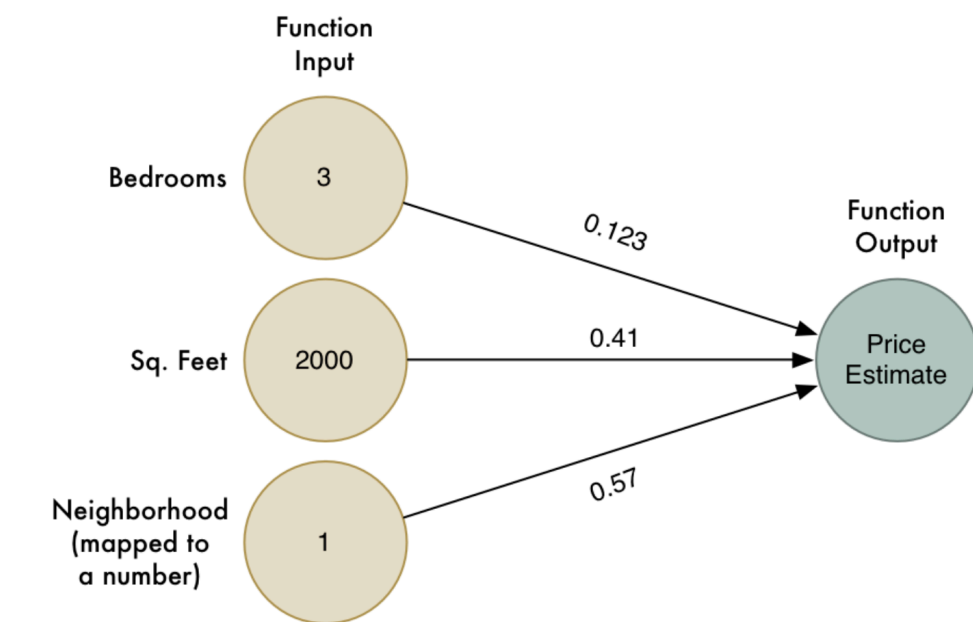# Machine Learning - Part 3

## Apr. 15, 2025

# Recap question:

A real estate agent computes their price estimate by pretending to be a neuron in a neural network as in:



Which of the following is a valid expression for the price estimate computed by a neural network?

A) $\tanh(0.123 \times 3 + 0.41 \times 2000 + 0.57 \times 1)$

B) $\tanh(3^{0.123} + 0.41 \times 2000 \times 0.57 + 1)$

# Recap question:

The answer is A!

For each neuron, we multiply the value coming in from each arrow by its corresponding weight. We add these together and then apply an **activation function** to it (here, tanh). This becomes the **value** of that neuron.

# Machine Learning - Part 3

## April 15, 2025

By the end of this lecture, you will be able to:

1. Discuss the cost of machine learning
2. Define a convolutional neural network
3. Explain how to use neural networks for facial detection

# https://playground.tensorflow.org/

# Under- and Over-fitting examples

# Deep Neural Networks (DNN)

## ORIGINAL CONTRIBUTION

# Multilayer Feedforward Networks are Universal Approximators

KURT HORNIK

Technische Universität Wien

MAXWELL STINCHCOMBE AND HALBERT WHITE

University of California, San Diego

**Abstract**—*This paper rigorously establishes that standard multilayer feedforward networks with as few as one hidden layer using arbitrary squashing functions are capable of approximating any Borel measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many hidden units are available. In this sense, multilayer feedforward networks are a class of universal approximators.*

Many researchers did not believe that neural networks are useful, until the ImageNet Challenge in 2012.

More than 14M images divided into 20 000 categories.
Task: perform classification!



mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

# ImageNet Challenge (2012)

Ranking of the best results from each team

**Error (5 predictions)**



**AlexNet**

# ImageNet Challenge

**What are the drivers of the improvements?**

1. Better hardware
2. Better training algorithms

Together, they allow us to train larger and larger networks, using more parameters, which gives better performance.

# The number of weights increases



16 weights (since there are 16 arrows)!

# The number of weights increases

# The number of weights increases

# The number of weights increases

## Shrinking deep learning's carbon footprint

Through innovation in software and hardware, researchers move to reduce the financial and environmental costs of modern artificial intelligence.

Kim Martineau | MIT Quest for Intelligence
August 7, 2020

The training of GPT-3 cost $4.6 million and 355 years in computing time on a single computer. Training smaller models than GPT-3 releases around 626,000 pounds of $CO_2$.

# Factors that affect the performance of a Neural Network:

1. Type of data (spiral data set vs the other ones)
2. The architecture of the neural network (number of hidden layers, number of neurons in each hidden layer and the choice of activation function)
3. The method used to train the network from the training data (which we will not discuss)

# Remaining topics in the class

1. Image processing and facial recognition
2. Natural Language Processing (NLP) and the Turing test
3. Generative Adversarial Networks (GAN) and Adversarial Attacks on Neural Networks

Facial recognition and image processing

# Image processing and facial recognition



Let's first see how computers can recognize the hand-written digit "8".

= [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 12, 0, 11, 39, 137, 37, 0, 152, 147, 84, 0, 0, 0, 0, 0, 1, 0, 0, 0, 41, 160, 250, 255, 235, 162, 255, 238, 206, 11, 13, 0, 0, 0, 0, 16, 9, 9, 150, 251, 45, 21, 184, 159, 154, 255, 23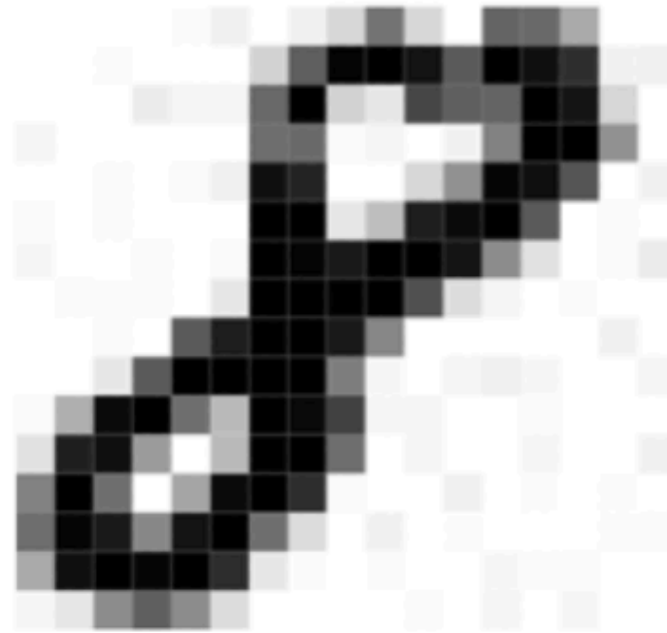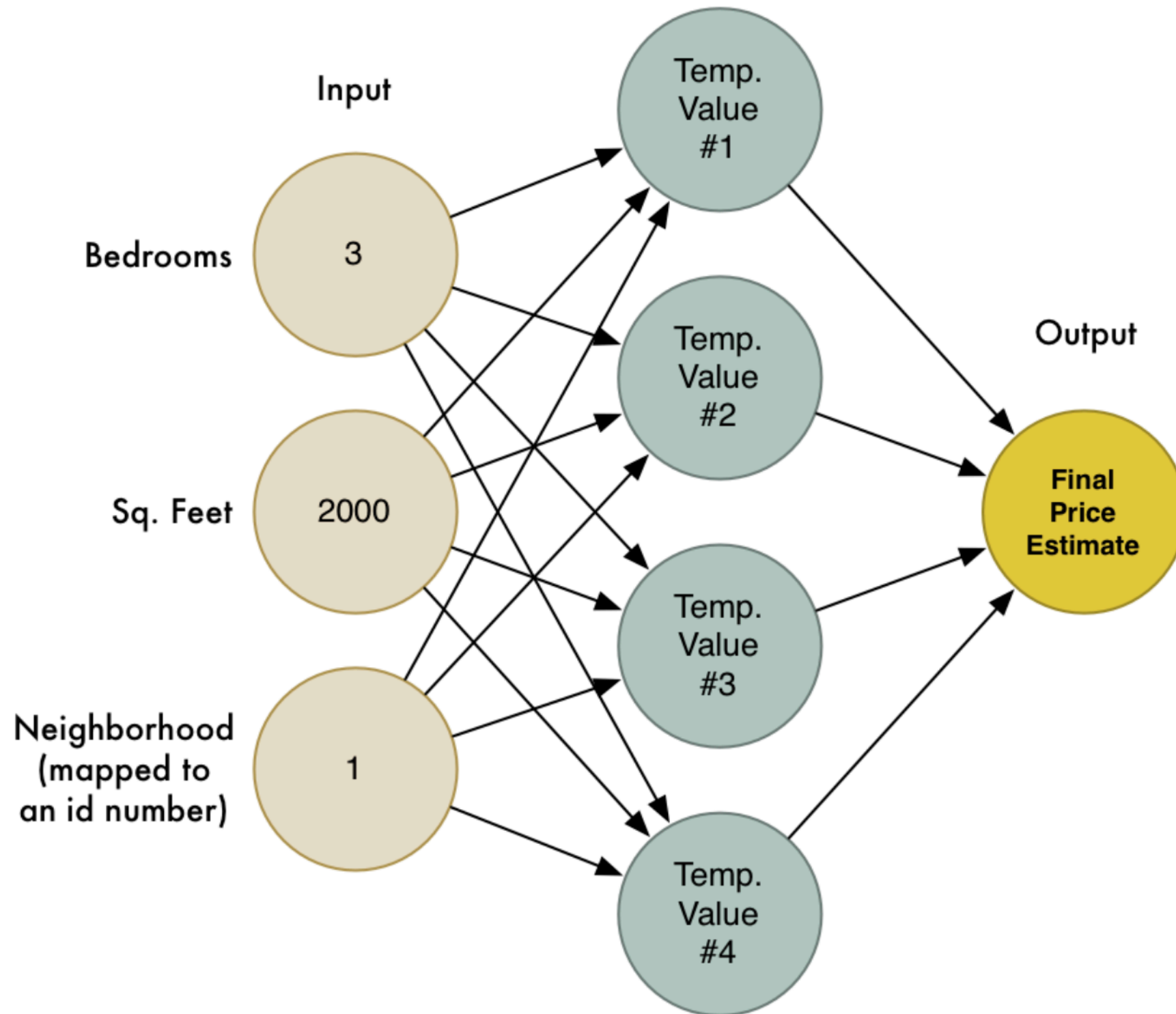3, 40, 0, 0, 10, 0, 0, 0, 0, 0, 145, 146, 3, 10, 0, 11, 124, 253, 255, 107, 0, 0, 0, 0, 3, 0, 4, 15, 236, 216, 0, 0, 38, 109, 247, 240, 169, 0, 11, 0, 1, 0, 2, 0, 0, 0, 253, 253, 23, 62, 224, 241, 255, 164, 0, 5, 0, 0, 6, 0, 0, 4, 0, 3, 252, 250, 228, 255, 255, 234, 112, 28, 0, 2, 17, 0, 0, 2, 1, 4, 0, 21, 255, 253, 251, 255, 172, 31, 8, 0, 1, 0, 0, 0, 0, 0, 4, 0, 163, 225, 251, 255, 229, 120, 0, 0, 0, 0, 0, 11, 0, 0, 0, 0, 21, 162, 255, 255, 254, 255, 126, 6, 0, 10, 14, 6, 0, 0, 9, 0, 3, 79, 242, 255, 141, 66, 255, 245, 189, 7, 8, 0, 0, 5, 0, 0, 0, 0, 26, 221, 237, 98, 0, 67, 251, 255, 144, 0, 8, 0, 0, 7, 0, 0, 11, 0, 125, 255, 141, 0, 87, 244, 255, 208, 3, 0, 0, 13, 0, 1, 0, 1, 0, 0, 145, 248, 228, 116, 235, 255, 141, 34, 0, 11, 0, 1, 0, 0, 0, 1, 3, 0, 85, 237, 253, 246, 255, 210, 21, 1, 0, 1, 0, 0, 6, 2, 4, 0, 0, 0, 6, 23, 112, 157, 114, 32, 0, 0, 0, 0, 2, 0, 8, 0, 7, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]

We can treat an 18x18 pixel image as a list of 18x18=324 numbers, containing the greyscale level of each pixel.

# Designing a neural network

# Designing a neural network



Too small!

# Designing a neural network

Input Nodes

Intermediate Values

Outputs

**Likelihood is an 8!**

**Likelihood not an 8!**

# Training the neural net



When we feed in an "8", we'll tell it the probability the image is an "8" is 100% and the probability it's not an "8" is 0%. Vice versa for the counter-example images.

# Testing the neural net



**Test Image #1** → **Prediction from our network**

100% an "8"!

**Test Image #2** → **Prediction from our network**

100% not an "8"!

The good news is that our "8" recognizer does work well on simple images where the letter is right in the middle of the image.

# Testing the neural net

Test
Image #1

Prediction from
our network

No idea!?!

Test
Image #2

Prediction from
our network

What is this?!@

The bad news is that our "8" recognizer *totally fails* to work when the letter isn't perfectly centered in the image. Just the slightest position change ruins everything.

# Testing the neural net

Test Image #1 → Prediction from our network: No idea!?!

Test Image #2 → Prediction from our network: What is this?!@

The bad news is that our "8" recognizer *totally fails* to work when the letter isn't perfectly centered in the image. Just the slightest position change ruins everything.
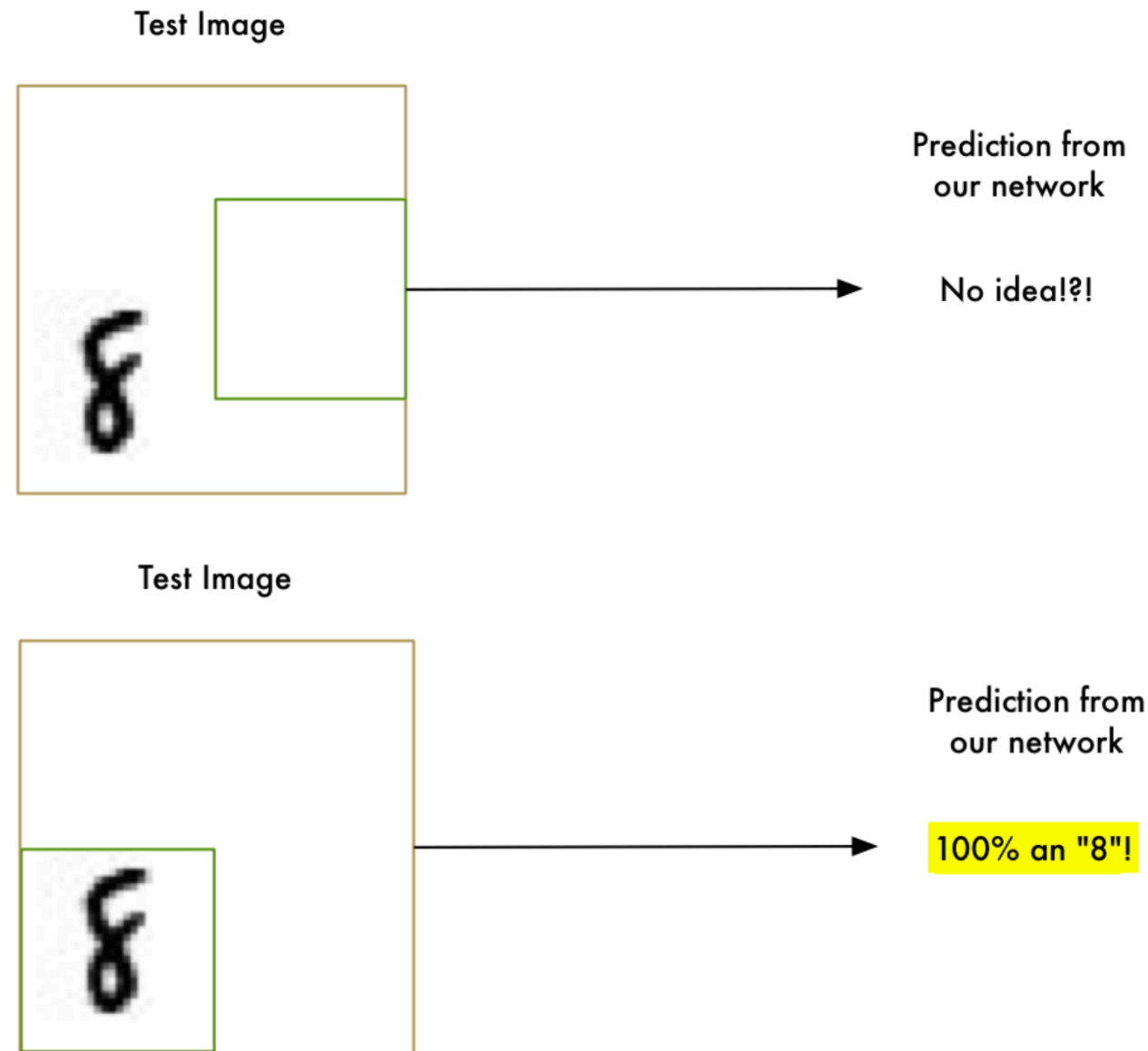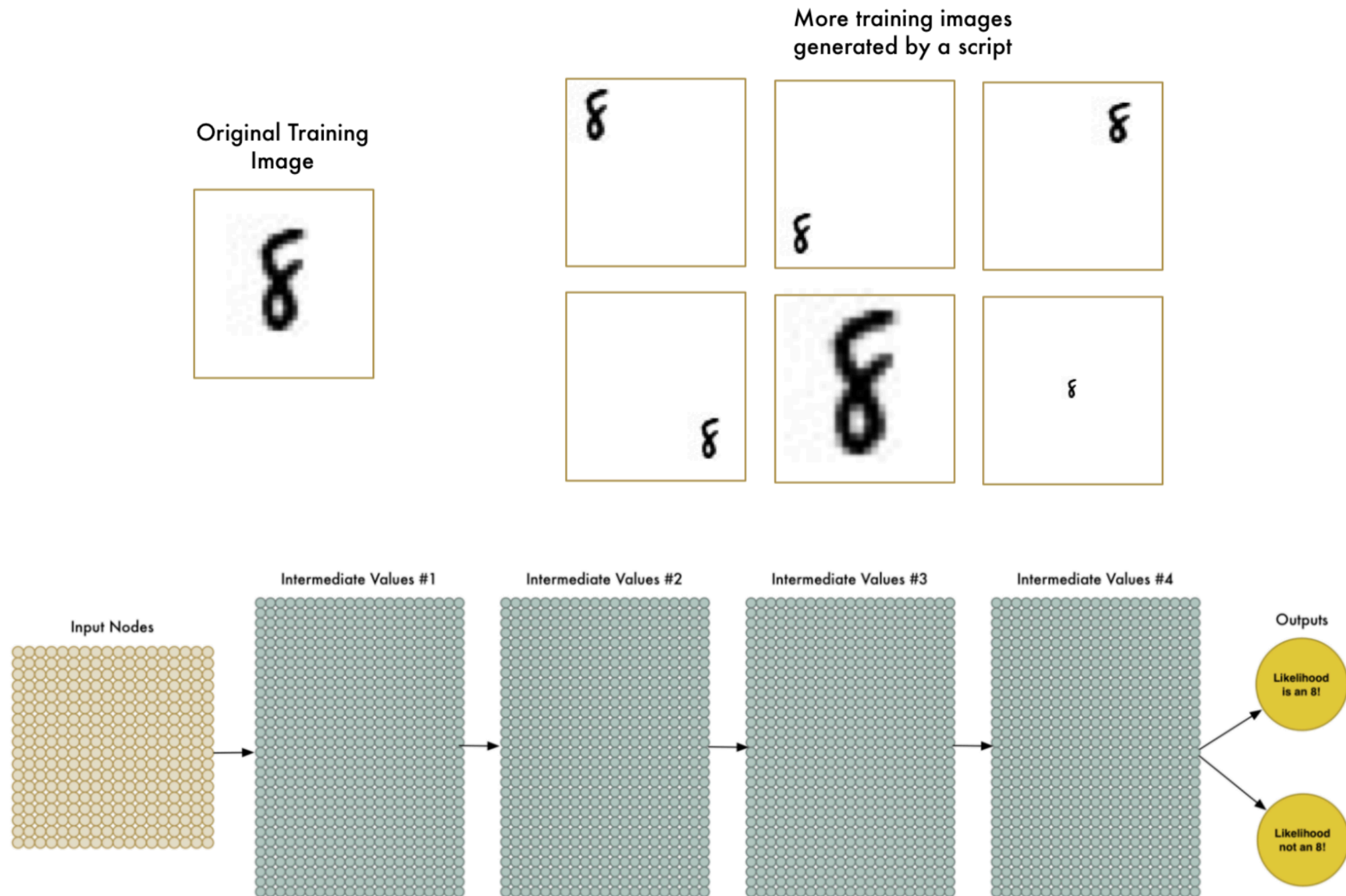
How can we mitigate this?

# Idea #1: Searching with a sliding window



Very inefficient!

# Idea #2: More data and a deep neural net



Very inefficient!
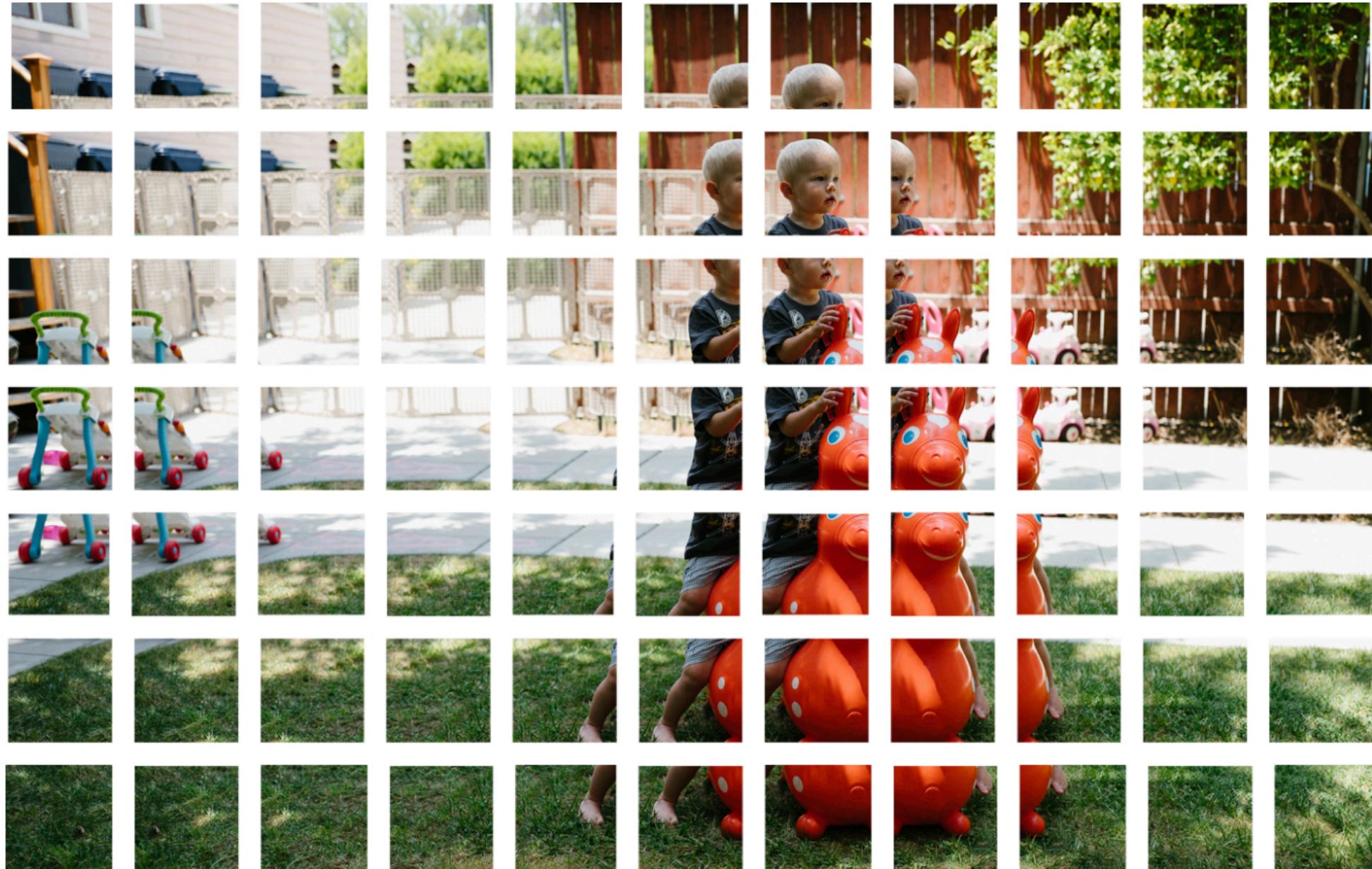
# The solution: convolution

For now, our neural network thinks that an "8" in a different part of the image is an entirely different thing. It doesn't understand that moving an object around in the picture doesn't make it something different.

We need to give our neural network understanding of **translation invariance**. This can be done using a process called **convolution**.
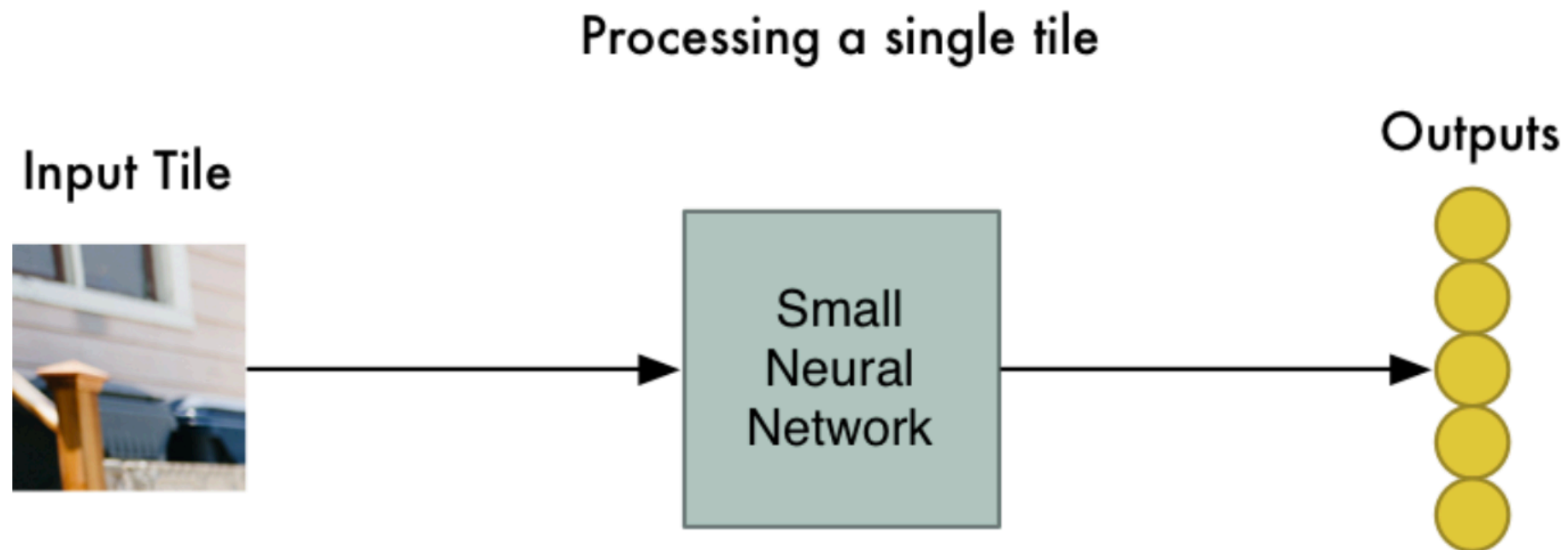
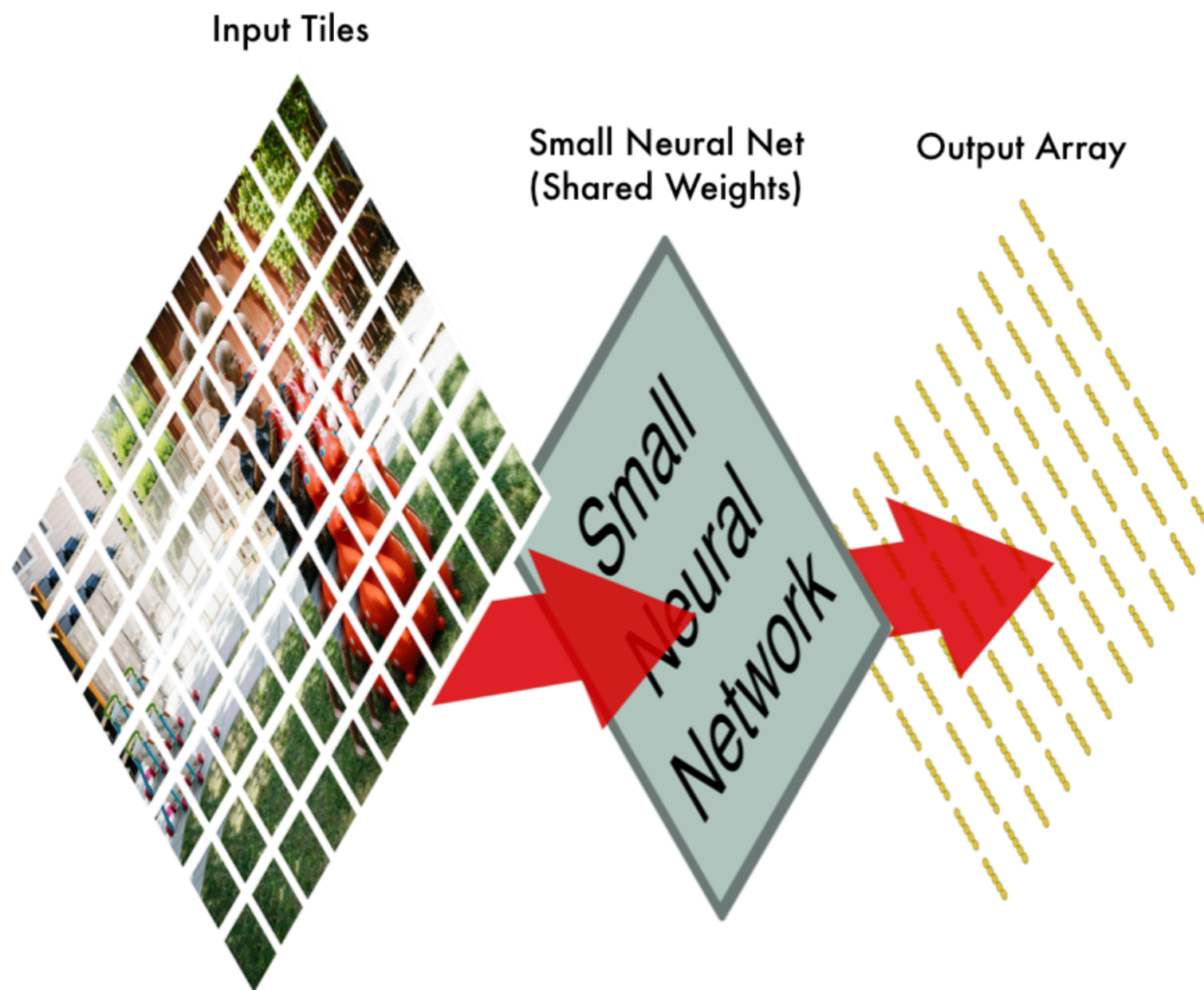**Example:** Find a face in the following picture

**Step 1:** Break the image into overlapping image tiles

**Step 2:** Feed each image tile into a small neural network. We'll keep the same neural network weights for every single tile in the same original image. In other words, we are treating every image tile equally.
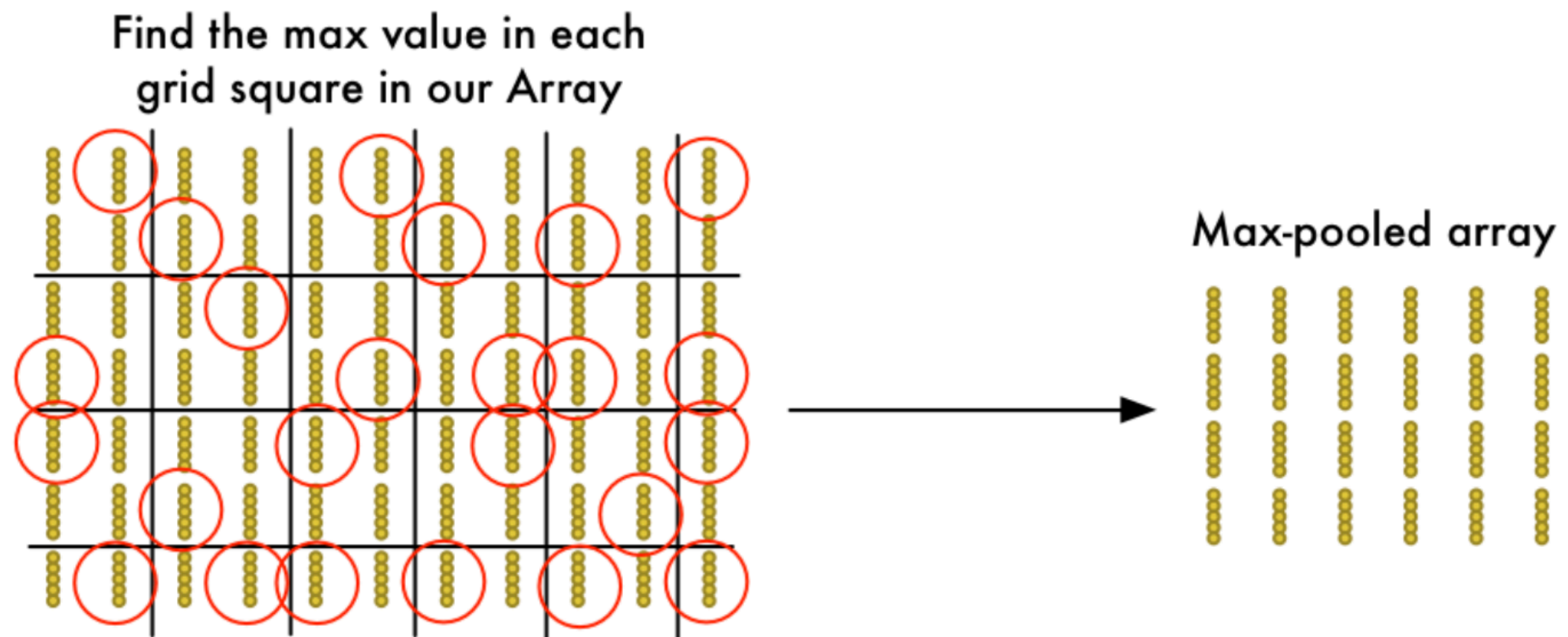


Processing a single tile

**Step 3:** Save the results from each tile into a new array in the same arrangement as the original image.
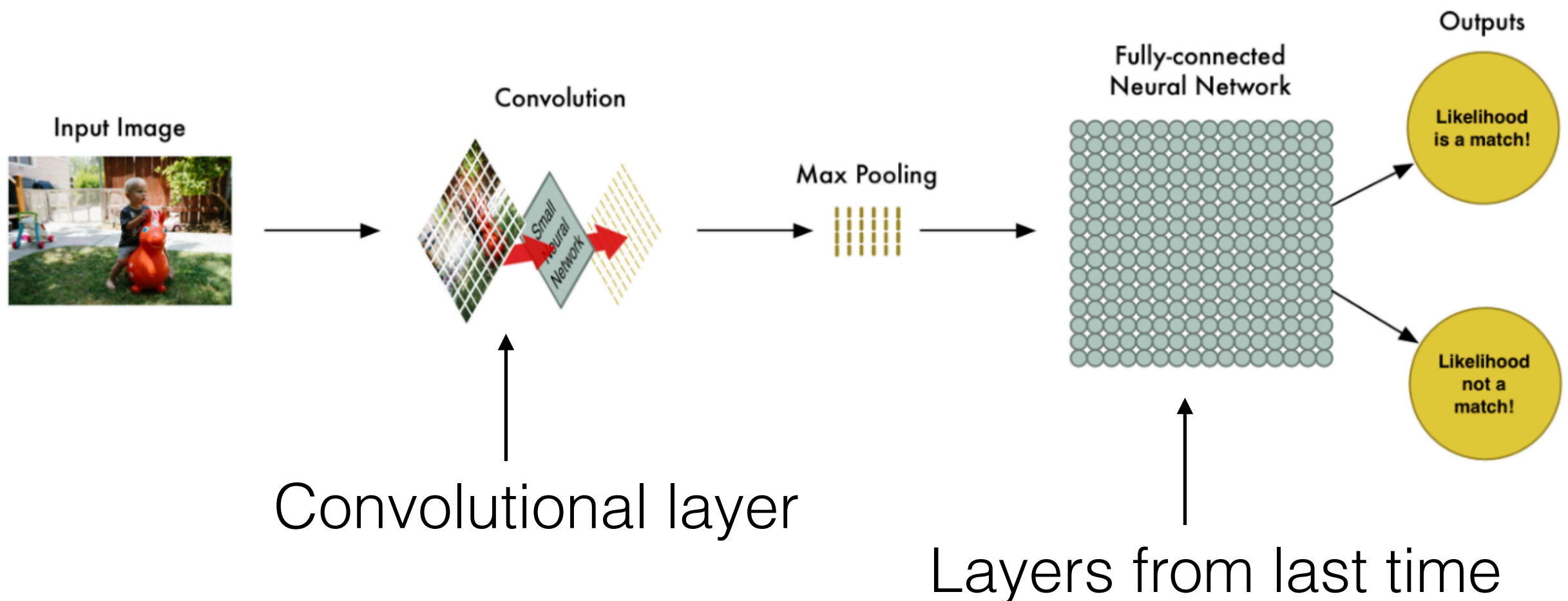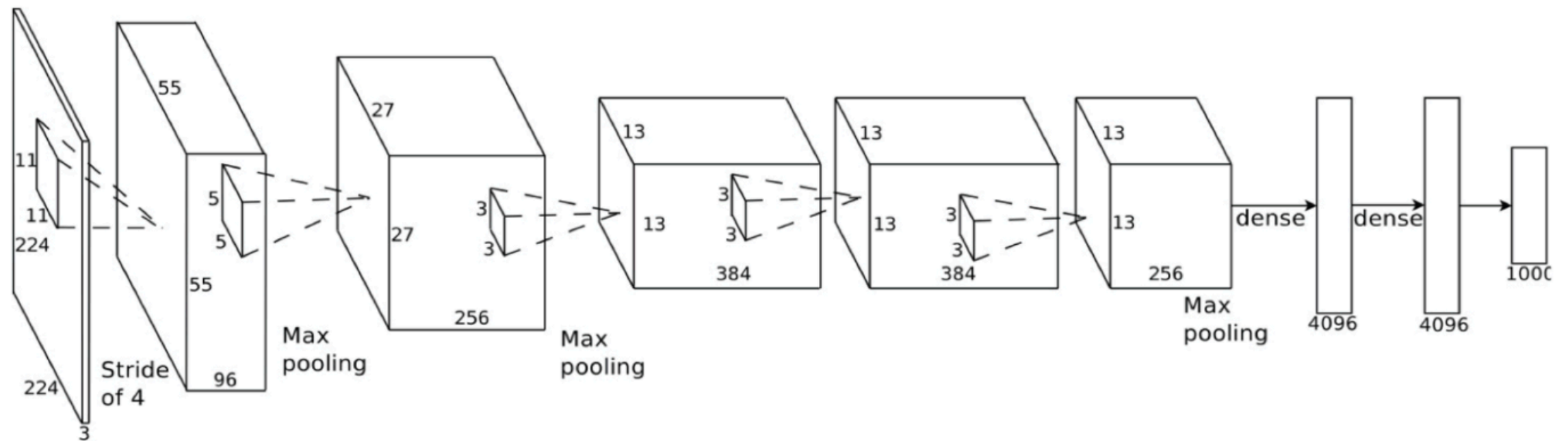
**Step 4:** Downsampling: The result of Step 3 was an array that maps out which parts of the original image are the most interesting. But that array is still pretty big. To reduce its size, we *downsample* it using an algorithm called *max pooling*.

We look at each 2x2 square of the array and keep the biggest number.



Find the max value in each grid square in our Array

Max-pooled array

**Final step:** Make a prediction. We can use that small array as input into *another neural network*. This final neural network will decide if the image is or isn't a match.
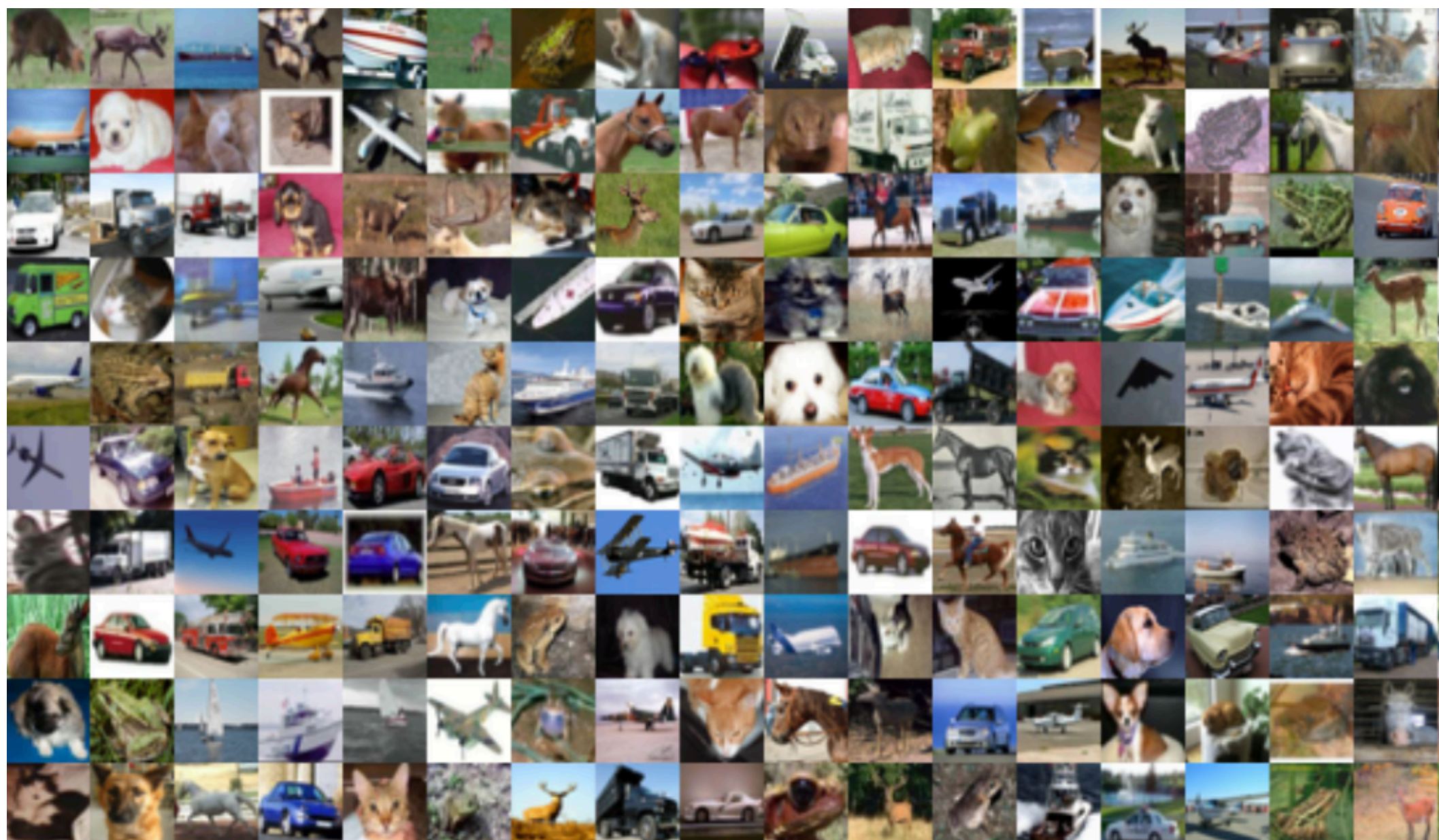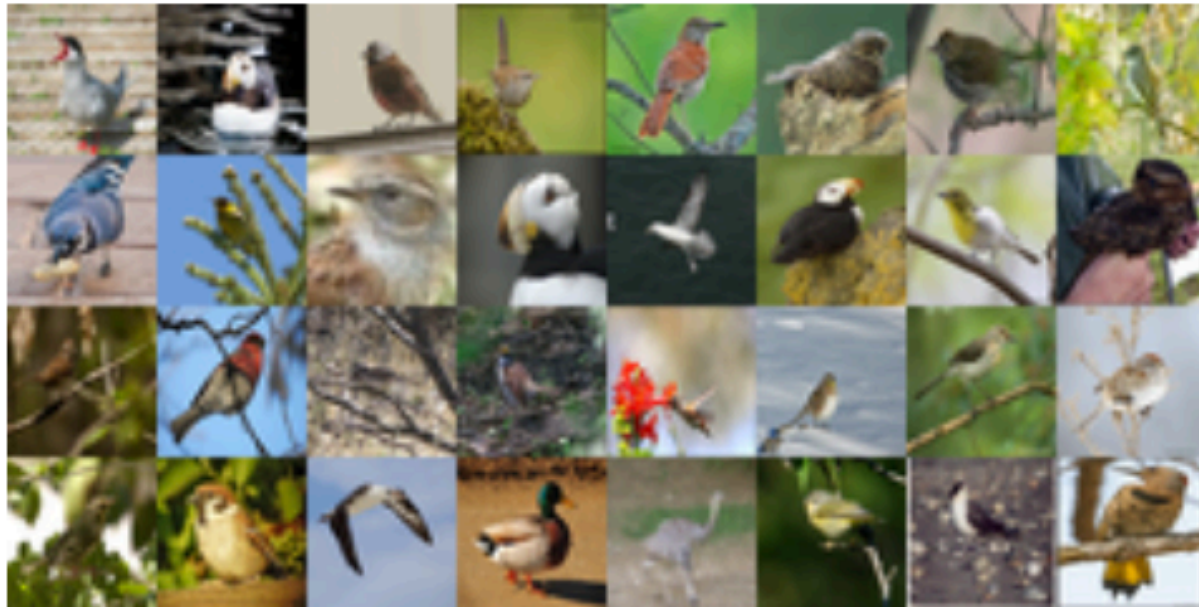
When solving problems in the real world, these steps can be combined and stacked as many times as we want to address tasks of different complexity.

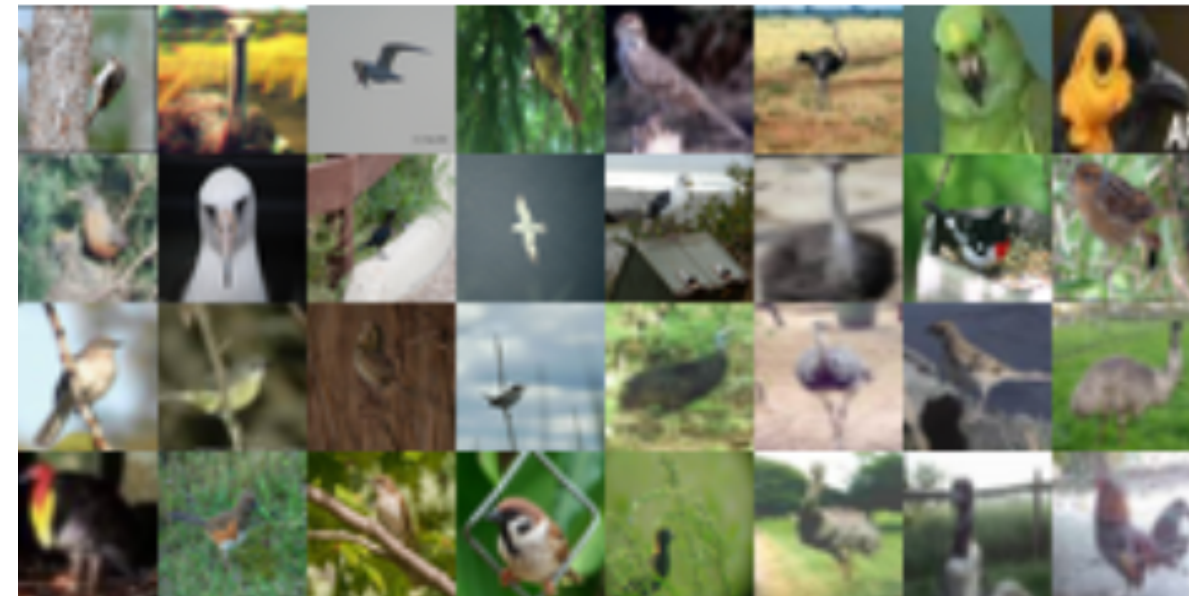# Building a bird classifier

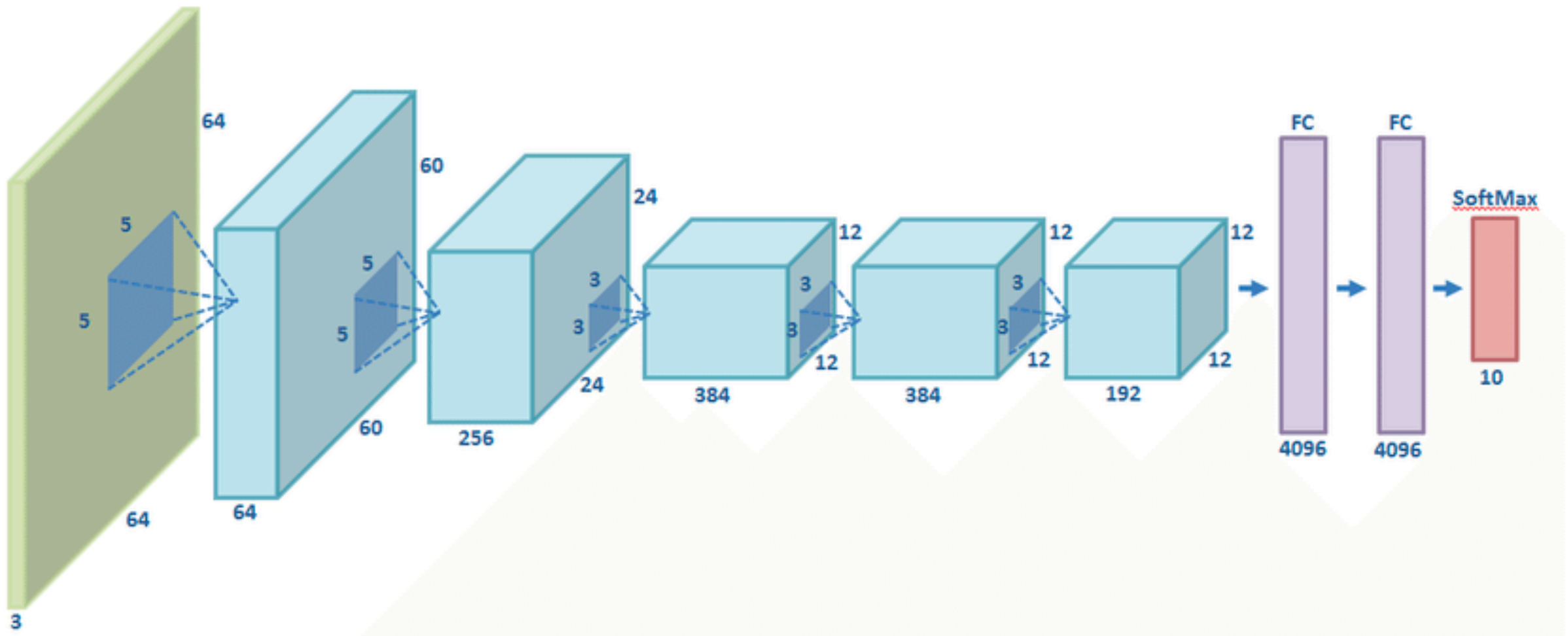# True positives

# True negatives

# False positives

# False negatives

# Results for 15,000 Validation Images

*(6000 images are birds, 9000 images are not birds)*

|  | Predicted 'bird' | Predicted 'not a bird' |
|---|---|---|
| **Bird** | 5,450 <br> *True Positves* | 550 <br> *False Negatives* |
| **Not a Bird** | 162 <br> *False Positives* | 8,838 <br> *True Negatives* |

| **Precision** <br> *If we predicted 'bird', how often was it really a bird?* | 97.11% <br> *(True Positives ÷ All Positive Guesses)* |
|---|---|
| **Recall** <br> *What percentage of the actual birds did we find?* | 90.83% <br> *(True Positives ÷ Total Birds in Dataset)* |

# AlexNet (16% error on ImageNet)

# VGG16 (7.3% error on ImageNet)



conv1

conv2

conv3

conv4

conv5

fc6    fc7    fc8

$1 \times 1 \times 4096$    $1 \times 1 \times 1000$

$14 \times 14 \times 512$

$28 \times 28 \times 512$

$7 \times 7 \times 512$

$56 \times 56 \times 256$

$112 \times 112 \times 128$

$224 \times 224 \times 64$

convolution+ReLU

max pooling

fully connected+ReLU

# ResNet (3.6% error on ImageNet)

# Convolutional layers extract features from an image



Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier

**Example:** Facebook automatically tags the people in your photos that you only tagged a couple of times before.

# Facial recognition



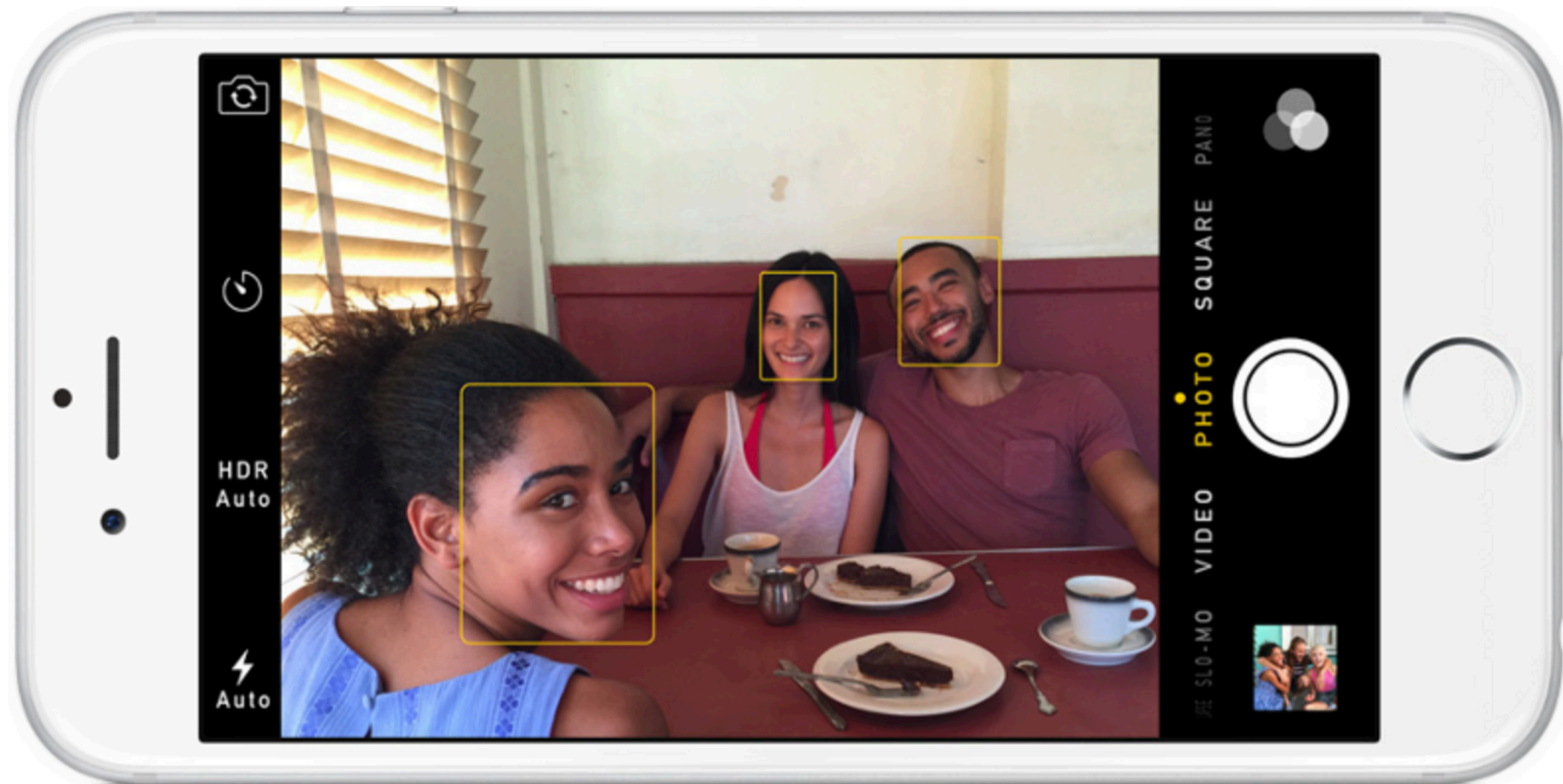Chad Smith (a famous rock musician)     Will Ferrell (a famous actor)

# Steps for face recognition

1. Look at the image and find all the faces in it.
2. Focus on each face and be able to understand that even if a face is turned in a weird direction or in bad lighting, it is still the same person.
3. Highlight the unique features of the face that can be used to distinguish it from other people - for example, how big the eyes are or how long the face is, etc.
4. Compare these unique features of a person with the features of other people you know to determine the name of the person.
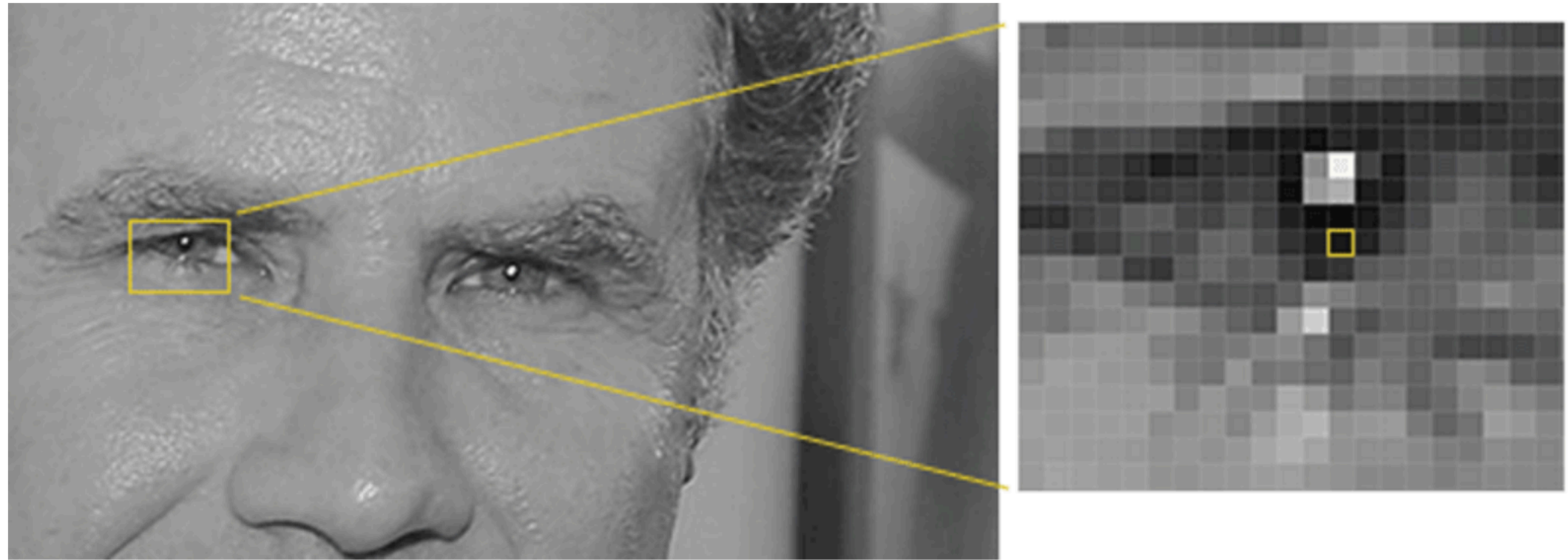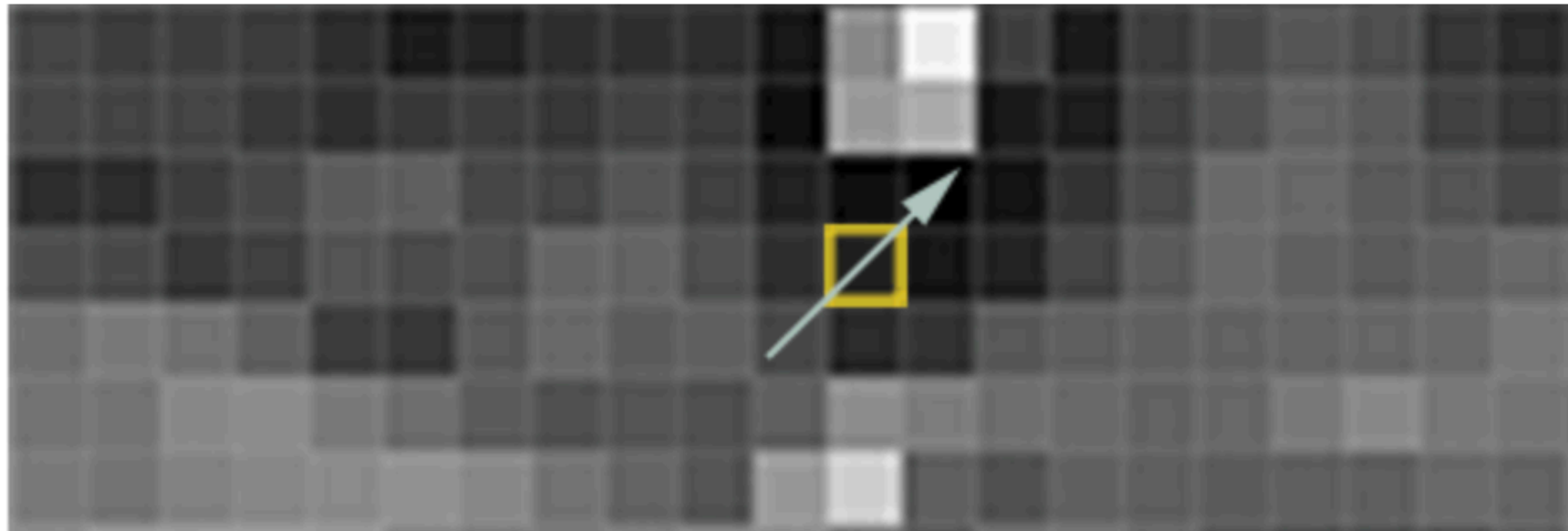
# Step 1: Finding all faces

Face detection became the mainstream in the early 2000s when Paul Viola and Michael Jones invented a way to detect faces that was fast enough to work on cheap cameras.

However, there are now many more reliable solutions. We are going to use a method discovered in 2005 called histogram of oriented gradients (HOG).
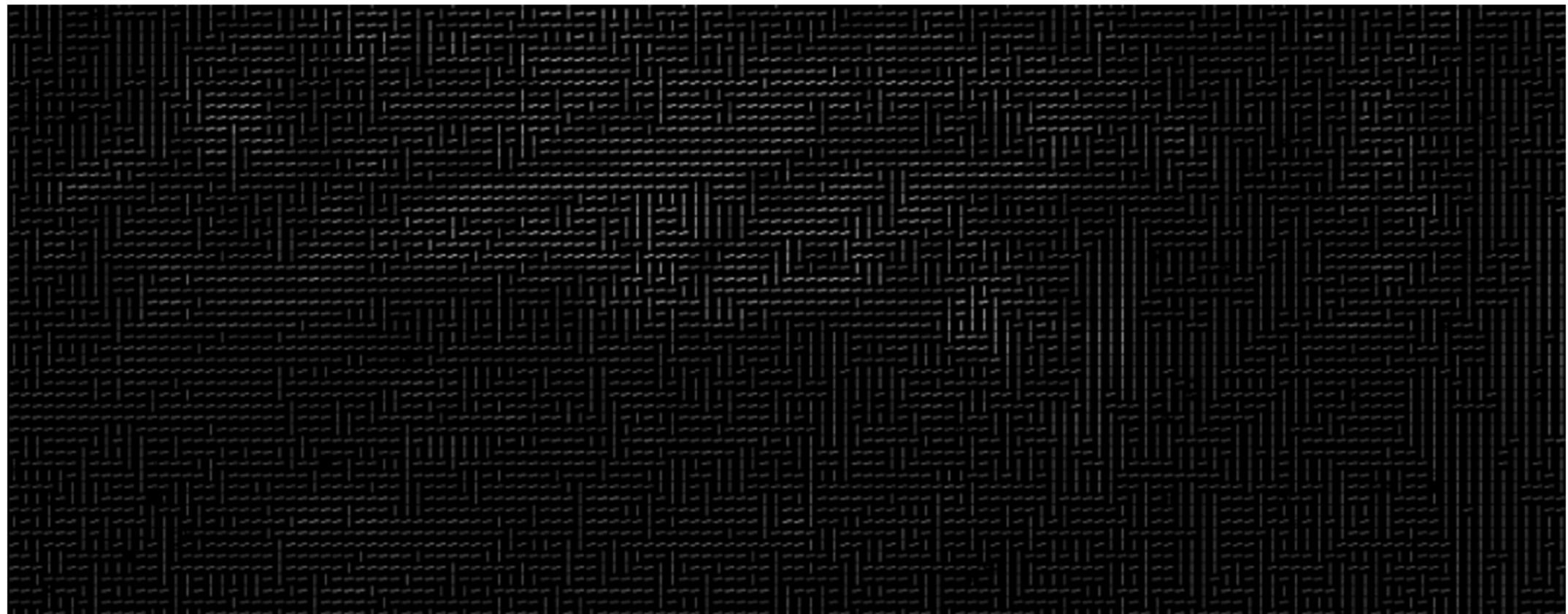
We look at each individual pixel in our image sequentially. For each individual pixel, its immediate environment should be considered.

Our goal is to highlight how dark the current pixel is in comparison with the pixels directly adjacent to it. Then draw an arrow showing the direction in which the image becomes darker.

If we repeat this process for each individual pixel in the image, then, in the end, each pixel will be replaced by an arrow. These arrows are called *gradients*, and they show the flow from light to dark throughout the image.
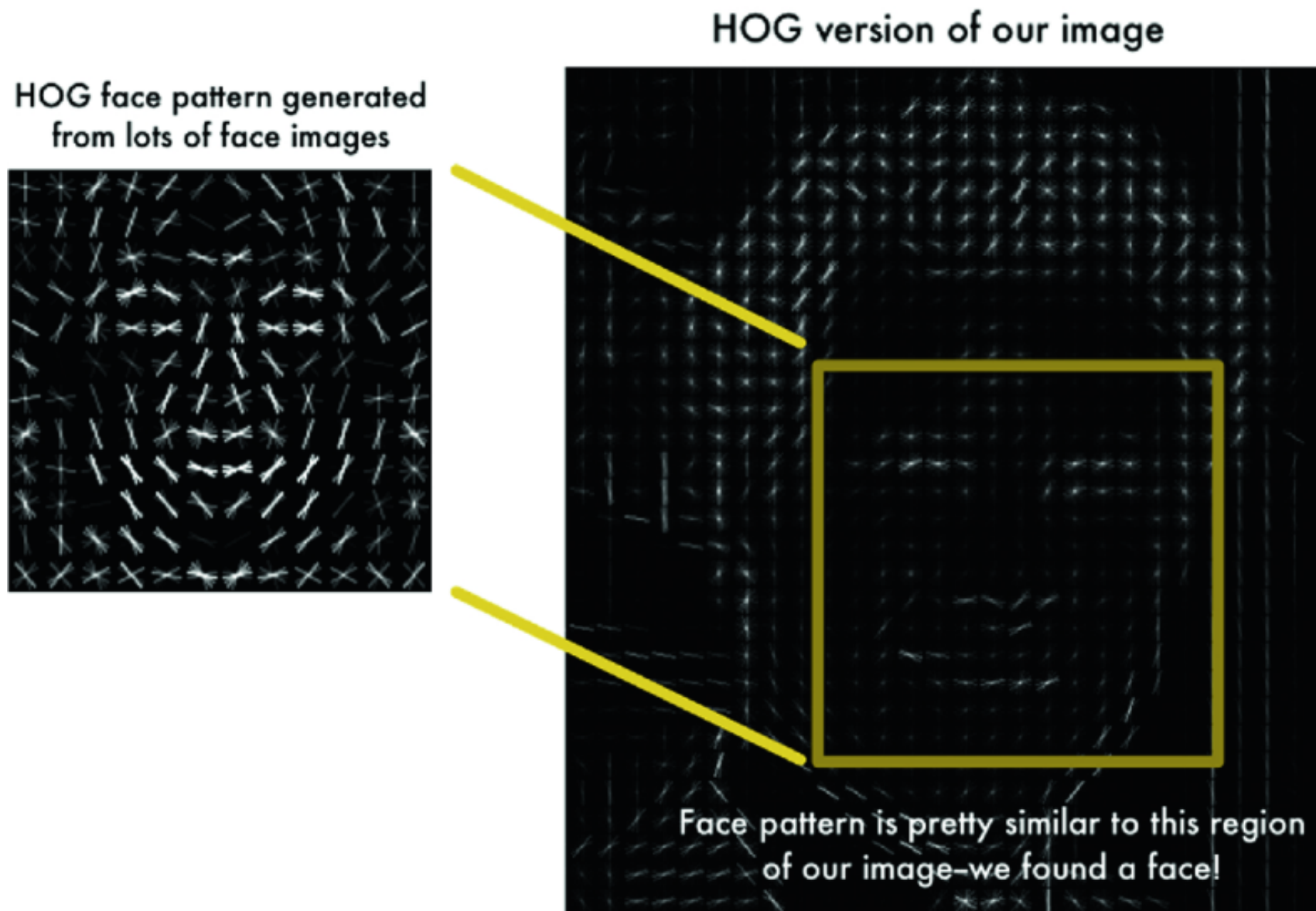
When we analyze pixels directly, dark and light images of the same person will have very different pixel intensities. But if we consider only the directional changes in brightness, then both dark and light images will have exactly the same representation.

However, preserving the gradient for each individual pixel gives us way too many details, so we need a way to keep only the basic structure of the image.
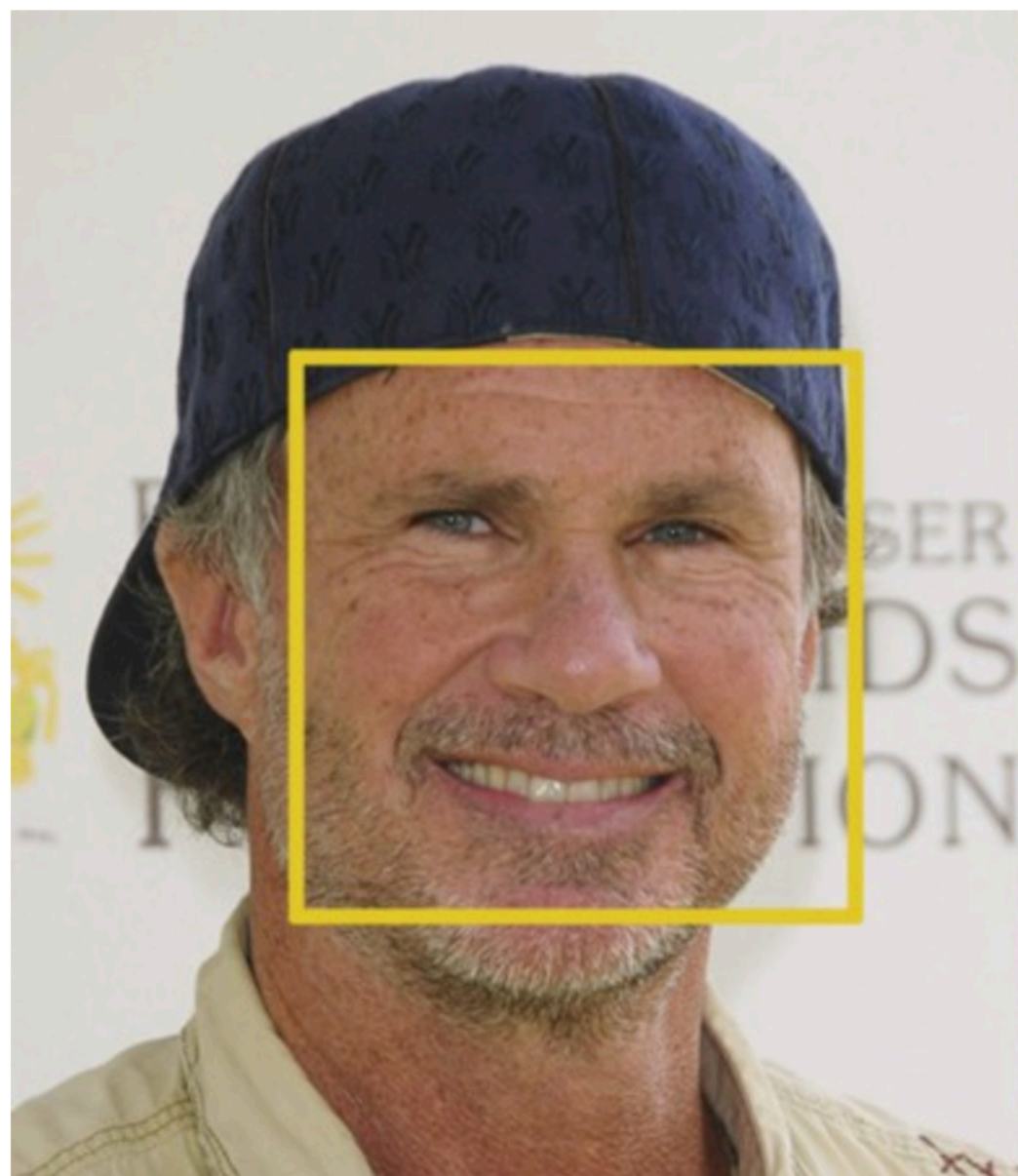
# Solution

1. Break the image into small squares of 16x16 pixels in each.
2. In each square, you should calculate how many gradient arrows are shown in each main direction (i.e. how many arrows are directed up, up right, right, etc.).
3. Replace the square under consideration by an arrow with the direction prevailing in this square.

HOG face pattern generated from lots of face images

HOG version of our image

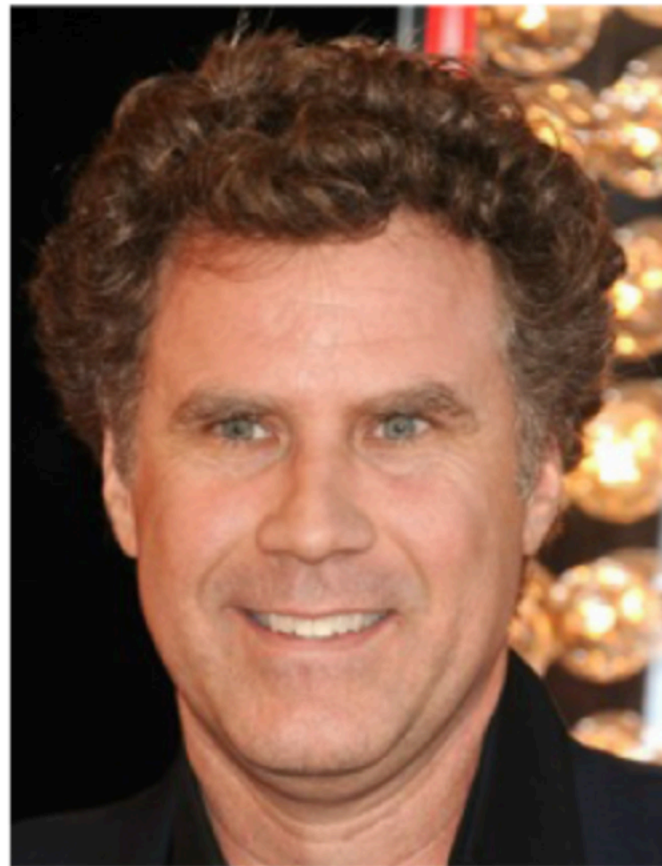Face pattern is pretty similar to this region of our image—we found a face!

To detect faces in this HOG image, all that is required is to find a part of the image that is most similar to the known HOG structure obtained from the group of faces used for training.
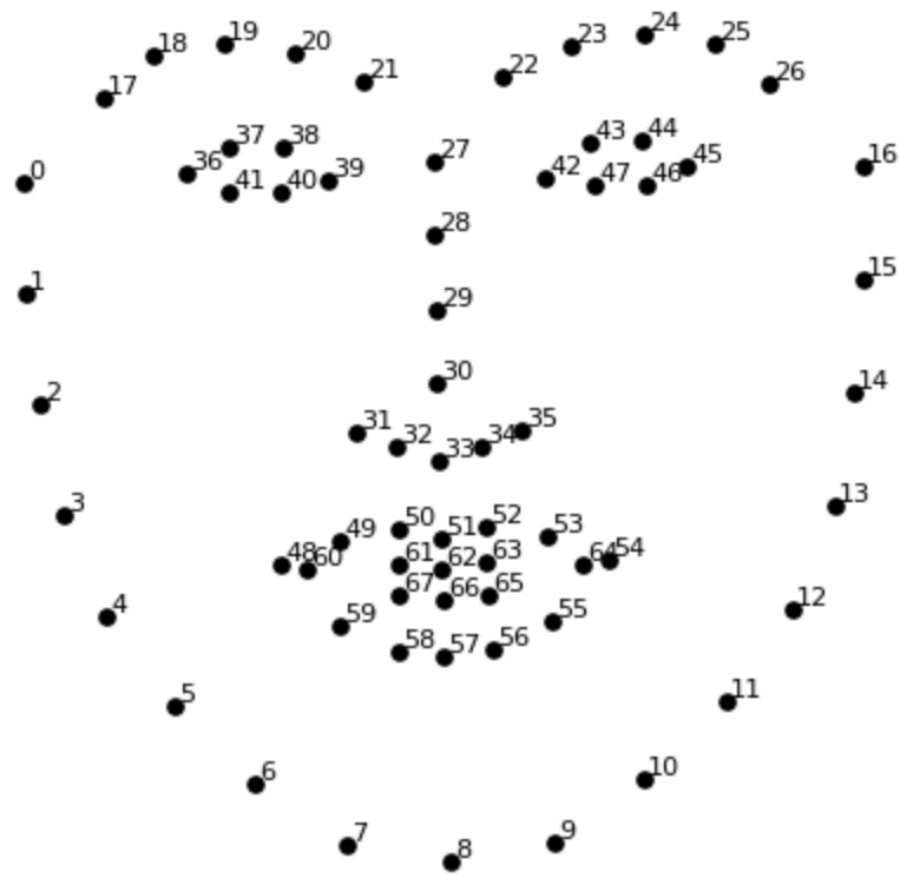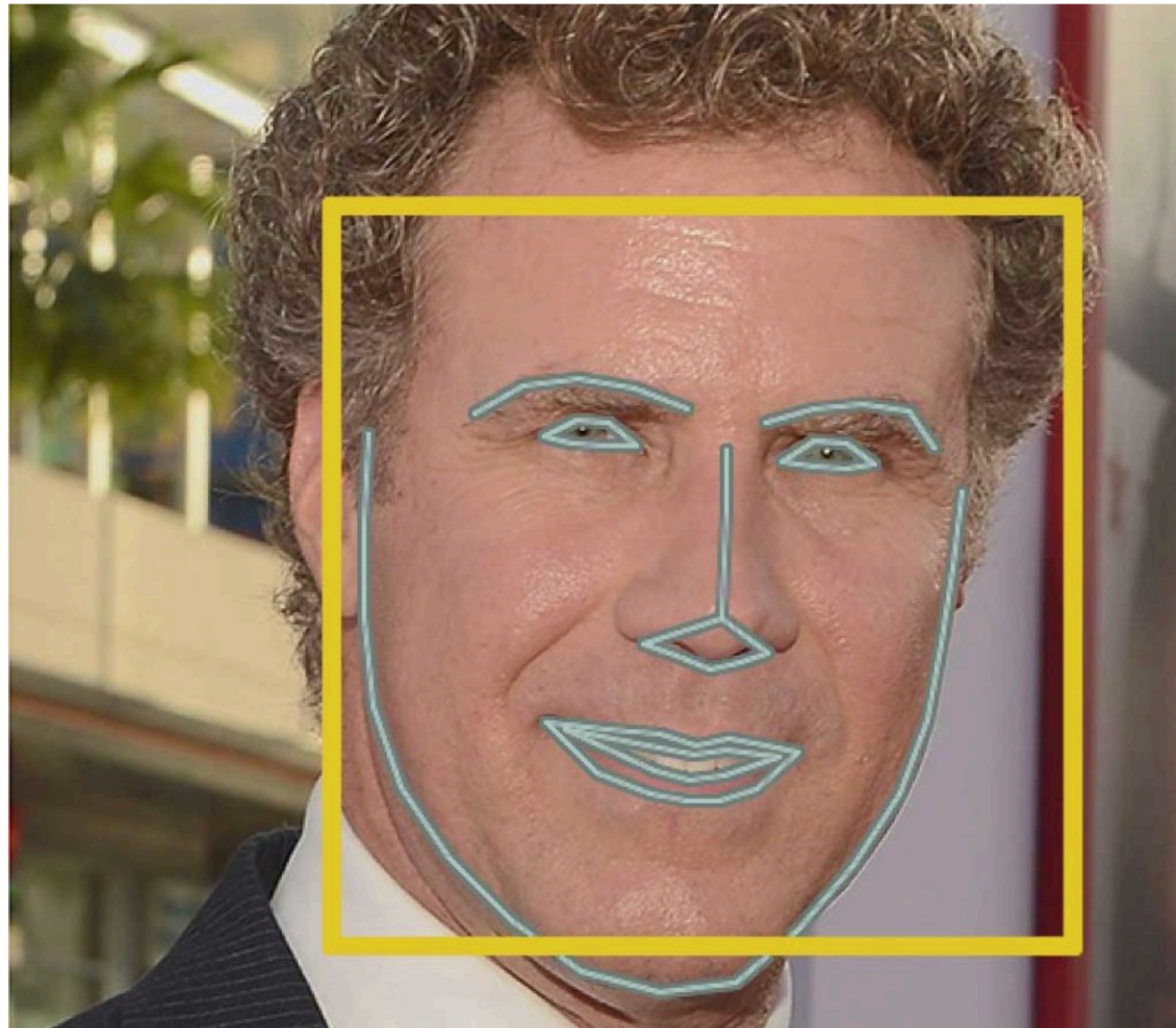
# Step 2: Location and display of faces



People can easily see that both images refer to the same person, but computers will consider them as faces of two different people.

To address this issue, we will transform each image so that the eyes and lips are always in the same place on the image.



We use a face landmark estimation algorithm called "assessment of anthropometric points."

The main idea is that 68 specific points (marks) are found on each face - the protruding part of the chin, the outer edge of each eye, the inner edge of each eyebrow, etc.
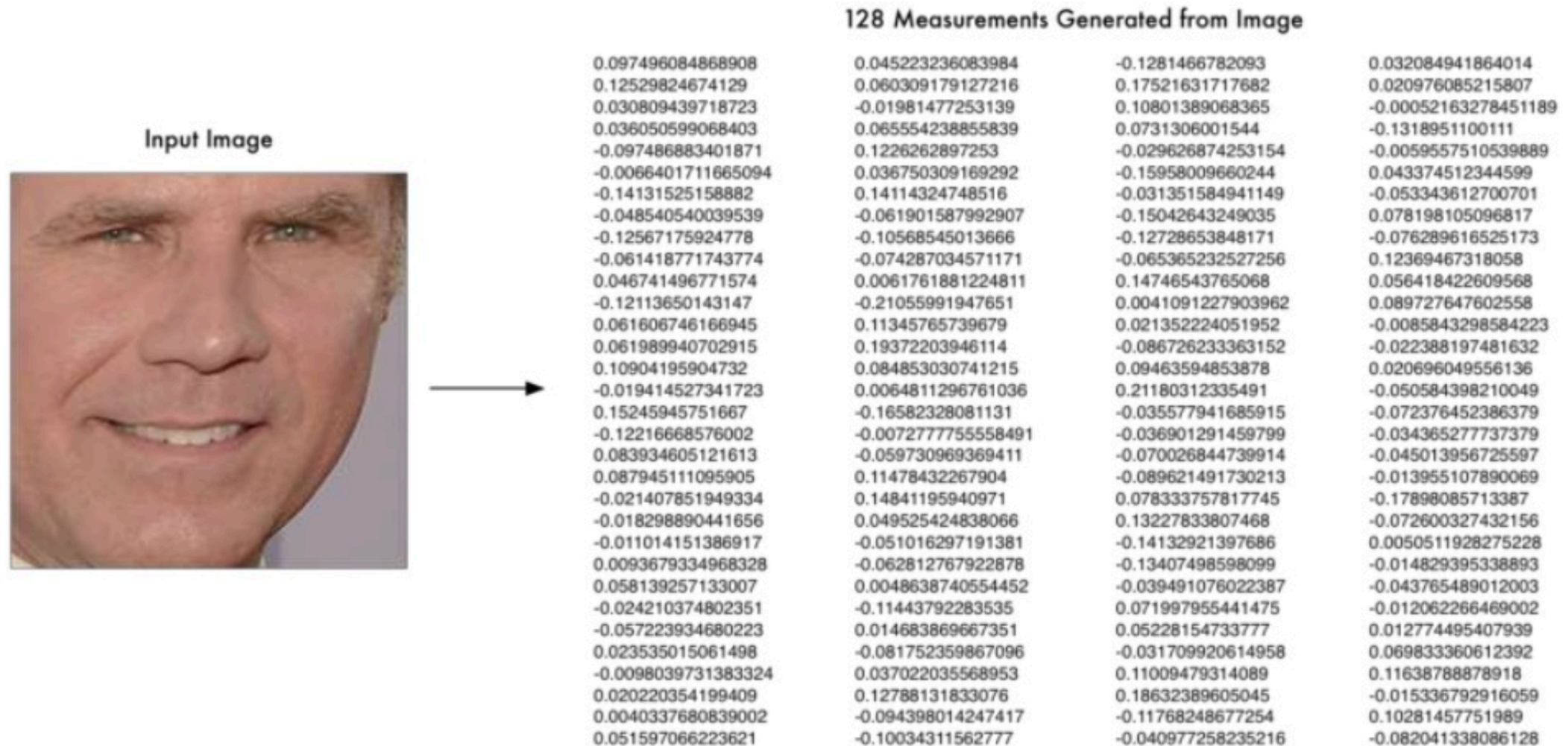
Now that we know where the eyes and mouth are, we'll just rotate, scale and move the image so that the eyes and mouth are centered as best as possible.
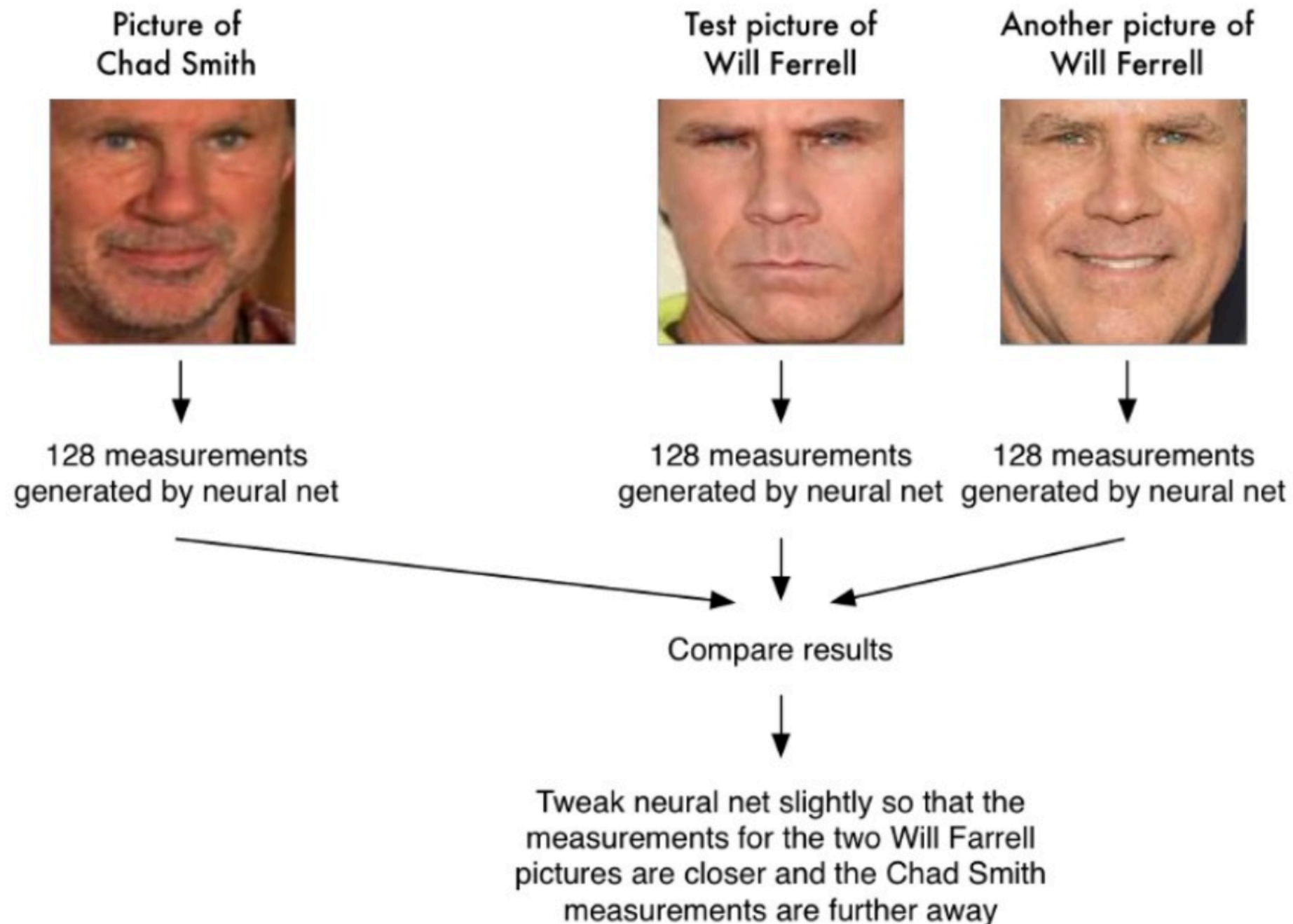
# Steps for face recognition

1. Look at the image and find all the faces in it.
2. Focus on each face and be able to understand that even if a face is turned in a weird direction or in bad lighting, it is still the same person.
3. Highlight the unique features of the face that can be used to distinguish it from other people - for example, how big the eyes are or how long the face is, etc.
4. Compare these unique features of a person with the features of other people you know to determine the name of the person.

# Step 3: Encoding Faces



Train a neural network to generate 128 characteristics for each face.

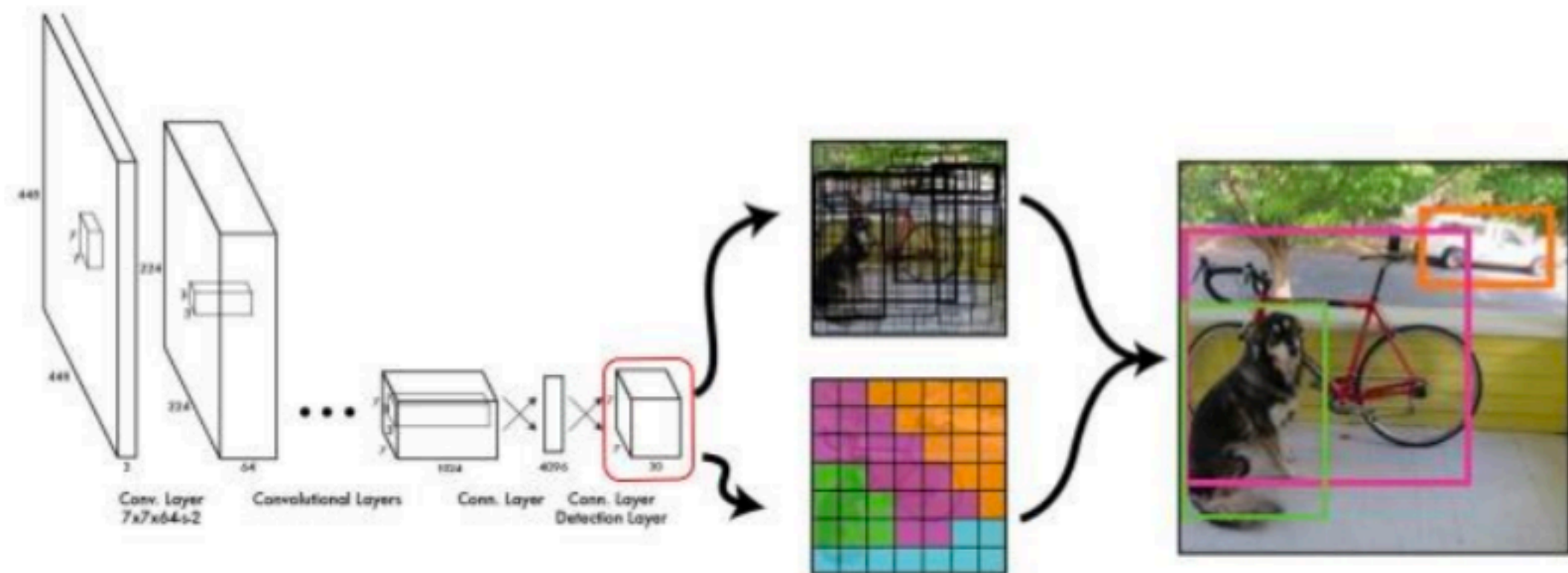# A single "triplet" training step

# Step 4: Finding a person's name after encoding a face

This step is the easiest one in the pipeline. We compare a given face encoding with a list of known face encodings and find the best match. The result gives the name of the person.

All steps can be combined into a single Neural Network!

# YOLO: You Only Look Once

# You Only Look Once:
# Unified, Real-Time Object Detection

Joseph Redmon*[†], Santosh Divvala*[†], Ross Girshick[¶], Ali Farhadi*[†]

University of Washington*, Allen Institute for AI[†], Facebook AI Research[¶]

http://pjreddie.com/yolo/

## Abstract

We present YOLO, a new approach to object detection. Prior work on object detection repurposes classifiers to perform detection. Instead, we frame object detection as a regression problem to spatially separated bounding boxes and associated class probabilities. A single neural network predicts bounding boxes and class probabilities directly from full images in one evaluation. Since the whole detection pipeline is a single network, it can be optimized end-to-end directly on detection performance.
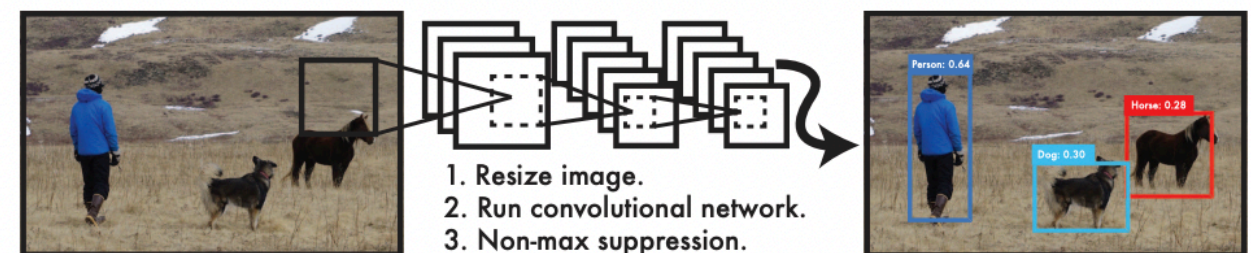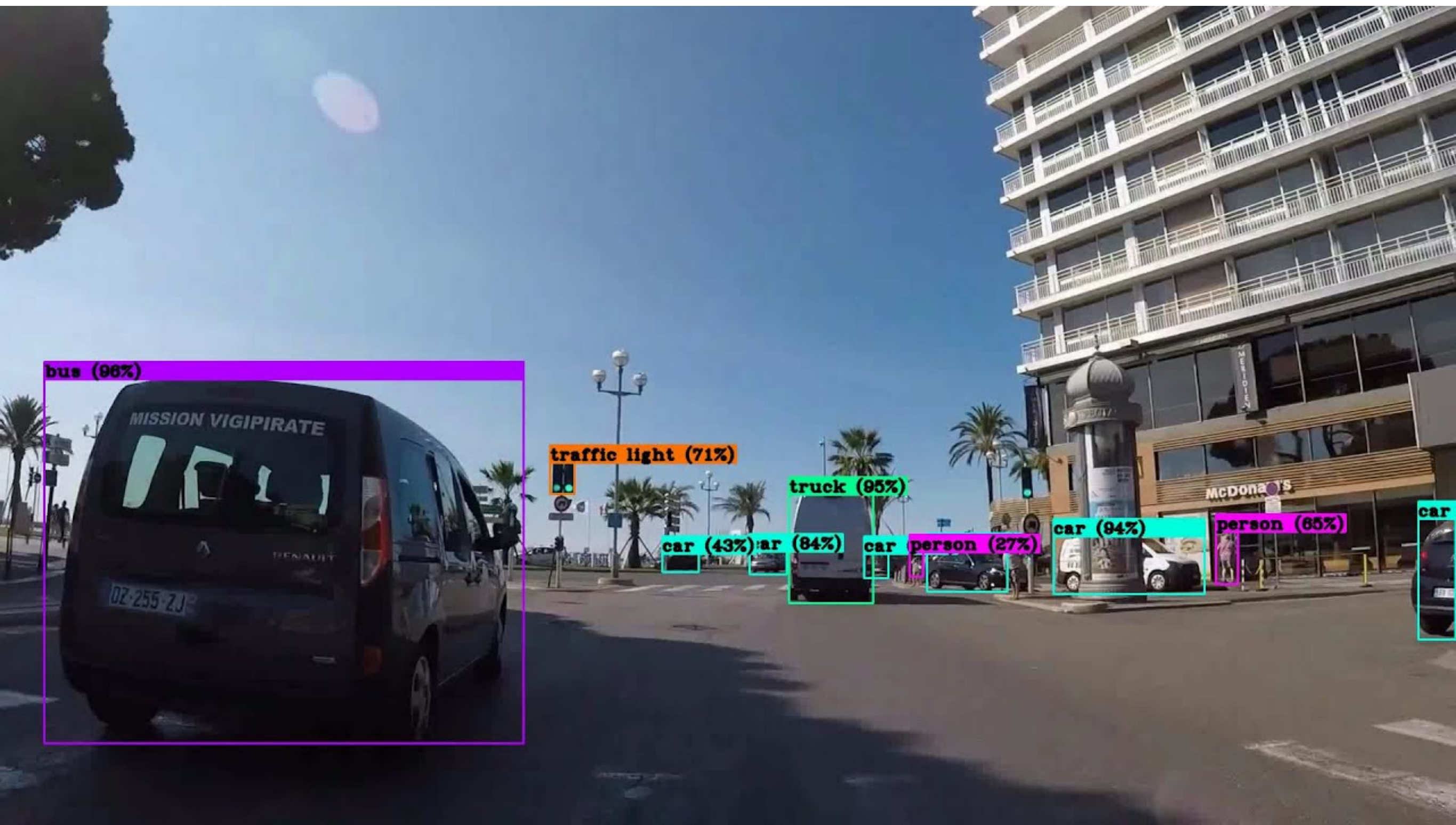
**Figure 1: The YOLO Detection System.** Processing images with YOLO is simple and straightforward. Our system (1) resizes the input image to $448 \times 448$, (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model's confidence.