



Consiglio Nazionale delle Ricerche



Istituto di Scienza e Tecnologie  
dell'Informazione "A. Faedo"



# DEEP LEARNING TOOLS FOR IMAGE CLASSIFICATION AND RETRIEVAL

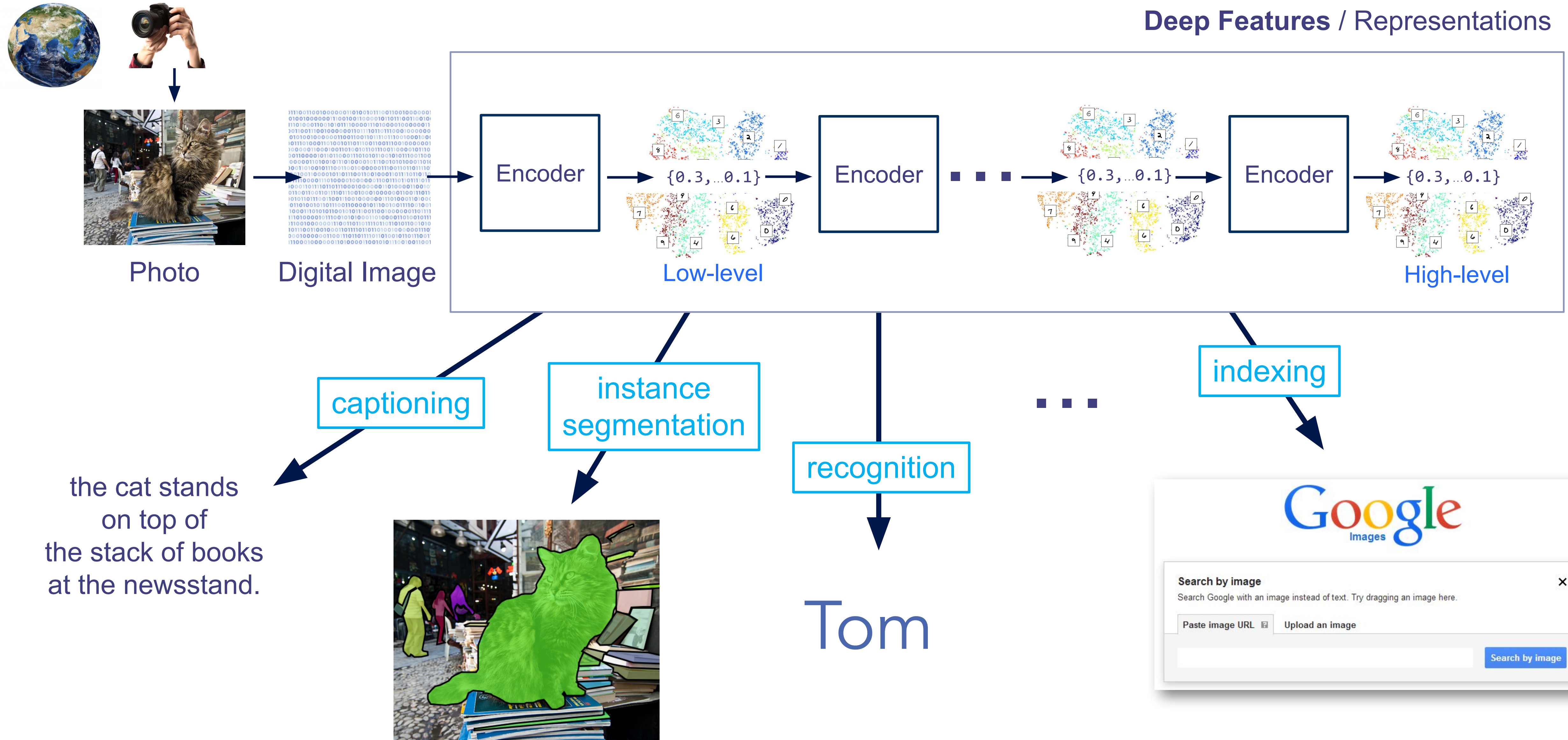
---

Fabio Carrara, Nicola Messina

Researcher @ ISTI-CNR  
[fabio.carrara@isti.cnr.it](mailto:fabio.carrara@isti.cnr.it)

PostDoc @ ISTI-CNR  
[nicola.messina@isti.cnr.it](mailto:nicola.messina@isti.cnr.it)

# The Overall Picture



# Outline

---

- Deep Learning for Images
  - Introduction: Image Features/Representation
  - Intuition: Why we use Deep Features
  - From Light to Bits
  - Storing and Sharing Images
- Image Classification
- Image Retrieval

# Introduction: Image Features / Representations

# An Image, Many Tasks!



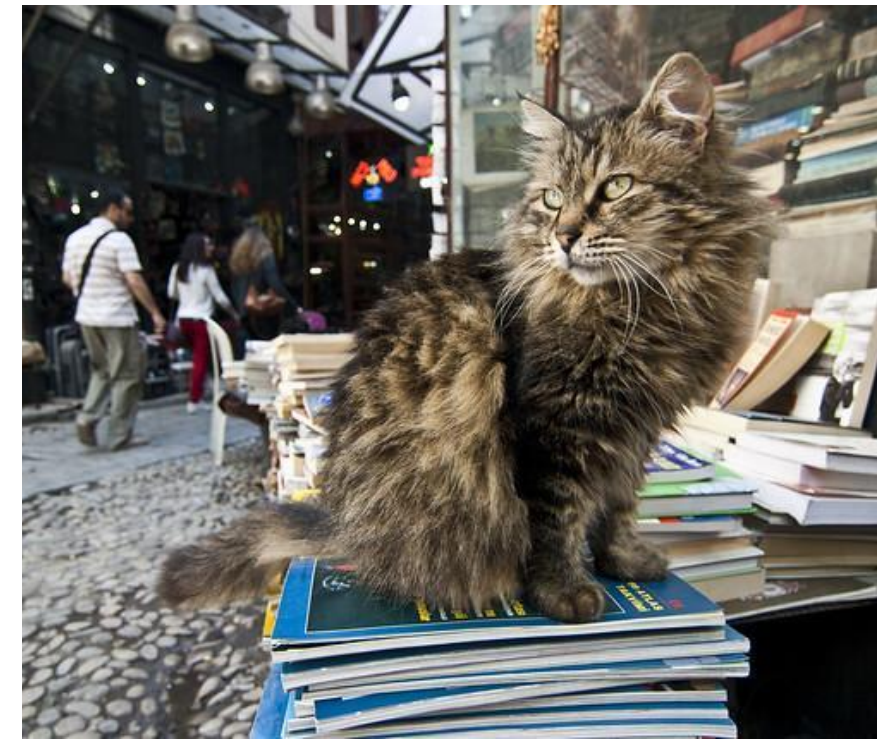
What we want to do:

- captioning
- similarity
- segmentation
- detection
- recognition
- classification
- retrieval
- ...

# From Real World to Bits



photography



digitalization



(compression)



Camera

Photo

Bits

Analysis are performed on a digital image

A digital image (usually) is a digital photo of the real world.  
(not the real world, not an analog photo)

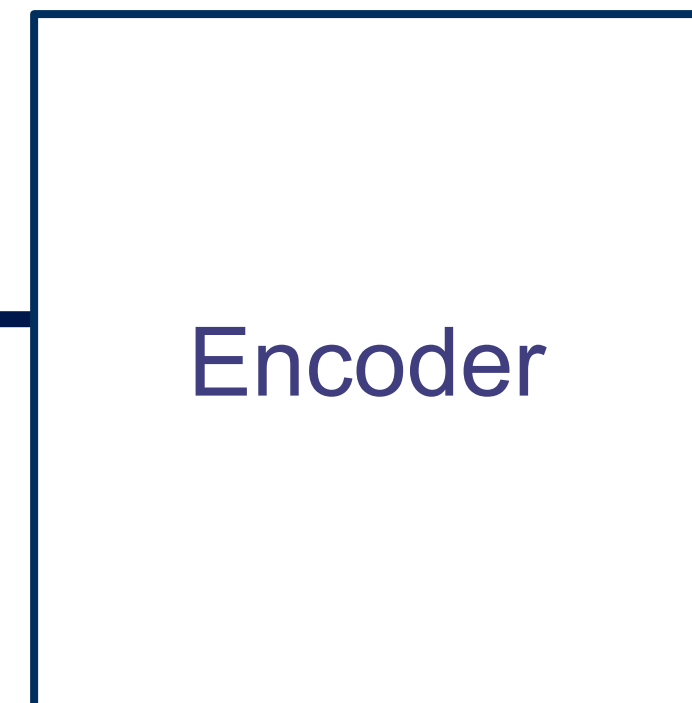
The bits representing the digital images have **very low semantic**.

# Representation / Features



000001  
1100101  
000011  
100000  
001001  
000001  
010110  
001100  
001101  
101110  
1111011  
110010  
1111011  
010110110011001100100000111010001101000  
011010010110110011000010111001001111001001  
1000111010101100101011100110010000001101111  
1110100001011100101010001101000011010010111  
111010000001100110110110110110110110110101010

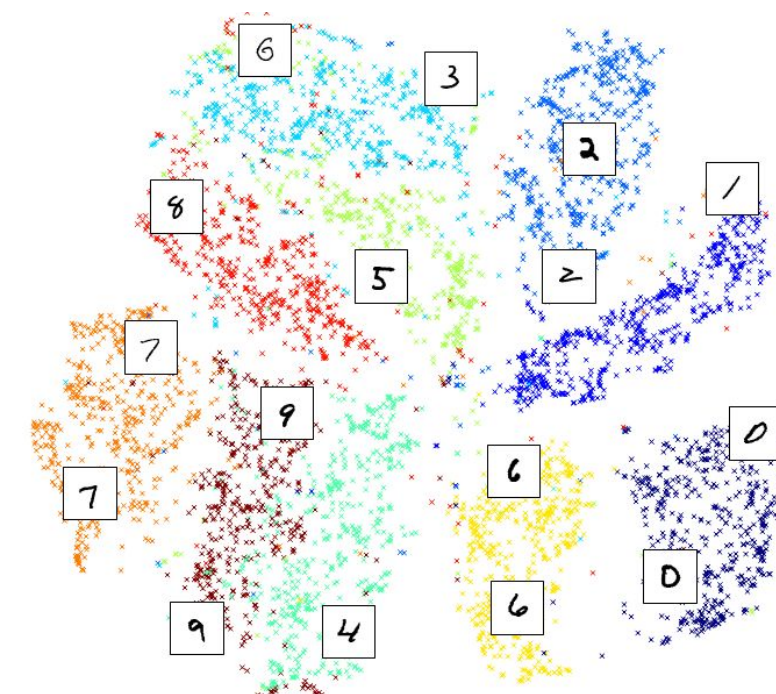
Digital Image



feature  
extractor

A representation  
(a feature)

{0.3, 0.5 ... 0.1}



in a **latent space**

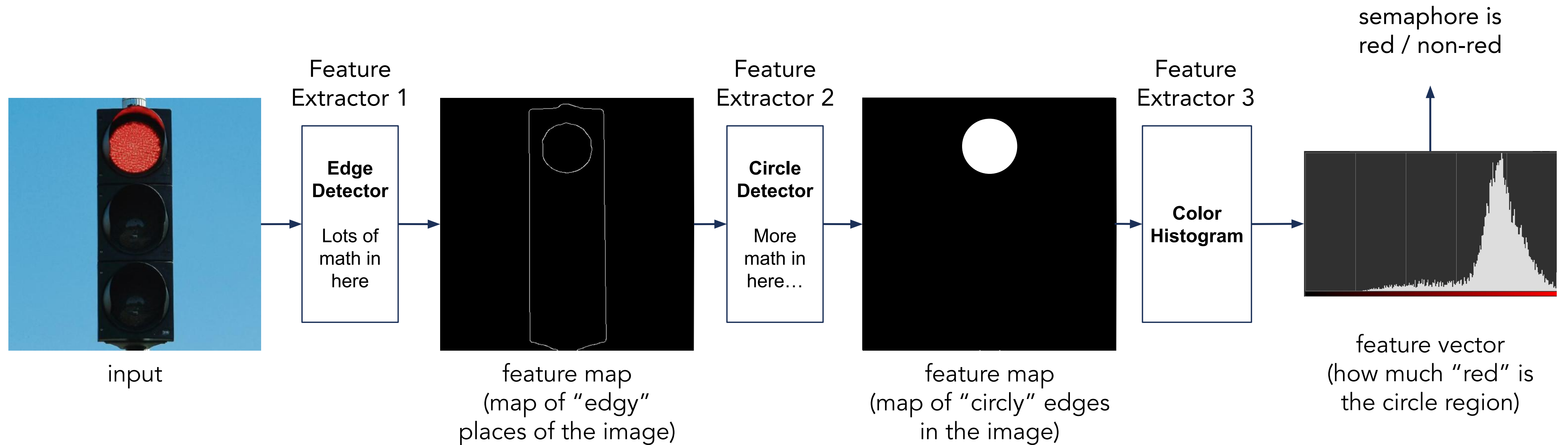
- captioning
- similarity
- segmentation
- detection
- recognition
- classification
- retrieval
- ...

An encoder (or feature extractor) takes an input (e.g., an image) and produce a representation (or descriptor) that is used (in place of the image) for the specific task.

# Handcrafted Features

Task: find out automatically if semaphore is red

Hierarchy of features you may manually define to solve the task:





# Features / Representation

Before Deep Learning, handcrafted features/representations:

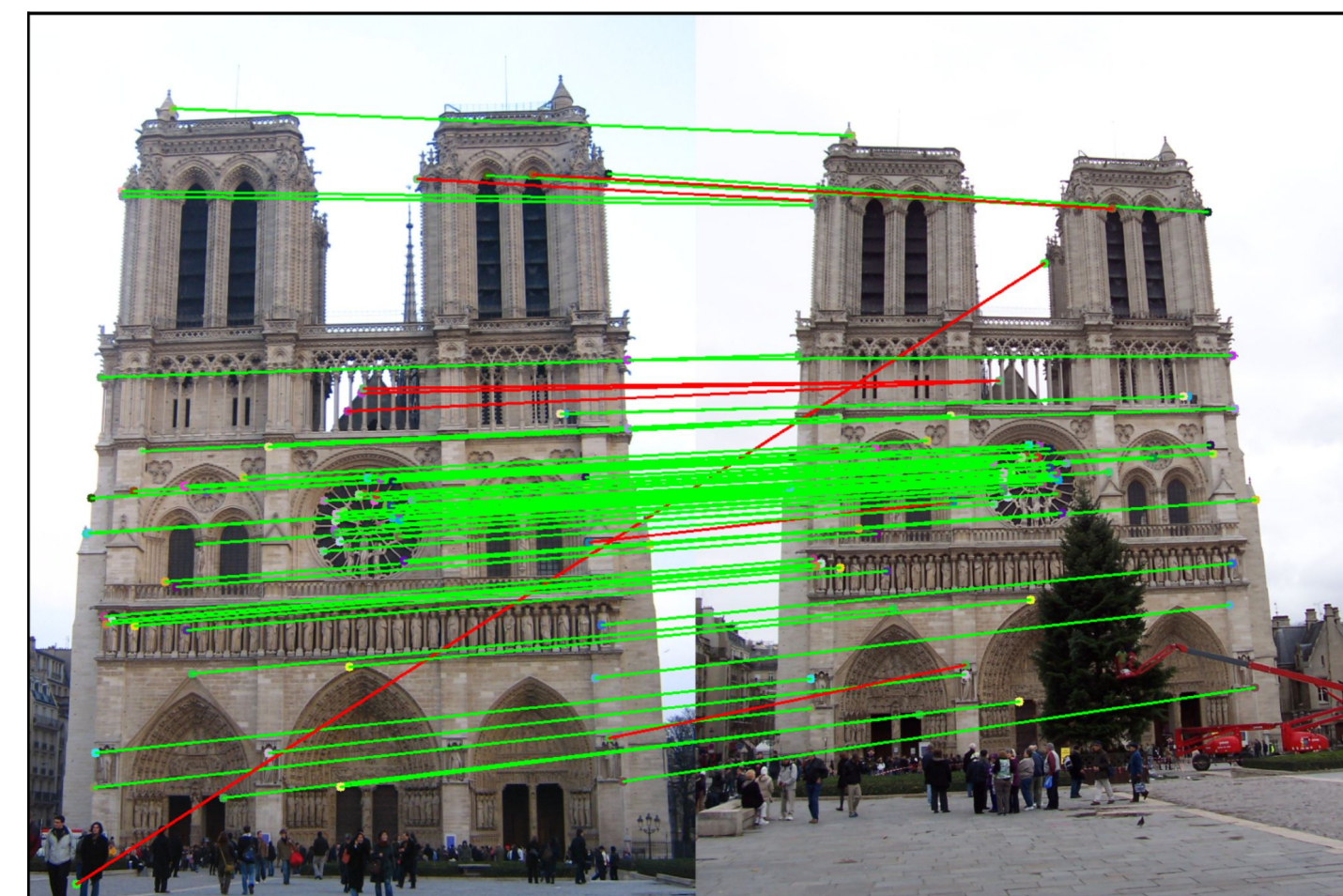
Global Features:

- color, edge, texture etc...



Local Features:

- representation of interest points/regions
- for image stitching or object recognition



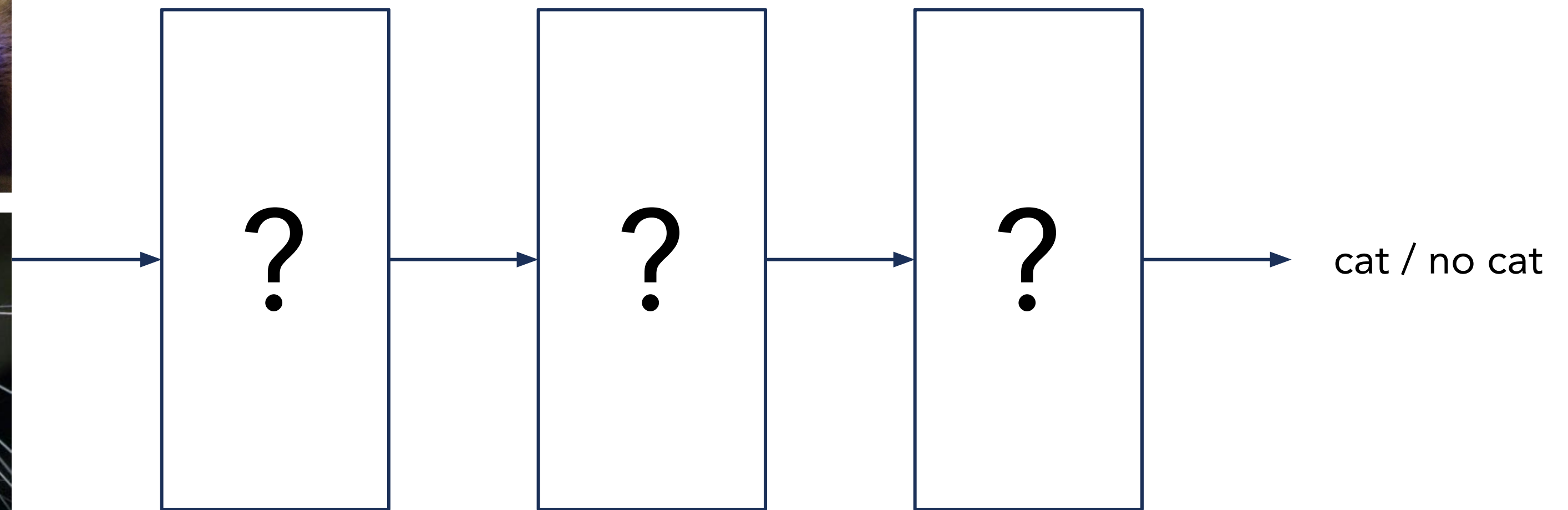
# Handcrafted Features

Task: find out automatically if image contains a cat

Hierarchy of features: ???

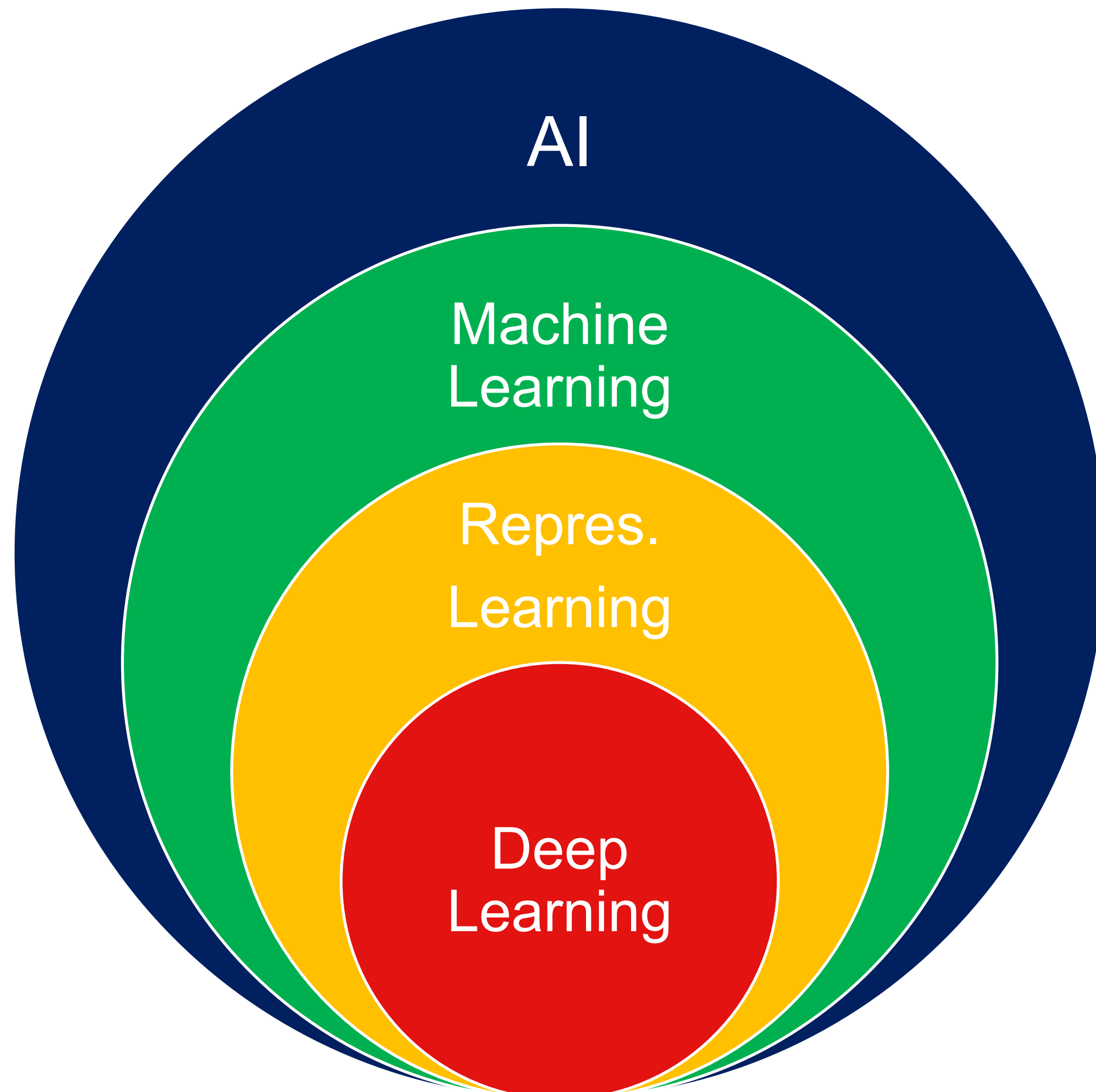


Define good and robust features for this task manually is way too hard.



You **LEARNT** to recognize a cat by examples and experience. You are an expert cat-recognizer. What do you think are the elements that let you recognize a cat?

# Deep Learning (from Nature)



## Representation learning methods that

“allow a machine to be fed with raw data and to automatically discover the representations needed for detection or classification.”

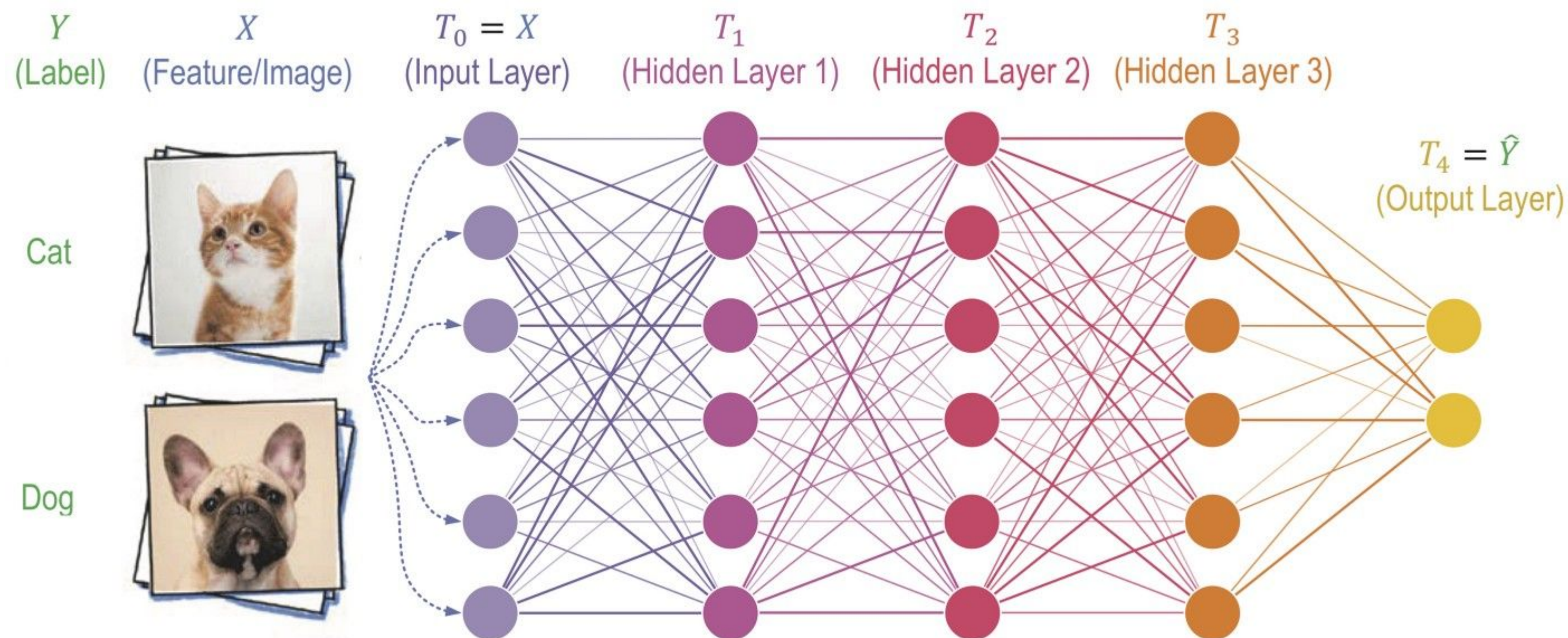
## Deep-learning are representation learning methods

- with multiple levels of representation, obtained by
- composing simple but non-linear modules that each
- transform the representation at one level into a representation at a higher, slightly more abstract level.

# Deep Learning & Artificial Neural Networks

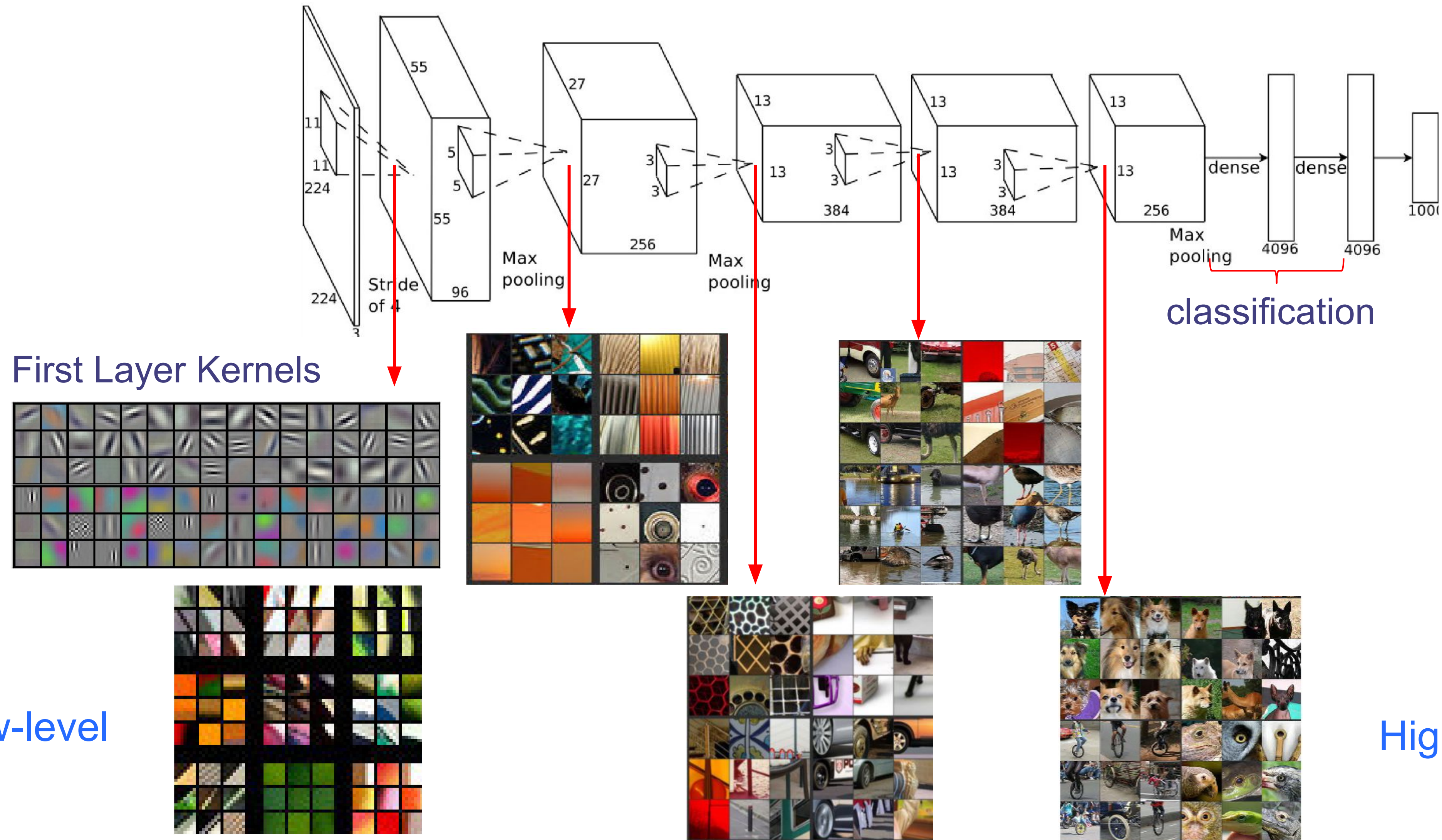
Deep Learning Models are often Artificial Neural Networks:

- Loosely-inspired by biological learning in mammal brains
- An artificial neuron can learn to recognize a pattern
- Several neurons are organized in layers
- Each layer can be thought as a learnable feature extractor
- Several layers, one feeding on the output of the previous one, form an artificial neural network
- Given inputs and a desired outputs, in the training phase neurons adapt to align the network output to the desired one



# Multiple Levels Of Abstraction

AlexNet, 2012, Trained on a Classification task of 1,000 classes.

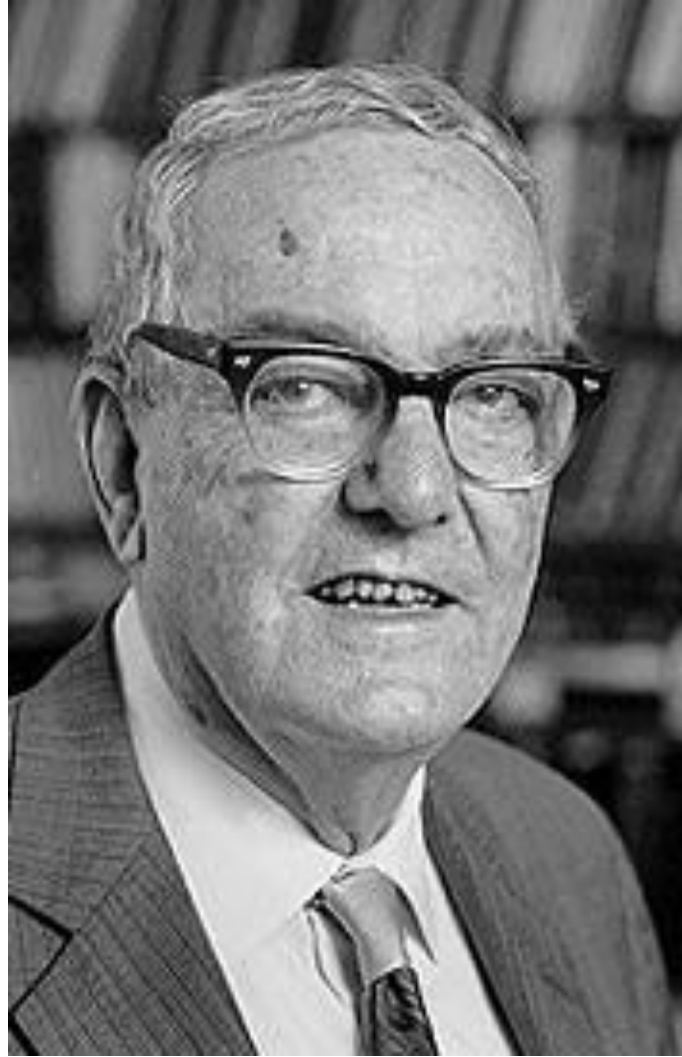




# Intuition: Why we use Deep Learning Features?

# Intuition is Recognition

---



“Intuition is nothing more and nothing less than recognition”  
Herbert Simon, Turing Award 1975 and the Nobel 1978

Simon defined intuition as the recognition of patterns stored in memory.



“There is really no difference between the physician recognising a particular disease from a facial expression and a little child learning, pointing to something and saying doggie. The little child has no idea what the clues are but he just said, he just knows this is a dog without knowing why he knows”.

Daniel Kahneman, Nobel Prize 2002

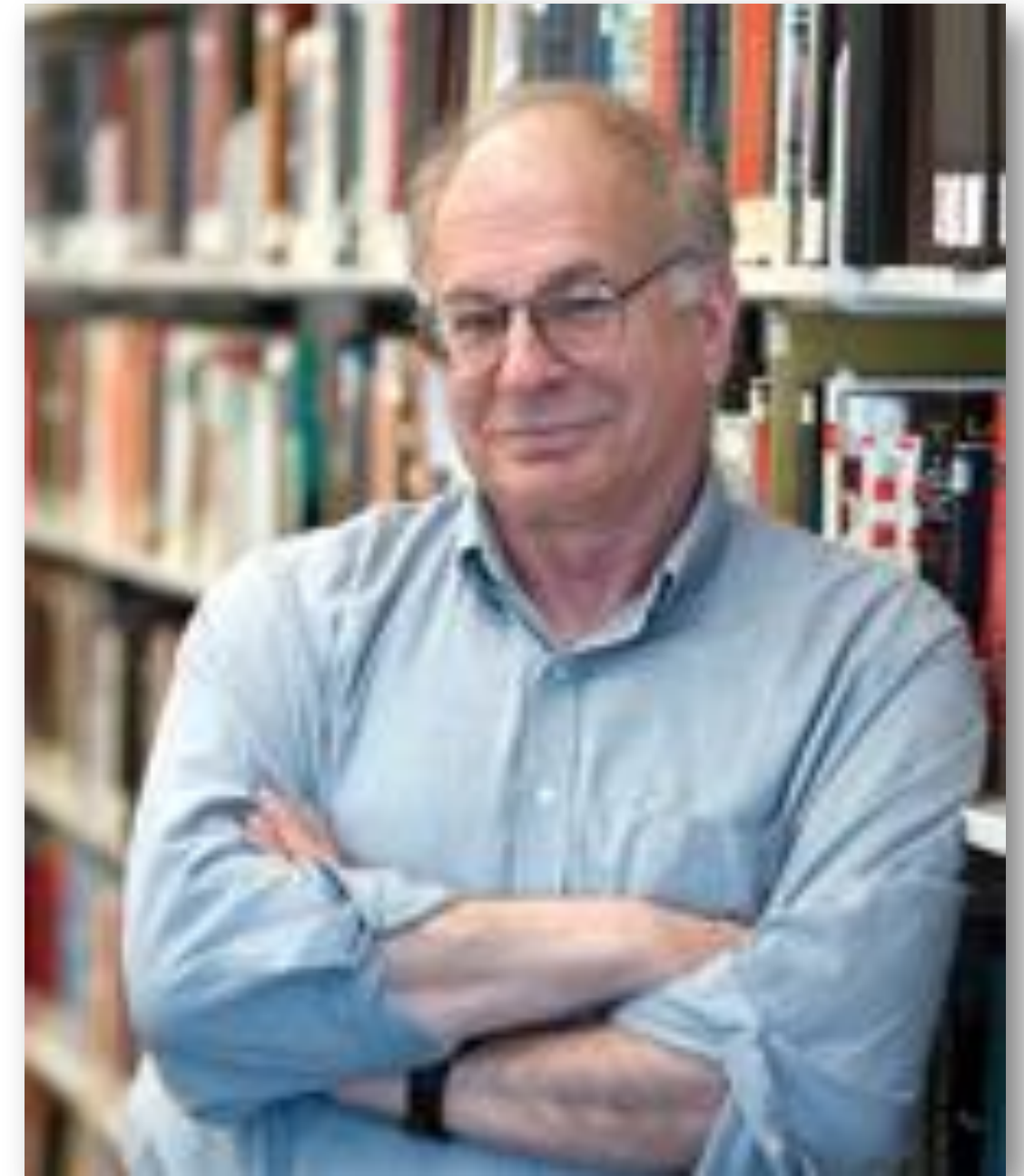


# Thinking, Fast and Slow

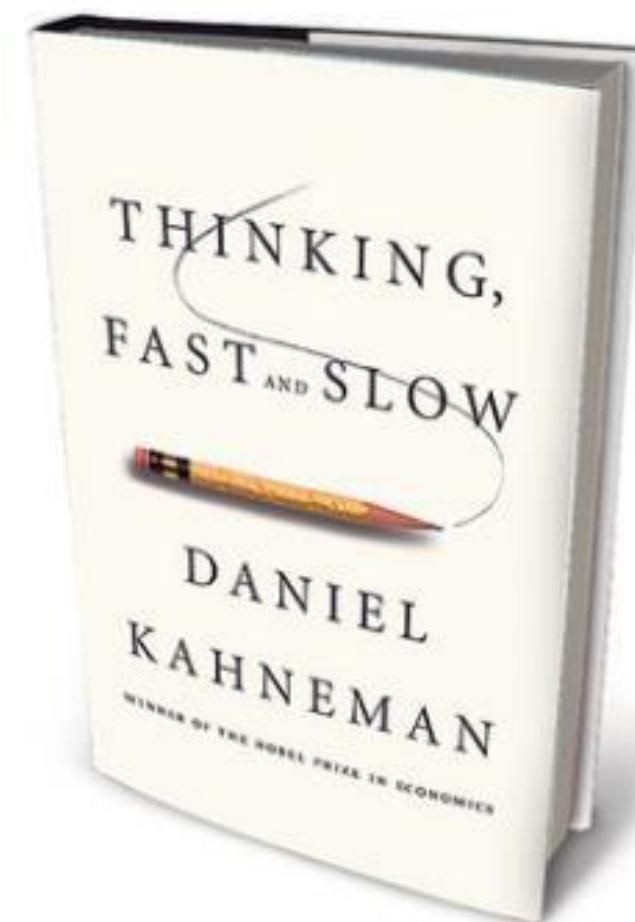
## Daniel Kahneman

- Psychologist
- Nobel Prize in Economic Sciences in 2002 (shared with Vernon L. Smith)

*for having integrated insights from psychological research into economic science, especially concerning human judgment and decision-making under uncertainty*



*Thinking, Fast and Slow (2011)*







What did you see?

What is she going to do?

# Thinking Fast

---



As surely and quickly as you saw that

- the young woman's hair is dark, you knew she is angry.

What you saw extended into the future.

- You sensed that this woman is about to say some very unkind words, probably in a loud and strident voice.

# Thinking Fast

---



You did not intend

- to assess her mood, or
- to anticipate what she might do.

Your reaction to the picture did not have the feel of something you did.

It just happened to you.

It was an instance of fast thinking.

# Thinking, Fast and Slow

---

$$17 \times 24 = ?$$

What came to your mind?





123

586

12.609

# Thinking Slow

---

$$17 \times 24 =$$

- this is a multiplication problem
- you knew that you could solve it
- you would be quick to recognize that both 12,609 and 123 are implausible
- you would not be certain at first that answer is not 568

# Thinking Slow

---



*A precise solution  
did not come to mind*



*You felt you could choose  
whether or not  
to engage in the computation*

# Thinking, Fast and Slow

If you engaged the computation, you proceeded through a sequence of steps.



you retrieved from memory  
the cognitive program  
then you implemented it



Carrying out the computation was a strain.



You felt the burden of:

- holding much material in memory,
- as you needed to keep track of where you were and of where you were going,
- while holding on to the intermediate result.

# Thinking, Fast and Slow

---

Kahneman describes two different ways the brain forms thoughts:

- System 1: Fast, automatic, frequent, emotional, stereotypic, subconscious
- System 2: Slow, effortful, infrequent, logical, calculating, conscious

*Deep Learning  $\approx$  Thinking Fast*

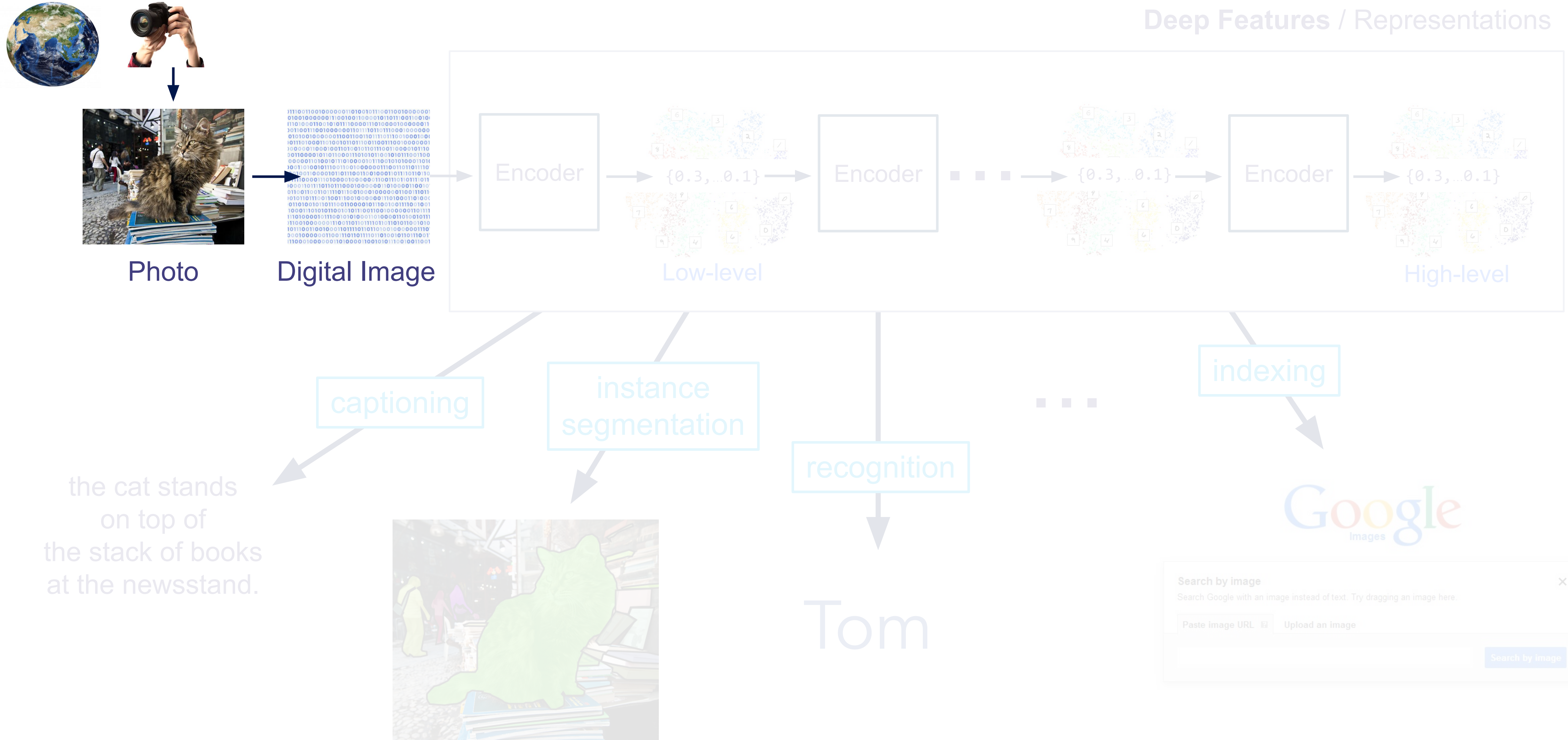
# Thinking, Fast and Slow

---

$$17 \times 24 = 408$$

# From Light to Bits

# From Light to Bits





# From Light to Bits



photography



digitalization



(compression)



Acquisition Device  
(Camera, Scanner)

Photo

Bits

# Exposure Time / Shutter Speed

- . The amount of time the image sensor is exposed to light
- . Side effects with moving objects



# Aperture

- Opening through which light travels
- Affects the depth of field



f/1.8

f/2.8

f/4.0

f/5.6



f/8

f/11

f/16

f/22

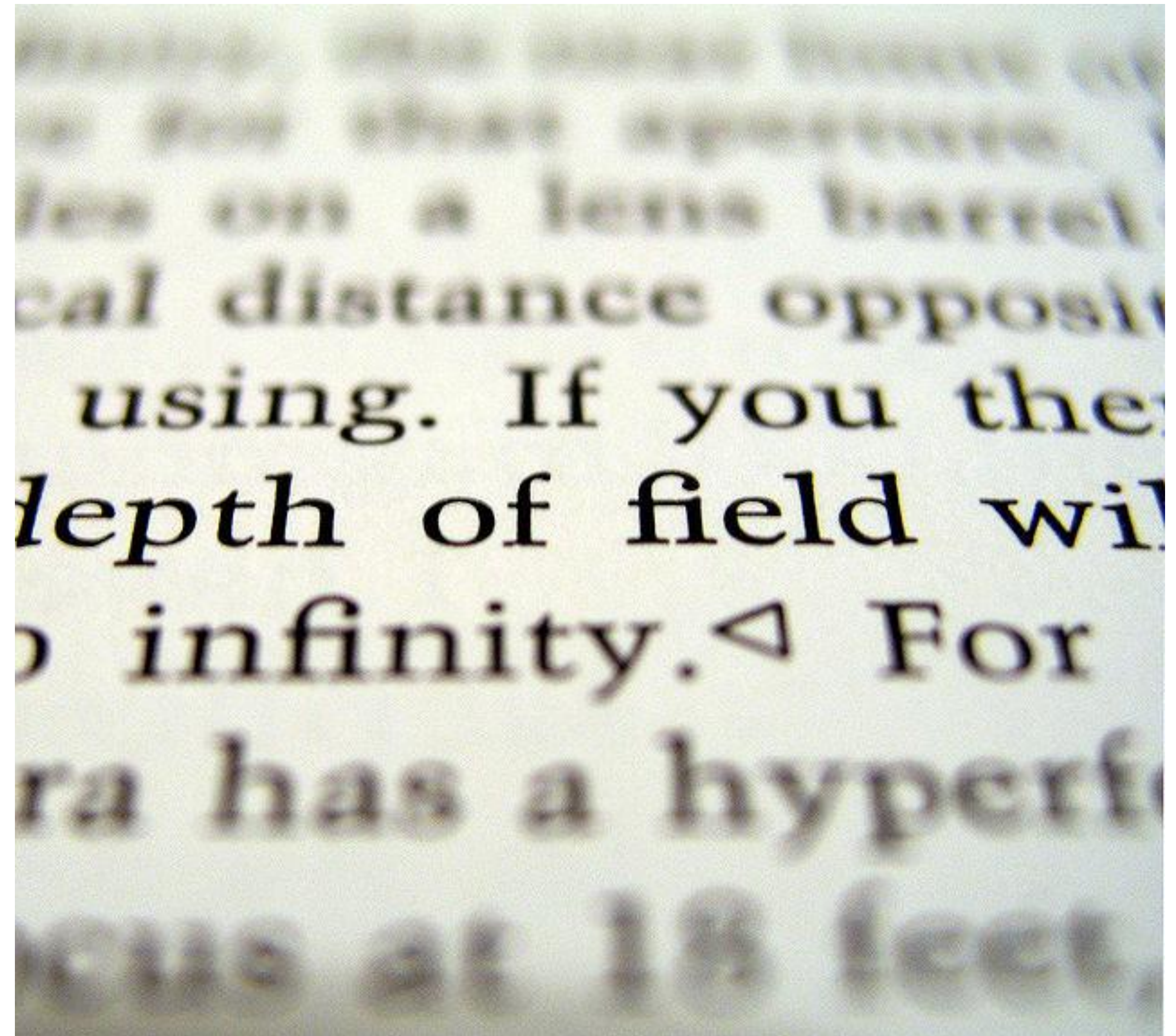
# Aperture / Depth of fields

- Same flowers with different apertures resulting in distinct depth of fields



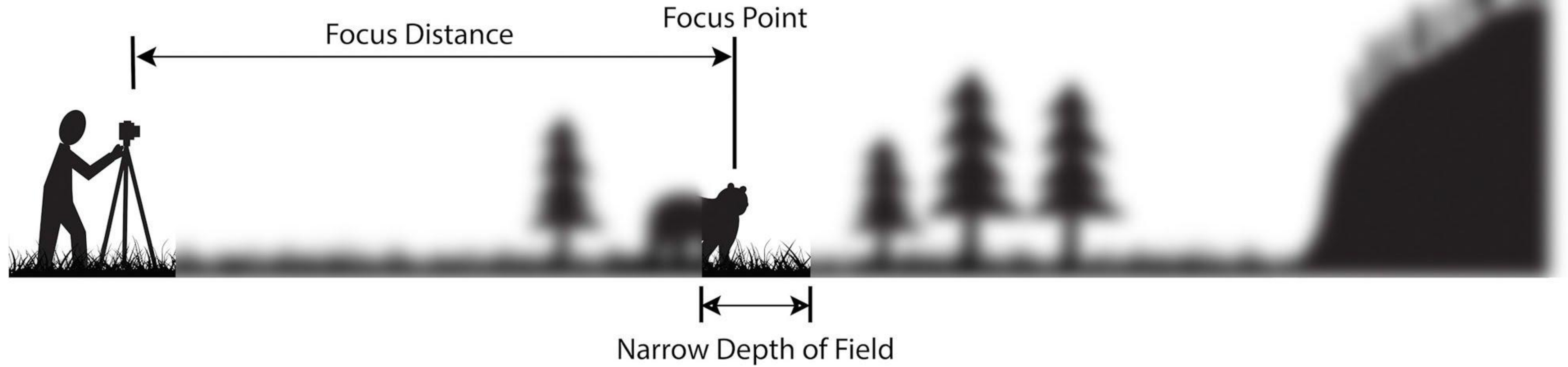
# Depth of Field

- The distance between the nearest and farthest objects in a scene that appear acceptably sharp in an image
- Precise focus is possible at only one distance; at that distance, a point object will produce a point image.
- When this circular spot is sufficiently small, it is indistinguishable from a point, and appears to be in focus; it is rendered as “acceptably sharp”.

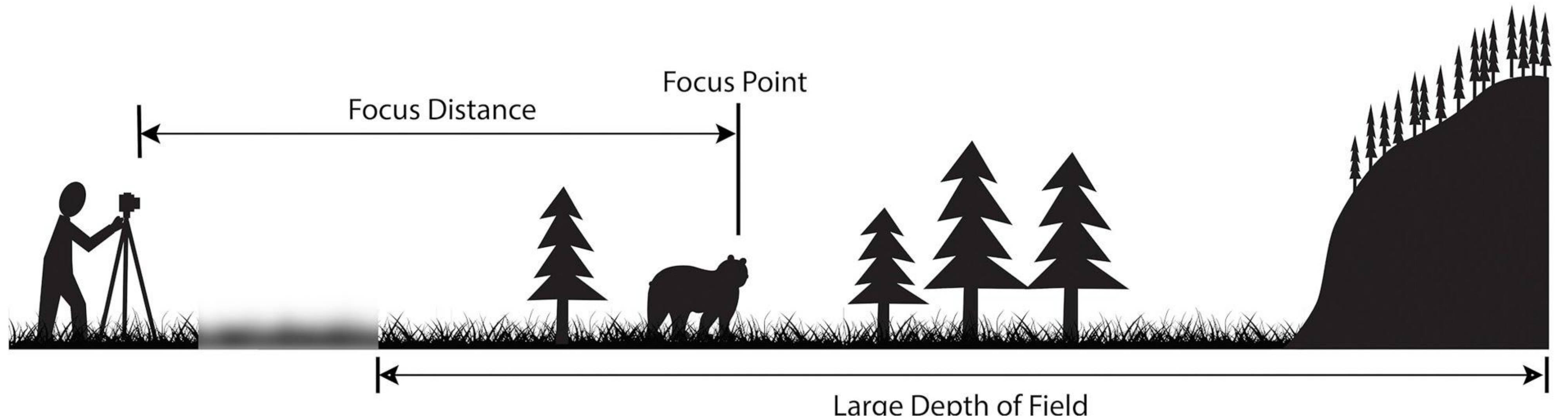


# Depth of Field

Big Aperture,  
Narrow Depth  
of Field



Small  
Aperture,  
Large Depth  
of Field



# ISO

- In Digital Photography ISO measures the sensitivity of the image sensor.



ISO 100

ISO 200

ISO 400

ISO 800

ISO 1600

ISO 100

ISO 3200

CLEAN IMAGE

NOISY IMAGE







Low ISO sensitivity, slow shutter speed.



High ISO sensitivity, fast shutter speed.

# From Light to Bits



photography



digitalization



(compression)

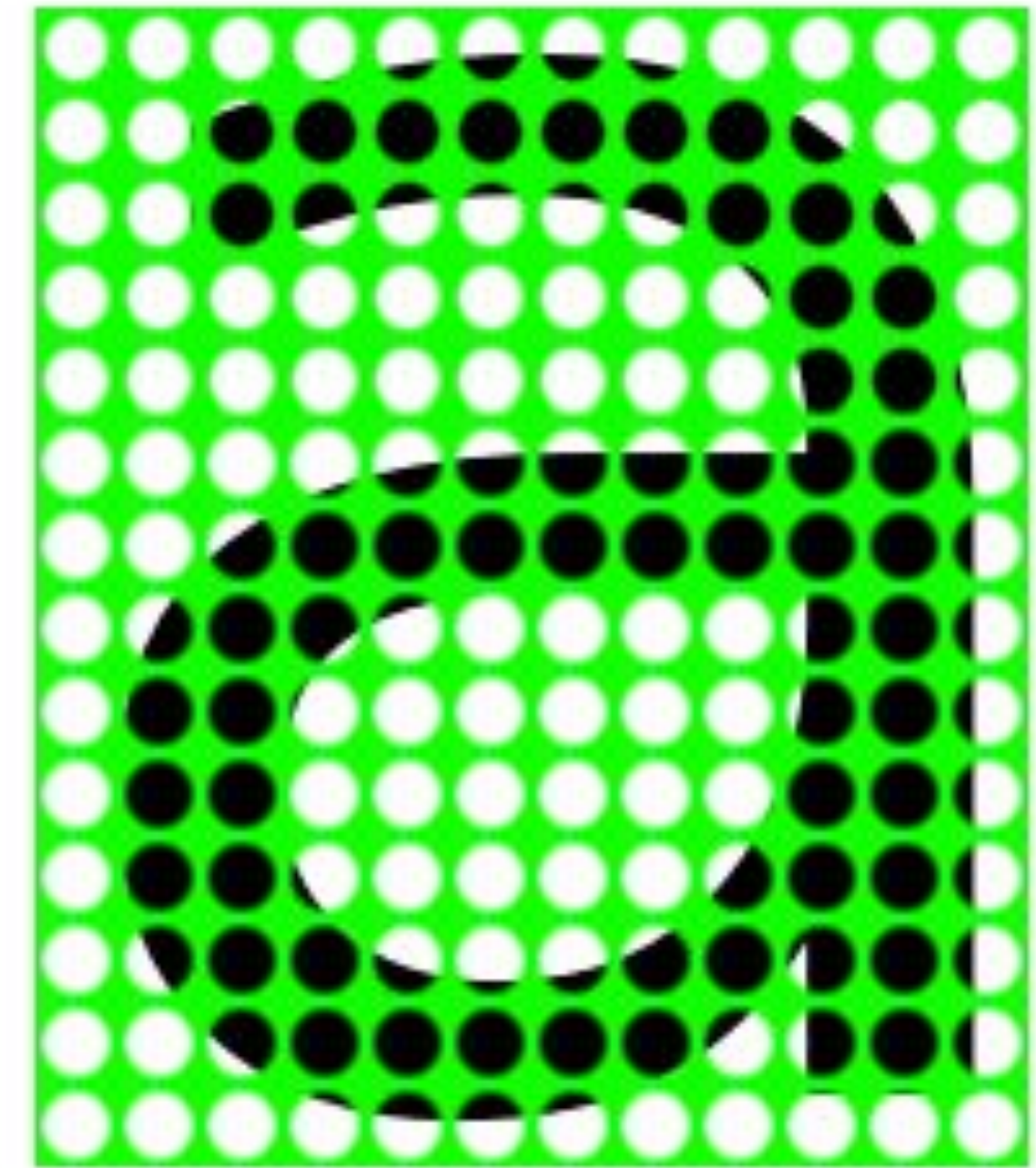
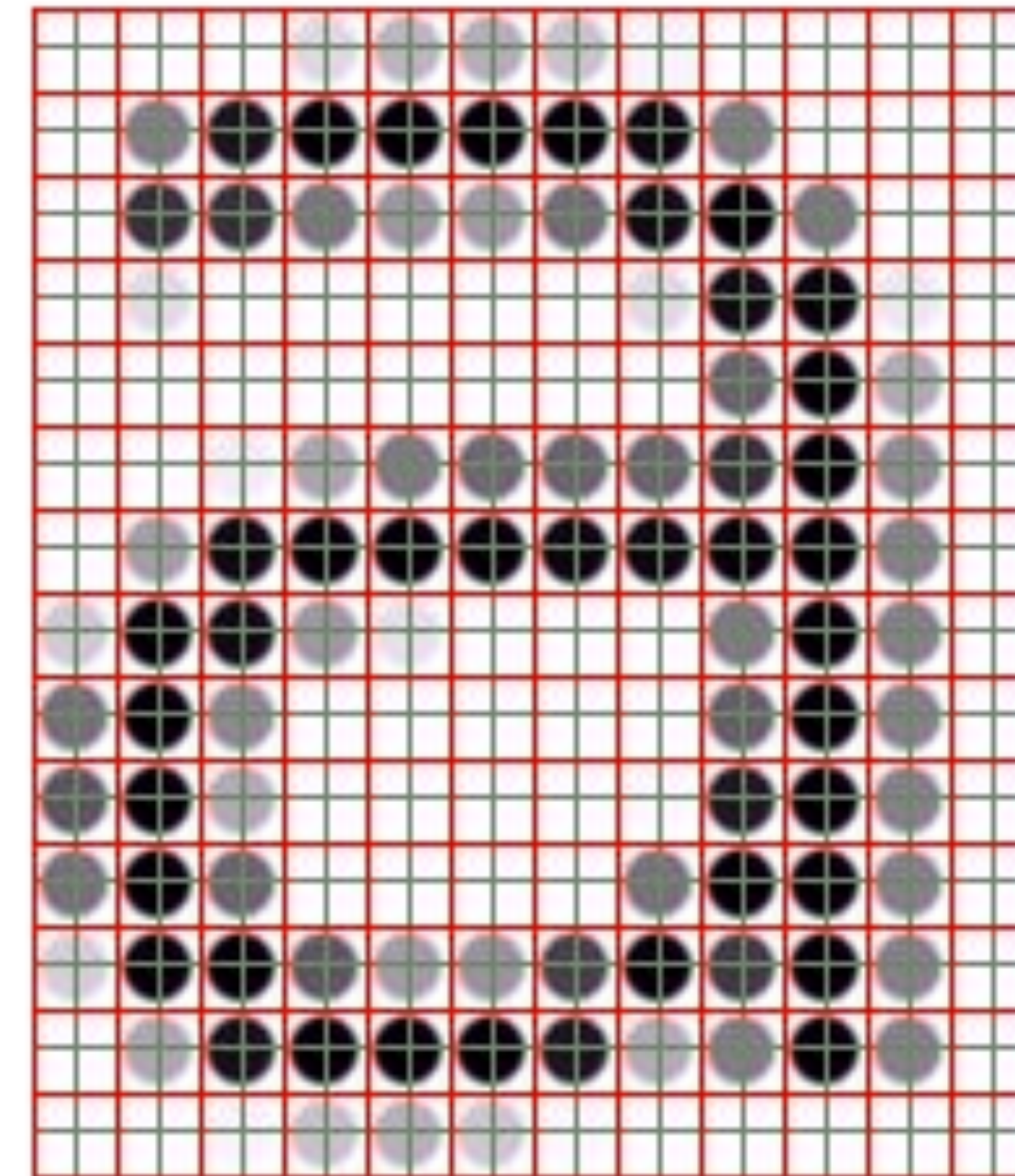
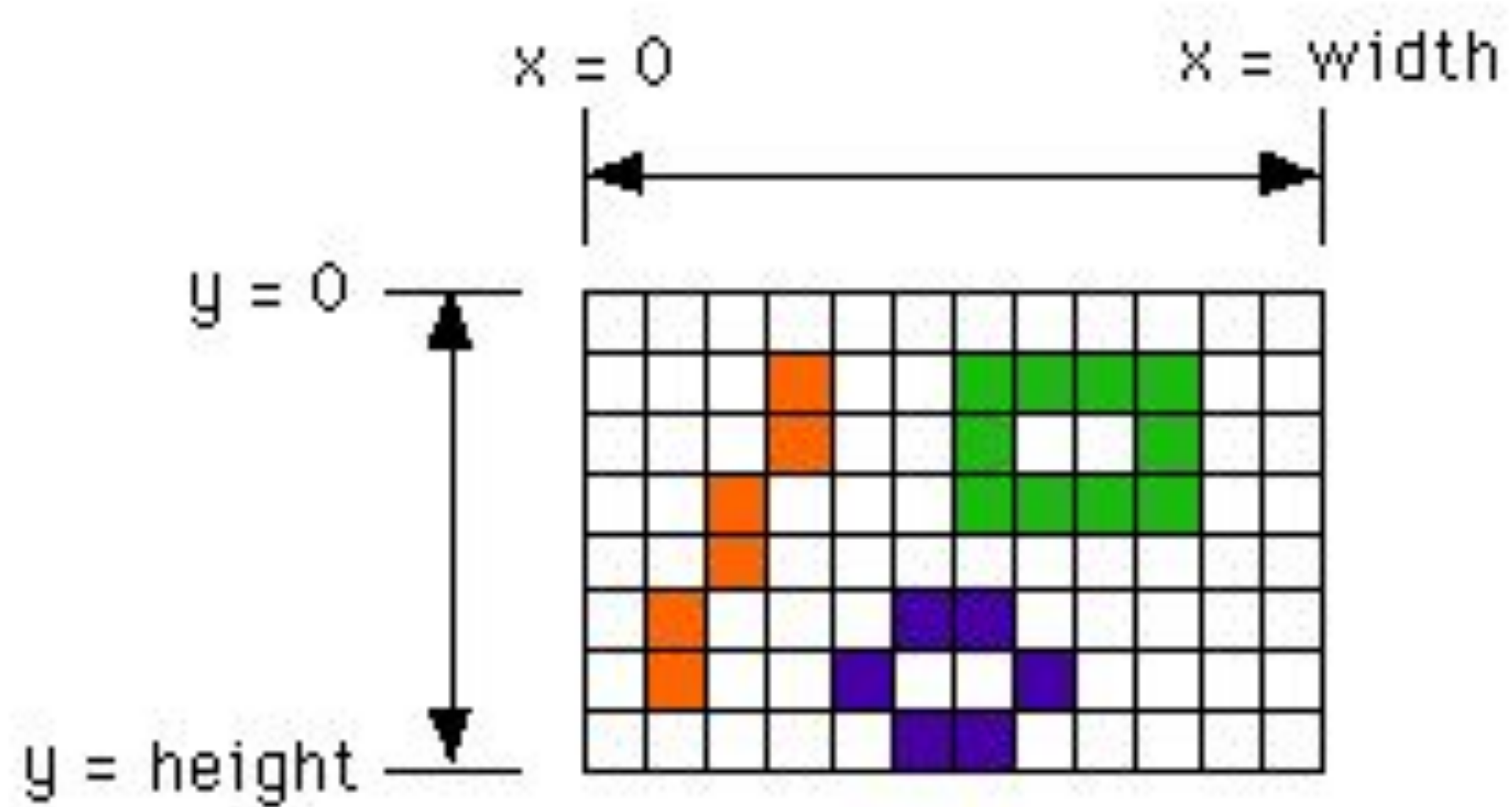


Camera

Photo

Bits

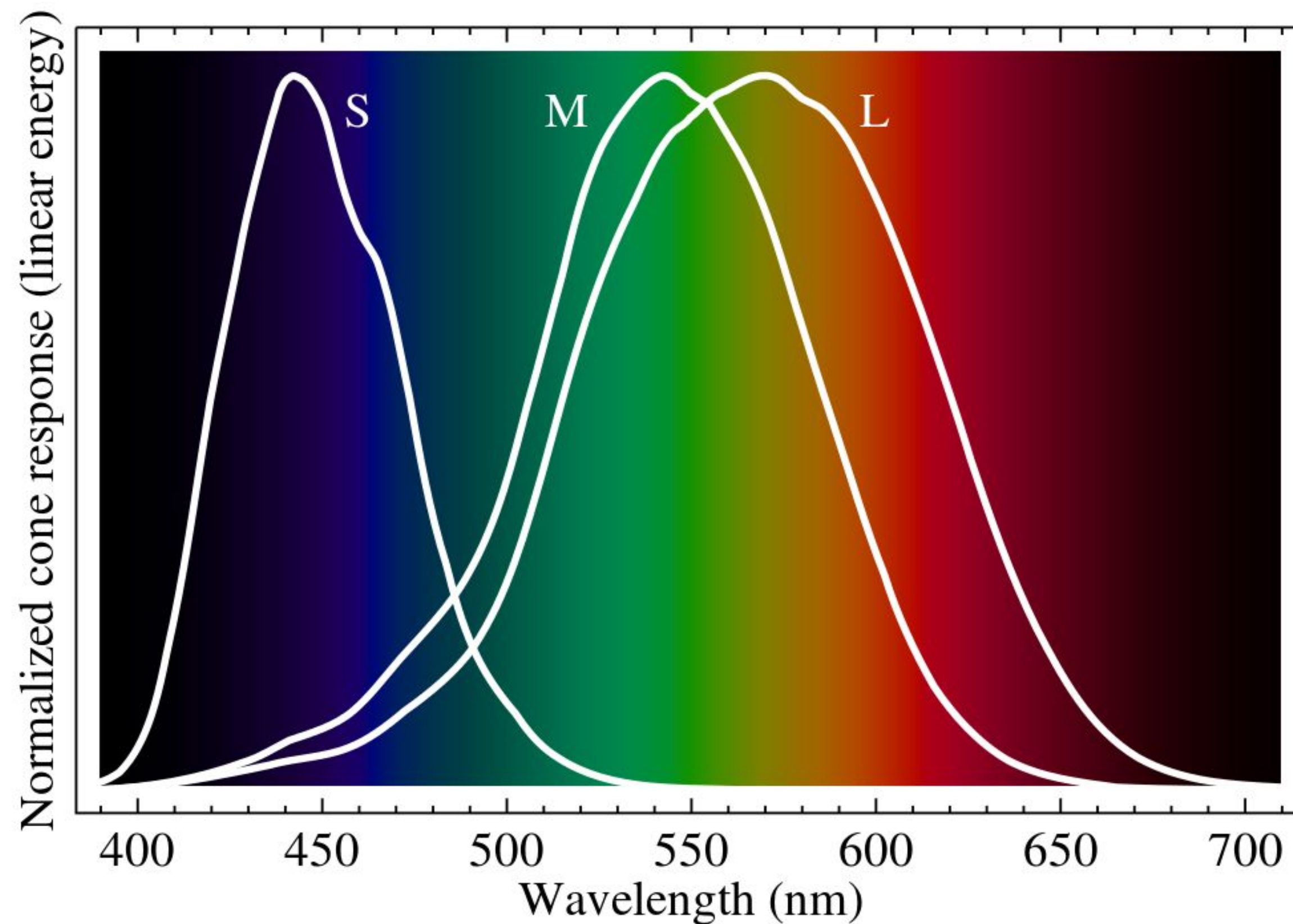
# Still Images



From the GIMP software documentation

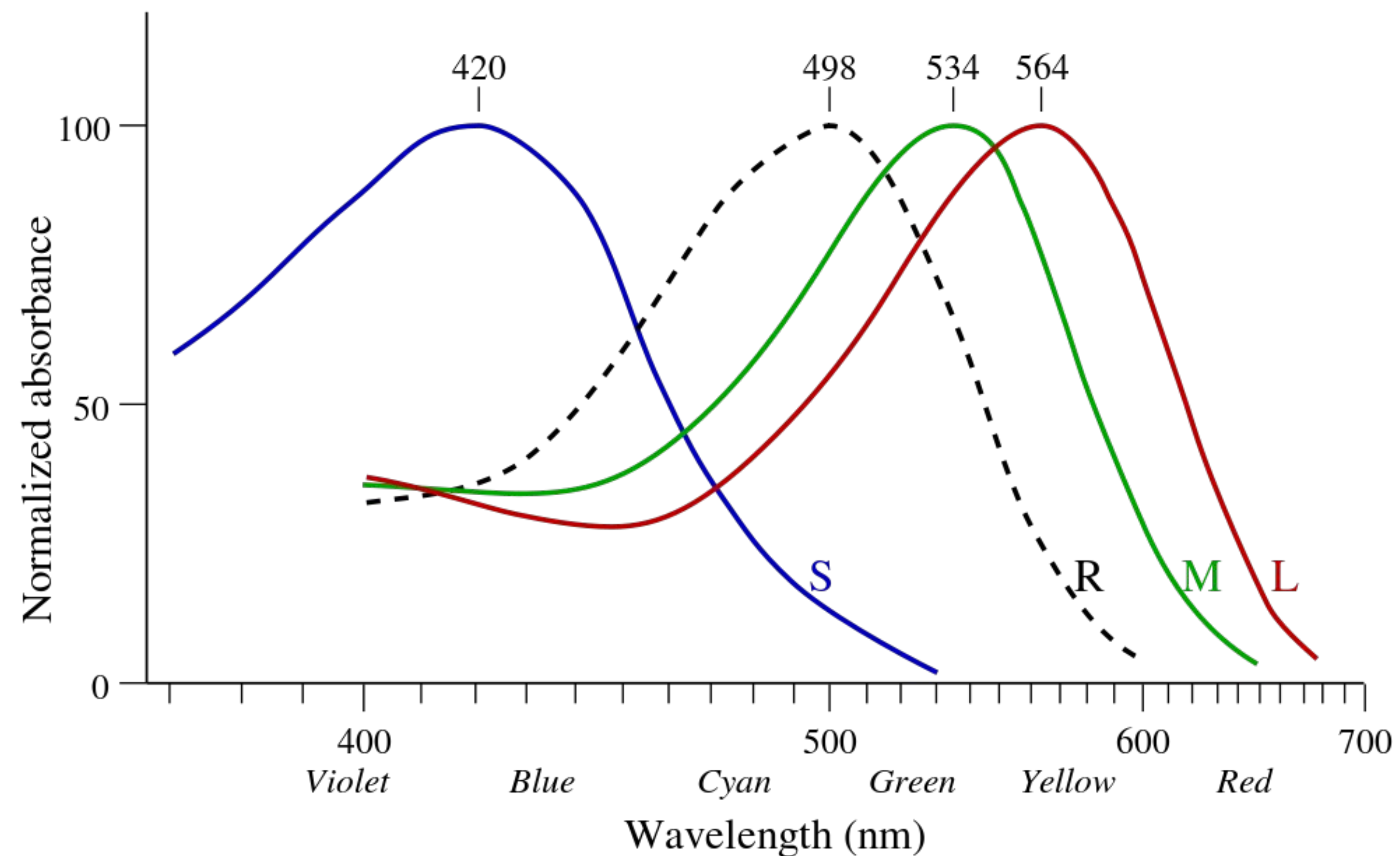
# Trichromacy

- Humans have 3 types of sensors (cones): L, M and S.
- Colors are perceived by the interaction of at least 2 types of cones
- N. of perceived color between 1 and 10 mln



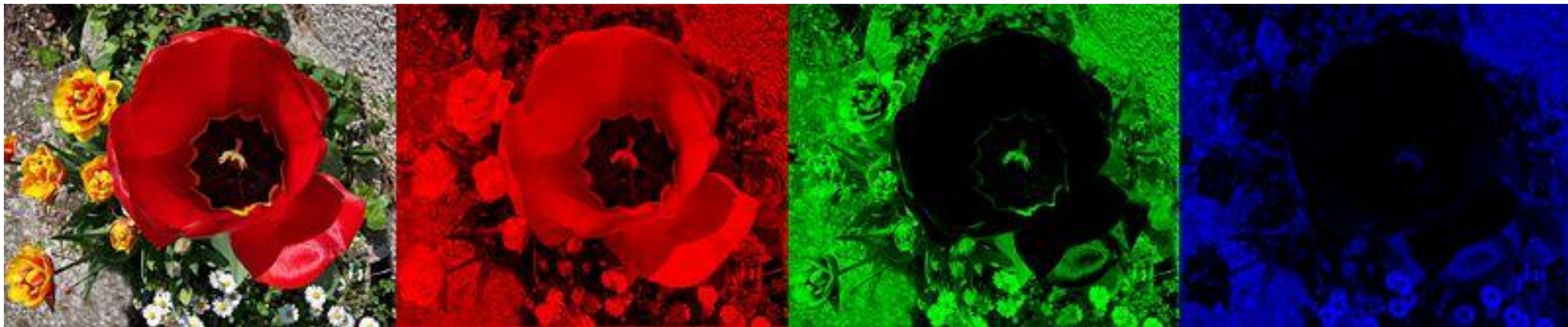
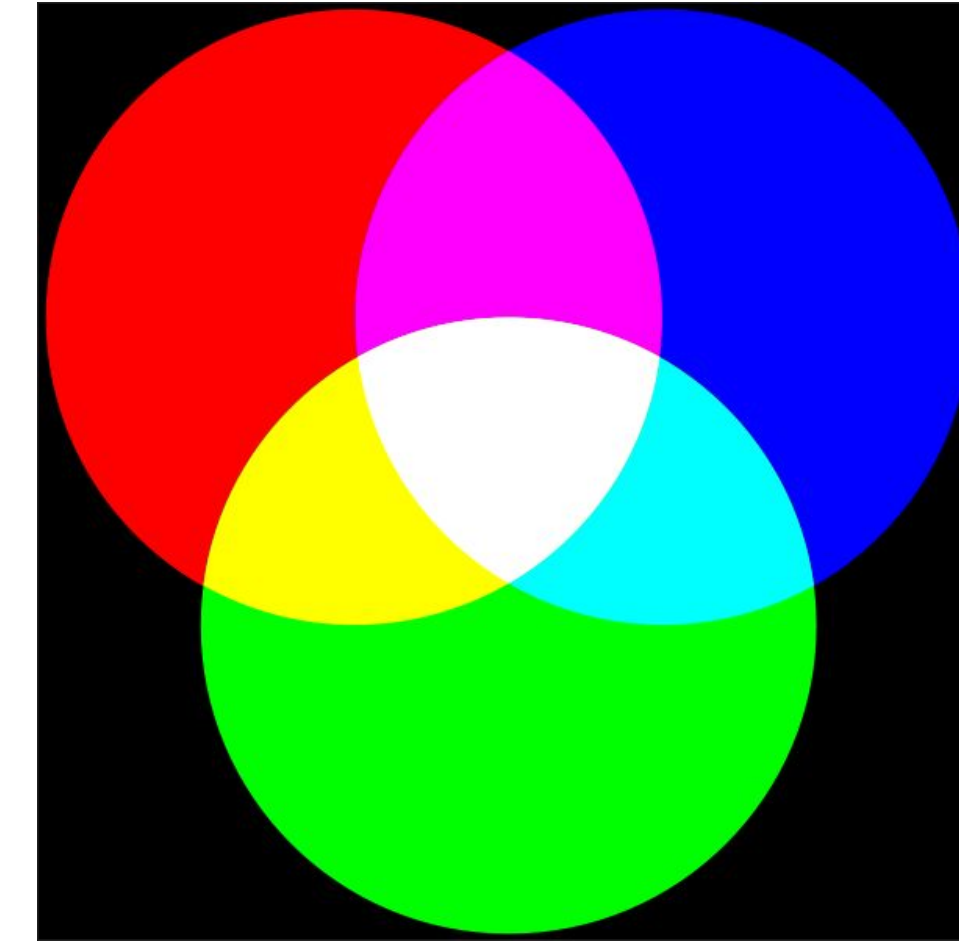
# Rod Cells

- Are in the retina of the eye
- Are sensitive to less intense light than the others (black dotted line)
- Concentrated on the outer edges
- Used in peripheral vision and in low light



# RGB

- Additive color model
- Based on human perception
- R, G, and B levels vary between devices
- Color management is needed
- Color spaces are used to ensure consistency



# Additive / Subtractive

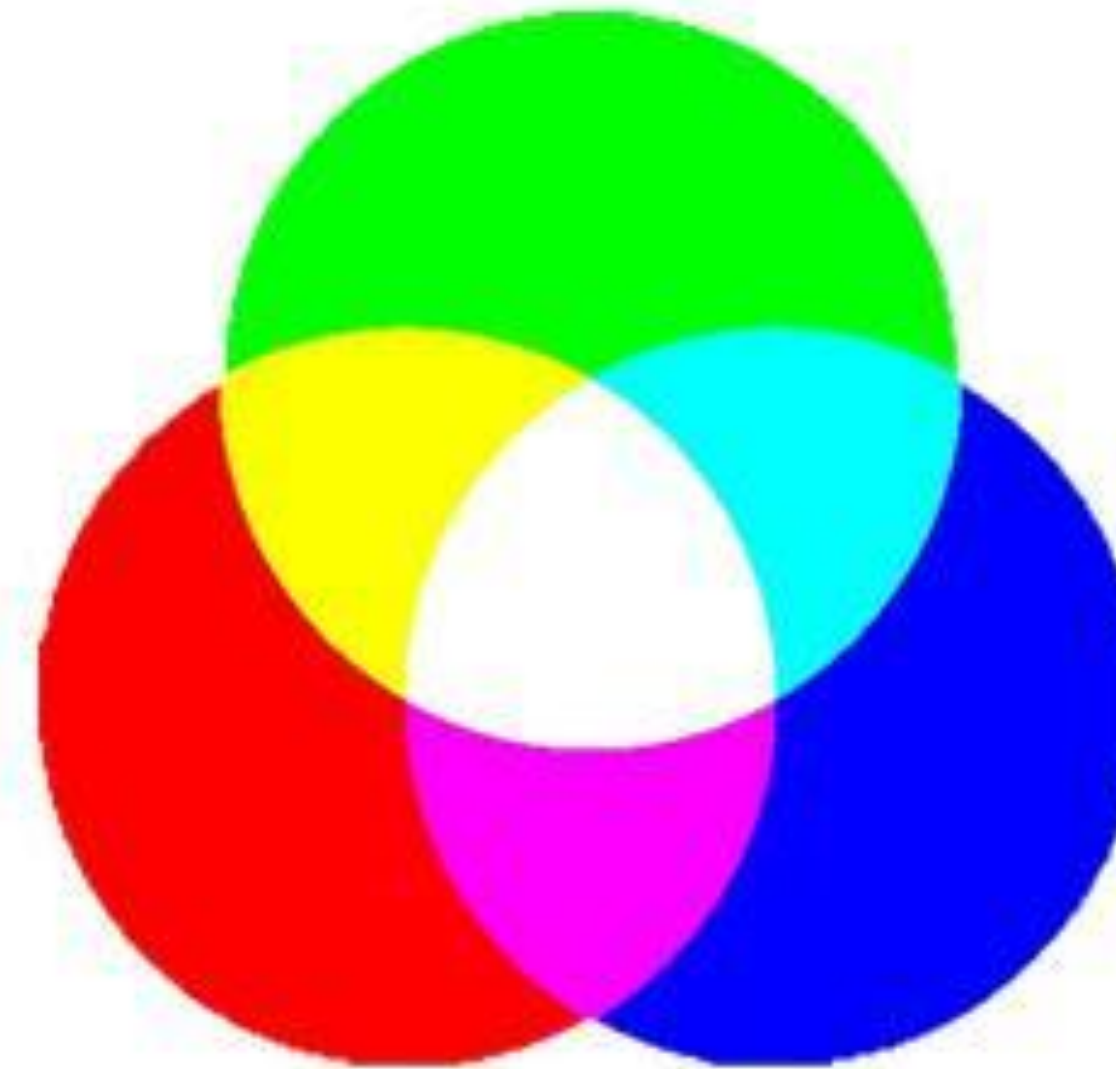
## RGB

Additive  
Color



*mixing light*

**RED GREEN BLUE**



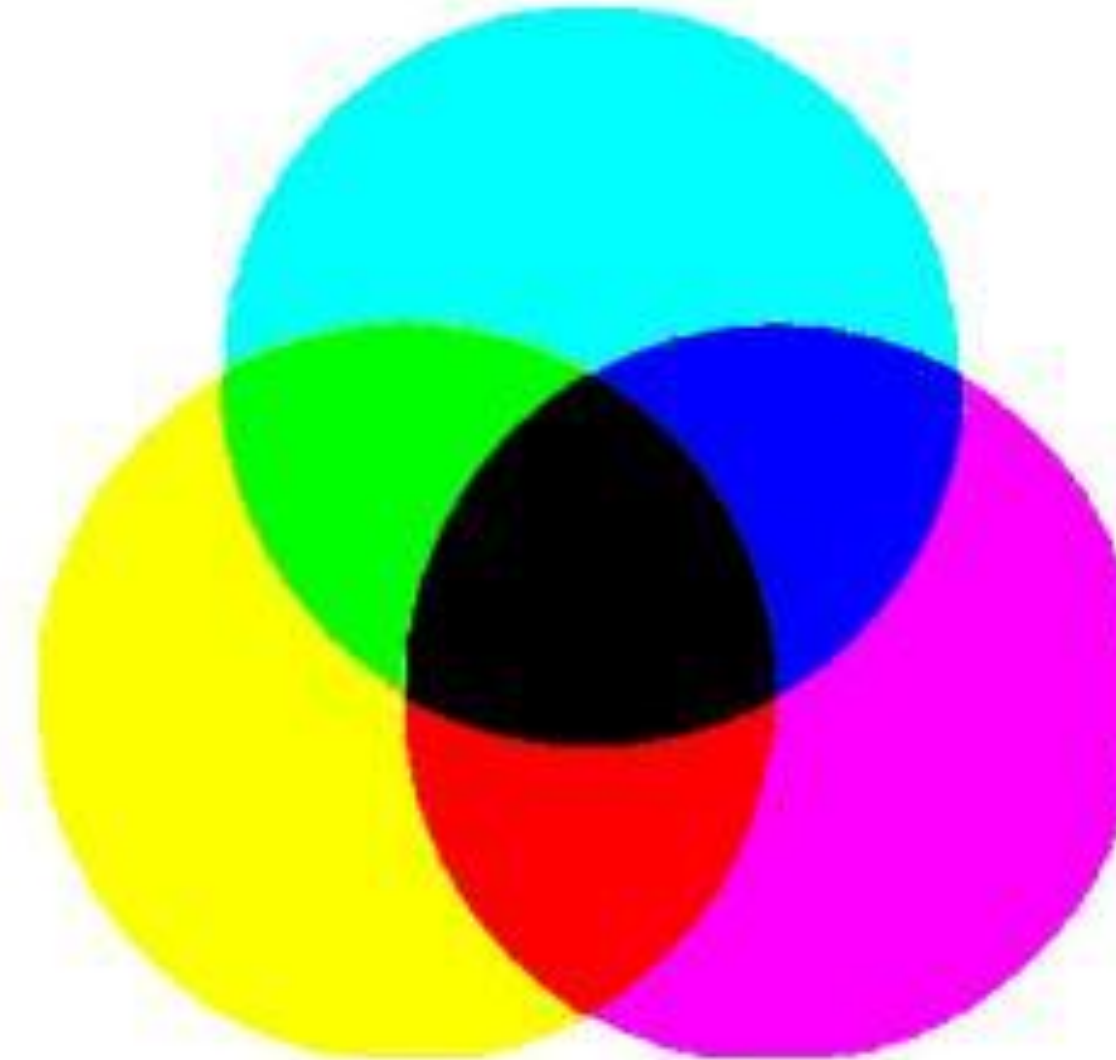
## CMYK

Subtractive  
Color



*mixing ink*

**CYAN MAGENTA YELLOW**



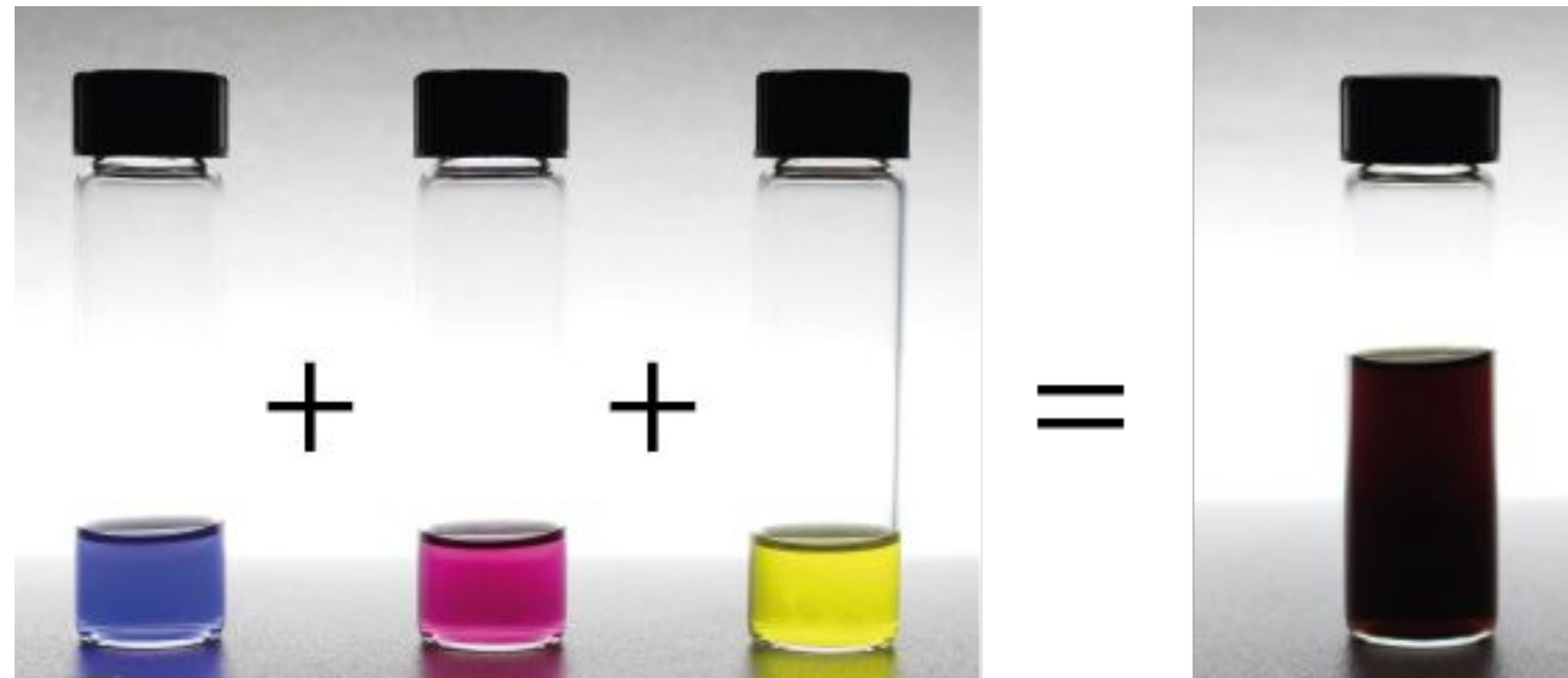
# Additive Color Mode

---

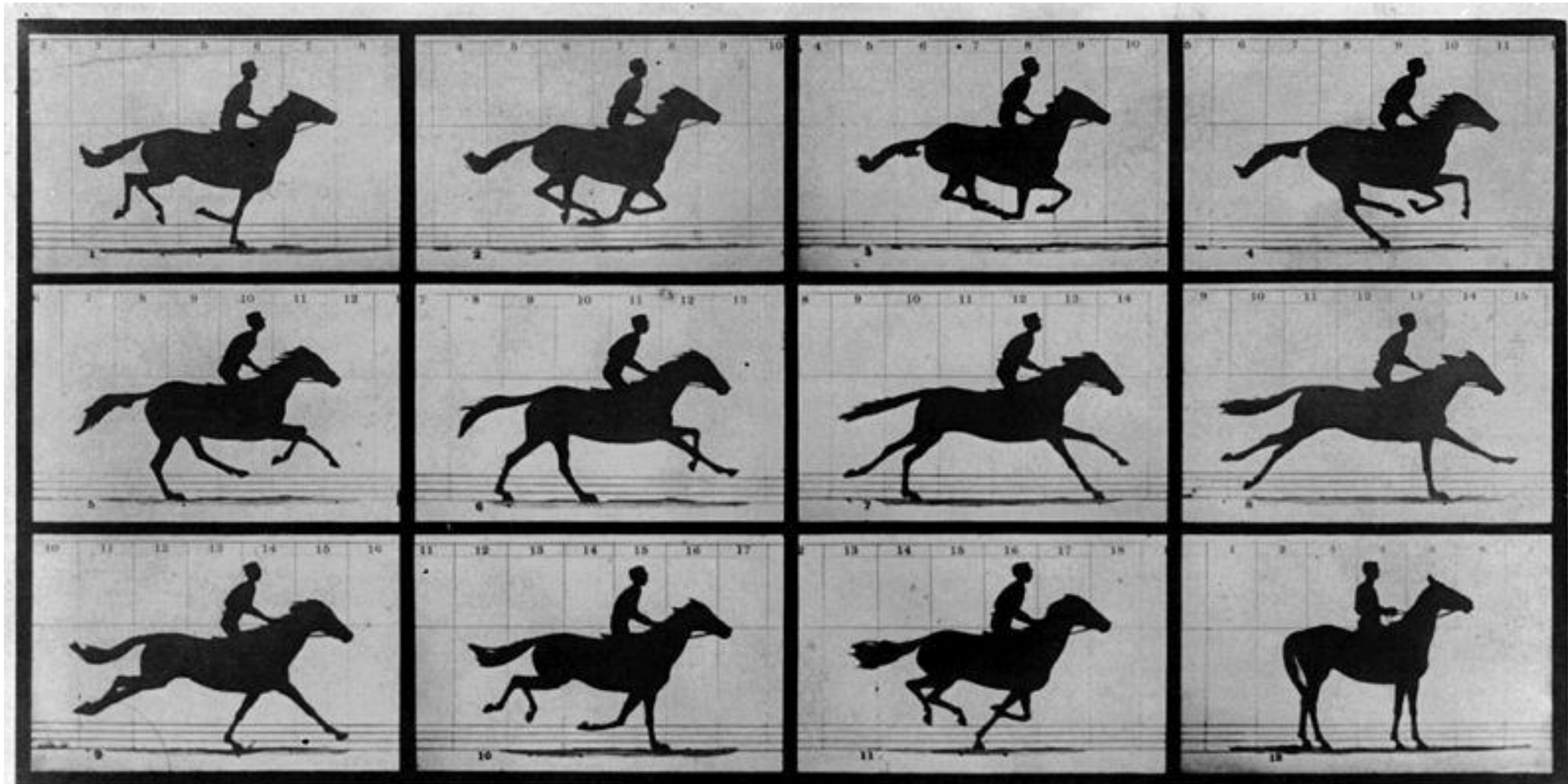




# Subtractive Color Mode



# 1878, Muybridge



Copyright, 1878, by MUYBRIDGE.

MORSE'S Gallery, 417 Montgomery St., San Francisco.

## THE HORSE IN MOTION.

Illustrated by  
MUYBRIDGE.

AUTOMATIC ELECTRO-PHOTOGRAPH

"SALLIE GARDNER," owned by LELAND STANFORD; running at a 1.40 gait over the Palo Alto track, 19th June, 1878.

The negatives of these photographs were made at intervals of twenty-seven inches of distance, and about the twenty-fifth part of a second of time; they illustrate consecutive positions assumed in each twenty-seven inches of progress during a single stride of the mare. The vertical lines were twenty-seven inches apart; the horizontal lines represent elevations of four inches each. The exposure of each negative was less than the two-thousandth part of a second.

# 1878: The horse in motion

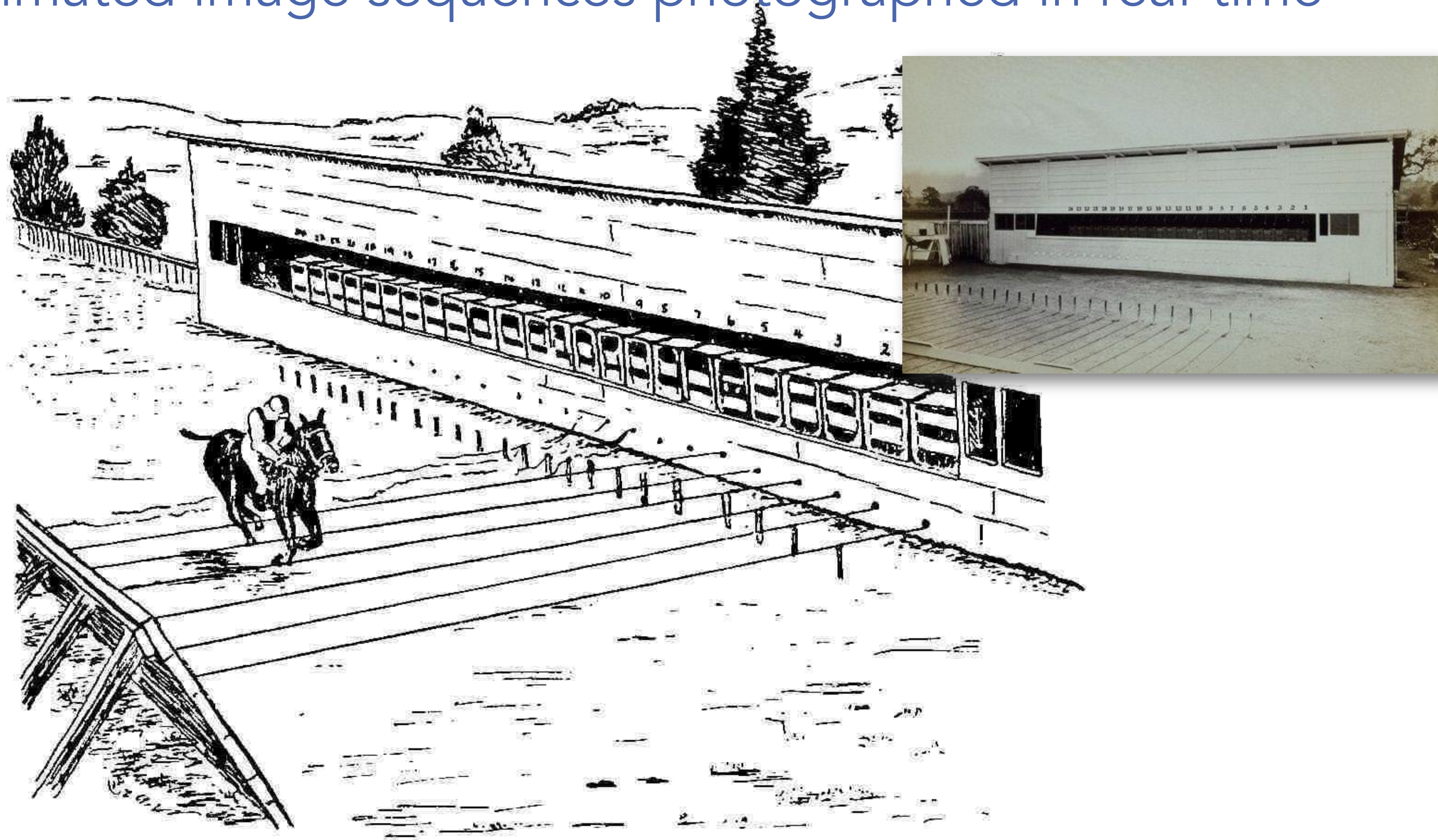
---



first animated image sequences photographed in real-time

# Moving pictures

- first animated image sequences photographed in real-time



# Reel

---

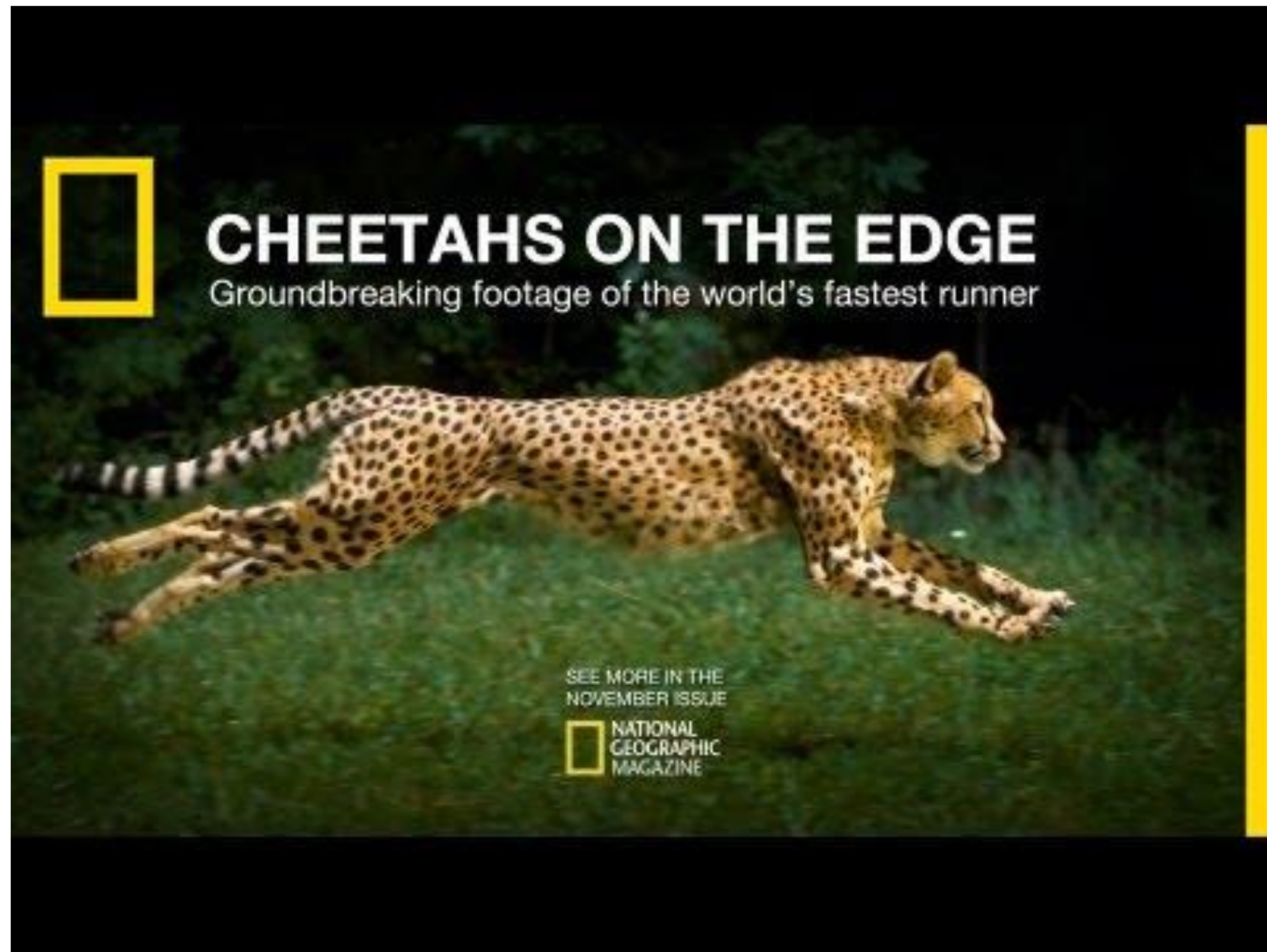


# How many frames per seconds?

---

- Frame Per Seconds (FPS):
  - 24: traditional 35 mm sound film starting from 1930
  - 25: PAL (EU TV)
  - 29.97: NTSC
  - 48: The Hobbit: An Unexpected Journey  
(accused to breaks the suspension of disbelief)
  - 50/60: HDTV
  - 72: experimental
  - 90/100: GoPro
  - 120: UHD TV
  - 144/240: Gaming monitors
  - 300: Tested by BBC for sports broadcasts

# FPS: recording vs display



FPS Recording  $>$  FPS Play  
slow motion video



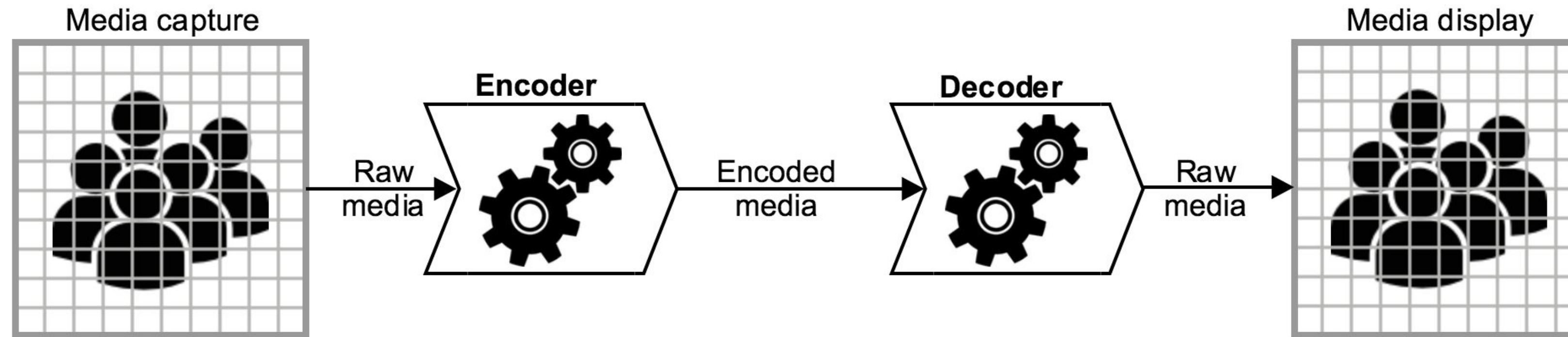
FPS Play  $>$  FPS Recording  
time lapse

# Storing and Sharing Images



# Codec

- Codec = coder-decoder (or compressor-decompressor)
- Software for encoding or decoding a data stream



- Lossy codecs (Quality vs Compression):
  - Mostly based on human perception, trying to achieve better human perceived quality at a given predefined compression
  - Examples: MPEG-2, MPEG-1 and 2 layer III (MP3), AAC, MPEG-2, H.264, DivX
- Lossless codecs (Original data can be perfectly reconstructed)
  - Mostly based on statistical modelling
  - Examples: LZW, FLAC, PNG, TIFF, etc.
  - Transcoding between lossy formats results in loss of data

# Image File Formats

IMAGE FORMAT	AVAILABLE COLORS	COMPRESSION	FILE SIZE	BEST FOR
RAW	Billions	No	Very big (<10MB)	Editing
JPEG	16,1 million	Lossy	Small (<1MB)	Websites and storage
GIF	256	Lossless	Small (<1MB)	Animation
PNG	16,1 million + transparency	Lossless	Big (<3MB)	Websites, editing, storage
TIFF	Variable	Variable	Big (<3MB)	Editing and printing
BMP	Variable	Lossless	Big (<3MB)	-

# Sharing Image Collections

---

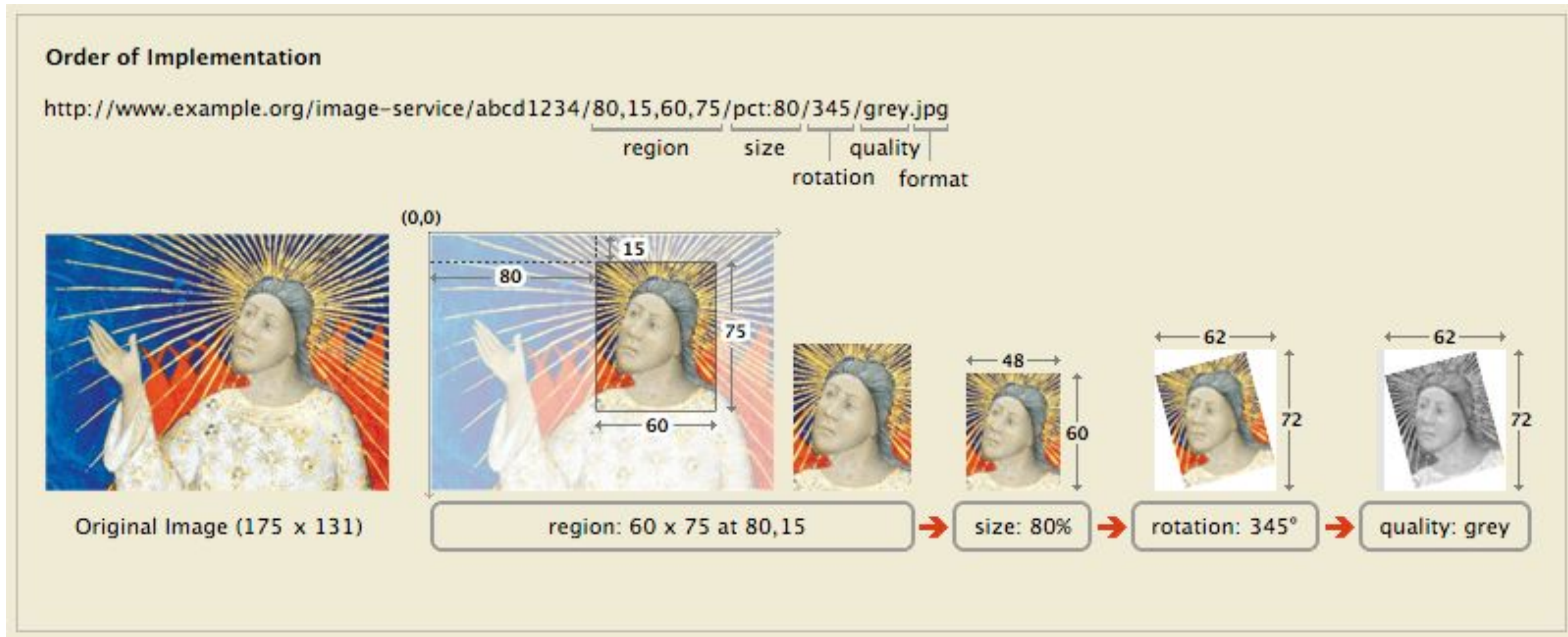
- IIIF - International Image Interoperability Framework (<https://iiif.io>)
- Set of open standards for delivering high-quality, attributed digital objects (images, audio/visual) online at scale.
- Born to “facilitate systematic reuse of image resources in digital image repositories maintained by cultural heritage organizations.”



International  
Image  
Interoperability  
Framework

# IIIF Image API

- Define how image data can be served and accessed (URL-based API)



# IIIF Presentation API

---

- . Provides the information necessary to allow a rich, online viewing environment for compound digital objects to be presented to a human user.
- . Manifest file in JSON format containing:
  - . Descriptive metadata like labels, rights and other information
  - . Links to Images and AV resources
  - . Ordering of Images in sequences and table of contents

- . Demos

<https://uv-v3.netlify.app/>

<https://demos.biblissima.fr/chateauroux/osd-demo/>

# More on IIIF

---

- <https://iiif.io/get-started/>
- <https://training.iiif.io/>
- IIIF-compliant image servers: <https://iiif.io/get-started/image-servers/>
- IIIF-compliant image viewers: <https://iiif.io/get-started/iiif-viewers/>



Consiglio Nazionale delle Ricerche



Istituto di Scienza e Tecnologie  
dell'Informazione "A. Faedo"



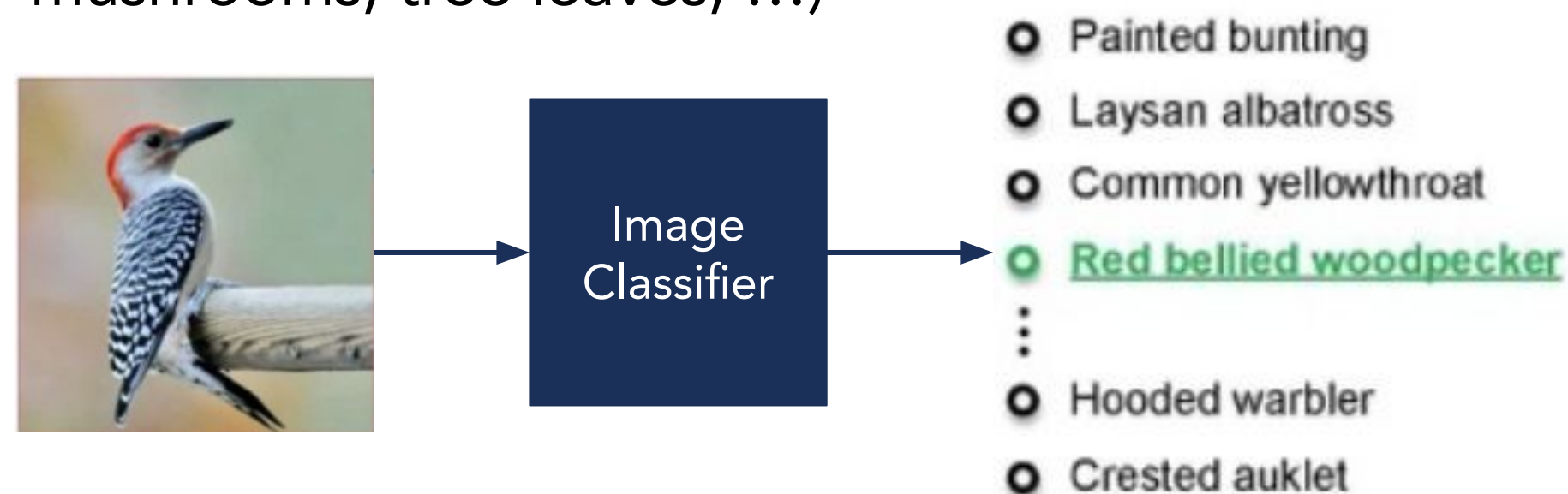
# IMAGE CLASSIFICATION

---

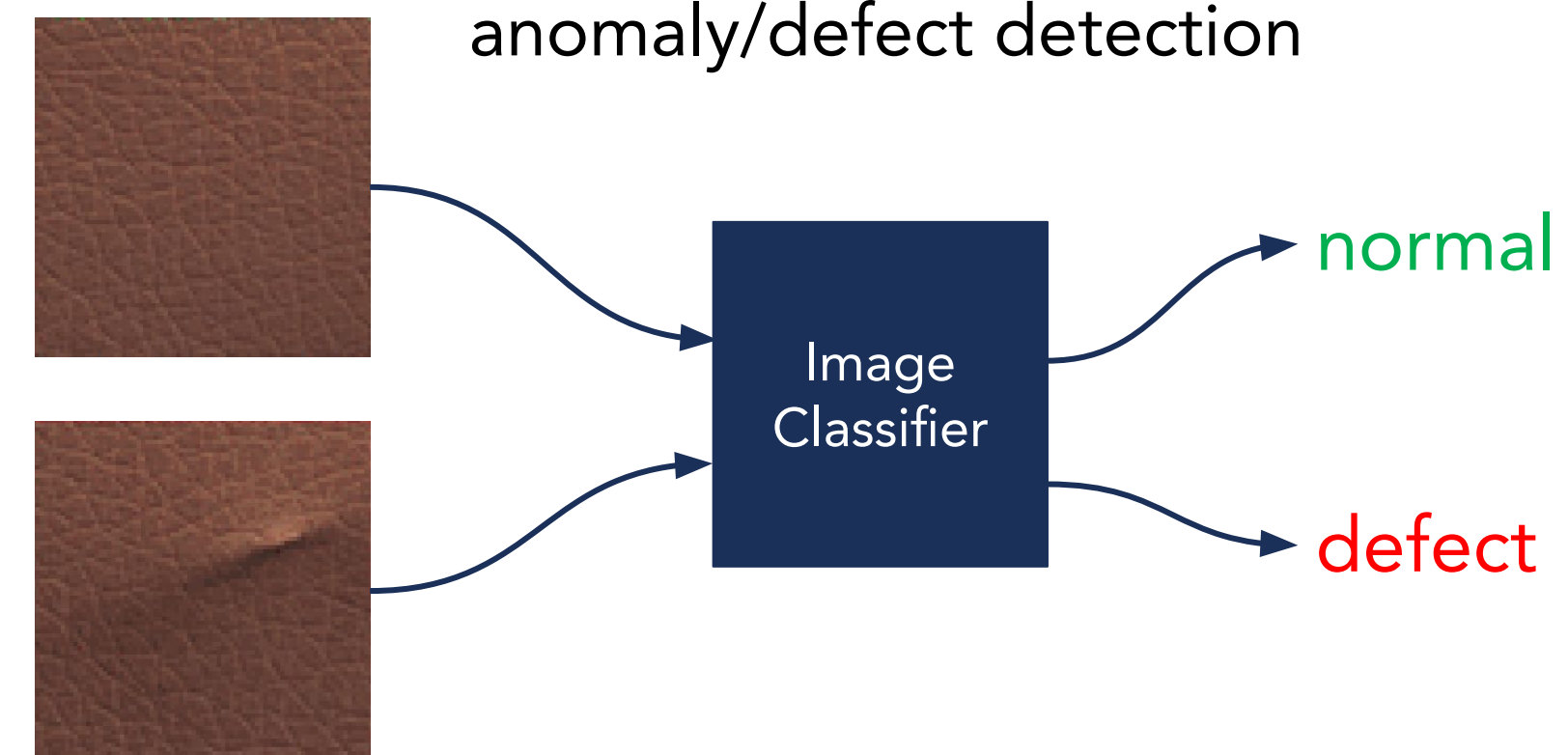
# Image Classification

- Automatically assign an image to categories or classes of interest depending on its visual content only. (No metadata available).
- Several problems can be framed as image classification. E.g.:

fine-grained animal/object recognition (birds, mushrooms, tree leaves, ...)

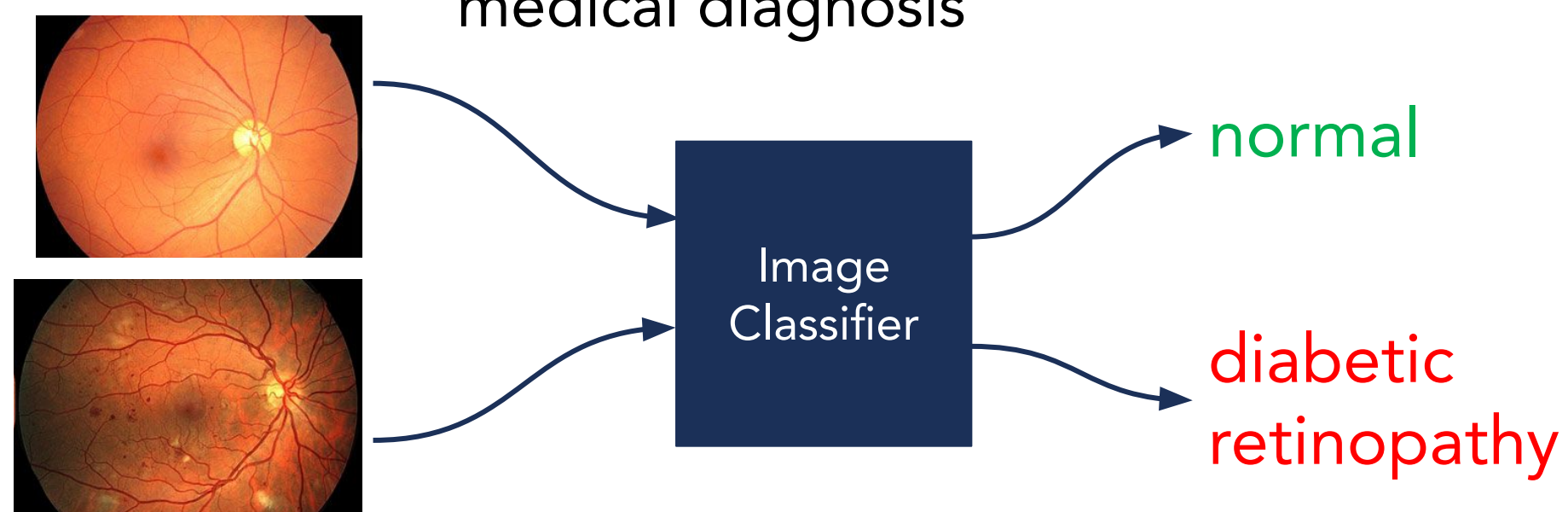


anomaly/defect detection

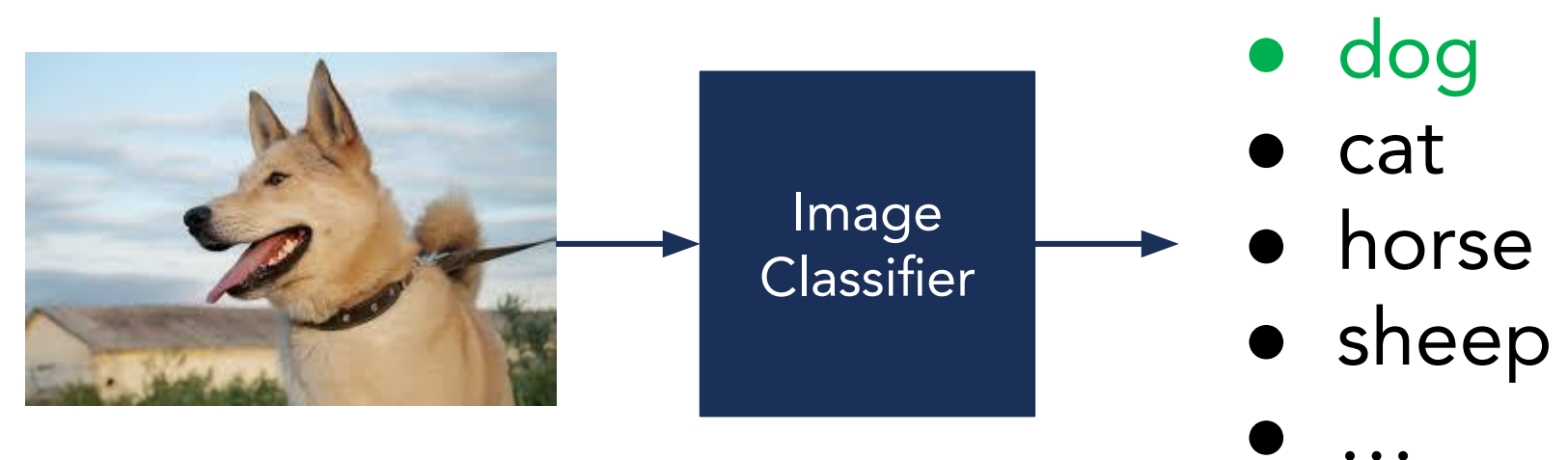


and many more....

medical diagnosis



generic object recognition





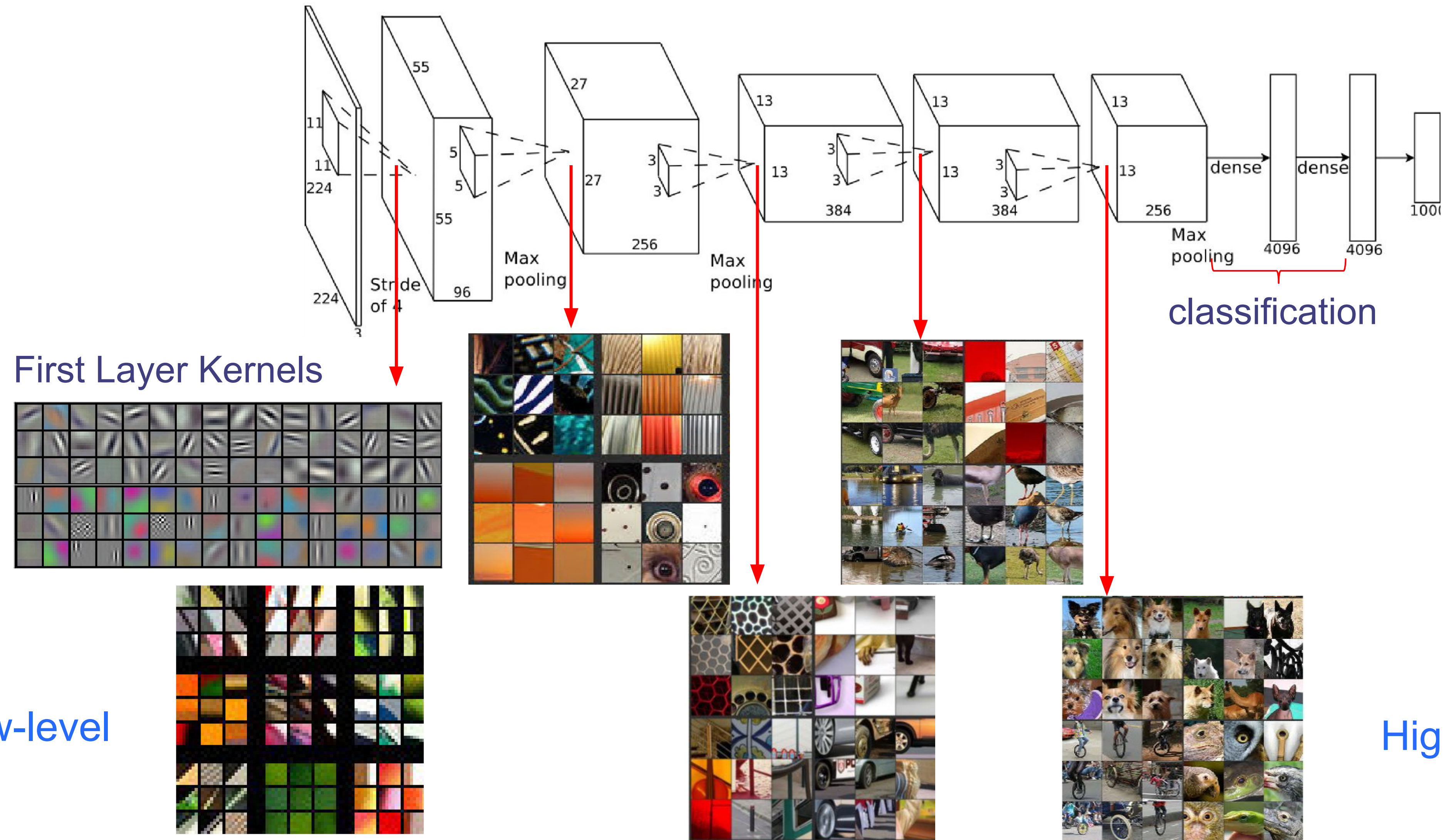
# Deep Learning for Image Classification

---

- Understand relationships between input images and categories is often difficult to do manually (hand-crafted features).
- Most solutions use Deep Learning (specifically, artificial neural networks) to learn the mapping between image and category.
- Networks must be trained on a set of images to learn the specific mapping; neurons change connections to learn patterns during training.
- Once trained, the network is frozen and used to classify new images.

# Multiple Levels Of Abstraction

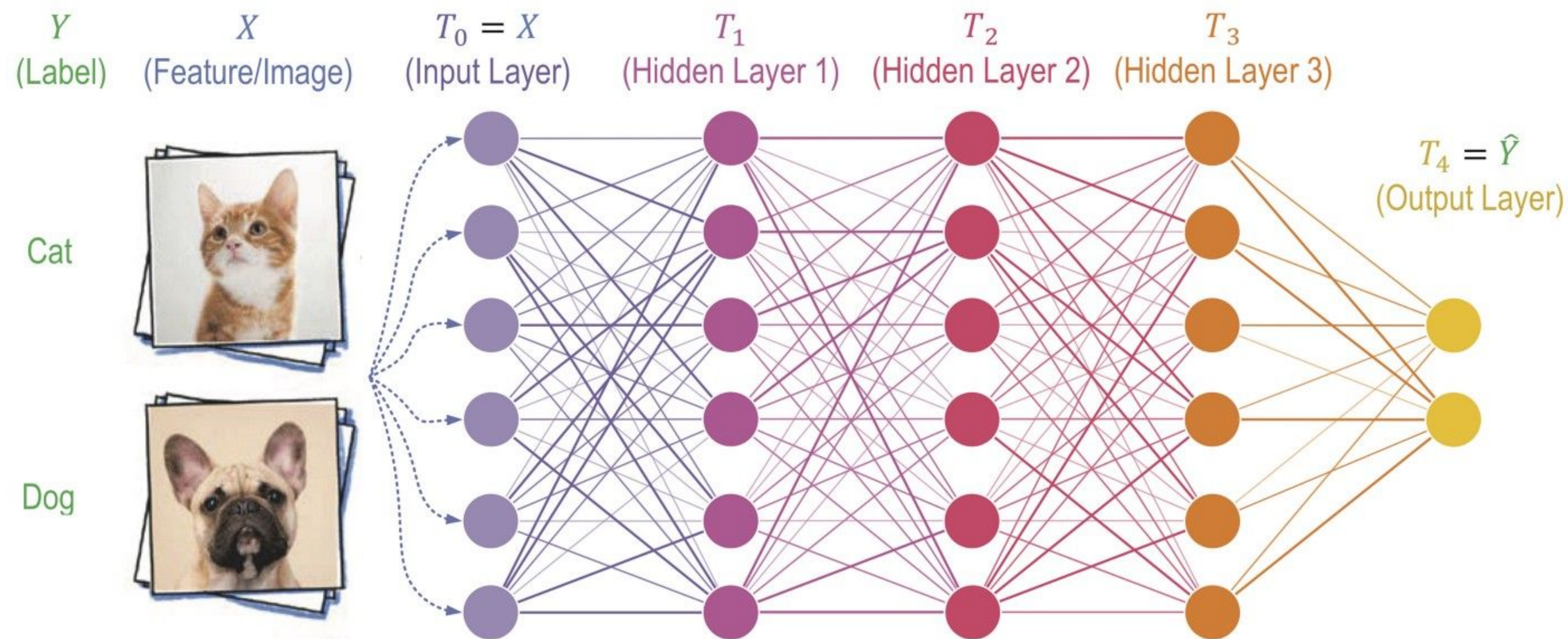
AlexNet, 2012, Trained on a Classification task of 1,000 classes.



# Closed-set vs Open-set Classification

## • Closed-set classification

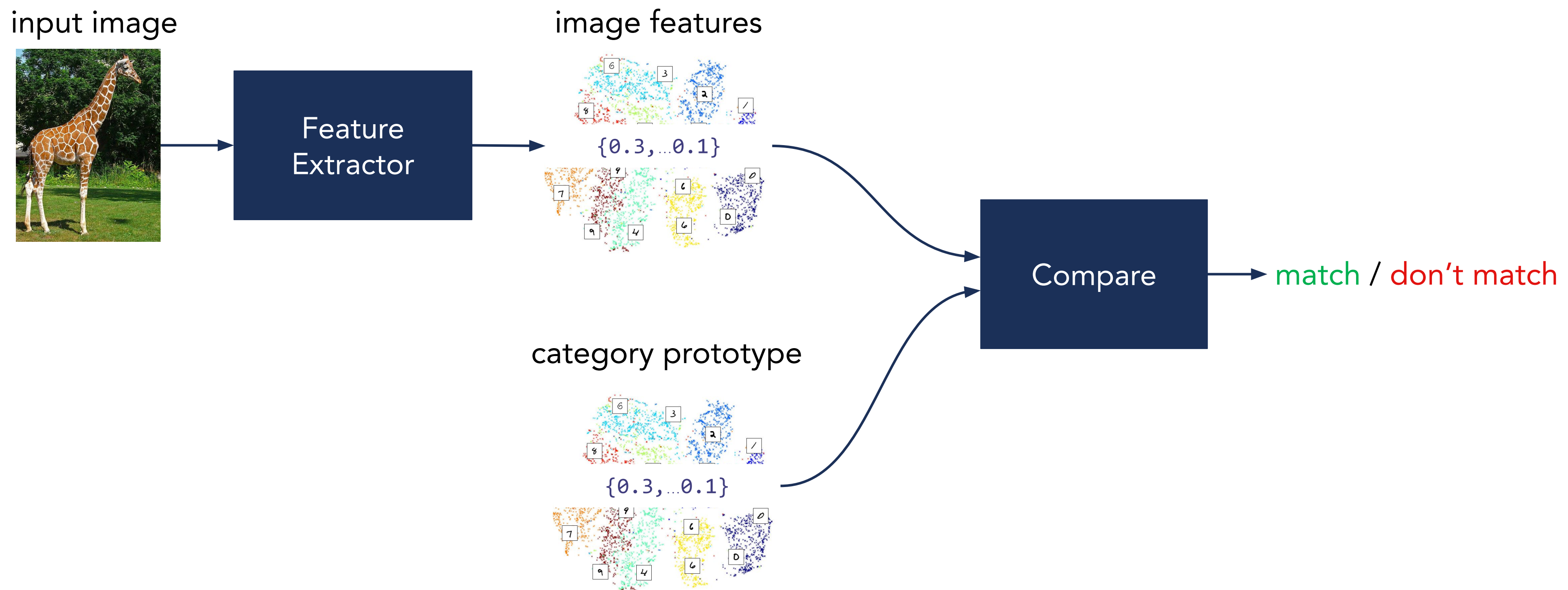
- We want to classify the input image into one of  $N$  predefined categories.
- The classifier guesses the best one out of the  $N$  categories.
- The classifier is trained on examples belonging to the  $N$  categories.



# Closed-set vs Open-set Classification

## • Open-set classification

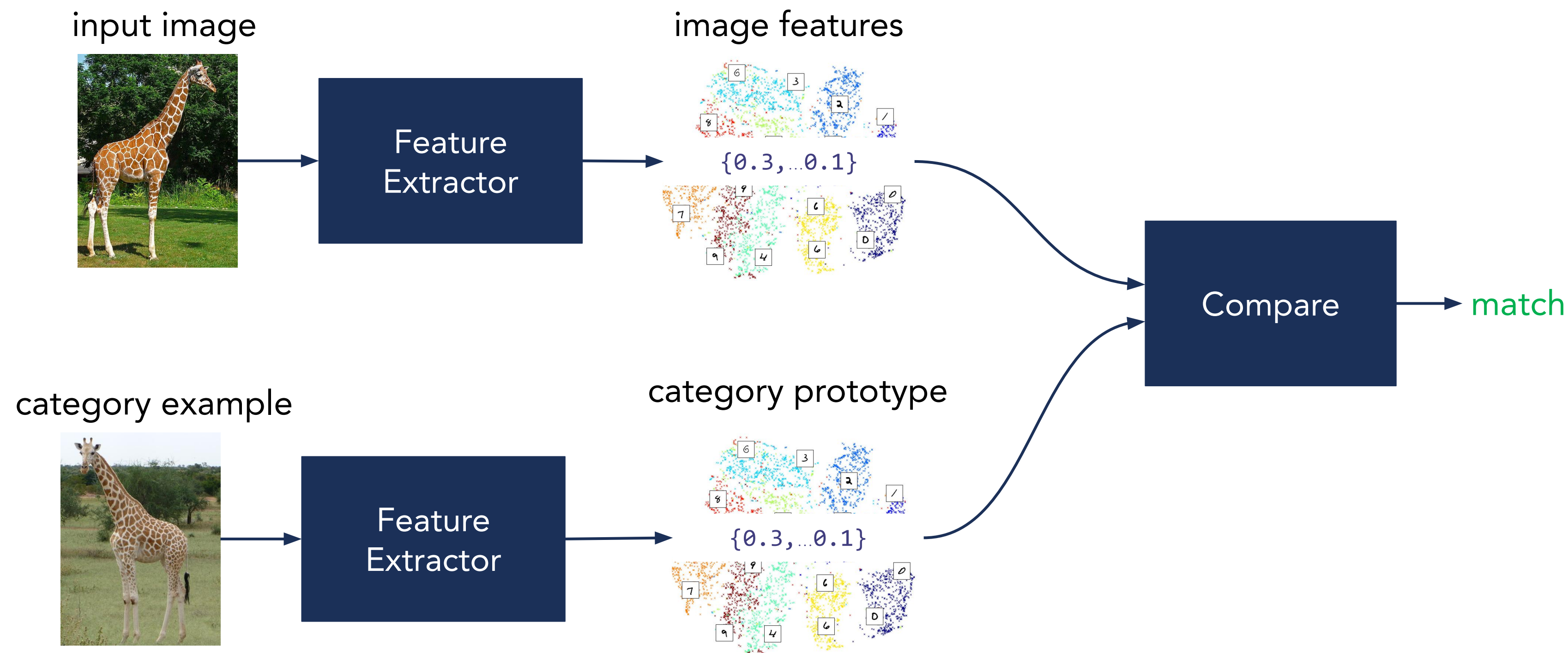
- We do not have a predefined set of categories
- The classifier is trained to extract generic image features/representations (a string of numbers!)
- Image features are compared with novel category “prototypes” to check if they match



# Closed-set vs Open-set Classification

## • Open-set classification

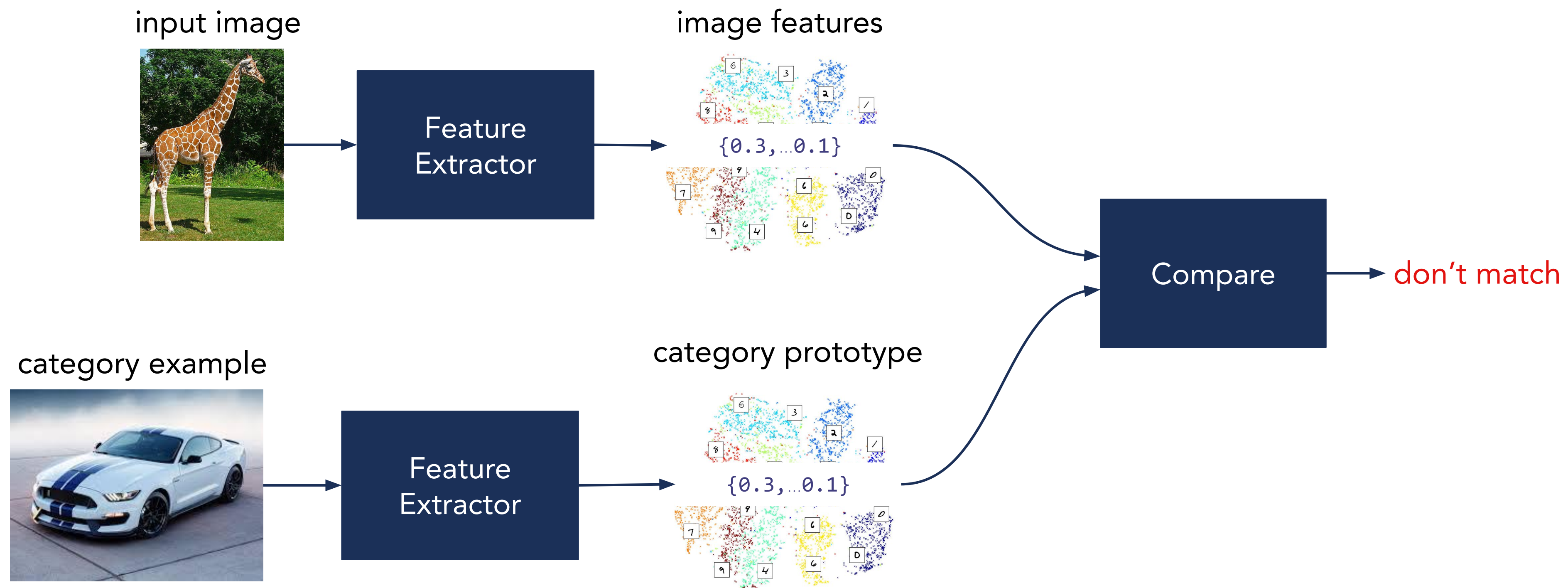
- We do not have a predefined set of categories
- The classifier is trained to extract generic image features/representations (a string of numbers!)
- Image features are compared with novel category “prototypes” to check if they match



# Closed-set vs Open-set Classification

## • Open-set classification

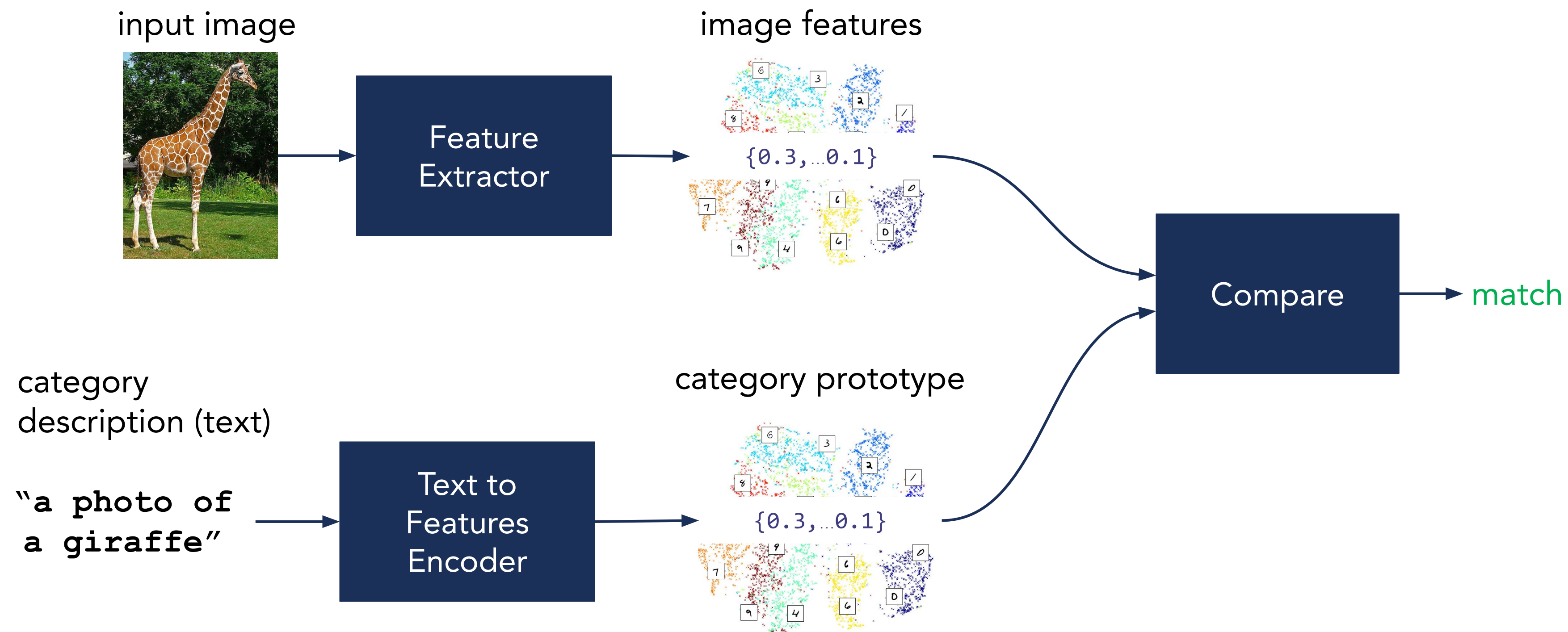
- We do not have a predefined set of categories
- The classifier is trained to extract generic image features/representations (a string of numbers!)
- Image features are compared with novel category “prototypes” to check if they match



# Closed-set vs Open-set Classification

## • Open-set classification

- We do not have a predefined set of categories
- The classifier is trained to extract generic image features/representations (a string of numbers!)
- Image features are compared with novel category “prototypes” to check if they match





Consiglio Nazionale delle Ricerche



Istituto di Scienza e Tecnologie  
dell'Informazione "A. Faedo"



# IMAGE RETRIEVAL

---



# Nicola Messina

---

- M.Sc. Computer Engineering @ UniPi (Feb. 2018)
- PhD in Information Engineering @ UniPi (May 2021)
- PostDoc @ Istituto di Scienza e Tecnologie dell'Informazione (ISTI) - CNR
  - AIMH (Artificial Intelligence for Media and Humanities) Lab

[nicola.messina@isti.cnr.it](mailto:nicola.messina@isti.cnr.it)

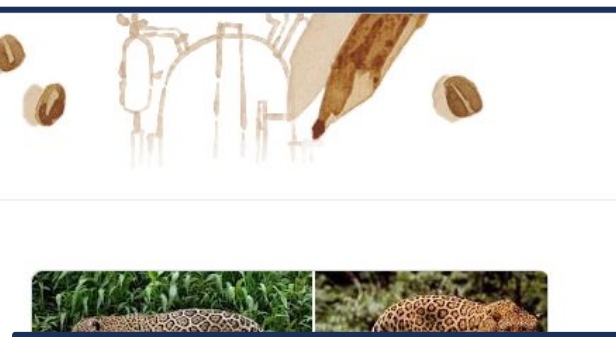
# Outline

---

- The Multimedia Information Retrieval. Why do we care?
- A brief look into textual retrieval (and his limitations)
- Image Retrieval
- Image Representations
- Deep Learning to obtain powerful representations
- Text-to-image retrieval

# Motivation

Google search results for "jaguar animal". The top result is a Wikipedia page titled "Jaguar - Wikipedia". The snippet reads: "The jaguar (Panthera onca) is a large cat species and the only living member of the genus Panthera native to the Americas. With a body length of up to 1.85 ...". Below the snippet are details: Family: Felidae, Kingdom: Animalia, Genus: Panthera, Order: Carnivora. There are also "Ricerche correlate" (related searches) for "jaguar animal wikipedia", "black panther animal", "black jaguar", "pantanal jaguar", "leopard animal", and "animals". A video section shows "Jaguar: The True King of the Jungle" and "Jaguar facts: and How They Compare to Leopards | Animal ...".



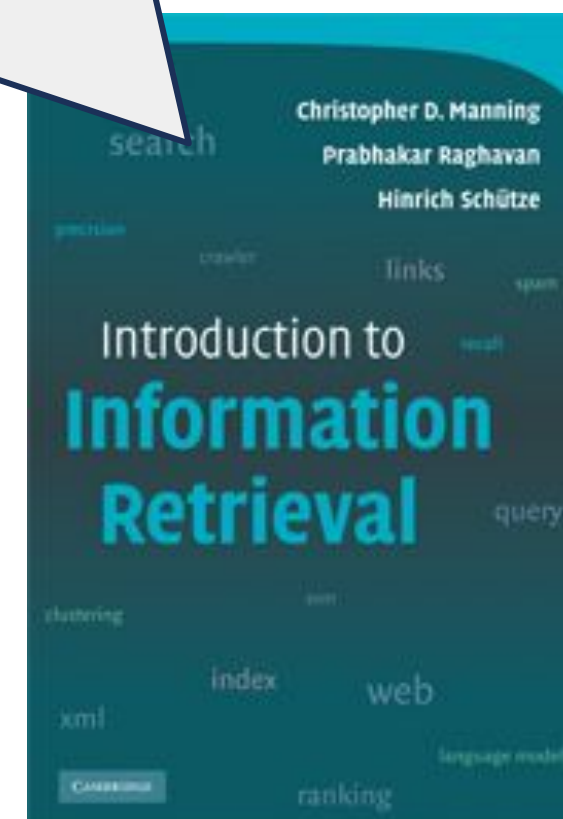
Screenshot of the Treccani website. The search bar contains "DANTE ALIGHIERI". The search results show "1169 RISULTATI" and "Tutti i risultati". The main entry is for "Dante Alighieri" under the category "ENCICLOPEDIA ON LINE". The text describes him as a poet from Florence, born in 1265 and died in 1321. It mentions his family, his father Alighiero, and his work "The Divine Comedy". A small image of Dante in a red robe is shown. Below the main entry, there are sections for "Dante" and "Alighieri, Dante" with further details and a "Mostra tutto" button.

Screenshot of a YouTube search for "cooking tutorial". The search results show several videos, including "Gordon Ramsay's Top Basic Cooking Skills | Ultimate Cookery Course FULL EPISODE", "20 Recipes You Should Learn In Your 20s • Tasty", and "Quick & Easy Recipes With Gordon Ramsay". The interface includes navigation menus, filters, and a list of video thumbnails with titles and durations.

Screenshot of a Google search for "the original hamlet manuscript". The search results show several images of historical manuscripts, including "Hamlet - Shakespeare in q...", "Hamlet Q1 - Wikipedia", "File:Hamlet, Shakespeare, 1...", "Shakespeare, Hamlet", and "File:Hamlet, Shakespeare, 1...". The interface includes navigation menus, filters, and a list of image thumbnails with titles and sources.

# Information Retrieval

“Information retrieval (IR) is finding material (usually documents) of an unstructured nature that satisfies an information need from within large collections (usually stored on computers).”



<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>

# Multimedia Information Retrieval

## *Multimedia* Information Retrieval

- *multimedia*: “two or more different media” and refers to different *modes* of information consumption
  - listening, seeing, reading, watching, smelling etc
- *multimedia information retrieval*: we want the query and the retrieved documents to have possibly different modalities

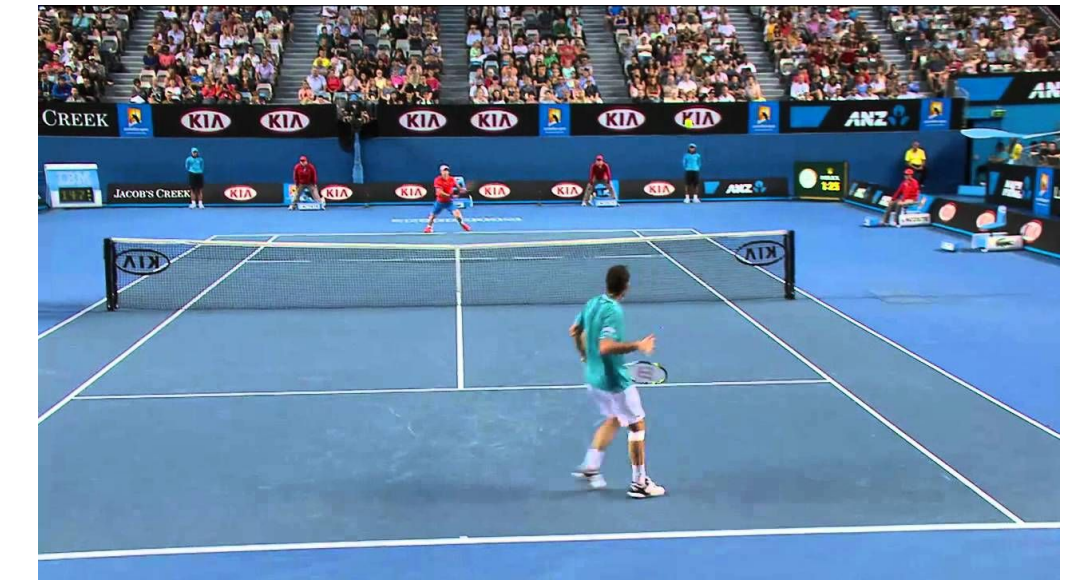


Milton Keynes's Peace Pagoda

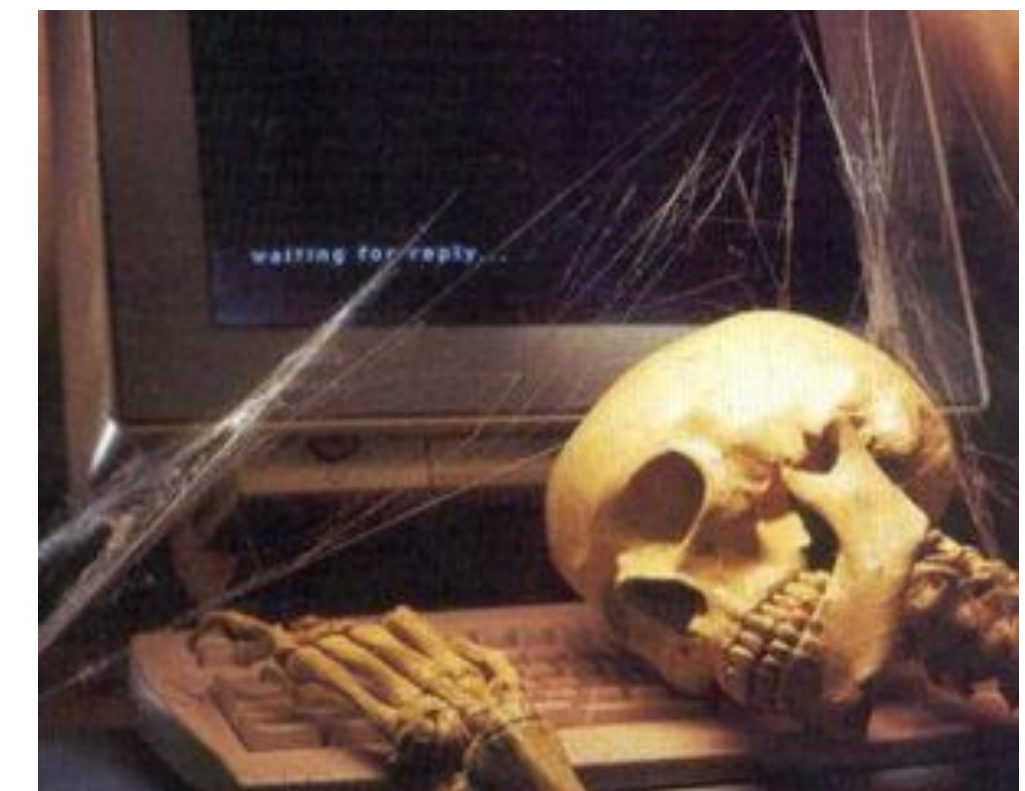
Built by the monks and nuns of the Nipponzan Myohoji, this was the first Peace Pagoda...

# Challenges in Multimedia Information Retrieval

- The content to retrieve must be encoded in some way
  - Exploit metadata
  - Directly exploit the content (Content-Based Information Retrieval)
- Large scale scenarios
  - Billions of items to be retrieved in few milliseconds
  - E.g.: Google

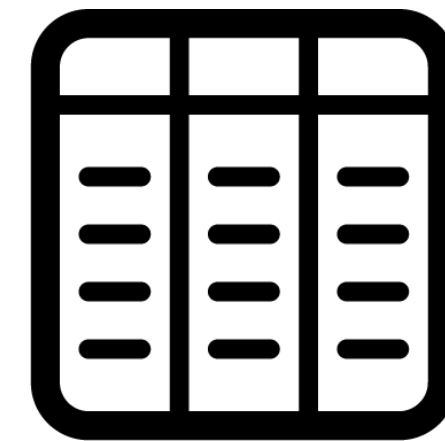
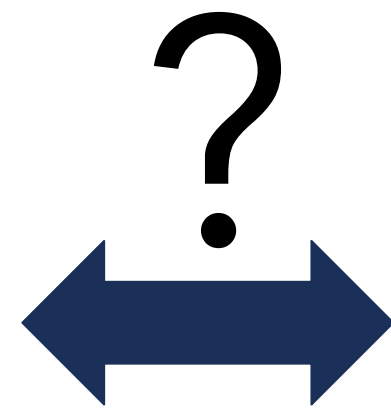
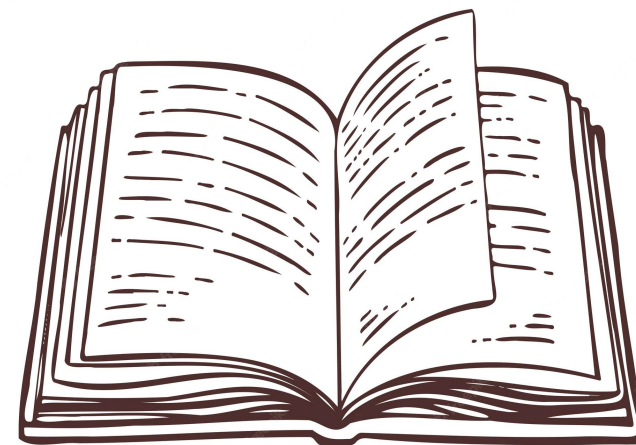
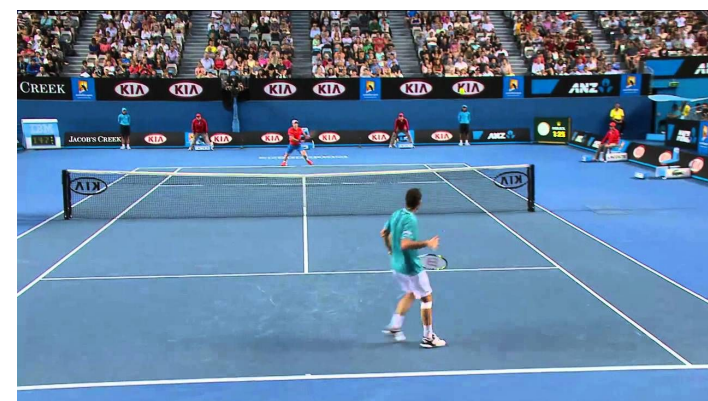


filename: "tennis.jpg"  
size: 1280x720  
timestamp: 22 gen 2012



# Challenges in Multimedia Information Retrieval

- Multimedia data (e.g. images, videos) are not structured data
  - Not possible to use standard relational databases (e.g. MySQL)
  - An image cannot be directly stored as nice tabular data



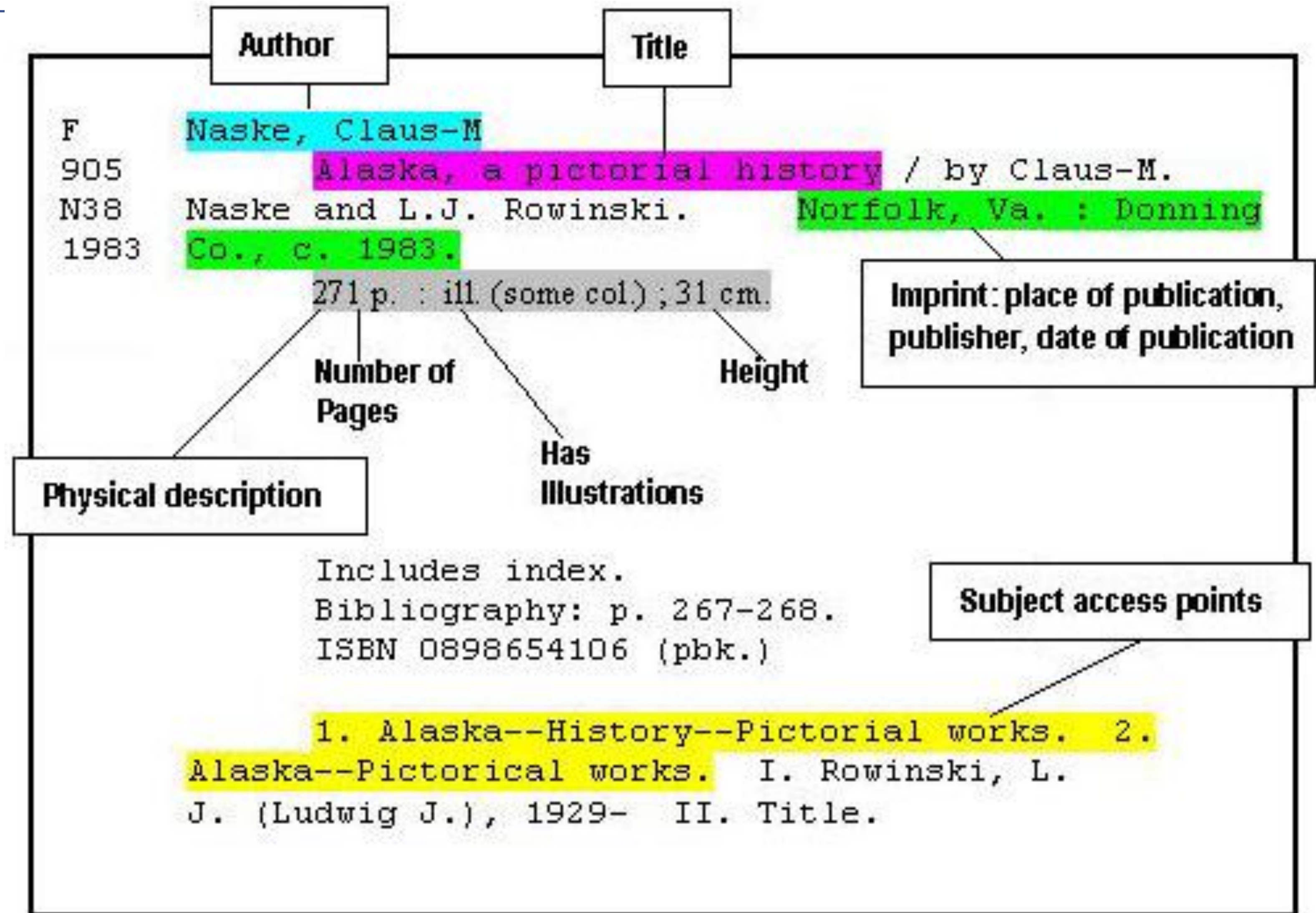
unstructured

structured

Image	date	size	pixels(?)
...			
cat.jpg	15/06/22	300x200	???
...			

# The analog way: Catalog card

A catalog card is an individual entry in a library catalog containing bibliographic information, including author's name, book title, etc...





# The importance of the Content

- . In general, we need to search also the content
- . Google indexes the content, so that we can search for something even in the content, not only in metadata
- . For books (and textual documents in general) this is not too much difficult
  - o Each document is a set of words
  - o Find out how many words match with the query
  - o We will see in few slides...

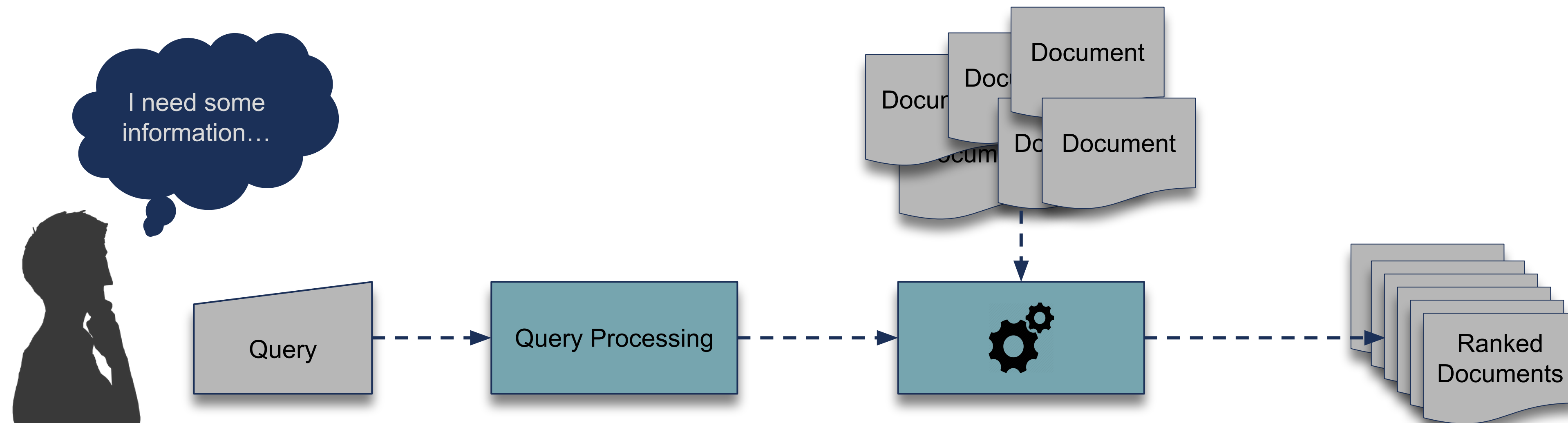
The screenshot shows a Google search for "deep blue scacchi". The search bar contains the query "deep blue scacchi" with "deep blue" highlighted in a red box. Below the search bar, there are navigation options: "Tutti", "Immagini", "Shopping", "Video", "Libri", and "Altro". The "Libri" option is selected. The search results are displayed in two columns. The first result is for the book "Scacchi per Principianti: Fai la Tua Mossa! Impara a Giocare ..." by Fiorenza Marino, published in 2021. The search results for this book show "CONTENUTO TROVATO ALL'INTERNO - PAGINA 10" and a snippet of text: "Infatti, nel 1989, Garry Kasparov sconfisse il computer 'Deep Thought' di IBM in una partita di 6 scontri. Anche la versione successiva, 'Deep Blue' fu affidata a Kasparov nel 1996. Ma nella rivincita del 1997, Deep Blue, ...". The words "Deep Blue" and "Deep Blue" are highlighted in red boxes. Below the snippet are buttons for "Anteprima", "Altre edizioni", and "Paperpile". The second result is for the book "Informatica - Pagina 435" by G. Michael Schneider and Judith L. Gersting, published in 2013. The search results for this book show "CONTENUTO TROVATO ALL'INTERNO - PAGINA 435" and a snippet of text: "Nel maggio 1997 l'attenzione internazionale si concentrò su una storica partita di scacchi tra il campione mondiale Garry Kasparov e il computer di IBM per giocare a scacchi dal nome Deep Blue IBM Blue Gene/L, discusso nel Capitolo 3 ...". The words "Deep Blue" and "Deep Blue" are highlighted in red boxes. Below the snippet are buttons for "Anteprima", "Altre edizioni", and "Paperpile".

# Information Retrieval System

In the more general terms, the main components are

- . A query
- . A database of documents among which we want to search

The output is a subset of documents, the ones relevant to the given query, ranked by decreasing relevance



# Textual retrieval

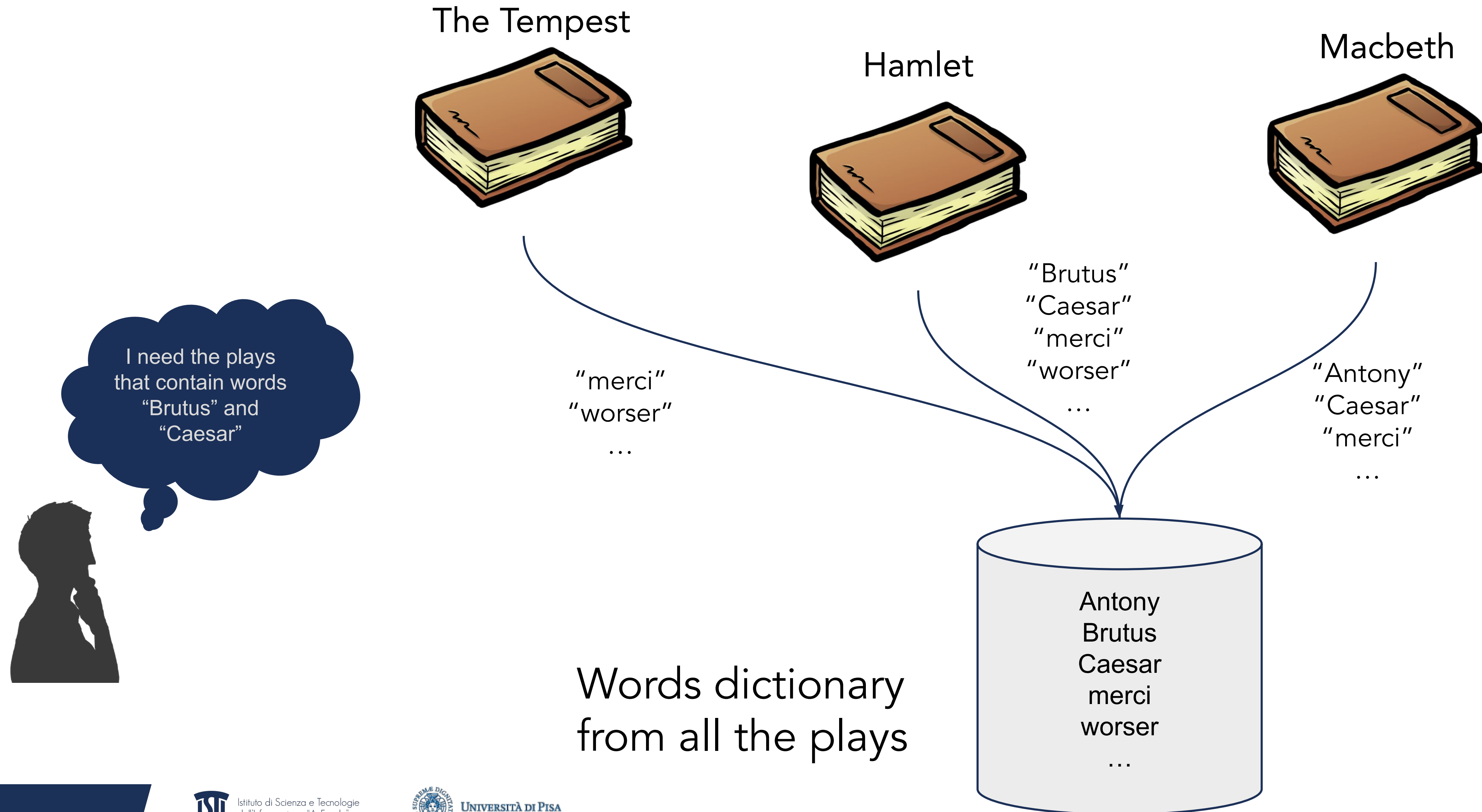
Historically, documents were only textual

- . 1950s: textual document retrieval
- . 1980s: multimedia documents acquired interest
  - Difficulty of processing “non-textual documents”
  - Medium mismatch problem, or semantic gap
    - E.g.: How to match an image with a text that describes it?

There is ... a machine called the Univac ... whereby letters and figures are coded as a pattern of magnetic spots on a long steel tape. By this means the text of a document, preceded by its subject code symbol, can be recorded ... the machine ... automatically selects and types out those references which have been coded in any desired way at a rate of 120 words a minute

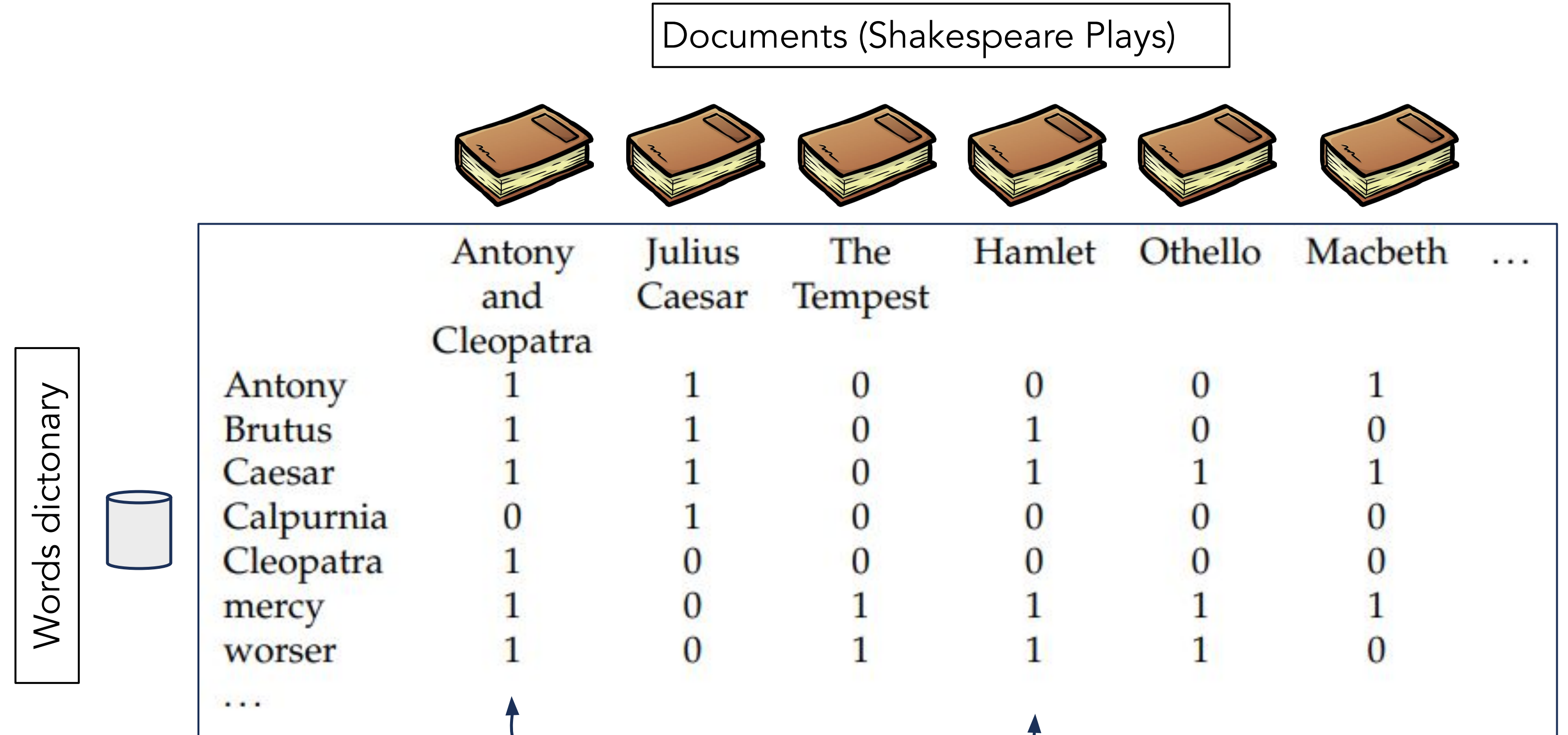
— J. E. Holmstrom, 1948

# A simple retrieval model for texts



# The simple boolean model

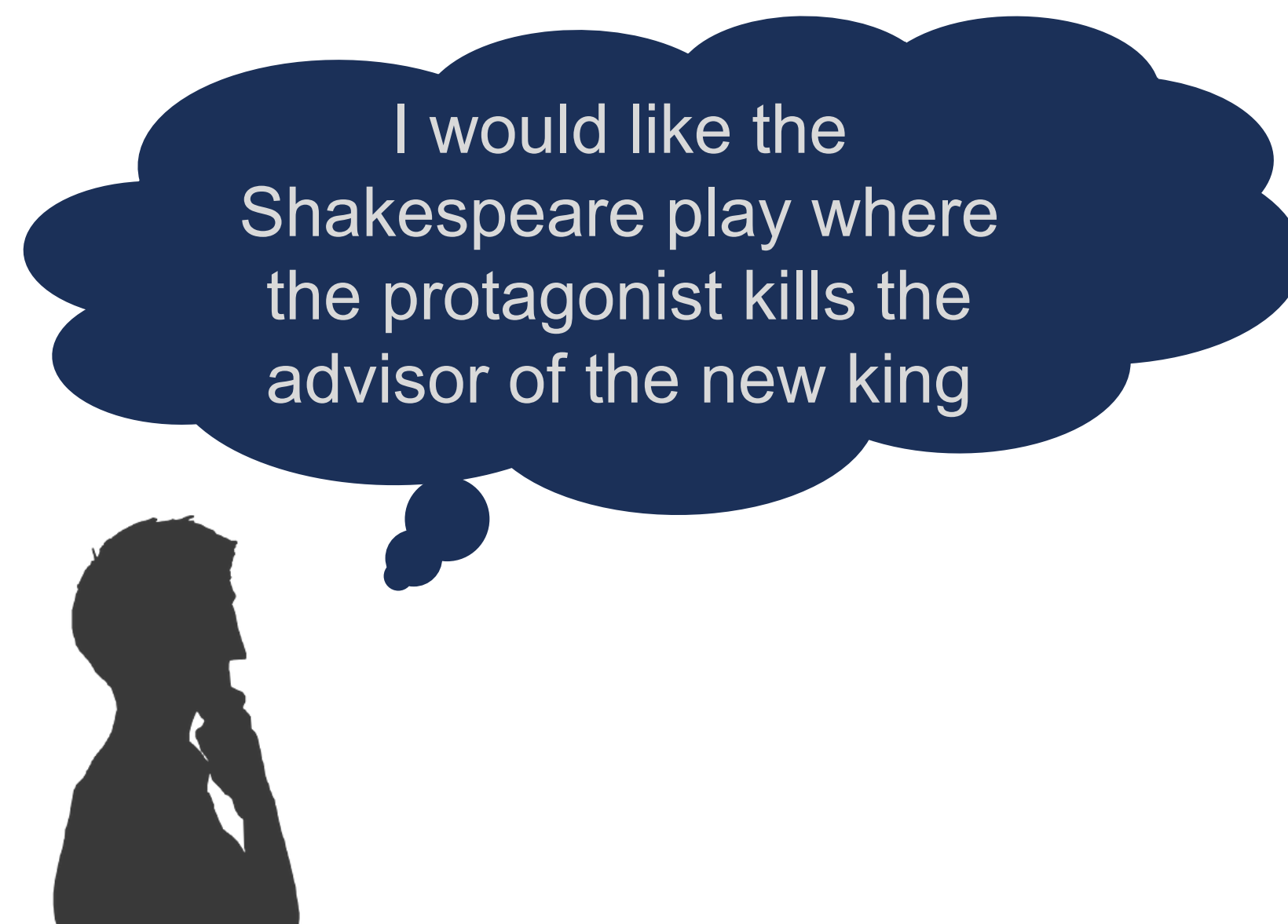
- We create an incidence matrix that tells us which word appears in each of the documents



- To answer the query Brutus AND Caesar AND NOT Calpurnia
  - We take the vectors for Brutus, Caesar and Calpurnia, complement the last, and then do a bitwise AND:  $110100 \text{ AND } 110111 \text{ AND } 101111 = 100100$
  - The answers for this query are thus Antony and Cleopatra and Hamlet

# Can be improved but...

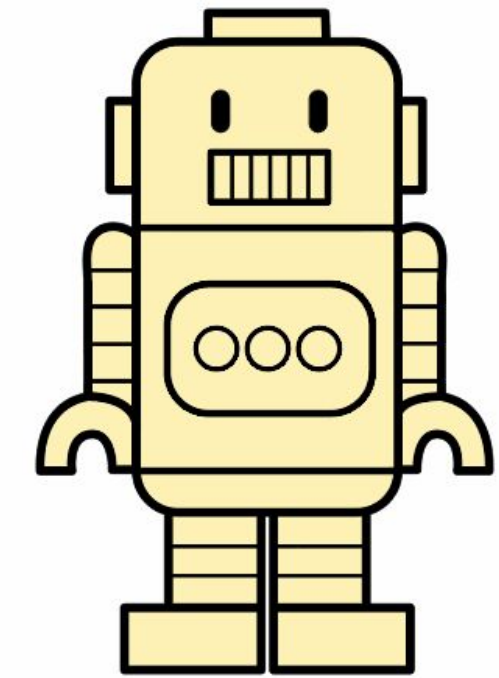
- This model could be improved by leveraging word occurrences
  - the more frequent a word, the more important it is
  - a word appearing too much in a corpus is not too significative (e.g, "the")
- However, what about a query like:



# Problems

---

- The system works at a very low level, very “robotic”
  - Only exact word matching
  - What about synonyms? And what about the context?
- This method does not scale to other media objects
  - Internet is full of less structured media objects that need to be stored and efficiently retrieved
    - Images, videos, audio
    - We cannot search these unstructured objects using exact match



# How to handle images, videos, ...?

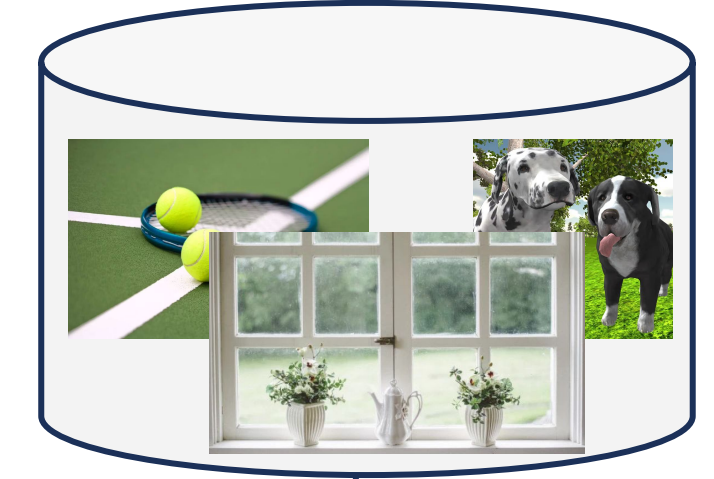
- Frame the problem as a text retrieval problem
  - For images, in principle we could rely on some metadata, e.g. the alternative text associated to the image HTML "alt" tag
  - For videos, use to the video description provided by the user that uploaded it



```
<span id="imageBlockEDPOverlay"></span>
<li class="image" item itemNo0 maintain-height
selected" style="cursor: pointer;">
  <span class="a-list-item">
    <span class="a-declarative" data-action="main-image-click" data-
main-image-click="{}">
      <div id="imgTagWrapperId" class="imgTagWrapper" style="height:
500px;">
        <img alt="Doritos Tortilla Chips, Nacho Cheese, 1.75-Ounce
Large Single Serve Bags (Pack of 64)" src="https://images-
na.ssl-images-amazon.com/images/I/
71Br1LeeJGL_SX486_SX679_P1bundle-
64_TopRight_0_0_SX486_SX679_CR_0_0_486_679_SH20.jpg" data-
old-hires="https://images-na.ssl-images-amazon.com/images/I/
71Br1LeeJGL_SL1366.jpg" class="a-dynamic-image a-stretch-
```



Extract text  
*"a girl playing tennis"*



Extract texts

back to text retrieval

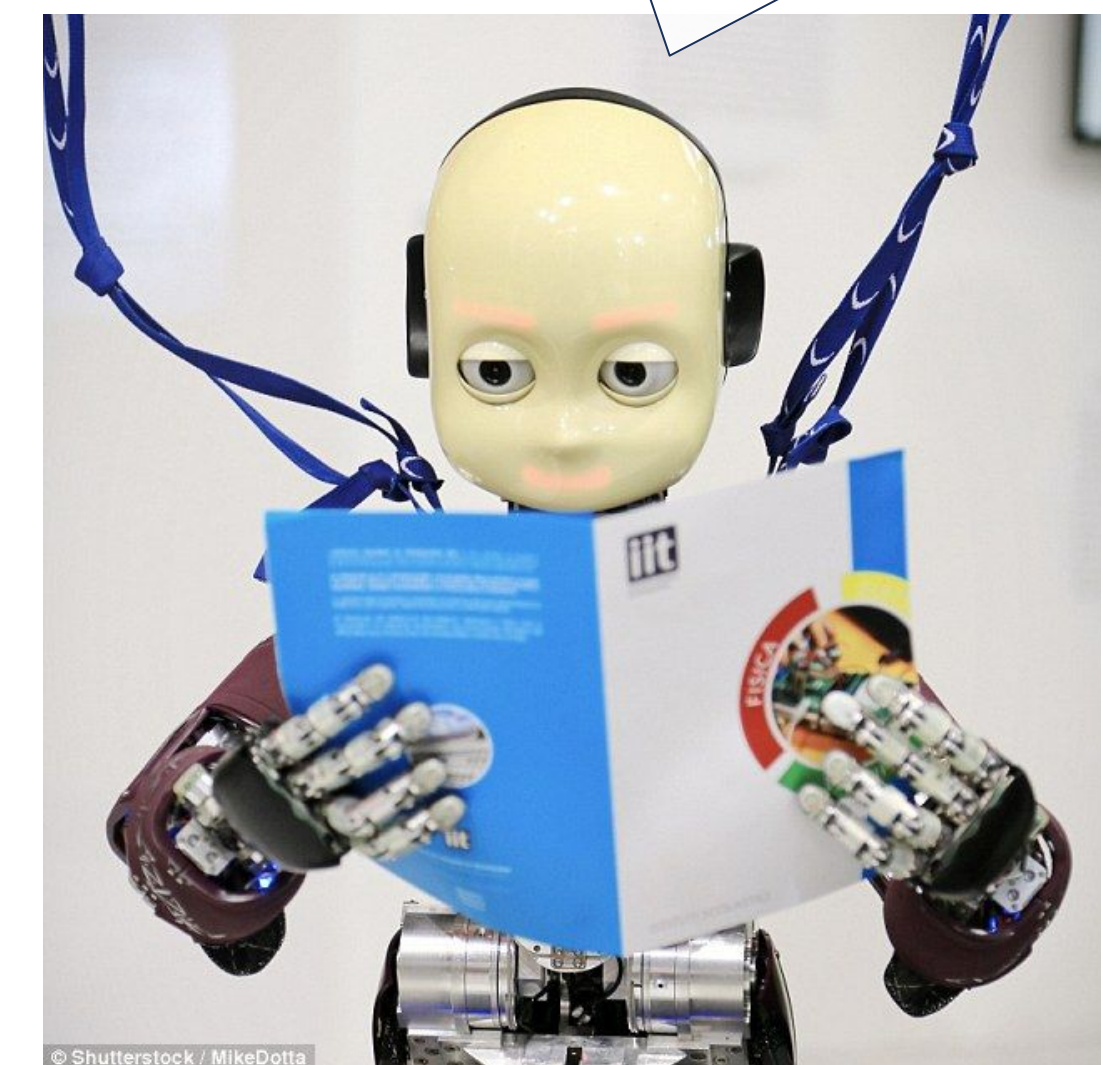
	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
Antony	1	1	0	0	0	1	
Brutus	1	1	0	1	0	0	
Caesar	1	1	0	1	1	1	
Calpurnia	0	1	0	0	0	0	
Cleopatra	1	0	0	0	0	0	
mercy	1	0	1	1	1	1	
worser	1	0	1	1	1	0	
...							



# How to handle images, videos, ...?

- But...
  - Often these data is not available, or contain errors
  - Not representative of the whole multimedia element
- We often need to rely only on the content

Looking at the content → high level understanding  
→ certain degree of intelligence



# Image Retrieval

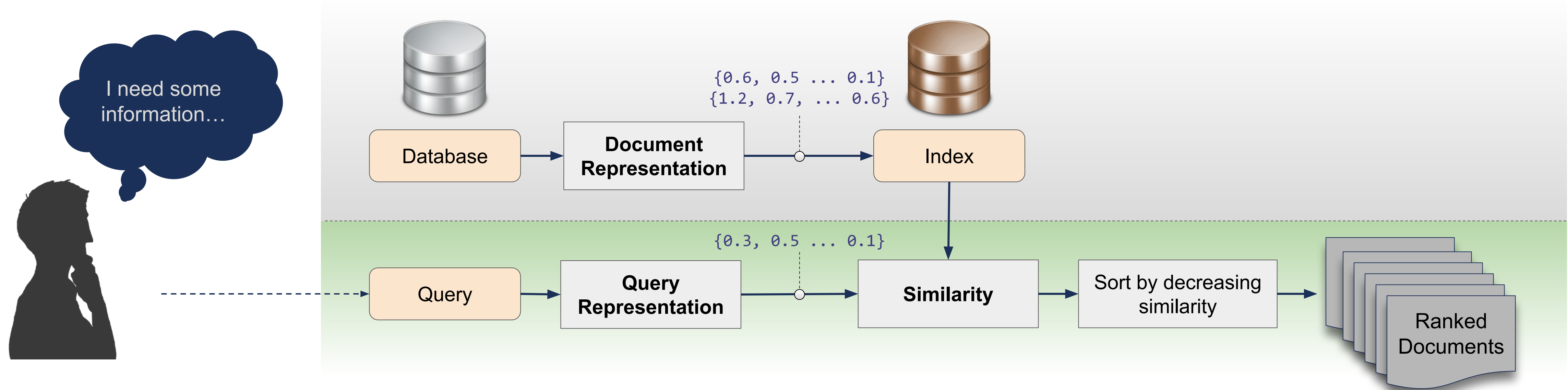
- Find images similar to the one given as a query
- Query by example



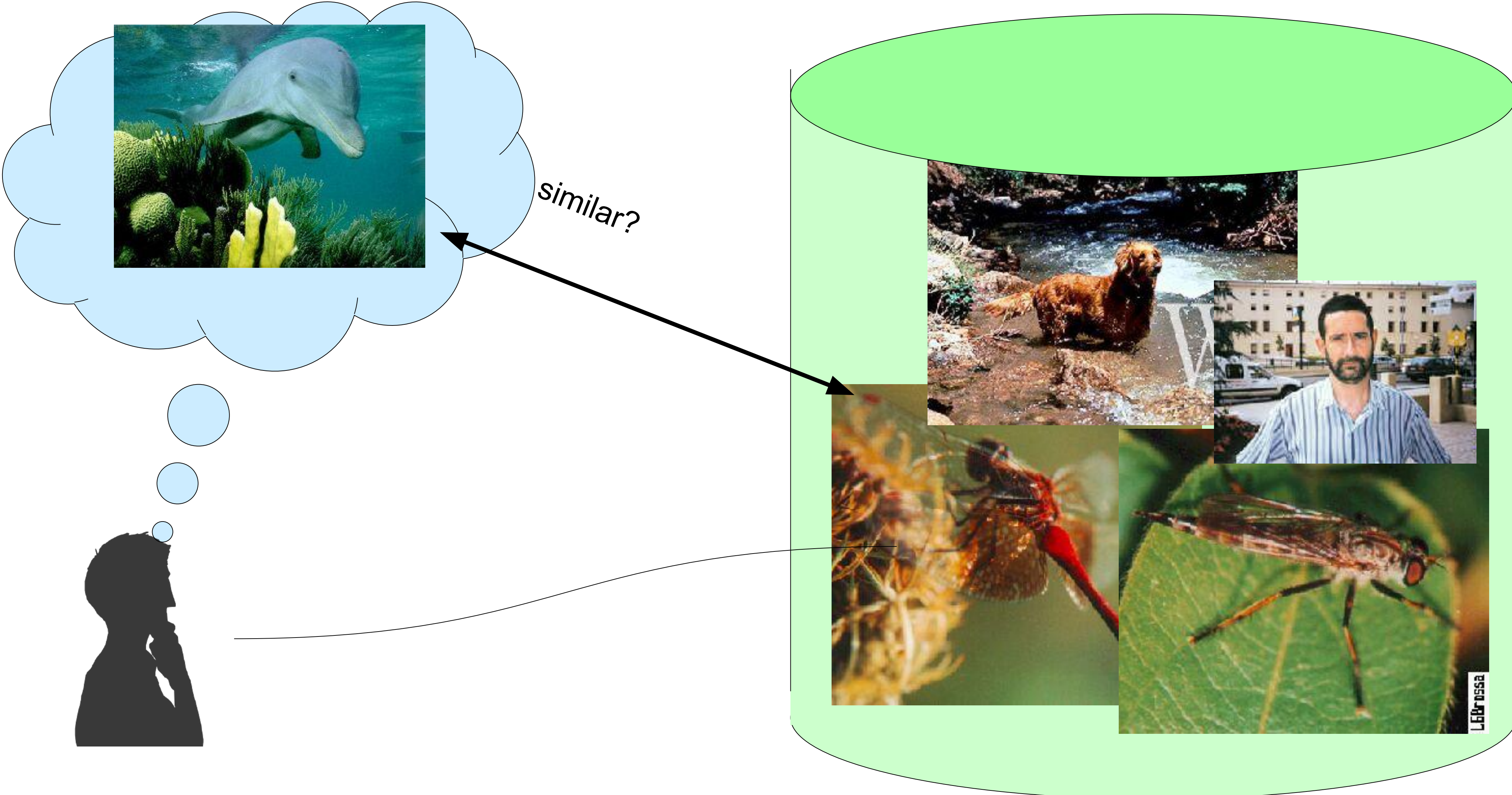
- We cannot define a precise matching criterion between images, but...
- We can quantify in some way their similarities

# Paradigm shift: similarity search

- Encode a media object (a text, an image, or a video) into some numerical representation (or feature)
- Measure the similarity between the representation of the query and the representation of the documents in the database

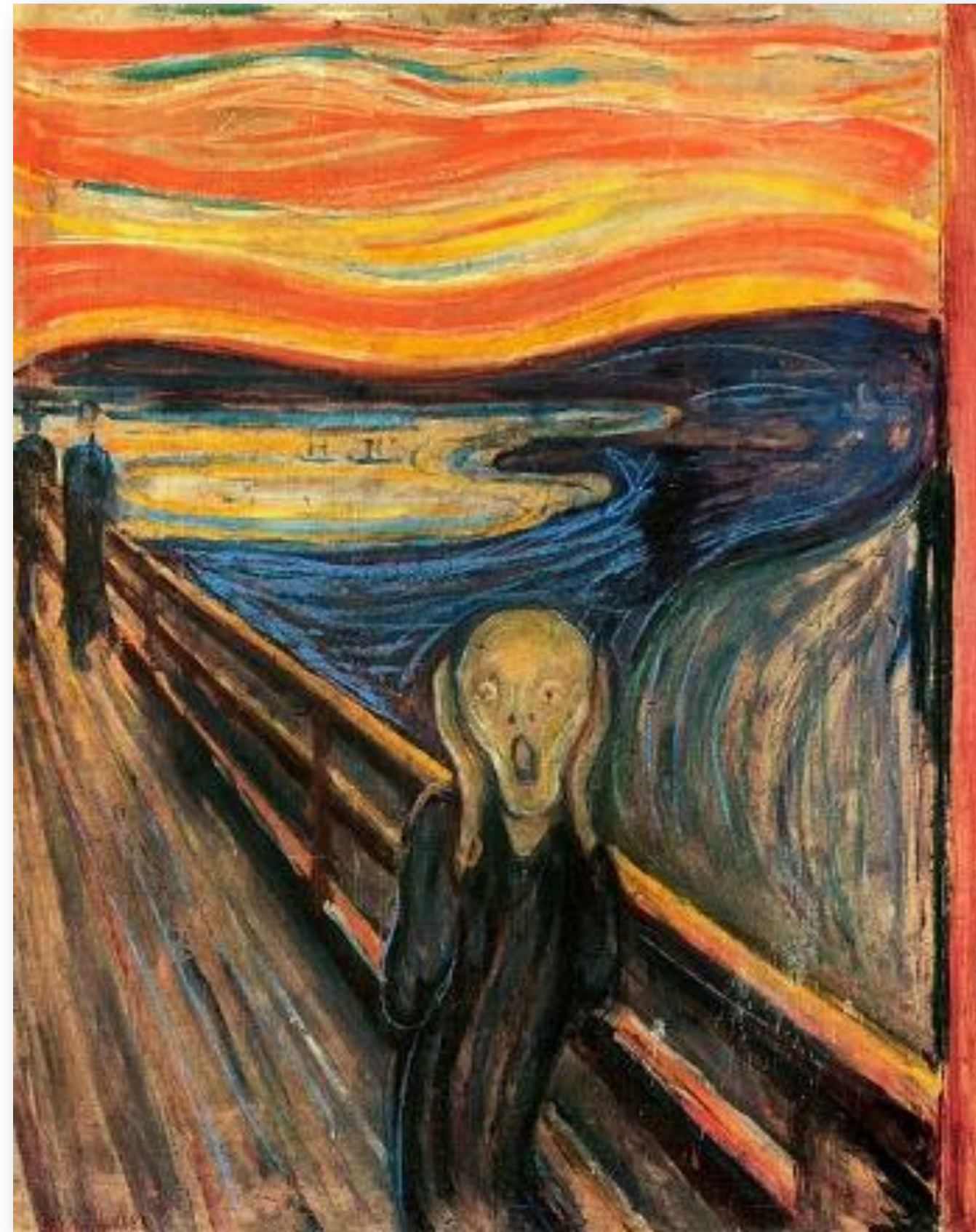


# Image Retrieval: what does “similar” mean?



# What's This?

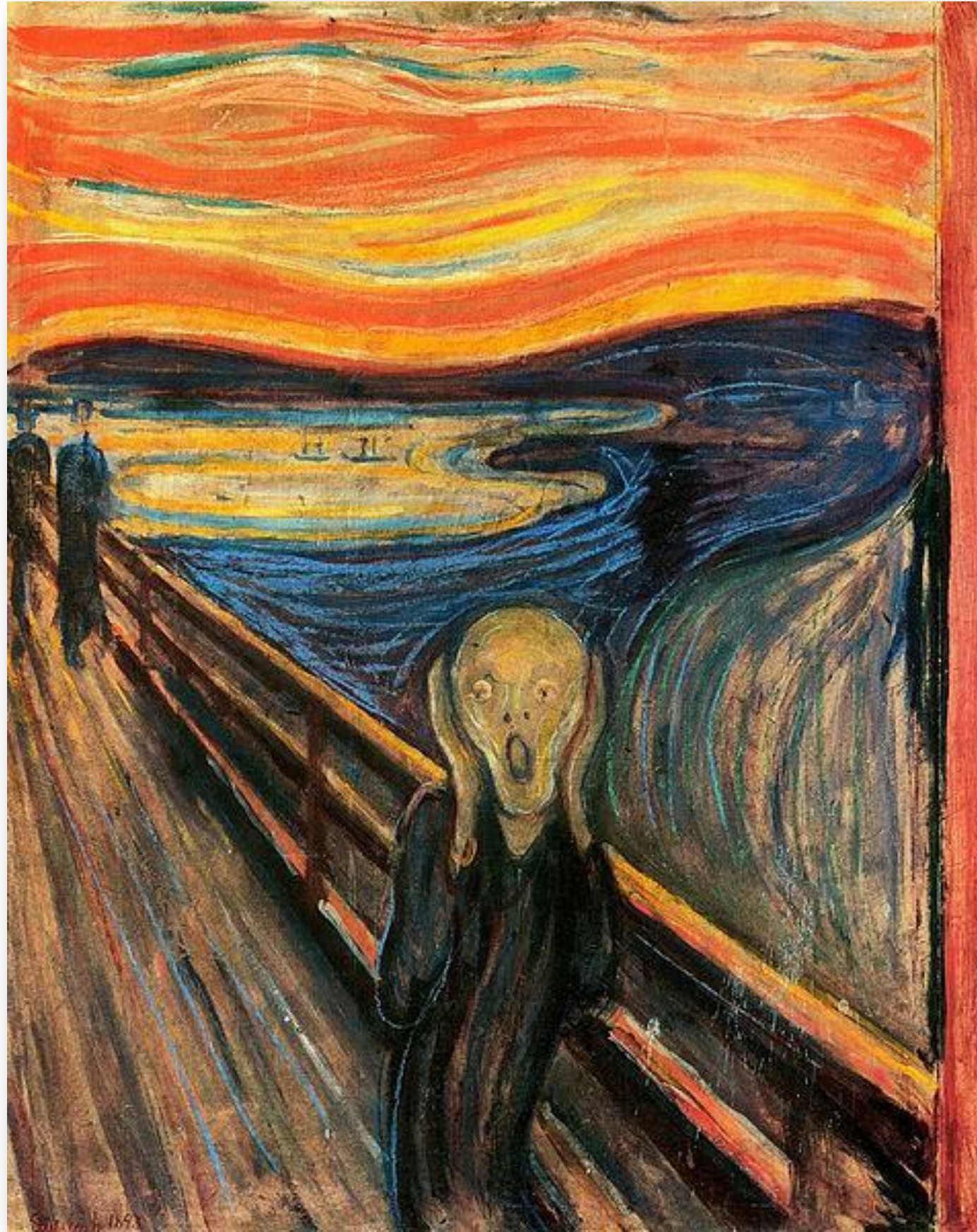
---



What's this?

# The Scream, Edvard Munch

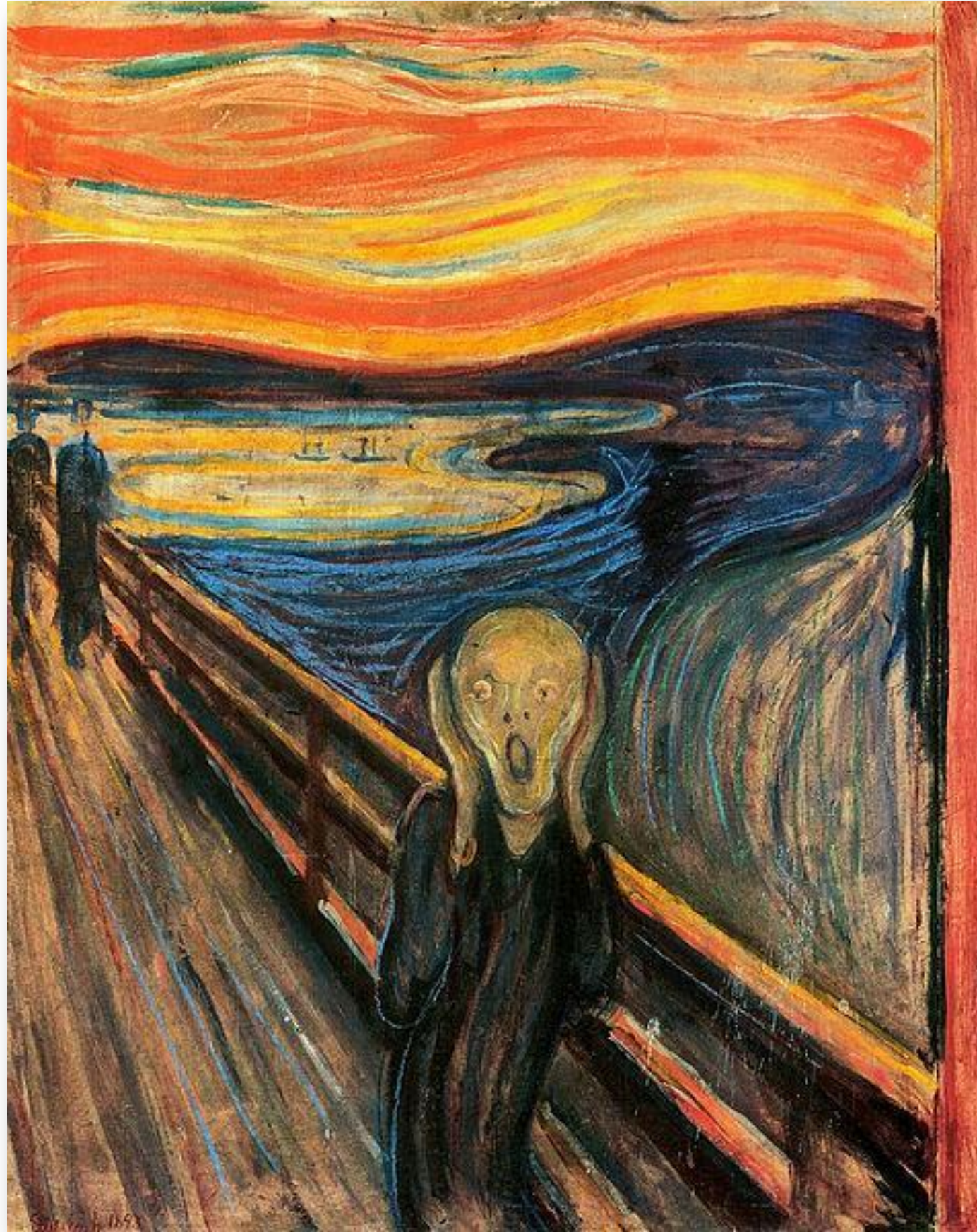
---



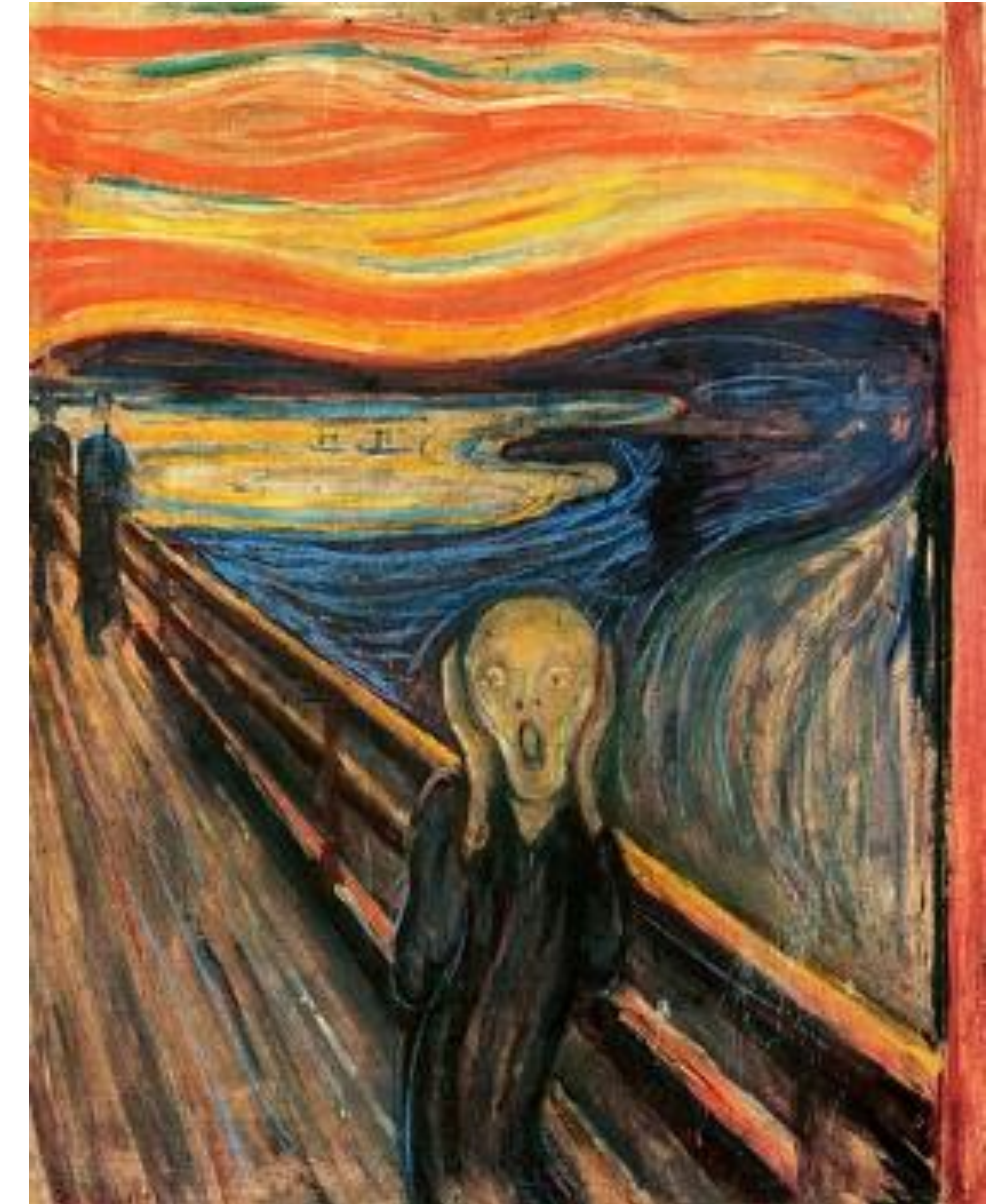
- The file at

[http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)

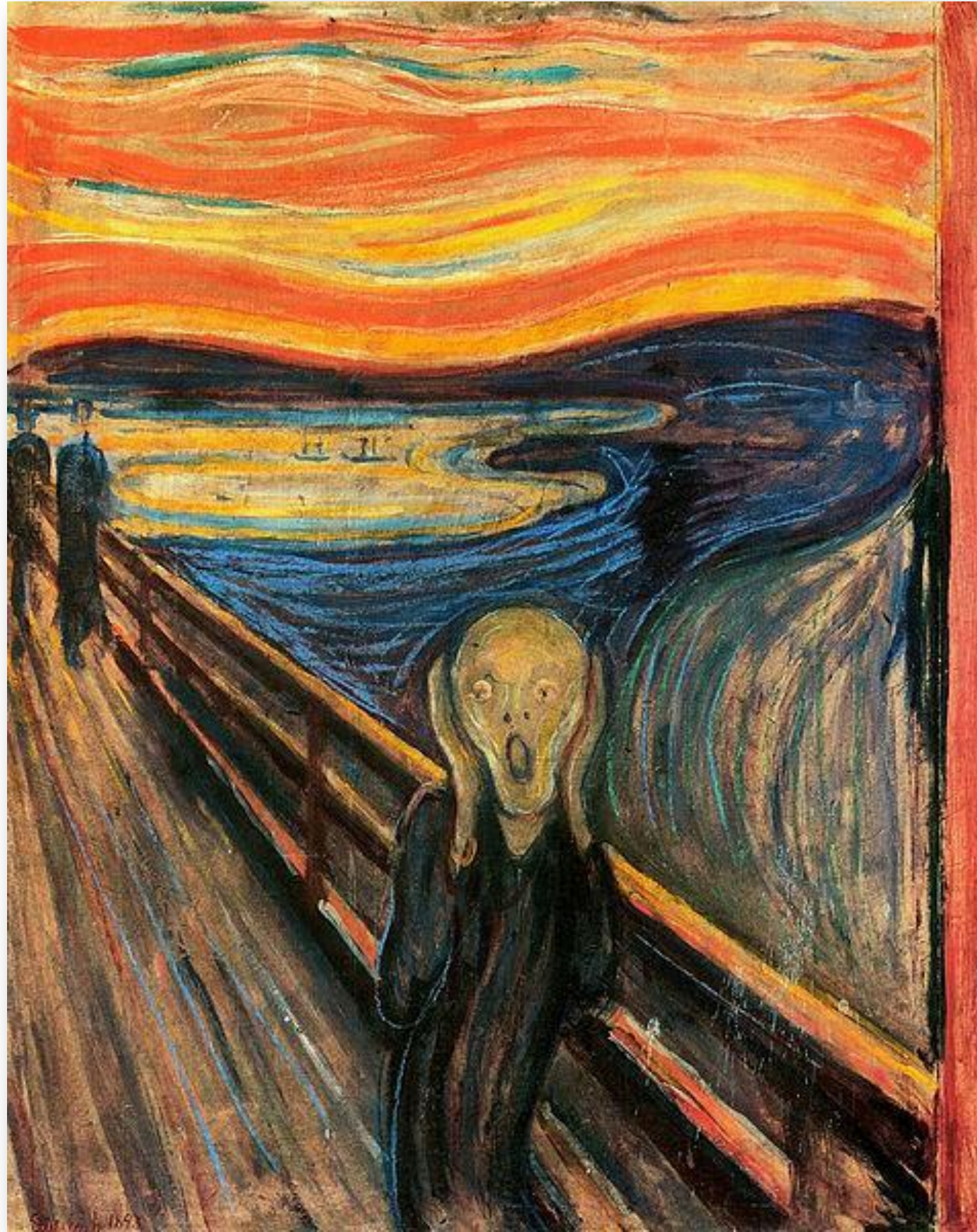
# The Scream, Edvard Munch



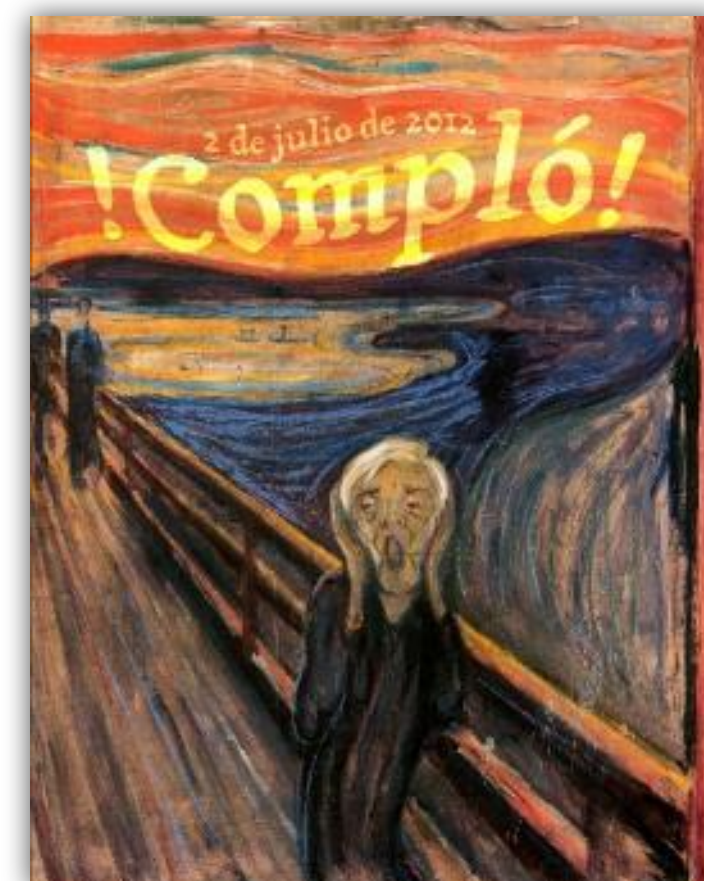
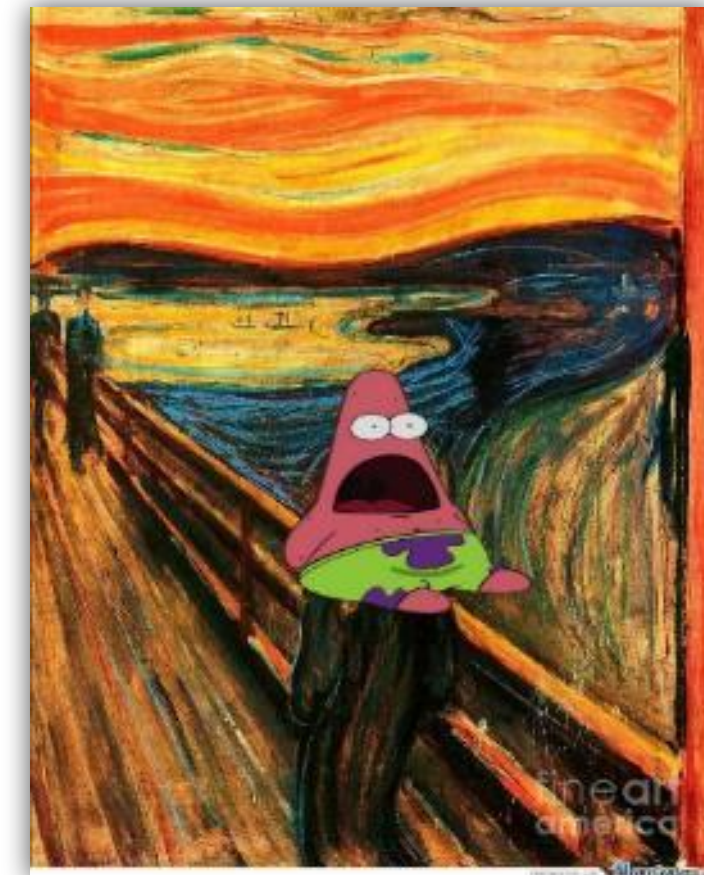
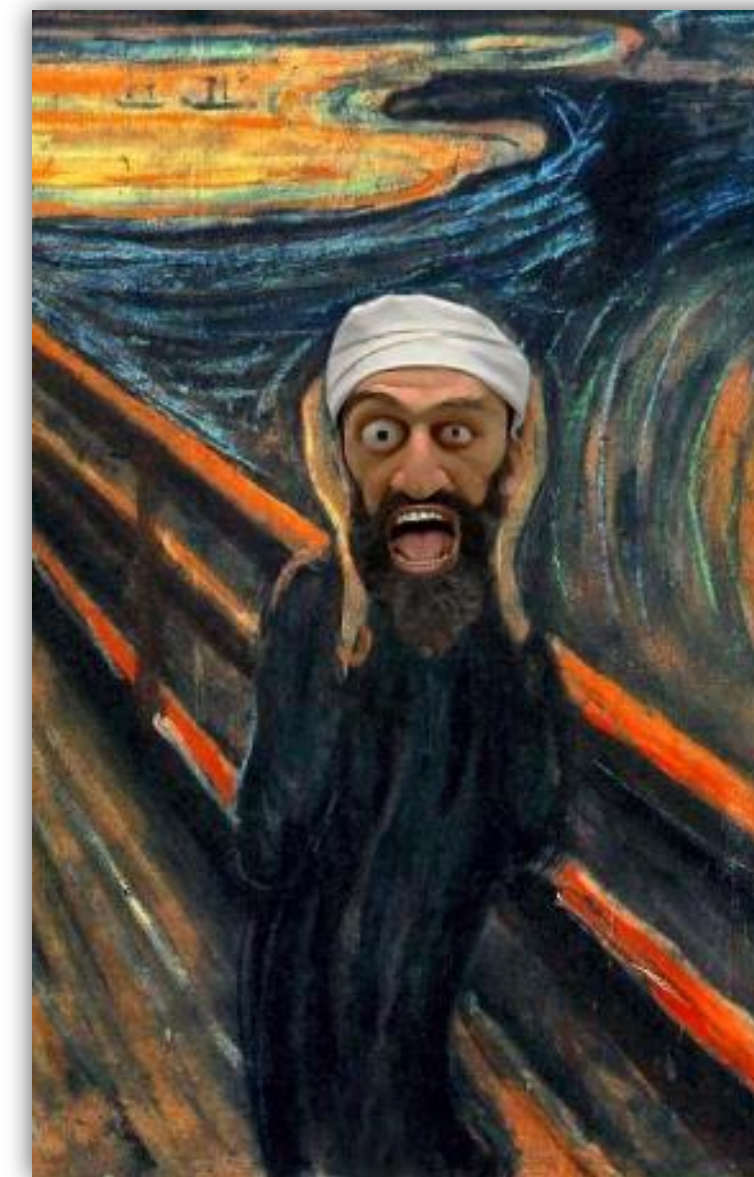
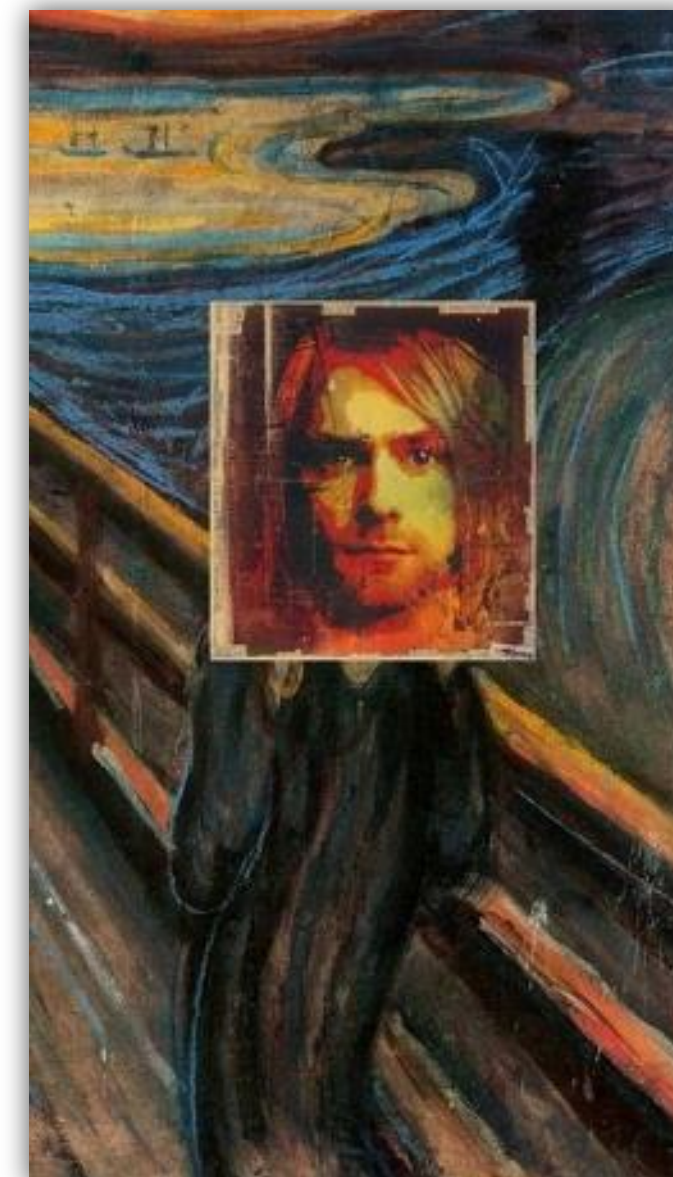
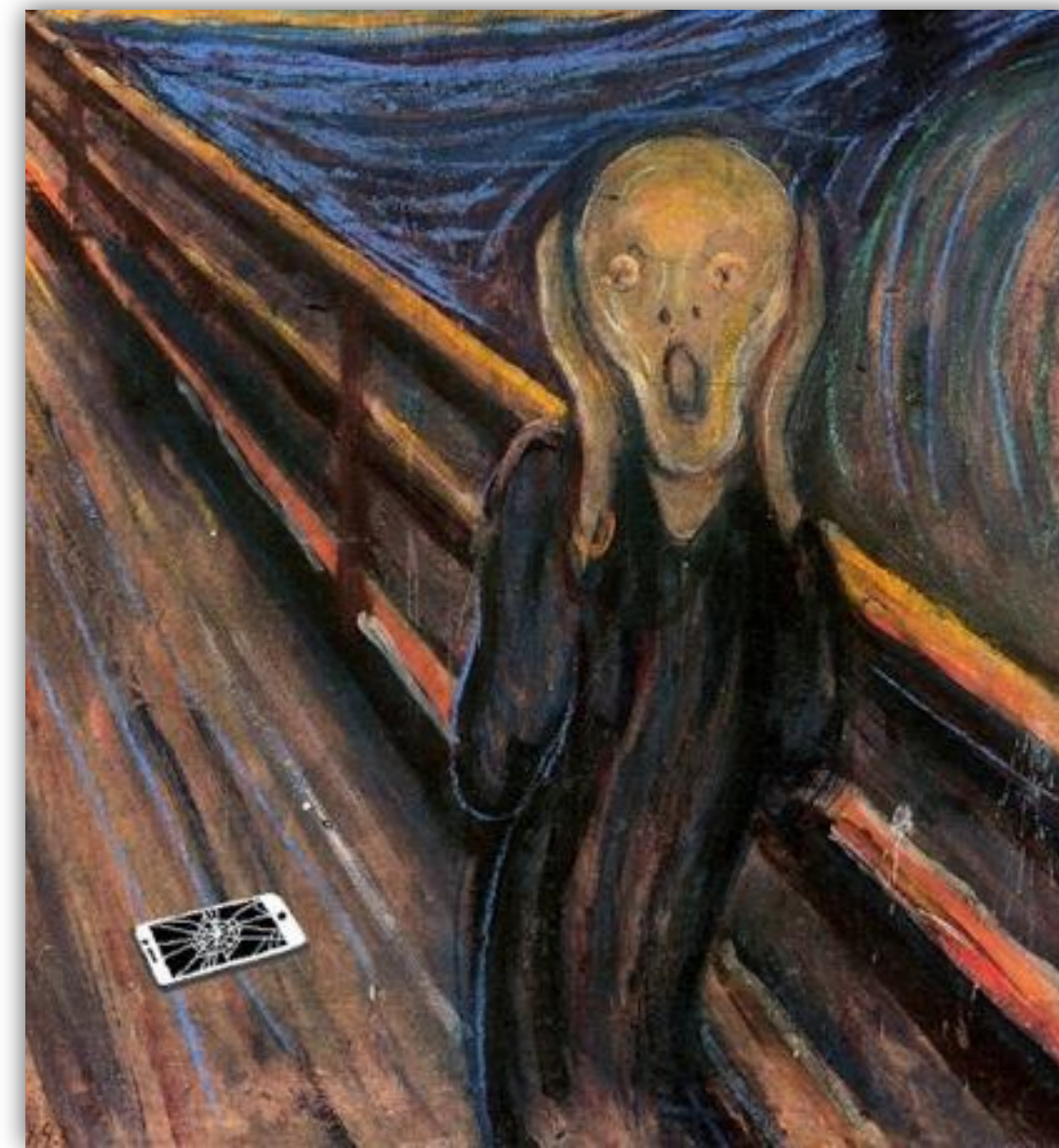
- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture **as**



# The Scream, Edvard Munch

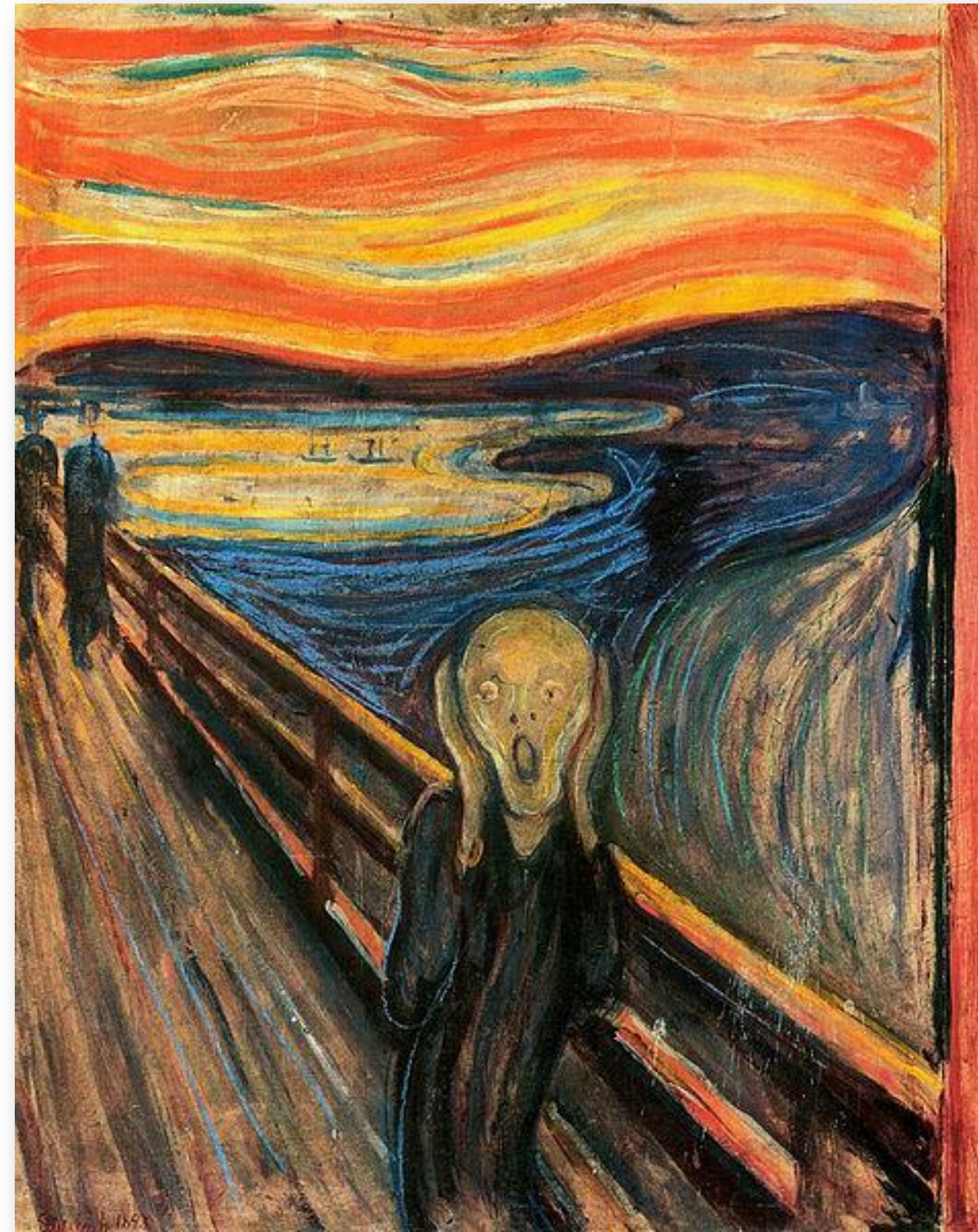


- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same **as**

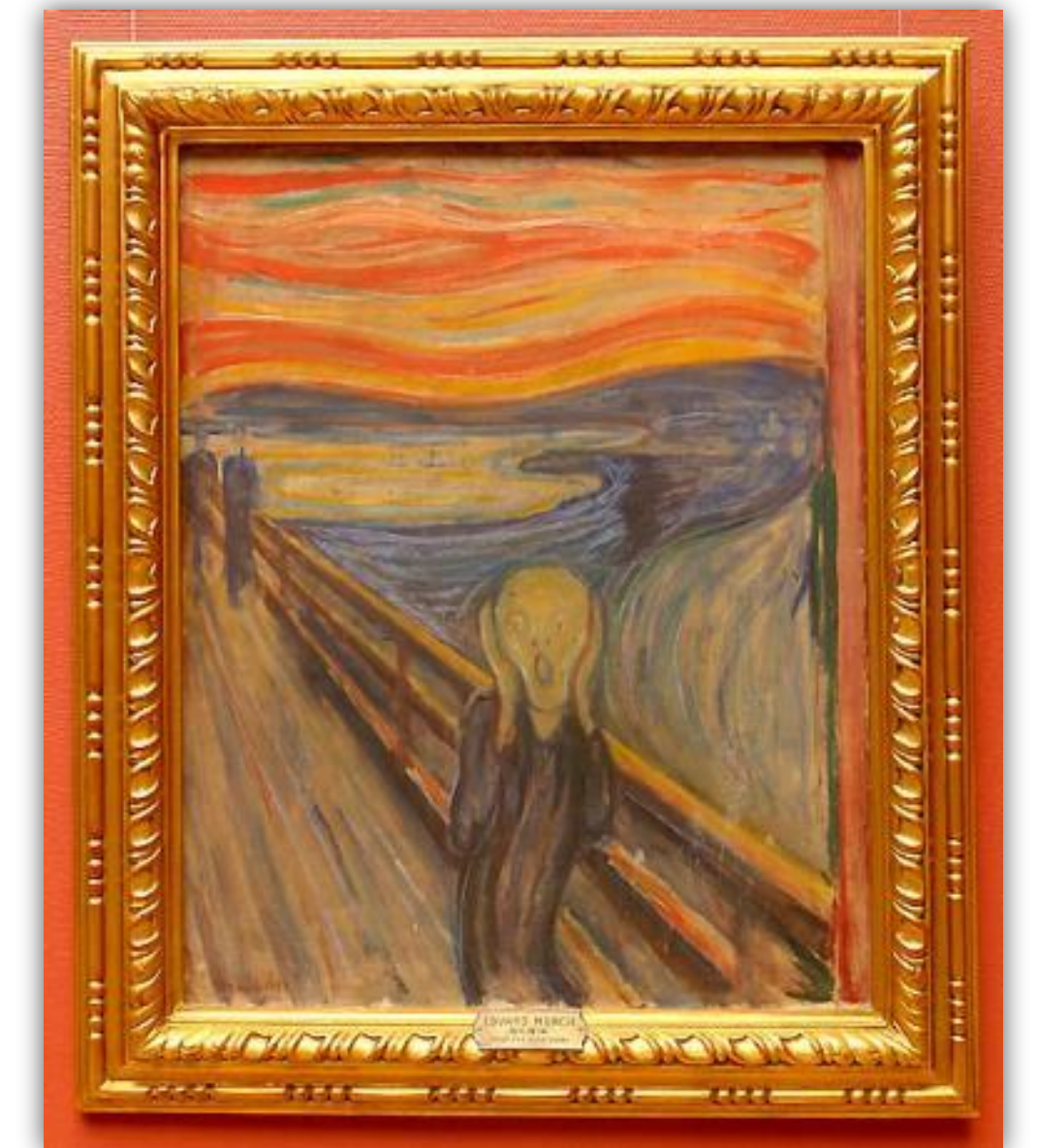




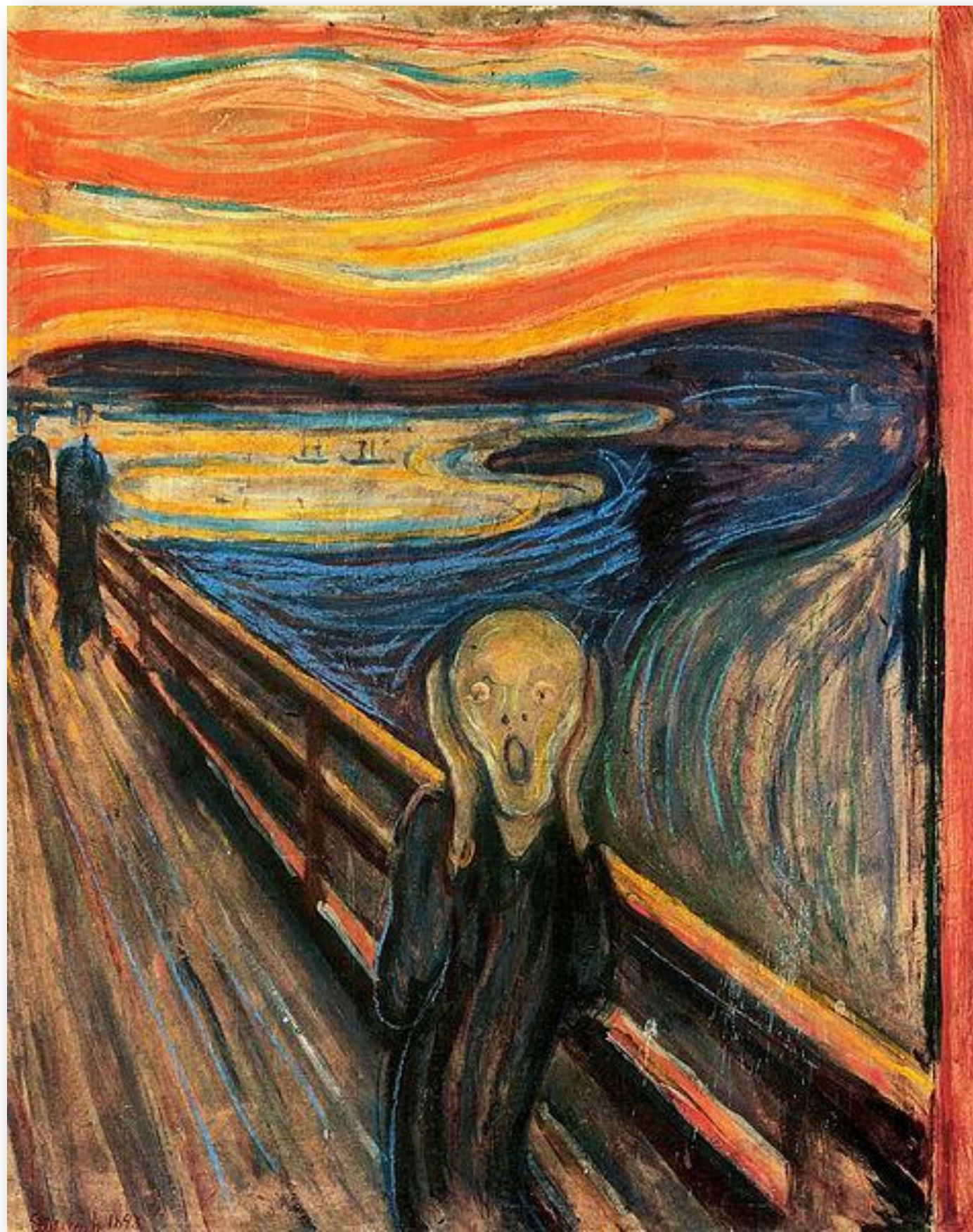
# The Scream, Edvard Munch



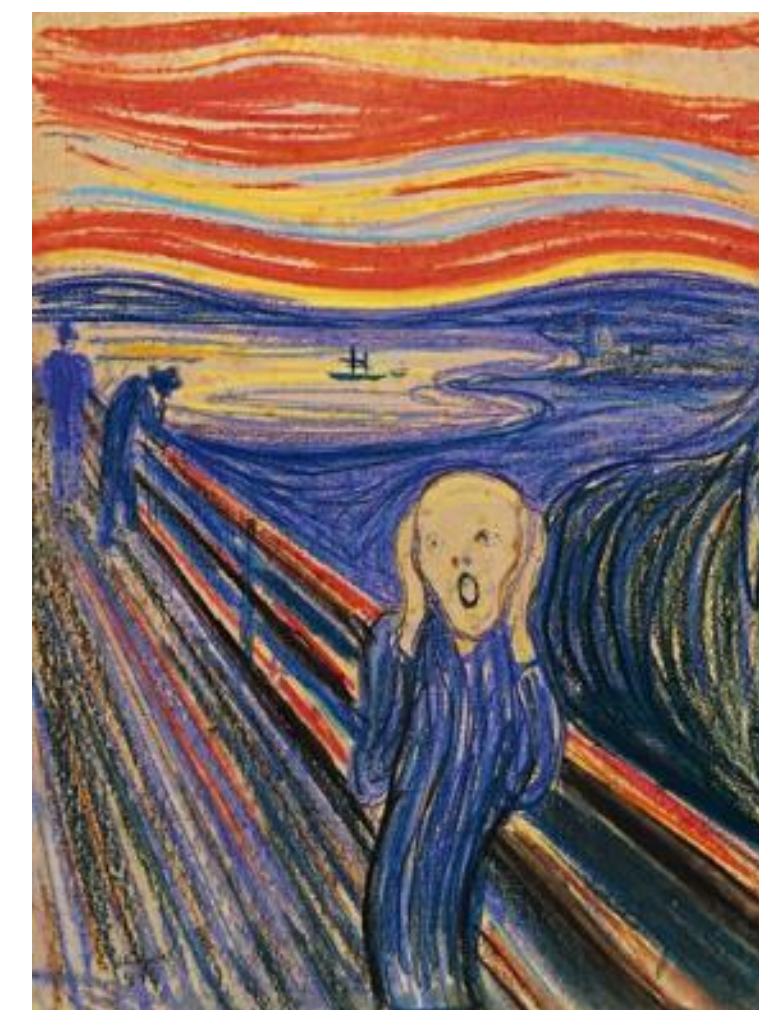
- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo **as**



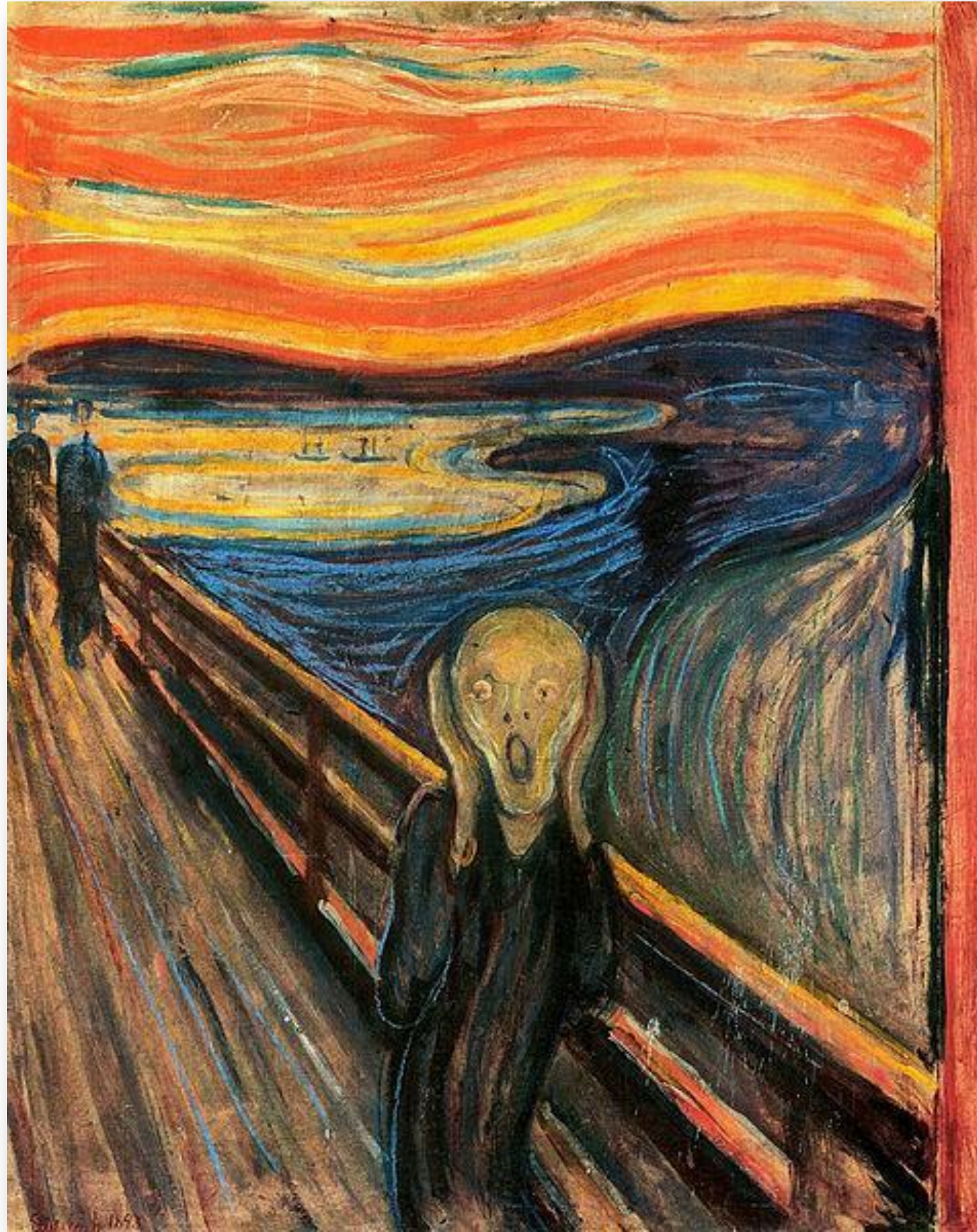
# The Scream, Edvard Munch



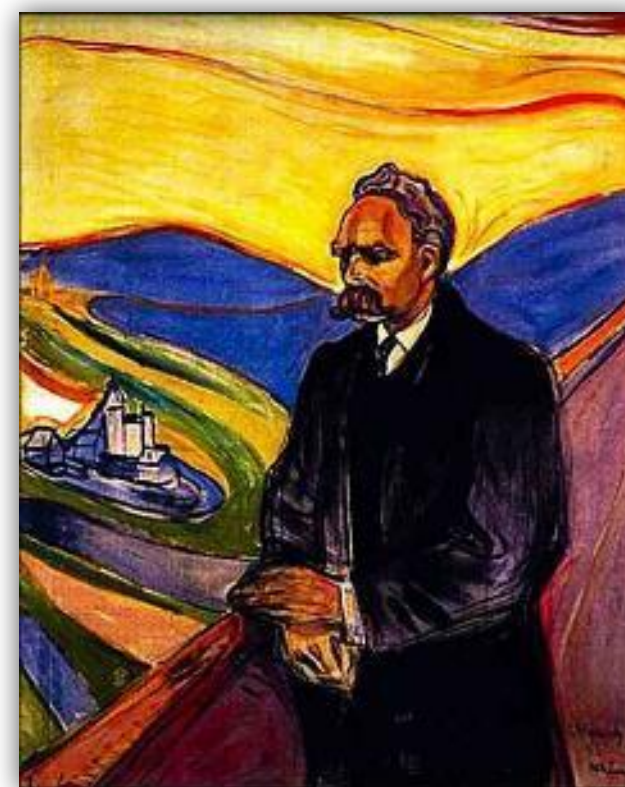
- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream" by Edvard Munch **as**



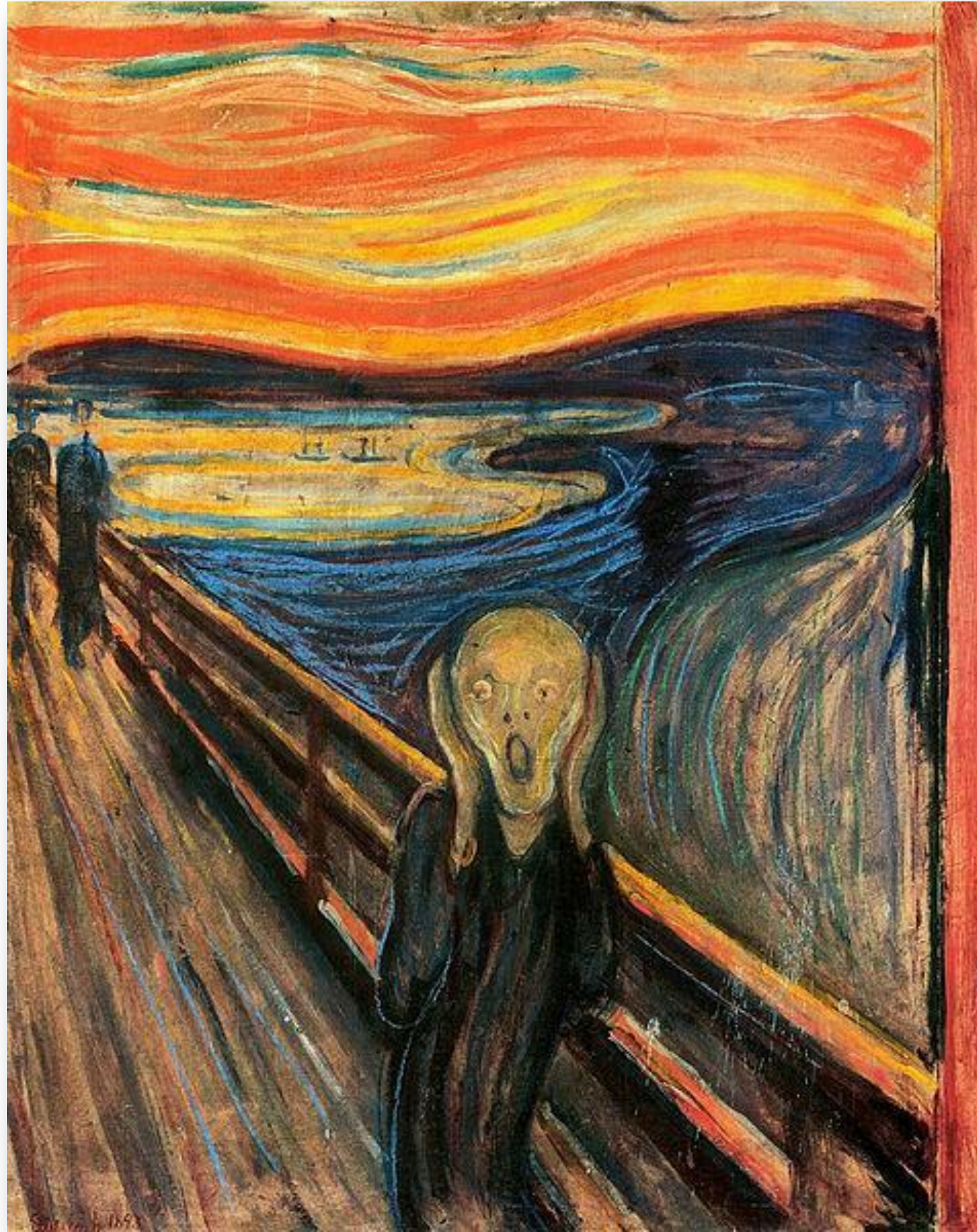
# The Scream, Edvard Munch



- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch **as**

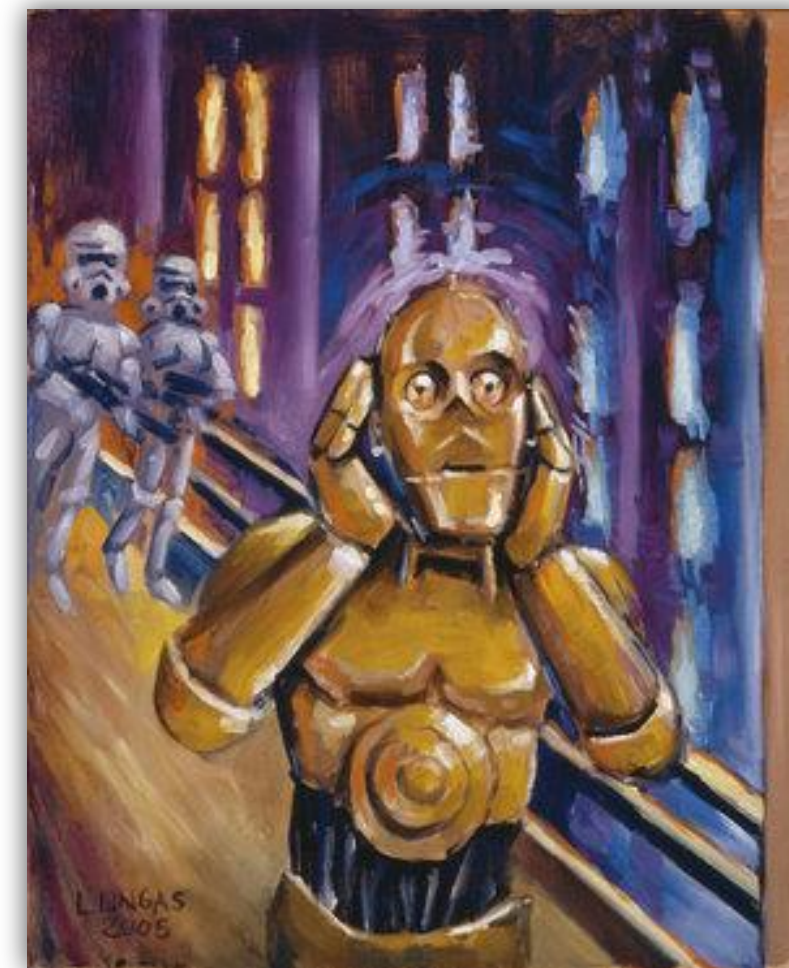


# The Scream, Edvard Munch

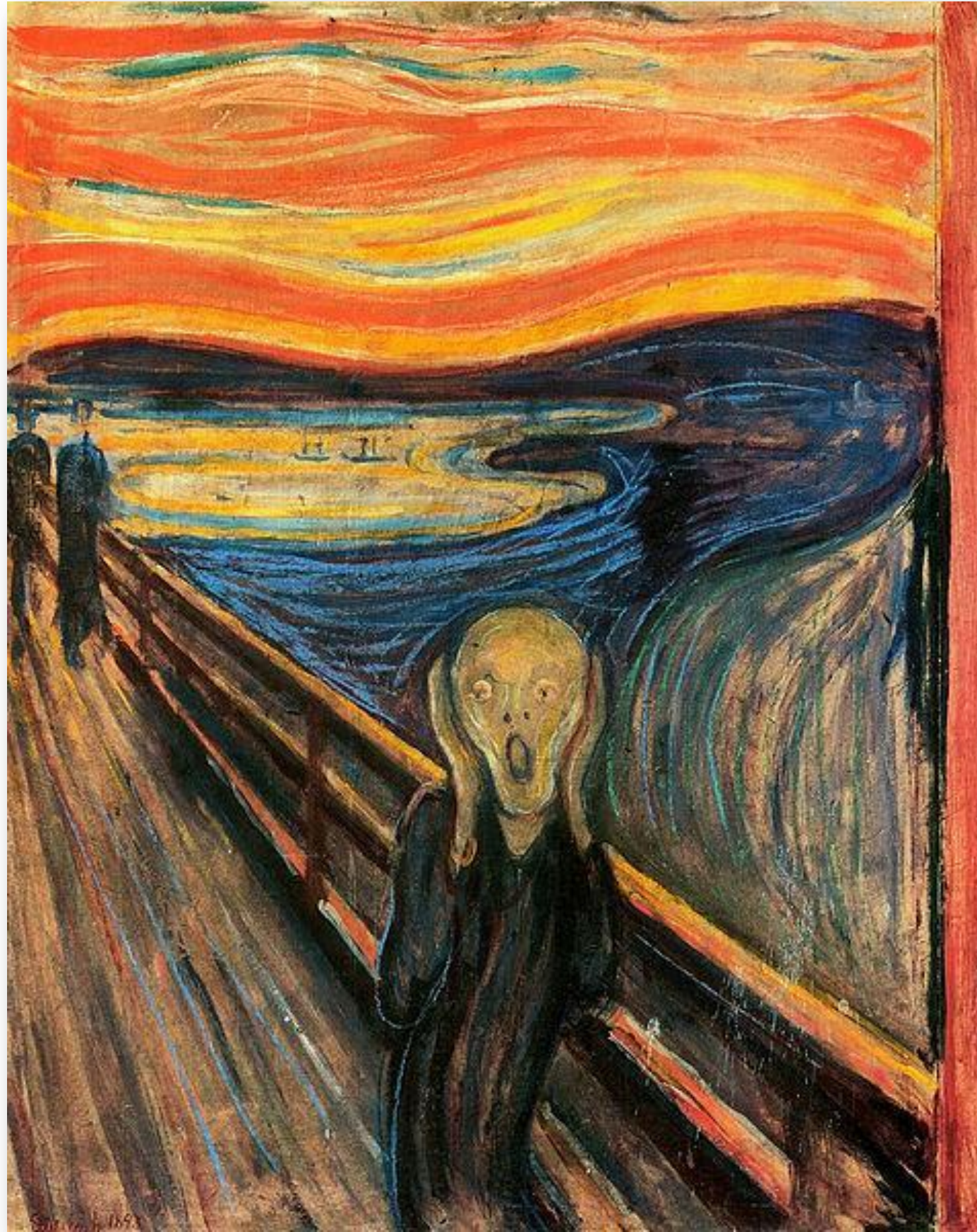


- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists

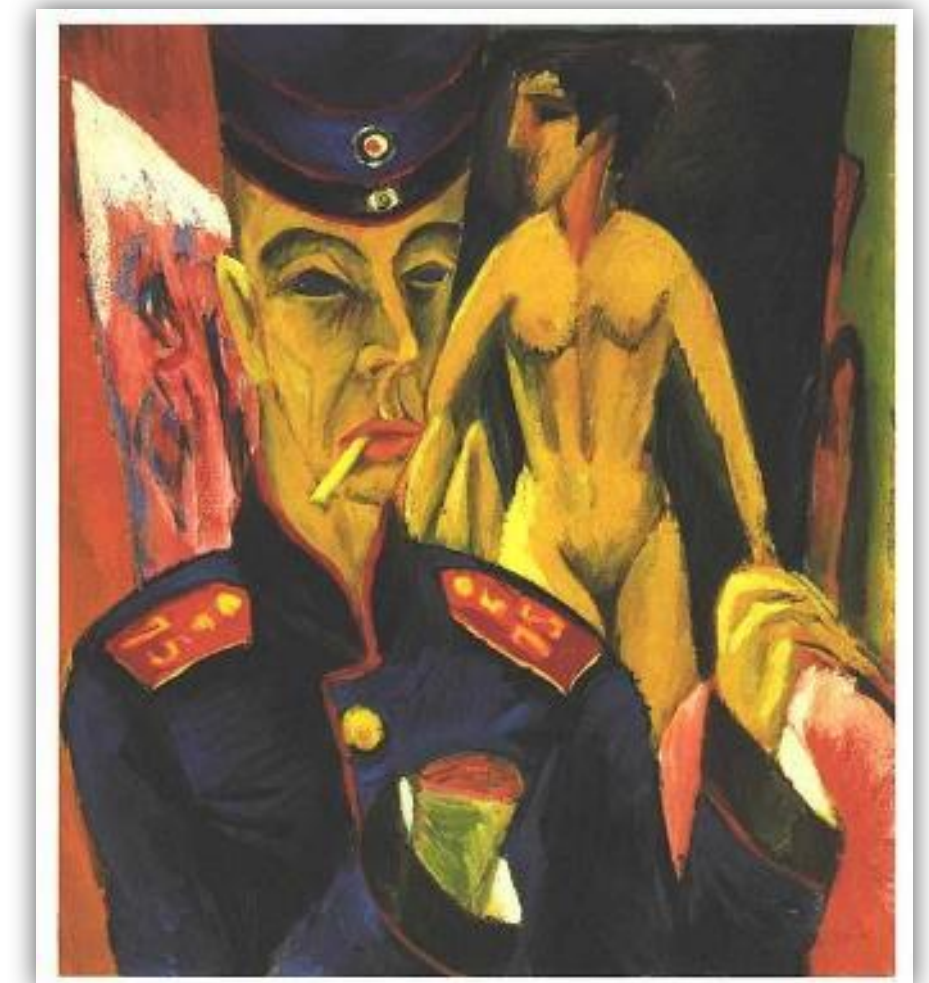
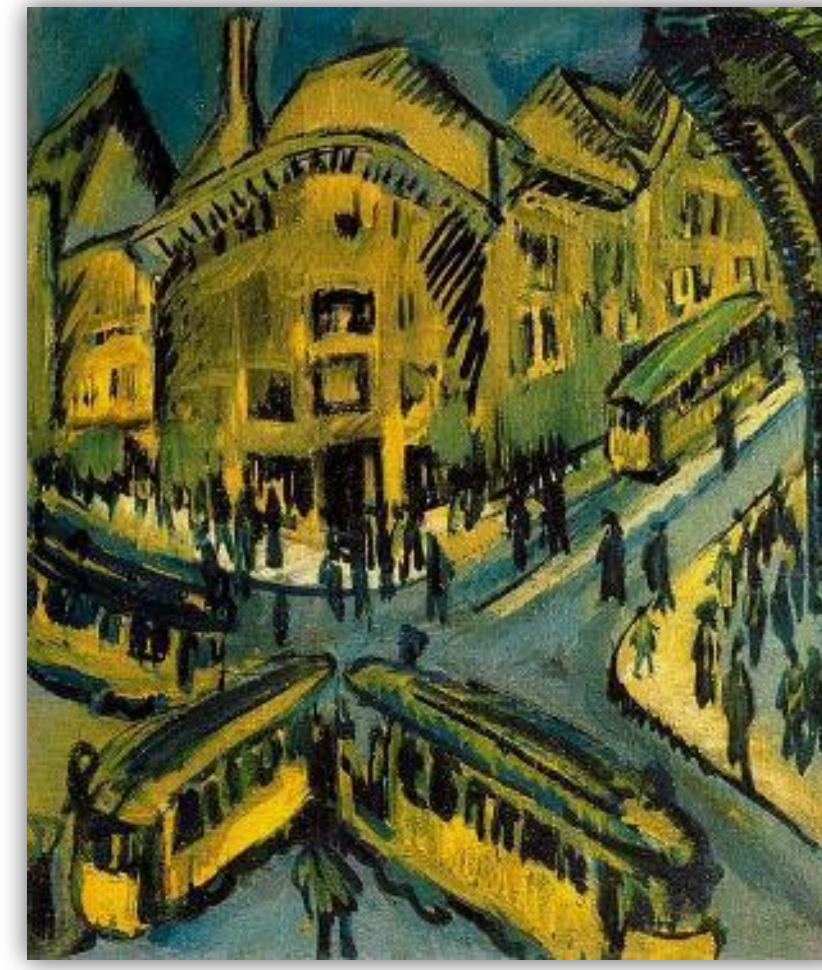
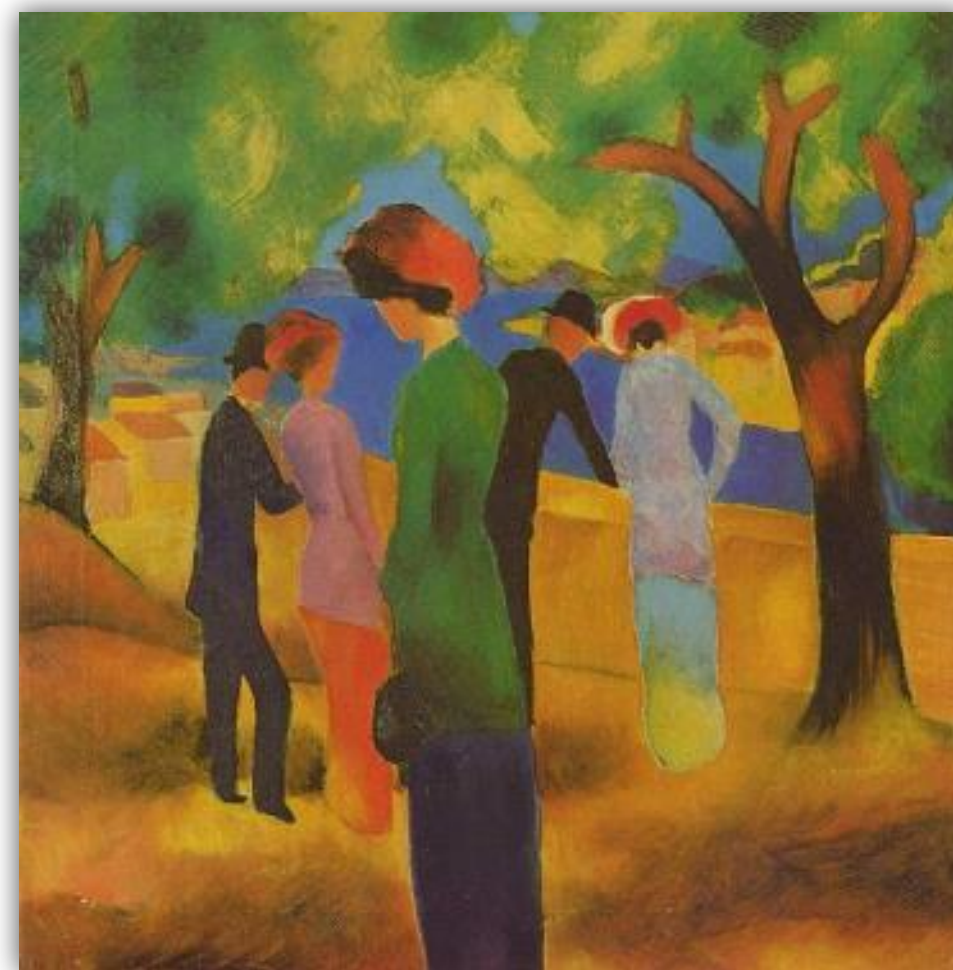
as



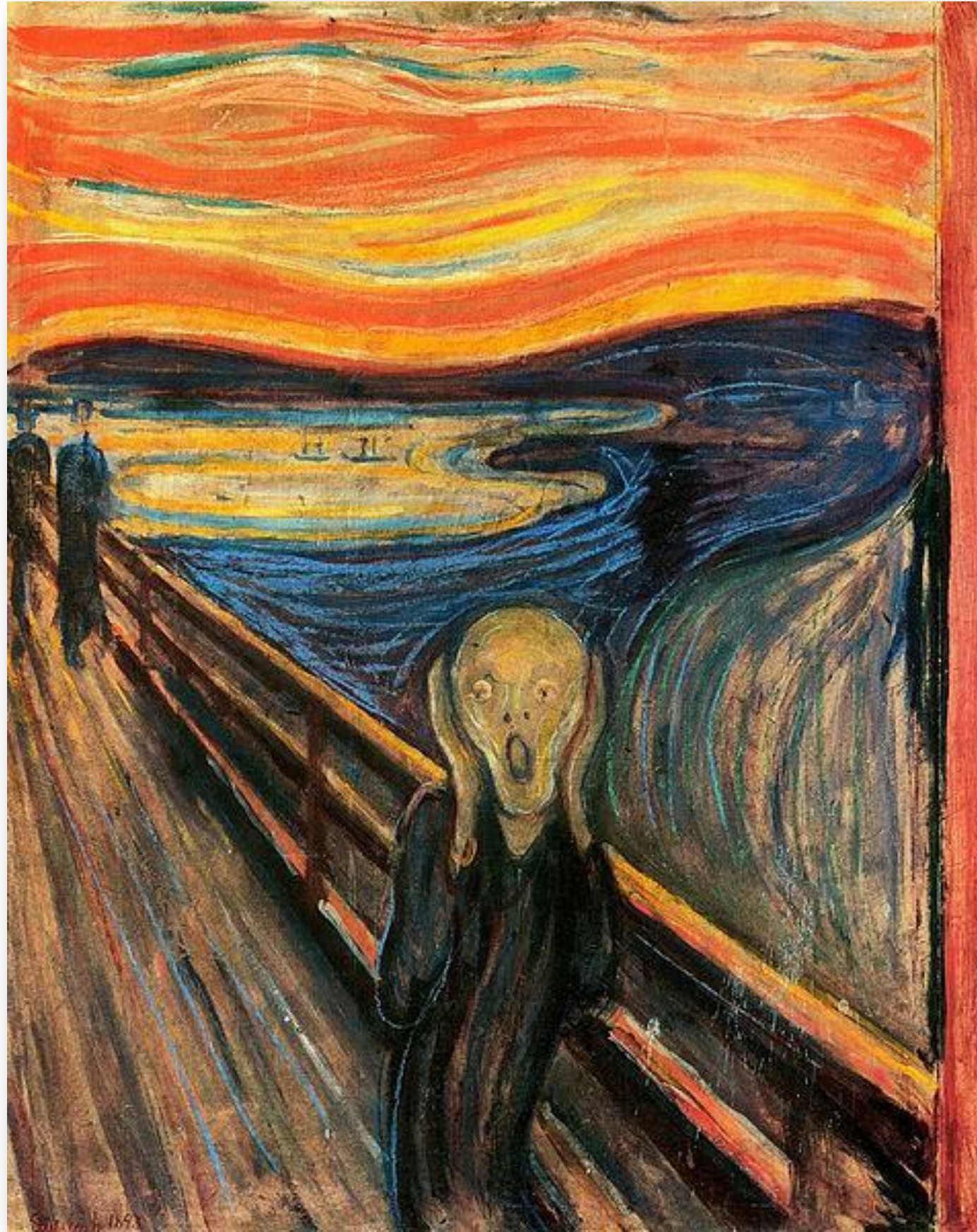
# The Scream, Edvard Munch



- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting **as**



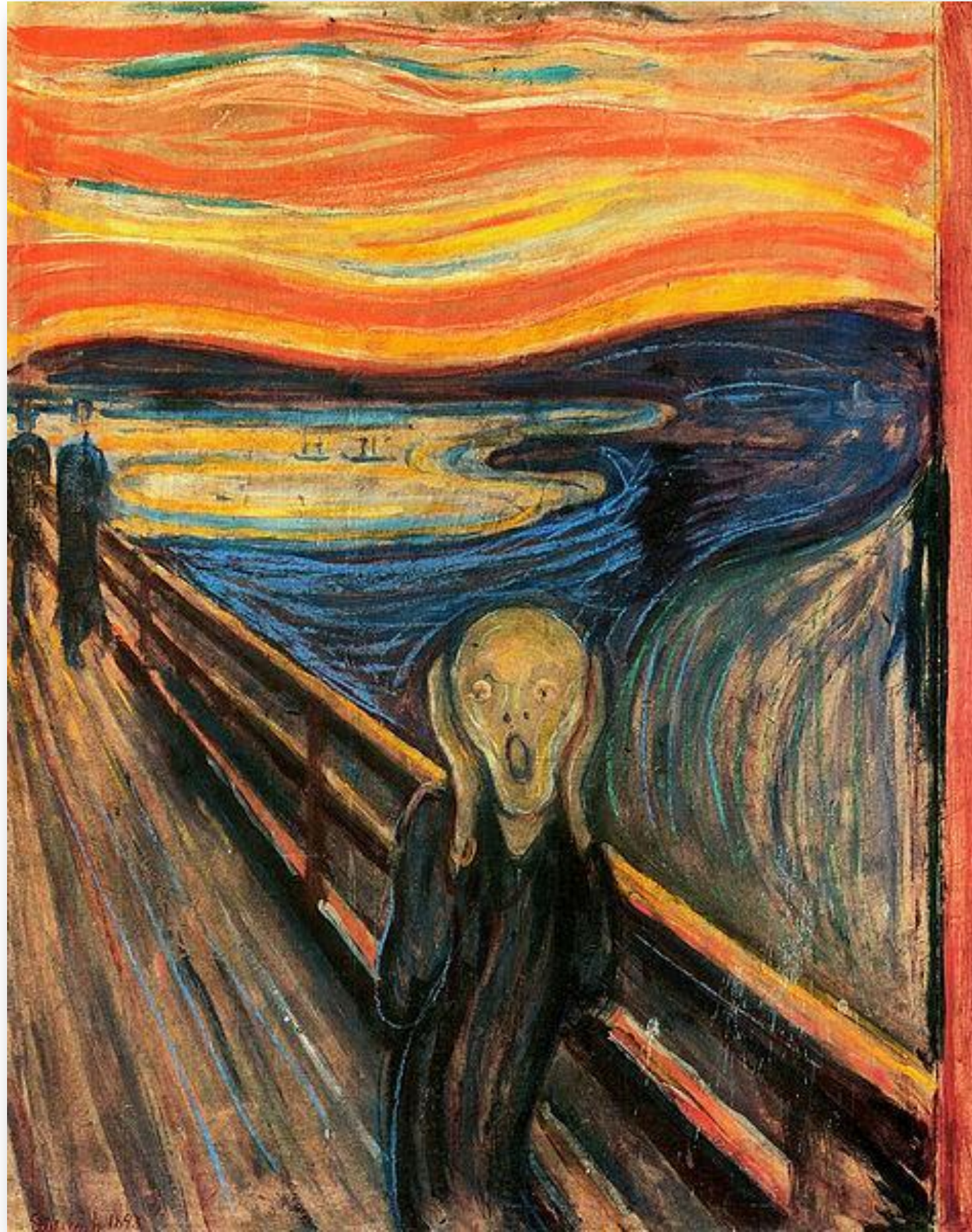
# The Scream, Edvard Munch



- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting
- A painting **as**



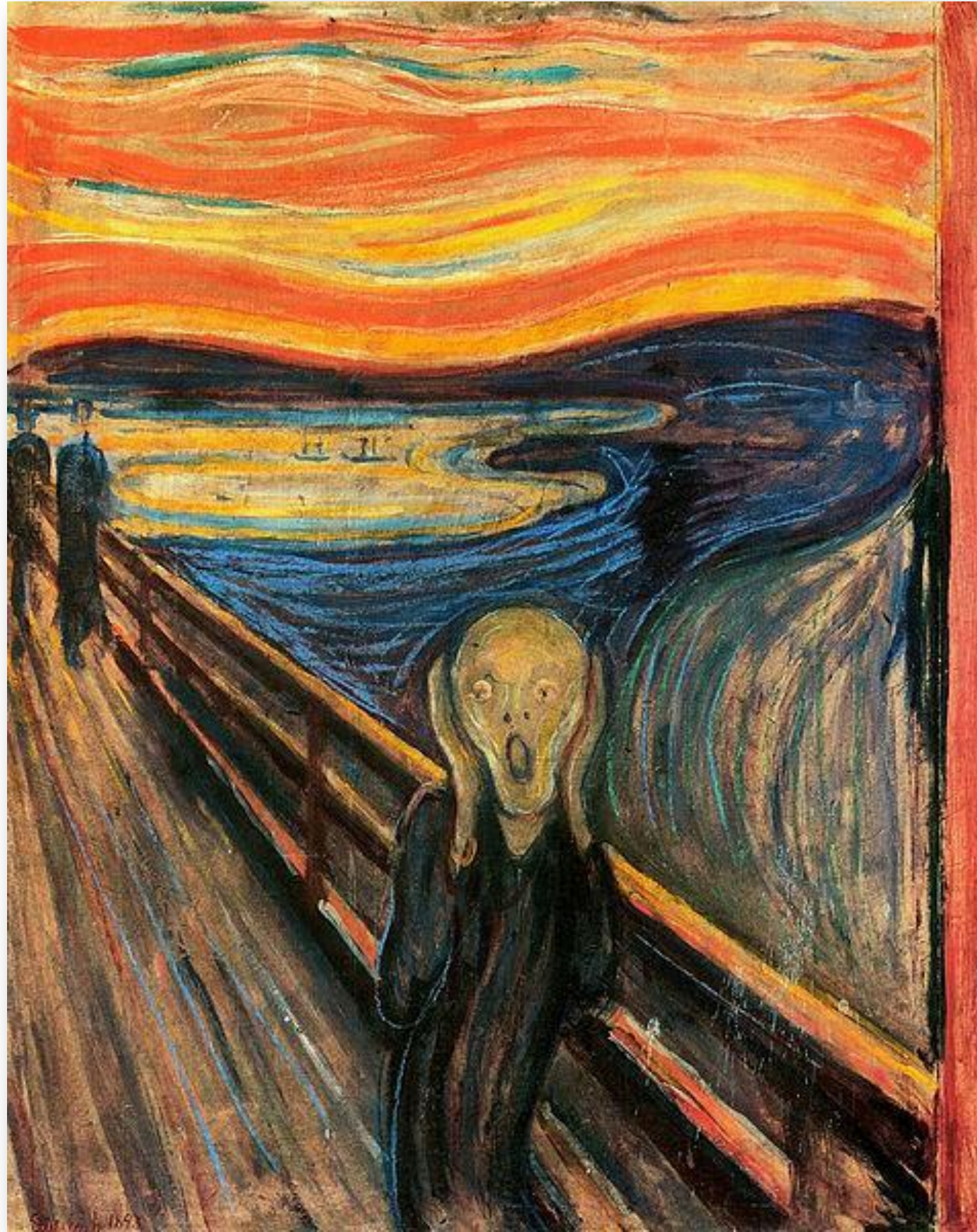
# The Scream, Edvard Munch



- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting
- A painting
- An hand made object **as**



# The Scream, Edvard Munch



- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting
- A painting
- An hand made object
- **An artificial object**  
being the product of intentional human manufacture



# The Scream, Edvard Munch



Low-level

- **The file** at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- **One** of the files of the same picture
- Almost the **same**
- A picture of **the object** at National Gallery, Oslo
- **One of** “The Scream”s by Edvard Munch
- **A** painting by Edvard Munch
- **One of** “The Scream”s by various artists
- **An** expressionist painting
- **A** painting
- **An** hand made object
- **An** artificial object  
being the product of intentional human manufacture

High-level

# The Scream, Edvard Munch



Matching

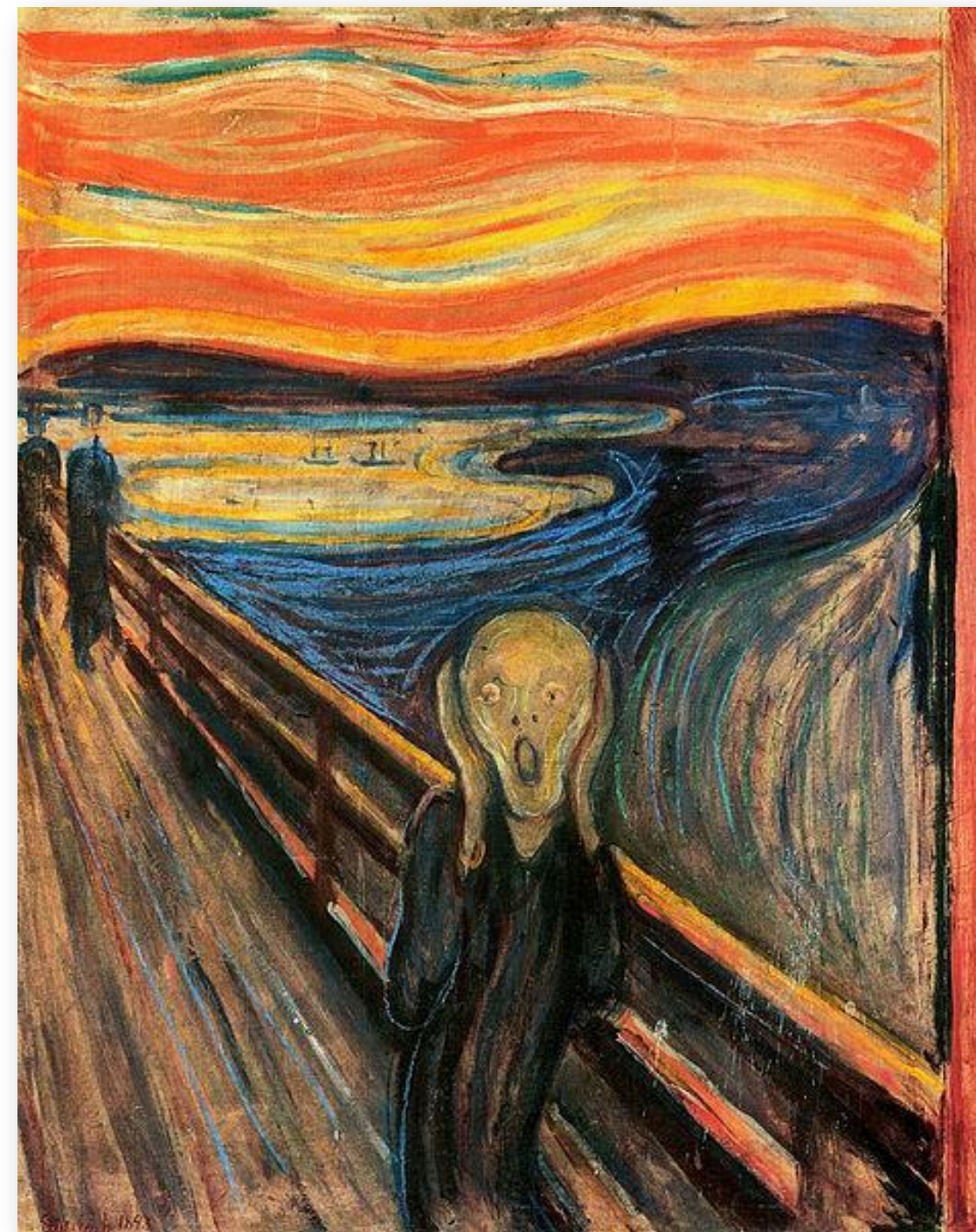
Recognition

Classification

- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting
- A painting
- An hand made object
- An artificial object  
being the product of intentional human manufacture

# Similarity between different representations

It is possible to define a concept of similarity for these different levels of abstraction



instance similarity  
(very low abstraction)



semantic similarity  
(very high abstraction)





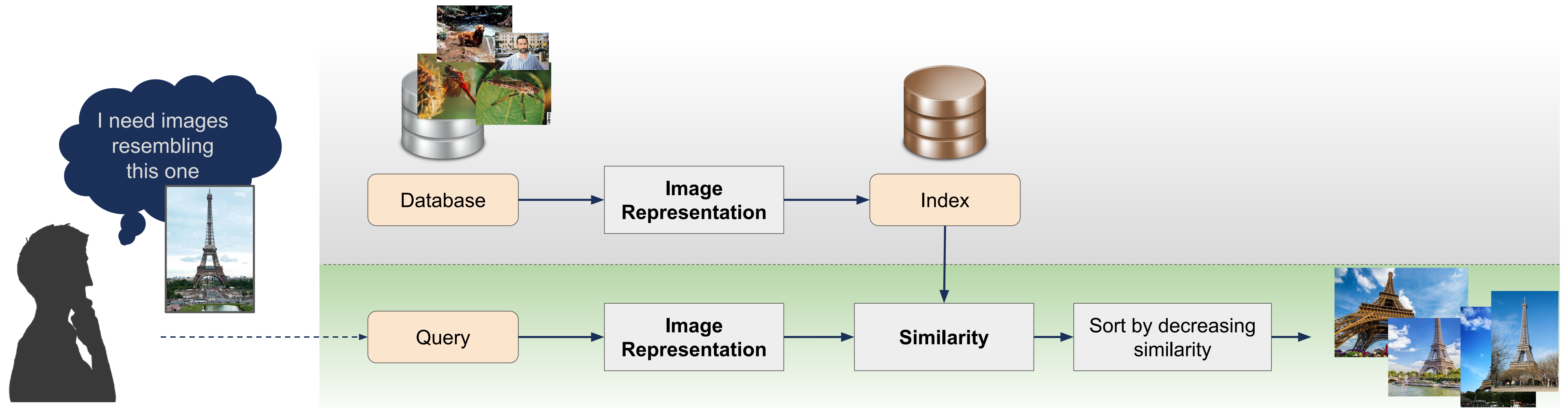
- The file at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- One of the files of the same picture
- Almost the same
- A picture of the object at National Gallery, Oslo
- One of "The Scream"s by Edvard Munch
- A painting by Edvard Munch
- One of "The Scream"s by various artists
- An expressionist painting
- A painting
- An hand made object
- An artificial object being the product of intentional human manufacture

How can we associate a **representation** to each one of these abstraction levels?

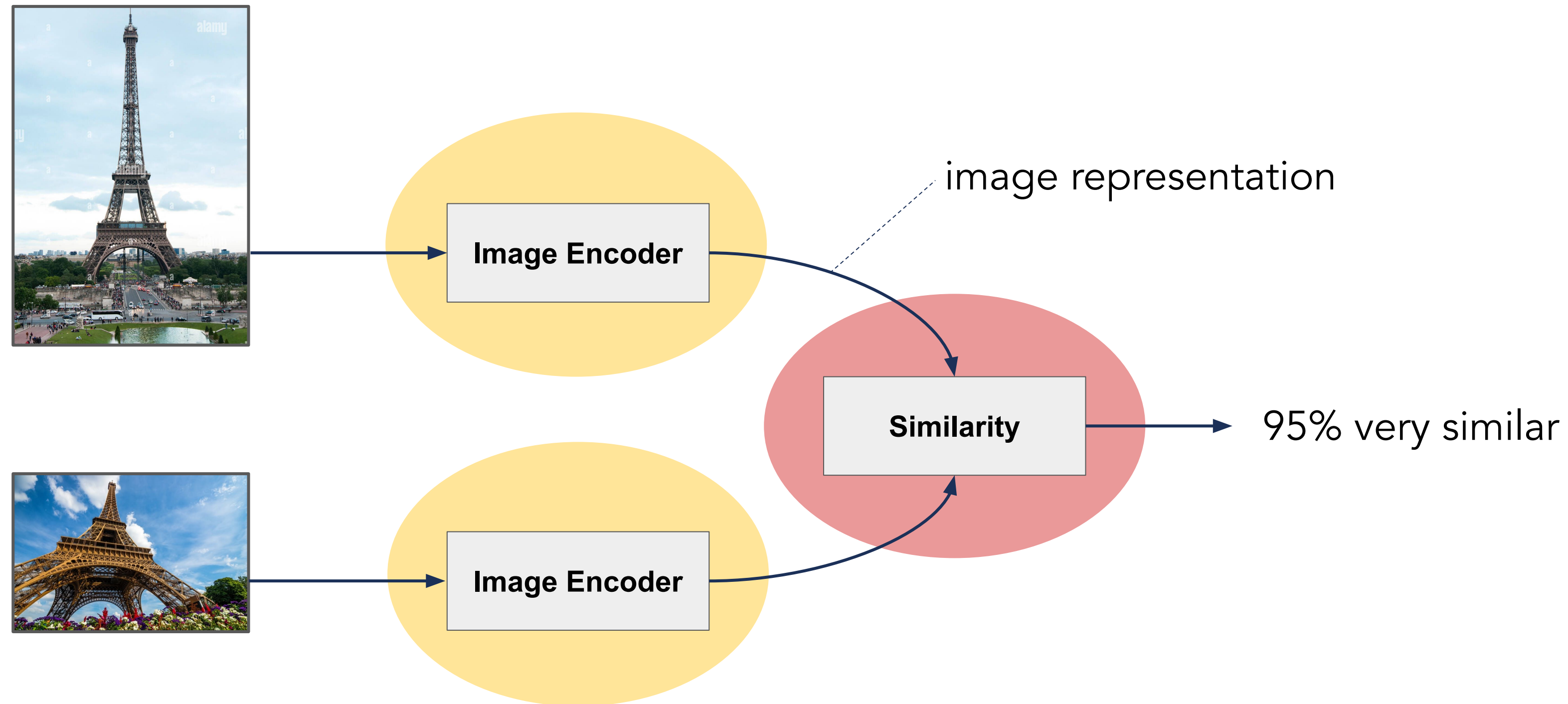


# Image Representations

# Image retrieval setup



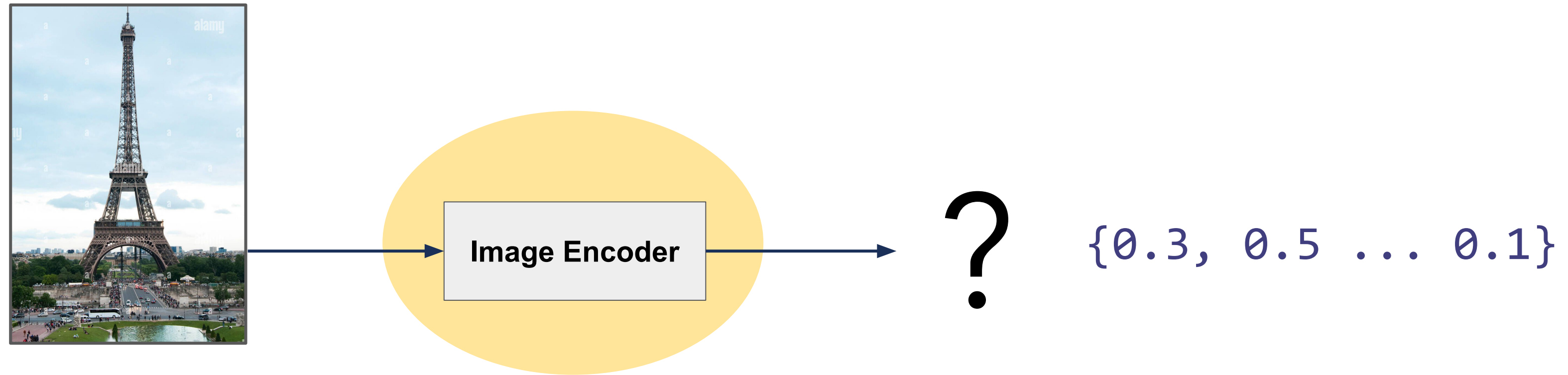
# Representations and similarities



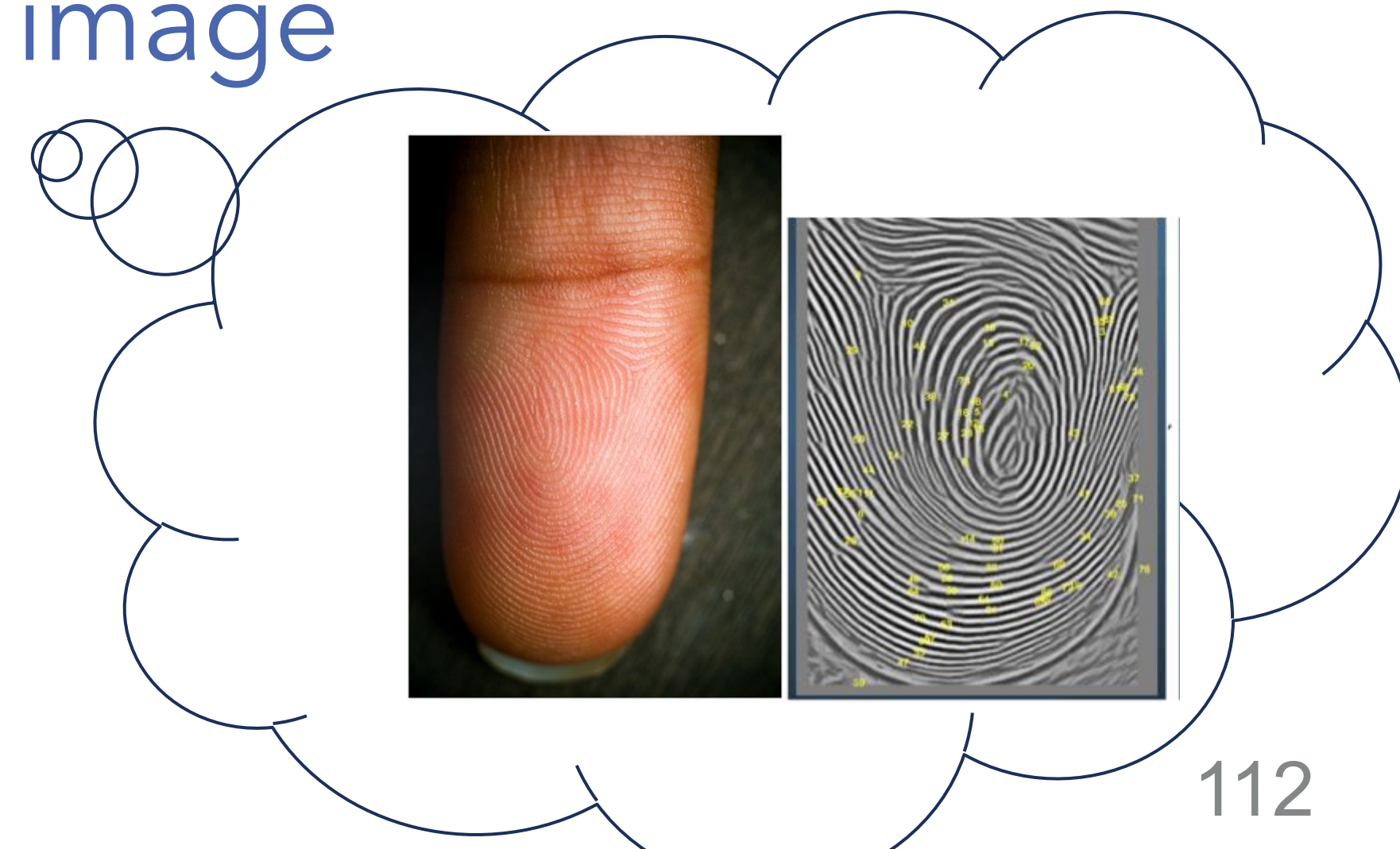
All boils down to

- Finding an image encoder that outputs good representations
- Measuring the similarity between these two representations

# Representation / Features



An image is converted into a set of numbers, called vector  
It can be considered the "fingerprint" of that image



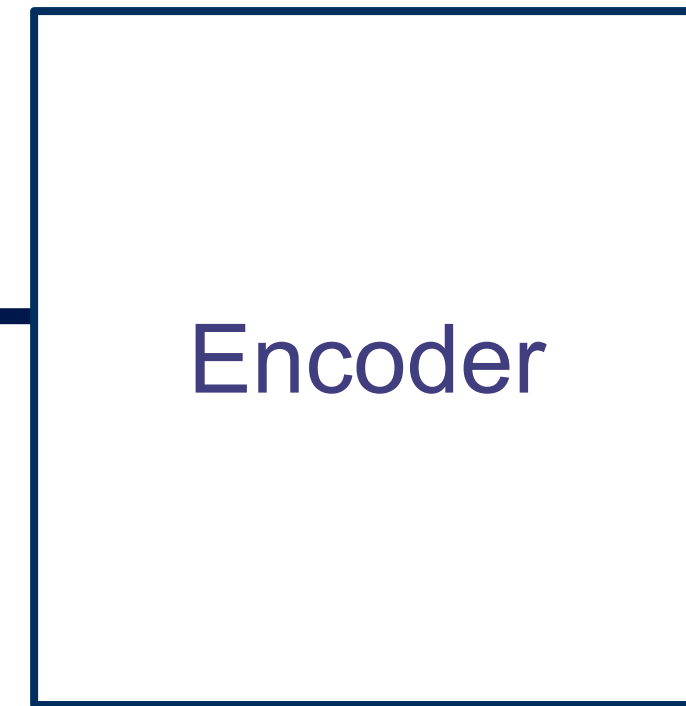


# Representation / Features



Digital Image

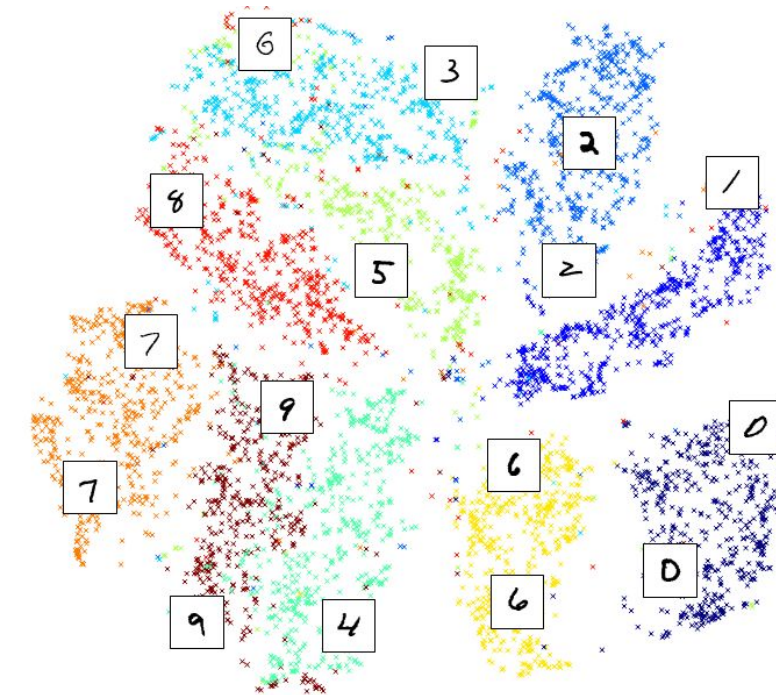
000001  
1100101  
000011  
100000  
001001  
000001  
010110  
001100  
001101  
101110  
1111011  
110010  
0111011  
0111011  
010110111001100110000011101000110100  
01101001011011100110000101110010011110010  
100011101010110010101110011001000001101111  
1110100001011100101010001101000011010010111  
1110100000011100110110111011011011010101010



feature  
extractor

A representation  
(a feature)

{0.3, 0.5 ... 0.1}



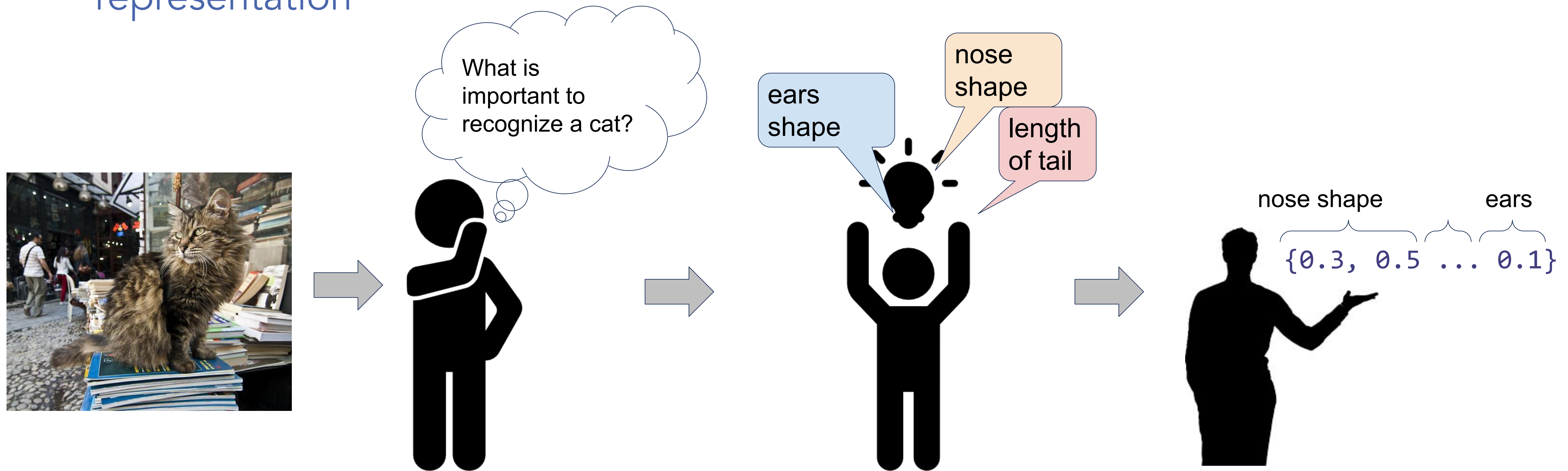
in a **latent space**



An encoder (or feature extractor) takes an input (e.g., an image) and produce a representation (or descriptor) that is used (in place of the image) for the specific task.

# Handcrafted features

- Before Deep Learning, handcrafted features/representations:
  - Human is always in the loop
  - He chooses what is important in the picture for creating a numerical representation



# Representation / Features

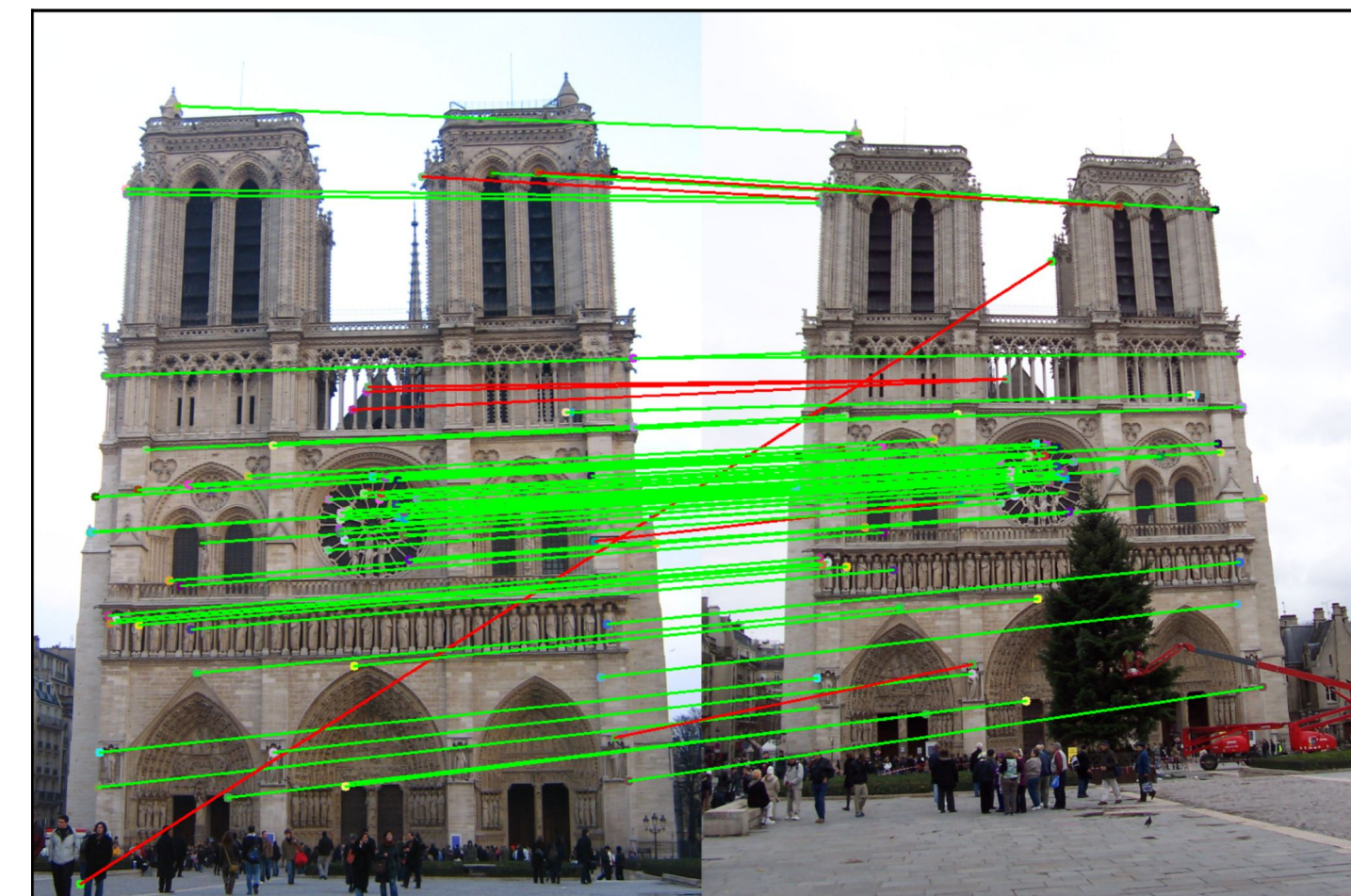
## Global Features:

- color, edge, texture etc...

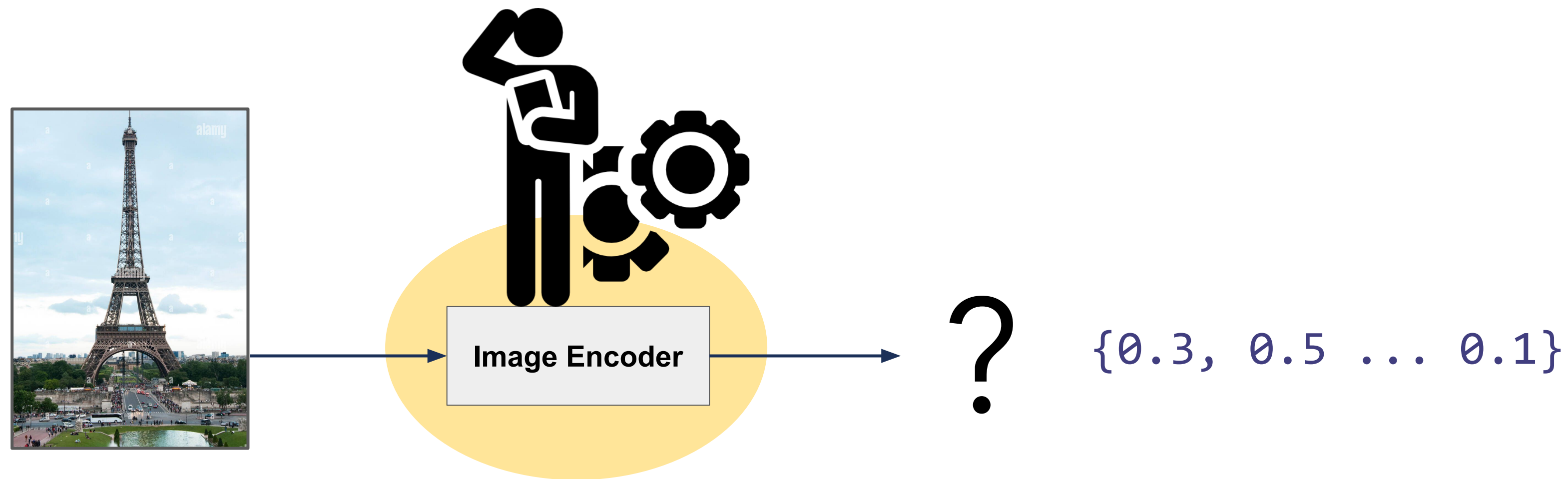


## Local Features:

- representation of interest points/regions
- for image stitching or object recognition

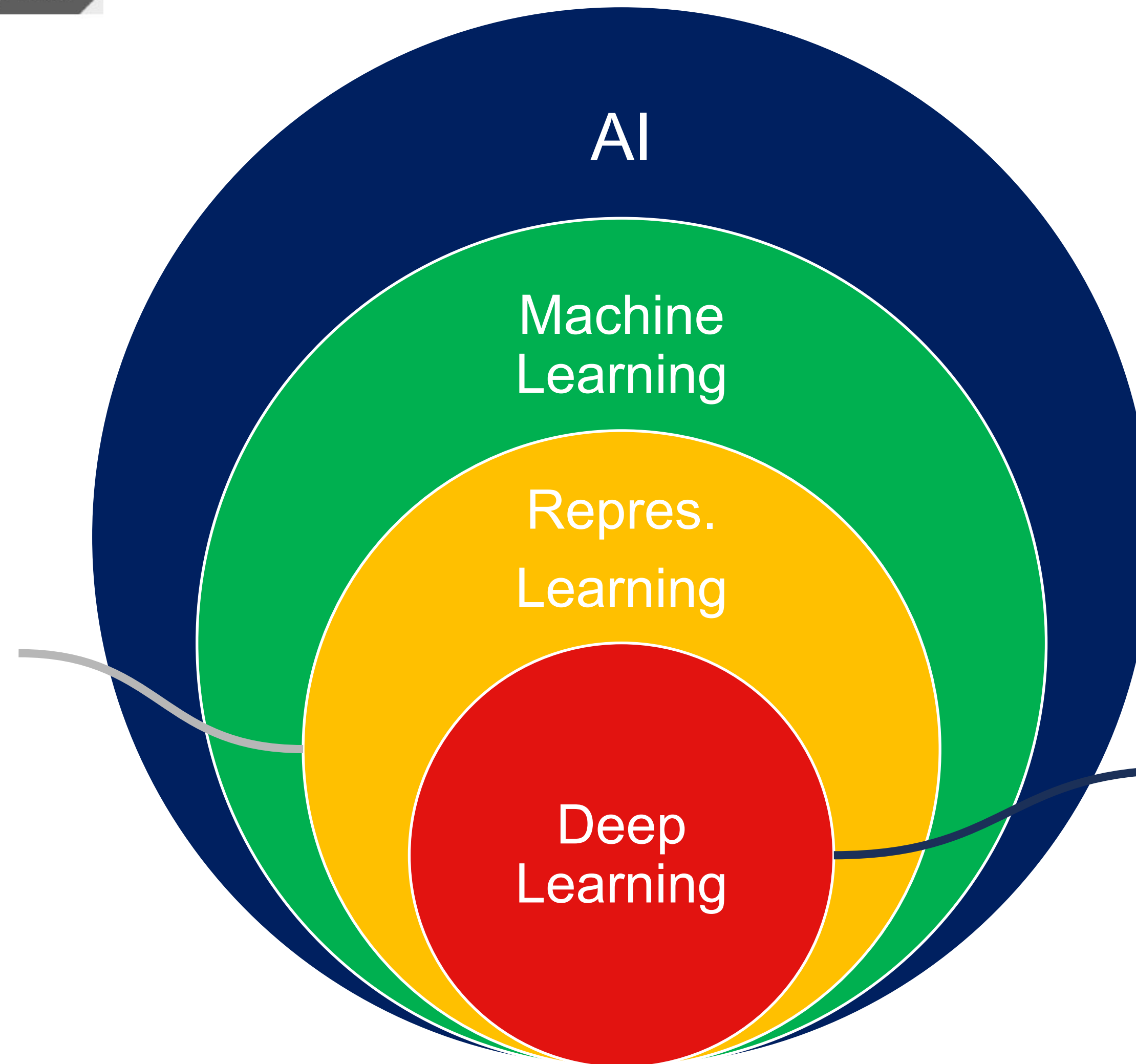


# Handcrafted features



- Good for low level features
  - Colors, shapes, important keypoints
- But how can we define a priori the characteristics of the image that enable us to recognize complex entities (e.g., a “cat”, or a “tower”)?

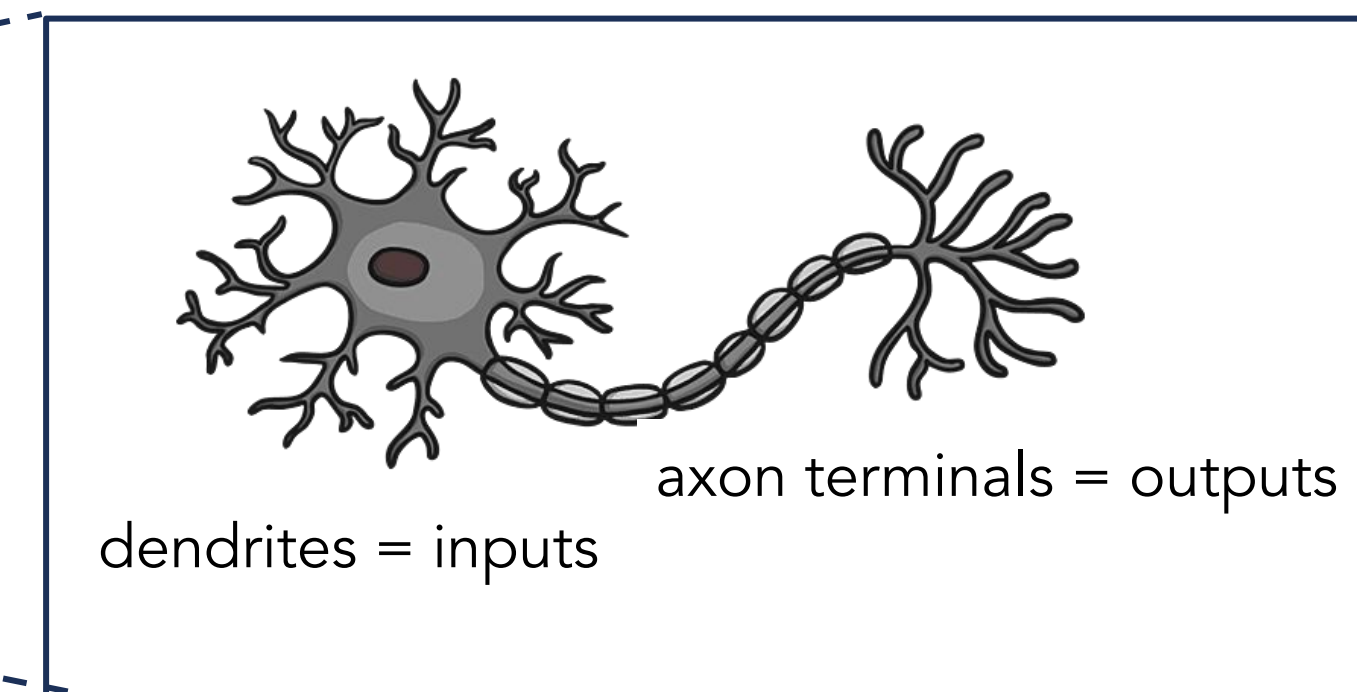
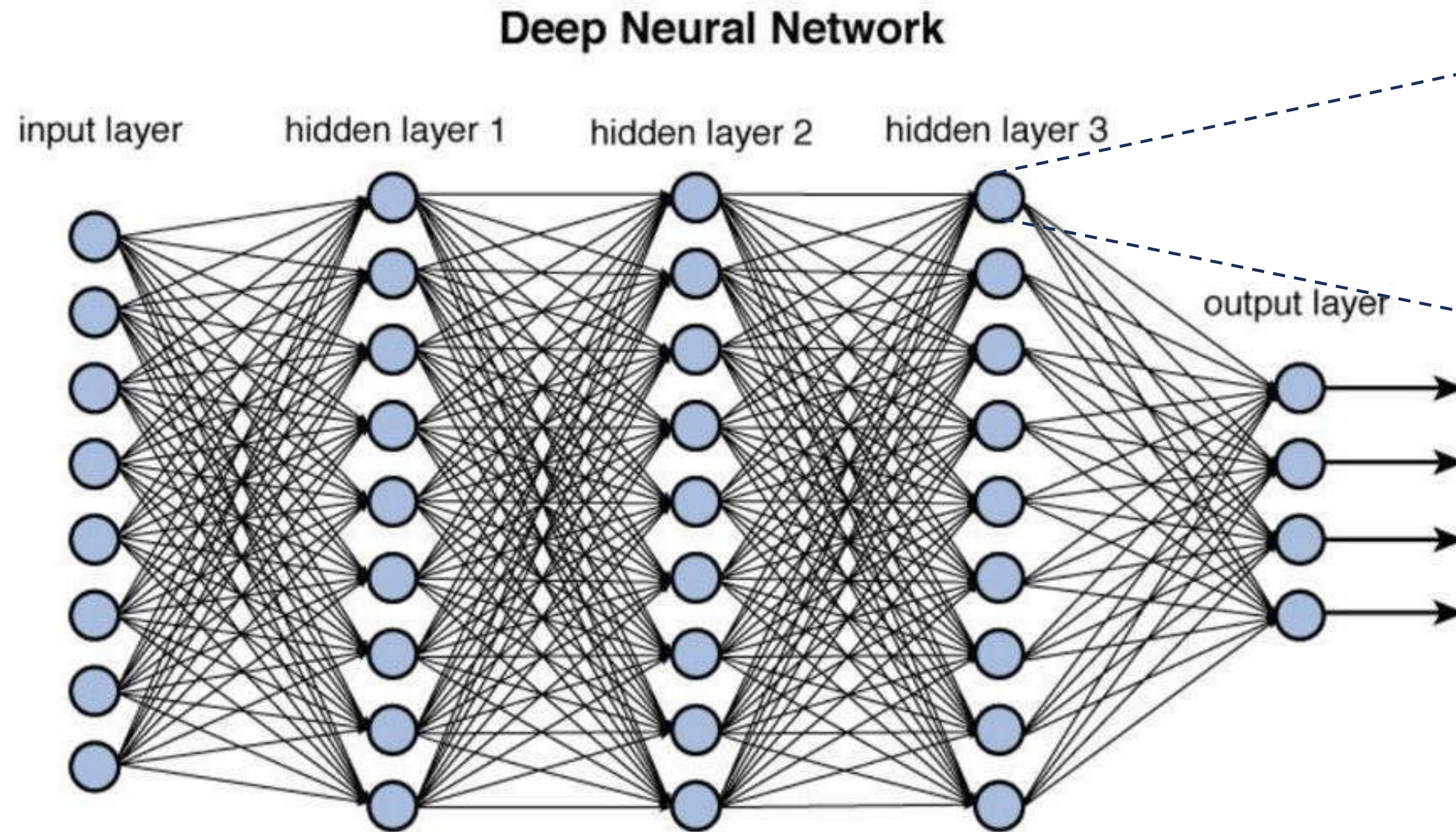
# Deep learning (from Nature)



Representations are learned from data!  
No more handcrafted!

Representations are learned inside multiple layers of a deep neural network

# A Deep Neural Network

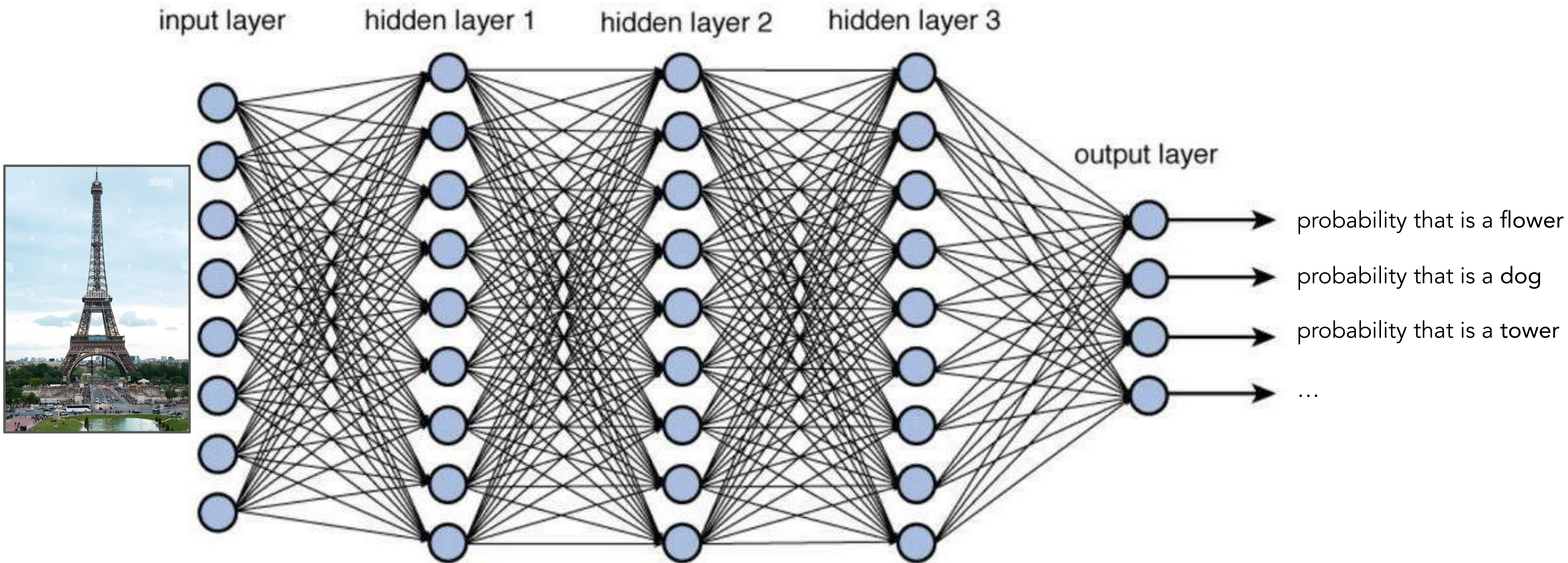


Every neuron in each layer is connected to all (or some) of the neurons of the previous layer

. Feed Forward Neural Networks

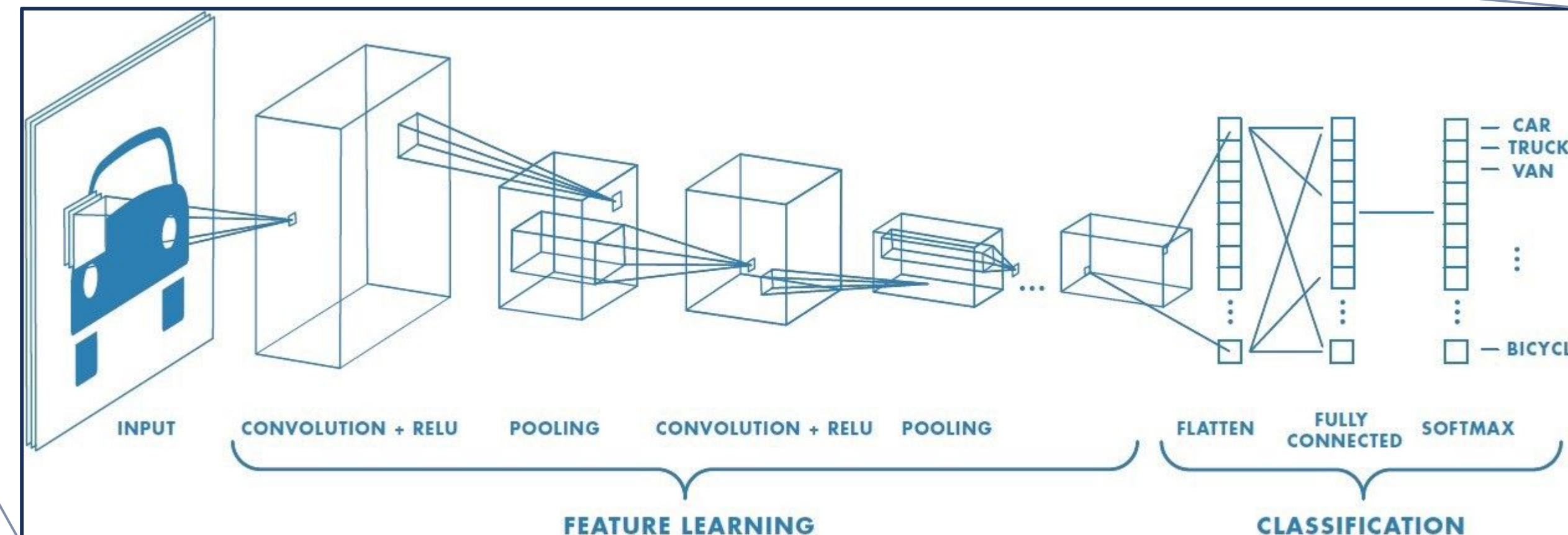
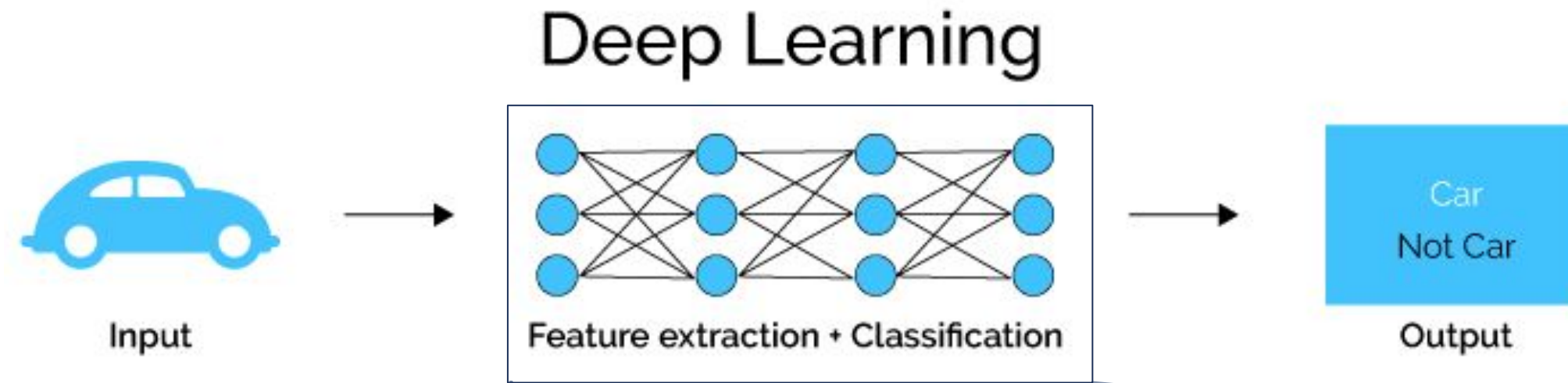
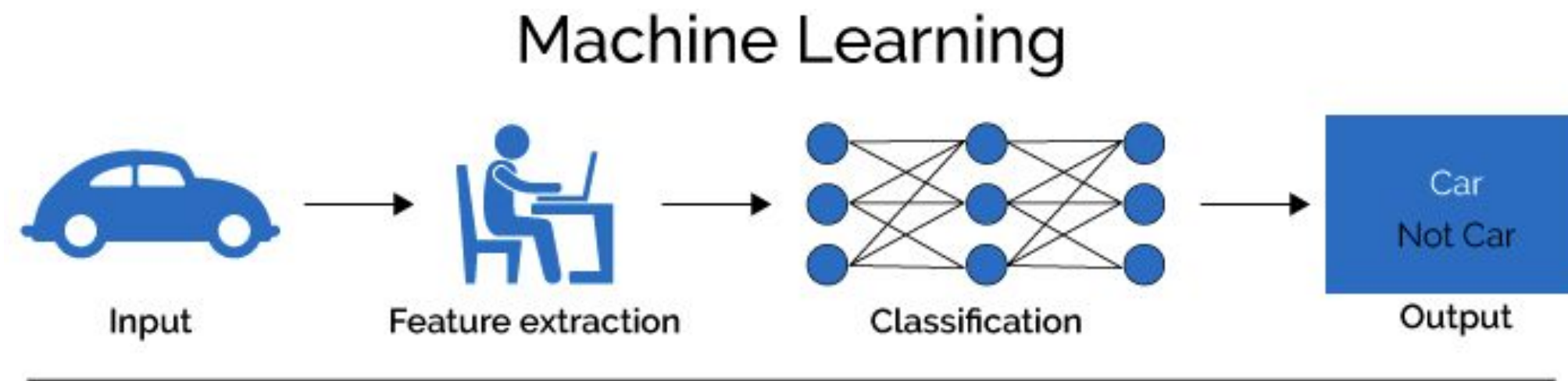
# A Deep Neural Network

Deep Neural Network



- . The input could be an image
- . The output the class of the object contained in the image (e.g., a "cat")

# Deep Learning vs standard Machine Learning





# Deep learning (from Nature)

## Representation learning methods:

allow a machine to be fed with raw data and to automatically discover the representations needed for classification.

## Deep-learning are representation learning methods

- with multiple levels of representation, obtained by
- composing simple modules that
- transform the representation at one level into a representation at a higher, slightly more abstract level.

# The Scream, Edvard Munch



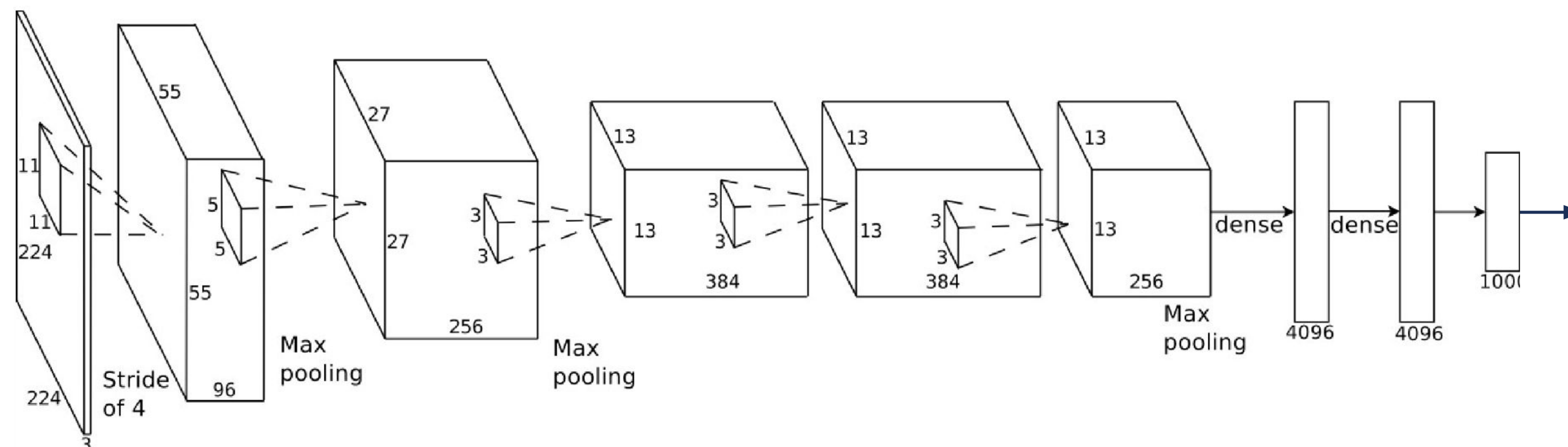
Low-level

- **The file** at [http://upload.wikimedia.org/.../475px-The\\_Scream.jpg](http://upload.wikimedia.org/.../475px-The_Scream.jpg)
- **One** of the files of the same picture
- Almost the **same**
- A picture of **the object** at National Gallery, Oslo
- **One of** “The Scream”s by Edvard Munch
- **A** painting by Edvard Munch
- **One of** “The Scream”s by various artists
- **An** expressionist painting
- **A** painting
- **An** hand made object
- **An** artificial object  
being the product of intentional human manufacture

High-level

# The training procedure

A deep network should be trained before being used  
How is training performed?



car? no, is a cat!  
dog? no, is a cat!  
lion? no, is a cat!  
cat? yes!

The network in this case learns like small childs (direct supervision)

# Supervised learning... Sometime it fails

Direct supervision does not always work

- works generally well but it is often overused
- does not work easily with logical or mathematical reasoning
- intuition vs reasoning



# A nice side-effect

---

After the training is completed, we can observe the rise of nice representations from the intermediate layers!

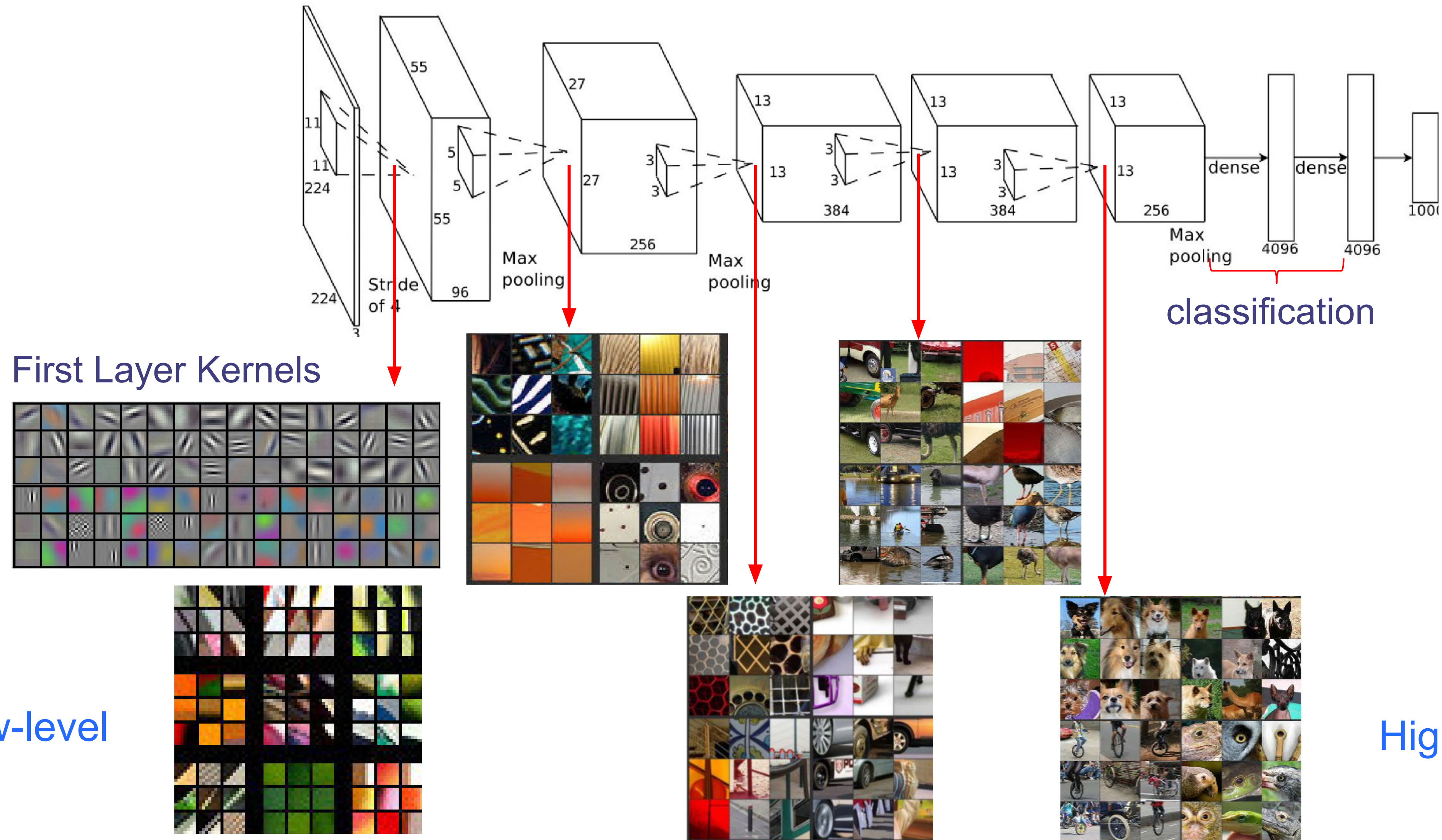
The first layers learnt low level details (textures)

The higher layers learnt very semantic details (faces, objects)

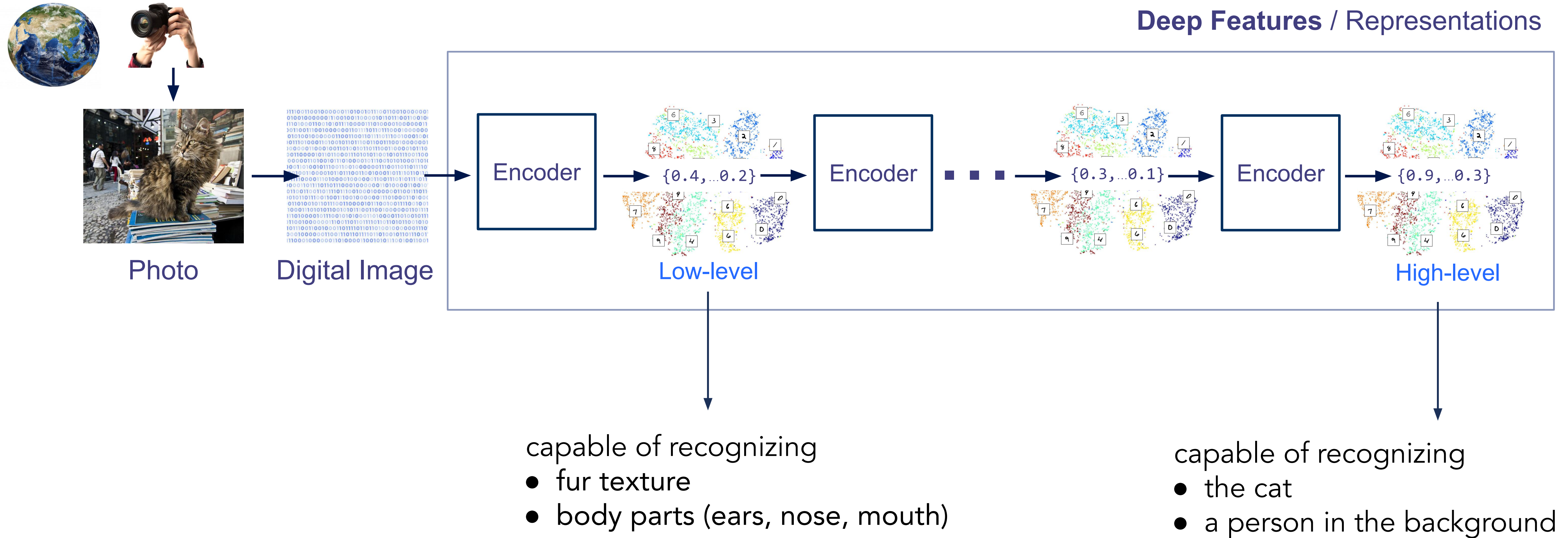


# Multiple Levels Of Abstraction

AlexNet, 2012, Trained on a Classification task of 1,000 classes.



# The Overall Picture

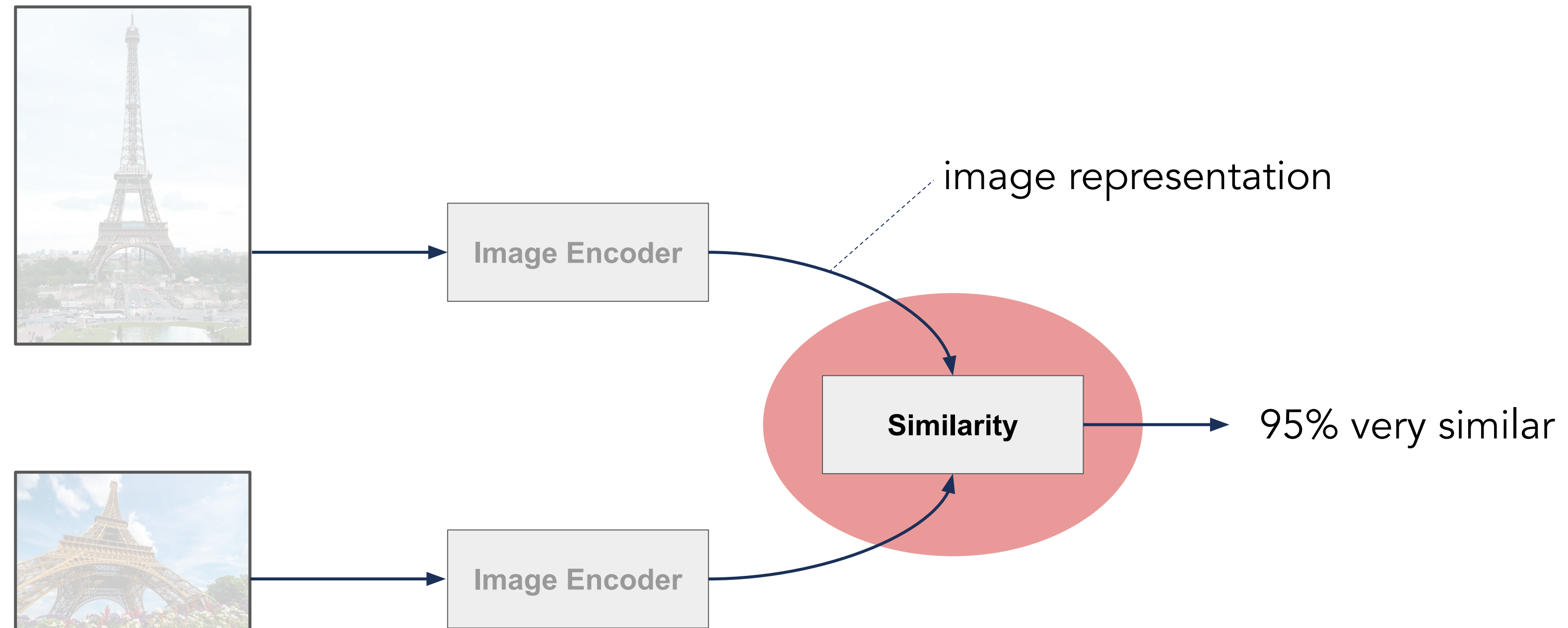


- 
- Define a similarity between representations
  - Similarity between representations as a way to measure the similarity between different images
    - of course, using representations from different levels bring to a different idea of relevance
      - low-level (instance) similarity
      - high-level (semantic) similarity



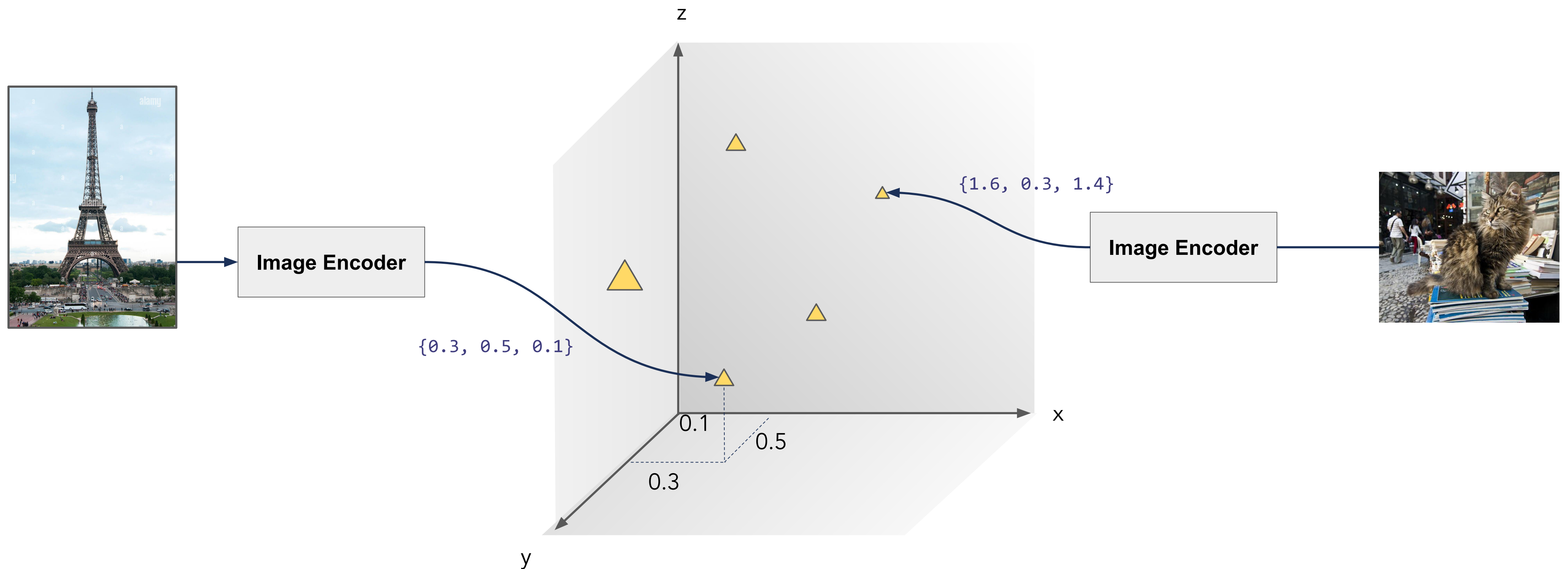
# Similarity between Representations

# Representations and similarities



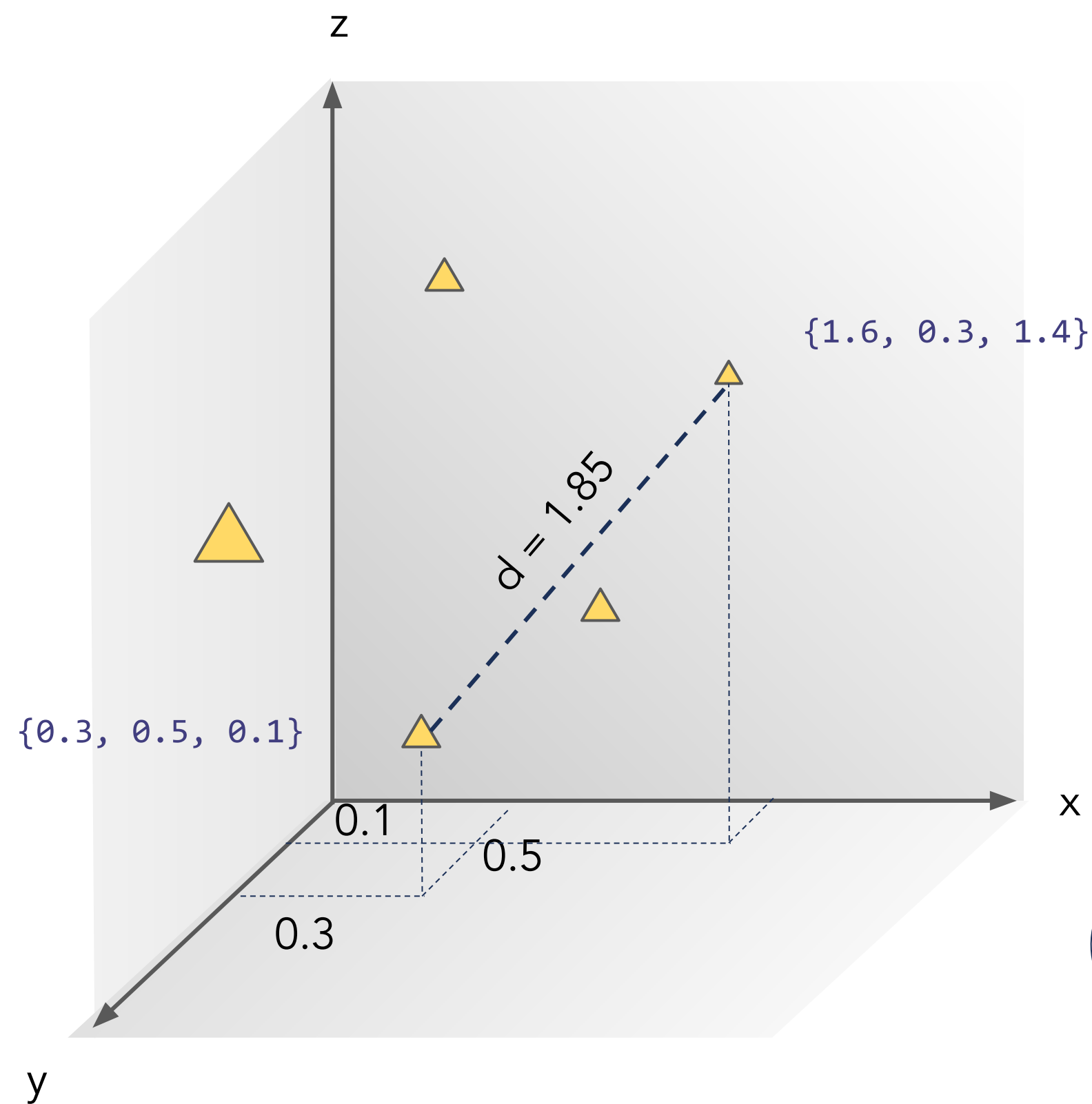
# Representations in the space

- The representations are list of numbers  $\{0.3, 0.5, 0.1\}$
- They can be represented in a cartesian space



# Distance between representations

- We can define a distance between representations
- Usually, Euclidean distance, a.k.a. Pythagorean Theorem



$$d = \sqrt{(1.6 - 0.3)^2 + (0.3 - 0.5)^2 + (1.4 - 0.1)^2} \simeq 1.85$$

This is how in the game of boules we measure distances between the balls.



# Euclidean distance (in 2D)

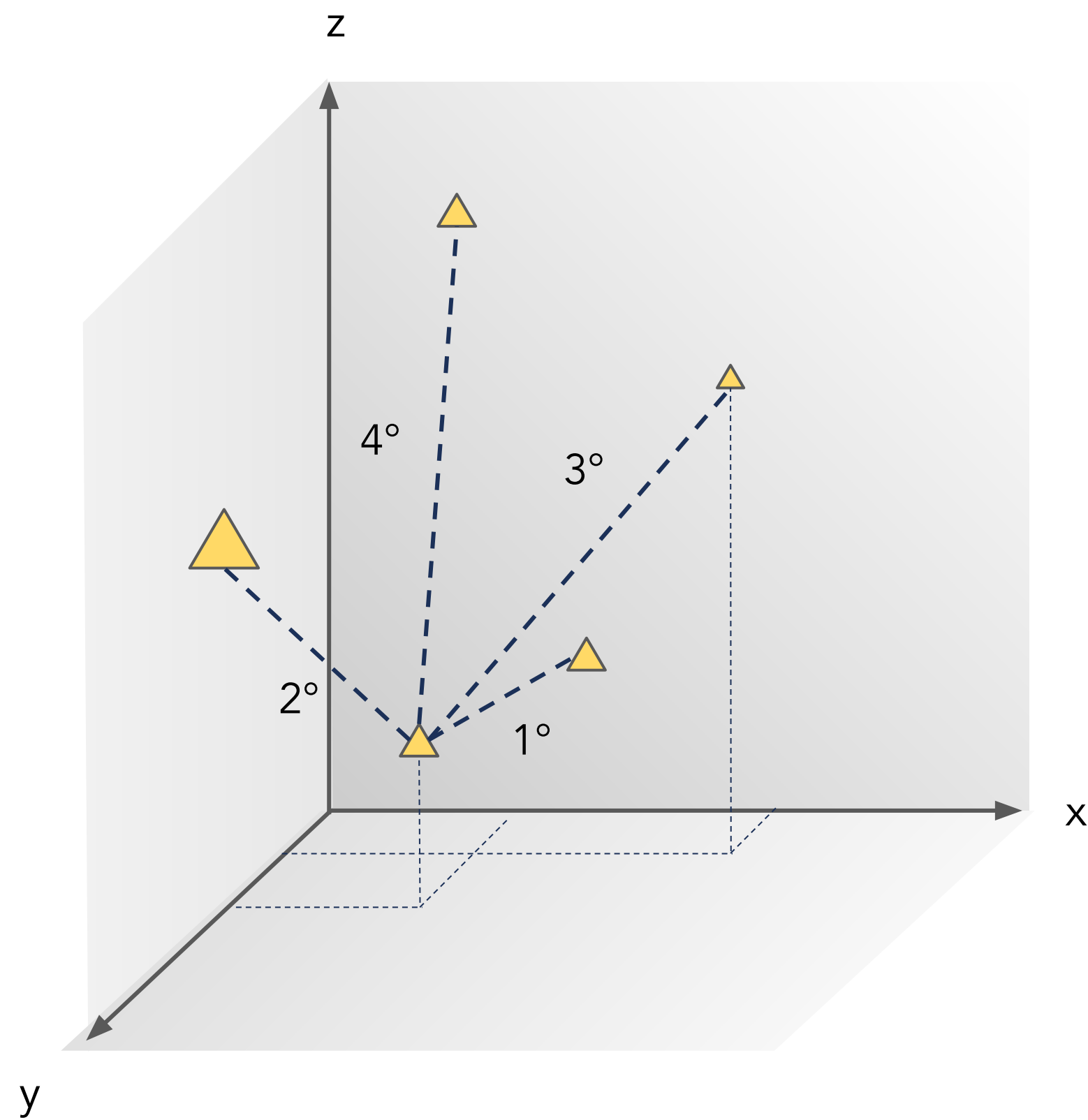
## Euclidean Distance

$$\text{Euclidean}(A, B) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



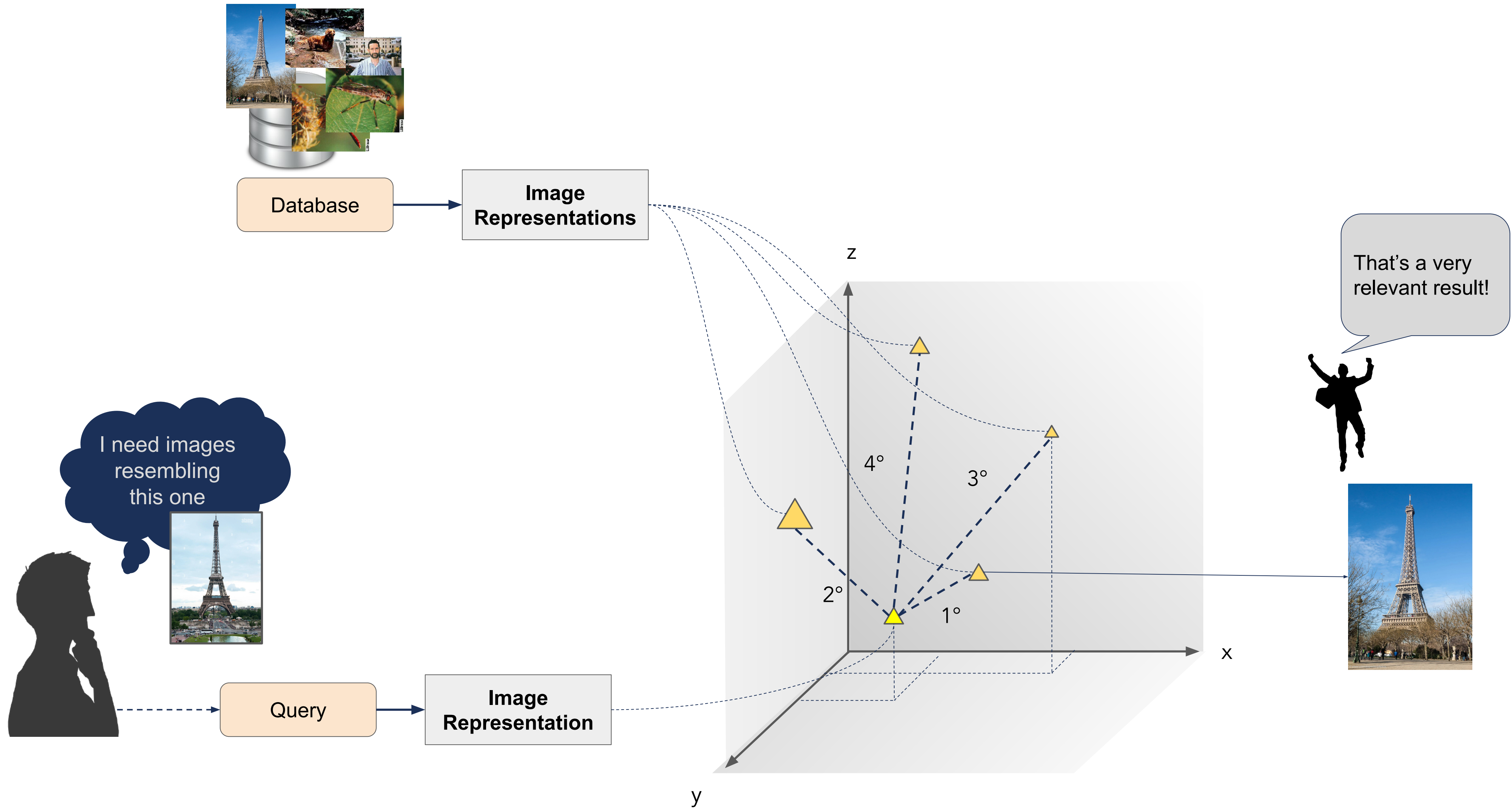
# Similarity

- . We can define the similarity as the opposite of the distance
- . The more distant, the less similar they are

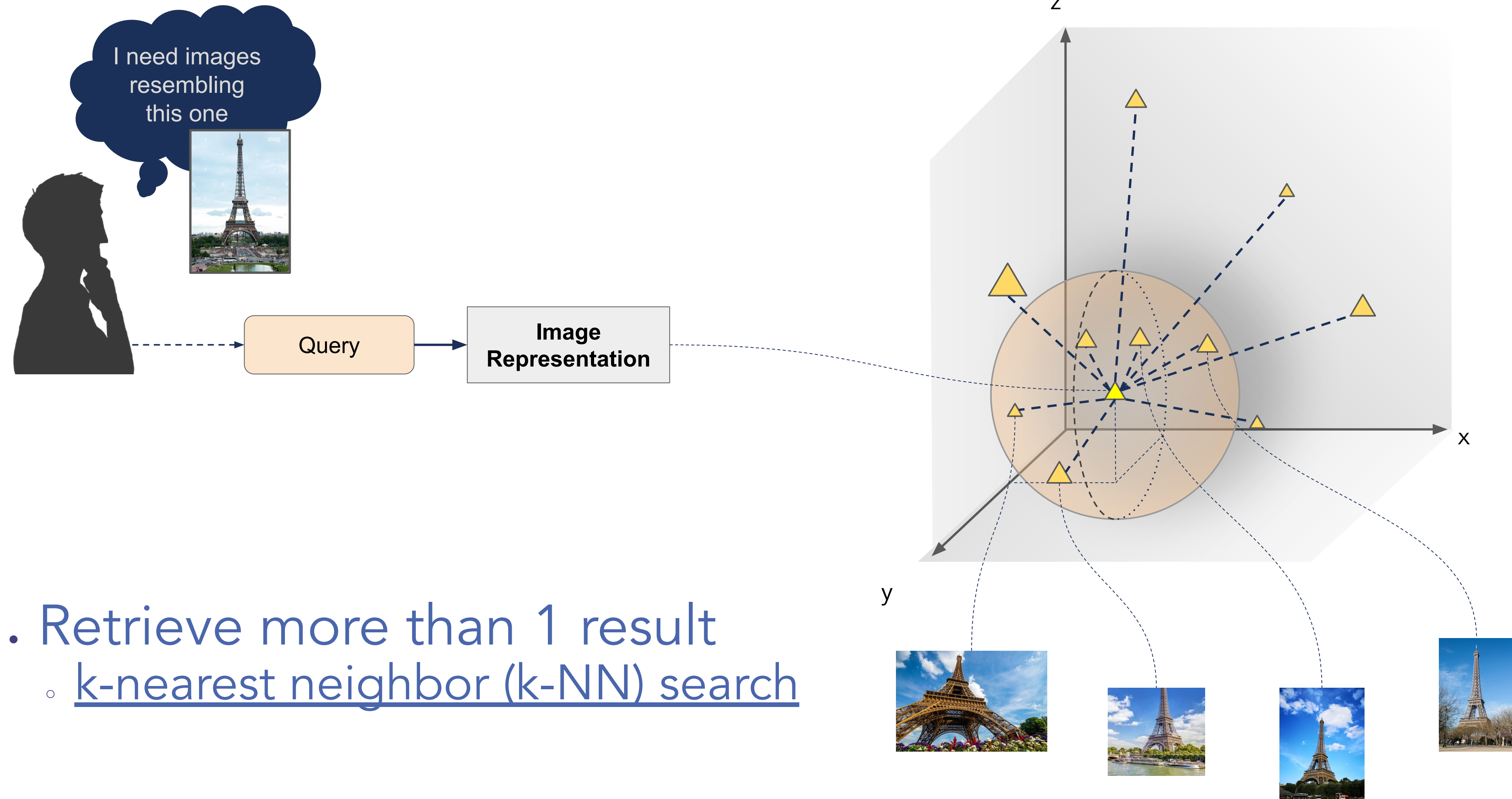


Again, this is actually the method used for assigning points in the game of boules! You have to find out which bowl is the nearest to the target ball

# Back to image retrieval



# Scaling up this idea

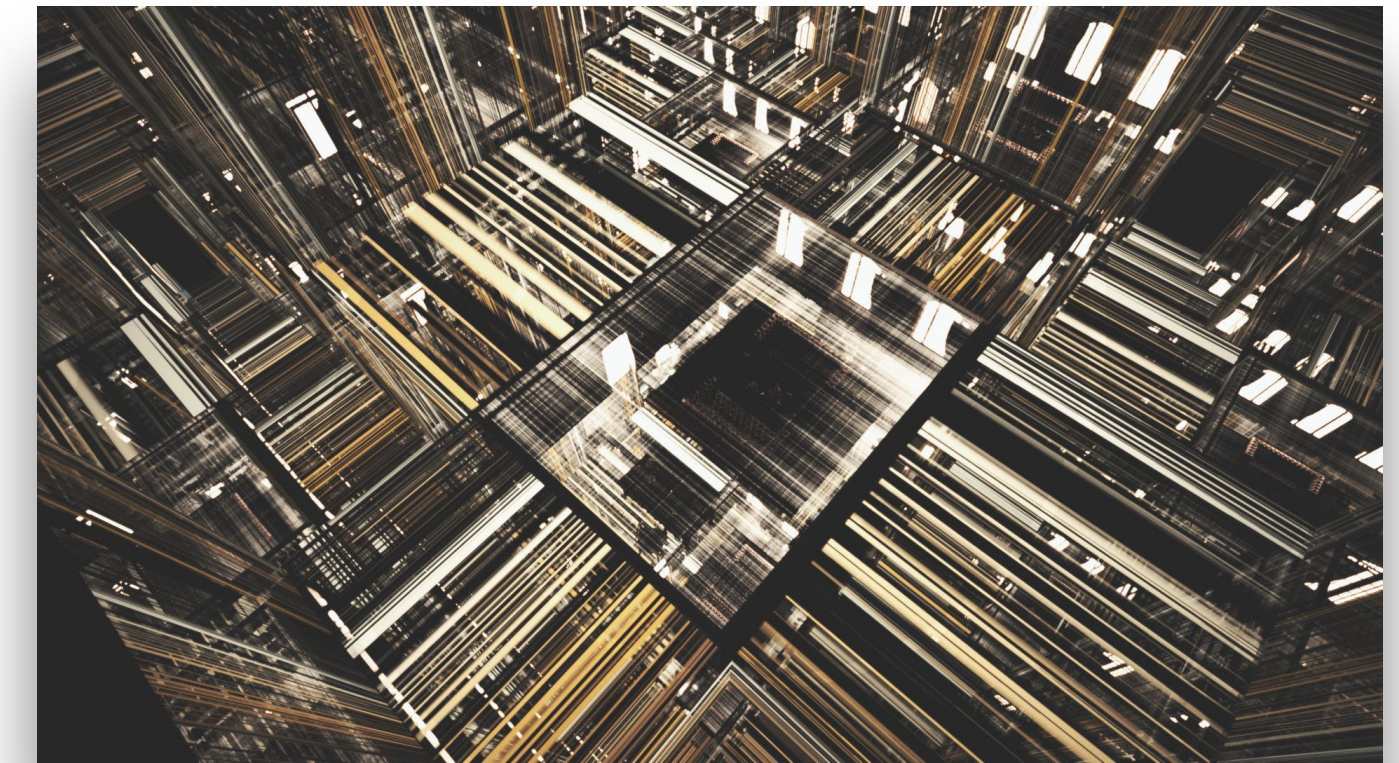




# Scaling up this idea

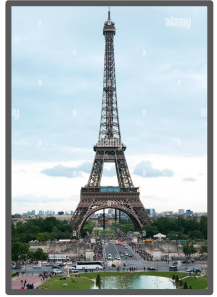
Usually, the features have 512 or more dimensions, not just 3

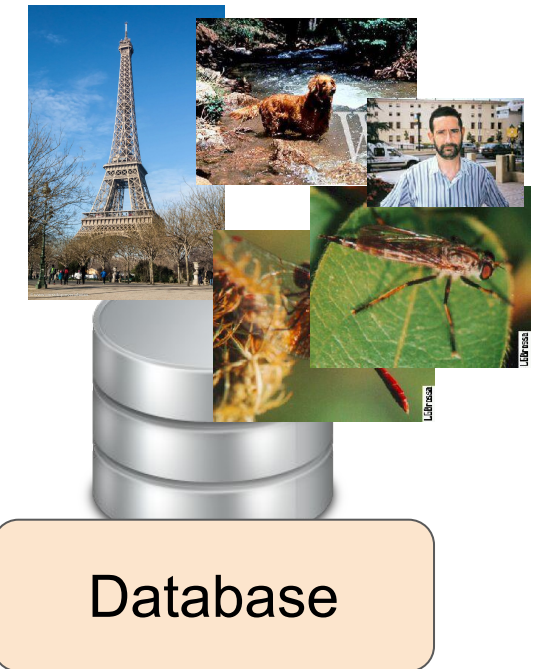
- Difficult to visualize for us, that we live only in a 3-dimensional world
- Mathematically, this is possible without loss of generality
  - The Euclidean distance is still defined
- High-dimensional representations carry more information



$$d = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

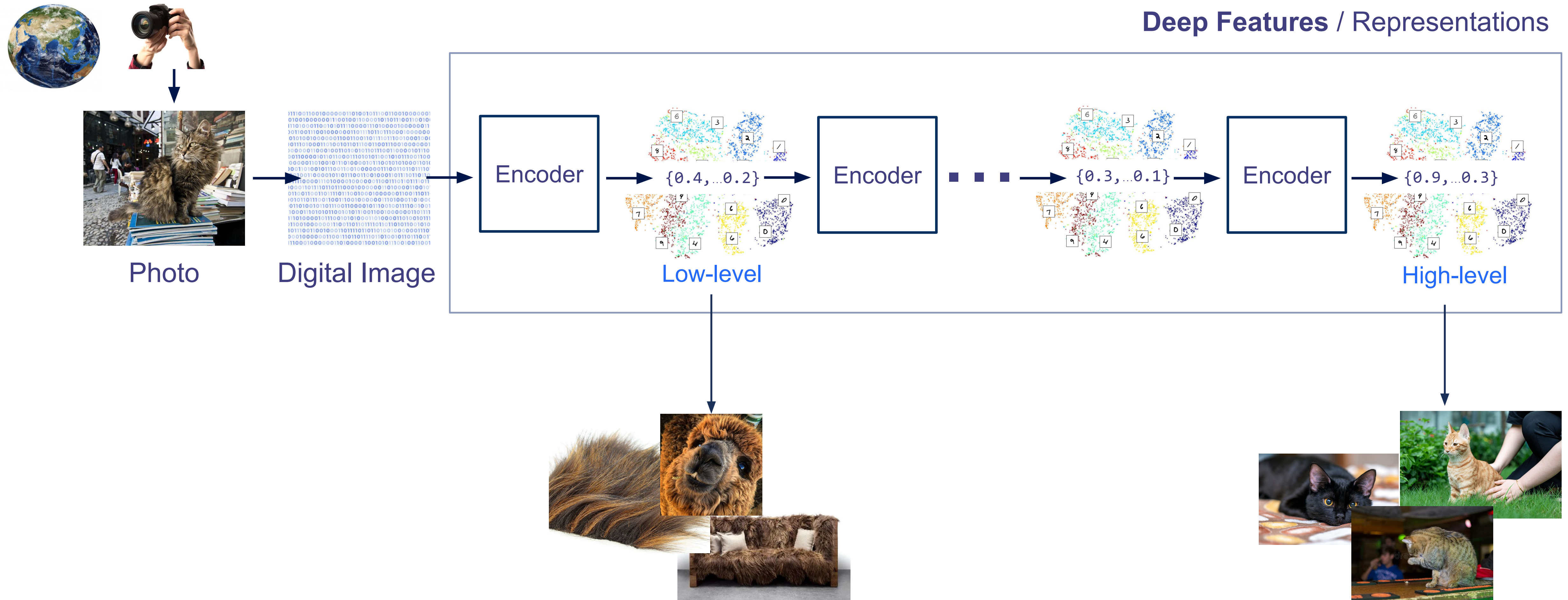
# Operatively

1. Compute representations from all the images in the database  $D = \{I_1, I_2, \dots, I_n\}$
2. Compute representation from the query image  $I_q$  
3. Compute the Euclidean distances between  $q$  and all the images in  $D$
4. K-nearest-neighbor search: sort these distances in decreasing order (or, in other words, by increasing similarity) and take the first  $k$  results



# Different representations, different similarities

By taking representations at different layers of the deep network, we give a different meaning to our similarity measure



# Elements of text-to-image search

# Using text as a query

- As of now, we used an image as query for retrieving other images
  - image → image
- What about using another modality as a query?
  - text → image

image → image



text → image

*“A football player kicked the ball”*



# Google image search

The screenshot shows a Google search interface with the query "a football player kicking the ball" in the search bar. Below the search bar are navigation tabs for "Tutti", "Immagini", "Video", "Notizie", "Shopping", and "Altro". There are also icons for camera, voice search, and a magnifying glass. On the right side, there are settings, a grid icon, a profile picture, and "Raccoglie" and "SafeSearch" options. Below the search bar are filters for "drawing", "vecteezy", "art", and "sports". The search results are displayed in a grid of 24 images, each with a small thumbnail icon in the bottom-left corner and a caption below it. The captions include sources like "vectorstock.com", "photos.com", "dreamstime.com", "shutterstock.com", "123rf.com", "istockphoto.com", "dribbble.com", "nonprofitlawblog.com", "pinterest.com", "depoaltphotos.com", and "vecteezy.com".

# Advantages

- In many cases it is not convenient to search using an image as a query
  - Imagine to always search images in Google using other images
- The natural language is natively less ambiguous than an image



I see a cat with a long fur over some books



I see a cat near a library in the street

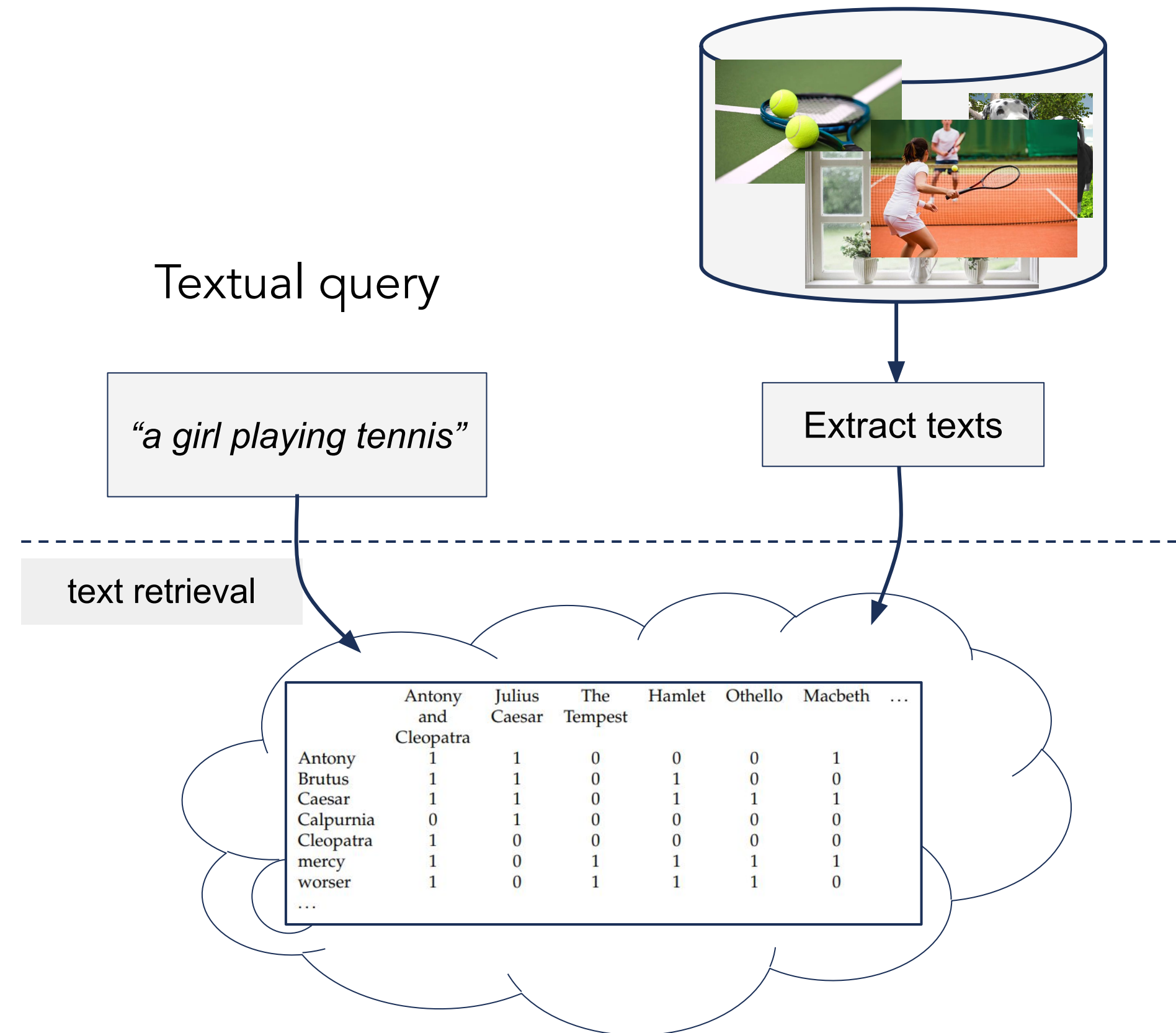


I see some persons in the street walking behind a cat



# A trivial solution

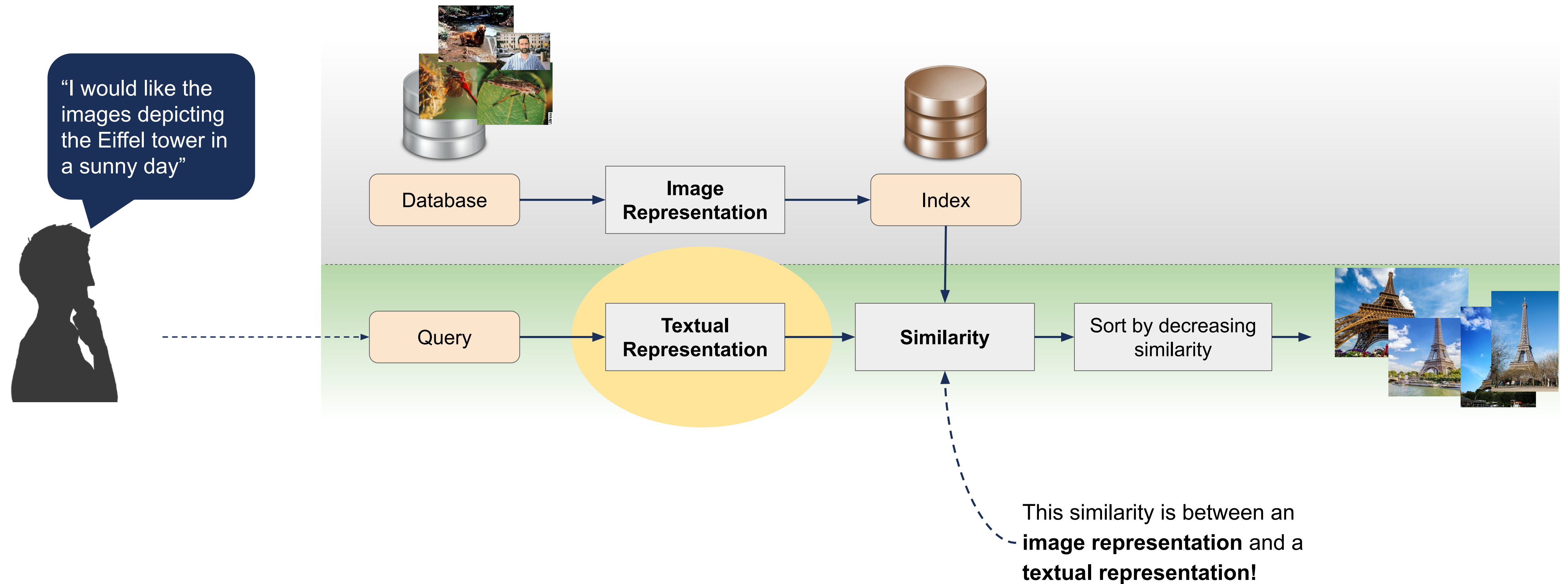
Use the textual metadata associated to the image (e.g., in the *alt* text) to perform a textual search



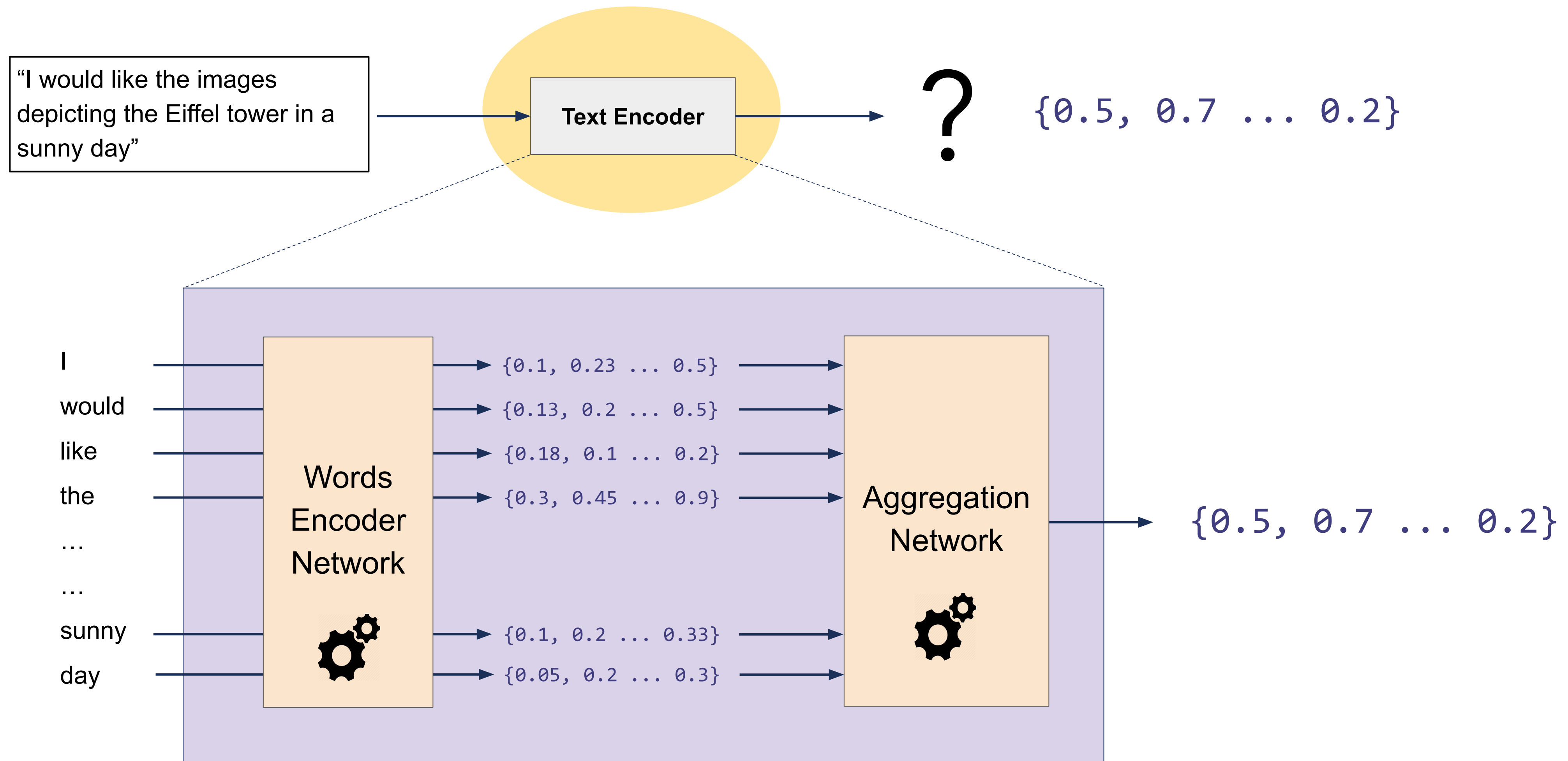
- Google partially works in this way
- What if we don't have textual descriptions for images?



# Text-to-image similarity search



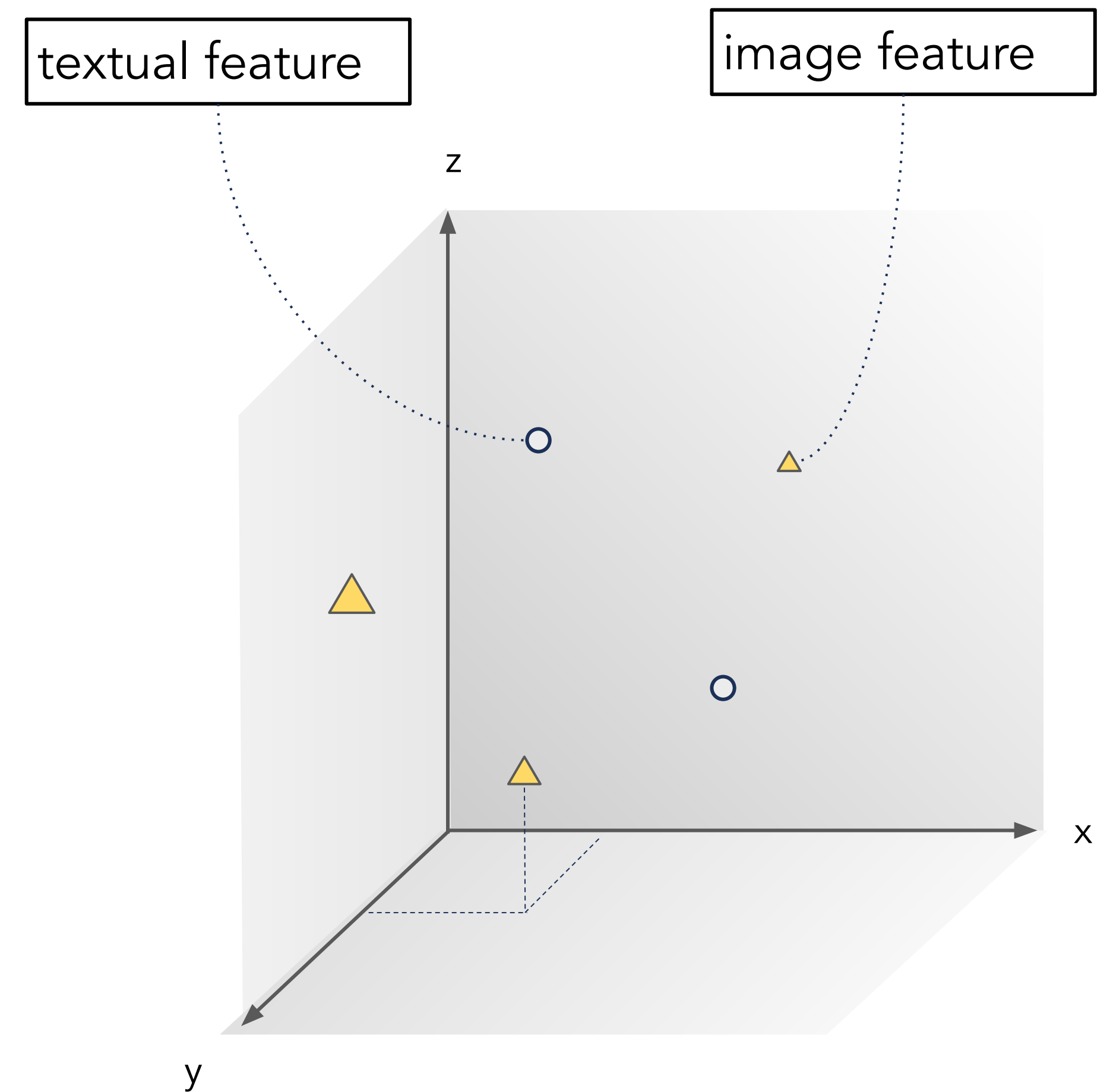
# A deep network for texts!



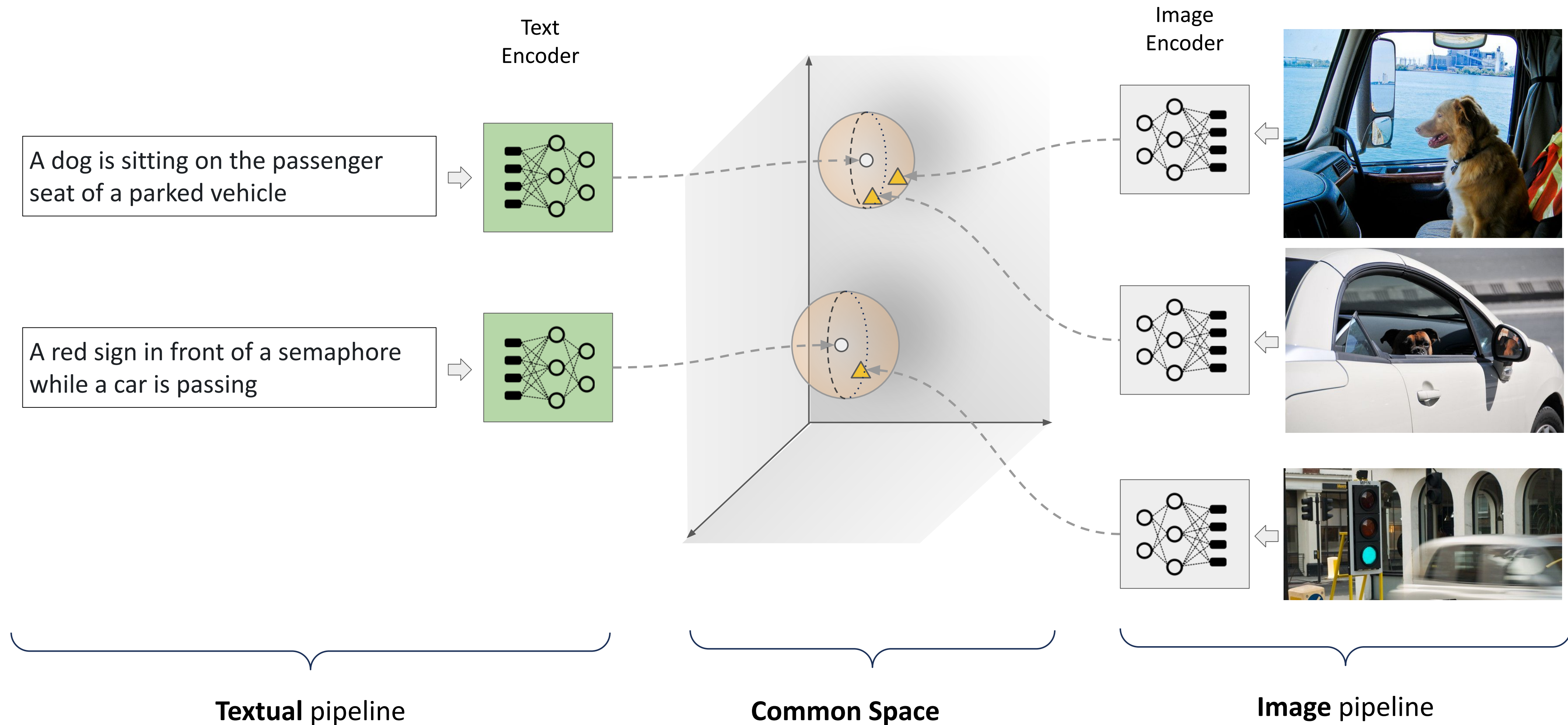
# The similarity computation

If the image representations and the textual representations have the same dimensions, they can be compared in the same space!

- We reuse the similarity framework developed before for image-image searches
- We can compute the Euclidean distance between textual and image representations

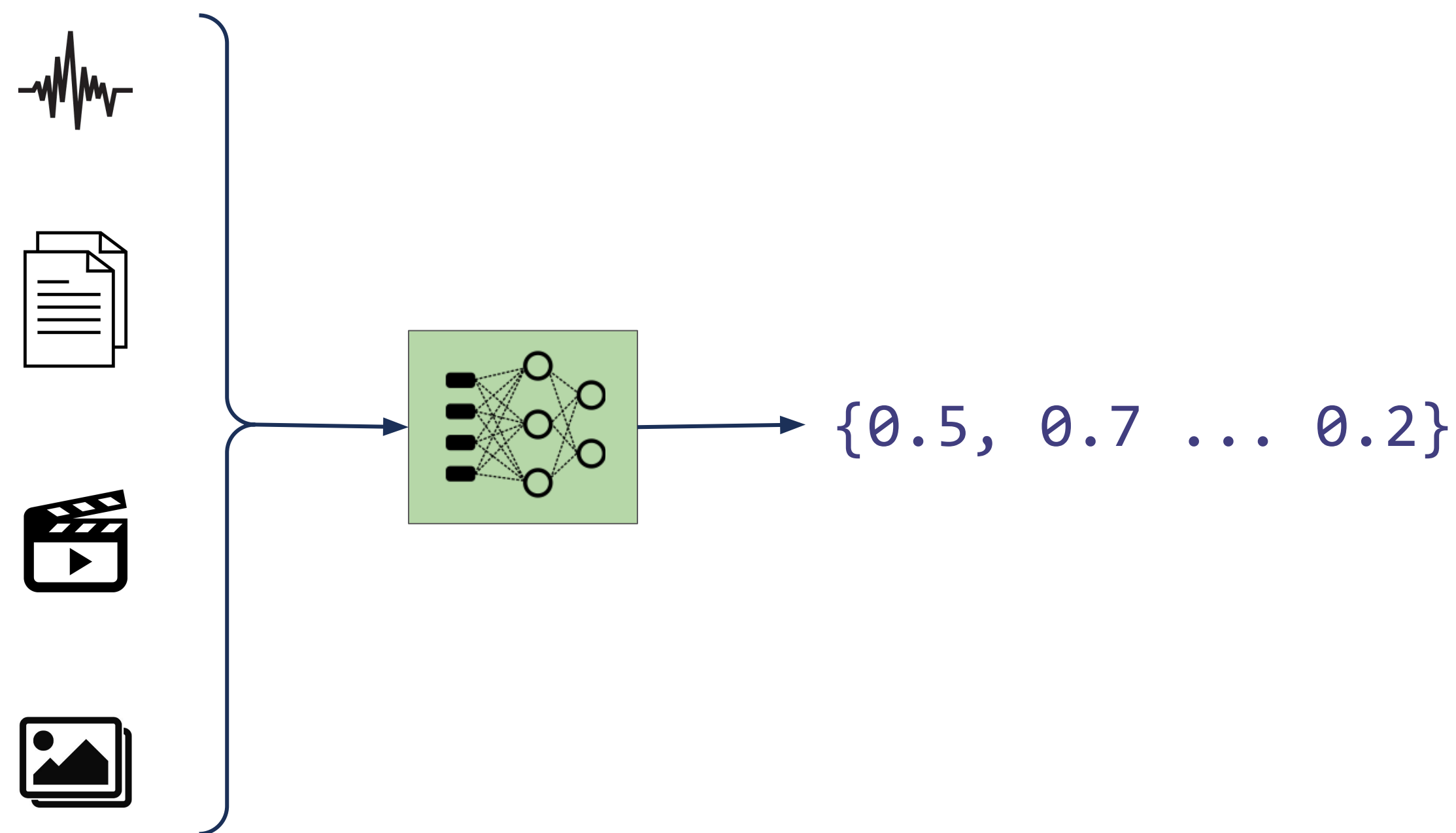


# The common visual-textual space



# Generalizing to multiple modalities

Every multimedia object can be converted in a numerical representation using deep neural networks



Using k-NN search in a common space, we can easily perform:

text  $\rightarrow$  image

text  $\rightarrow$  text

text  $\rightarrow$  video

image  $\rightarrow$  text

audio  $\rightarrow$  text

text  $\rightarrow$  audio, video

# The end!

## Questions?

