



SpotOn

SpotOn

**Progress #2**



# Agenda

- Recap
- Architecture
  - CNN
    - One-stage v.s. Two-stage CNN
  - ViT
- Update the model's latency when using GPU
- Re-Identification
  - Re-Identification model selection
  - Re-Identification model experimental
- Camera view selection

# Recap

# Detection models experimental

## Feedback & Solutions



### Privacy when deploying



### Stakeholder

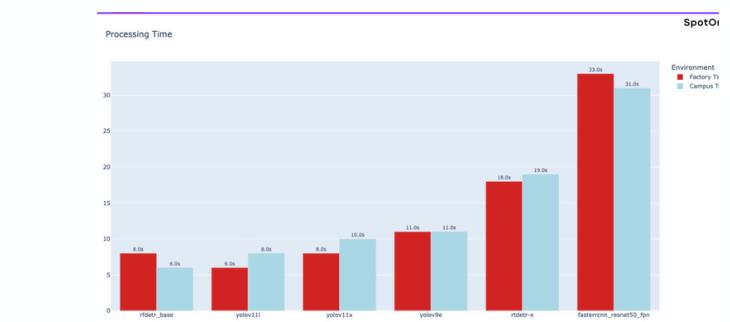
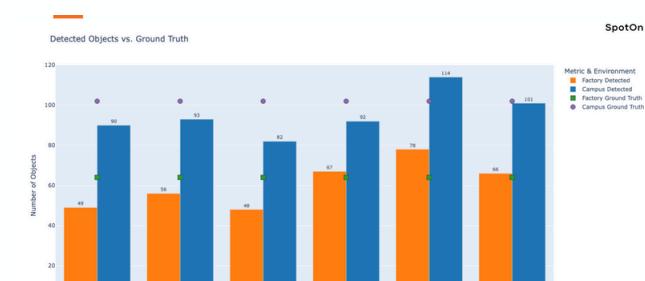
• Stakeholders involved in the project include the data providers, distributors of the Multi-Target Multi-Camera (MTMMC) dataset used for training and evaluating the AI components of this project (4.5.3).

• Stakeholders: They have an interest in ensuring the dataset is used responsibly and ethically, strictly according to the terms and conditions outlined in the usage agreement signed by the project team. This includes preventing misuse or unauthorised redistribution of the data.

• Privacy Considerations: All stakeholders agree to all stipulations in the dataset usage license. Compliance will be documented, and proper attribution/citation will be provided in all relevant project outputs.

### Model experimental

- One stage detector (Bounding boxes / Classes in a single pass)
  - yolov11
  - yolov11x
  - yolov9e
- Transformer based (Self-attention mechanisms)
  - rtdetr-x
  - rfdetr-l
- Two stage detector (Rois → Bounding boxes / Classes)



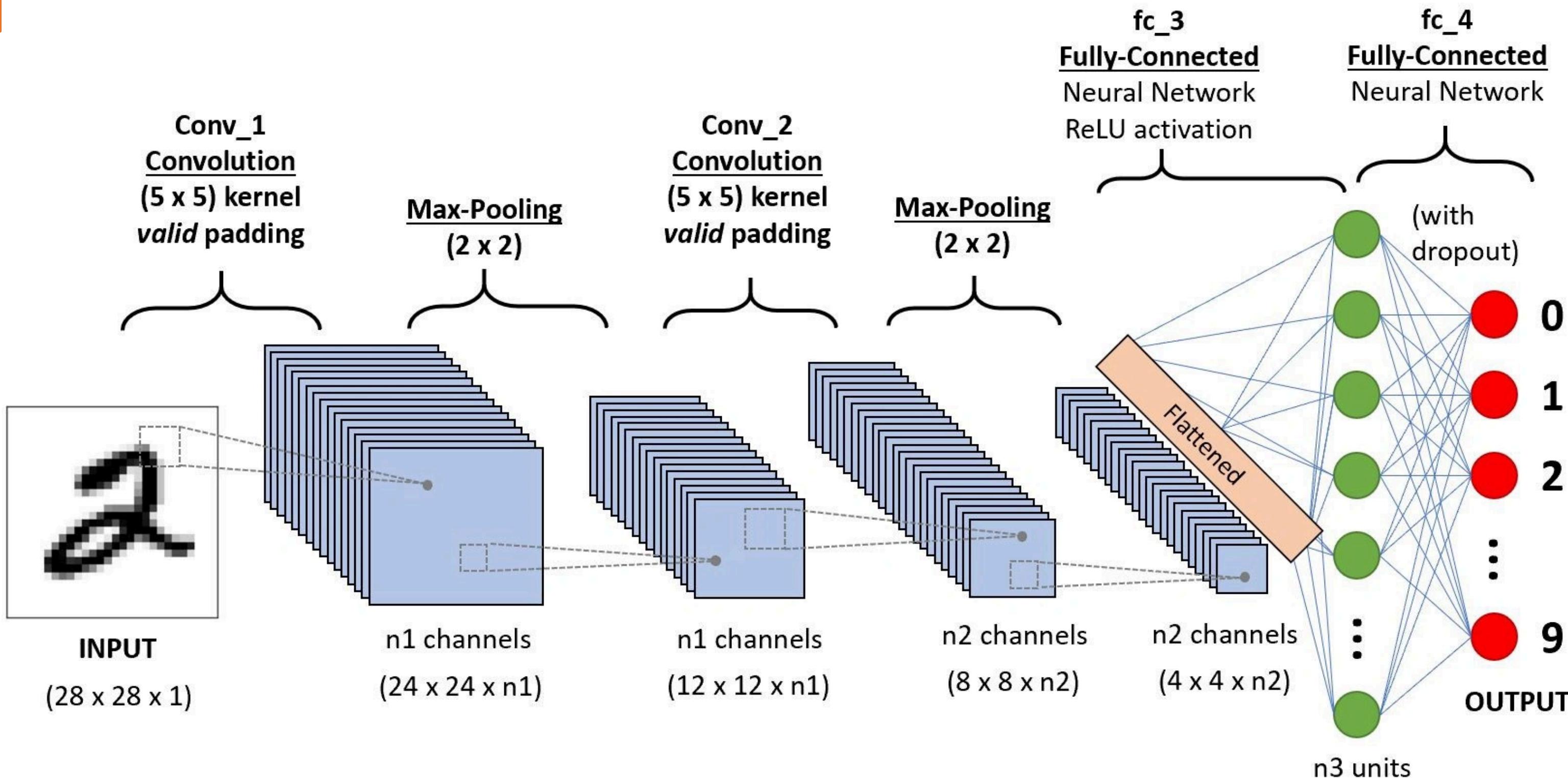
### Impus



# Architecture (Simple term)

# CNN (1989)

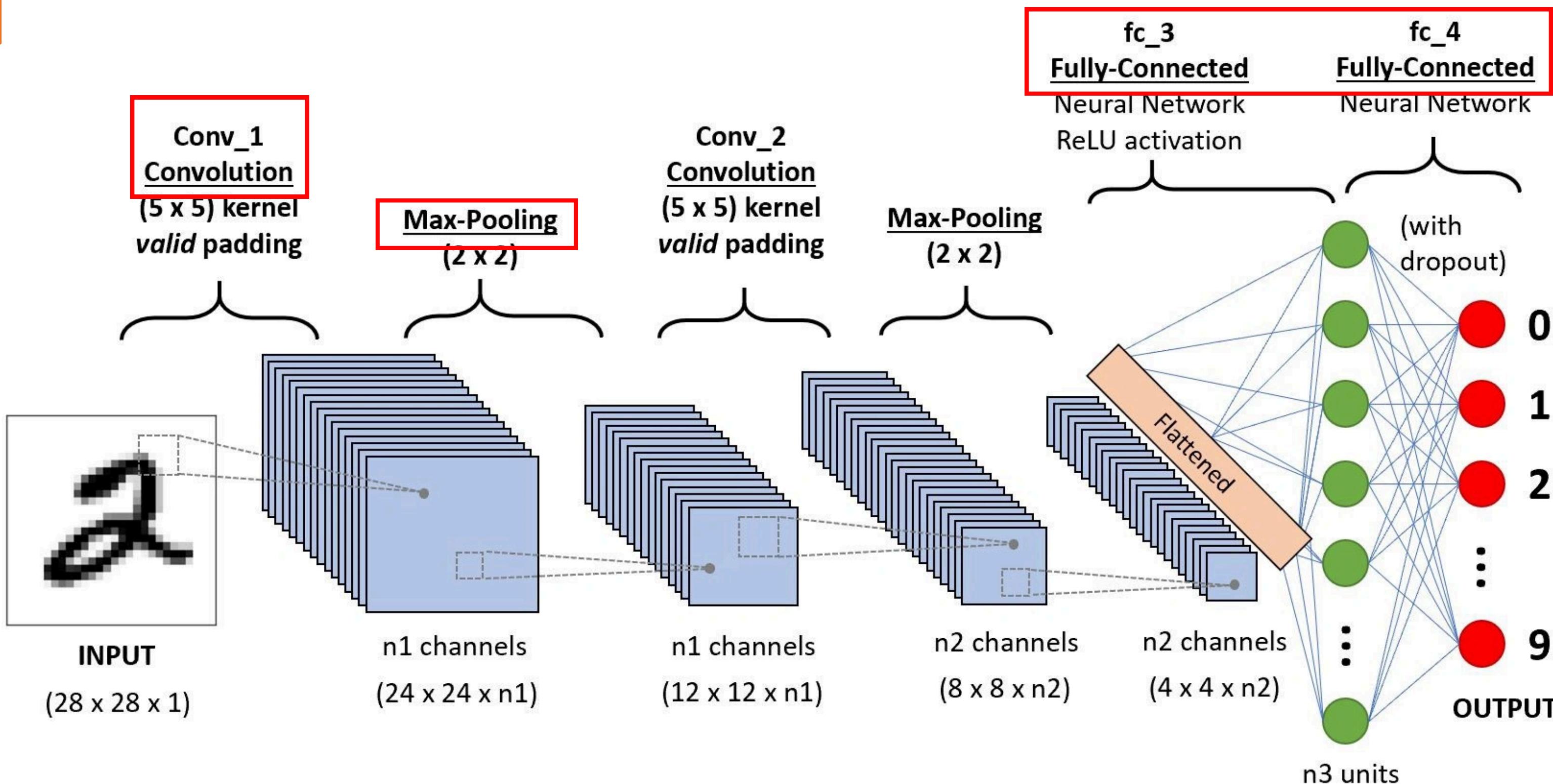
SpotOn



# CNN (1989)

SpotOn

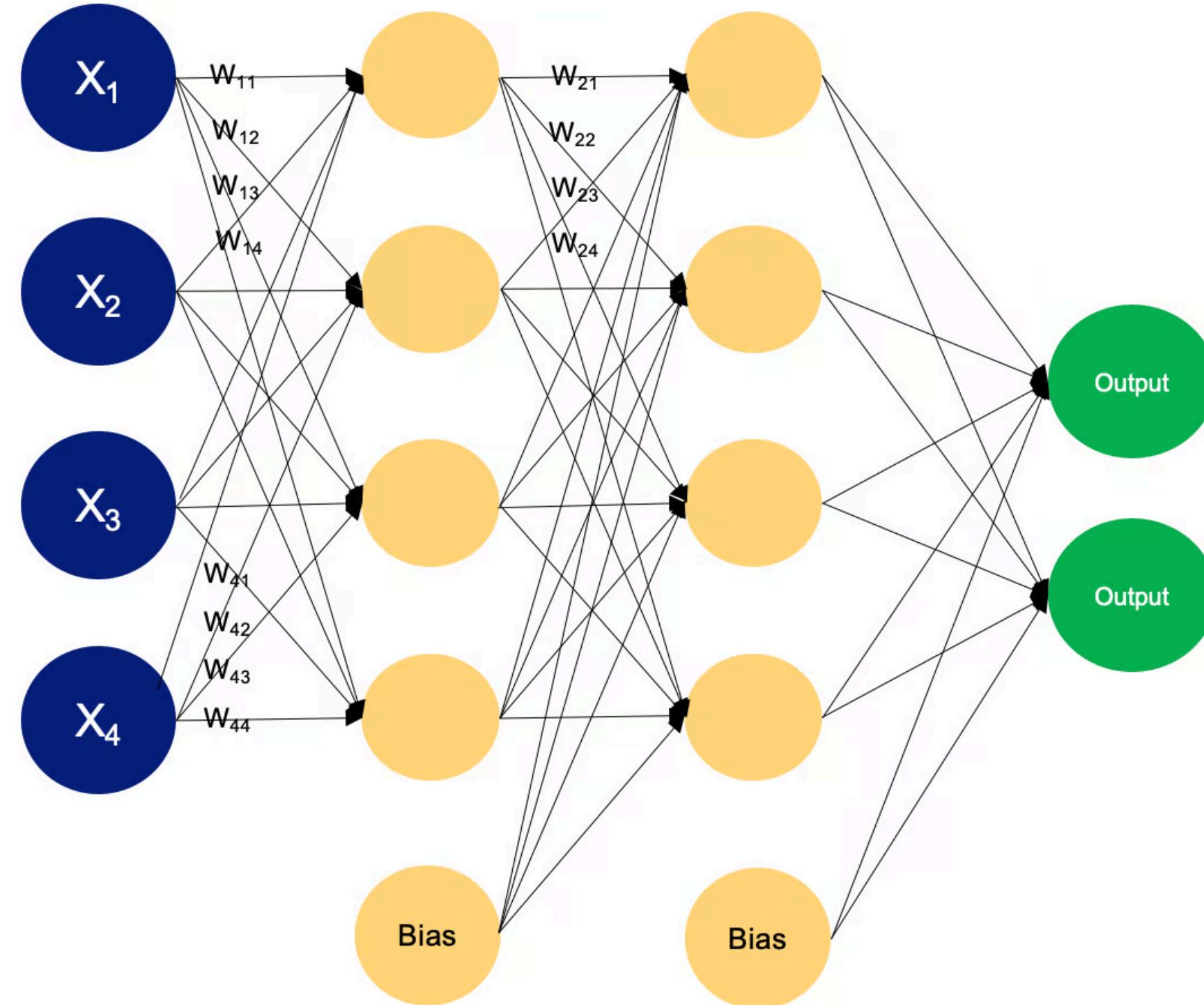
MLP



# MLP

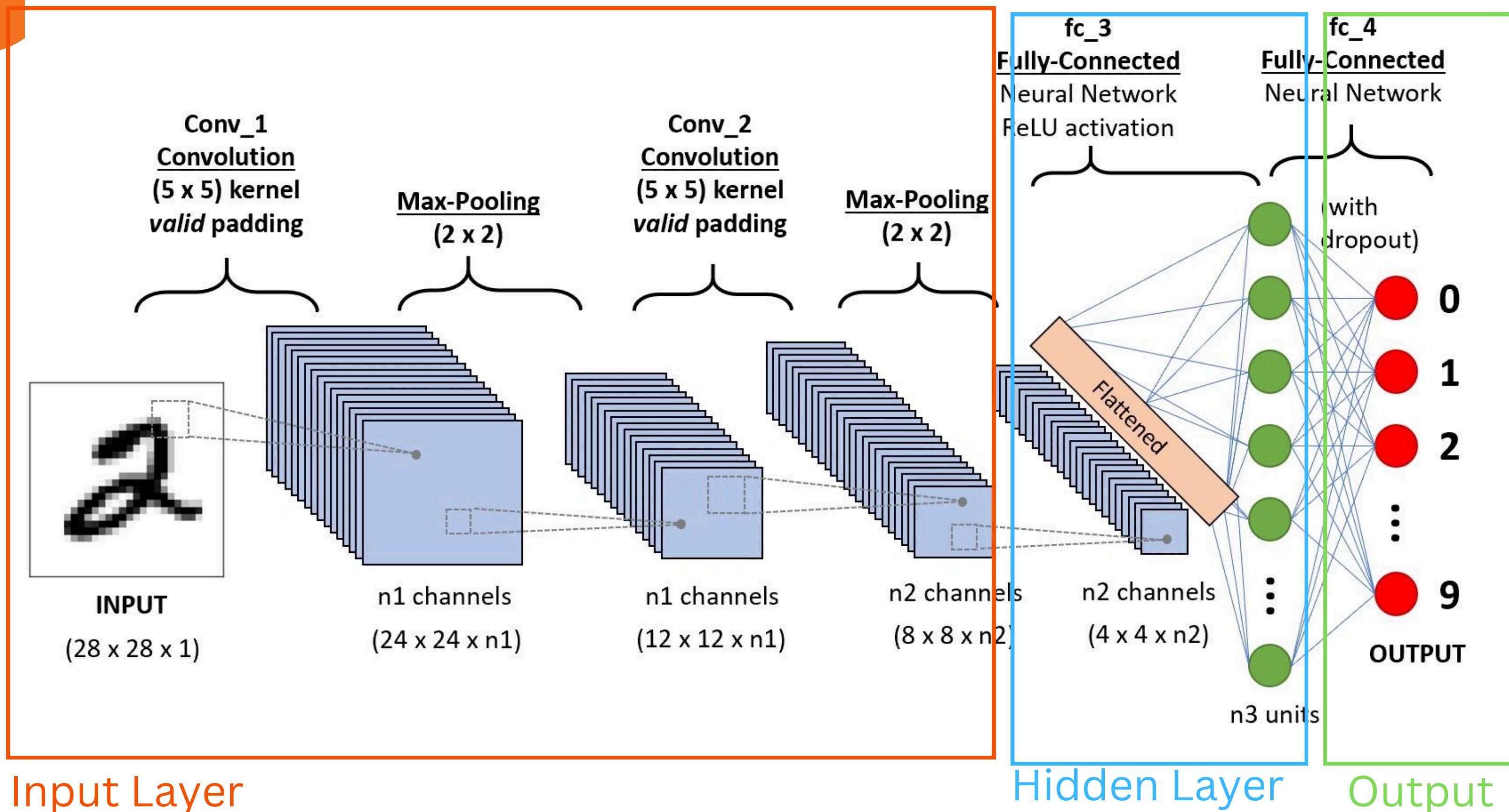
SpotOn

Inner layer | Hidden layers | Outer layer

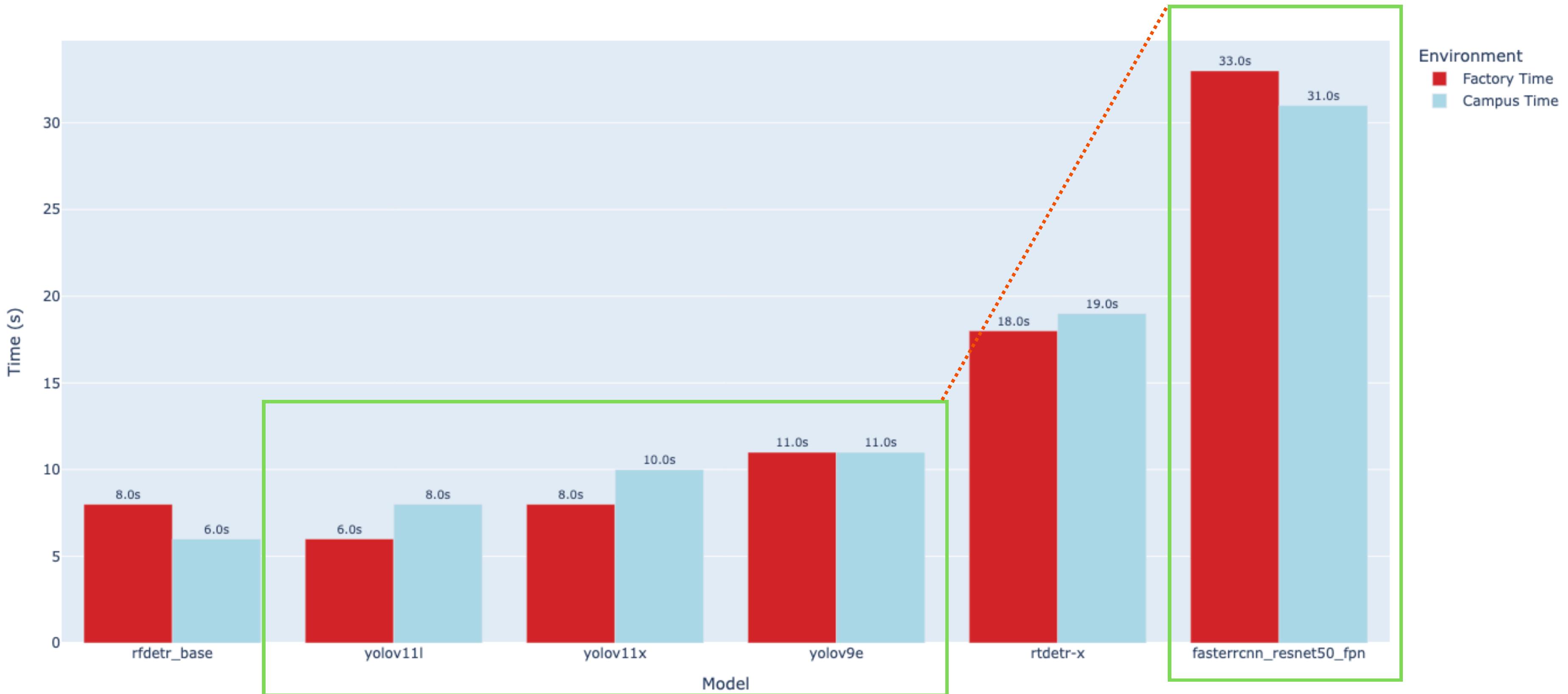


# MLP

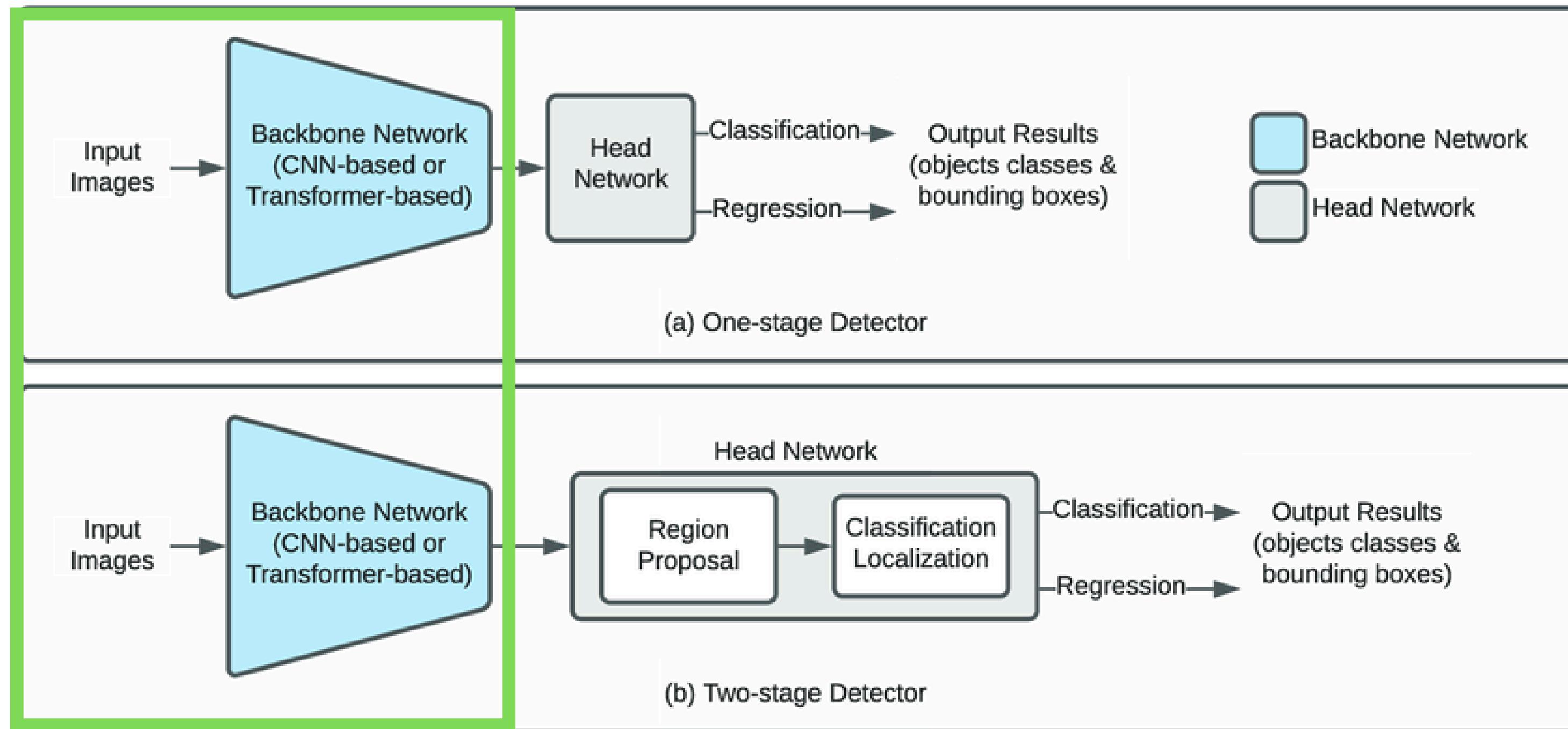
SpotOn



## Processing Time

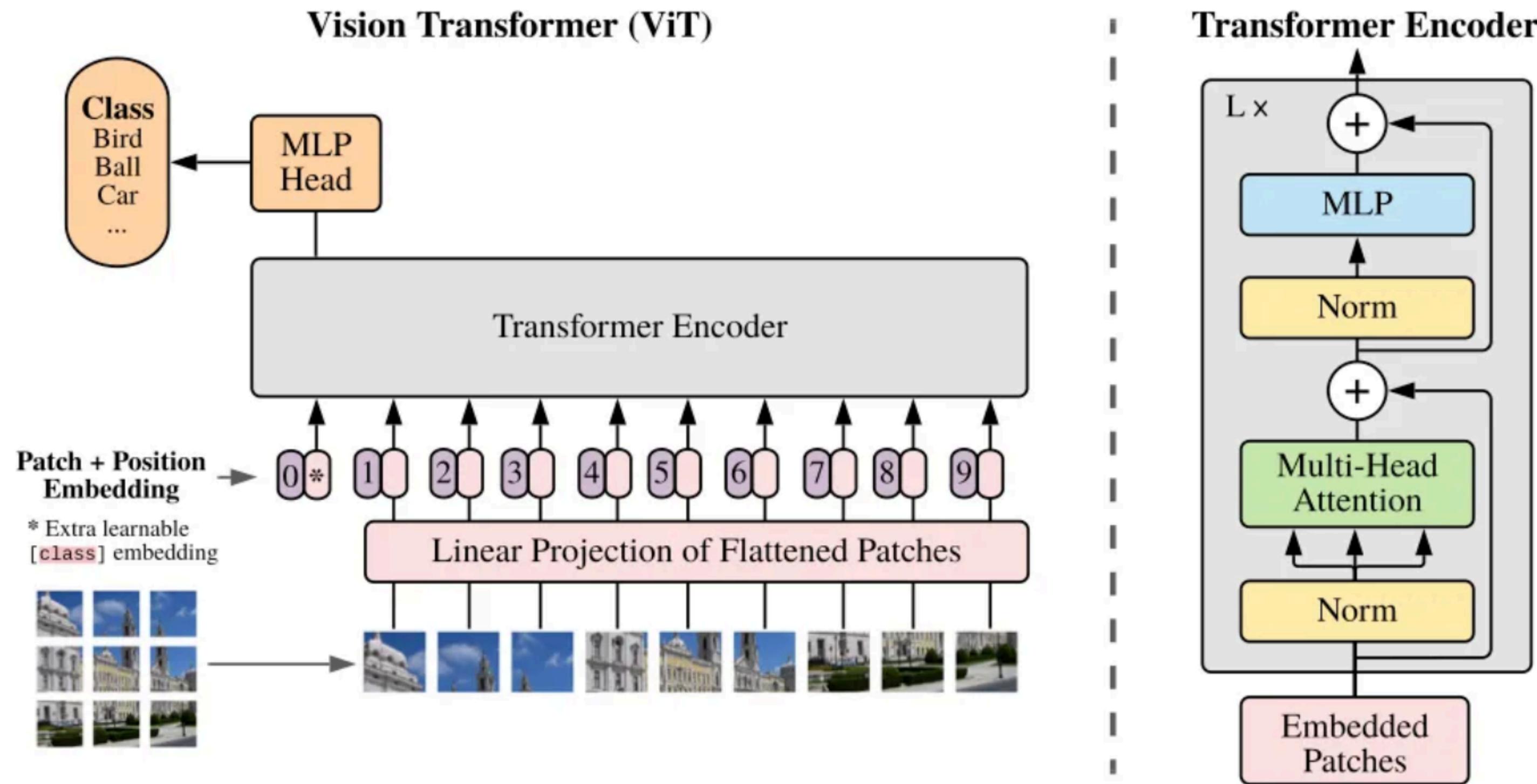


# One/Two Stage



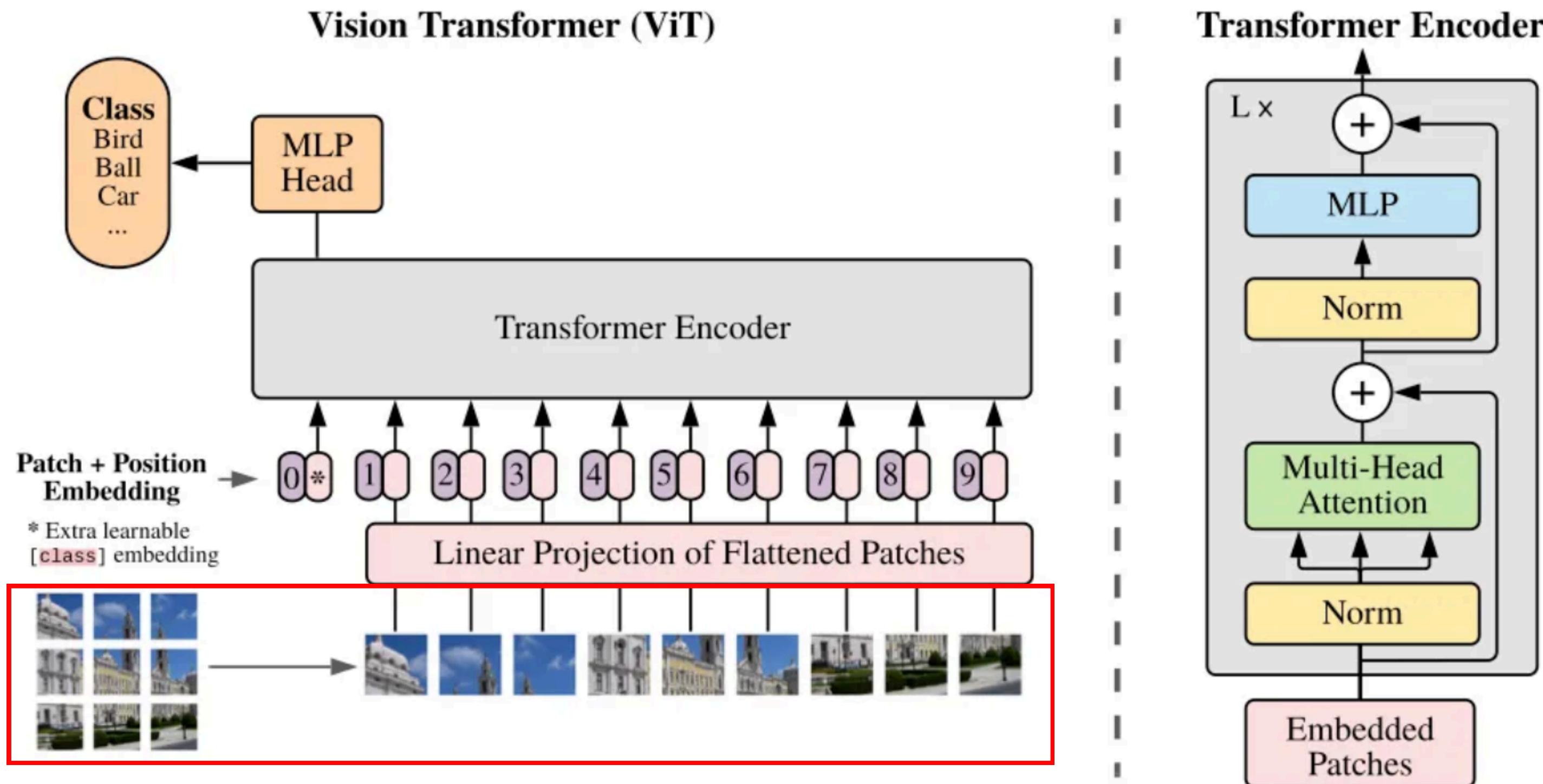
# ViT (2020)

SpotOn

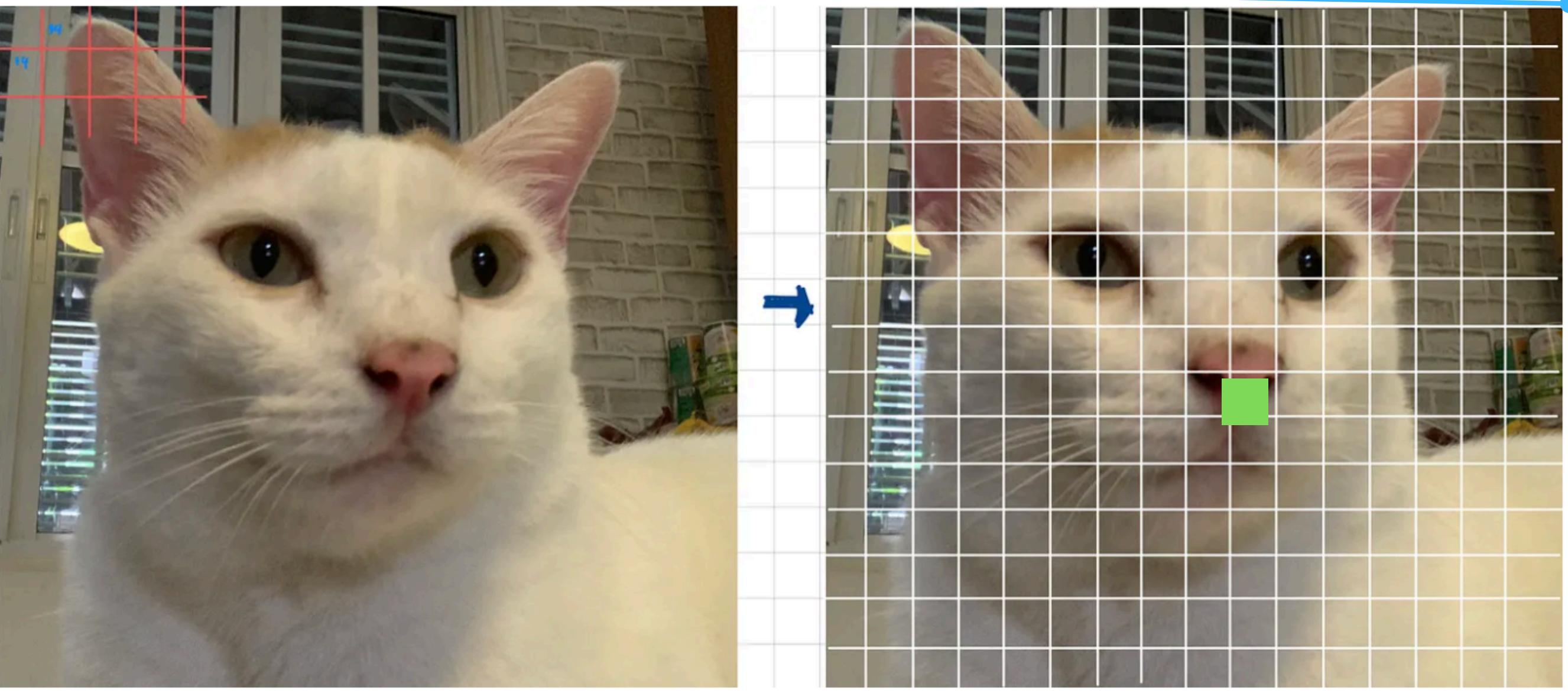


# ViT (2020)

SpotOn



# 1. Split an image into patches



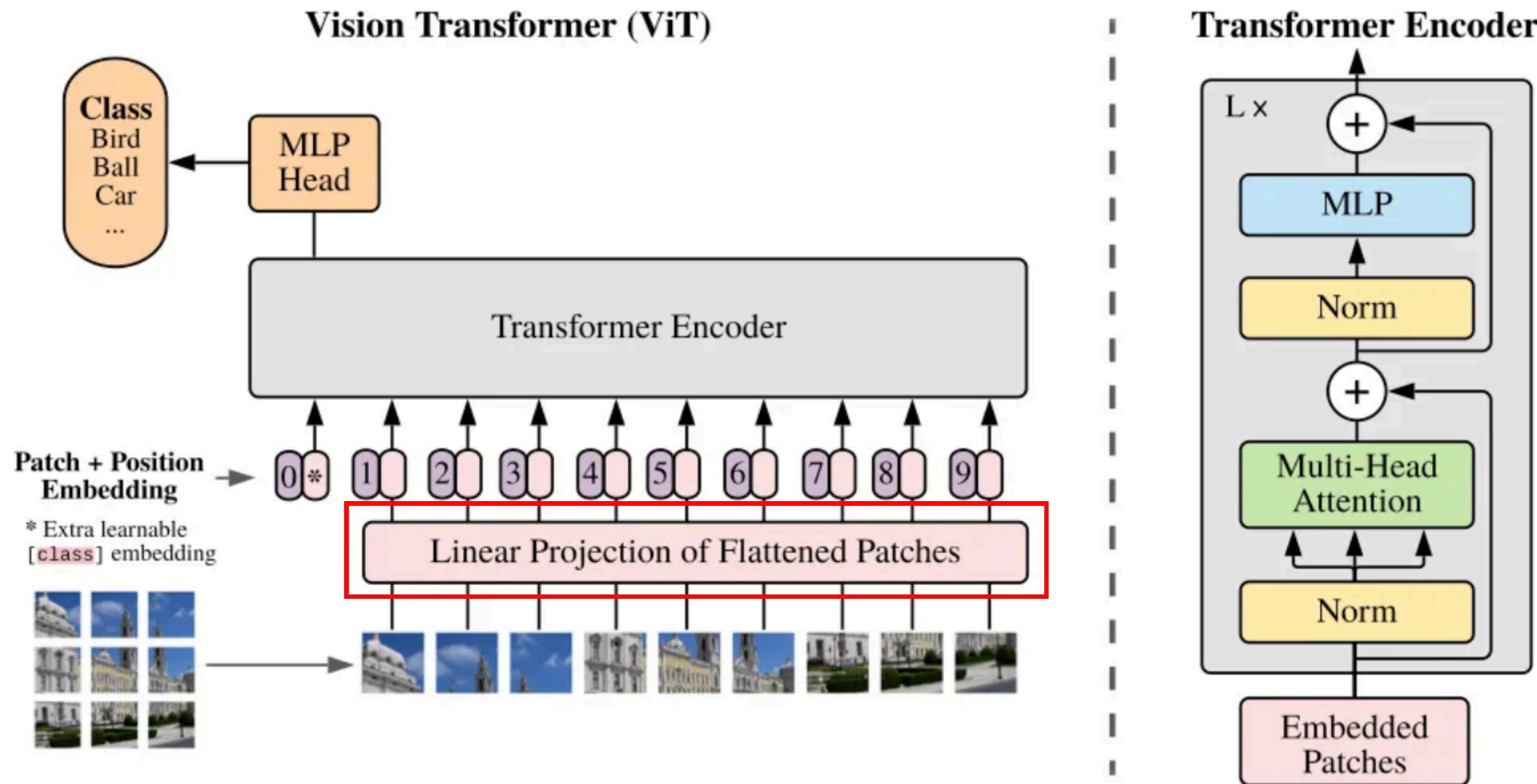
224 x 224 x 3 (RGB)  
(H x W x C)



N (number) =  $16 \times 16 = 256$   
P (patch size) =  $14 \times 14 = 196$

# ViT (2020)

SpotOn



## 2. Flattening Each Patch

2 (height) x 2 (width) x 3 (channels - R, G, B)

- Red Channel (Layer 0):

```
[[ R1, R2 ],  
 [ R3, R4 ]]
```

Example values: [[10, 20], [30, 40]]

- Green Channel (Layer 1):

```
[[ G1, G2 ],  
 [ G3, G4 ]]
```

Example values: [[51, 61], [71, 81]]

- Blue Channel (Layer 2):

```
[[ B1, B2 ],  
 [ B3, B4 ]]
```

Example values: [[92, 82], [72, 62]]

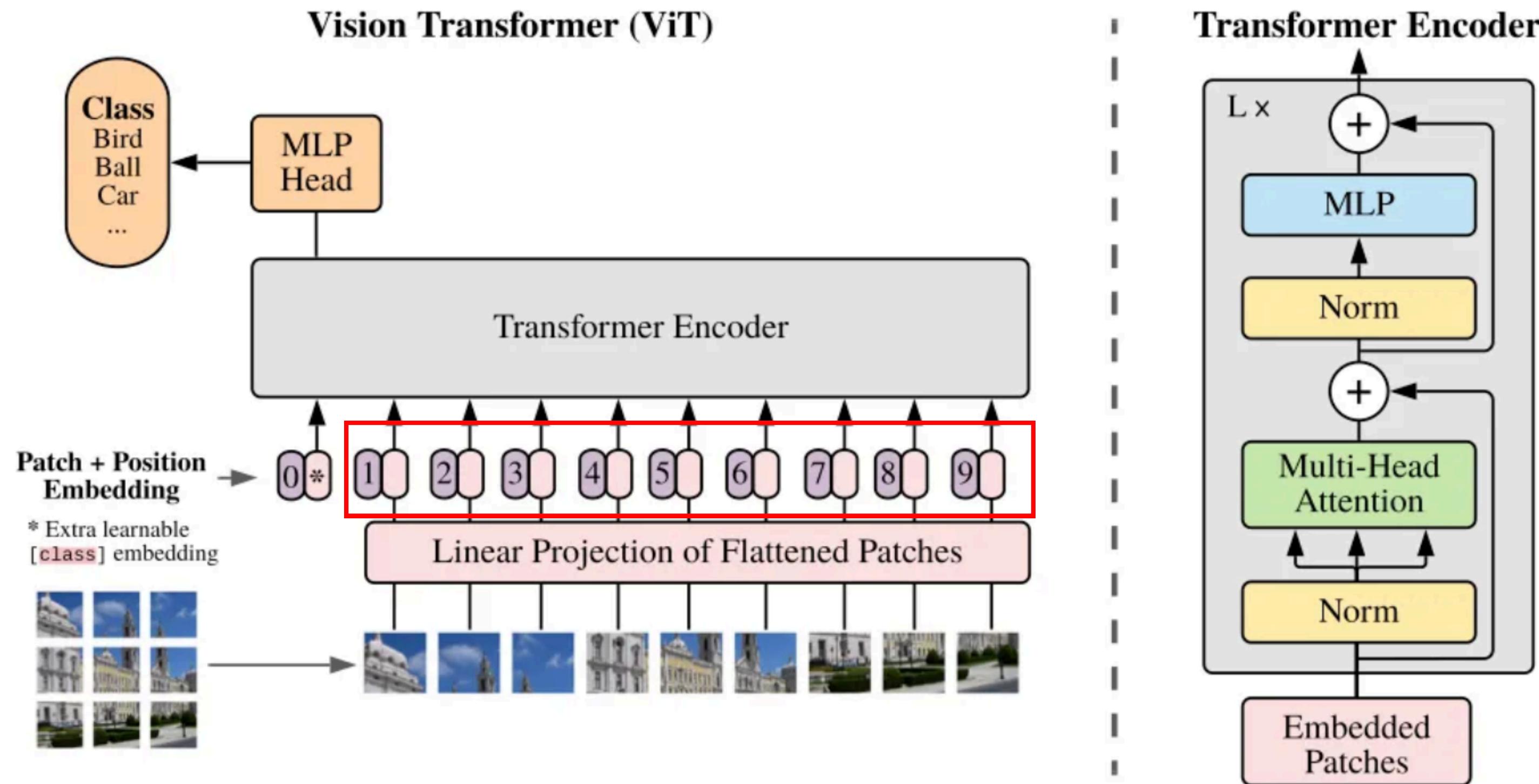
Flattened Vector = [[R1,G1,B1], [R2,G2,B2], [R3,G3,B3], [R4,G4,B4]]

Flattened Vector = [10, 51, 92, 20, 61, 82, 30, 71, 72, 40, 81, 62]

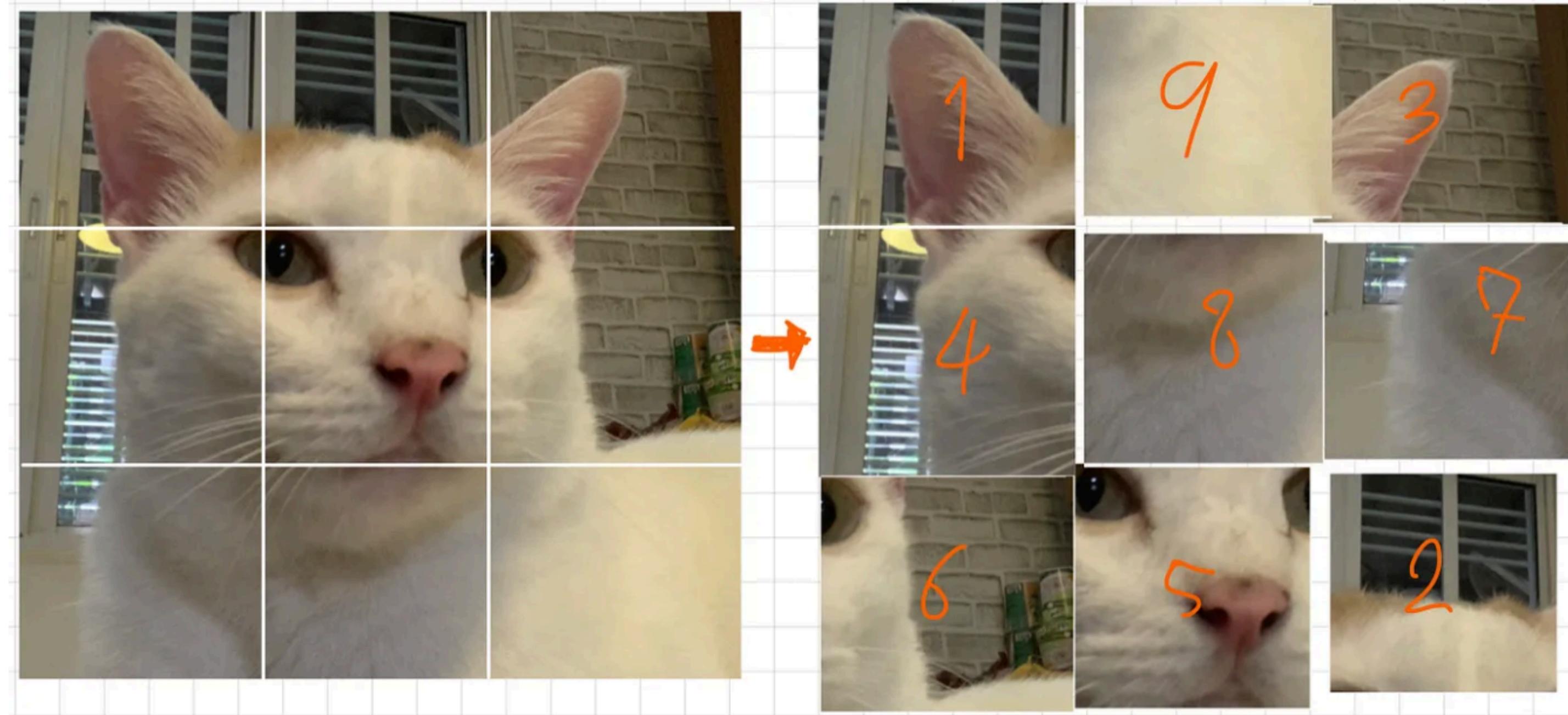
```
[  
  [vector_for_patch_1], // e.g., [10, 51, 92, ..., 62]  
  [vector_for_patch_2], // e.g., [15, 55, 95, ..., 65]  
  [vector_for_patch_3], // e.g., [12, 52, 92, ..., 60]  
  ...  
  ... N times (e.g., 256 times) ...  
  ...  
  [vector_for_patch_N] // e.g., [18, 58, 98, ..., 68]  
]
```

# ViT (2020)

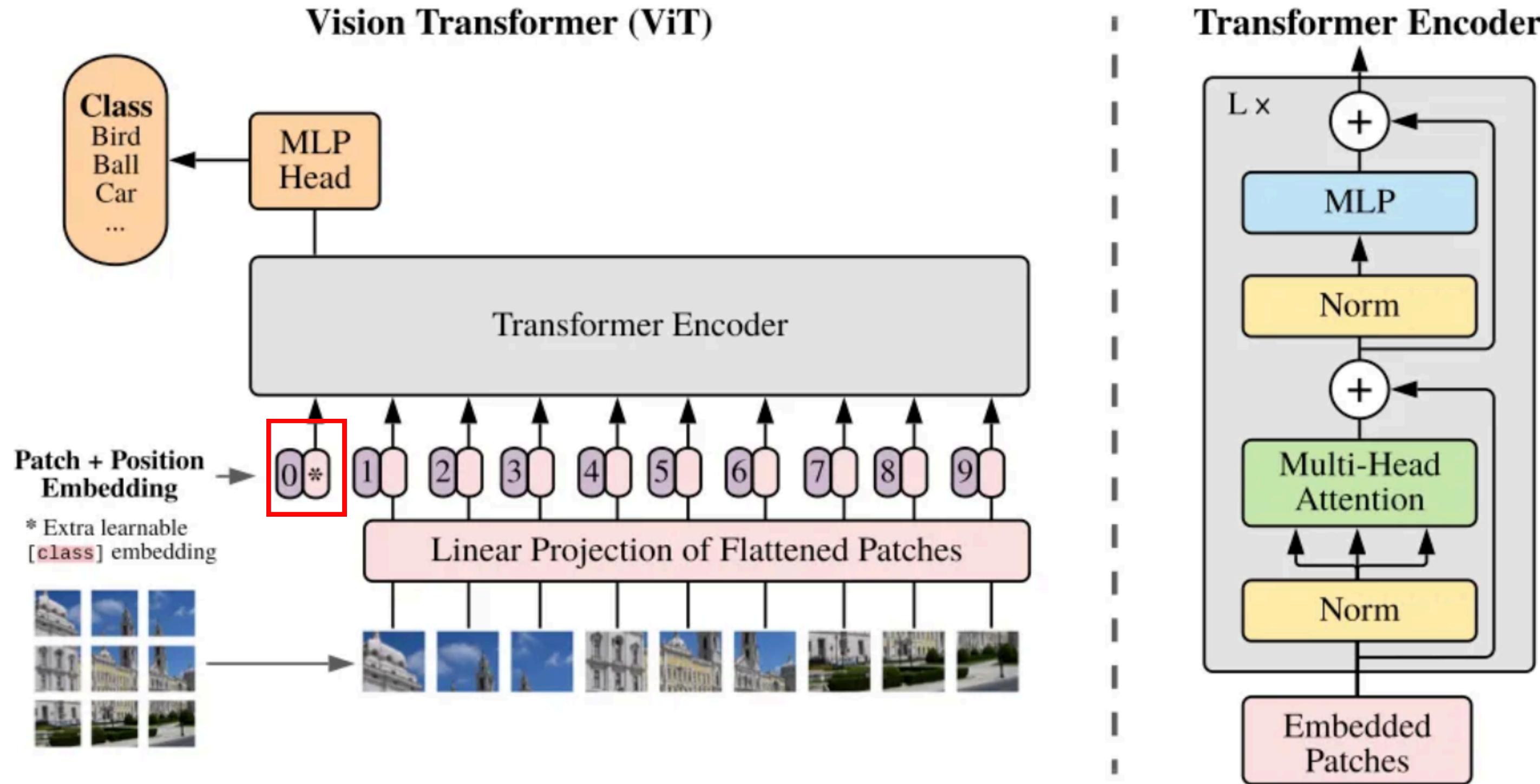
SpotOn



### 3. Add positional embeddings

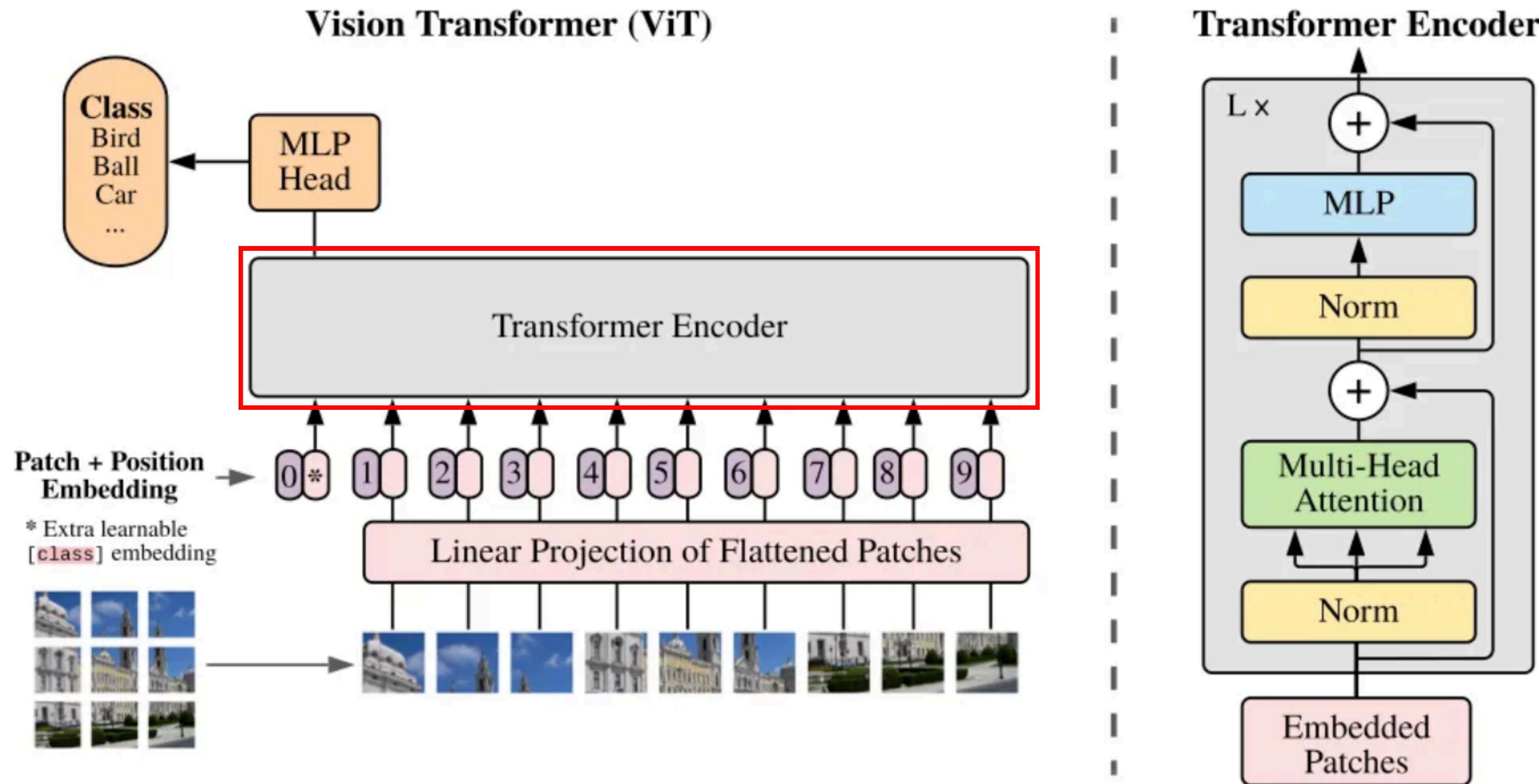


# 4. Add Class token for classification



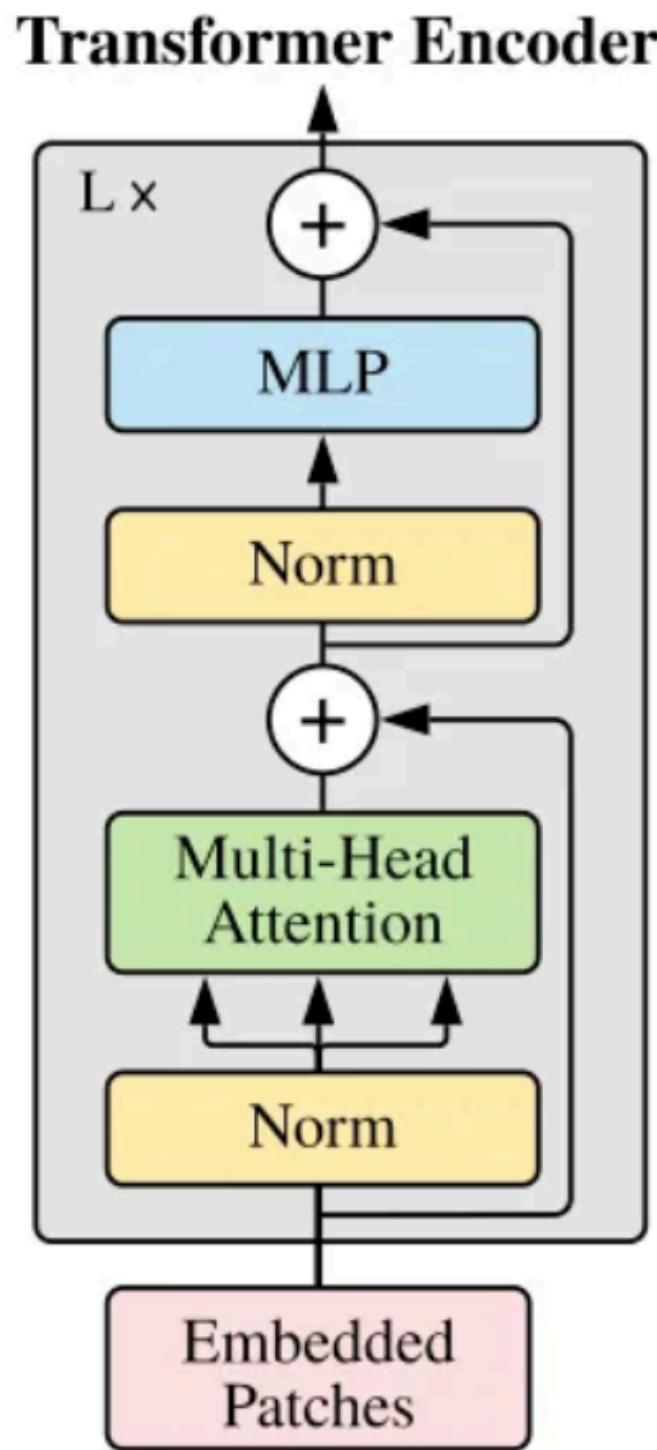
# ViT (2020)

SpotOn



# 5. Feed the sequence as an input to standard transformer encoder

SpotOn



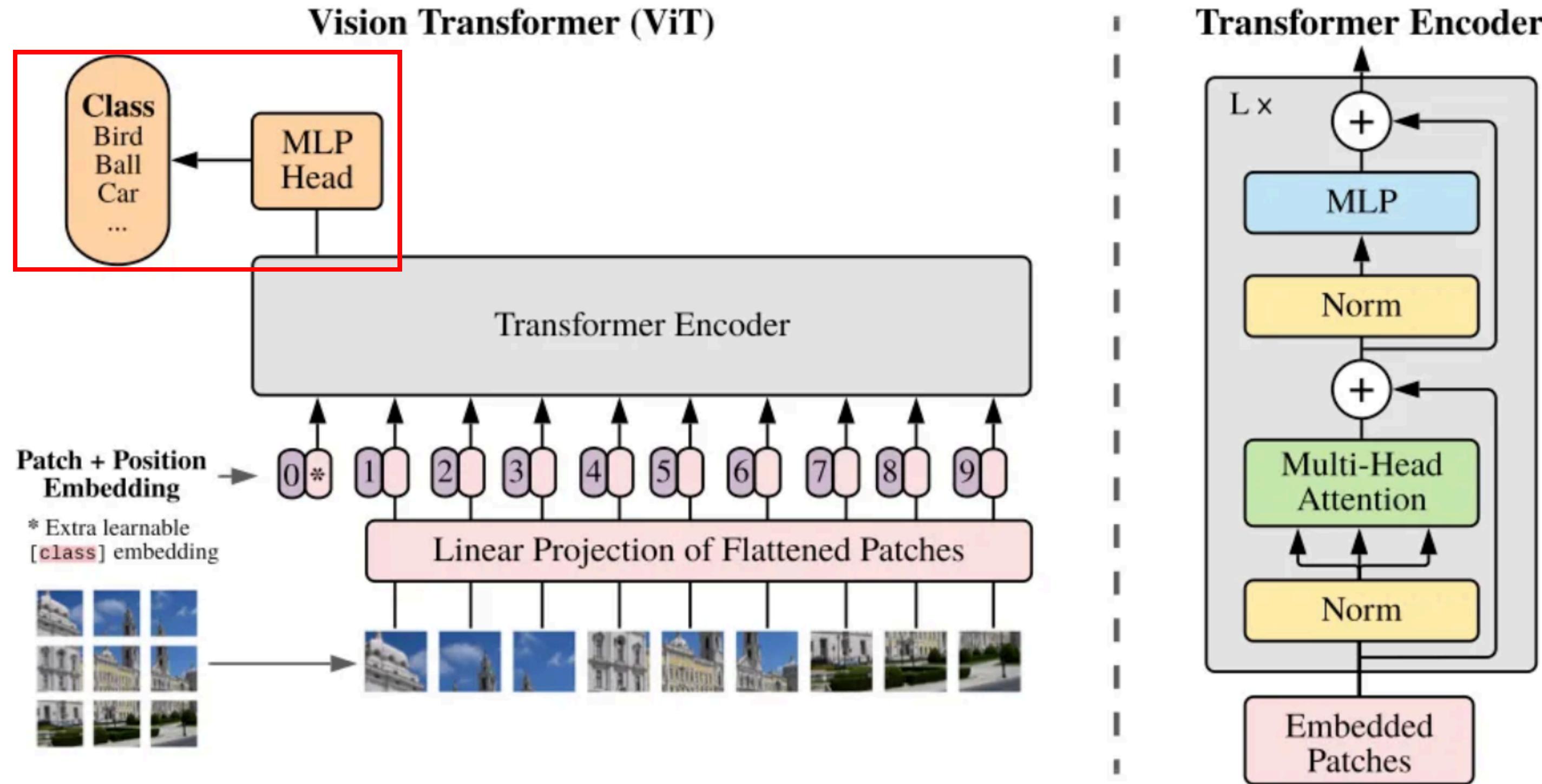
## Multi-Head Attention

1. Look at all patches
2. Calculate attention score (important)
3. Repeat multiple times in parallel

This allows the model to understand the context of the whole image e.g. A patch of a cat's ear can learn it's related to a patch of the cat's eye

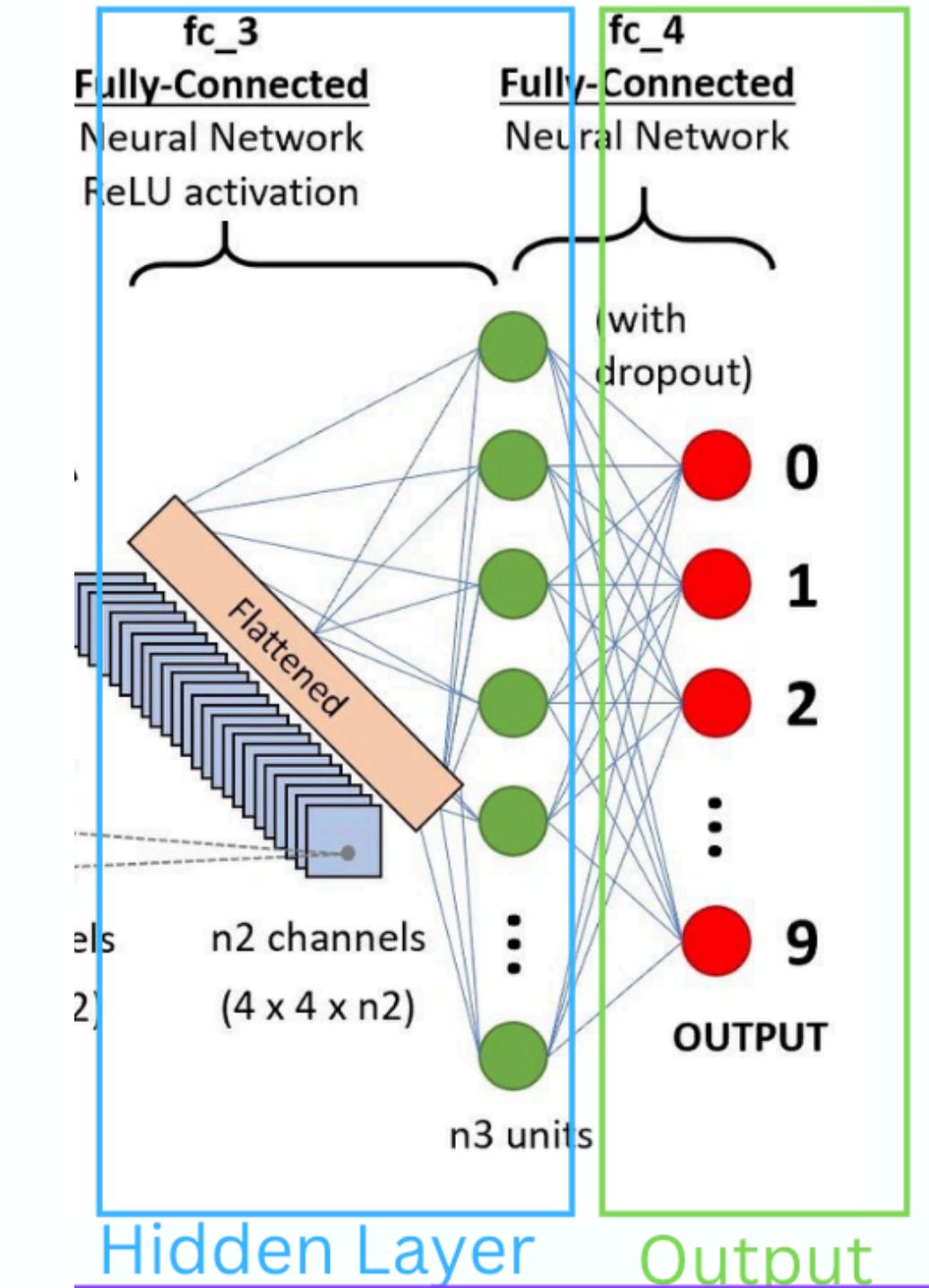
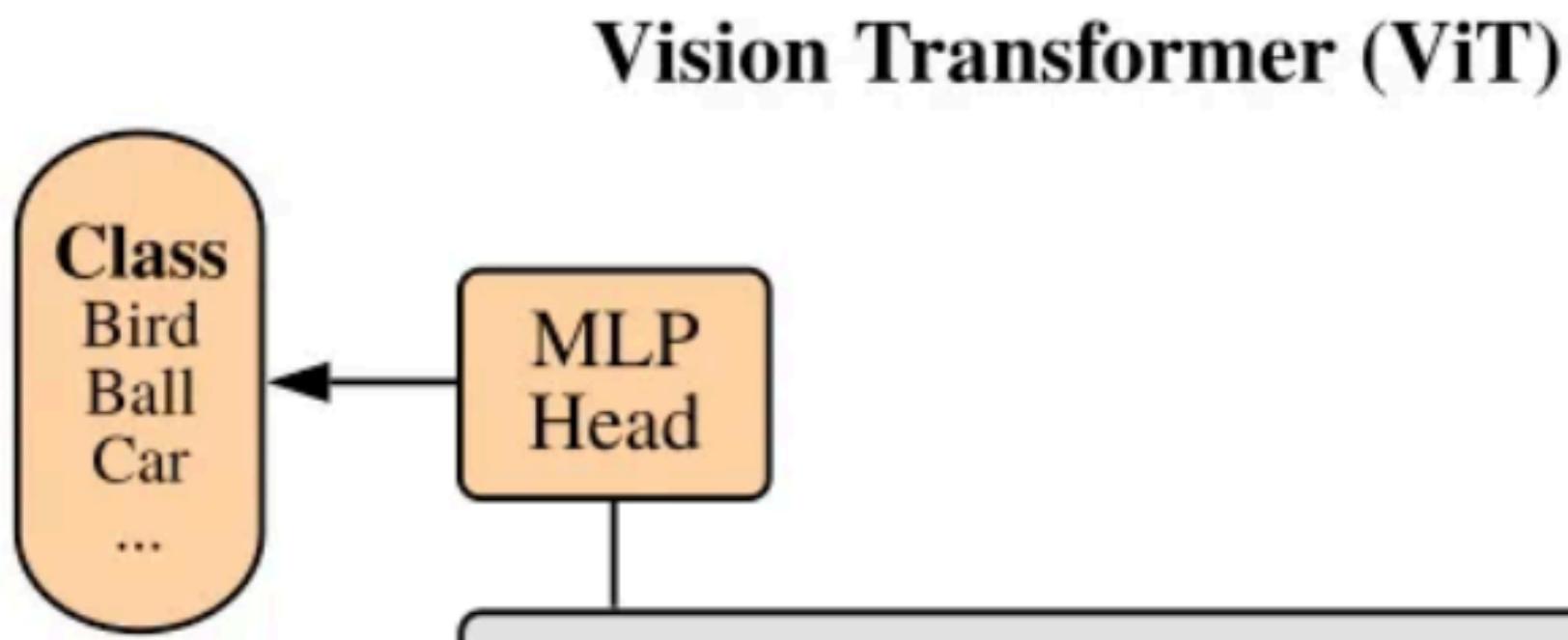
# Classify

SpotOn



# Classify

SpotOn



-is-the-convolutional-neural-network-archite

Update the model's  
latency when using GPU

6s

# Re-Identification

# Model selection

## Benchmarks

[Add a Result](#)

These leaderboards are used to track progress in Person Re-Identification

Trend	Dataset	Best Model	Paper	Code	Compare
	Market-1501	st-ReID(RE, RK)			<a href="#">See all</a>
	DukeMTMC-reID	DenseLL			<a href="#">See all</a>
	MSMT17	CLIP-ReID (with re-ranking)			<a href="#">See all</a>
	Occluded-DukeMTMC	KPR + SOLIDER			<a href="#">See all</a>
	Market-1501-C	TransReID			<a href="#">See all</a>
	MARS	B-BOT + OSM + CL Centers* (Re-rank)			<a href="#">See all</a>
	CUHK03 labeled	Weakly Pre-training (ResNet101+RK)			<a href="#">See all</a>
	CUHK03	Proposed SGGNN			<a href="#">See all</a>
	CUHK03 detected	Top-DB-Net + RK			<a href="#">See all</a>
	PRID2011	B-BOT + Attention and CL Loss*			<a href="#">See all</a>

# Find a benchmarks

<https://paperswithcode.com/task/person-re-identification>

## Benchmarks

[Add a Result](#)

These leaderboards are used to track progress in Person Re-Identification

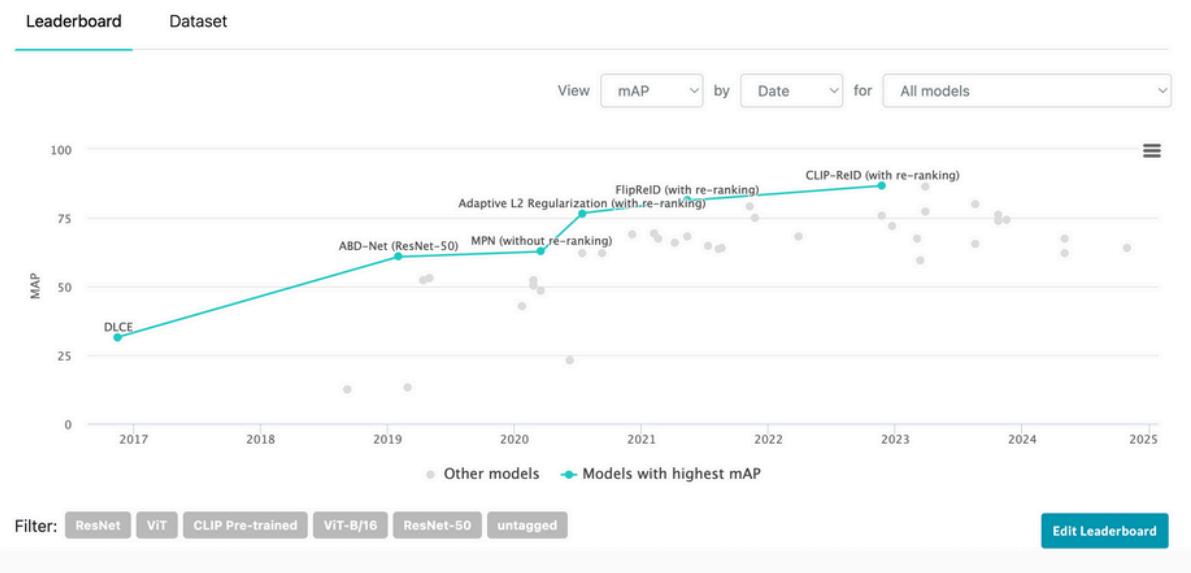
Trend	Dataset	Best Model	Paper	Code	Compare
	Market-1501	st-ReID(RE, RK)	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	DukeMTMC-reID	DenseLL	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	MSMT17	CLIP-ReID (with re-ranking)	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	Occluded-DukeMTMC	KPR + SOLIDER	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	Market-1501-C	TransReID	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	MARS	B-BOT + OSM + CL Centers* (Re-rank)	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	CUHK03 labeled	Weakly Pre-training (ResNet101+RK)			<a href="#">See all</a>
	CUHK03	Proposed SGGNN	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	CUHK03 detected	Top-DB-Net + RK	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>
	PRID2011	B-BOT + Attention and CL Loss*	<a href="#">Paper</a>	<a href="#">Code</a>	<a href="#">See all</a>

# Filter similar dataset

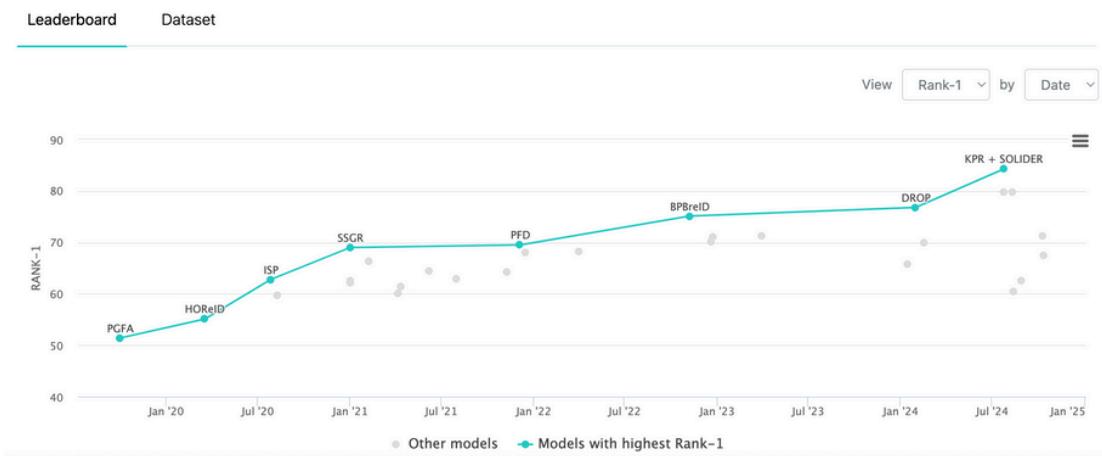
All of them are similar btw.

# Research on interesting models

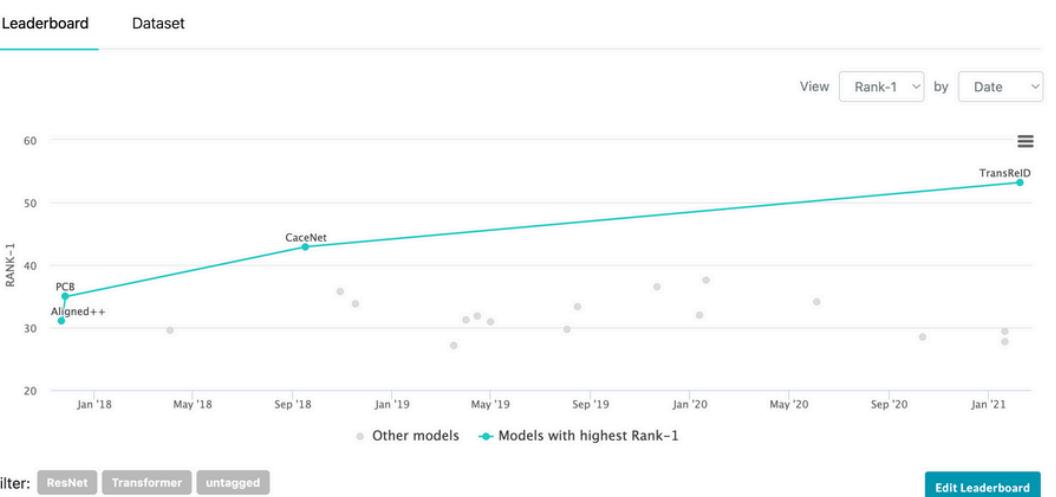
Person Re-Identification on MSMT17



Person Re-Identification on Occluded-DukeMTMC



Person Re-Identification on Market-1501-C



- General Accuracy: **CLIP-ReID**
- High Speed/Low Latency: **OSNet**
- Occlusion: **KPR + SOLIDER**

Compare  
the models

- General Accuracy: **CLIP-ReID**
- High Speed/Low Latency: **OSNet**
- Occlusion: **KPR + SOLIDER**
  - Too much complex
  - Too much resource consuming

Pick the  
models

# Model experimental

Comparison Results (Threshold: 0.7)  
OSNet: 0.5357 (Different Person)  
CLIP: 0.3113 (Different Person)

c1\_000352.jpg  
Person 3



c3\_000000.jpg  
Person 2

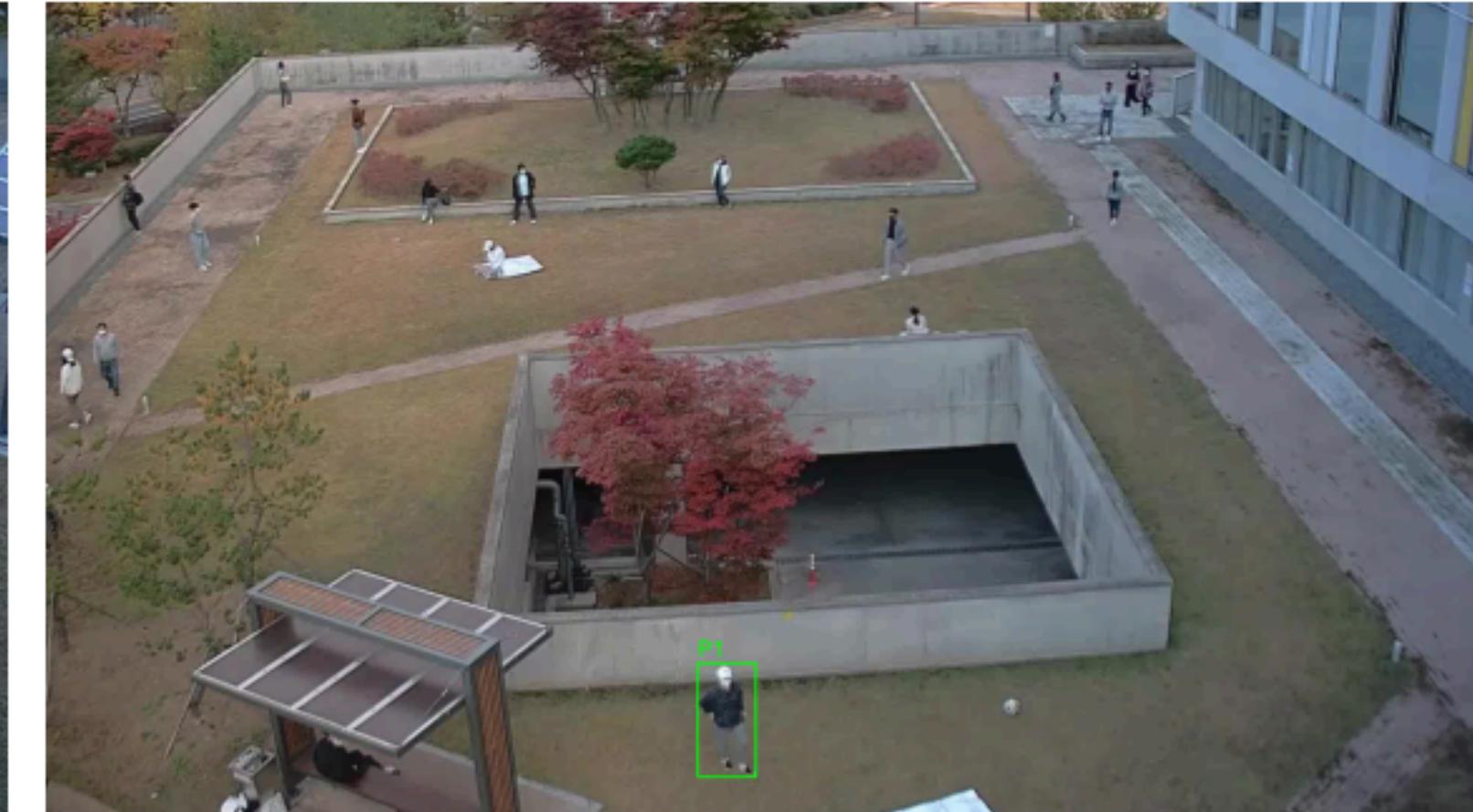


Comparison Results (Threshold: 0.7)  
OSNet: 0.7089 (Same Person)  
CLIP: 0.6850 (Different Person)

c1\_000352.jpg  
Person 3



c3\_000000.jpg  
Person 1



Comparison Results (Threshold: 0.7)  
OSNet: 0.8078 (Same Person)  
CLIP: 0.8346 (Same Person)

c1\_000217.jpg  
Person 1



c1\_000352.jpg  
Person 1



Comparison Results (Threshold: 0.7)  
OSNet: 0.6874 (Different Person)  
CLIP: 0.5773 (Different Person)

c1\_000217.jpg  
Person 1



c1\_000352.jpg  
Person 2



Comparison Results (Threshold: 0.7)  
OSNet: 0.7854 (Same Person)  
CLIP: 0.7958 (Same Person)

c1\_000217.jpg  
Person 2



c1\_000352.jpg  
Person 3

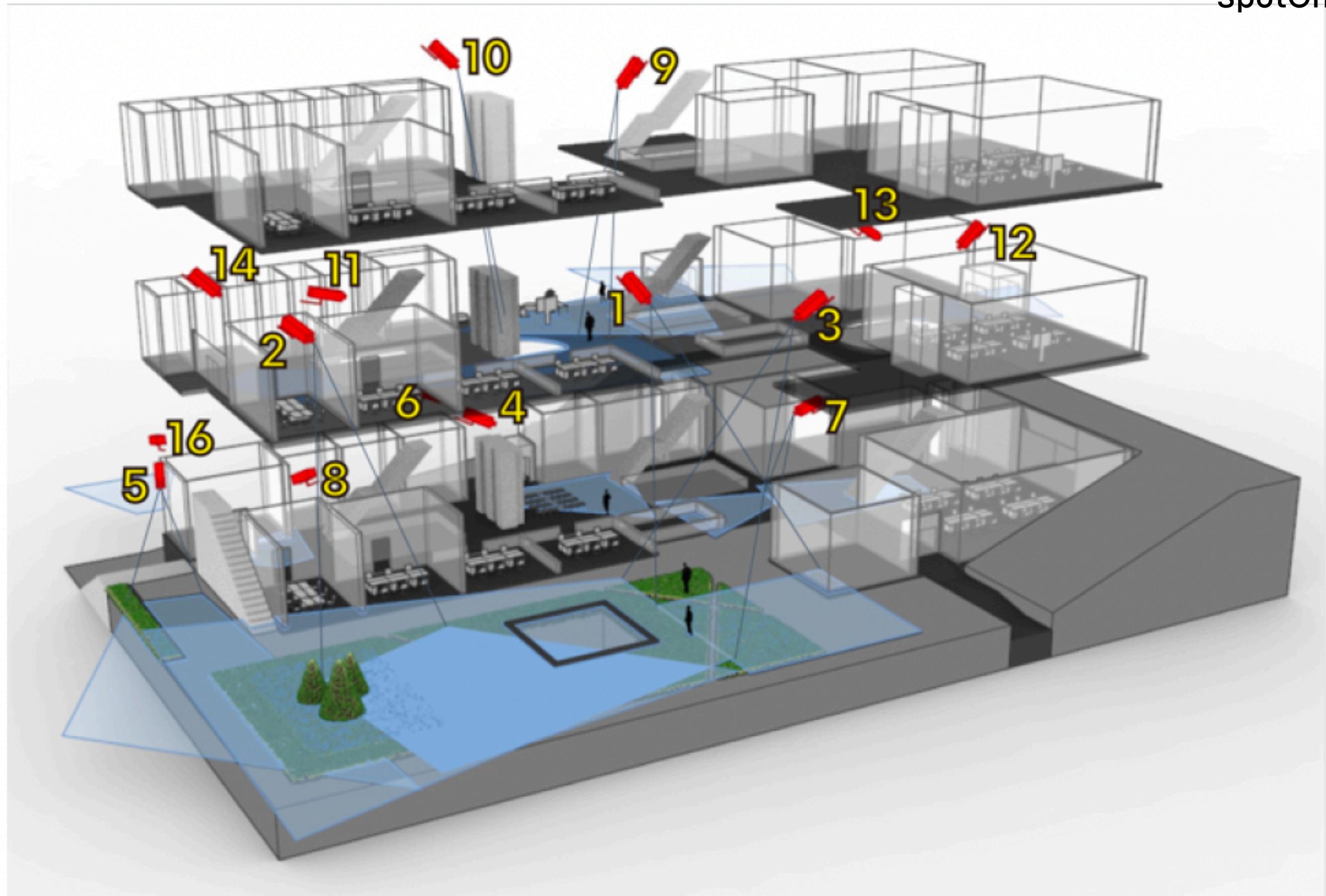


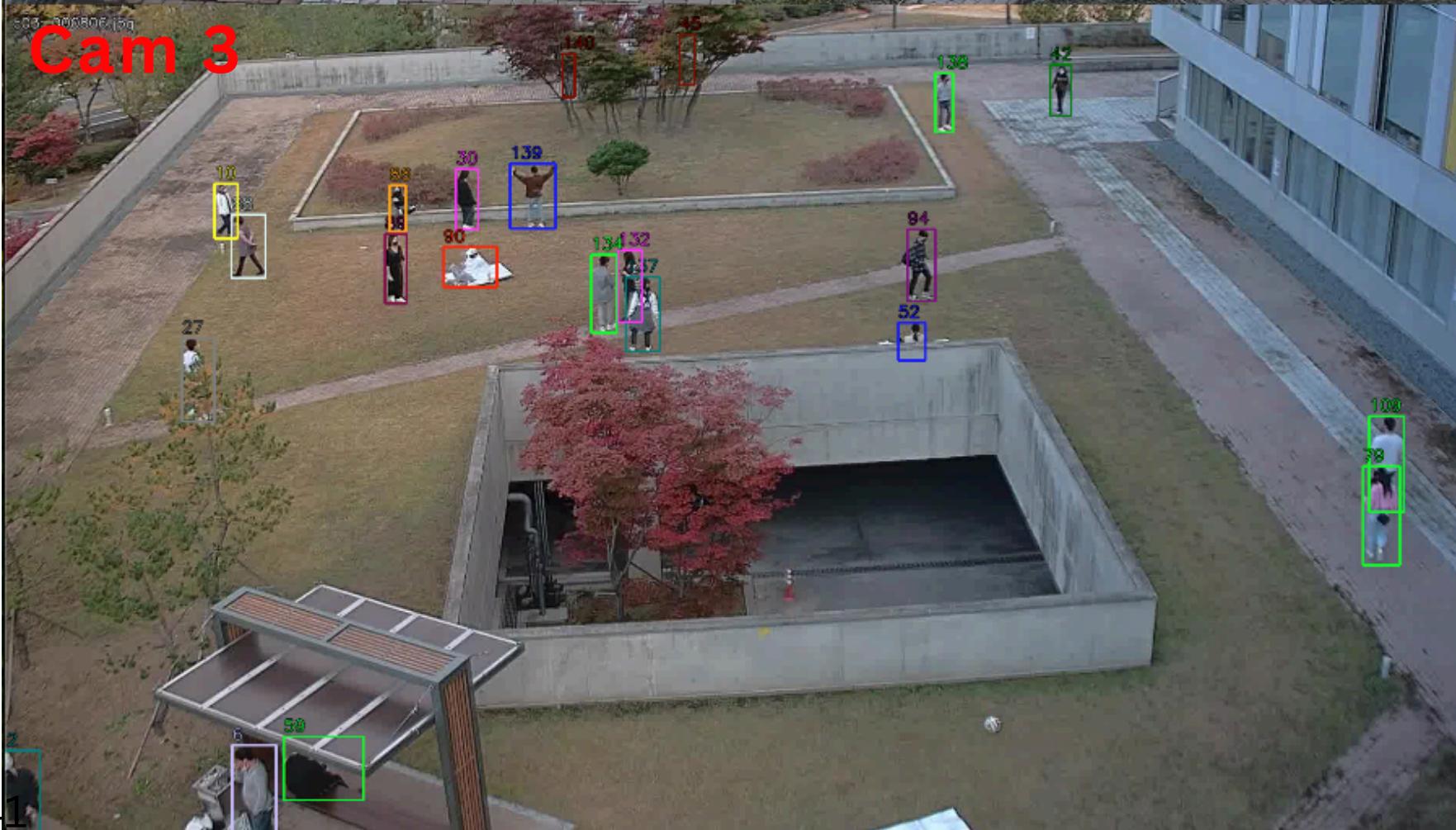
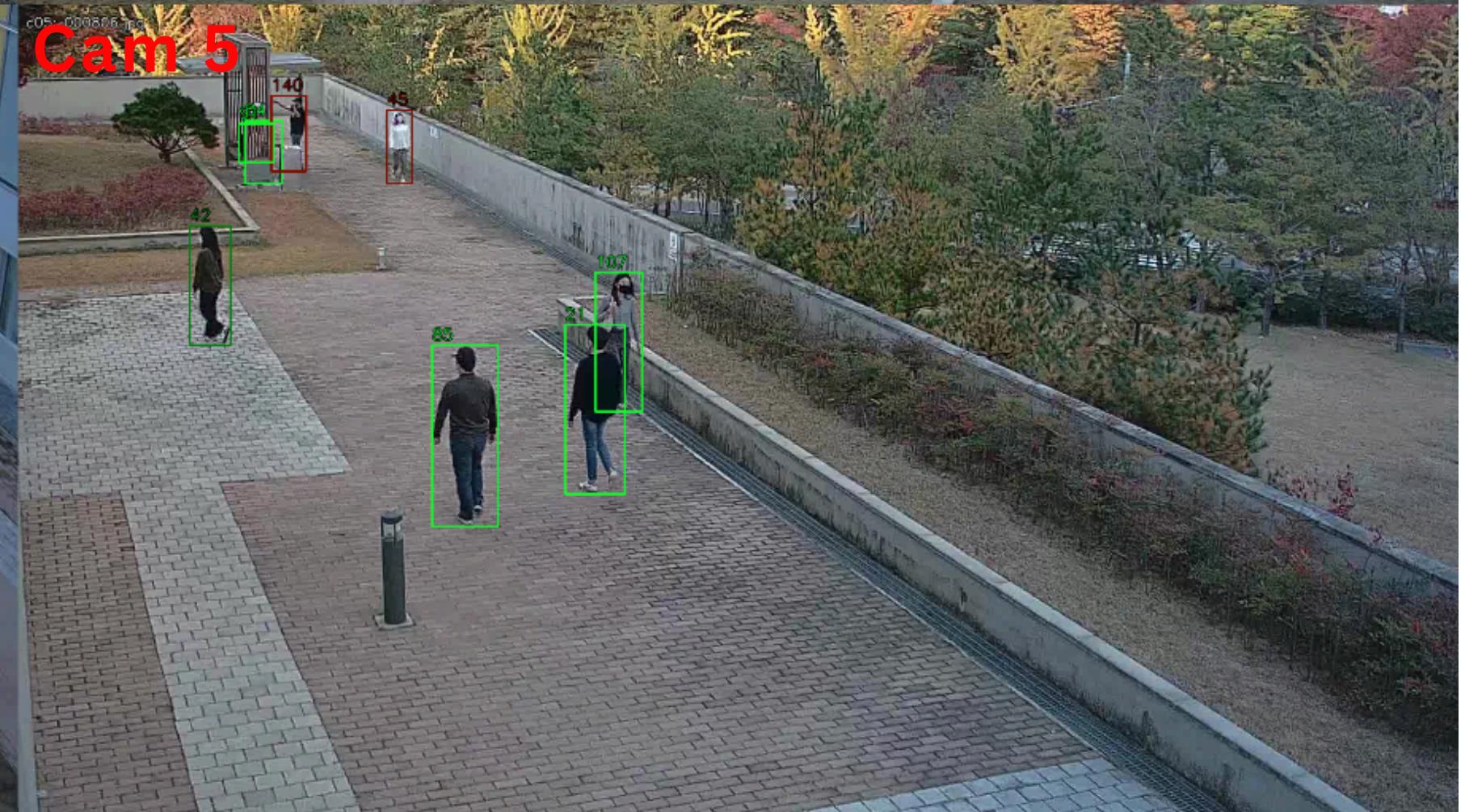
# Camera view selection

# Campus

Selected Cams:

- 1
- 2
- 3
- 5



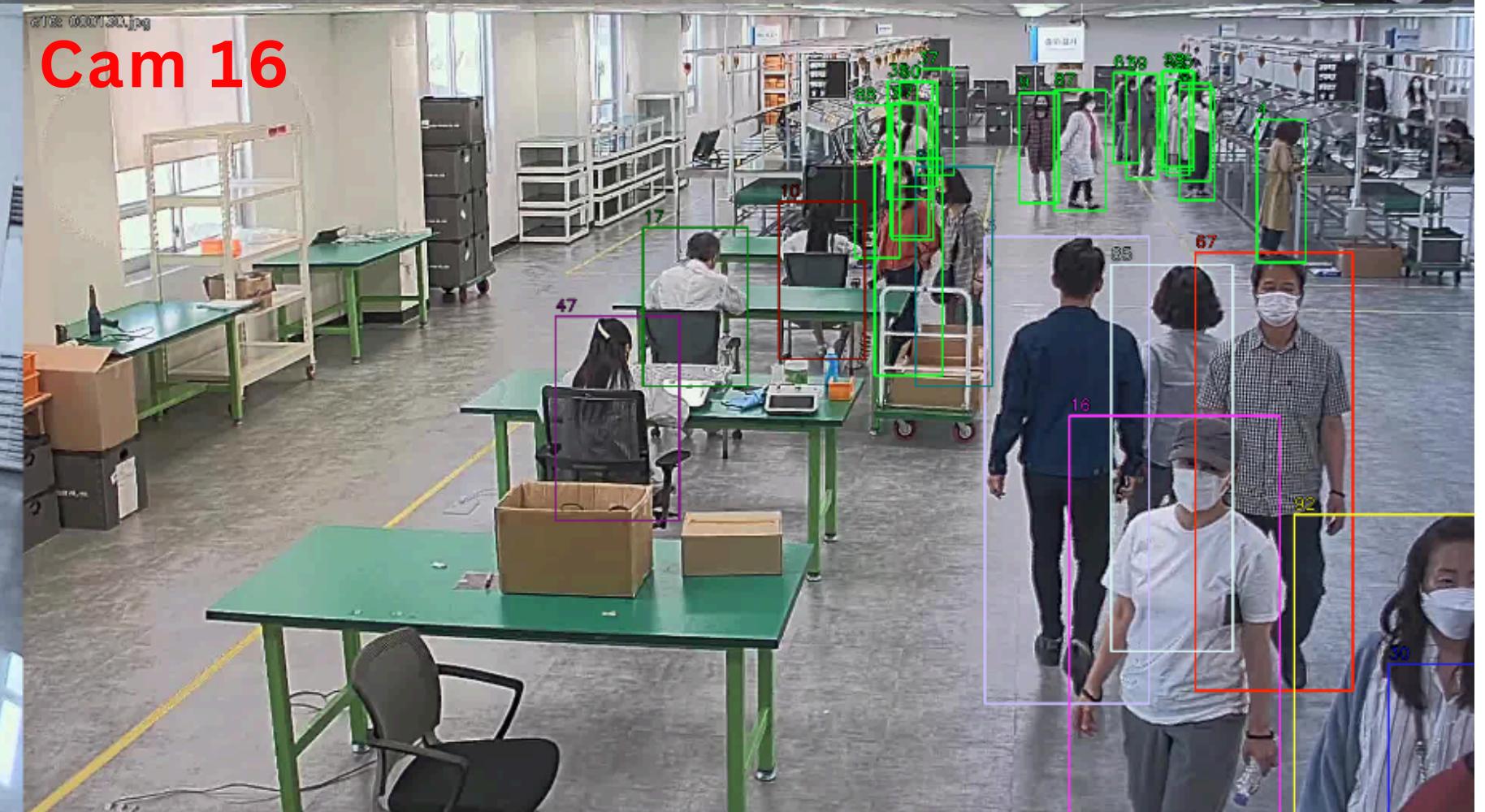


# Factory

Selected Cams:

- 9
- 12
- 13
- 16



**Cam 9****SpotOn****Cam 13****Cam 16**

# What's next

- MLOPs designing