

# COURSE: APPLICATIONS OF ARTIFICIAL INTELLIGENCE (AAI)



AMP in Business Analytics 2020

INDIAN SCHOOL OF BUSINESS

Vamshi Ambati

Vamshi\_Ambati@isb.edu;

vamshi@predera.com

Session 4: Computer Vision

March 1, 2020

# USE CASES OF APPLICATION AREAS

- Natural Language Processing -> IE, IR, MT, NLU, NLG, Dialog -> Chatbot, Customer complaint translation
- Speech Processing -> SR (S2T), SS (T2S), SG (S2S) -> Voice Assistant (Siri), Education, IVR
- Computer Vision -> IU, IG, OR, OD, Video {U,G,R,D} -> OCR, Face Recognition, Self-driving, Entertainment
- Predictive Analytics -> Prediction, Recommender, Forecasting Systems -> E-commerce , Fraud detection
- Robotics -> Locomotion, Mechanics, Sensors, Planning -> Medical Bots, Agriculture, Military
- Agent Systems -> Planning, Reinforcement Learning, Space Optimization -> AlphaGo

# THIS COURSE

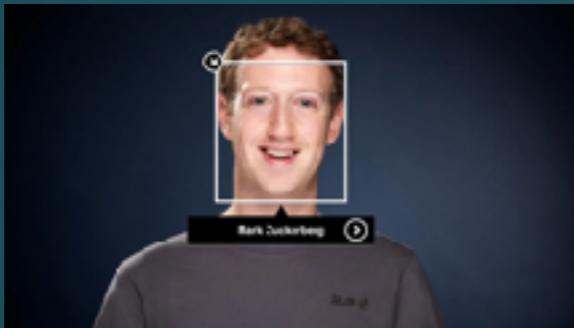
- Session 1: Machine Translation
- Session 2: Dialog Systems (Chatbot)
- Session 3: Recommender Systems
- Session 4: Computer Vision (Image | Video)

# AGENDA

- Introduction
- Image Processing
  - Traditional Approaches
  - Deep Learning
  - Case study - nutrition label extraction
- Video Understanding
  - Traditional Approaches
  - Deep Learning
- Summary & Challenges

# 1. INTRODUCTION

# Introduction

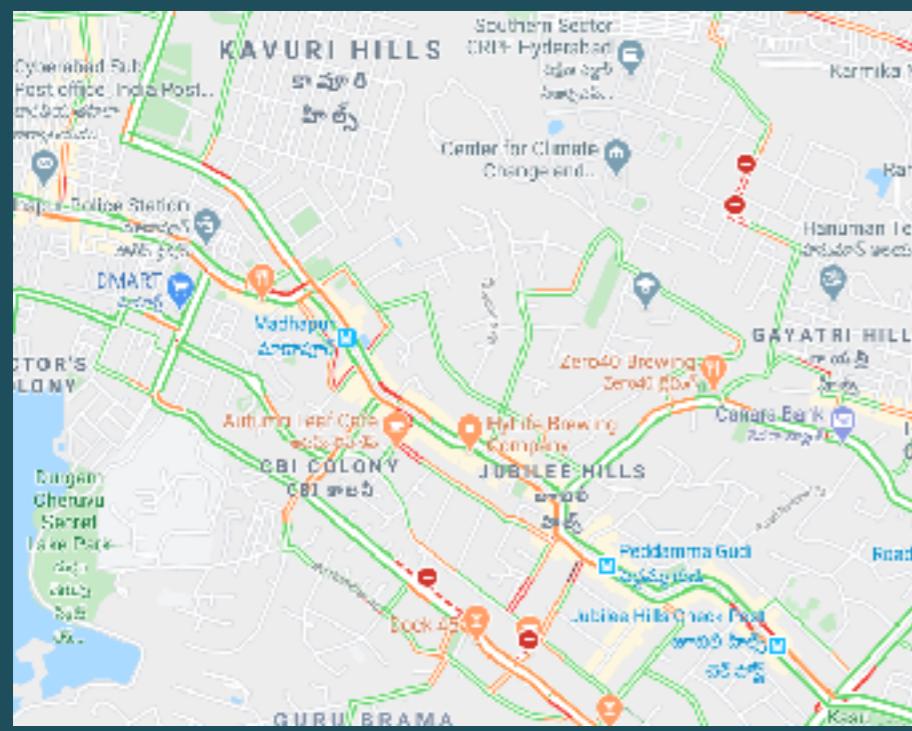
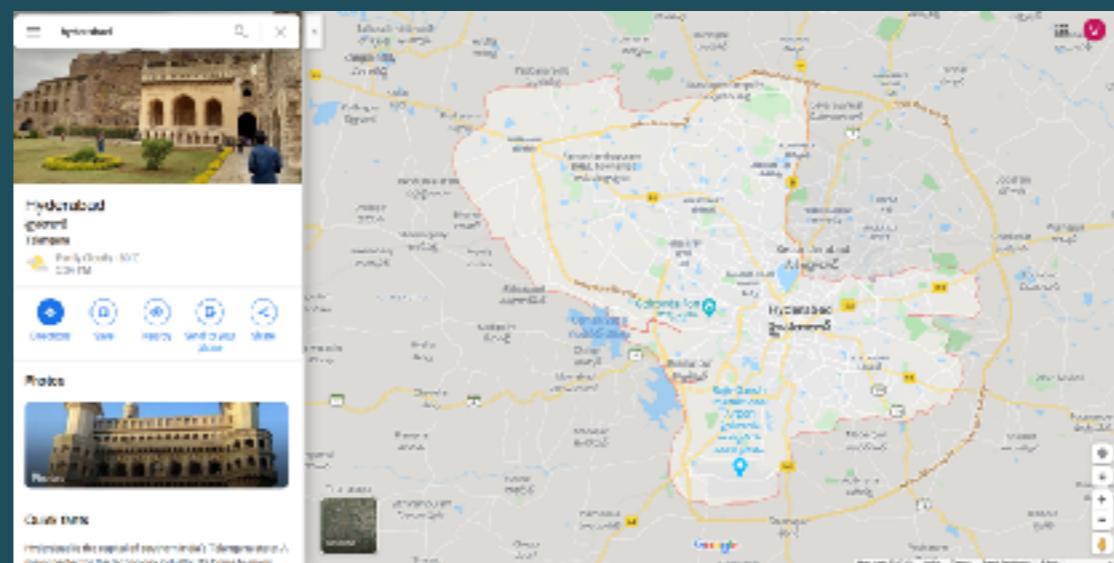
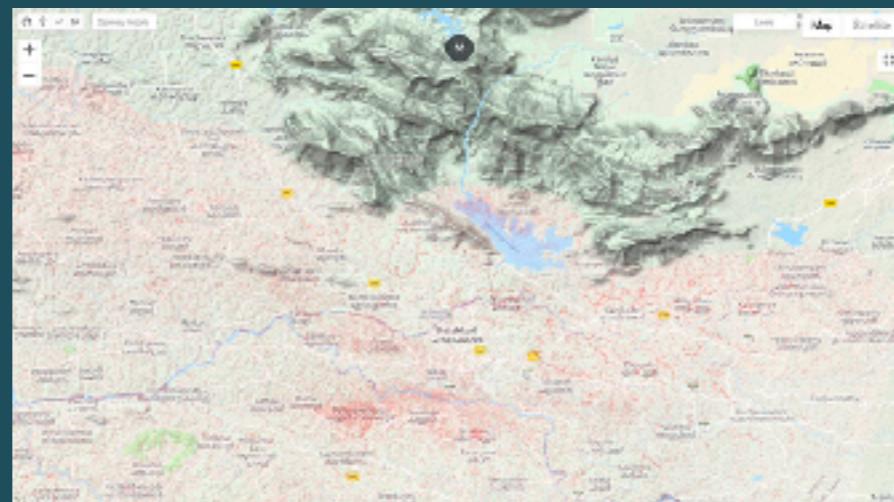


Walmart						
Save money. Live better.						
C 0000 0000 - 0000 MANAGER ELANA BARBERET 241 WILMINGTON DR 00000000						
NEW BRITAIN, CT 06051-0000 STN 00115 OPR 005004 TEL 44 0000 61801 INTL TEL 004123181866 3.83 X FLORRY BIRBY 004123181446 3.83 X SARLOKARNEE S 004123181219 4.83 X J-S SQUEAK 004499800224 8.83 X BUNNY THERAPY 004123181595 1.83 X DOG TREAT 004123181851 2.83 X PUP PCH 1 002310011062 6.83 X PUP PCH 1 002310011062 6.83 X COUPON 2100 002310010709 1.00 X BUNNY KERMIT 004123181837 P 3.83 X PRENUO DOBBO 004532859555 P 3.83 X J-COONES 001488844101 P 8.83 X KATE CLEOPATRA 003144333382 P 3.83 X COLLIE 00449980063687 3.83 X CALICO 005104494062 P 2.83 X BB TOW HWT 003103105844 26.33 X SPONGEBOB 0015594679436 2.83 X STICKERLAB 001886479436 2.83 X STICK STICKLAB 0015594679429 6.23 X STICK STICKLAB 0015594679435 6.23 X STICK STICKLAB 001886479429 6.23 X STICK STICKLAB 0015594679419 6.23 X BLIND BEARD 004534859569 P 6.23 X GREAT VALUE 007874253191 P 9.87 X LIFTON 001200011224 P 4.48 X DRY DOG 002310011035 12.44 X SUBTOTAL 93.62 X TAX 1 6.750 4.59 X TOTAL 98.21 X VISA TEND 58.21 X						



Covid-19 from lung CT scans : <https://www.wired.com/story/chinese-hospitals-deploy-ai-help-diagnose-covid-19/>

# Google Maps



# History of Image Processing

- Origin at Bell Laboratories, American Jet Propulsion Laboratory, MIT
- Early purpose was to improve quality of images, owing to low resolution cameras
- Applications in space exploration and radiology
- Early success in 1964 through understanding surface of moon by improving quality of images sent back by spacecrafts

# History of image processing

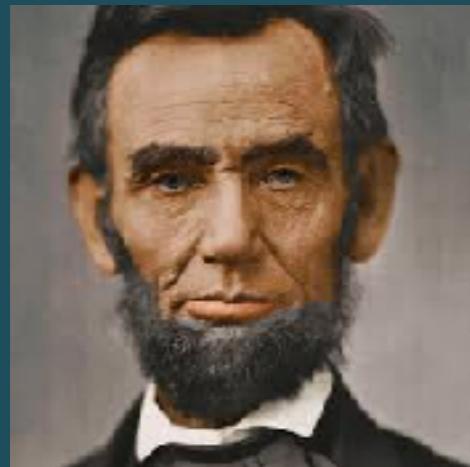
- 1970s- breakthroughs in computing, image compression
- 1980s- Contour modeling, texture understanding
- 1990-now: Machine learning and Deep learning
- Real time object detection, video analytics
- Today - security, retail, defence, art, healthcare, public policy, entertainment, education, space exploration, environment etc.

# Challenges

- Images of varied quality
- Size of objects of interest in images varies
- Background noise
- Varied colours
- Orientation
- Labeling
- Computing



# Cross section of an image

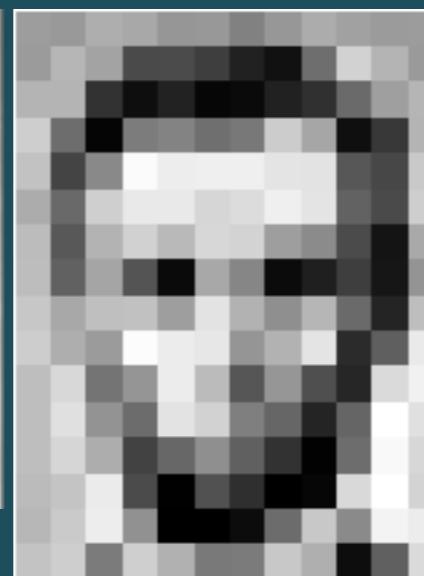
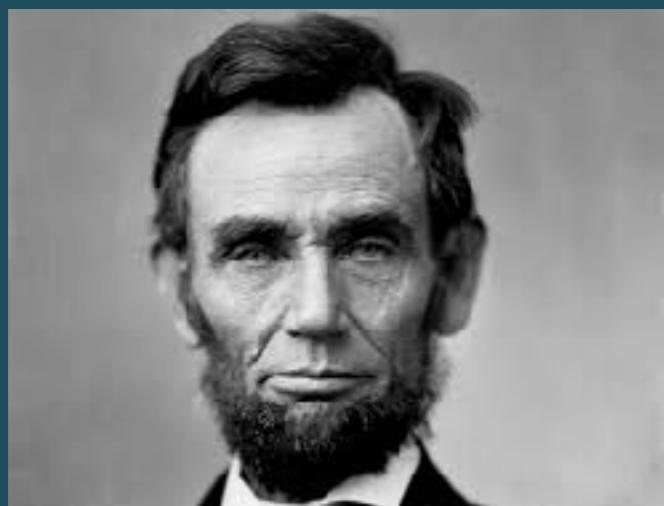


Resolution 1024X1024

Size = 1024X1024 \* 3 {R,G,B}

Pixel value range: 0-127 or 0-255

165	187	203	58	7	
14	120	205	20*	98	159
253	142	120	251	47	147
67	100	92	241	23	185
209	118	124	27	59	201
210	236	105	139	13	210
35	78	99	37	4	14
115	04	34	11	15	156
32	59	231	223	74	



157	152	154	158	150	152	129	121	172	161	155	150
156	155	154	152	153	156	153	157	153	150	152	154
150	146	150	14	14	14	5	16	23	48	166	151
205	166	6	124	121	111	126	204	166	16	66	120
154	48	117	267	272	263	209	229	207	21	13	211
172	146	567	239	230	214	226	229	229	90	54	230
188	45	170	306	308	216	311	168	130	76	56	156
159	31	165	54	56	159	154	11	20	62	22	356
199	146	191	193	156	257	179	143	108	185	56	190
205	174	186	267	236	209	148	138	208	45	55	254
180	216	138	149	236	167	46	150	79	26	218	281
153	224	147	104	327	213	137	102	56	191	268	221
193	214	173	54	193	143	26	50	2	165	549	215
187	146	323	73	1	86	97	6	6	217	154	211
189	205	227	184	6	9	15	108	200	148	143	236
185	206	133	207	177	130	121	200	175	79	56	219

# Video- stream of images



Zoetrope



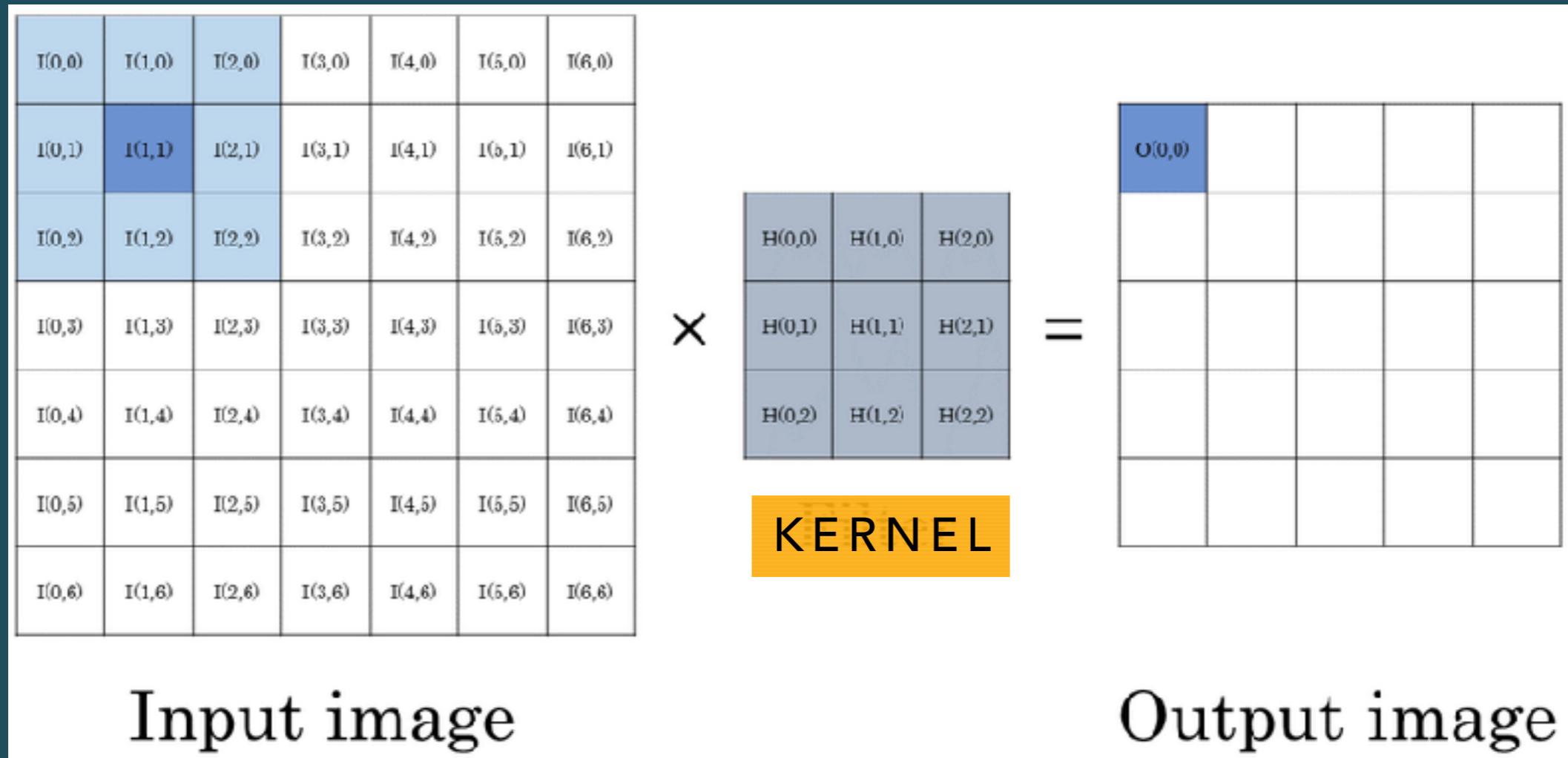
Flipbook

Video has frame rates of 24 to 30 per second

# Image Processing Tasks

- Image augmentation
- Image enhancement/ blurring
- Noise removal
- Boundary detection

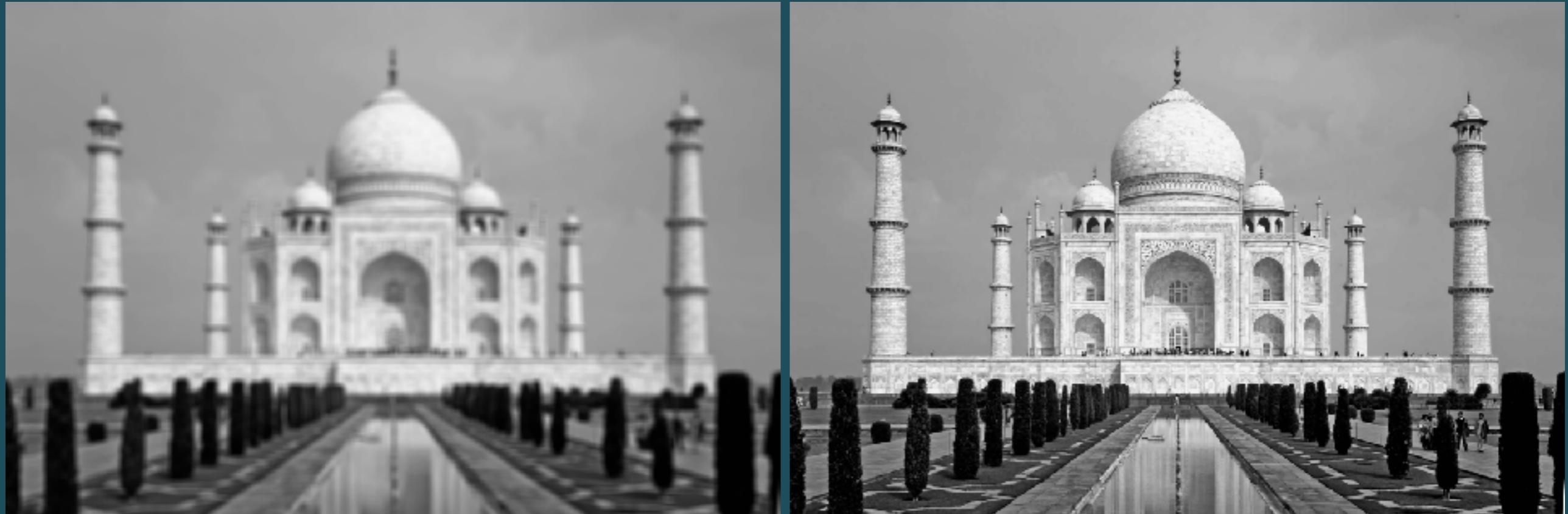
# Convolution for Image Processing



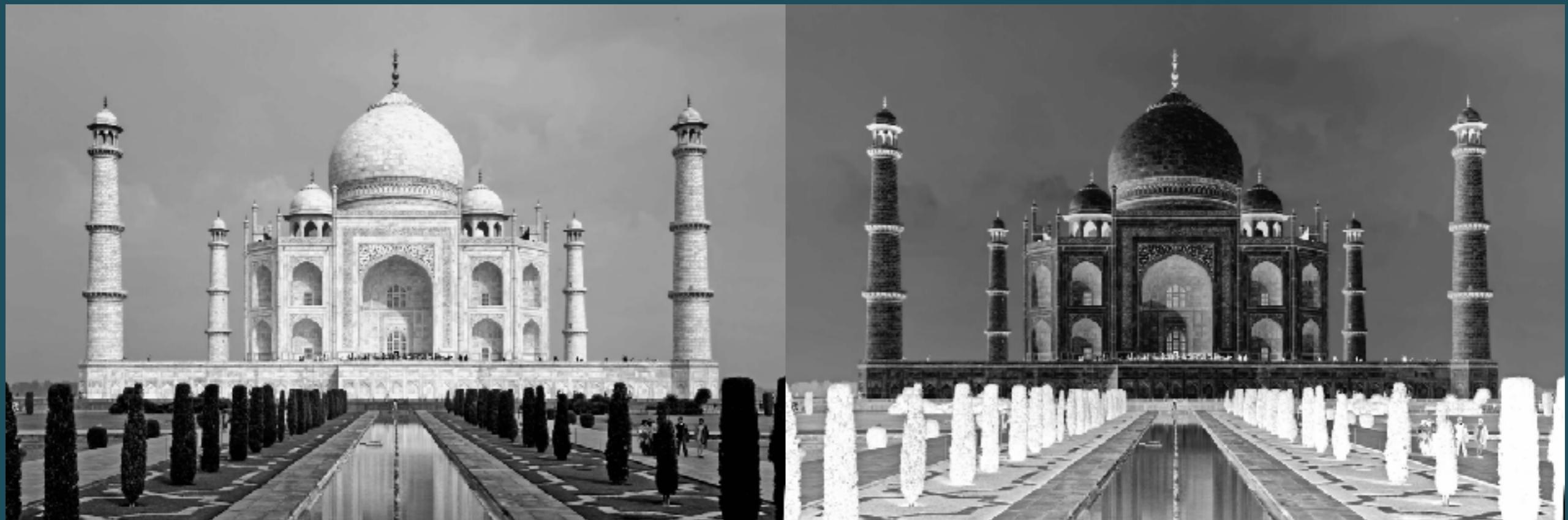
Convolution provides a way of 'multiplying together' two arrays of numbers to produce a third

# Image enhancement: De-blurring

- Laplacian kernel



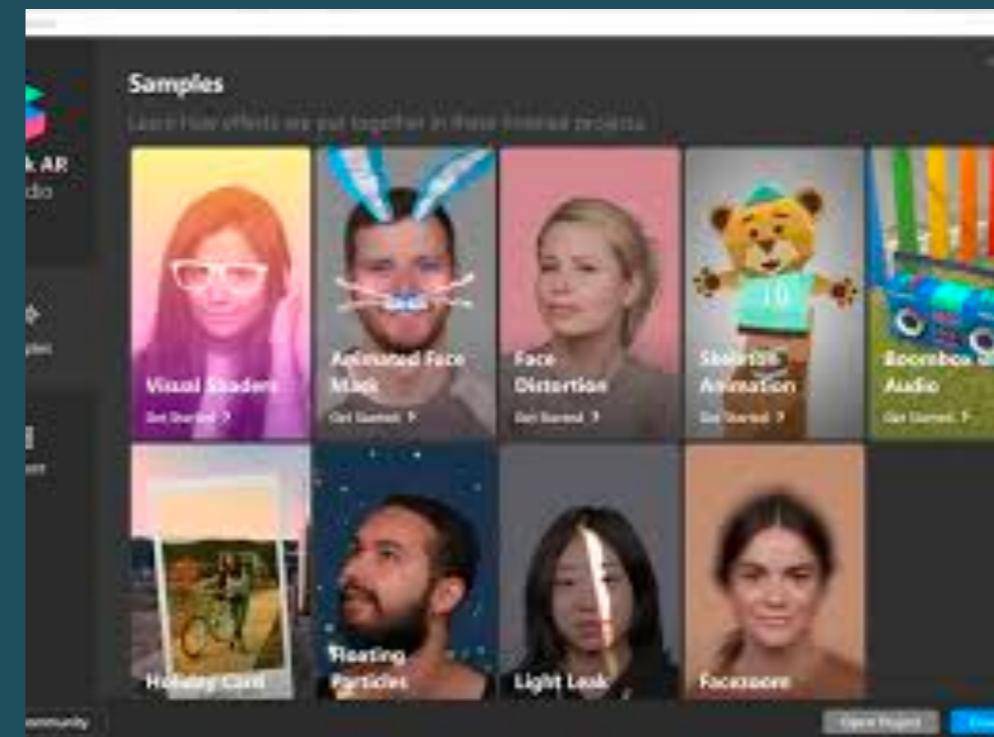
# Image enhancement: Contrast Improvement



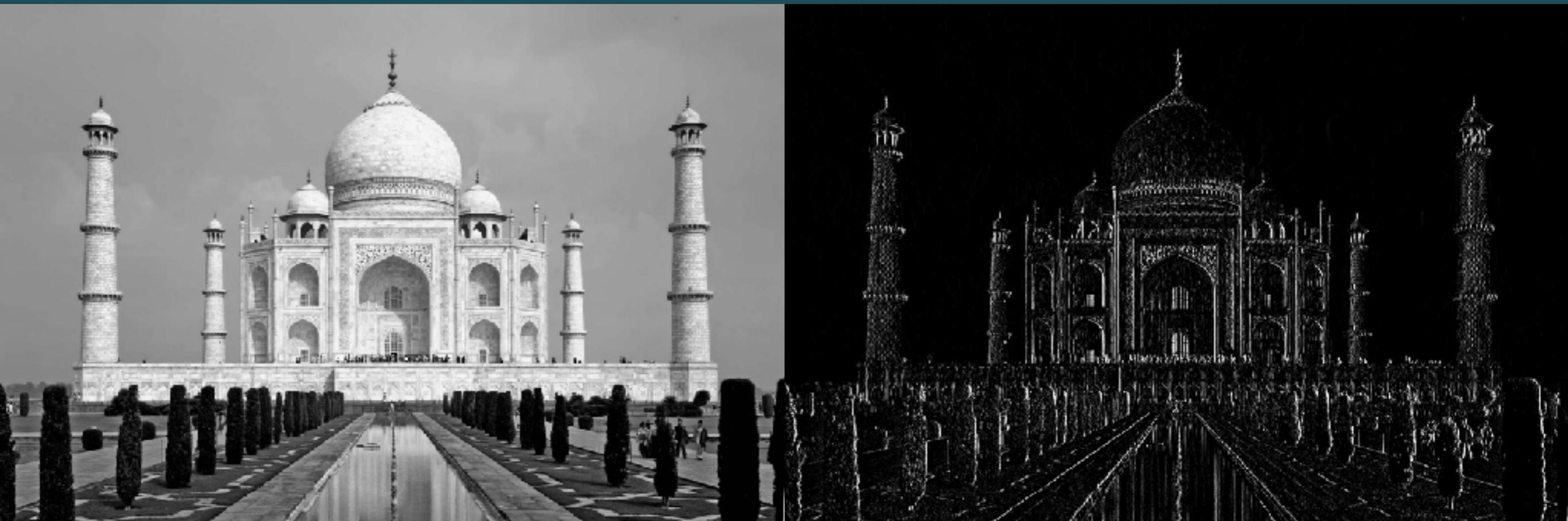
# Image enhancement: Noise removal



# More complex filters



# Edge detection



# Edge detection

## Sobel Filter/Kernel

X – Direction Kernel

-1	0	1
-2	0	2
-1	0	1

Y – Direction Kernel

-1	-2	-1
0	0	0
1	2	1

# Going beyond Image Tasks (Math gets hairy)

## Image Processing tasks

- Image augmentation
- Image enhancement/  
blurring
- Noise removal
- Boundary detection

## Computer vision

- Face recognition
- Object detection
- Text detection- OCR
- Image classification
- Video processing

# Recognition as Classification

## Face Recognition

- Identify human
- Segment face
- Align position
- Classify face

## OCR

- Quality of page
- Extract text block
- Extract word block
- Extract letters
- Classify letter

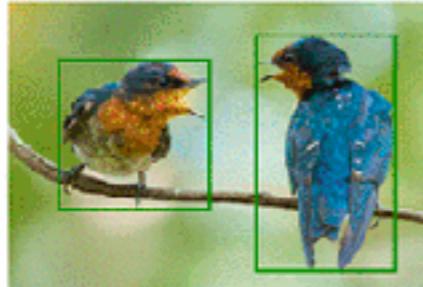
# Image Classification as Supervised ML

**Task type** [Info](#)

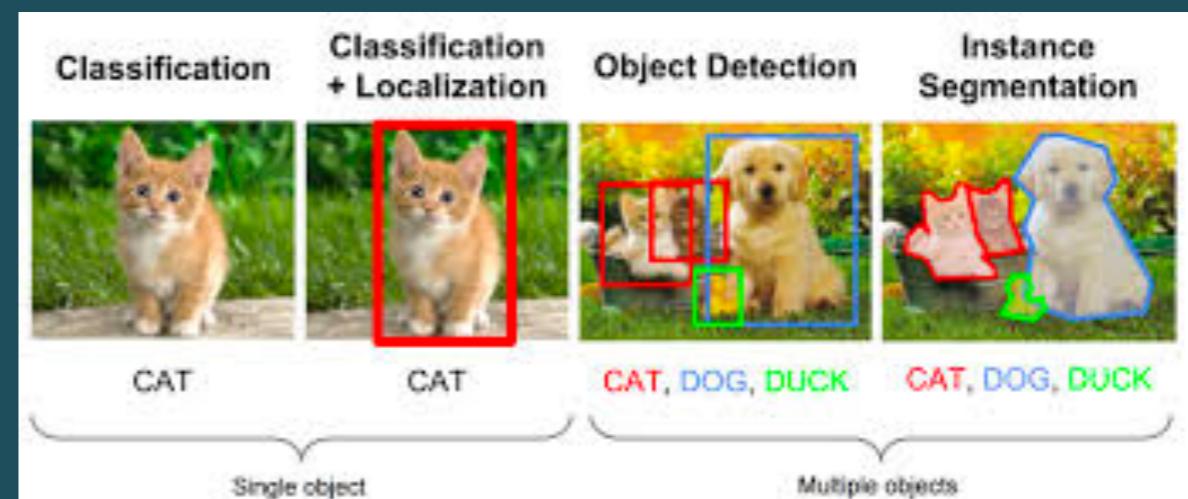
**Task category**  
Select the type of data being labeled to view available task templates for it or select 'Custom' to create your own.

**Image**

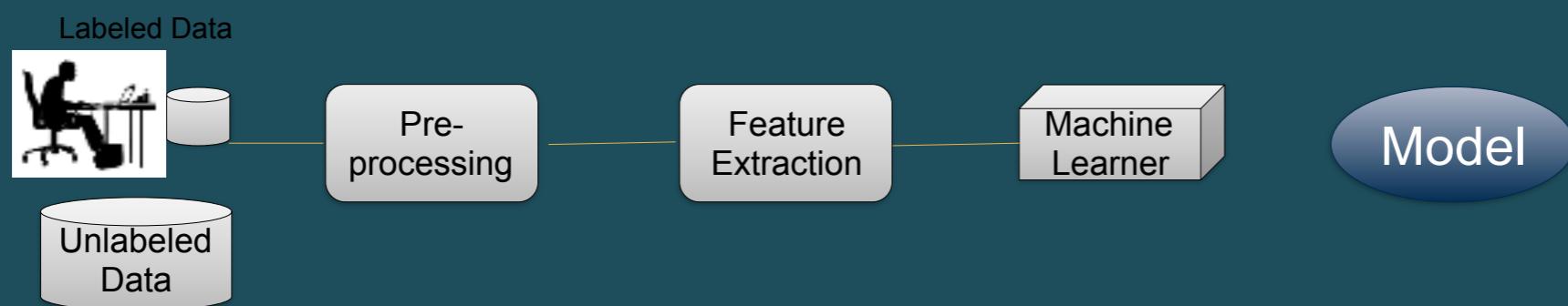
**Task selection**  
Select the task that a human worker will perform to label objects in your dataset.

- Image classification**  
Get workers to categorize images into specific classes. [Info](#)  
  
 Basketball  
 Soccer
- Bounding box**  
Get workers to draw bounding boxes around specified objects in your images. [Info](#)  

- Semantic segmentation**  
Get workers to draw pixel level labels around specific objects and segments in your images. [Info](#)  

- Label verification**  
Get workers to verify existing labels in your dataset. [Info](#)  
 Correct label  
 Incorrect label  

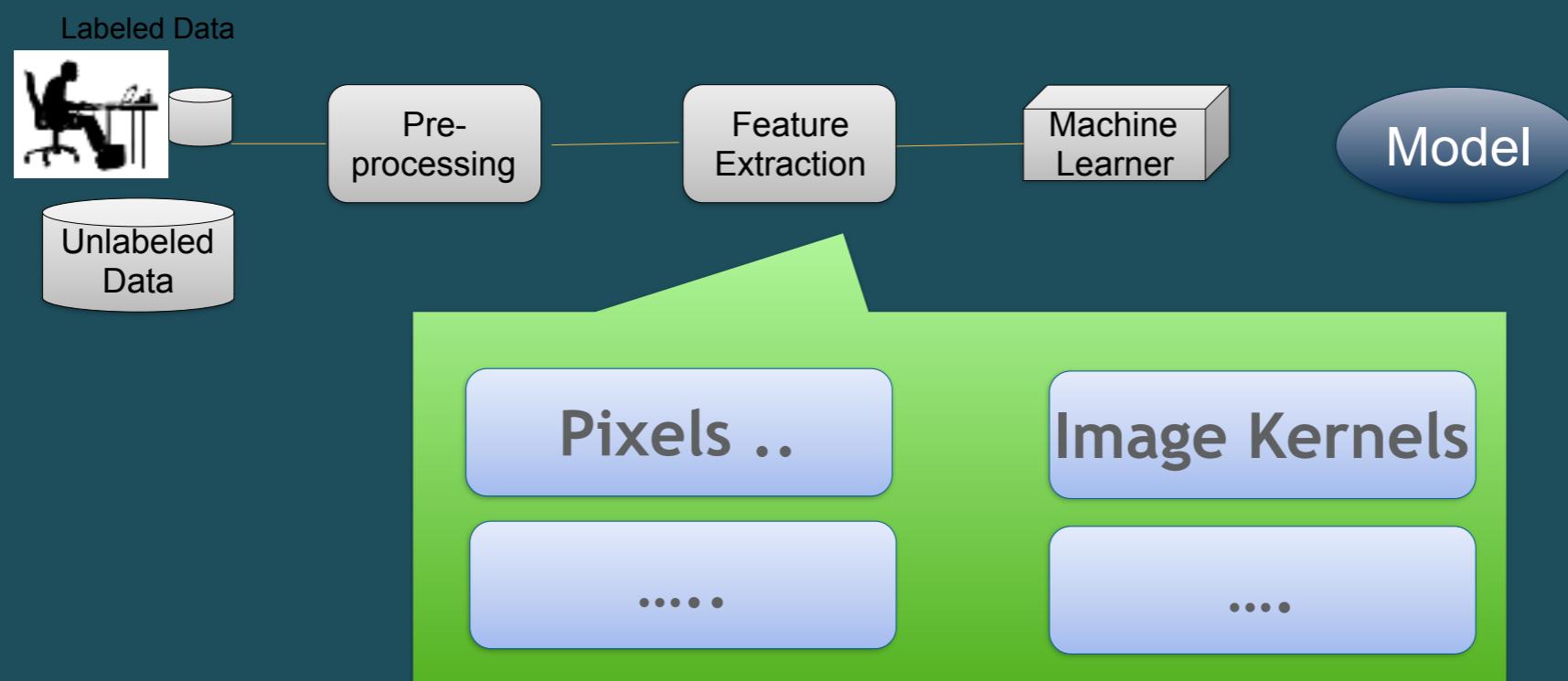



# MACHINE LEARNING FOR FACE RECOGNITION

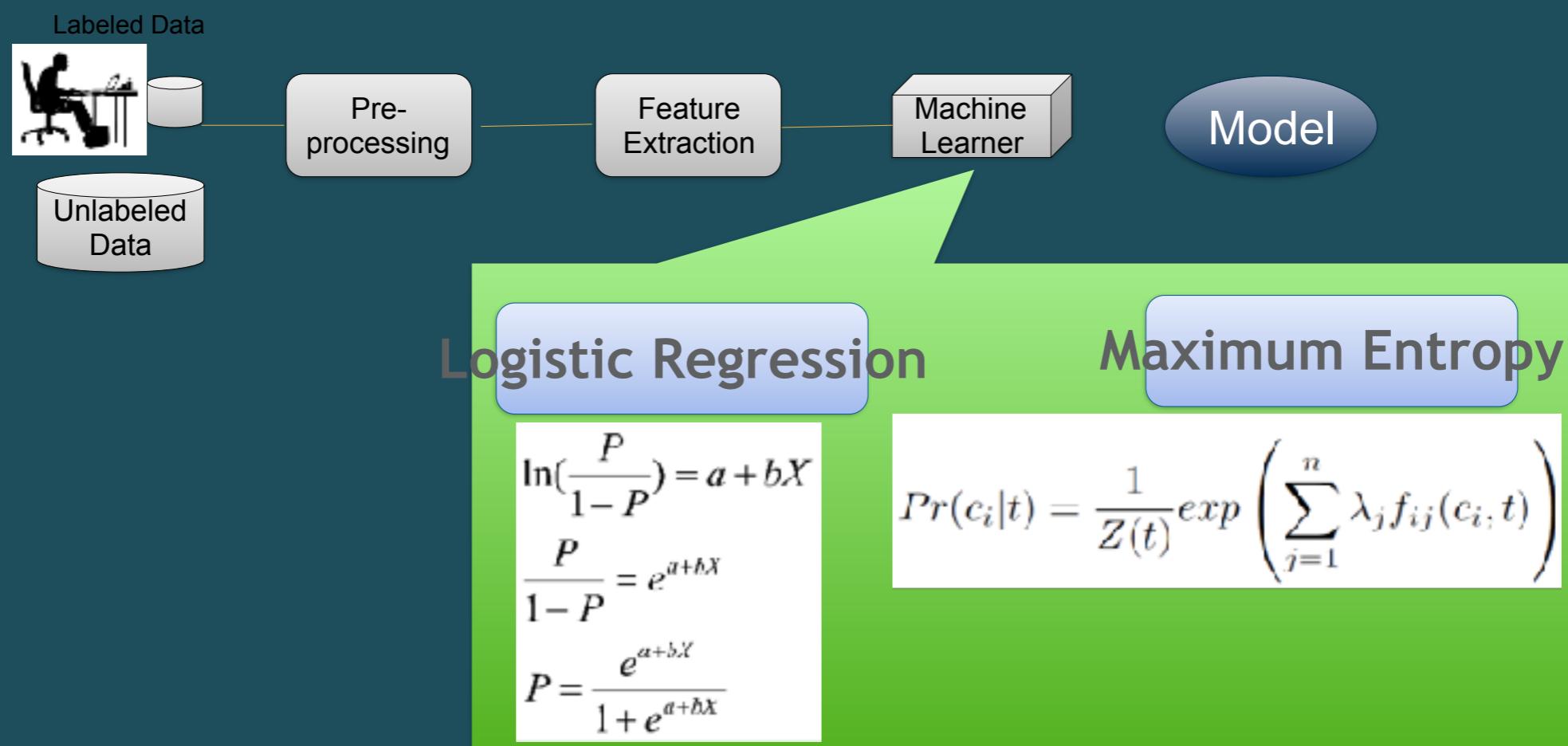


Nancy

# MACHINE LEARNING FOR FACE RECOGNITION

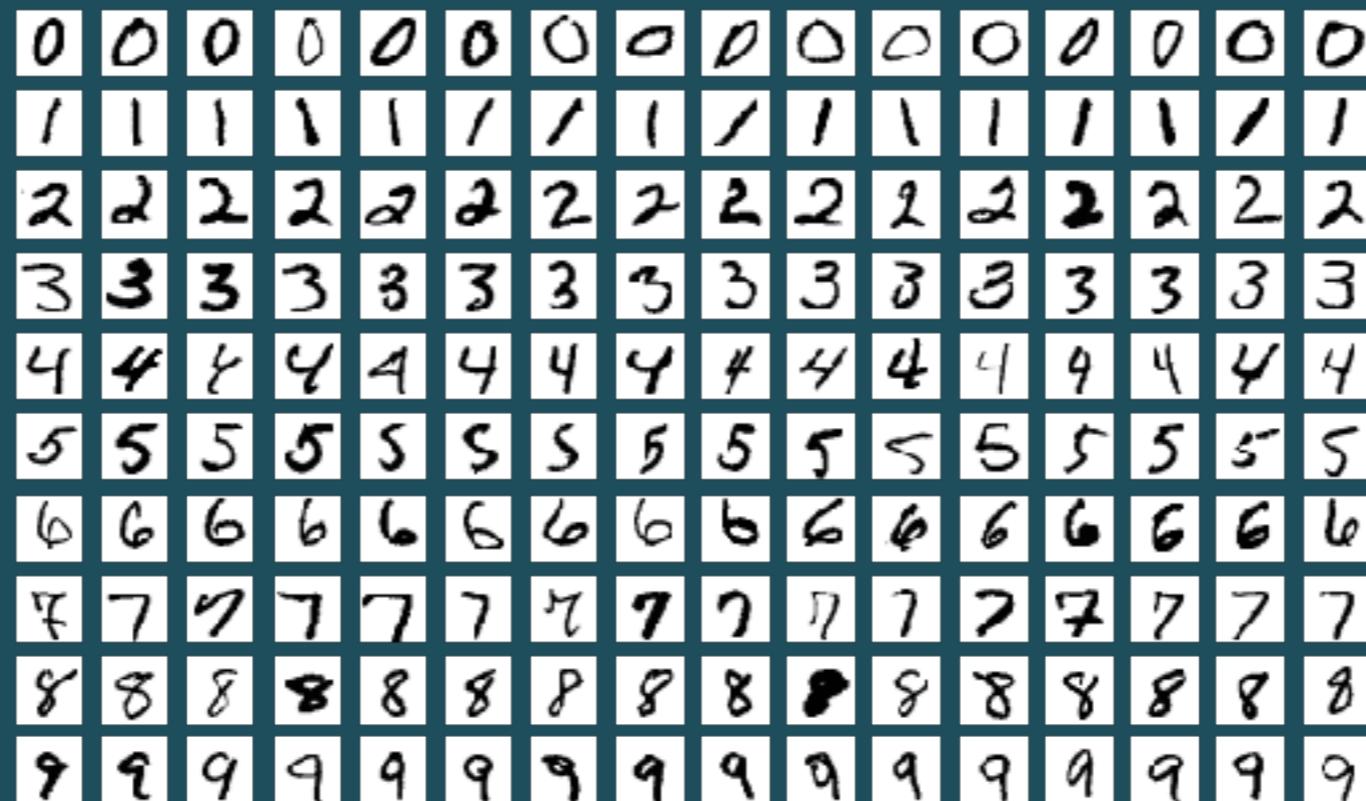


# MACHINE LEARNING FOR FACE RECOGNITION



# Image Classification Dataset (MNIST)

- Modified National Institute of Standards and Technology database
- Identify handwritten digits
- 28\*28, 60,000 images



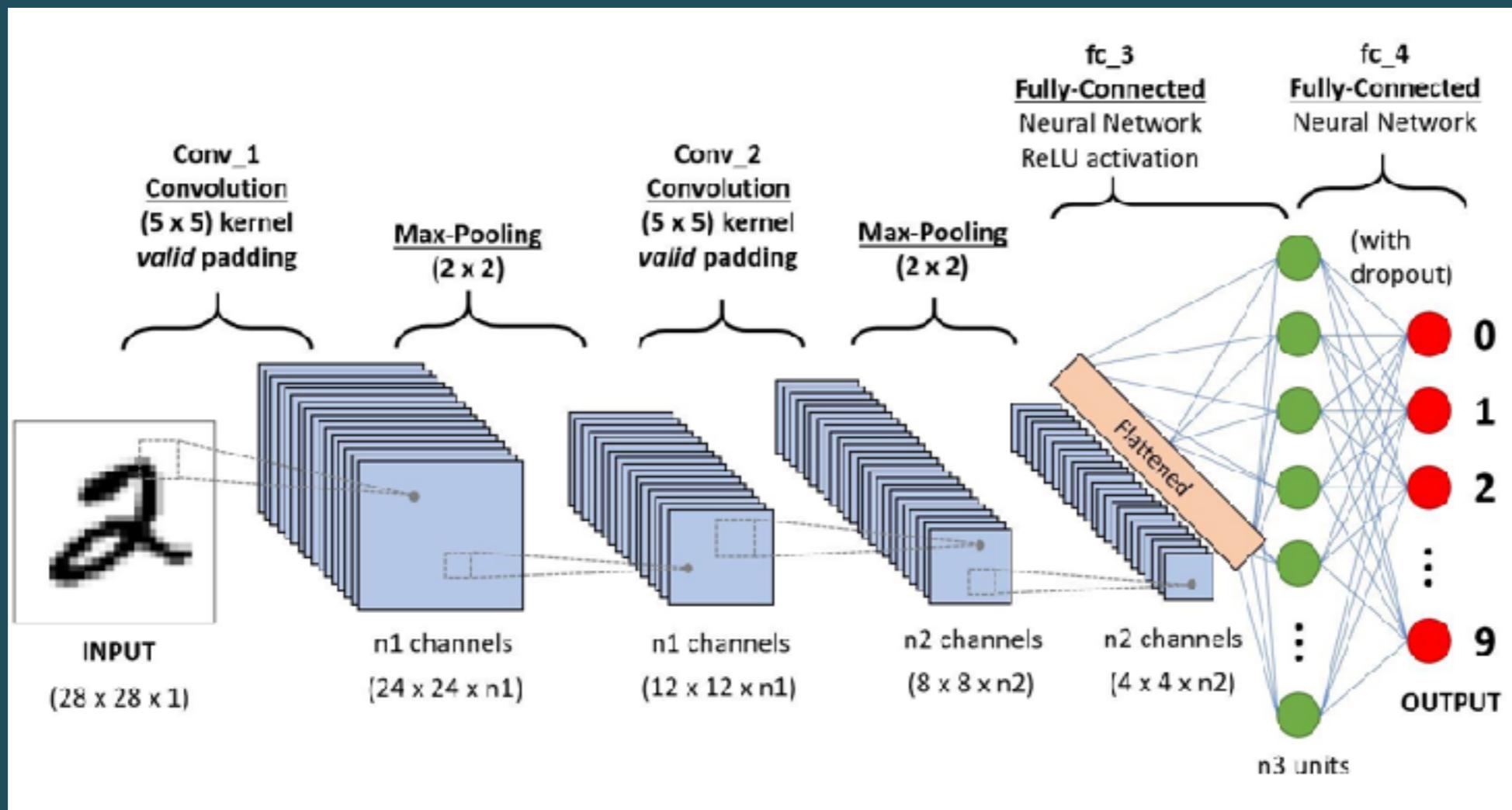
# MNIST Classification

- Machine learning way
    - Each pixel as a feature

pixel_42	pixel_43	pixel_44	pixel_45	pixel_46	pixel_47	pixel_48	pixel_49	pixel_50	pixel_51	pixel_52	pixel_53	pixel_54	pixel_55	pixel_56	pixel_57	pixel_58	pixel_59	pixel_60	pixel_61	pixel_62	pixel_63	target
11	0	1	12	7	0	0	2	14	5	10	12	0	0	0	0	6	13	10	0	0	0	0
1	15	16	6	0	0	0	0	1	16	16	6	0	0	0	0	0	11	16	10	0	0	0
16	16	5	0	0	0	0	3	13	16	16	11	5	0	0	0	0	3	11	16	9	0	0
0	0	1	10	8	0	0	0	8	4	5	14	9	0	0	0	7	13	13	9	0	0	3
15	16	13	16	1	0	0	0	0	3	15	10	0	0	0	0	0	2	16	4	0	0	4
0	0	4	16	9	0	0	0	5	4	12	16	4	0	0	0	9	16	16	10	0	0	5
13	16	13	16	3	0	0	0	7	16	11	15	8	0	0	0	1	9	15	11	3	0	6
0	16	5	0	0	0	0	0	9	15	1	0	0	0	0	0	13	5	0	0	0	0	7
16	8	10	13	2	0	0	1	15	1	3	16	8	0	0	0	11	16	15	11	1	0	8
0	3	0	9	11	0	0	0	0	0	9	15	4	0	0	0	9	12	13	3	0	0	9
16	5	1	11	3	0	0	0	12	12	10	10	0	0	0	0	1	10	13	3	0	0	0
0	5	16	15	0	0	0	0	0	4	16	14	0	0	0	0	0	1	13	16	1	0	1
0	0	15	0	0	0	0	0	9	16	15	9	8	2	0	0	3	11	8	13	12	4	2
0	0	2	15	4	0	0	1	5	6	13	10	6	0	0	2	12	12	13	11	0	0	3
10	15	16	14	0	0	0	0	0	1	16	10	0	0	0	0	0	10	15	4	0	0	4
0	5	16	3	0	0	0	1	5	15	13	0	0	0	0	4	15	16	2	0	0	0	5
15	16	9	9	14	0	0	0	3	14	9	2	16	2	0	0	0	7	15	16	11	0	0
8	12	14	8	3	0	0	0	0	10	13	0	0	0	0	0	11	9	0	0	0	0	7

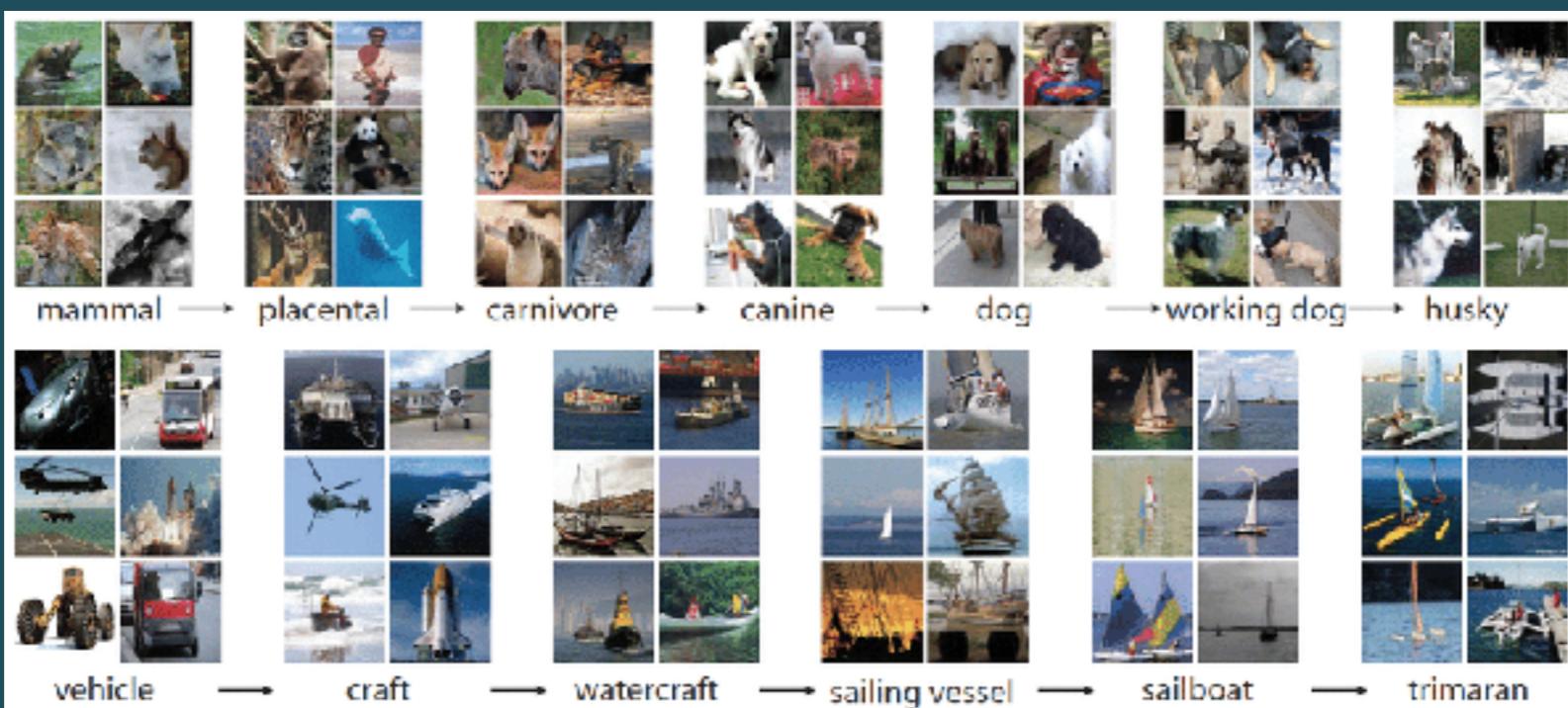
# MNIST classification

- Deep learning way- Convolutional Neural Networks

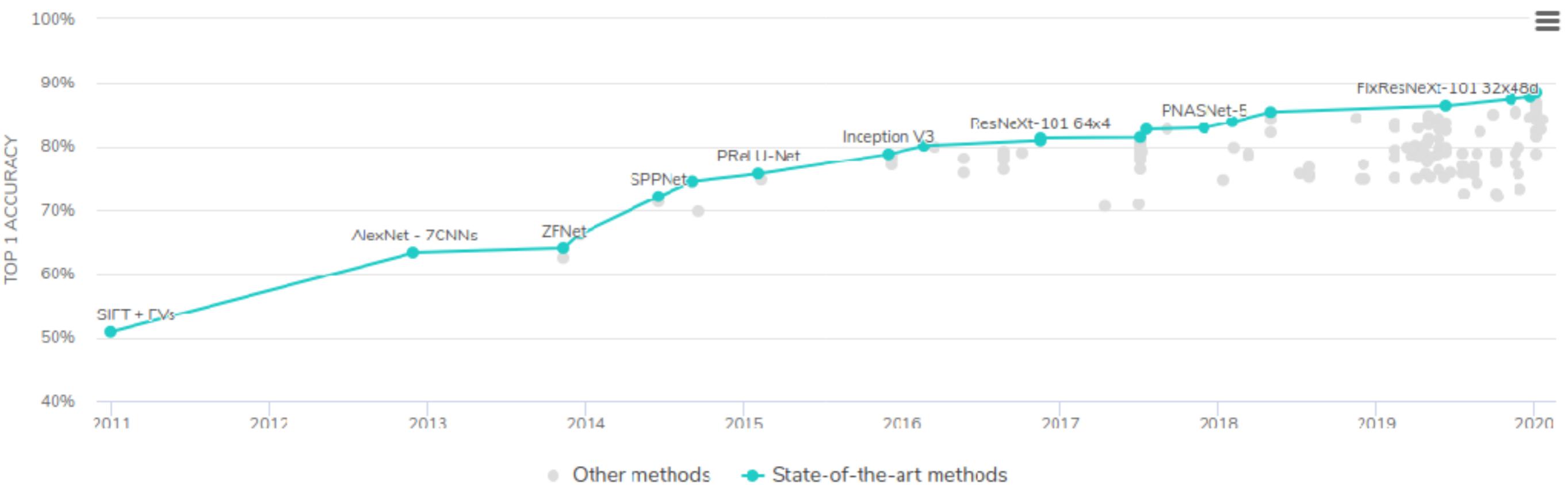


# ImageNet

- 14,197,122 images, > 20,000 categories
- Images hand annotated through Amazon Mechanical Turk
- ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
  - A subset of it (1000 class dataset) is used to report results



# ImageNet results



# ImageNet architectures

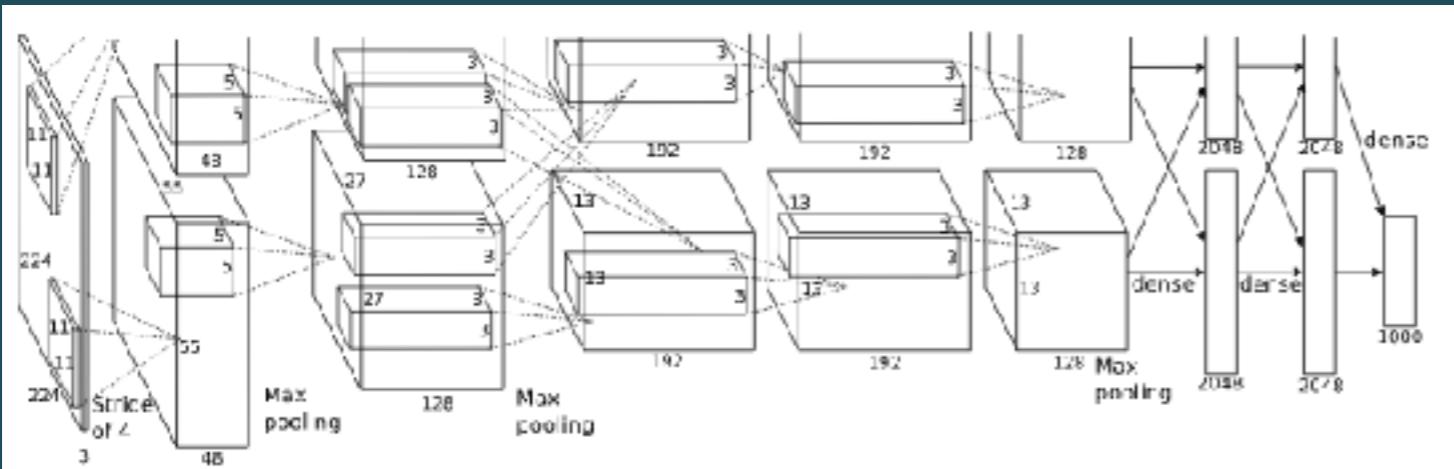


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

Alexnet

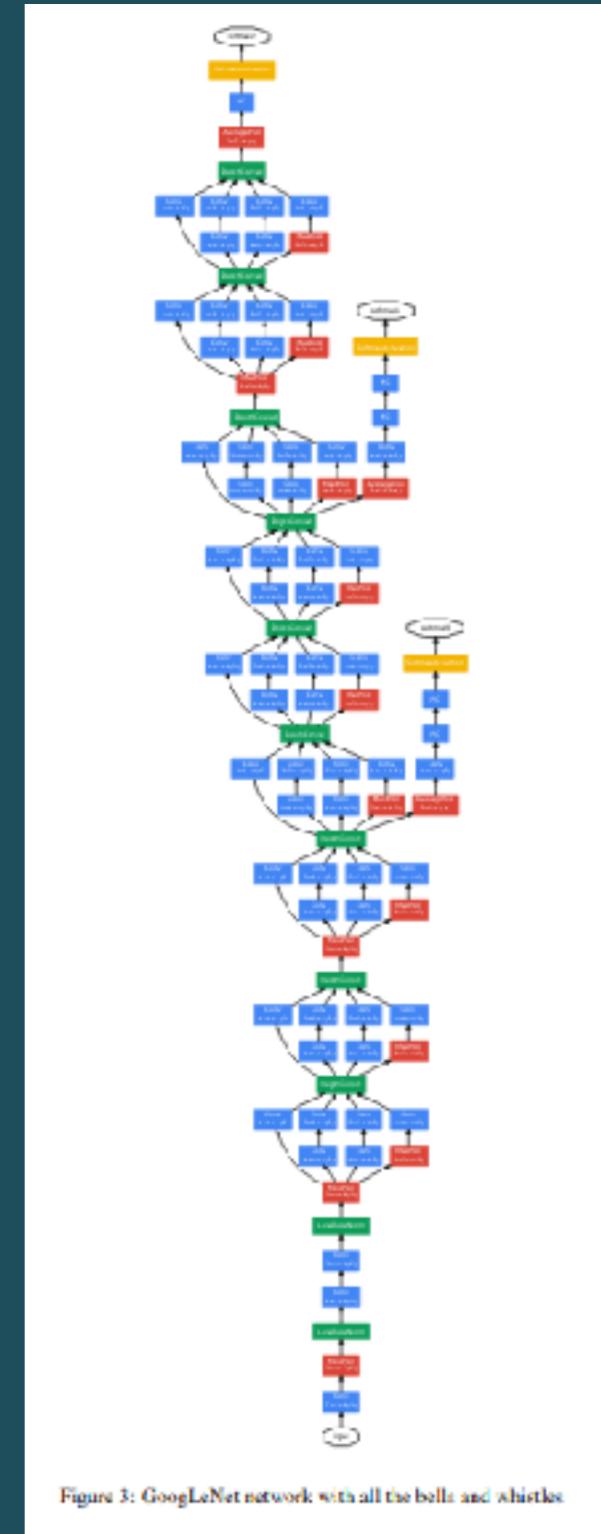
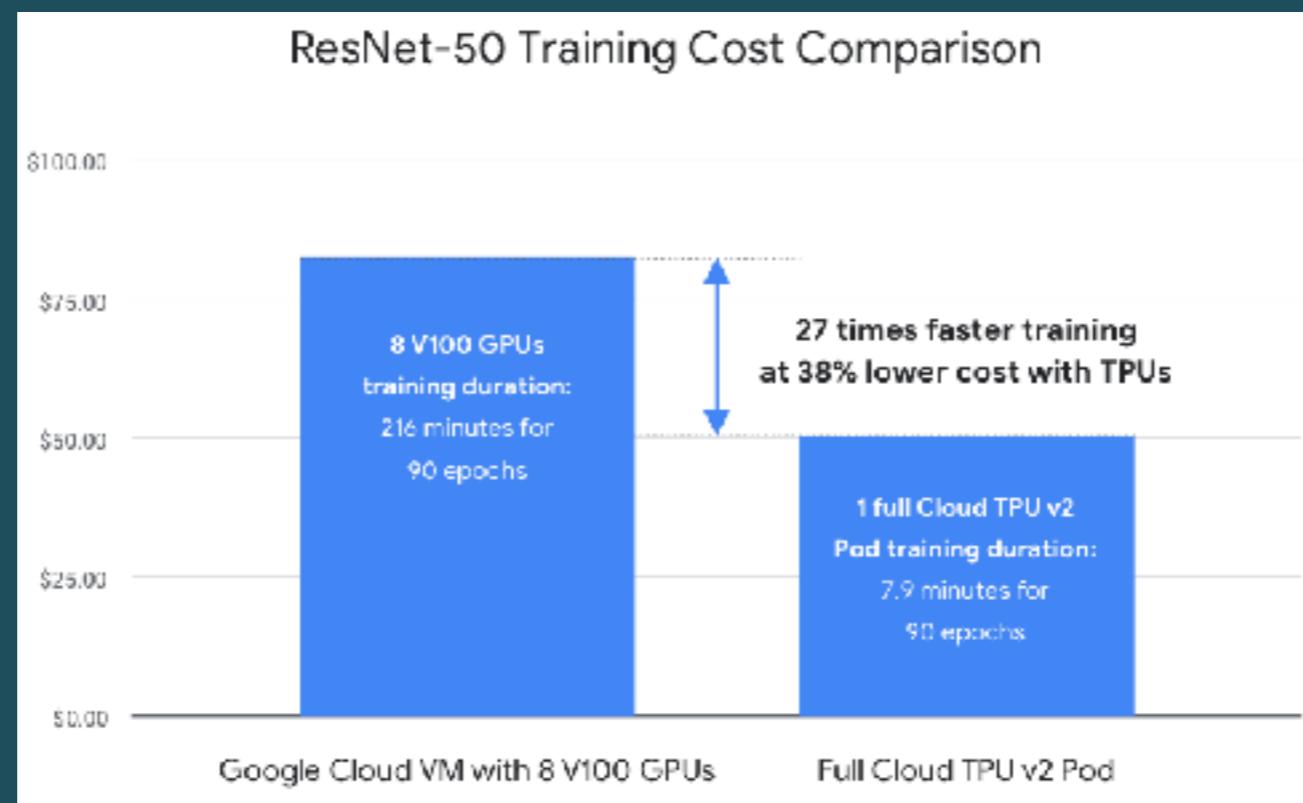


Figure 3: GoogLeNet network with all the bells and whistles

Inception

# Compute power is getting cheaper

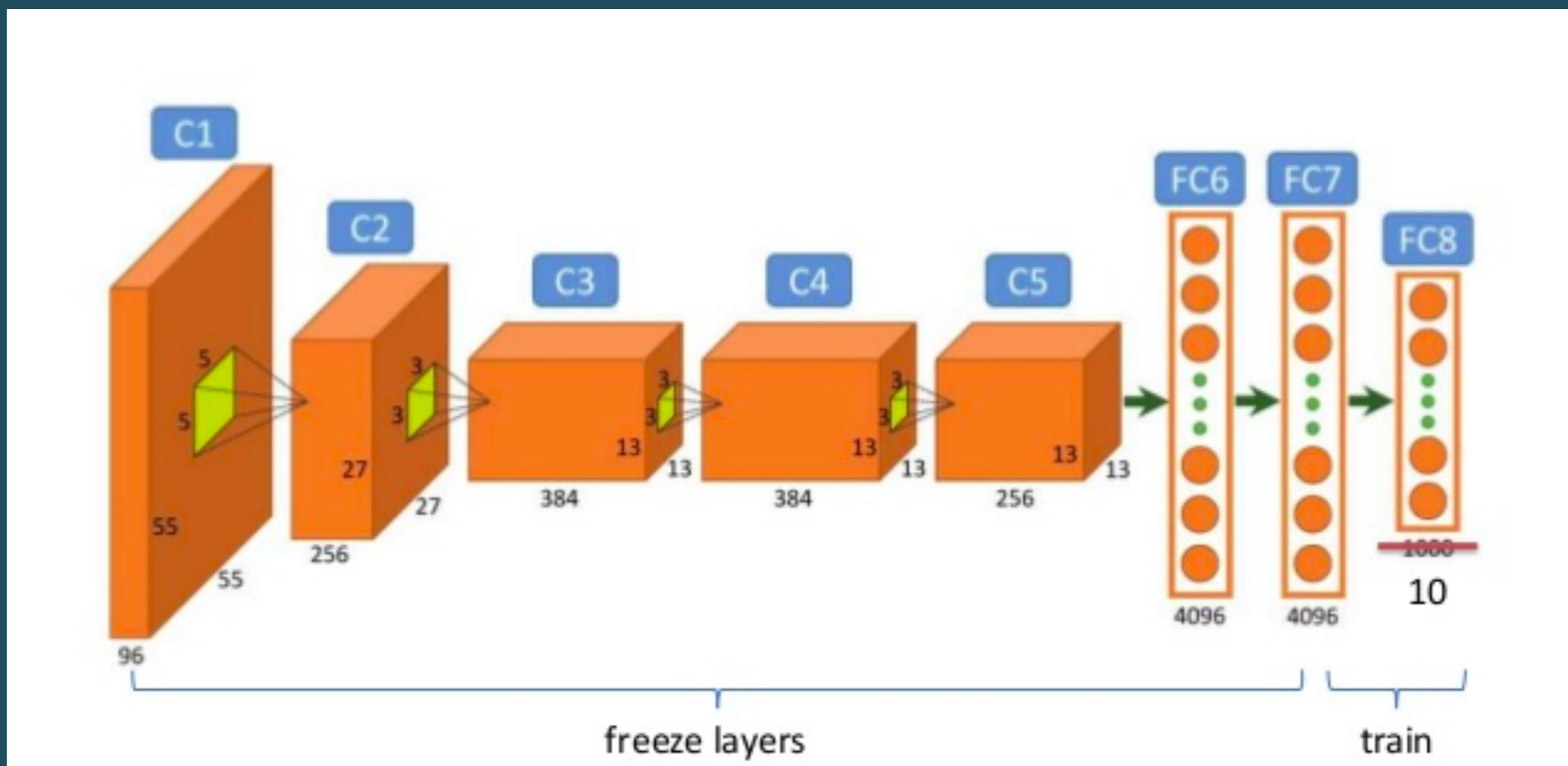
- GPU- faster matrix computation
- Tensor Processing Unit (TPU)- Integrated Circuit developed by Google for faster training of machine learning/ deep learning models



# Building Classifiers for Every Object and Domain

- Creating label sets for training is expensive manual effort
  - 1M types of objects
- Building new models per domain and managing them is infeasible
- Can we reuse and leverage existing effort to create new models?

# Transfer learning



# Transfer learning- approach

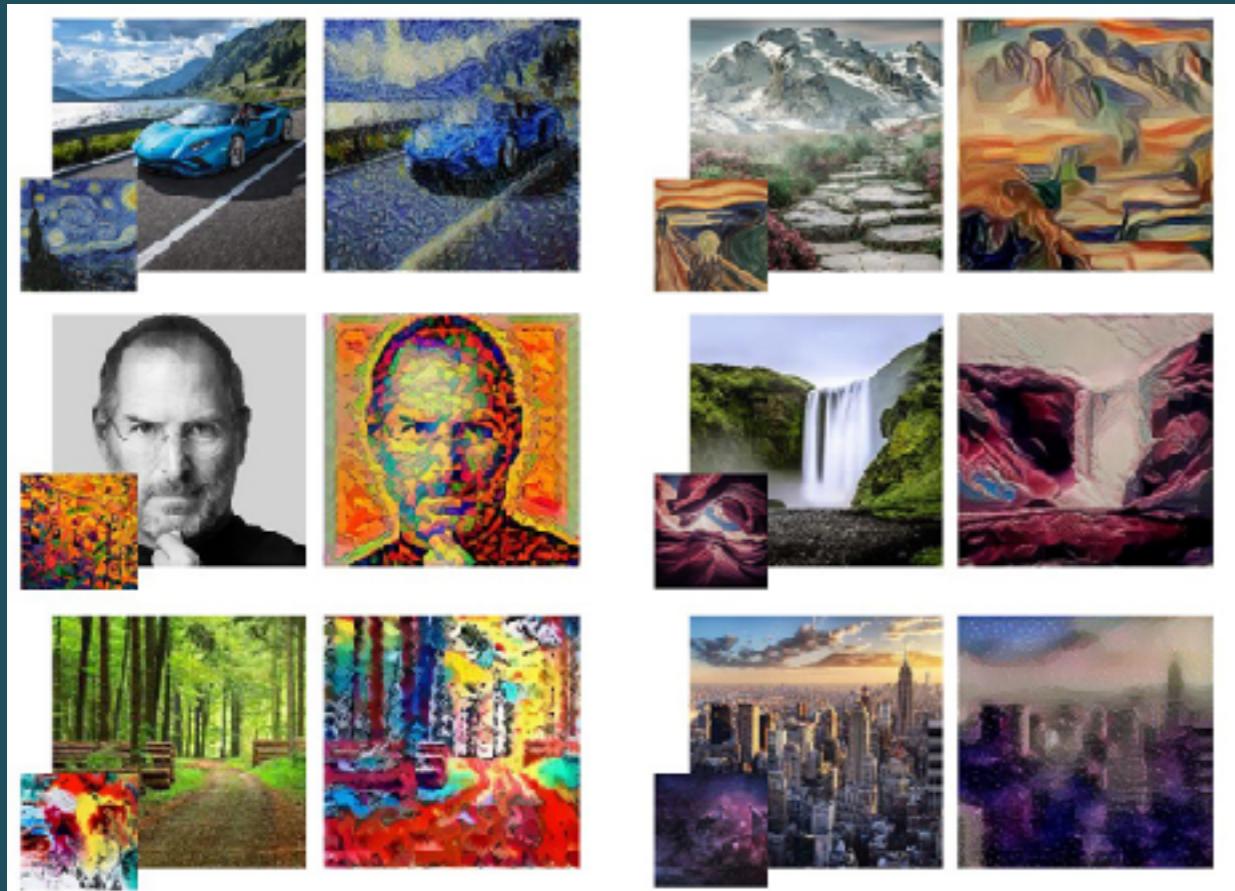
- Bring your dataset of images and classes
- Use any pretrained model- Resnet50, Inception, VGGNet etc.
- Feed the images through the pretrained model and collect the transformations from penultimate layer
- Use them as features and train a classifier
- You could also unfreeze the hidden layers of pretrained model and train the entire model- expensive computation

# New challenges in Image Processing

- Learning with less data
- Generating data (GAN)
- Yolo
- 3D reconstruction from 2D (improving maps)

# Generating data

## Neural style transfer

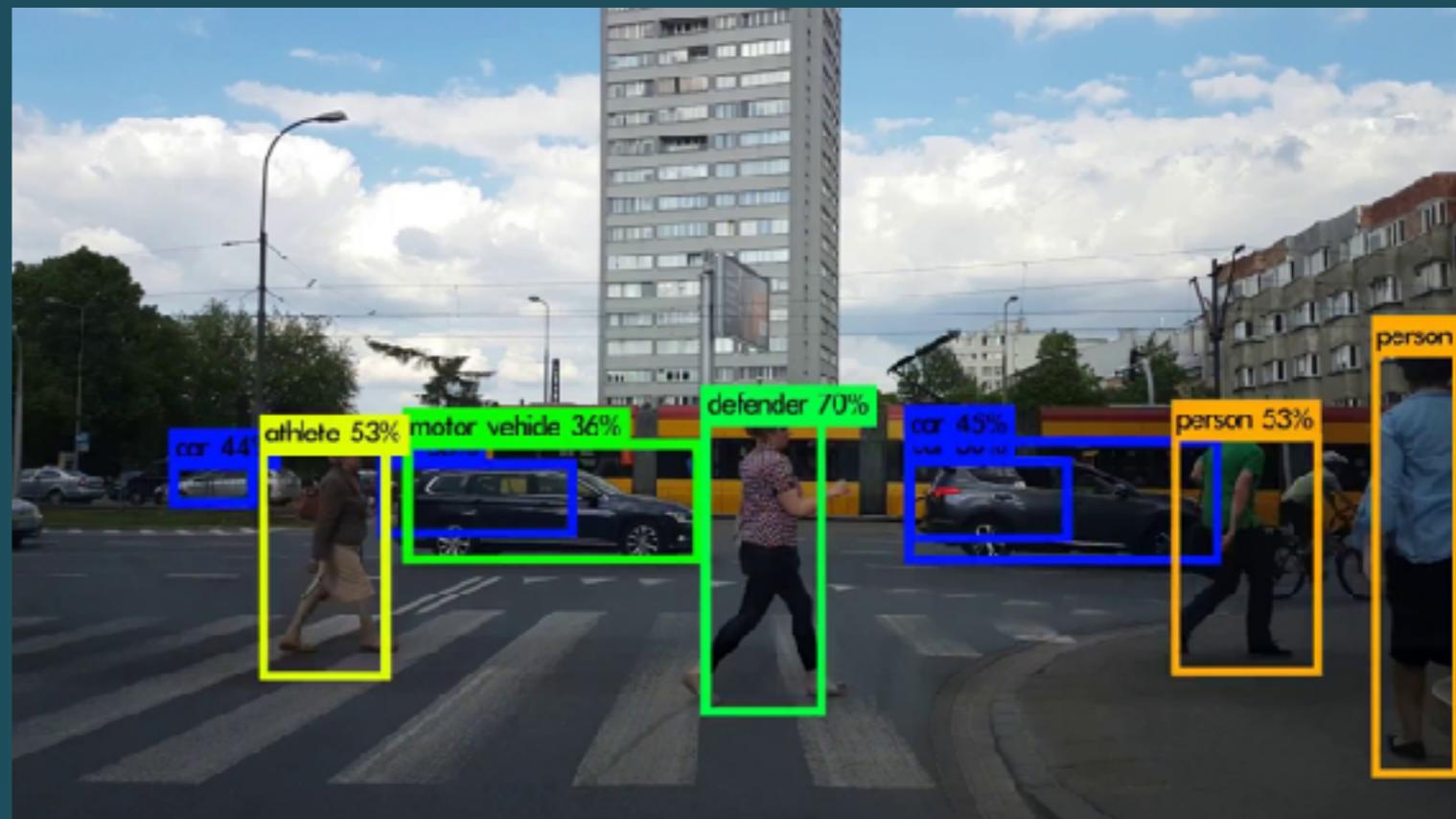
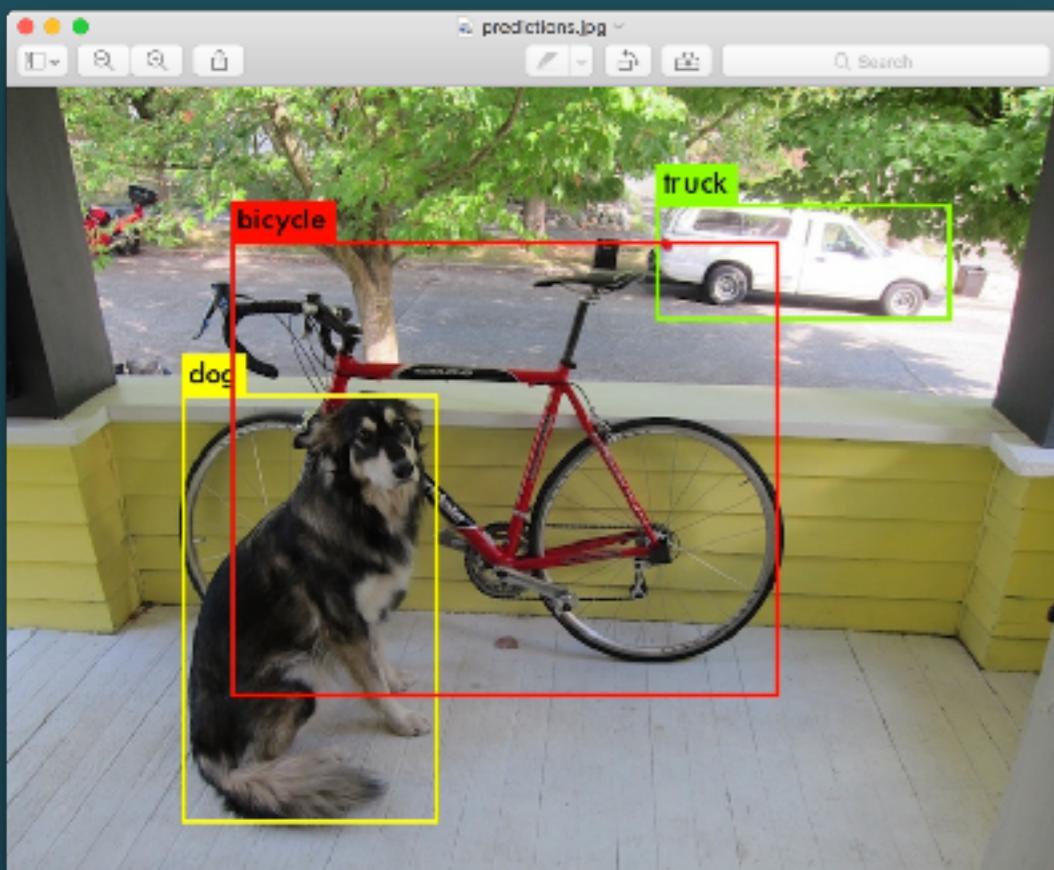


Generated images from images of celebrities using Generative Adversarial Networks (GAN)



# Object detection

- You Only Look Once (YOLO)



# Detecting objects from multiple angles

- Capsule networks



# Use case: Google Maps

- Extract information from Street View images to update locations and landmarks



Avenue des Sapins

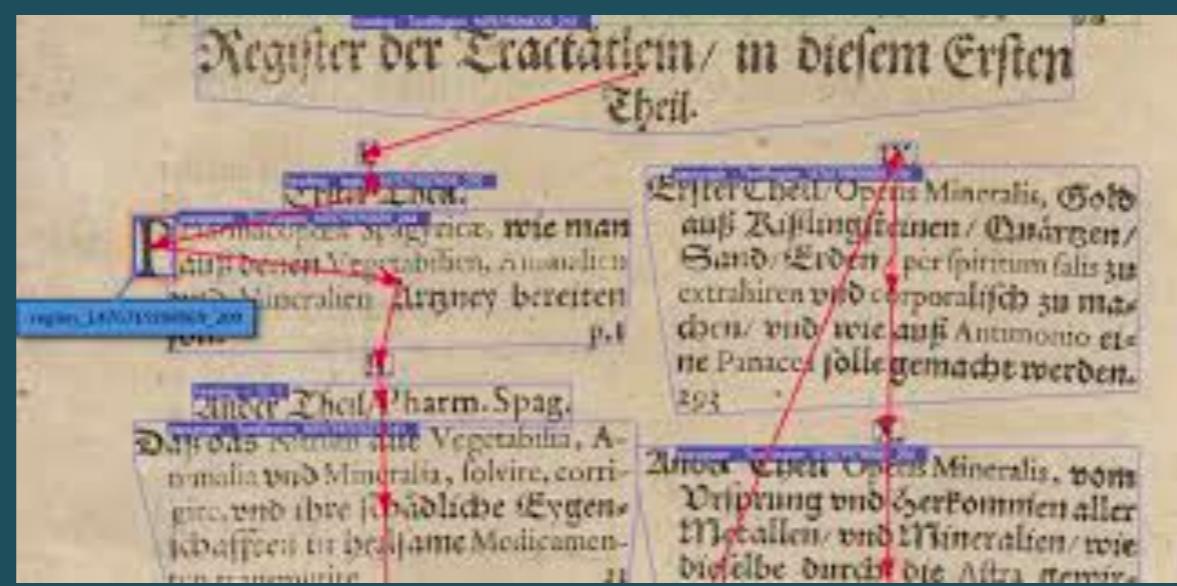
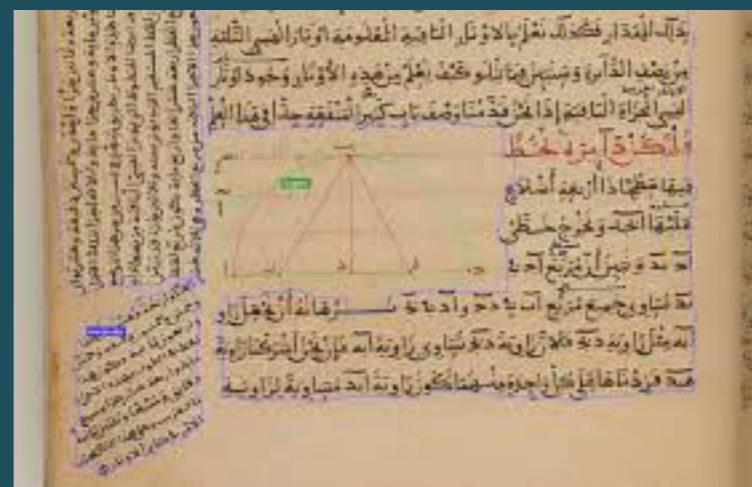
# Enhancing Google Maps 3D vs. 2D

- Seamless street view panoramas, image stitching



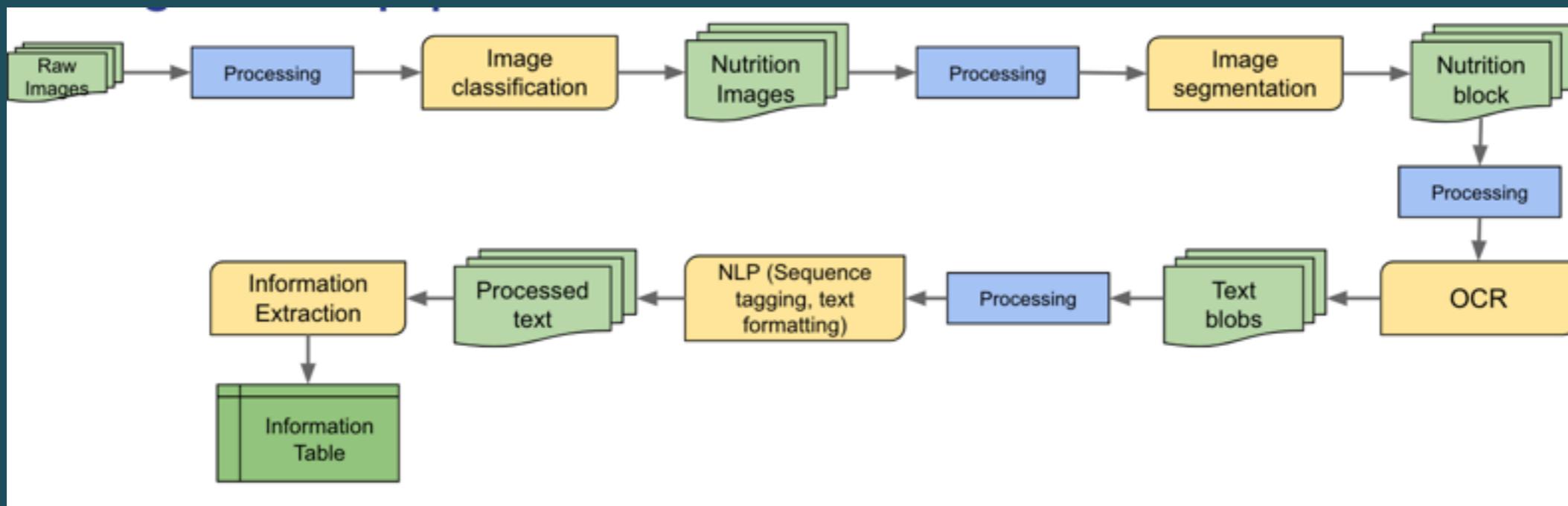
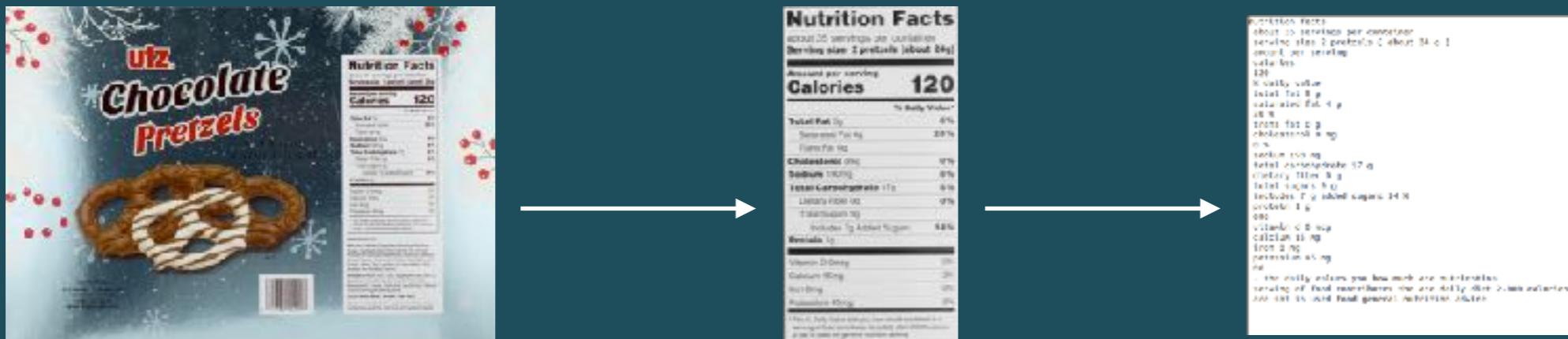
# Optical Character Recognition (OCR)

- Extract handwritten, printed or typed text present in images of documents, scans or scenes



# Retail Case Study

- Extract information from packaged good labels, and make the products searchable



# 3. VIDEO UNDERSTANDING

(CREDIT: STANFORD LECTURES)

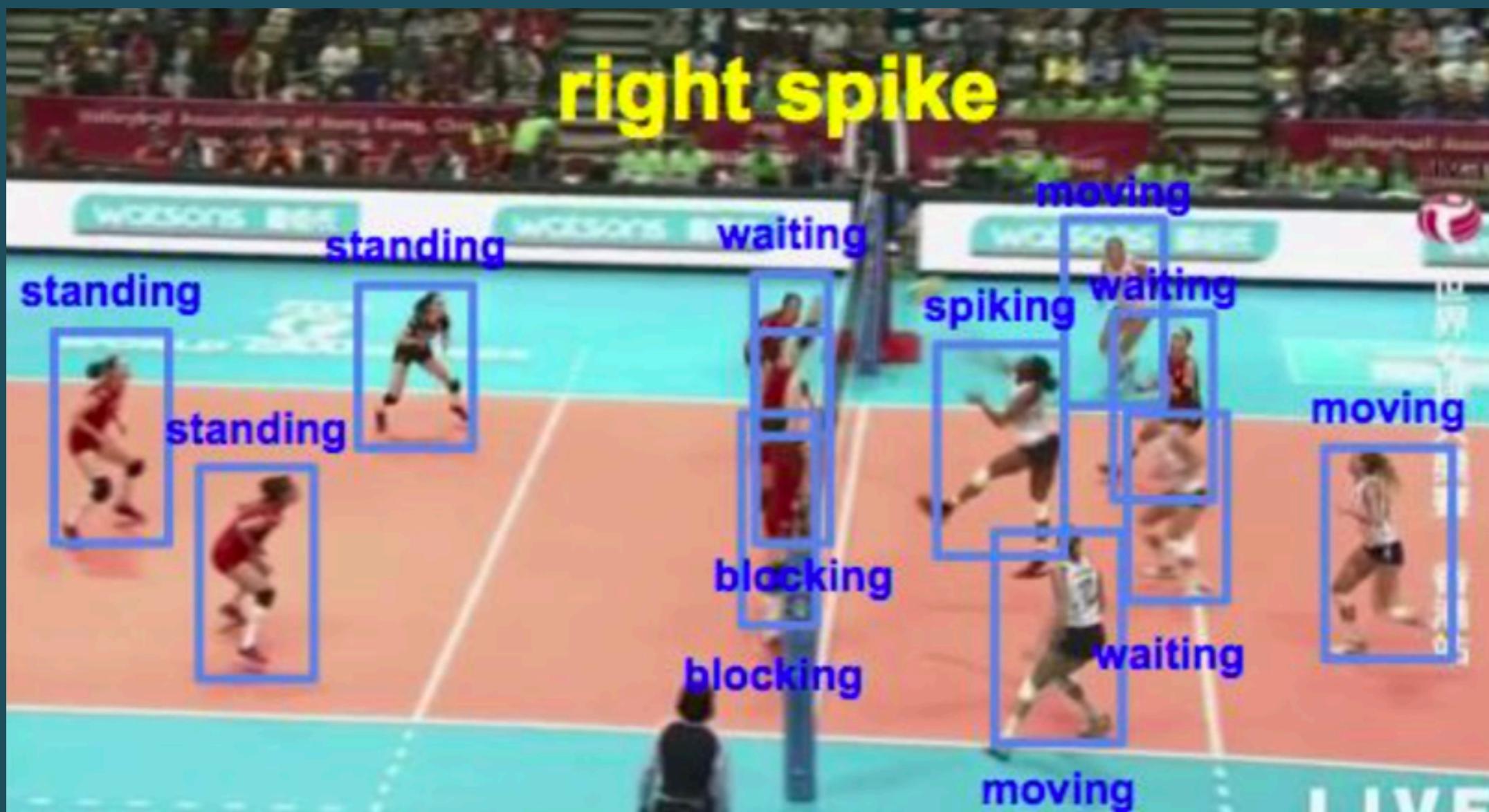
# ROBOTS



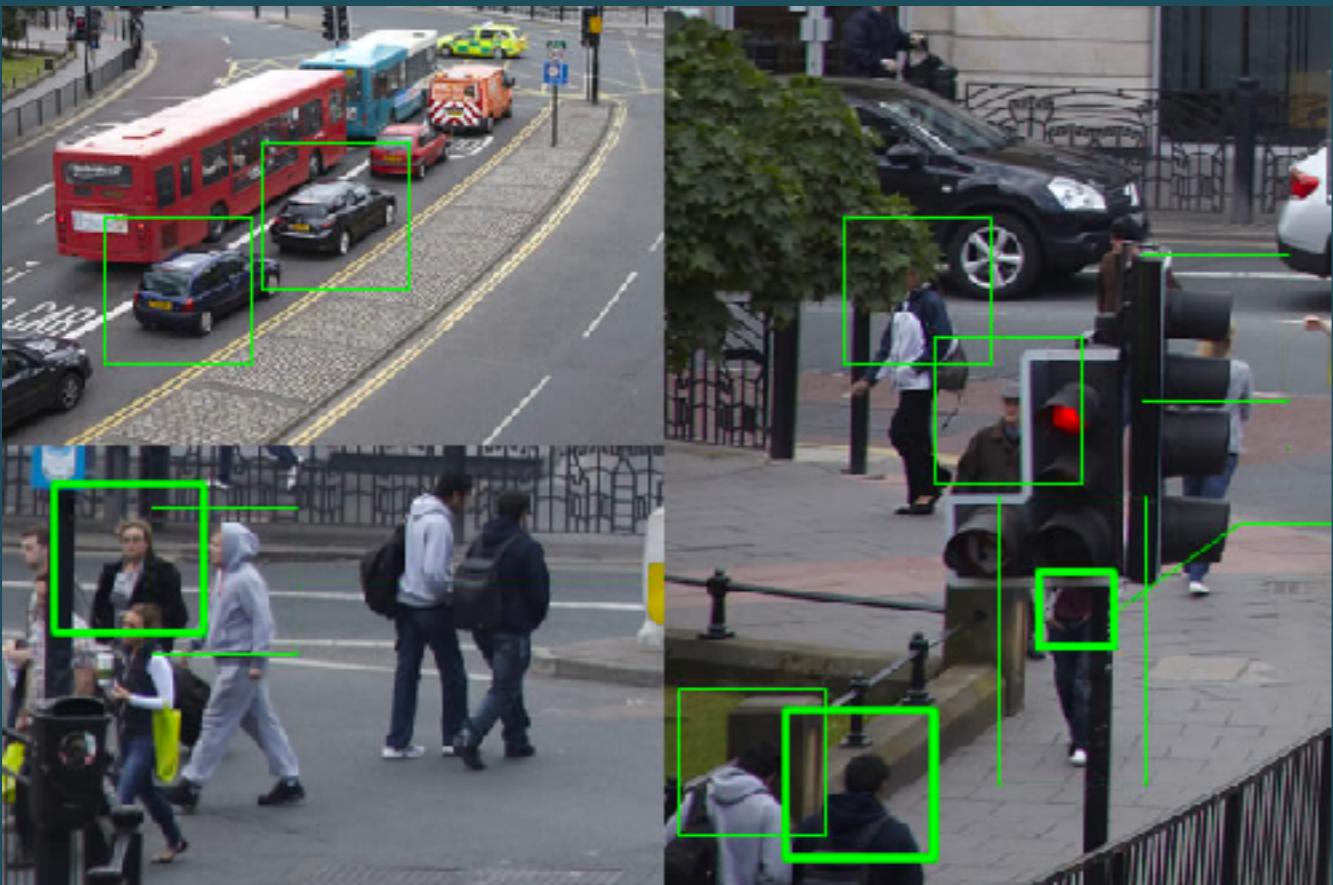
# SELF DRIVING CARS



# GAME ANALYSIS



# SECURITY



# VIDEO UNDERSTANDING TASKS

- Video tagging (analysis)
- Video editing (improvements)
- Video generation (gaming/VR)
- Video QA (teaching and human task automation)
- Video Retrieval

# DATASETS

- UCF101 ([https://www.crcv.ucf.edu/  
data/UCF101.php](https://www.crcv.ucf.edu/data/UCF101.php))
  - 13320 videos from youtube
  - 101 actions
- Sports 1M
  - 1,133,157 videos
  - 487 sports labels



# DATASETS

- Youtube 8M

- <https://research.google.com/youtube8m/>



237K  
Human-verified  
Segment Labels

1000  
Classes

5.0  
Avg. Segments /  
Video

## Entities

Games (788288)

Video game (539945)

Vehicle (415890)

Concert (378135)

Musician (286532)

Cartoon (236948)

Performance art (203343)

Car (200813)

Dance (181579)

Guitar (156226)

String instrument (144667)

Food (135357)

Football (130835)

Musical ensemble (125668)

Music video (116098)

Animal (107788)

Animation (98140)

Motorsport (93443)

Pet (90779)

Racing (84258)

Recipe (75819)

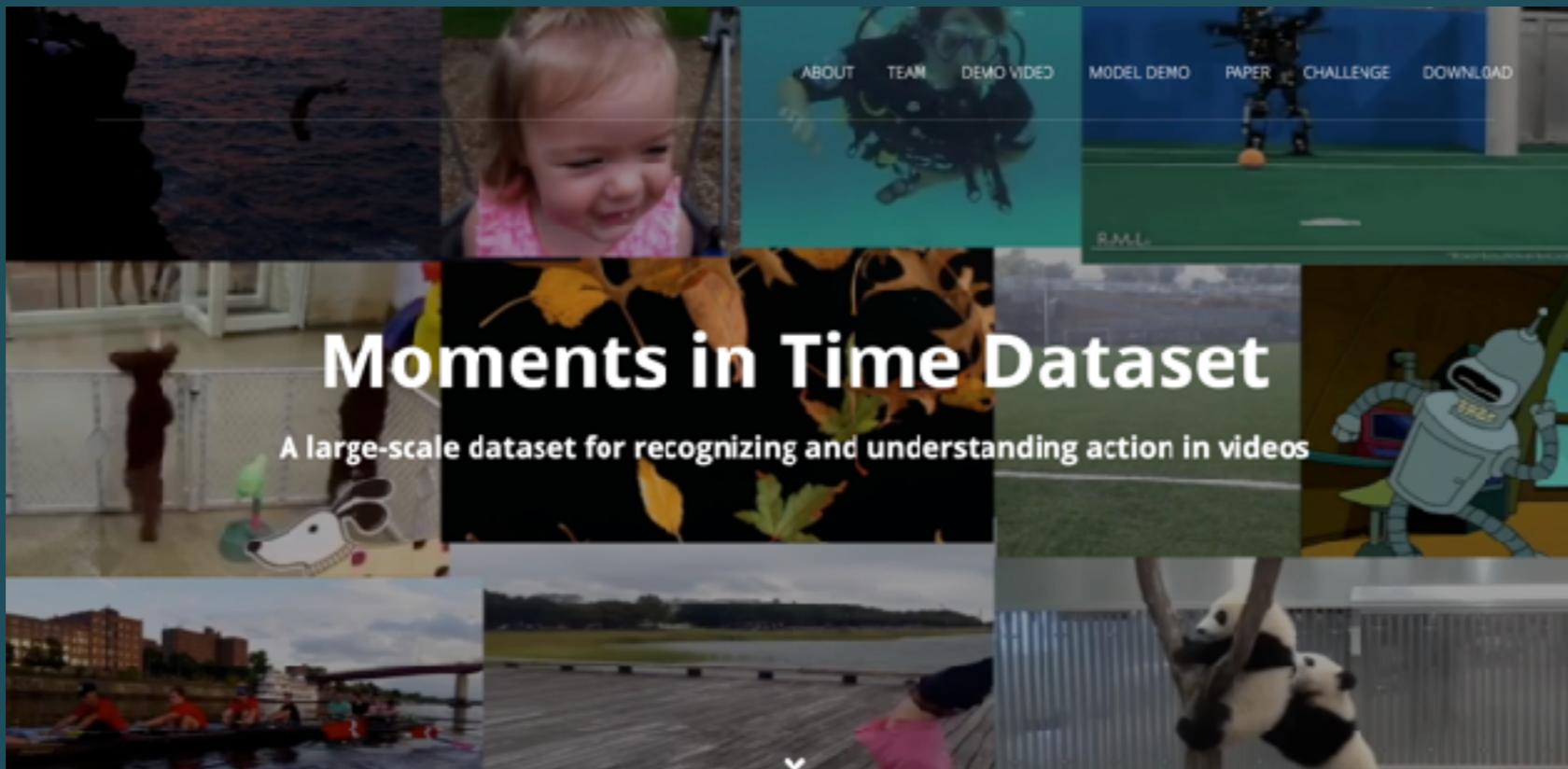
# APPROACH 1: VIDEO FRAME ANALYSIS

- Atomic visual actions
  - Pose and objects
  - Does not describe interaction with objects



# APPROACH 2: VIDEO SEGMENT ANALYSIS

- Video classification (not detection)
  - 339 verbs were tagged (including humans and other actions)
- Moments in Time dataset (<http://moments.csail.mit.edu/>)

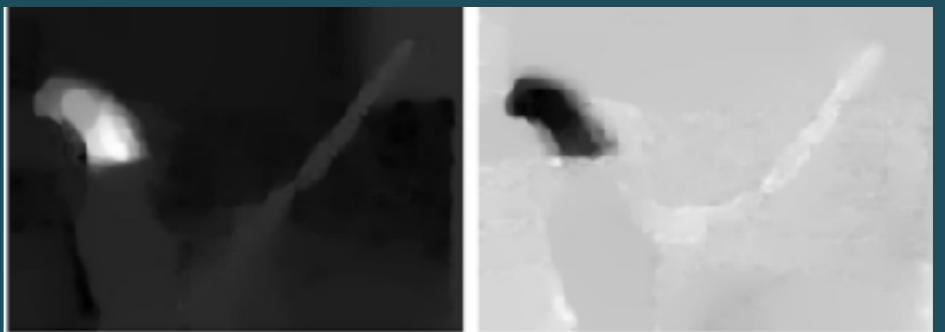
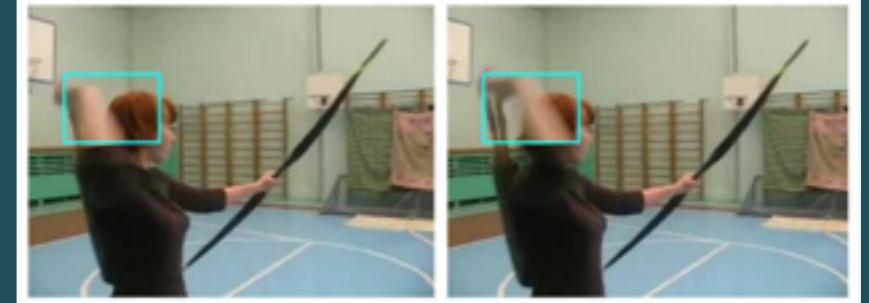


# VIDEO ANALYTICS FRAMEWORK

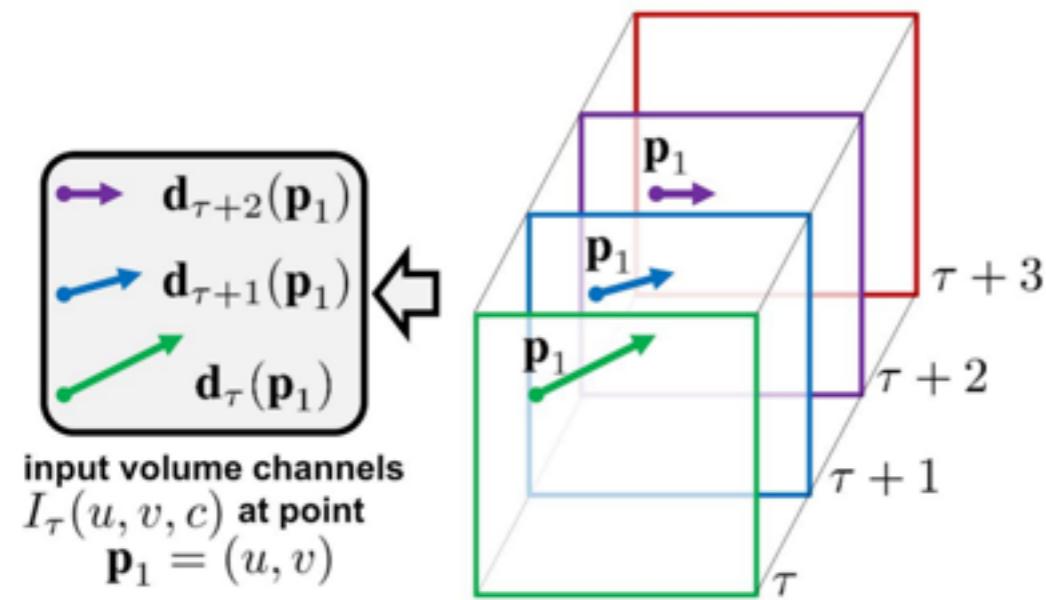
- Labeled Datasets
- Sequence modeling modules
- Temporal reasoning modules
- Action recognition algorithm
  - Representative action
  - Miscellaneous actions

# TRADITIONAL MACHINE LEARNING APPROACH

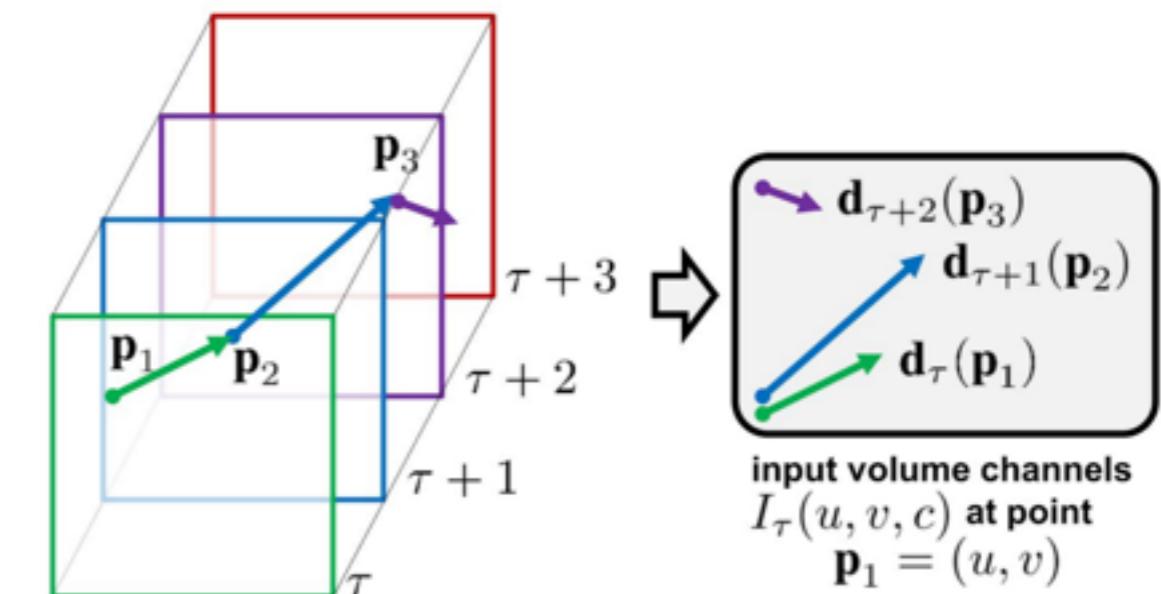
- What is the flow of action?
  - Optical flow algorithm
  - Trajectory flow algorithms



## 1) Optical flow

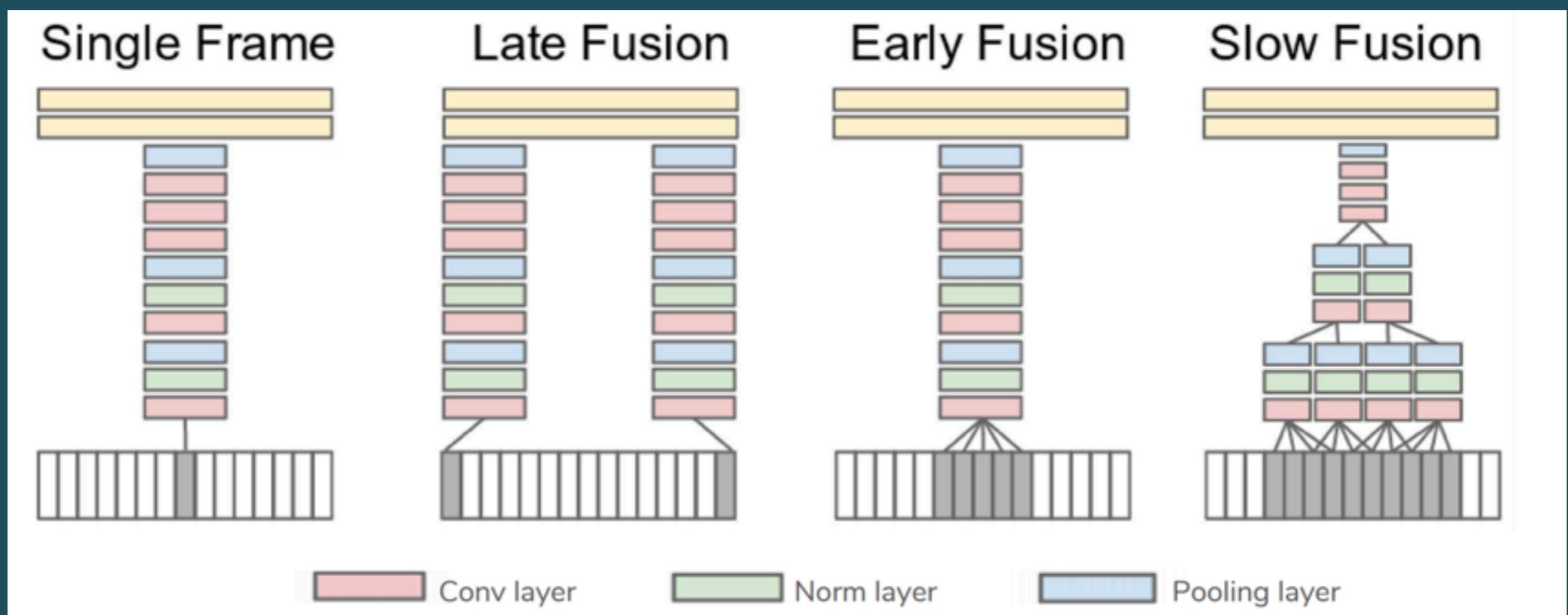


## 2) Trajectory stacking



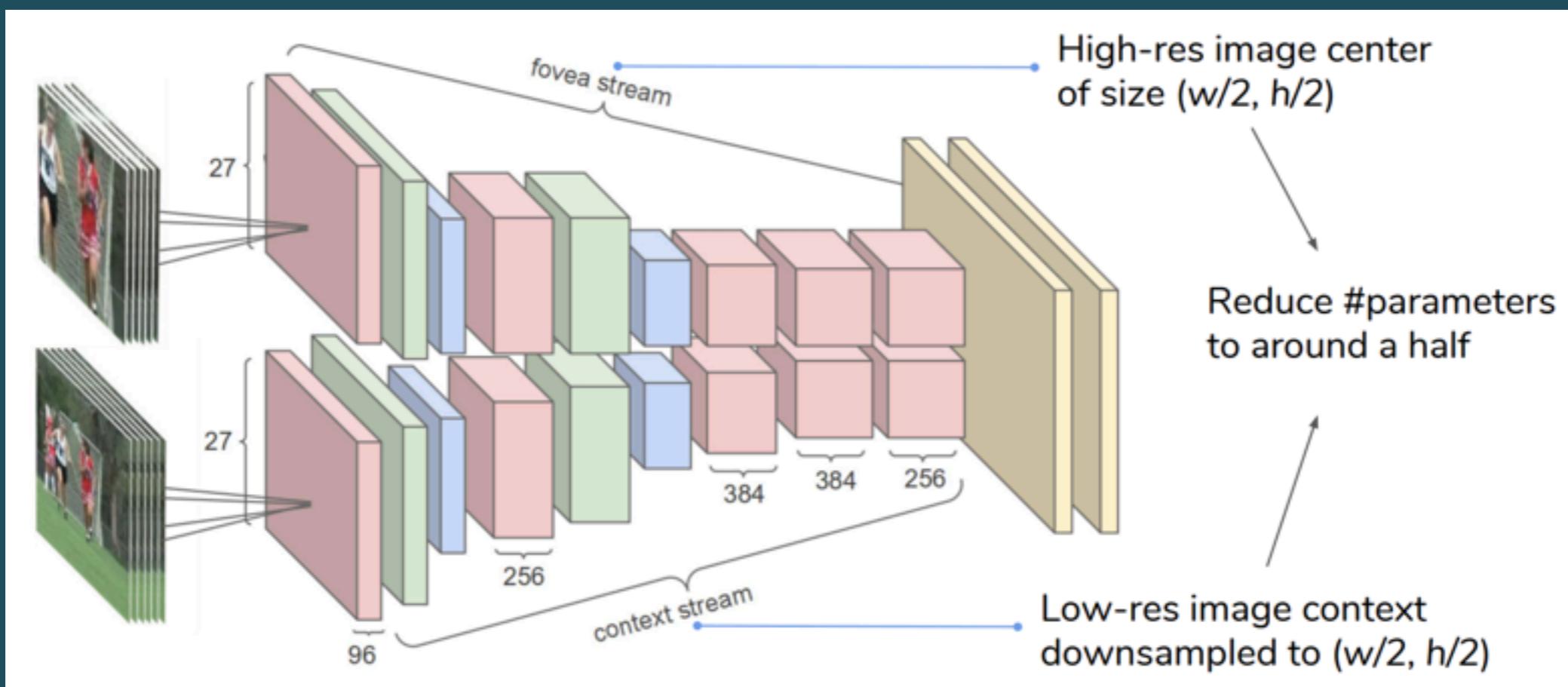
# DEEP LEARNING APPROACH

- Convolution Neural Networks
  - Architecture to capture temporal aspects (combining frames)



# DEEP LEARNING APPROACH

- Convolution Neural Networks
  - Reduce spatial dimension to reduce computational cost



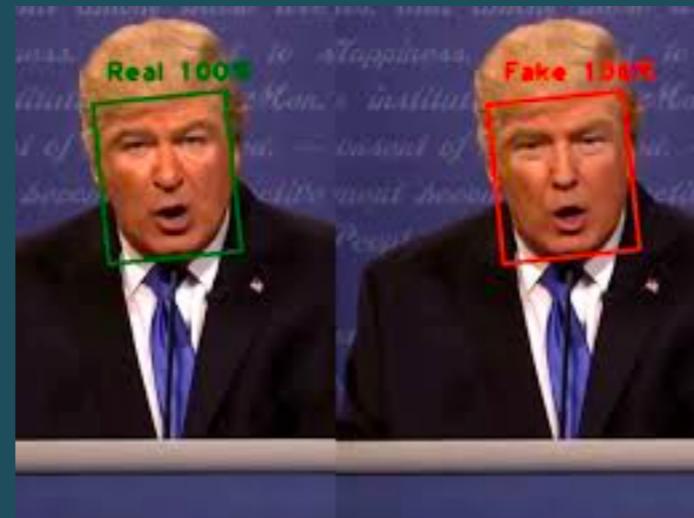
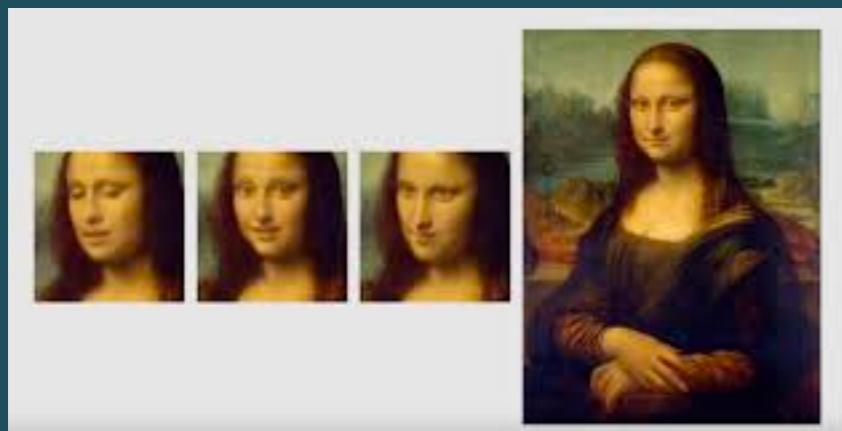
# EVALUATING VIDEO RECOGNITION TASK

- Video Classification
  - Accuracy metrics
- Video querying
  - Provide a verb/task query
  - Retrieve top videos
  - Measure HIT ratios

# CHALLENGES

- New architectures to handle long video segments (E.g summarizing highlights of a cricket game)
- Techniques and Hardware improvements to reduce cost of computation
- Labeling datasets through crowdsourcing and semi-supervised learning
- 3D video analysis with spatial interactions between all objects in the frame
- Separating the motion part of video from the static part of video to improve analysis (two-stream)

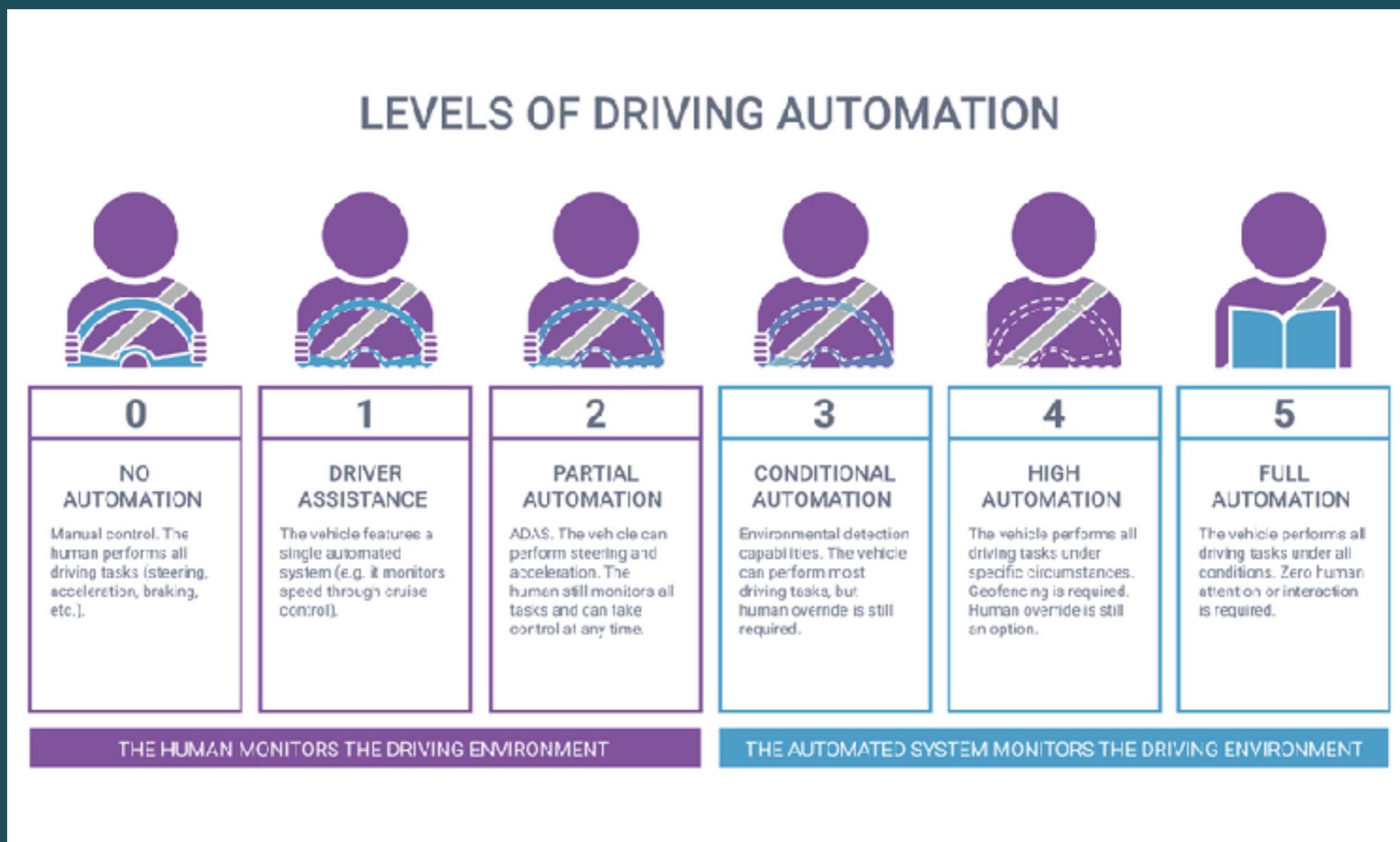
# SCARY REALITY AHEAD



# Self driving cars

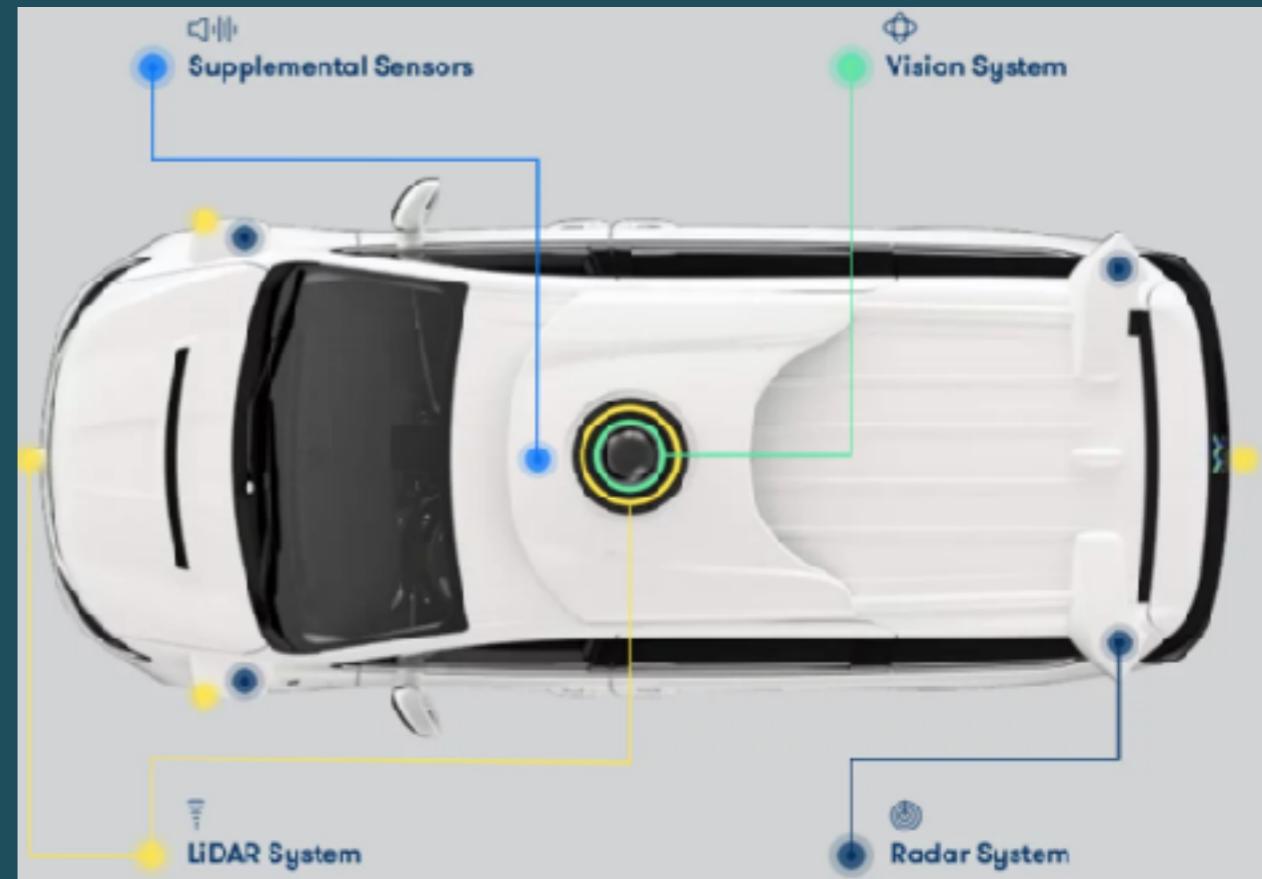
- Defense Advanced Research Projects Agency (DARPA)-2004 challenge to build all terrain self driving cars
- Industry
  - Waymo
  - Tesla autopilot
  - Uber
- Modules
  - Perception
  - Localization
  - Planning
  - Control

# Use case- Self driving cars



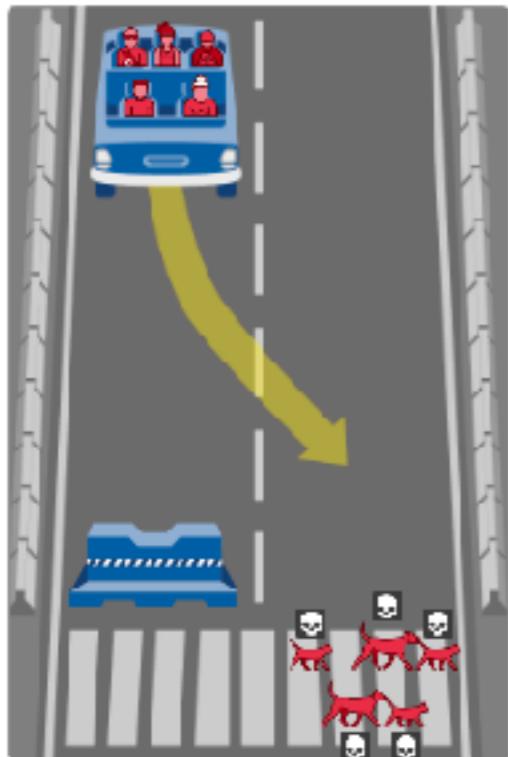
# Waymo

- Trained on 5 billion miles of simulated driving and 5 million miles of on-road driving experience
- Components
  - Video sensor
  - Traffic signal detection
  - Lane tracking
  - Radar

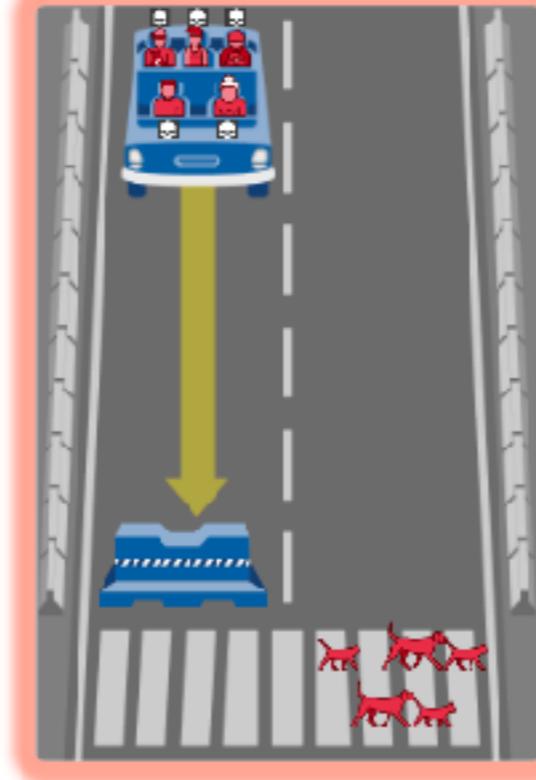


# Ethics

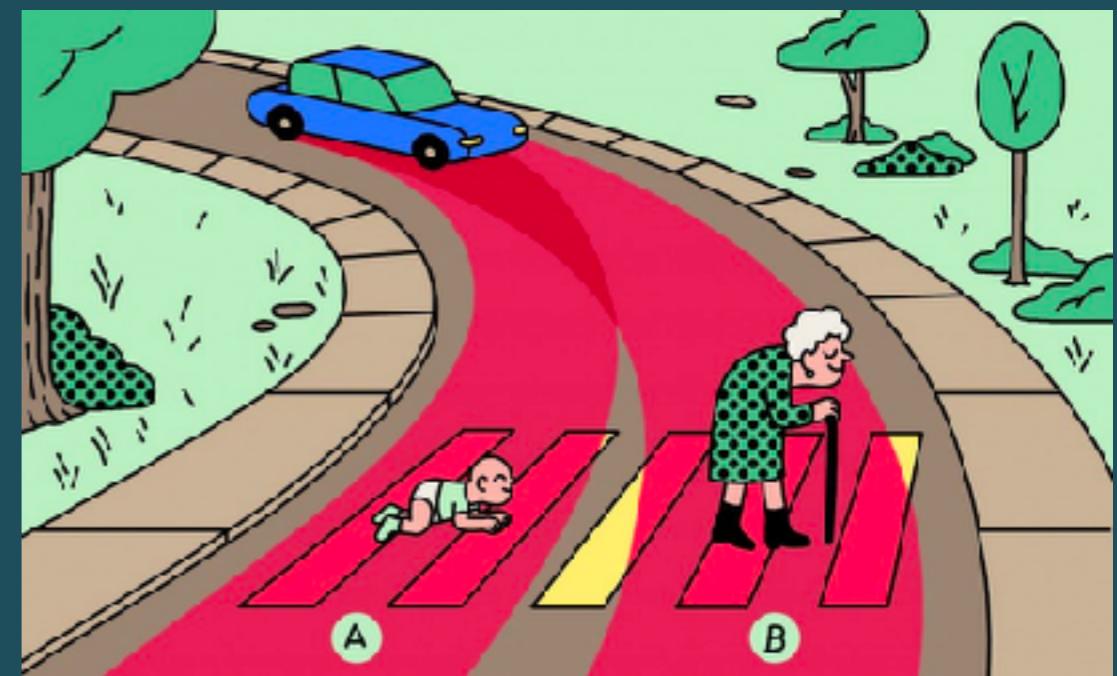
What should the self-driving car do?



Show Description



Show Description



- [https://www.reddit.com/r/MachineLearning/comments/f8wsyg/nd\\_yolo\\_creator\\_joseph\\_redmon\\_stopped\\_cv\\_research/](https://www.reddit.com/r/MachineLearning/comments/f8wsyg/nd_yolo_creator_joseph_redmon_stopped_cv_research/)
- [https://en.wikipedia.org/wiki/List\\_of\\_self-driving\\_car\\_fatalities](https://en.wikipedia.org/wiki/List_of_self-driving_car_fatalities)
- Privacy

# FROM HERE

- Open source libraries
  - OpenCV
  - Tensorflow
  - FastAI
- Image Datasets
  - MNIST Dataset
  - Youtube Video Dataset
  - Amazon dataset
- Video Datasets
  - Youtube Dataset
  - AVA dataset
  - Sports 1M

THANKS