
CS 232:
Artificial Intelligence

Spring 2024

Prof. Carolyn Anderson
Wellesley College

Reminders

- Reading for Friday: Chiang (2023)
- My help hours Friday: 3:30-4:30
- Lepei's help hours Thursday: 3:45-5:45
- Bonus late day option:
 - Attend Psych 216: Psychology of Language from 11:20-12:35 this Friday for a discussion of large language models

Upcoming Talks

BABSON COLLEGE

Renowned AI Ethics Pioneer is Coming to Babson!

6:00 PM | April **2** 2024 | Winn Auditorium



Dr. Rumman Chowdhury

Named by Forbes as one of the five key people shaping AI - Dr. Chowdhury is the former Director of Machine Learning, Ethics, Transparency, and Accountability team at Twitter and now CEO and co-founder of Humane Intelligence.

Join us for an enlightening session that explores the intersection of AI, ethics, policy, and entrepreneurship.



Butler Institute for Free Enterprise Through Entrepreneurship

Upcoming Talks

WELLESLEY CS CLUB — ALUMNAE EVENT

CAREERS IN → TECH

Explore different roles within the industry from a panel of Wellesley alums. Learn what it means be a designer, engineer, data scientist, project manager, language scientist, founder, and how to get started!

Wednesday, April 3
6:00 pm in Sci L031
RSVP bit.ly/W-panel-24

? nt101, ec116 | Accomodations accessibility@wellesley.edu



Christine Doran

Clockwork Language
Verified email at clockworklanguage.com
Corpus linguistics evaluation dialogue



Catherine Chen

PhD Candidate at Brown University



Dr. Rachel Lomasky

DIRECTOR OF MACHINE LEARNING | MANIFOLD

Dr. Rachel Lomasky is Director of Machine Learning at Manifold, where she helps clients train and productionalize their machine learning algorithms.

Prior to Manifold, she was co-founder and Chief Data Officer of WEVO Conversion, a platform for digital marketers that uses AI to improve websites and search

Upcoming Talks

TBA : around lunch time on Thursday, April 18th

Homework 7: Art Generation Competition

future tense

“When Robot and Crow Saved East St. Louis”

A new short story about a disease surveillance robot whose social programming gets put to the test.

BY ANNALEE NEWITZ DEC 29, 2018 • 5:50 AM



Lisa Larson-Walker

Goal: develop a set of illustrations for this short story using generative AI

Computer Vision

color images are tensors



0	3	2	5	4	7	6	0	8
0	3	2	5	4	7	6	0	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

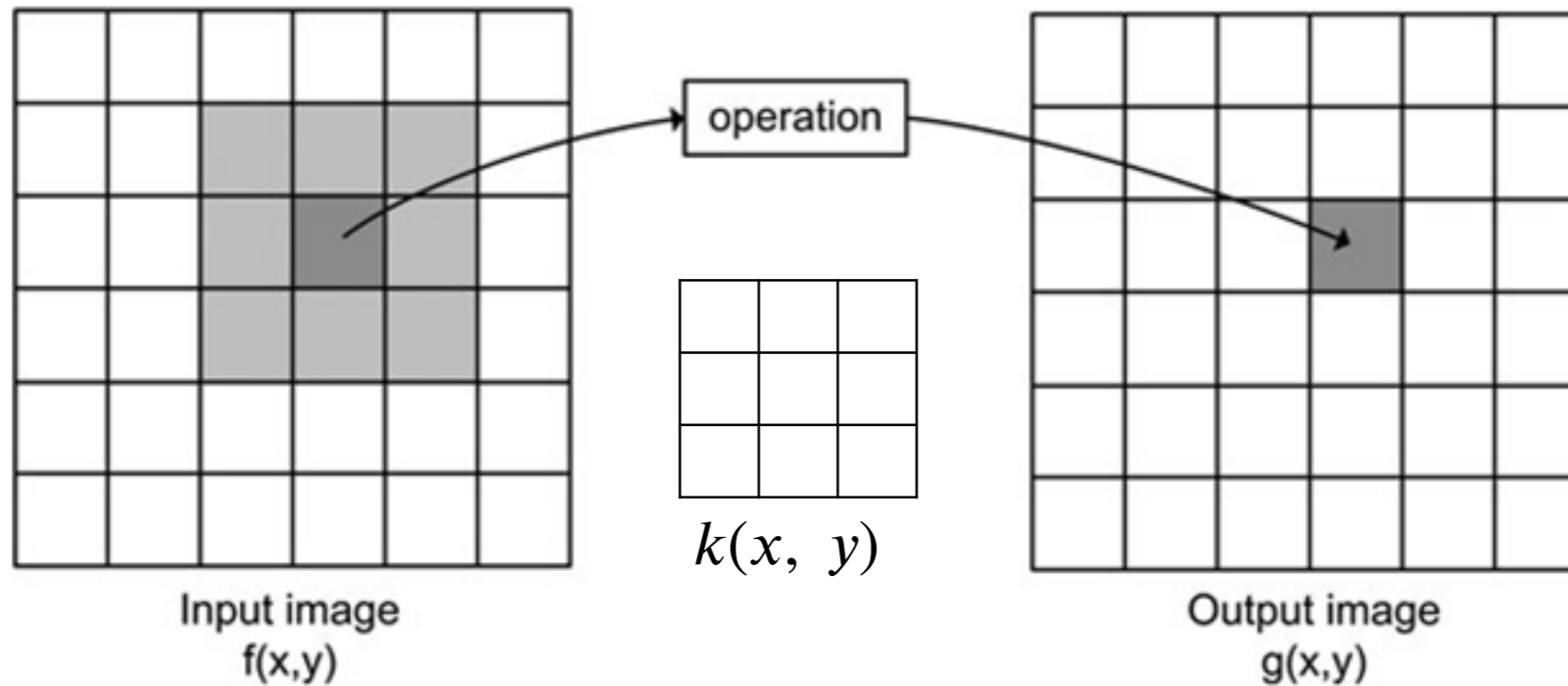
channel x height x width

Channels are usually RGB: Red, Green, and Blue

Other color spaces: HSV, HSL, LUV, XYZ, Lab, CMYK, etc

Convolutional Neural Networks

Convolution Operator



$$g(x, y) = \sum_v \sum_u k(u, v) f(x - u, y - v)$$

Convolution operation

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

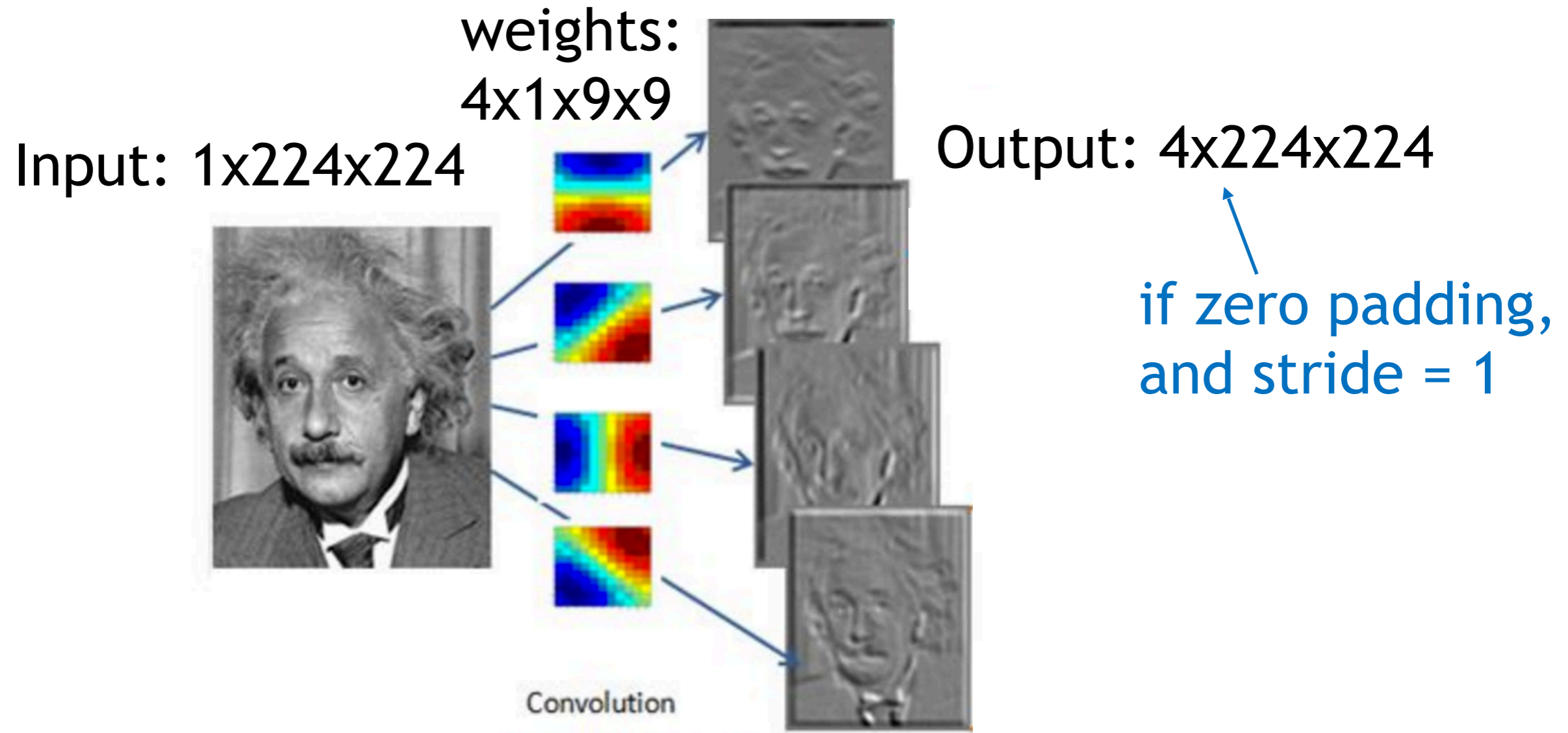
4		

Convolved
Feature

demo:

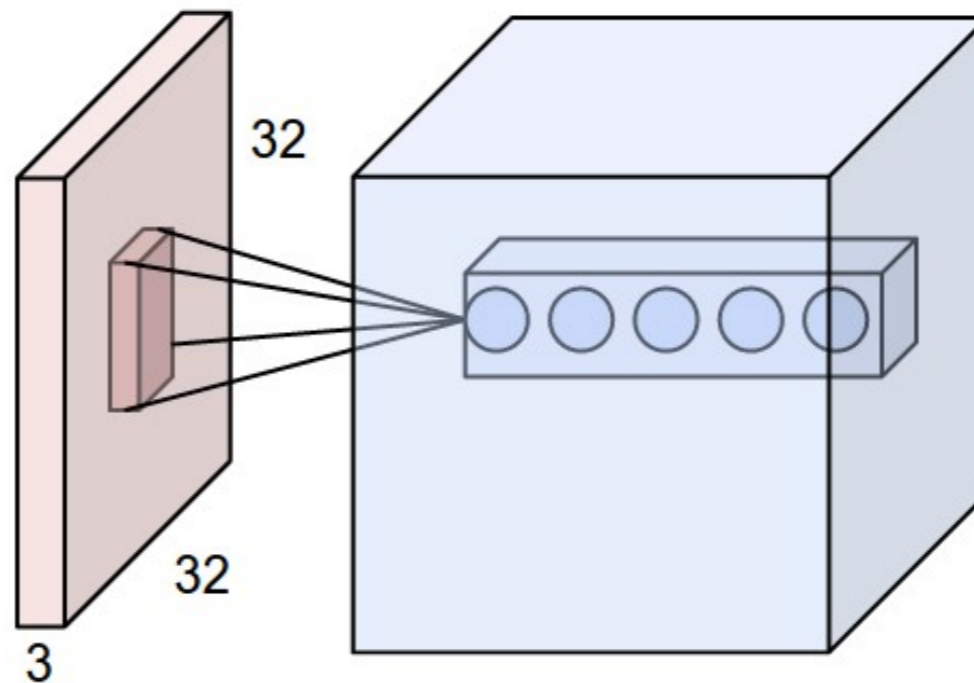
<http://setosa.io/ev/image-kernels/>

Convolutional Layer (with 4 filters)

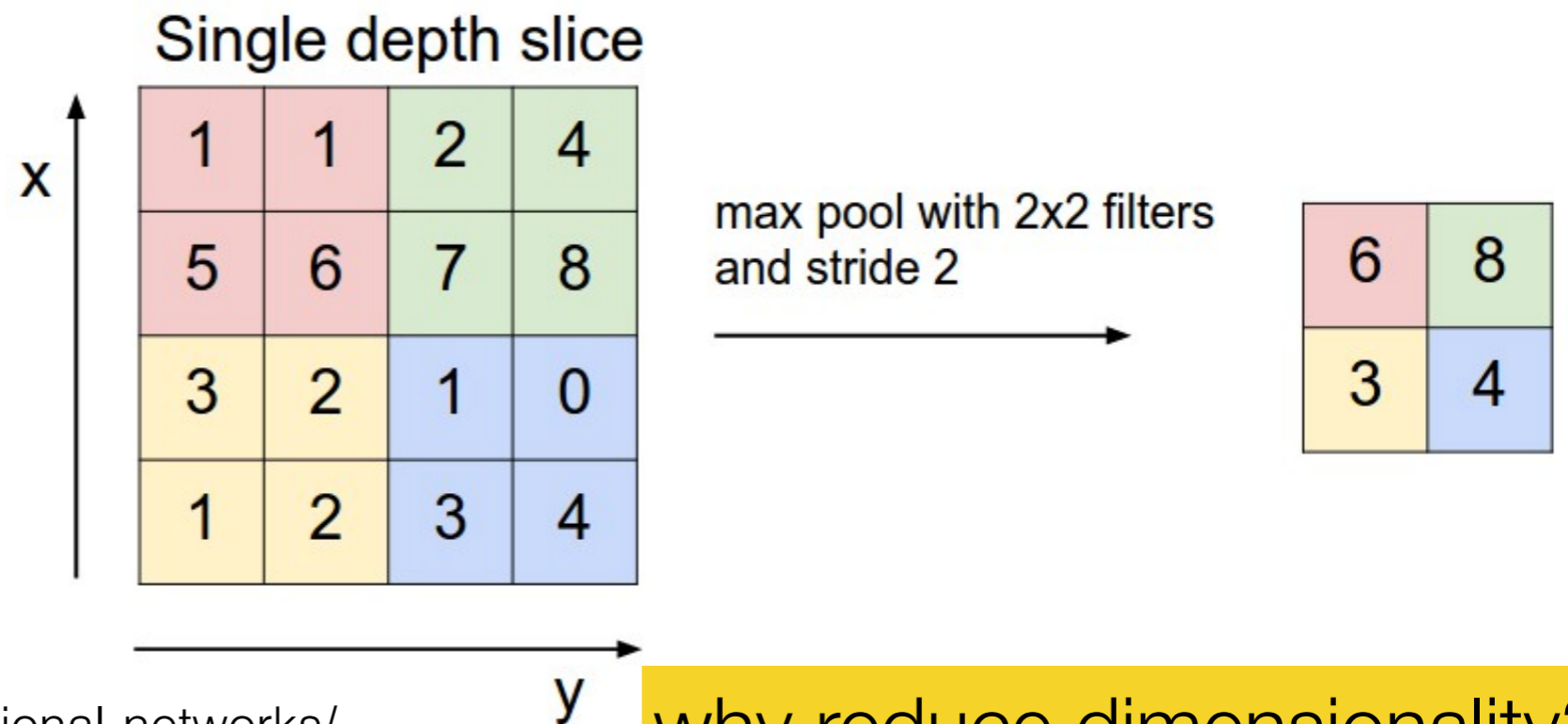


pooling layers also used to reduce dimensionality

Convolutional Layers:
slide a set of small filters over the image



Pooling Layers:
reduce dimensionality of representation



Alexnet

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

the paper that started the
deep learning revolution!

image classification

Classify an image into 1000 possible classes:

e.g. Abyssinian cat, Bulldog, French Terrier, Cormorant,
Chickadee,
red fox, banjo, barbell, hourglass, knot, maze, viaduct, etc.

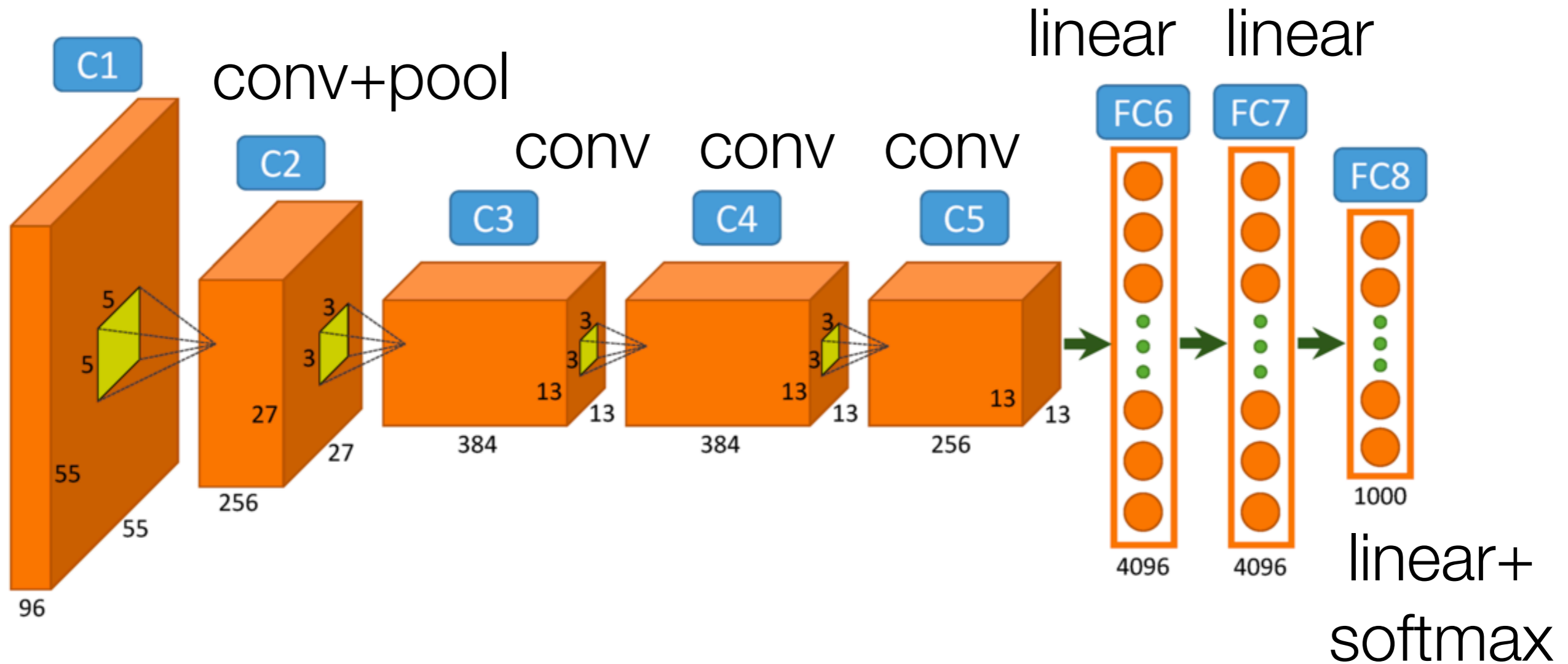


cat, tabby cat (0.71)
Egyptian cat (0.22)
red fox (0.11)
.....

train on the ImageNet
challenge dataset,
~1.2 million images

Alexnet

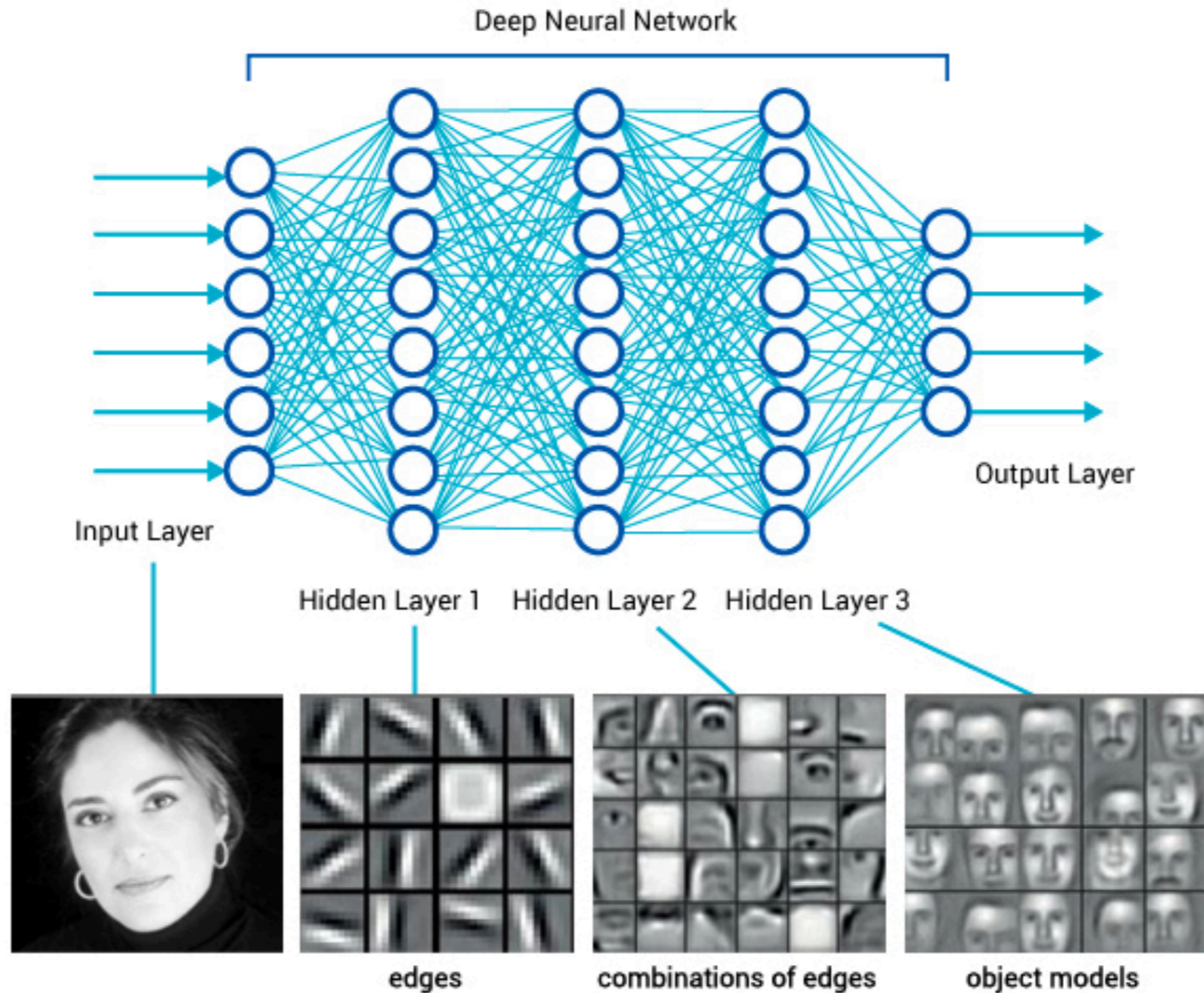
conv+pool



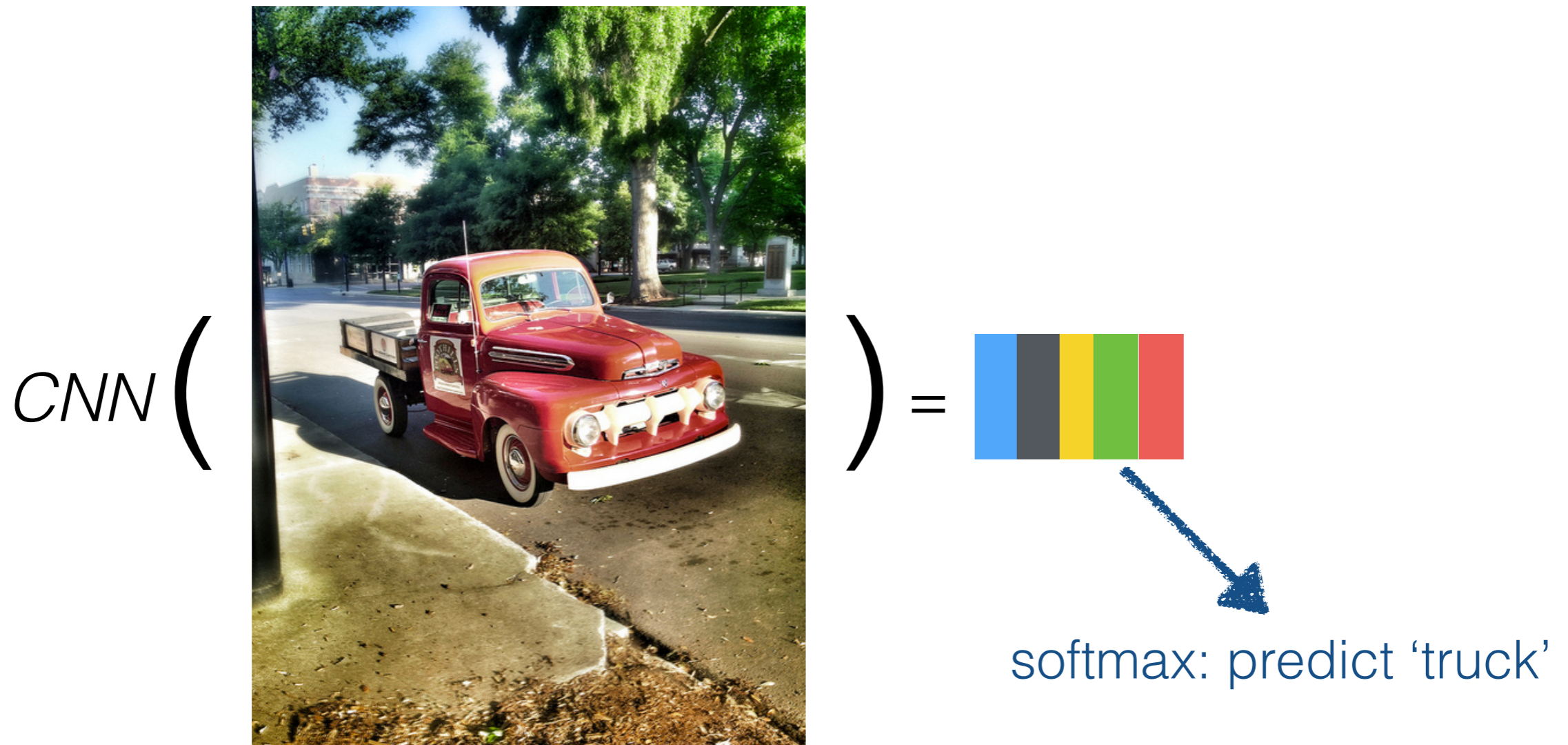
<https://www.saagie.com/fr/blog/object-detection-part1>

Slides adapted from Mohit Iyyer

What is happening?



at the end of the day, we generate a fixed size vector from an image and run a classifier over it

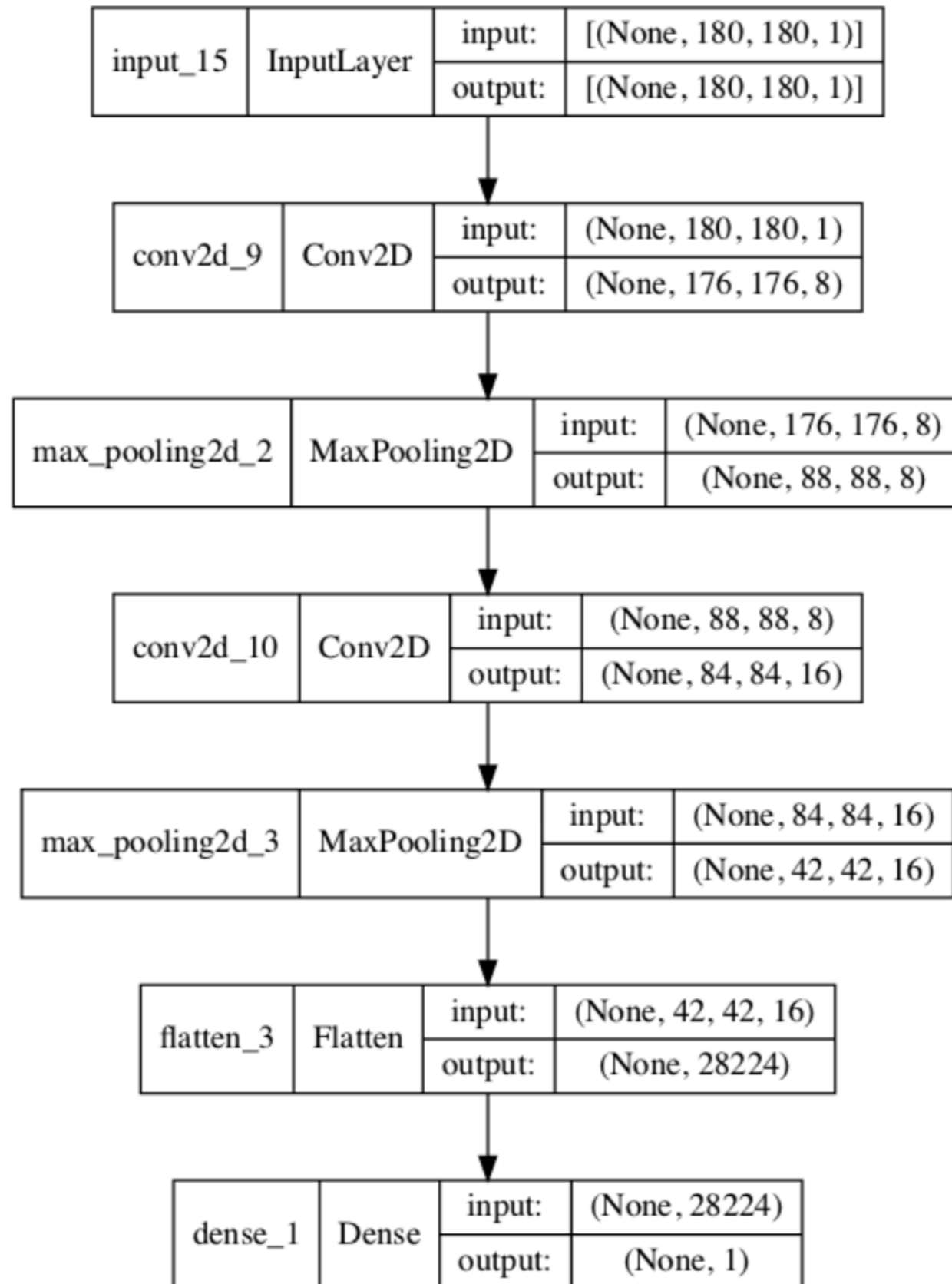


Adding More Layers

```
def make_model(input_shape, num_classes):
    inputs = keras.Input(shape=input_shape)
    x = layers.Conv2D(8, (5, 5), activation='relu', strides=1)(inputs)
    x = layers.MaxPooling2D((2, 2))(x)
    x = layers.Conv2D(16, (5, 5), activation='relu', strides=1)(x)
    x = layers.MaxPooling2D((2, 2))(x)
    x = layers.Flatten()(x)
    if num_classes == 2:
        activation = "sigmoid"
        units = 1
    else:
        activation = "softmax"
        units = num_classes
    outputs = layers.Dense(units, activation=activation)(x)
    return keras.Model(inputs, outputs)
```

```
model = make_model(input_shape=image_size+(1,), num_classes=2)
keras.utils.plot_model(model, show_shapes=True)
```

New Architecture



Auto-Encoders

Auto-encoders

Auto-encoders are a class of neural networks that do not require labeled data.

Supervised NNs: predict the **output** given the **input**.

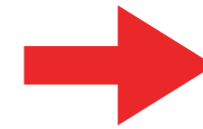
Auto-encoders: predict the **input** given the **input**.

Key idea: select features by **reducing then increasing** dimensionality.

Normal NN goes:



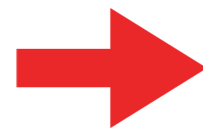
shutterstock.com · 451203238



Auto-encoder goes:

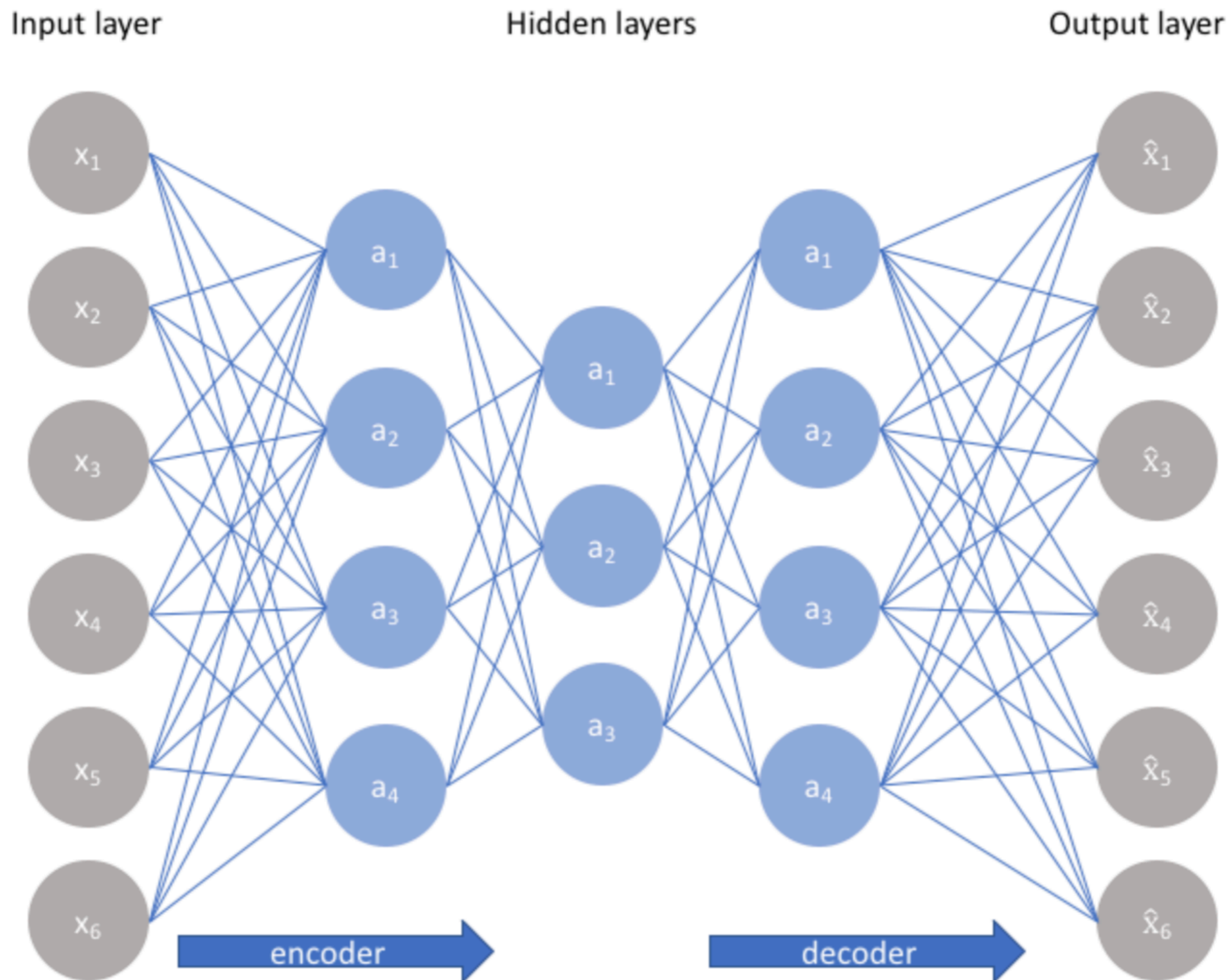


shutterstock.com · 451203238



shutterstock.com · 451203238

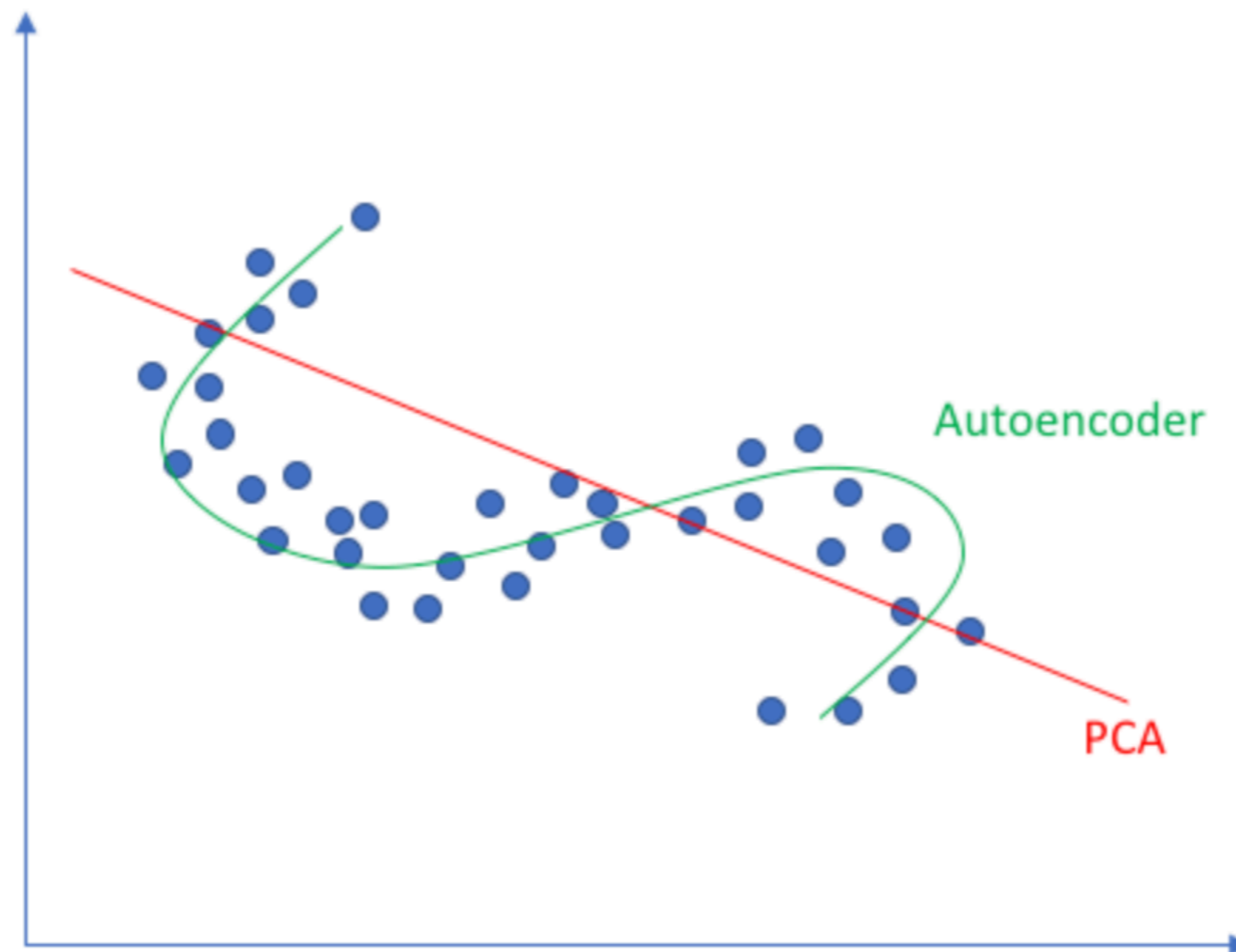
Auto-Encoder Architecture



Auto-Encoders as Dimensionality Reduction

Auto-encoders are a more powerful form of dimensionality reduction than traditional techniques like PCA, because they can learn nonlinear transformations.

Linear vs nonlinear dimensionality reduction



Encoder

Model: "sequential_6"

Layer (type)	Output Shape	Param #
conv2d_9 (Conv2D)	(None, 32, 32, 32)	320
dropout_18 (Dropout)	(None, 32, 32, 32)	0
conv2d_10 (Conv2D)	(None, 16, 16, 64)	18496
dropout_19 (Dropout)	(None, 16, 16, 64)	0
conv2d_11 (Conv2D)	(None, 8, 8, 128)	73856
dropout_20 (Dropout)	(None, 8, 8, 128)	0
flatten_3 (Flatten)	(None, 8192)	0
dense_6 (Dense)	(None, 128)	1048704

=====
Total params: 1,141,376
Trainable params: 1,141,376
Non-trainable params: 0
=====

Decoder

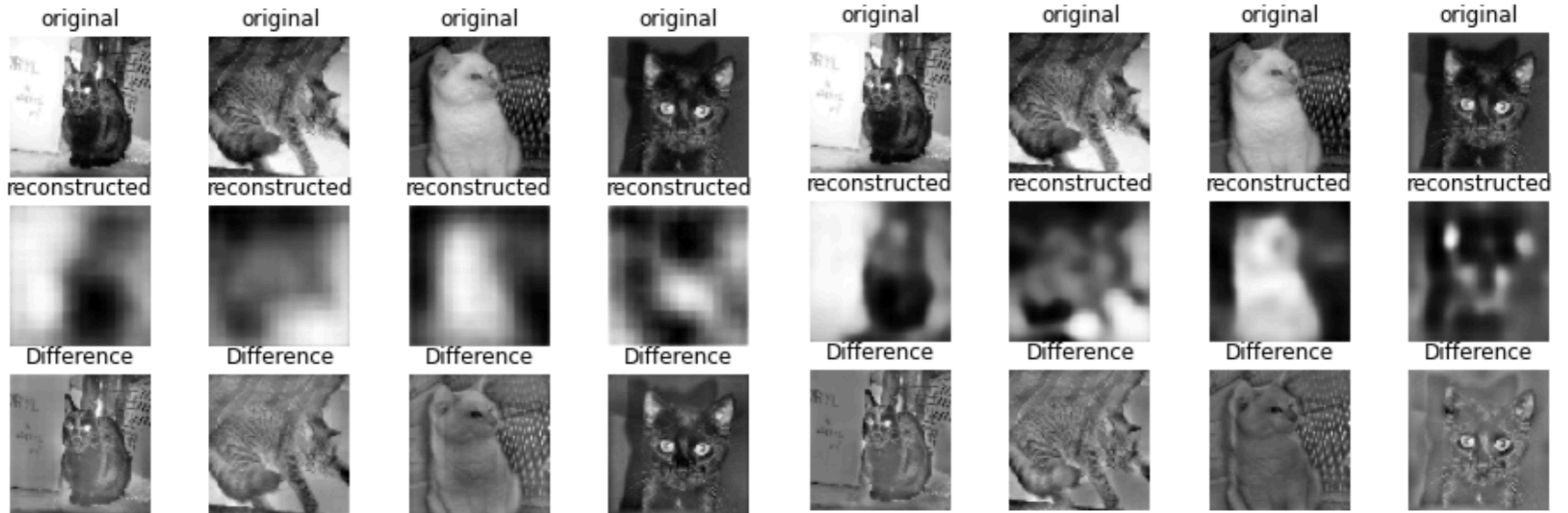
Layer (type)	Output Shape	Param #
dense_7 (Dense)	(None, 8192)	1056768
reshape_3 (Reshape)	(None, 8, 8, 128)	0
conv2d_transpose_12 (Conv2D Transpose)	(None, 16, 16, 128)	147584
dropout_21 (Dropout)	(None, 16, 16, 128)	0
conv2d_transpose_13 (Conv2D Transpose)	(None, 32, 32, 64)	73792
dropout_22 (Dropout)	(None, 32, 32, 64)	0
conv2d_transpose_14 (Conv2D Transpose)	(None, 64, 64, 32)	18464
dropout_23 (Dropout)	(None, 64, 64, 32)	0
conv2d_transpose_15 (Conv2D Transpose)	(None, 64, 64, 1)	289

Total params: 1,296,897
Trainable params: 1,296,897
Non-trainable params: 0

Input, Output, Difference

Epoch 1

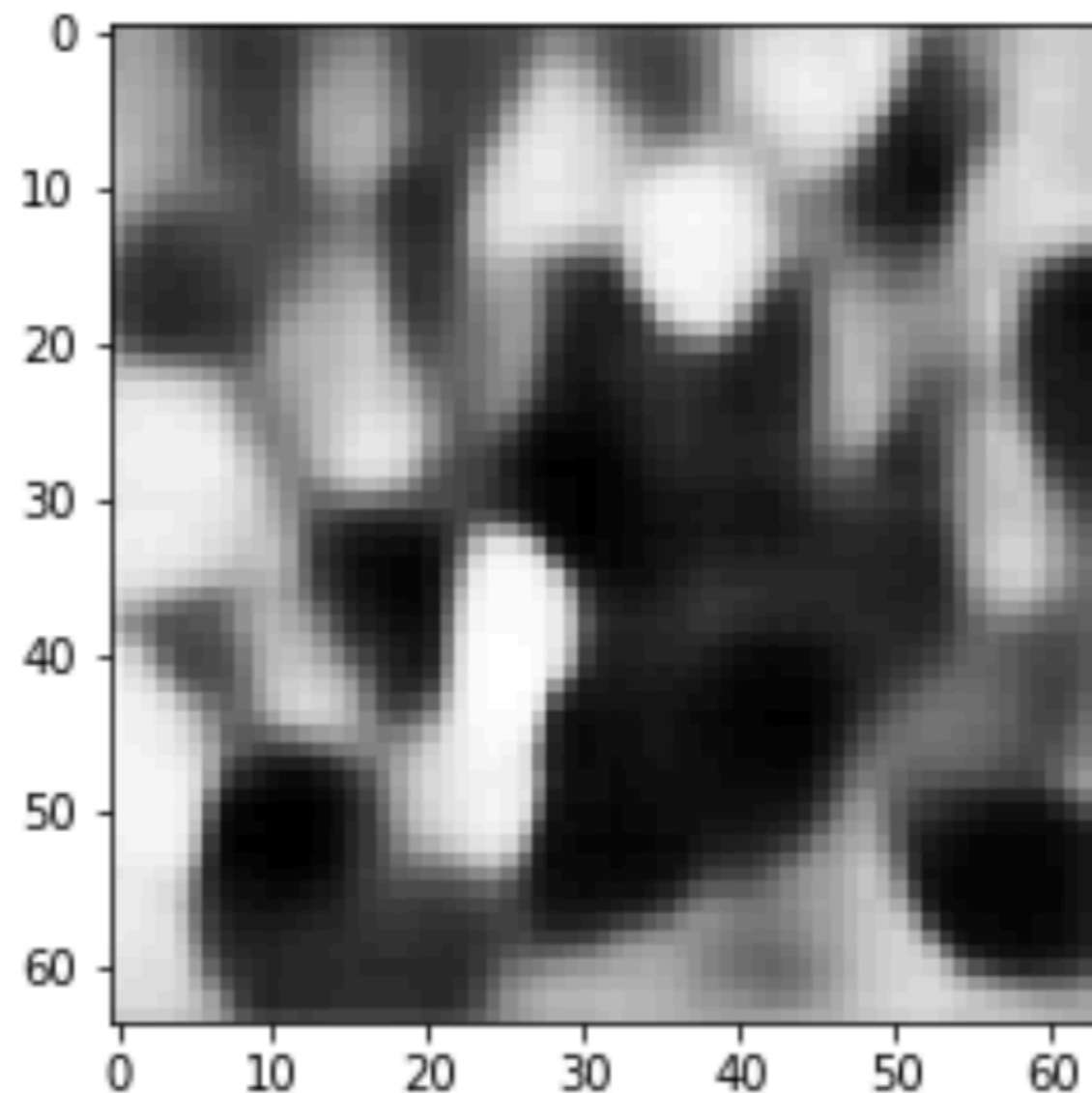
Epoch 10



Using Decoder to Generate

Input noise to the decoder to make it hallucinate a cat:

```
x = autoencoder.decoder(np.random.randn(1, 128)).numpy()  
plt.imshow(x[0, :, :, 0], cmap='gray')
```



Stable Diffusion

Stable Diffusion

3 components:

1. VAE: an auto-encoder to map images to a latent space
2. U-Net: an architecture that learns to denoise images
3. CLIP: a text-encoder to allow multi-modal input

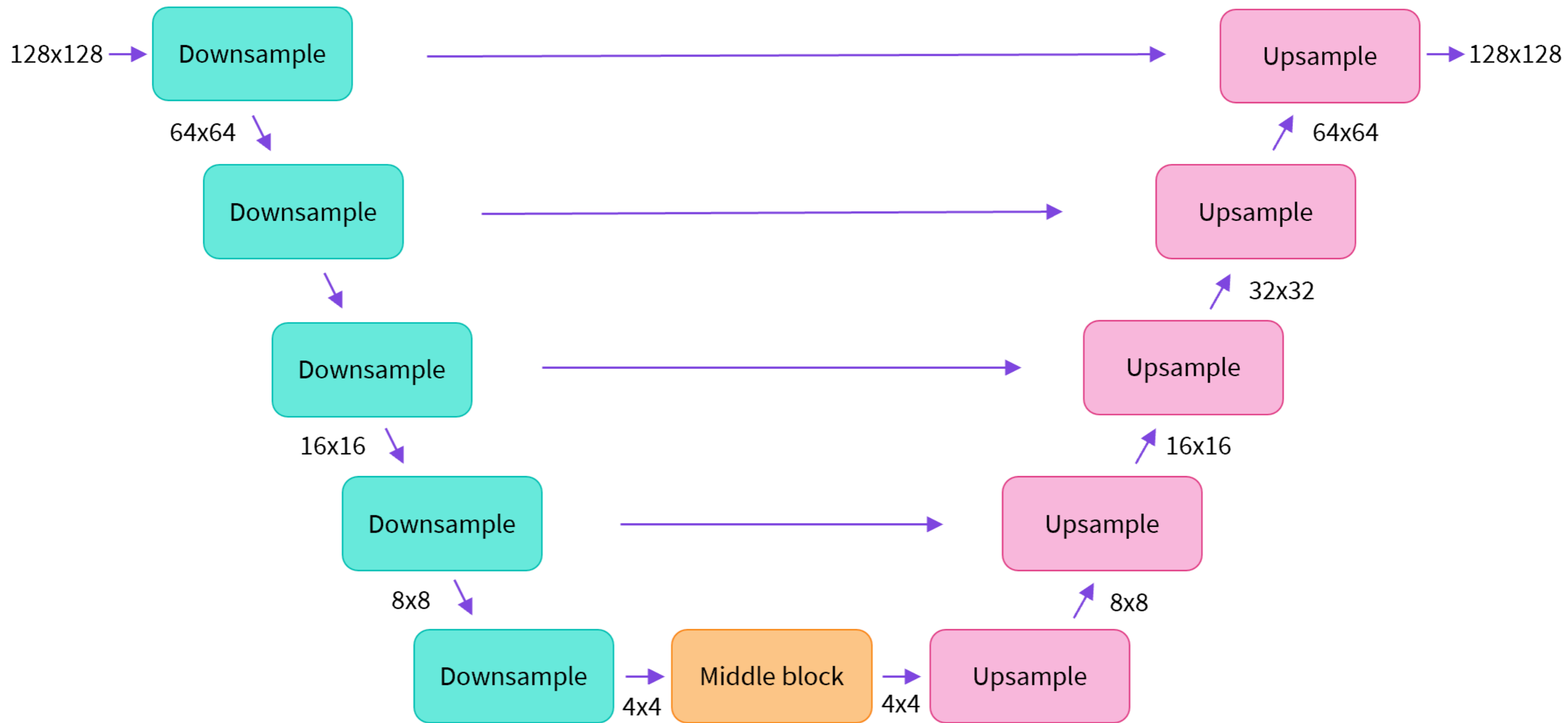
VAE: variational autoencoder

VAE is an encoder / decoder model.

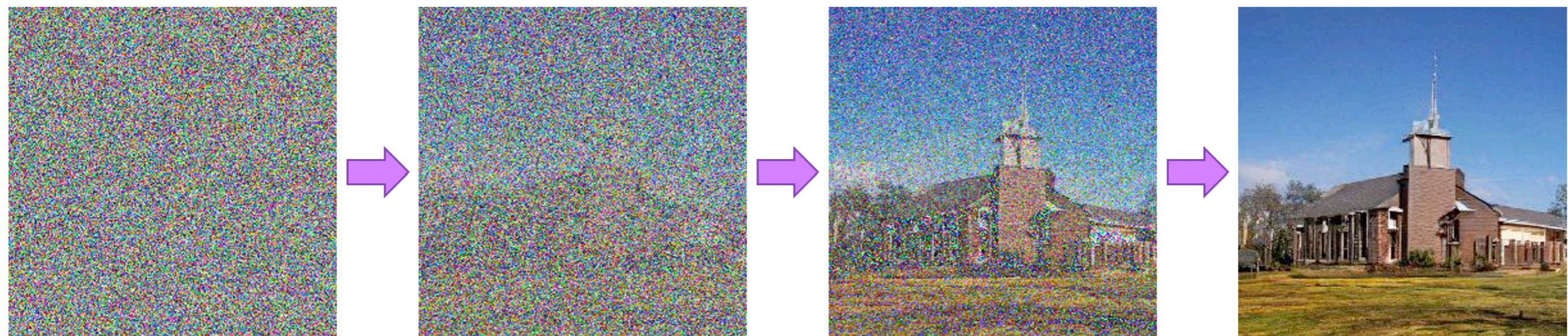
The encoder maps an input image (pixels) to a lower-dimension latent space.

The decoder takes the output of the model and maps it back to an image in pixels.

U-Net model (auto-encoder)



Iteratively Denoising



CLIP: a text encoder for multi-modal input

Objective: given a batch of text and image inputs, predict the correct image-text pairings.

Learning Transferable Visual Models From Natural Language Supervision

Alec Radford^{*1} Jong Wook Kim^{*1} Chris Hallacy¹ Aditya Ramesh¹ Gabriel Goh¹ Sandhini Agarwal¹
Girish Sastry¹ Amanda Askell¹ Pamela Mishkin¹ Jack Clark¹ Gretchen Krueger¹ Ilya Sutskever¹

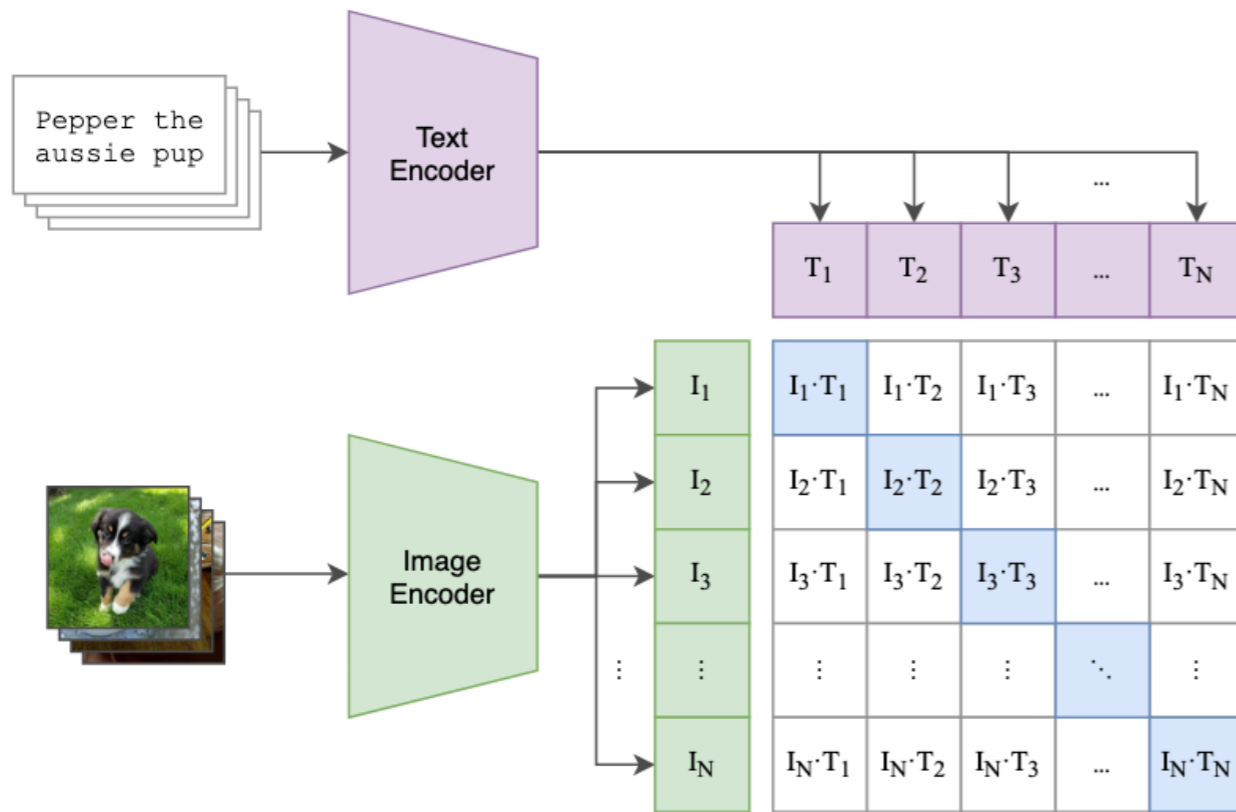
Abstract

State-of-the-art computer vision systems are trained to predict a fixed set of predetermined object categories. This restricted form of supervision limits their generality and usability since additional labeled data is needed to specify any other visual concept. Learning directly from raw text about images is a promising alternative which leverages a much broader source of supervision. We demonstrate that the simple pre-training task

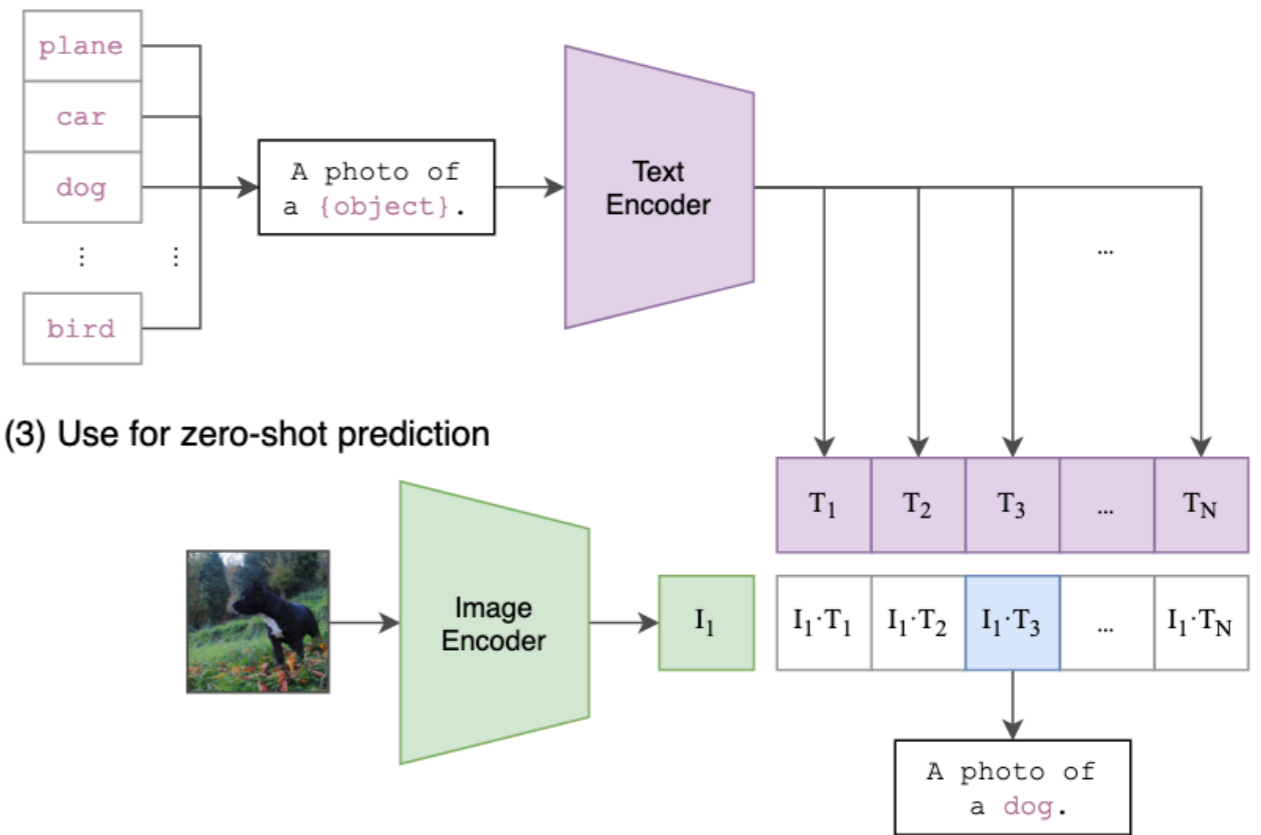
Task-agnostic objectives such as autoregressive and masked language modeling have scaled across many orders of magnitude in compute, model capacity, and data, steadily improving capabilities. The development of “text-to-text” as a standardized input-output interface (McCann et al., 2018; Radford et al., 2019; Raffel et al., 2019) has enabled task-agnostic architectures to zero-shot transfer to downstream datasets removing the need for specialized output heads or dataset specific customization. Flagship systems like GPT-3 (Brown et al., 2020) are now competitive across many tasks with bespoke models while requiring little to no dataset

CLIP: a text encoder for multi-modal input

(1) Contrastive pre-training



(2) Create dataset classifier from label text



Stable Diffusion: putting the pieces together

