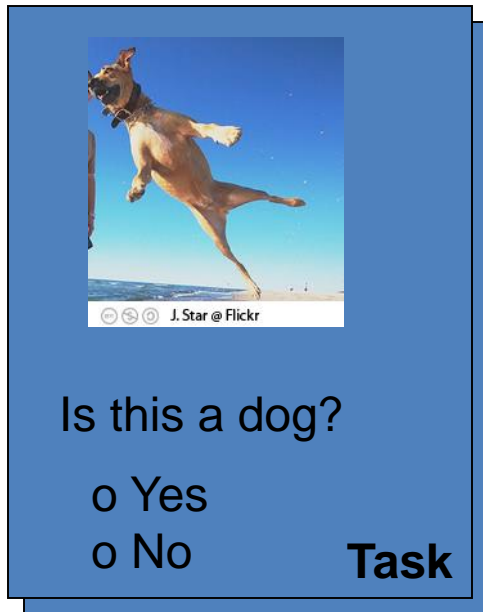# CS 1674: *Intro to Computer Vision*
# Recent Topics
# (Unsupervised Learning)

Prof. Adriana Kovashka

University of Pittsburgh

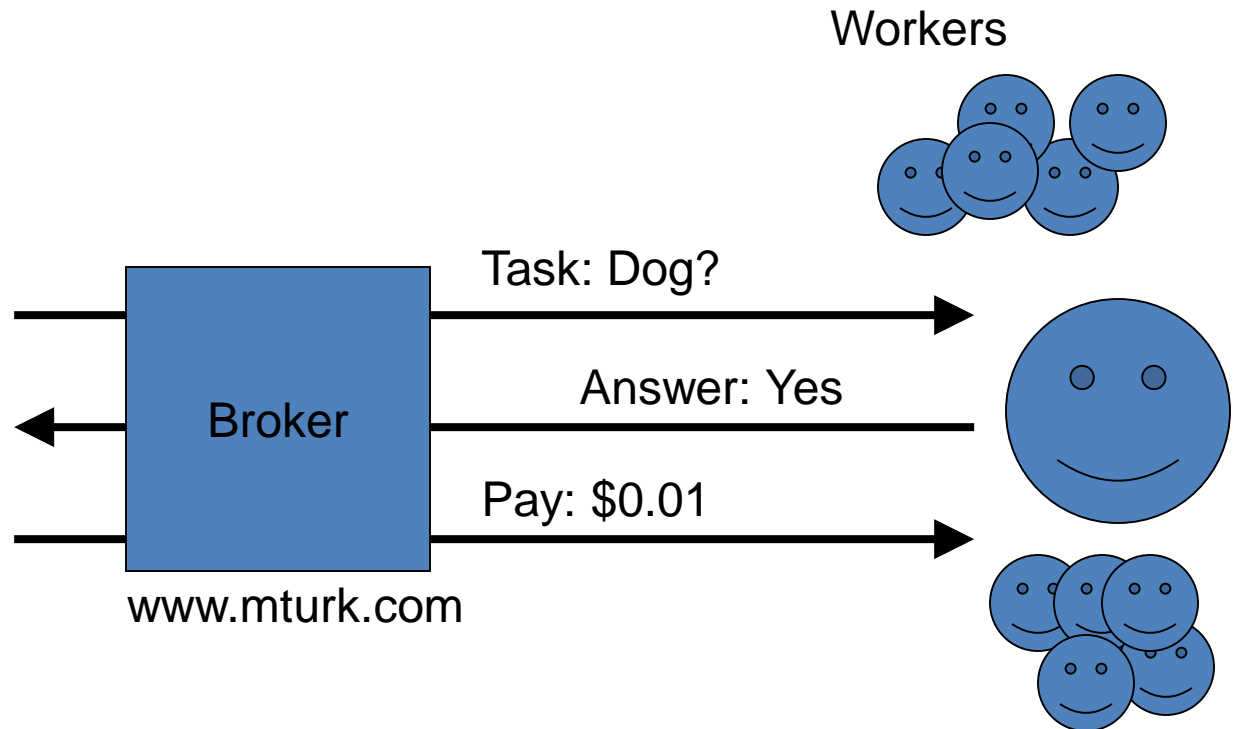December 4, 2018

# Motivation

- So far we've assumed access to plentiful labeled data

- How is this data obtained?

# Crowdsourcing

Workers

Is this a dog?

o Yes
o No

**Task**

$0.01

Broker

www.mturk.com

Task: Dog?

Answer: Yes

Pay: $0.01

J. Star @ Flickr

Alex Sorokin

# Crowdsourcing via games

- The ESP Game
  - Two-player online game
  - Partners don't know each other and can't communicate
  - Objective of the game: type the same word
  - The only thing in common is an image

von Ahn and Dabbish, "Labeling Images with a Computer Game," CHI 2004

# THE **ESP GAME**

## PLAYER 1



**GUESSING: CAR**

**GUESSING: HAT**

**GUESSING: KID**

**SUCCESS!**
**YOU AGREE ON CAR**

## PLAYER 2



**GUESSING: BOY**

**GUESSING: CAR**

**SUCCESS!**
**YOU AGREE ON CAR**

von Ahn and Dabbish, "Labeling Images with a Computer Game," CHI 2004

# Motivation

- So far we've assumed access to plentiful labeled data
- **What if we have limited or no labeled data?**
- One approach: learn from unlabeled data (unsupervised learning)
  - Mine for interesting patterns (discovery)
  - Use supervision (labels) inherent in the data (self-supervised learning)
- Another approach (not discussed): carefully choose which data to label
  - Active learning, human-in-the-loop

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data**: $(x, y)$
x is data, y is label

**Goal**: Learn a *function* to map x -> y

**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.

## Unsupervised Learning

**Data**: x
Just data, **no labels!**

**Goal**: Learn some underlying hidden *structure* of the data

**Examples**: Clustering, dimensionality reduction, feature learning, density estimation, etc.
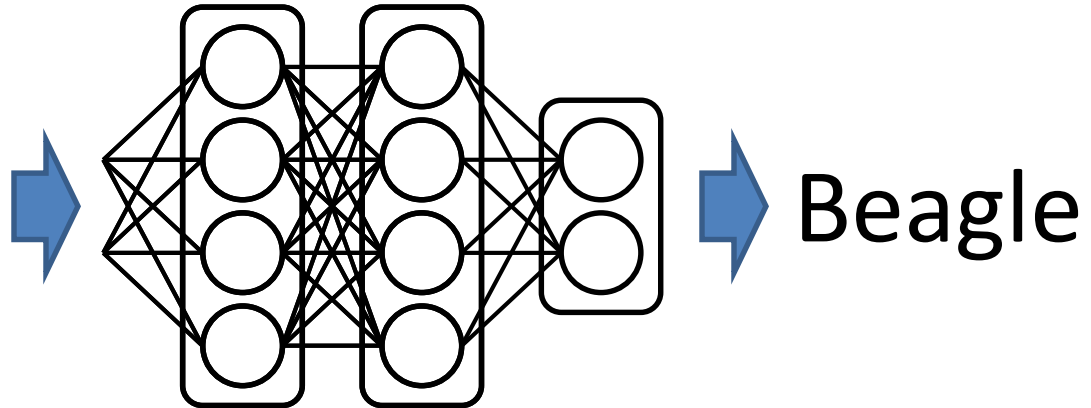
Serena Young

# Plan for this last lecture

- Self-supervised learning
  - For images
  - For video

- Visual discovery
  - Discovering style-specific elements

- Generative models (sep. from recognition)
  - Theory/technique
  - Applications
  - Generating synthetic training data

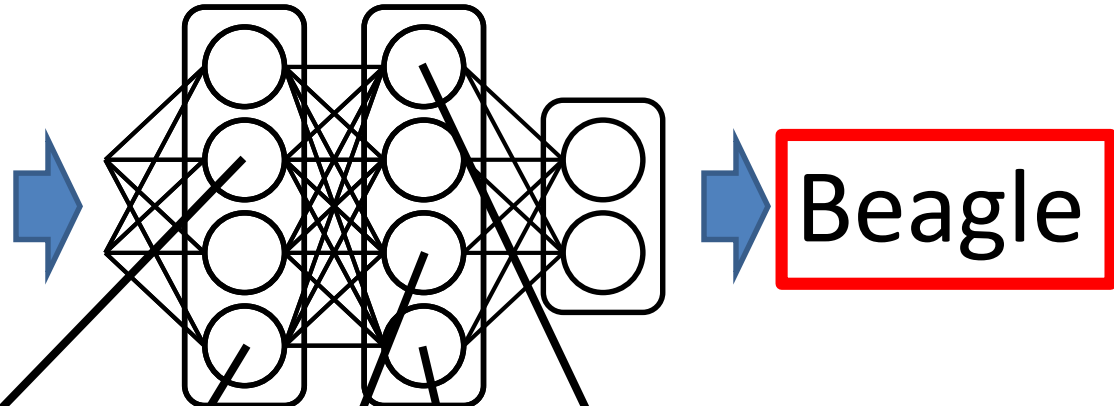# Unsupervised Visual Representation Learning by Context Prediction

Carl Doersch, Alexei Efros and Abhinav Gupta

ICCV 2015

# ImageNet + Deep Learning



Beagle

- Image Retrieval
- Detection (RCNN)
- Segmentation (FCN)
- Depth Estimation
- ...

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

# ImageNet + Deep Learning



Materials?

Parts?

Pose?

Geometry?

Boundaries?

*Do we even need this task?*

*Do we ever need semantic labels?*

Beagle

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

# Context as Supervision
[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk, but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would

Deep Net

# Context Prediction for Images

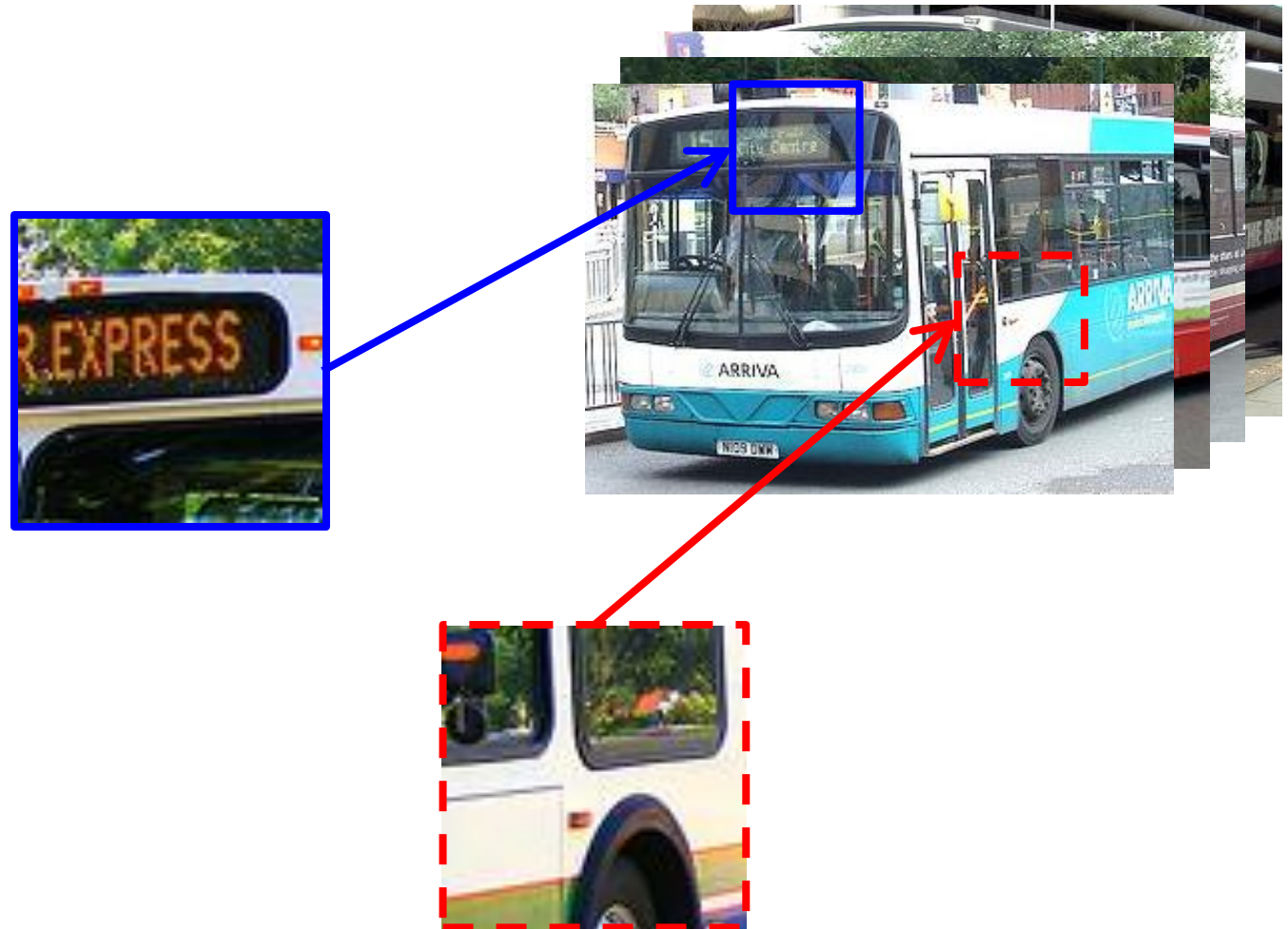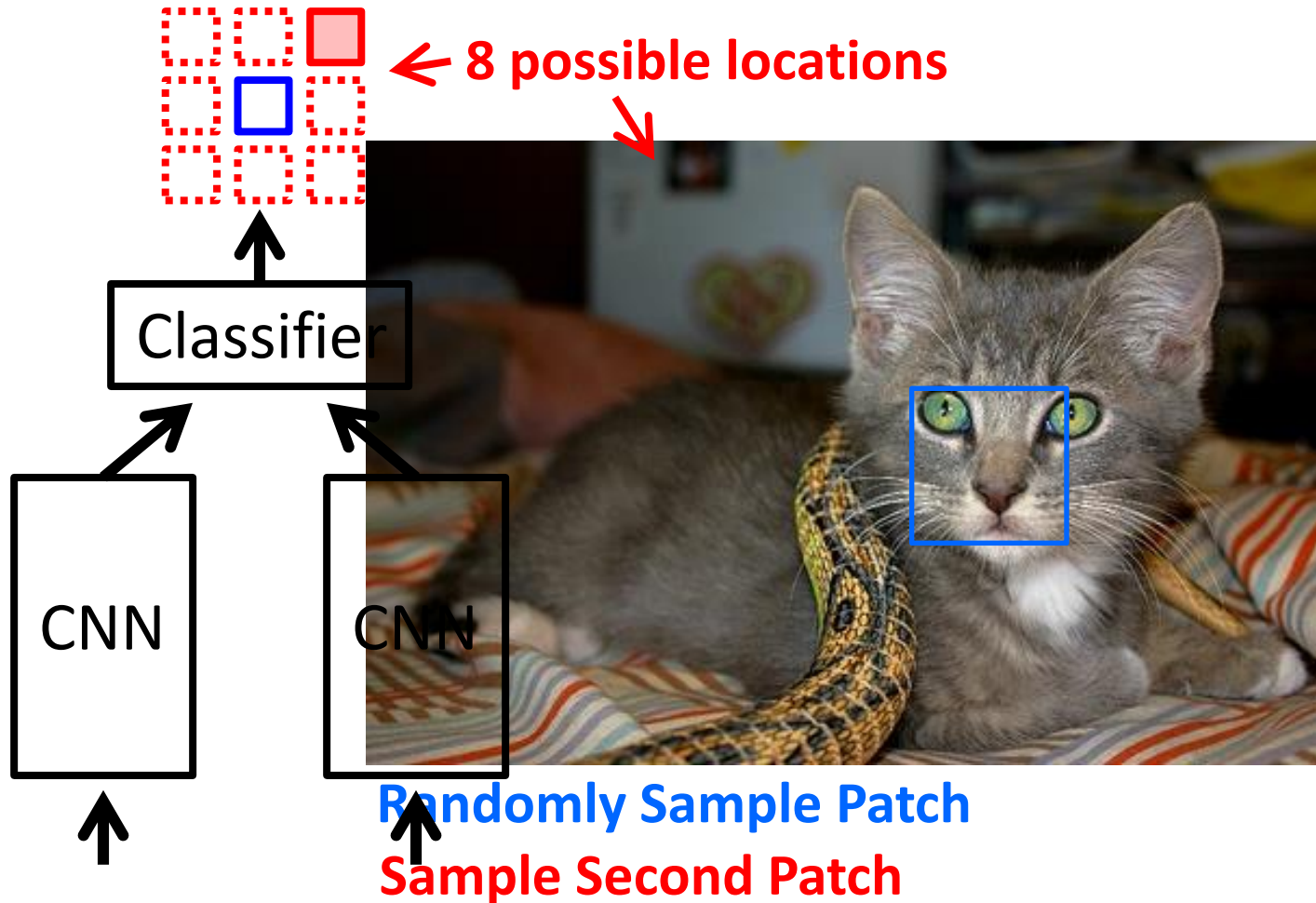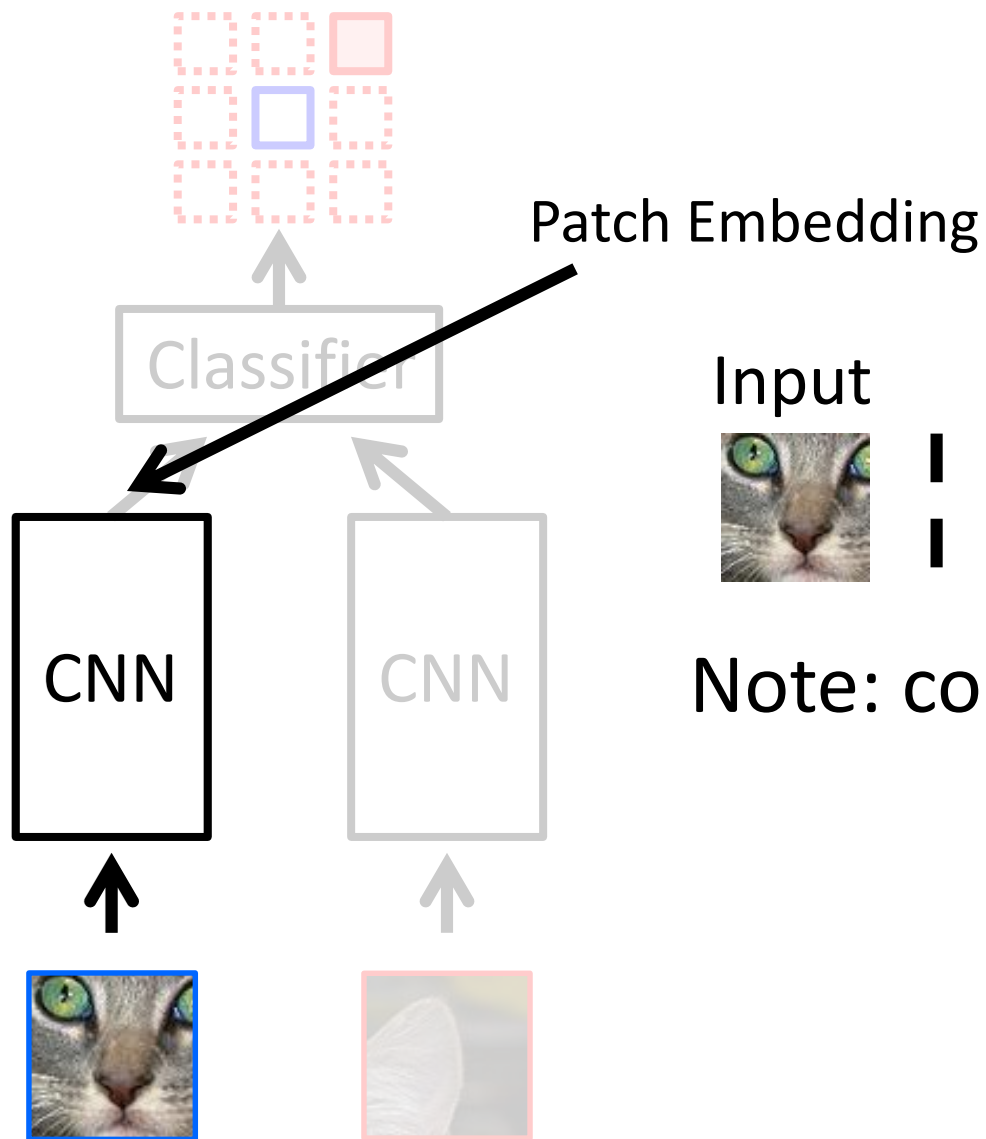# Semantics from a non-semantic task



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

# Relative Position Task



← **8 possible locations**

Classifier

CNN    CNN

**Randomly Sample Patch**

**Sample Second Patch**

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

Patch Embedding

Classifier

CNN

CNN

Input

Nearest Neighbors

Note: connects *across* instances!

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

# Architecture



Softmax loss

Fully connected

Fully connected

| Fully connected | | Fully connected |
| Max Pooling | | Max Pooling |
| Convolution | | Convolution |
| Convolution | | Convolution |
| Convolution | | Convolution |
| LRN | | LRN |
| Max Pooling | | Max Pooling |
| Convolution | | Convolution |
| LRN | | LRN |
| Max Pooling | | Max Pooling |
| Convolution | Tied Weights | Convolution |

Patch 1

Patch 2

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

# What is learned?

Input        Ours        ImageNet AlexNet



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

# Pre-Training for R-CNN



1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

Pre-train on relative-position task, w/o labels

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

[Girshick et al. 2014]

# VOC 2007 Performance
## (pretraining for R-CNN)



% Average Precision

54.2 — ImageNet Labels

46.3 — Ours

40.7 — No Pretraining

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

# Shuffle and Learn: Unsupervised Learning using Temporal Order Verification

Ishan Misra, C. Lawrence Zitnick, and Martial Hebert

ECCV 2016

Fig. 1: **(a)** A video imposes a natural temporal structure for visual data. In many cases, one can easily verify whether frames are in the correct temporal order (shuffled or not). Such a simple sequential verification task captures important spatiotemporal signals in videos. We use this task for unsupervised pre-training of a Convolutional Neural Network (CNN). **(b)** Some examples of the automatically extracted positive and negative tuples used to formulate a classification task for a CNN.

Misra et al., "Shuffle and Learn: Unsupervised Learning using Temporal Order Verification", ECCV 2016

Fig. 2: **(a)** We sample tuples of frames from high motion windows in a video. We form positive and negative tuples based on whether the three input frames are in the correct temporal order. **(b)** Our triplet Siamese network architecture has three parallel network stacks with shared weights upto the `fc7` layer. Each stack takes a frame as input, and produces a representation at the `fc7` layer. The concatenated `fc7` representations are used to predict whether the input tuple is in the correct temporal order.

Misra et al., "Shuffle and Learn: Unsupervised Learning using Temporal Order Verification", ECCV 2016

**Table 2:** Mean classification accuracies over the 3 splits of UCF101 and HMDB51 datasets. We compare different initializations and finetune them for action recognition.

| Dataset | Initialization | Mean Accuracy |
|---|---|---|
| UCF101 | Random | 38.6 |
| | (Ours) Tuple verification | **50.2** |
| HMDB51 | Random | 13.3 |
| | UCF Supervised | 15.2 |
| | (Ours) Tuple verification | **18.1** |

Misra et al., "Shuffle and Learn: Unsupervised Learning using Temporal Order Verification", ECCV 2016

# Plan for this last lecture

- Self-supervised learning
  - For images
  - For video
- Visual discovery
  - Discovering style-specific elements
- Generative models
  - Theory/technique
  - Applications
  - Generating synthetic training data

# What Makes Paris Look like Paris?

Carl Doersch, Saurabh Singh, Abhinav Gupta,  Josef Sivic, Alexei Efros

# One of these is from Paris
# Raise your hand if...

...this is Paris



Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

# Raise your hand if...



We showed 20 subjects:

- 100 Random Street View Images

- 50 from Paris

- They classified Paris vs non-Paris

- Accuracy: 79%

How do they know?

We showed 20 subjects:

- 100 Random Street View Images

- 50 from Paris

- They classified Paris non-Paris

- Accuracy: 79%

How do they know?

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

# Our Goal:

*Given a large geo-tagged image dataset,*
*we automatically discover **visual elements***
*that characterize a geographic location*

*Why might this be a useful task?*

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

# Our Hypothesis

- The visual elements that capture Paris:
  - Frequent: Occur often in Paris
  - Discriminative: Are not found outside Paris

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

London
Prague
Paris
New York     Boston
Milan
San
Tokyo
Barcelona
Francisco   Philadelphia
Mexico City

Sao
Paulo

- Positive Set
- Negative Set

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

# Step 1: Nearest Neighbors for Every Patch
## Using normalized correlation of HOG features as a distance metric

patch          nearest neighbors



Paris
Not Paris

# Step 2: Find the Parisian Clusters by Sorting



patch    nearest neighbors

Sort by # Paris Neighbors

patch    nearest neighbors

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

# Paris: A Few Top Elements



Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

Doersch et al., "What Makes Paris Look Like Paris?", SIGGRAPH 2012

Elements from Prague          Elements from London          Elements from Barcelona

# In the U.S.



Elements from San Francisco



Elements from Boston

# Plan for this last lecture

- Self-supervised learning
  - For images
  - For video

- Visual discovery
  - Discovering style-specific elements

- Generative models (sep. from recognition)
  - Theory/technique
  - Applications
  - Generating synthetic training data

# Generative Models



Training data ~ $p_{data}(x)$



Generated samples ~ $p_{model}(x)$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

Serena Young

# Generative Models



Training data ~ $p_{data}(x)$                    Generated samples ~ $p_{model}(x)$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

Addresses density estimation, a core problem in unsupervised learning
**Several flavors:**
- Explicit density estimation: explicitly define and solve for $p_{model}(x)$
- Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

# Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models can be used to enhance training datasets with diverse synthetic data
- Generative models of time-series data can be used for simulation

# Generative Adversarial Networks

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise. Learn transformation to training distribution.

Q: What can we use to represent this complex transformation?

Serena Young

# Generative Adversarial Networks

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise. Learn transformation to training distribution.

Q: What can we use to represent this complex transformation?

A: A neural network!

Output: Sample from
 training distribution

Generator
Network

Input: Random noise        z

Serena Young

# Training GANs: Two-player game

**Generator network**: try to fool the discriminator by generating real-looking images
**Discriminator network**: try to distinguish between real and fake images

Serena Young

# Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

**Generator network**: try to fool the discriminator by generating real-looking images
**Discriminator network**: try to distinguish between real and fake images



Real or Fake

Discriminator Network

Fake Images (from generator)

Real Images (from training set)

Generator Network

Random noise          z

Serena Young

# Adversarial Networks Framework



Ian Goodfellow

# Adversarial Networks Framework



Generator
$x \sim G(z)$

Discriminator
Real vs. Fake

[Goodfellow et al. 2014]

Jun-Yan Zhu

# Training GANs: Two-player game

**Generator network**: try to fool the discriminator by generating real-looking images
**Discriminator network**: try to distinguish between real and fake images

Train jointly in **minimax game**

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Serena Young

# Training GANs: Two-player game
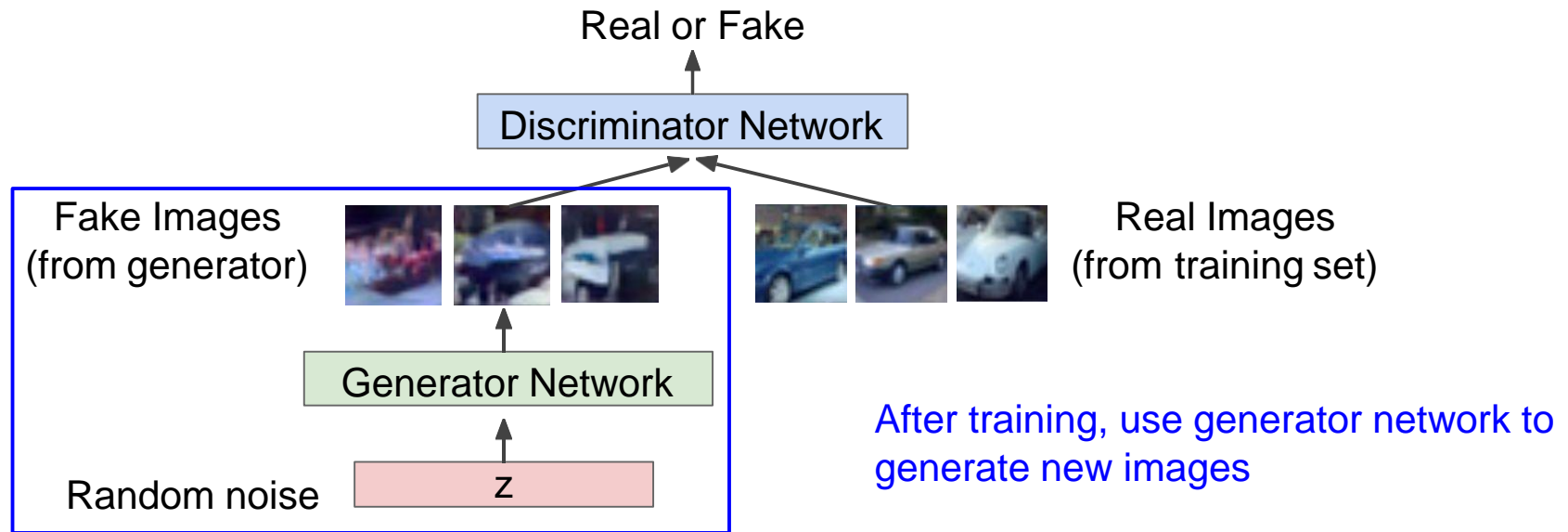
**Generator network**: try to fool the discriminator by generating real-looking images
**Discriminator network**: try to distinguish between real and fake images

Train jointly in **minimax game**

Discriminator outputs likelihood in (0,1) of real image

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output
for real data x

Discriminator output for
generated fake data G(z)

Serena Young

# Training GANs: Two-player game

**Generator network**: try to fool the discriminator by generating real-looking images
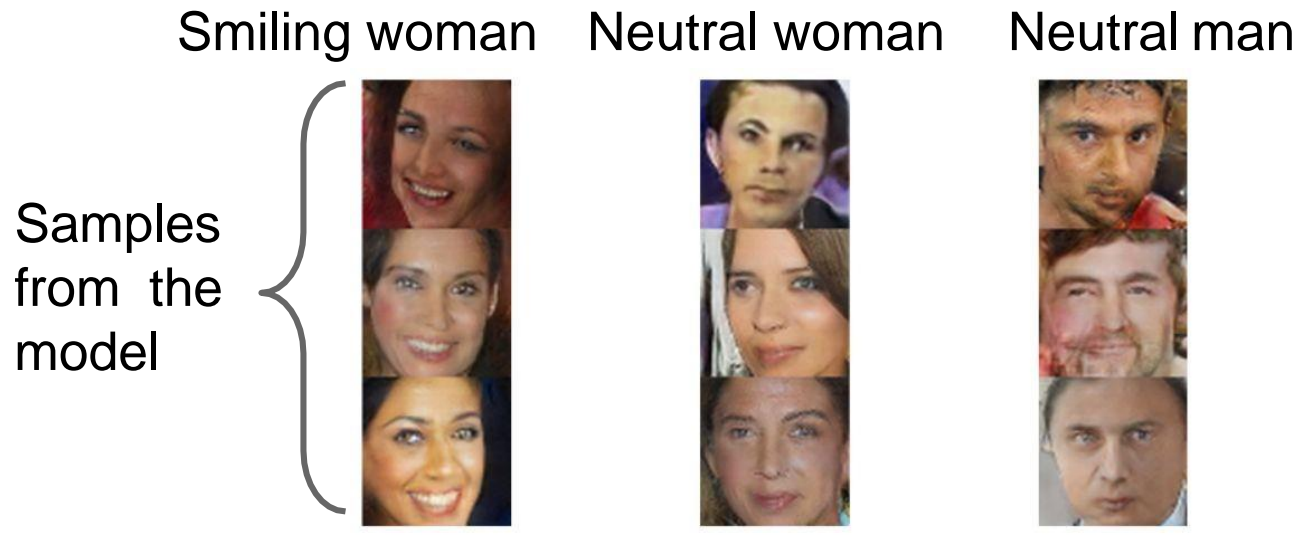**Discriminator network**: try to distinguish between real and fake images

Train jointly in **minimax game**

Discriminator outputs likelihood in (0,1) of real image

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log \underbrace{D_{\theta_d}(x)}_{} + \mathbb{E}_{z \sim p(z)} \log(1 - \underbrace{D_{\theta_d}(G_{\theta_g}(z))}_{}) \right]$$

Discriminator output
for real data x

Discriminator output for
generated fake data G(z)

- Discriminator ($\theta_d$) wants to **maximize objective** such that D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)
- Generator ($\theta_g$) wants to **minimize objective** such that D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)

Serena Young

# Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:
1. **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

# Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

**Generator network**: try to fool the discriminator by generating real-looking images
**Discriminator network**: try to distinguish between real and fake images

Real or Fake

Discriminator Network

Fake Images
(from generator)

Generator Network

Random noise          z

Real Images
(from training set)

After training, use generator network to generate new images

Serena Young

# Generative Adversarial Nets

**Samples from the model look amazing!**

Radford et al, ICLR 2016



Serena Young

# Generative Adversarial Nets: Interpretable Vector Math

Radford et al, ICLR 2016

Smiling woman    Neutral woman    Neutral man

Samples from the model



Serena Young

# Generative Adversarial Nets: Interpretable Vector Math

Radford et al, ICLR 2016



Smiling woman    Neutral woman    Neutral man

Samples from the model

Average Z vectors, do arithmetic

Serena Young

# Generative Adversarial Nets: Interpretable Vector Math

Smiling woman     Neutral woman     Neutral man

Smiling Man

Samples from the model

Average Z vectors, do arithmetic

**−**     **+**     **=**

Serena Young

# Generative Adversarial Nets: Interpretable Vector Math

Glasses man     No glasses man     No glasses woman

Radford et al,
ICLR 2016

Serena Young

# Generative Adversarial Nets: Interpretable Vector Math

Glasses man     No glasses man     No glasses woman

Woman with glasses



Serena Young

# Celebrities Who Never Existed



*Karras et al., "Progressive Growing of GANs for Improved Quality, Stability, and Variation", ICLR 2018*

# StarGAN



| Input | Blond hair | Gender | Aged | Pale skin | Input | Angry | Happy | Fearful |

Choi et al., "StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation", CVPR 2018

# GANs for Privacy (Action Detection)



*Ren et al., "Learning to Anonymize Faces for Privacy Preserving Action Detection", arxiv 2018*

# Artificial Fashion: vue.ai

# Edges → Images

Input        Output        Input        Output        Input        Output



Edges from [Xie & Tu, 2015]

Pix2pix / CycleGAN

# *Sketches* → Images

Input    Output    Input    Output    Input    Output



Trained on Edges → Images

Data from [Eitz, Hays, Alexa, 2012]

Pix2pix / CycleGAN

**#edges2cats**    [Christopher Hesse]



@gods_tail



Ivy Tasi @ivymyt



@matthematician



Vitaly Vidmirov @vvid

https://affinelayer.com/pixsrv/

Pix2pix / CycleGAN

# Changing artistic style



Input    Monet    Van Gogh    Cezanne    Ukiyo-e

Pix2pix / CycleGAN

# Changing seasons