



Carnegie Mellon University  
Language Technologies Institute

# Vision, Robotics, Hardware...

Yonatan Bisk and Emma Strubell

Some slides stolen from 11-777

# Benchmarking on Pi Zero

	<b>Delta</b>	<b>Remark</b>
USB Ethernet dongle	+104 mA	No network connection, only a USB Ethernet dongle inserted into the USB OTG port.
USB Ethernet dongle + link	+180 mA	After inserting an Ethernet cable.
USB Ethernet dongle + iperf	+244 mA	Average. See chart below for details.
2.4 GHz WiFi	+11 mA	Connected to 2.4 GHz access point.
2.4 GHz WiFi + iperf	+187 mA	Average. See chart below for details.
Logitech RF dongle for keyboard/mouse	+29 mA	
Seagate USB HDD (idle)	+255 mA	Stable current after initial insertion that peaked at 1.06A. See chart below for details.
Seagate USB HDD (iozone)	+554 mA	Average. See chart below for details
HDMI cable	+7 mA	This is only after inserting the cable, not about enabling/disabling HDMI. More on that below.

**Processor clocks can also be modified!**

# Operating Systems and Libraries

Ubuntu 18.04 and 20.04

Python 3.6+

PyTorch 1.09+

Torchvision 0.9.0

Huggingface transformers 4.9.2+

These are not required.

You might find several others very useful:

1. TFLite
2. ONNX
3. Detectron2
4. ROS

But they are not officially supported by the course  
so you will have limited TA support with debugging.

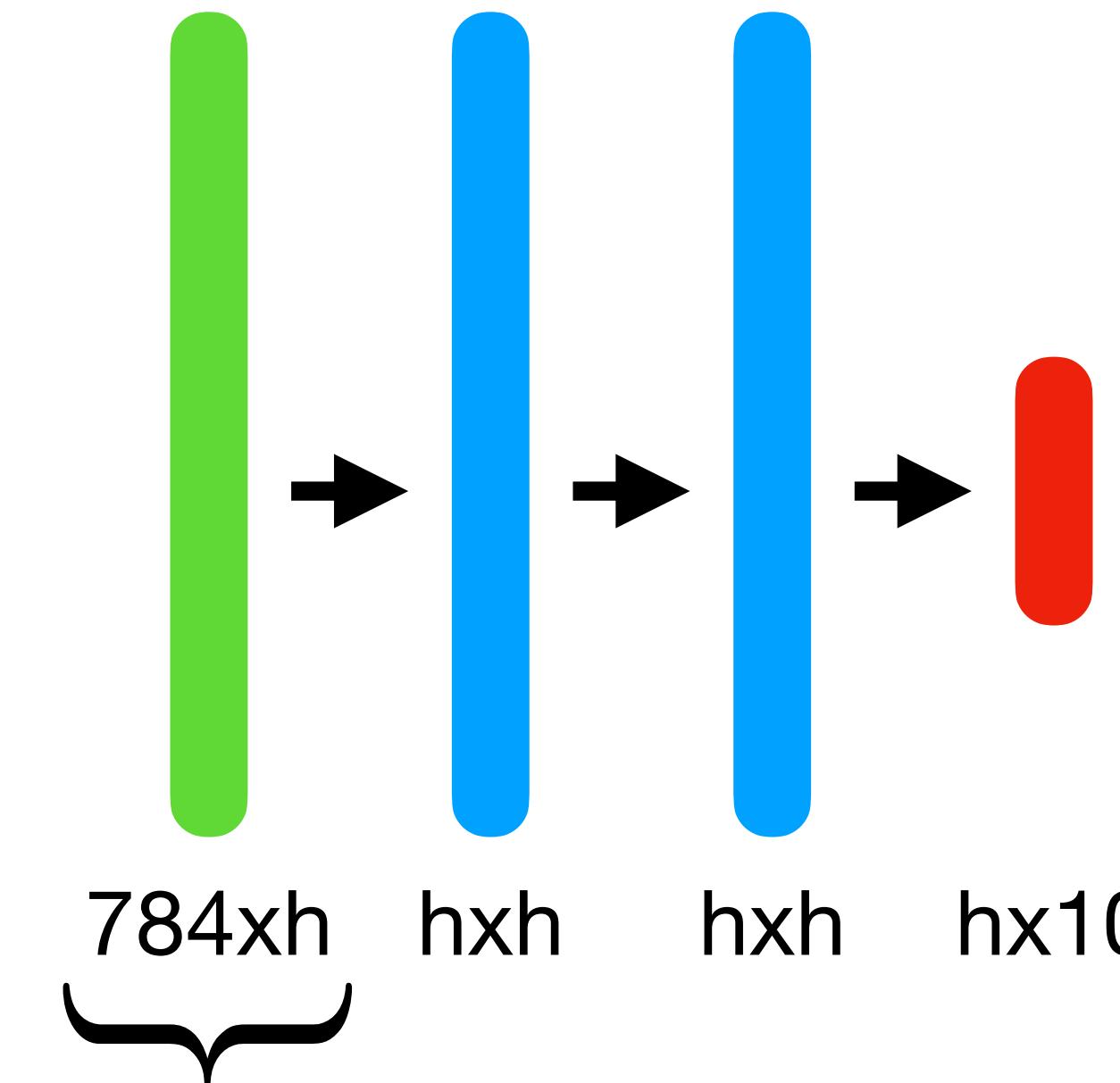
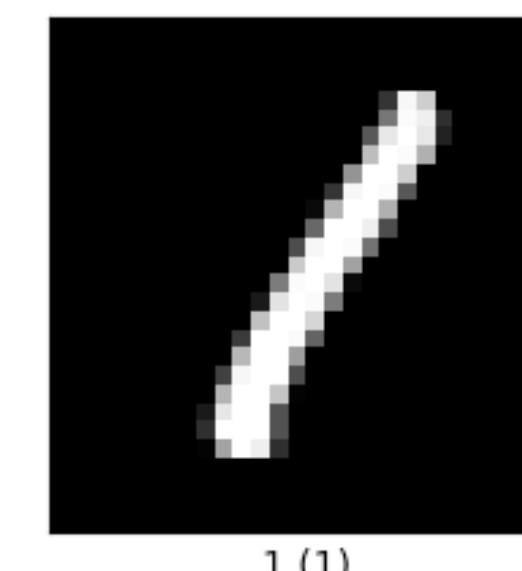
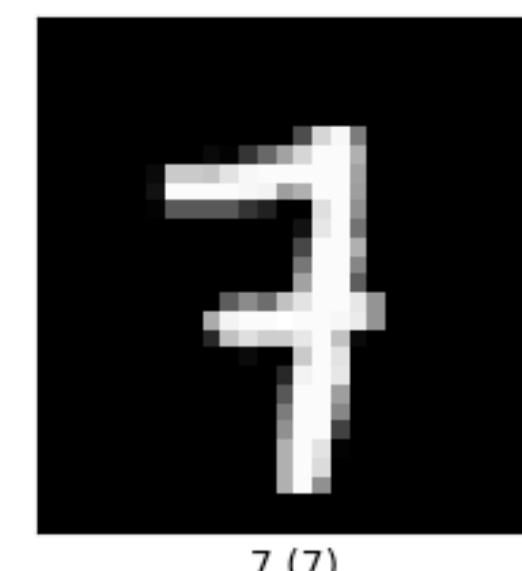
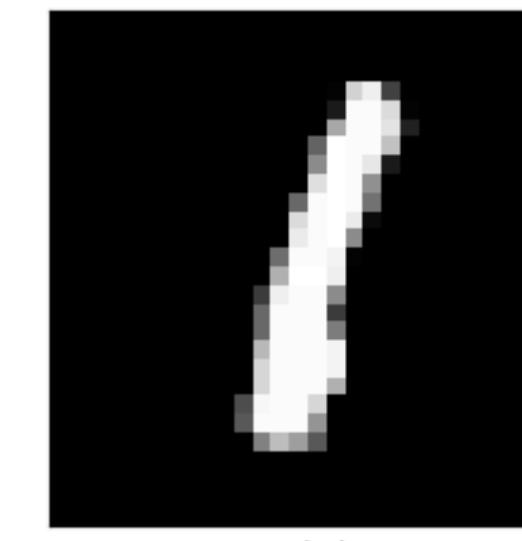
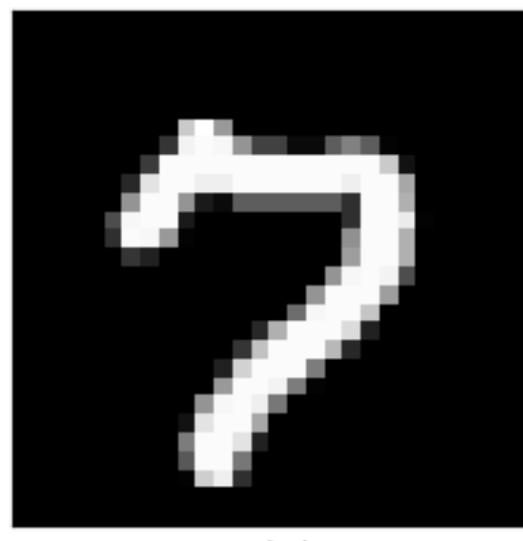
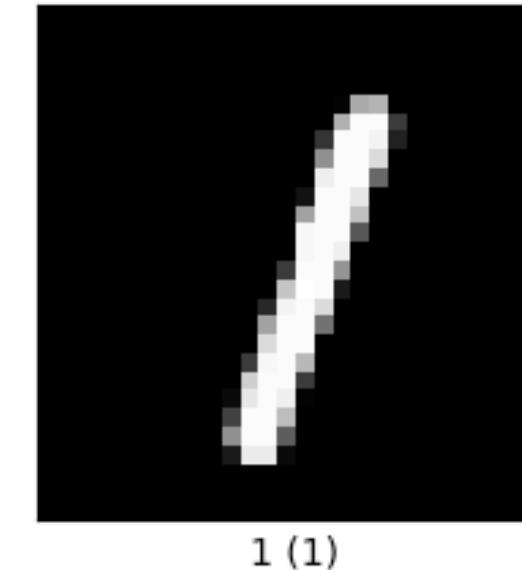
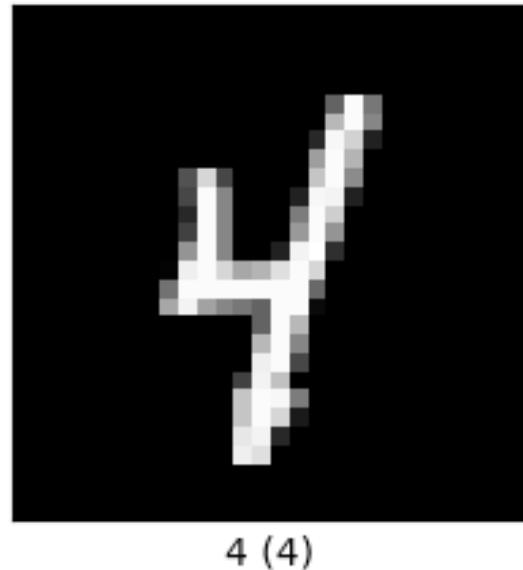
# Keywords

- Feature Representations
- Input/Output dimensionality
- CNNs: Layers (norms, pooling, etc)
  - 1D, 2D, 3D
- Transformers: (Self-Attention Blocks, Auxiliary Loss, Masking, Auto-regressive)
- Mapping, real valued control, ...

# Questions

- How much information does a pixel, a spectrogram, a word contain?
- How much do I care about parameter size vs model size?
- Does my task require a large vocabulary (why?), high resolution (why?), ...
- What interactions are required in the input?  
(Which pixels or words need to talk to each other?)

# Lab 1: Computer Vision

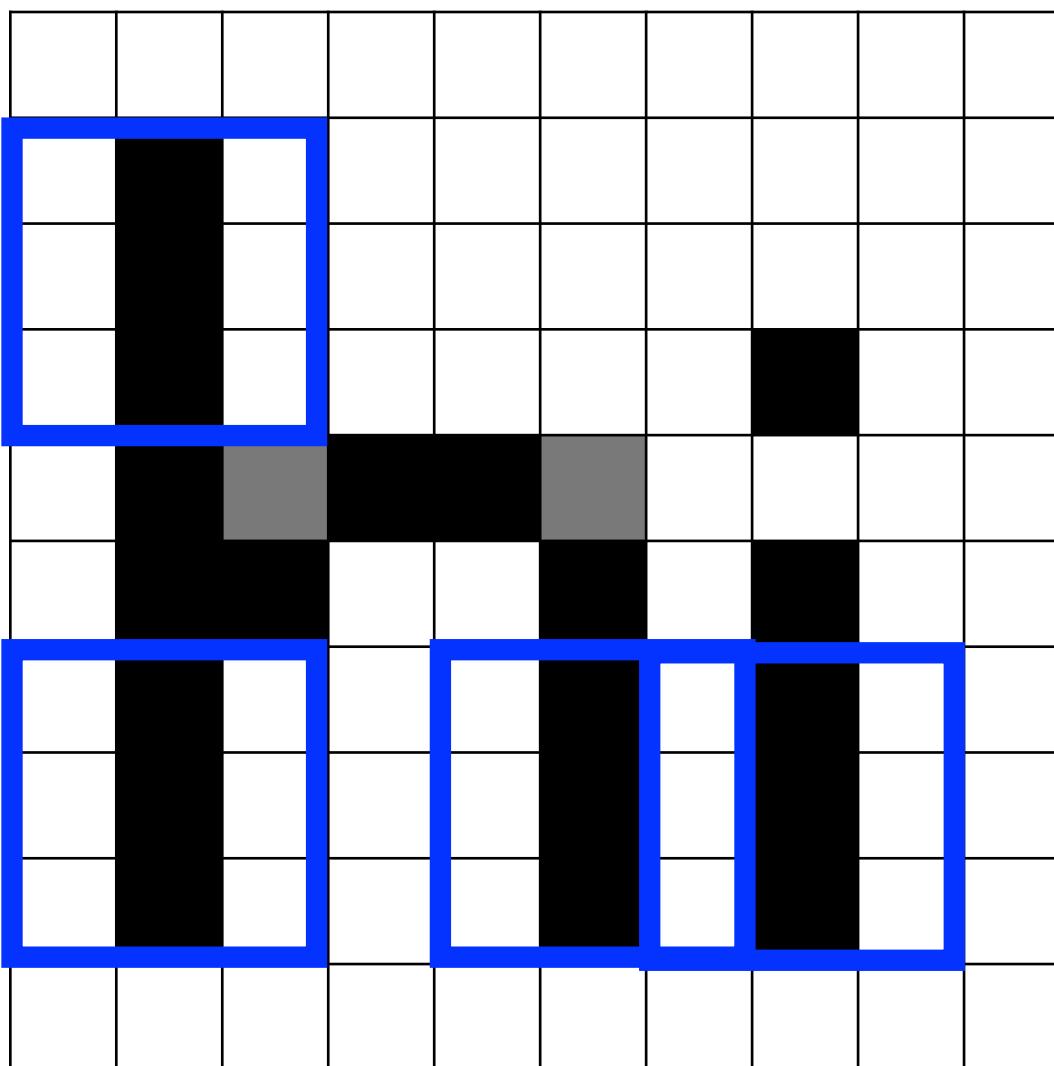


Each hidden dimension is a function of *all* input pixels

Let's make some assumptions:

# Parameter Efficient, Local Processing

Learning a code book  
Small, inexpensive, and reusable parameters



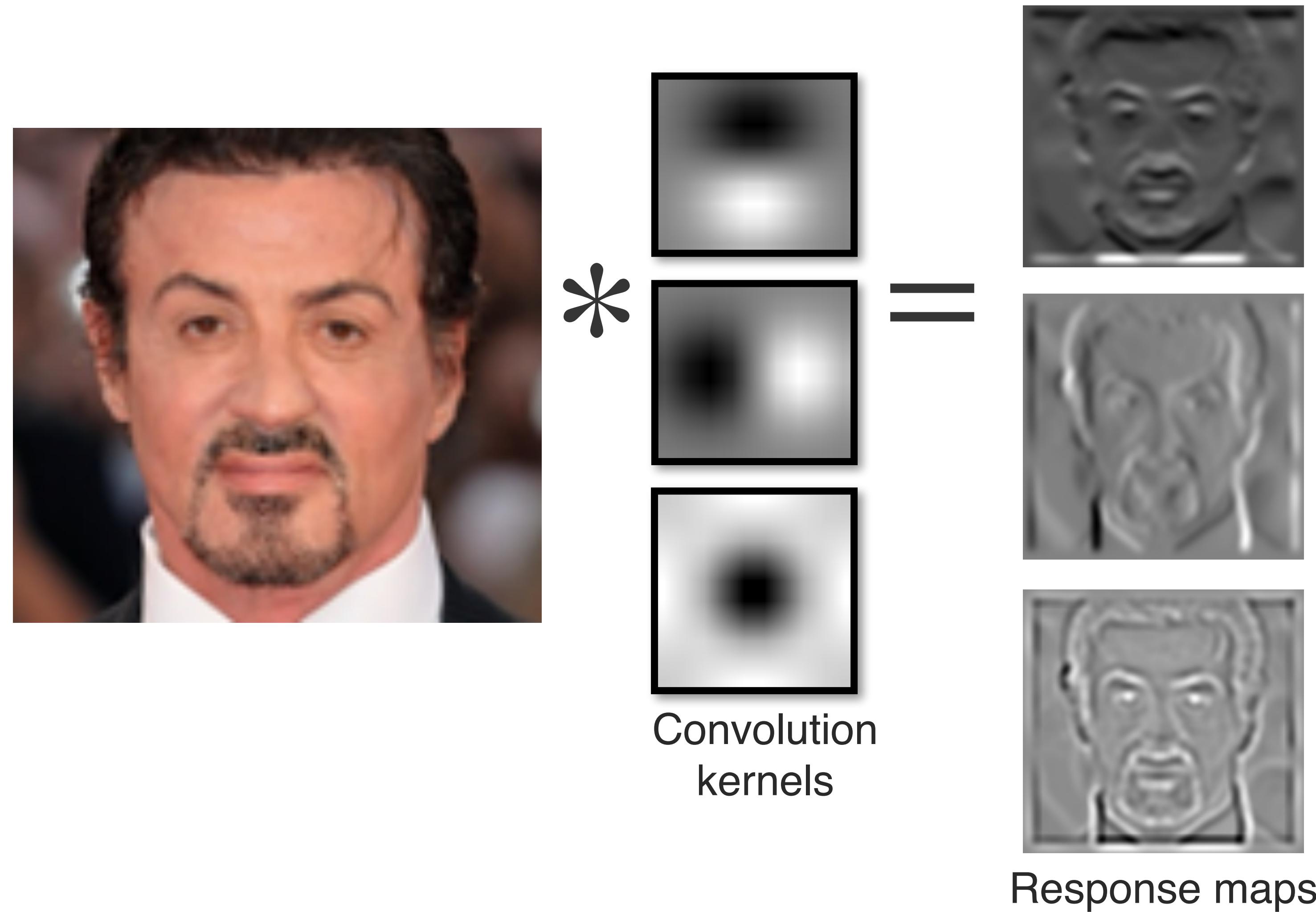
FF 10x10xh Run once

CNN 3x3xh Run 81x? 🤔

Run 16x?

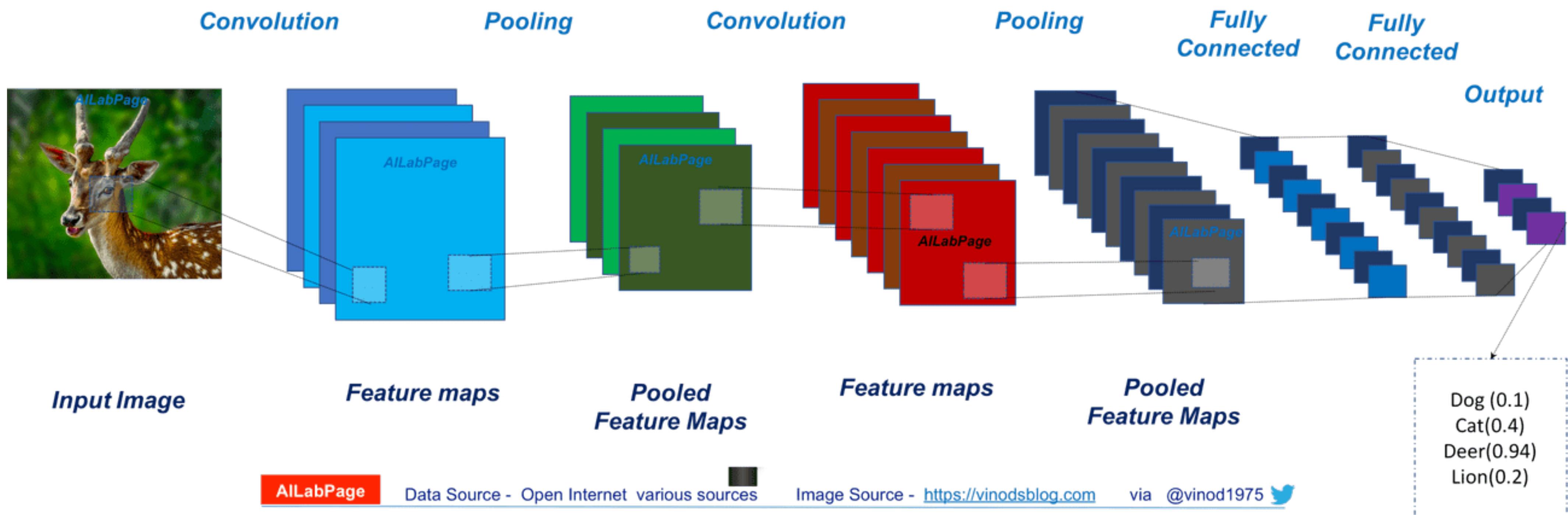
Learning efficiency – “lots of small images”

# Convolution Kernels



# Object Descriptors

## Convolutional Neural Network (CNN)



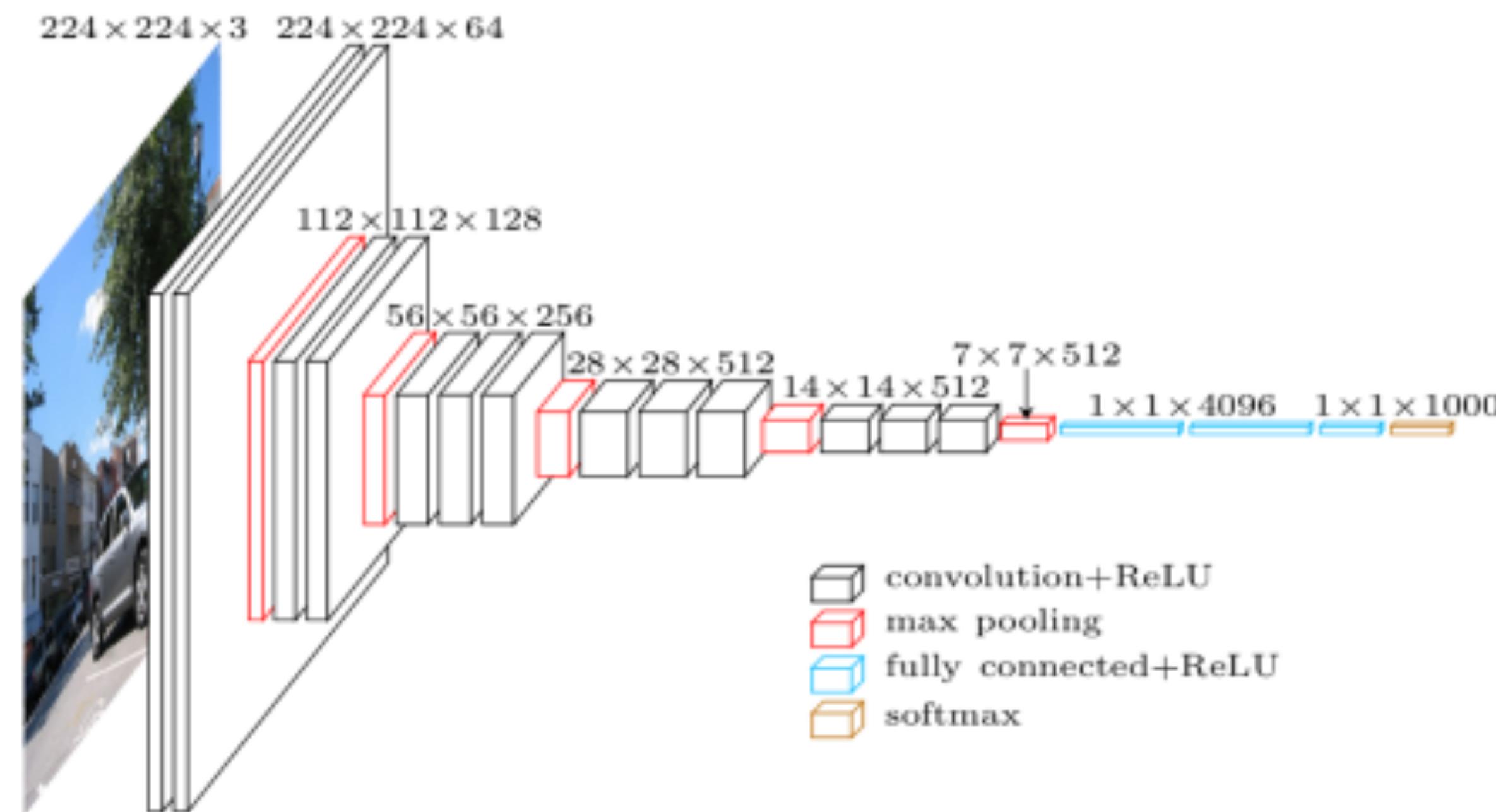
# One representation, lots of tasks



# VGGNet model

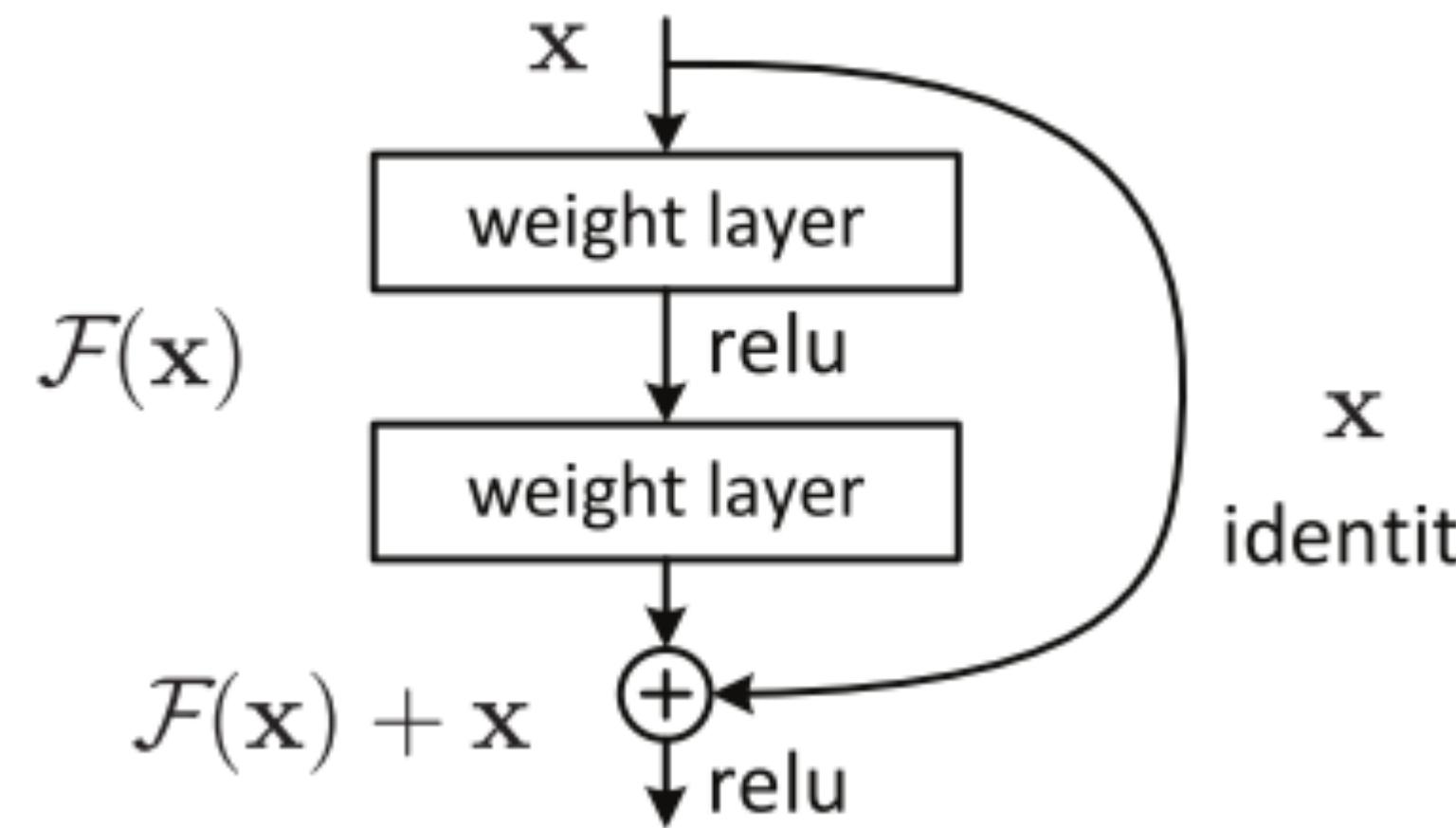
Used for object classification task

- 1000-way classification task
- 138 million parameters



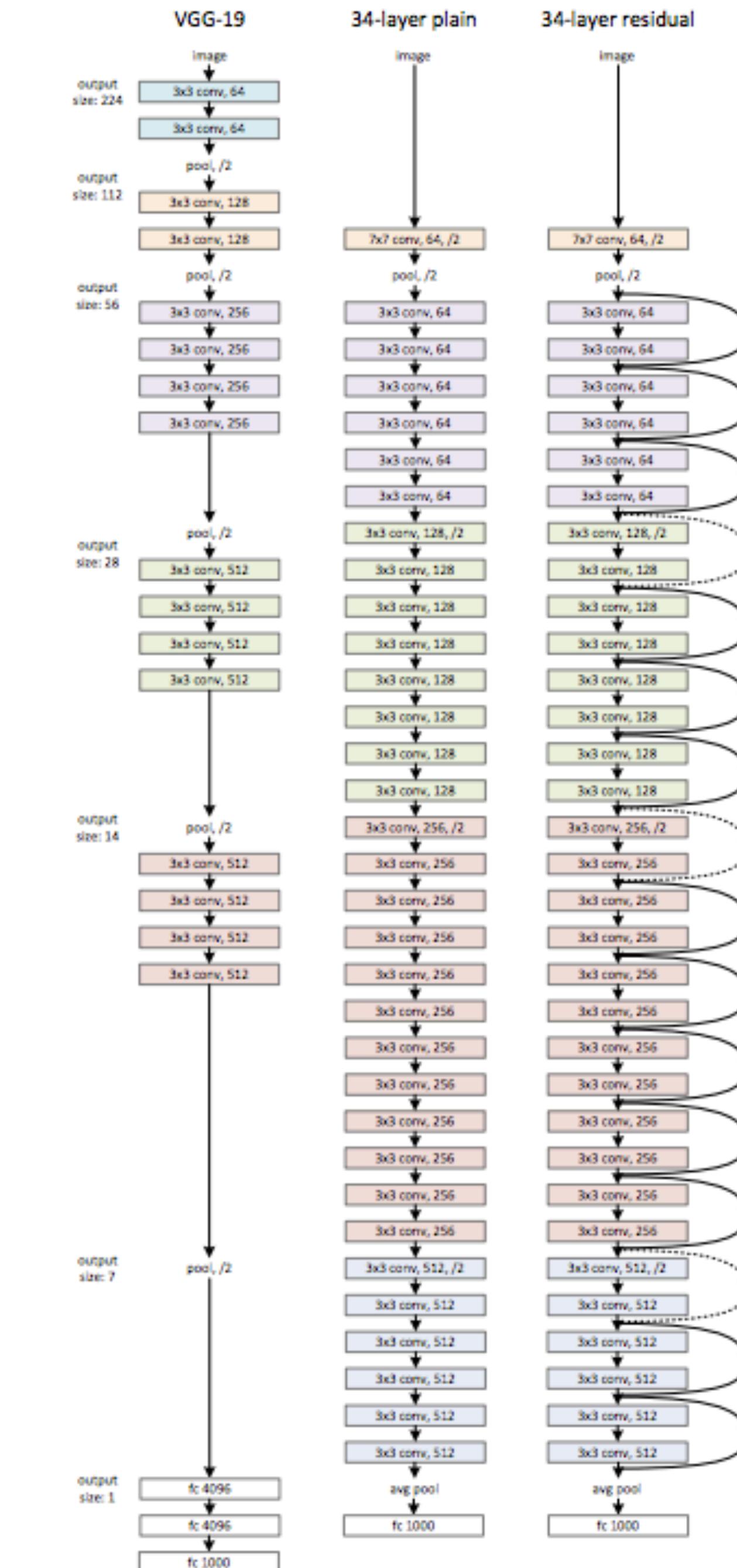
# Residual Networks

## Adding residual connections

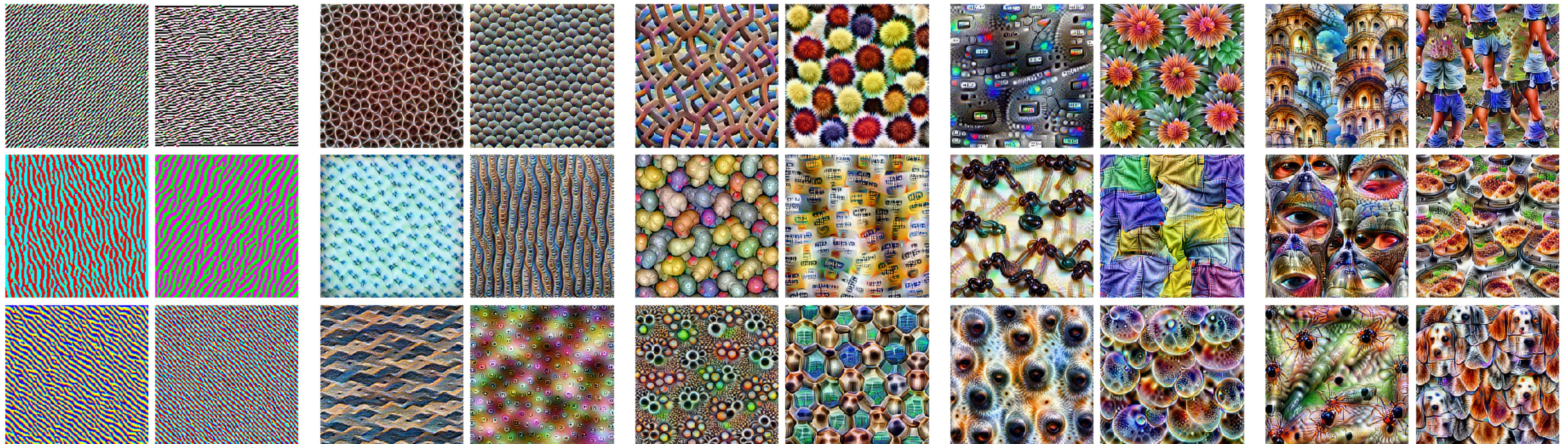


ResNet (He et al., 2015)

- Up to 152 layers!



# Hierarchical Concepts



Edges (layer conv2d0)

Textures (layer mixed3a)

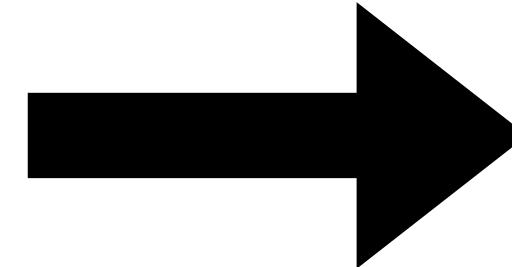
Patterns (layer mixed4a)

Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)



# Discriminative Labels Select Features



“Dog”

What will the model learn is the meaning of dog?

**Dataset Examples** show us what neurons respond to in practice



**Optimization** isolates the causes of behavior from mere correlations. A neuron may not be detecting what you initially thought.



Baseball—or stripes?  
*mixed4a, Unit 6*

Animal faces—or snouts?  
*mixed4a, Unit 240*

Clouds—or fluffiness?  
*mixed4a, Unit 453*

Buildings—or sky?  
*mixed4a, Unit 492*



# Object recognition



Input Image



# Object Detection (and Segmentation)



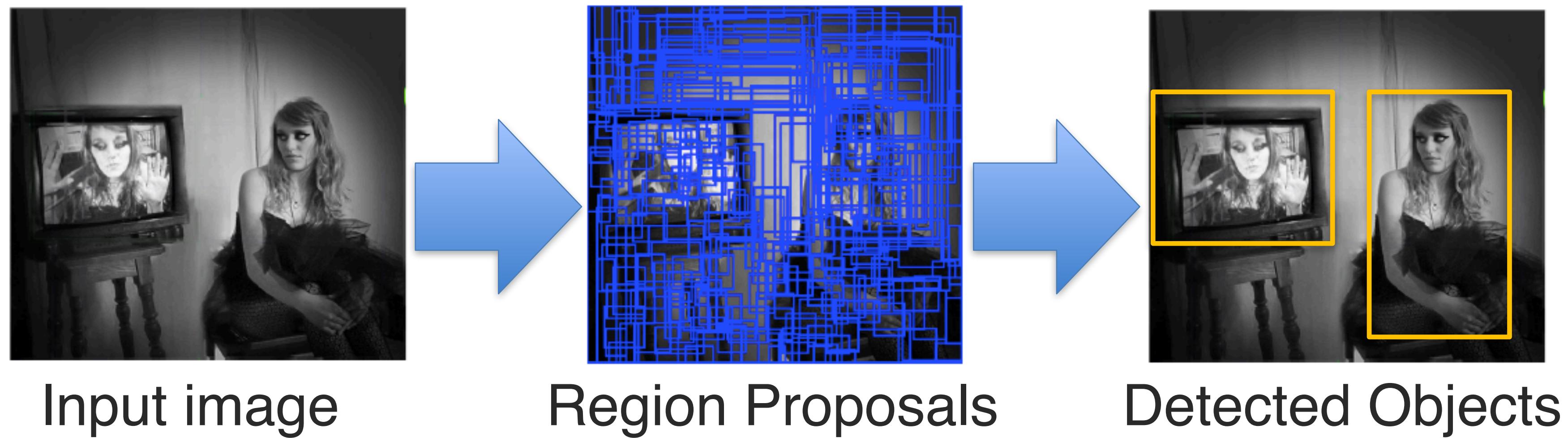
Input image



Detected Objects

One option: Sliding window

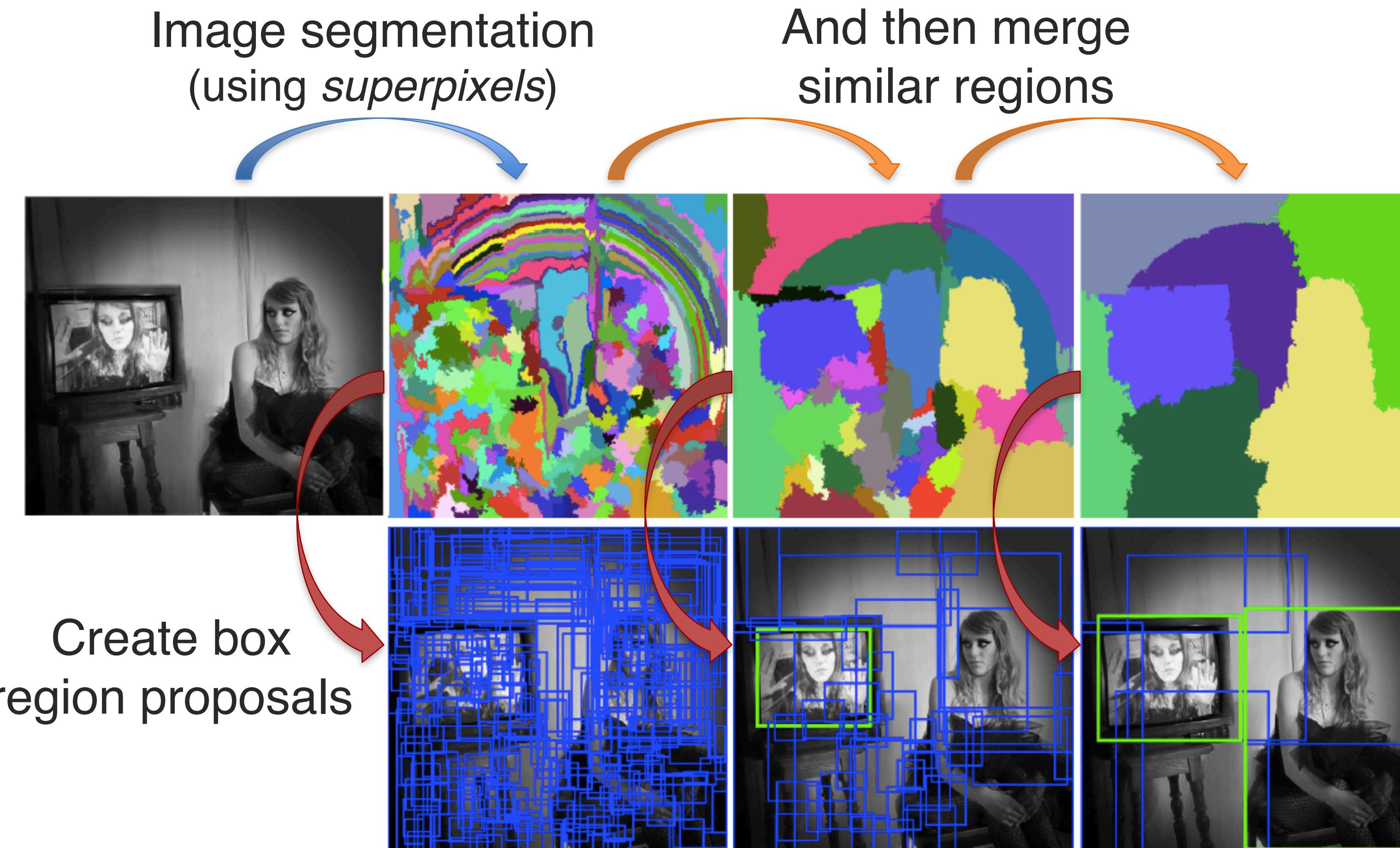
# Object Detection (and Segmentation)



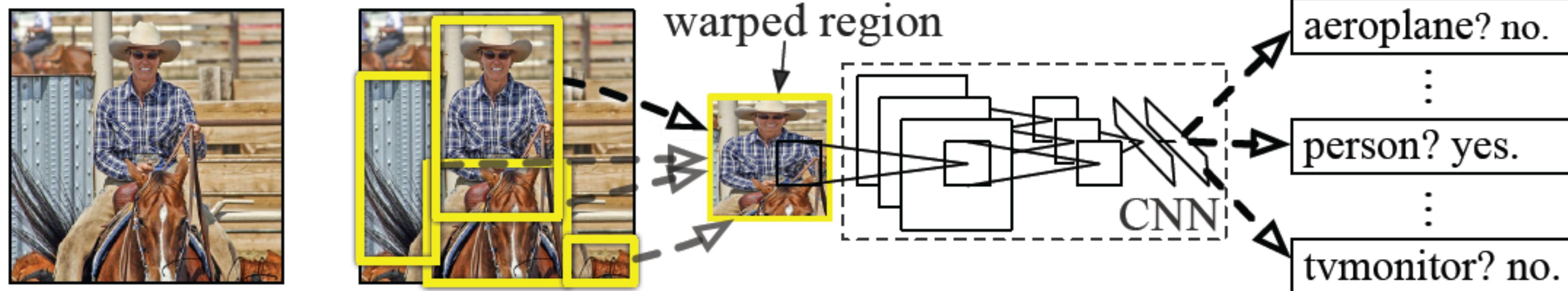
A better option: Start by Identifying hundreds of region proposals and then apply our CNN object detector

*How to efficiently identify region proposals?*

# Selective Search [Uijlings et al., IJCV 2013]



# R-CNN [Girshick et al., CVPR 2014]

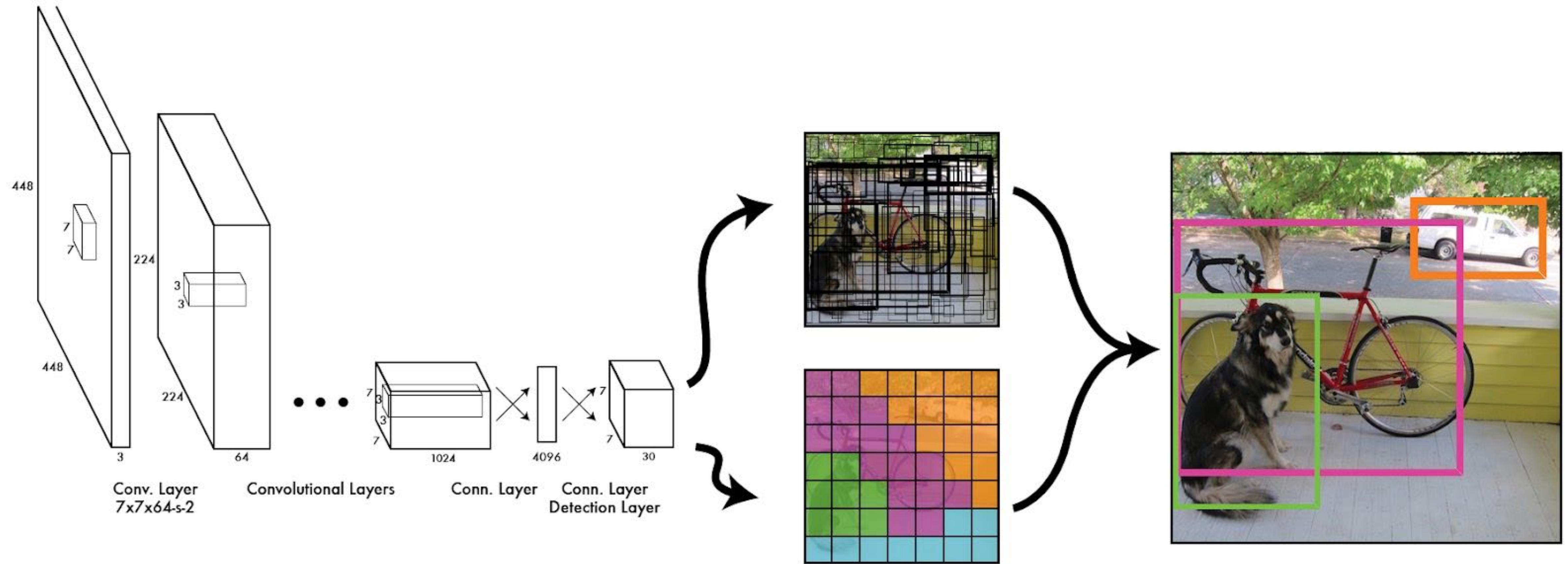


- Select ~2000 region proposals → Time consuming!
- Warp each region
- Apply CNN to each region → Time consuming!

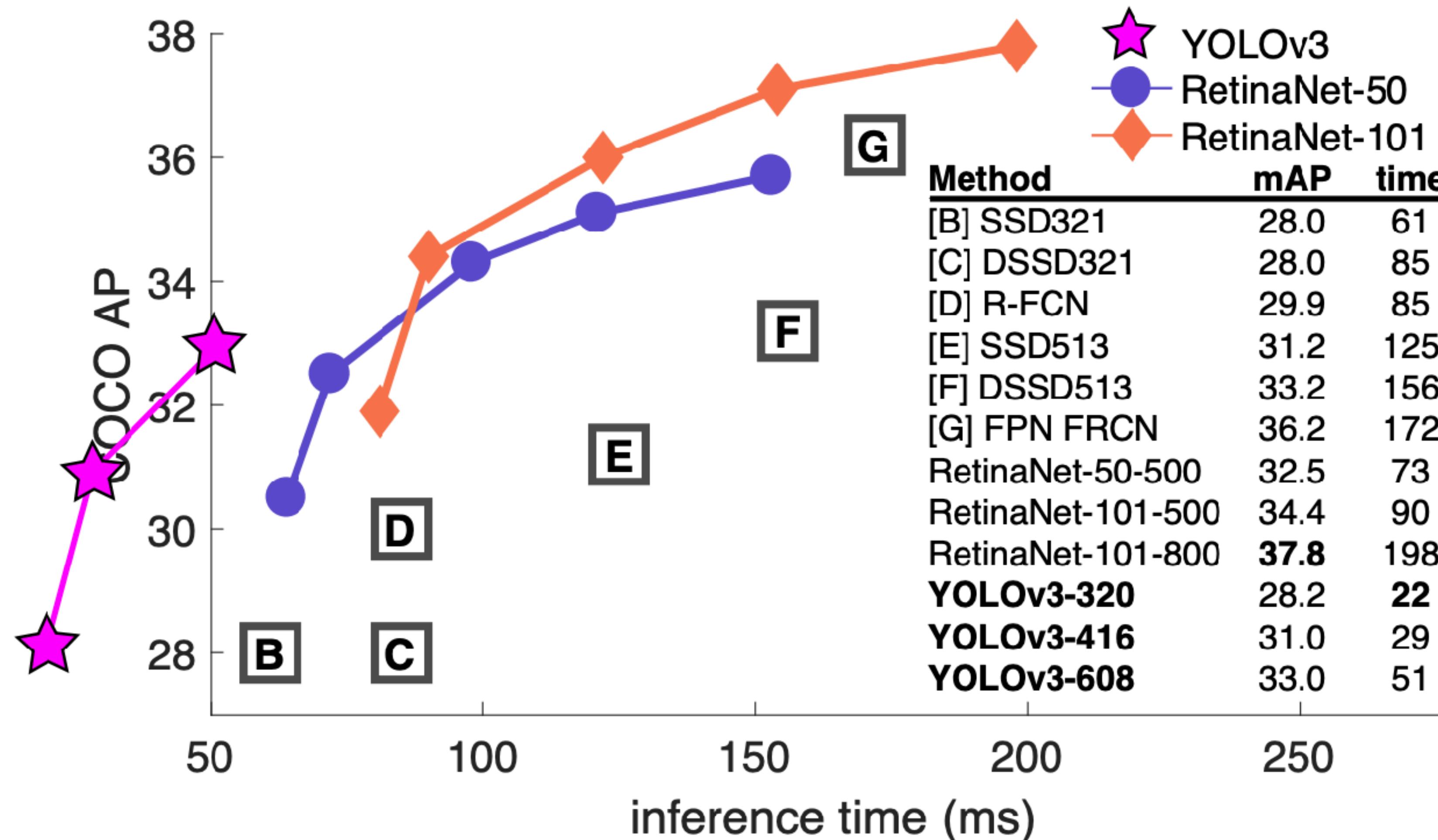
**Fast R-CNN:** Applies CNN only once, and then extracts regions

**Faster R-CNN:** Region selection on the Conv5 response map

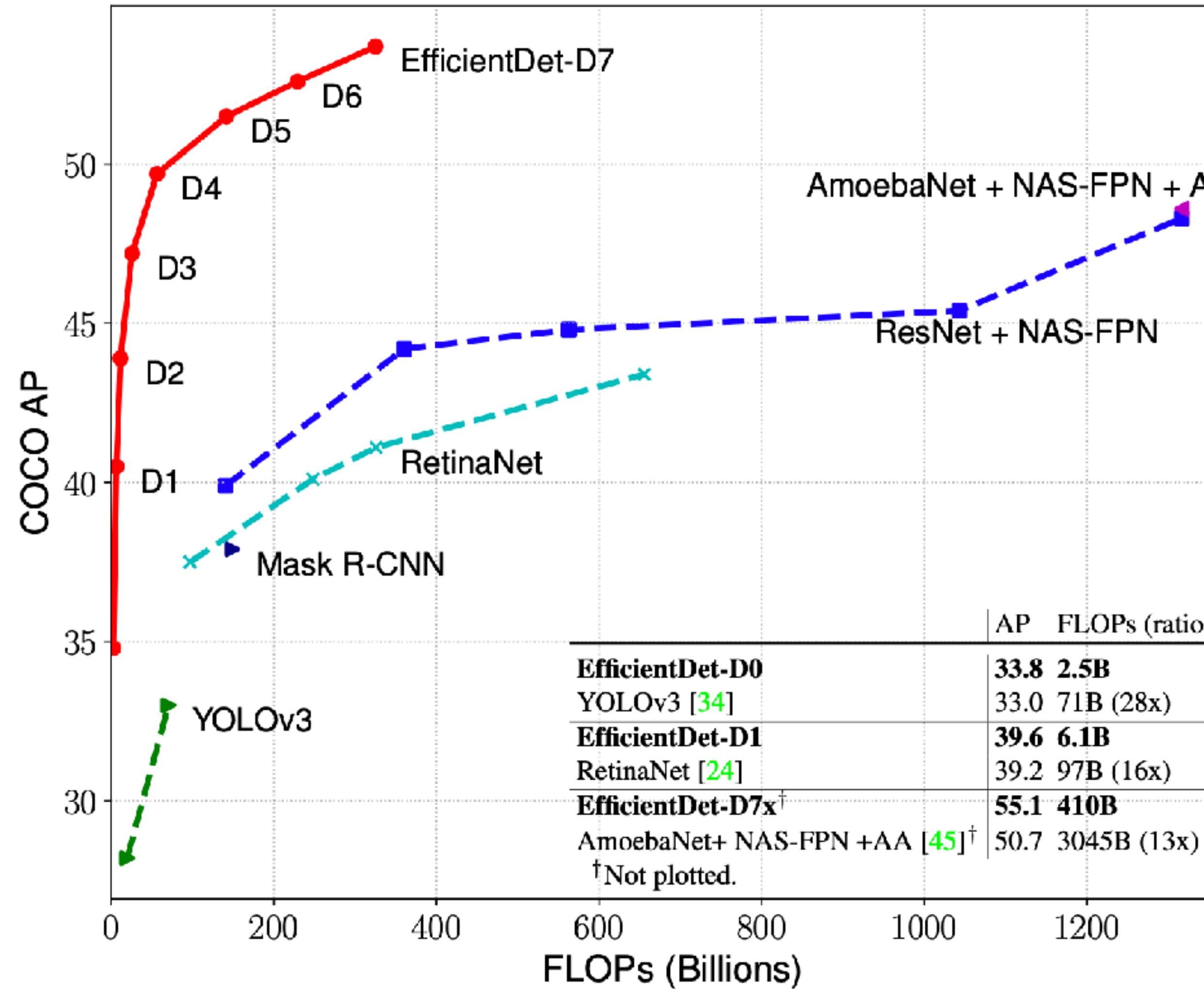
# YOLO [cvpr 2016]



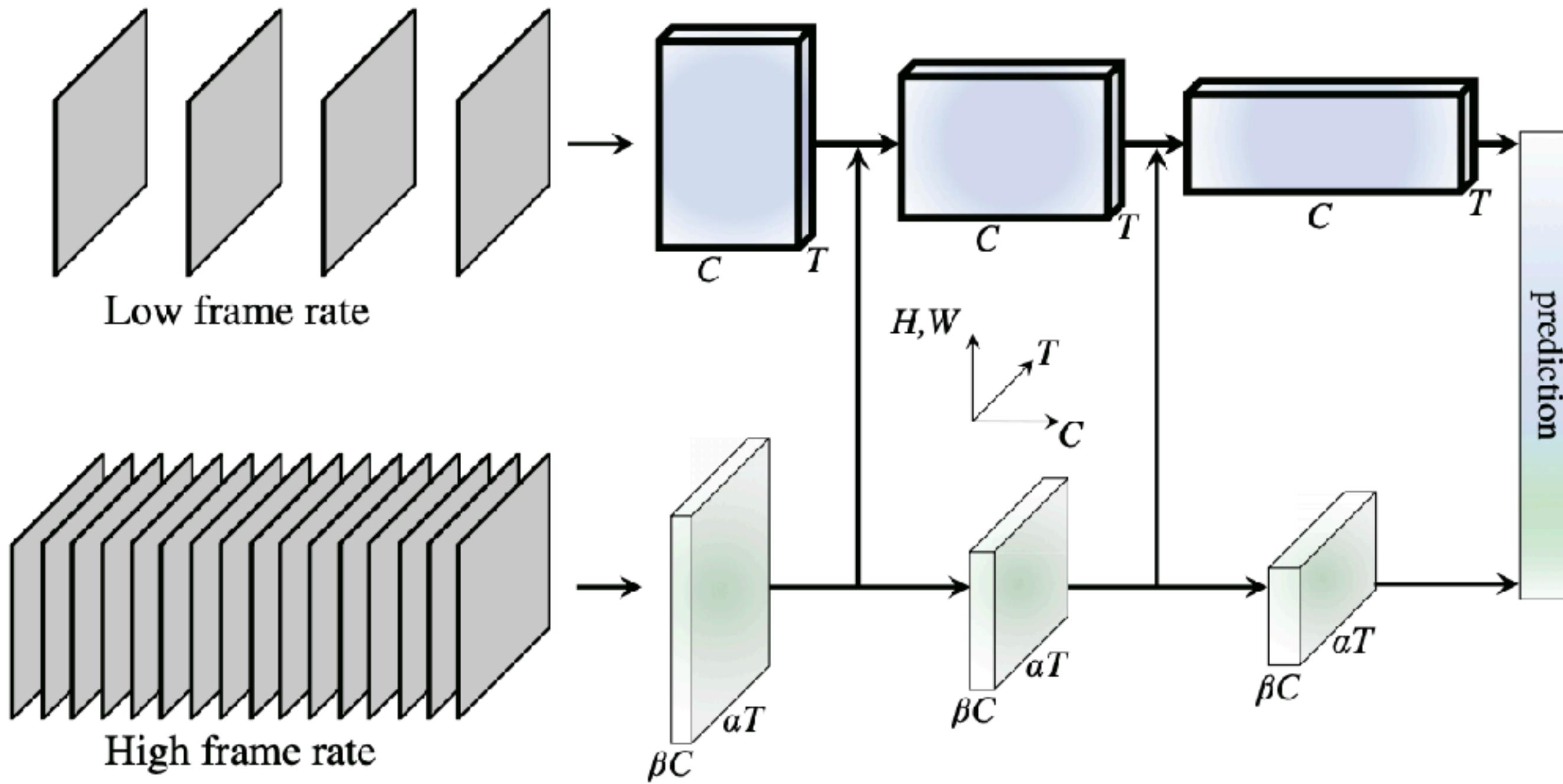
# Trade-off Between Speed and Accuracy



# Trade-off Between Speed and Accuracy



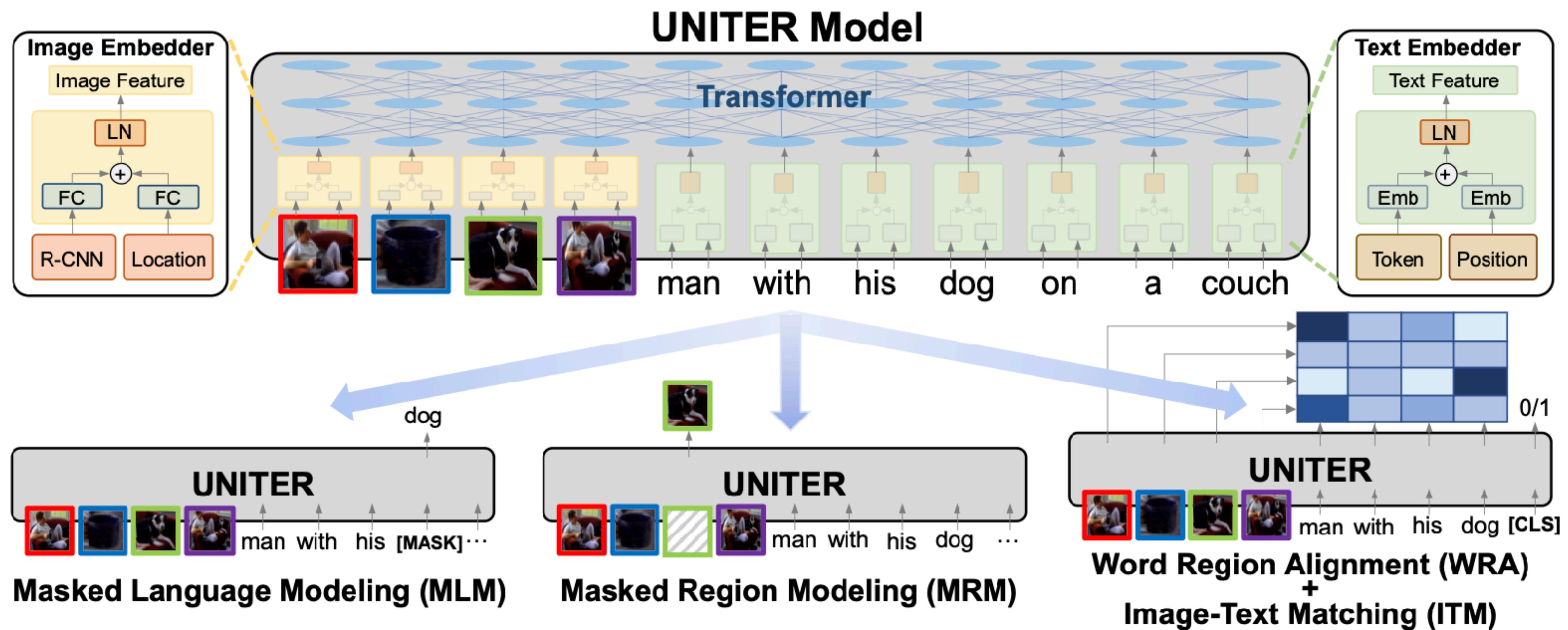
# SlowFast Networks for Video Recognition



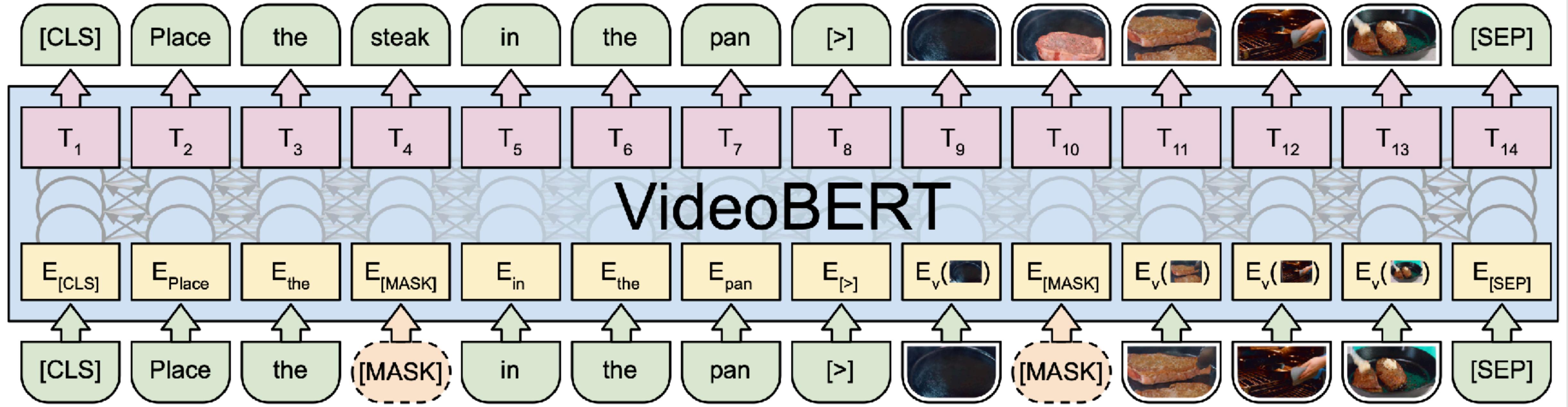
stage	Slow pathway	Fast pathway	output sizes $T \times S^2$
raw clip	-	-	$64 \times 224^2$
data layer	stride 16, 1 <sup>2</sup>	stride 2, 1 <sup>2</sup>	Slow : $4 \times 224^2$ Fast : $32 \times 224^2$
conv <sub>1</sub>	1×7 <sup>2</sup> , 64 stride 1, 2 <sup>2</sup>	$\frac{5 \times 7^2}{2}, 8$ stride 1, 2 <sup>2</sup>	Slow : $4 \times 112^2$ Fast : $32 \times 112^2$
pool <sub>1</sub>	1×3 <sup>2</sup> max stride 1, 2 <sup>2</sup>	1×3 <sup>2</sup> max stride 1, 2 <sup>2</sup>	Slow : $4 \times 56^2$ Fast : $32 \times 56^2$
res <sub>2</sub>	$\begin{bmatrix} 1 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 1^2, 8 \\ 1 \times 3^2, 8 \\ 1 \times 1^2, 32 \end{bmatrix} \times 3$	Slow : $4 \times 56^2$ Fast : $32 \times 56^2$
res <sub>3</sub>	$\begin{bmatrix} 1 \times 1^2, 128 \\ 1 \times 3^2, 128 \\ 1 \times 1^2, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 3 \times 1^2, 16 \\ 1 \times 3^2, 16 \\ 1 \times 1^2, 64 \end{bmatrix} \times 4$	Slow : $4 \times 28^2$ Fast : $32 \times 28^2$
res <sub>4</sub>	$\begin{bmatrix} 3 \times 1^2, 256 \\ 1 \times 3^2, 256 \\ 1 \times 1^2, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 3 \times 1^2, 32 \\ 1 \times 3^2, 32 \\ 1 \times 1^2, 128 \end{bmatrix} \times 6$	Slow : $4 \times 14^2$ Fast : $32 \times 14^2$
res <sub>5</sub>	$\begin{bmatrix} 3 \times 1^2, 512 \\ 1 \times 3^2, 512 \\ 1 \times 1^2, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{bmatrix} \times 3$	Slow : $4 \times 7^2$ Fast : $32 \times 7^2$
	global average pool, concat, fc		
	# classes		



# UNITER



# VideoBERT



How do we get visual words now?!?



# Clusters define visual “words” — 12<sup>4</sup> clusters

Original



ASR:

*“So it's really up to you cut it off get this and nice slices.”*

Centroids



Top verbs: cut, prepare, make

Top nouns: orange, lemon, tomato

Input text

*“Put the pizza into oven.”*

Retrieved centroid



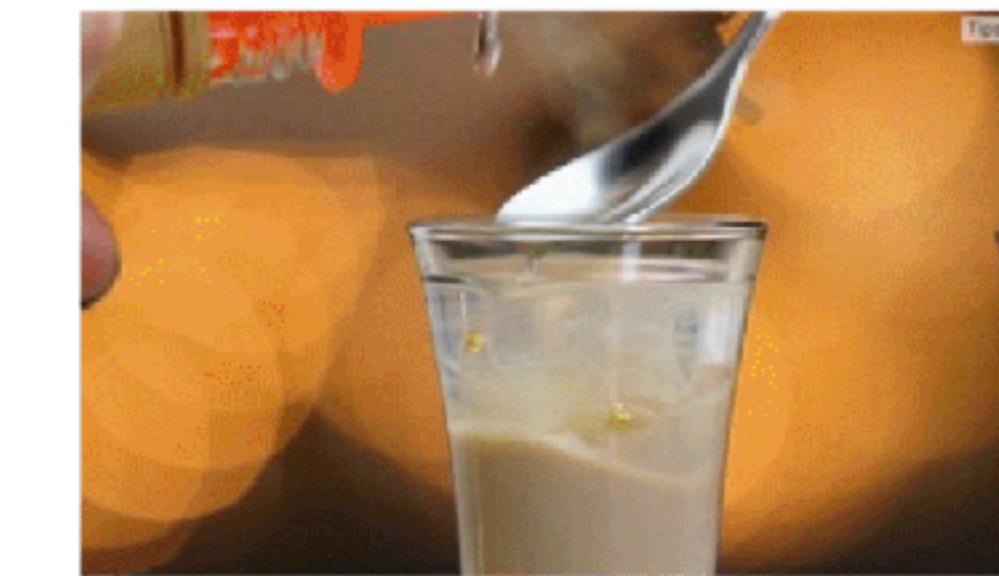
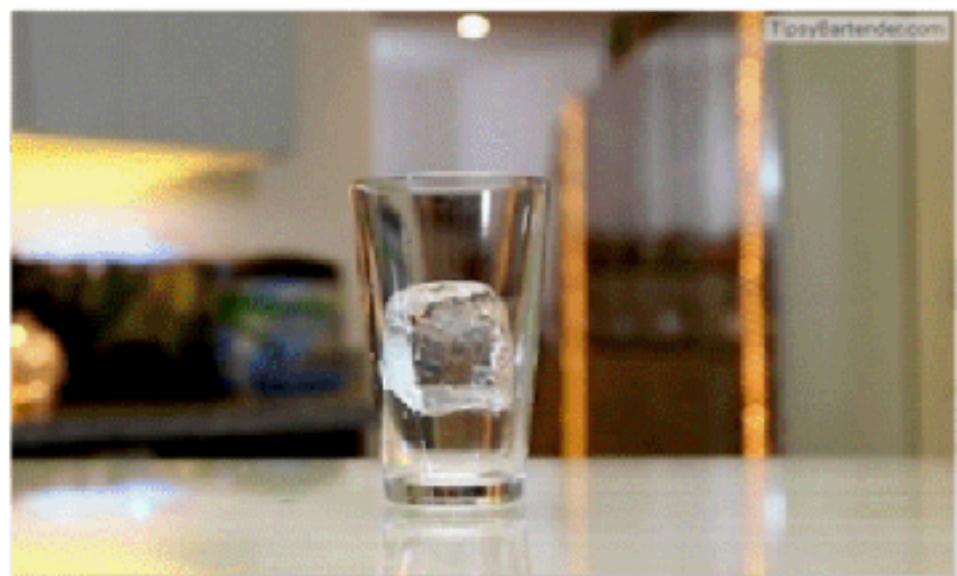
Retrieved centroid



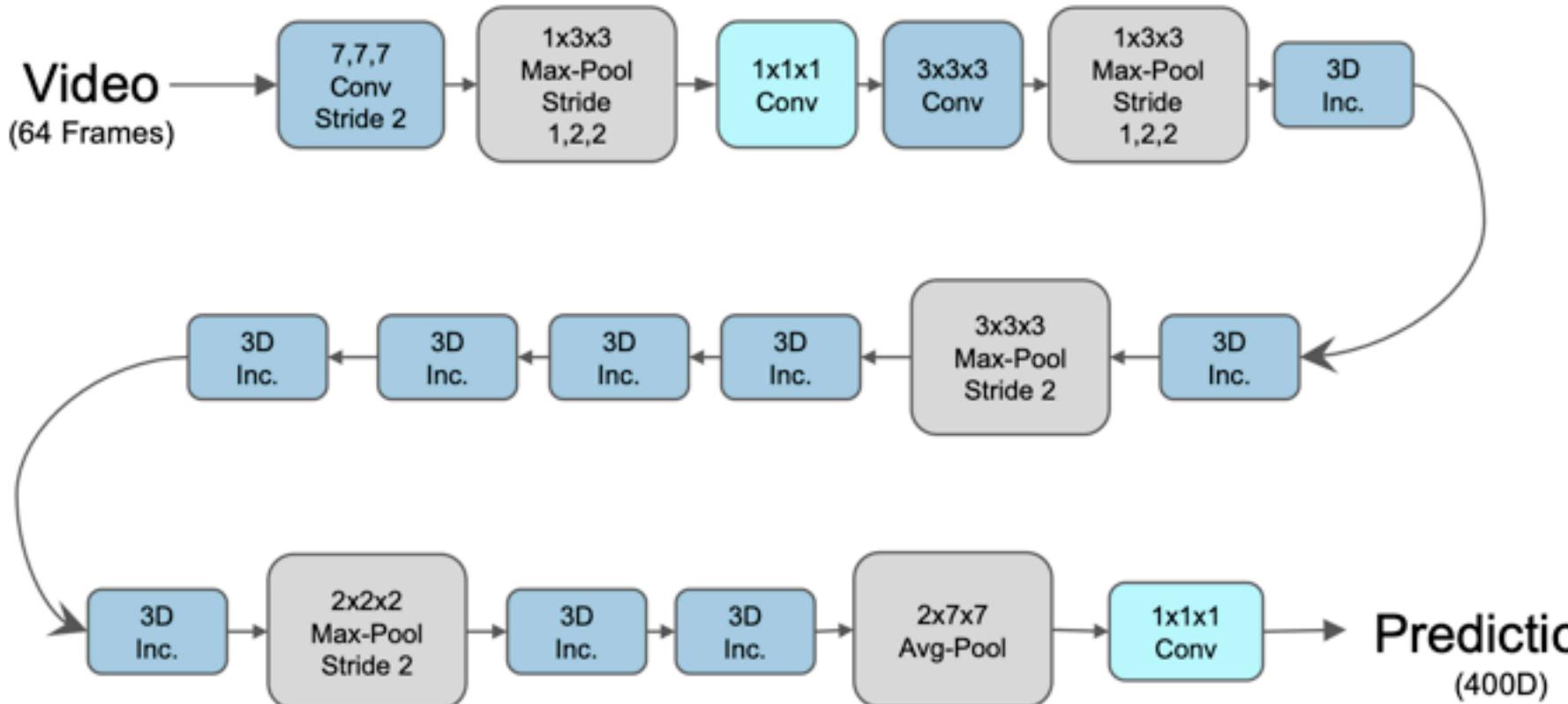
*“Put the cookies into oven.”*



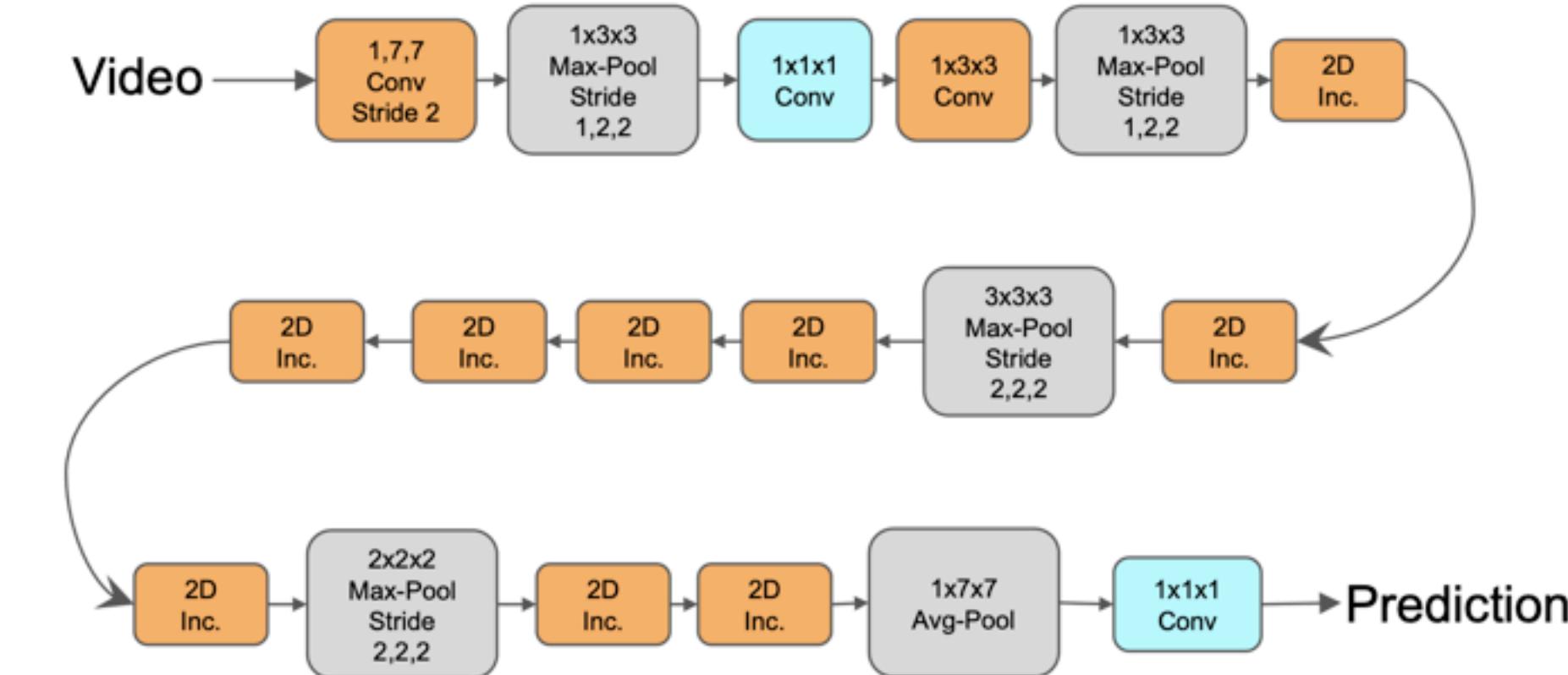
# Video-to-Video



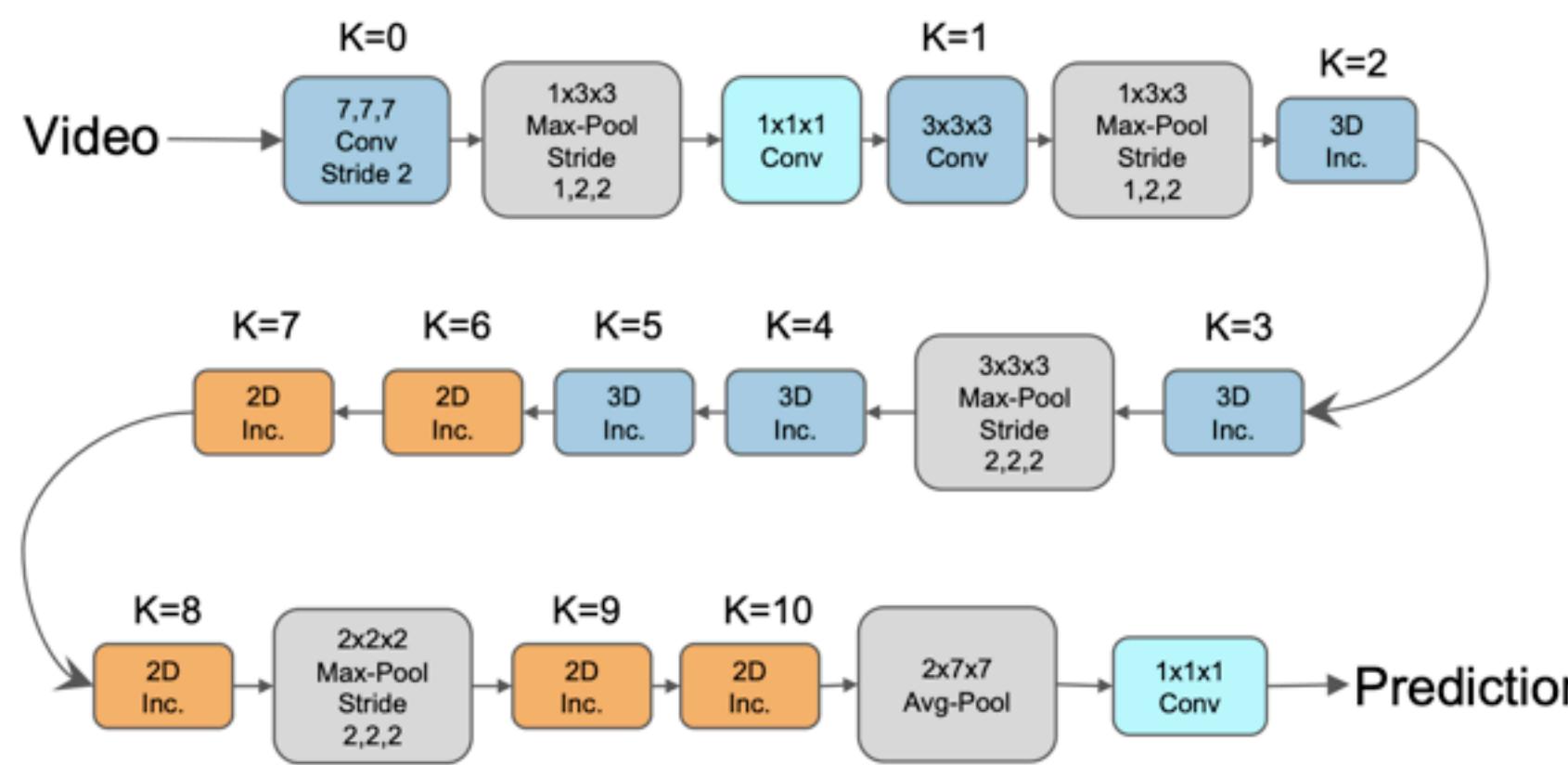
# Video Features



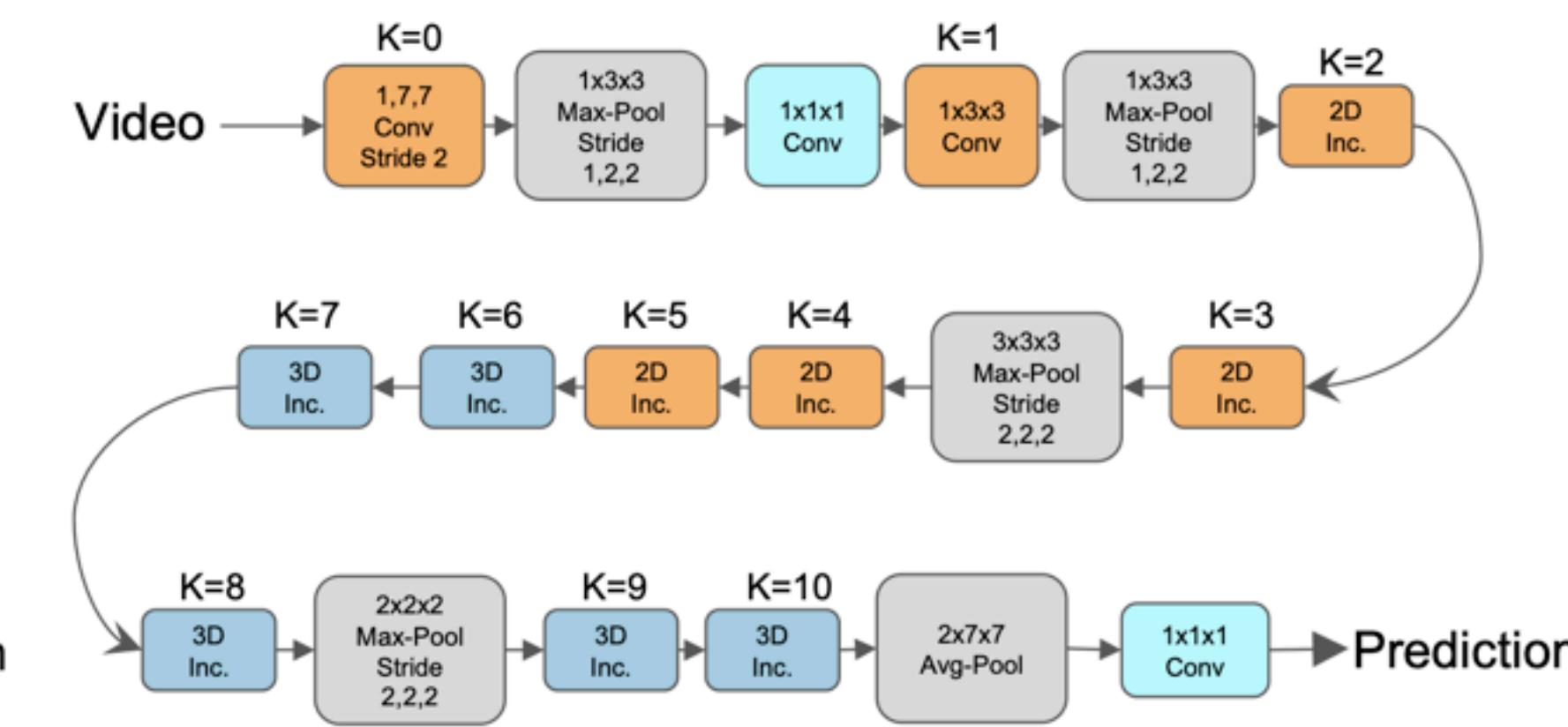
(a) I3D



(b) I2D

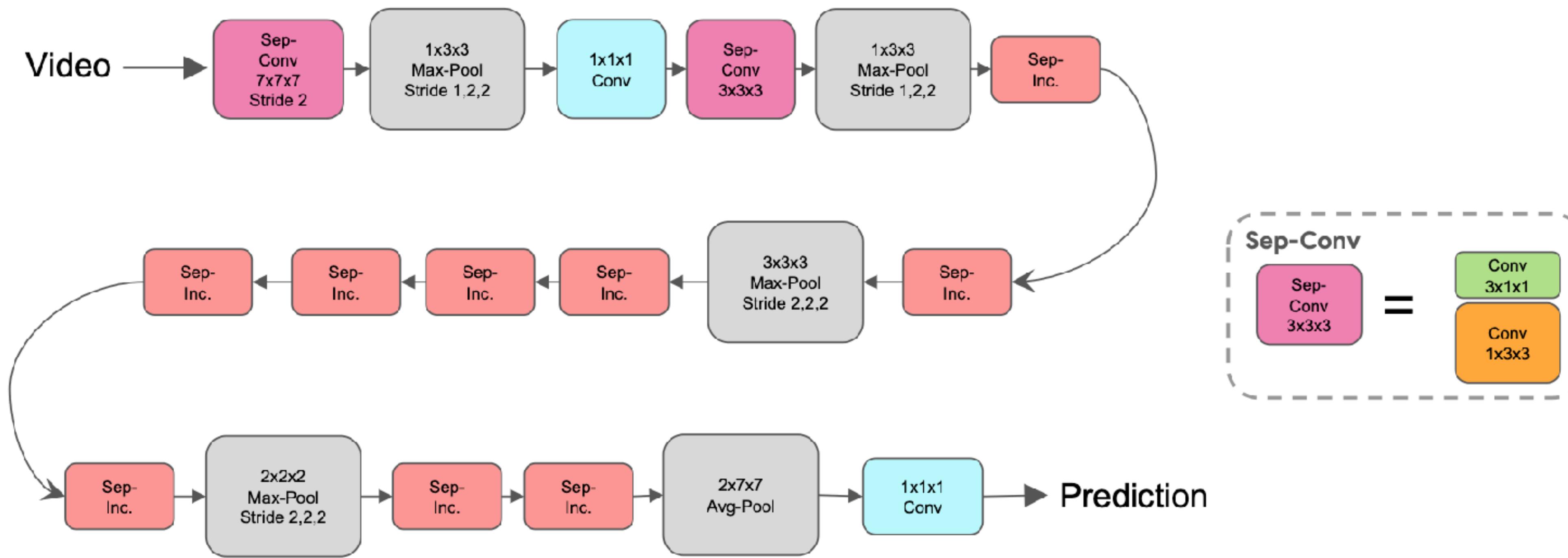


(c) Bottom-heavy I3D



(d) Top-heavy I3D

# S3D – Separable 3D CNN



$$k_t \times k \times k$$

$$1 \times k \times k$$

$$k_t \times 1 \times 1$$

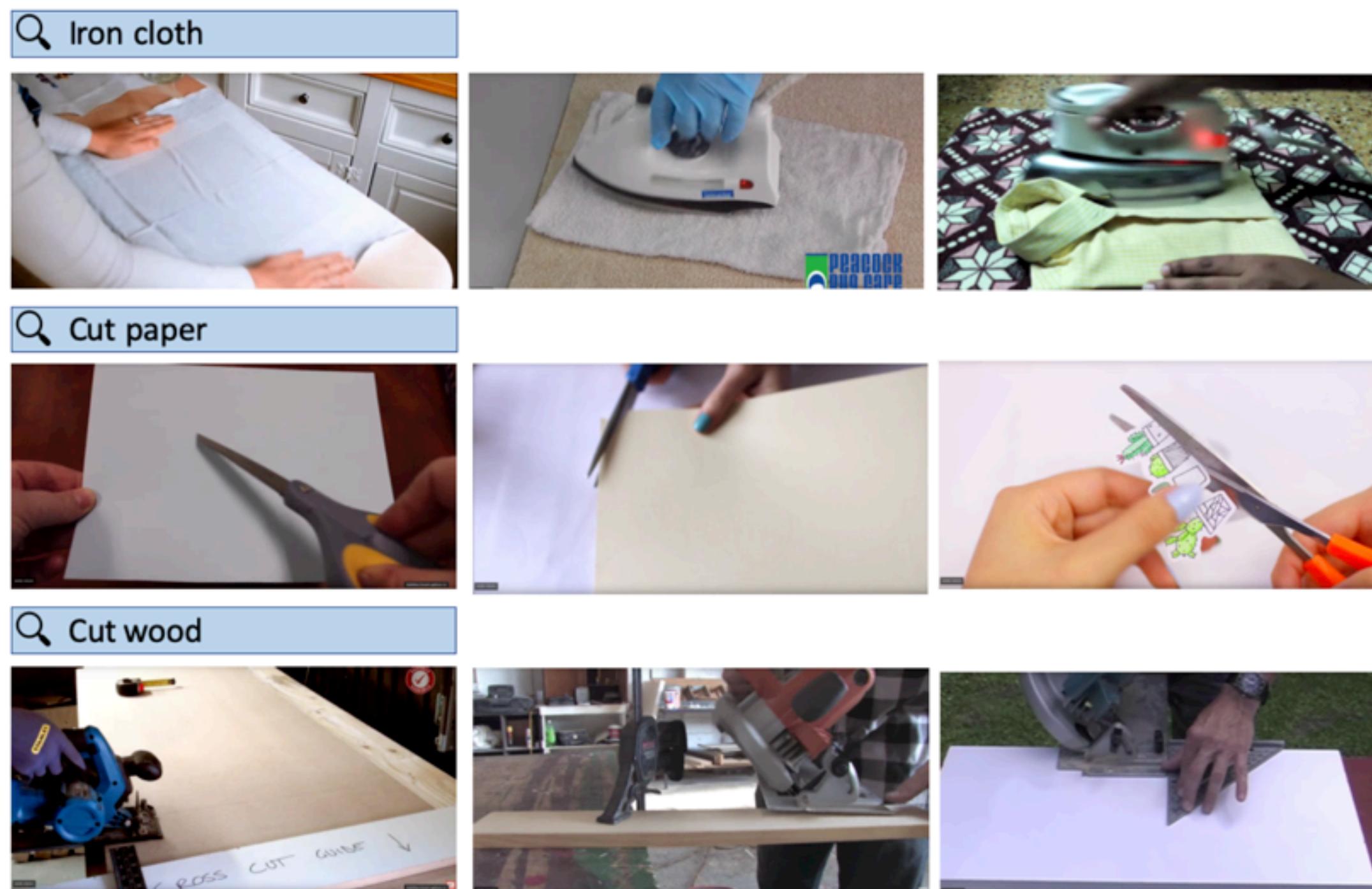
$k_t$

Width of the filter in time

# Data (Something-Something)



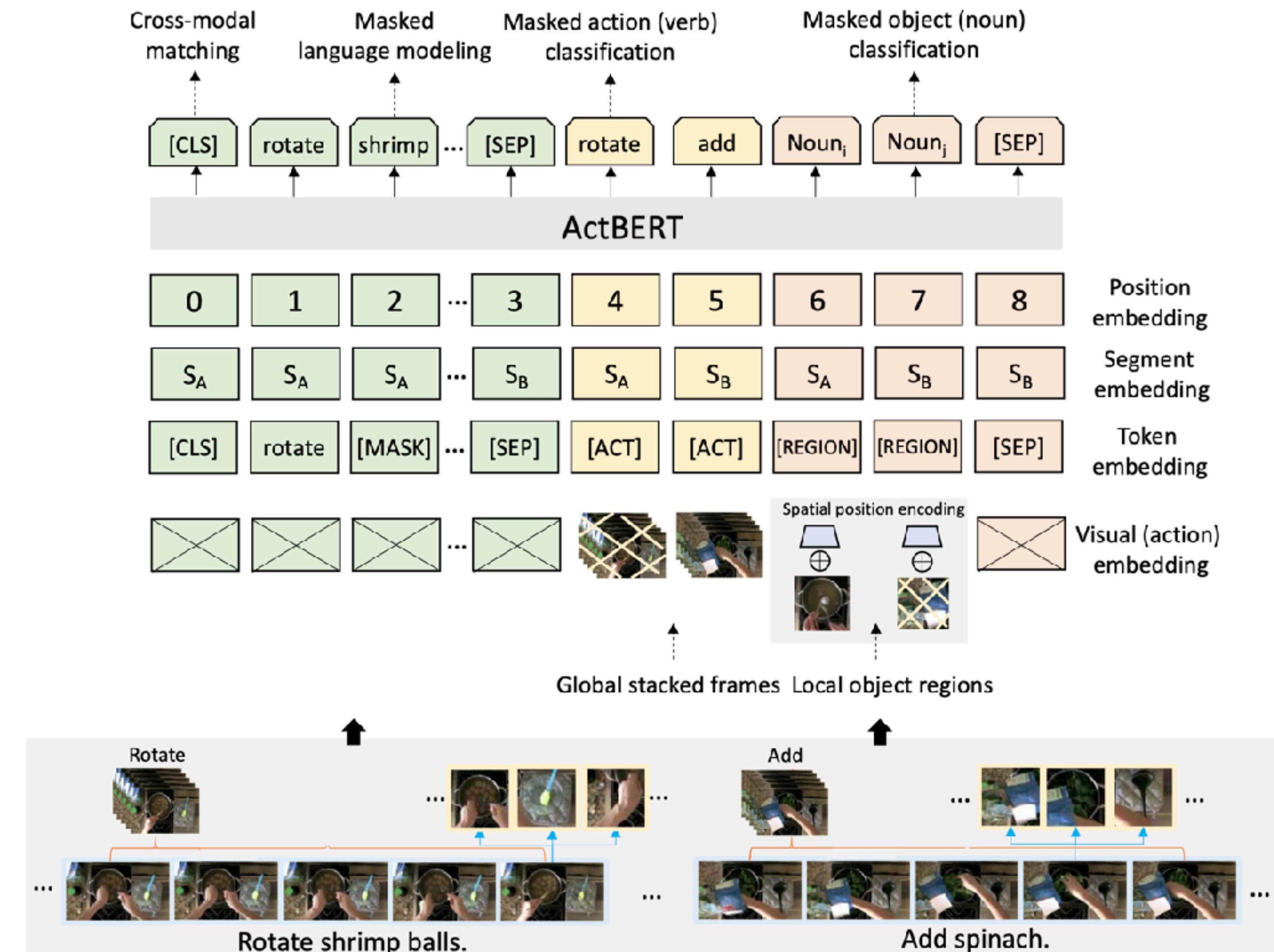
# HowTo100M



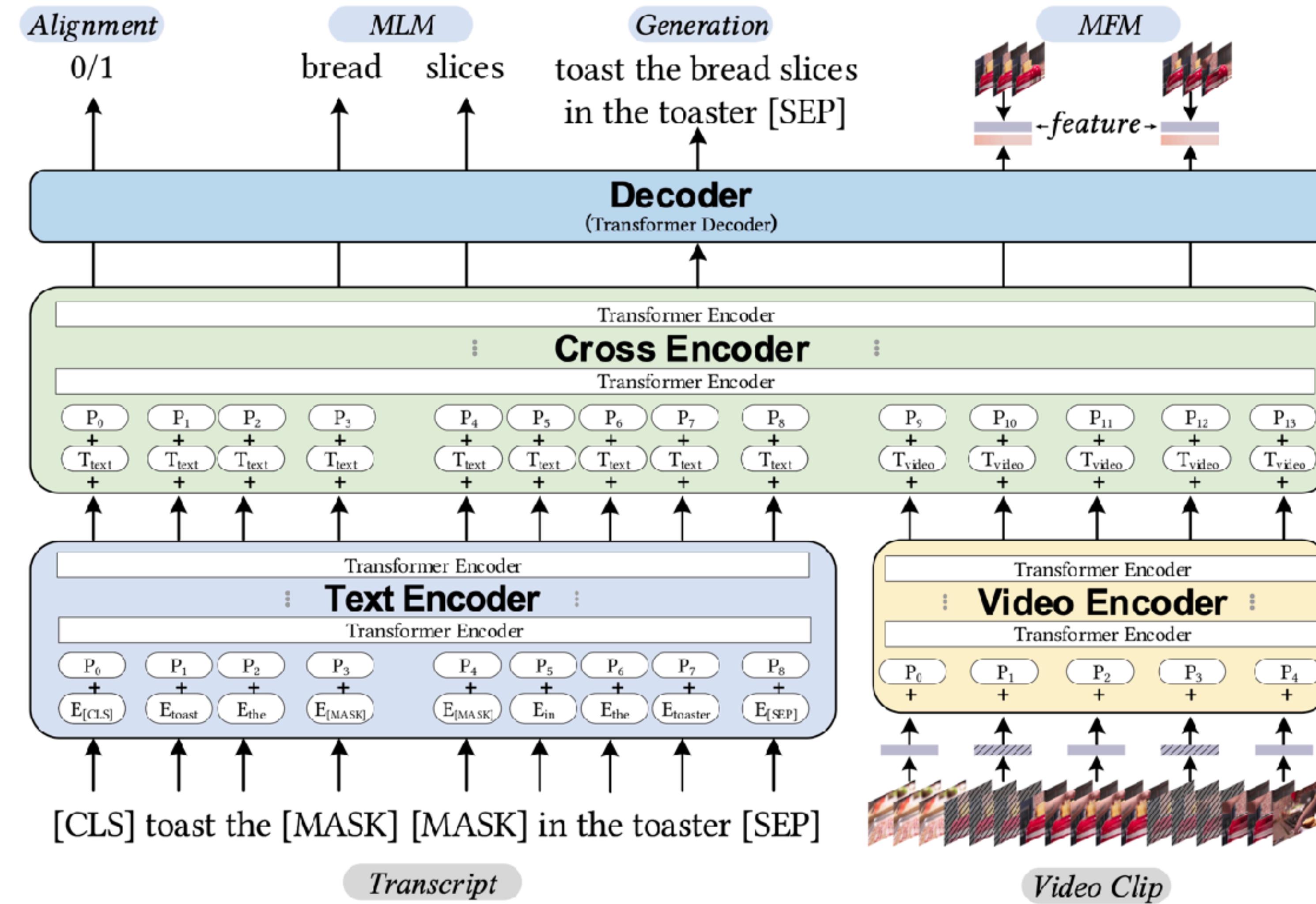
Category	Tasks	Videos	Clips
Food and Entertaining	11504	497k	54.4M
Home and Garden	5068	270k	29.5M
Hobbies and Crafts	4273	251k	29.8M
Cars & Other Vehicles	810	68k	7.8M
Pets and Animals	552	31k	3.5M
Holidays and Traditions	411	27k	3.0M
Personal Care and Style	181	16k	1.6M
Sports and Fitness	205	16k	2.0M
Health	172	15k	1.7M
Education and Communications	239	15k	1.6M
Arts and Entertainment	138	10k	1.2M
Computers and Electronics	58	5k	0.6M
Total	23.6k	1.22M	136.6M

Table 2: Number of tasks, videos and clips within each category.

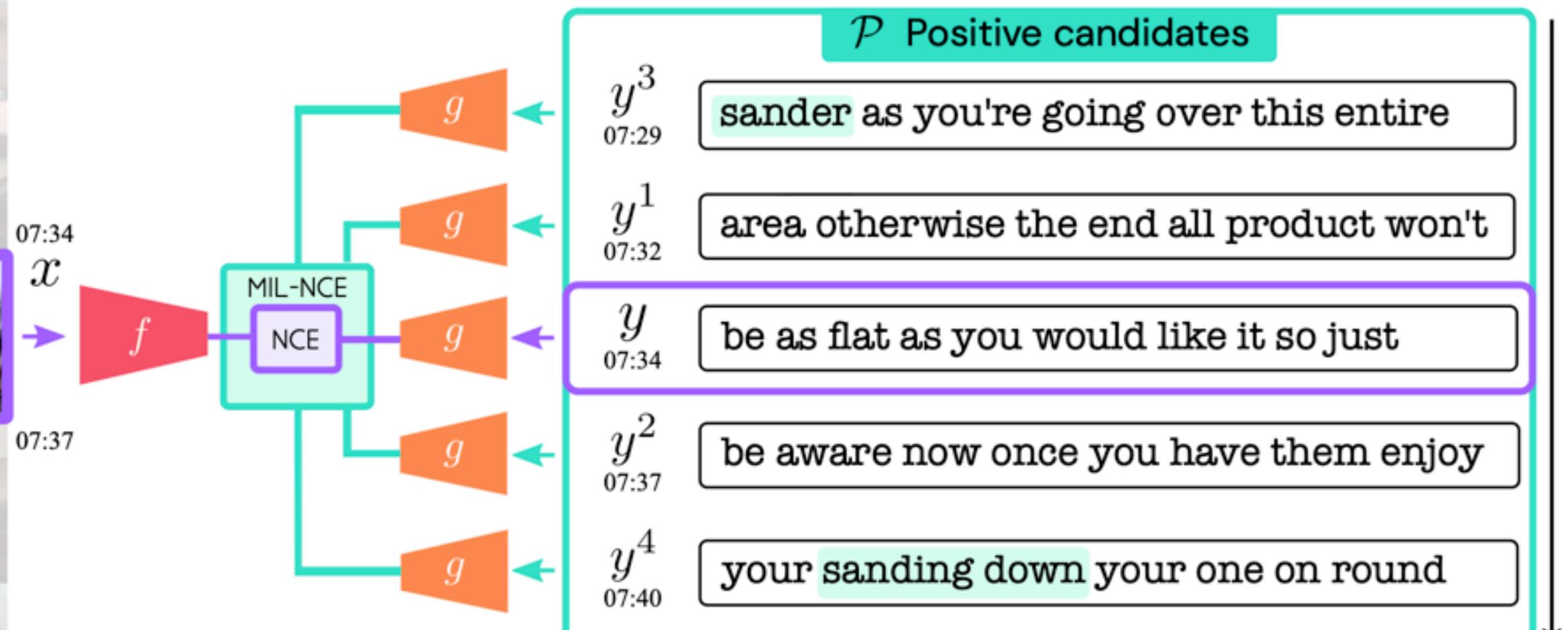
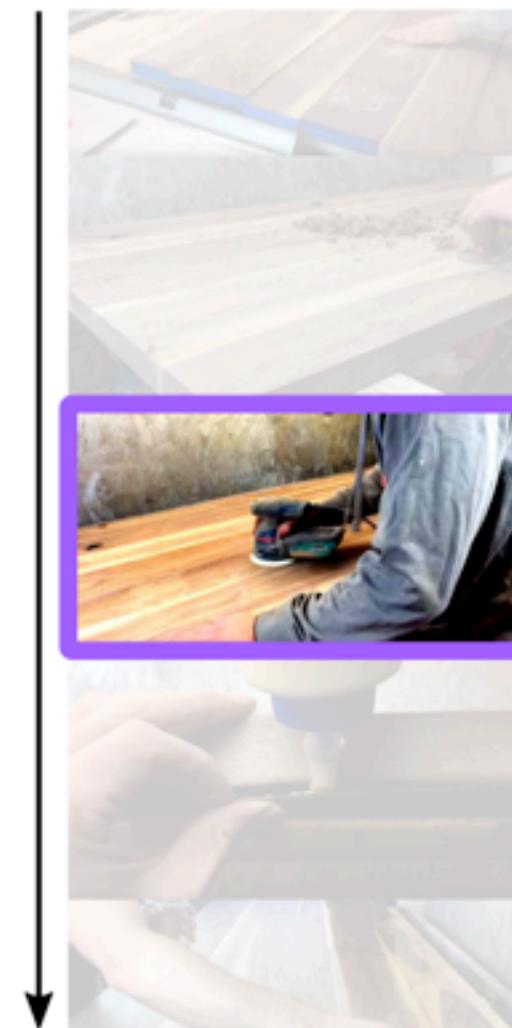
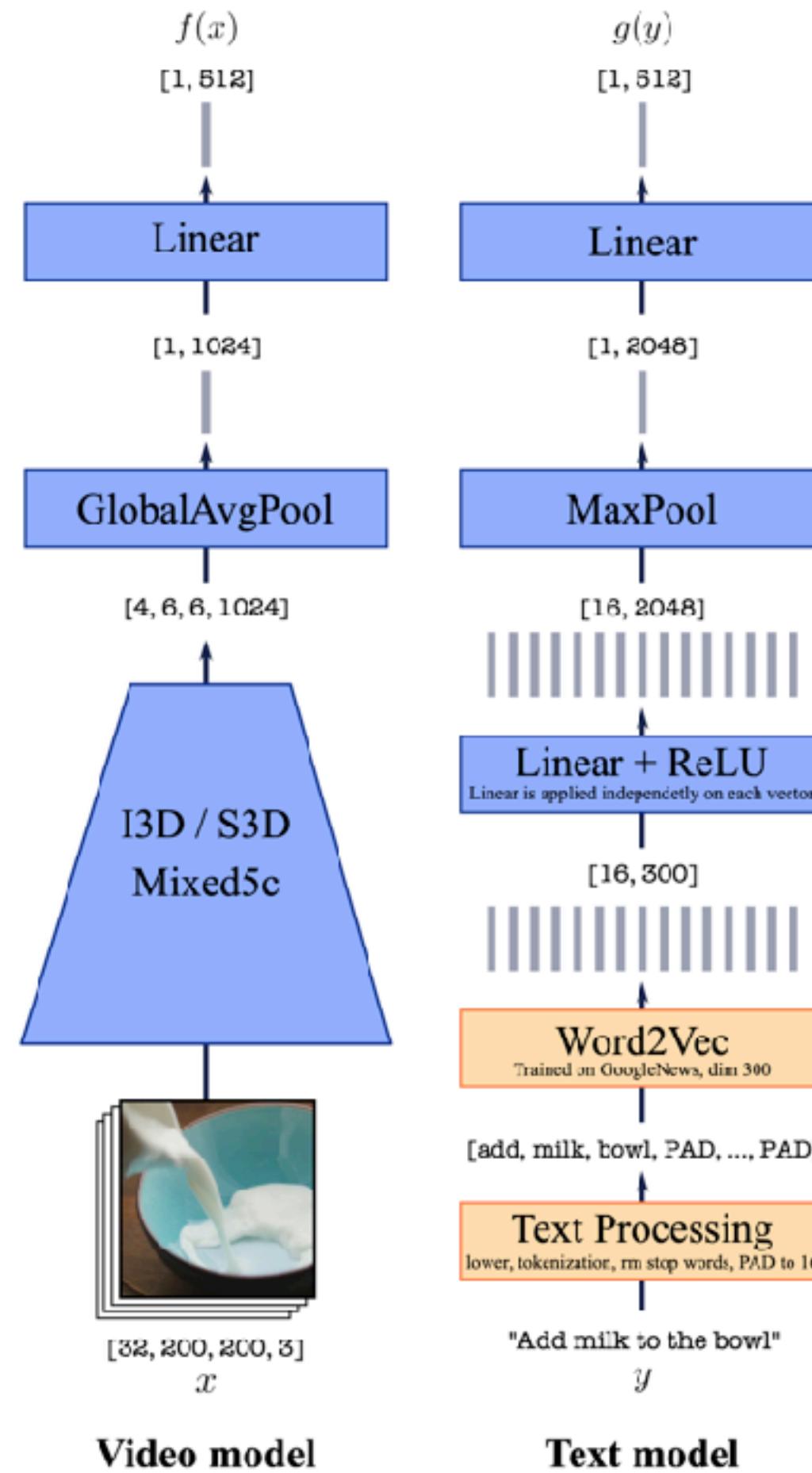
# ActBERT



# UniViLM



# End-to-End Training



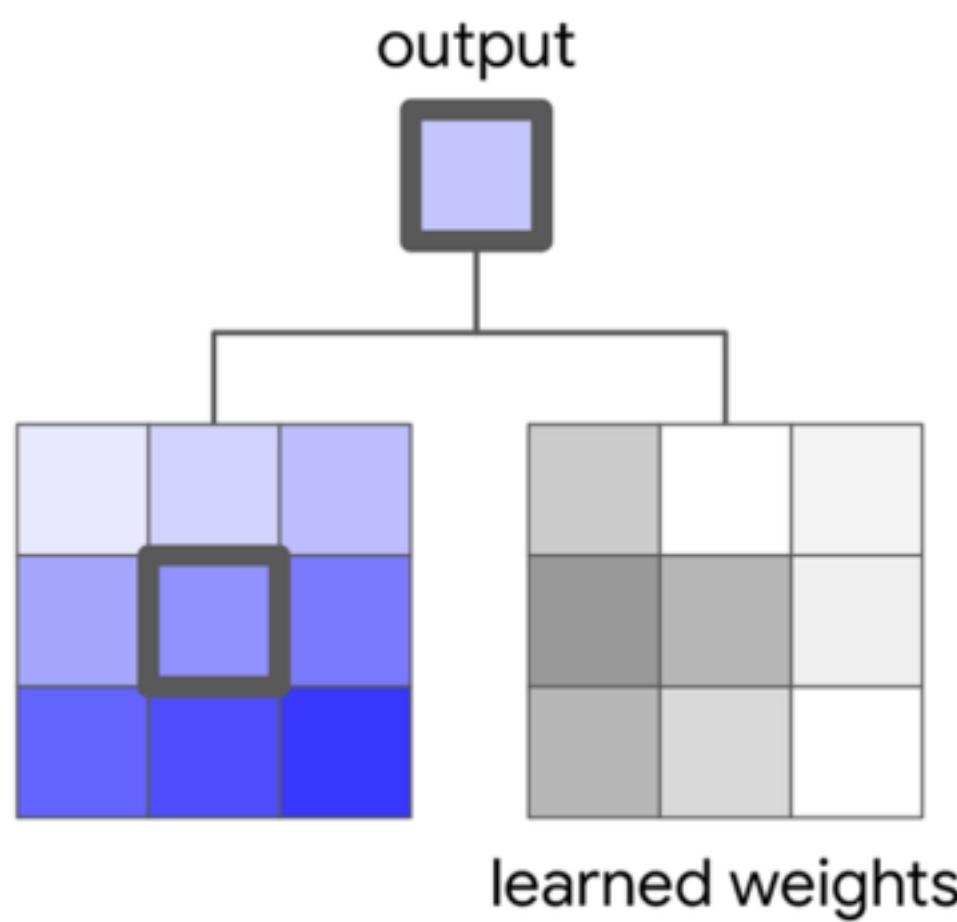
# End-to-End Training

Input word vectors	YR10	MR10
BERT wo. stop words	19.0	19.7
BERT w. stop words	17.6	23.9
Word2Vec	<b>35.0</b>	<b>31.8</b>

[https://github.com/antoine77340/MIL-NCE\\_HowTo100M](https://github.com/antoine77340/MIL-NCE_HowTo100M)

Implementation	Epochs	Total batch size	Accelerator	CPU cores	Training input size	R@1	R@5	R@10	Median Rank
TPU Tensorflow [1]	<b>2645</b>	<b>8192</b>	64 x Cloud TPU v3 128Gb	64 x N.A.	32 frames at 200x200	15.1	38.0	51.2	10
TPU Tensorflow [1]	<b>206</b>	<b>512</b>	4 x Cloud TPU v3 128Gb	4 x N.A.	16 frames at 200x200	8.2	24.6	36.2	23
This implementation	<b>150</b>	<b>512</b>	2 x 4 Tesla V100 32Gb	2 x 40	16 frames at 224x224	7.4	21.5	32.2	29
This implementation	<b>300</b>	<b>1024</b>	4 x 4 Tesla V100 32Gb	4 x 40	16 frames at 224x224	8.8	24.3	34.6	23

# Replacing a CNN w/ Self-Attention



Convolution

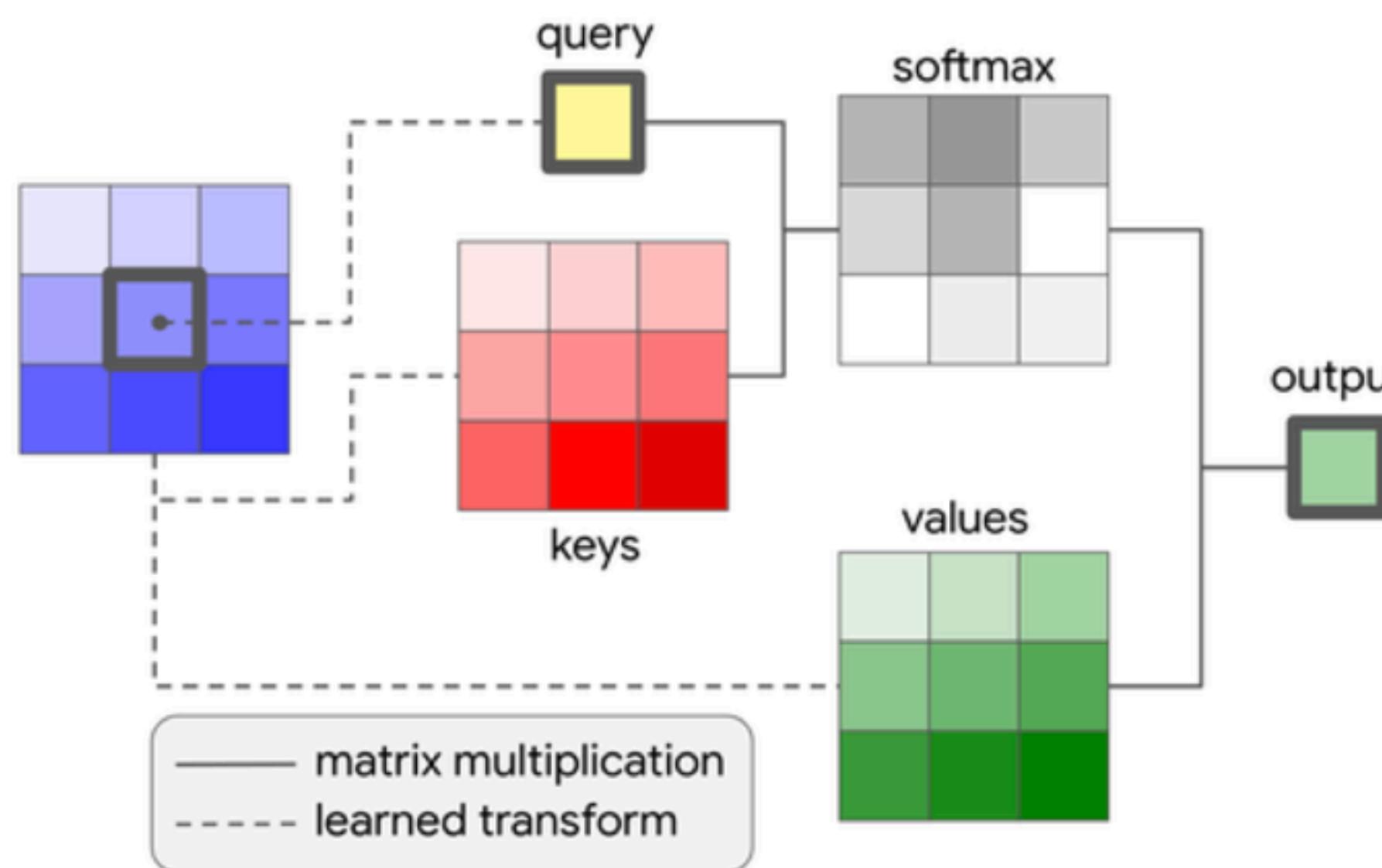


Figure 3: An example of a local attention layer over spatial extent of  $k = 3$ .

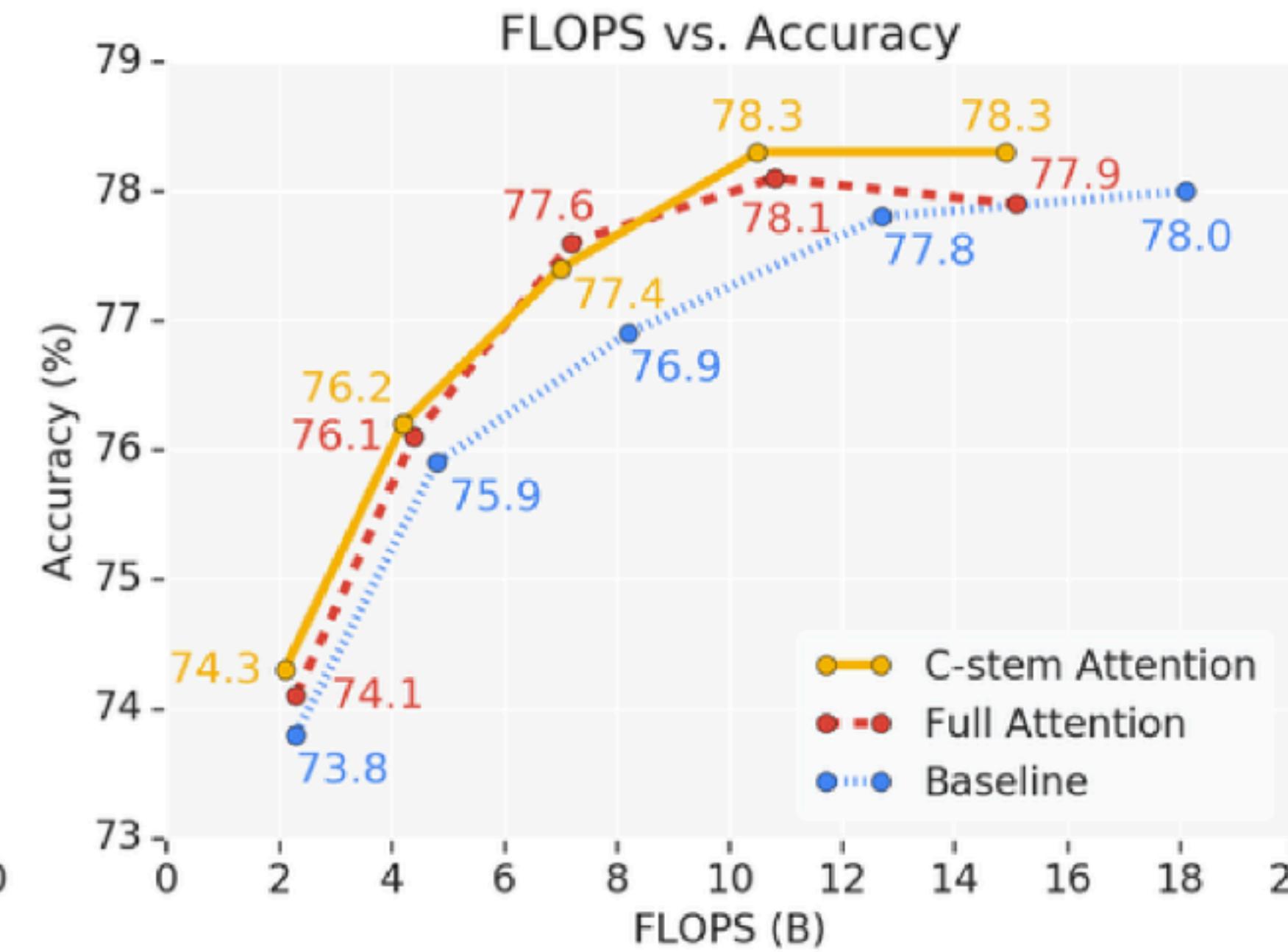
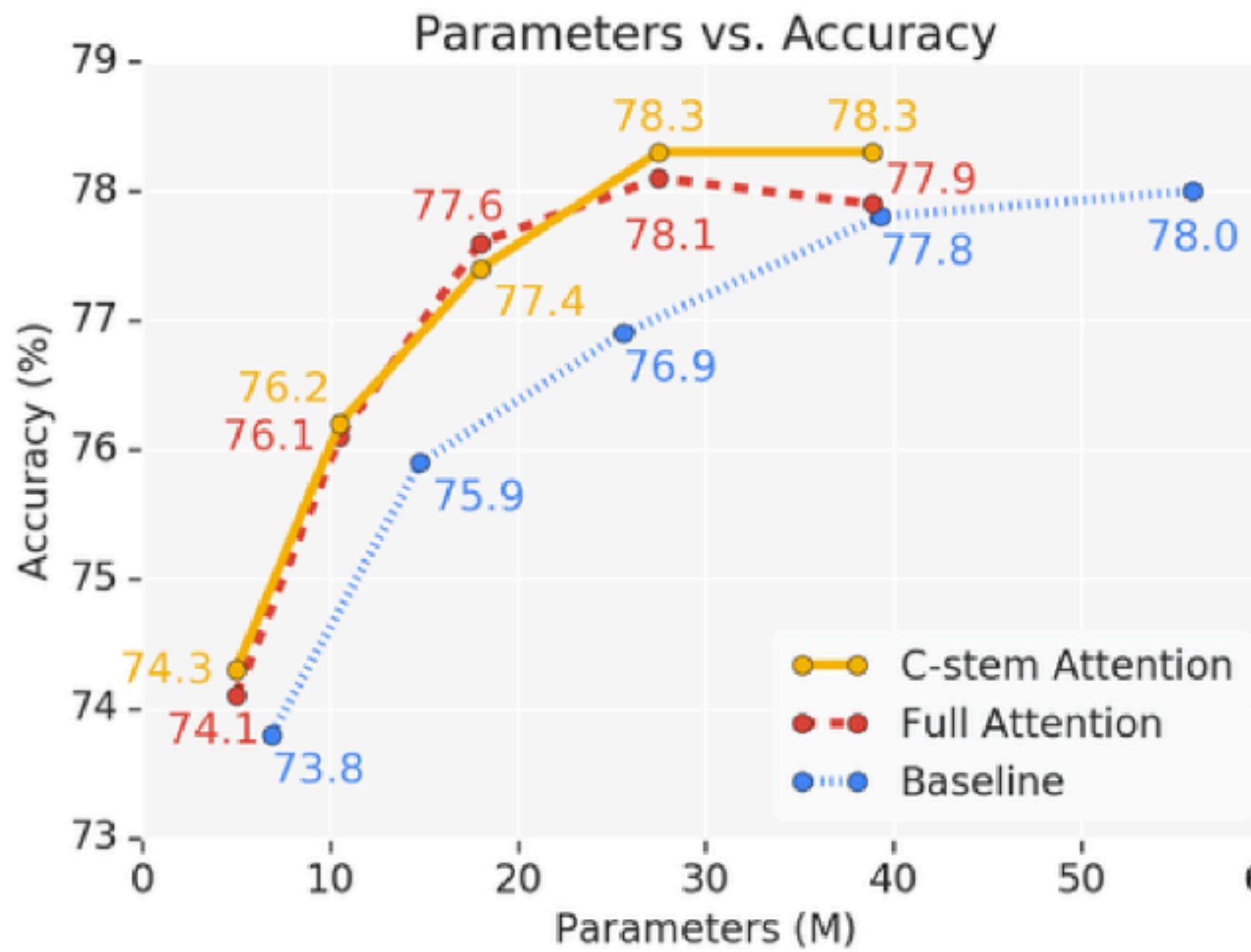
-1, -1	-1, 0	-1, 1	-1, 2
0, -1	0, 0	0, 1	0, 2
1, -1	1, 0	1, 1	1, 2
2, -1	2, 0	2, 1	2, 2

Parameters vs Computation

How do they change with receptive field?

# Parameters/FLOPs vs Accuracy

ResNet 50 and k=7, 8 attention heads

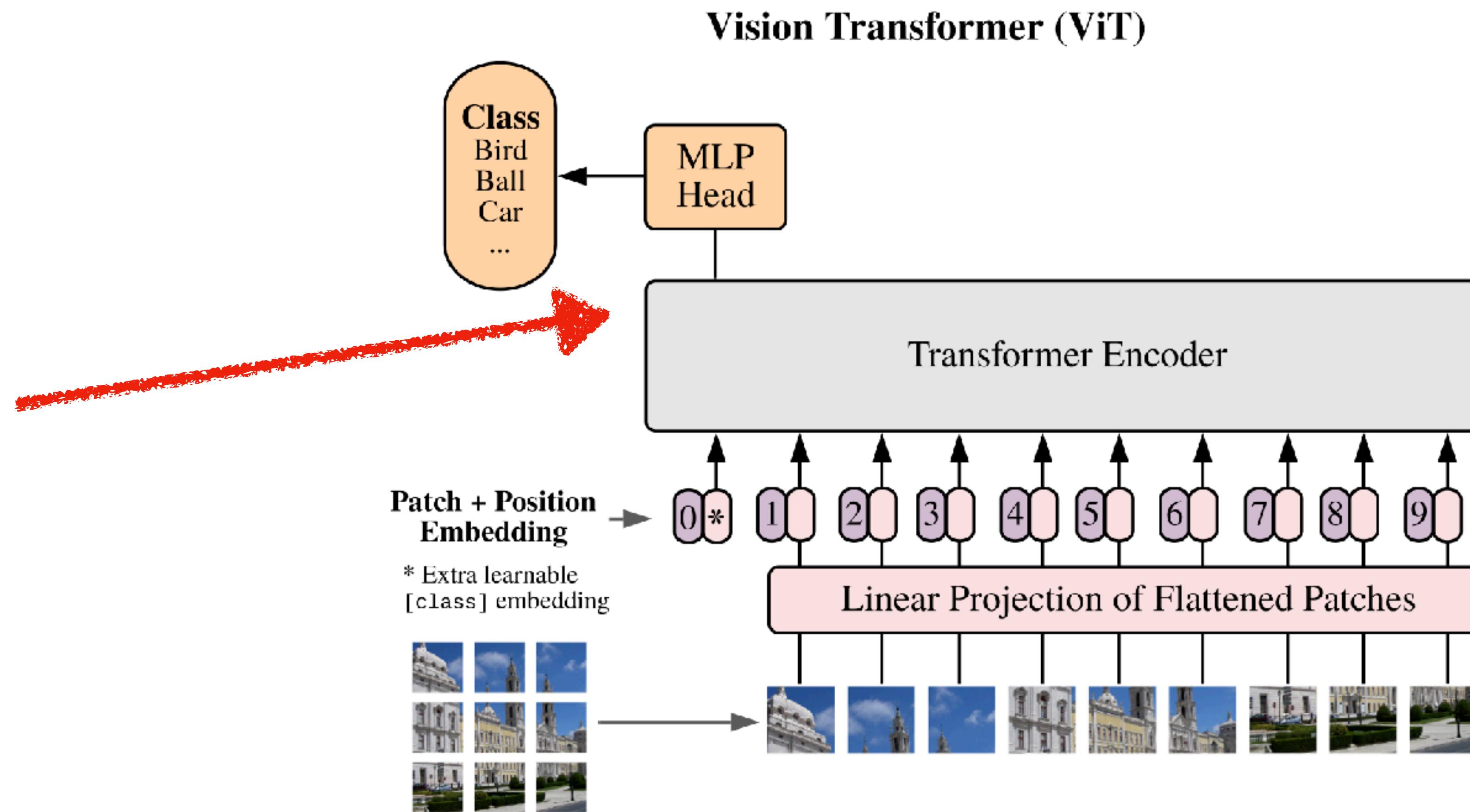


Spatial Extent ( $k \times k$ )	FLOPS (B)	Top-1 Acc. (%)
$3 \times 3$	6.6	76.4
$5 \times 5$	6.7	77.2
$7 \times 7$	7.0	77.4
$9 \times 9$	7.3	77.7
$11 \times 11$	7.7	77.6

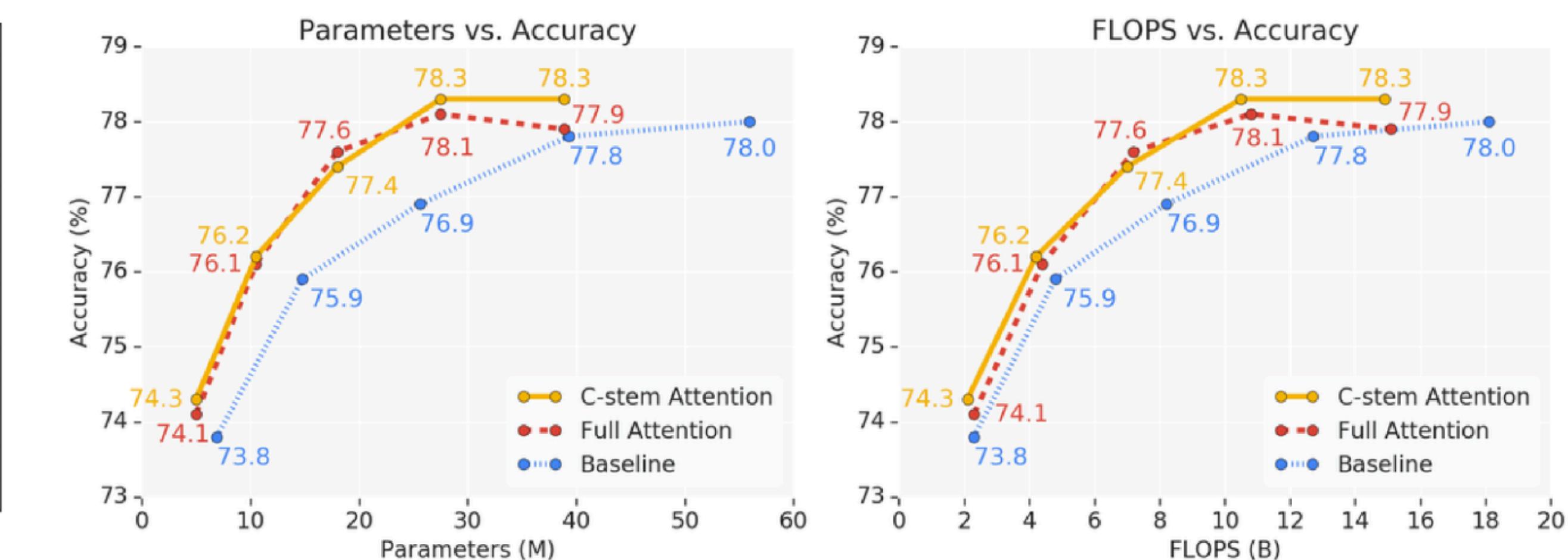
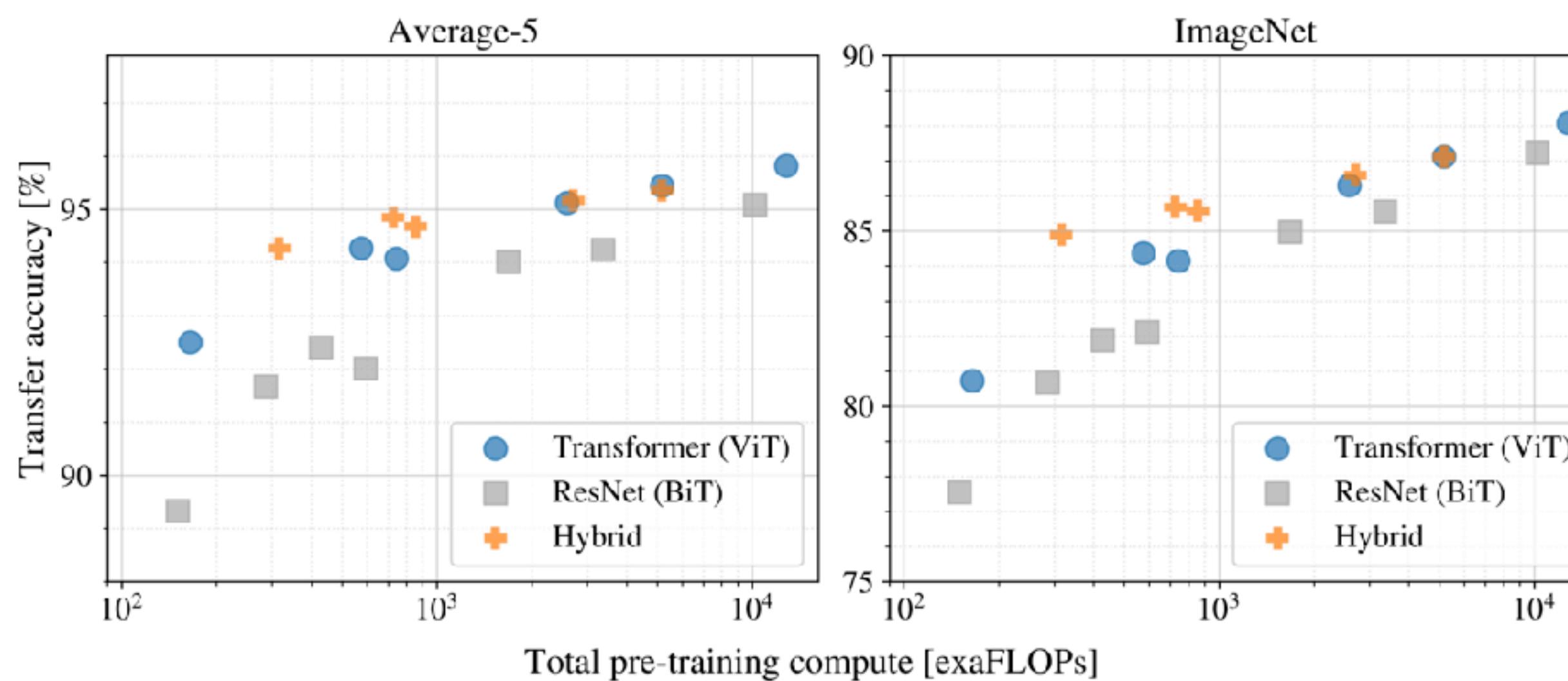
ViT



# “Visual” Features



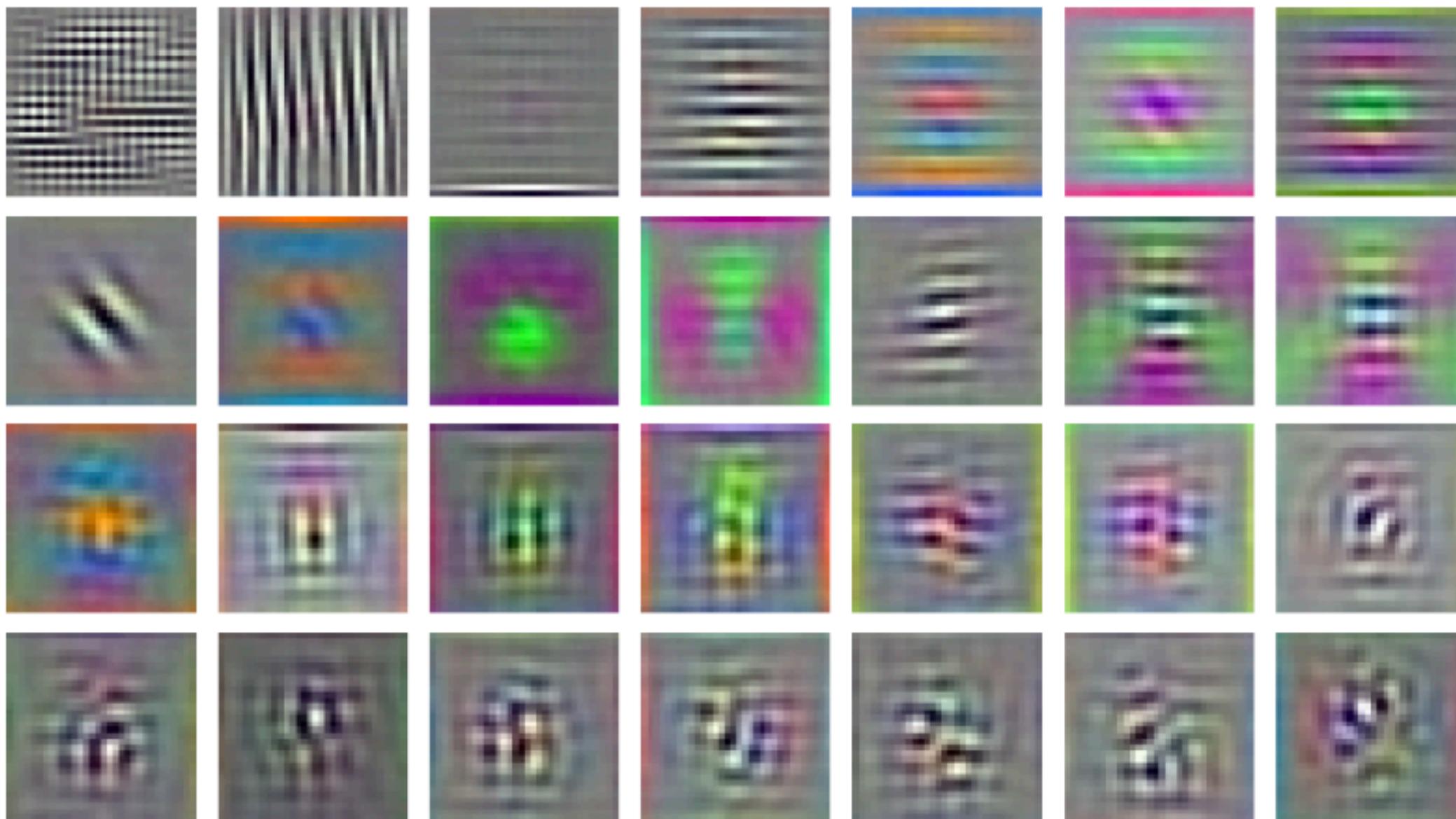
# Curves



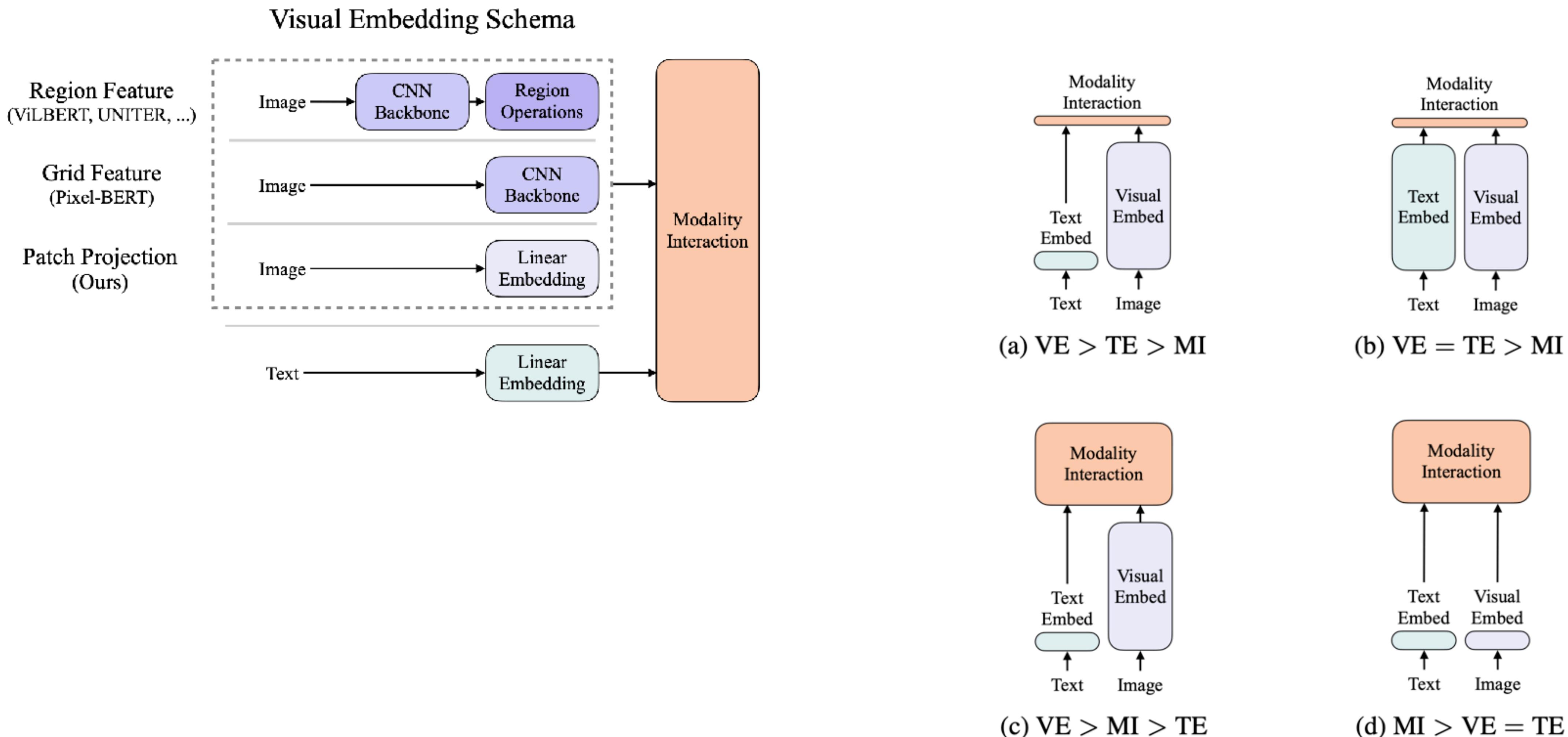
# Which is the best?

ImageNet	ImageNet ReaL	CIFAR-10	CIFAR-100	Pets	Flowers	exaFLOPs
80.73	86.27	98.61	90.49	93.40	99.27	164
84.15	88.85	99.00	91.87	95.80	99.56	743
84.37	88.28	99.19	92.52	95.83	99.45	574
86.30	89.43	99.38	93.46	96.81	99.66	2586
87.12	89.99	99.38	94.04	97.11	99.56	5172
88.08	90.36	99.50	94.71	97.11	99.71	12826
77.54	84.56	97.67	86.07	91.11	94.26	150
82.12	87.94	98.29	89.20	93.43	97.02	592
80.67	87.07	98.48	89.17	94.08	95.95	285
81.88	87.96	98.82	90.22	94.17	96.94	427
84.97	89.69	99.06	92.05	95.37	98.62	1681
85.56	89.89	99.24	91.92	95.75	98.75	3362
87.22	90.15	99.34	93.53	96.32	99.04	10212
84.90	89.15	99.01	92.24	95.75	99.46	315
85.58	89.65	99.14	92.63	96.65	99.40	855
85.68	89.04	99.24	92.93	96.97	99.43	725
86.60	89.72	99.18	93.64	97.03	99.40	2704
87.12	89.76	99.31	93.89	97.36	99.11	5165

# Filters

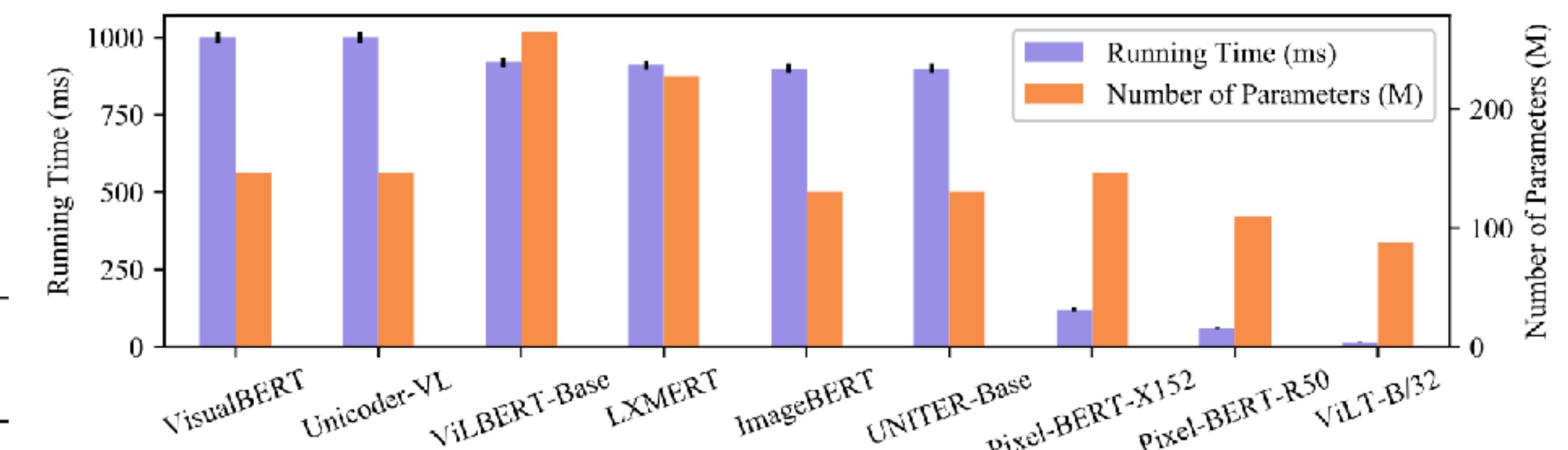


# What about ROIs? What about Different Pipelines?

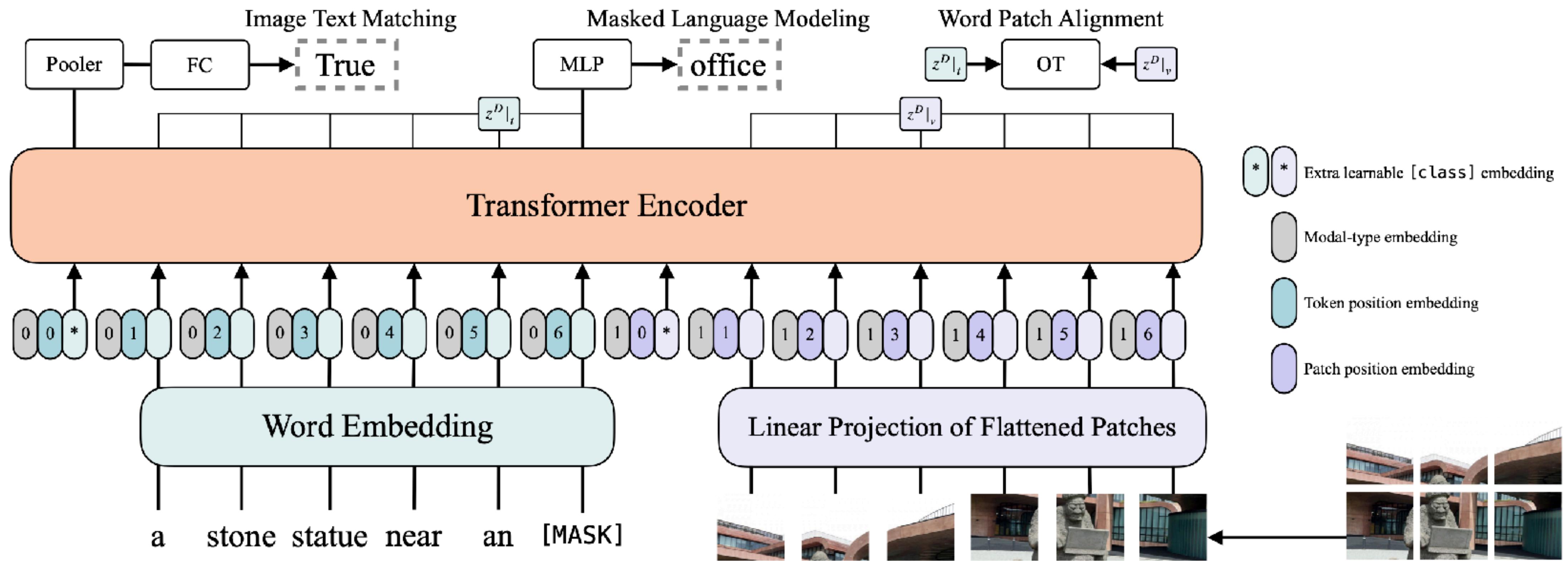


# Faster Inference? What about training the features?

Visual Embed	Model	Time (ms)	VQAv2 test-dev	NLVR2 dev	NLVR2 test-P
Region	w/o VLP SOTA	~900	70.63	54.80	53.50
	ViLBERT	~920	70.55	-	-
	VisualBERT-Base	~1000	70.80	67.40	67.00
	LXMERT	~910	72.42	74.90	74.50
	UNITER-Base	~900	72.70	75.85	75.80
	OSCAR-Base <sup>†</sup>	~900	73.16	78.07	78.36
	VinVL-Base <sup>†‡</sup>	~1000	75.95	82.05	83.08
Grid	Pixel-BERT-X152	~120	74.45	76.50	77.20
	Pixel-BERT-R50	~60	71.35	71.70	72.40
Linear	ViLT-B/32	~15	70.34	74.56	74.66
	ViLT-B/32 <sup>④</sup>	~15	70.94	75.24	76.21



# Transformer “figures it out”



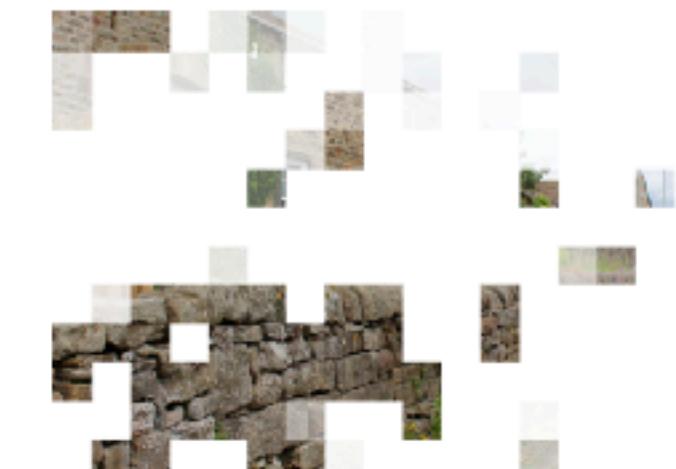
# Alignment



a display of **flowers** growing out and over the retaining **wall** in front of **cottages** on a **cloudy** day.



**flowers**



**wall**



**cottages**



**cloudy**



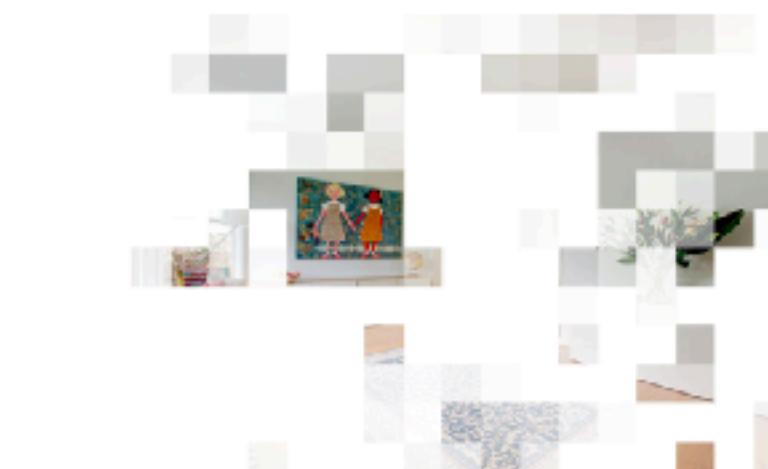
a room with a **rug**, a **chair**, a **painting**, and a **plant**.



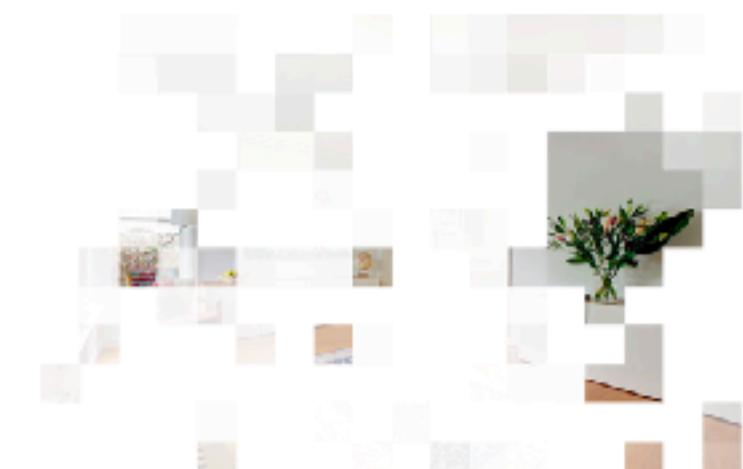
**rug**



**chair**



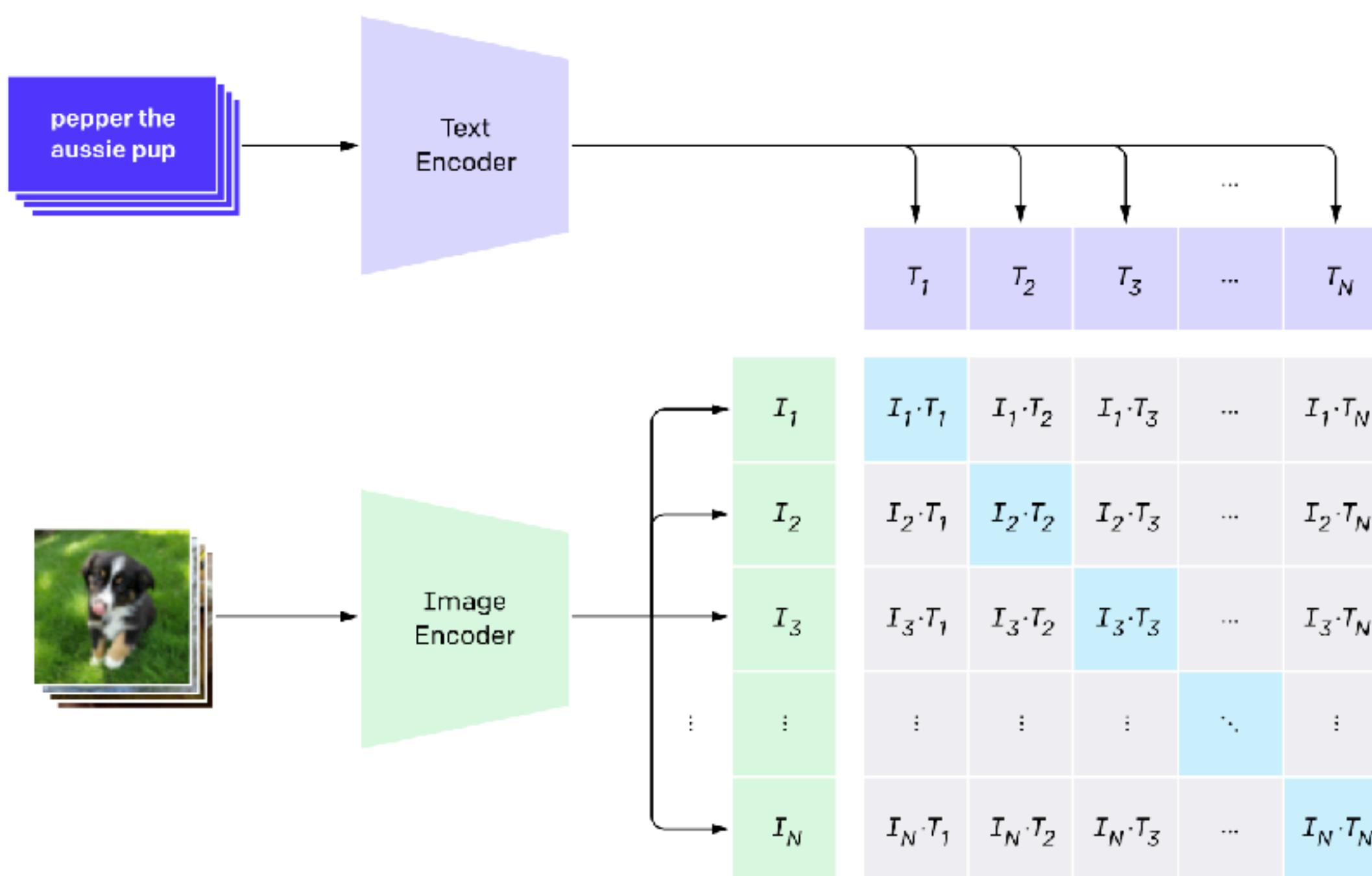
**painting**



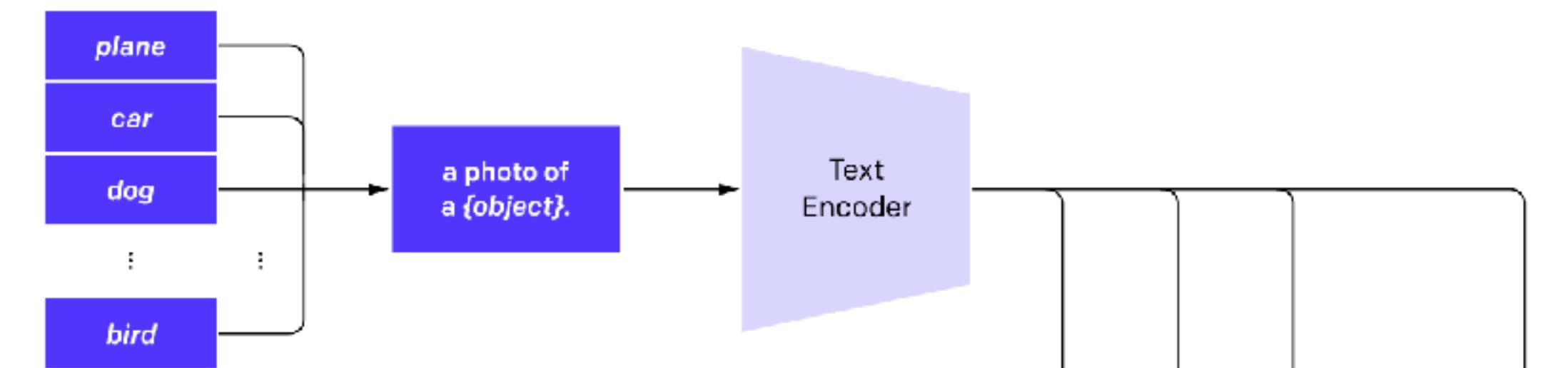
**plant**

# CLIP – 400 million (image, text) pairs

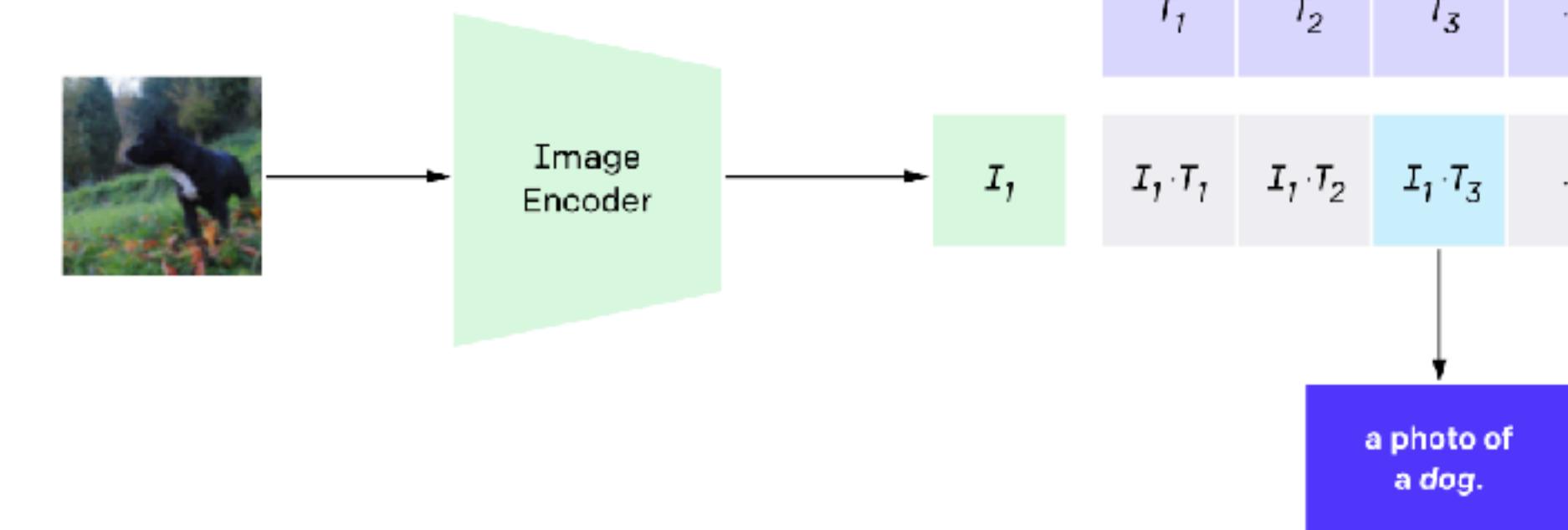
## 1. Contrastive pre-training



## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction



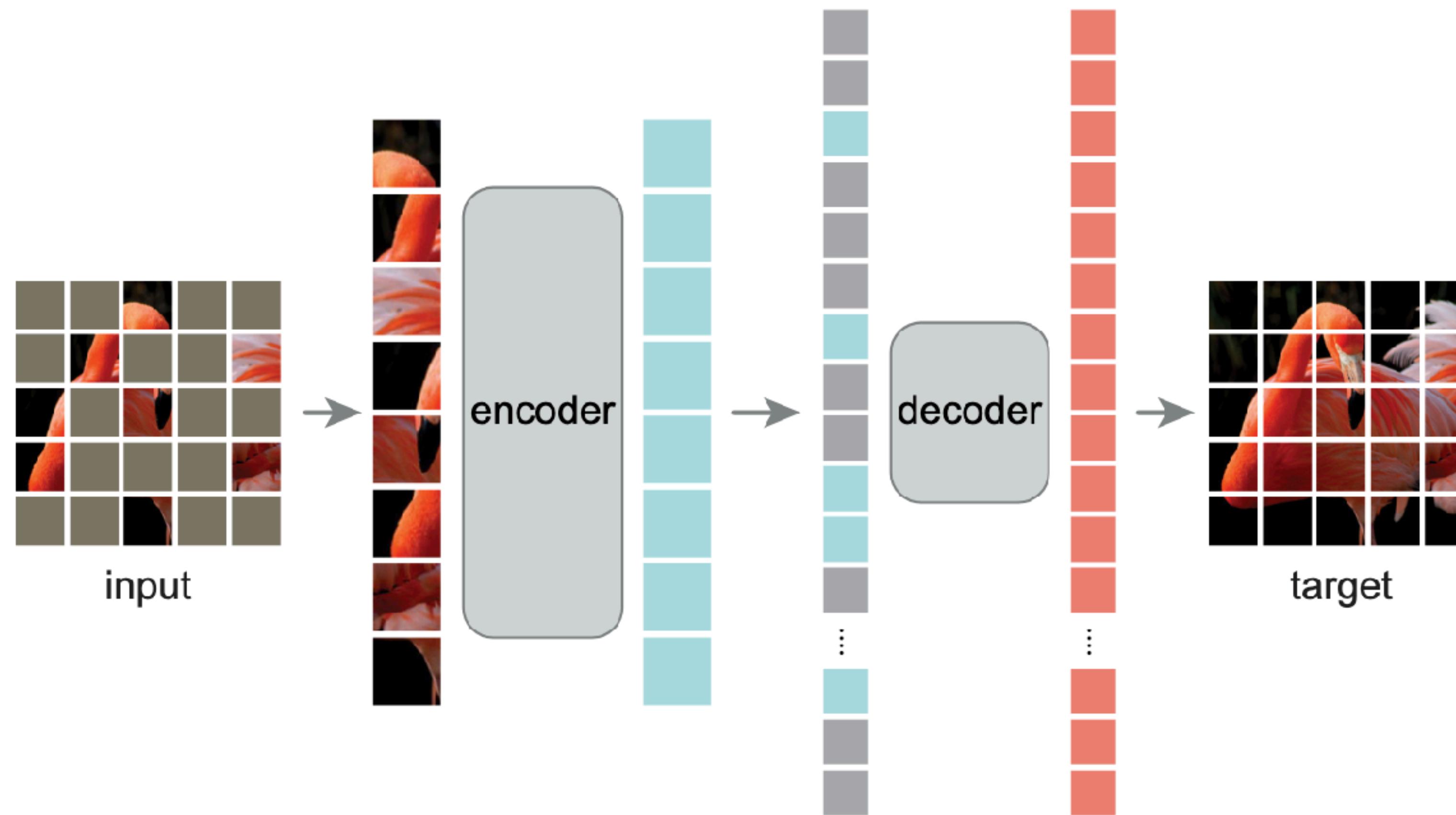
# So should we Transformer?



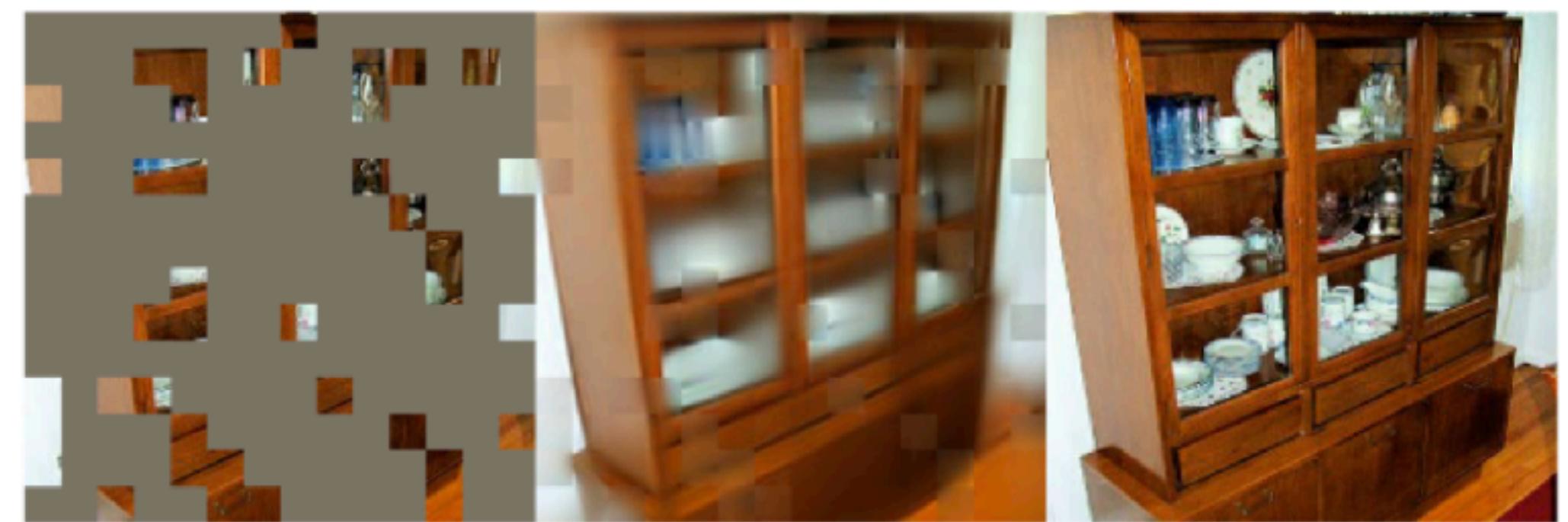
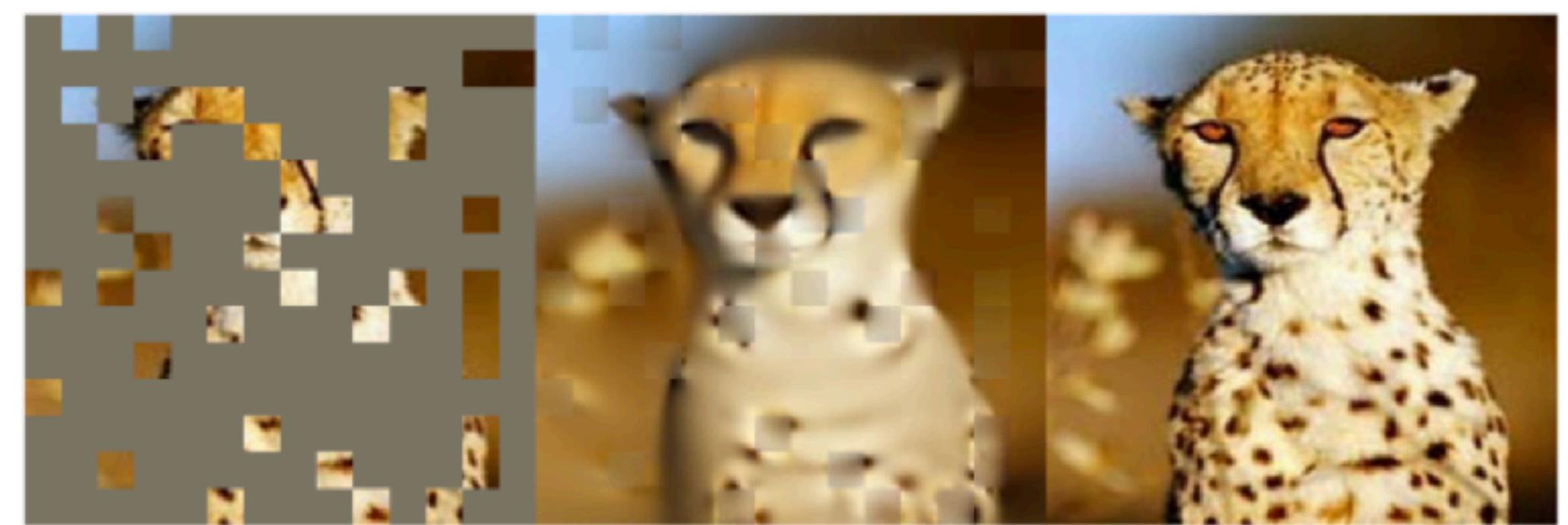
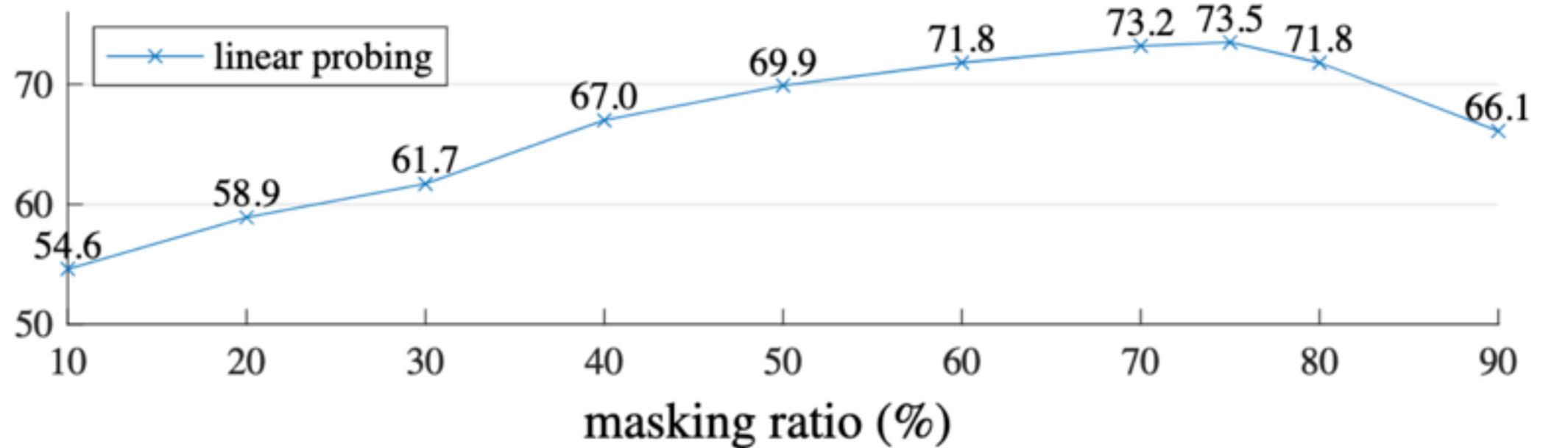
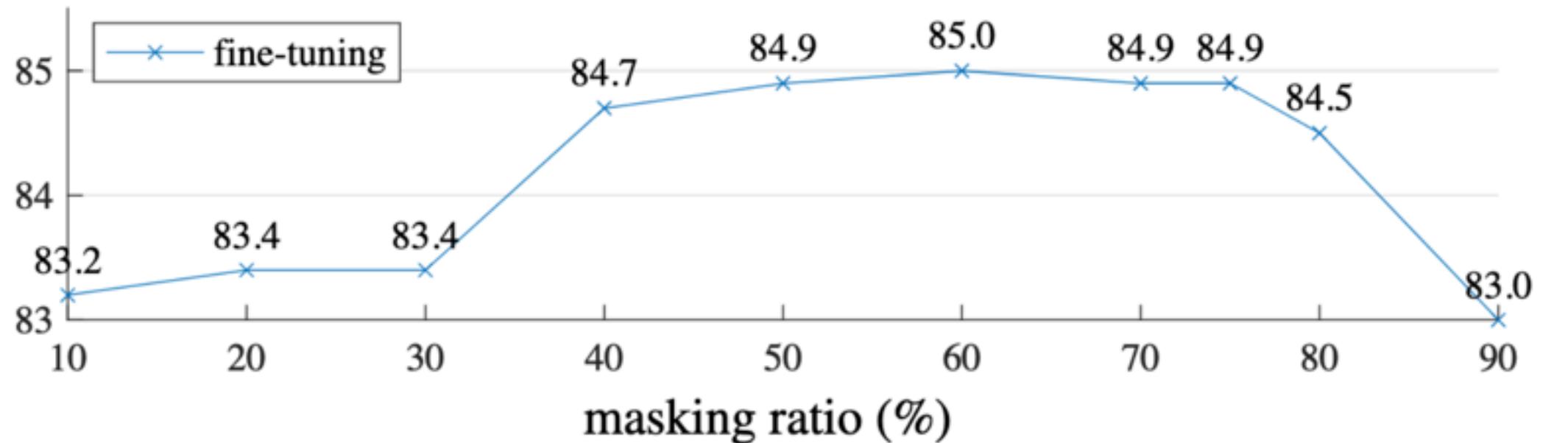
		Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech101	Flowers	MINIST	FER2013	STL10*	EuroSAT	RESISC45	GTSRB	KITTI	PCAM	UCF101	Kinetics700	CLEVR	HatefulMemes	SST	ImageNet	Country211	
LM RN50		81.3	82.8	61.7	44.2	69.6	74.9	44.9	85.5	71.5	82.8	85.5	91.1	96.6	60.1	95.3	93.4	84.0	73.8	70.2	19.0	82.9	76.4	51.9	51.2	65.2	76.8	65.2	
CLIP-RN	50	86.4	88.7	70.3	56.4	73.3	78.3	49.1	87.1	76.4	88.2	89.6	96.1	98.3	64.2	96.6	95.2	87.5	82.4	70.2	25.3	82.7	81.6	57.2	53.6	65.7	72.6	73.3	
	101	88.9	91.1	73.5	58.6	75.1	84.0	50.7	88.0	76.3	91.0	92.0	96.4	98.4	65.2	97.8	95.9	89.3	82.4	73.6	26.6	82.8	84.0	60.3	50.3	68.2	73.3	75.7	
	50x4	91.3	90.5	73.0	65.7	77.0	85.9	57.3	88.4	79.5	91.9	92.5	97.8	98.5	68.1	97.8	96.4	89.7	85.5	59.4	30.3	83.0	85.7	62.6	52.5	68.0	76.6	78.2	
	50x16	93.3	92.2	74.9	72.8	79.2	88.7	62.7	89.0	79.1	93.5	93.7	98.3	98.9	68.7	98.6	97.0	91.4	89.0	69.2	34.8	83.5	88.0	66.3	53.8	71.1	<b>80.0</b>	81.5	
	50x64	94.8	94.1	78.6	77.2	81.1	90.5	67.7	88.9	82.0	94.5	95.4	98.9	98.9	71.3	99.1	97.1	92.8	90.2	69.2	40.7	83.7	89.5	75.0	55.0	<b>81.2</b>	83.6		
CLIP-ViT	B/32	88.8	95.1	80.5	58.5	76.6	81.8	52.0	87.7	76.5	90.0	93.0	96.9	<b>99.0</b>	69.2	98.3	97.0	90.5	85.3	66.2	27.8	83.9	85.5	61.7	52.1	66.7	70.8	76.1	
	B/16	92.8	96.2	83.1	67.8	78.4	86.7	59.5	89.2	79.2	93.1	94.7	98.1	<b>99.0</b>	69.5	99.0	97.1	92.7	86.6	67.8	33.3	83.5	88.4	66.1	57.1	70.3	75.5	80.2	
	L/14	95.2	98.0	87.5	77.0	<b>81.8</b>	<b>90.9</b>	69.4	89.6	82.1	<b>95.1</b>	<b>96.5</b>	99.2	<b>99.2</b>	<b>72.2</b>	<b>99.7</b>	<b>98.2</b>	<b>94.1</b>	<b>92.5</b>	<b>64.7</b>	42.9	85.8	<b>91.5</b>	72.0	<b>57.8</b>	<b>76.2</b>	<b>80.8</b>	83.9	
	L/14-336px	<b>95.9</b>	97.9	87.4	<b>79.9</b>	<b>82.2</b>	<b>91.5</b>	<b>71.6</b>	89.9	83.0	<b>95.1</b>	<b>96.0</b>	99.2	<b>99.2</b>	<b>72.9</b>	<b>99.7</b>	<b>98.1</b>	<b>94.9</b>	<b>92.4</b>	69.2	46.4	85.6	<b>92.0</b>	73.0	60.3	<b>77.3</b>	<b>80.5</b>	85.4	
	B/8	74.3	92.5	76.5	59.7	62.0	62.5	55.7	84.4	71.2	93.0	93.3	91.7	98.2	57.2	97.1	97.3	85.5	80.0	73.8	12.4	83.1	74.4	47.6	55.7	53.4	76.9		
EfficientNet	B/1	74.2	93.2	77.2	61.3	62.6	62.5	56.1	84.7	74.2	93.4	93.6	92.4	98.3	57.0	97.5	96.8	84.5	75.9	75.5	12.5	82.7	74.7	48.5	44.3	54.5	54.4	78.6	
	B/2	75.8	93.6	77.9	64.4	64.0	63.2	57.0	85.3	73.5	93.9	93.5	92.9	98.5	56.6	97.7	96.9	84.4	76.4	73.1	12.6	84.3	75.1	49.4	42.6	55.4	55.2	79.7	
	B/3	77.4	94.0	78.0	66.5	64.4	66.0	59.3	85.8	73.1	94.1	93.7	93.3	98.5	57.1	98.2	97.3	85.0	75.8	76.1	13.4	83.3	78.1	50.9	45.1	53.8	54.8	81.0	
	B/4	79.7	94.1	78.7	70.1	65.4	66.4	60.4	86.5	73.4	94.7	93.5	93.2	98.8	57.9	98.6	96.8	85.0	78.3	72.3	13.9	83.1	79.1	52.5	46.5	54.4	55.4	82.9	
	B/5	81.5	93.6	77.9	72.4	67.1	72.7	68.9	86.7	73.9	<b>95.0</b>	94.7	94.5	98.4	58.5	98.7	96.8	86.0	78.5	69.6	14.9	84.7	80.9	54.5	46.6	53.3	56.3	83.7	
	B/6	82.4	94.0	78.0	73.5	65.8	71.1	68.2	87.6	73.9	<b>95.0</b>	94.1	93.7	98.4	60.2	98.7	96.8	85.4	78.1	72.7	15.3	84.2	80.0	54.1	51.1	53.3	57.0	84.0	
	B/7	84.5	94.9	80.1	74.7	69.0	77.1	<b>72.3</b>	87.2	76.8	<b>95.2</b>	94.7	95.9	98.6	61.3	99.1	96.3	86.8	80.8	<b>75.8</b>	16.4	85.2	81.9	56.8	51.9	54.4	57.8	84.8	
	B/8	84.5	95.0	80.7	75.2	69.6	76.8	<b>71.5</b>	87.4	77.1	<b>94.9</b>	95.2	96.3	98.6	61.4	99.2	97.0	87.4	80.4	70.9	17.4	85.2	82.4	57.7	51.4	55.8	55.8	85.3	
EfficientNet Noisy Student	B/0	78.1	94.0	78.6	63.5	65.5	57.2	53.7	85.6	75.6	93.8	93.1	94.5	98.1	55.6	98.2	97.0	84.3	74.0	71.6	14.0	83.1	76.7	51.7	55.7	55.0	78.5		
	B/1	80.4	95.1	80.2	66.6	67.6	59.6	53.7	86.2	77.0	91.6	94.4	95.1	98.0	56.1	98.6	96.9	84.3	73.1	67.1	14.5	83.9	79.9	54.5	46.1	54.3	54.9	81.1	
	B/2	80.9	95.3	81.3	67.6	67.9	60.9	55.2	86.3	77.7	<b>95.0</b>	94.7	94.4	98.0	55.5	98.8	97.3	84.6	74.6	71.7	10.0	14.6	82.9	80.1	55.1	46.1	54.1	55.3	82.2
	B/3	82.6	95.9	82.1	68.6	68.8	60.6	55.4	86.5	77.2	<b>95.0</b>	94.8	95.2	98.1	56.0	99.1	96.5	85.0	70.5	69.5	15.1	83.1	81.8	56.8	45.1	55.7	52.0	83.8	
	B/4	85.2	95.6	81.0	72.5	69.7	56.1	52.6	87.0	78.7	<b>94.8</b>	95.2	95.3	98.2	56.0	99.3	95.3	84.8	61.9	64.8	16.0	82.8	83.4	59.8	43.2	55.3	53.0	85.4	
	B/5	87.6	96.3	82.4	75.3	71.6	64.7	64.8	87.8	79.6	<b>95.5</b>	95.6	96.6	98.8	60.9	99.4	96.1	87.0	68.5	73.7	16.4	83.5	86.4	61.6	46.3	53.4	55.8	85.8	
	B/6	87.3	97.0	83.9	75.8	71.4	67.6	65.6	87.3	78.5	<b>95.2</b>	96.4	97.2	98.6	61.9	99.5	96.6	87.2	62.7	74.7	17.6	84.2	85.5	61.0	49.6</td				

# Self-Supervised Vision Transformers

## Masked Autoencoders Are Scalable Vision Learners



# What are we learning?

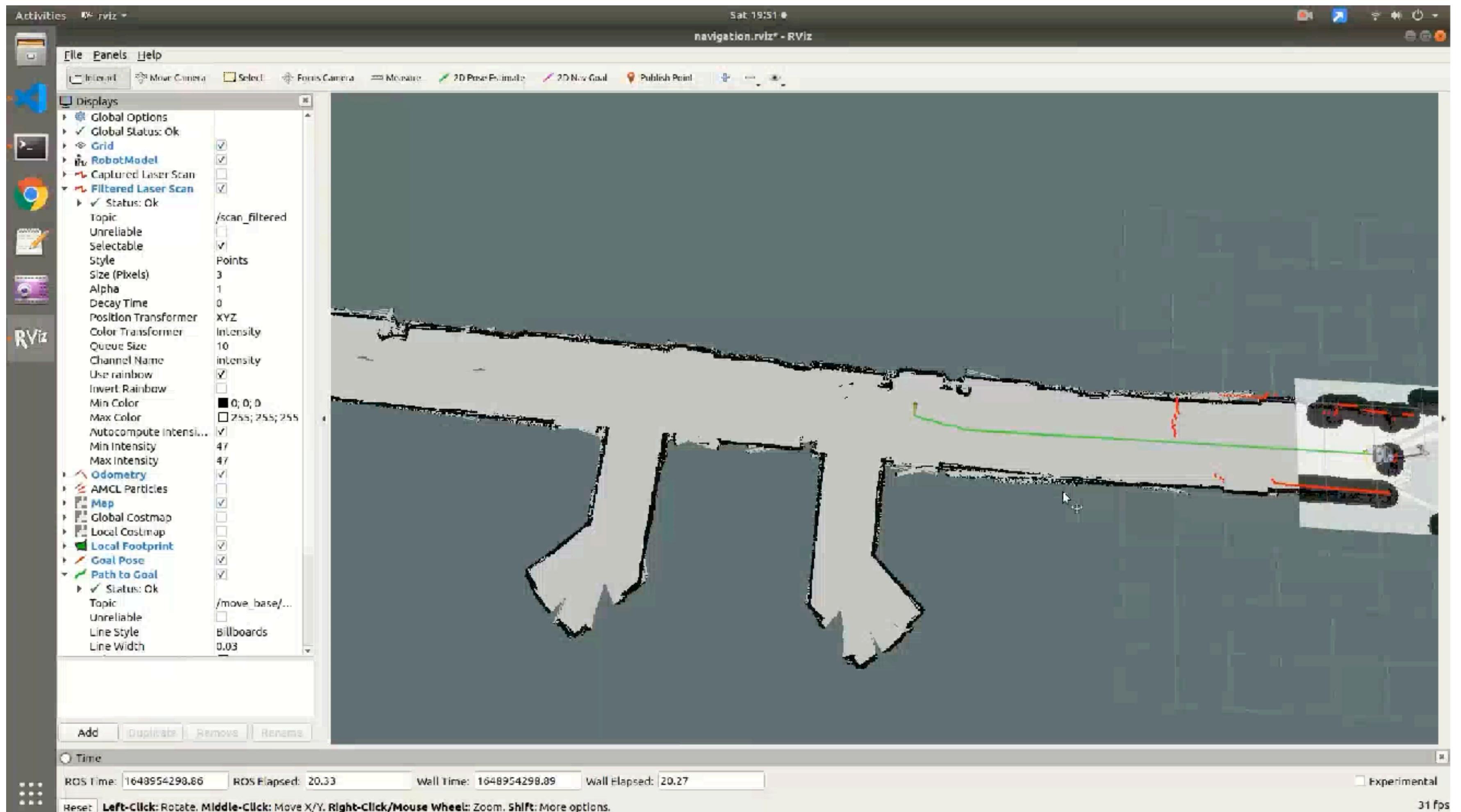


# What even is a Robot?

Are you seriously going to try and cram a robotics class into this lecture? What do you have? Like 5m left?

# Mapping

- Sensors
- Location estimation
- Map updates

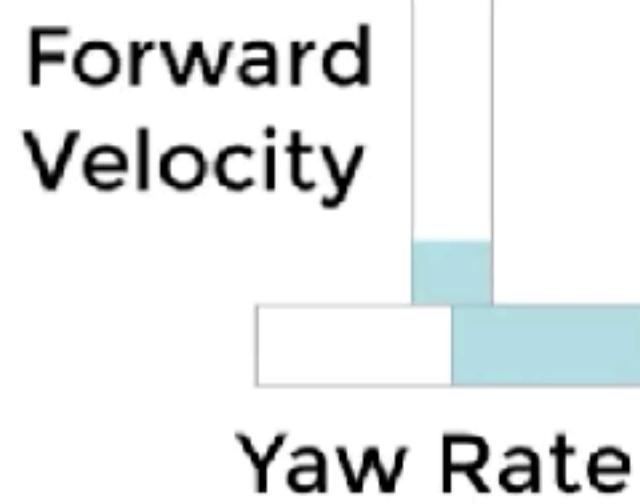


# Real Valued Predictions

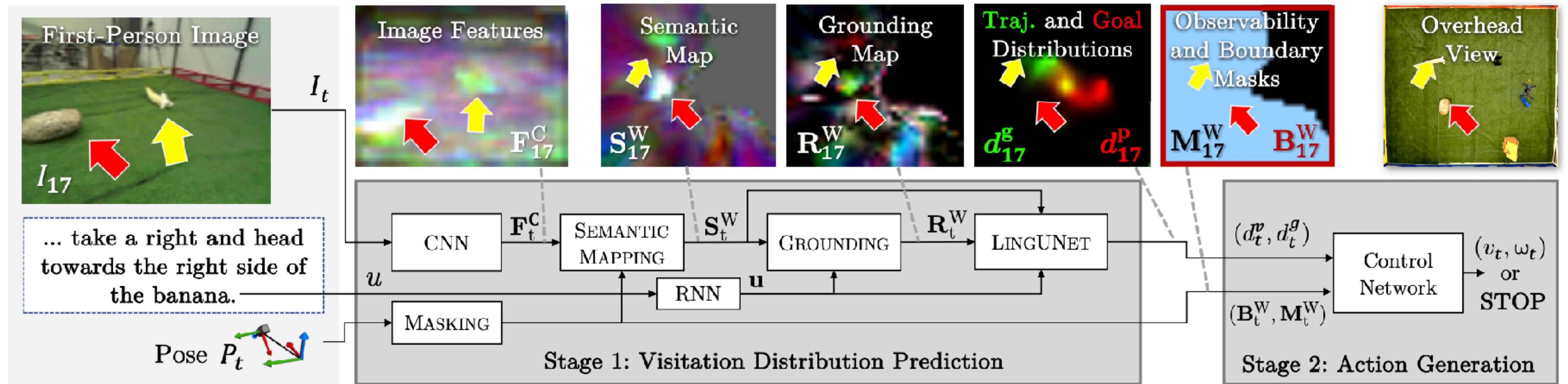
Once near the rear of the gorilla, turn right and head towards the rock stopping once near it



Continuous  
velocity  
commands



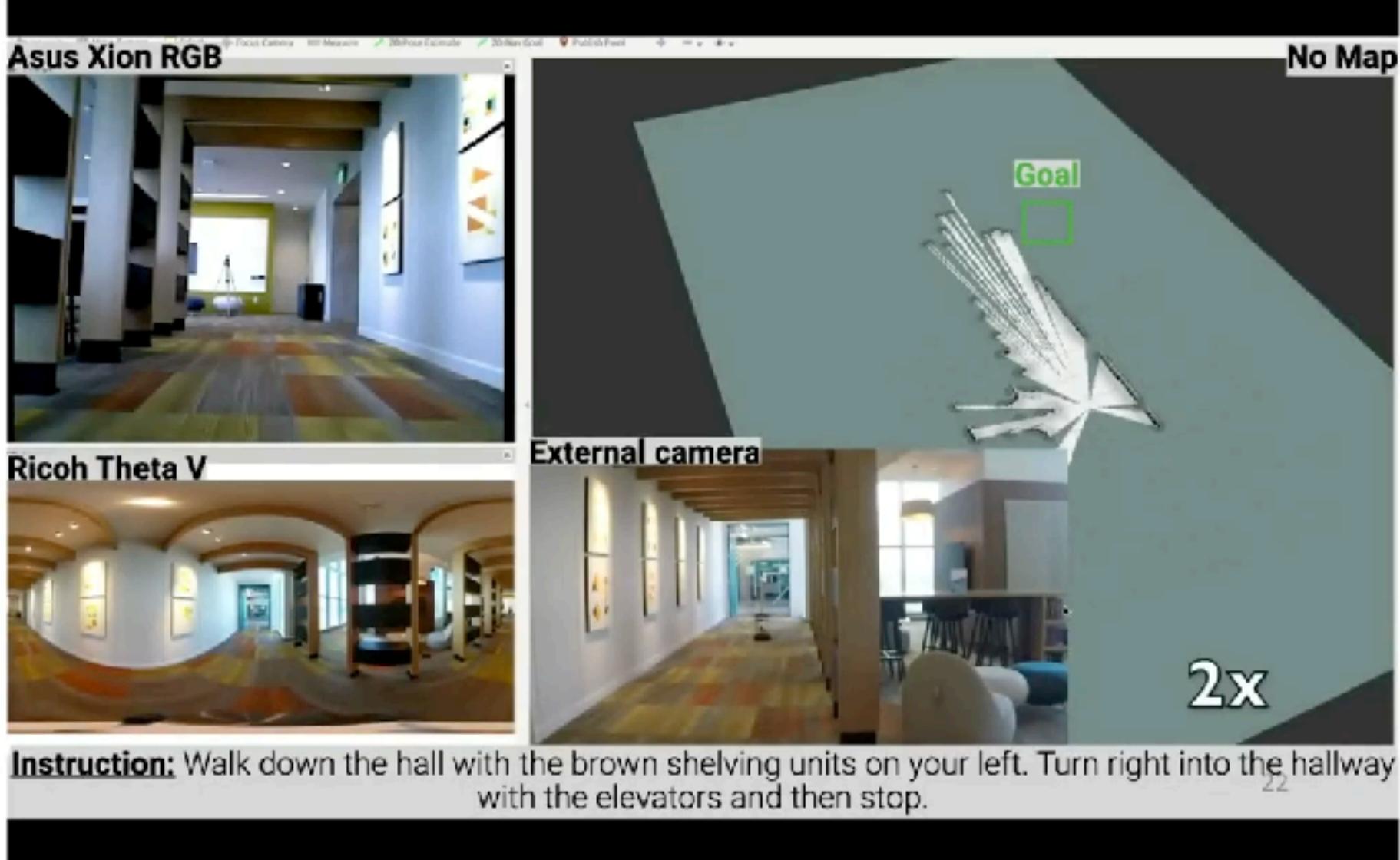
# Lots of individual models / algorithms



# Planning (algorithms)



Conference on  
Robot Learning



- Inverse Kinematics – What sequence of joint movements is required to get to a target position (e.g. gripper)
- A\*, Dijkstra – What's the shortest path from point A to B
- PDDL – Given pre- and post-conditions, what sequence of actions are required to complete a task?

