

# CS230: Lecture 3

# Various Deep Learning Topics

Kian Katanforoosh, Andrew Ng

## Today's outline

We will learn how to:

- Analyse a problem from a deep learning approach
- Choose an **architecture**
- Choose a **loss** and a **training strategy**

- I. Day'n'Night classification
- II. Face Recognition
- III. Art generation
- IV. Object detection
- V. Image Segmentation

## Day'n'Night classification (warm-up)

**Goal:** Given an image, classify as taken “during the day” (0) or “during the night” (1)

**1. Data?**

10,000 images

Split? Bias?

**2. Input?**



Resolution?

(64, 64, 3)

**3. Output?**

$y = 0$  or  $y = 1$

Last Activation?

sigmoid

**4. Architecture ?**

Shallow network should do the job pretty well

**5. Loss?**

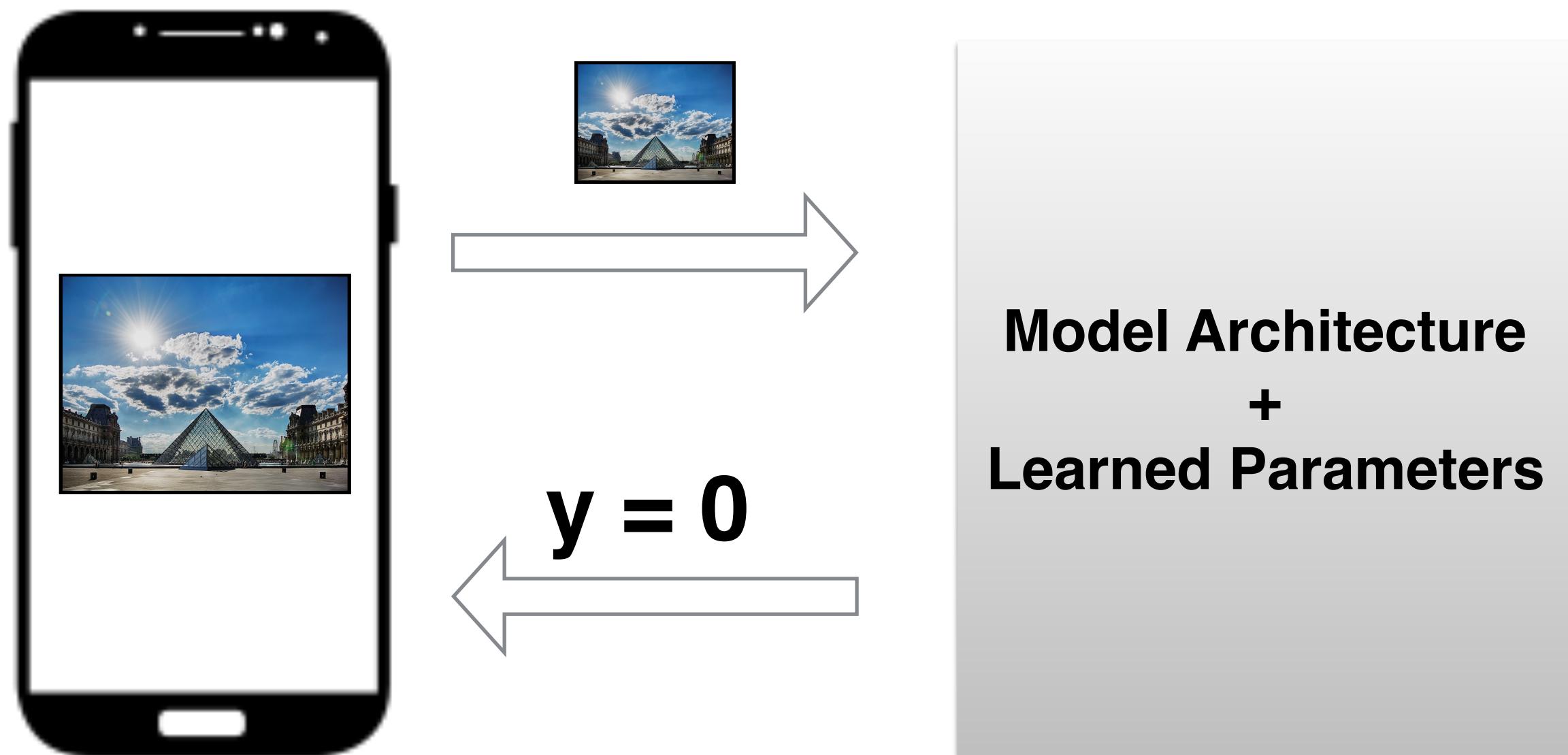
$$L = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

Easy

# Server-based or on-device?

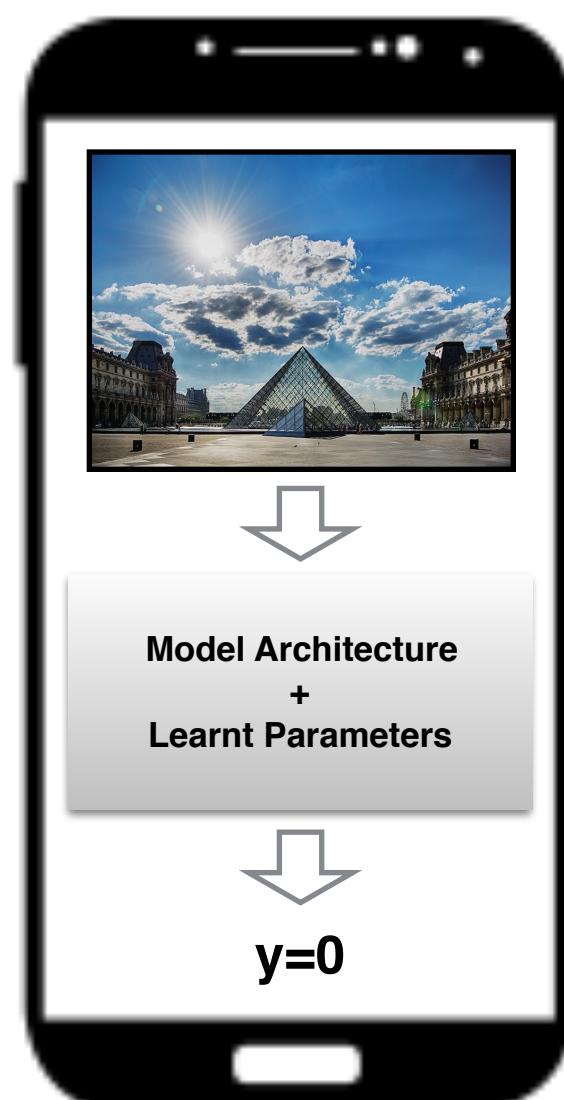
## Server-based

- + App is light-weight



## On-device

- + Faster predictions



# Face Recognition

**Goal:** A school wants to use Face Verification for validating student IDs in facilities (dinning hall, gym, pool ...)

## 1. Data?

Picture of every student labelled with their name



Bertrand

## 2. Input?



Resolution?  
(412, 412, 3)

## 3. Output?

$y = 1$  (it's you)  
or  
 $y = 0$  (it's not you)

# Face Recognition

**Goal:** A school wants to use Face Verification for validating student IDs in facilities (dinning hall, gym, pool ...)

## 4. What architecture?

Simple solution:



compute distance  
pixel per pixel  
if less than threshold  
then  $y=1$



database image

input image

Issues:

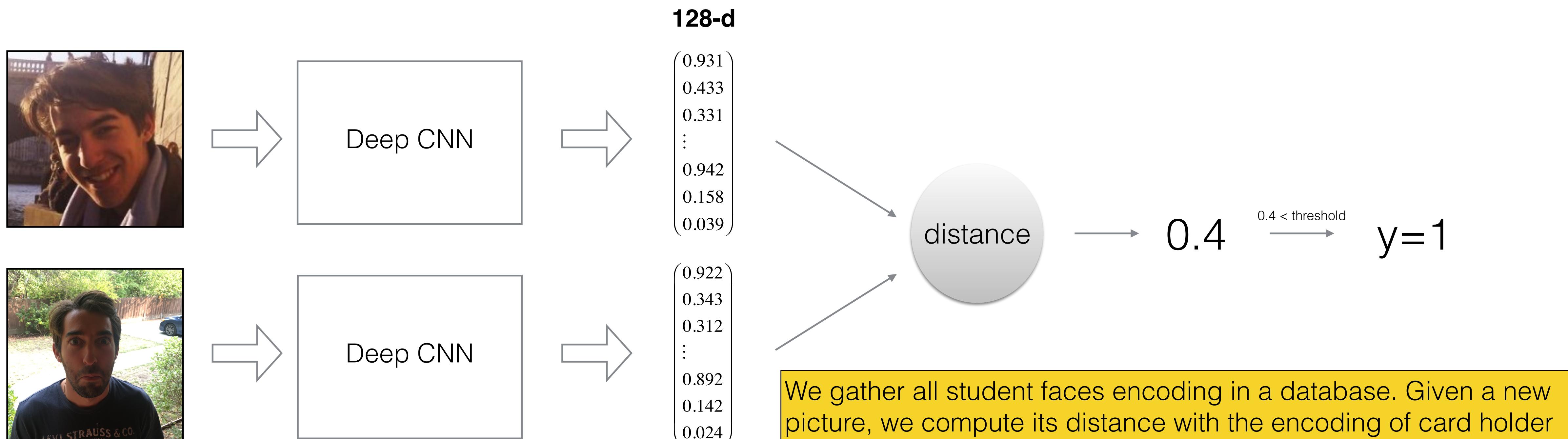
- Background lighting differences
- A person can wear make-up, grow a beard...
- ID photo can be outdated

# Face Recognition

**Goal:** A school wants to use Face Verification for validating student IDs in facilities (dinning hall, gym, pool ...)

## 4. What architecture?

Our solution: encode information about a picture in a vector



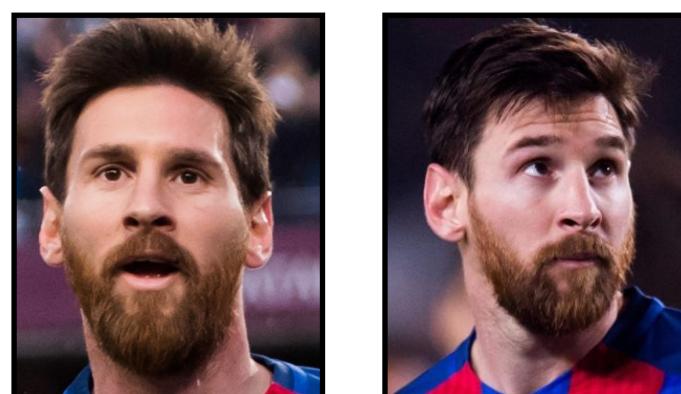
# Face Recognition

**Goal:** A school wants to use Face Verification for validating student IDs in facilities (dinning hall, gym, pool ...)

## 4. Loss? Training?

We need more data so that our model understands how to encode:  
Use public face datasets

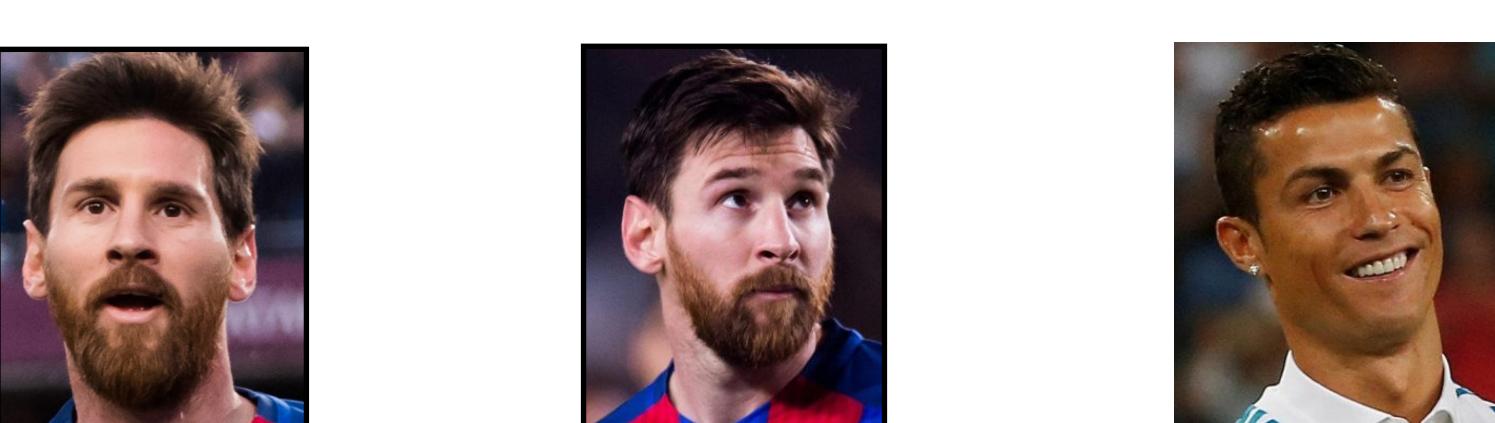
What we really want:



similar encoding



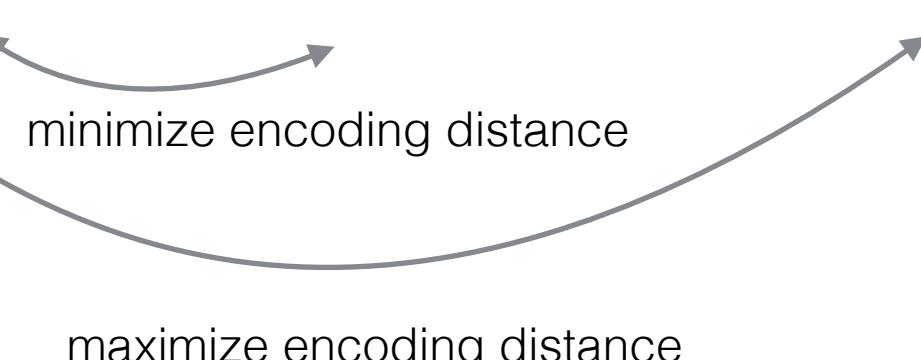
different encoding



anchor

positive

negative



So let's generate triplets:

$$L = \left\| \text{Enc}(A) - \text{Enc}(P) \right\|_2^2 - \left\| \text{Enc}(A) - \text{Enc}(N) \right\|_2^2 + \alpha$$

# Face Recognition

**Goal:** A school wants to use Face Identification for recognize students in facilities  
(dinning hall, gym, pool ...)

## K-Nearest Neighbors

**Goal:** You want to use Face Clustering to group pictures of the same people on your smartphone

## K-Means Algorithm

Maybe we need to detect the faces first?

# Art generation (Neural Style Transfer)

**Goal:** Given a picture, make it look beautiful

## 1. Data?

Let's say we have  
any data



## 2. Input?



content  
image

## 3. Output?



style  
image

generated  
image

## Art generation (Neural Style Transfer)

### 4. Architecture?

We want a model that **understands images** very well  
We load an **existing model trained on ImageNet** for example



When this image forward propagates, we can get information about its content & its style by inspecting the layers.

*Content<sub>C</sub>*  
*Style<sub>S</sub>*

### 5. Loss?

$$L = \|\text{Content}_C - \text{Content}_G\|_2^2 + \|\text{Style}_S - \text{Style}_G\|_2^2$$

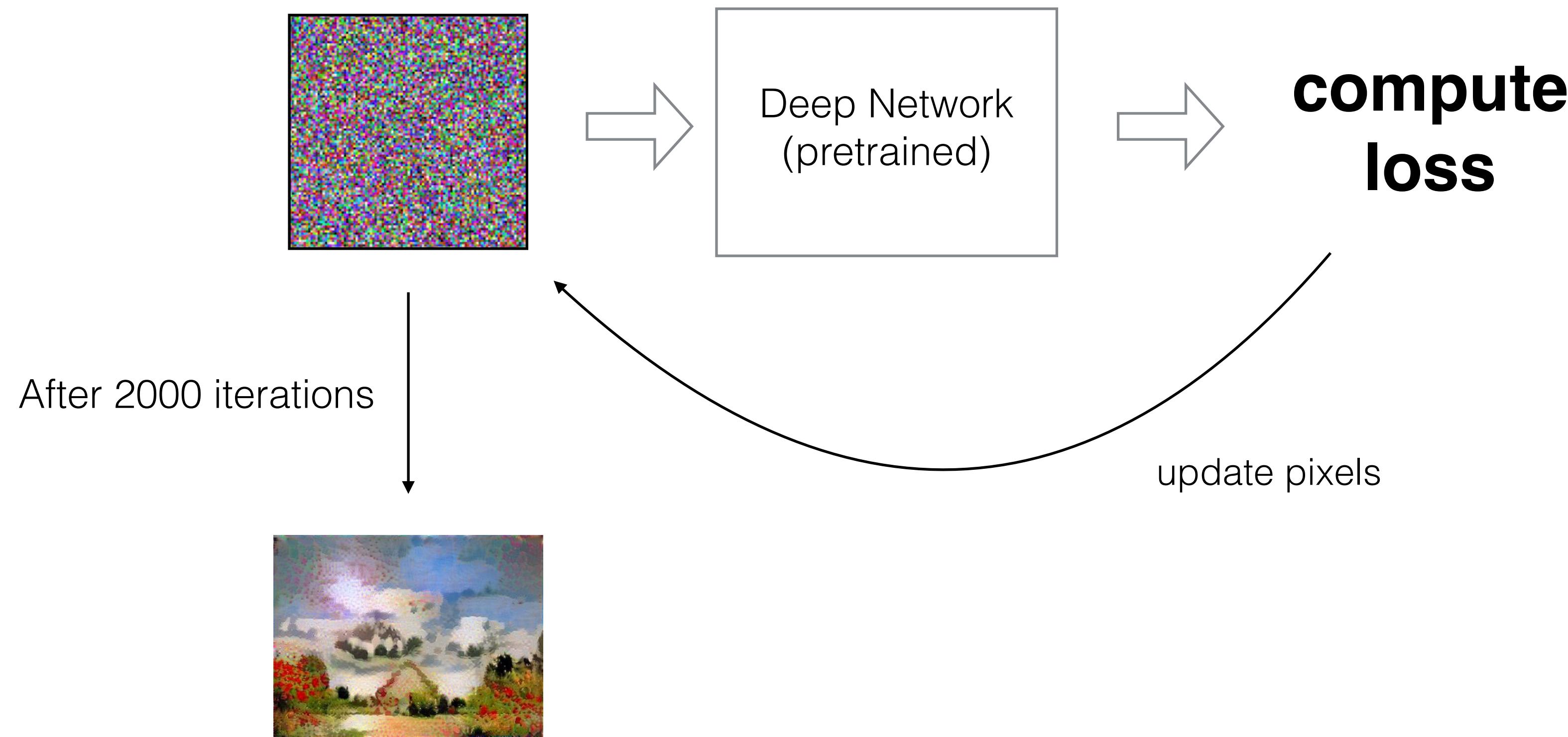


We are not learning parameters by minimizing L. We are learning an image!

# Art generation (Neural Style Transfer)

## Correct Approach

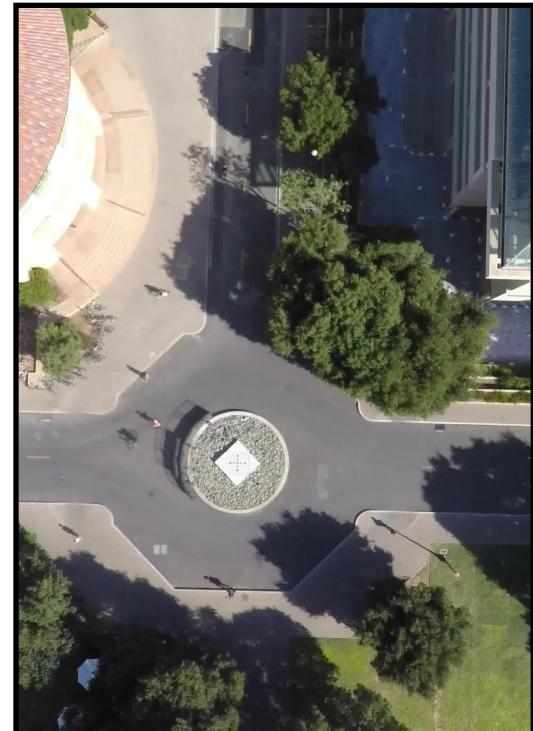
$$L = \left\| Content_C - Content_G \right\|_2^2 + \left\| Style_S - Style_G \right\|_2^2$$



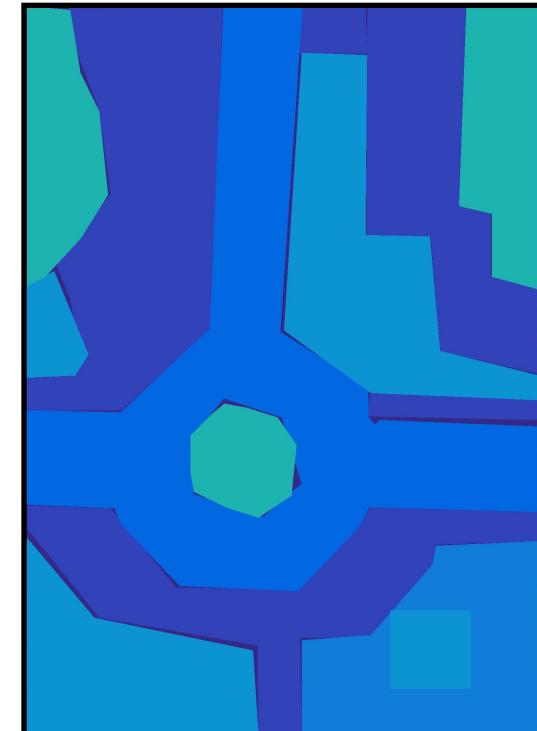
# Image Segmentation

**Goal:** Separate the foreground from the background on a picture

**1. Data?**

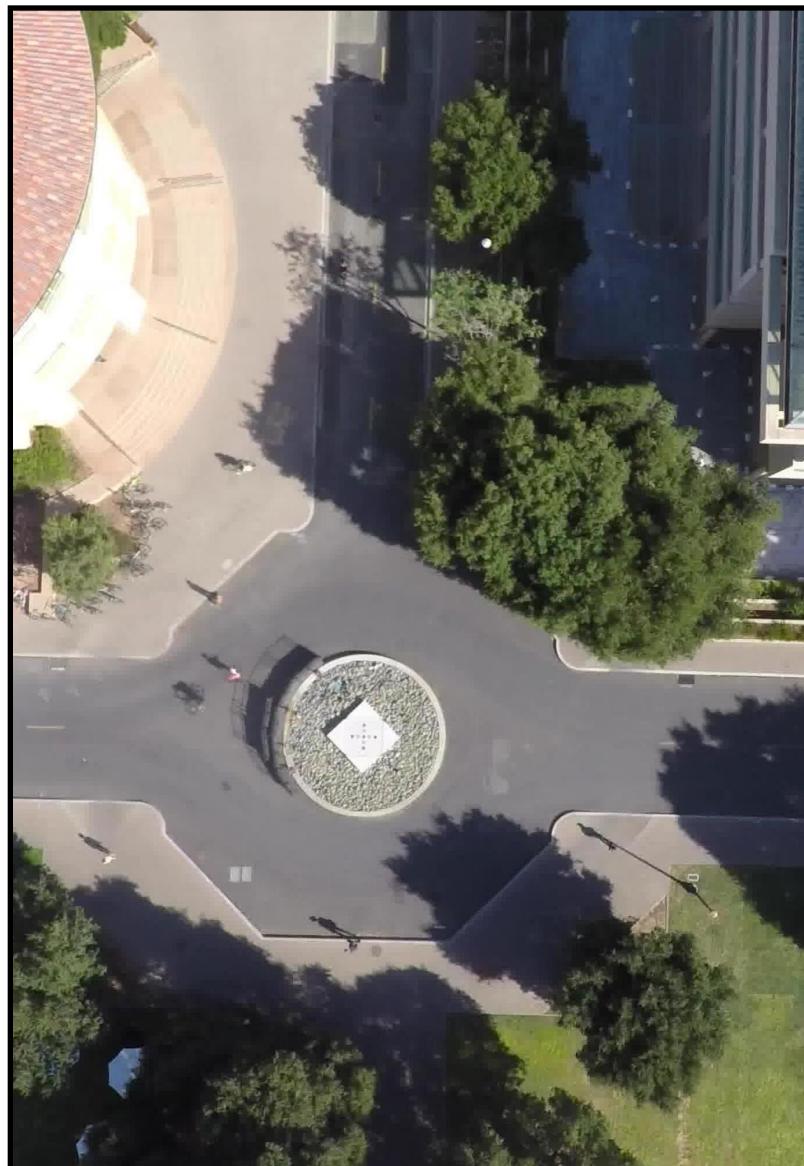


image



labels

**2. Input?**

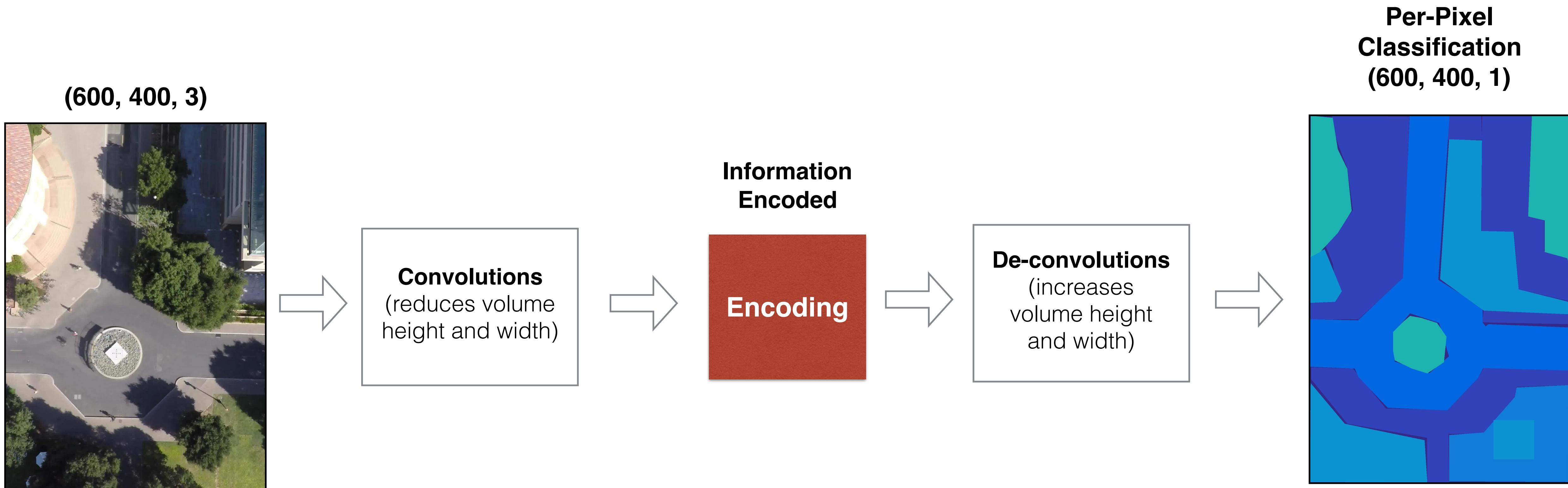


**3. Output?**



# Image Segmentation

## 4. Architecture?



# Image Segmentation

## 4. Loss?

pixel-wise cross-entropy

$$L = \sum_{pixels} \sum_{classes} y \log(\hat{y})$$

The diagram shows the calculation of pixel-wise cross-entropy loss. It consists of two vectors: a vertical column vector  $y$  and a horizontal row vector  $\log(\hat{y})$ . The formula  $L = \sum_{pixels} \sum_{classes} y \log(\hat{y})$  is written above them. A large curly brace under the summation over pixels groups the two vectors together. Two arrows point from the first element of  $y$  to the first element of  $\log(\hat{y})$ , and another arrow points from the last element of  $\log(\hat{y})$  back to the first element of  $y$ .

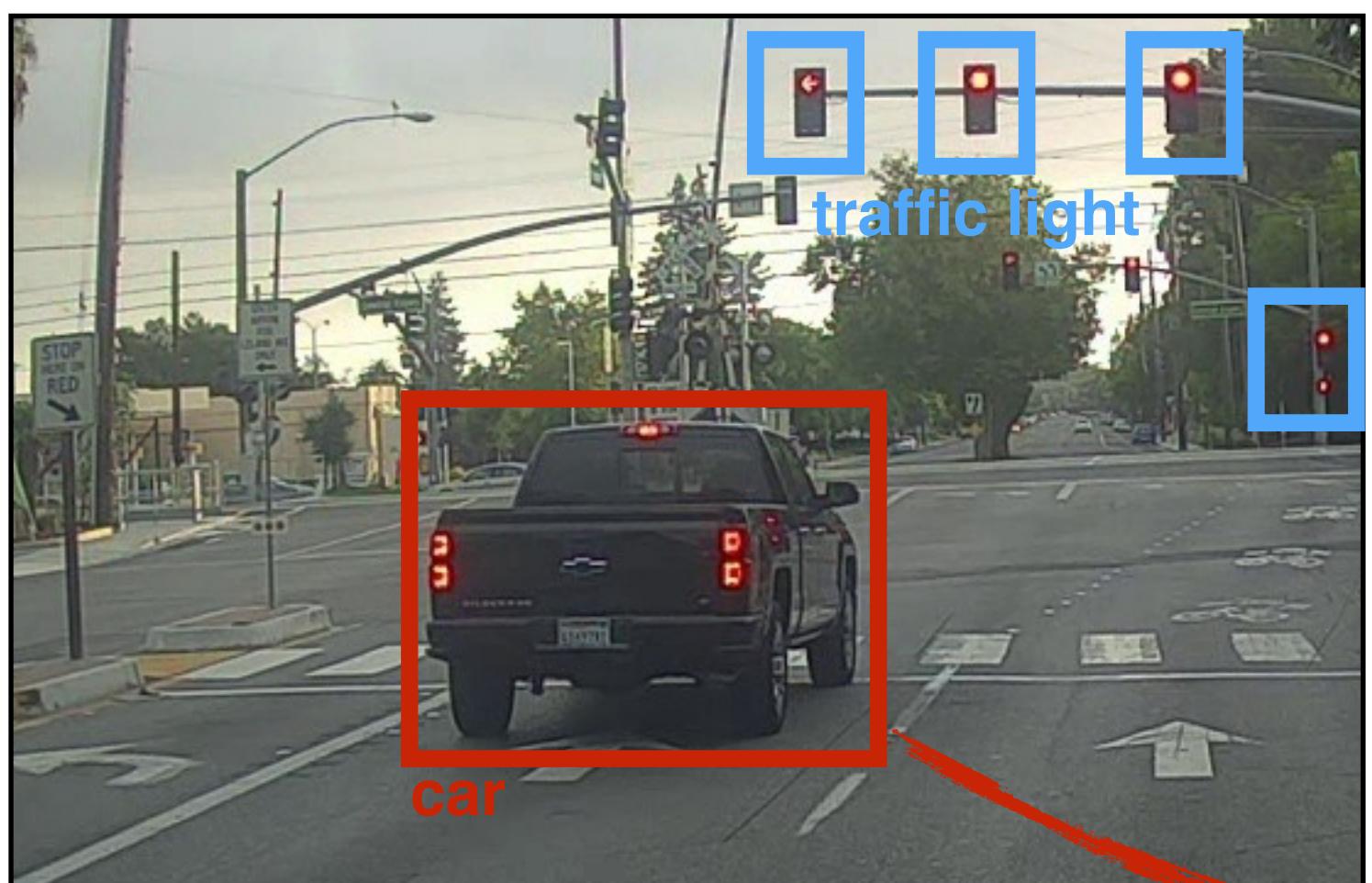
0	0.02
1	0.93
0	0.04
:	:
0	0.07
0	0.11
0	0.09

# Object Detection

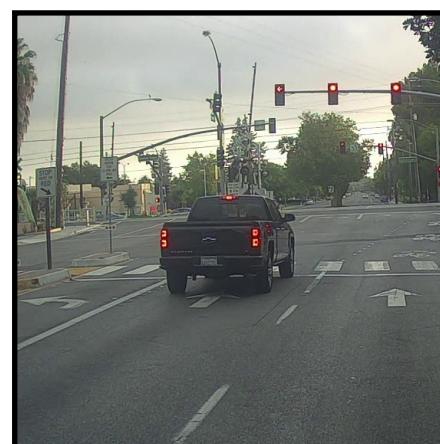
**Goal:** Find objects in images

## 1. Data?

Very large set of labelled images



## 2. Input?



## 3. Output?

$$y_1 = (b_x, b_y, b_h, b_w, p_c, c)$$

$$y_2 = (b_x, b_y, b_h, b_w, p_c, c)$$

$$\dots$$
  
$$y_k = (b_x, b_y, b_h, b_w, p_c, c)$$

Problem: size of output varies  
1. Use a mask?  
2. Change the output of the model

$$y = (b_x, b_y, b_h, b_w, p_c, c)$$

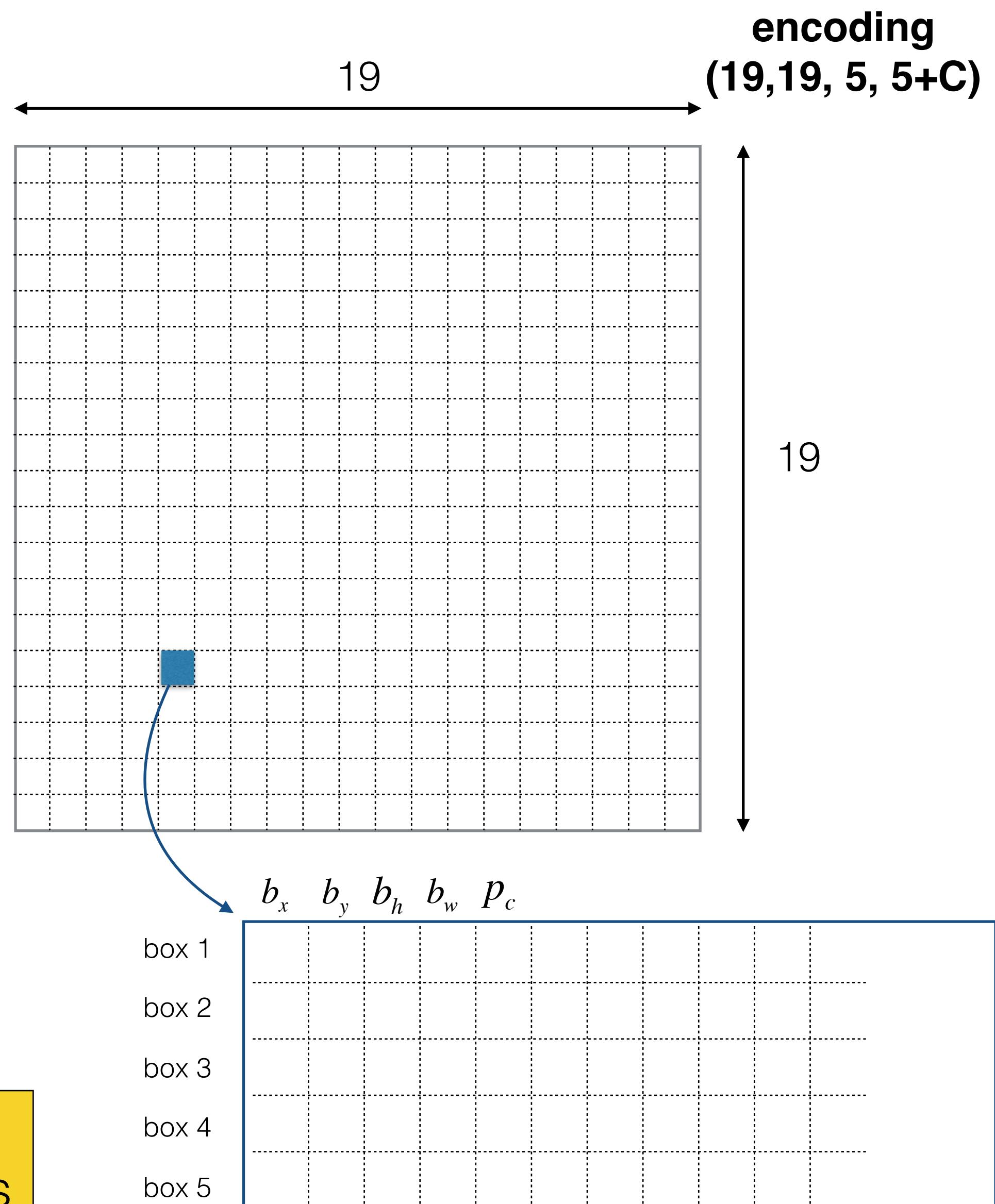
# Object Detection

## 4. Architecture?



Preprocessed image  
(608, 608, 3)

Deep CNN  
reduction  
factor: 32



We have a lot of boxes  
We select the most likely ones using thresholding and other methods

# Object Detection

## 5. Loss?

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

# Visual Question Answering

**Goal:** Find objects in images

## 1. Data?

Very large set of labelled images



## 2. Input?



## 3. Output?

$$y_1 = (b_x, b_y, b_h, b_w, p_c, c)$$

$$y_2 = (b_x, b_y, b_h, b_w, p_c, c)$$

$$\dots$$
  
$$y_k = (b_x, b_y, b_h, b_w, p_c, c)$$

Problem: size of output varies

1. Use a mask?
2. Change the output of the model

$$y = (b_x, b_y, b_h, b_w, p_c, c)$$