# Deep Learning 1

2025-2026 – Pascal Mettes

## Lecture 8

*From supervised to unsupervised deep learning*

# Previous lecture

| Lecture | Title |
| --- | --- |
| 1 | Intro and history of deep learning |
| 3 | Deep learning optimization I |
| 5 | Convolutional deep learning |
| 7 | Graph deep learning |
| 9 | Multi-modal deep learning |
| 11 | What doesn't work in deep learning |
| 13 | Q&A |

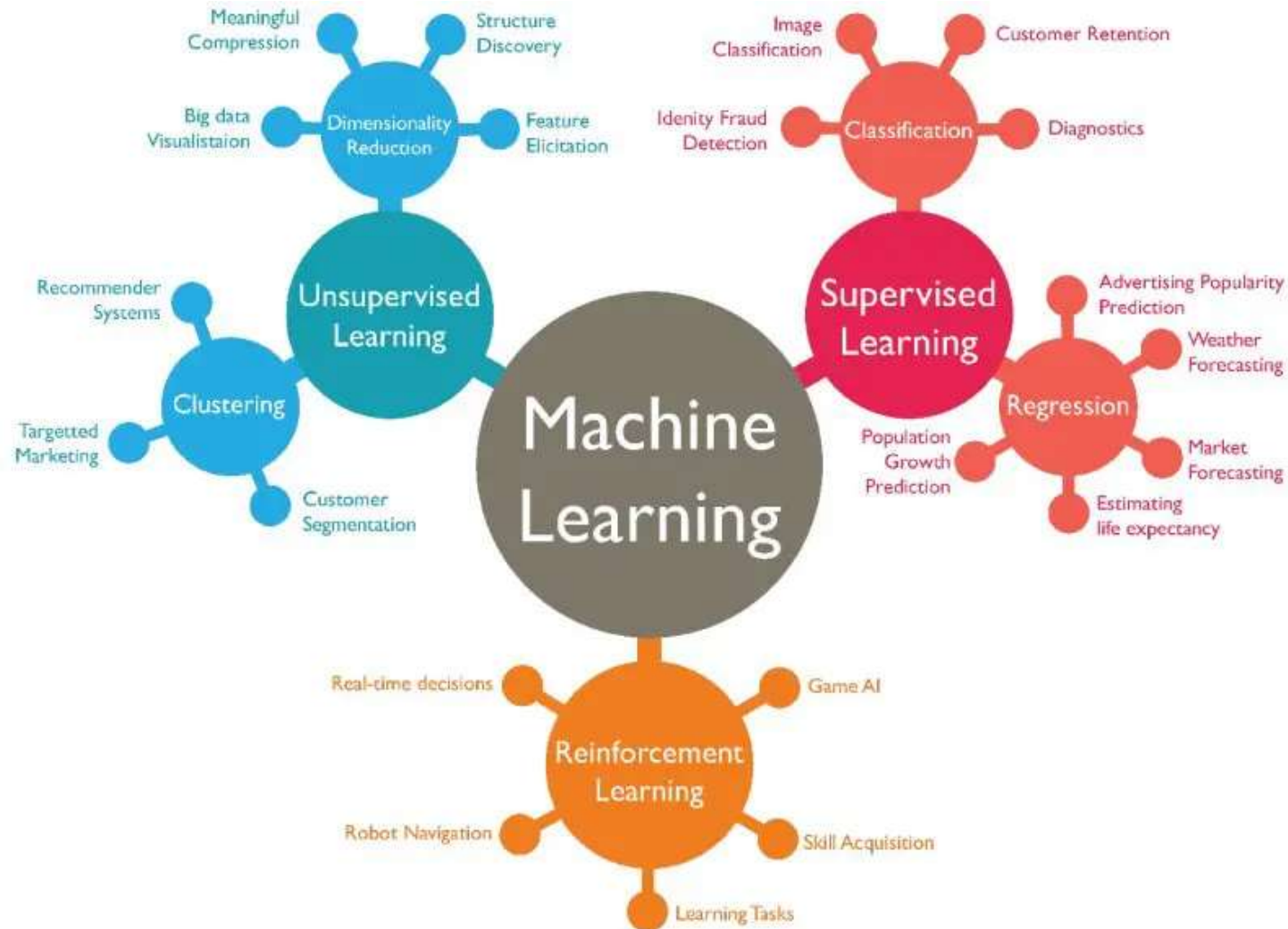| Lecture | Title |
| --- | --- |
| 2 | AutoDiff |
| 4 | Deep learning optimization II |
| 6 | Attention-based deep learning |
| 8 | From supervised to unsupervised deep learning |
| 10 | Generative deep learning |
| 12 | Non-Euclidean deep learning |
| 14 | Deep learning for videos |

# This lecture

Self-supervised learning for vision
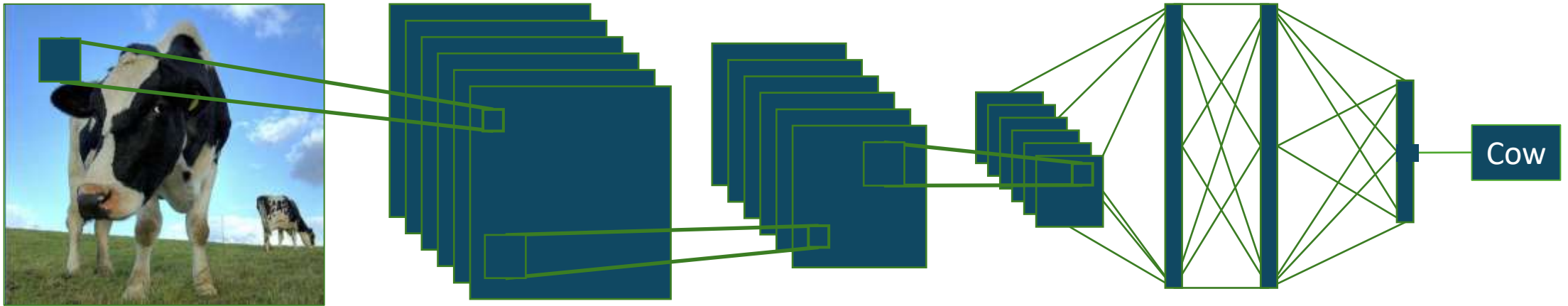
Self-supervised learning for language

In between supervised and self-supervised learning

# Traditional pillars of machine learning

# Strength and weakness of supervision in DL

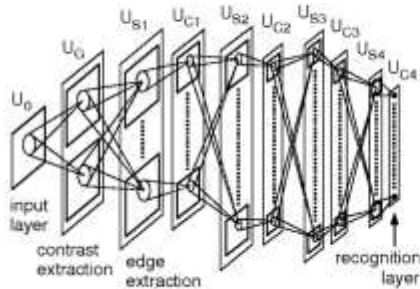Supervision makes it possible to propagate signals back to train networks.



Cow

Classification labels no longer the backbone of latest models, why?

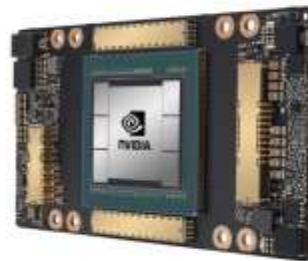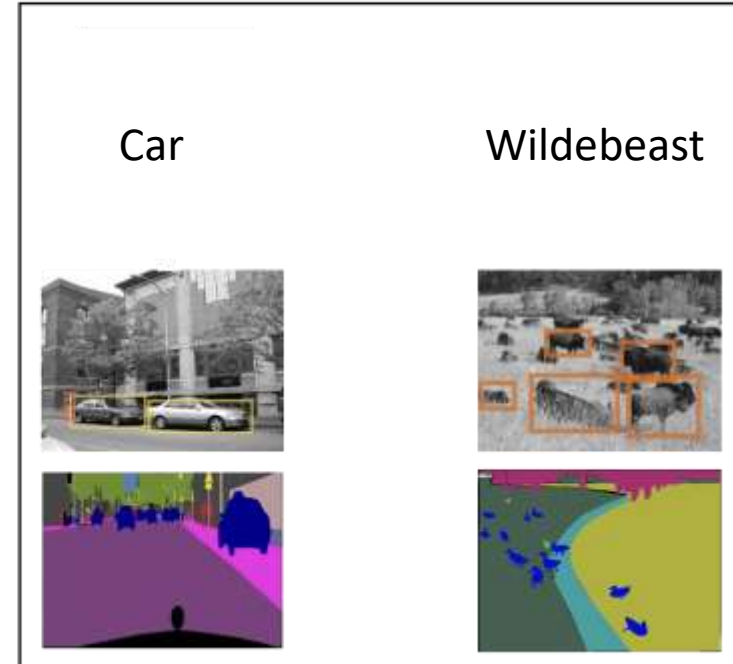# Self-supervised learning

# Data as fuel for deep learning

# The bottleneck of data



Images are often cheap

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

Supervised Learning

Car                Wildebeast

# The two stages of deep learning



Feature learning

Classifier learning

Cow

The final layer requires labels, but is that also true for all other layers?
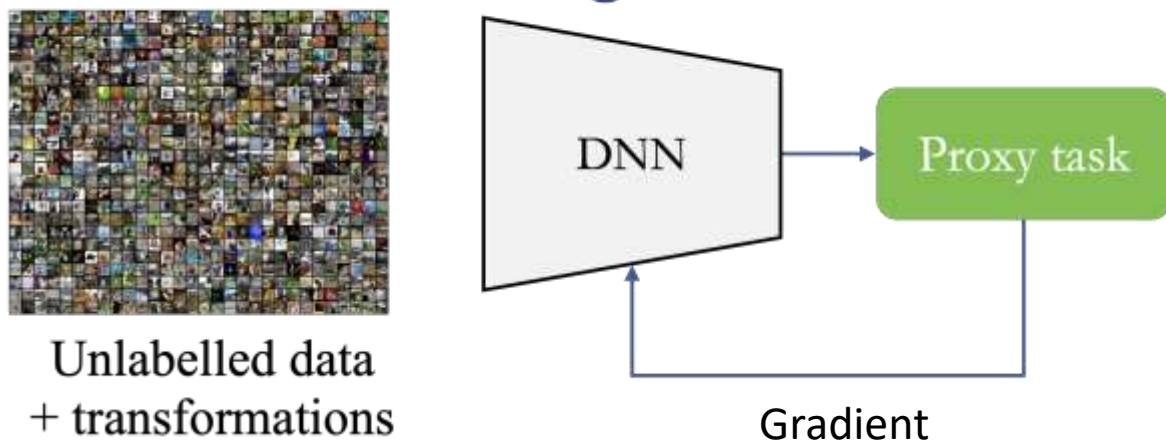
# Solving the problem of expensive annotations: self



Self-supervision

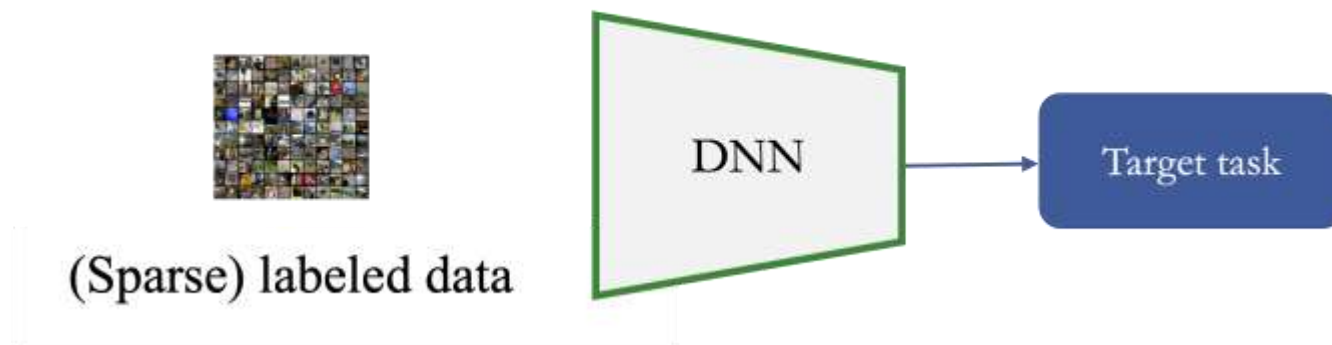Extract a supervisory signal from the raw data

# Main idea of self-supervised learning

**Phase 1: Pretraining**



Unlabelled data
+ transformations

DNN

Proxy task

Gradient

**Phase 2: Downstream tasks**

(Sparse) labeled data
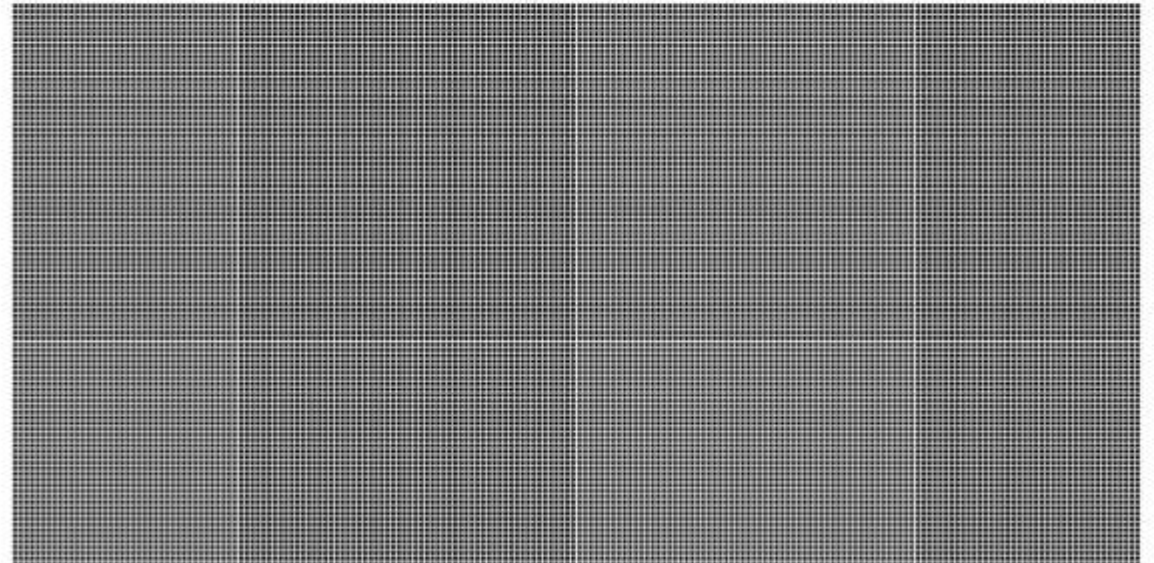
DNN

Target task

# Why do we want self-supervised learning?

# Reason 1: Scalability



50K·
1M ▬
1B

ImageNet
~1 million annotated images

Instagram
~50 billion images floating about
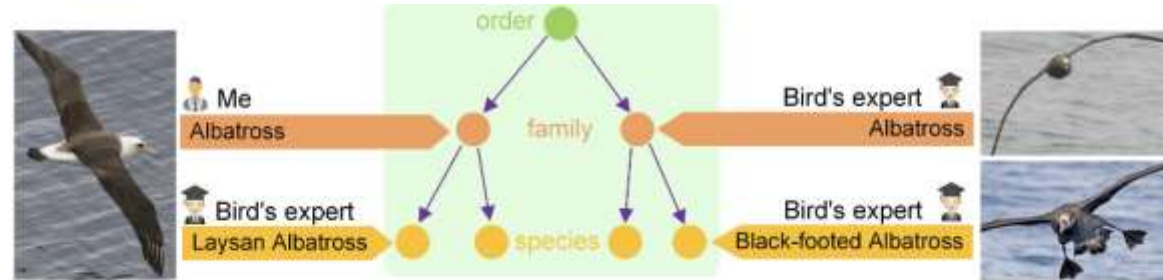
The web is filled with unanontated data.
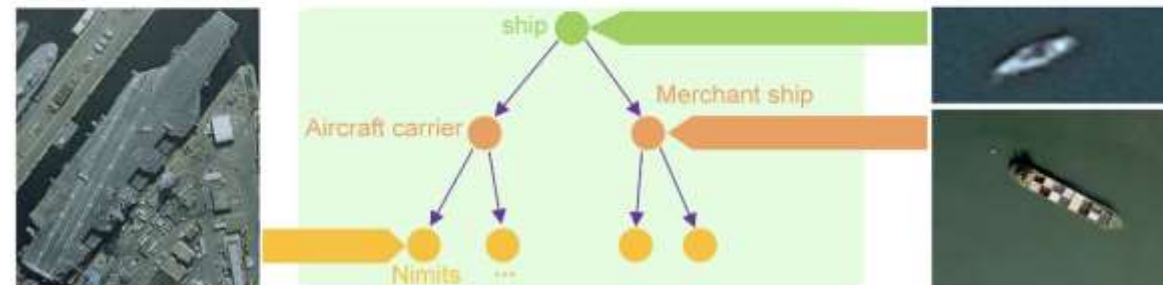
# Reason 2: Generalizability



We want models that generalize to many domains and shifts.

# Reason 3: Label are not perfect



(a) Differences in domain knowledge and interference from the image occlusion.
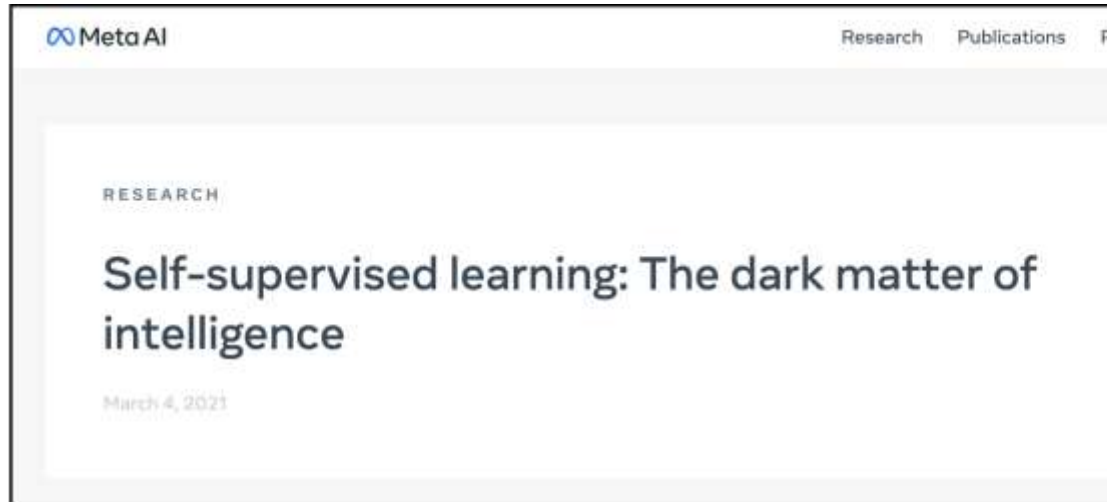
(b) Large variations of image resolutions.
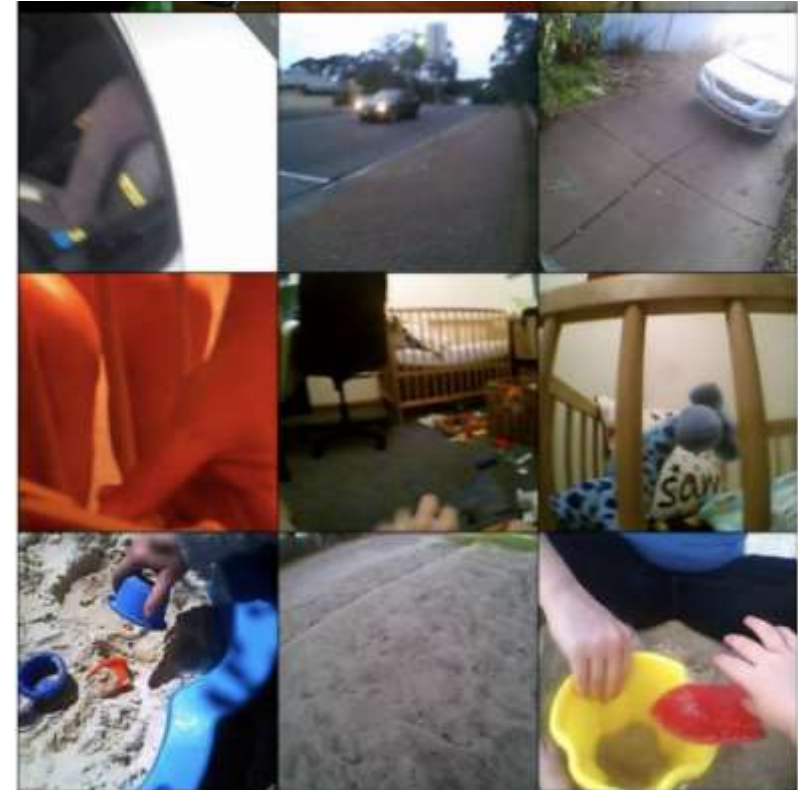
Chen et al. 2022

Labels can be ambiguous, biased, or simply wrong.

# Reason 4: Humans are self-supervised



Still a lot of lessons from human learning that can be transferred.

# How do we train deep networks without labels?

# Self-supervied visual learning

The first popular domain for self-supervised learning.

Main idea: exploit the structure of images and videos to learn without labels.

Goal is not to develop new algorithms, but borrow losses from supervised learning and think of your own loss functions.

# Early attempt: relative positioning



← 8 possible locations

Word2Vec

Motivated from NLP

Classifier

CNN    CNN

**Randomly Sample Patch**
**Sample Second Patch**

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

# Learning by coloring



Grayscale image: L channel

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

Concatenate (L,ab)

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

# Visual results



The representations for learning color can be used to initialize a new network.

# Learning by rotations



90° rotation    270° rotation    180° rotation    0° rotation    270° rotation

Assumption: if we know the object, we understand which rotation is most natural.

# Modern approach: contrastive self-supervision



NPID

SimCLR

The contrastive loss for positive pairs i,j:

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

with $z_i$, $z_j$ embeddings for images $i$ and $j$,

$\tau$ a temperature, *sim() is the dot-product*

"non-parametric" softmax

Enforces image uniqueness and augmentation invariance.

# How to train with contrastive losses

# Self-supervised learning is conservative supervised learning

Contrastive supervised learning:

Pull samples of same class together, push others away.

Contrastive self-supervised learning:

Pull augmented versions of same sample together, push others away.

# Self-supervised learning is learning augmentation invariance



(a) Original    (b) Crop and resize    (c) Crop, resize (and flip)    (d) Color distort. (drop)    (e) Color distort. (jitter)

(f) Rotate {90°, 180°, 270°}    (g) Cutout    (h) Gaussian noise    (i) Gaussian blur    (j) Sobel filtering

We want augmentations of a sample to lead to the same embedding representation.

# Self-supervised video learning

Videos have a super strong extra signal to learn from: time.

How can time be used for pretext tasks?

<span style="color:red">Predict temporal order.</span>

<span style="color:red">Predict whether video is played in reverse or not.</span>

<span style="color:red">Predict alignment between video and audio.</span>

# Examples of self-supervised video learning



Wang et al. (2020): Predict pace.



Basura et al. (2017): Predict odd-one-out.



Recasens et al. (2021): Narrow to broad prediction.



Xu et al. (2019): Predict clip order.

# Break

# Self-supervised learning on other modalities



[Akbari et al. NeurIPS 2021]



[Gong et al. AAAI 2022]



[Yang et al. ICCV 2023]

# How good are self-supervised models?
# A DINOv2 study

DINOv2 (2024): Contrastive learning + patch-level masking + tricks + 142M dataset.

# How good are self-supervised models?
## A DINOv2 study



(a)    (b)    (c)    (d)

Supervised networks struggle when being deployed in new settings,
self-supervised networks thrive in such settings.

# Self-supervised learning for language

What if our data is a collection of sentences?

*"The quick brown fox jumps over the lazy dog."*

I.e., how can we train Large Language Models on the internet?

# Masked Language Modelling

Main idea is simple: remove some tokens and predict them.

Set of classification labels:    All tokens.

Targets:    Tokens that were removed.

Just like self-supervised visual learning, the problem falls back to a standard classification setup, but now with "free labels".

# Masked Language Modelling

Standard setting: Sample 15% of tokens and replace with [MASK].

*"The quick brown [MASK] jumps over the [MASK] dog."*

Modified MLM: Sample 15% of tokens. Replace 80% with [MASK], 10% with random token, and 10% left unchanged.

*"The quick brown [oven] jumps over the [maybe] dog."*

*"The quick brown [fox] jumps over the [lazy] dog."*

# Modified Masked Language Modelling

For 80% of the sampled tokens, we simply need to predict the masked input.

For 10% of the tokens, the model needs to figure out that the word needs to be replaced.

For the remaining 10%, the model needs to figure out to do nothing.

Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

BERT

Randomly mask 15% of tokens

1  2  3  4  5  6  7  8  •••  512

[CLS]  Let's  stick  to  [MASK]  in  this  skit

Input

[CLS]  Let's  stick  to improvisation in  this  skit

https://jalammar.github.io/illustrated-bert/

# Next Sentence Prediction

Main idea: Given two sentences, predict whether the first follows the second.



How to automatically generate "labels" for this setting?

# MLM in vision: masked autoencoding

Many ideas from one domain are inspiration for the next domain.



"Masked Autoencoders Are Scalable Vision Learners" He et al. CVPR 2022

# Revisiting RLHF



**Step 1**

**Collect demonstration data and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

**Step 2**

**Collect comparison data and train a reward model.**
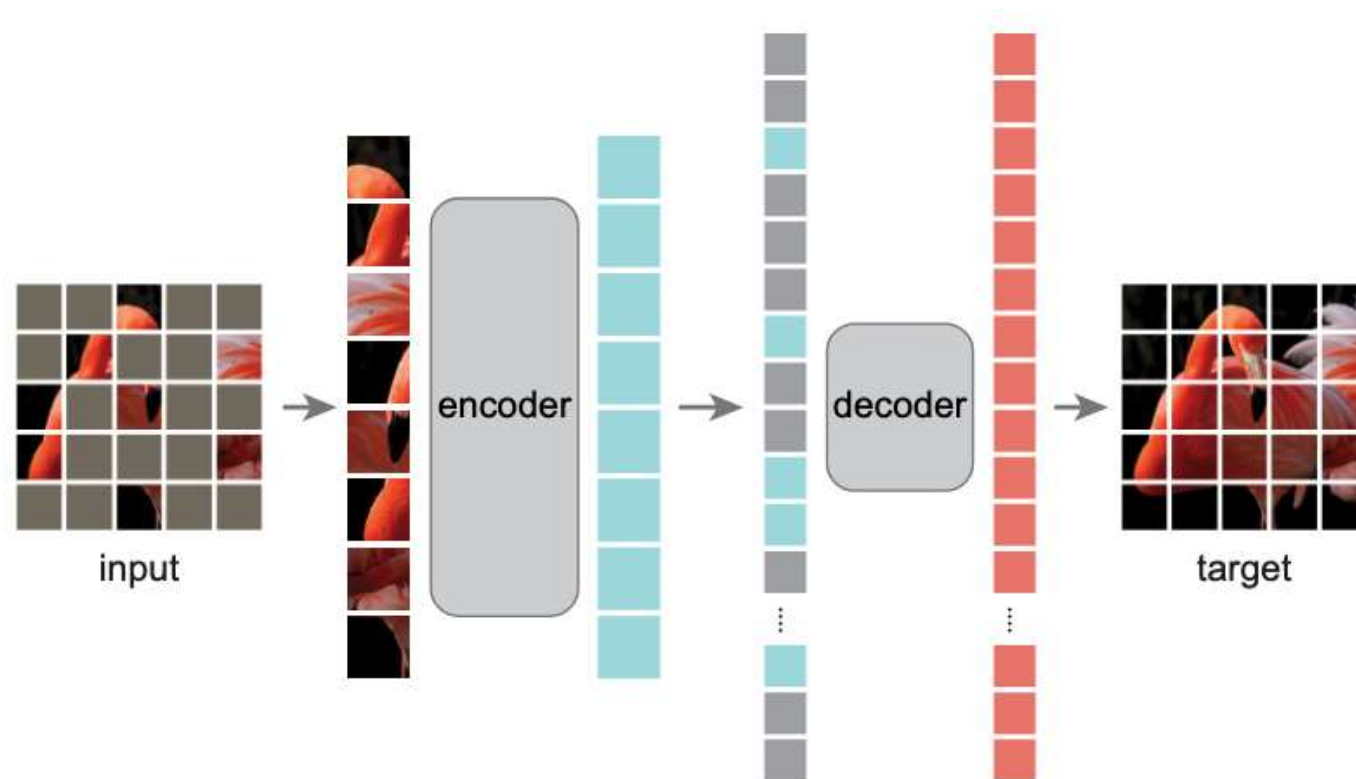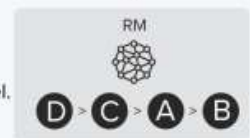
A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

A. In reinforcement learning, the agent is...

B. Explain rewards...

C. In machine learning...

D. We give treats and punishments to teach...

A labeler ranks the outputs from best to worst.

D > C > A > B

This data is used to train our reward model.

RM

D > C > A > B

**Step 3**

**Optimize a policy against the reward model using the PPO reinforcement leaning algorithm.**

A new prompt is sampled from the dataset.

Write a story about otters.

The PPO model is initialized from the supervised policy.

PPO

The policy generates an output.

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

# RLHF from a supervision perspective

Self-supervision (GPT) is not enough for Large Language Models (ChatGPT).

*Prompt: Explain why we need water to survive as humans.*

*(Imaginary) GPT: Explain why humans die when not given any water.*

What is going on here?

We lack alignment with the intent of the user input.

This needs human supervision.

## Step 1

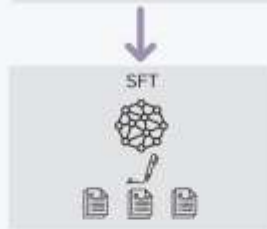**Collect demonstration data and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

## Step 2

**Collect comparison data and train a reward model.**
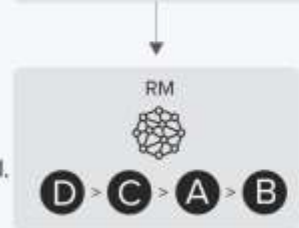
A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

**A**
In reinforcement learning, the agent is...

**B**
Explain rewards...

**C**
In machine learning...

**D**
We give treats and punishments to teach...

A lableler ranks the outputs from best to worst.

D > C > A > B

This data is used to train our reward model.

RM

D > C > A > B

## Step 3

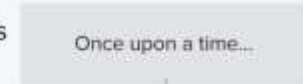**Optimize a policy against the reward model using the PPO reinforcement leaning algorithm.**

A new prompt is sampled from the dataset.

Write a story about otters.

The PPO model is initialized from the supervised policy.

PPO

The policy generates an output.

Once upon a time...

The reward model calculates a reward for the output.

RM

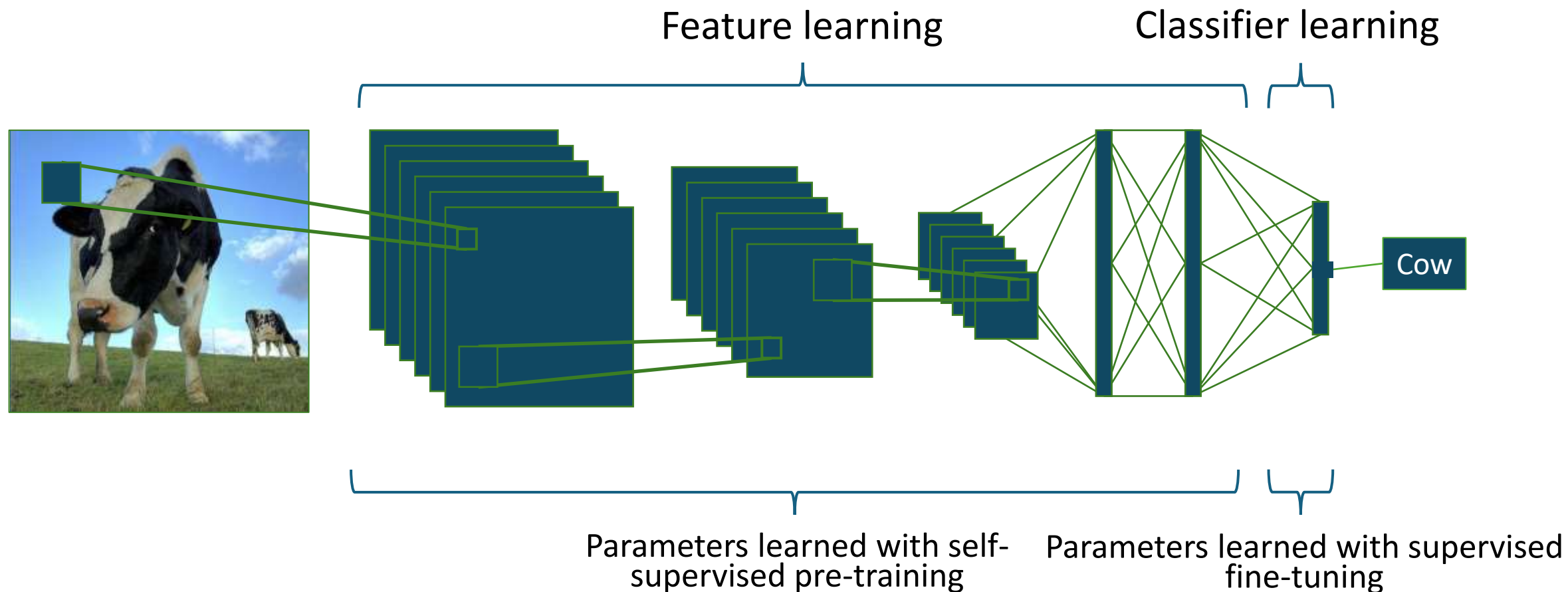The reward is used to update the policy using PPO.

$r_k$

LLM = Transformer + self-supervised pre-training + human aligned tuning
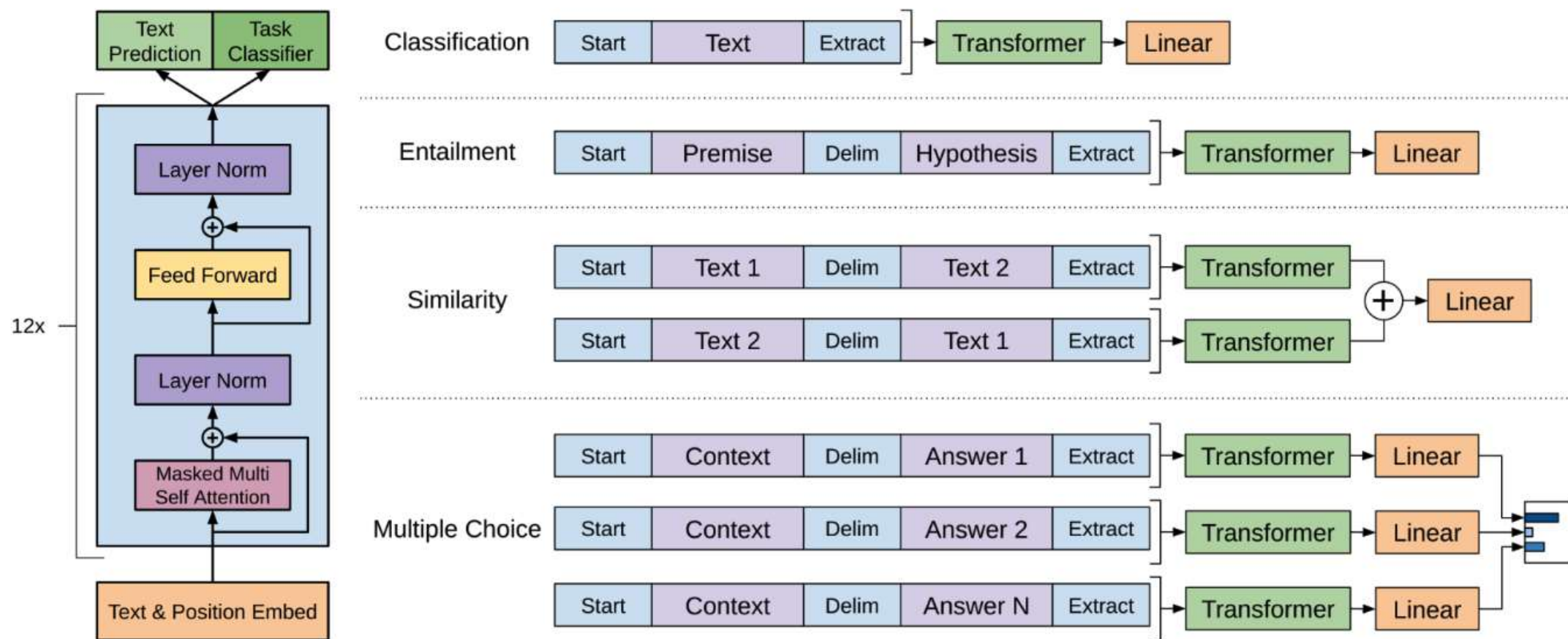
Discussion: why are they so good? And is this all you need for deep learning?

# Is there something in between supervised and self-supervised learning?

# Classical setup: pre-training and fine-tuning



Feature learning

Classifier learning

Cow

Parameters learned with self-supervised pre-training

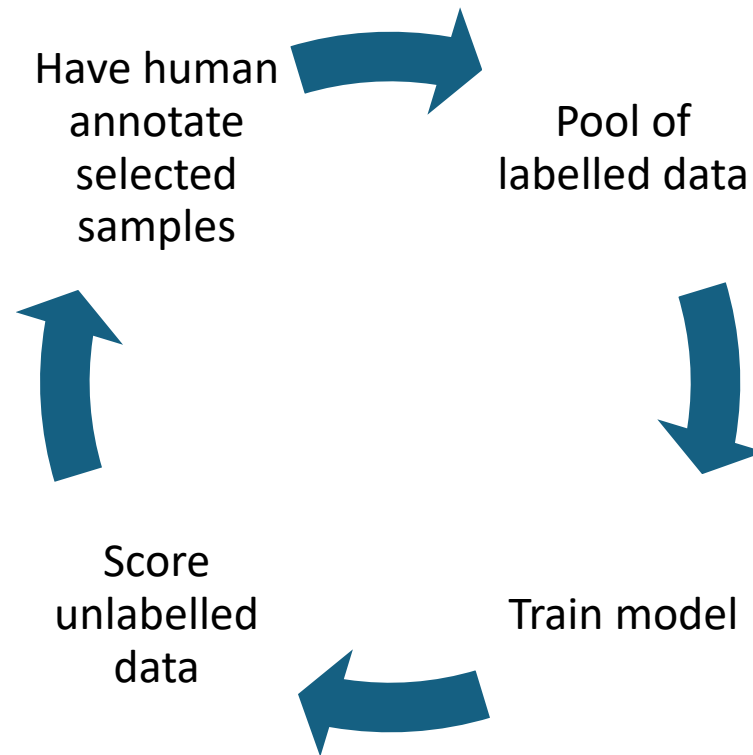Parameters learned with supervised fine-tuning

# Examples of fine-tuning in language

# Active learning

Assume we only have an unlabelled training set. Is it possible to simultaneously label samples and train a model?
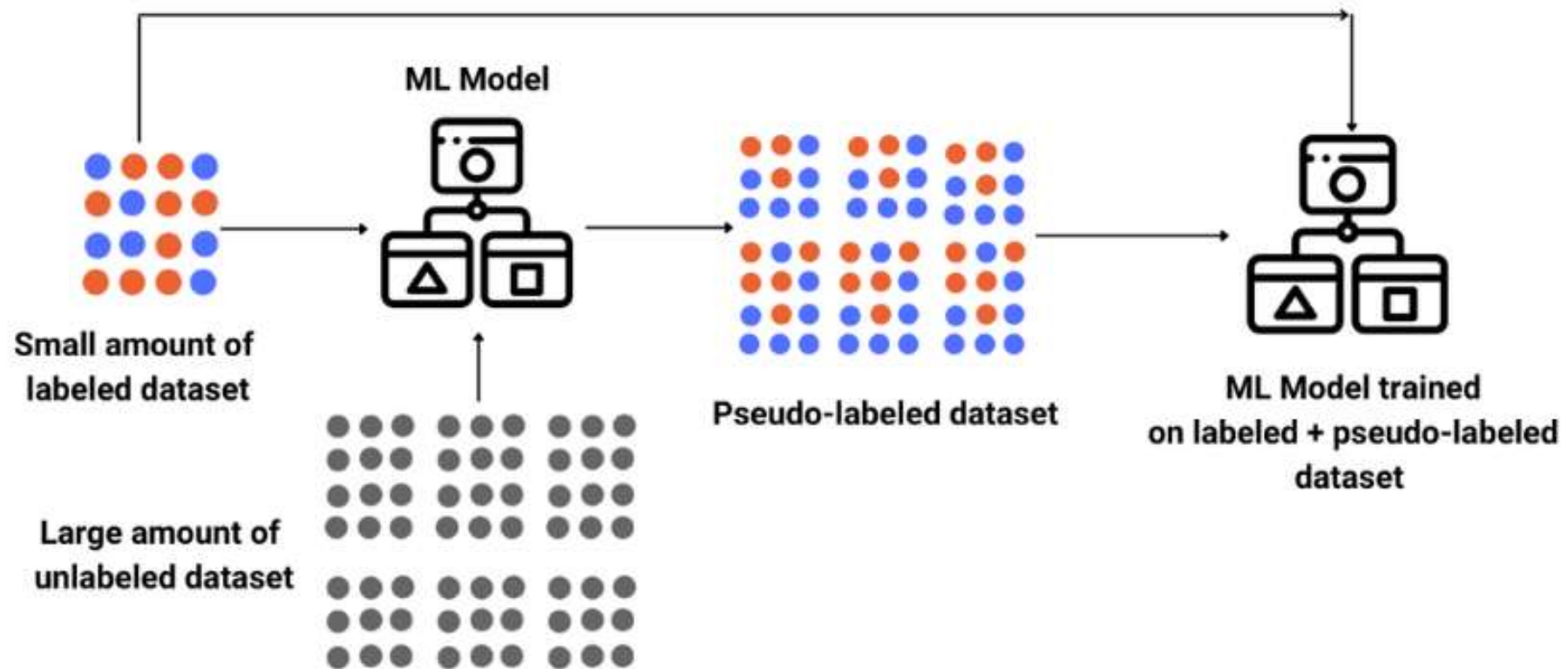
# Which samples to select?

**Random:** randomly select (ignore scoring).

**Most uncertain:** closest to the decision bounary, or lowest norm in embedding space (second to last layer), or highest likelihood entropy.

**Group-based metrics:** Uniformity over classes to avoid biases.

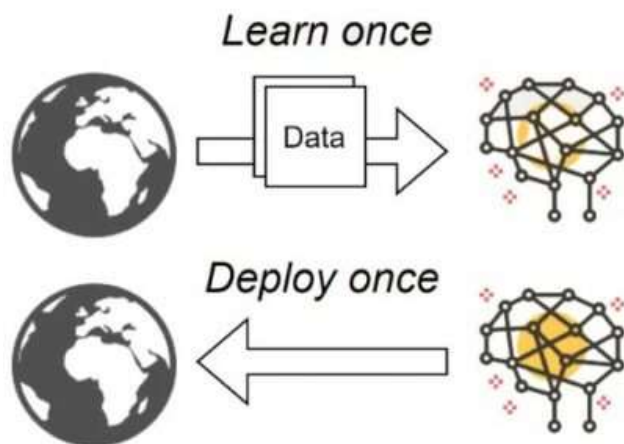**Mix:** Combine a mix of X% random and (100-X)% uncertain+group.

# Semi-supervised learning

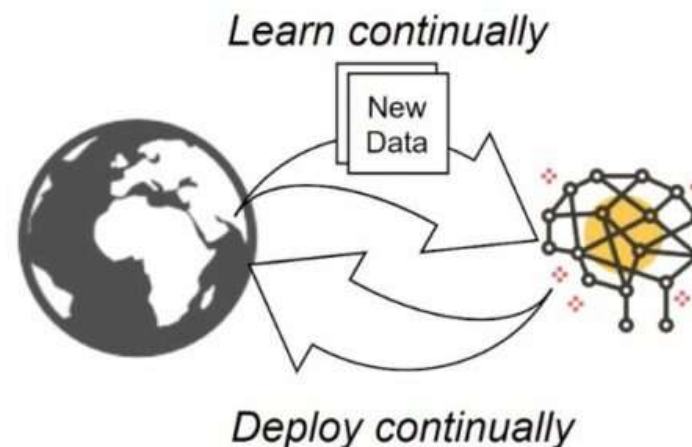# Continual learning

In the real world, there is no such thing as a static dataset.



Shockingly, we can't just train on new data! Much more in lecture 11.

# Is self-supervised learning truly without supervision?

# My view on supervised vs self-supervised learning

**Supervised learning**

**Self-supervised learning**

Label by sample.

Label by rule.

Invariance defined
at global semantic level.

Invariance defined
at geometric or local semantic level.

Self-supervised learning is conservative supervised learning from pre-defined invariances.

# Next lecture

| Lecture | Title | Lecture | Title |
|---------|-------|---------|-------|
| 1 | Intro and history of deep learning | 2 | AutoDiff |
| 3 | Deep learning optimization I | 4 | Deep learning optimization II |
| 5 | Convolutional deep learning | 6 | Attention-based deep learning |
| 7 | Graph deep learning | 8 | From supervised to unsupervised deep learning |
| 9 | Multi-modal deep learning | 10 | Generative deep learning |
| 11 | What doesn't work in deep learning | 12 | Non-Euclidean deep learning |
| 13 | Q&A | 14 | Deep learning for videos |

# Learning and reflection

TODO

# Thank you