



# *Self-Supervised Learning*

*COMPSCI 760*  
*2024 Semester 1*

*Olivier Graffeuille*  
[ogra439@aucklanduni.ac.nz](mailto:ogra439@aucklanduni.ac.nz)



## Outline

1. Motivation
2. Background
3. Methods for SSL
  1. What is a generative model?
  2. Generative models for SSL
  3. Contrastive methods for SSL



## Outline

- 1. Motivation**
2. Background
3. Methods for SSL
  1. What is a generative model?
  2. Generative models for SSL
  3. Contrastive methods for SSL



## Issues with Supervised Learning

- Supervised learning models have had great success in modelling various tasks in recent years
- However, these systems require *massive amounts of carefully labelled* data
  - Practically, it is impossible to label everything in the world
  - Some things don't have sufficient data (e.g. low-resource languages)
  - Supervised learning models tend to have narrow intelligence



## Human vs Machine Learning

- It's clear that humans are able to learn things with far fewer examples
  - A self-driving car must drive off a cliff in a simulation thousands of times before learning not to
  - Children recognising animals
- We seem to have a lot of "background" knowledge about the world: "Common sense"?
- This ability to learn quickly comes from millions of years of "pre-training" via evolution



## **Self-Supervised Learning**

- Allows models to learn from unlabelled data
  - Substantially more unlabelled data
  - Important to learn subtle knowledge and less common concepts
- Self-Supervised Learning is a Machine Learning “Paradigm”
  - Other topics recently covered including data streams, domain generalisation etc. can be thought of as “settings” of ML



# Motivation

“Self-Supervised Learning is one of the most promising ways to build such background knowledge and approximate a form of ‘common sense’ in AI systems”

- Yann LeCun



# Self-Supervised Learning

## Outline

1. Motivation
2. *Background*
3. Methods for SSL
  1. What is a generative model?
  2. Generative models for SSL
  3. Contrastive methods for SSL



## Background

# **Self-Supervised Learning (SSL)**

- Learning by observation?
- We still want to make a prediction!
  - But without a human-defined label



## Background

# Self-Supervised Learning (SSL)

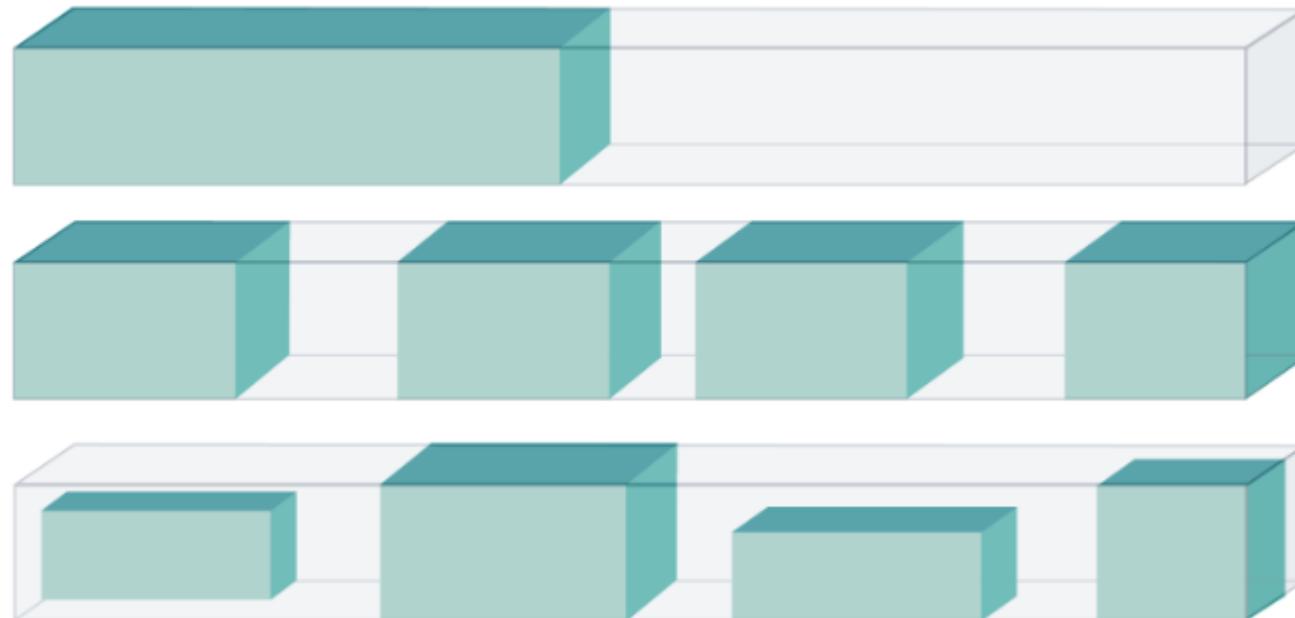
Learns to predict a part of the data from another part of the data.  
Two steps:

1. Obtain “labels” from the data itself “semi-automatically”
  - Could be incomplete, transformed, distorted, or corrupted etc.
2. Predict something about the data from other parts of the data
  - E.g. the machine learns to ‘recover’ all or some parts of its original input
  - “Repairing” the data

# Background

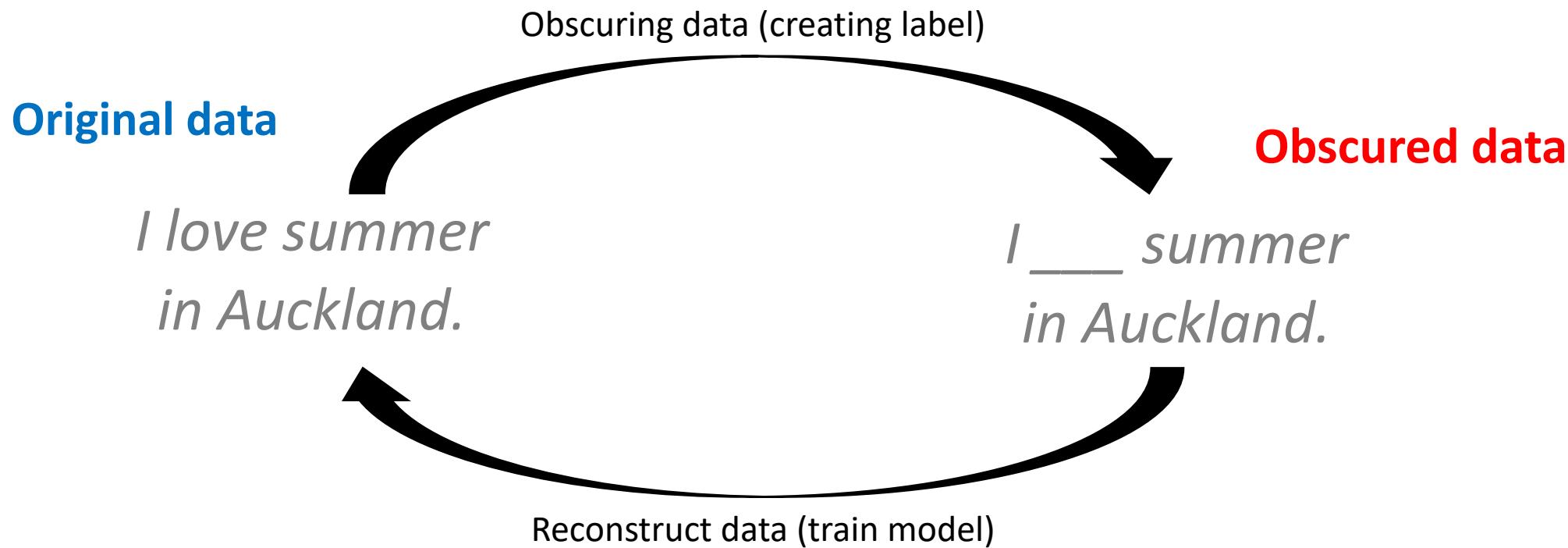
## Obtaining Labels

- One recipe: hide part of the label



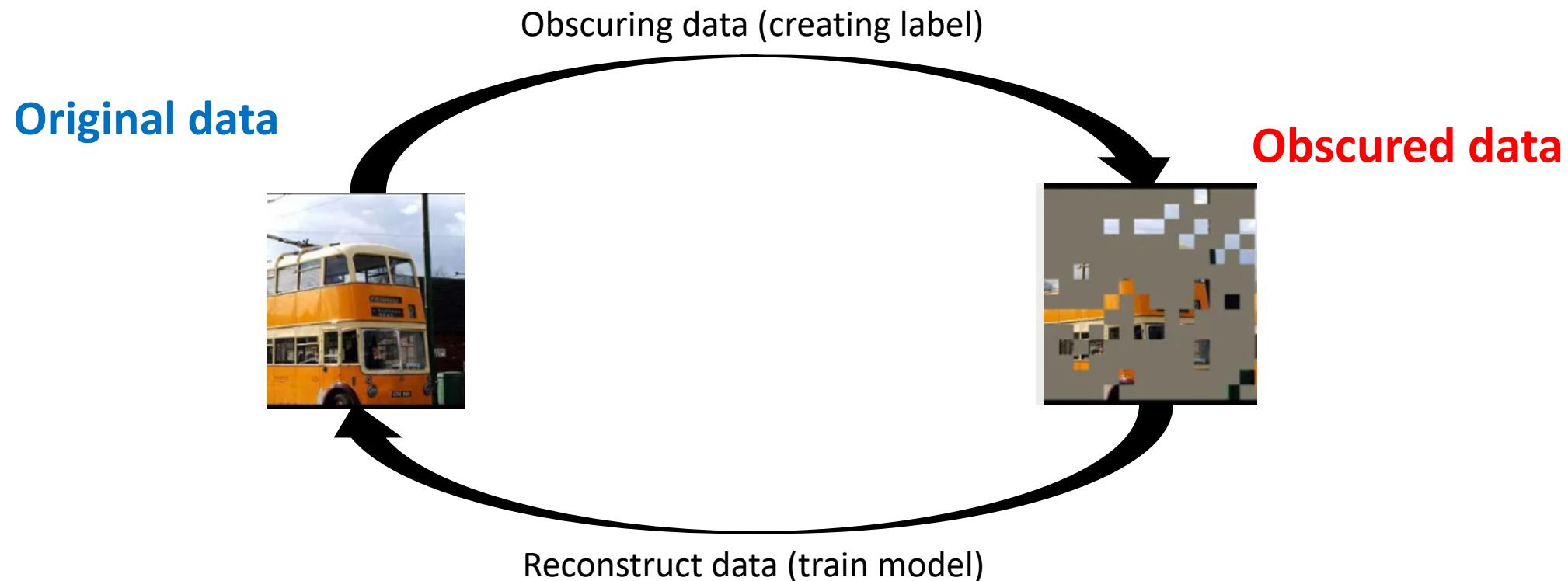
## Background

# Obtaining Labels: Examples



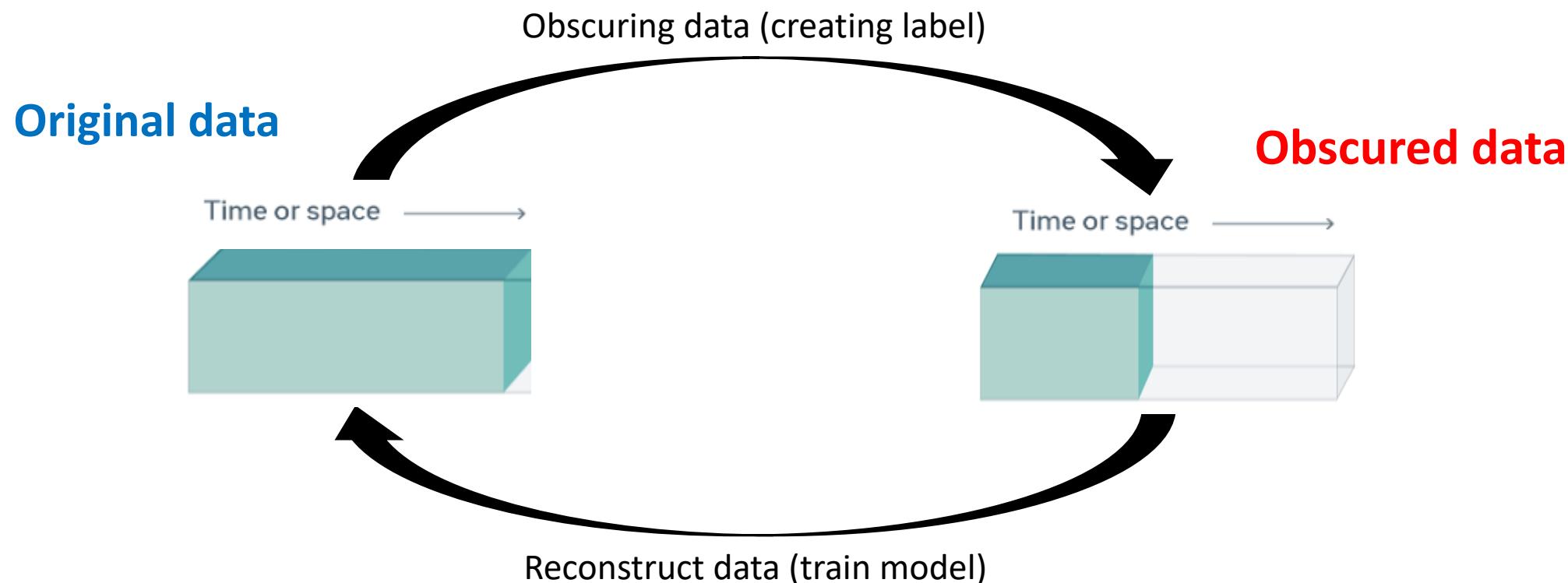
## Background

# Obtaining Labels: Examples



## Background

# Obtaining Labels: Examples





## Background

# Unsupervised Learning?

- If we are learning from unlabelled data, is this not just unsupervised learning?
  - E.g. clustering, dimensionality reduction, anomaly detection
- “Self-Supervised Learning” is used to imply that we are still learning from supervision
  - SSL is not unsupervised at all
  - We are still making predictions and learning from many labels
    - Arguably, we are making richer predictions

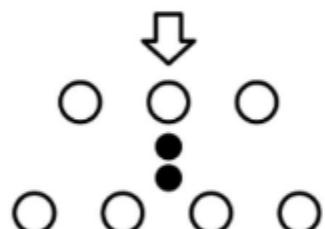
# Background

# Unsupervised Learning?

Supervised  
implausible labels

**"COW"**

Target



Input

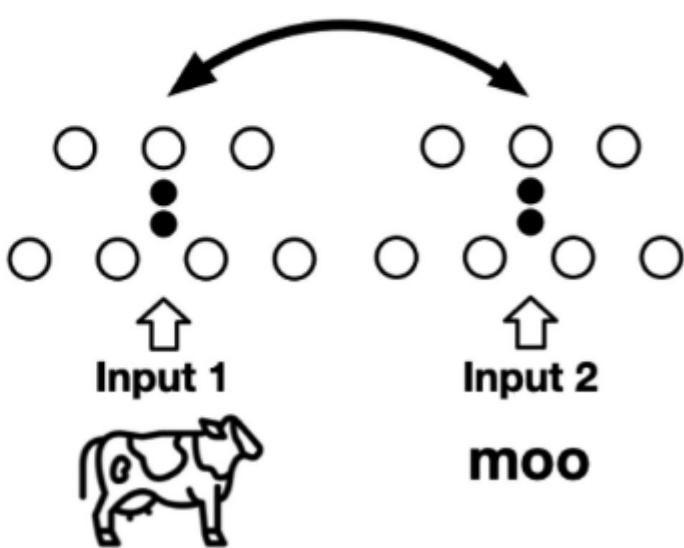


Unsupervised  
limited power

Input



Self-supervised  
derives label from a  
co-occurring input to  
related information



Input 2

moo



## Background

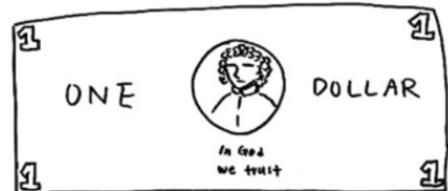
# Obtaining Labels

- One of the challenges of supervised learning is human bias in labels
- There is still a human bias in SSL labels, since humans are defining the labelling process. Do we want the model to:
  - Fill in gaps?
  - De-noise?
  - Predict the future?

# Background

## Why is this useful?

- This is not just useful to repair data
- It is useful for learning **meaningful representations** from huge quantities of data



High-level representation  
of a \$1 bill



Low-level representation  
of a \$1 bill



# Background

## Why is this useful?

- This is not just useful to repair data
- It is useful for learning **meaningful representations** from huge quantities of data
  - If we can predict the future of a sentence, we have probably learnt the representations required to label that sentence

“Despite its flaws, I have come to \_\_\_\_ this product.” Self-prediction

:

“Despite its flaws, I have come to love this product.” -> Positive sentiment Classification

- SSL can have more feedback signals from labels than supervised



# Self-Supervised Learning

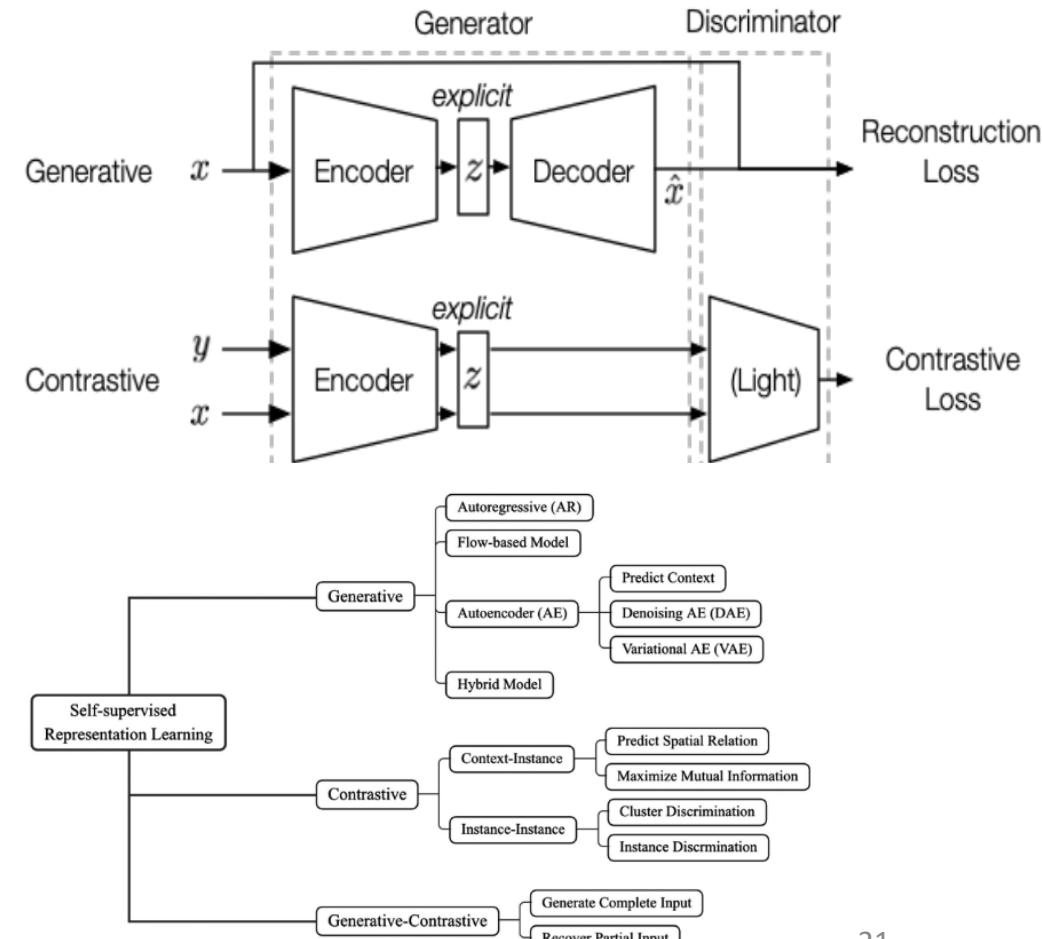
## Outline

1. Motivation
2. Background
- 3. *Methods for SSL***
  1. What is a generative model?
  2. Generative models for SSL
  3. Contrastive methods for SSL

# Self-Supervised Learning Methods

## Self-Supervised Learning Approaches

- Generative Methods
  - Reconstruct the original data
  - Typically encoder/decoder architecture
- Contrastive Methods
  - Compare embeddings of input data
  - Typically encoder architecture
- Generative-contrastive methods





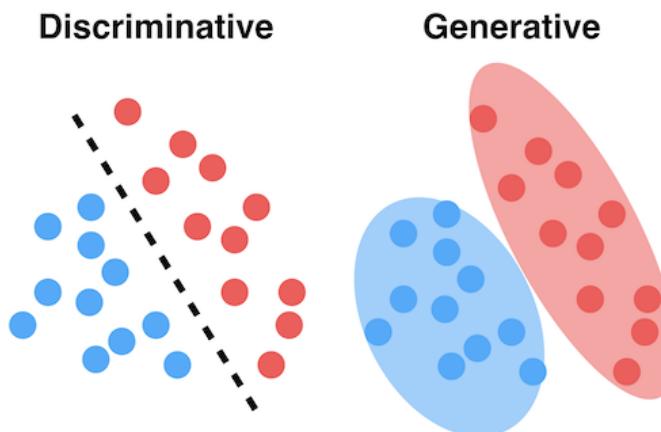
# Self-Supervised Learning

## Outline

1. Motivation
2. Background
3. Methods for SSL
  1. *What is a generative model?*
  2. Generative models for SSL
  3. Contrastive methods for SSL

# Generative vs. Discriminative Methods

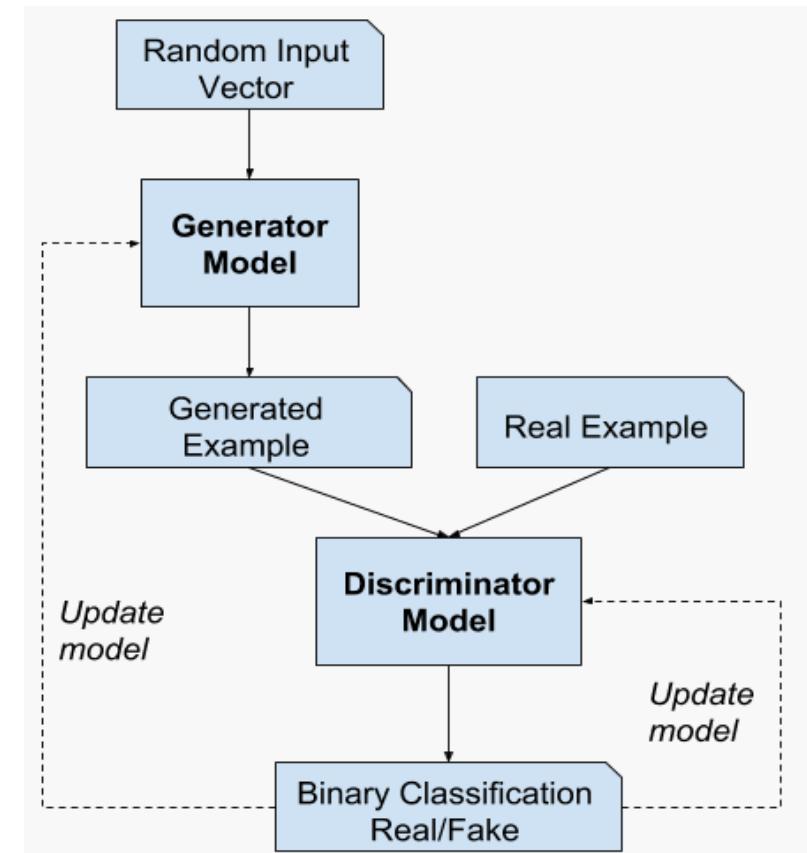
- You will be used to discriminative models:  
learn conditional probability  $P(y|X)$
- Generative models achieve a different goal:  
learn joint probability  $P(X,y)$
- In other words: learn the data probability density directly



# Generative Models

## Generative Adversarial Networks (GANs)

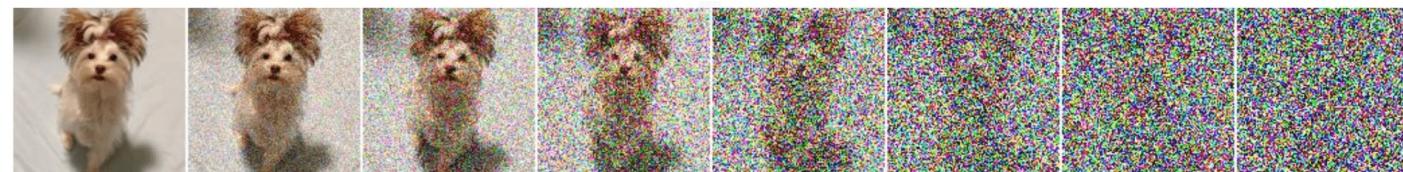
- Generator
  - Generate fake samples, trying to fool the discriminator
- Discriminator
  - Discriminating real and fake samples
- Disadvantages:
  - Difficult to train, unstable
  - Mode collapse



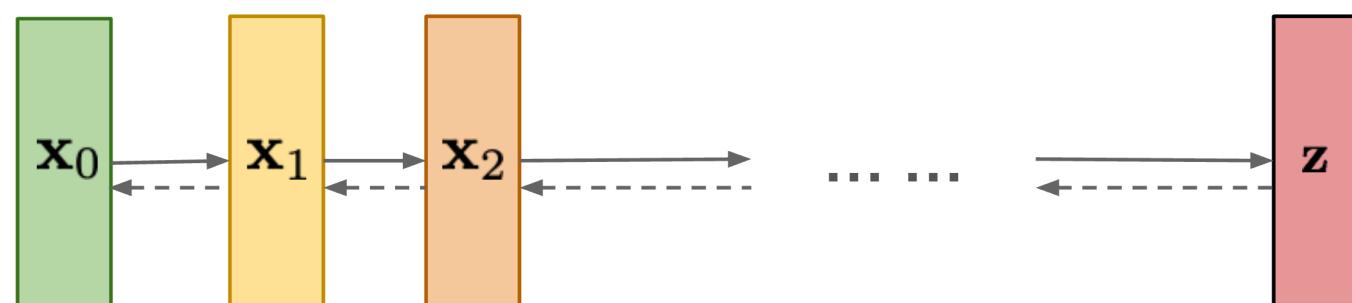
# Generative Models

## Diffusion Models

1. Add some noise to data, train model to remove noise
2. Repeat many times
3. Input noise
4. ???



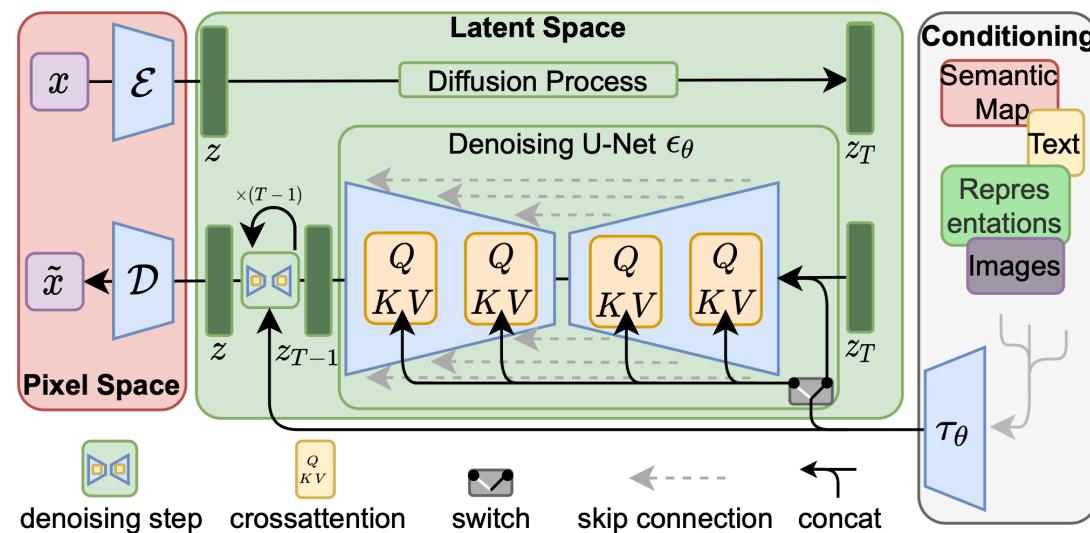
**Diffusion models:**  
Gradually add Gaussian  
noise and then reverse



# Generative Models

## Diffusion Models

- Overall architecture involves more steps
- Conditioning: allows us to “direct” the generating process





THE UNIVERSITY OF  
**AUCKLAND**  
Te Whare Wānanga o Tamaki Makaurau  
NEW ZEALAND

# Generative Models

## Diffusion Models

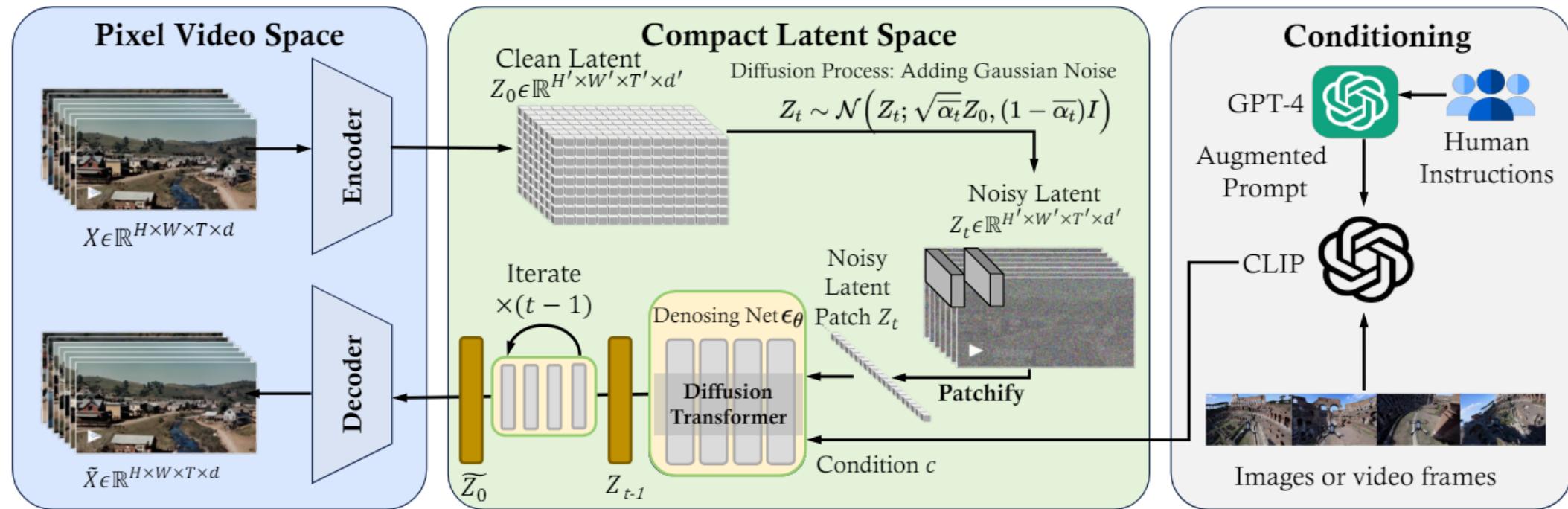


<https://stability.ai/blog/stable-diffusion-public-release>



# Generative Models

## Video Diffusion Models



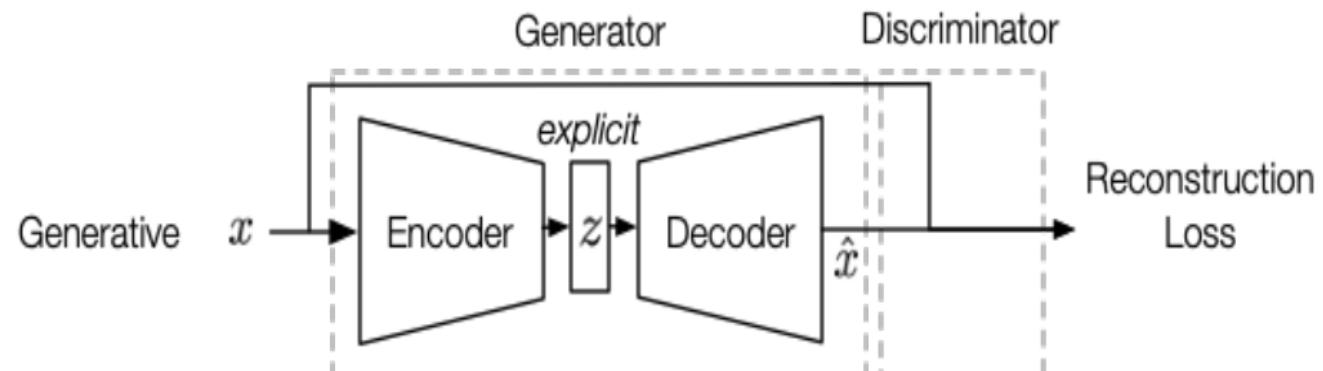


## Outline

1. Motivation
2. Background
3. Methods for SSL
  1. What is a generative model?
  - 2. *Generative models for SSL***
  3. Contrastive methods for SSL

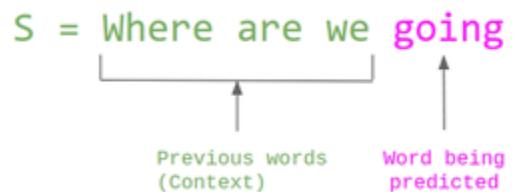
# Generative Methods

- Reconstructs a part of the input data
- More exactly:
  - Train an encoder to encode input  $x$  into an explicit vector, and a decoder to reconstruct  $x$  from  $z$



## Autoregressive models

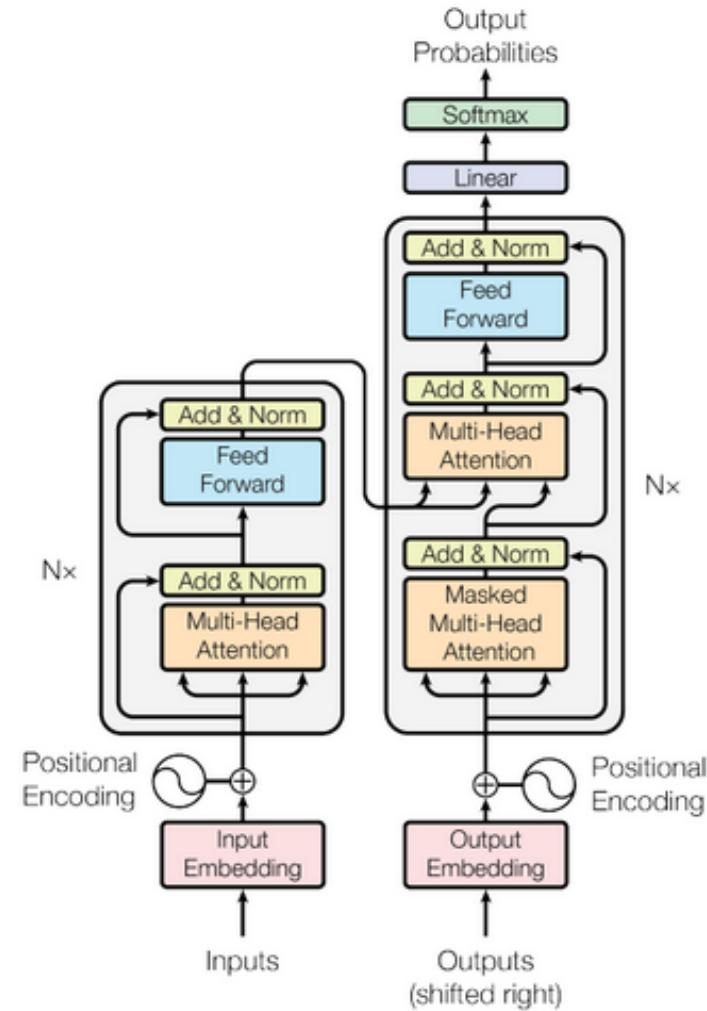
- A feed-forward network that predicts future values based on past values
- RNN, LSTM and transformers (and timeseries)



# Generative Methods

## Transformers, GPTs

- Most famous examples of autoregressive models
  - Trained on a large portion of the internet
  - Trained to predict the next word
- Fine-tuned with Reinforcement Learning with Human Feedback (RLHF)
  - To behaves more like a useful chatbot, rather than just completing text

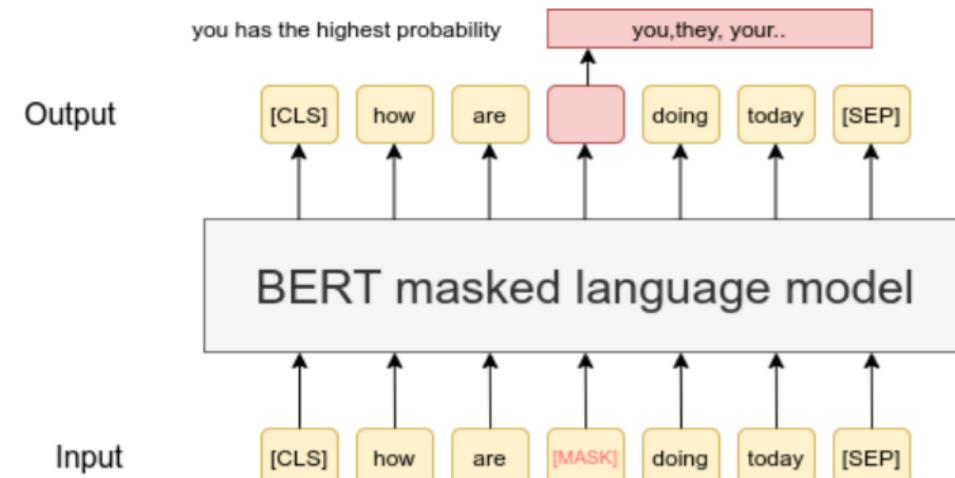


# Generative Methods

## Masked Language Models (MLM)

E.g. BERT

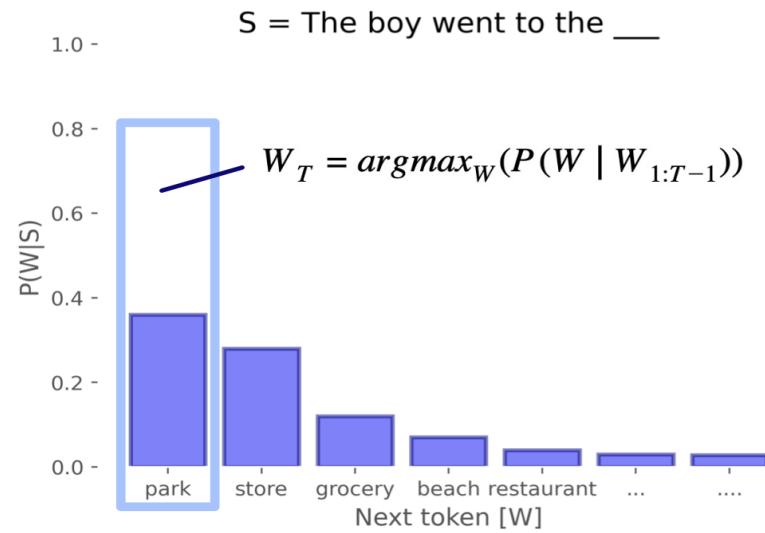
- Mask a part of the input with special [MASK] token
- Get model to determine missing token



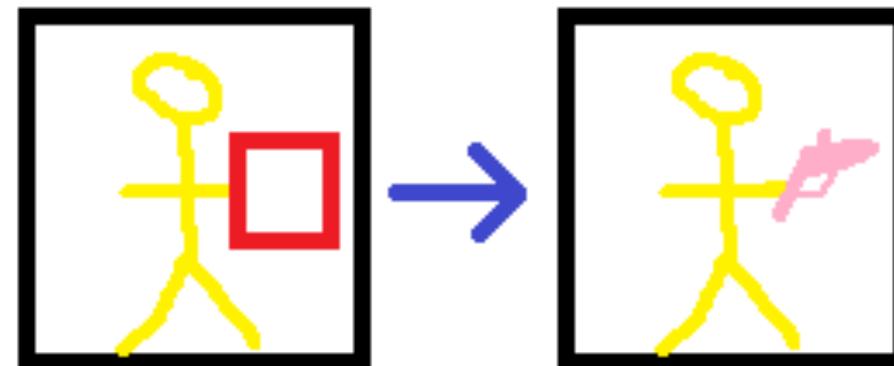
# Generative Methods

## SSL for Language vs Vision Tasks

- Let's think of how we generate text versus images:



Language Generation



Vision Generation



## SSL for Language vs Vision Tasks

- Language tasks are much better able to model uncertainty
  - We can give a distribution over every possible word (maybe 50,000 possibilities?)
  - We cannot generate a distribution over every possible image ( $2^{512 \times 512 \times 3}$  possibilities)



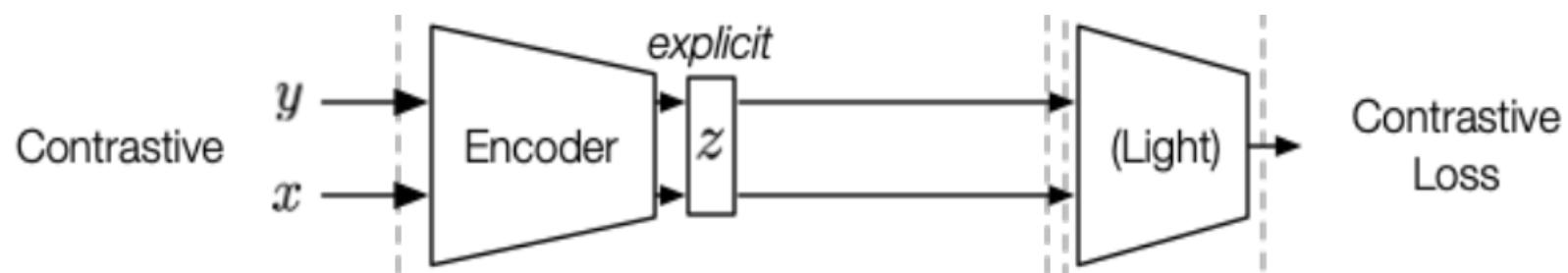
# Self-Supervised Learning

## Outline

1. Motivation
2. Background
3. Methods for SSL
  1. What is a generative model?
  2. Generative models for SSL
  3. *Contrastive methods for SSL*

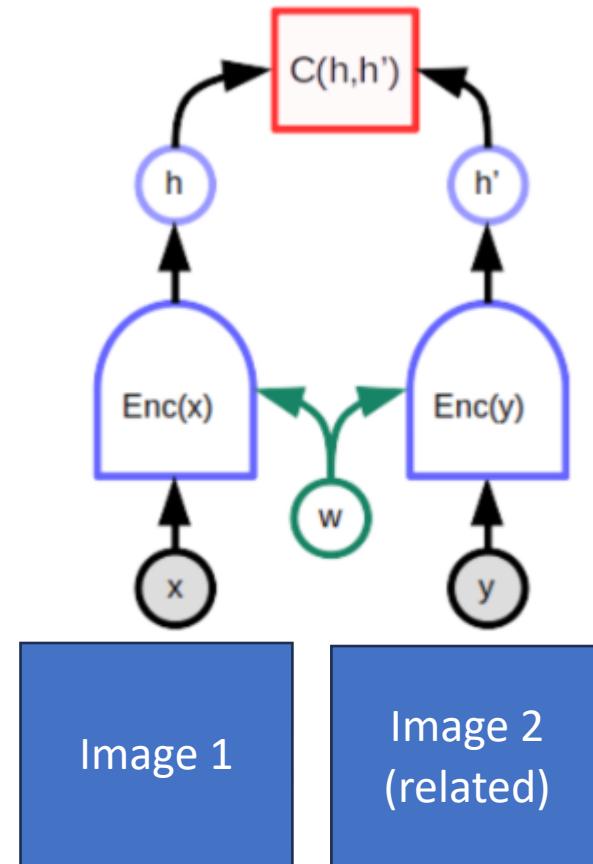
# Contrastive Methods

- Train an encoder to encode input  $x$  into an explicit vector  $z$  to measure similarity
  - No longer generating hidden part of the input
  - Instead, we predict a hidden property (latent variables) about the input
- The goal is still to learn meaningful hidden representations



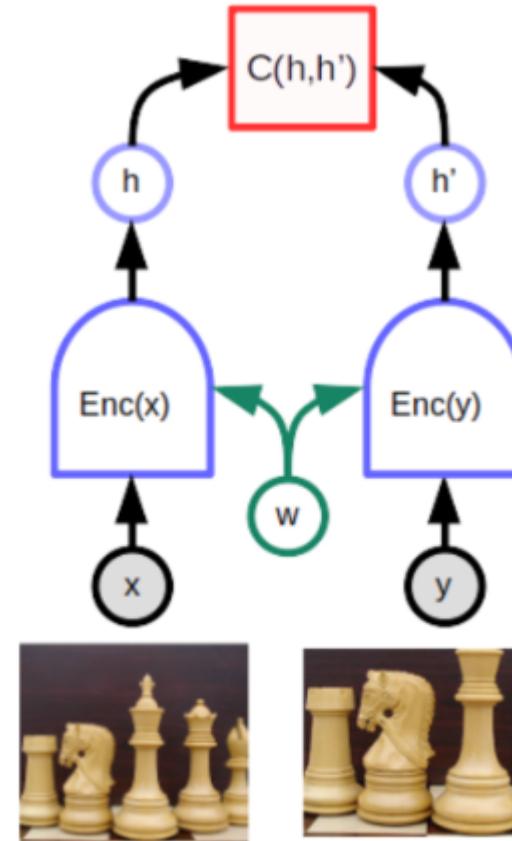
# Origins: Siamese Networks

- Pass two *related* inputs through two identical neural networks
- Compare the generated latent embeddings generated by each network



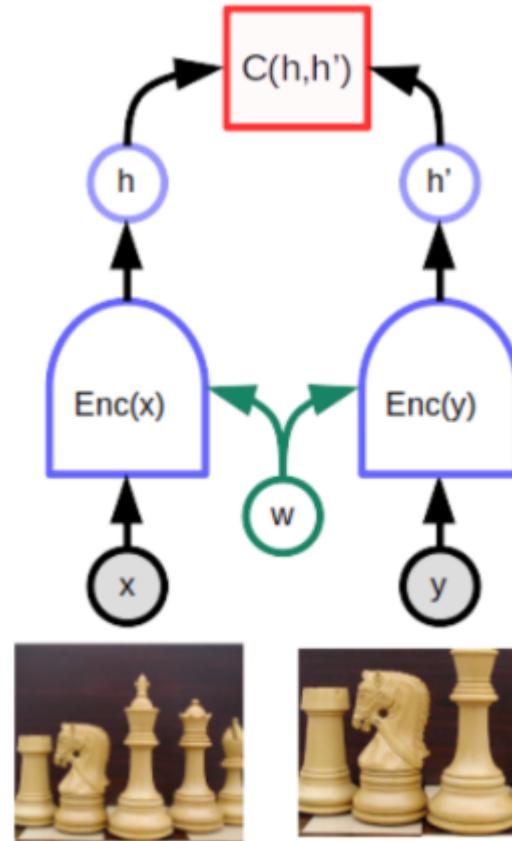
## How do choose related inputs?

- Remember, we have no labels
- We can find two related inputs by transforming an image into two inputs
  - (most common) two crops of same image
  - Changing colour, noise, etc...
  - Transformation is important
- We have a problem...



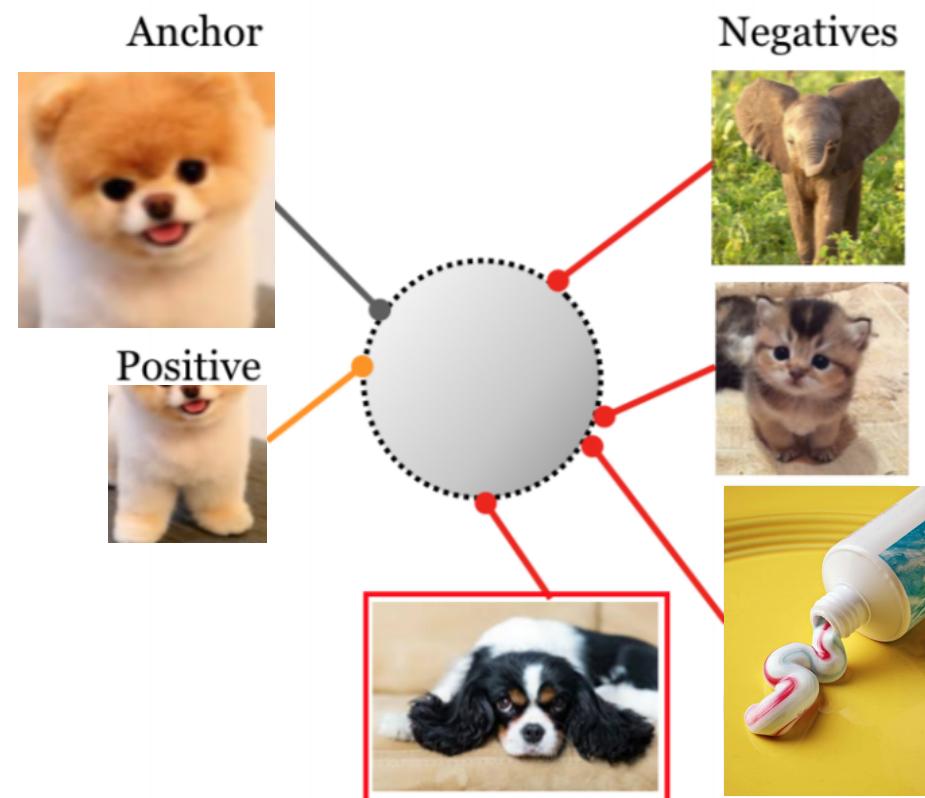
# The Problem: Mode Collapse

- Nothing is stopping the network from generating a trivial embedding
  - Ignoring the input
  - Output the same embedding to get a good loss score
- This is called mode collapse
  - Think back to GANs...



# The Solution: Contrastive Learning

- Take two *unrelated* inputs and make sure they generate different embeddings
- How to find unrelated inputs?
  - Remember: still no labels
  - We just transform any other image

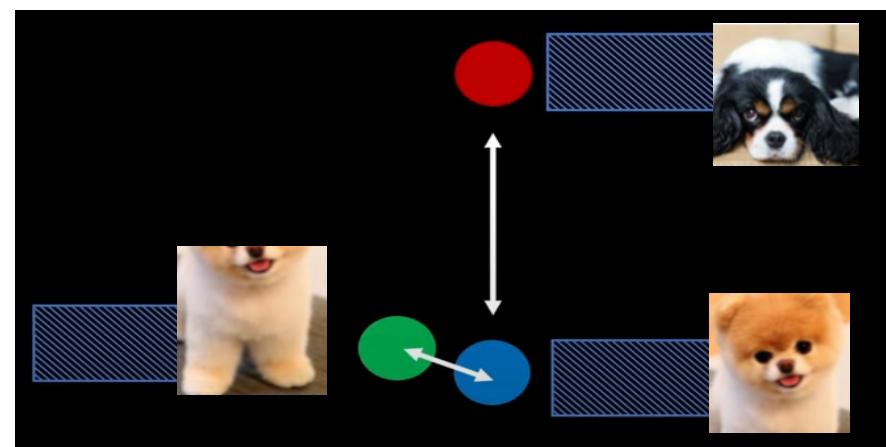
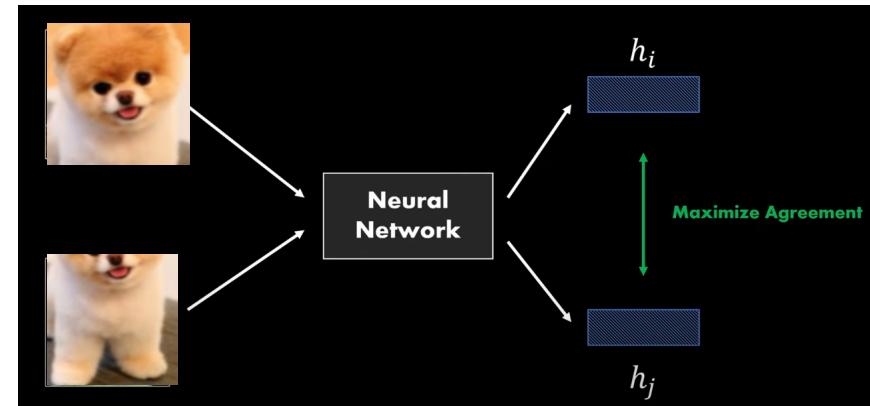


# Contrastive Methods

## The Solution: Contrastive Learning

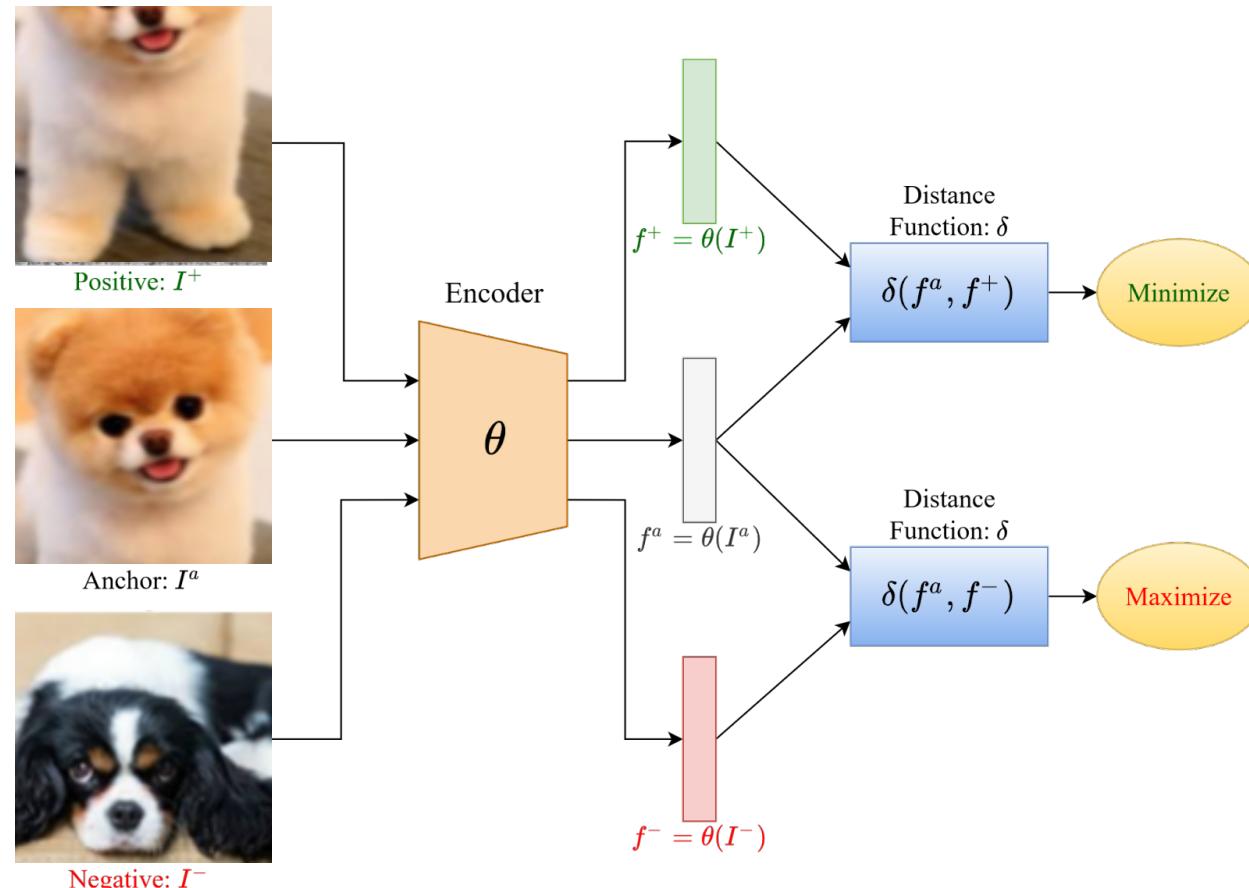
Simple implementation:

1. Get anchor data point
  - Choose positive sample
  - Choose negative sample
2. Calculate Euclidean distance between output space of input pairs
3. Minimise distance for related pairs
4. Maximise distance for unrelated pairs



# Contrastive Methods

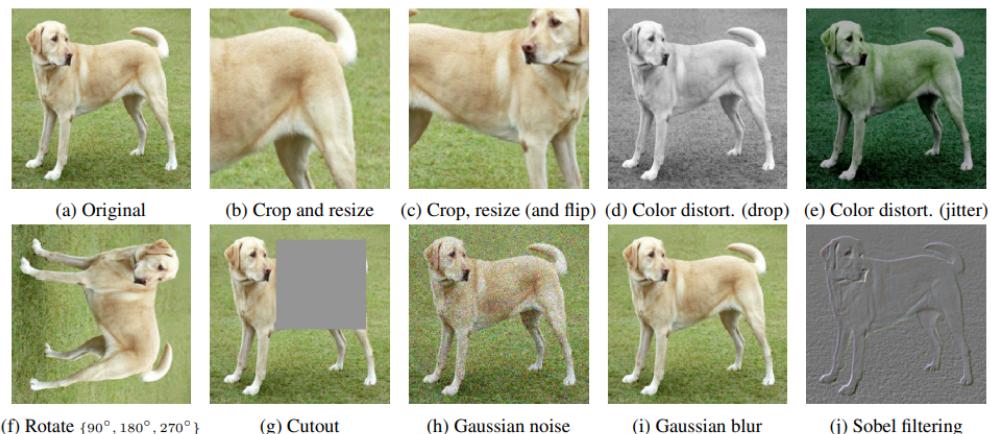
## The Solution: Contrastive Learning



# Contrastive Methods

## SimCLR (Simple Contrastive Learning of visual Representations)

- For positive and negative examples, generate many transformations of each image
  - Focus on transformations: not just cropping
  - Focus on hard positive samples



# Contrastive Methods

## SimCLR (Simple Contrastive Learning of visual Representations)

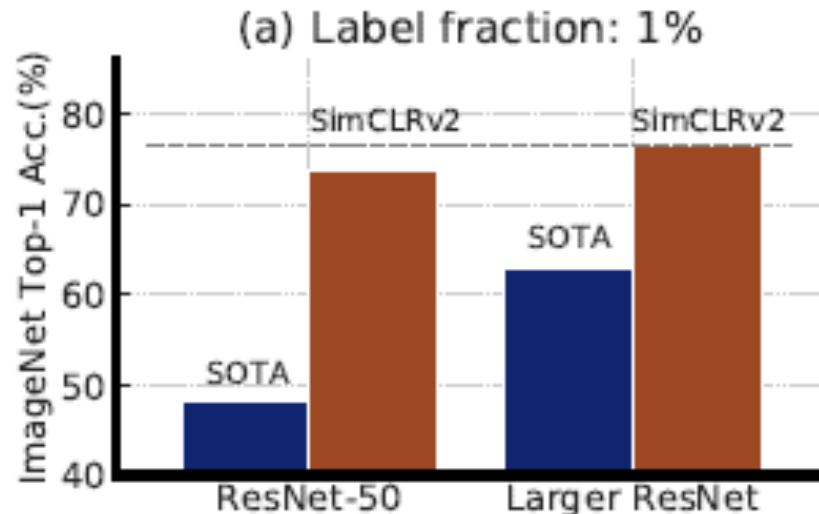
- What transformations to use?



# Contrastive Methods

## SimCLRv2

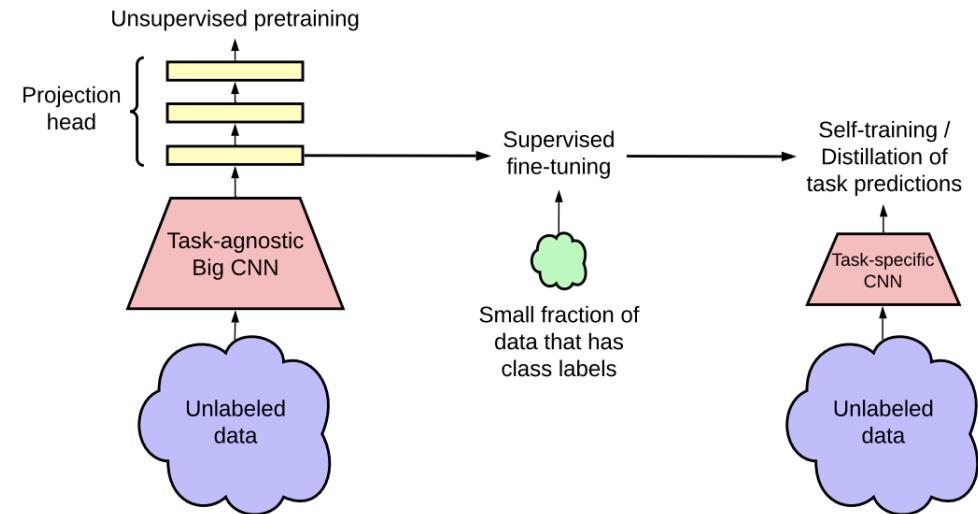
- Develop SimCLR further by applying it to supervised downstream tasks
- Works very well



# Contrastive Methods

## SimCLRV2

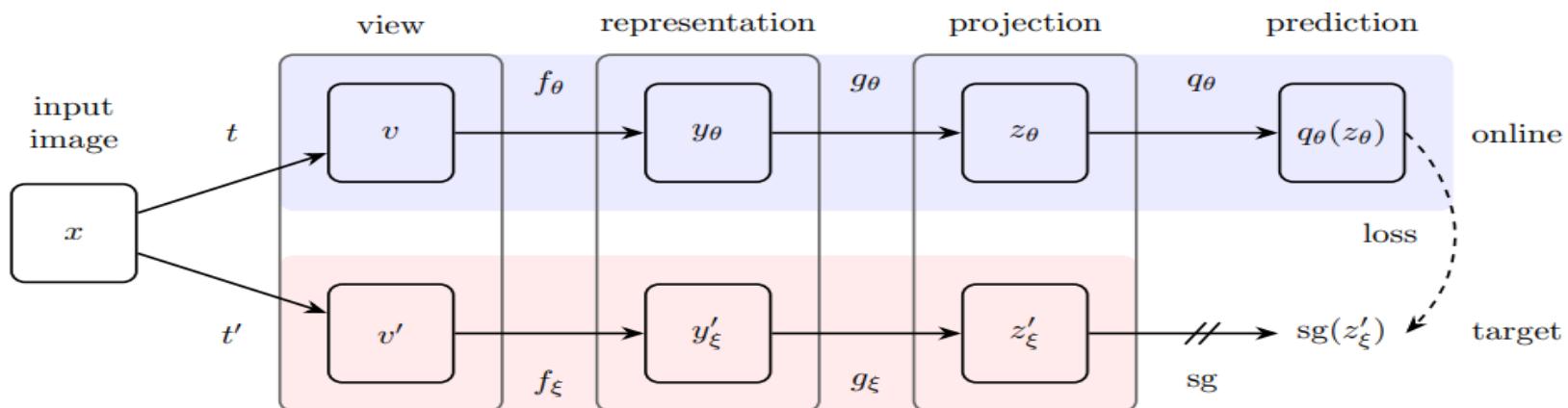
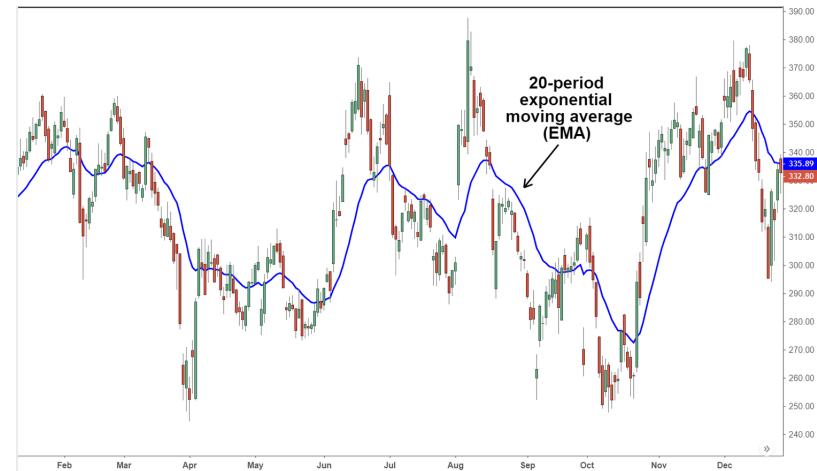
- Phase 1
  - Initial self-supervised pre-training (same as SimCLR with projection head)
- Phase 2
  - A supervised fine-tuning for specific tasks
- Phase 3
  - Unsupervised self-distillation further improves performance



# Contrastive Methods

## BYOL (Bootstrap Your Own Learning)

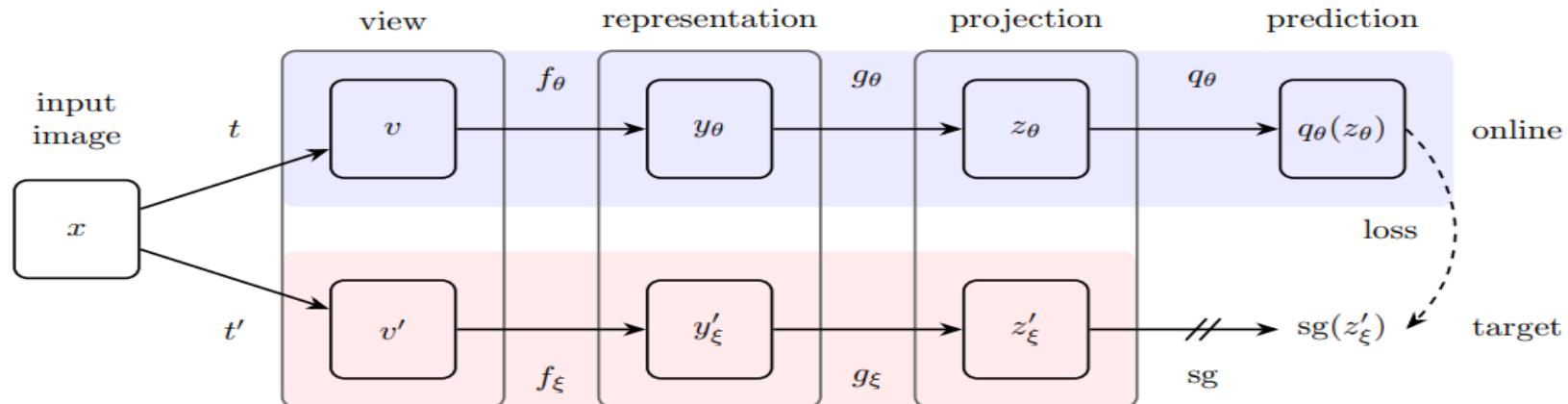
- Uses two models
  - Online model is a typical neural network, trained continuously
  - Target model is more stable: exponential moving average



# Contrastive Methods

## BYOL (Bootstrap Your Own Learning)

- No negative examples to stop collapse
- Collapse seems to be stopped by stability from target network



# Contrastive Methods

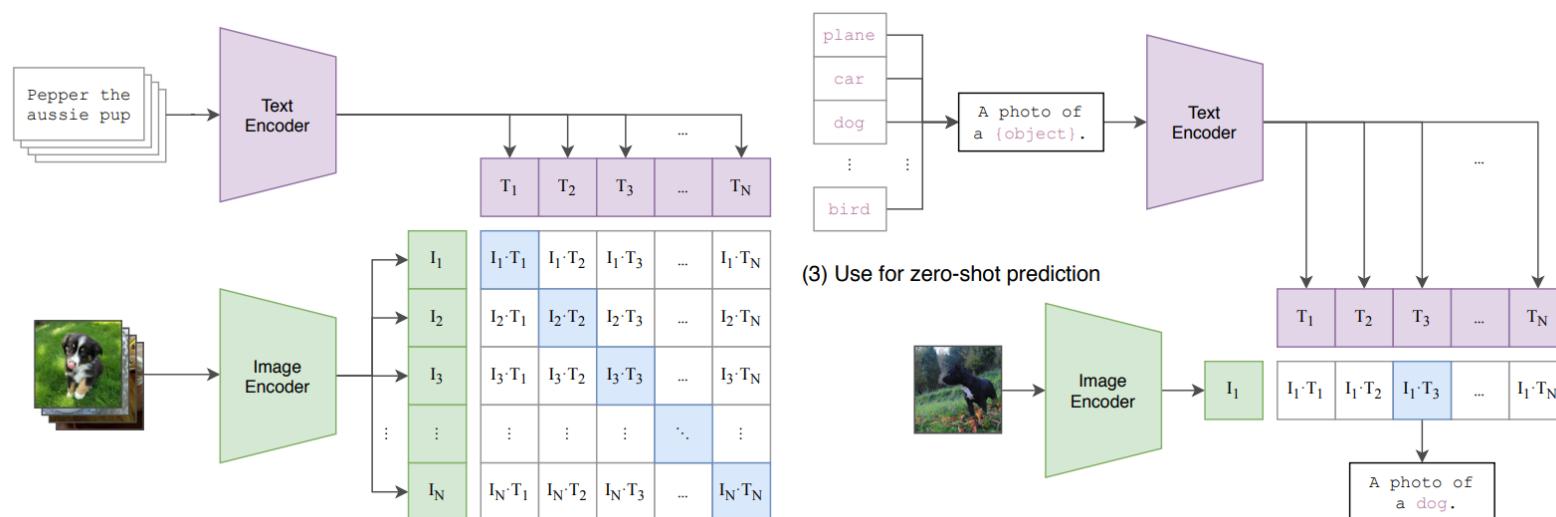
## Emerging Properties in Self-Supervised Vision Transformers

- Training Vision Transformers with Self-Supervised Learning
  - Attention tokens learn informative class features
  - Applying kNN to learnt representation achieves 78.3% top-1 accuracy on ImageNet, without finetuning



# Multi-Modal Representations

- CLIP aligns text representations with visual representations
- A different form of supervision using image captions





## Some thoughts

- SSL has been behind all the recent AI hype
- Huge quantities of unlabelled data have allowed for more “intelligent” systems
- Existing systems are (presumably) learning from most of human knowledge
  - How will this scale in the future?

# Announcements

- We'll announce which research paper readings are examinable next week (recap session)
- SET evaluations
  - Set evaluations are out for a couple of weeks
  - Please give us useful feedback (including what you liked!)