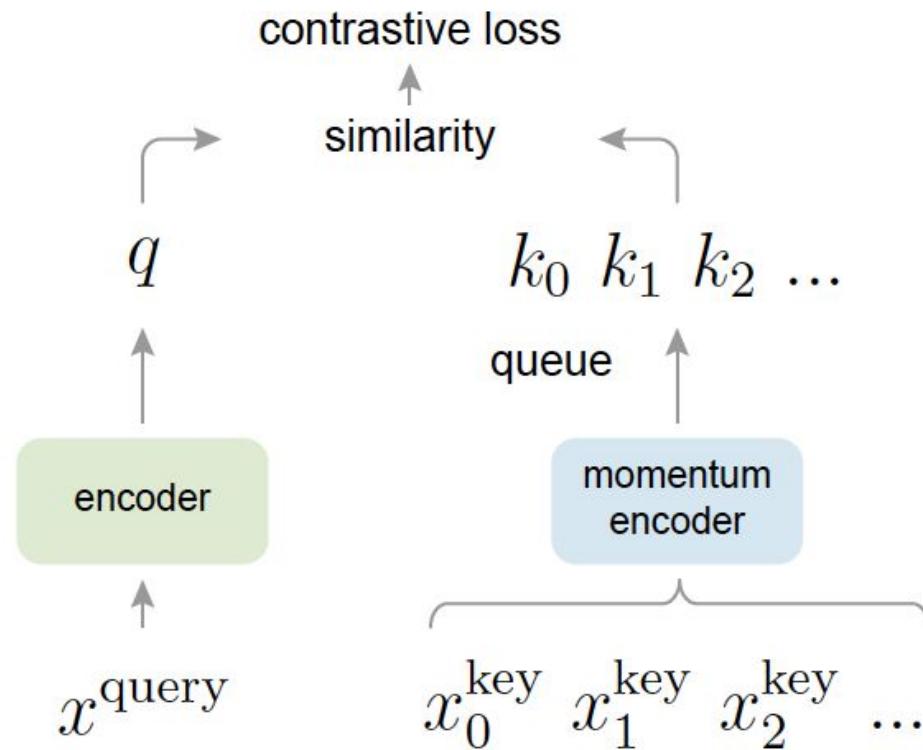


# Momentum Contrast (MoCo)

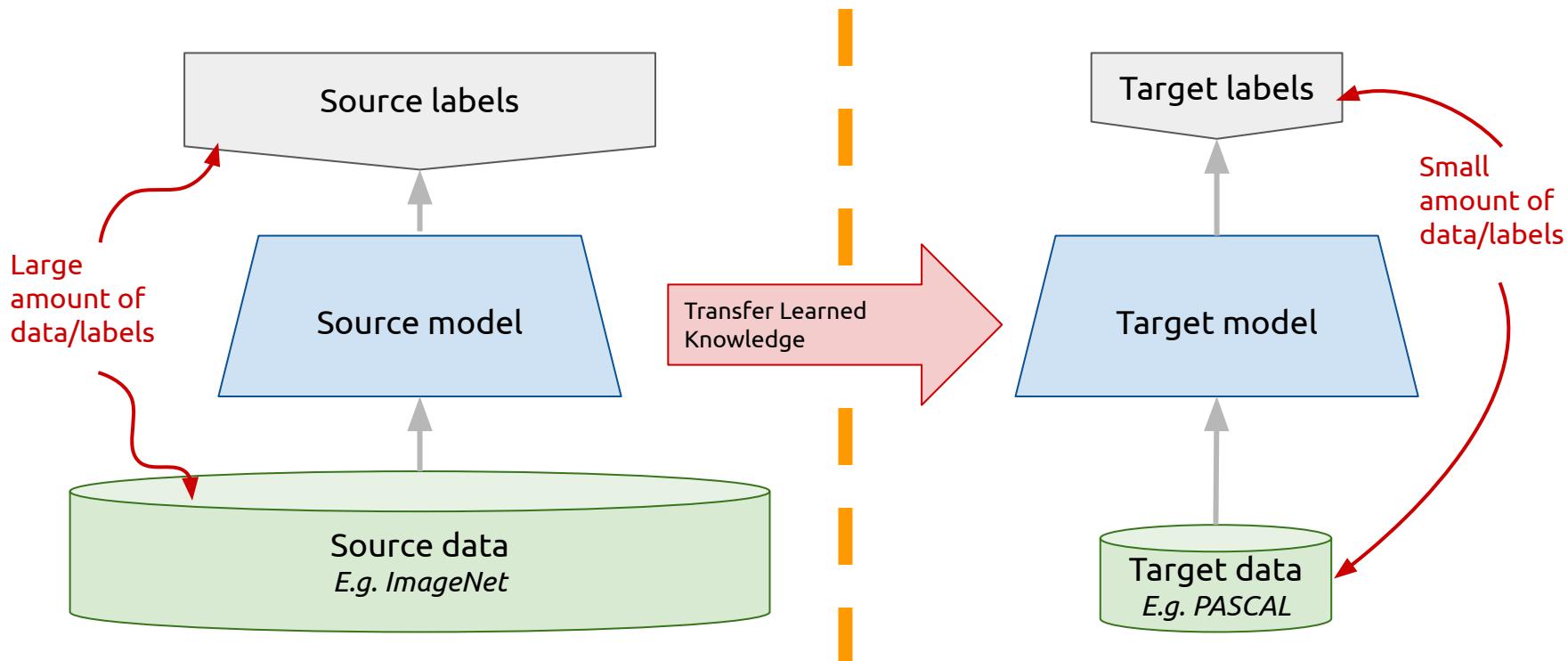
Learn a contrastive loss by building training batches between:

- Encoded features of a query image.
- A running queue of encoded keys.

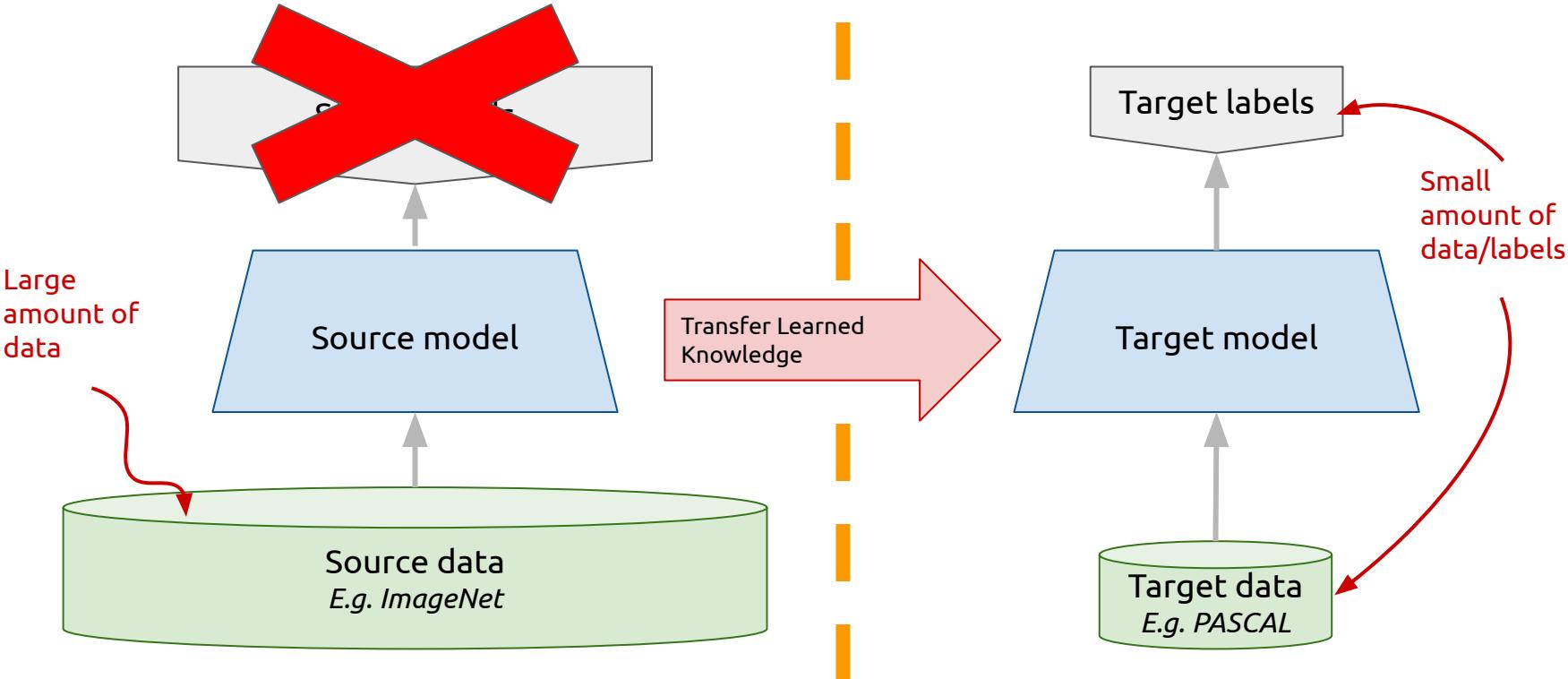
**Losses over representations**  
(instead of pixels) facilitate abstraction.



# Transfer Learning from ImageNet labels (super. IM-1M)



# Transfer Learning without ImageNet labels (MoCo IN-1M)



# Fine-tune task (super. IN-1M vs MoCo IN-1M)

**Transferring to Pascal VOC object detection with Faster R-CNN model:**  
Unsupervised MoCo training outperformed supervised ImageNet pre-training.

pre-train	AP <sub>50</sub>	AP	AP <sub>75</sub>
random init.	60.2	33.8	33.1
super. IN-1M	81.3	53.5	58.8
<b>MoCo IN-1M</b>	<b>81.5 (+0.2)</b>	<b>55.9 (+2.4)</b>	<b>62.6 (+3.8)</b>

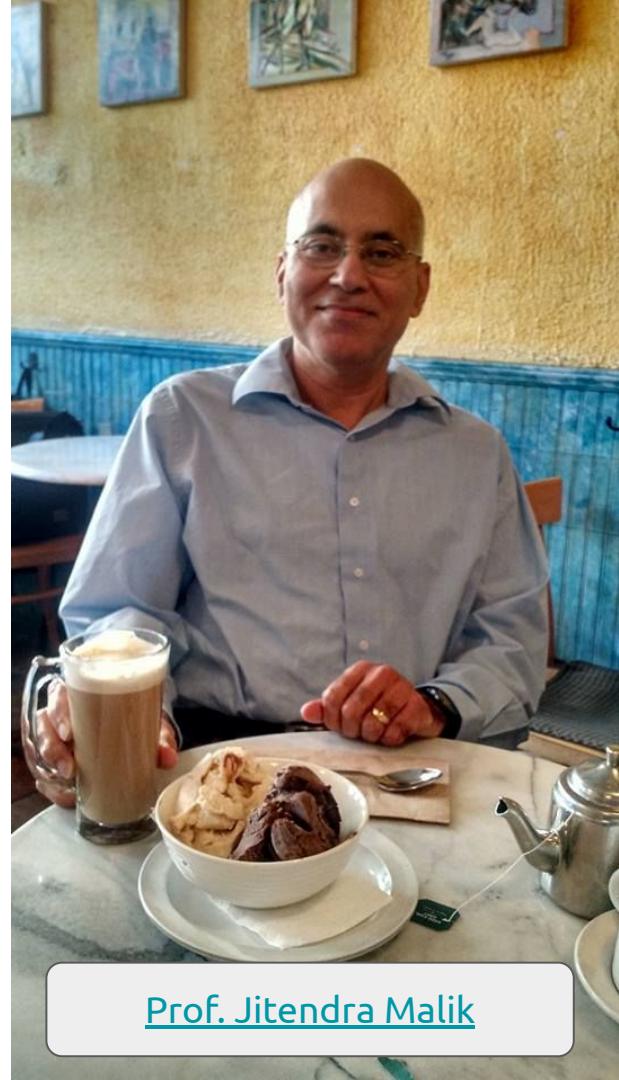
## The Gelato Bet (Sept 23, 2014)

*"If, by the first day of autumn (Sept 23) of 2015, a method will exist that can match or beat the performance of R-CNN on Pascal VOC detection, without the use of any extra, human annotations (e.g. ImageNet) as pre-training, Mr. Malik promises to buy Mr. Efros one (1) gelato (2 scoops: one chocolate, one vanilla)."*

Who won The Gelato Bet ?



Source: Ankesh Anand, "Contrastive Self-Supervised Learning" (2020)"



Prof. Jitendra Malik

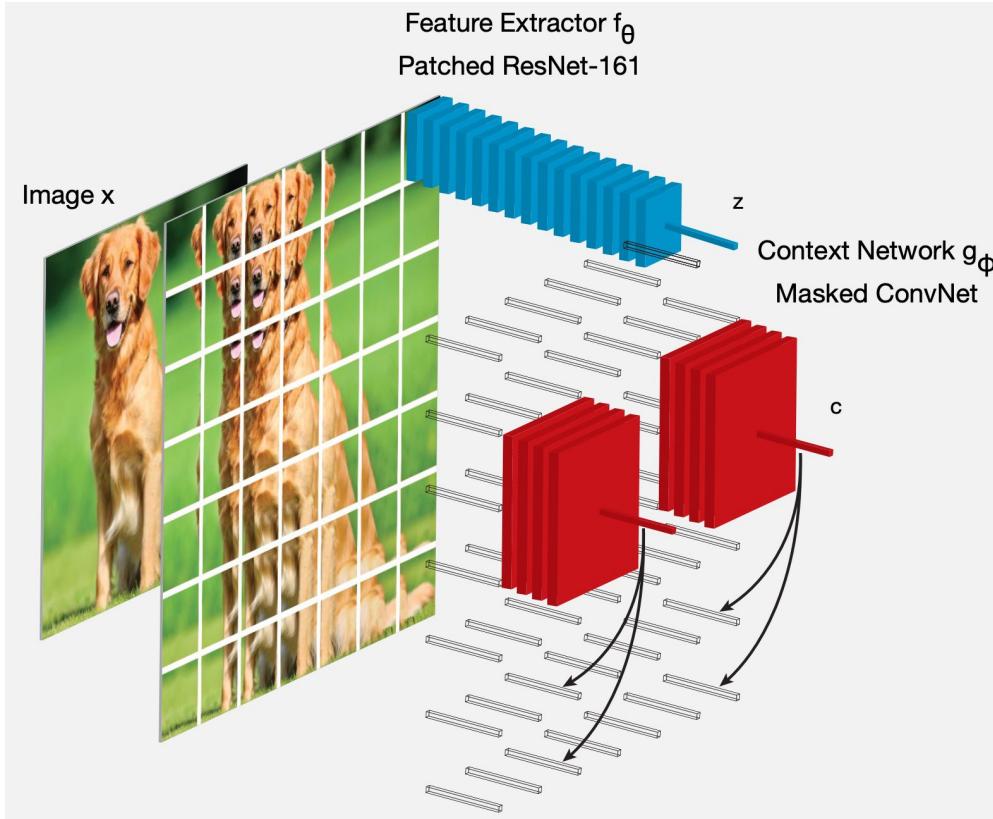
# The Revolution will not be Supervised



Prof. Alexei A. Efros



# Spatial Prediction Task



- 1) Image is divided in patches.
- 2) Each patch is encoded independently (in blue).
- 3) Features in the upper half of the images are aggregated with a context network (in red).
- 4) Rows of context vectors must predict the “blue” vectors below.

# Correlated views by data augmentation



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(i) Gaussian blur

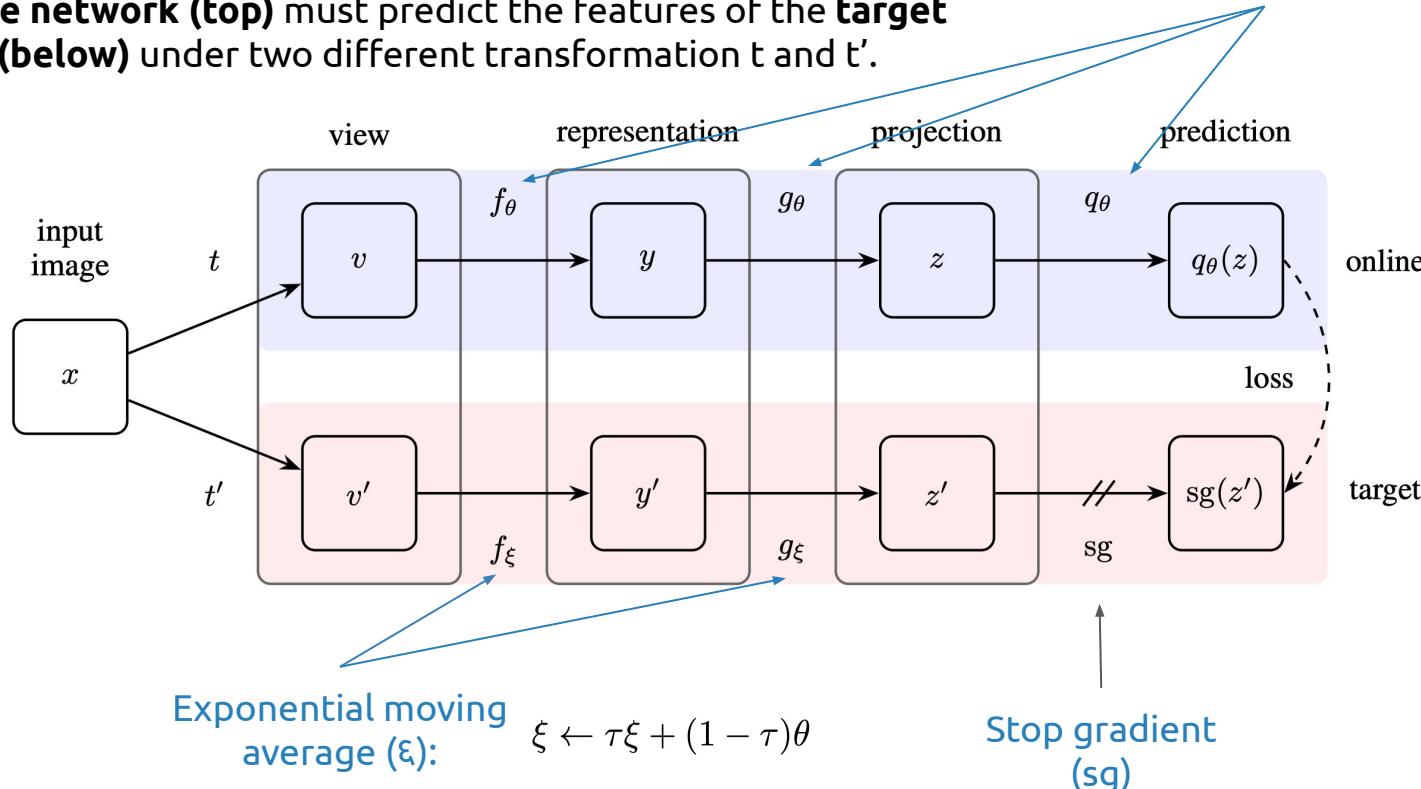
Original image cc-by: Von.grzanka

#SimCLR Chen, Ting, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. "[A simple framework for contrastive learning of visual representations.](#)" ICML 2020. [\[tweet\]](#)

# Predict transformed features

Network parameters  
being trained ( $\theta$ )

The **online network (top)** must predict the features of the **target network (below)** under two different transformation  $t$  and  $t'$ .

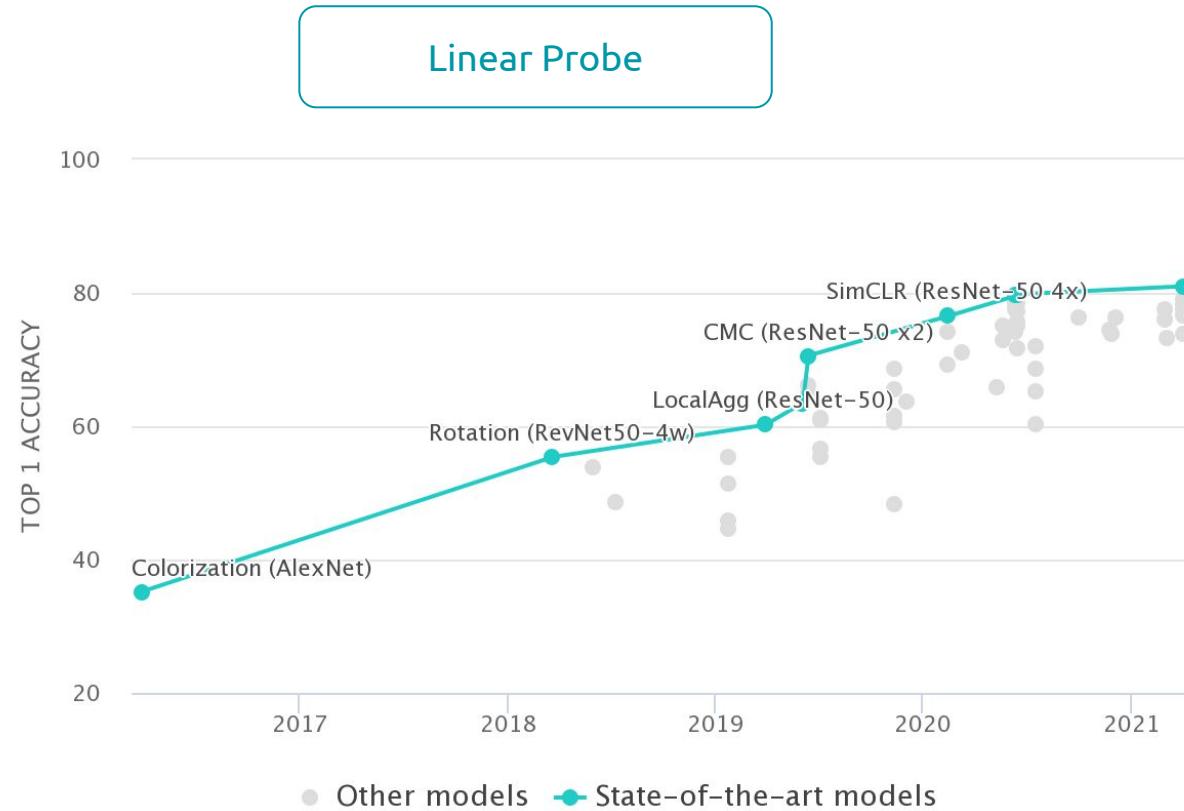


# Cluster features

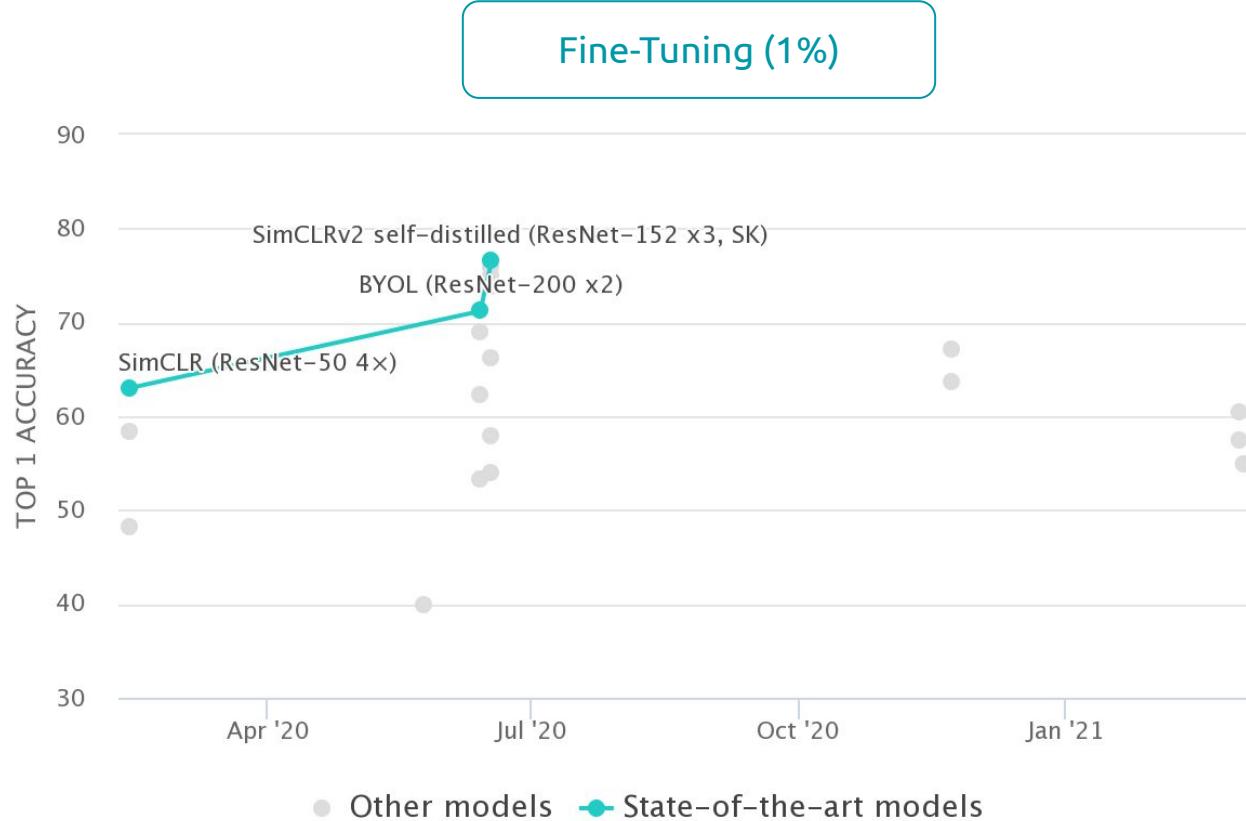
Patches from different views (transformations) of a source image are independently assigned to clusters of images. The cluster assignments must be consistent over time.

**Pretext task:** predicting the cluster of one version of the image with the other version:

# State of the Art

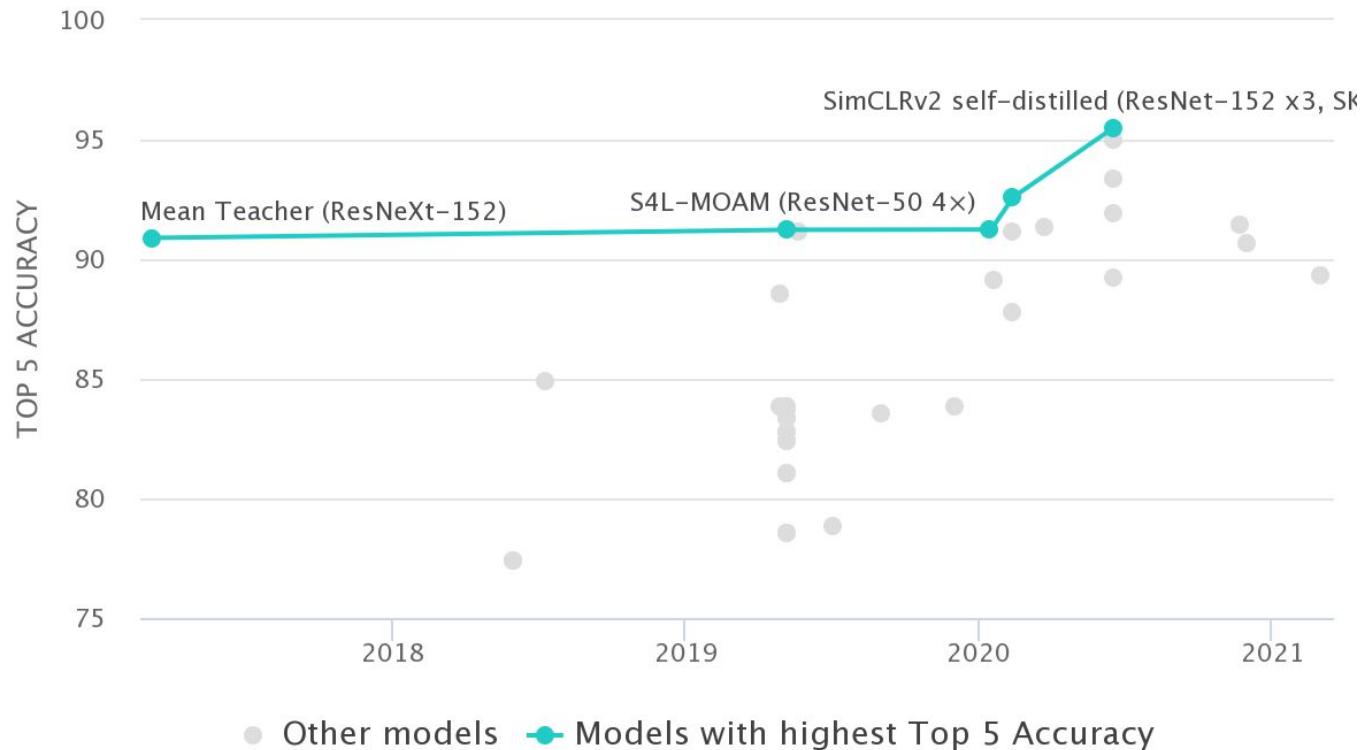


# State of the Art



# State of the Art

Fine-Tuning (10%)



# Triplet Ranking Loss (reminder)

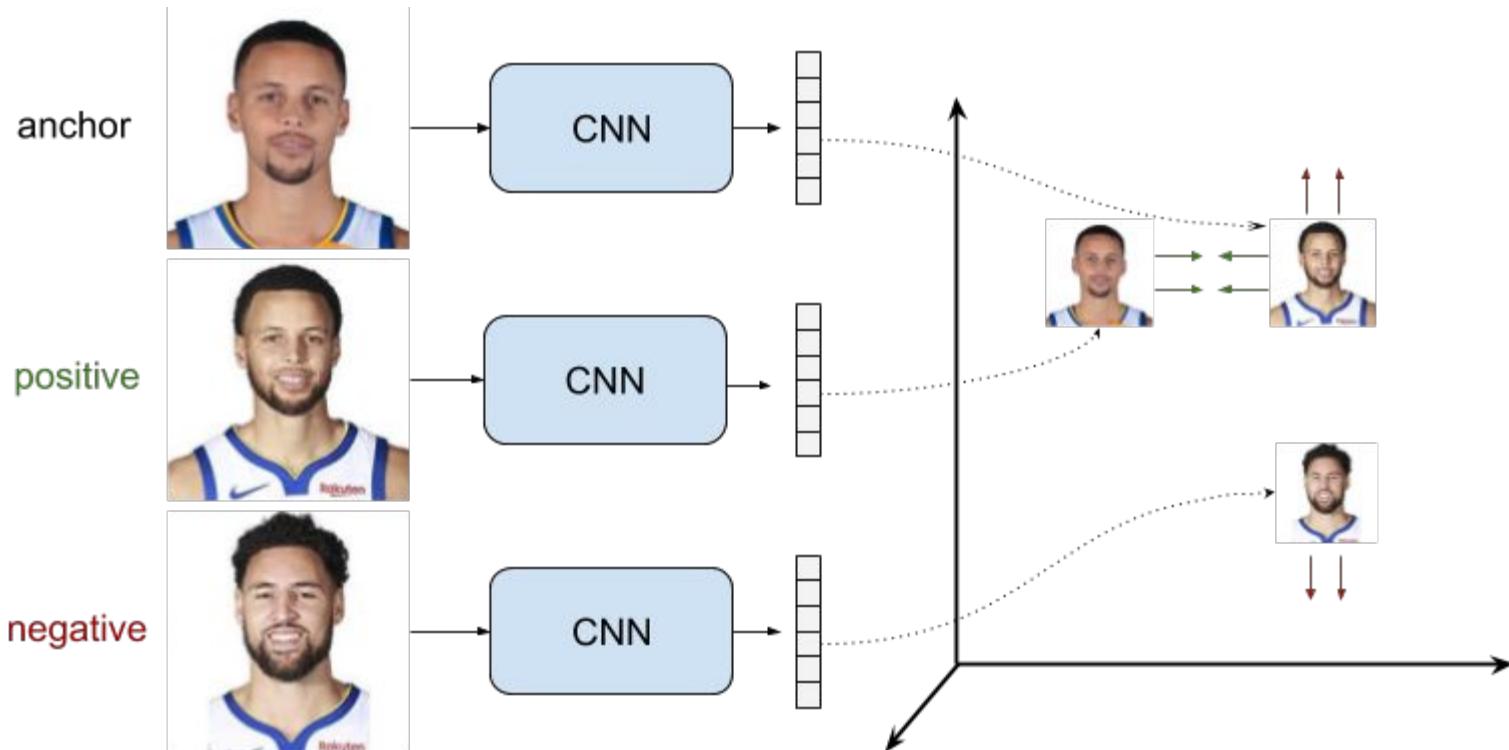


Figure: Raul Gómez, "[Understanding Ranking Loss, Contrastive Loss, Margin Loss, Triplet Loss, Hinge Loss and all those confusing names](#)" (2019)

# Triplet Ranking Loss = Contrastive Loss (too)

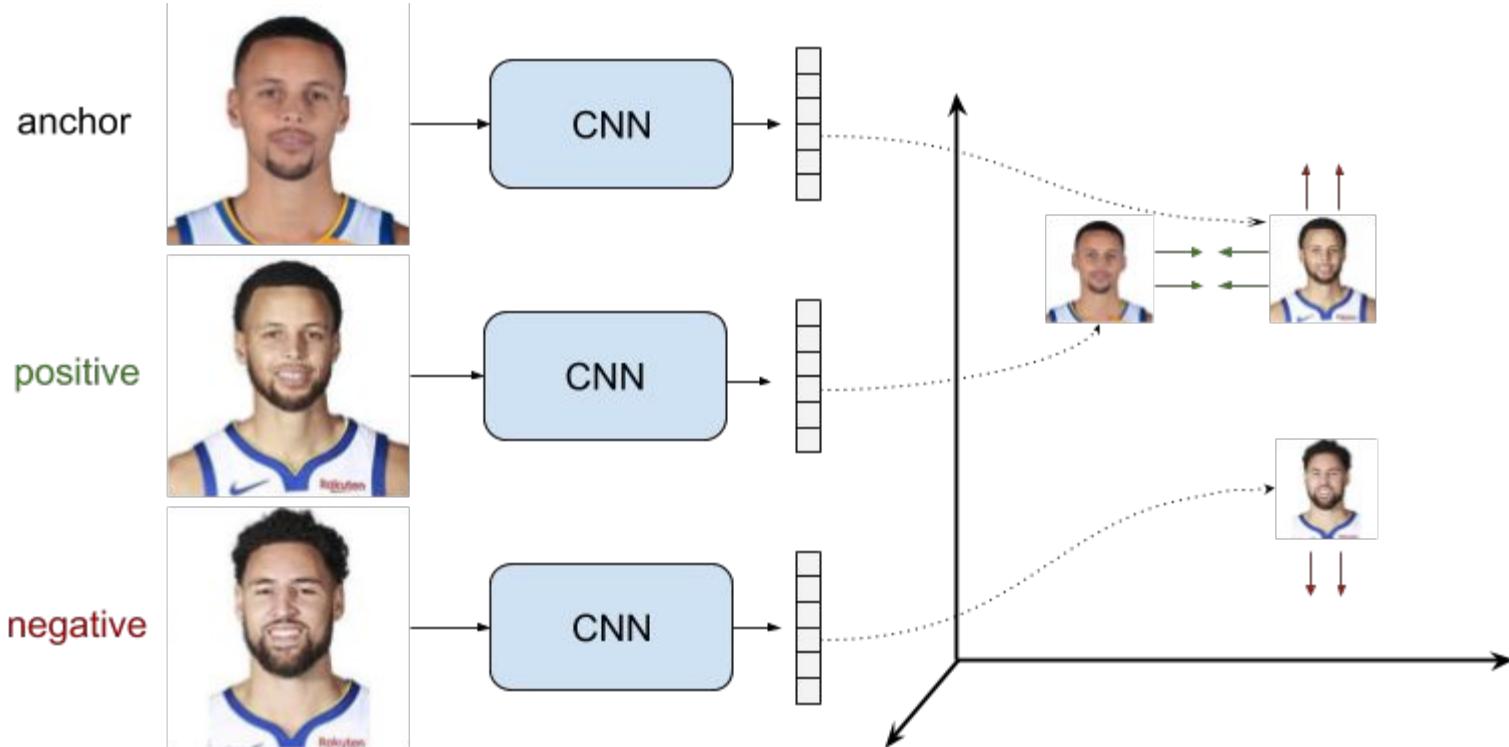
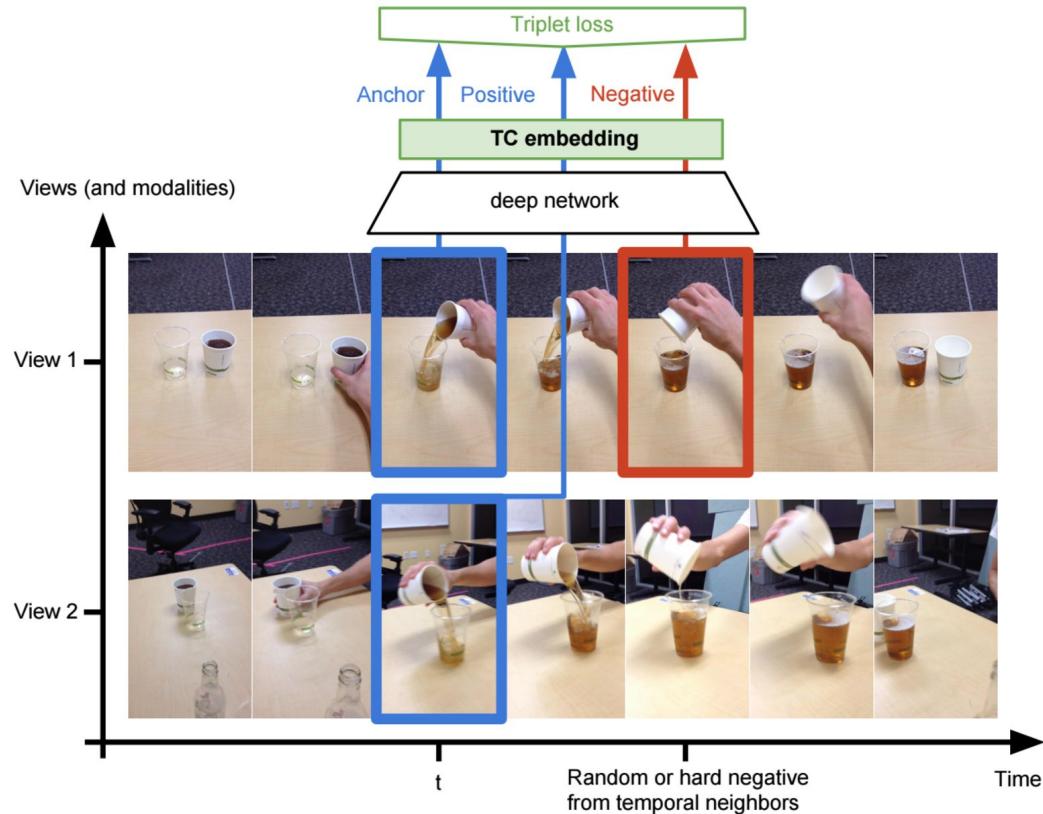


Figure: Raul Gómez, "[Understanding Ranking Loss, Contrastive Loss, Margin Loss, Triplet Loss, Hinge Loss and all those confusing names](#)" (2019)

# Multi-camera + Temporal ordering

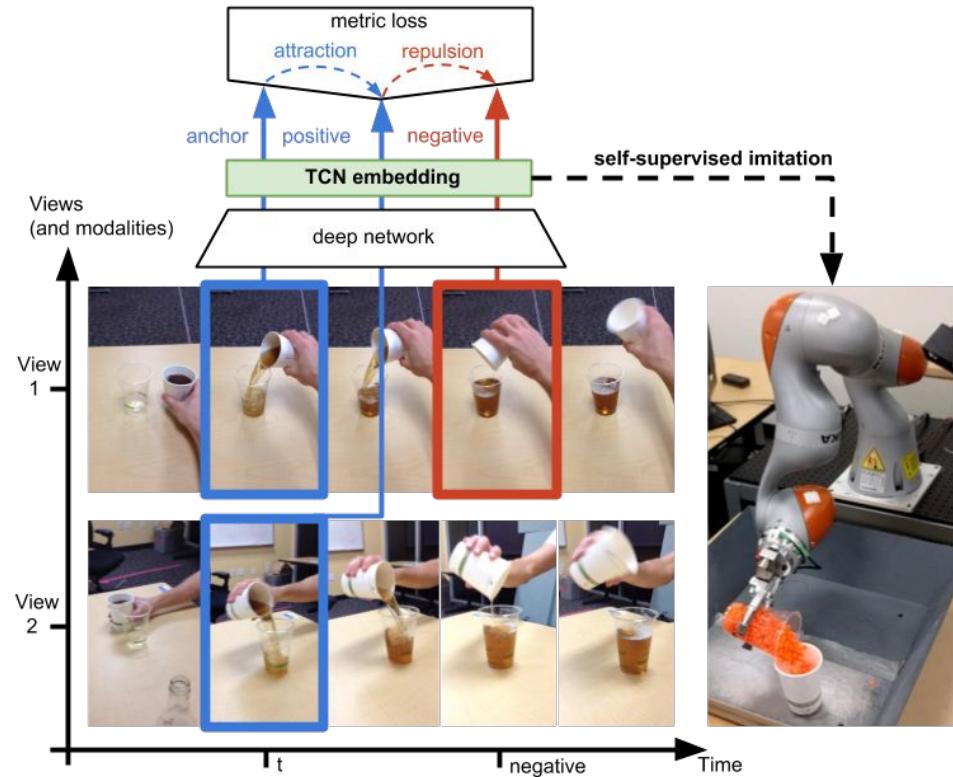
Pre-Training: Learn features with a triplet loss by:

- **Positive**: Same time-stamp from different cameras, to achieve viewpoint invariance.
- **Negative**: Other frames from the same camera, to be sensitive to temporal ordering



# Multi-camera + Temporal ordering

Inference: Imitation learning for a robotic arm.



#TCN Sermanet, Pierre, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. ["Time-contrastive networks: Self-supervised learning from video."](#) ICRA 2018.

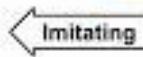
## Learning to imitate, from video, without supervision



3rd-person observation

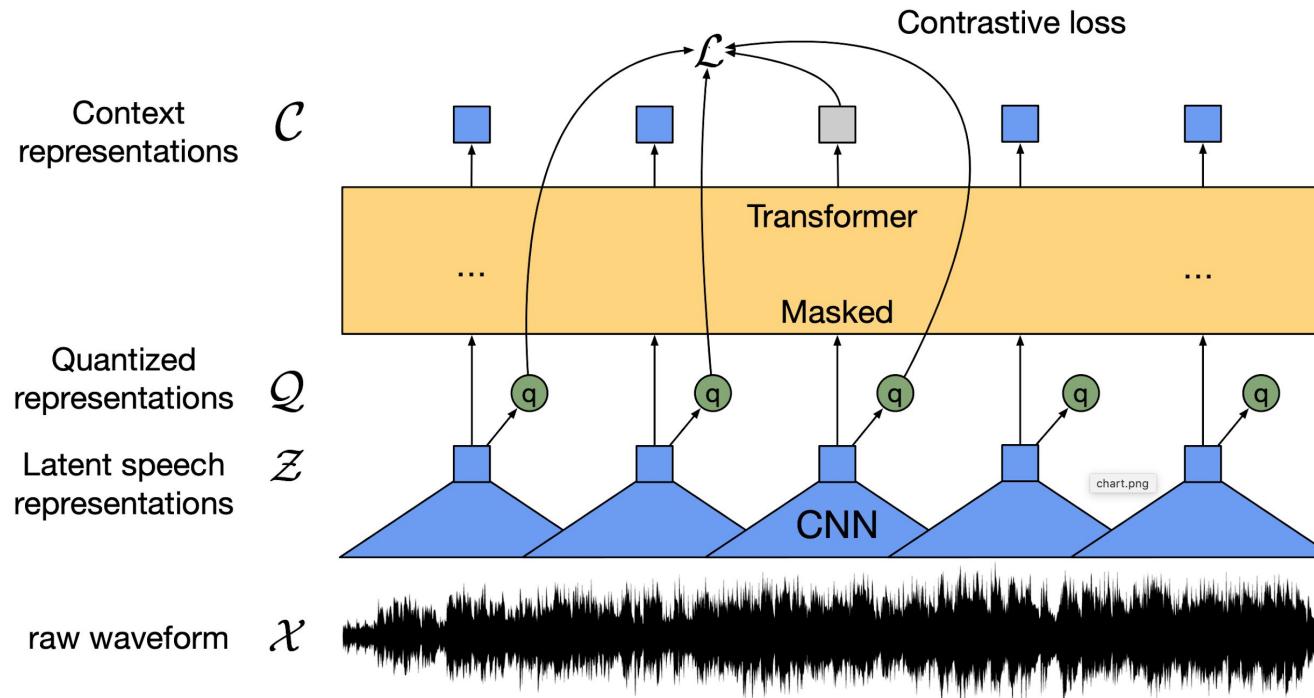


Learned policy



# Masking + Contrastive loss for speech

Predict whether a sample belong to the future of the same sequence.

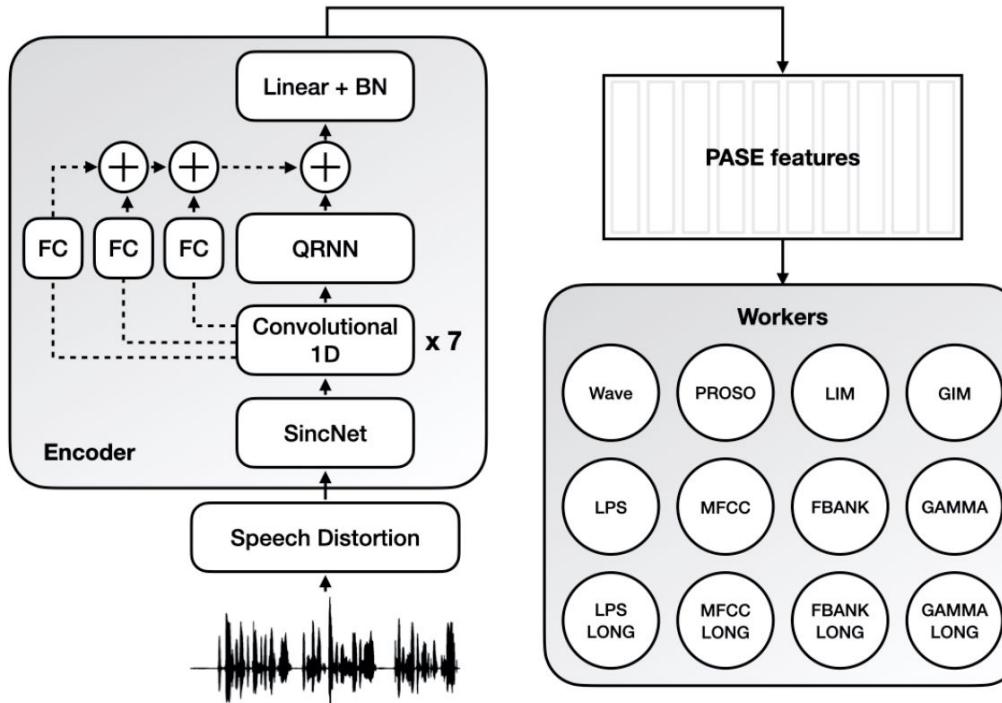


#Wav2Vec Baevski, A., Zhou, H., Mohamed, A., & Auli, M. (2020). [wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations](#). arXiv preprint arXiv:2006.11477. [\[tweet\]](#) [\[code\]](#)

# Outline

1. Transfer Learning
2. Representation Learning
3. Unsupervised Learning
4. Self-supervised Learning
5. Predictive methods
6. Contrastive methods
7. **Self-supervised Learning @ UPC**

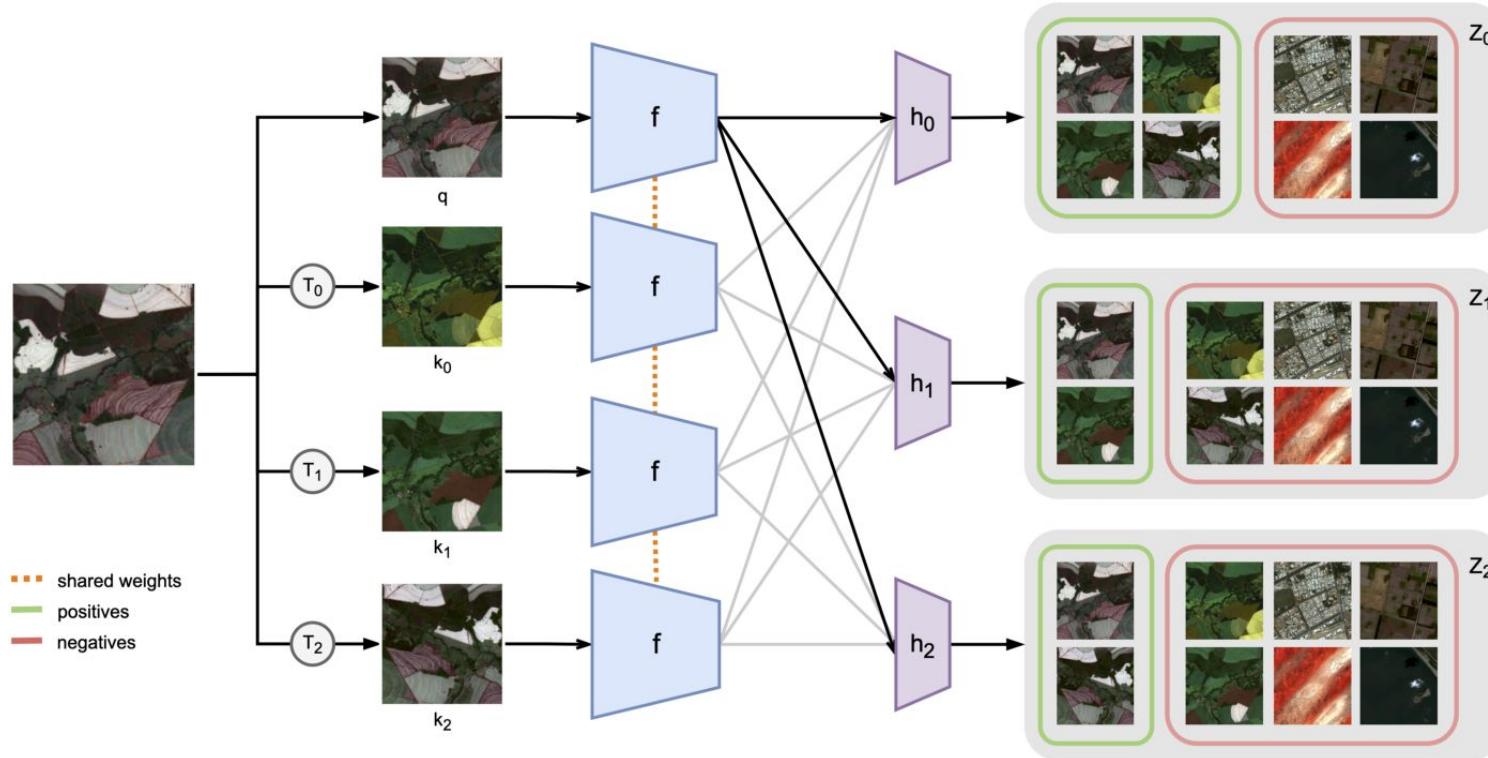
# Predictive Method: Handcrafted Workers



#PASE Pascual, S., Ravanelli, M., Serrà, J., Bonafonte, A., & Bengio, Y. (2019). [Learning problem-agnostic speech representations from multiple self-supervised tasks](#). INTERSPEECH 2019. [\[code\]](#)  
Figure: Gerard I. Gallego, ["End-to-end Speech Translation with Self-supervised Speech Representations"](#). UPC TelecomBCN 2020.



# Contrastive Methods: Seasonal Contrast



# Bonus Track

Lecture 16

## Self-Supervised Audio-Visual Learning



Xavier Giro-i-Nieto  
Associate Professor  
Universitat Politècnica de Catalunya  
[@DocXavi](#)  
[xavier.giro@upc.edu](mailto:xavier.giro@upc.edu)

[course site]

Xavier Giró, "Self-Supervised  
Audio-Visual Learning".

UPC TelecomBCN DLAI, 2020.

# Homework for December 3



Víctor Campos, “**Towards RL that scales: Autonomous acquisition and transfer of knowledge**”.

UPC TelecomBCN DRL, 2020.



"Inefficient Data Efficient".  
CVPR Deepvision 2020.



"The InfoNCE loss in self-supervised learning" [[slides](#)] [[tweet](#)]  
NeurIPS SSL Workshop 2020.

# Outline

1. Transfer Learning
2. Representation Learning
3. Unsupervised Learning
4. Self-supervised Learning
5. Predictive methods
6. Contrastive methods
7. Self-supervised Learning @ UPC

# Suggested talks



Demis Hassabis (Deepmind)



Aloysha Efros (UC Berkeley)

**Use audio and visual features**

facebook research

What can be learnt by watching and listening to videos?

- Good representations
  - Visual features
  - Audio features
- Intra- and cross-modal retrieval
  - Aligned audio and visual embeddings
- "What is making the sound?"
  - Learn to localize objects that sound

"Objects that Sound", Arandjelović and Zisserman, ICCV 2017 & ECCV 2018

The diagram illustrates a dual-path network architecture. It starts with two inputs: a 'single frame' (image) and a '1 s' (audio). Both inputs pass through their respective 'subnetworks' (labeled 'visual subnetwork' and 'audio subnetwork'). The outputs of these subnetworks are then compared in a 'correspond' module, which outputs a 'yes/no' classification.

Andrew Zisserman (Oxford/Deepmind)



Yann LeCun (Facebook AI)

# Related lecture

## CS294-158 Deep Unsupervised Learning

### Lecture 7 Self-Supervised Learning



Pieter Abbeel, Xi (Peter) Chen, Jonathan Ho, Aravind Srinivas, Alex Li, Wilson Yan  
UC Berkeley

# Suggested reading



Alex Graves, Kelly Clancy, "[Unsupervised learning: The curious pupil](#)". Deepmind blog 2019.

# Related events & publications

## Events

- [Self-supervised Learning: What's Next](#). ECCV Workshop 2020.
- [Self-supervised Learning - Theory and Practice](#). NeurIPS Workshop 2020.

## Publications

- Le-Khac, Phuc H., Graham Healy, and Alan F. Smeaton. "[Contrastive representation learning: A framework and review.](#)" IEEE Access (2020).
- Purushwarkam, S., & Gupta, A. [Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases](#). NeurIPS 2020.
- YHH Tsai, Y Wu, R Salakhutdinov, LP Morency, "[Self-Supervised Learning from a Multiview Perspective](#)". ICLR 2021

# Transfer Learning: Video-lectures

**DEEP LEARNING FOR COMPUTER VISION**  
Master Course UPC TelecomBCN, Fall 2016

Instructors: [List of 6 names]

Organizers: [List of logos]

+ info: [TelecomBCN.DeepLearning.Barcelona](http://TelecomBCN.DeepLearning.Barcelona)

Day 2 Lecture 5

**Transfer learning and domain adaptation**

UPC UNIVERSITAT POLITÈCNICA DE CATALUNYA  
BANCO DE CATALUÑA  
Departament de Teoria del Senyal i Comunicacions

[Kevin McGuinness \(UPC DLCV 2016\)](#)

**DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE**  
Master Course UPC ETSETB TelecomBCN Barcelona, Autumn 2018

Instructors: [List of 6 names]

Organizers: [List of logos]

Supporters: Google Cloud, GitHub Education

+ info: <http://bit.ly/dlai2018>

Day 5 Lecture 2

**Transfer learning and domain adaptation**

#DLUPC

Many slides from:

Ramon Morros

Kevin McGuinness

Eric Arazo

[course site]

[Ramon Morros \(UPC DLAI 2018\)](#)

# Video lectures on Unsupervised Learning



**The manifold hypothesis**

The data distribution lie close to a low-dimensional manifold

Example: consider image data

- Very high dimensional (1,000,000)
- A randomly generated image will almost certainly not look like any real-world scene
  - The space of images that occur in nature is almost completely empty
- Hypothesis: real world images lie on a smooth, low-dimensional manifold
  - Mahalanobis distance is a good measure of similarity

Similar for audio and text



UPC  
UNIVERSITAT POLITÈCNICA DE CATALUNYA  
BANDELONTECH  
Departament de Teoria del Senyal i Comunicacions



**DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE**

Master Class UPC ETSIIETB Telecommunications Barcelona, Autumn 2017

Instructors



Organizers



Supporters



+ info: <http://daii.deeplearning.barcelona>

[course site]



Xavier Giro-i-Nieto  
[savent.giro@upc.edu](mailto:savent.giro@upc.edu)

Associate Professor  
Universitat Politècnica de Catalunya  
Technical University of Catalonia



Kevin McGuinness, [UPC DLCV 2016](#)

Xavier Giró, [UPC DLAI 2017](#)

# Questions ?

## Undergradese

What undergrads ask vs. what they're REALLY asking

