

Transfer Learning: Self-Supervised Learning

Danna Gurari

University of Colorado Boulder

Spring 2022



Review

- Last lecture topic:
 - Visual dialog applications
 - Visual dialog dataset
 - Visual dialog evaluation
 - Mainstream 2017 challenges: baseline approaches
- Assignments (Canvas)
 - Lab assignment 4 due Friday
 - Final project proposal due next week
- Questions?

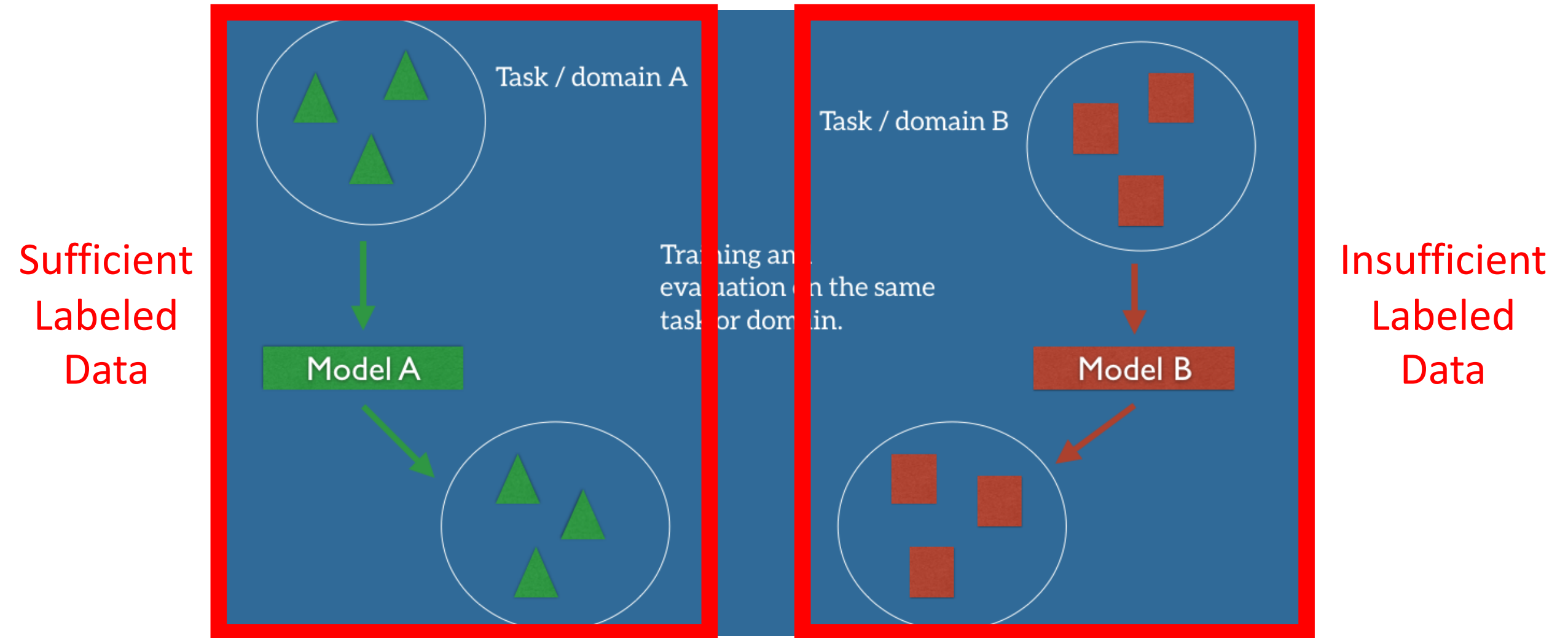
Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- Generative-based methods
- Generative adversarial networks
- Context-based methods

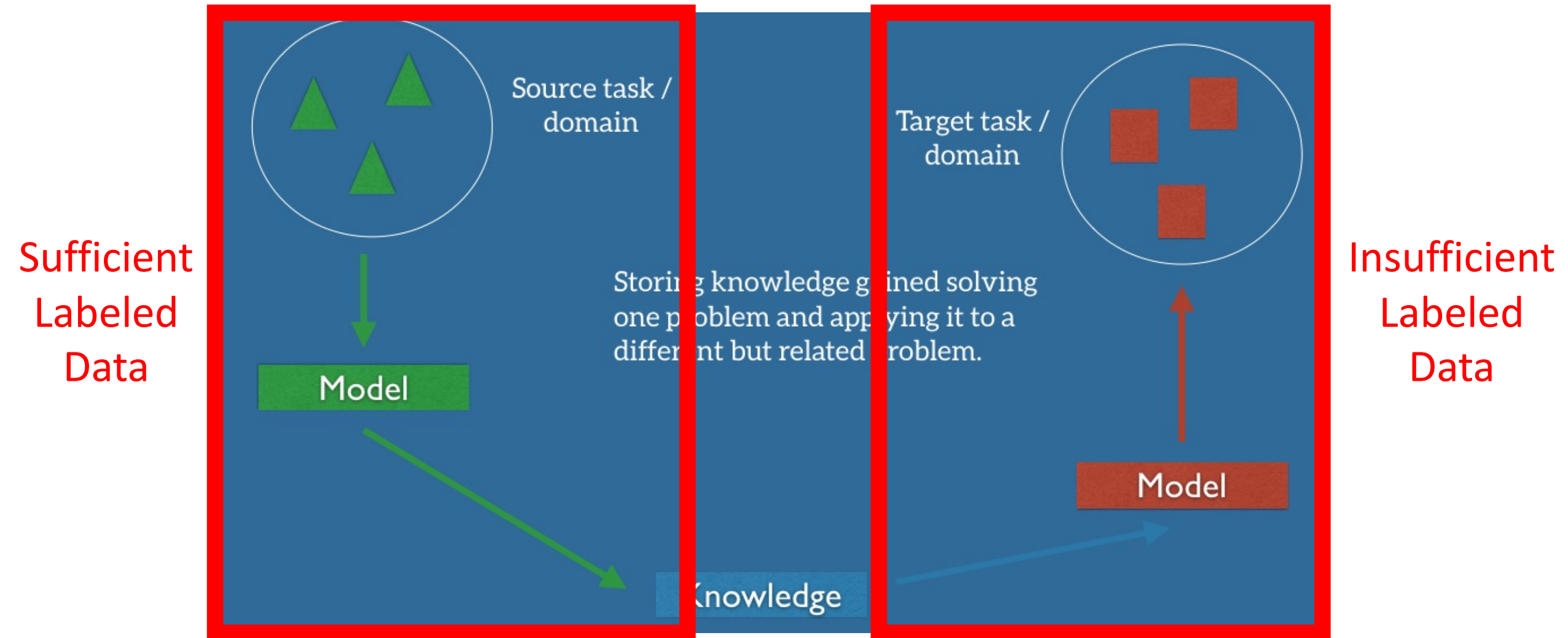
Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- Generative-based methods
- Generative adversarial networks
- Context-based methods

Rather than Learn Solution from Scratch For Each Task/Domain Pair... (Problem for B)



Idea: Improve the Learning for Conditions Not Observed During Training



Transfer Learning When Data Sampling Changes (e.g., Sentiment Classification)



News (formal and lengthy)



Tweets (informal and brief)

Transfer Learning When Feature Space Changes (e.g., Sentiment Classification in Different Language)

★★★★★ Cool charger

By Tiffany on March 30, 2015

Verified Purchase

Bought this for my Galaxy phone and I have to say, this is a pretty cool USB cord! :) I like the lights in the cord as it puts off a cool glowing effect in my room at night and it makes it much easier to see, thanks for the great product!

★★★★★ Definitely buying more.

By Krystal Willingham on March 28, 2015

Verified Purchase

I was impressed with how bright the lights on the cable are. It works amazing and as described. I received earlier than expected so that made me very happy. So far is working like a charm and I can't wait to buy a few more.

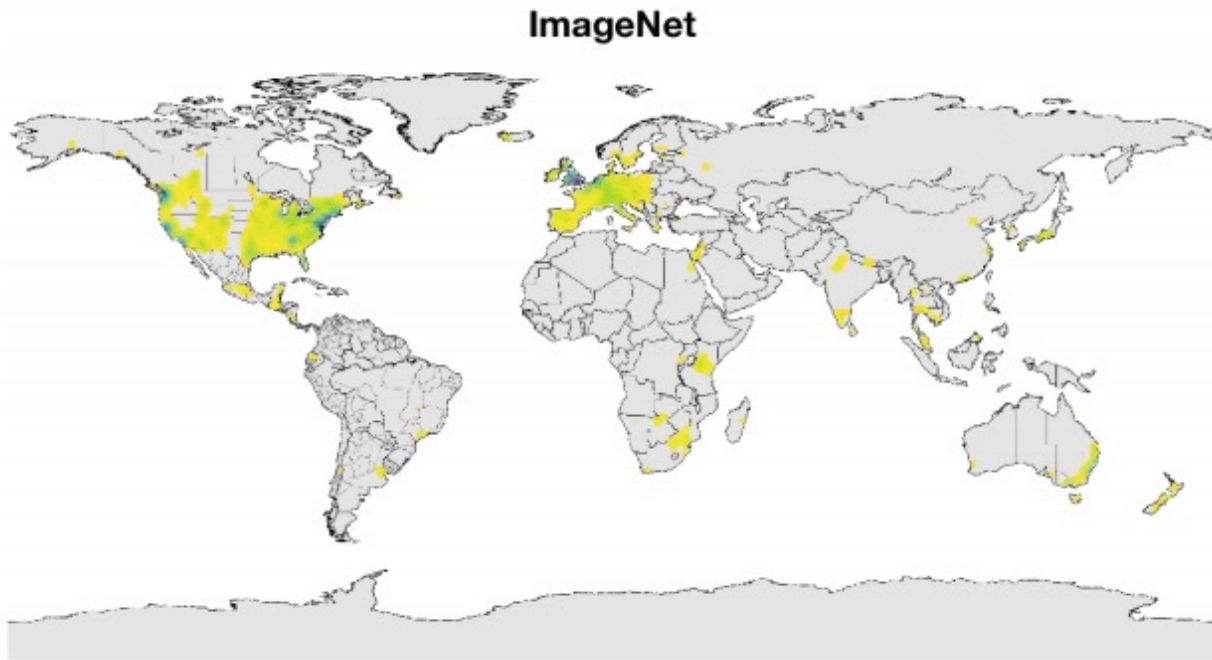
★★★★★ Spot It In the Crowd

By Heather-Joan Carls on March 29, 2015

Verified Purchase

Such a cool product. I was so happy with how bright the lights on the cable are. It shipped super fast. The light shuts off when the charging is complete, so that's super helpful. I don't have to keep checking.

Transfer Learning When Target Categories Change (e.g., Items in Low Income Household vs ImageNet)



Zhao et al. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. 2017.



Ground truth: Soap Nepal, 288 \$/month

Azure: food, cheese, bread, cake, sandwich
Clarifai: food, wood, cooking, delicious, healthy
Google: food, dish, cuisine, comfort food, spam
Amazon: food, confectionary, sweets, burger
Watson: food, food product, turmeric, seasoning
Tencent: food, dish, matter, fast food, nutriment

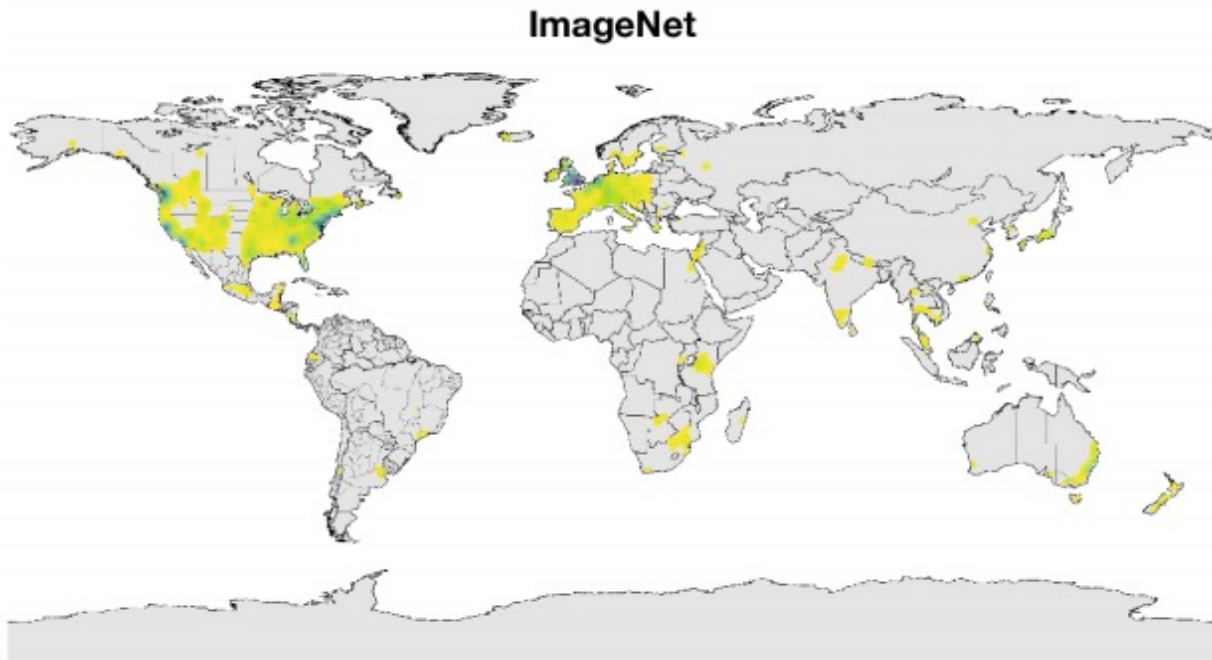


Ground truth: Soap UK, 1890 \$/month

Azure: toilet, design, art, sink
Clarifai: people, faucet, healthcare, lavatory, wash closet
Google: product, liquid, water, fluid, bathroom accessory
Amazon: sink, indoors, bottle, sink faucet
Watson: gas tank, storage tank, toiletry, dispenser, soap dispenser
Tencent: lotion, toiletry, soap dispenser, dispenser, after shave

DeVries et al. Does object recognition work for everyone? CVPR workshops, 2019.

Transfer Learning When Limited Data Available (e.g., Items in Low Income Household vs ImageNet)



Zhao et al. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. 2017.



Ground truth: Soap Nepal, 288 \$/month

Azure: food, cheese, bread, cake, sandwich
Clarifai: food, wood, cooking, delicious, healthy
Google: food, dish, cuisine, comfort food, spam
Amazon: food, confectionary, sweets, burger
Watson: food, food product, turmeric, seasoning
Tencent: food, dish, matter, fast food, nutriment

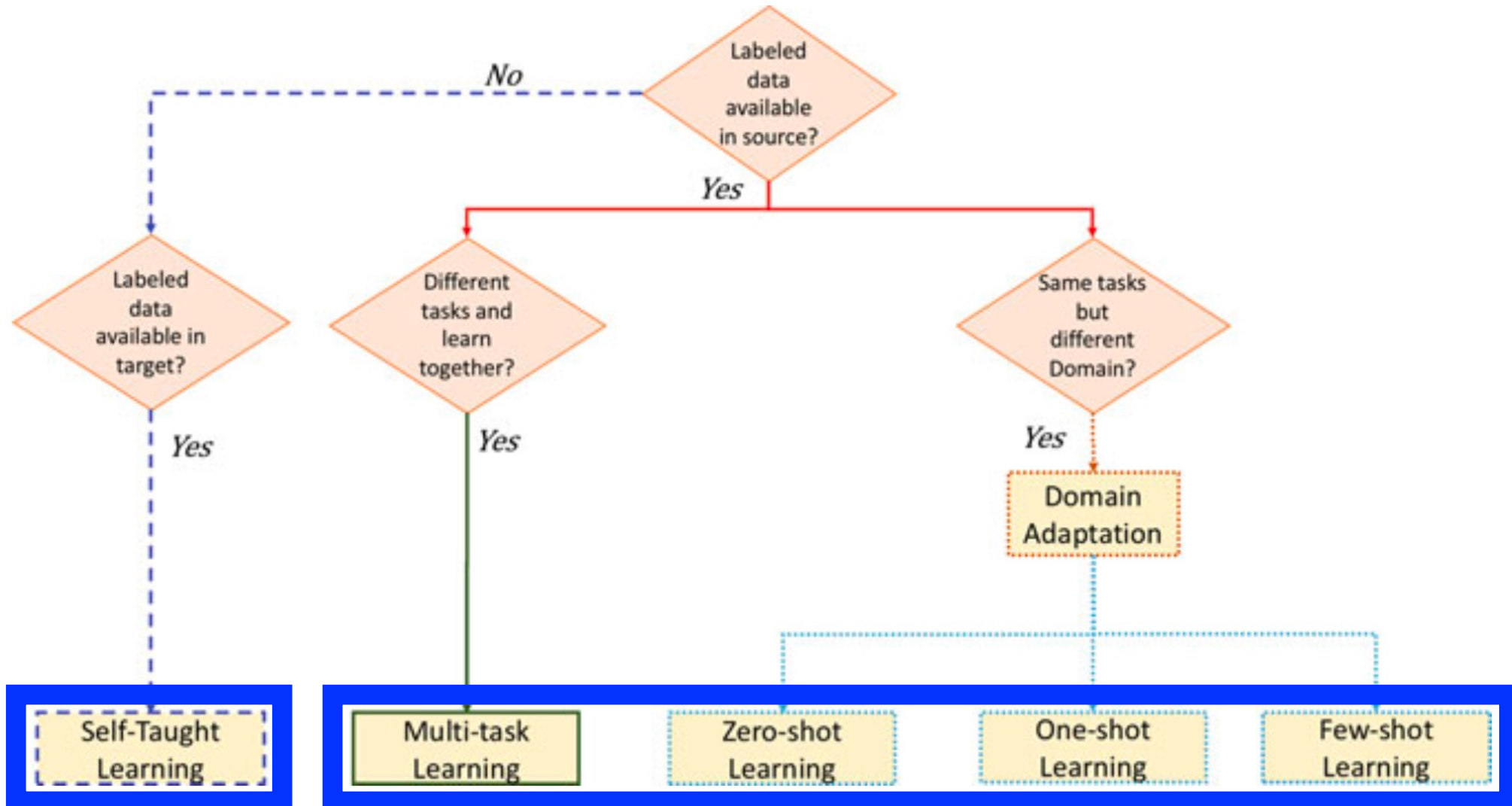


Ground truth: Soap UK, 1890 \$/month

Azure: toilet, design, art, sink
Clarifai: people, faucet, healthcare, lavatory, wash closet
Google: product, liquid, water, fluid, bathroom accessory
Amazon: sink, indoors, bottle, sink faucet
Watson: gas tank, storage tank, toiletry, dispenser, soap dispenser
Tencent: lotion, toiletry, soap dispenser, dispenser, after shave

DeVries et al. Does object recognition work for everyone? CVPR workshops, 2019.

Transfer Learning Approaches



Today's
focus

Next
lecture

Transfer Learning: Key Challenges

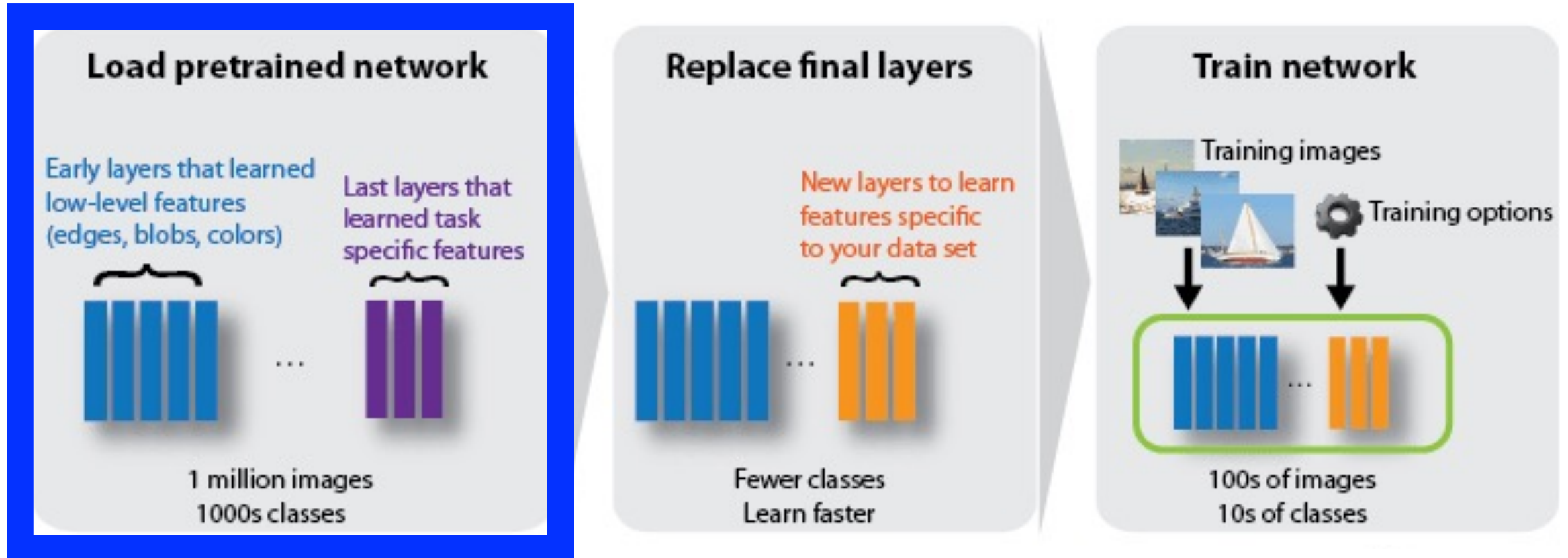
- **What to transfer?** i.e., what knowledge generalizes
- **How to transfer?**
- **When to transfer?** i.e., transferring knowledge can harm performance

Today's Topics

- Transfer learning definition
- **Overview of self-supervised learning**
- Generative-based methods
- Generative adversarial networks
- Context-based methods

Goal: Create Generalizable Features

Key observation: features from a pretrained network can be useful for other datasets/tasks



Intuition: How Do Humans Learn?

With Supervision

Learn from instruction



Unsupervised

Learn from experience



Today's
scope

<https://pixabay.com/en/toddler-learning-book-child-423227/>

<https://www.maxpixel.net/Father-Child-Family-Dad-Baby-Daughter-3046495>

Self-Supervised Learning: Data Gives Supervision

- Relatively Cheap
- Can Collect Data Fast



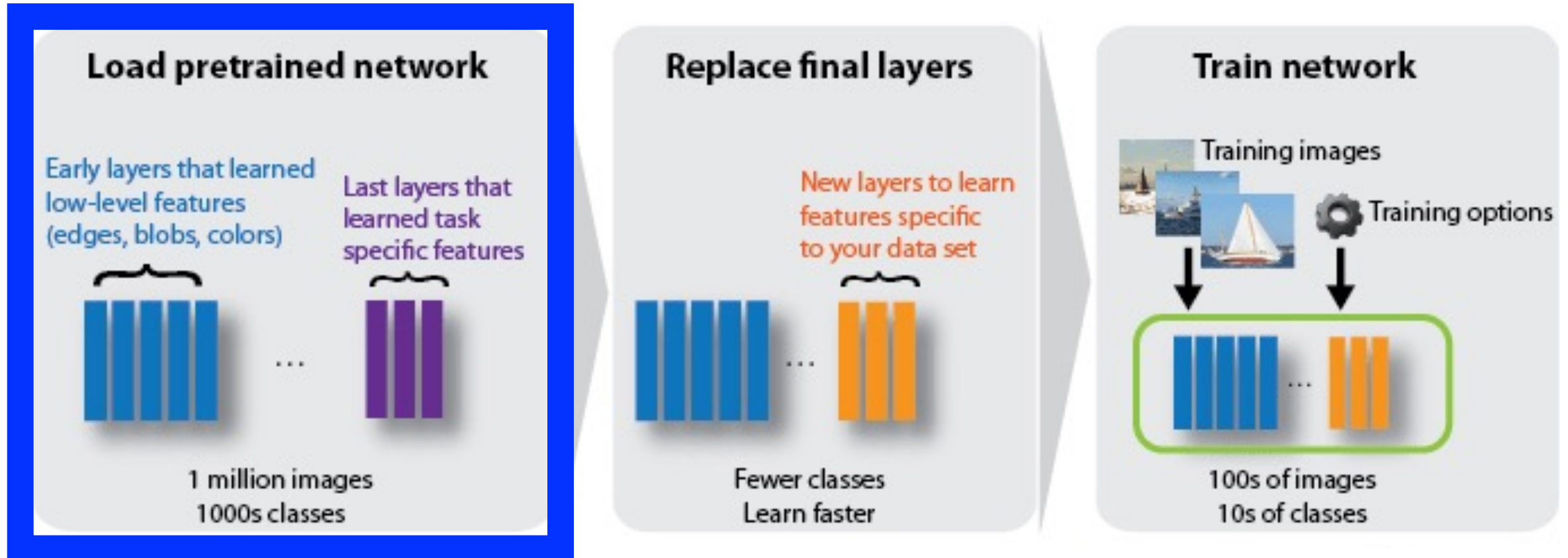
<https://lovevery.com/community/blog/child-development/the-surprising-learning-power-of-a-household-mirror/>



<https://www.rockettes.com/blog/how-to-use-the-mirror-in-dance-class/>

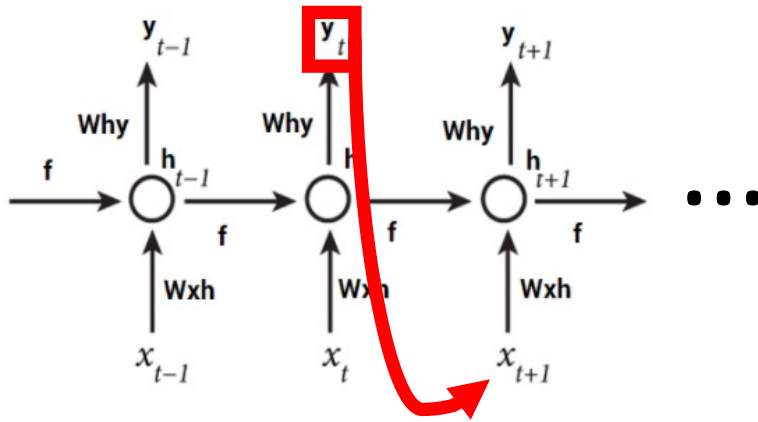
Self-Supervised Learning: Data Gives Supervision

Approach: create features that are useful for other datasets/tasks



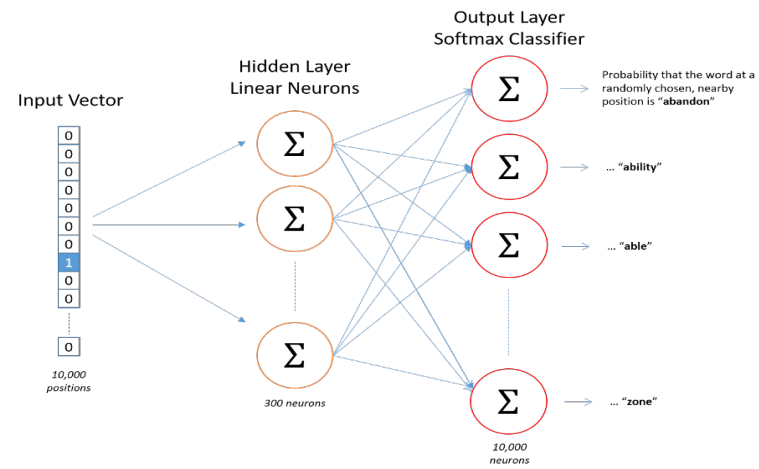
Self-Supervised Learning Methods Already Covered in This Course (Many NLP Methods)

Character prediction with RNNs



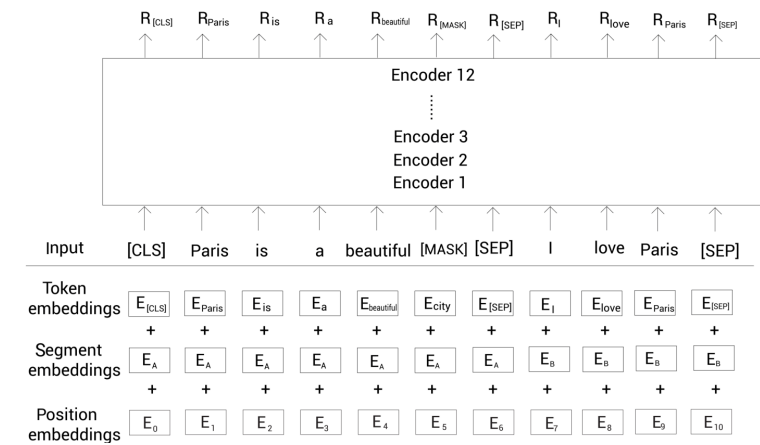
<https://www.analyticsvidhya.com/blog/2017/12/introduction-to-recurrent-neural-networks/>

Word embeddings
(e.g., word2vec; predict nearby word for given word)



<https://towardsdatascience.com/word2vec-skip-gram-model-part-1-intuition-78614e4d6e0b>

Transformers
(e.g., BERT and LXMERT with masking)



https://static.packt-cdn.com/downloads/9781838821593_ColorImages.pdf

Next: self-supervised learning methods
explored in computer vision to learn
visual features which generalize

Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- **Generative-based methods**
- Generative adversarial networks
- Context-based methods

Generative-based Methods

- **Autoencoder**: predict self
- **Colorization**: convert grayscale to color
- **Video prediction**: predict future frames

Generative-based Methods

- **Autoencoder**: predict self
- Colorization: convert grayscale to color
- Video prediction: predict future frames

Image Autoencoder Architecture

- Learn to copy the input to the output

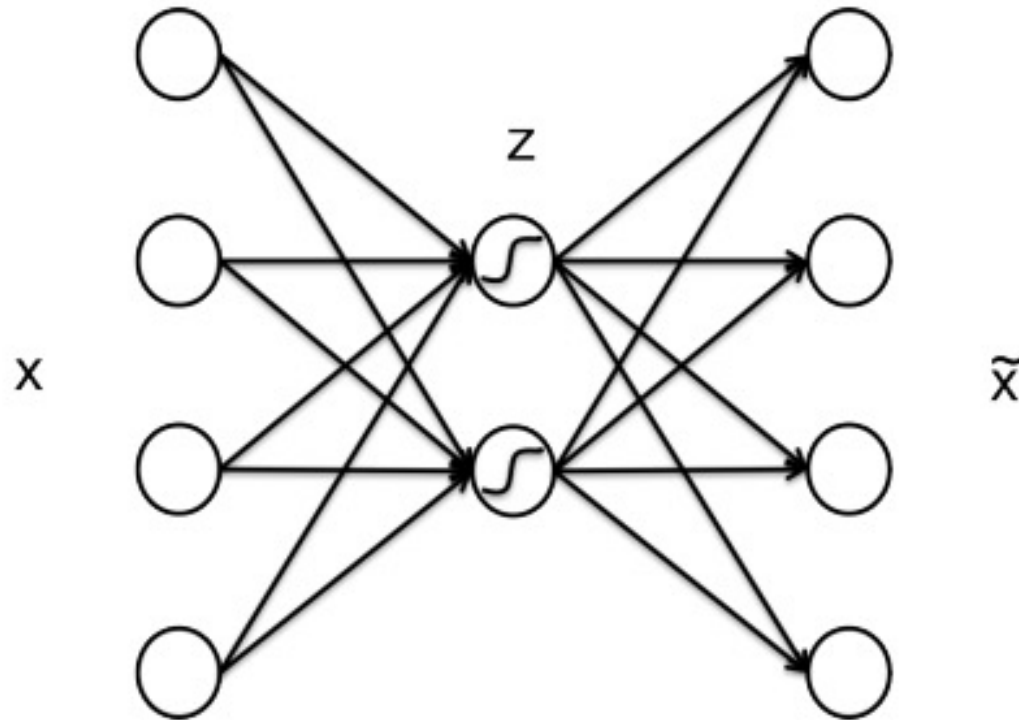


Image Autoencoder Architecture

- Consists of two parts:
 - **Encoder:** compresses inputs to an internal representation
 - **Decoder:** tries to reconstruct the input from the internal representation

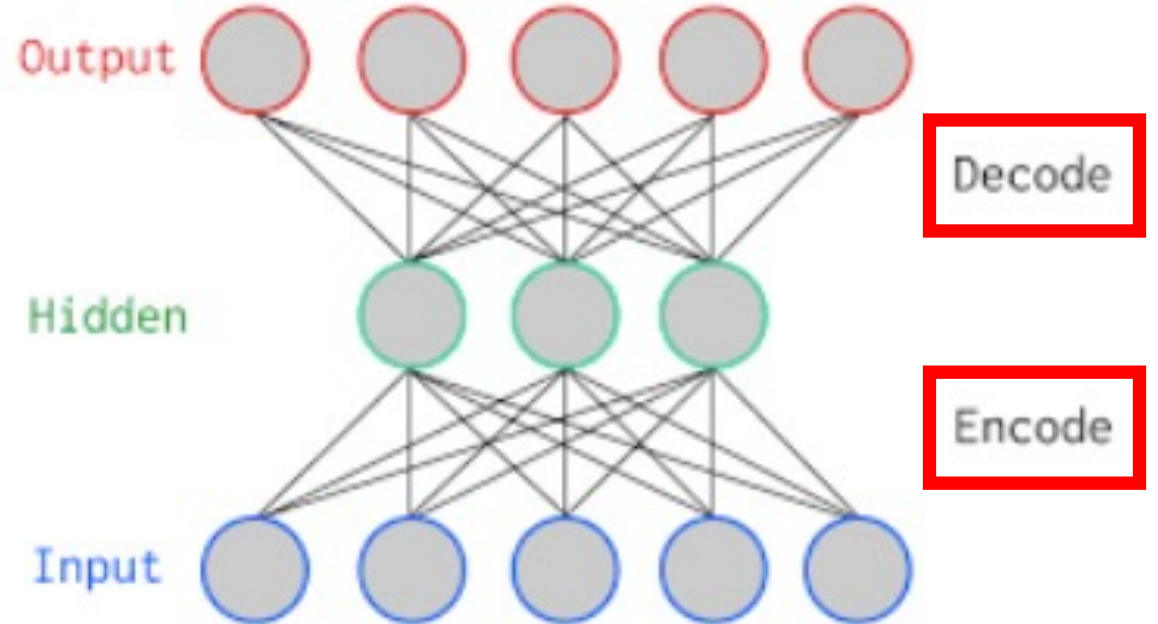


Image Autoencoder Architecture

- Given this input 620 x 426 image (264,120 pixels):



- What would a perfect autoencoder predict?
 - Itself
- What number of nodes are in the final layer?
 - 264,120

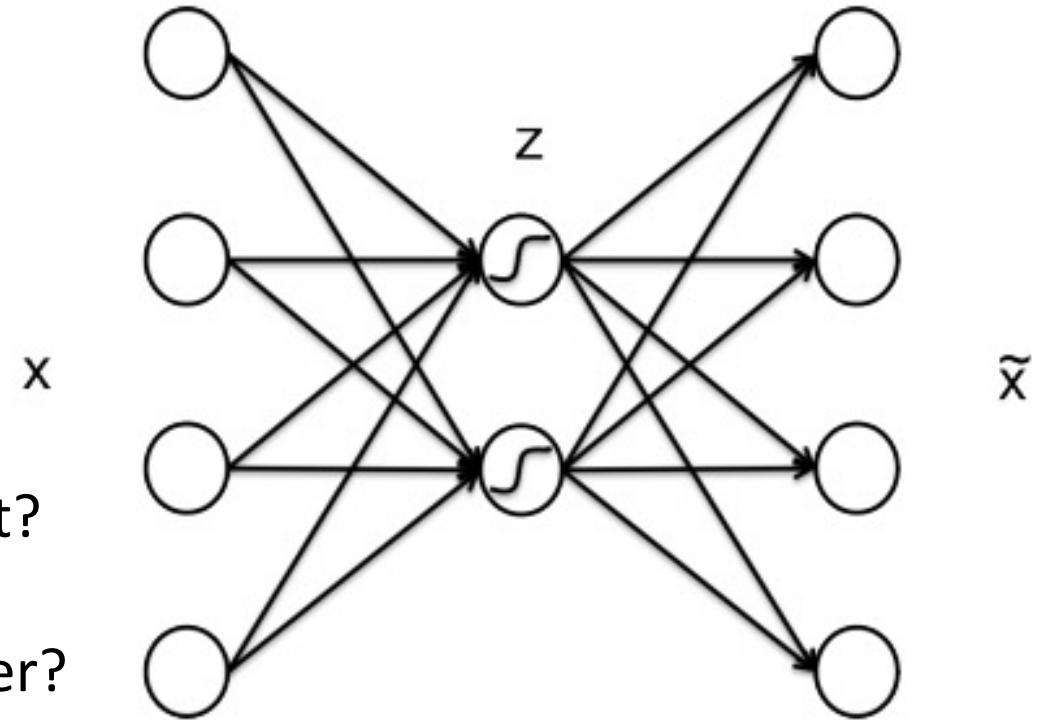


Image Autoencoders

- Intuition: which number sequence is easier to remember?
 - **A:** 30, 27, 22, 11, 6, 8, 7, 2
 - **B:** 30, 15, 46, 23, 70, 35, 106, 53, 160, 80, 40, 20, 10, 5
- **B:** need learn only two rules
 - If even, divide by 2
 - If odd, multiply by 3 and add 1

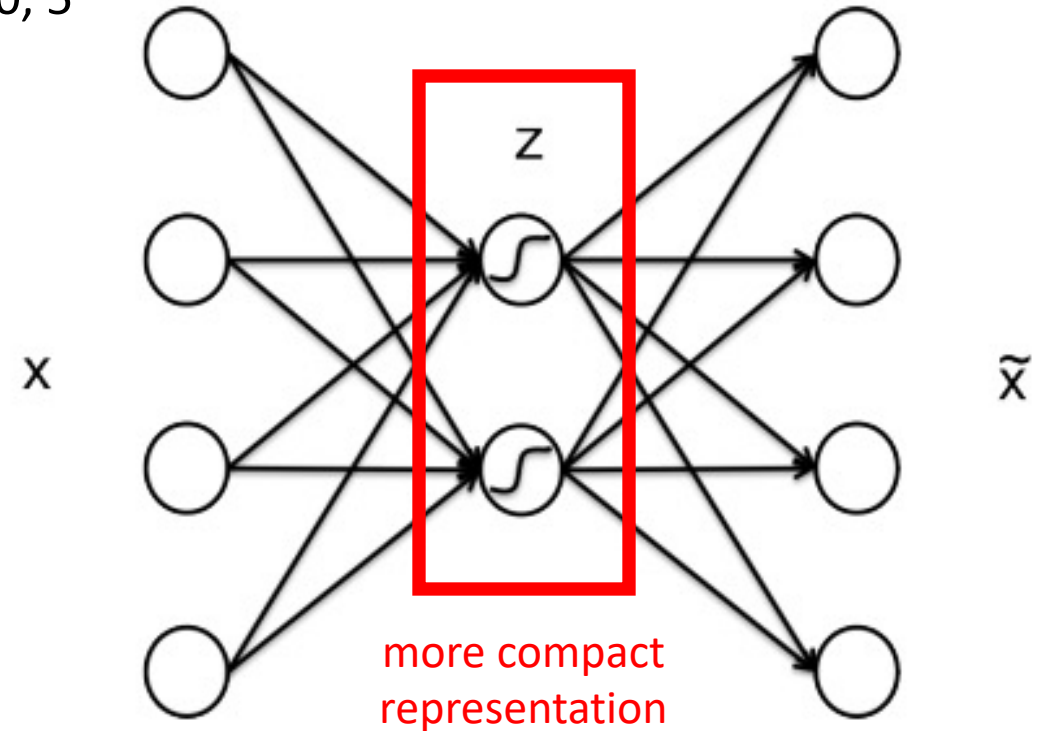


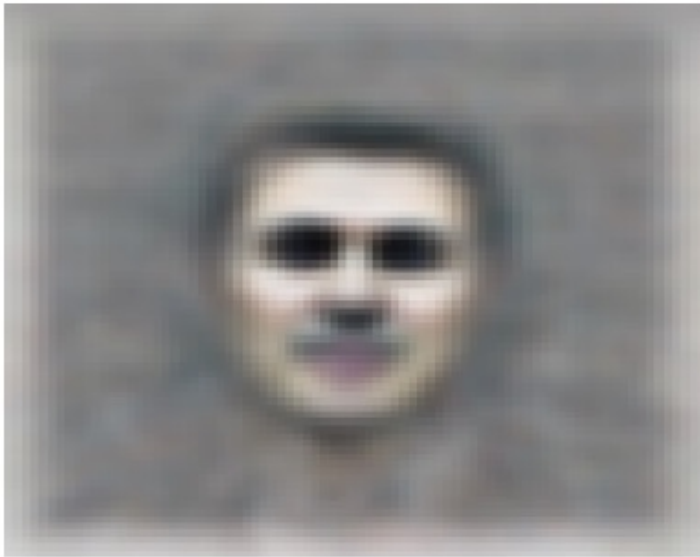
Image Autoencoder Training

Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through network to make prediction
2. **Backward pass:** using predicted output, calculate error gradients backward
3. Update each weight using calculated gradients

Image Autoencoder Features

- e.g., features learned include:



human face



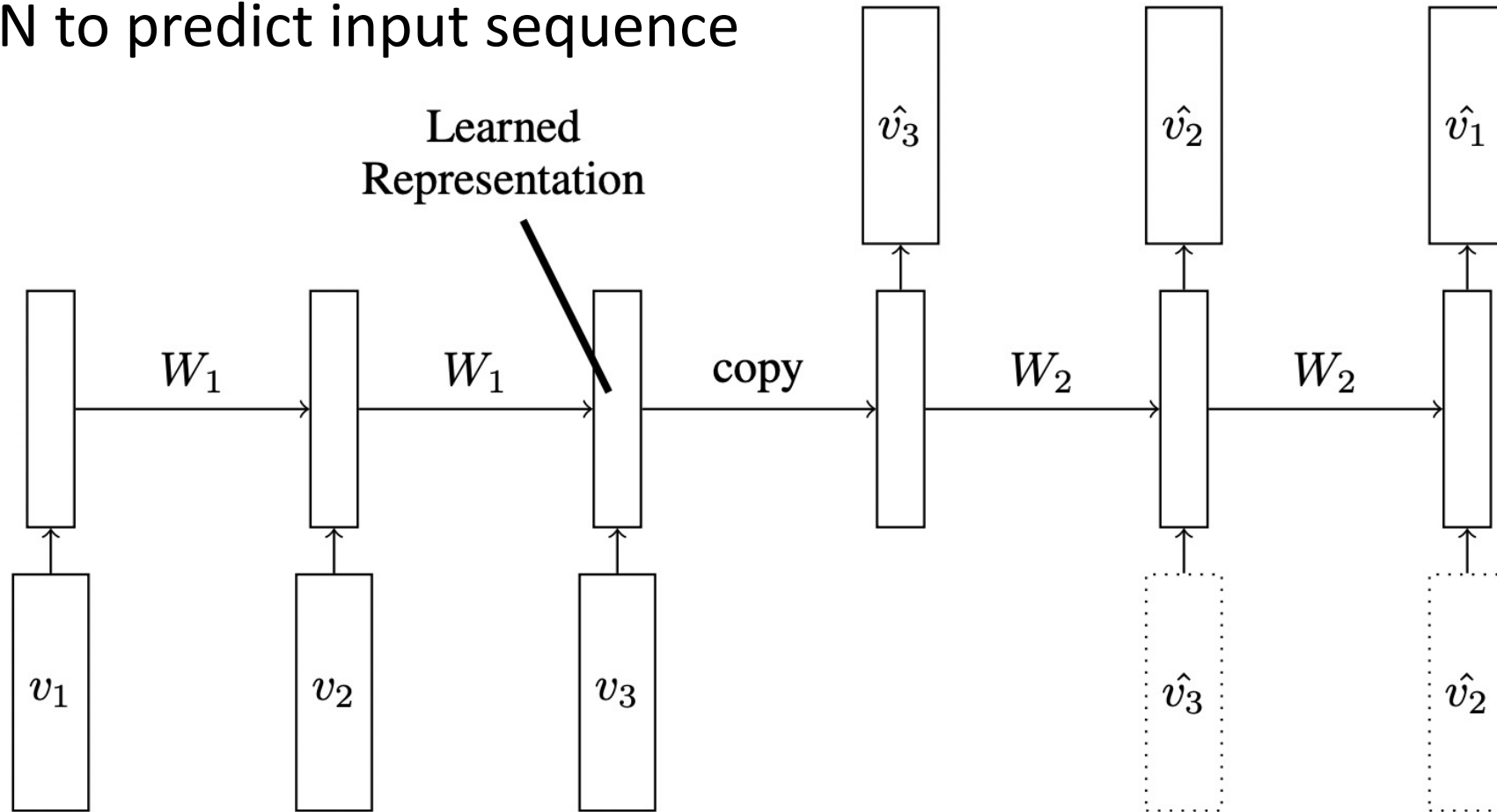
cat face



human body

Video Autoencoder

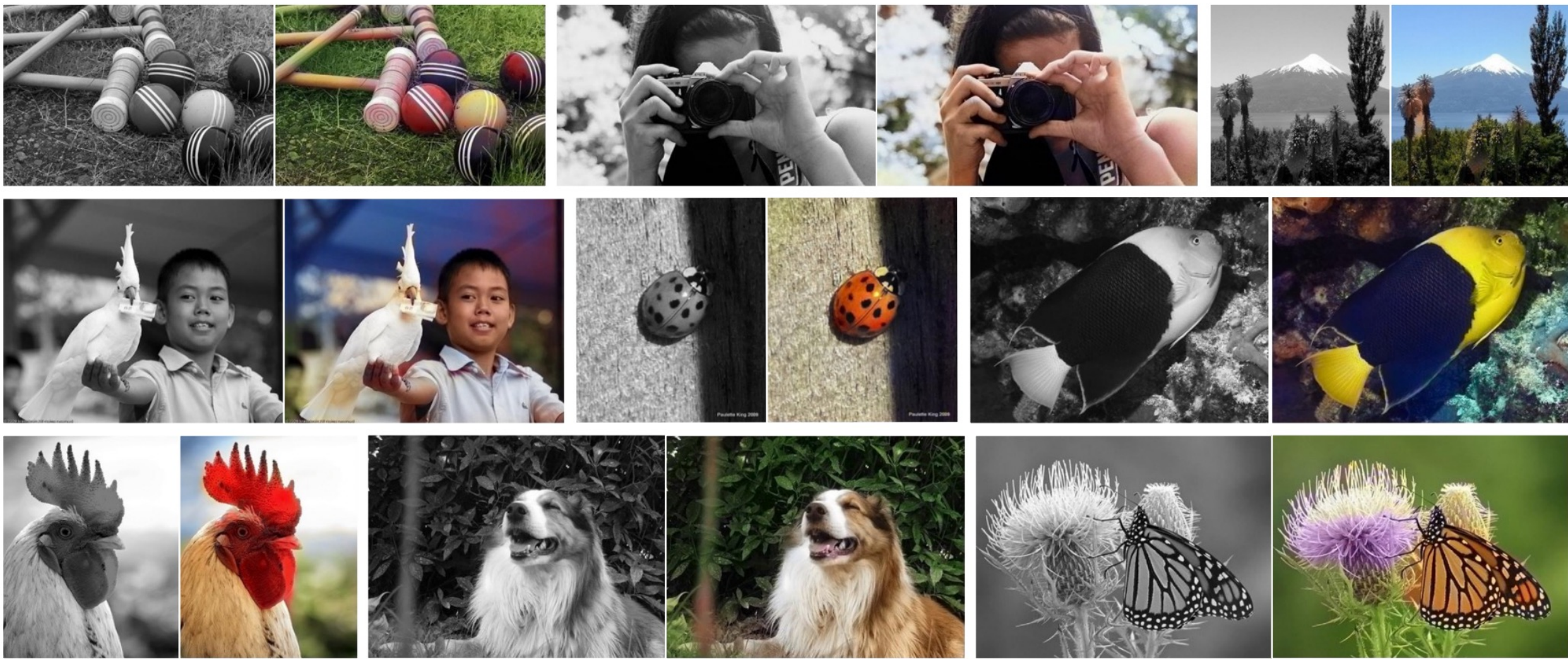
- Train RNN to predict input sequence



Generative-based Methods

- Autoencoder: predict self
- **Colorization**: convert grayscale to color
- Video prediction: predict future frames

Colorization: *Plausible* Coloring Results



Colorization: *Plausible* Coloring Results



Figure Sources: [https://www.flickr.com/photos/applesnpearsau/12197380673/in/photostream/;](https://www.flickr.com/photos/applesnpearsau/12197380673/in/photostream/)
https://commons.wikimedia.org/wiki/File:JACQUES_VILET_-_1982,_Les_Fruits_du_Jardin.jpg

Image Colorization Architecture

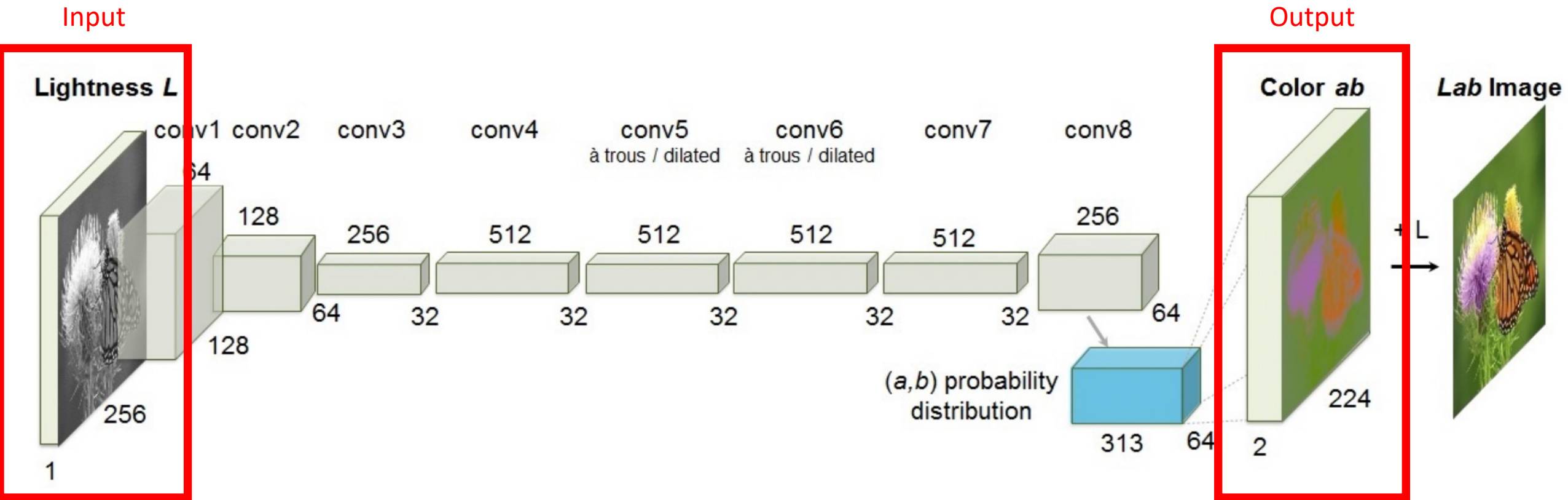
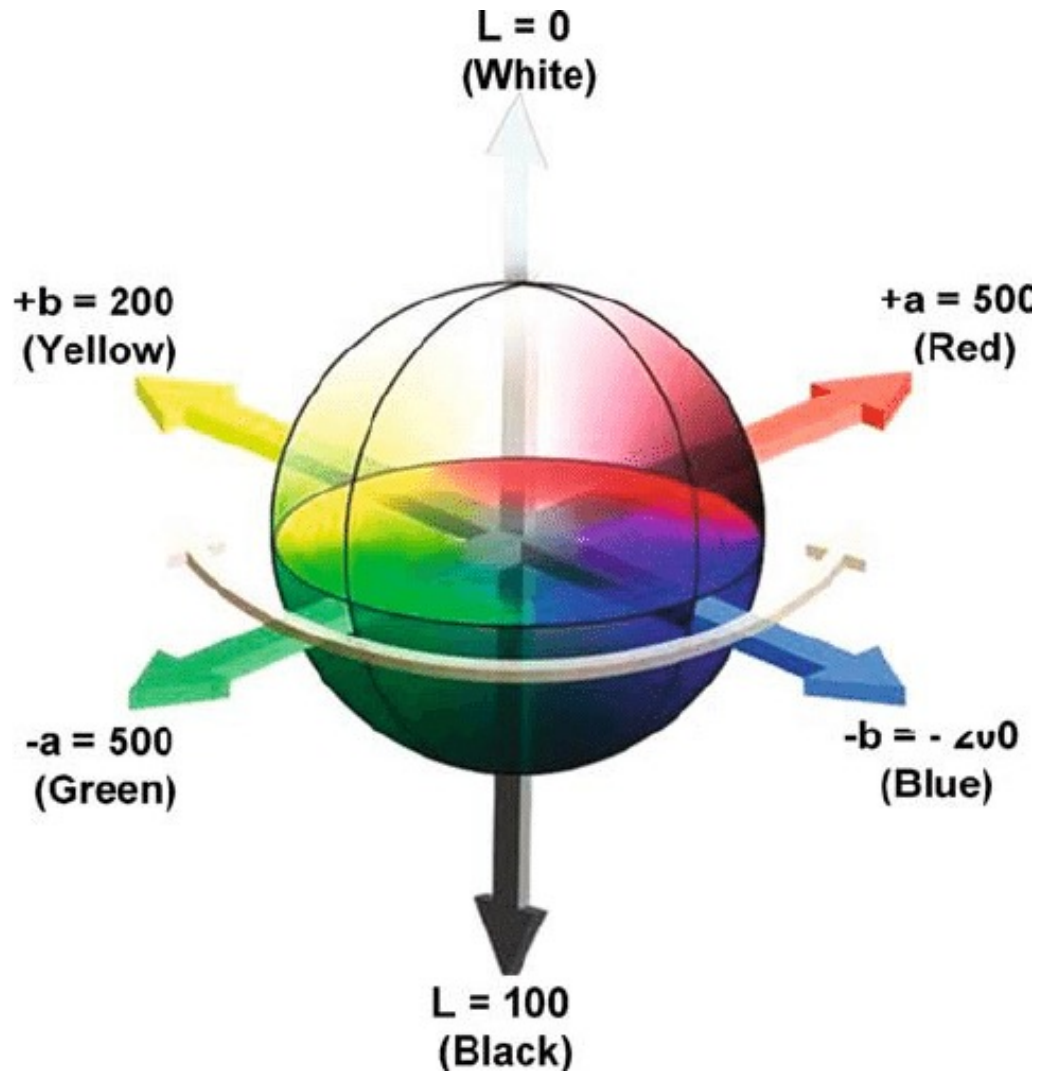
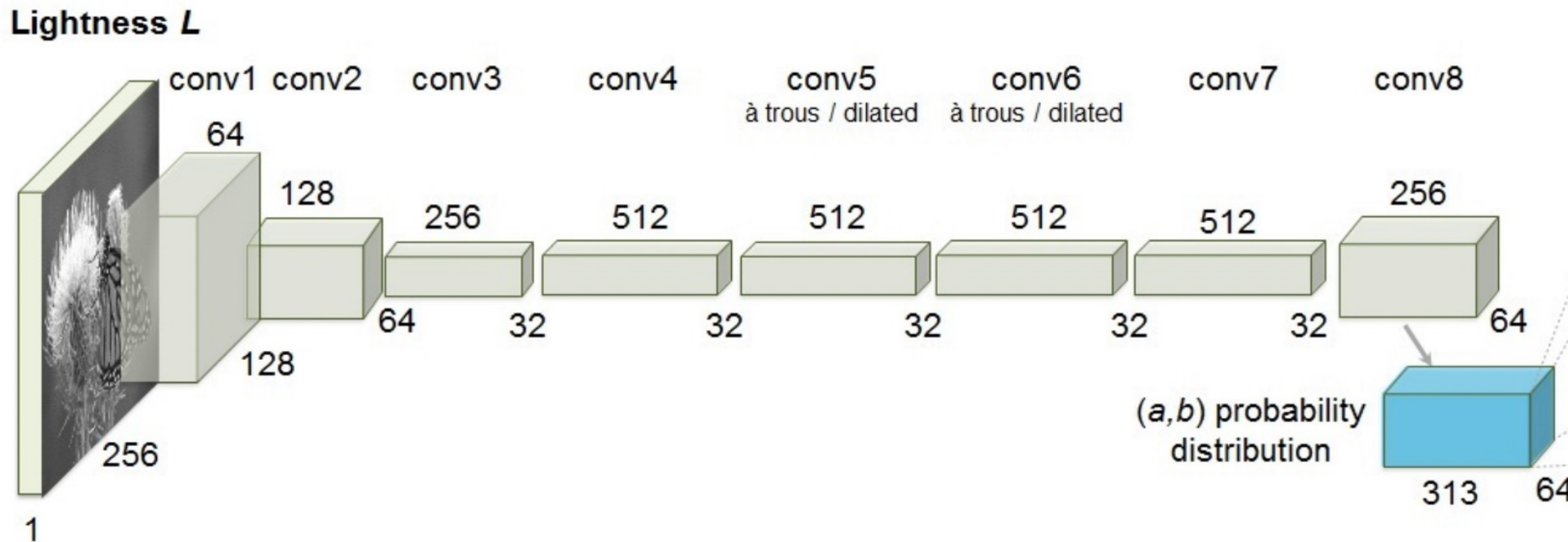


Image Colorization Architecture: CIE *Lab* Color



L indicates grayscale information whereas a and b represent colors

Image Colorization Architecture



Create image by combining predicted a and b channels with the L channel

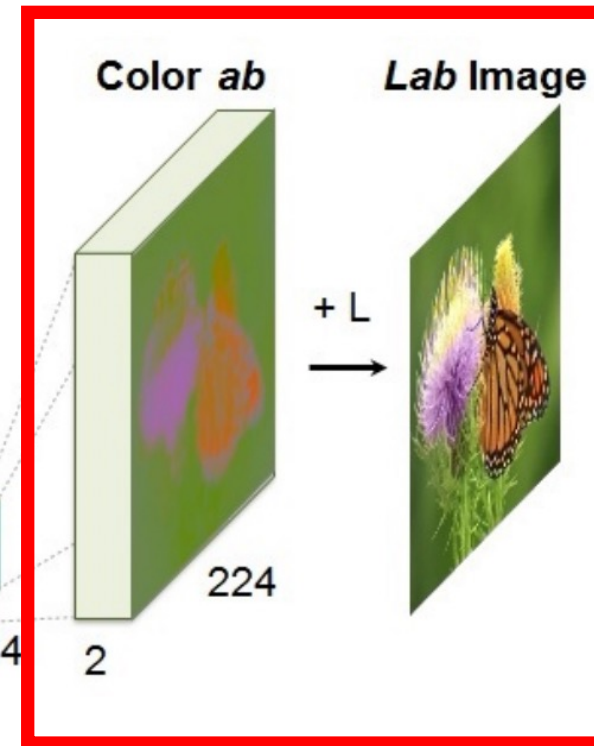


Image Colorization Architecture



Grayscale image: L channel

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

L

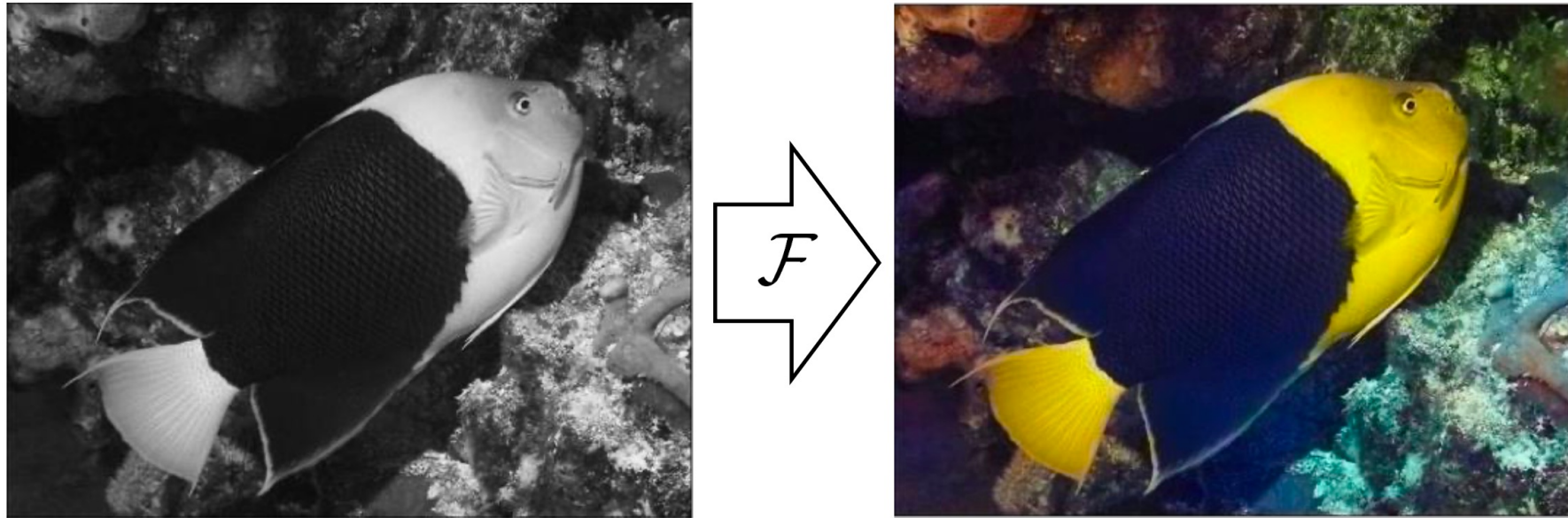


Color information: ab channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

ab

Image Colorization Architecture



Grayscale image: L channel

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

Concatenate (L,ab)

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

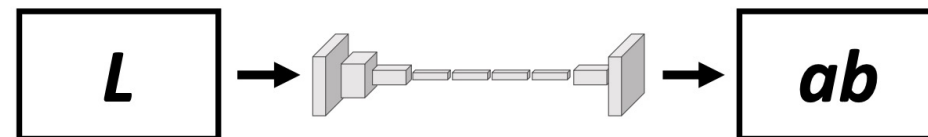


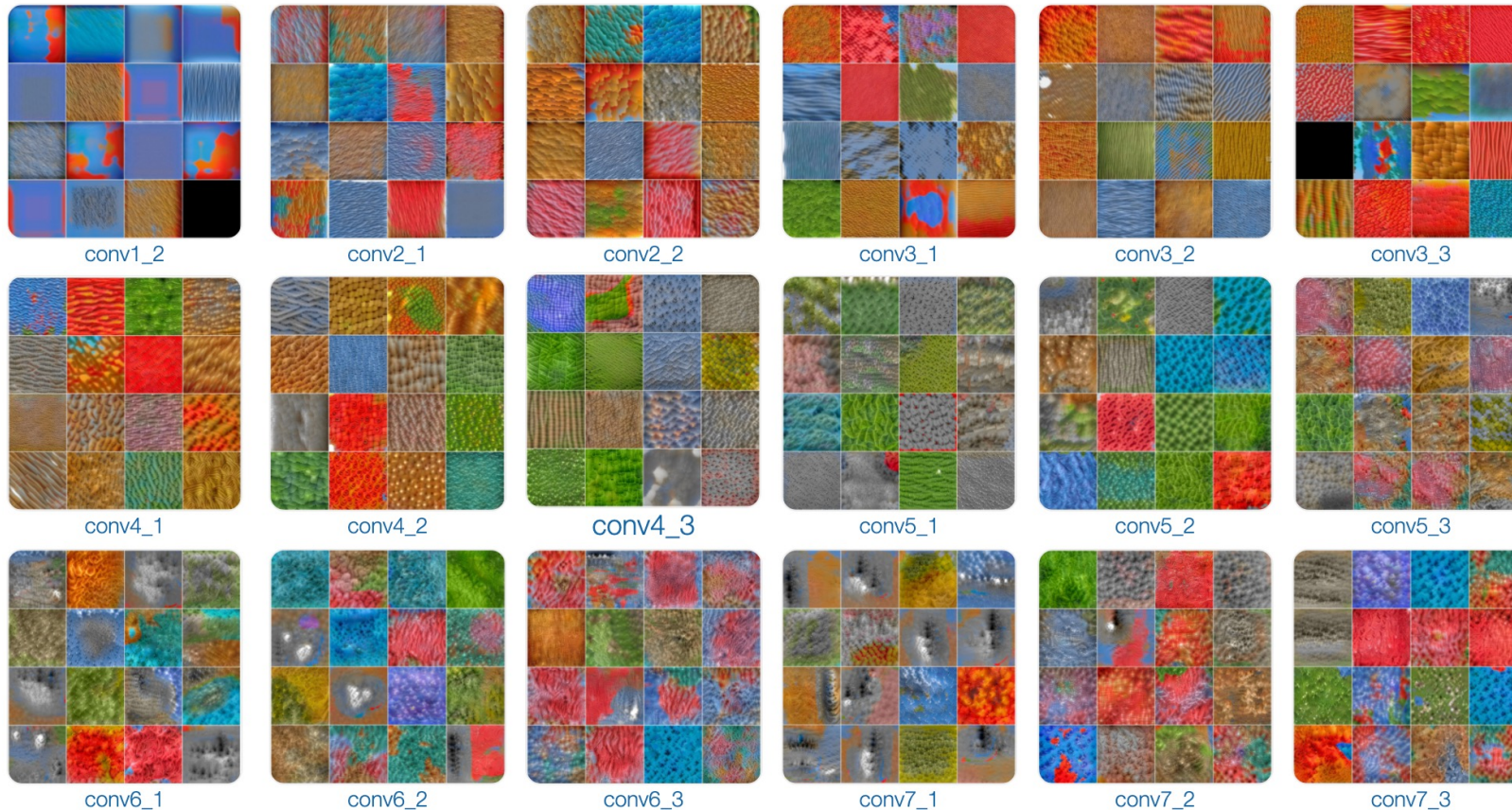
Image Colorization Training

For 1.3 million ImageNet images, repeat until stopping criterion met:

1. **Forward pass:** propagate training data through network to make prediction
2. **Backward pass:** using predicted output, calculate error gradients backward
3. Update each weight using calculated gradients

Image Colorization Features

Task requires understanding an image at the pixel and semantic-level

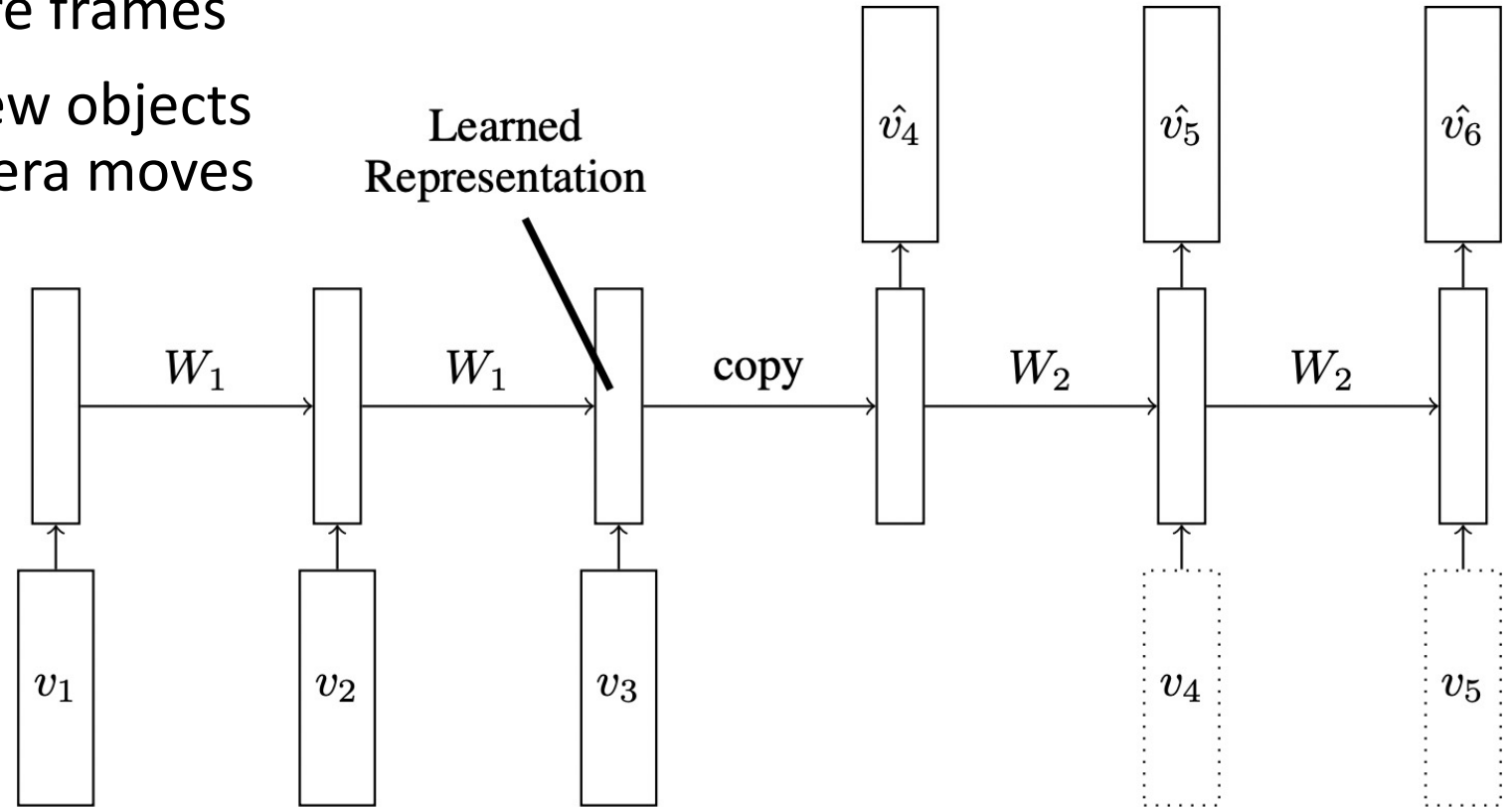


Generative-based Methods

- Autoencoder: predict self
- Colorization: convert grayscale to color
- **Video prediction**: predict future frames

Video Prediction

- Train RNN to predict future frames
- Limitations: identifying new objects and background as a camera moves



What type of features might be learned?

Generative-based Methods

- Autoencoder: predict self
- Colorization: convert grayscale to color
- Video prediction: predict future frames

Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- Generative-based methods
- **Generative adversarial networks**
- Context-based methods

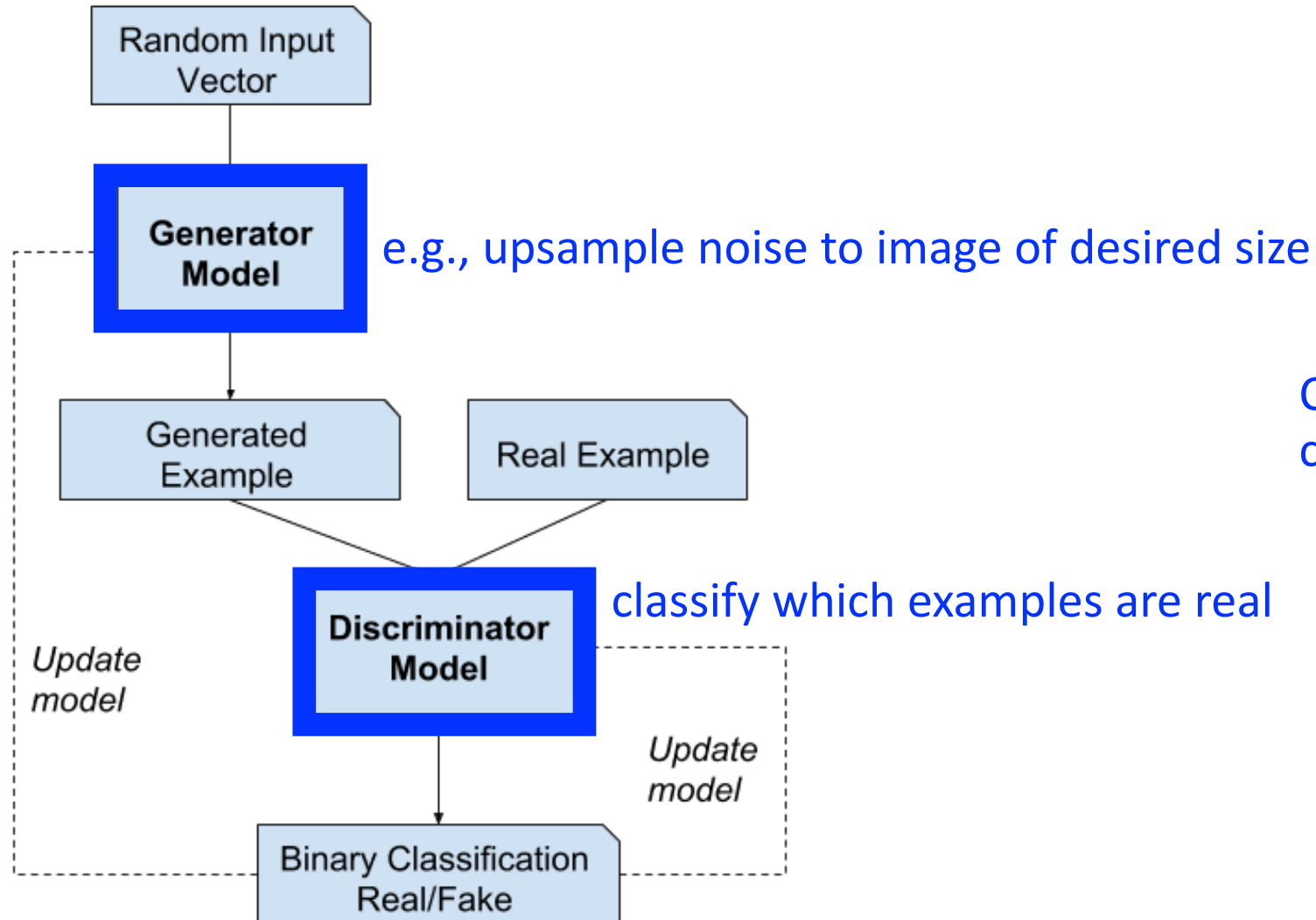
Generative adversarial networks

- Generative adversarial networks (GANs)
- Context encoder

Generative adversarial networks

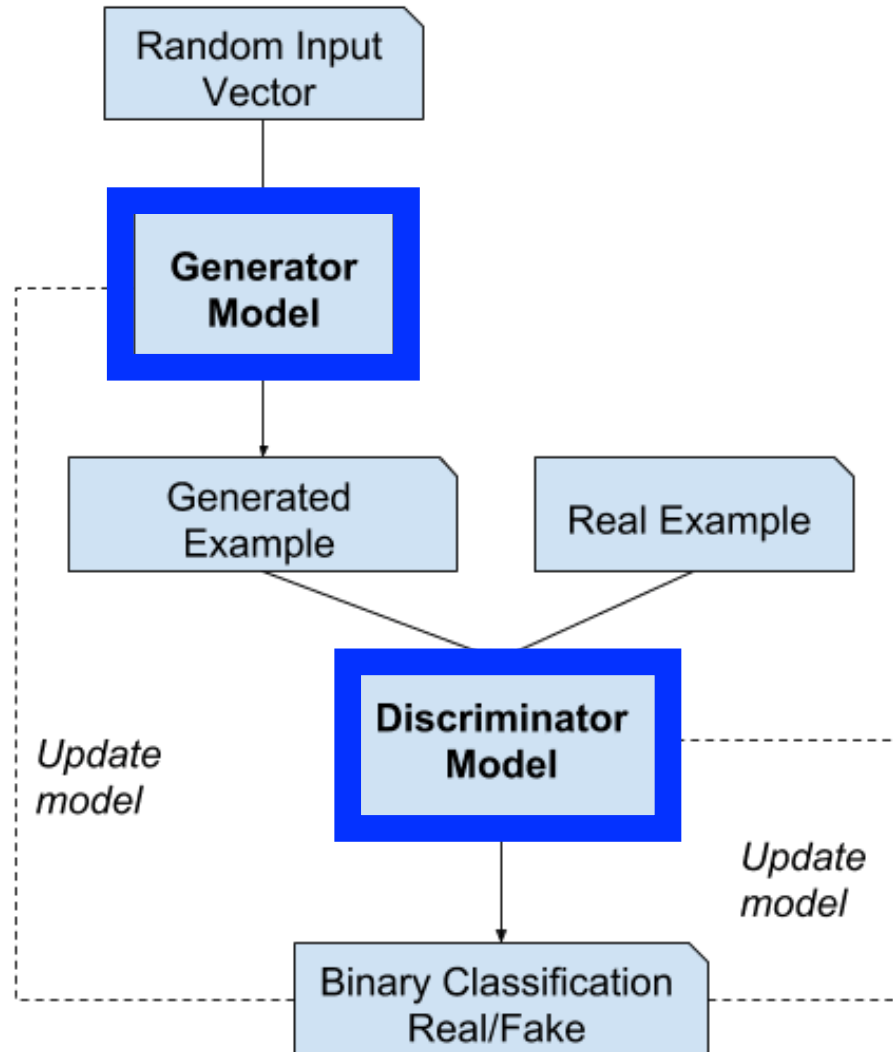
- Generative adversarial networks (GANs)
- Context encoder

GAN: Basic Architecture



Consists of two models that compete against each other

GAN: Training



The two models are iteratively trained separately

- Train discriminator using fake and real images
- Train generator using just fake images and penalize it when the discriminator recognizes images are fake

GAN: Discriminator Loss Function

Discriminator tries to minimize classification error

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

Discriminator wants a value of 1 for real images

Discriminator wants a value of 0 for fake images

Real image

Input noise

GAN: Generator Loss Function

Generator tries to maximize classification error

$$J^{(G)} = -J^{(D)}$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Want the discriminator to mistakenly arrive at a value of 1 for fake images

Input noise

DGANs: GANs that Use Convolutional Layers



Bedrooms generated by observing over 3M bedroom images

DGANs: GANs that Use Convolutional Layers



What objects does it learn to generate?

DGANs: GANs that Use Convolutional Layers



What objects may it not have learned to generate?

DGANs: GANs that Use Convolutional Layers



Faces generated by observing over 3M images of 10K people

DGANs: GANs that Use Convolutional Layers



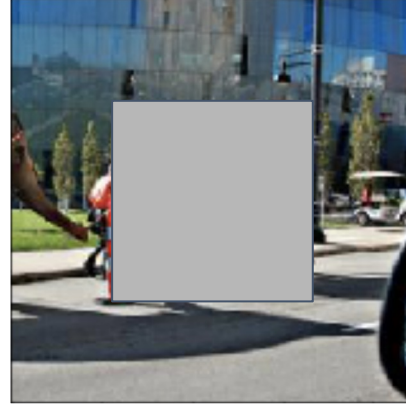
What does it generate poorly or not all?

Generative adversarial networks

- Generative adversarial networks (GANs)
- Context encoder

Task: Hole Filling

- What might fit into this hole?

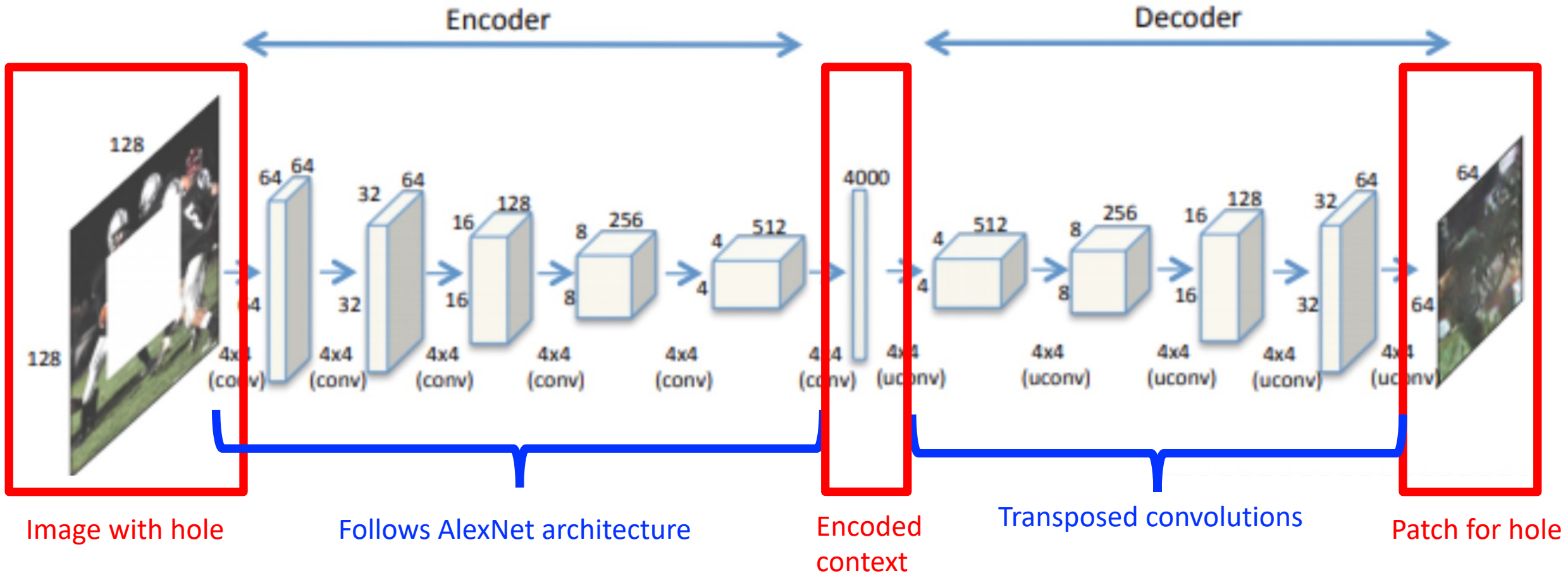


- Many items may plausibly fit into the hole:

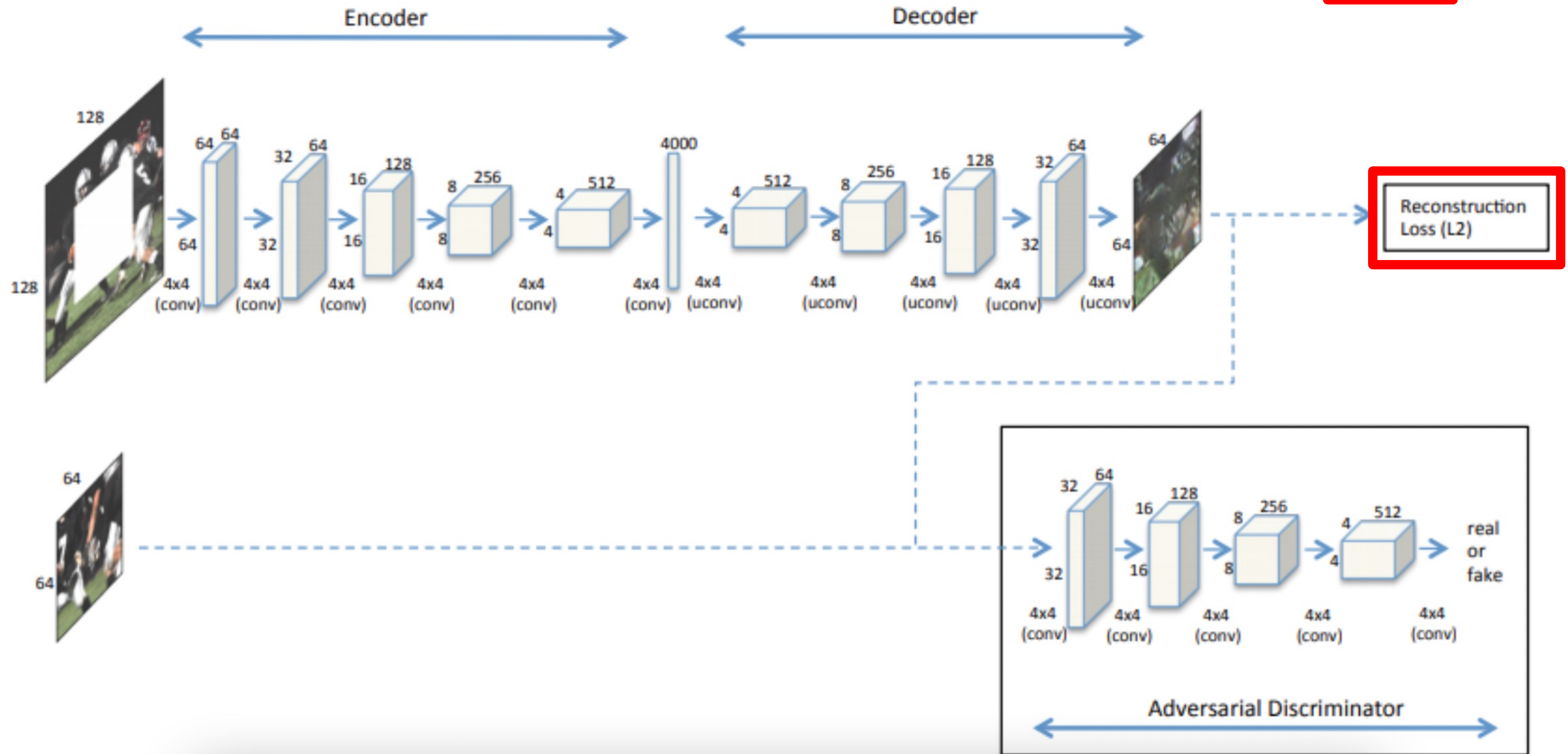


- Challenge: have up to 1 known ground truth region per hole

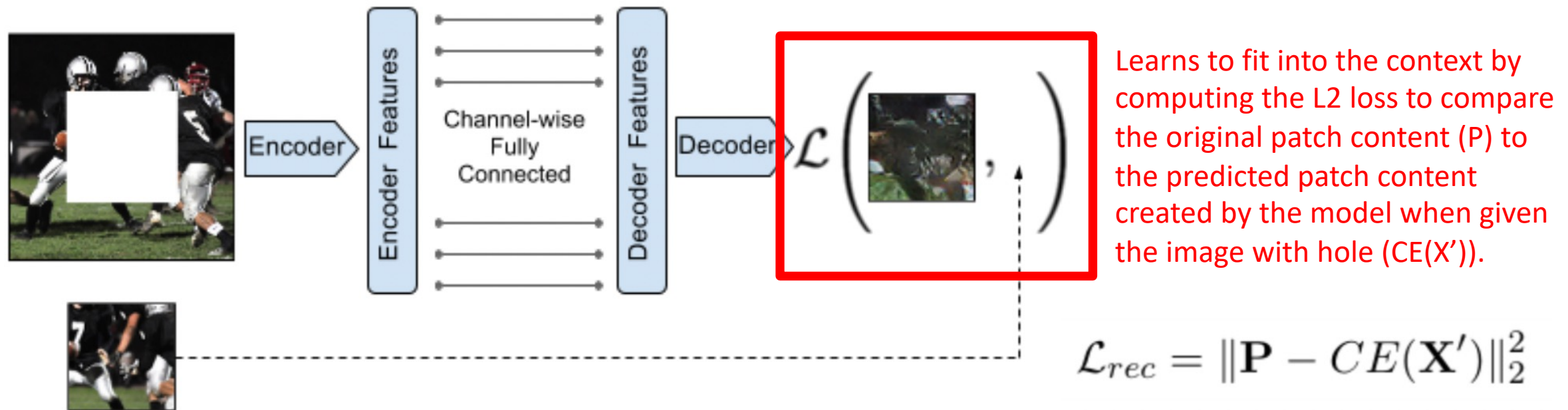
Architecture



Training: Loss Functions ($\mathcal{L} = \lambda_{adv}\mathcal{L}_{adv} + \lambda_{rec}\mathcal{L}_{rec}$)



Training: Reconstruction Loss (i.e., Self-Supervised Learning Approach)



Training: Reconstruction Loss (i.e., Self-Supervised Learning Approach)



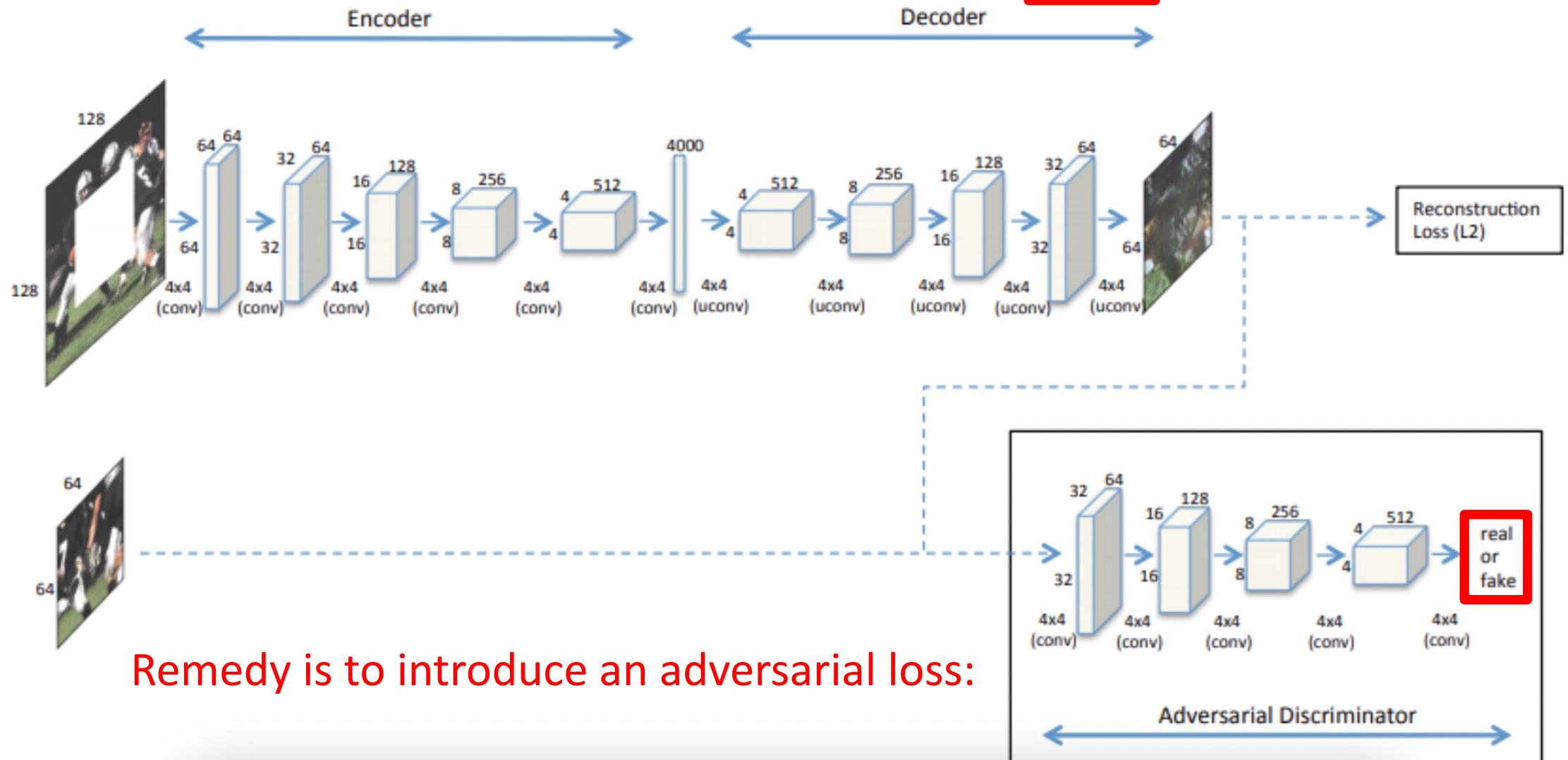
(a) Input context



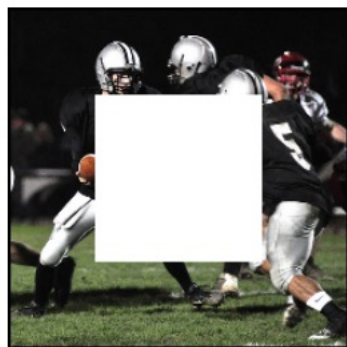
(c) Context Encoder
(L_2 loss)

Why might training with this loss function alone lead to blurry results?
- It averages the multiple plausible inpaintings for a hole

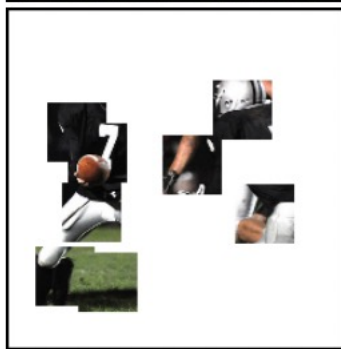
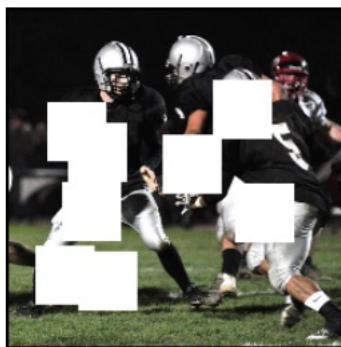
Training: Loss Functions ($\mathcal{L} = \lambda_{adv} \mathcal{L}_{adv} + \lambda_{rec} \mathcal{L}_{rec}$)



Training: Datasets



(a) Central region



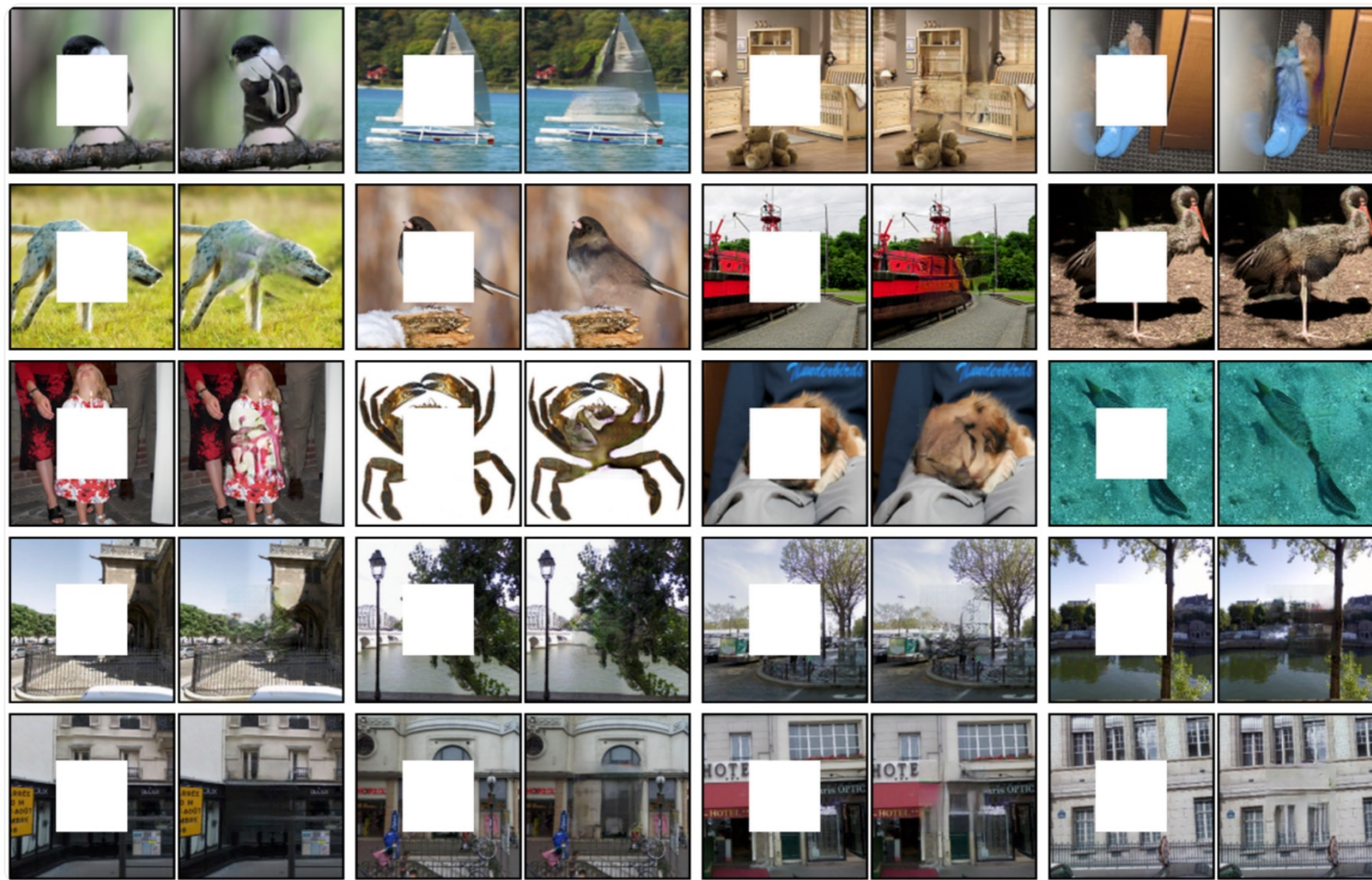
(b) Random block



(c) Random region

Training completed on ImageNet (all 1.2M and a 100K subset) for three hole types

Results: https://www.cs.cmu.edu/~dpathak/context_encoder/



What type of features might be learned?

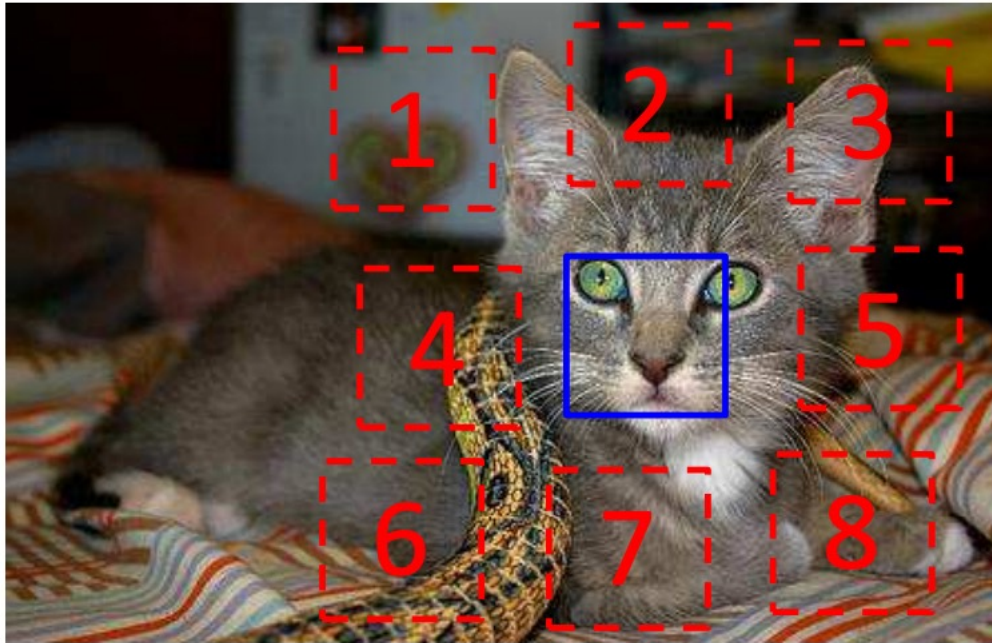
Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- Generative-based methods
- Generative adversarial networks
- Context-based methods

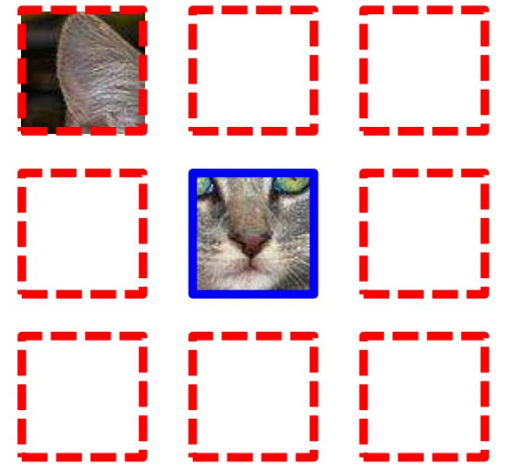
Context-based Methods

- **Spatial context**: predict relative positions of image patches
- **Timing context**: predict relative positions of video frames
- **Similarity context**: clustering

Spatial Context: Predict Image Index Per Patch

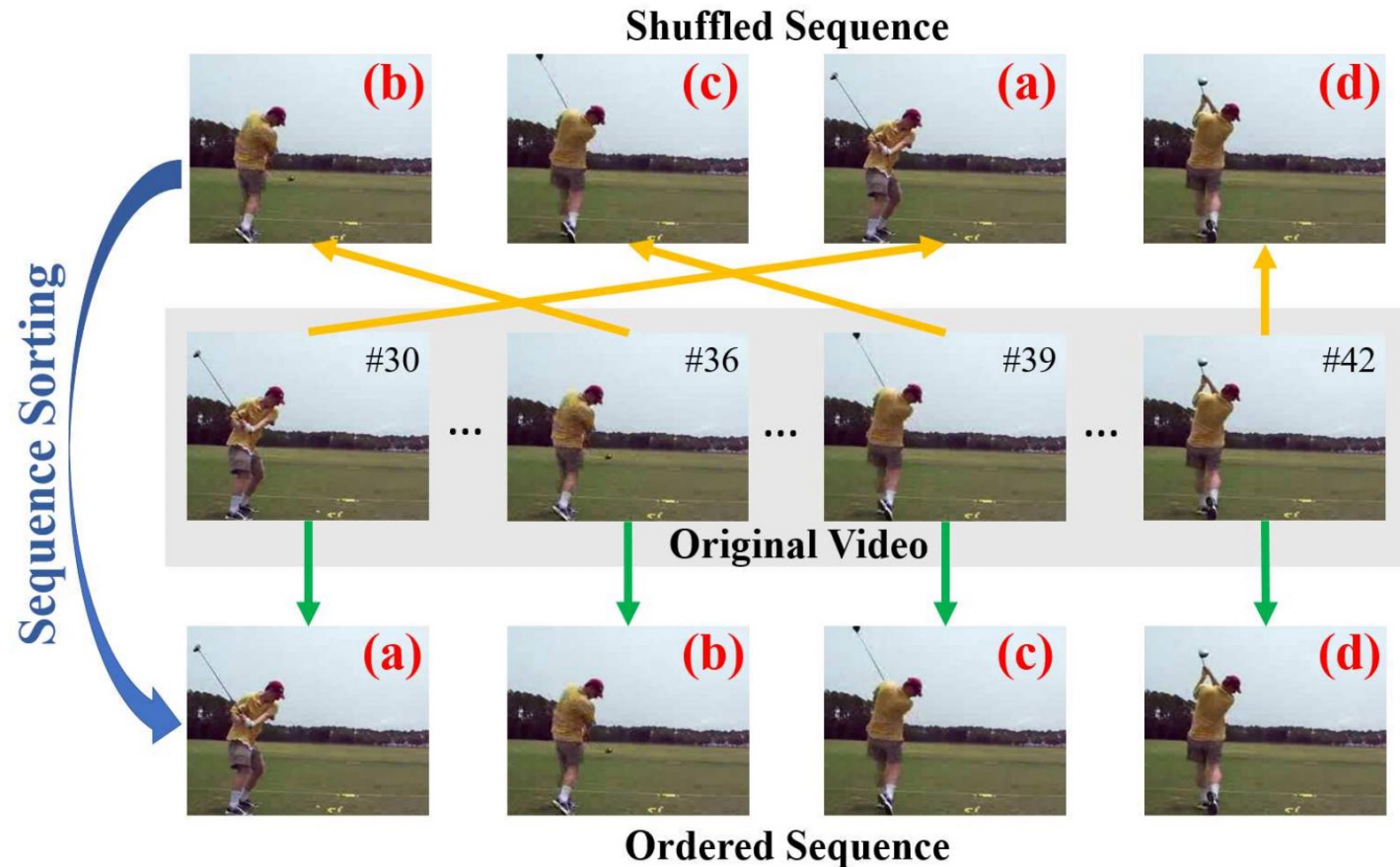


Example:



What type of features might be learned?

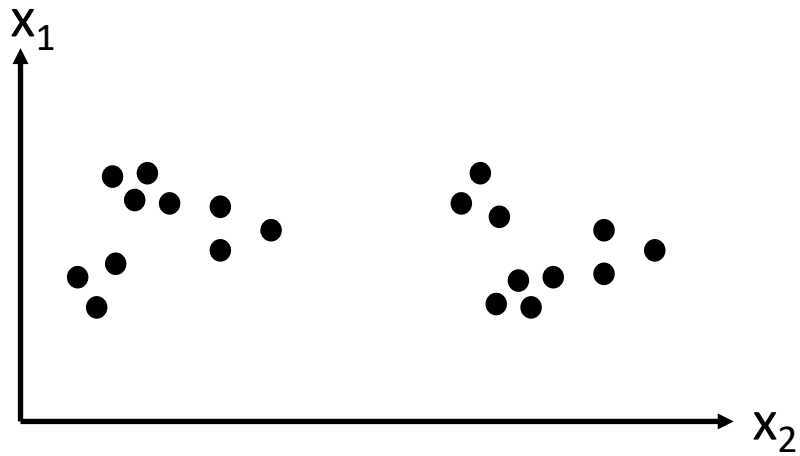
Timing Context : Predict Order of Video Frames



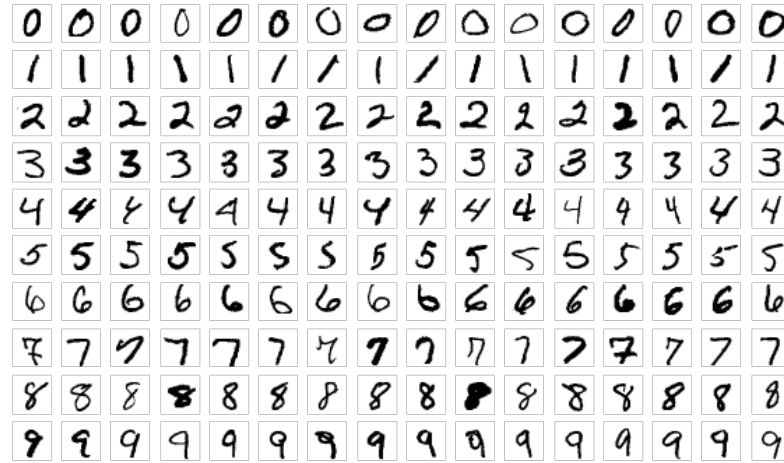
What type of features might be learned?

Similarity Context: Predict Clusters

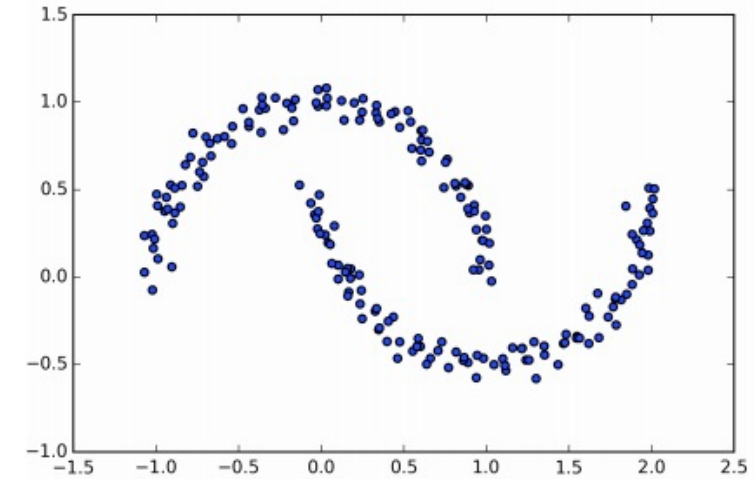
A.



B.



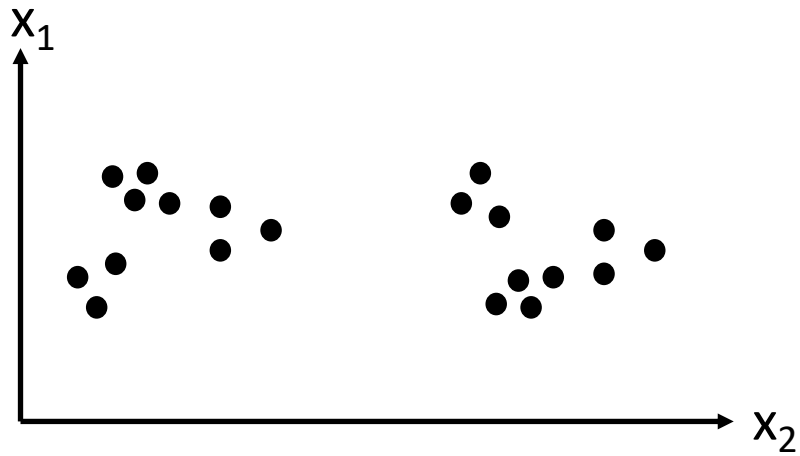
C.



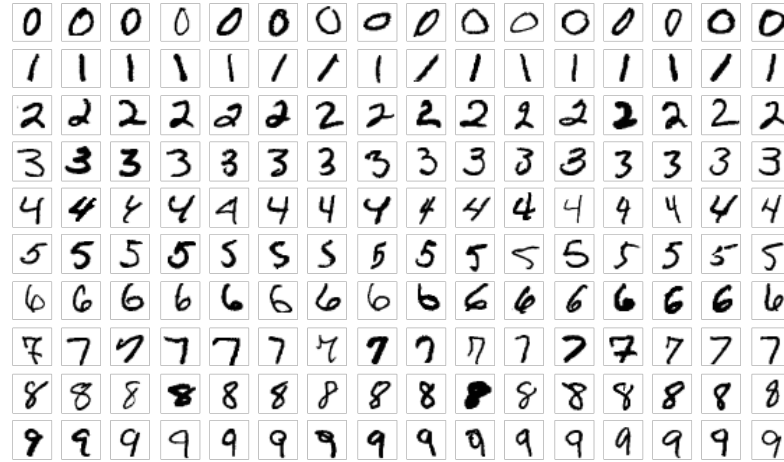
CNNs are trained to identify cluster assignments OR to recognize whether images belong to the same cluster

Clustering

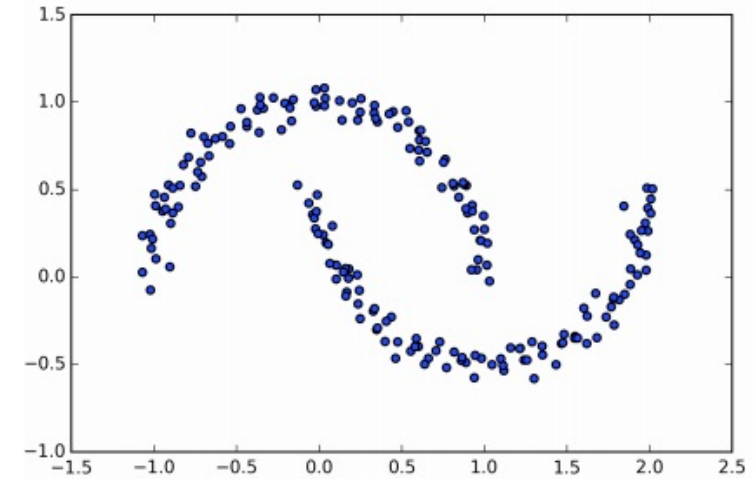
A.



B.



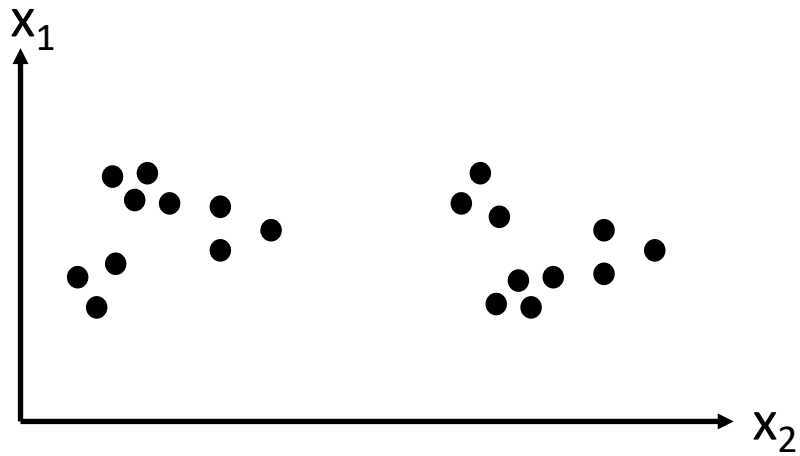
C.



Create groupings so entities in a group will be similar to each other and different from the entities in other groups.

Clustering: Key Questions

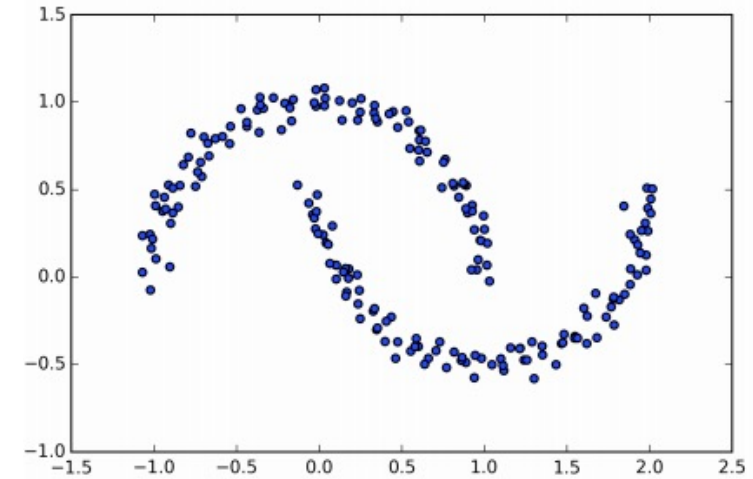
A.



B.

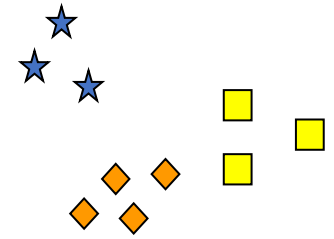
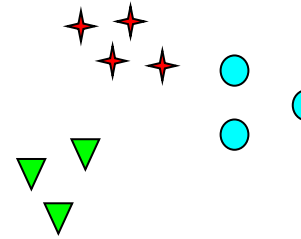
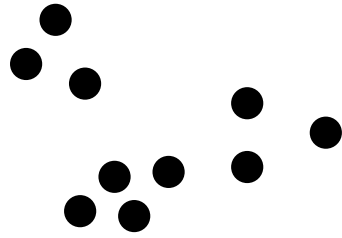
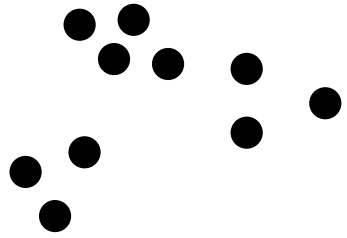


C.

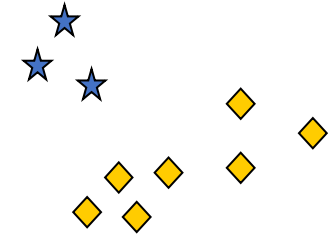
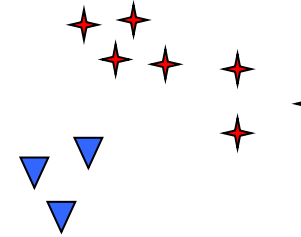
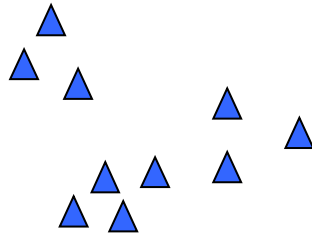
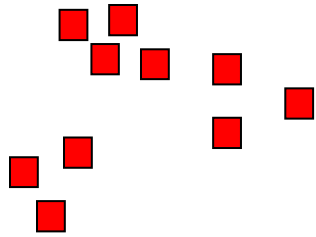


- How many data clusters to create?
- What “algorithm” to use to partition the data?

Clustering: How Many Clusters to Create?



Six Clusters



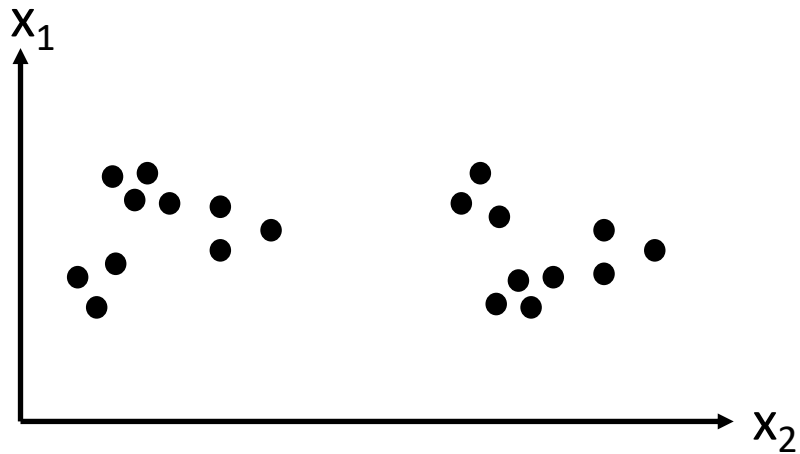
Two Clusters

Four Clusters

Number of clusters can be ambiguous.

Clustering

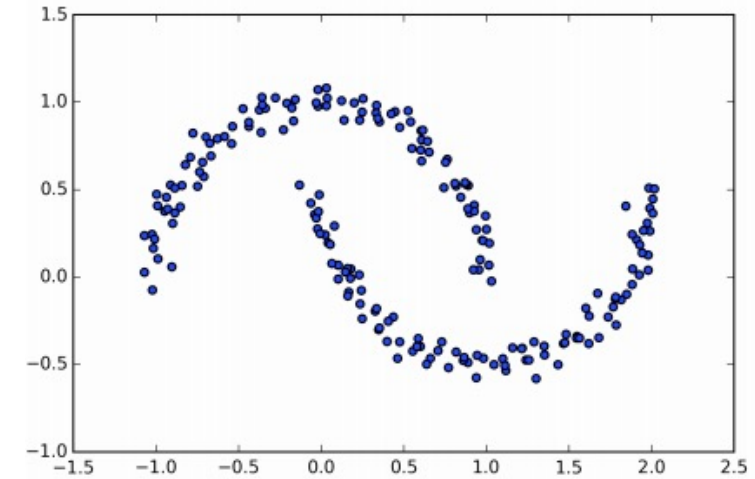
A.



B.



C.



Create groupings so entities in a group will be similar to each other and different from the entities in other groups.

What type of features might be learned?

Context-based Methods: How Might Such Methods Be Used in the NLP Field?

- **Spatial context**: predict relative positions of image patches
- **Timing context**: predict relative positions of video frames
- **Similarity context**: clustering

Today's Topics

- Transfer learning definition
- Overview of self-supervised learning
- Generative-based methods
- Generative adversarial networks
- Context-based methods

The image features a dark gray background with a large, faint, circular glow in the center. A white film strip border, consisting of a series of rectangular sprocket holes, frames the entire scene. In the center of the glow, the words "The End" are written in a white, elegant, cursive script font. The text has a slight drop shadow, giving it a three-dimensional appearance as if it's floating within the scene.

The End