

# Novice CV/VLM Workshop

Image Processing, Computer  
Vision, and Your First Vision  
Language Model

Ming Hao

Data Scientist @ FWD



# Table of Contents

1

Introduction to  
Computer Vision

2

Image Processing with  
OpenCV

Hands-on

3

Using Object  
Detection Models

Hands-on

4

Introduction to  
Vision Language  
Models

Hands-on

5

Deploying  
pretrained VLM to  
FastAPI

Hands-on

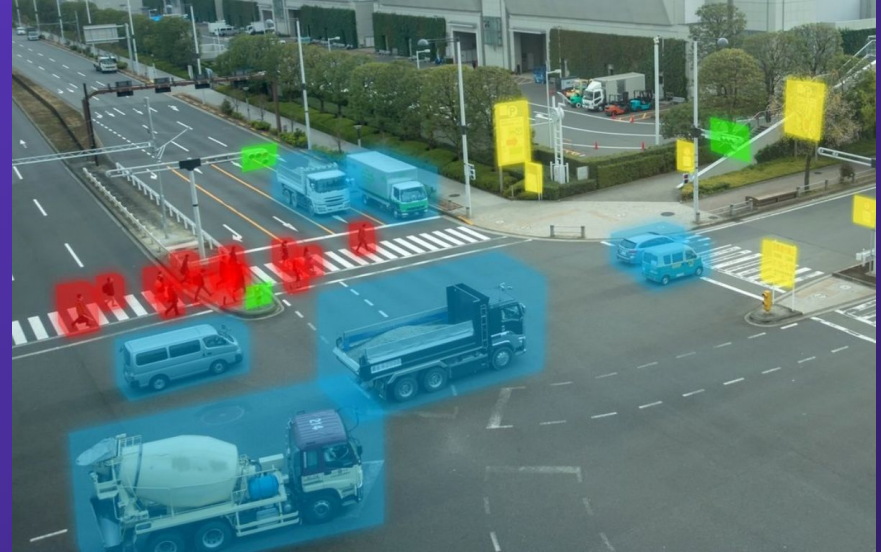
# 1

## Introduction to Computer Vision

Fundamentals of Computer Vision

# What is Computer Vision?

- A branch of AI that teaches computers to see and understand visual data in a way that mimics human vision
- Utilizes algorithms to analyse and interpret visual inputs



Adapted from [viso.ai](https://viso.ai)

# Analog Comparison

## Human Vision vs Computer Vision



Tomato



Eye



Brain



Tomato

Result



Tomato



Sensor



Algorithm

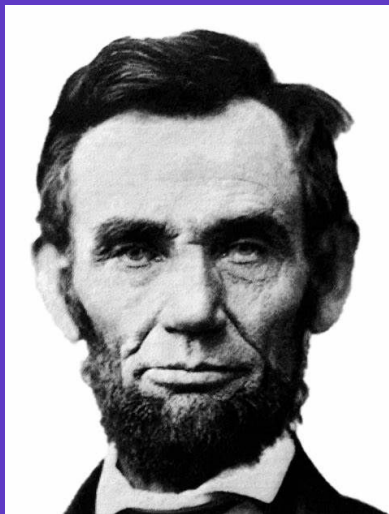


84 % Tomato  
15% Apple  
1% Peach

Result

# What Computers See - Pixel Values

What we see



Grayscale  
(Black and White)

What computers see

167	163	174	168	160	162	129	161	172	163	166	166
166	182	163	74	76	62	33	17	110	210	180	164
180	180	60	14	34	6	10	33	48	106	169	181
206	109	6	124	131	111	120	204	166	16	66	180
194	68	137	261	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	186	216	211	198	139	76	20	169
189	97	166	84	10	168	134	11	31	62	22	148
199	168	191	193	198	227	178	143	182	106	36	190
206	174	195	252	236	231	149	178	228	43	96	234
190	216	116	149	236	187	86	190	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	90	2	109	249	216
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
196	206	123	207	177	121	123	200	176	13	96	218

167	163	174	168	160	162	129	161	172	163	166	166
166	182	163	74	76	62	33	17	110	210	180	164
180	180	60	14	34	6	10	33	48	106	169	181
206	109	6	124	131	111	120	204	166	16	66	180
194	68	137	261	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	186	216	211	198	139	76	20	169
189	97	166	84	10	168	134	11	31	62	22	148
199	168	191	193	198	227	178	143	182	106	36	190
206	174	195	252	236	231	149	178	228	43	96	234
190	216	116	149	236	187	86	190	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	90	2	109	249	216
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
196	206	123	207	177	121	123	200	176	13	96	218

Array of Pixel values

Colored Images could be coded in RGB (Red, Green, Blue), HSL (Hue, Saturation, Lightness), etc, which would be channels (depth) of pixel value arrays.

# 2

## Image Processing with OpenCV

Hands On



# Notebook



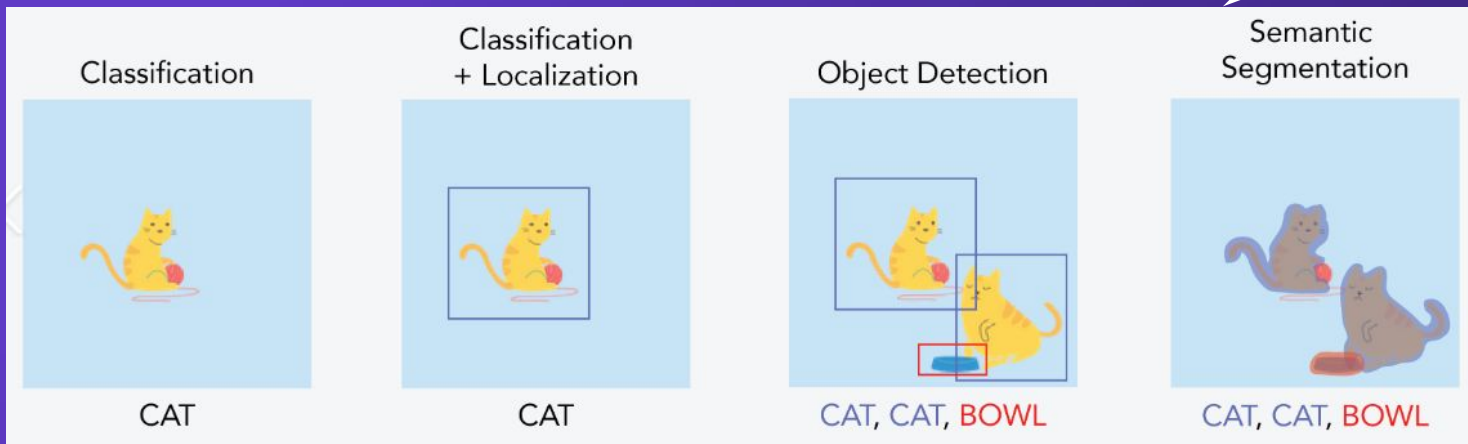
# 3

## Using Object Detection Models

Hands On

# Common Problem Types in Computer Vision

Complexity



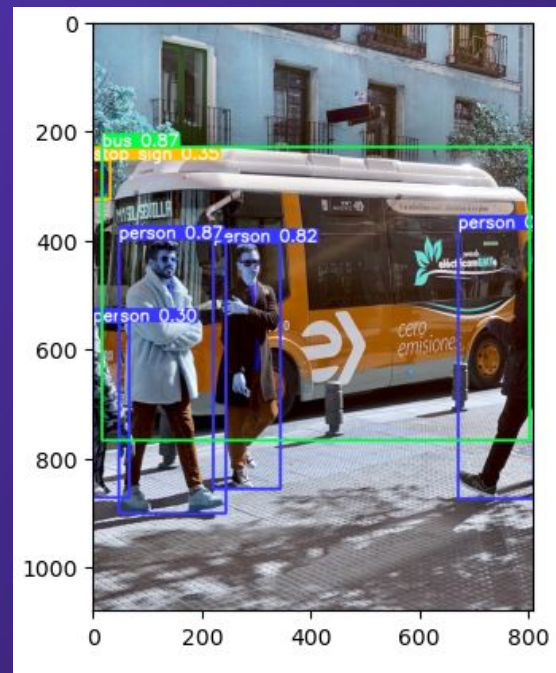
Features	Classification	Bounding Box	Multiple objects	Segment Mask
Applications	Image Tagging, Scene Understanding, Content moderation		Robot navigation, Facial Recognition Self Driving Cars	Medical Image Diagnosis Aerial Image processing

# What is Object Detection?

Localizing and identifying objects within an image.

## Tasks:

Where	<b>Localization:</b> Pinpointing the bounding box (rectangular frame) around the detected object.
What	<b>Classification:</b> Assigning a class label to the object (e.g. bus, person, dog).
Score	<b>Confidence Score:</b> Indicating the model's certainty in the detection and classification.

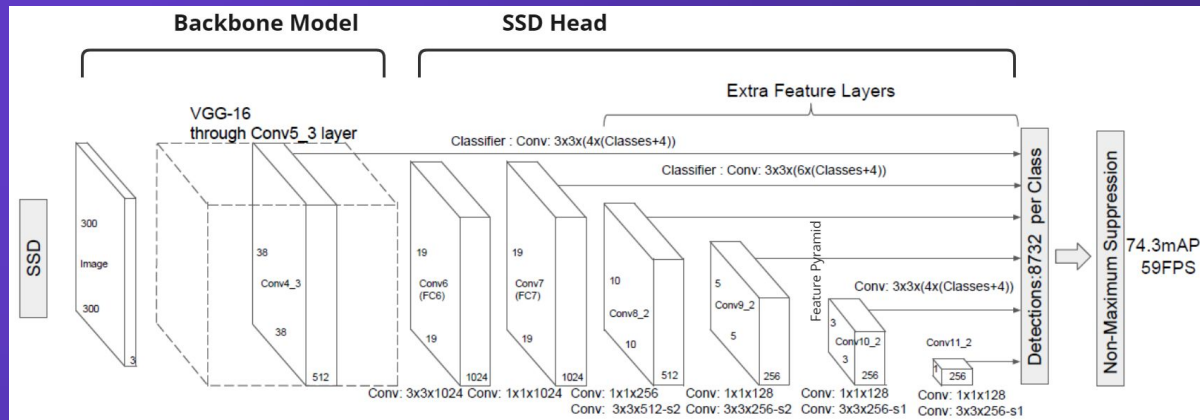


# Single-Shot Detector (SSD)

SSD is just one technique;  
explore the CV/VLM educational  
content for more!

## Two Components

1. **Backbone model:** Pre-trained image CNN (Resnet, VGG)
2. **SSD Head:** More convolutional layers added to the backbone, whereby the outputs are bounding boxes and classes of the objects in the spatial locations.

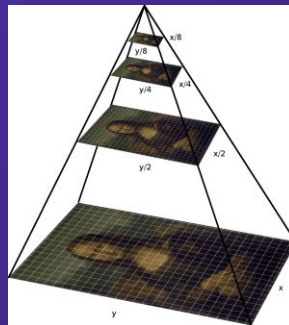


Wei Liu et al. in the paper SSD: Single Shot MultiBox Detector.

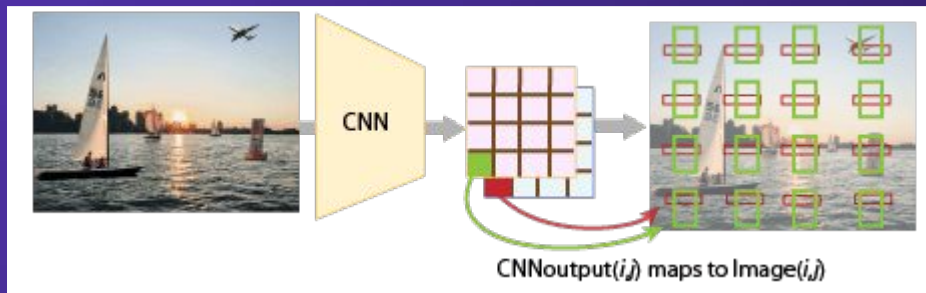
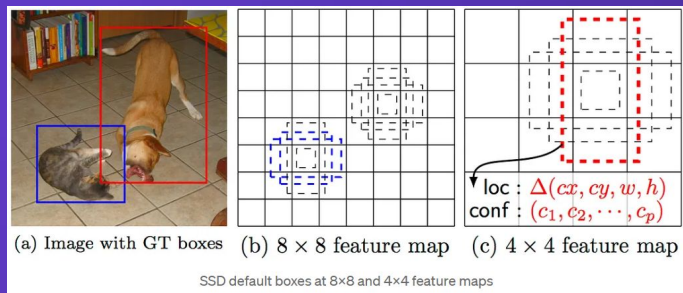
# Single-Shot Detector (SSD)

## Innovations

- Applies various feature map grid cell sizes (e.g.  $8 \times 8$ ,  $6 \times 6$ , ...,  $1 \times 1$ ) to detect objects of different sizes [image pyramid]
- Anchor Boxes



This is all done in the SSD Head network!



# YOLO v8 Architecture

This is another popular Object Detection algorithm.

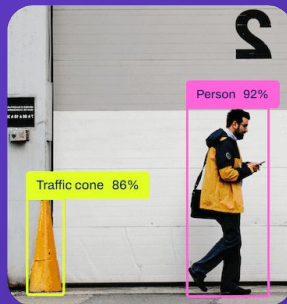
YOLOv8 has also integrated other submodules

- Classify models pretrained on the ImageNet dataset.
- Detect, Segment and Pose models pretrained on the COCO dataset (with track mode).

Classify



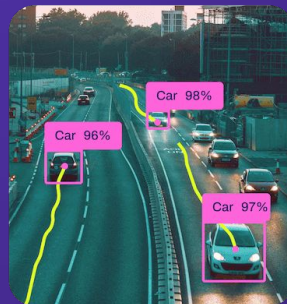
Detect



Segment



Track



Pose



- Latest version: YOLO v8
- Adapted from Brief summary of YOLOv8 model structure · [GitHub](#)

# Notebook



# 4

## Introduction to Multimodal Models

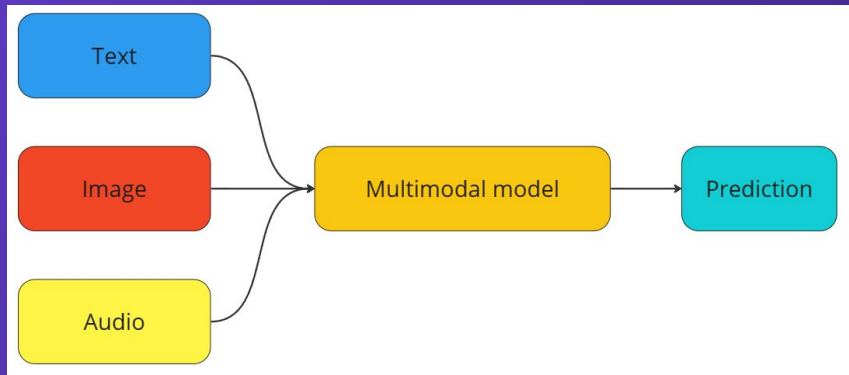
What are multimodal models?

# Multimodal models

## Definition

Multimodal models are models that integrates information from multiple modalities (e.g. text, image, videos, audio, gestures) to create a unified representation.

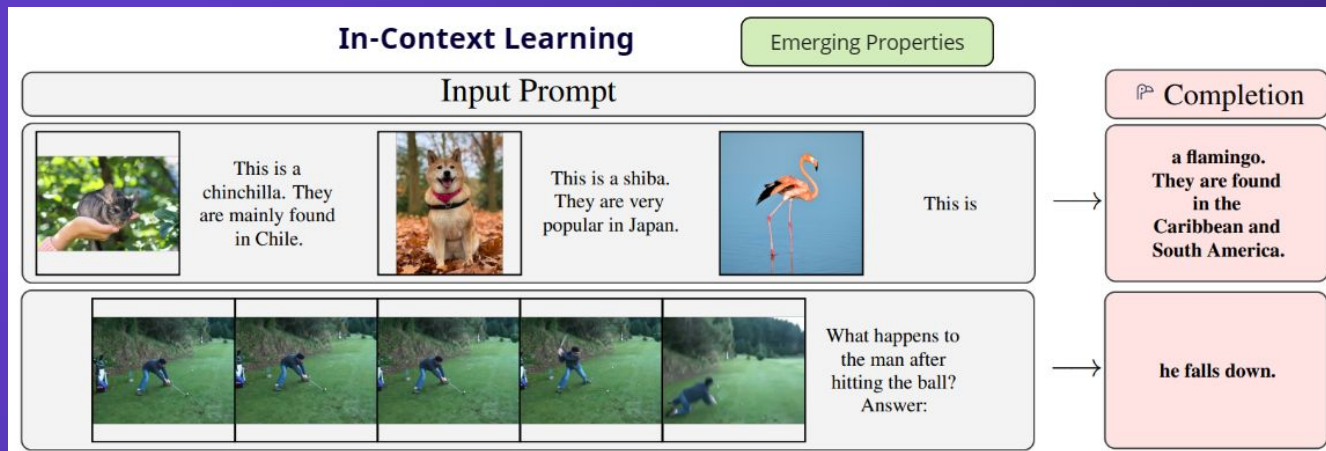
By leveraging different types of data, multimodal models can perform tasks that involve complex interactions between modalities and thus understand and reason about the world in more comprehensive way.



# Multimodal models and their importance

The key importance of multimodal models in computer vision are:

1. Richer Understanding
2. Better Performance
3. Facilitate Emerging Properties



(left) The emerging properties of pre-training on web-scale interleaved image-text data: multimodal in-context-learning.VLM (Flamingo)

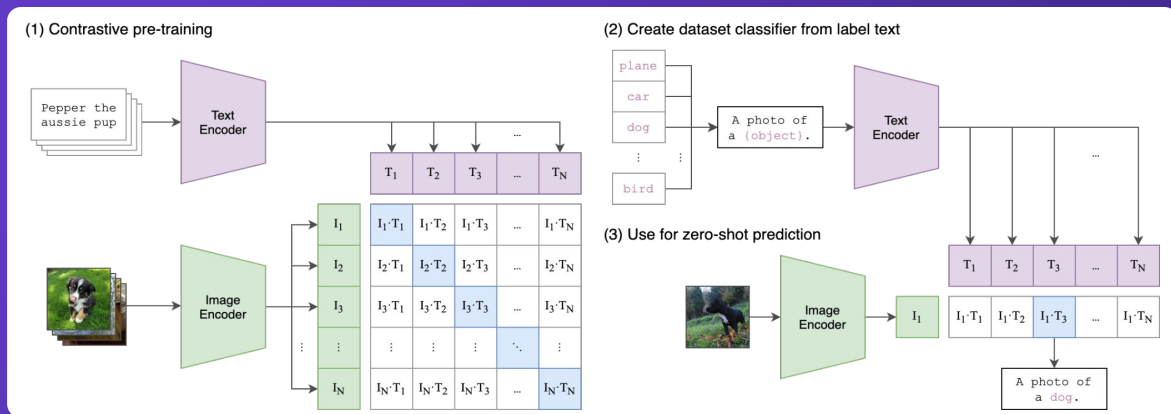
Adapted from [Alayrac et al. \(2022\)](#) via

# CLIP (Contrastive Language-Image Pre-training)

## Connecting Text and Image

"CLIP is a neural network trained on a variety of (image, text) pairs. It can be instructed in natural language to predict the most relevant text snippet, given an image, without directly optimizing for the task, similar to the zero-shot capabilities of GPT-2 and 3. We found CLIP matches the performance of the original ResNet50 on ImageNet "zero-shot" without using any of the original 1.28M labeled examples, overcoming several major challenges in computer vision."

## Approach



- [openai/CLIP: CLIP \(Contrastive Language-Image Pretraining\)](https://github.com/openai/CLIP), Predict the most relevant text snippet given an image (github.com)
- [CLIP: Connecting text and images \(openai.com\)](https://openai.com/research/clip)

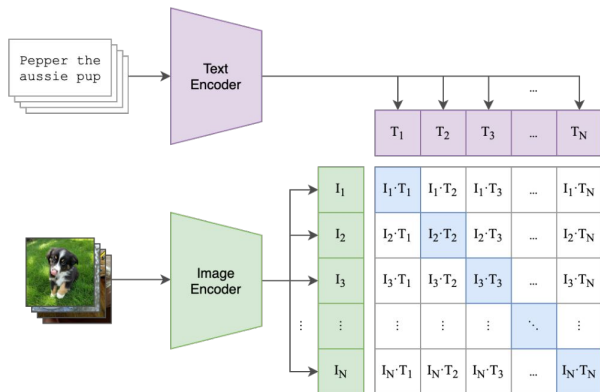
# CLIP (Contrastive Language-Image Pre-training)

## Details

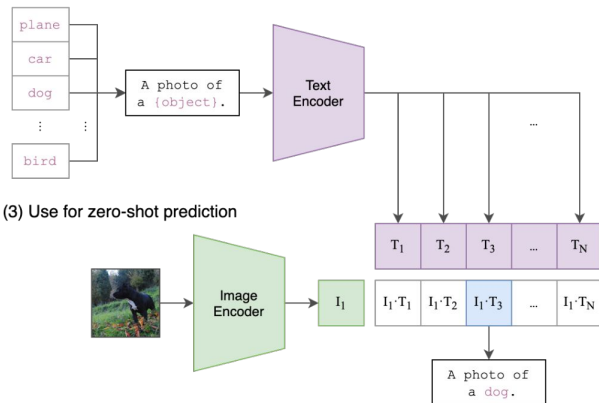
Given a batch of  $N$  (image, text) pairs, predict which of the  $N \times N$  possible (image, text) pairings across a batch actually occurred.

(Image, Text)  
Pairs

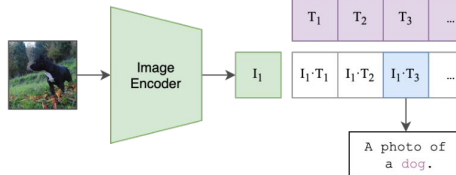
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Author constructed a new dataset of 400 million (image, text) pairs collected from a variety of publicly available sources on the Internet

Multi-model embedding space  
Maximize the cosine similarity of  
image vs text

Optimized for a symmetric cross  
entropy loss.

- [openai/CLIP: CLIP \(Contrastive Language-Image Pretraining\), github.com](https://github.com/openai/CLIP)
- [CLIP: Connecting text and images \(openai.com\)](https://openai.com/research/clip)

# Notebook

# 5

## Deploying pretrained VLM to FastAPI

Hands On



# Notebook