



Multimodal Generative AI 2025

Self-Supervised Learning



Vicky Kalogeiton

Lecture 3: CSC_52002_EP



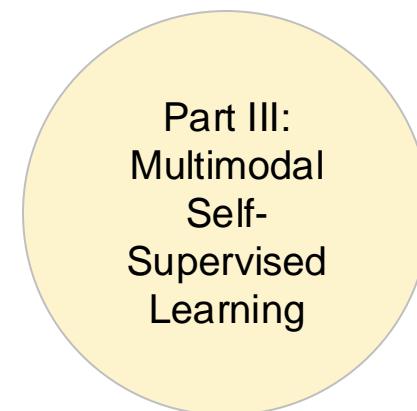
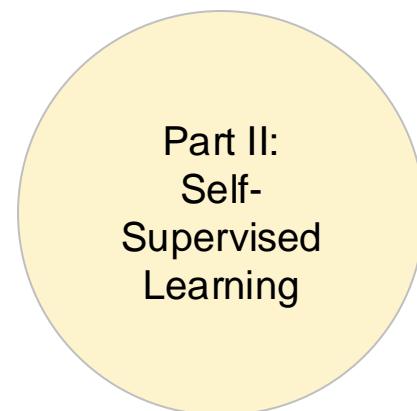
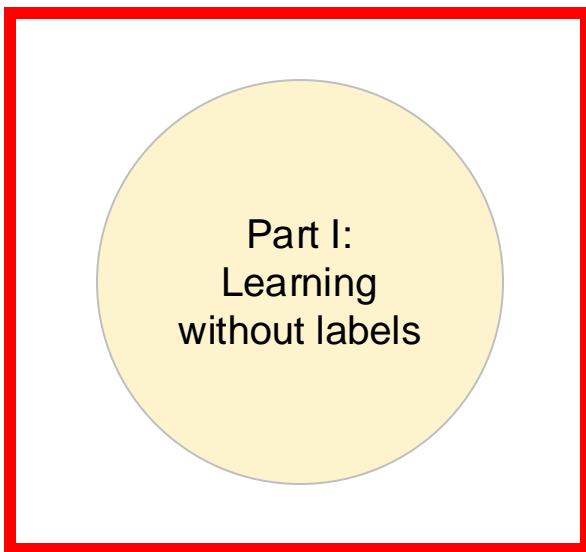
Today's lecture

Part I:
Learning
without labels

Part II:
Self-
Supervised
Learning

Part III:
Multimodal
Self-
Supervised
Learning

Today's lecture



Pop Quiz

- Types of learning that do not require labeled data?

Pop Quiz - Hint

- Types of learning that do not require labeled data?
 - Supervised Learning
 - Unsupervised Learning

Pop Quiz Answers

- Types of learning that do not require labeled data?
 - Supervised Learning
 - Semi-supervised Learning
 - Weakly-supervised Learning
 - Self-supervised Learning
 - Unsupervised Learning
- Active Learning

Part I: Outline

Learning without labels

- Learning from data:
 - Supervised Learning
 - Unsupervised Learning
 - Weakly-supervised Learning
 - Semi-supervised Learning
- Active Learning

Supervised Learning

Supervised Learning

- **Data:** (x, y) , where x is data, y is label
- **Goal:** Learn a function to map $x \rightarrow y$

Examples:

Classification,
regression,
object detection,
segmentation,
captioning...

Input: Image



Output: Assign Image to one of a fixed set of categories



Cat

Dog

Deer

Bird

Car

Unsupervised Learning

Supervised Learning

- **Data:** (x, y) , where x is data, y is label
- **Goal:** Learn a function to map $x \rightarrow y$

Examples:
Classification,
regression,
object detection,
segmentation,
captioning...

Unsupervised Learning

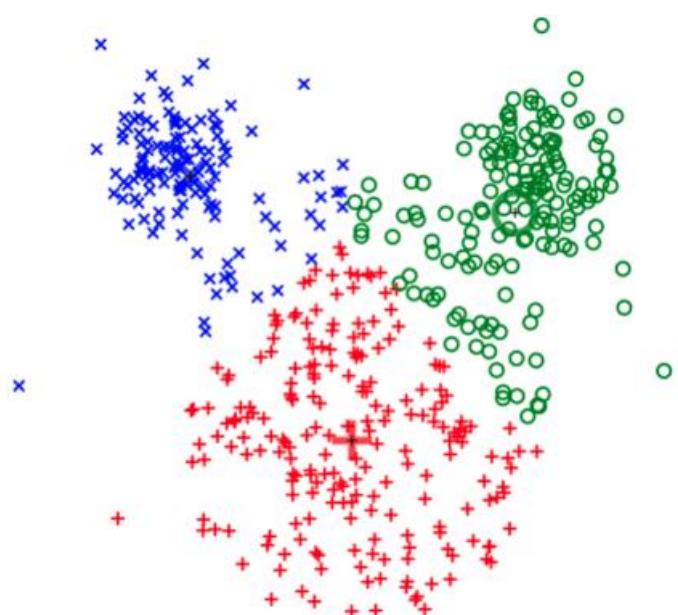
- **Data:** x , where x is data, **NO label!**
- **Goal:** Learn some underlying hidden structure of the data
- **Advantage:** No labelling

Examples:
Clustering
Dim reduction
Feature learning
Density estimation...

Unsupervised Learning

- **Data:** x , where x is data, **NO** label!
- **Goal:** Learn some underlying hidden structure of the data

Examples:
Clustering
Dim reduction
Feature learning
Density estimation...

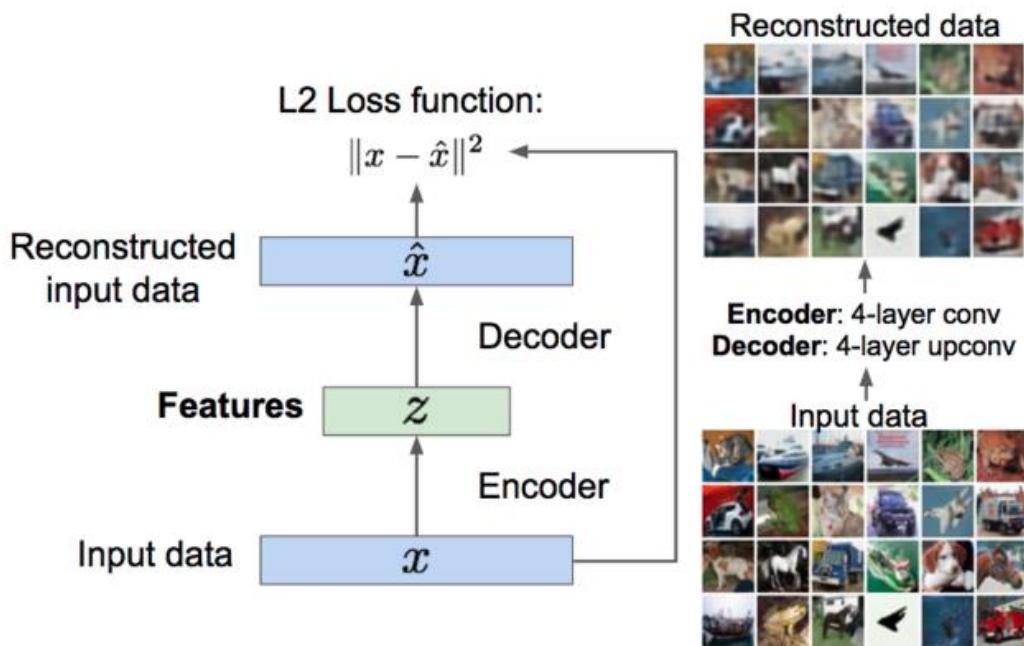


K-Means Clustering

Unsupervised Learning

- **Data:** x , where x is data, **NO** label!
- **Goal:** Learn some underlying hidden structure of the data

Examples:
 Clustering
 Dim reduction
 Feature learning
 Density estimation...



Feature Learning
 e.g. autoencoders

Weakly-Supervised Learning

Manual annotations are expensive



80-200s



40s



3s



(free)

Su et al., Crowdsourcing annotations for visual object detection, AAAI 2012

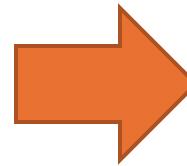
Deng et al., Scalable multi-label annotation, CHI 2014

Lin et al., Microsoft COCO: common objects in context, ECCV 2014

Fully supervised setting

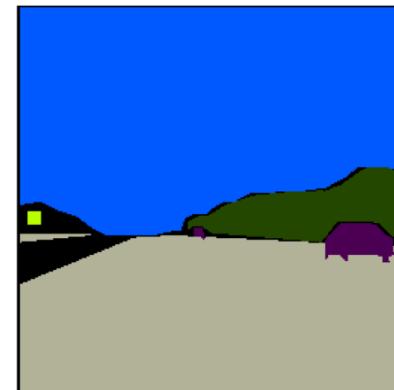
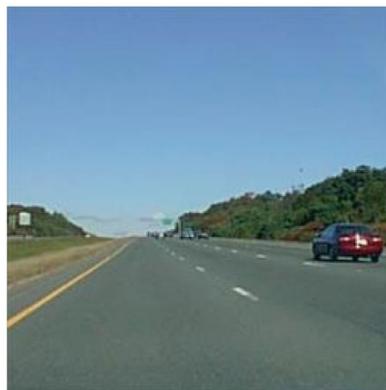


...



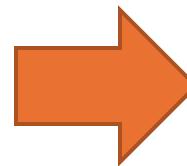
Object
detection
model

motorbike



- Car
- Road
- Sign
- Sky
- Tree

...



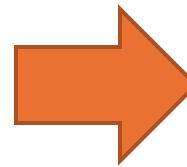
Semantic
segmentation
model

*Annotation to the **same** degree as outputs on test images*

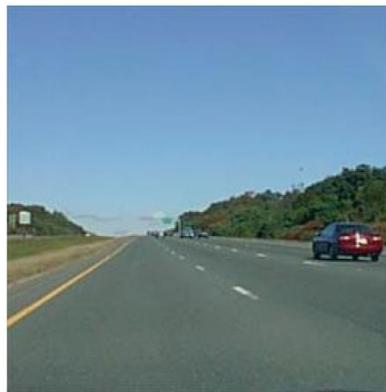
Weakly-supervised learning



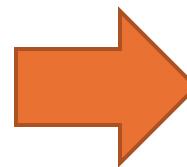
...



Object
detection
model



...

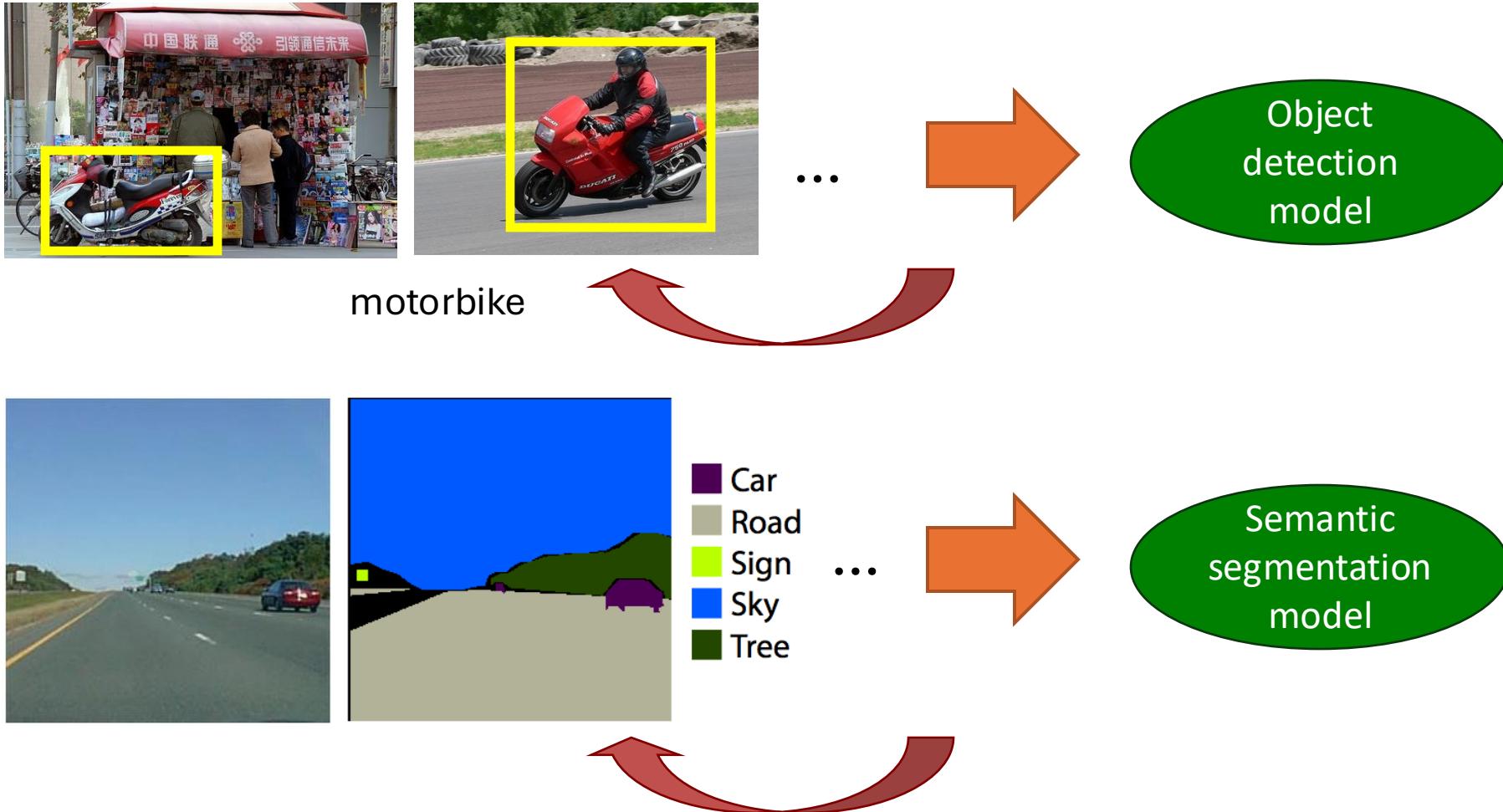


Semantic
segmentation
model

car, road, sign, sky, tree

*Annotation to a **lower** degree than outputs on test images*

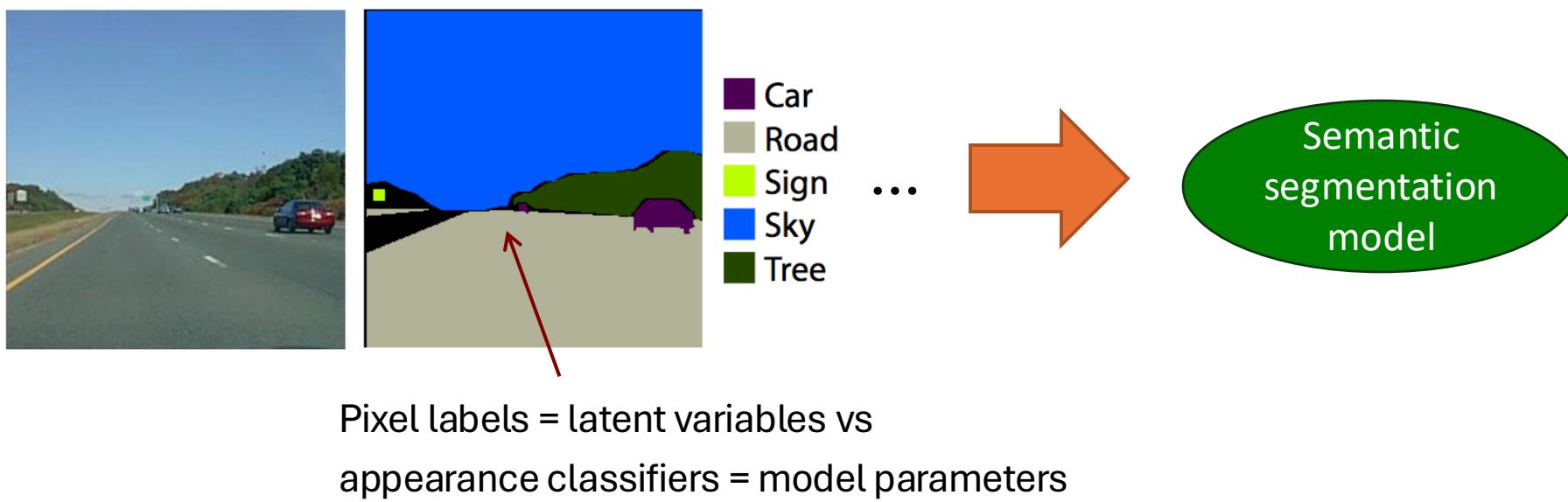
Weakly-supervised learning



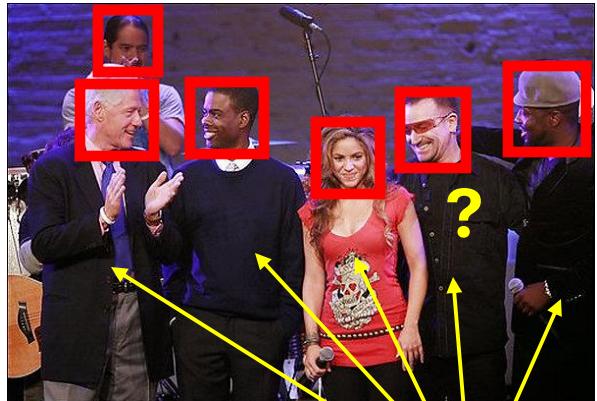
Core issue: need for auto-annotation

Weakly-supervised learning

- *Train with lower degree of supervision than necessary to fix latent (internal) variables*
- *Need to inject additional assumptions (bias) in the learner, e.g. appearance classes are compact, objects form recurring appearance patterns*



Learn face models from news captions



Former U.S. President Bill Clinton, comedian **Chris Rock**, musicians Shakira, Bono, and Wyclef jean participate ...

Correspondence ambiguity

Source: Yahoo-news



McChrystal and Gates raise tension in Afghanistan debate

Weakly supervised learning here means:

- Recover latent name-to-face assignments
- Learn face appearance models specific to person sub-classes

Semi-Supervised Learning

Supervised Learning

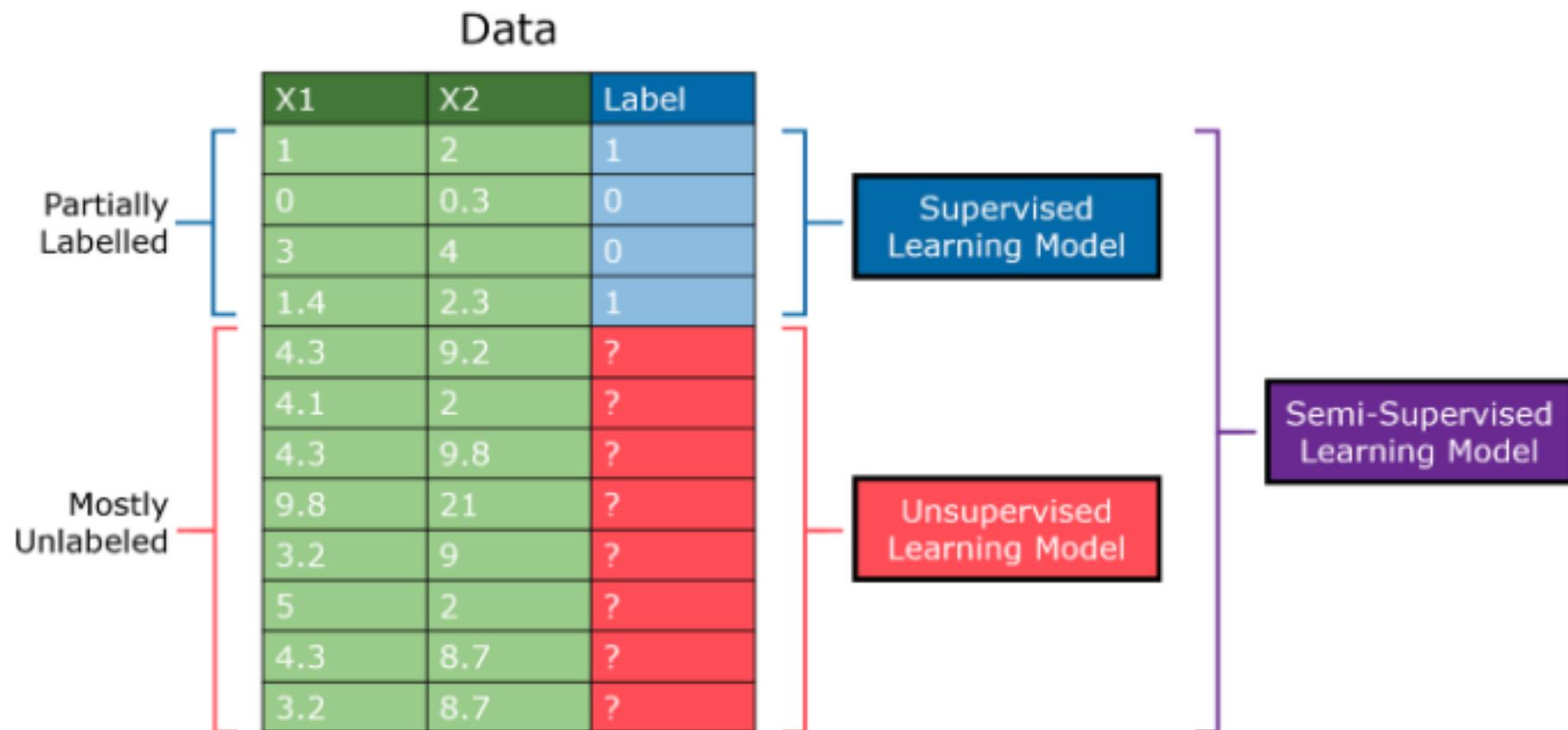
- **Data:** (x, y) , where x is data, y is label
- **Goal:** Learn a function to map $x \rightarrow y$

Examples:
Classification,
regression,
object detection,
segmentation,
captioning...

Semi-supervised Learning

- **Data:** labeled (x_1, y_1) , where x is data, y is the label and (x_2) data with **NO** label!
 - $x_2 \gg x_1 \rightarrow$ much more unlabeled data
- **Goal:** Learn a better prediction rule than based on labeled data alone
- **Advantage:** Partial labelling

Semi-supervised Learning



Semi-supervised Learning

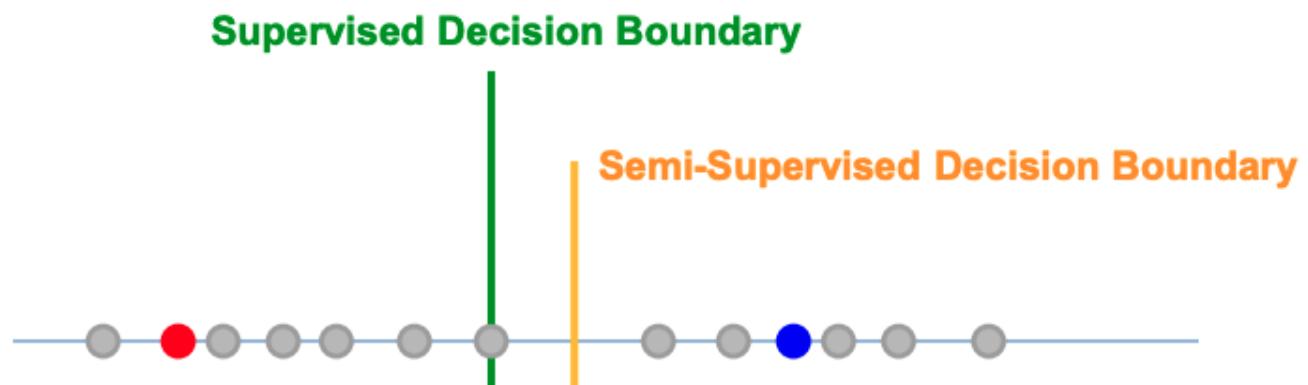
- Limited set of labelled sample data to train itself
→ partially-trained model
- Use partially-trained model to label the unlabeled data
→ pseudo-labeled data
- Combine labeled and pseudo-labeled data

Semi-supervised Learning:
Supervised + Unsupervised:

1. Classification to identify data assets
2. Clustering to group it into distinct parts

Can unlabeled data help?

- Positive labeled data
- Negative labeled data
- Unlabeled data



Assume each class is a coherent group (e.g. Gaussian)

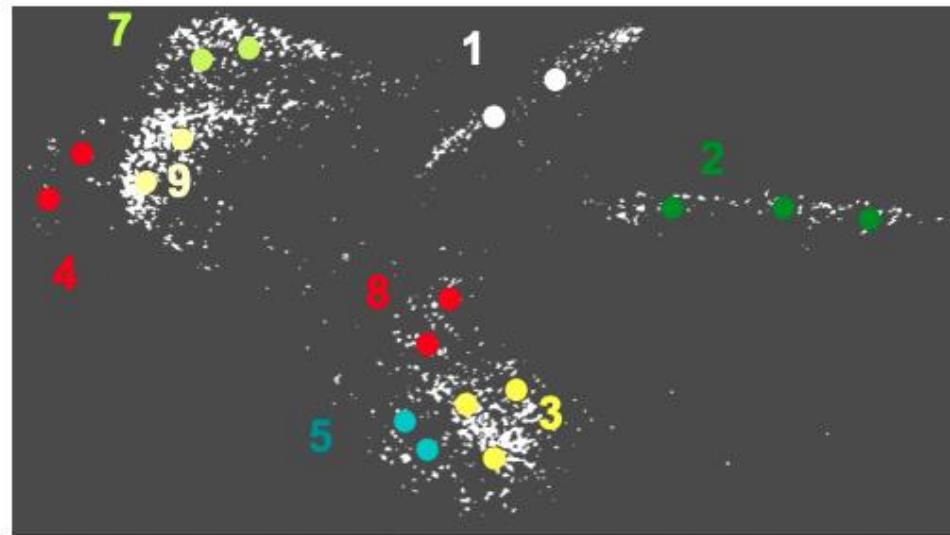
Then unlabeled data can help identify the boundary more accurately.

Can unlabeled data help?

Unlabeled Images

0	1	2	3	4	5	6	7	8	9
8	9	0	1	1	3	4	5	6	7
6	7	8	9	0	1	2	3	4	5

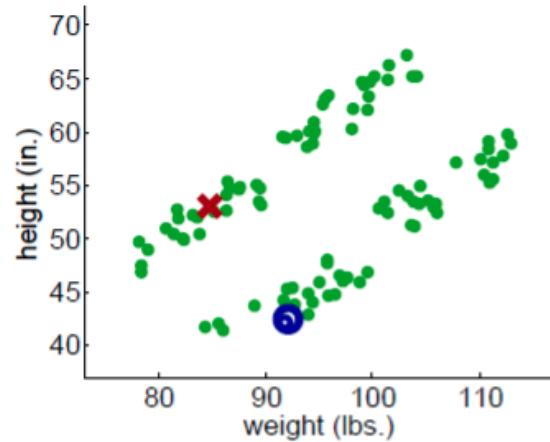
Labels "0" "1" "2" ...



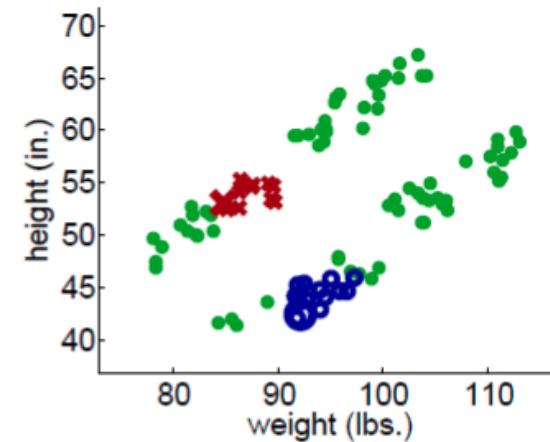
This embedding can be done by manifold learning algorithms

“Similar” data points have “similar” labels

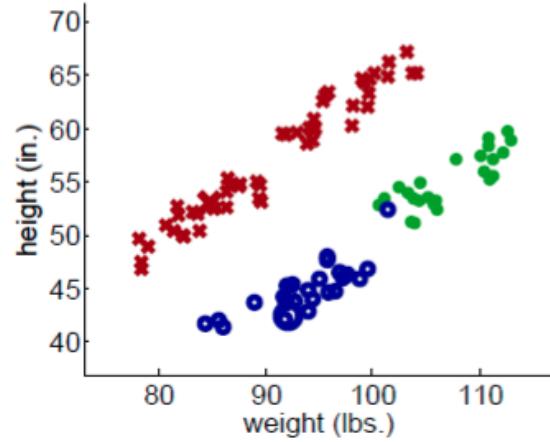
Propagating 1NN



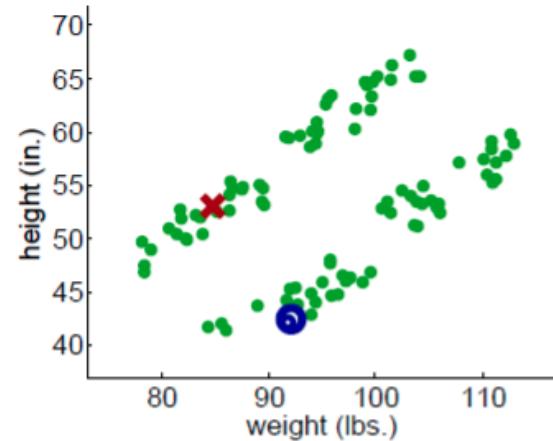
(a) Iteration 1



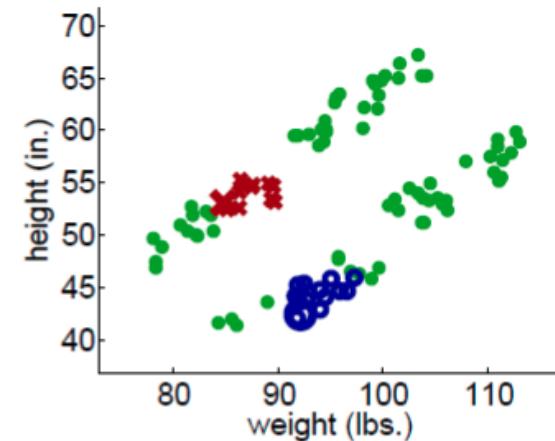
(b) Iteration 25



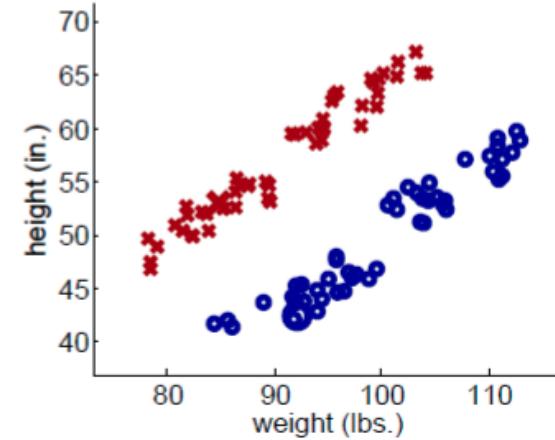
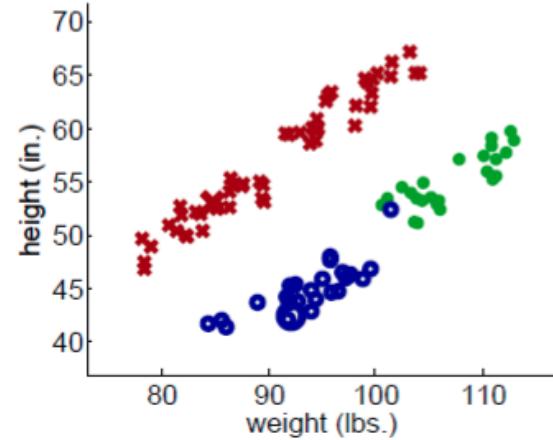
Pop Quiz



(a) Iteration 1

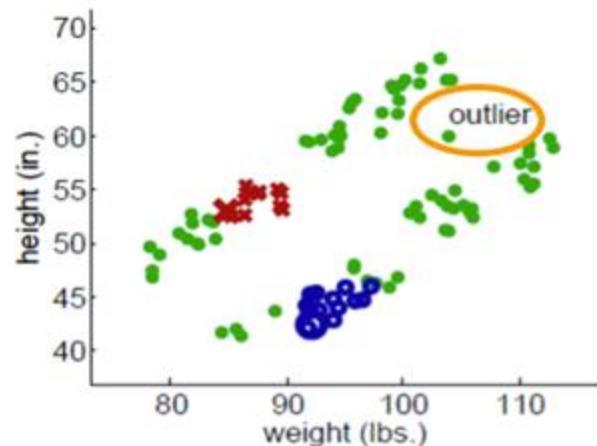


(b) Iteration 25

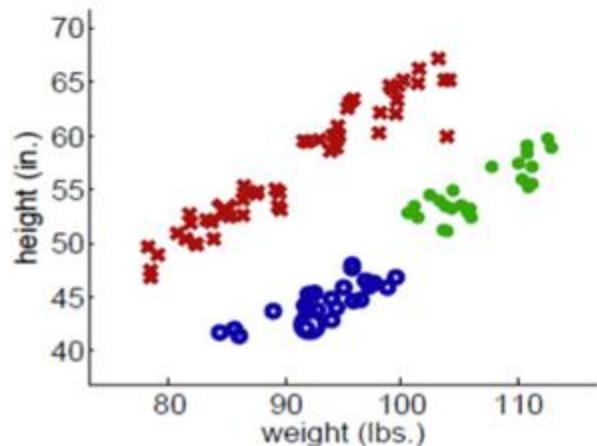


When the above won't work?

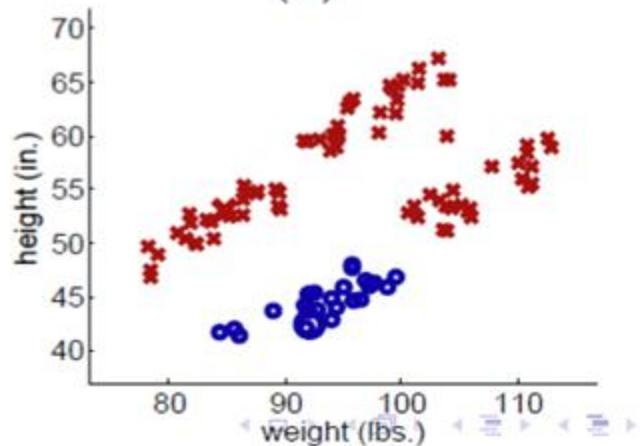
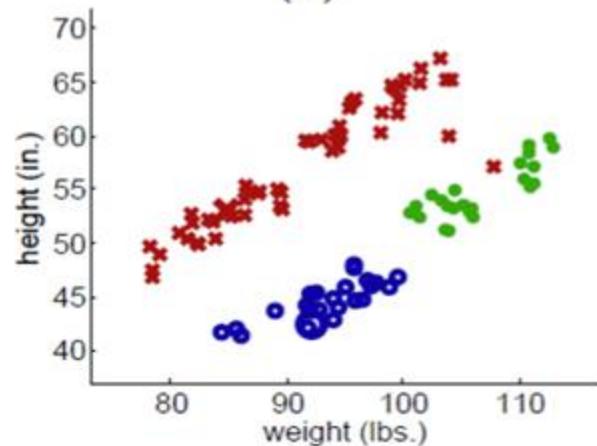
Propagating 1NN



(a)

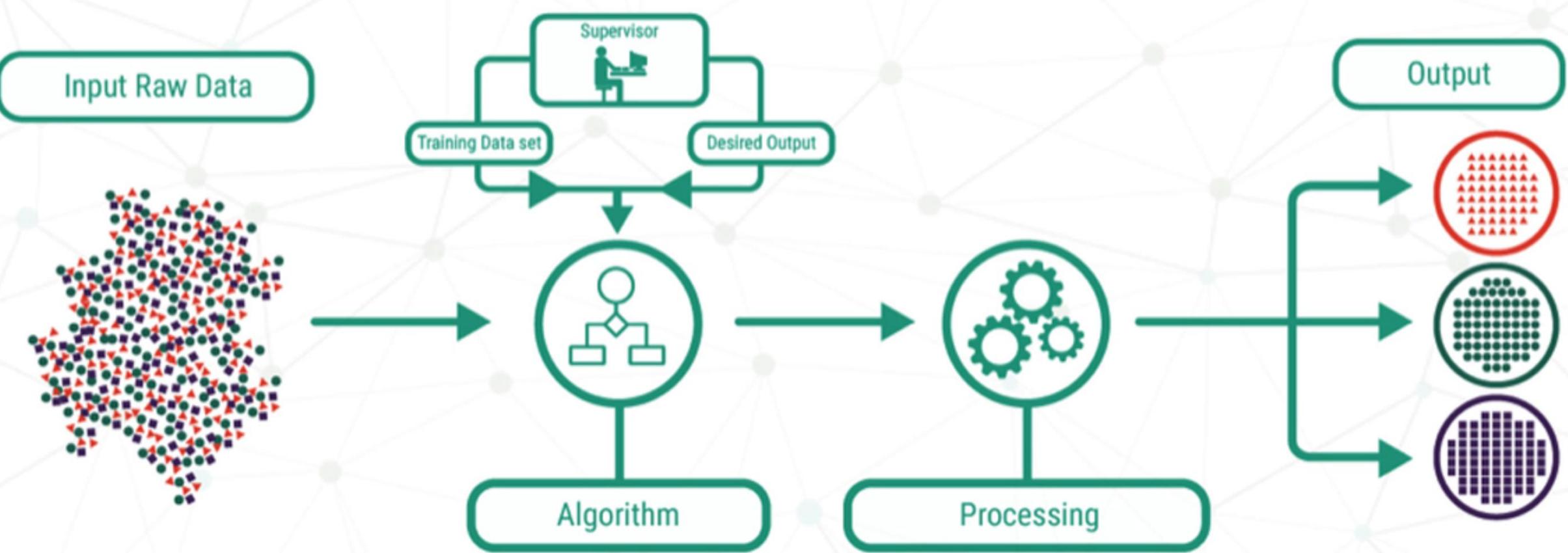


(b)

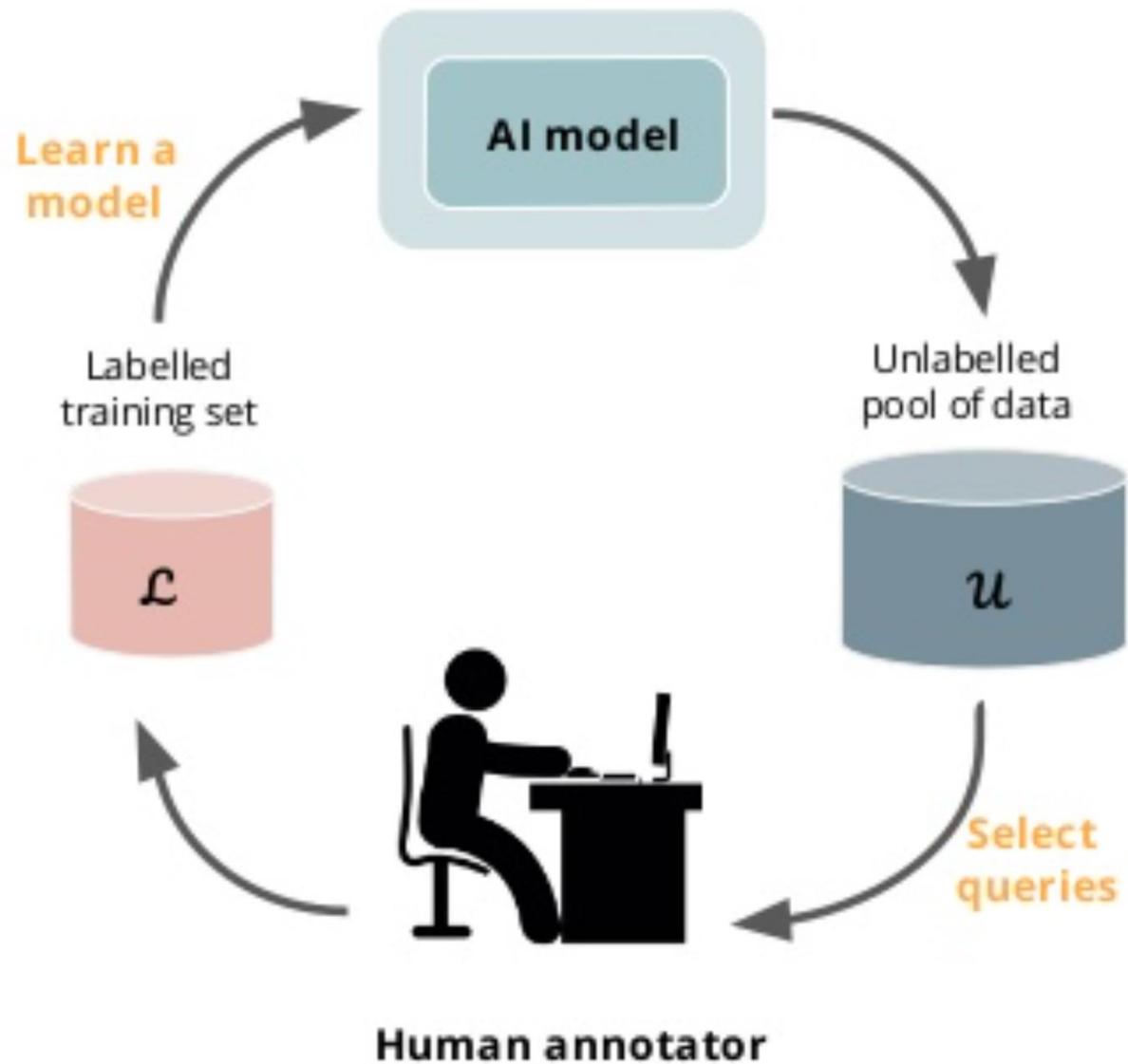


Active learning with human in the loop

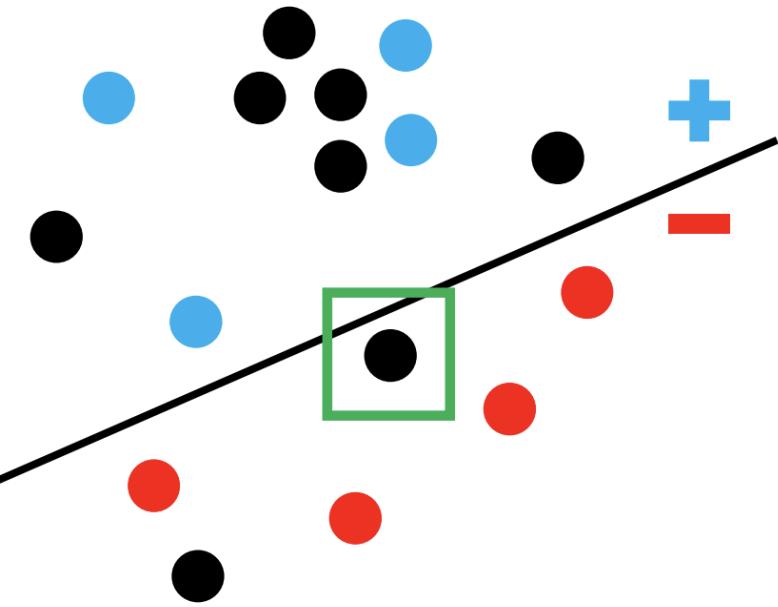
Fully supervised paradigm



Active learning



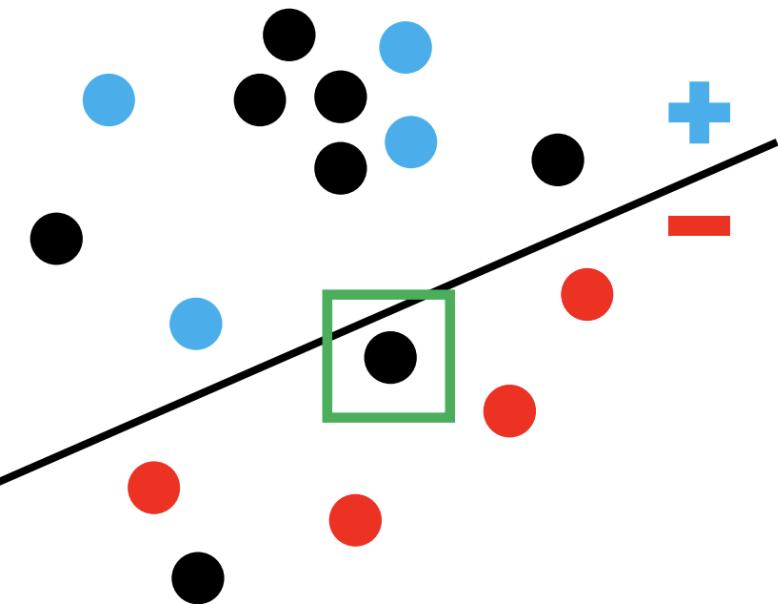
Active learning



(a) Uncertainty sampling
e.g. [Tong and Koller, 2002;
Kapoor et al., 2010]

[Tong 2002, Kapoor 2010, Jain 2009, Kovashka 11, Vijayanarasimhan 2009, Branson 2010]

Active learning



(a) Uncertainty sampling
e.g. [Tong and Koller, 2002;
Kapoor et al., 2010]

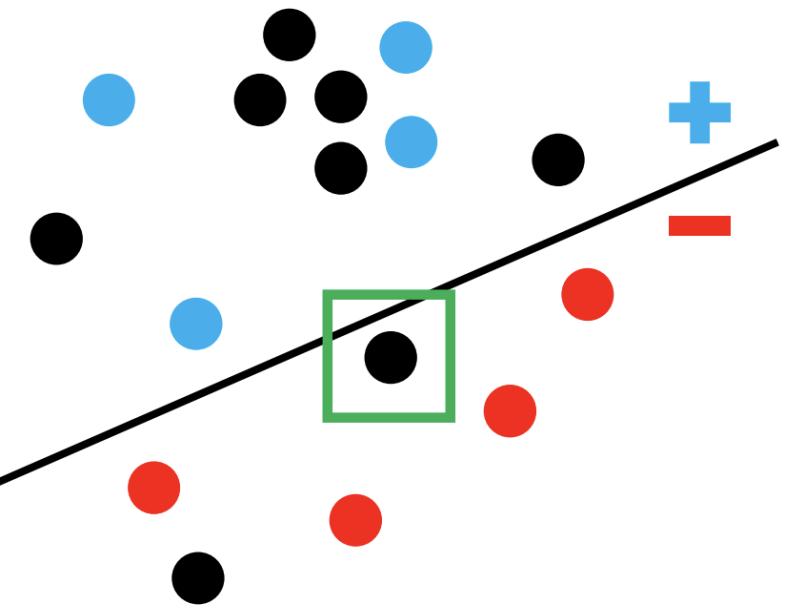
[Tong 2002, Kapoor 2010, Jain 2009, Kovashka 11, Vijayanarasimhan 2009, Branson 2010]

- Very efficient
- Labeling does not guarantee that will improve uncertainty on all images
- Annotated set is specific to one classifier

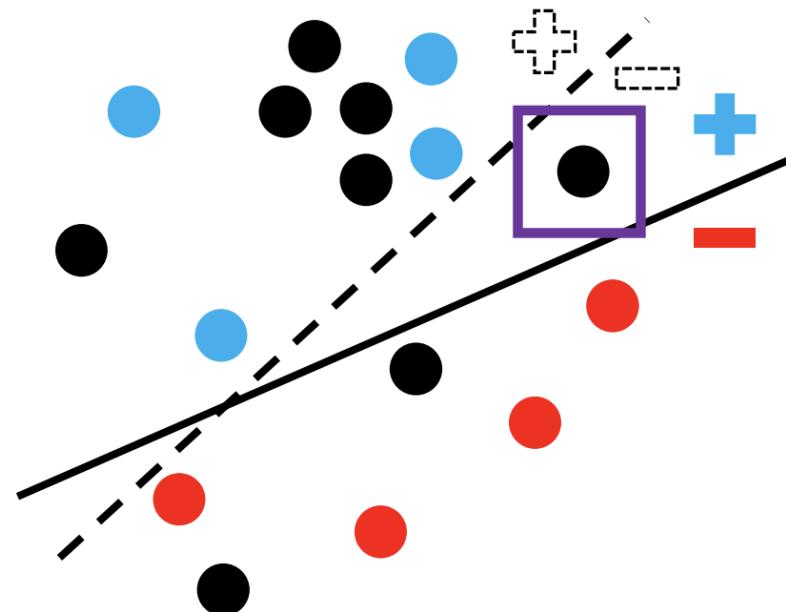
Alternatives

- Misclassification risk
- Measure expected entropy

Active learning



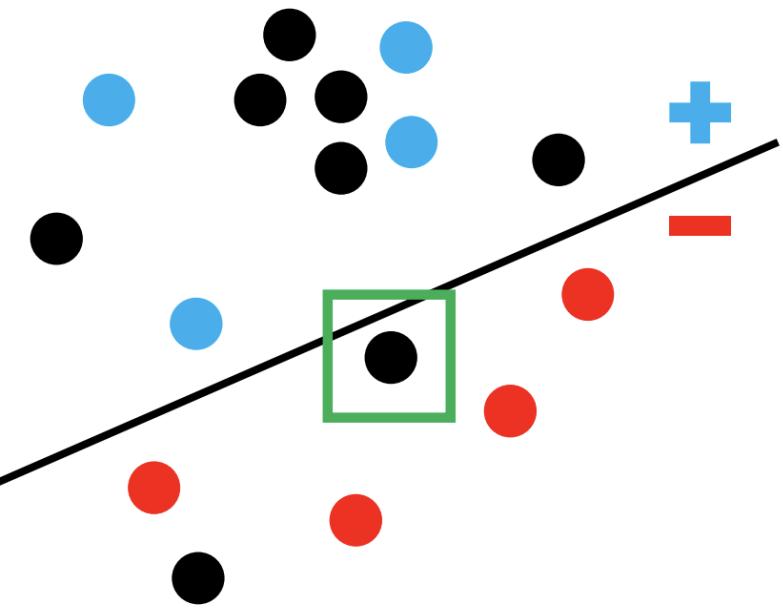
(a) Uncertainty sampling
e.g. [Tong and Koller, 2002;
Kapoor et al., 2010]



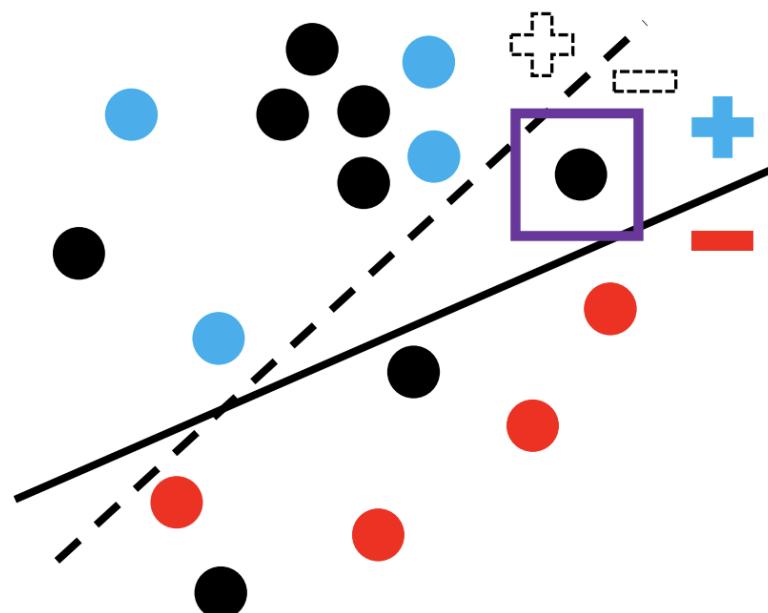
(b) Query by committee
e.g. [Seung et al., 1992;
Loy et al., 2012]

[Tong 2002, Kapoor 2010, Jain 2009, Kovashka 11, Vijayanarasimhan 2009, Branson 2010]

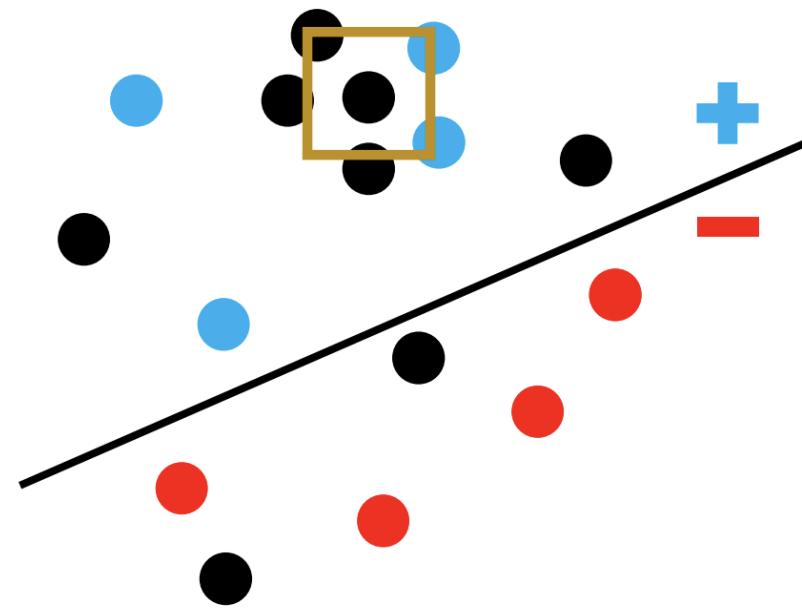
Active learning



(a) Uncertainty sampling
e.g. [Tong and Koller, 2002;
Kapoor et al., 2010]



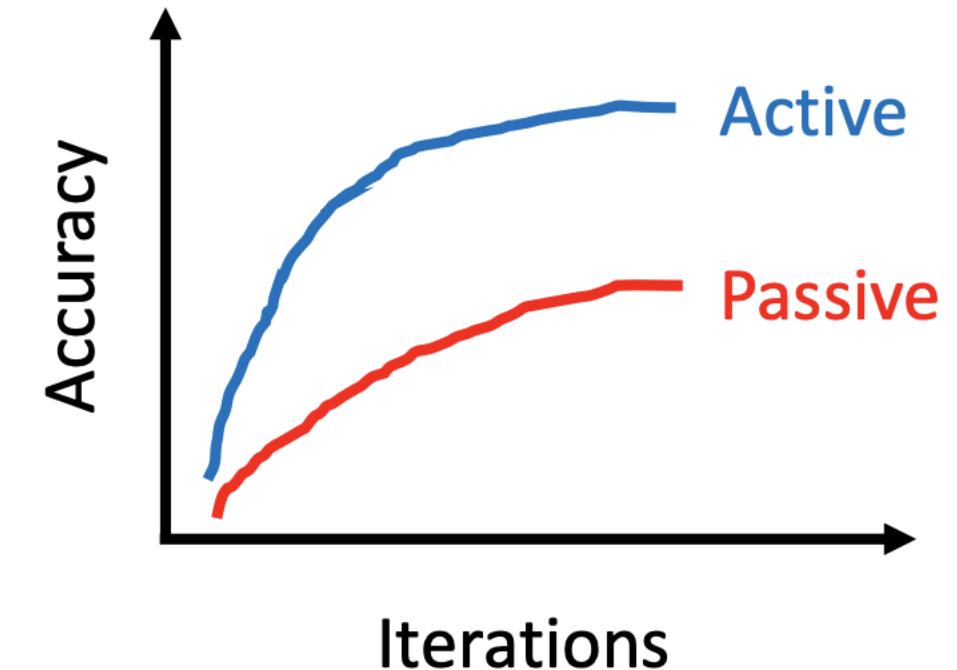
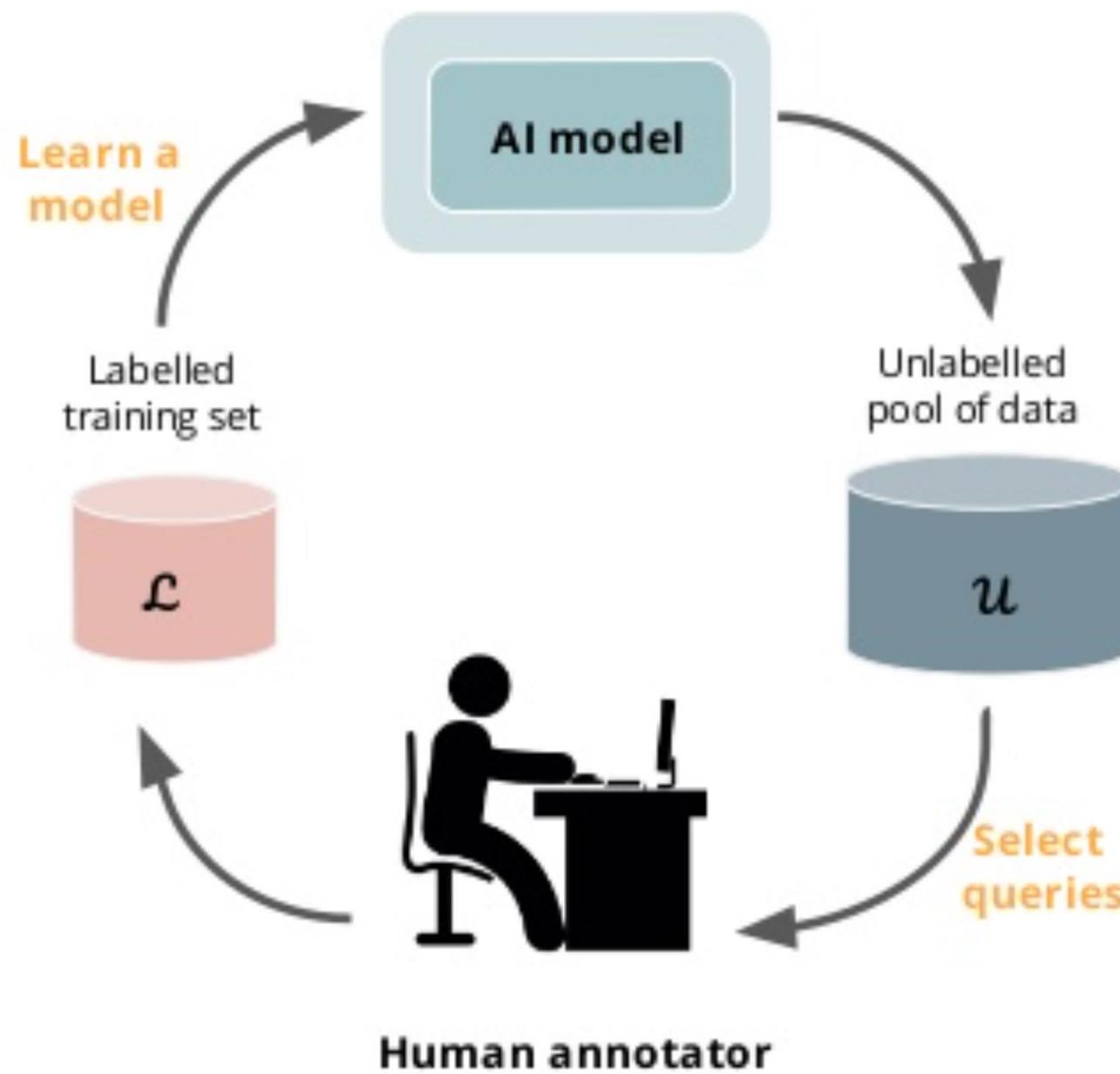
(b) Query by committee
e.g. [Seung et al., 1992;
Loy et al., 2012]



(c) Sampling from dense region
e.g. [Li and Guo, 2013]

[Tong 2002, Kapoor 2010, Jain 2009, Kovashka 11, Vijayanarasimhan 2009, Branson 2010]

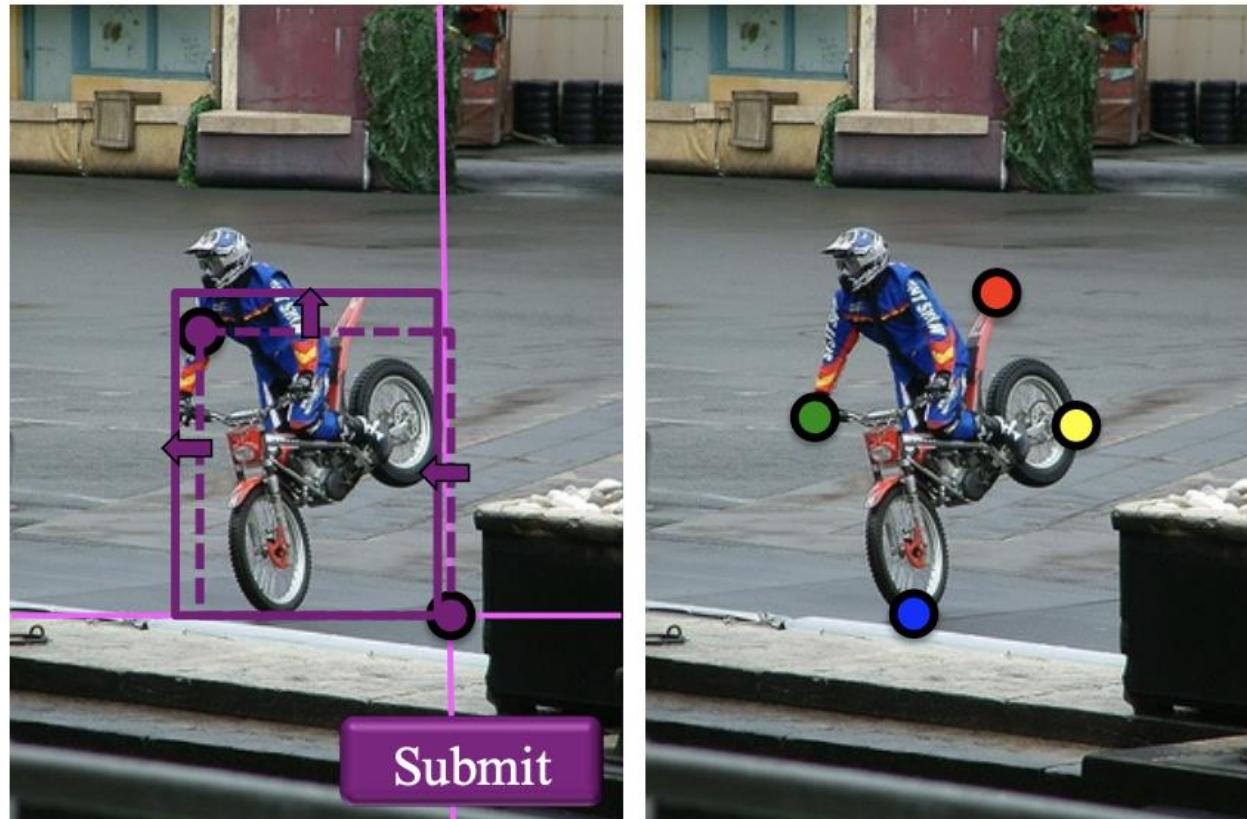
Active learning



- **Active vs. Passive:**
In practice the gap is very small
Selection is specific to one classifier
- **Generalization (?):**
Selection is based on one specific model

Active learning and annotation types

Clicking

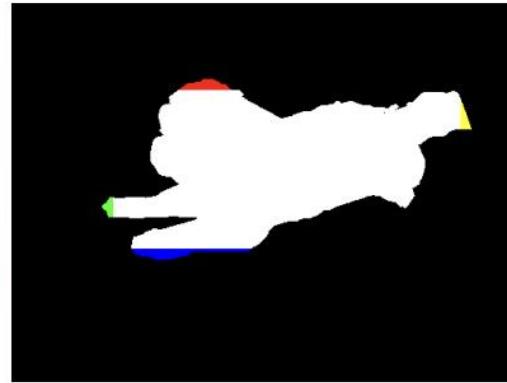


Clicking

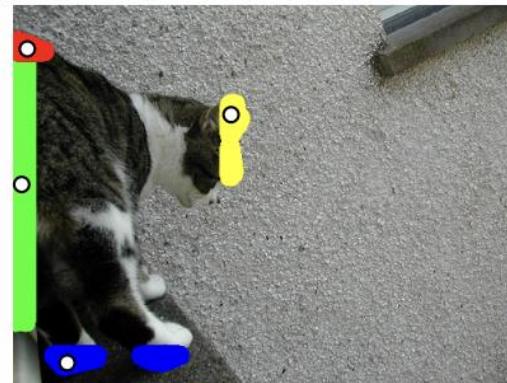
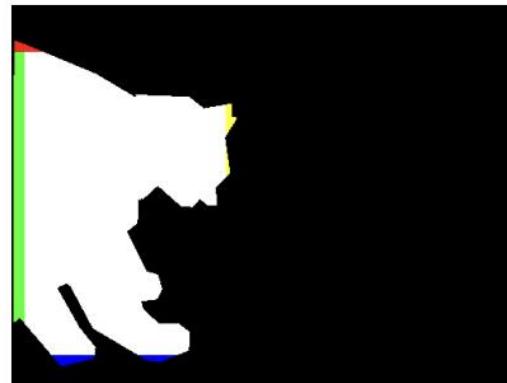
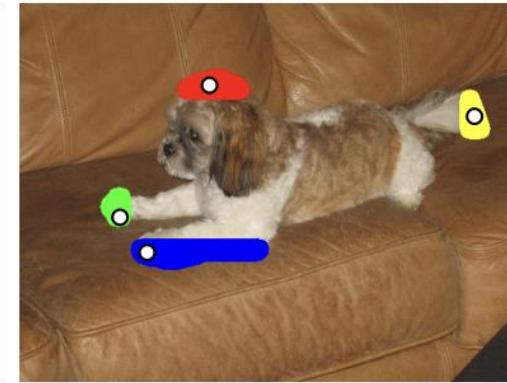
qualification image



extreme areas of segmentation mask



accepted areas of qualification image



Clicking

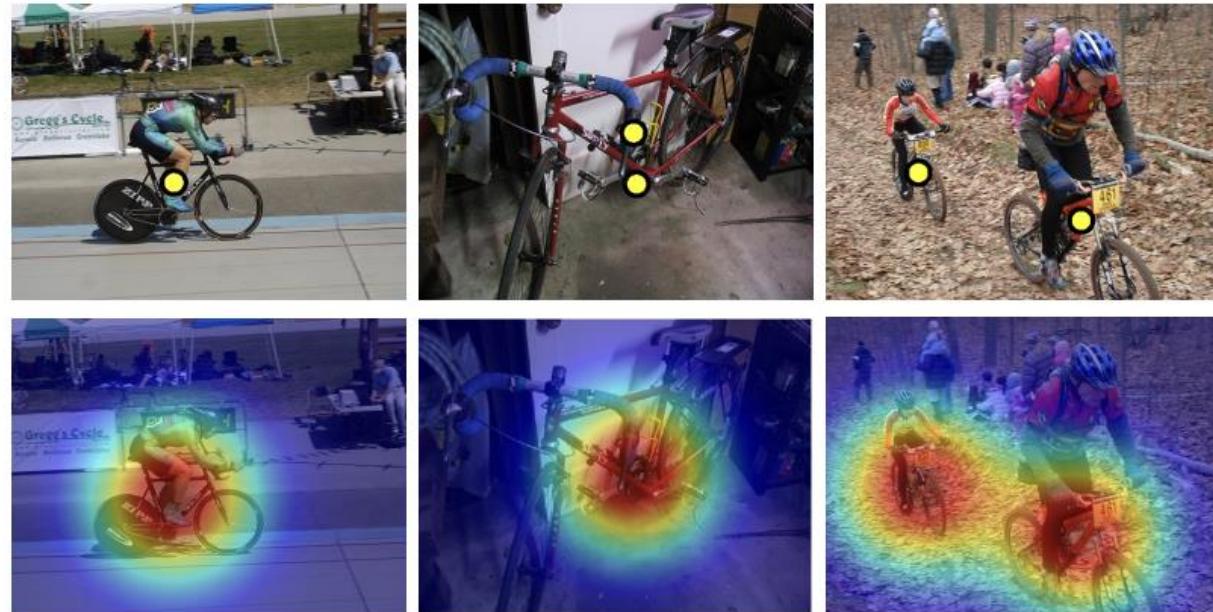
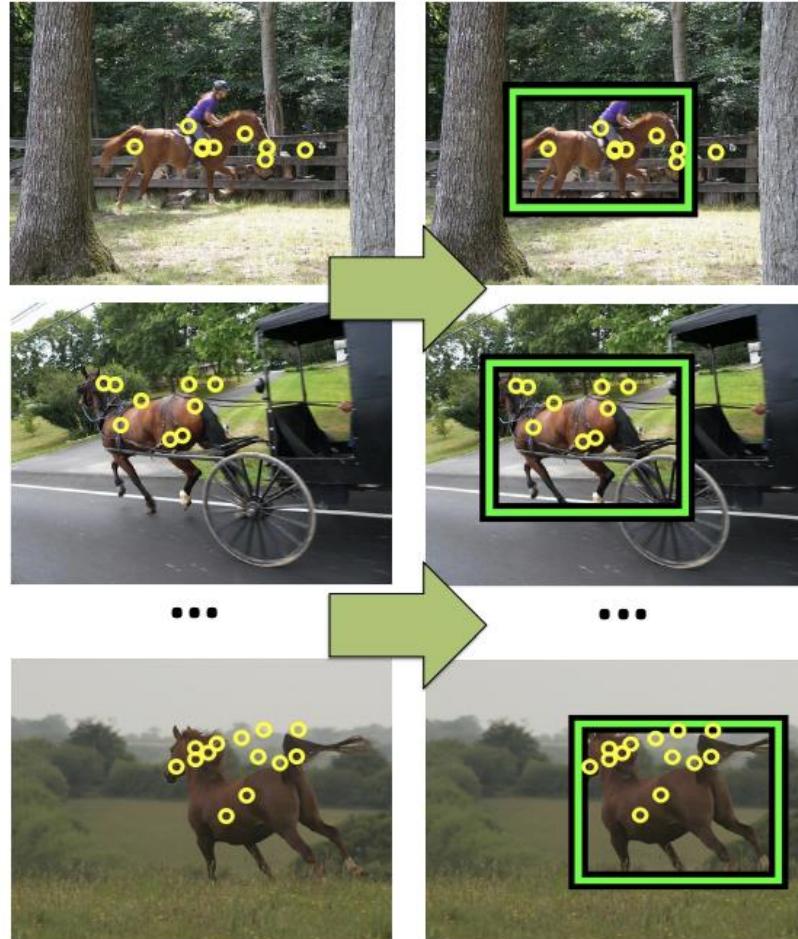


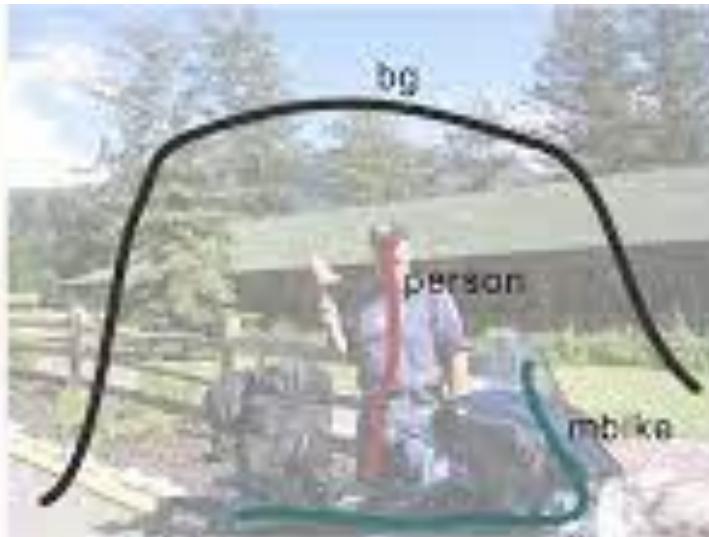
Figure 5. Box center score S_{bc} on bicycle examples. (left): One-click annotation. (middle): Two-click annotation on the same instance. (right): Two-click annotation on different instances. The values of each pixel in the heatmaps give the S_{bc} of an object proposal centered at that pixel.

Eye-tracking



- annotation time (1s per image)
- reduce annotation time by 6.8x
- simple annotation guidelines
- correct localizations in half images

Scribbles



More references

Actively selecting between different types of annotations (which images to label and at what level to label them)

e.g. [Vijayanarasimhan CVPR 2011]

Selecting batches of labels

e.g. [Liang 2014, Jain 2010, Vijayanarasimhan 2014]

Transfer learning

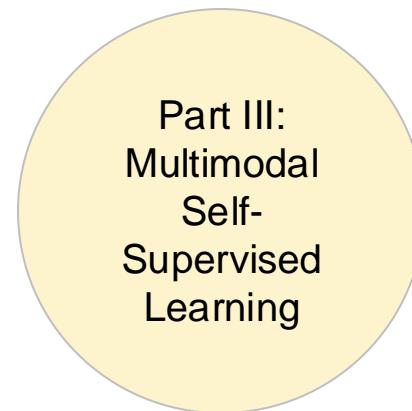
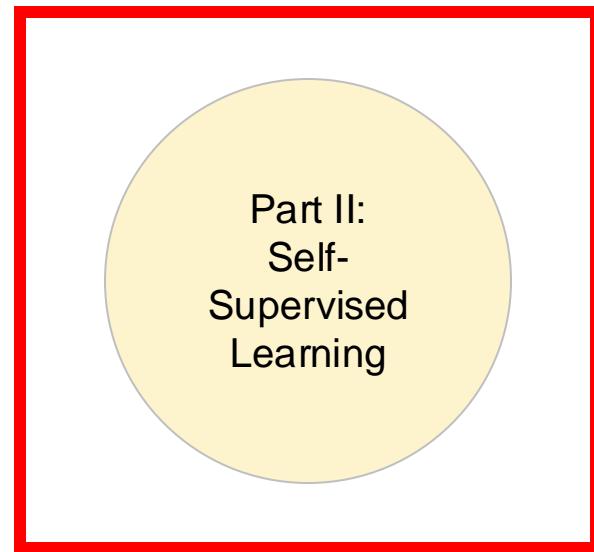
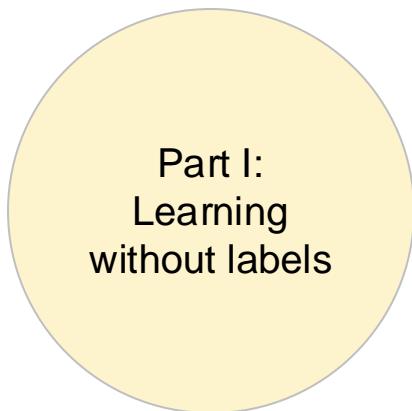
e.g. [Gavves 2015, Johns 2015]

Interactive annotation

e.g. [Branson ICCV 2011, Yao CVPR2012 ,Kovashka CVPR 2012, Wah CVPR 2013, Lad ECCV 2014,

Rubinstein ECCV 2012, Russakovsky CVPR 2015, Jain CVPR2016, Nagaraja ICCV2015,
Castrejon CVPR 2017, Acuna CVPR 2018, Benenson CVPR 2019, Ling CVPR 2019]

Today's lecture



Part II: Outline Self-Supervised Learning (SSL)



- Definition
- Motivation
- Benefits
- Methods
- Contrastive learning
- Masked Modeling
- Examples

Multimodal Self-supervised learning is everywhere

FINANCIAL TIMES

US COMPANIES TECH MARKETS CLIMATE OPINION WORK & CAREERS LIFE & ARTS HTSI

Artificial intelligence + Add to myFT

GPT-4 from OpenAI shows advances – ai moneymaking potential

Microsoft-backed group shifts towards showing less openness amid race to com AI systems

REUTERS®

World ▾ Business ▾ Markets ▾ Legal ▾ Breakingviews ▾ Technology ▾ Investigations

Disrupted

3 minute read · March 15, 2023 7:17 PM GMT+1 · Last Updated a month ago

Bar exam score shows AI can keep up with 'human law'

By Karen Sloan

VentureBeat

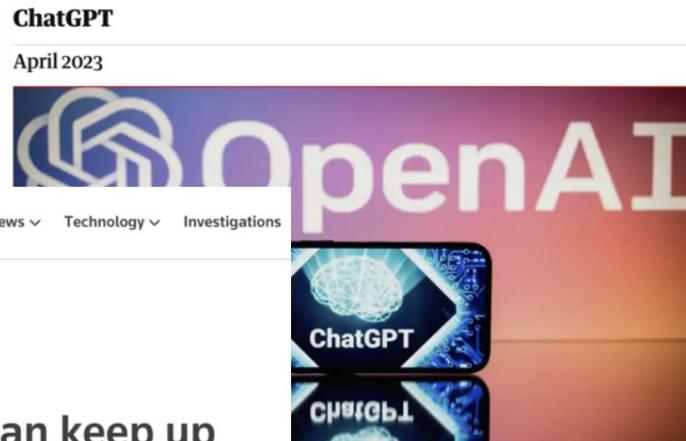
Why self-supervised learning is a medical AI game-changer

The Guardian

Support us →

News Opinion Sport Culture Lifestyle

World UK Coronavirus Climate crisis Environment Science Global development More



deVolkskrant

Topverhalen vandaag Opinie Cultuur & Media Podcasts Beter Leven Log in

ZES VRAGEN

Nieuwe 'turbo-versie' van ChatGPT is een stuk veelzijdiger en kan ook omgaan met plaatjes



GPT-4 The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

Slide credit: Y.Asano

Self-Supervised Learning

The ImageNet Challenge Story

Classification Results (CLS)



The ImageNet Challenge Story

- Strong supervision
- Features from networks trained on ImageNet can be used for other visual tasks, e.g. detection, segmentation, action recognition, fine grained visual classification
- To some extent, *any* visual task can be solved now by:
 - Construct a large-scale dataset labelled for that task
 - Specify a training loss and neural network architecture
 - Train the network and deploy
- What else can we do?

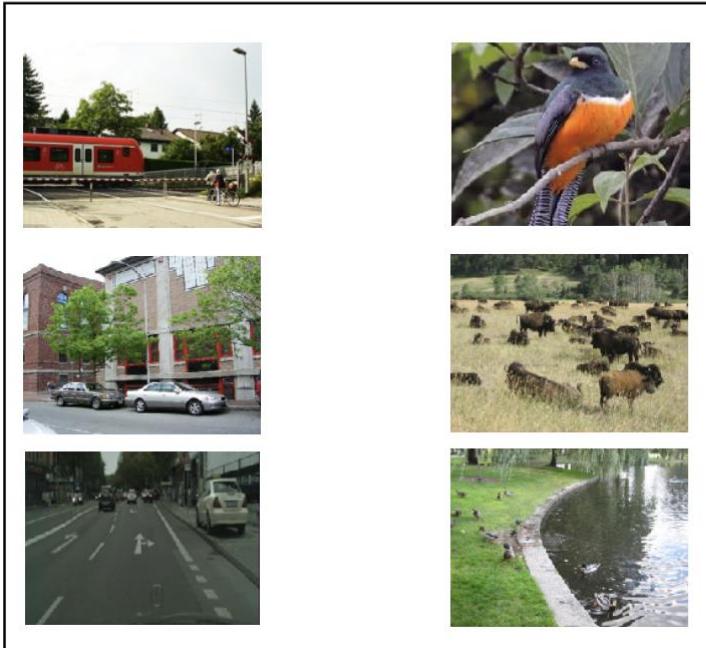
Self-supervised Learning

Motivation

- + Expense of producing a new dataset for each new task
- + Some areas are supervision-starved, e.g. medical data, where it is hard to obtain annotations
- + Availability of vast numbers of unlabelled images/videos
 - Facebook: one billion images uploaded per day
 - 300 hours of video are uploaded to YouTube every minute
- + How infants may learn ...

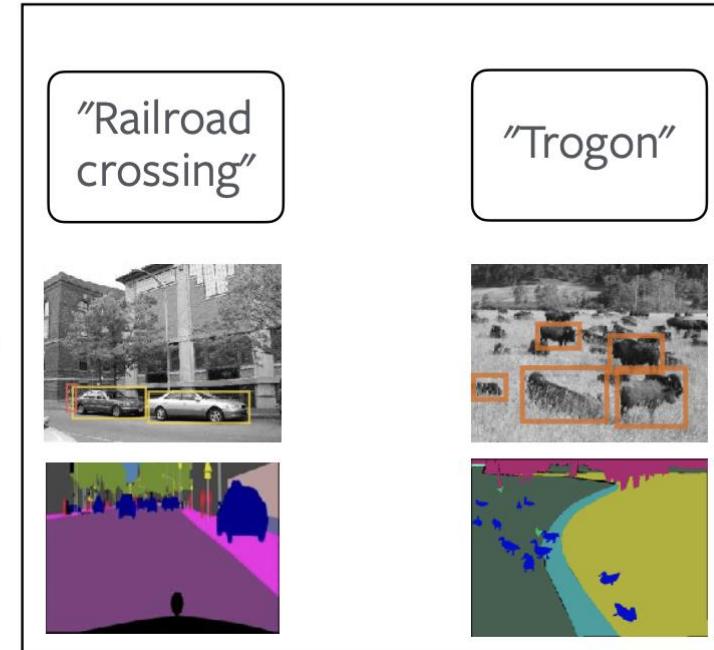
Self-supervised Learning Motivation

Images are often cheap



Supervised Learning

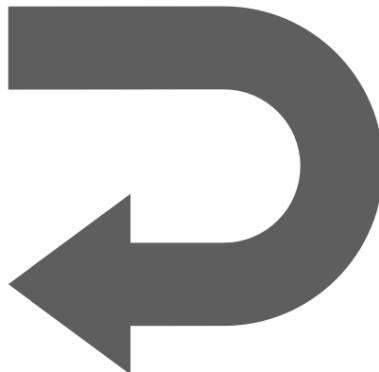
But manual annotations are expensive:
e.g. 30min per image / requiring experts



Self-supervised Learning

Motivation

Images are often cheap



Self-supervision

Extract a supervisory signal
from the raw data alone

Self-supervised Learning

Definition



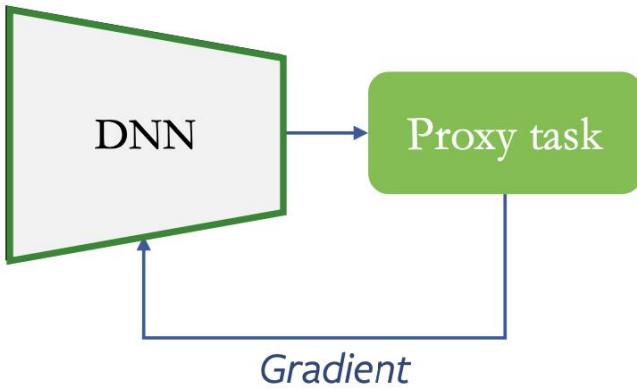
- A form of unsupervised learning where the data provides the supervision
- In general, withhold some part of the data, and task the network with predicting it
- The task defines a proxy loss, and the network is forced to learn what we really care about, e.g. a semantic representation, in order to solve it

Self-supervised Learning Procedure

Phase 1: Pretraining



Unlabelled data
+ transformations



Types:

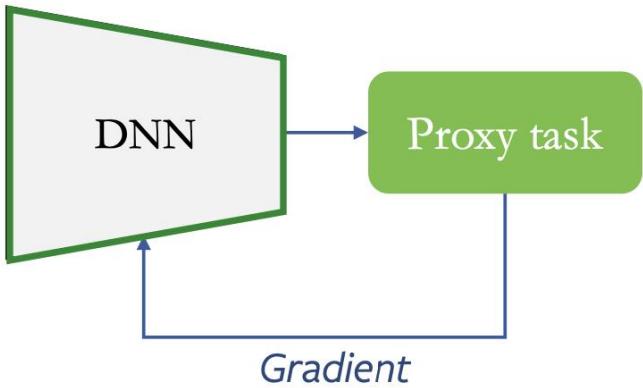
- Geometry based
- Clustering
- Contrastive
- Generative (partial/full)
- (more)

Self-supervised Learning Procedure

Phase 1: Pretraining



Unlabelled data
+ transformations



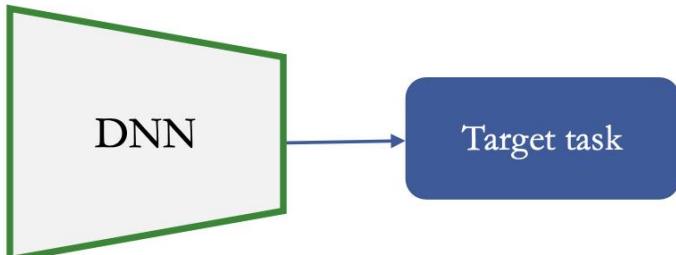
Types:

- Geometry based
- Clustering
- Contrastive
- Generative (partial/full)
- (more)

Phase 2: Downstream tasks



(Sparse) labeled data

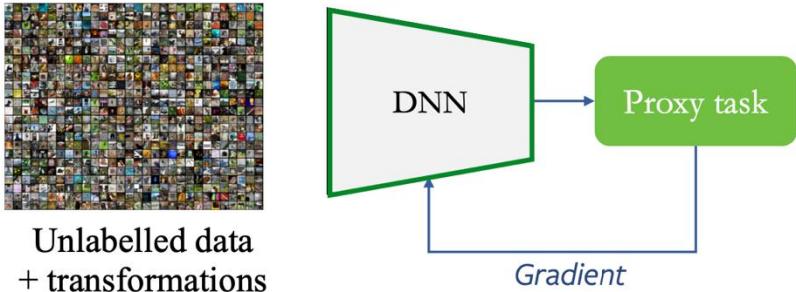


Types:

- Limited fine-tuning (e.g. linear layer)
- Finetuning (w/ full or fraction of database)

Self-supervised Learning Procedure

Phase 1: Pretraining



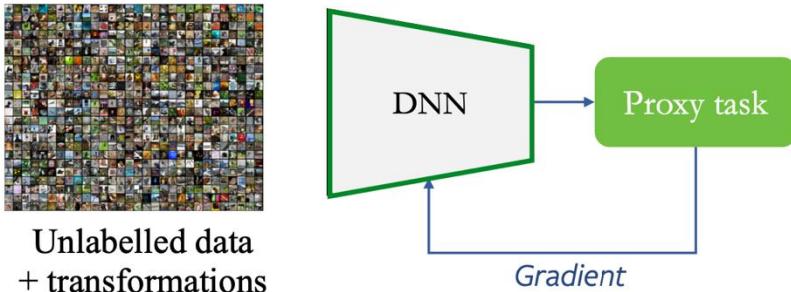
Phase 2: Downstream tasks



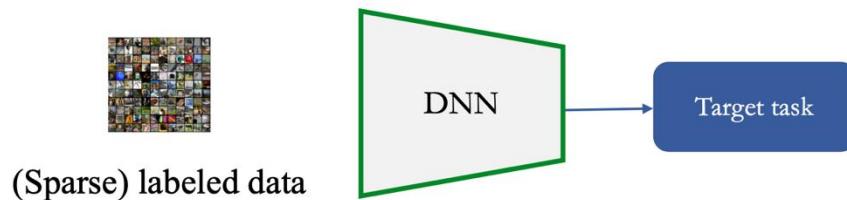
Representation Learning

Self-supervised Learning Procedure

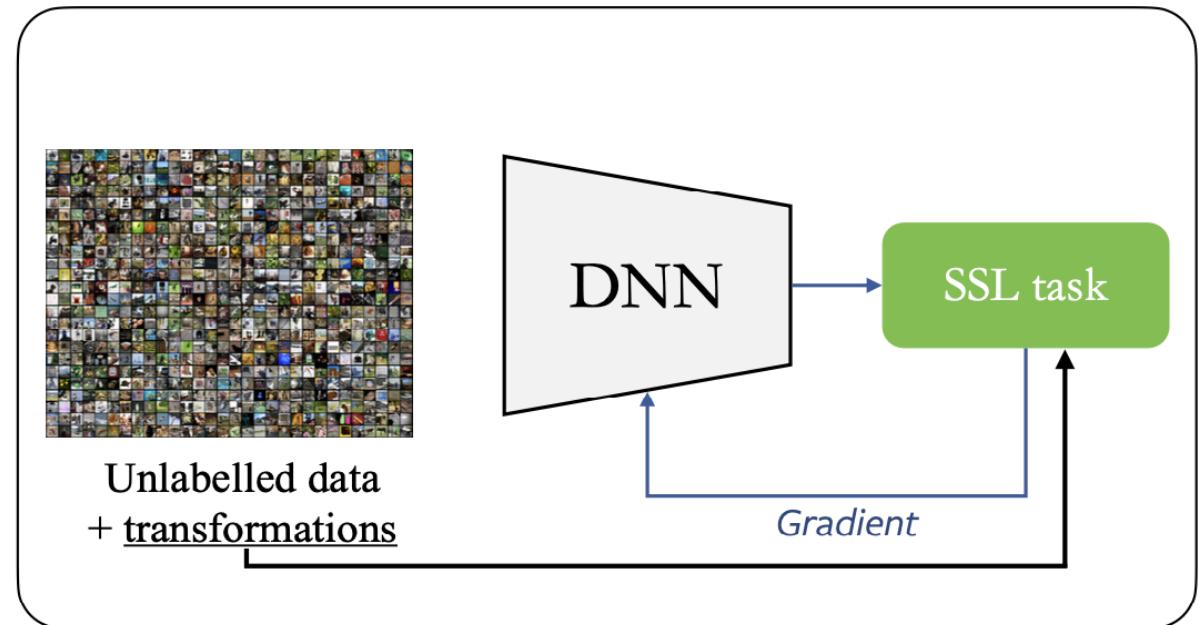
Phase 1: Pretraining



Phase 2: Downstream tasks



Representation Learning



Useful Self-supervised Learning, e.g.
 SSL object detection & segmentation
 SSL speaker detection, SSL dataset labelling etc..

10

Why SSL?

Why SSL?

1. Scalability



(above) x 50 = 1.2M images



$90\text{ms} * 1.2\text{M} = 30\text{h}$

Instagram: >50B images

Annotation is expensive, yet datasets
keep getting bigger

Why SSL?

2. Changing domains



Unclear when & what to relabel. Again, large costs just to "keep up"

Why SSL?

3. Accessibility & generalisability



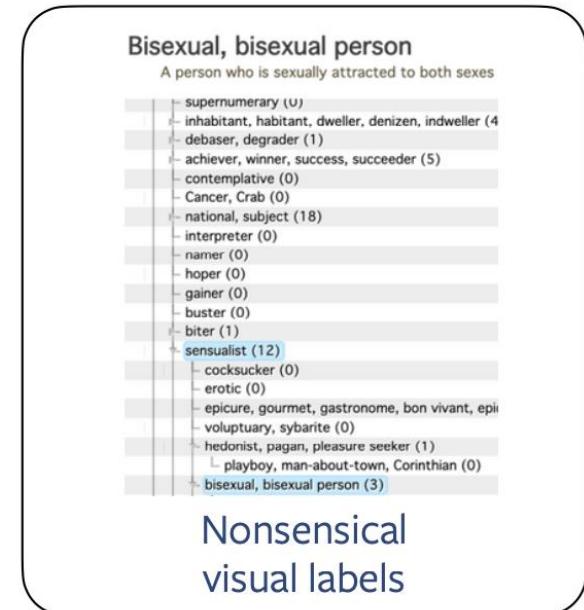
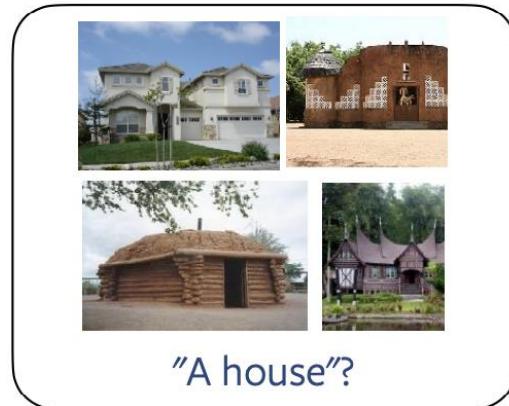
Query image		Top 3 nearest neighbours		
LP-Ubeidiya				Lower Palaeolithic-Gesher Benot Ya'aqov
ER-Caesarea				Early Roman-Caesarea
L-P-Tannur				Lower Palaeolithic- Tannur
Bx-En Gedi				Byzantine- En Gedi
Crus-Atilt				Crusader- Atilt
				Crusader- Arshaf



Pretrained models are very useful for a variety of tasks.

Why SSL?

4. Ambiguity of labels



Labels are ambiguous at best, discriminating and bias-propagating at worst.
Do we really wish to provide our models with these priors?

Why SSL?

5. Investigating the fundamentals of visual understanding

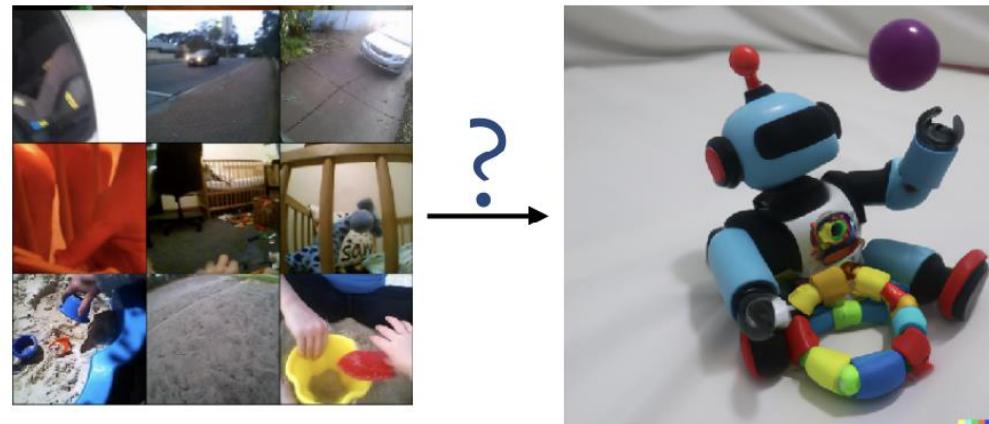
 MetaAI

RESEARCH

Self-supervised learning: The dark matter of intelligence

March 4, 2021

As babies, we learn how the world works largely by observation. We form generalized predictive models about objects in the world by learning concepts such as object permanence and gravity. Later in life, we observe the world, act on it, observe again, and build hypotheses to explain how our actions change our environment by trial and error.



What, if there are, are the limits of learning without labels?

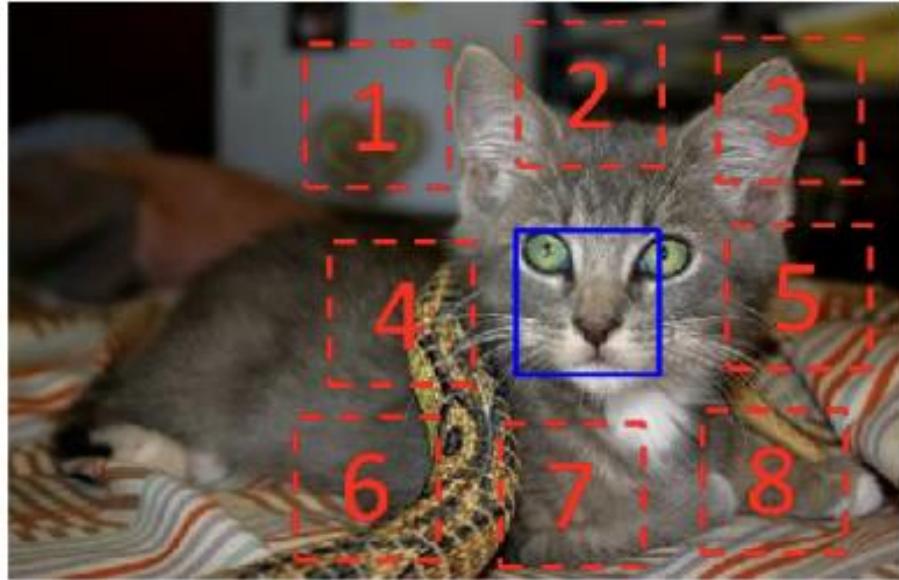
Pop Quiz

Can you think of examples?

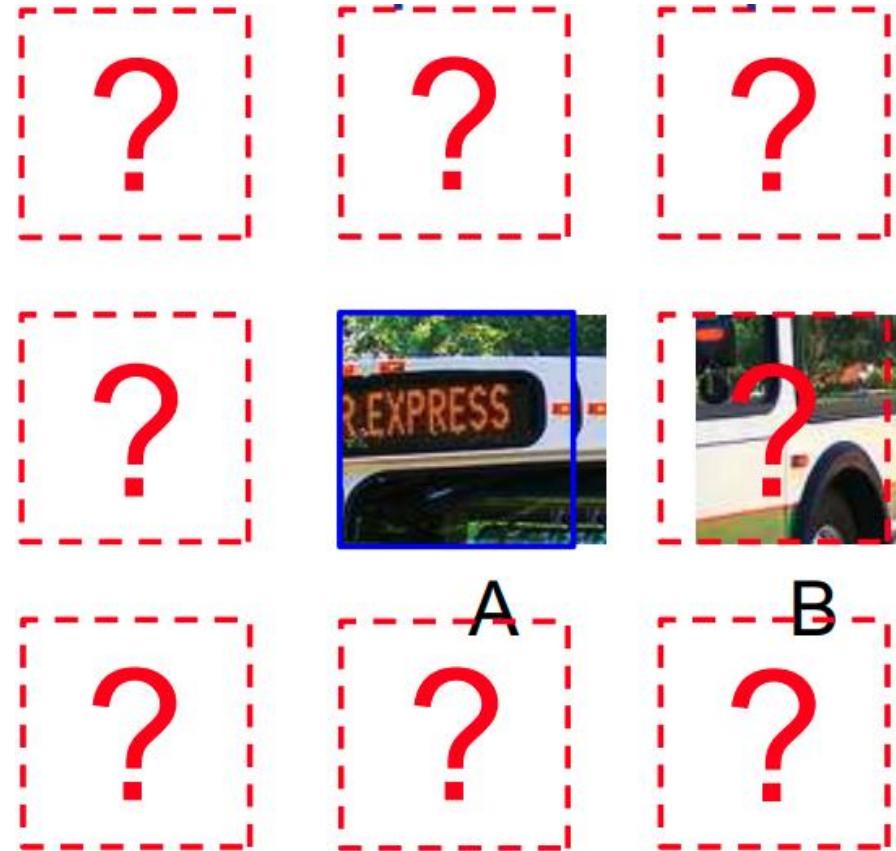
Methods for SSL

Self-supervised Learning

Relative positioning → Context prediction



$$X = (\begin{matrix} \text{cat eye} \\ \text{cat ear} \end{matrix}, \text{blanket}) ; Y = 3$$

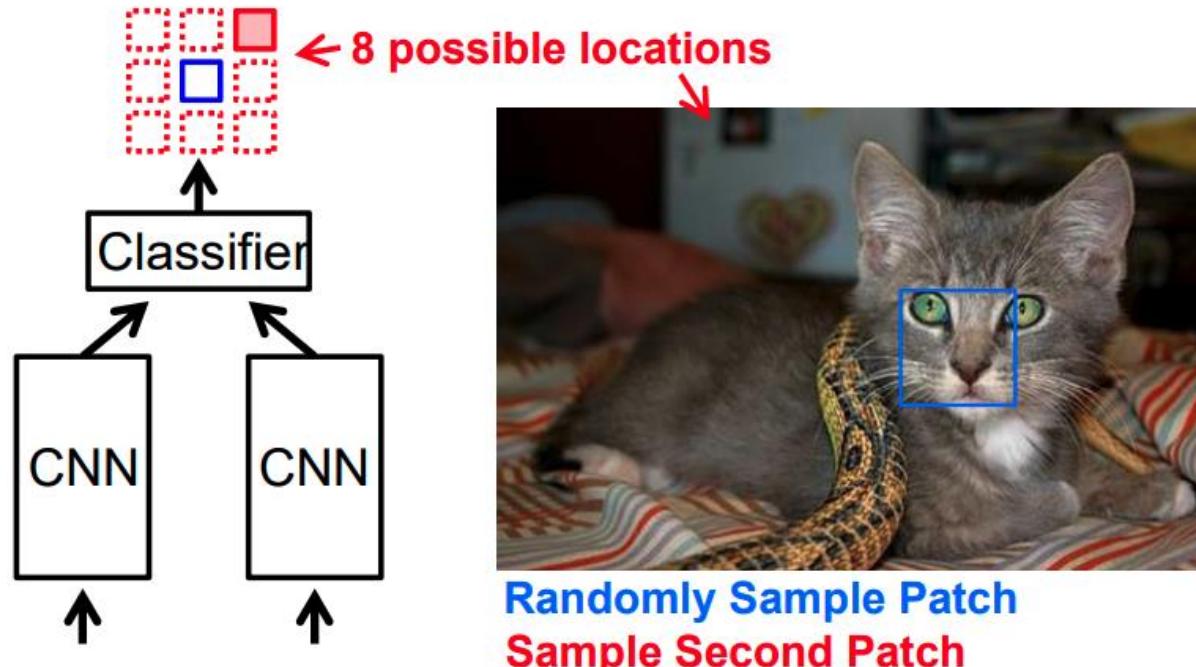


Self-supervised Learning

Relative positioning → Context prediction



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

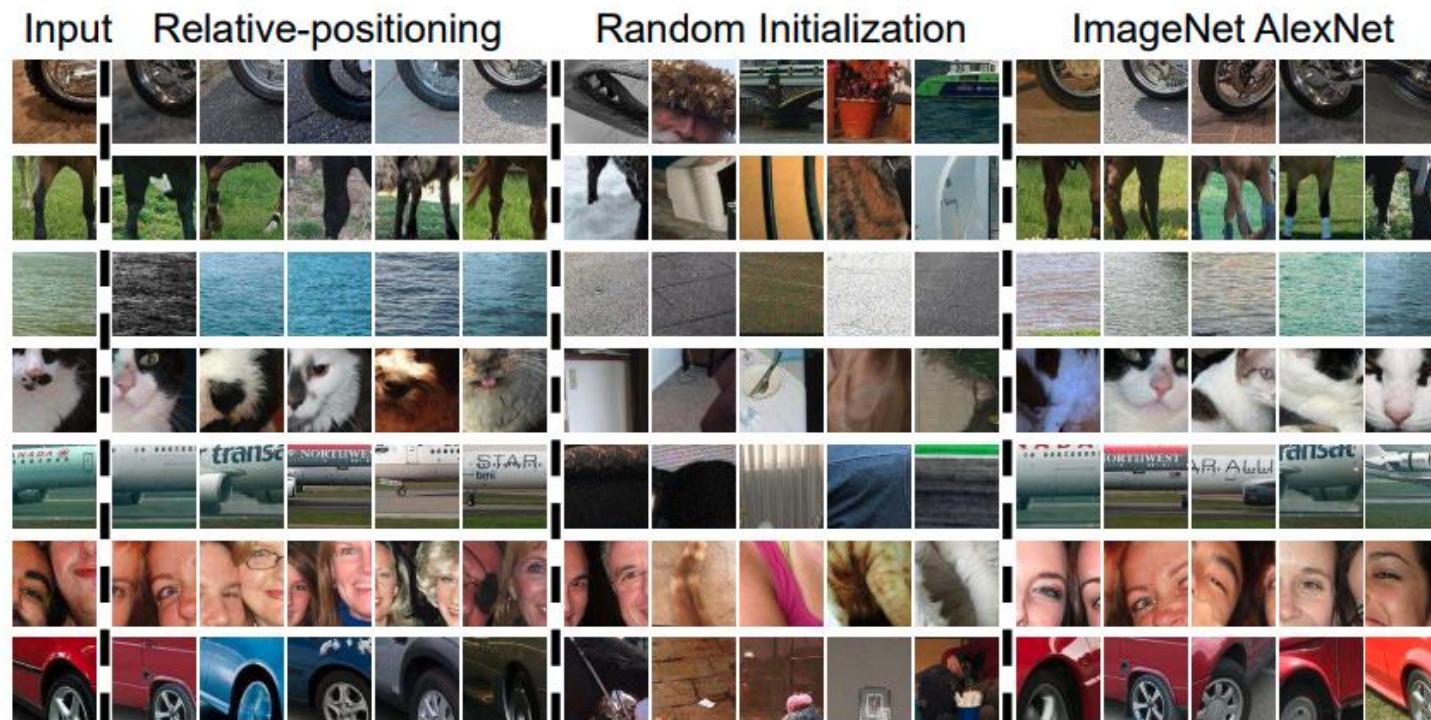
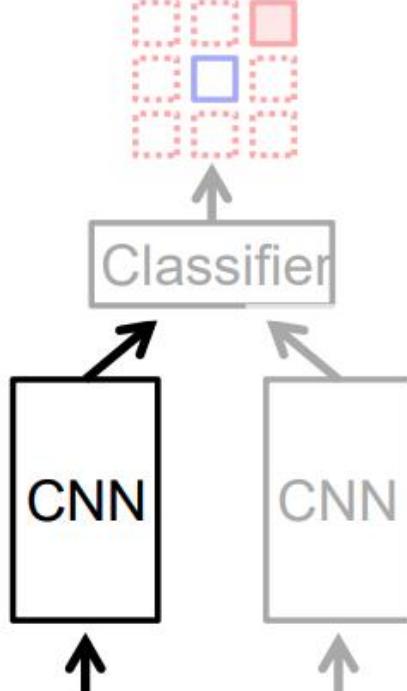


Unsupervised visual representation learning by context prediction,
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

Self-supervised Learning

Relative positioning → Context prediction

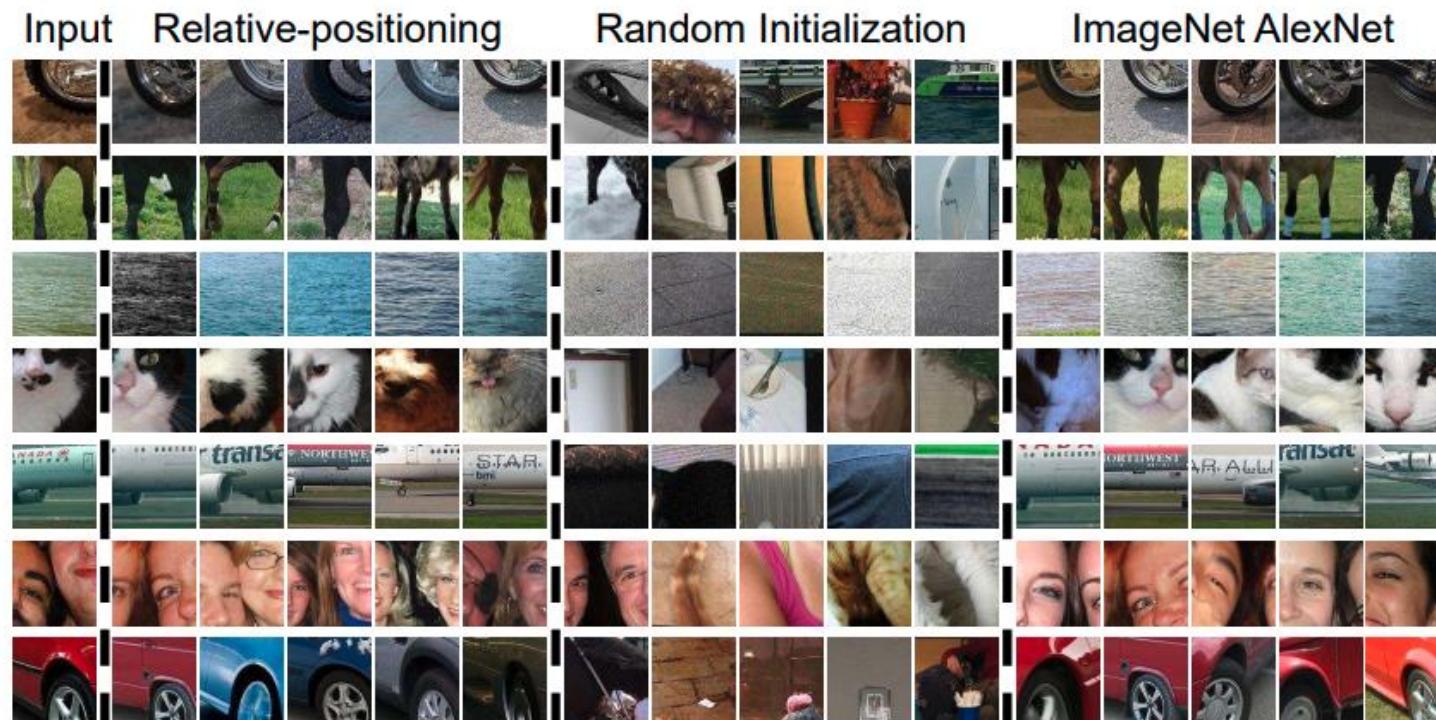
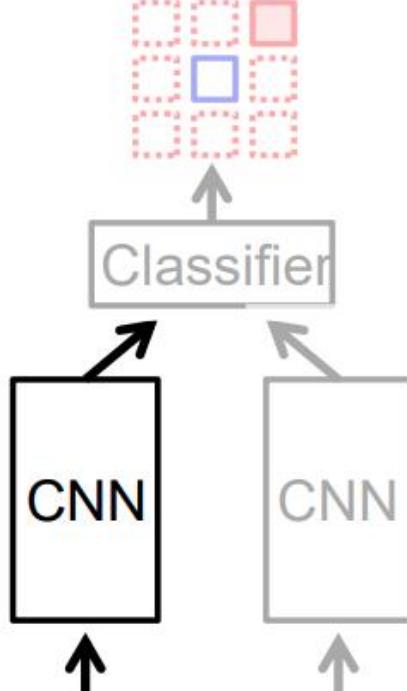
What is the network learning?



Self-supervised Learning

Relative positioning → Context prediction

What is the network learning?

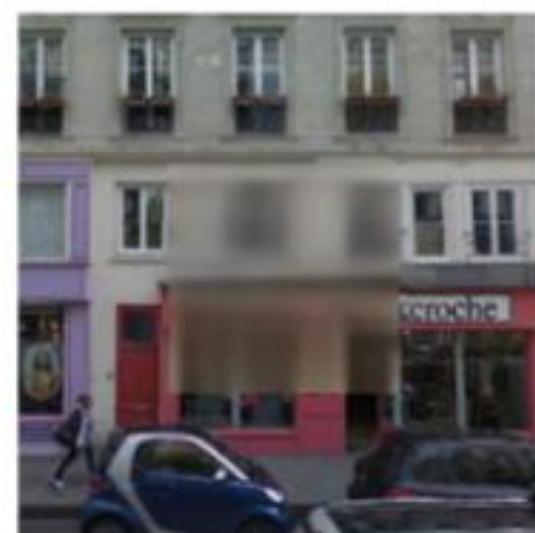


Self-supervised Learning

Context Encoder → Context prediction



(a) Input context



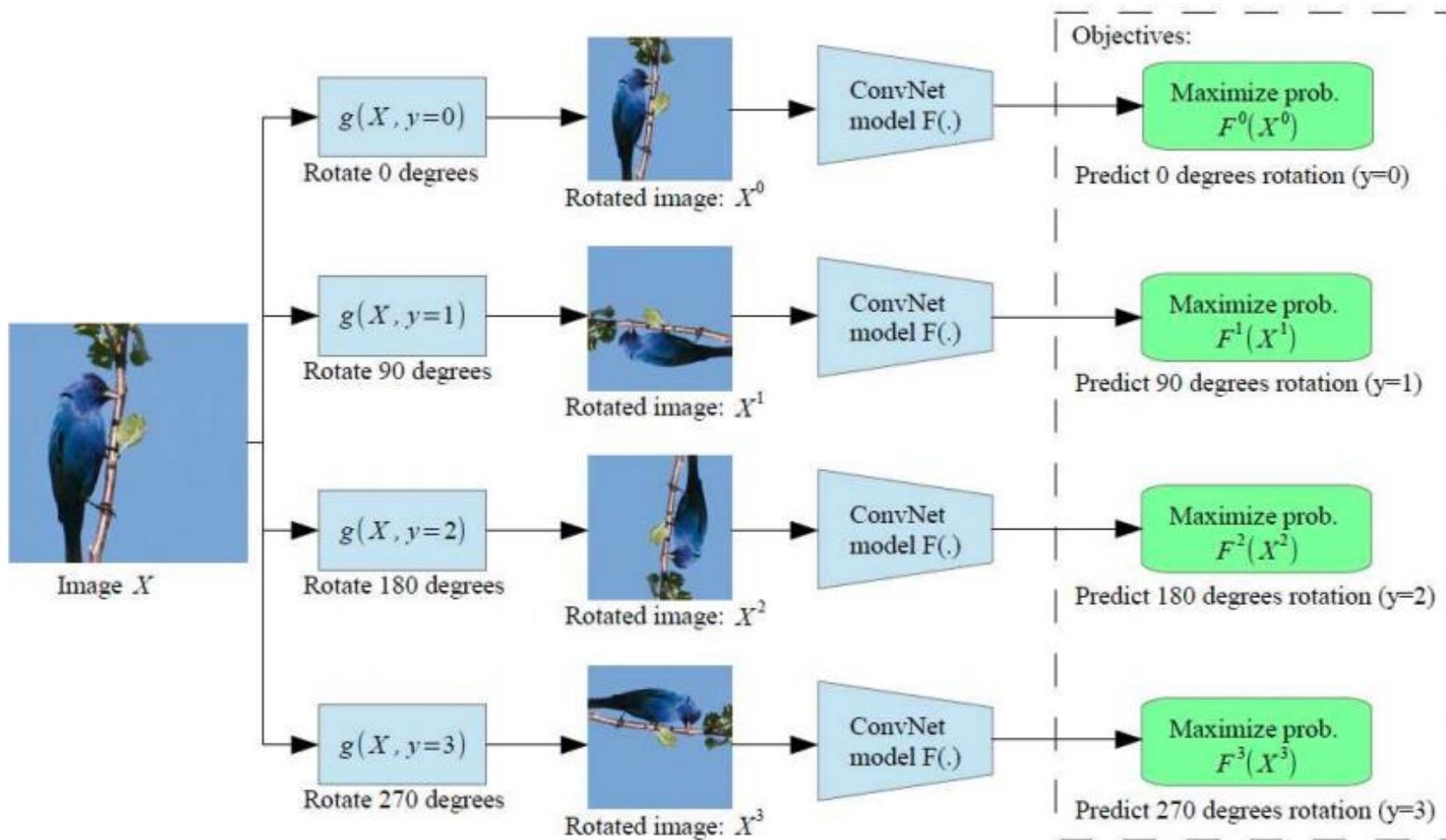
(c) Context Encoder
(L_2 loss)

https://www.researchgate.net/figure/Word2Vec-CBOW-and-Skip-gram-There-are-two-different-methods-in-the-Word2Vec-algorithm_fig2_320829283

Doersch et al. Unsupervised Visual Representation Learning by Context Prediction. ICCV 2015.

Pathak et al. Context Encoders: Feature Learning by Inpainting. CVPR 2016.

Self-supervised Learning Image Transformations



Self-supervised Learning

Image Transformations



But:



Self-supervised Learning Image Transformations

- AlexNet
- Closes gap between ImageNet and self-supervision

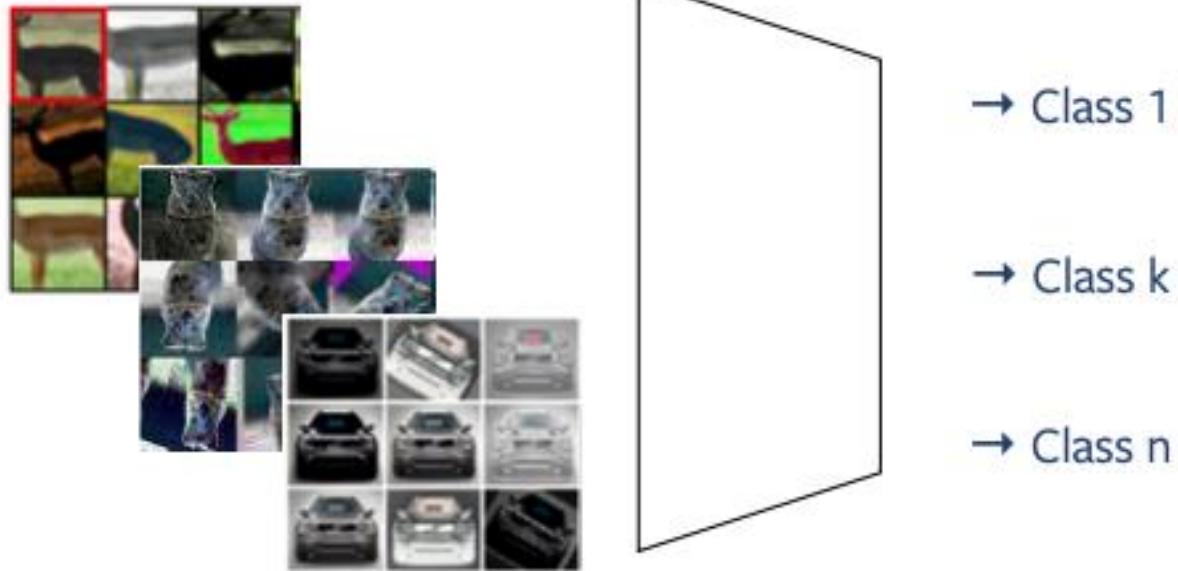
	PASCAL VOC Detection mAP
Random	43.4
Rel. Pos.	51.1
Colour	46.9
Rotation	54.4
ImageNet Labels	56.8

Contrastive learning

Self-supervised Learning

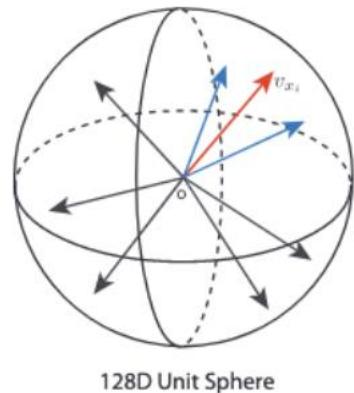
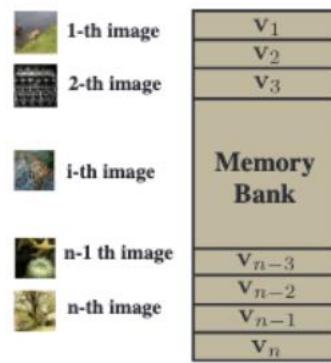
Image Uniqueness

- Exemplar CNN, precursor to contrastive learning

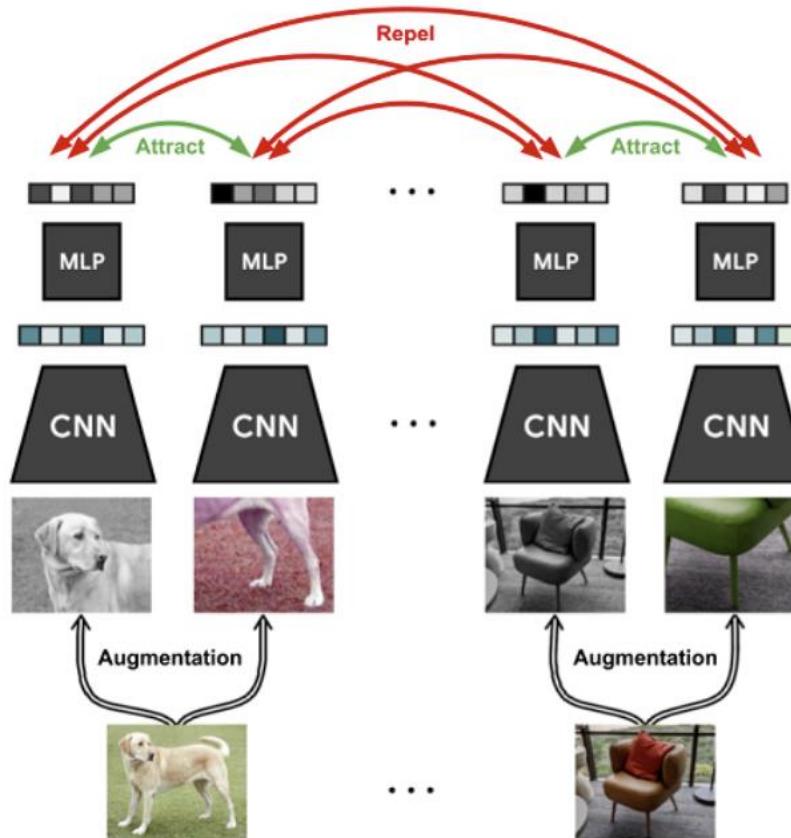


Uses image-uniqueness and enforces augmentation-invariance

Noise-contrastive self-supervised learning



NPID

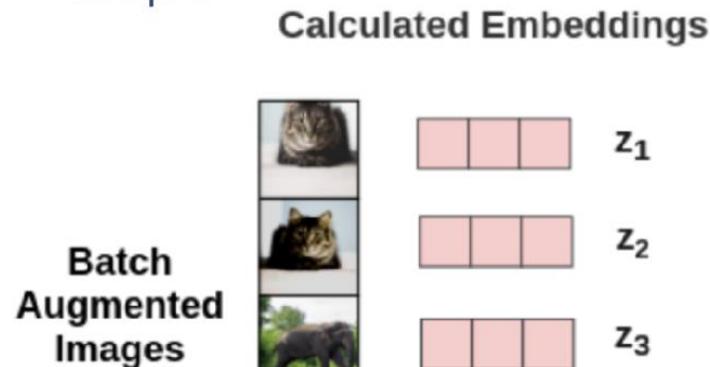


Wu et al. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. CVPR 2018

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

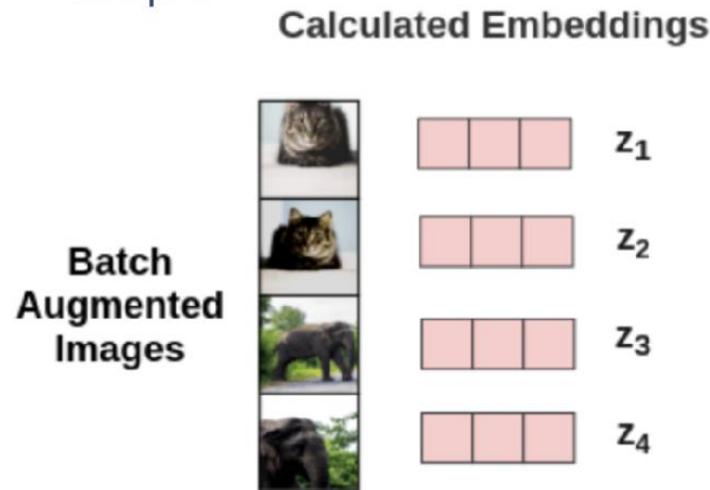
SimCLR

Step 1



SimCLR

Step 1



Step 2

Similarity Calculation of Augmented Images

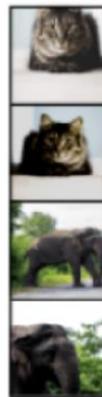
$$\text{similarity}(\underset{x_i}{\boxed{\text{tiger}}}, \underset{x_j}{\boxed{\text{tiger}}}) = \frac{\text{cosine similarity}(\underset{z_i}{\boxed{\text{pink}}}, \underset{z_j}{\boxed{\text{pink}}})}{\underset{\tau ||z_i|| ||z_j||}{\text{similarity}}}$$

$$s_{i,j} = \frac{z_i^T z_j}{(\tau ||z_i|| ||z_j||)}$$

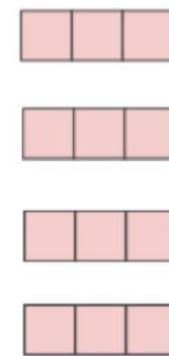
- τ is the adjustable temperature parameter. It can scale the inputs and widen the range [-1, 1] of cosine similarity
- $\|z_i\|$ is the norm of the vector.

SimCLR

Step 1
Batch Augmented Images



Calculated Embeddings



z_1

z_2

z_3

z_4

Step 2

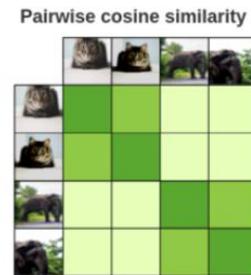
Similarity Calculation of Augmented Images

$$\text{similarity}(\underset{x_i}{\boxed{\text{cat}}}, \underset{x_j}{\boxed{\text{cat}}}) = \text{cosine similarity} \left(\underset{z_i}{\boxed{\text{---}}}, \underset{z_j}{\boxed{\text{---}}} \right)$$

$$s_{i,j} = \frac{z_i^T z_j}{(\tau \|z_i\| \|z_j\|)}$$

- τ is the adjustable temperature parameter. It can scale the inputs and widen the range [-1, 1] of cosine similarity
- $\|z_i\|$ is the norm of the vector.

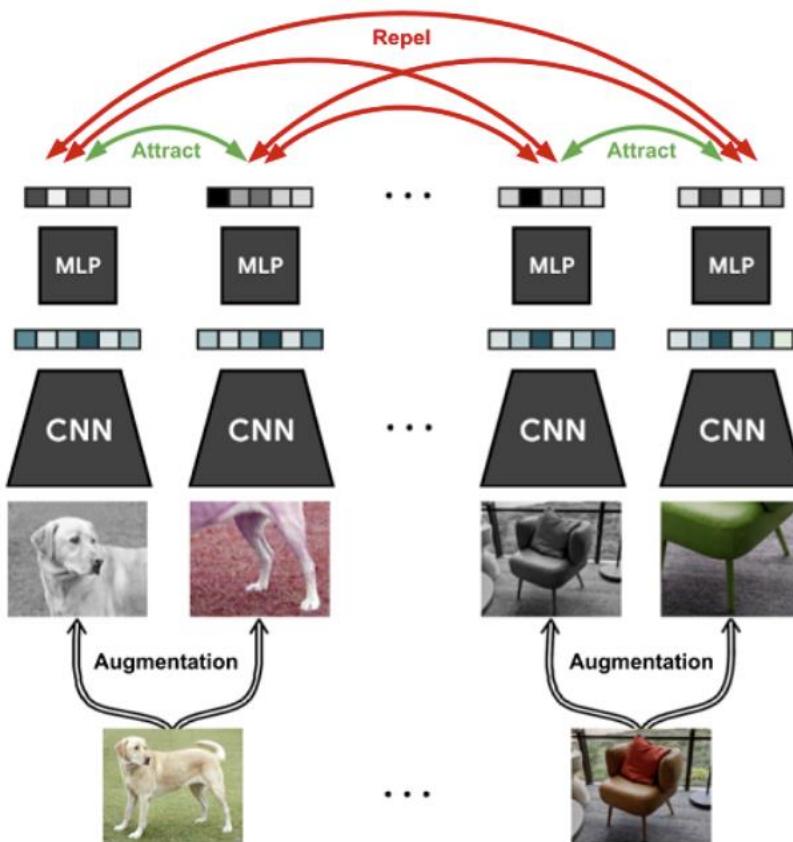
Step 3



Loss: relatively increase similarity for pairs, decrease rest

What happens if you only try to increase the diagonal?

SimCLR



SimCLR

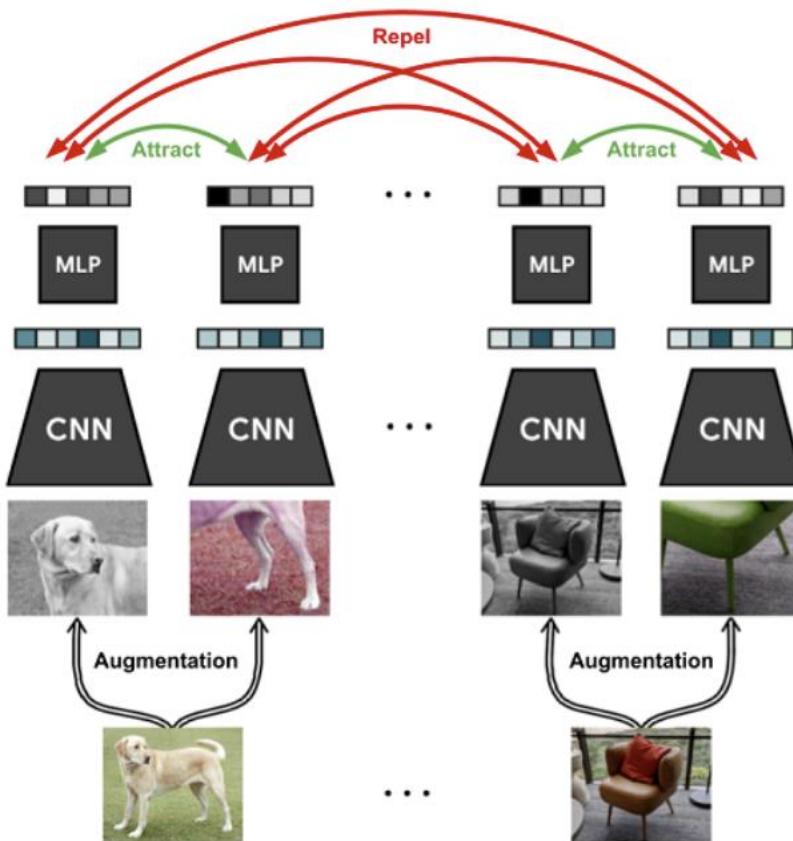
The contrastive loss for positive pairs i, j :

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

with $\mathbf{z}_i, \mathbf{z}_j$ embeddings for images i and j ,
 τ a temperature, $\text{sim}()$ is the dot-product

"non-parametric" softmax

SimCLR



SimCLR

The contrastive loss for positive pairs i, j :

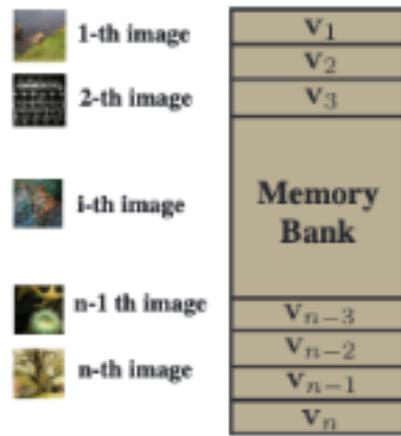
$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} [\mathbf{k} \neq i] \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)},$$

with $\mathbf{z}_i, \mathbf{z}_j$ embeddings for images i and j ,
 τ a temperature, $\text{sim}()$ is the dot-product

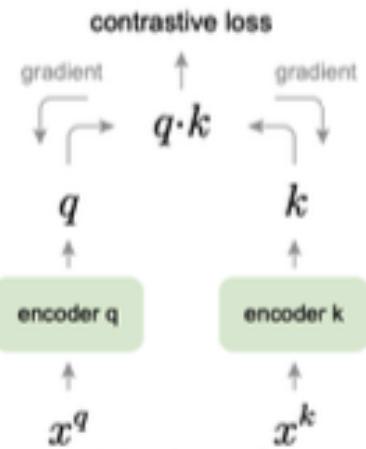
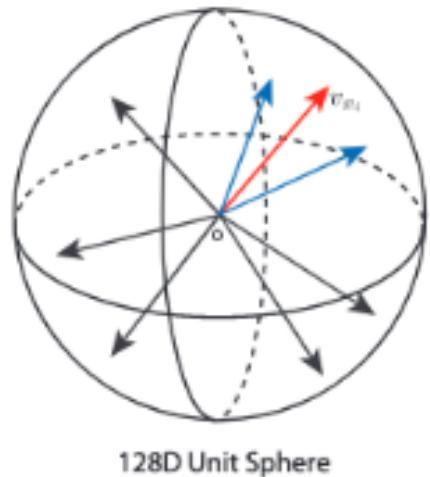
"non-parametric" softmax

Enforces image-uniqueness and
enforces augmentation-invariance

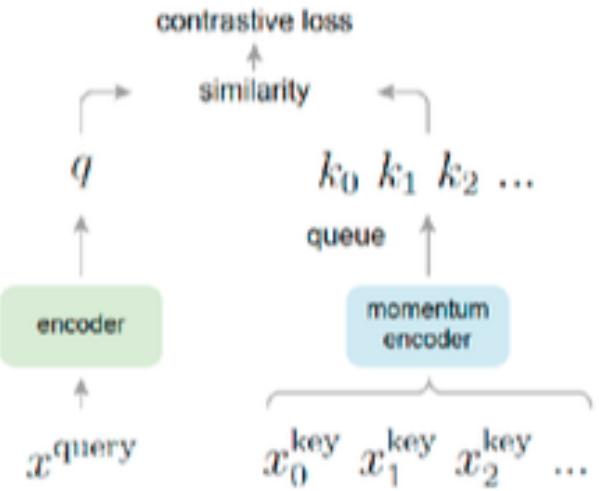
Modern Noise-contrastive self-supervised learning



NPID



SimCLR

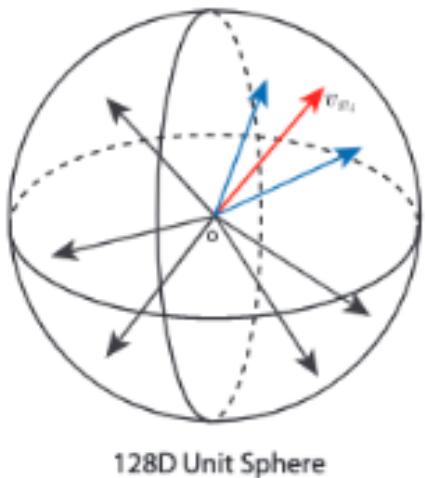
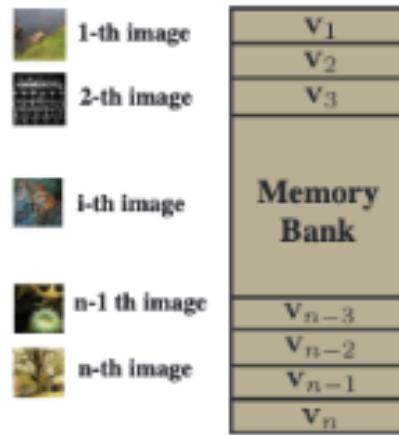


MoCo

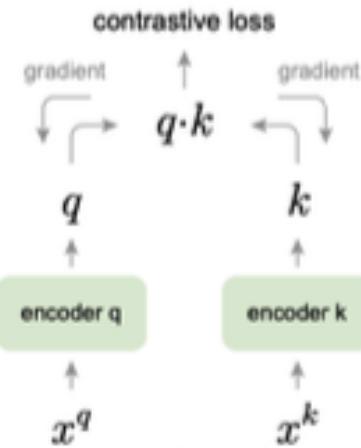
Wu et al. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. CVPR 2018
Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020
He et al. Momentum Contrast for Unsupervised Visual Representation Learning. CVPR 2020

Modern Noise-contrastive self-supervised learning

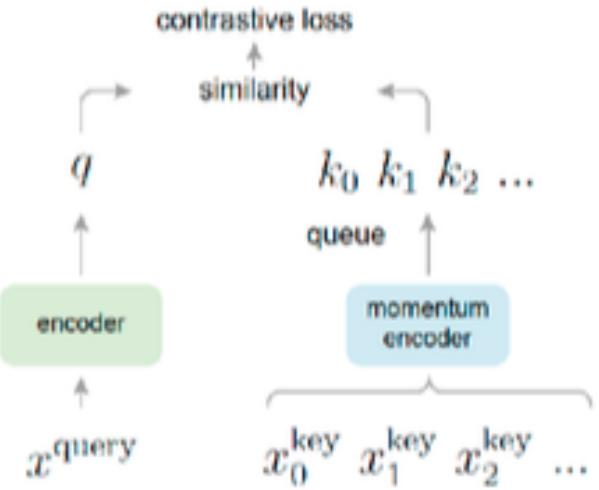
The start of large-scale & industrial self-supervised learning.
These works heavily rely on image augmentations.



NPID



SimCLR



MoCo

Momentum encoder:

```
# momentum update: key network
f_k.params = m*f_k.params+(1-m)*f_q.params
```

Wu et al. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. CVPR 2018

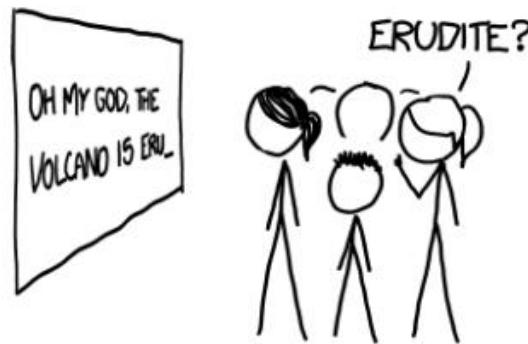
Chen et al. A Simple Framework for Contrastive Learning of Visual Representations. ICML 2020

He et al. Momentum Contrast for Unsupervised Visual Representation Learning. CVPR 2020

Masked Modeling

Language Modelling

Why "erudite" is not a good guess



Slide credit: Y.Asano

<https://thegradient.pub/understanding-evaluation-metrics-for-language-models/>
Radford et al. Improving Language Understanding by Generative Pre-Training . 2018

Language Modelling

Why "erudite" is not a good guess



Factor the probability of a datapoint (w_1, \dots, w_n):

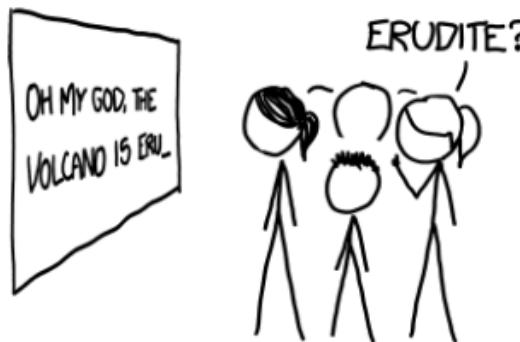
$$P(w_1, w_2, \dots, w_n) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2)\dots p(w_n|w_1, w_2, \dots, w_{n-1}) \\ = \prod_{i=1}^n p(w_i|w_1, \dots, w_{i-1})$$

Slide credit: Y.Asano

<https://thegradient.pub/understanding-evaluation-metrics-for-language-models/>
Radford et al. Improving Language Understanding by Generative Pre-Training . 2018

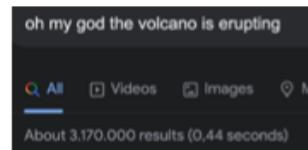
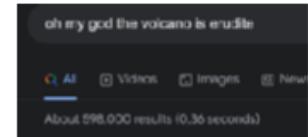
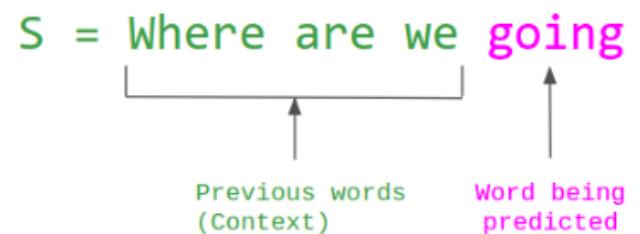
Language Modelling

Why "erudite" is not a good guess



Factor the probability of a datapoint (w_1, \dots, w_n):

$$P(w_1, w_2, \dots, w_n) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2)\dots p(w_n|w_1, w_2, \dots, w_{n-1}) \\ = \prod_{i=1}^n p(w_i|w_1, \dots, w_{i-1})$$



$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

Slide credit: Y.Asano

<https://thegradient.pub/understanding-evaluation-metrics-for-language-models/>
Radford et al. Improving Language Understanding by Generative Pre-Training . 2018

Language Modelling: GPT

Generative Pretrained Transformer (aka GPT)
simply does language modelling with a Transformer (decoder)

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

Language Modelling: GPT

Generative Pretrained Transformer (aka GPT)
simply does language modelling with a Transformer (decoder)

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

$$h_0 = UW_e + W_p$$

$$h_l = \text{transformer_block}(h_{l-1}) \forall i \in [1, n]$$

$$P(u) = \text{softmax}(h_n W_e^T)$$

$U = (u_k, \dots, u_1)$ is the context vector of tokens,
 n is the number of layers,
 W_e is the token embedding matrix,
 W_p is the position embedding matrix

in practice: "causal" (left-to-right) context via masking

Language Modelling: GPT

GPT-1,2,3: same loss. different training data and model sizes

GPT-1

117 million parameters

1.2 GB sized training dataset

Supervised finetuning afterwards

No release

<https://towardsdatascience.com/how-chatgpt-works-the-models-behind-the-bot-1ce5fca96286>

Language Modelling: GPT

GPT-1,2,3: same loss. different training data and model sizes

GPT-1

117 million parameters

1.2 GB sized training dataset

Supervised finetuning afterwards

No release

GPT-2

1.5 billion parameters

40 GB text training dataset

Often fine-tuned to perform specific tasks

Smaller version of the model was released
to the public open source

<https://towardsdatascience.com/how-chatgpt-works-the-models-behind-the-bot-1ce5fca96286>

Language Modelling: GPT

GPT-1,2,3: same loss. different training data and model sizes

GPT-1

117 million parameters

1.2 GB sized training dataset

Supervised finetuning afterwards

No release

GPT-2

1.5 billion parameters

40 GB text training dataset

Often fine-tuned to perform specific tasks

Smaller version of the model was released to the public open source

GPT-3

176 billion parameters

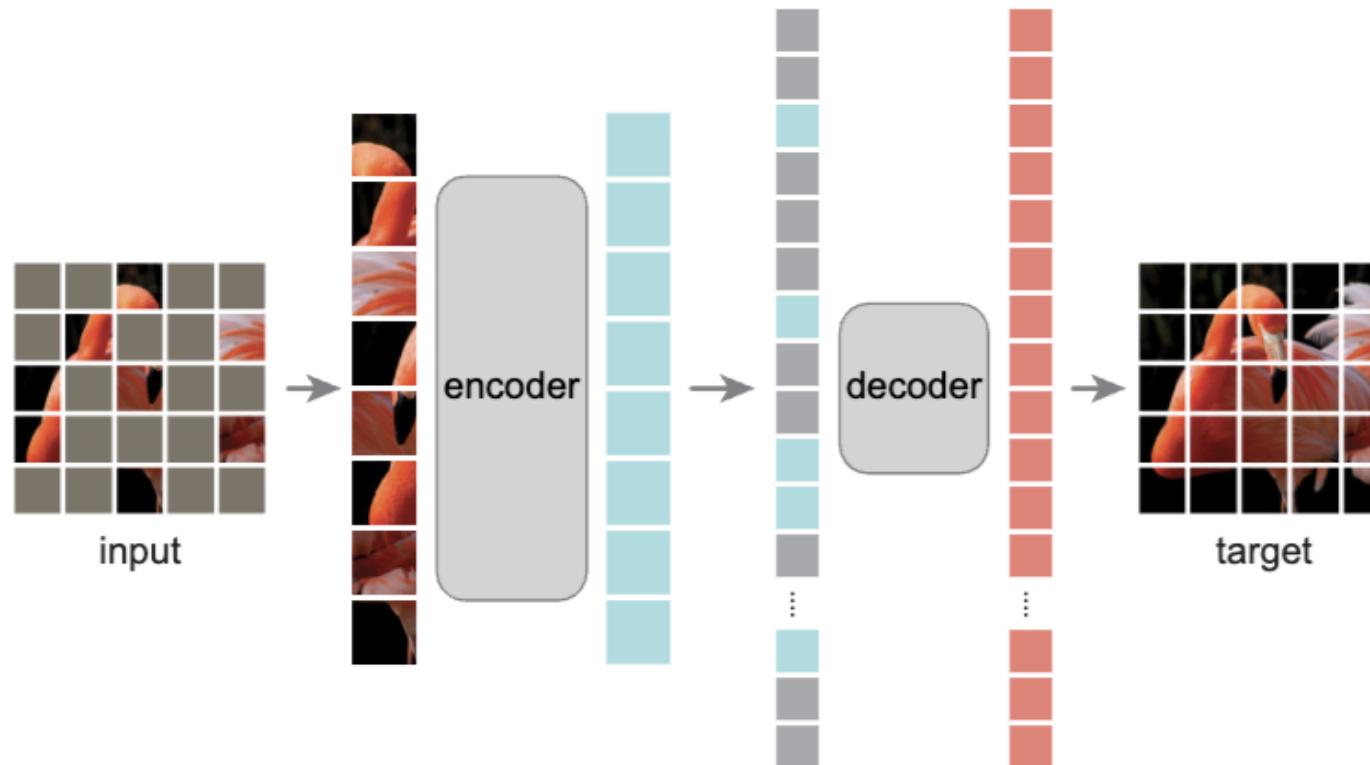
570 GB training dataset comprising of books, articles, websites, and more

Ability to perform most language tasks without additional tuning

Launched as an API service

Masked Image Modelling

Masked AutoEncoder (MAE) ~ same, but instead of words use image patches



He et al. Masked Autoencoders Are Scalable Vision Learners. CVPR'21
 Xie et al. SimMIM: A Simple Framework for Masked Image Modeling. ArXiv
 Dosovitskiy et al. An Image is Worth 16x16 Words: Transformers for Ima

Other examples

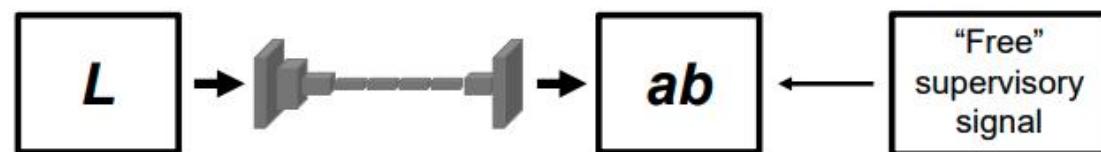
Self-supervised Learning Colorization

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)

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

Self-supervised Learning Colorization



Self-supervised Learning 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



Self-supervised Learning

Image Transformations

- Which image has the correct rotation?



Self-supervised Learning Image Transformations

- Which image has the correct rotation?



90° rotation



270° rotation



180° rotation



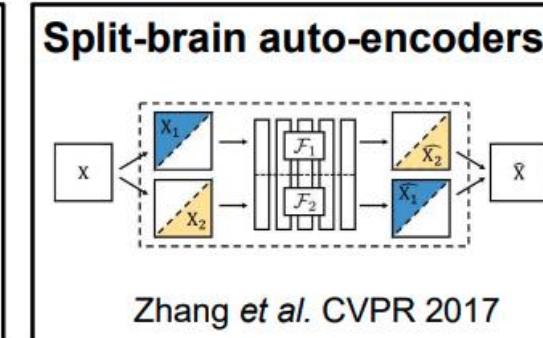
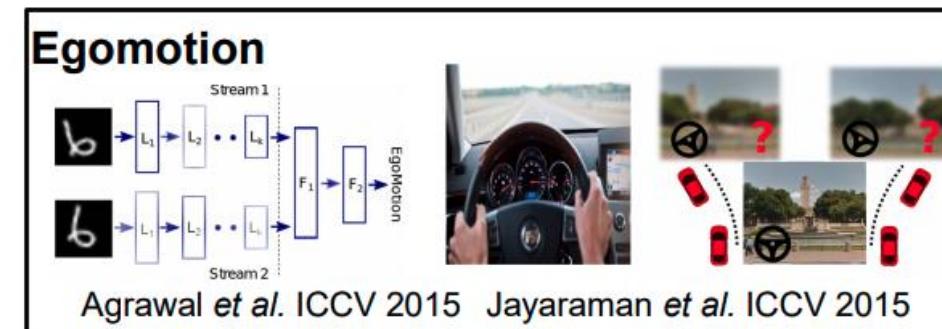
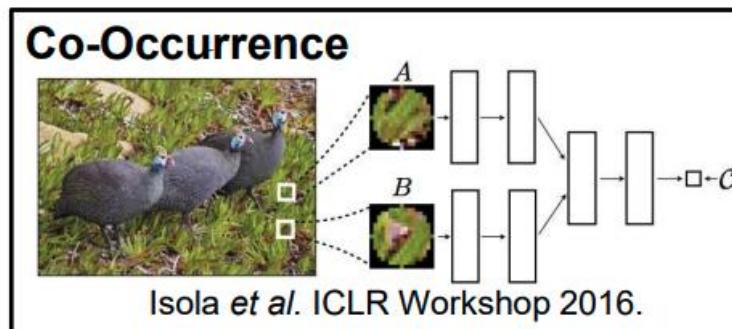
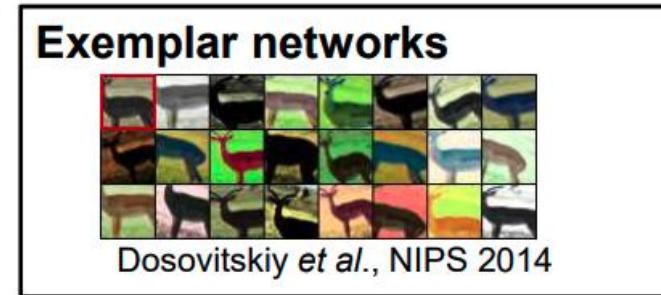
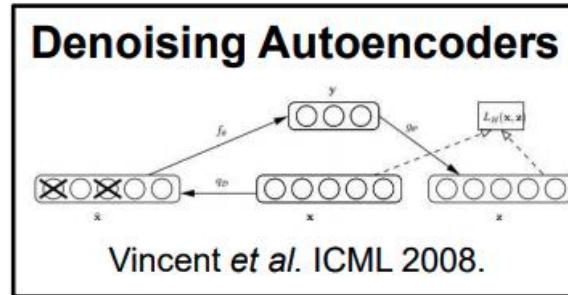
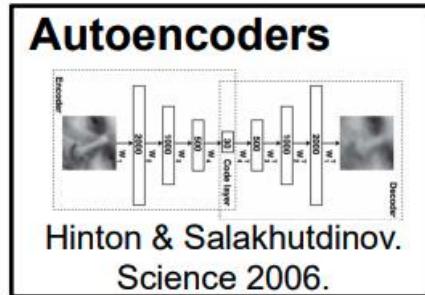
0° rotation



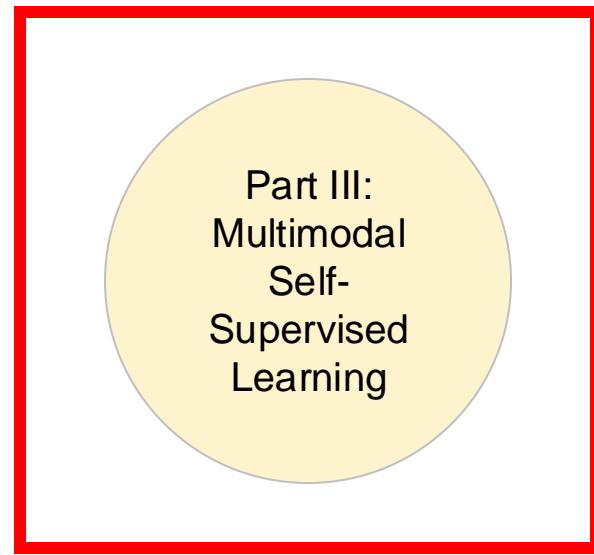
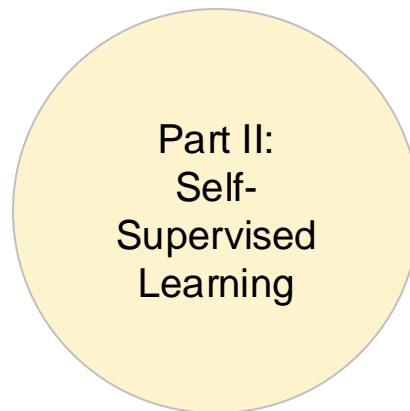
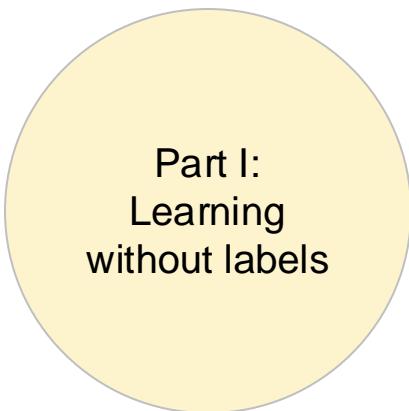
270° rotation

Figure 1: Images rotated by random multiples of 90 degrees (e.g., 0, 90, 180, or 270 degrees). The core intuition of our self-supervised feature learning approach is that if someone is not aware of the concepts of the objects depicted in the images, he cannot recognize the rotation that was applied to them.

Self-supervised Learning Images



Today's lecture

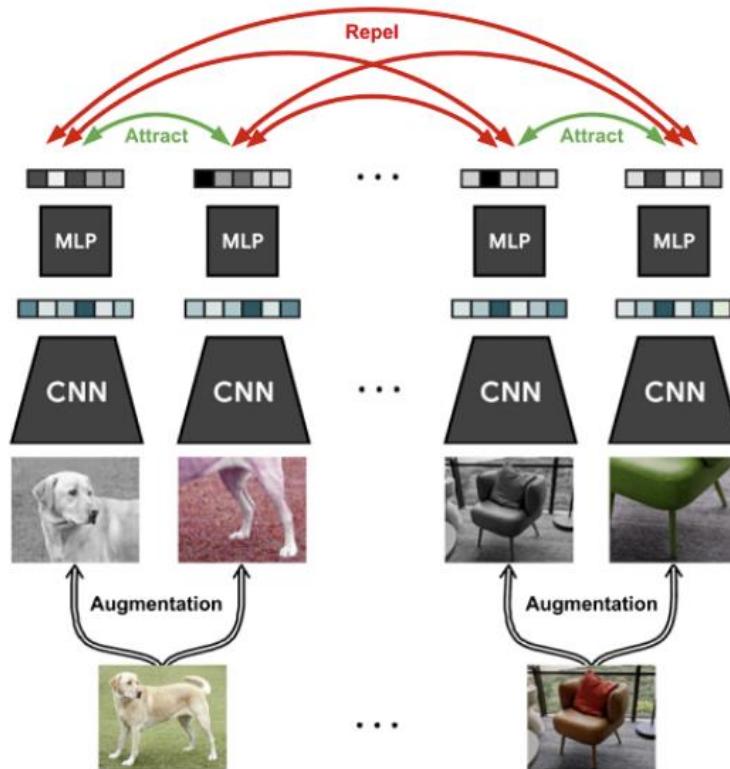


Part II: Outline Multimodal Self-Supervised Learning (SSL)

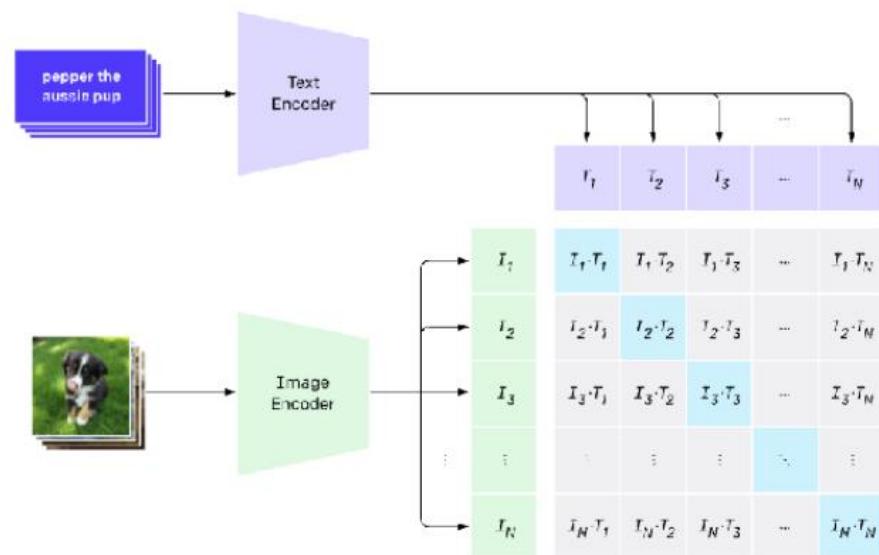


- CLIP
- Videos
- Sound

CLIP



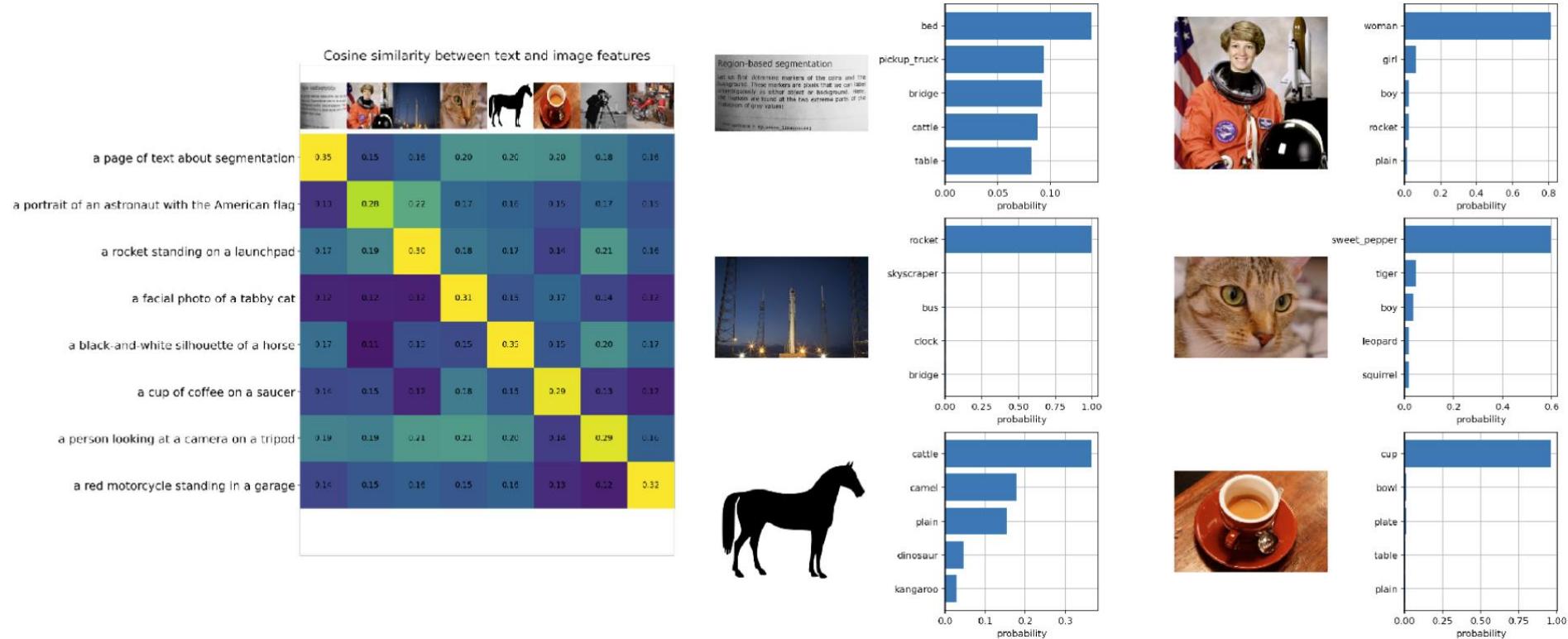
1. Contrastive pre-training



Radford et al. Learning Transferable Visual Models From Natural Language Supervision. ICML 2020
 but also: Desai & Johnson VirTex: Learning Visual Representations from Textual Annotations. CVPR 2021 etc.

CLIP

Zero-shot classification



https://colab.research.google.com/github/openai/clip/blob/master/notebooks/Interacting_with_CLIP.ipynb

Radford et al. Learning Transferable Visual Models From Natural Language Supervision. ICML 2020
but also: Desai & Johnson VirTex: Learning Visual Representations from Textual Annotations. CVPR 2021 etc.

When comparing pretrained image and language models, which one to adapt (more?)

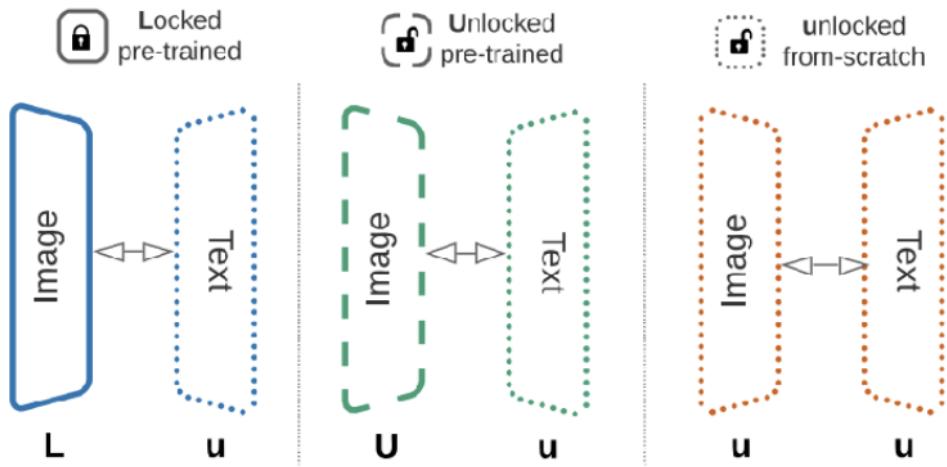


Figure 2. Design choices for contrastive-tuning on image-text data. Two letters are introduced to represent the image tower and text tower setups. L stands for locked variables and initialized from a pre-trained model, U stands for unlocked and initialized from a pre-trained model, u stands for unlocked and randomly initialized. Lu is named as “Locked-image Tuning” (LiT).

Method	ImgNet	ImgNet-v2	Cifar100	Pets
Lu	70.1	61.7	70.9	88.1
Uu	57.2	50.2	62.1	74.8
uu	50.6	43.3	47.9	70.3

Locking the image model is better.

Table 3: Zero-shot transfer results on ImageNet (variants).

Model	IN	IN-v2	IN-R	IN-A	ObjNet	ReaL
CLIP	76.2	70.1	88.9	77.2	72.3	-
ALIGN	76.4	70.1	92.2	75.8	72.2	-
BASIC	85.7	80.6	95.7	85.6	78.9	-
CoCa	86.3	80.7	96.5	90.2	82.7	-
LiT-g/14	85.2	79.8	94.9	81.8	82.5	88.6
LiT-e/14	85.4	80.6	96.1	88.0	84.9	88.4
LiT-22B	85.9	80.9	96.0	90.1	87.6	88.6

With only requiring one forward pass for getting image embeddings, can combine this with using a 22B parameter ViT

ALIGN

```
<figure class="wp-block-image size-large"></figure>
```



"motorcycle front wheel"



"thumbnail for version as of 21
57 29 june 2010"



"file frankfurt airport
skyline 2017 05 jpg"



"file london barge race 2 jpg"



"moustache seamless
wallpaper design"



"st oswalds way and shops"

ALIGN

```
<figure class="wp-block-image size-large"></figure>
```



"motorcycle front wheel"



"thumbnail for version as of 21
57 29 june 2010"



"file frankfurt airport
skyline 2017 05 jpg"



"file london barge race 2 jpg"



"moustache seamless
wallpaper design"



"st oswalds way and shops"

Novelty:

Start with very noisy dataset and:

- Filter based on images:
 - remove small ones, remove ones with >1k captions/alt texts
- Filter based on text:
 - alt-text with >10 occurrences are removed(e.g. "1920x10280")
 - too short or too long, or too rare
- Result:
 - dataset size ~2B (CLIP: 400M)

ALIGN

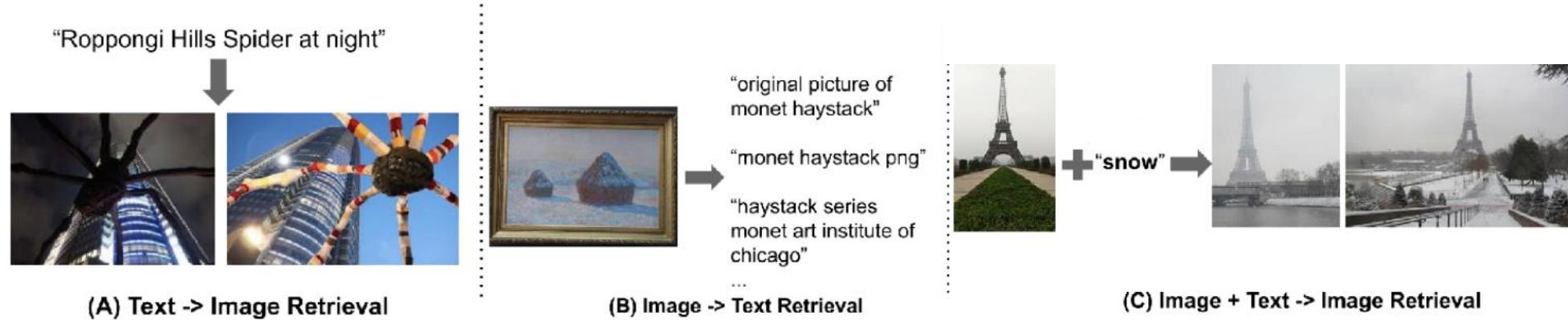


Image-text retrieval

Self-Supervised Learning from Videos

Self-supervised Learning Videos



- A temporal sequence of frames



Self-supervised Learning Videos

- A temporal sequence of frames



What can we use to define a proxy loss?

- Nearby (in time) frames are strongly correlated, further away may not be
- Temporal order of the frames
- Motion of objects (via optical flow)
- ...

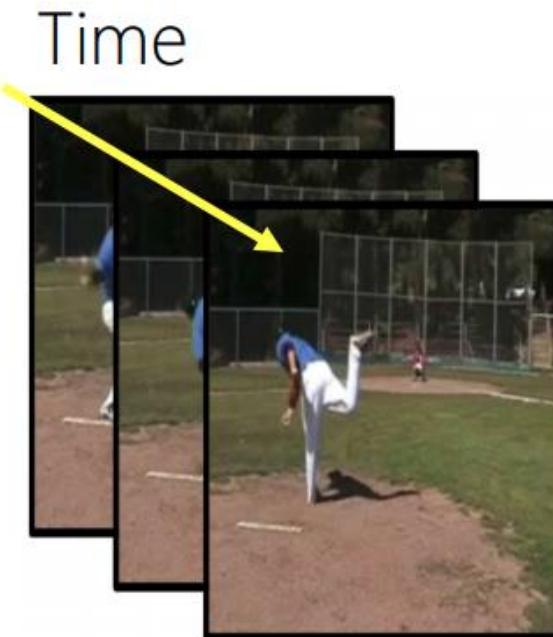
Temporal structure

Self-supervised Learning

Temporal structure in videos



- Shuffle and Learn: Unsupervised Learning using Temporal Order Verification



Self-supervised Learning

Temporal structure in videos

- Is this a valid sequence?



Self-supervised Learning

Temporal structure in videos

- Is this a valid sequence?



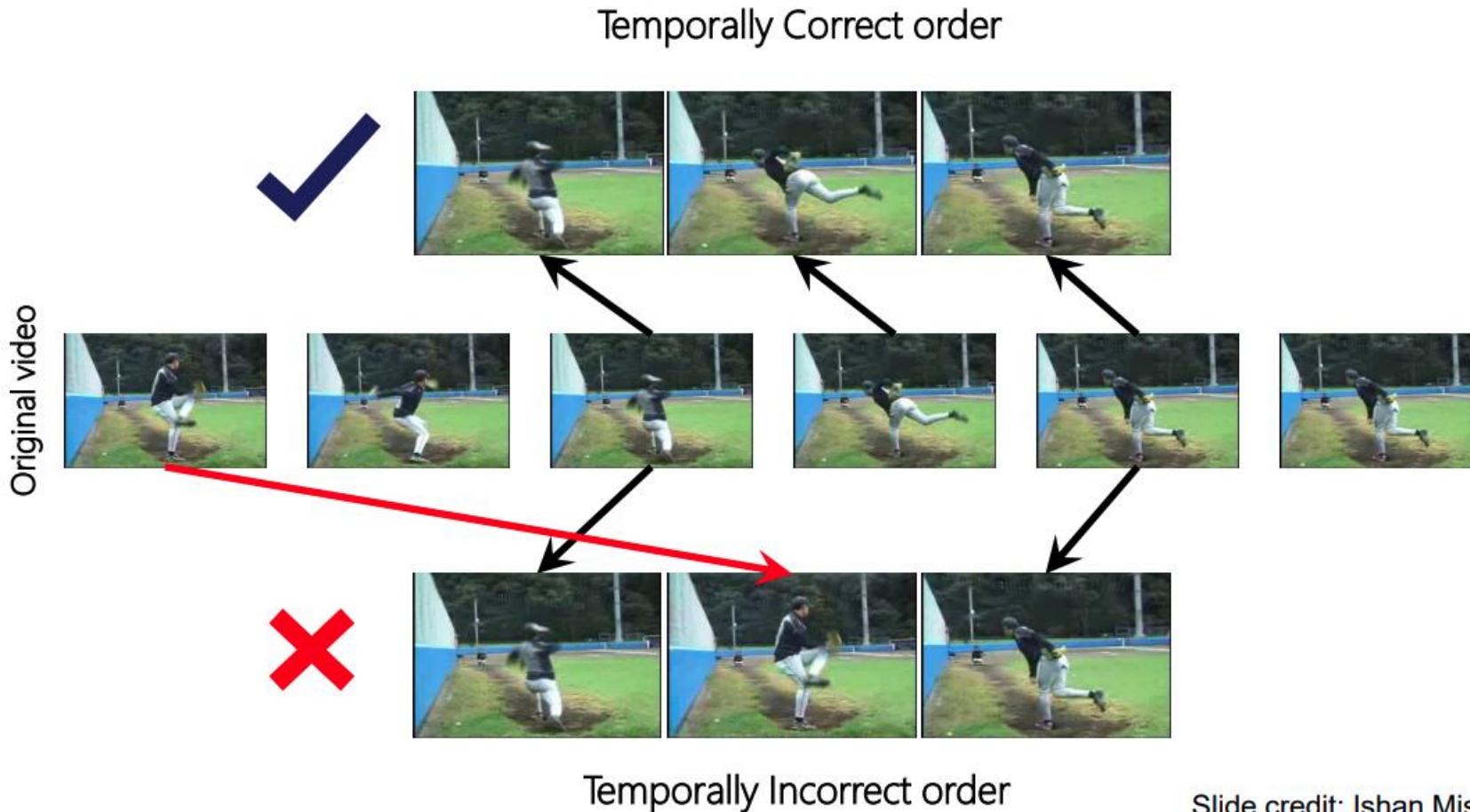
- Original video

Original video



Self-supervised Learning

Temporal structure in videos

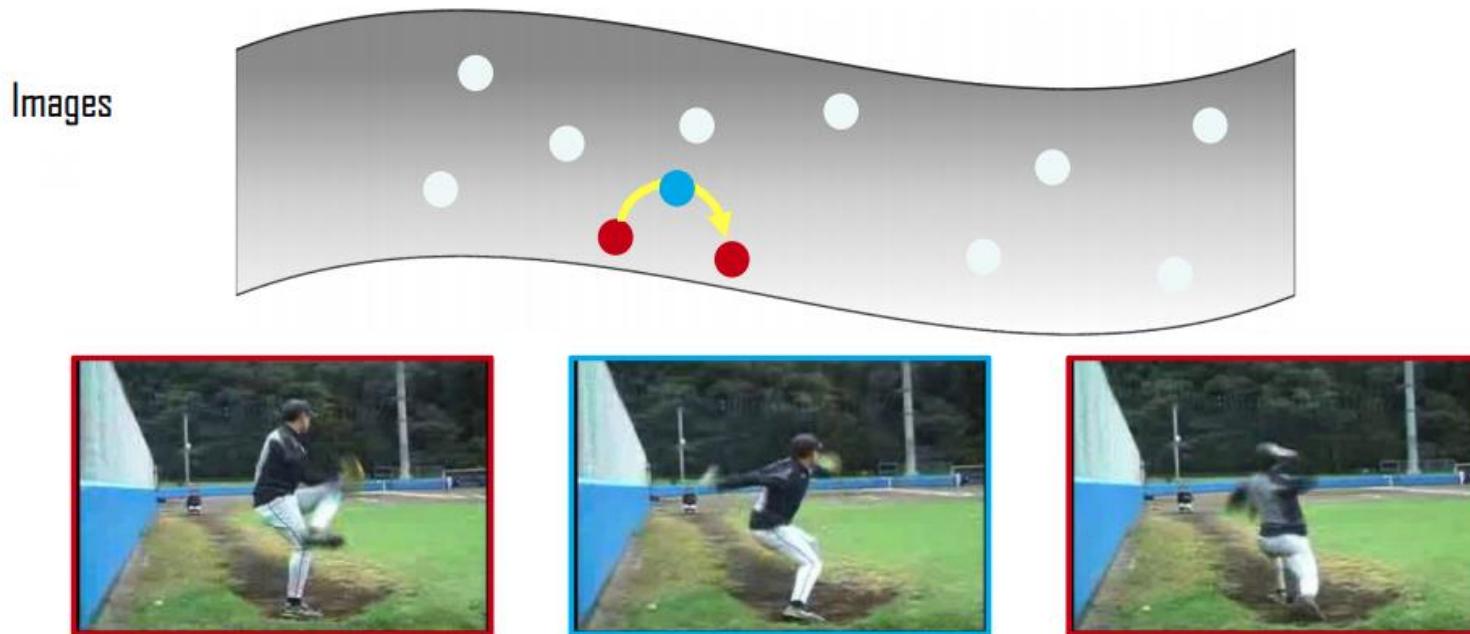


Slide credit: Ishan Misra

Self-supervised Learning

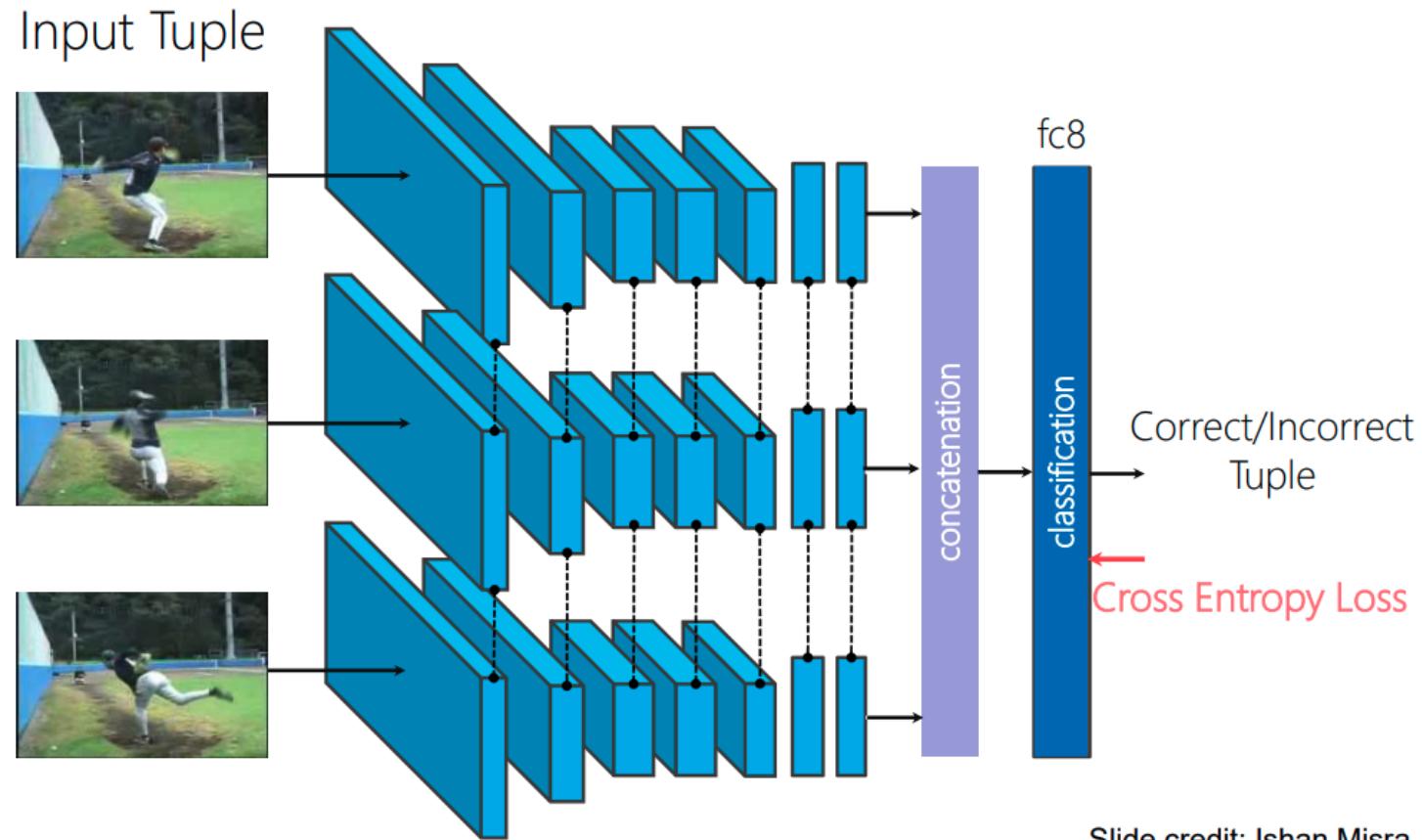
Temporal structure in videos

- Geometric View



Given a start and an end, can this point lie in between?

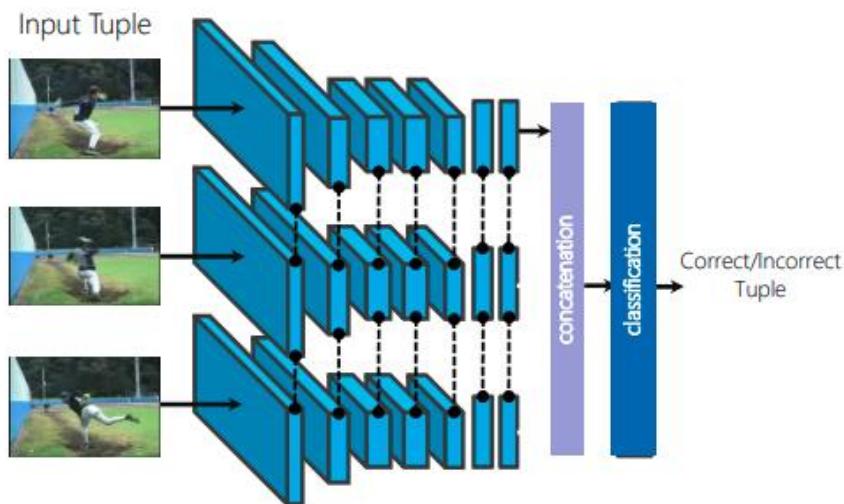
Self-supervised Learning Temporal structure in videos



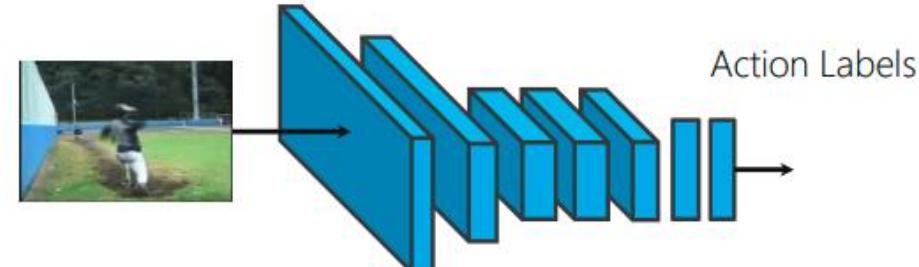
Self-supervised Learning Temporal structure in videos

- Fine-tuning setup

Self-supervised Pre-train



Test -> Finetune



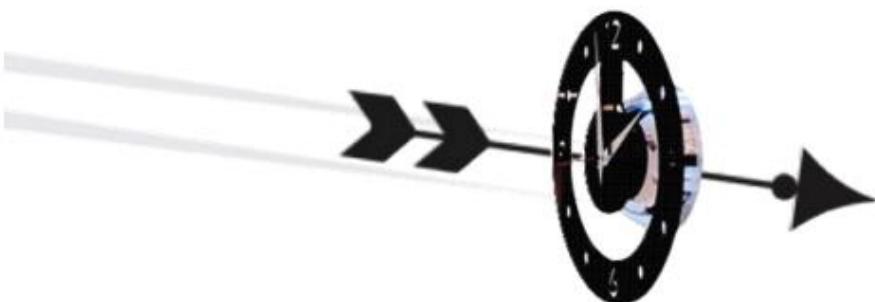
Learning the arrow of time

Self-supervised Learning

Learning the arrow of time



Task: predict if video playing forwards or backwards



Supervision:

Positive training samples: video clips playing forwards

Negative training samples: video clips playing backwards

Self-supervised Learning

Learning the arrow of time



- Strong cues

Semantic, face motion direction, ordering



Self-supervised Learning

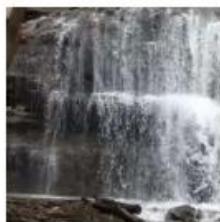
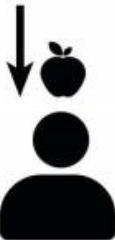
Learning the arrow of time



- Strong cues

'Simple' physics:

- gravity
- entropy
- friction
- causality



Self-supervised Learning

Learning the arrow of time

- Weak or no cues

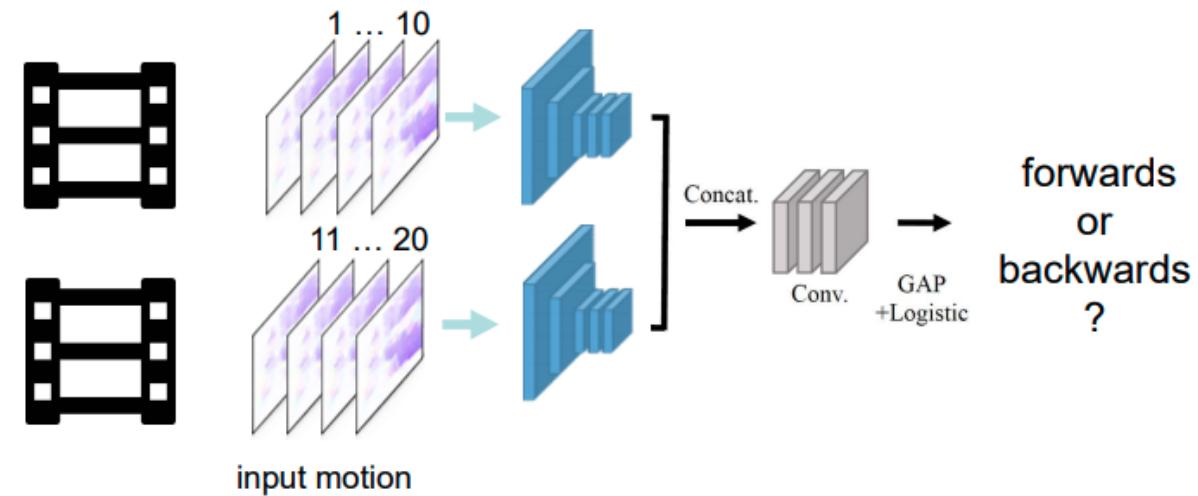
Symmetric in time, constant motion, repetitions



Self-supervised Learning

Learning the arrow of time

- Temporal class-activation Map Network



T-CAM Model:

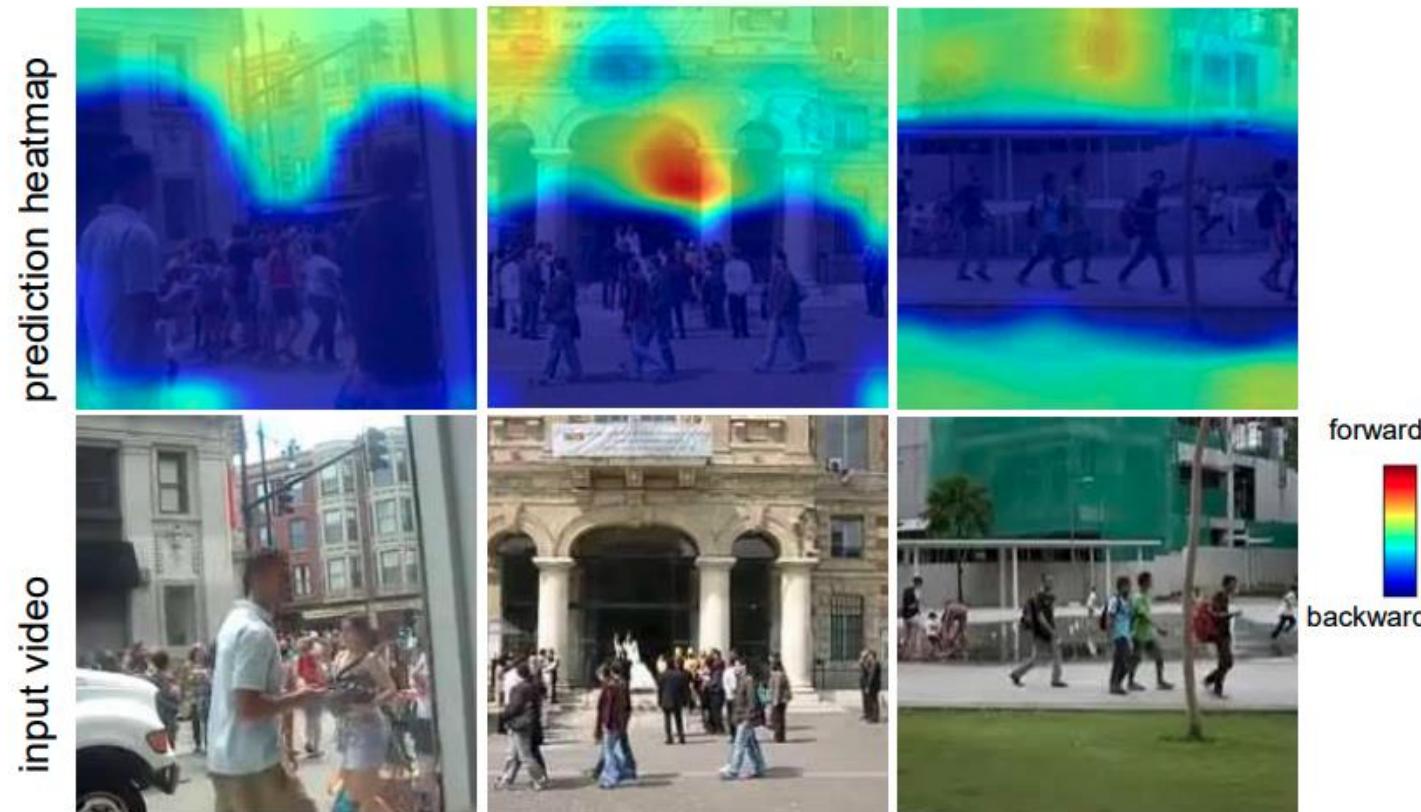
Input: optical flow in two chunks

Final layer: global average pooling to allow class activation map (CAM)

Self-supervised Learning

Learning the arrow of time

- Semantic motions

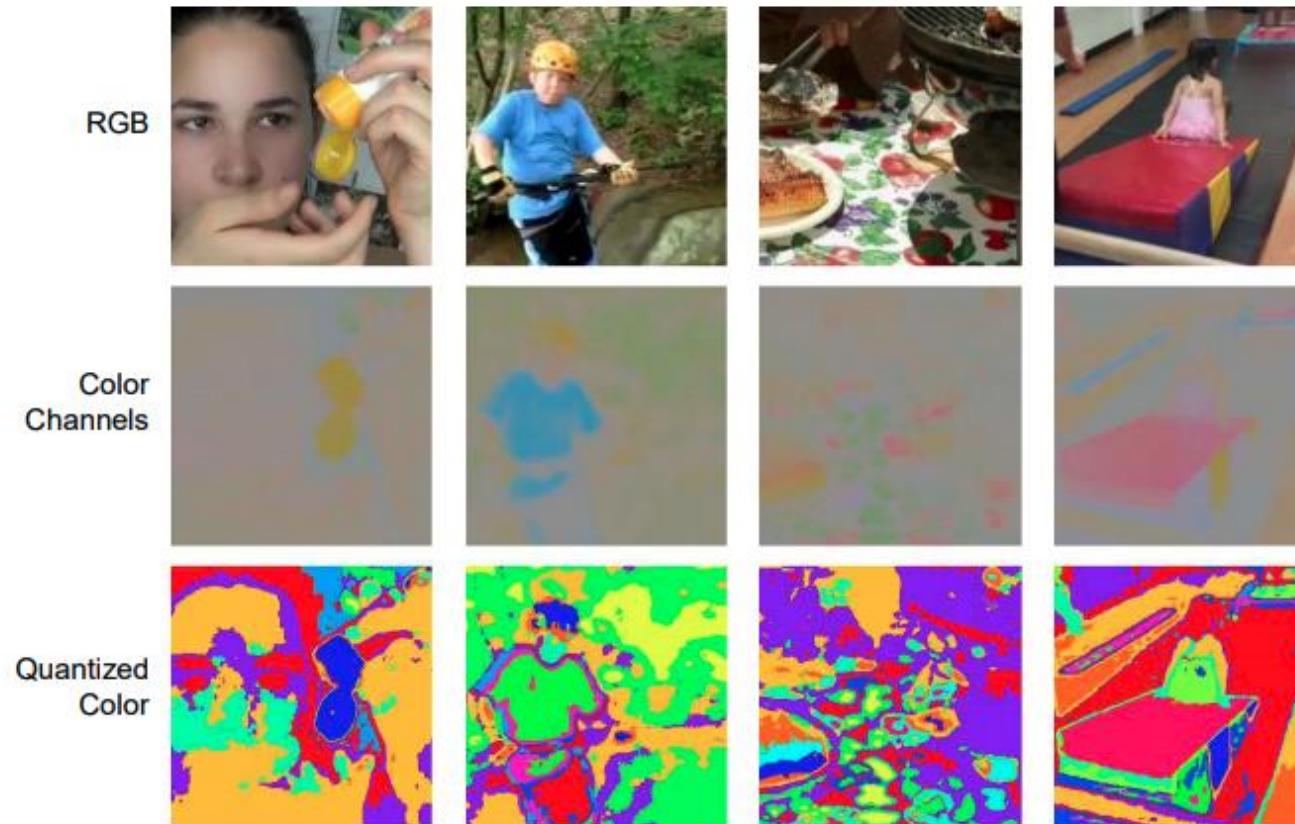


Tracking emerges by colorization

Self-supervised Learning

Tracking emerges by colorization

- Temporal coherence of color



Self-supervised Learning

Tracking emerges by colorization



- Self-supervised tracking

Task: given a color video ...

Colorize all frames of a gray scale version using a reference frame



Reference Frame



Gray-scale Video

Self-supervised Learning

Tracking emerges by colorization

- What is this color?



Self-supervised Learning

Tracking emerges by colorization

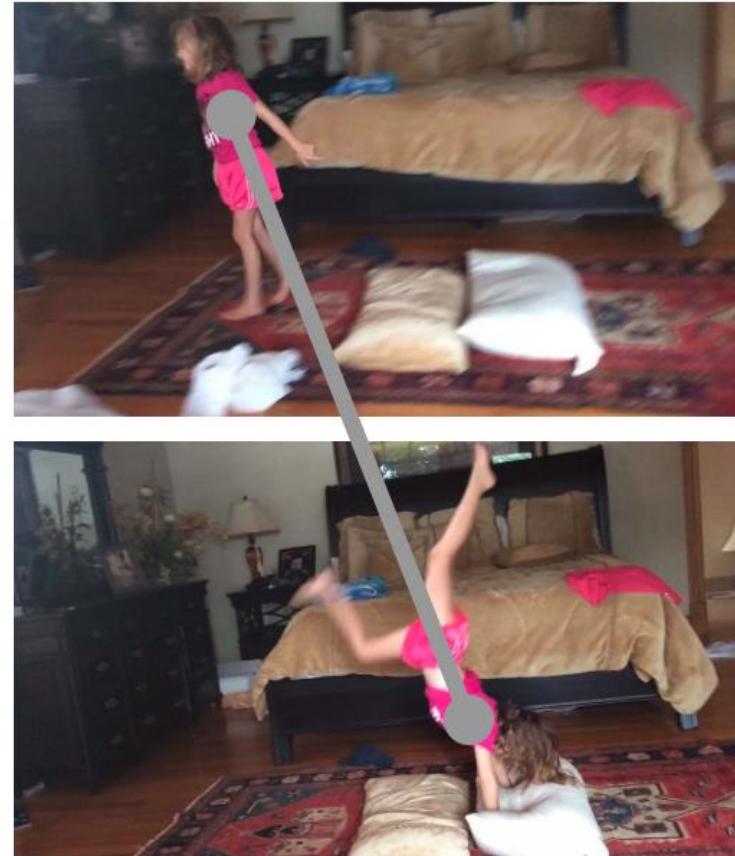
- Where to copy color from?



Self-supervised Learning

Tracking emerges by colorization

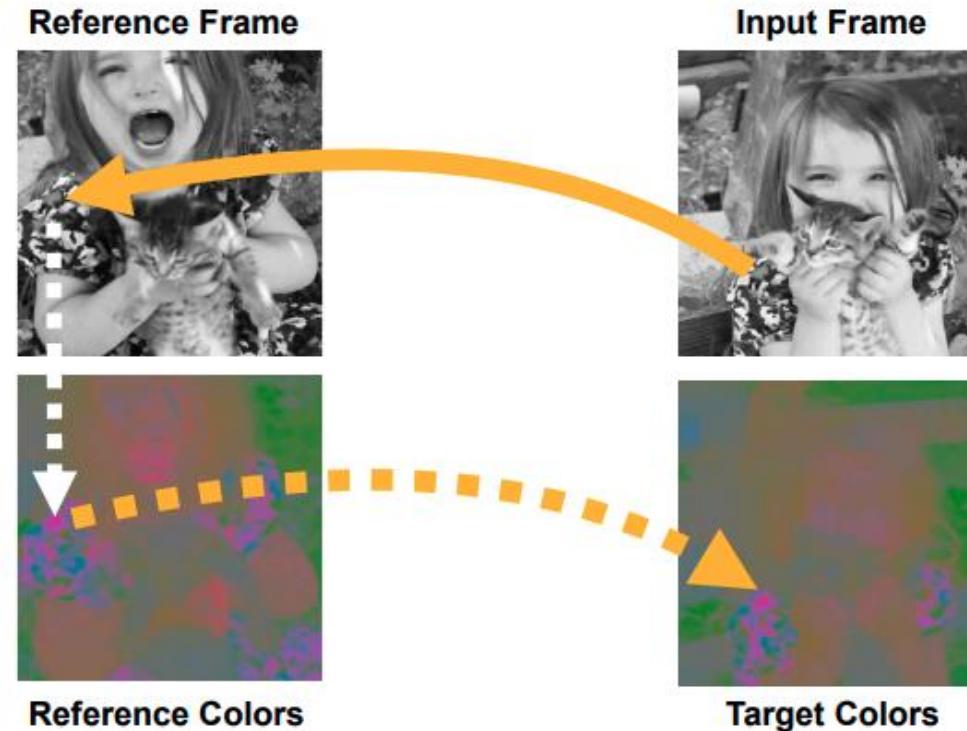
- Semantic correspondence



Self-supervised Learning

Tracking emerges by colorization

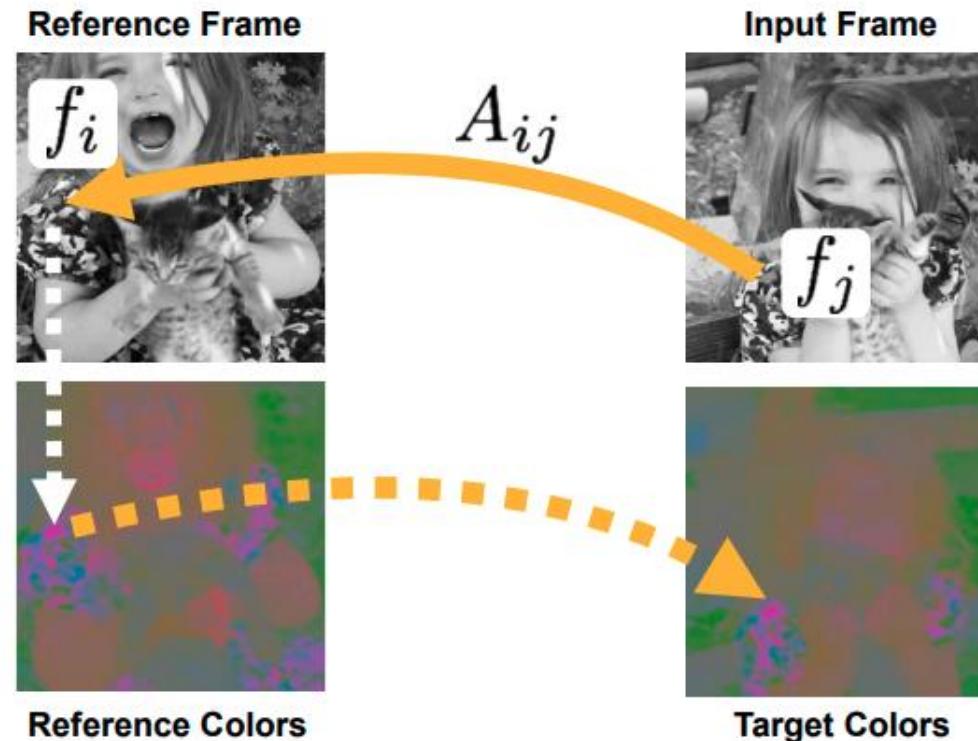
- Colorize by pointing



Self-supervised Learning

Tracking emerges by colorization

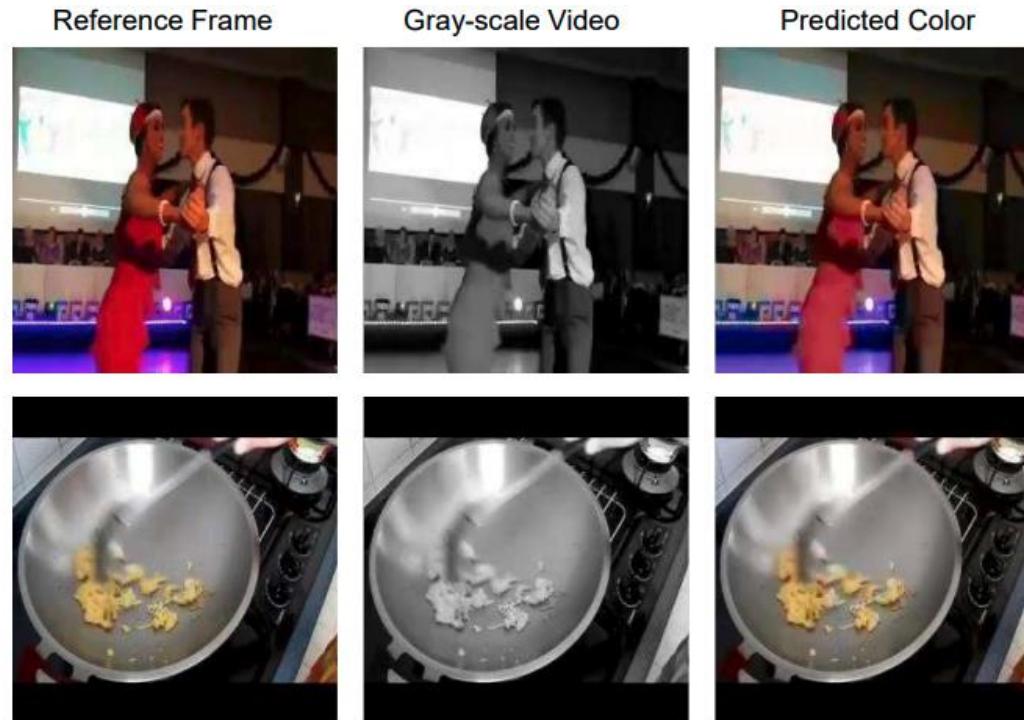
- Colorize by pointing



Self-supervised Learning

Tracking emerges by colorization

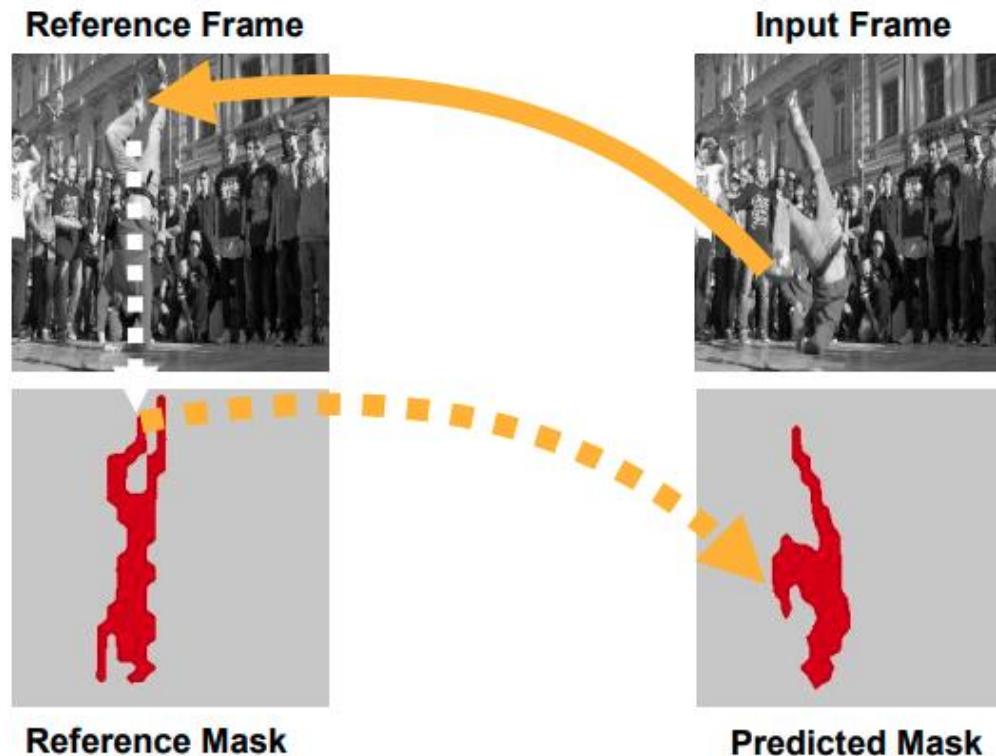
- Video colorization



Self-supervised Learning

Tracking emerges by colorization

- Tracking emerges!



Self-supervised Learning

Tracking emerges by colorization

- Segment tracking results

Only the first frame is given. Colors indicate different instances.



Videos with sound

Self-supervised Learning

Audio-Visual Co-supervision



- Sound and frames are:
 - Semantically consistent
 - Synchronized



Self-supervised Learning Audio-Visual Co-supervision



- Goal: use vision and sound to learn from each other
- Two types of proxy task:
 - Predict audio-visual correspondence
 - Predict audio-visual synchronization



Self-supervised Learning

Audio-Visual Co-supervision



- Train a network to predict if image and audio clip correspond

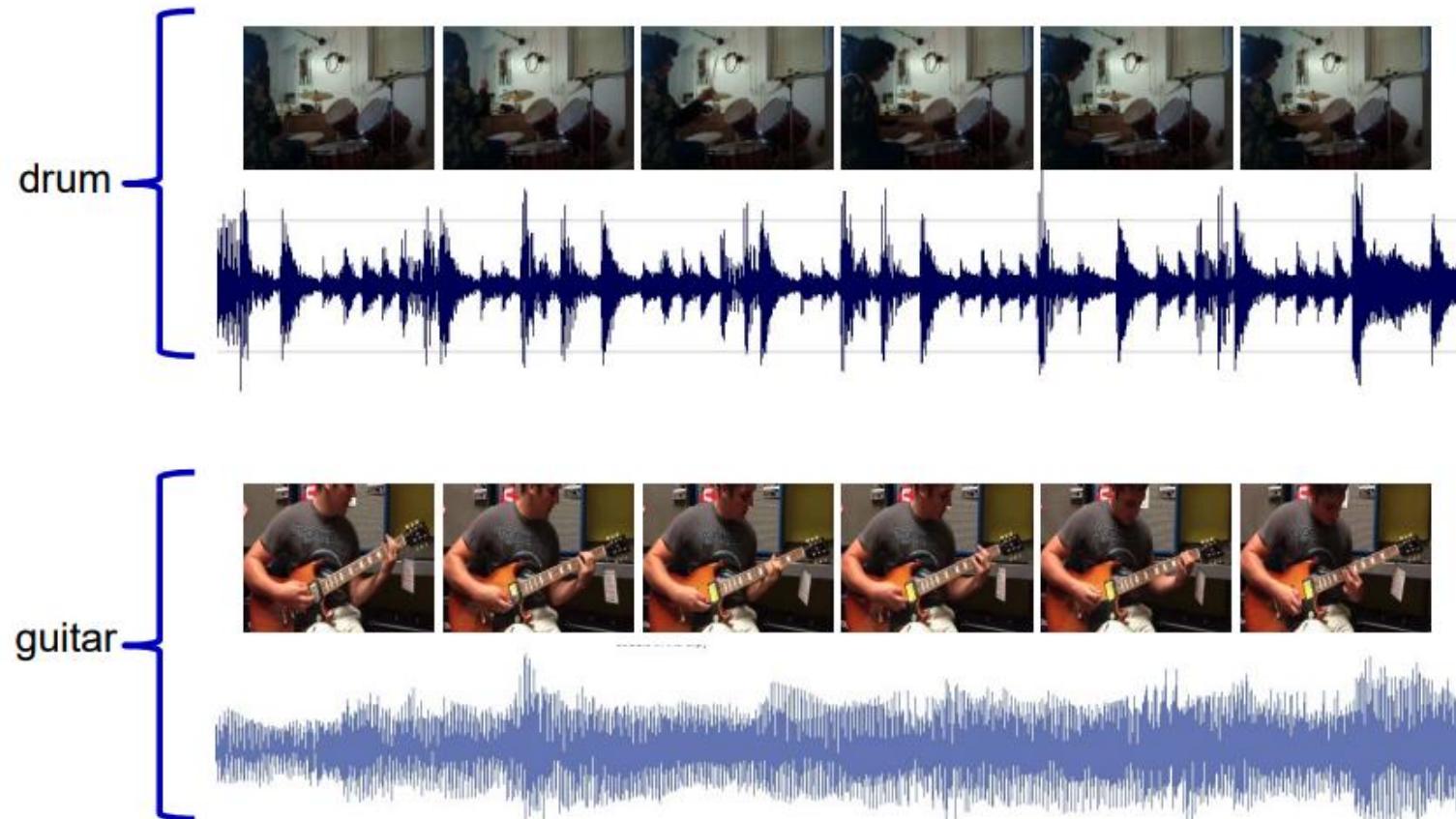


Correspond?



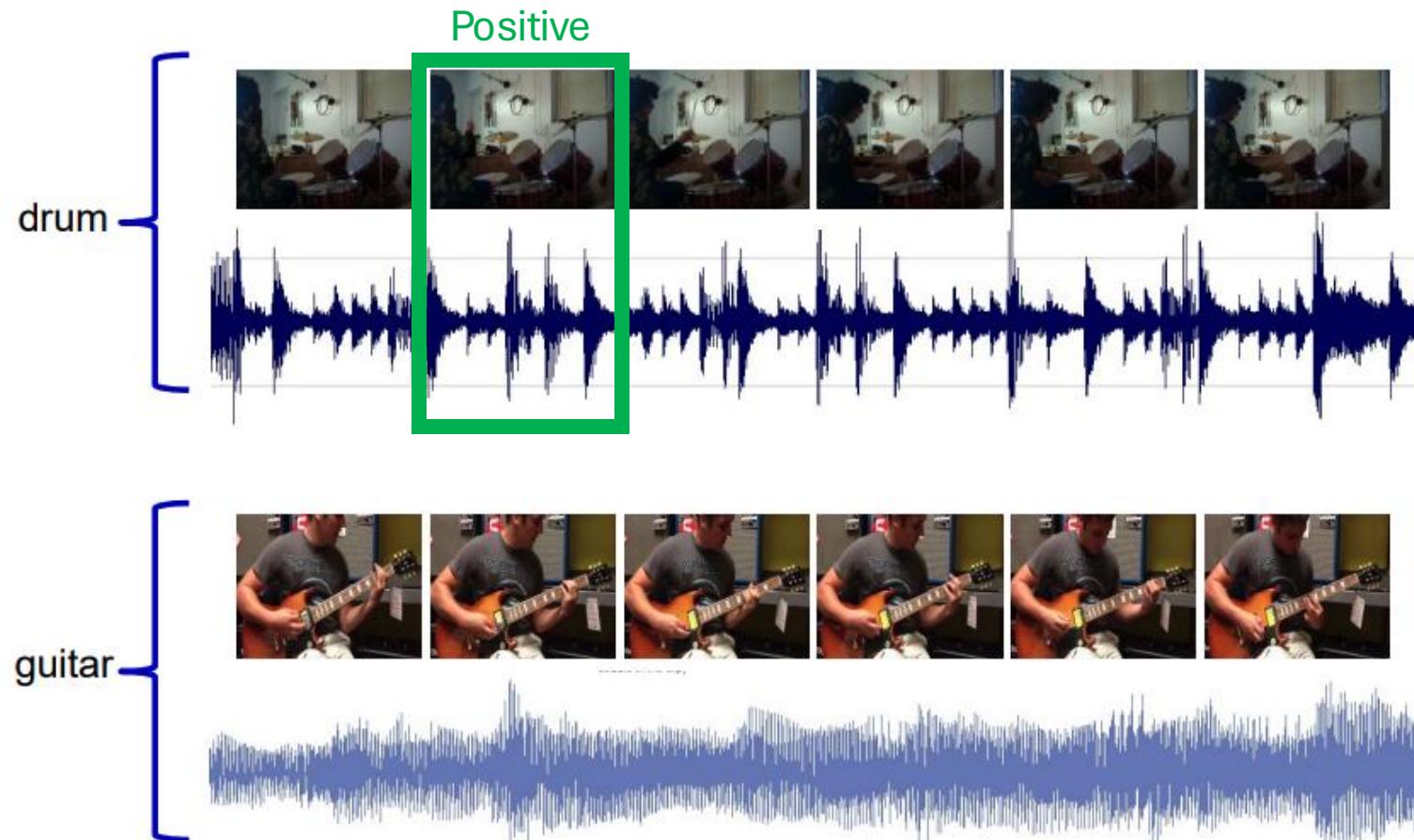
Self-supervised Learning Audio-Visual Co-supervision

- Audio-visual Correspondence



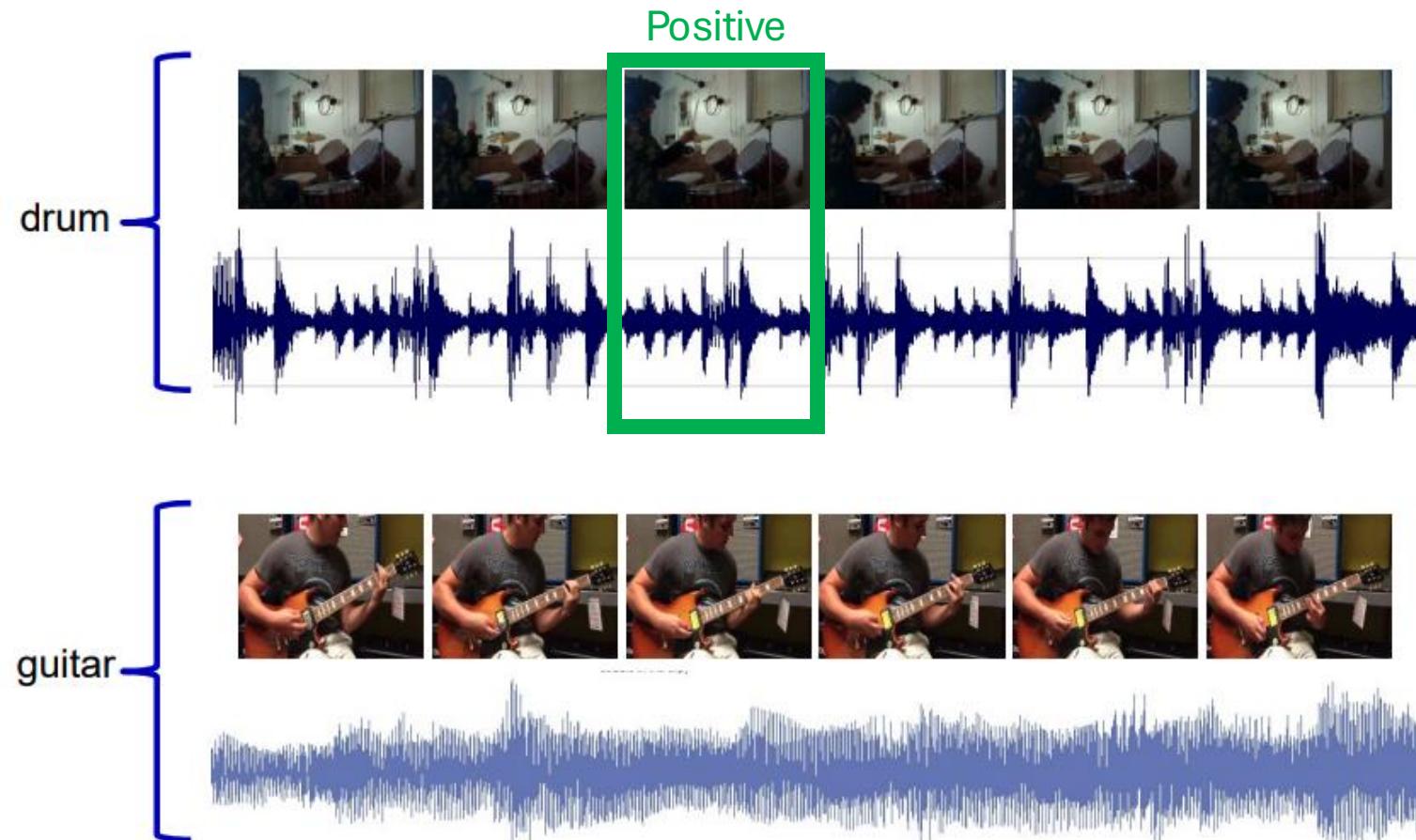
Self-supervised Learning Audio-Visual Co-supervision

- Audio-visual Correspondence



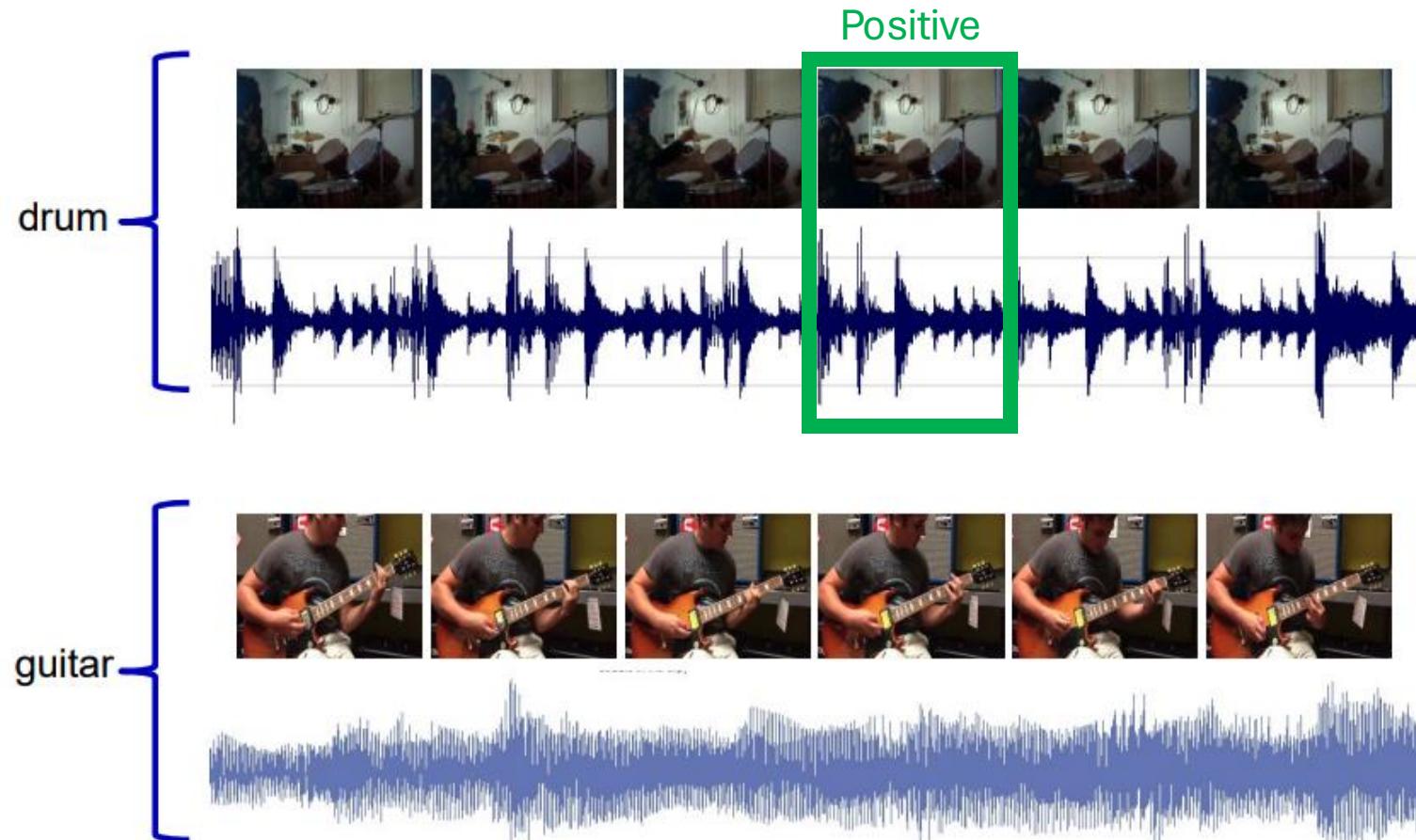
Self-supervised Learning Audio-Visual Co-supervision

- Audio-visual Correspondence



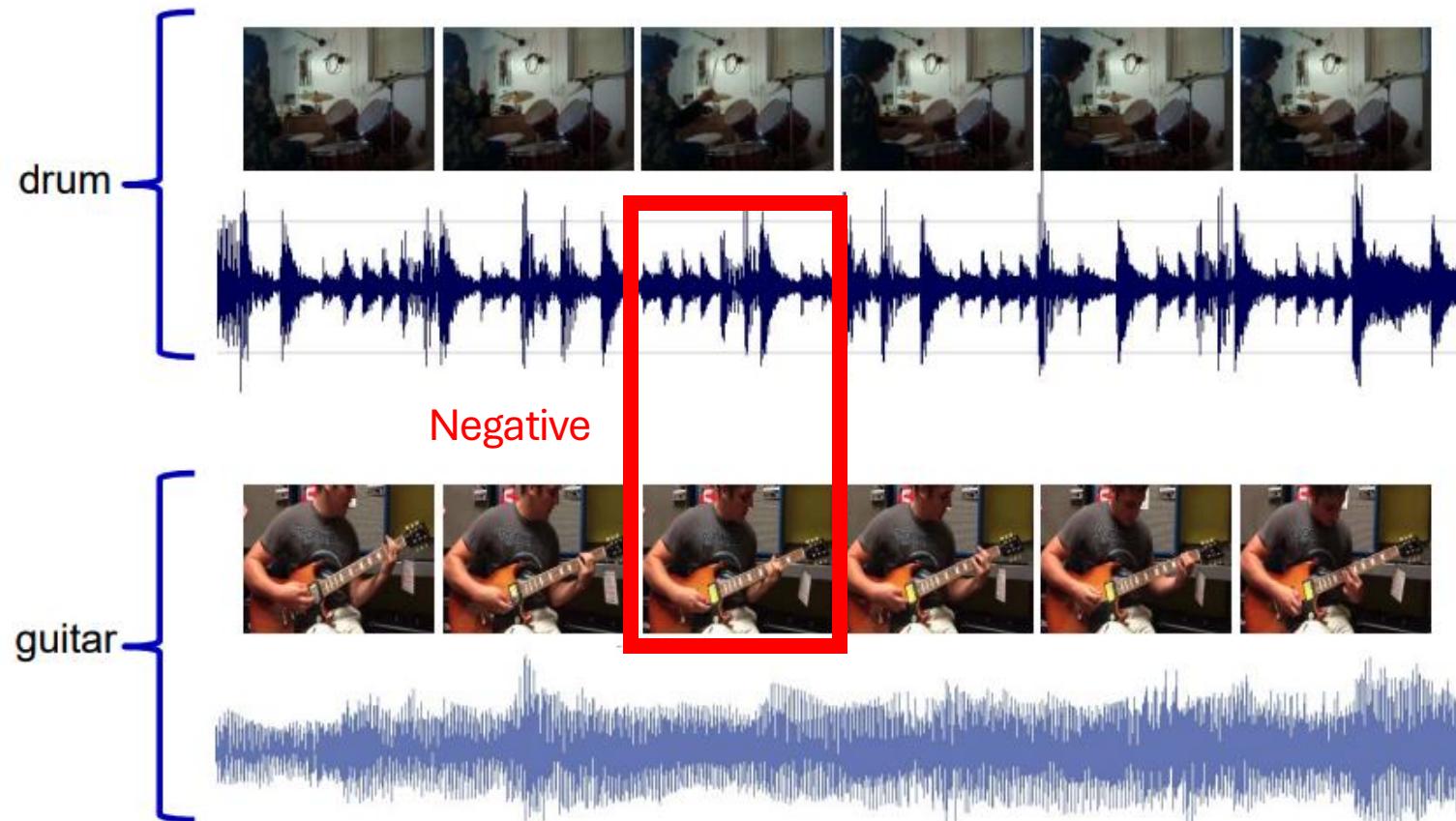
Self-supervised Learning Audio-Visual Co-supervision

- Audio-visual Correspondence



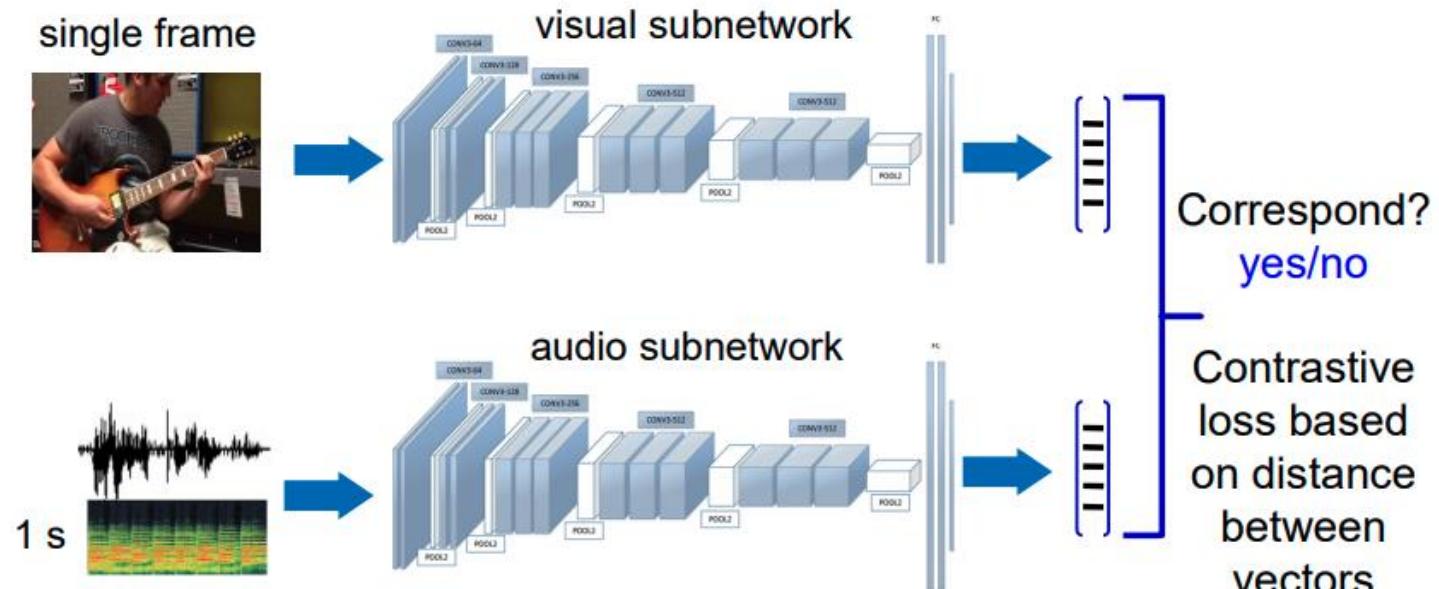
Self-supervised Learning Audio-Visual Co-supervision

- Audio-visual Correspondence



Self-supervised Learning Audio-Visual Co-supervision

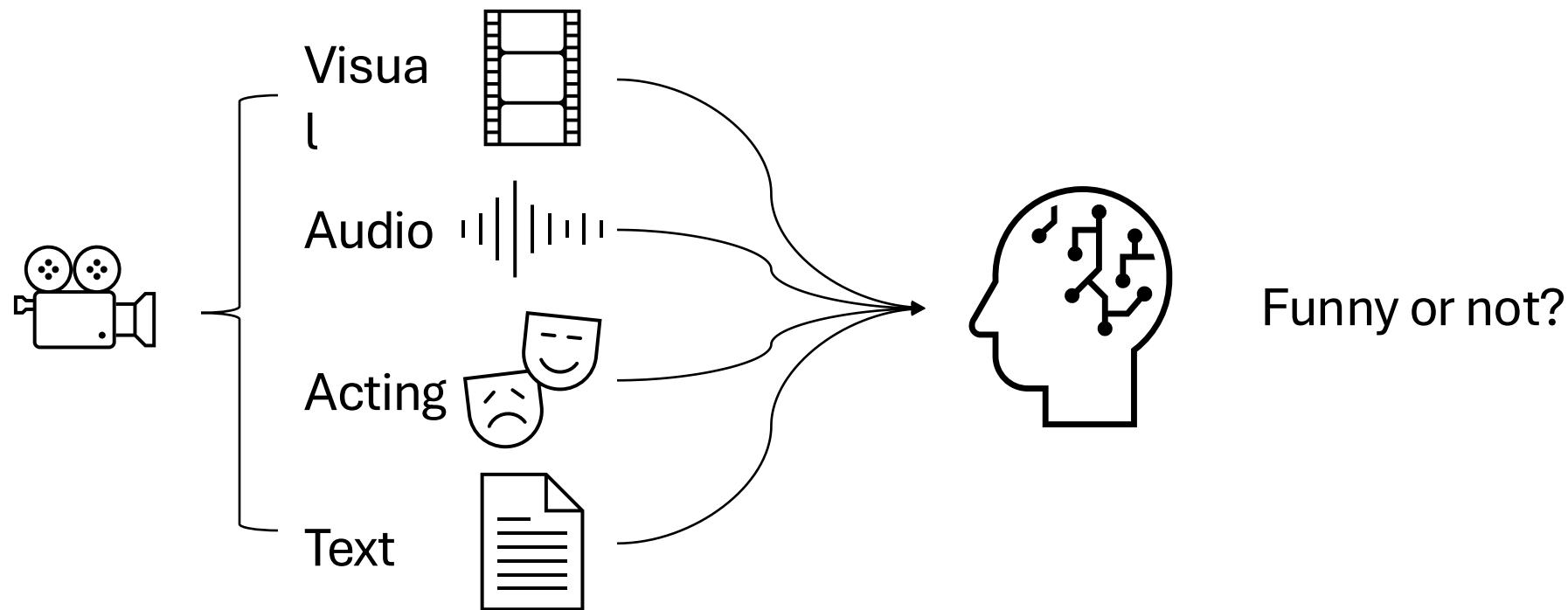
- Audio-visual Embedding (AVE-Net)



Distance between audio and visual vectors:

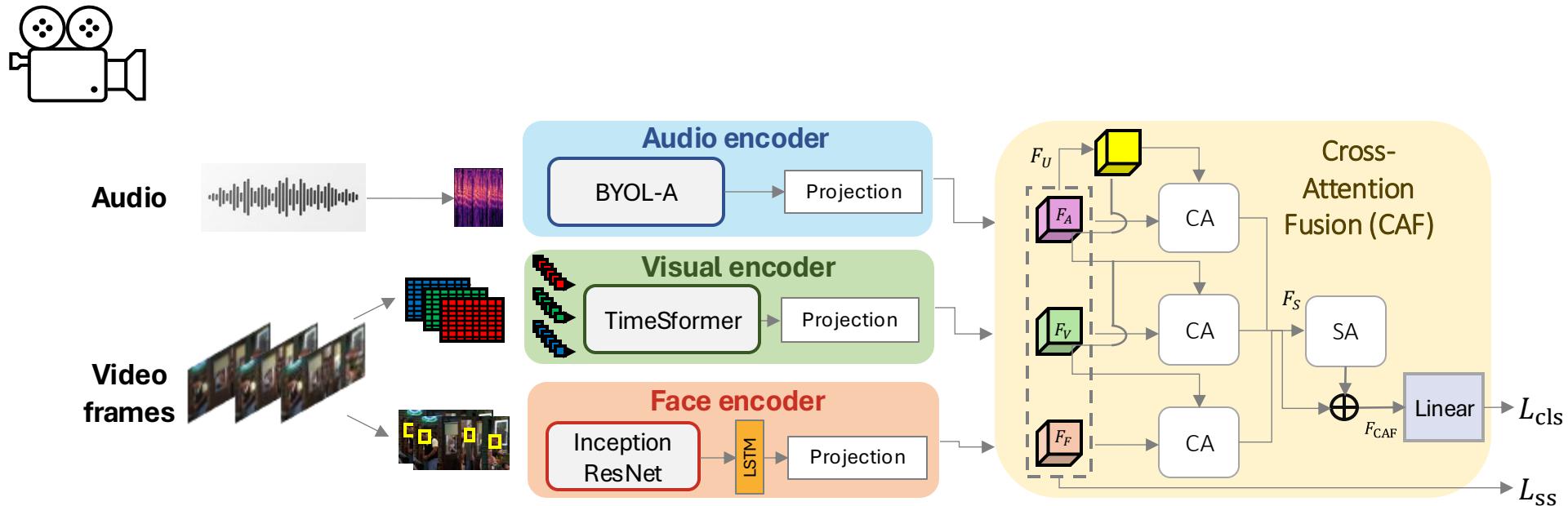
- Small: AV from the same place in a video (Positives)
- Large: AV from different videos (Negatives)

FunnyNet: Audiovisual learning of funny moments in videos



[Liu et al, ACCV 2022, IJCV 2024]

FunnyNet



FunnyNet-W

Sitcom with Canned Laughter



Examples of well classified
funny moments

FunnyNet-W

Sitcom without Canned Laughter



Examples of well classified
funny moments

Self-supervised learning categories

Thank you