



# Deep Learning

Oct 30-31<sup>st</sup>, 2025

EIT Oxford

# Lecture II – Computer Vision

EIT Oxford

# What can we solve with computer vision

---

- **Instance classification / regression**
- **Representation learning**
- **Semantic segmentation**
- Object detection
- Generation
- Super-resolution
- Inpainting
- Image + time (= **video**)
- Object tracking
- Summarization
- ...

# What is covered in other lectures

---

- **Fundamentals of Deep Learning – Day I**
  - Optimizers
  - Regularization
  - Training stability
  - Practical considerations
- **Deep generative modelling**

# Tasks

---

**Classification**



CAT

No spatial extent

**Semantic Segmentation**



GRASS, CAT,  
TREE, SKY

No objects, just pixels

**Object Detection**



DOG, DOG, CAT

Multiple Object

**Instance Segmentation**

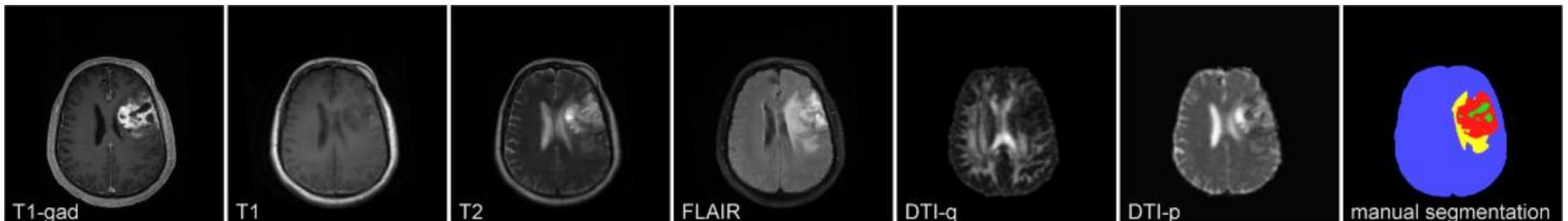
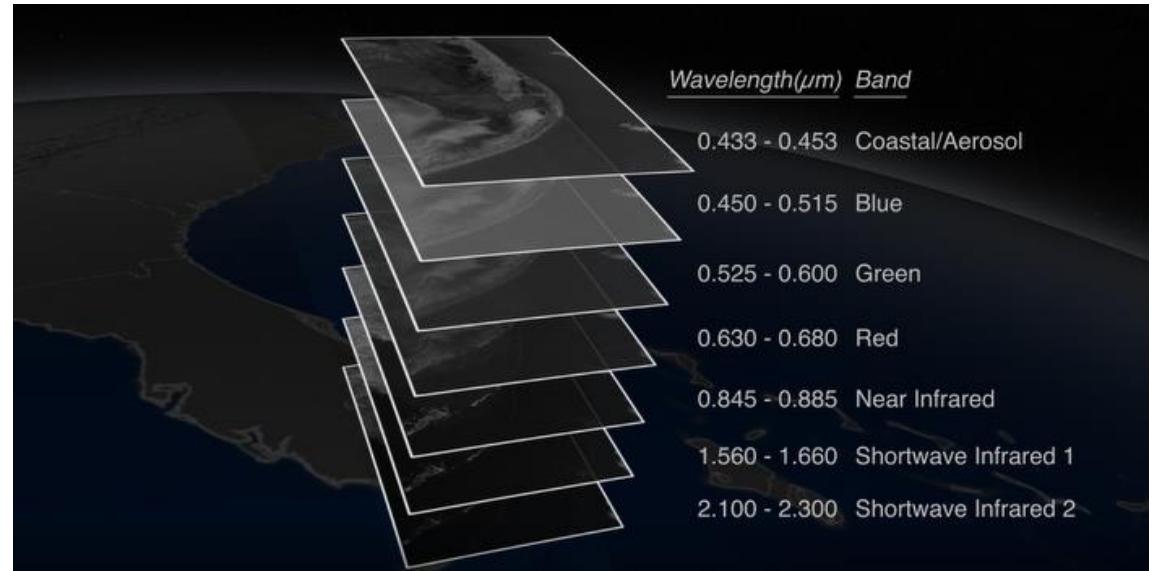
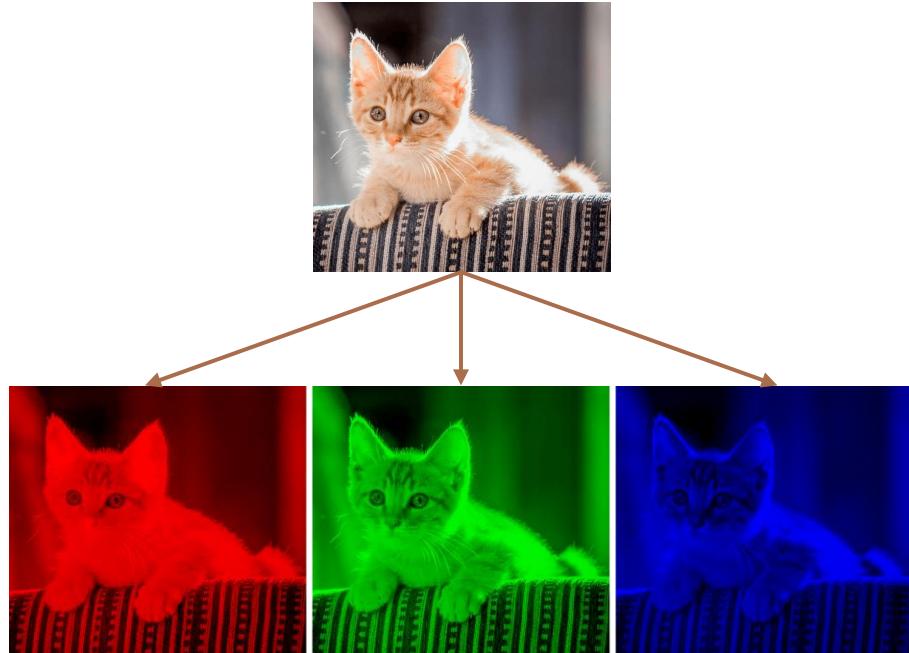


DOG, DOG, CAT

This image is CC0 public domain

Source: [http://cs231n.stanford.edu/slides/2020/lecture\\_12.pdf](http://cs231n.stanford.edu/slides/2020/lecture_12.pdf)

# Input channels vary depending on the task



# Intuitions from the visual system

106

*J. Physiol.* (1962), **160**, pp. 106-154  
With 2 plates and 20 text-figures  
Printed in Great Britain

## RECEPTIVE FIELDS, BINOCULAR INTERACTION AND FUNCTIONAL ARCHITECTURE IN THE CAT'S VISUAL CORTEX

BY D. H. HUBEL AND T. N. WIESEL

*From the Neurophysiology Laboratory, Department of Pharmacology  
Harvard Medical School, Boston, Massachusetts, U.S.A.*

(Received 31 July 1961)

What chiefly distinguishes cerebral cortex from other parts of the central nervous system is the great diversity of its cell types and inter-connexions. It would be astonishing if such a structure did not profoundly modify the response patterns of fibres coming into it. In the cat's visual cortex, the receptive field arrangements of single cells suggest that there is a degree of complexity far exceeding anything yet seen at lower

# Neural networks are inspired by the visual system

---

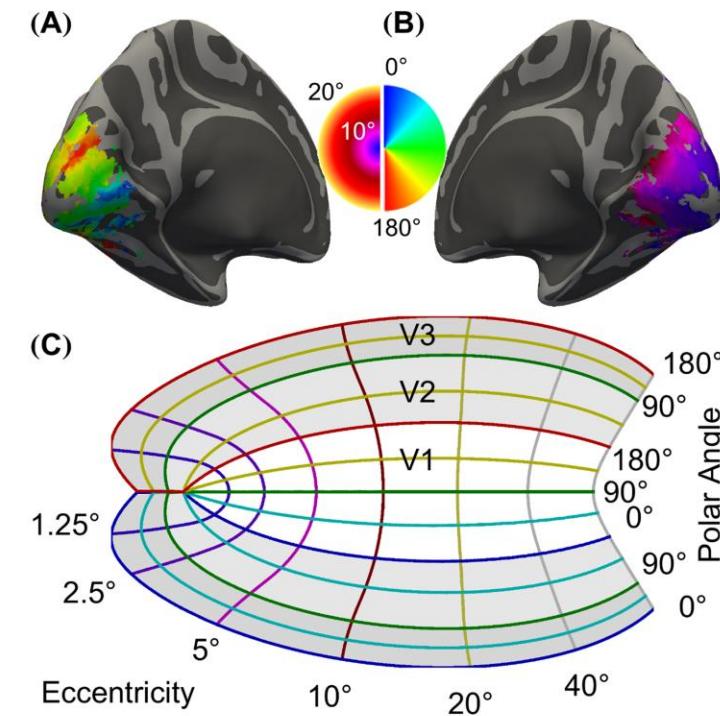
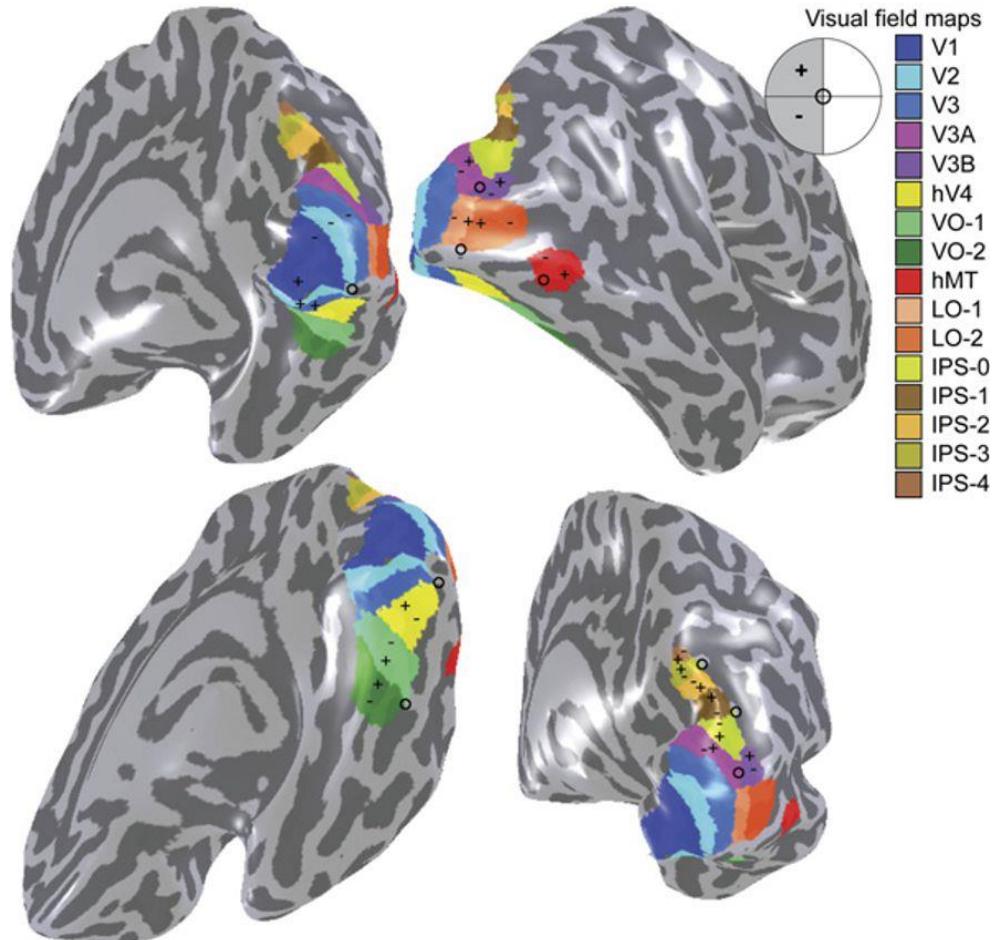


# Neural networks are inspired by the visual system

---

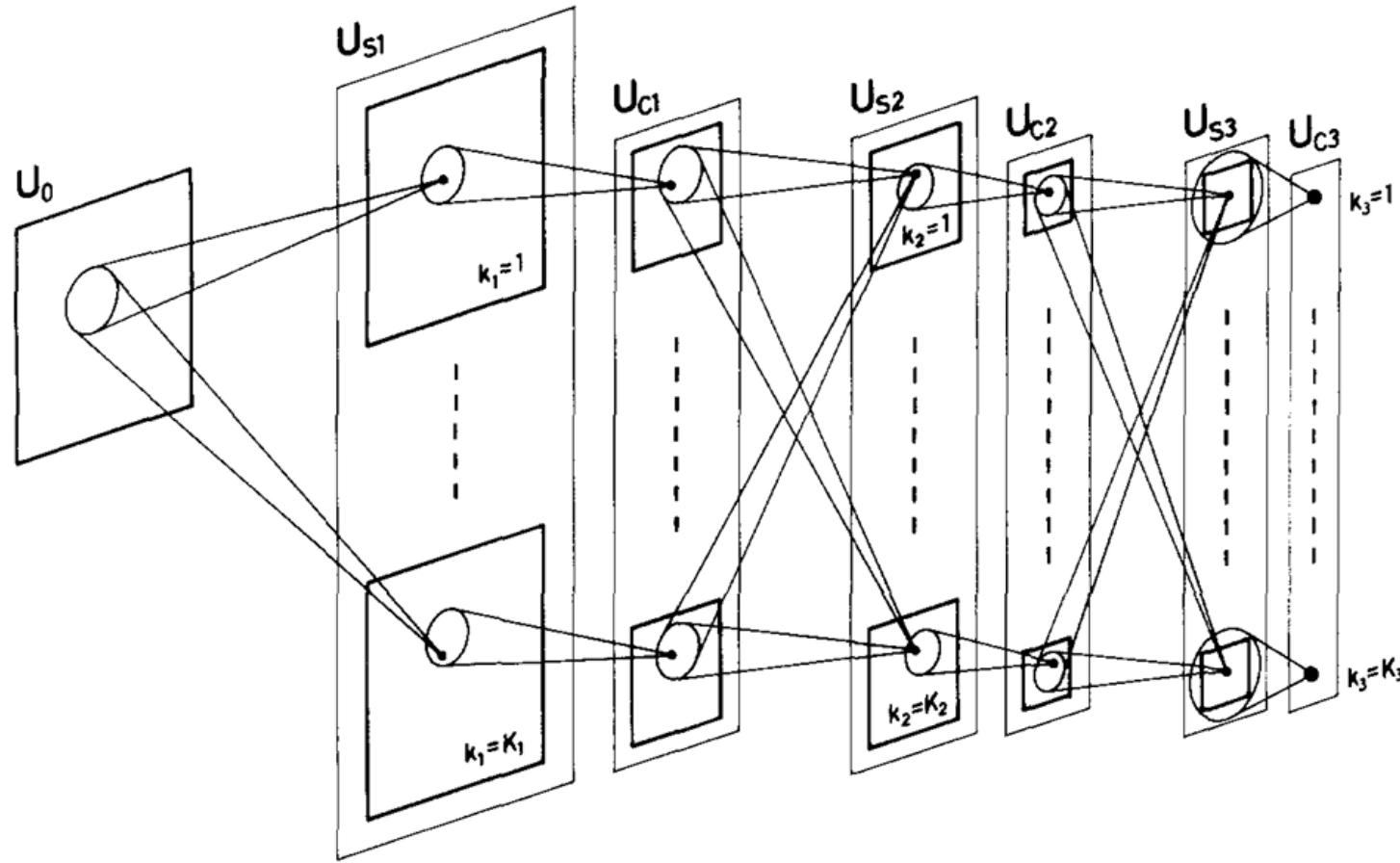
- Measured the electrical response in a cat's brain by stimulating it with simple patterns on a screen.
- Found that neurons in the early visual cortex are organized hierarchically:
  - Early layers include orientation-selective cells responding preferentially to simple patterns like edges and contours
  - Later layers respond to more complex patterns by combining the activations they receive.
- In the proposed model, neurons in the upper layers have a larger receptive field and are less sensitive to the position where the stimulus originates from.

# Organization of the human visual cortex



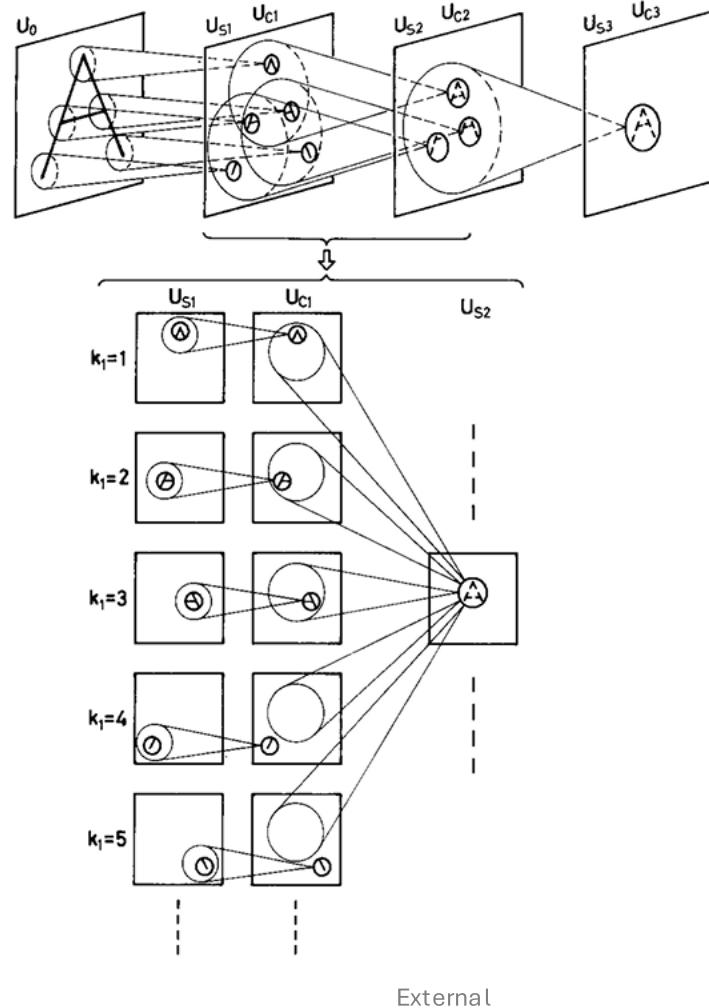
# Neocognitron – Fukushima (1980)

Schematic diagram illustrating the interconnections between layers in the neocognitron



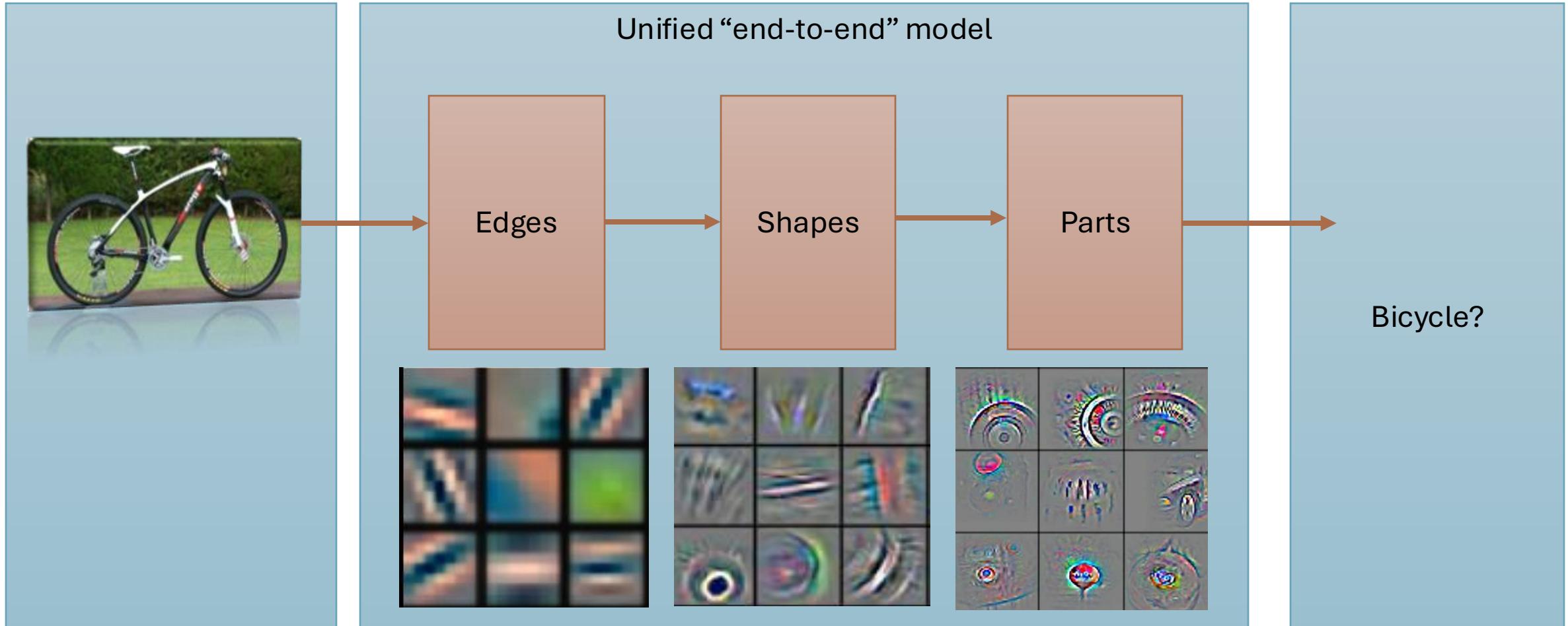
# Neocognitron – Fukushima (1980)

Interconnections between cells and the response of the cells after completion of self-organisation



# Multiple levels of abstraction

---



# Building blocks

---

- Invariance and equivariance
- Convolutions
- Skip / residual connections
- Pooling layers
- Receptive fields
- Upsampling layers
- 1x1 convolution

## Invariance

---

A function  $f[x]$  is **invariant** to a transformation  $t[]$  if:

$$f[t[x]] = f[x]$$

i.e., the function output is the same even after the transformation is applied.

# Example

---

e.g., Image classification

Image has been translated, but we want our classifier to give the same result



# Equivariance

---

A function  $f[x]$  is **equivariant** to a transformation  $t[]$  if:

$$f[t[x]] = f[x]$$

i.e., the output is transformed in the same way as the input

# Example

---

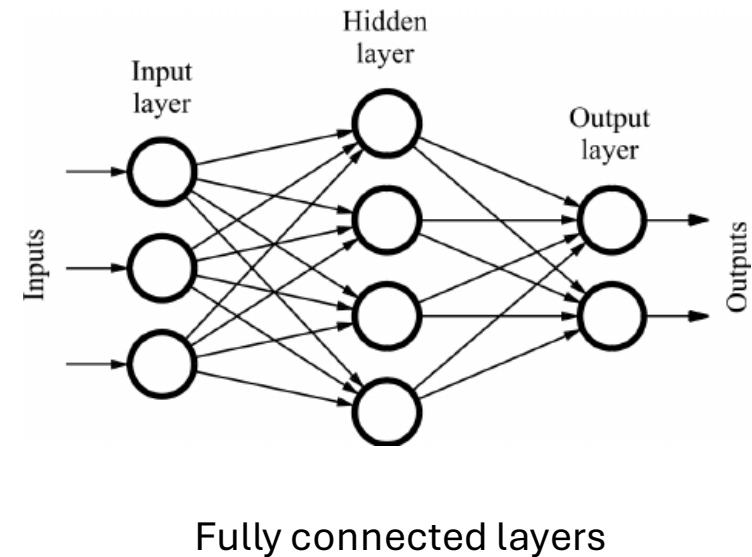
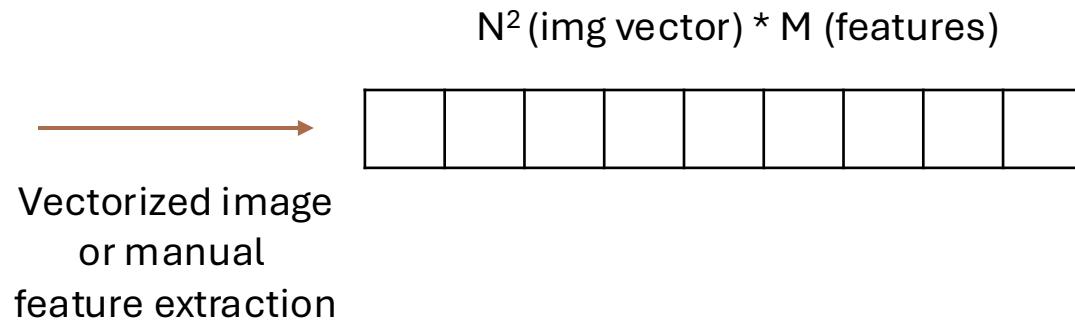
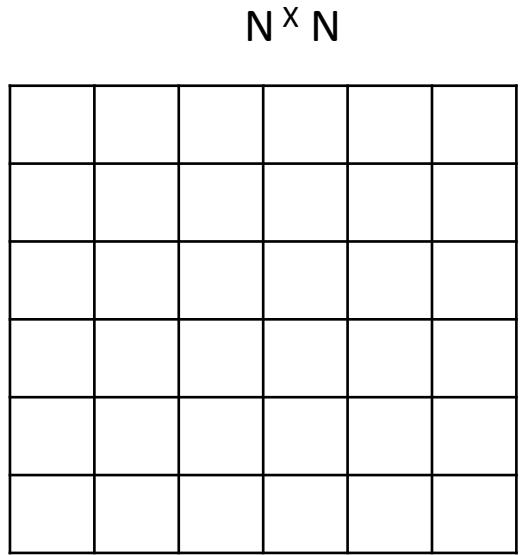
e.g., Image segmentation

Image has been translated and we want segmentation to translate with it

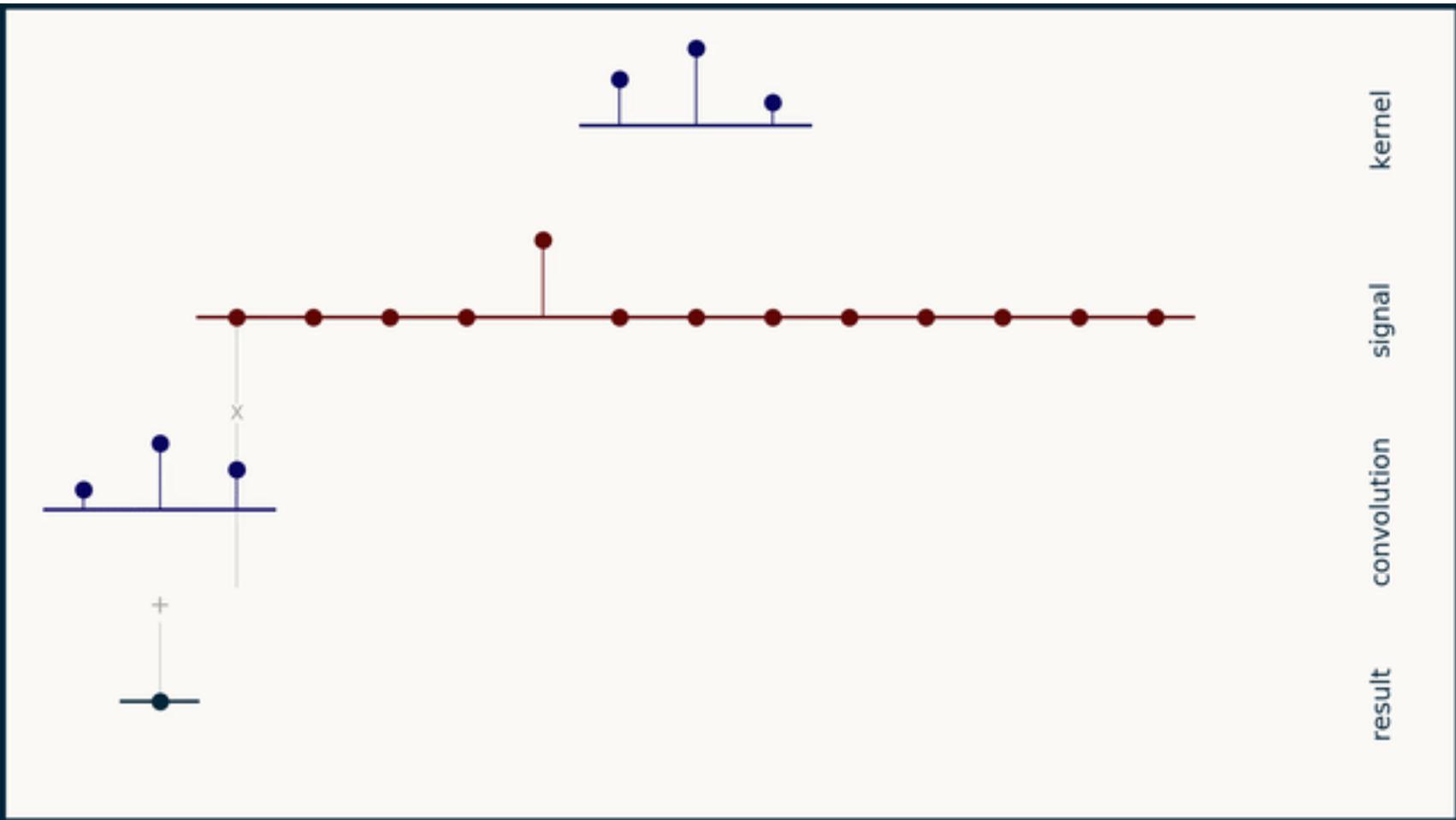


Exte.....

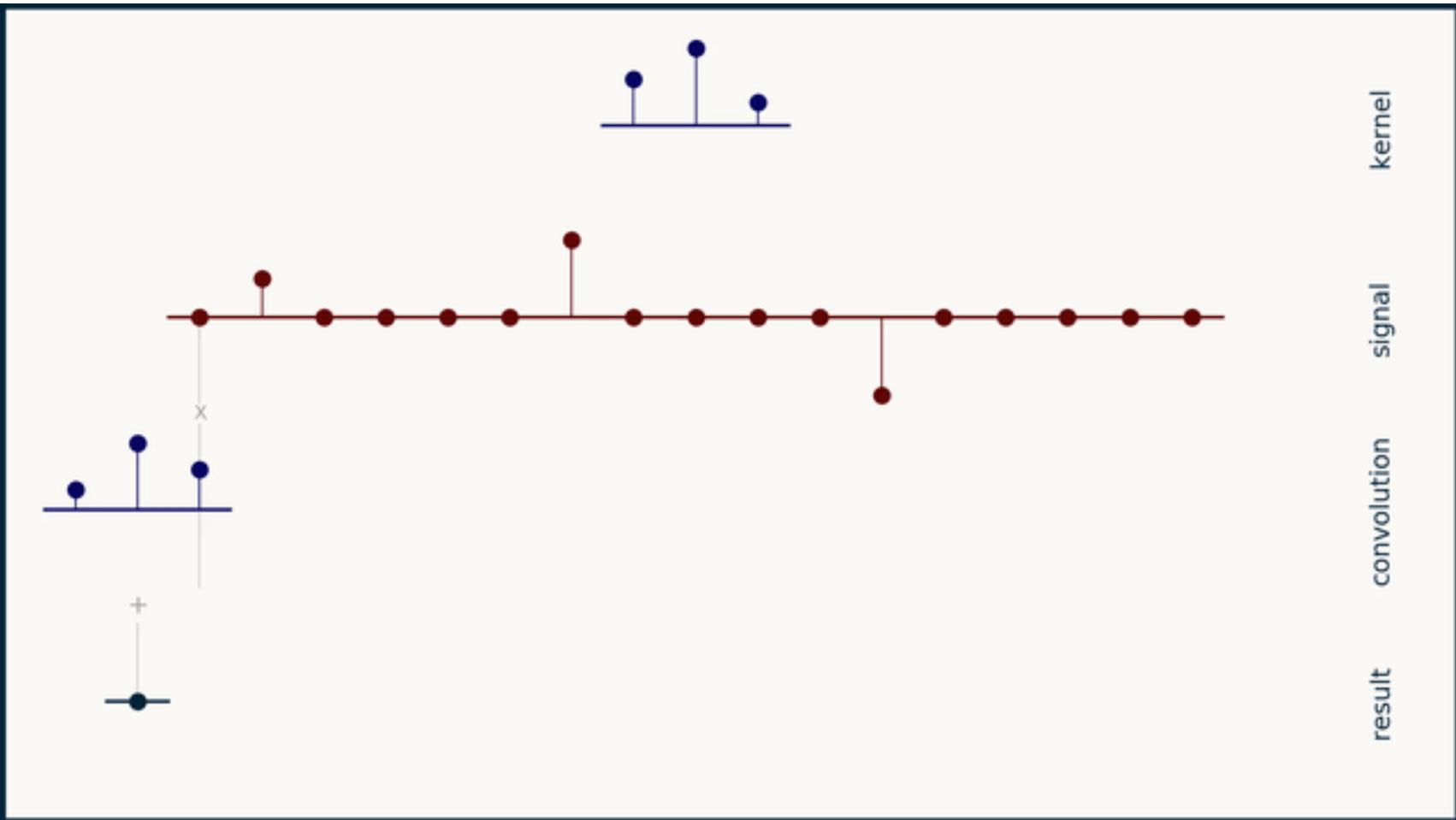
# How do we process an image with a FCN?



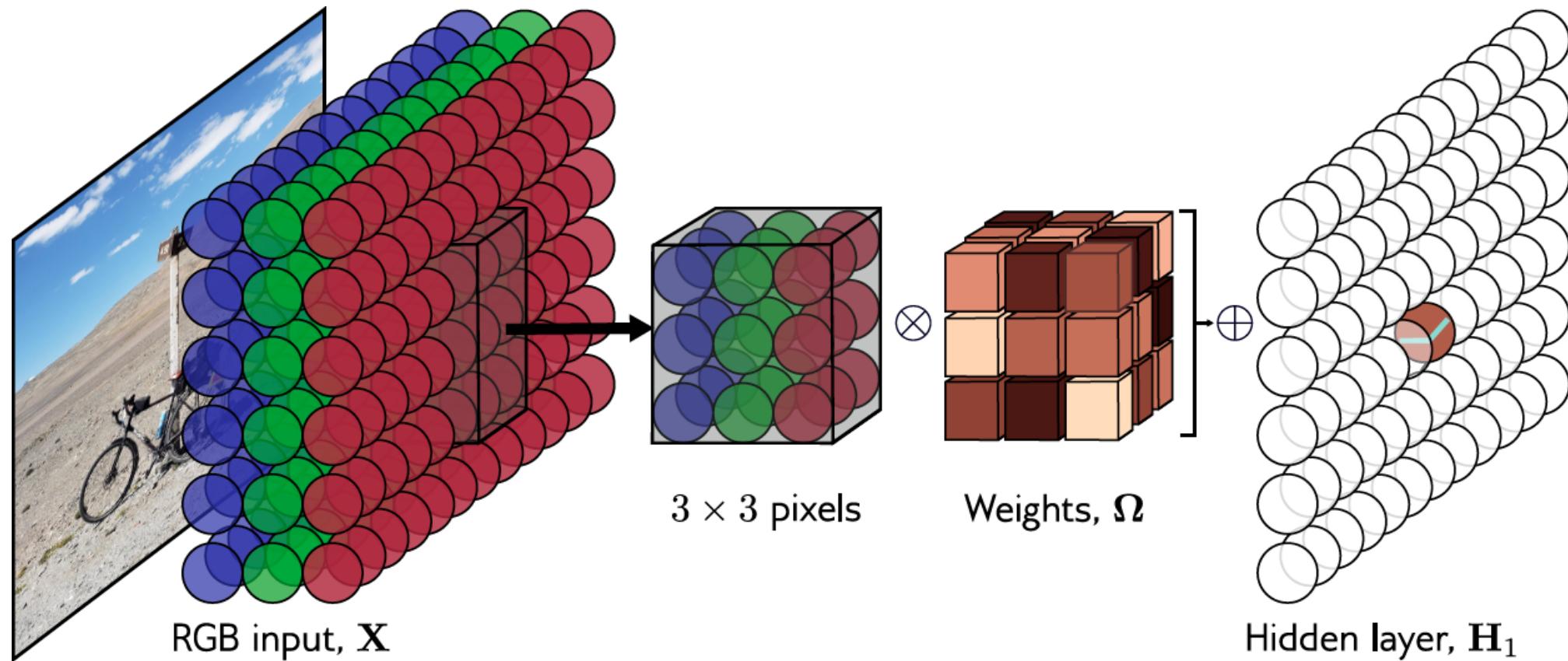
# Convolutions in 1D



# Convolutions in 1D



# CNNs



# Stride, kernel size, and dilation

---

- **Stride** = shift by k positions for each output
  - Decreases size of output relative to input
- **Kernel size** = weight a different number of inputs for each output
  - Combine information from a larger area
  - But kernel size 5 uses 5 parameters
- **Dilated convolutions** = intersperse kernel values with zeros
  - Combine information from a larger area
  - Fewer parameters

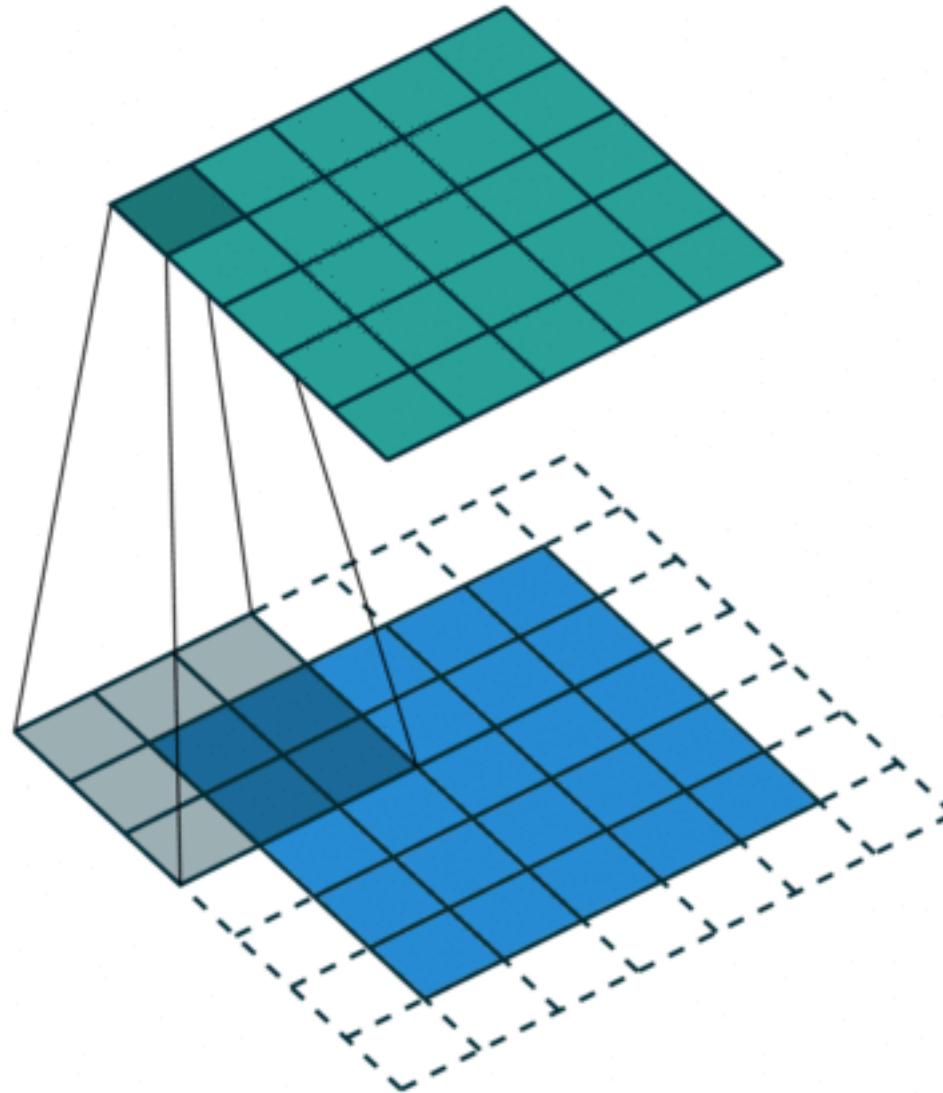
# Stride, kernel size, and dilation

---

Stride = 1

Kernel size = 3

Dilation = 1



## 2D convolution

---

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n)$$

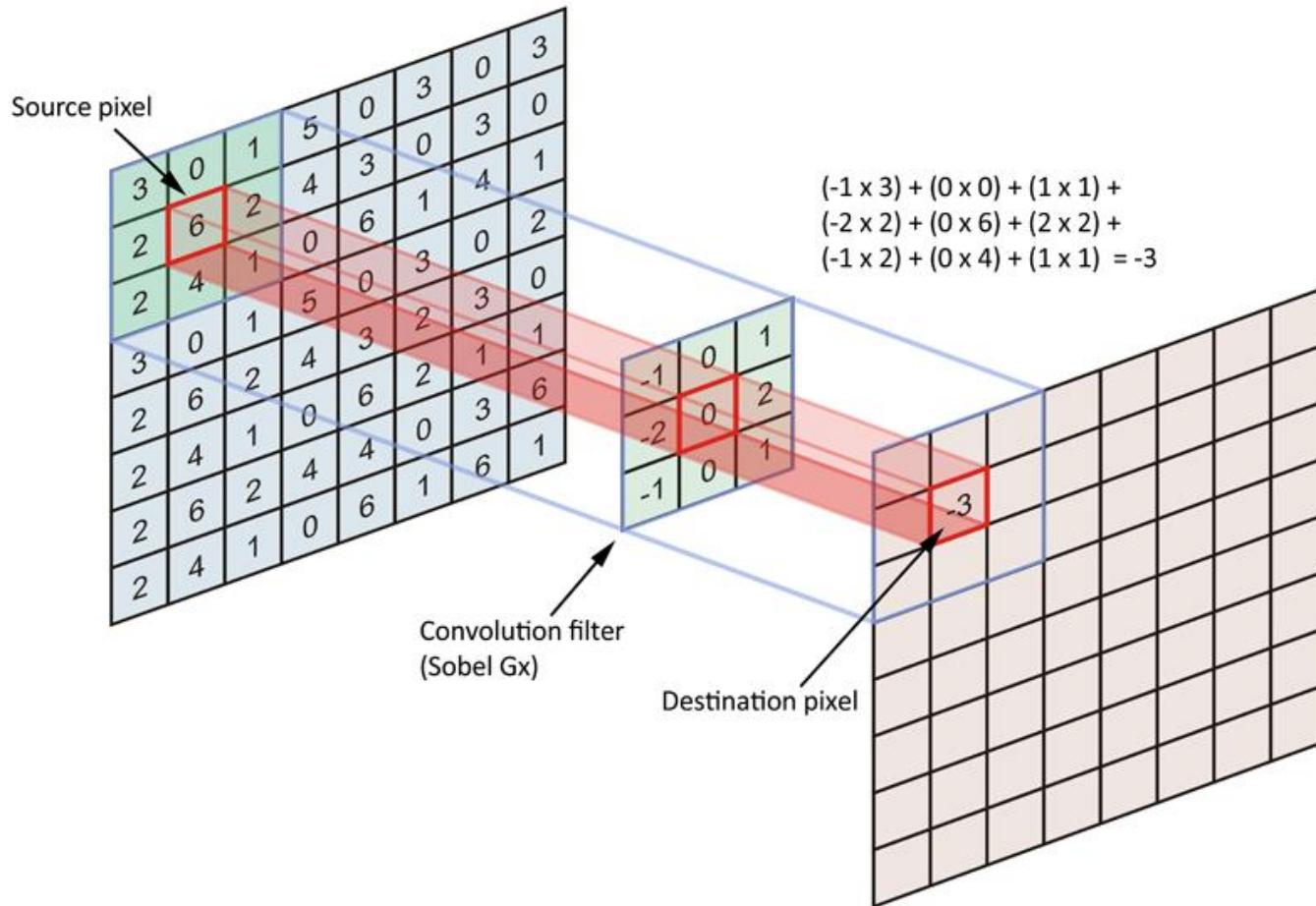


Convolution is commutative

$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$

Useful property since the range of valid values of  $m, n$  for an image is generally larger than that of the kernel.

# 2D convolution



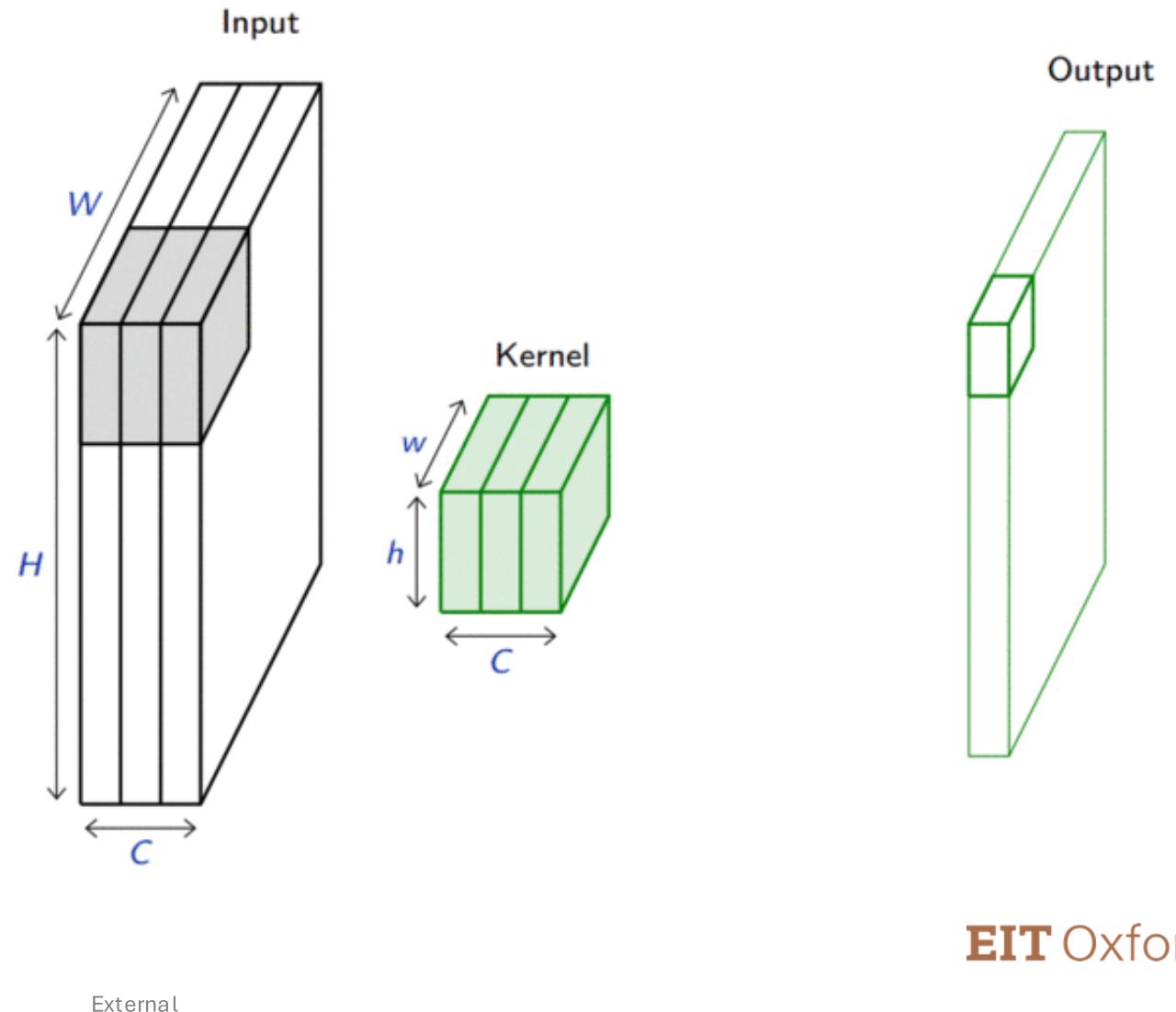
# Channels

---

- The convolutional operation averages together the inputs
- Plus passes through an activation function
- Loses information as a result (moving forward)
- Solution:
  - apply several convolutions and stack them in **channels**
  - Sometimes also called **feature maps**

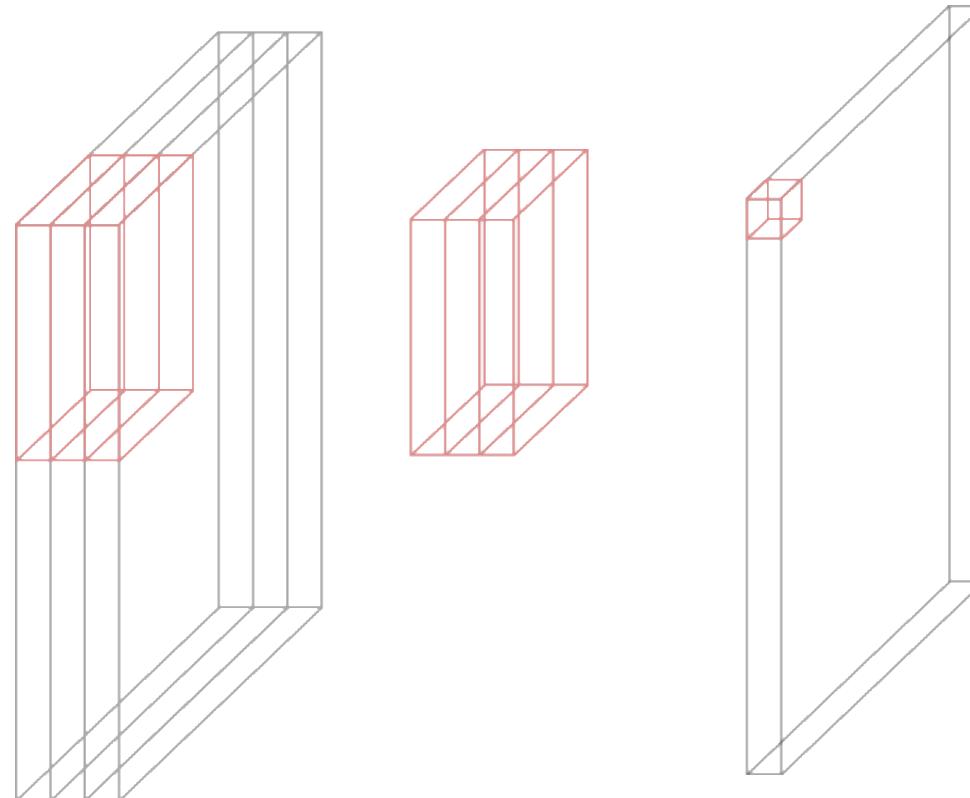
# 2D convolution layer

- Local spatial connectivity
- Weight sharing reduces the number of parameters and introduces translation invariance
- Non-linear activation functions are applied to each element individually



# Output of a convolution layer

---



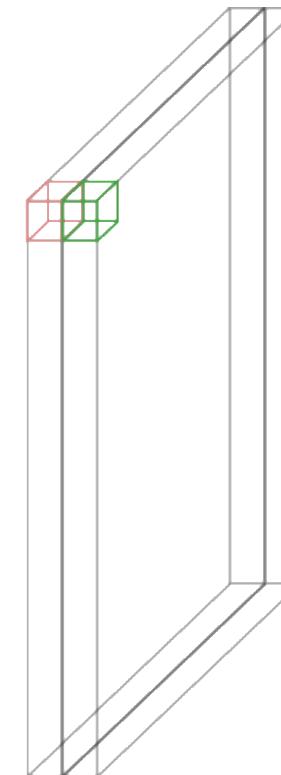
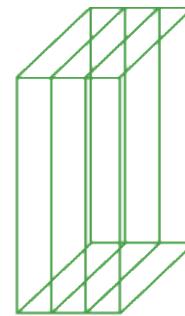
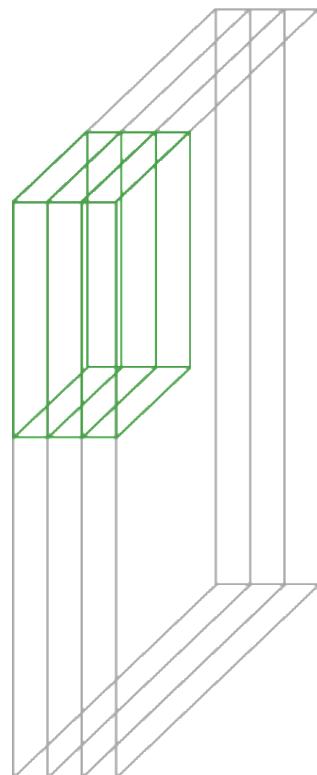
Learnable parameters



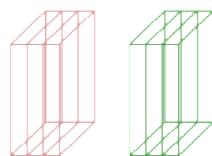
External

# Output of a convolution layer

---



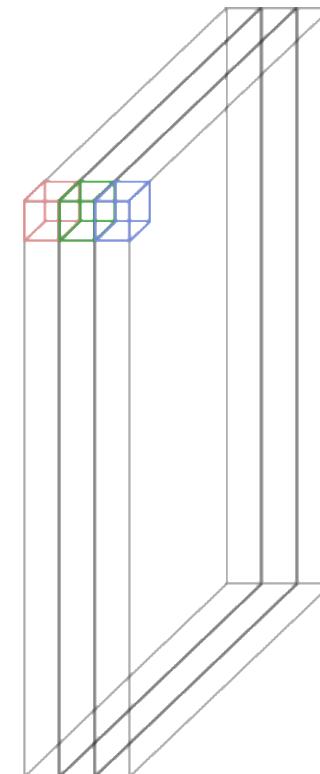
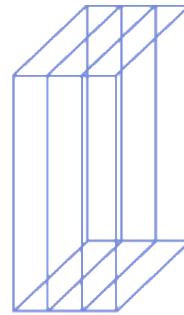
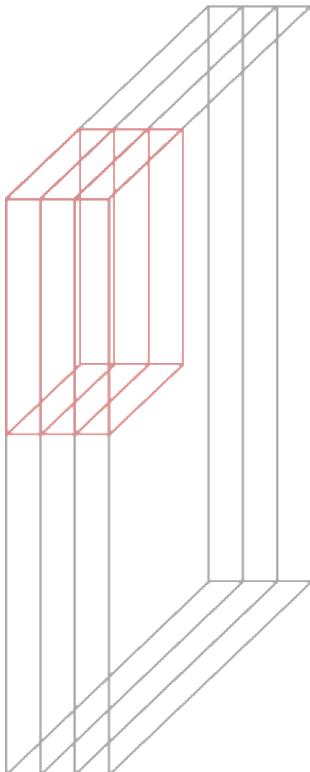
Learnable parameters



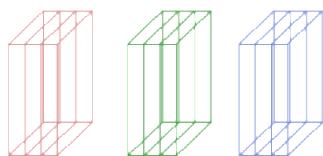
External

# Output of a convolution layer

---



Learnable parameters



External

# Examples of filters

---



$$* \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} =$$



# Examples of filters

---

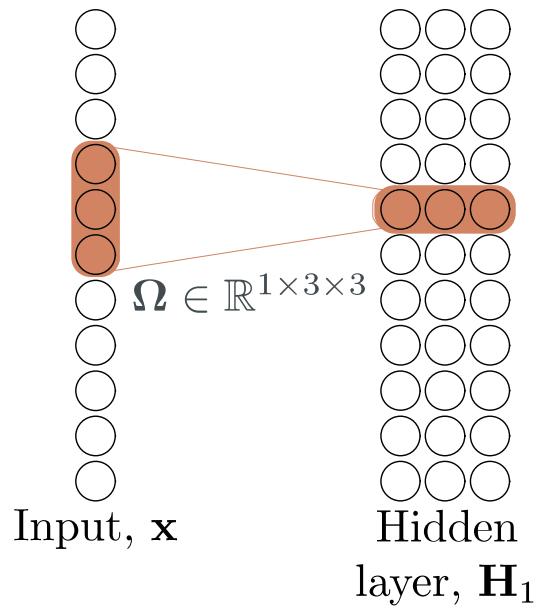


$$* \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} =$$



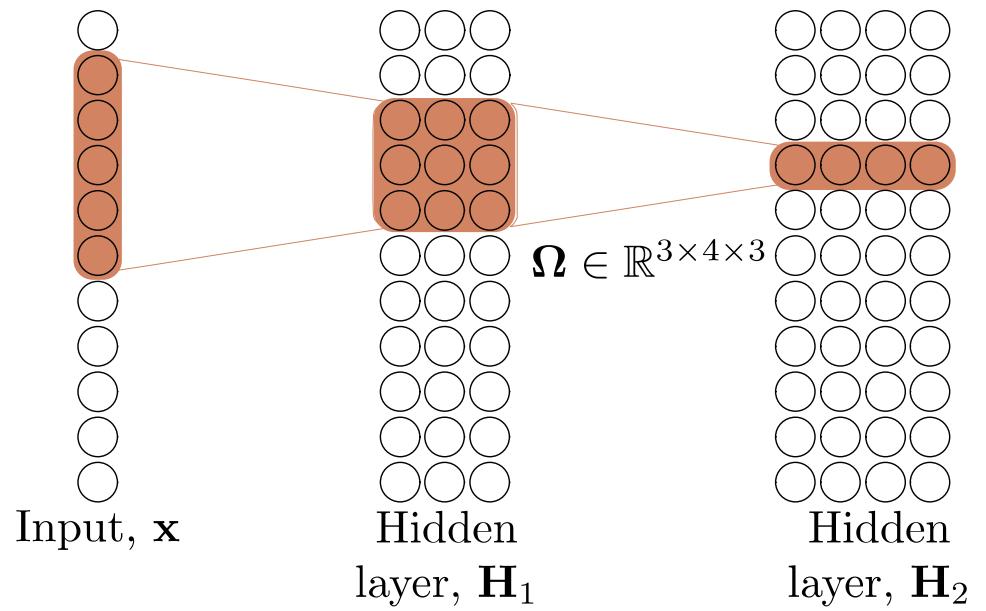
# Receptive fields

---



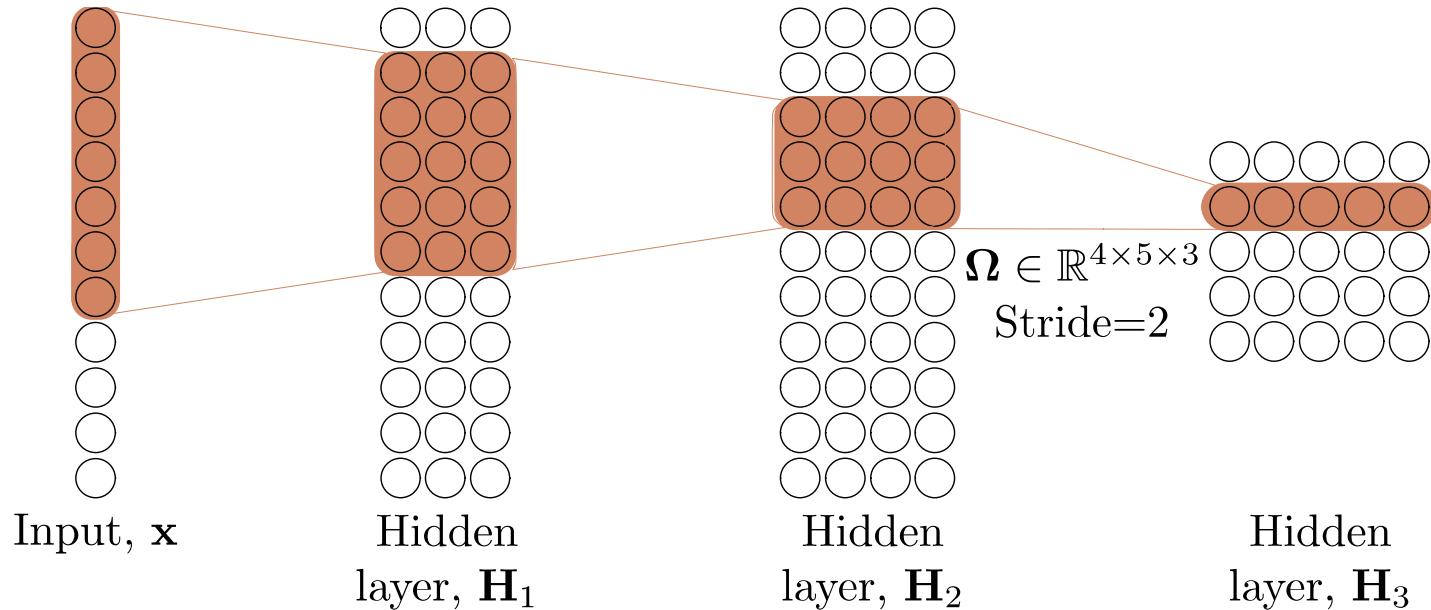
# Receptive fields

---



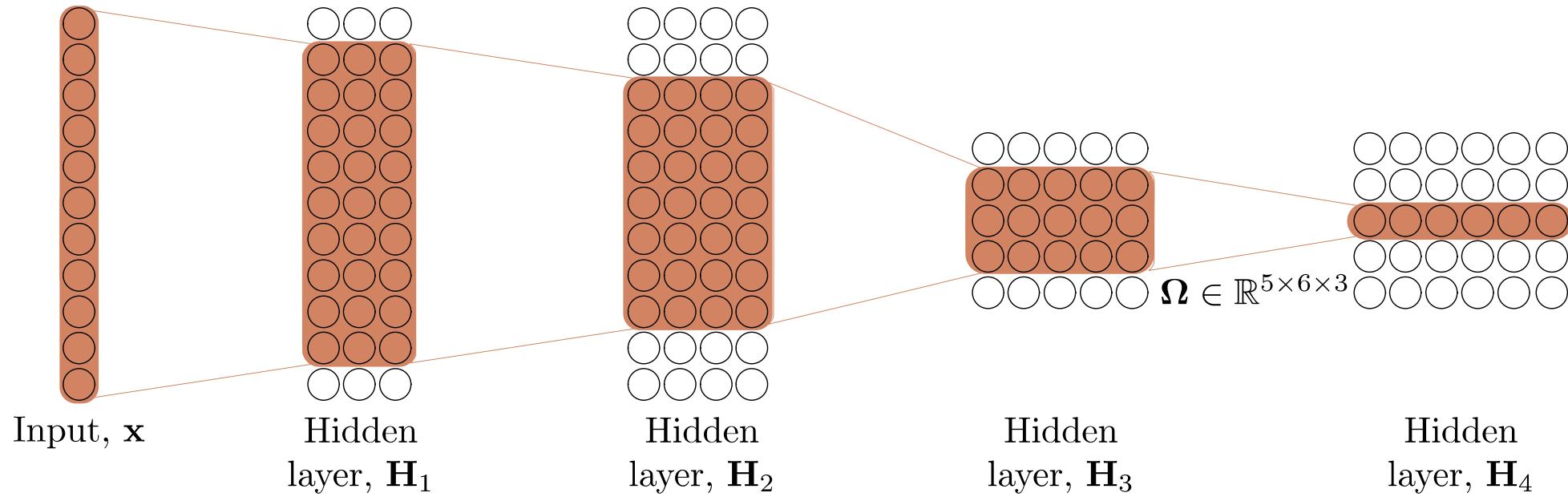
# Receptive fields

---

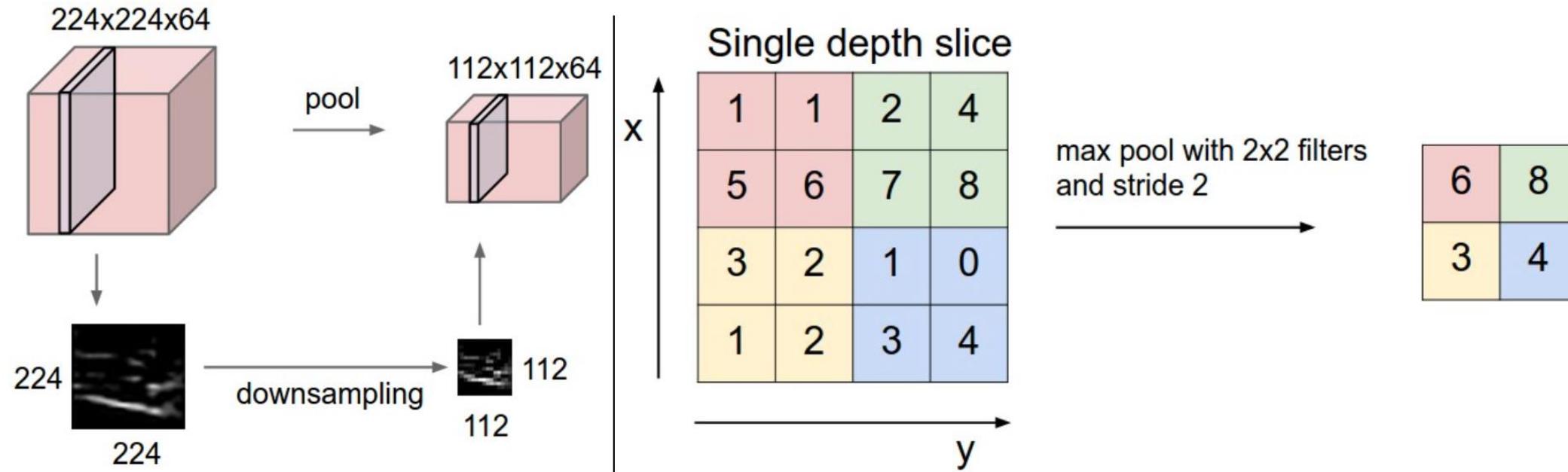


# Receptive fields

---



# Pooling layers



- Reduce dimensions of feature maps
- Contribute to invariance with respect to small translations in the input images

# Max-pooling layer

---

0	2	2	1
4	3	0	1
2	1	3	2
3	1	2	0

**2x2 Max Pooling**



4	2
3	3

# Average-pooling layer

---

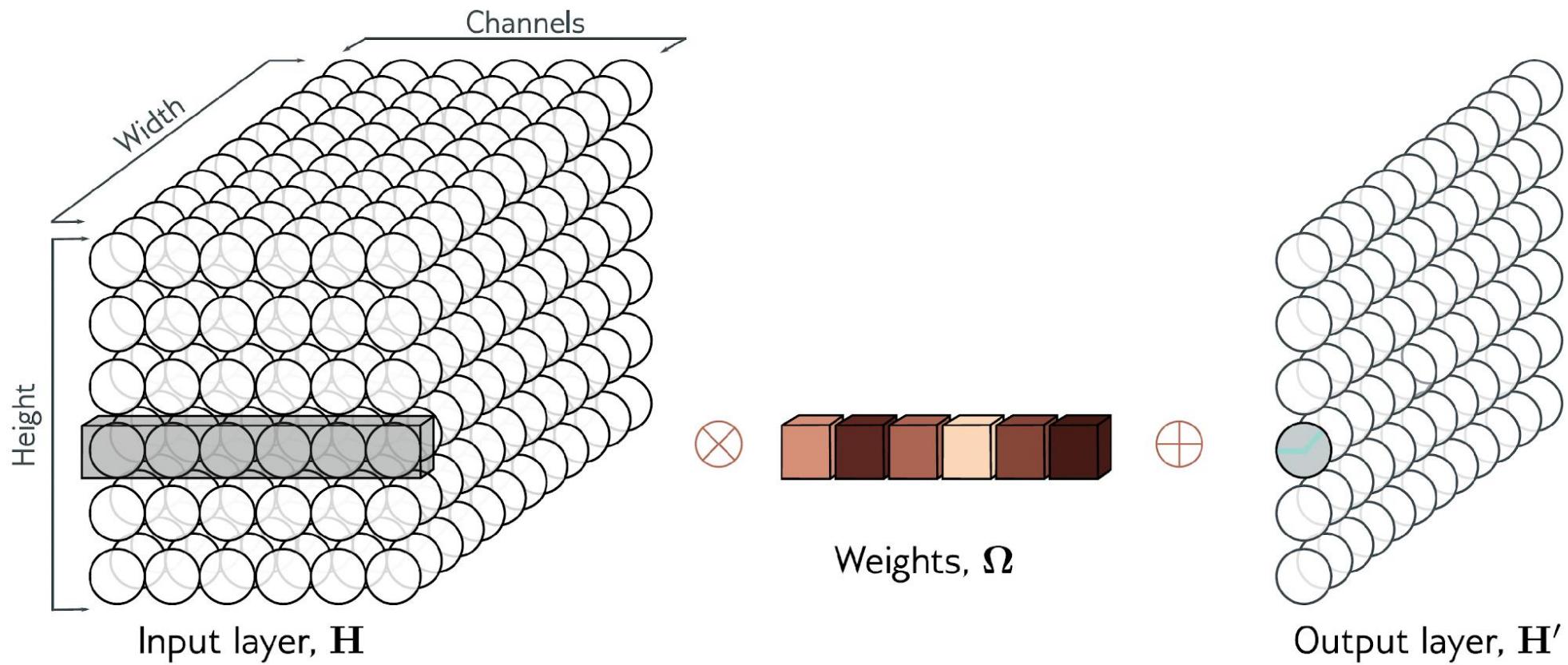
0	2	2	1
4	3	0	1
2	1	3	2
3	1	2	0

2x2 Avg Pooling



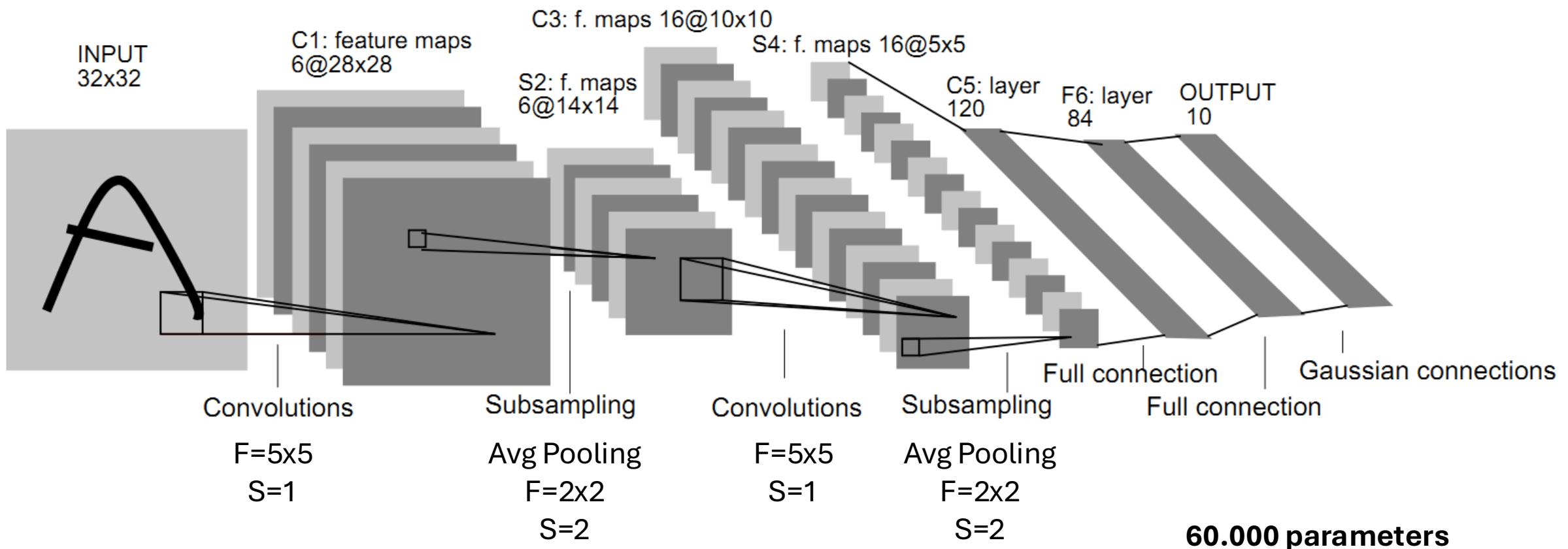
9/4	4/4
7/4	7/4

# 1x1 convolution

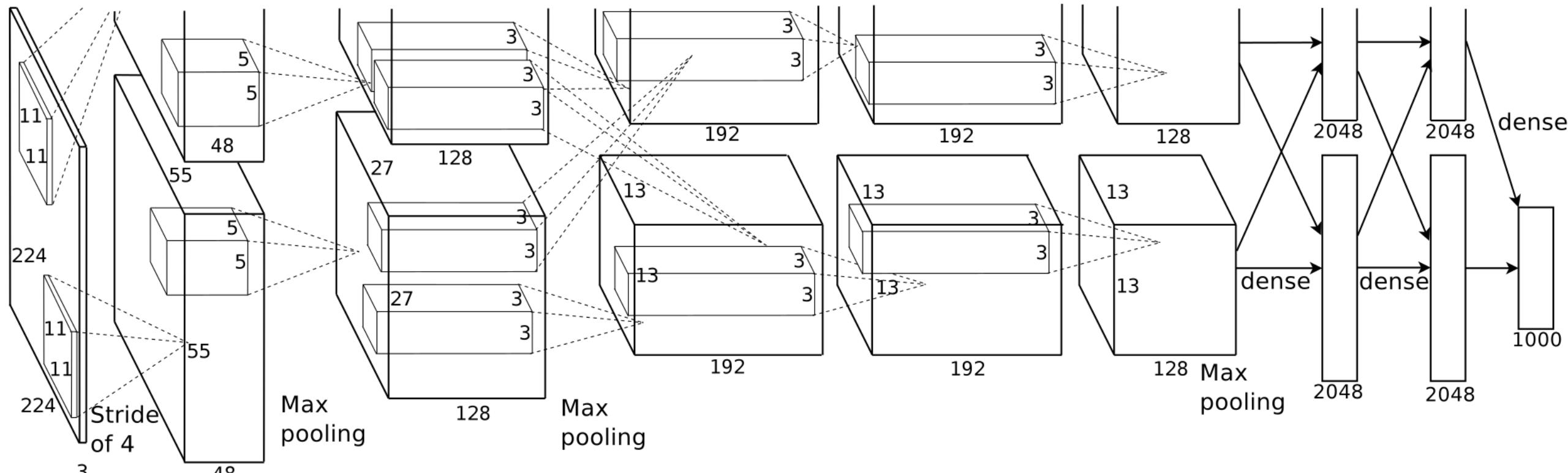


# LeNet-5 (1998)

$$O = \frac{(W - F + 2P)}{S} + 1$$



# AlexNet (2012)

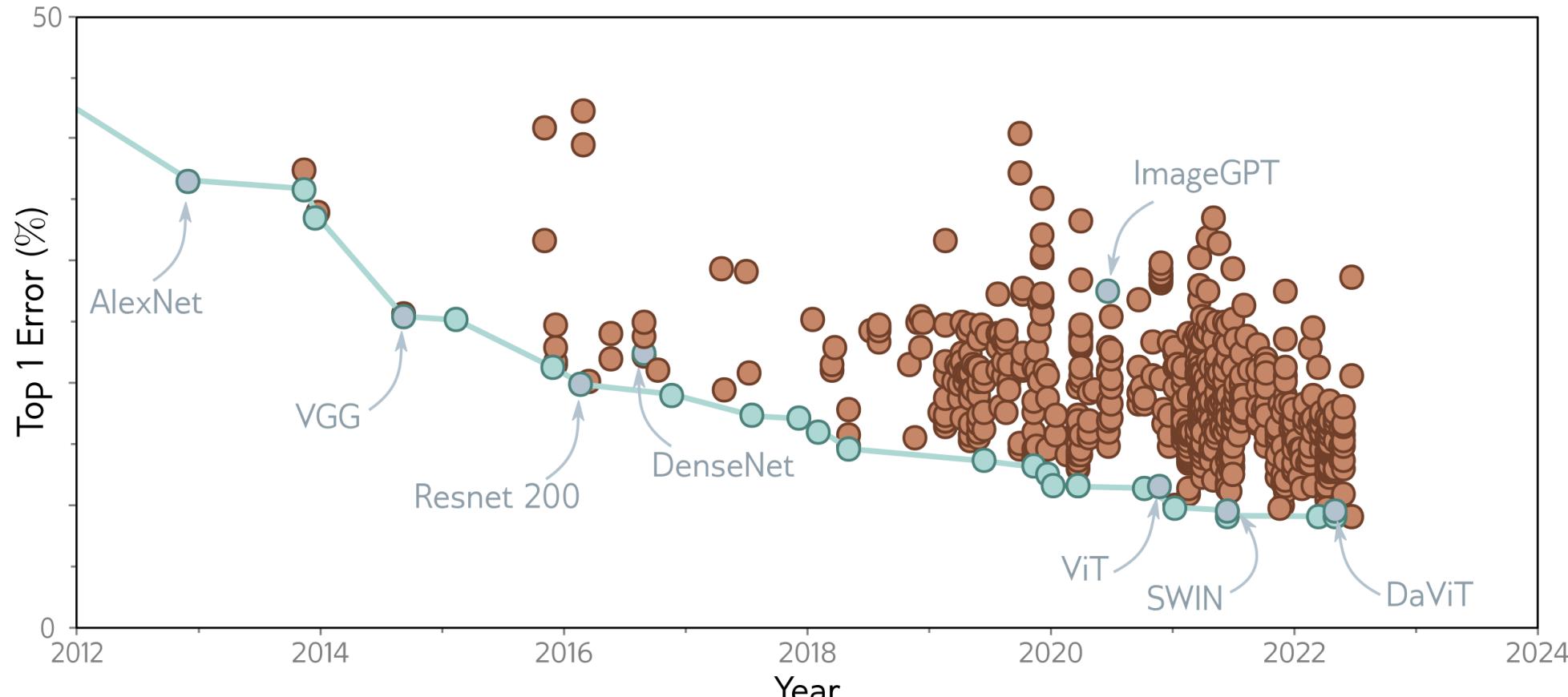


How many parameters does the first convolutional layer have?  $\rightarrow 11 \times 11 \times 3 \times 96 = 35.000$

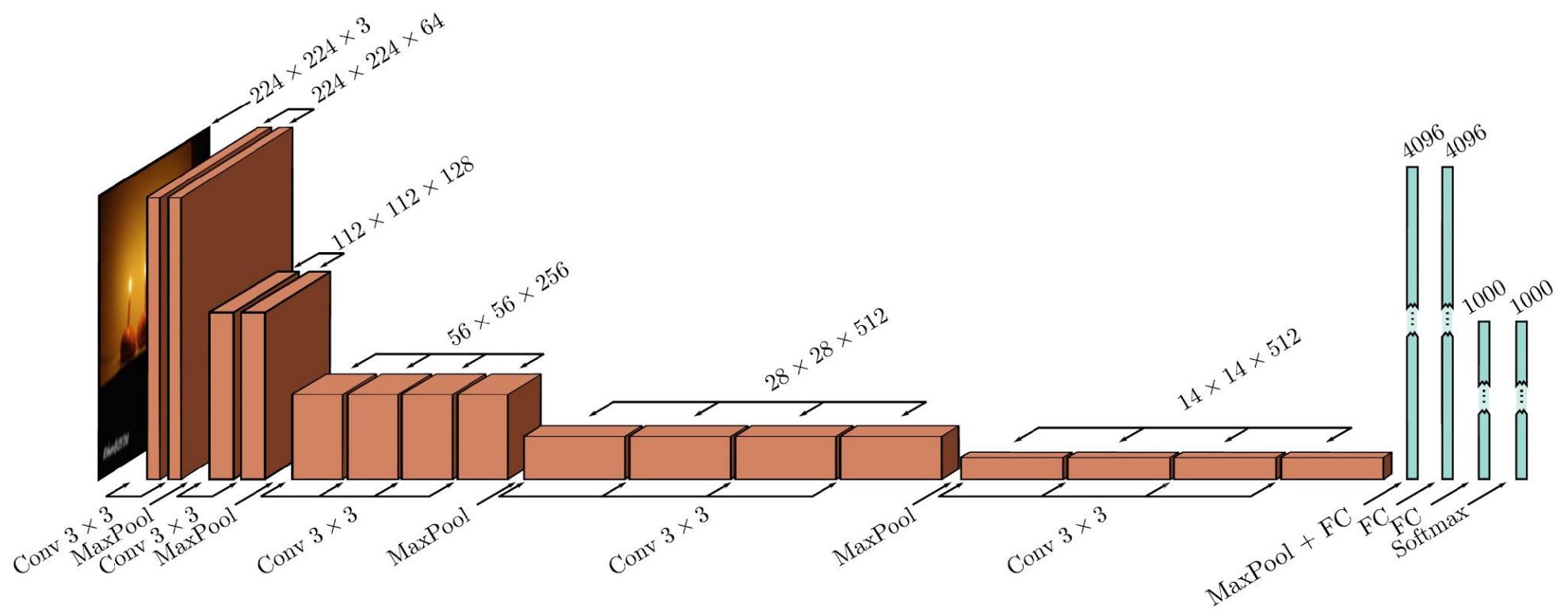
From the 60 million parameters, 58 million correspond to the dense layers

# Classical architectures for classification

ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

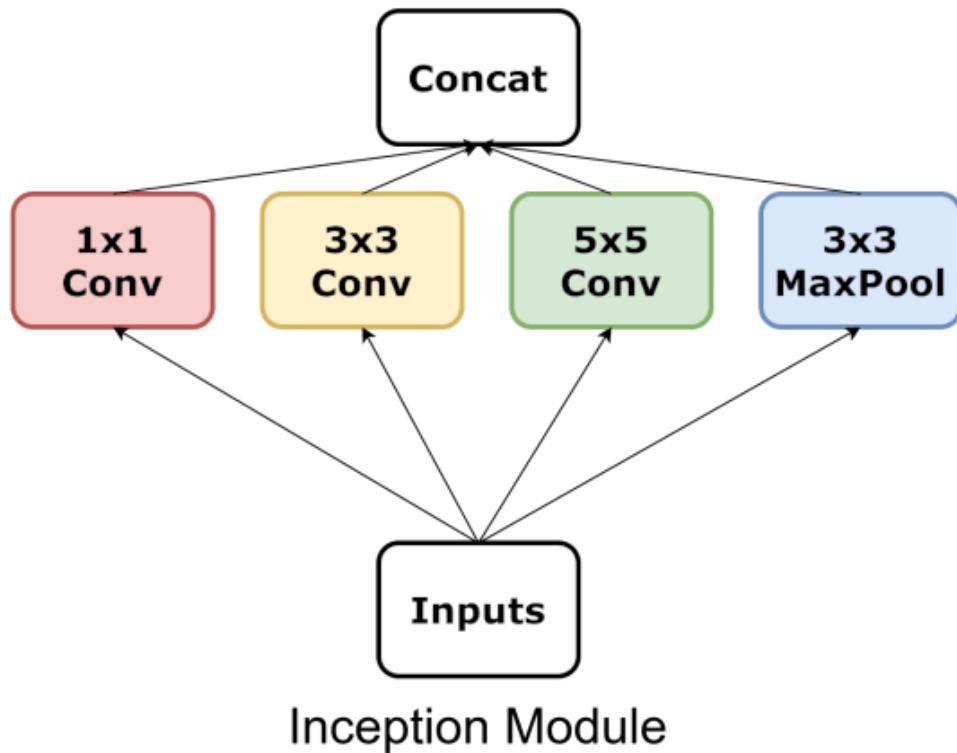


# VGG (2015)

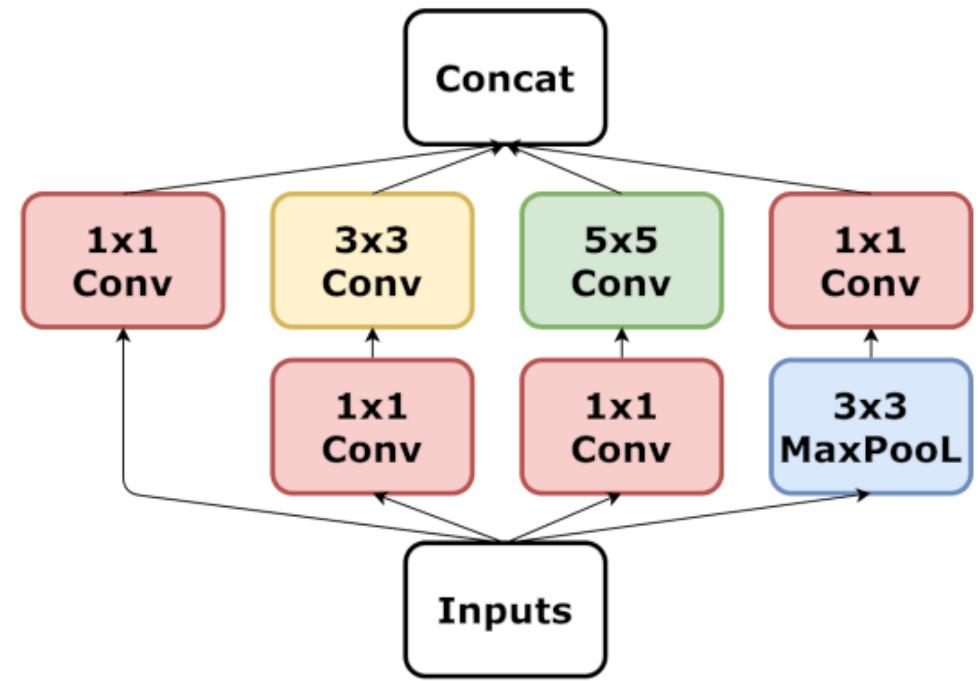


# GoogLeNet (2014)

Inception module – originally suffered from vanishing gradient problem



Inception Module

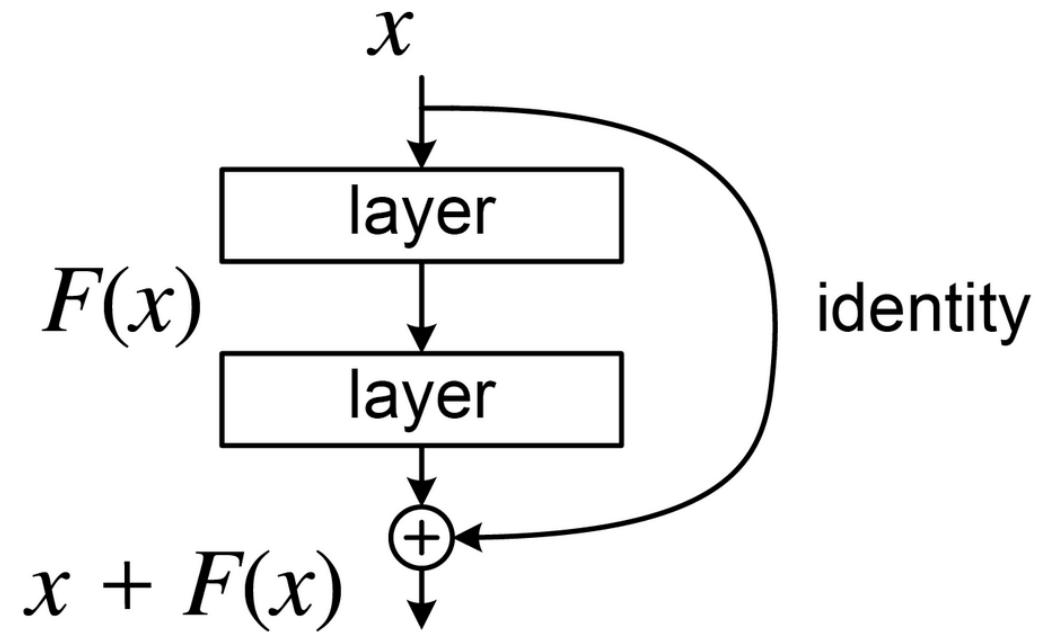


Inception Module with Dimension Reduction

# Skip / residual connections

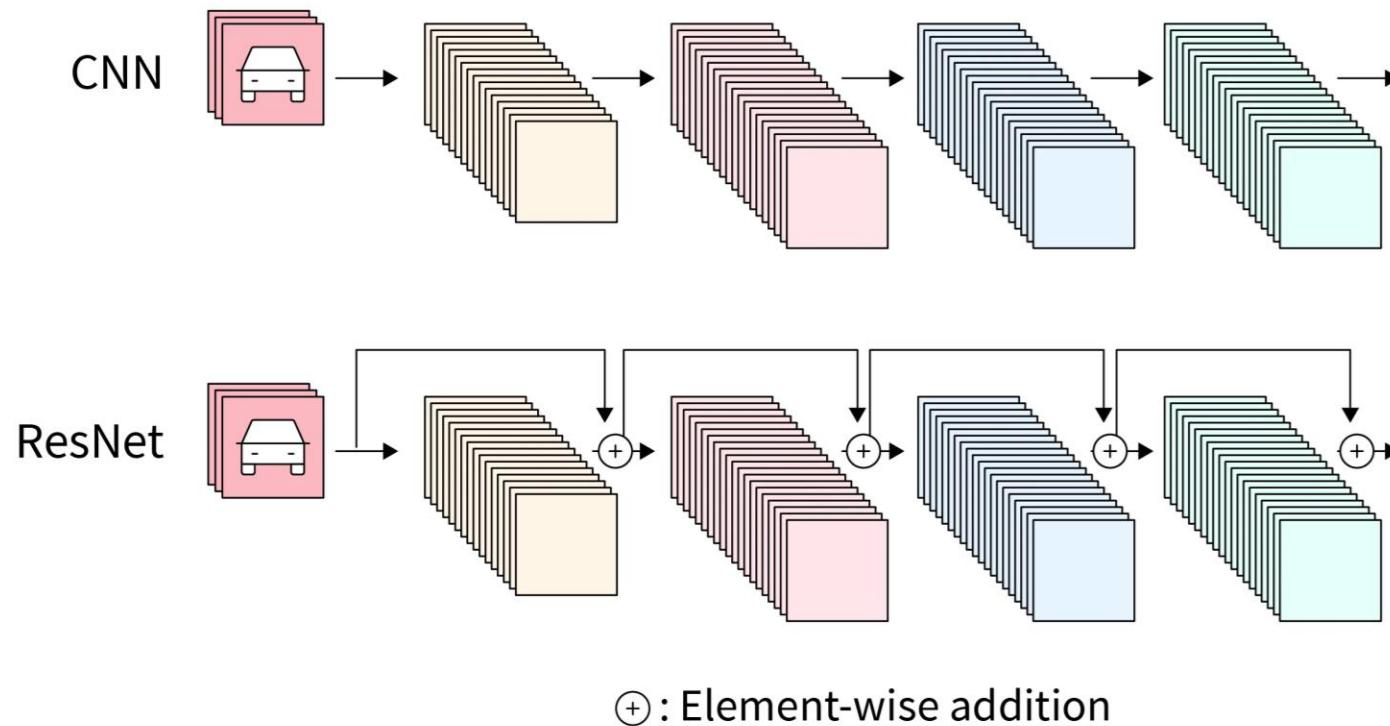
---

- Performs an identity mapping to connect the input of the subnetwork with its output
- All transformer architectures include residual connections. Very deep transformers cannot be trained without them.



# Skip / residual connections

---



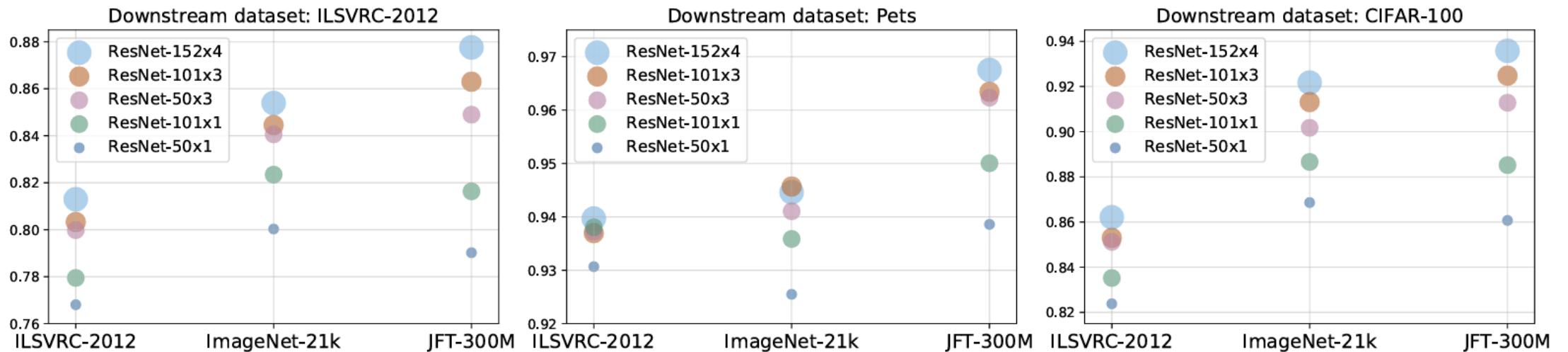
# Big Transfer: The value of pretraining

---

- **Scale wins (data × model).** Pre-training on larger supervised sets (ImageNet-21k, JFT-300M) and using bigger ResNet backbones yields better transfer than ImageNet-1k alone
- **Strong across data regimes, including few-shot.** The same pre-trained checkpoint transfers well from 1 example per class up to million-example datasets
- **A simple, general fine-tuning recipe (“BiT-HyperRule”).** One set of hyper-params works broadly: vanilla SGD, large batch ( $\approx 512$ ), learning rate decay by factor of 10 at 30/60/90% of training steps, and minimal augmentations (resize → random crop → flip). This reduces per-task tuning.
- **Pre-train once, adapt cheaply.** BiT emphasizes generalist representations: heavy compute is spent upstream; downstream fine-tuning is short and consistent across >20 datasets.
- **Benefits extend beyond classification.** Using BiT backbones improves detection when pre-training goes beyond ImageNet-1k (ImageNet-21k/JFT).
- **Bigger models + higher resolution help.** Very wide/deep ResNet-v2 variants (e.g., R152 $\times 4$ ) and higher crop resolutions further lift transfer accuracy.

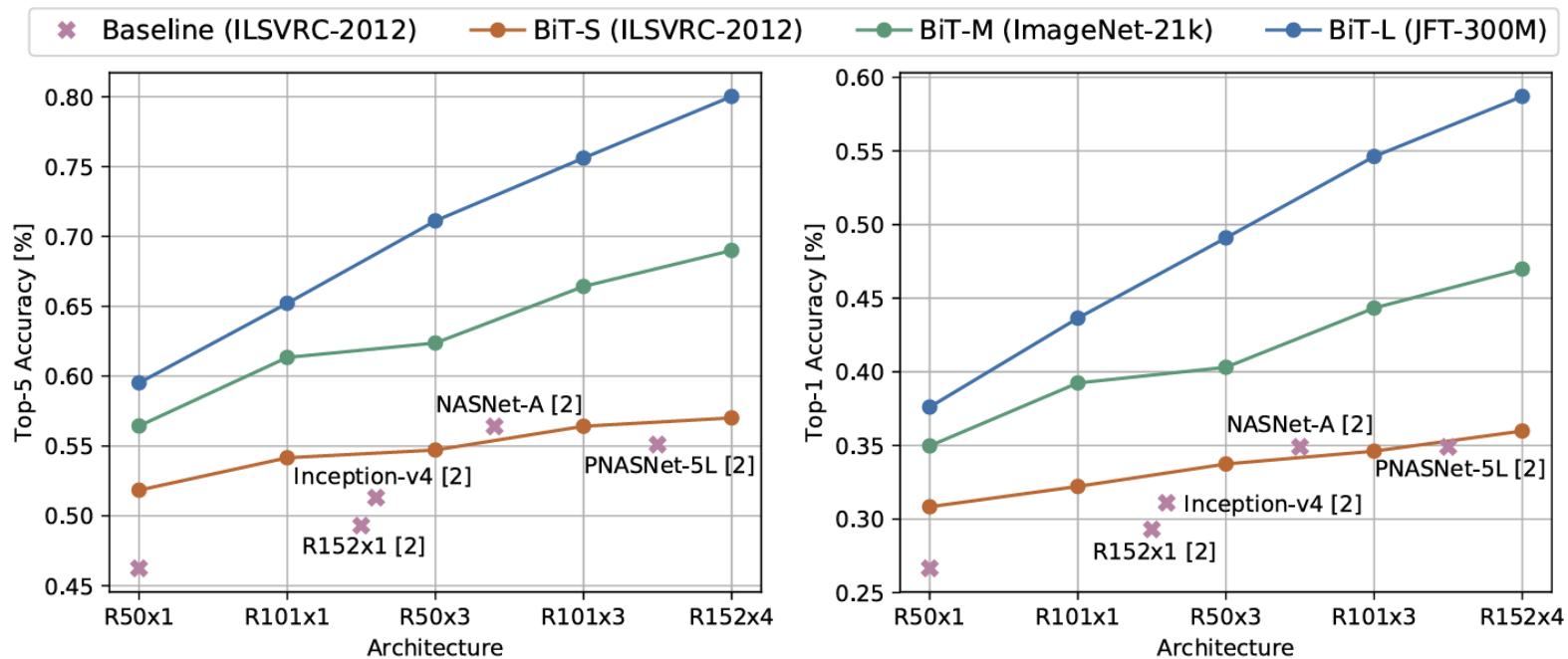
# Scale wins: classification

Effect of upstream data (shown on the x-axis) and model size on downstream performance.  
Exclusively using more data or larger models may hurt performance; both need to be increased in tandem.



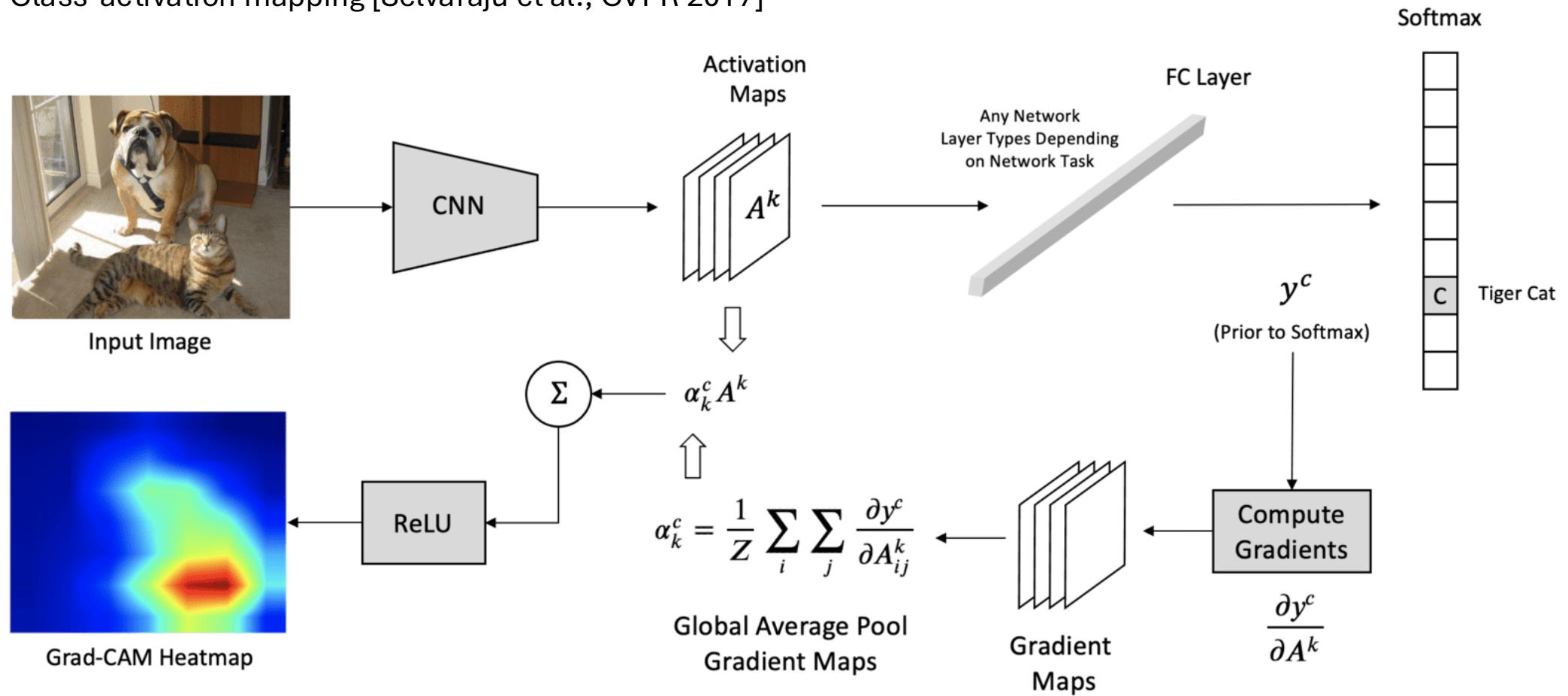
# Scale wins: object detection

Results on the original ObjectNet test set with resize and central crop augmentations.



# Interpreting trained models

Class-activation mapping [Selvaraju et al., CVPR 2017]



# Diagnosing failures through visual explanations

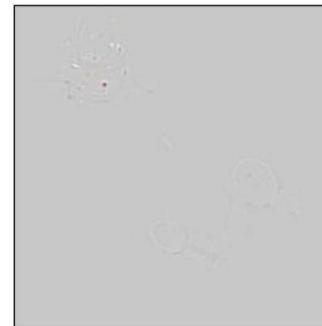
VGG-16 model



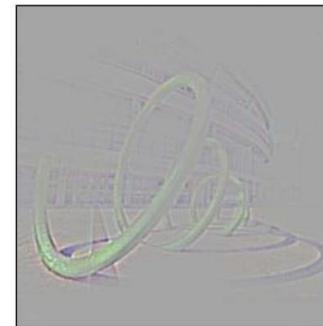
Ground truth: volcano



Ground truth: volcano



Ground truth: beaker



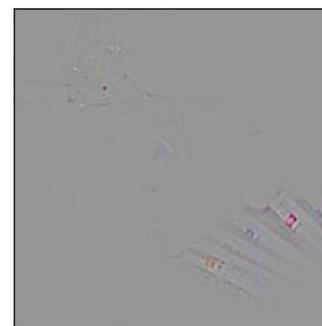
Ground truth: coil



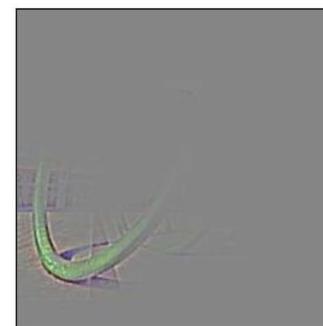
Predicted: sandbar



Predicted: car mirror



Predicted: syringe



Predicted: vine snake

# Vision transformers

---

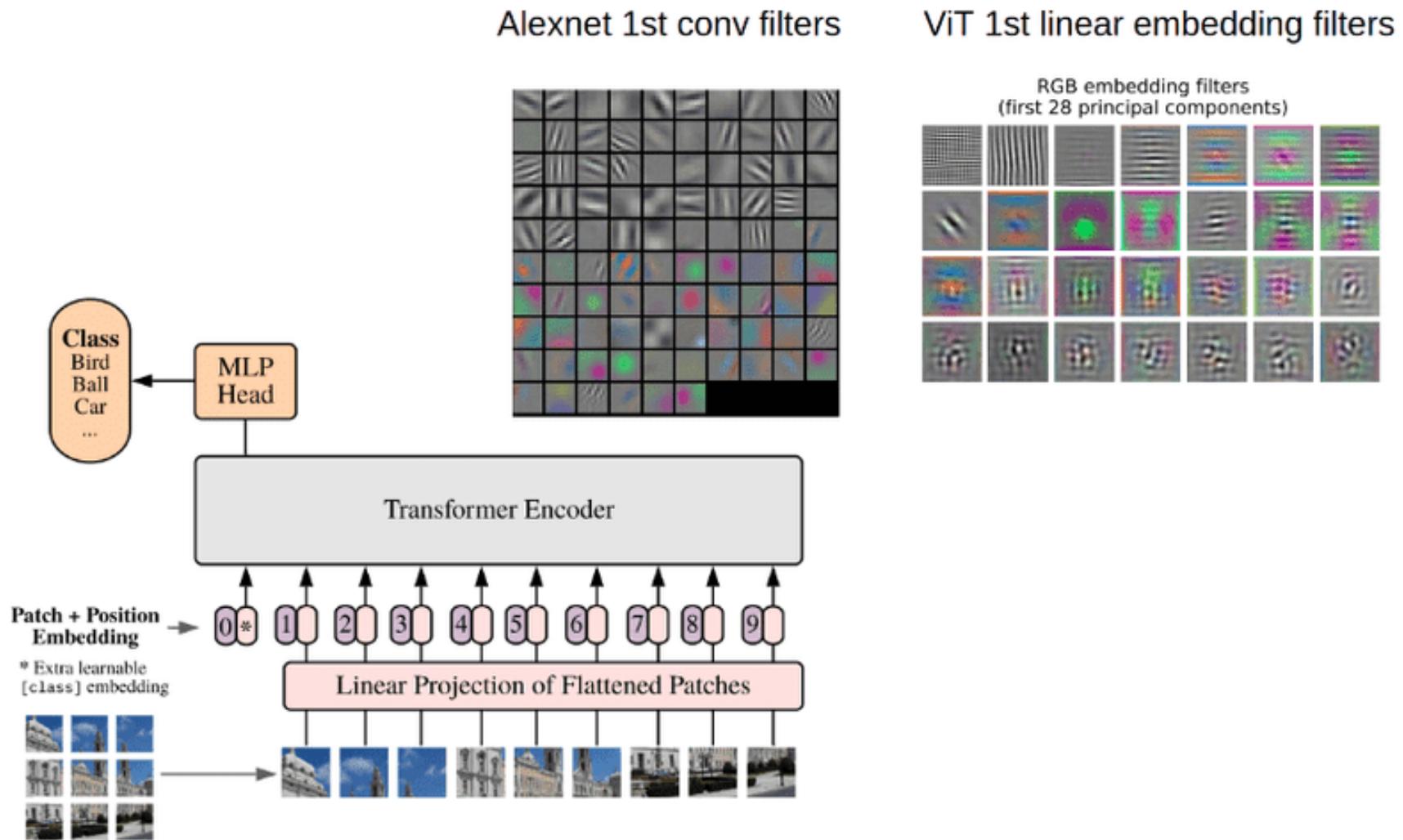
## AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy\*,†, Lucas Beyer\*, Alexander Kolesnikov\*, Dirk Weissenborn\*,  
Xiaohua Zhai\*, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby\*,†

\*equal technical contribution, †equal advising

Google Research, Brain Team  
`{adosovitskiy, neilhoulsby}@google.com`

# Vision transformers



# Vision transformers

---

Main claim: We don't need convolutions anymore!

Less image-specific inductive biases than CNNs:

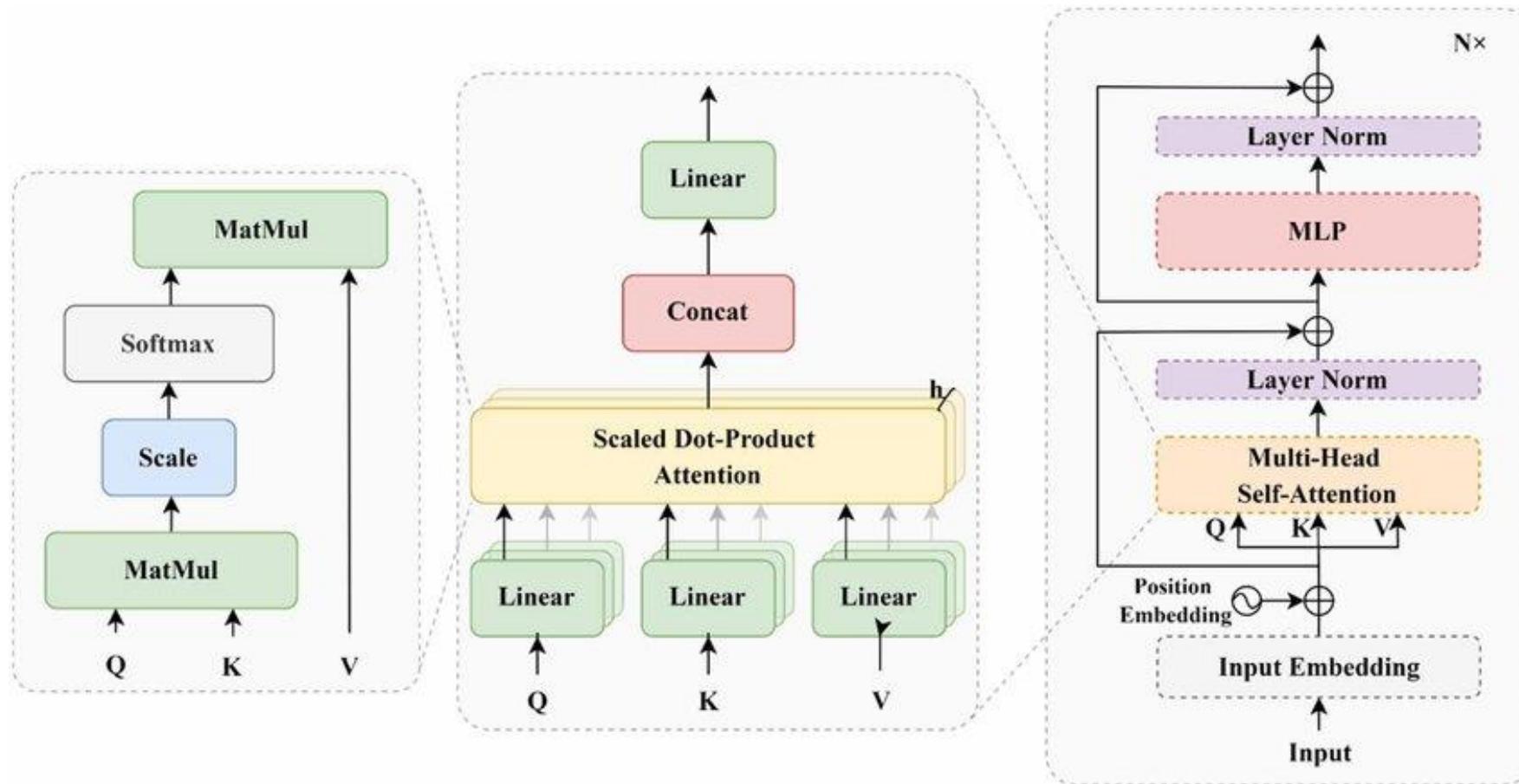
- **CNNs:** locality, 2D neighborhood structure & translation equivariance baked into each layer
- **ViTs:** only MLP layers are local and translation equivariant; attention layers are global

Reshape images  $x \in \mathbb{R}^{H \times W \times C}$  to flattened 2D-patch representations  $x_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$

( $P, P$ ) is the resolution of each image patch and  $N = \frac{WH}{P^2}$  the resulting number of patches

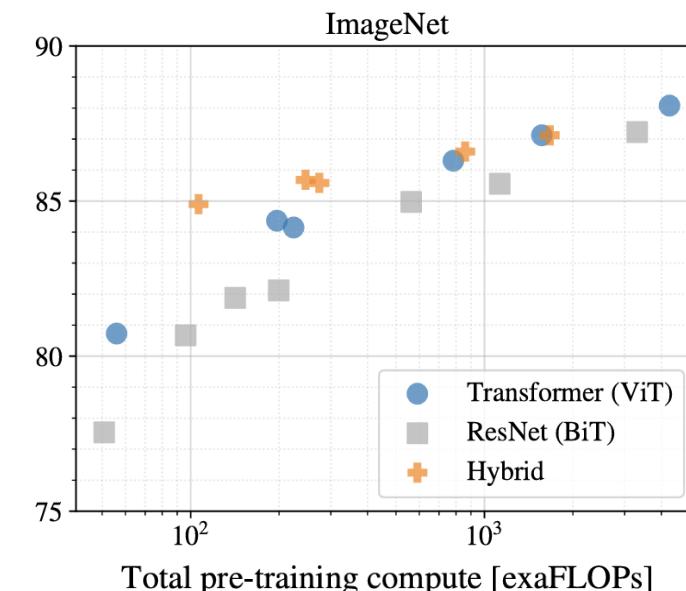
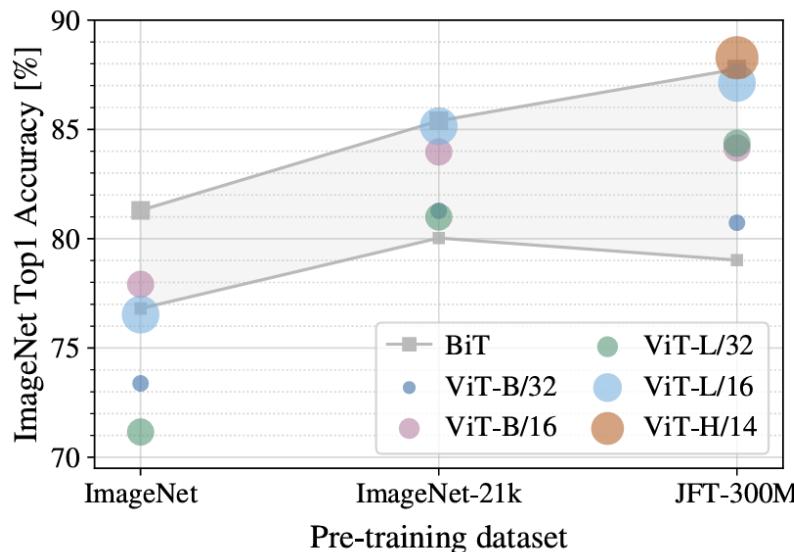
$$\begin{aligned}\mathbf{z}_0 &= [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, & \mathbf{E} &\in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \\ \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, & \ell &= 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, & \ell &= 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0)\end{aligned}$$

# Transformer block with MSA



# Transformers vs. CNNs

- **Compute/memory efficiency is favourable at scale.** With big pretraining, ViT attains SOTA with substantially less pretraining compute than strong CNN baselines (e.g., BiT), and large ViTs are **more memory-efficient** (larger per-core batch sizes) than comparable ResNets.
- **Model/patch scaling:** ViT comes in Base/Large/Huge variants; **smaller patches → longer sequences → higher compute but better accuracy.**
- **Hybrid option exists.** You can feed **CNN feature maps** into the Transformer instead of raw patches; this can help when data are smaller (otherwise ViTs are data hungry).



# The representation problem

---

**Can useful representations develop without labels?**

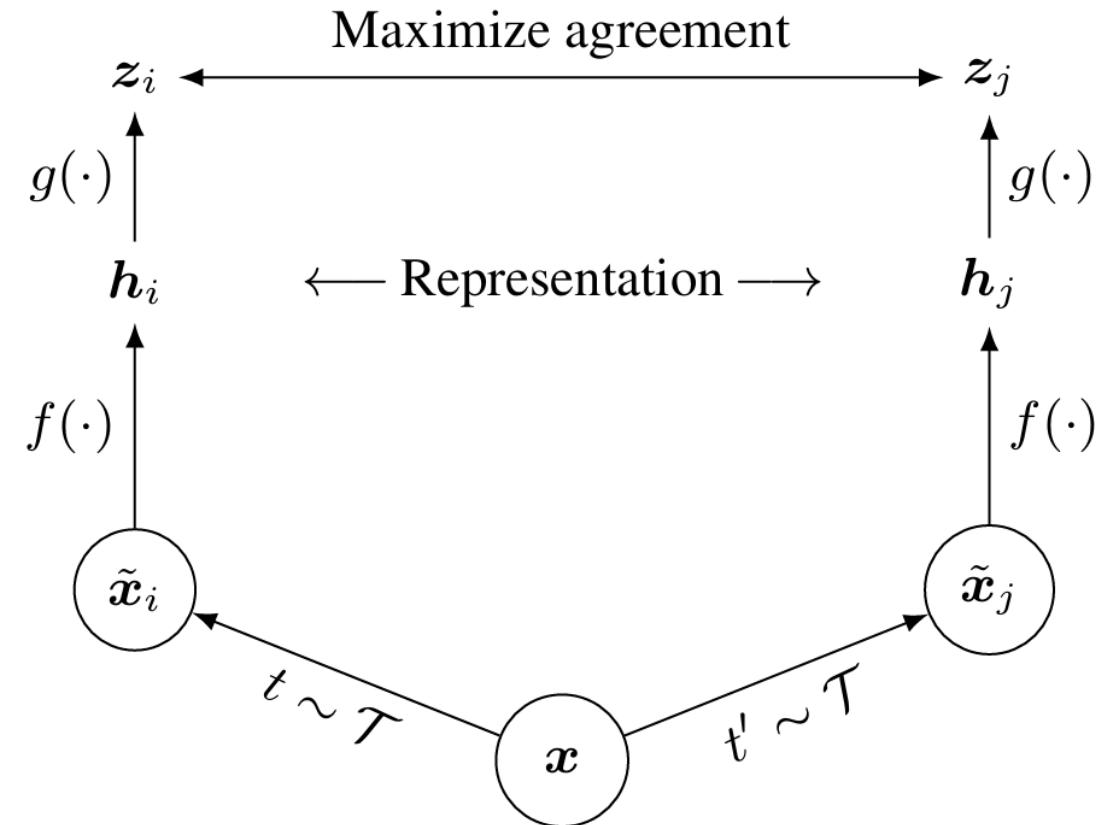
If so,

- We can learn more efficiently when we do gain access to labels
- We still learn useful things when label collection is impossible

# Can we beat supervision without labels?

- Leverage contrastive learning = SimCLR
- Key ingredients:
  - Composition of **data augmentations**
  - **Learnable non-linear projection** head to improve representation quality
  - **Large batch sizes**; long training
- Importance of positive and negative mining

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$



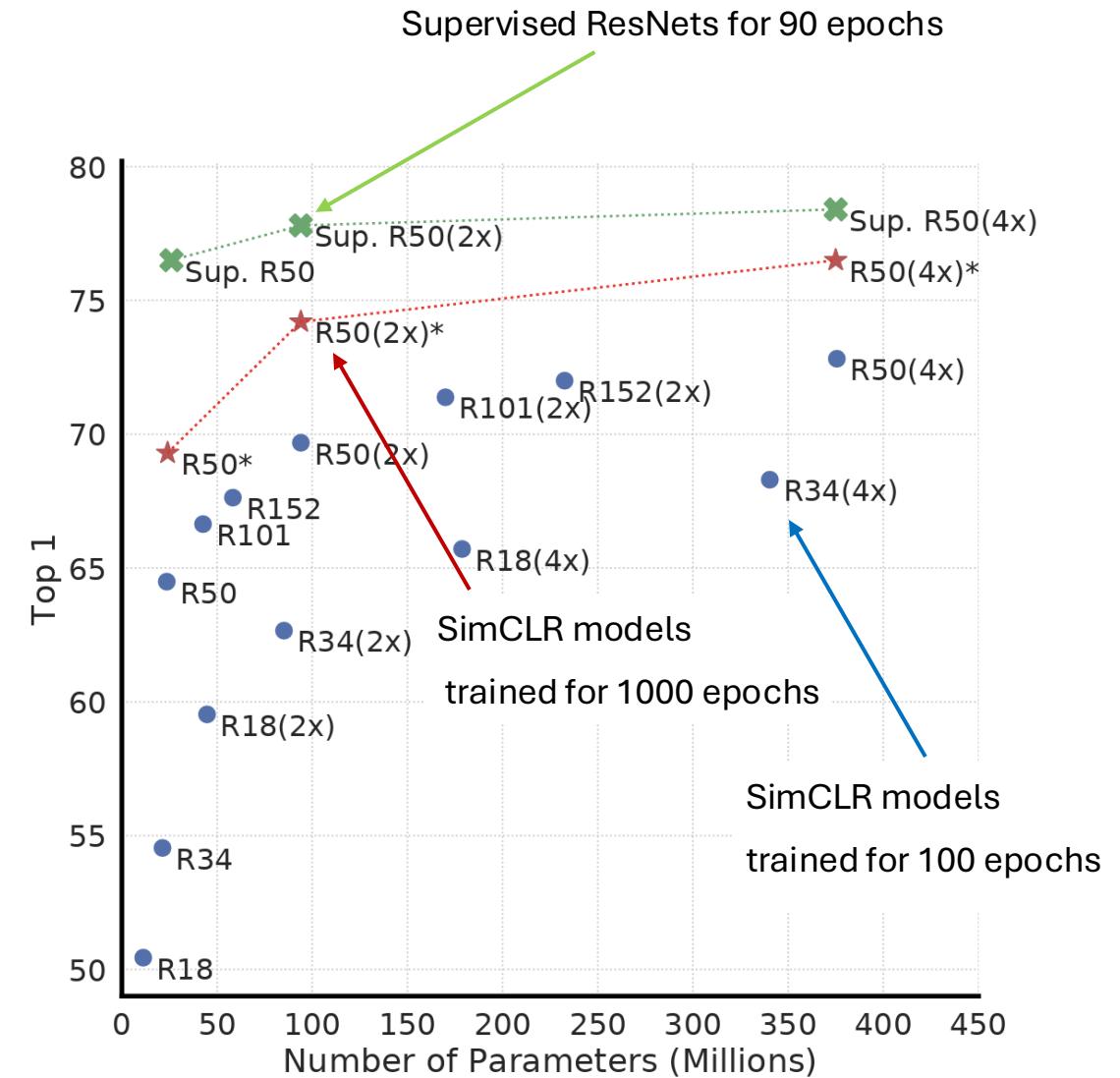
# SimCLR training

**Algorithm 1** SimCLR's main learning algorithm.

```

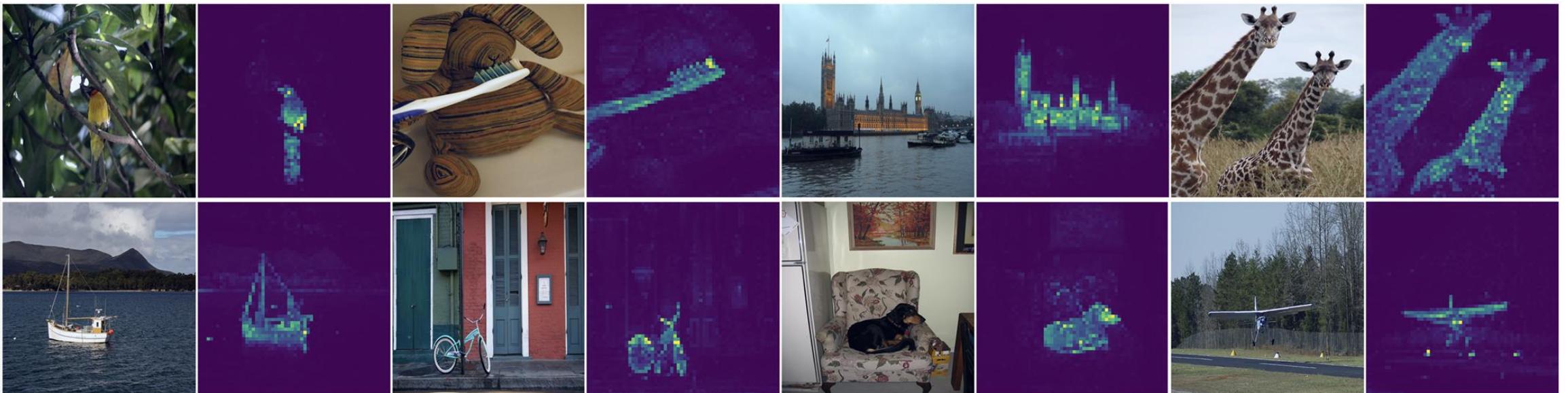
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
    for all  $k \in \{1, \dots, N\}$  do
        draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$ 
        # the first augmentation
         $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$ 
         $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation
         $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection
        # the second augmentation
         $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$ 
         $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation
         $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection
    end for
    for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity
    end for
    define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$ 
     $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$ 
    update networks  $f$  and  $g$  to minimize  $\mathcal{L}$ 
end for
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 

```



# Can we learn representations with ViTs?

Self-**distillation** with **no** labels instead of contrastive negatives [<https://arxiv.org/pdf/2104.14294>]

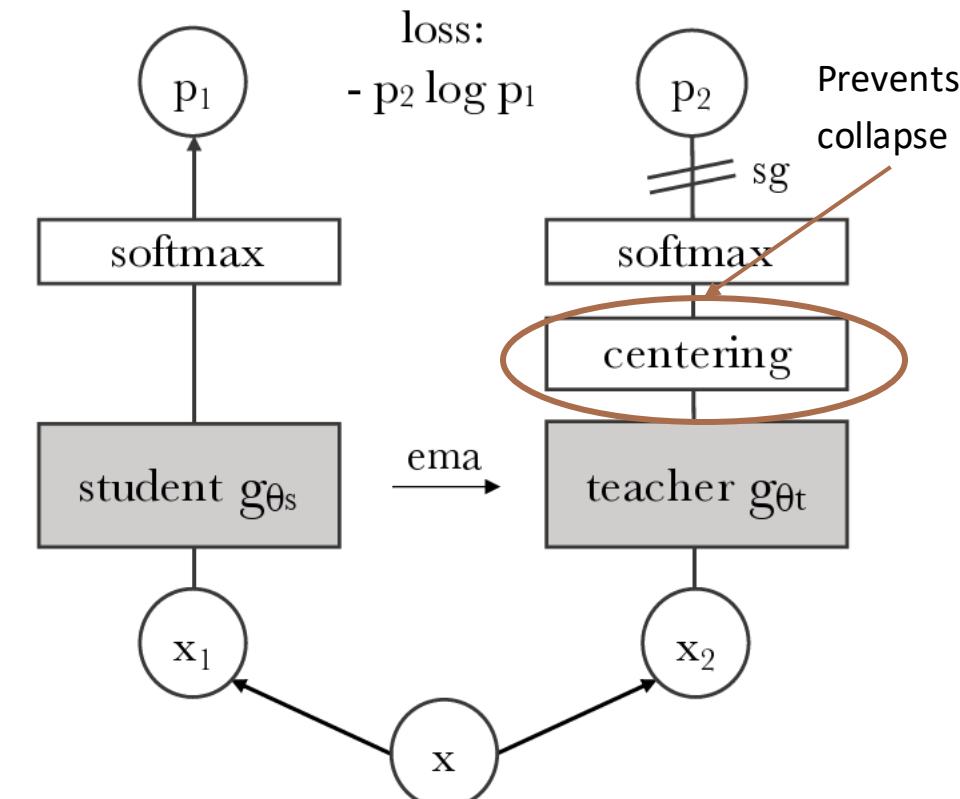


Self-attention from a ViT with 8x8 patches trained with no supervision

# DiNO v1

Architecture and training recipe:

- **Form of self-distillation** with student–teacher ViTs (sharing same architecture); **teacher = EMA** of student; teacher uses a copy of the student with **stop-gradient**.
- Train by **distribution matching**: cross-entropy between temperature-scaled softmax (teacher “sharpened”, student “softer”) + **centering** to prevent collapse.
- **Multi-crop views**: a couple of global crops + several local crops encourage locality/globality.
- Outcomes: strong **linear-probe** ImageNet (e.g., ViT-B  $\sim$ 80% top-1) and **object-like attention maps** without labels.



# Downstream tasks

## Image retrieval

Pretrain	Arch.	Pretrain	$\mathcal{R}_{\text{Ox}}$		$\mathcal{R}_{\text{Par}}$	
			M	H	M	H
Sup. [57]	RN101+R-MAC	ImNet	49.8	18.5	74.0	<b>52.1</b>
Sup.	ViT-S/16	ImNet	33.5	8.9	63.0	37.2
DINO	ResNet-50	ImNet	35.4	11.1	55.9	27.5
DINO	ViT-S/16	ImNet	41.8	13.7	63.1	34.4
DINO	ViT-S/16	GLDv2	<b>51.5</b>	<b>24.3</b>	<b>75.3</b>	51.6

## Copy detection

Method	Arch.	Dim.	Resolution	mAP
Multigrain [5]	ResNet-50	2048	$224^2$	75.1
Multigrain [5]	ResNet-50	2048	largest side 800	82.5
Supervised [69]	ViT-B/16	1536	$224^2$	76.4
DINO	ViT-B/16	1536	$224^2$	81.7
DINO	ViT-B/8	1536	$320^2$	<b>85.5</b>

## Semantic segmentation

*Supervised*



*DINO*



External

# DiNO v2

---

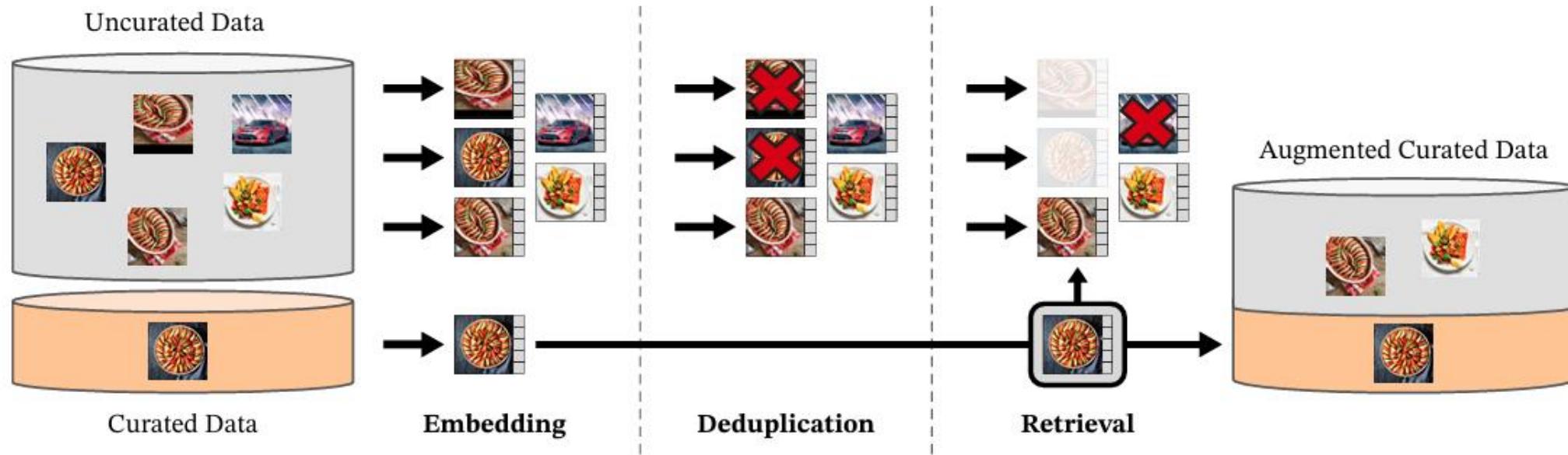
Learning Robust Visual Features without Supervision [<https://arxiv.org/pdf/2304.07193>]

Scaling recipe to “foundation-style” features → data is key

- **Scale everything:** larger ViTs (up to ~1B params teacher) and a **curated 142M-image dataset (LVD-142M)** built by retrieval + deduplication across diverse sources.
- **Stability/throughput tweaks** and **distillation:** train a huge teacher, **distill to smaller ViTs** for release.
- Result: **robust frozen features** that are SOTA/competitive across image-level and pixel-level tasks **without fine-tuning.**

# DiNO v2 – data curation

Retrieve from pool of uncurated data, images that are close to those in curated datasets

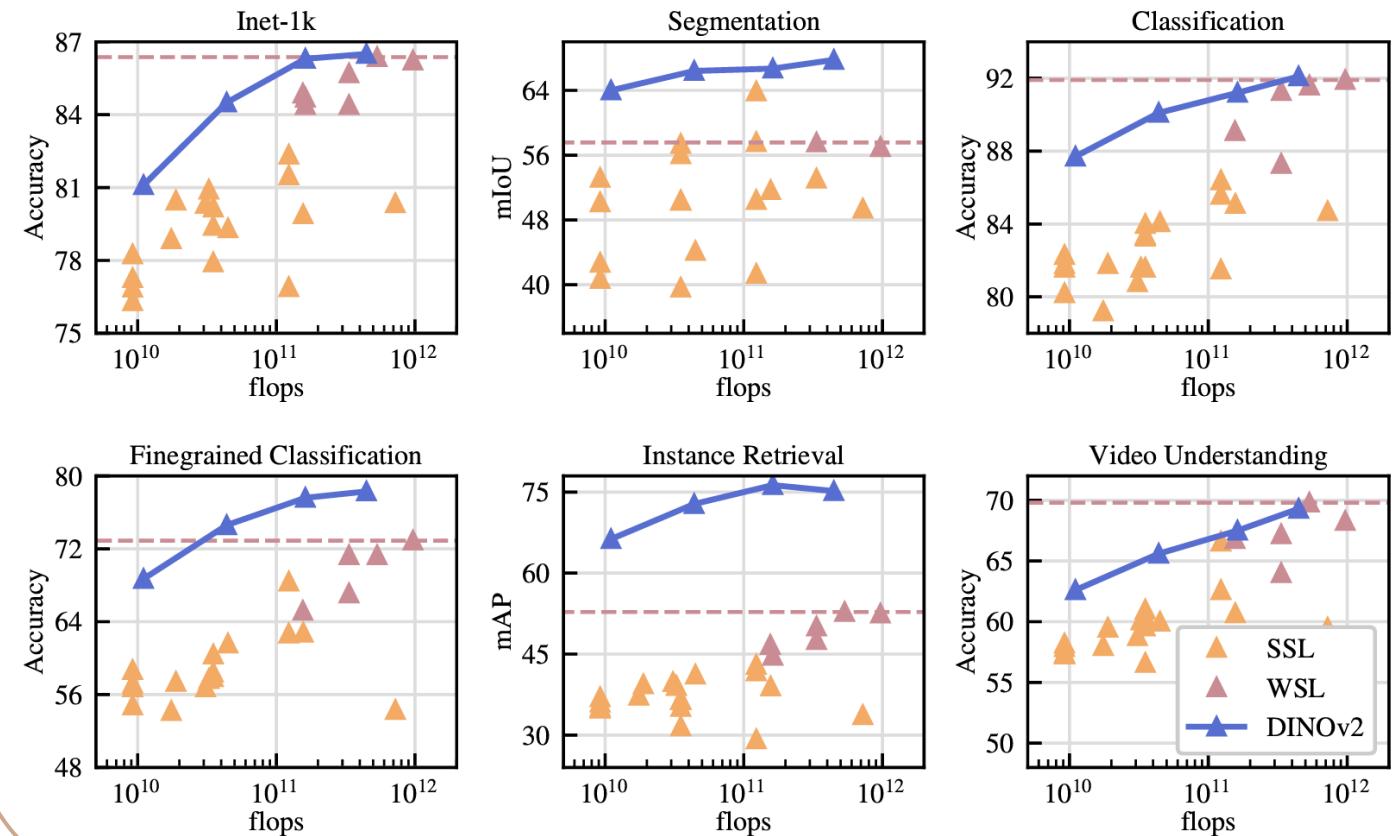


# DiNO v2 – representations & performance

## Visualisation of PCA components

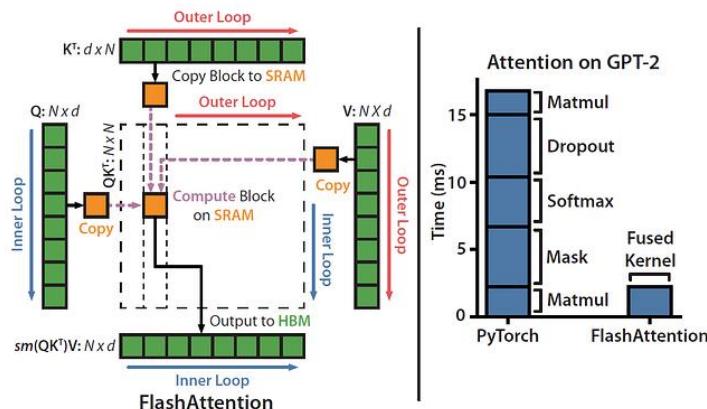
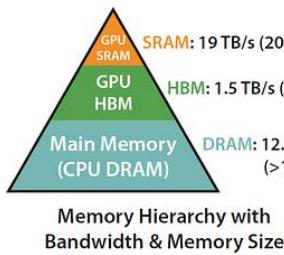


## Performance with scaling

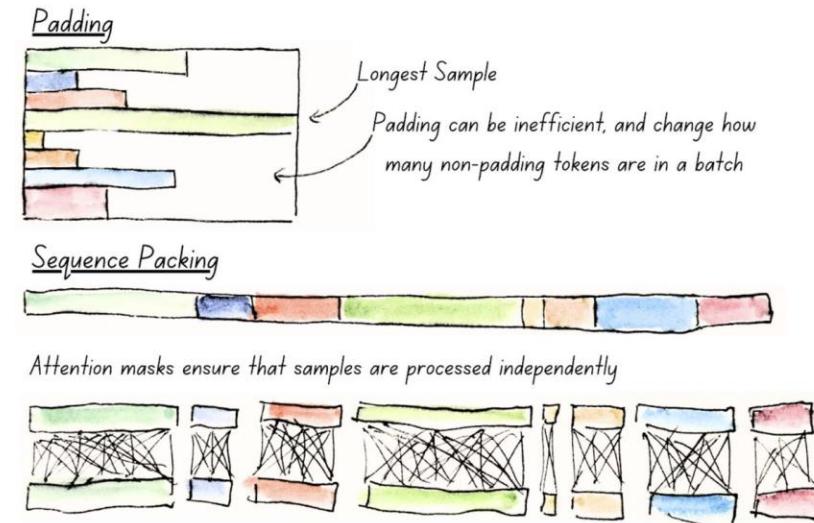


# DiNO v2 - efficiency

## Flash attention



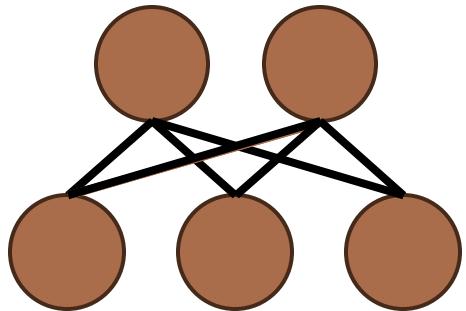
## Sequence packing



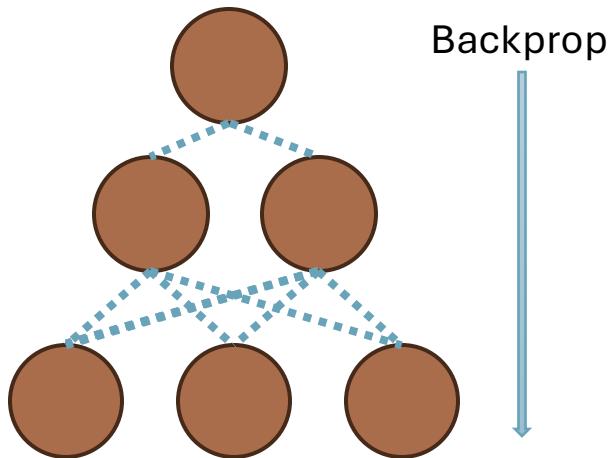
# Putting it all together

---

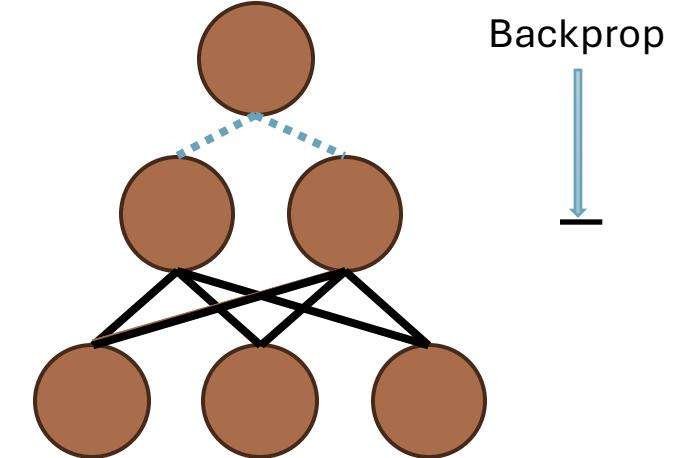
**Pre-training**



**Fine-tuning**



**Linear probing**



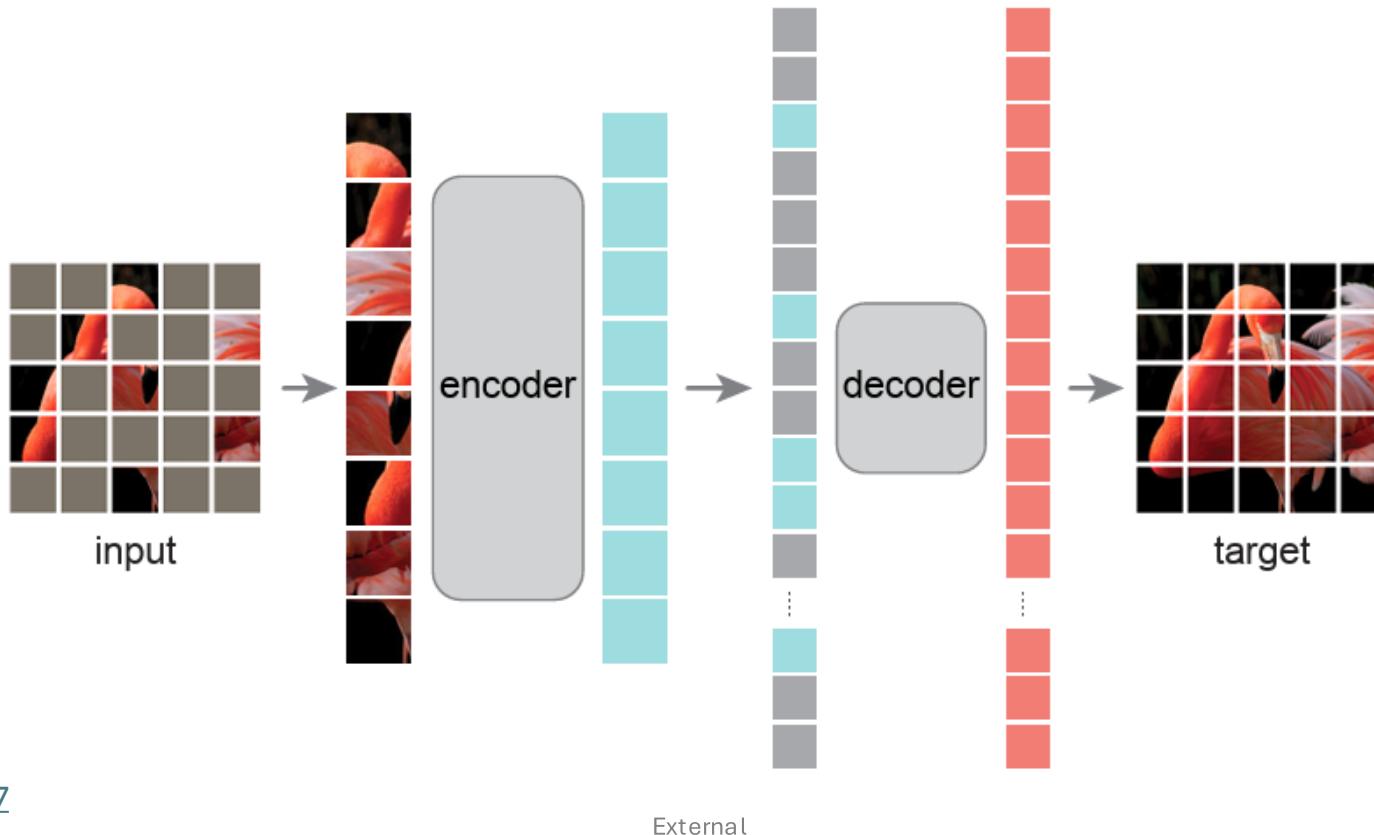
# Other approaches to self-supervision

## Masked Autoencoders Are Scalable Vision Learners

Kaiming He<sup>\*,†</sup> Xinlei Chen<sup>\*</sup> Saining Xie Yanghao Li Piotr Dollár Ross Girshick

\* equal technical contribution † project lead

Facebook AI Research (FAIR)



# From image to video

## VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training

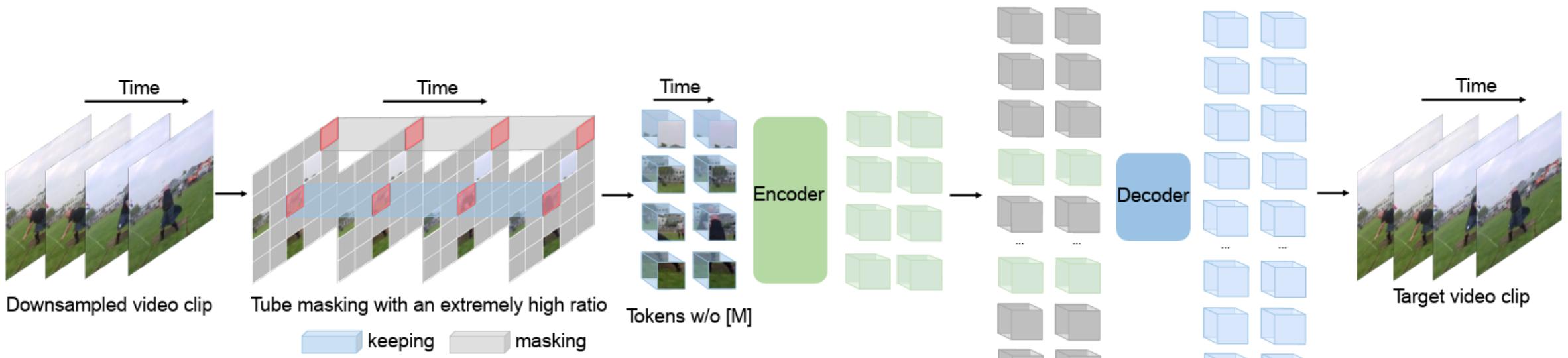
Zhan Tong<sup>1,2\*</sup> Yibing Song<sup>2</sup> Jue Wang<sup>2</sup> Limin Wang<sup>1,3†</sup>

<sup>1</sup>State Key Laboratory for Novel Software Technology, Nanjing University

<sup>2</sup>Tencent AI Lab

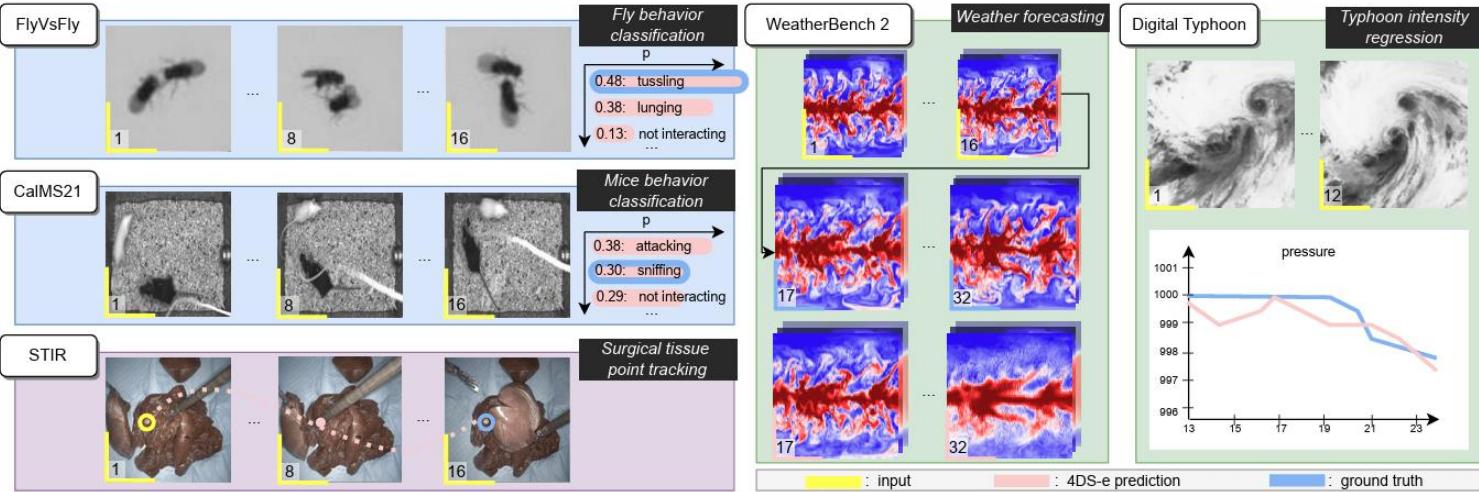
<sup>3</sup>Shanghai AI Lab

tongzhan@smail.nju.edu.cn {yibingsong.cv, arphid}@gmail.com lmwang@nju.edu.cn

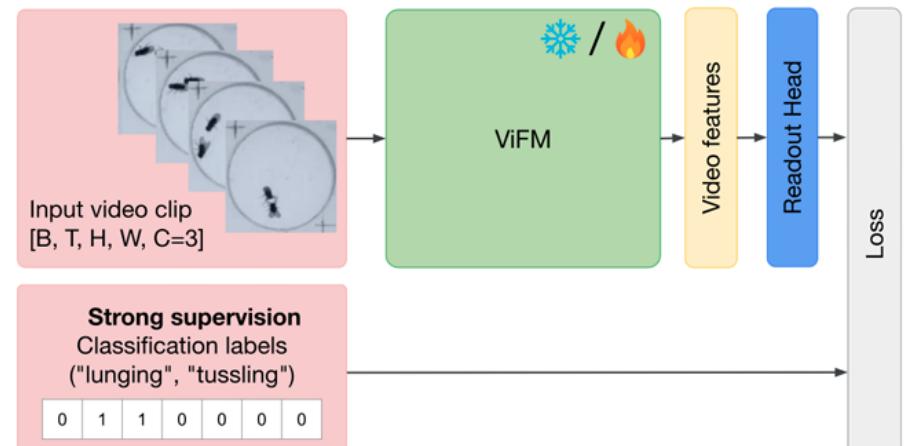


# Generalizing representations to scientific domains

## Tasks



## Evaluation framework

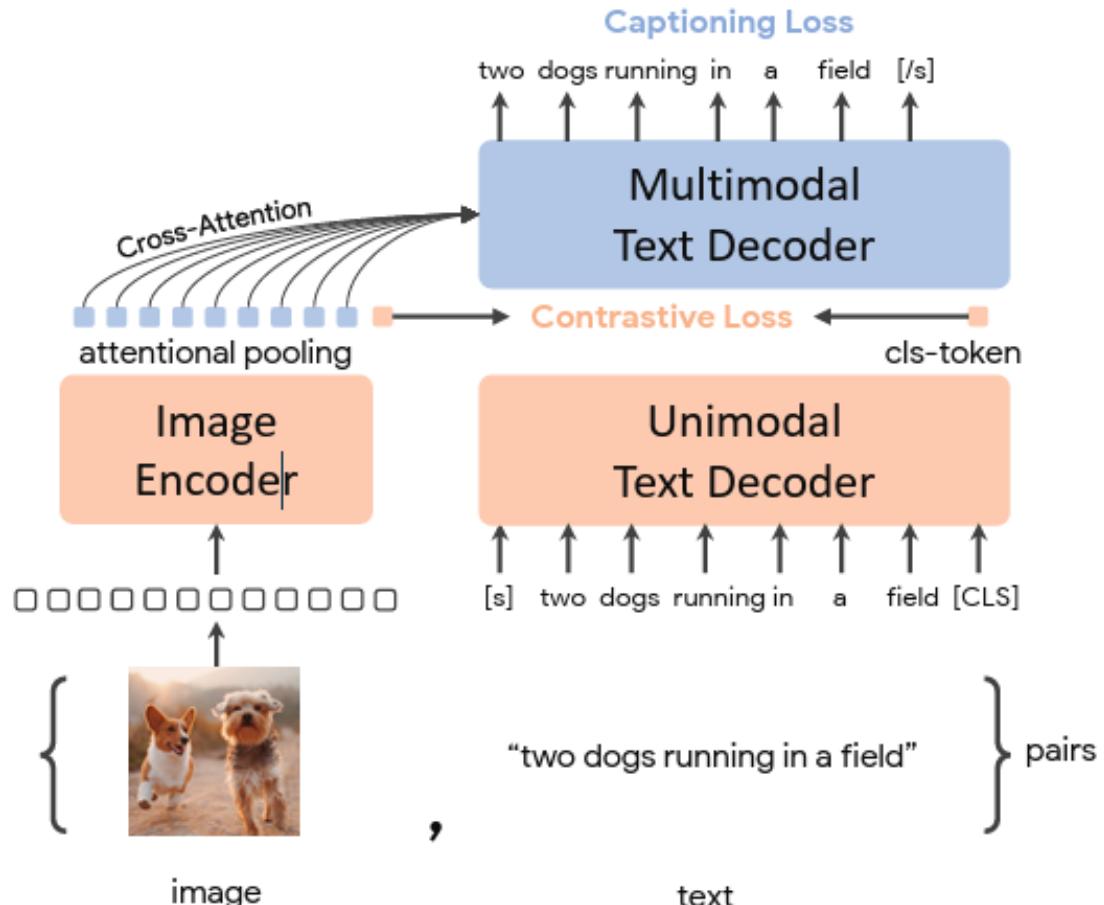


Zhao et al. ICML (2024)

Hasson et al. ICCV (2025)

External

# Obtaining image representations from text



**Algorithm 1** Pseudocode of Contrastive Captioners architecture.

```
# image, text.ids, text.labels, text.mask: paired {image, text} data
# con_query: 1 query token for contrastive embedding
# cap_query: N query tokens for captioning embedding
# cls_token_id: a special cls_token_id in vocabulary

def attentional_pooling(features, query):
    out = multihead_attention(features, query)
    return layer_norm(out)

img_feature = vit_encoder(image) # [batch, seq_len, dim]
con_feature = attentional_pooling(img_feature, con_query) # [batch, 1, dim]
cap_feature = attentional_pooling(img_feature, cap_query) # [batch, N, dim]

ids = concat(text.ids, cls_token_id)
mask = concat(text.mask, zeros_like(cls_token_id)) # unpad cls_token_id
txt_embs = embedding_lookup(ids)
unimodal_out = lm_transformers(txt_embs, mask, cross_attn=None)
multimodal_out = lm_transformers(
    unimodal_out[:, :-1, :], mask, cross_attn=cap_feature)
cls_token_feature = layer_norm(unimodal_out)[:, -1:, :] # [batch, 1, dim]

con_loss = contrastive_loss(con_feature, cls_token_feature)
cap_loss = softmax_cross_entropy_loss(
    multimodal_out, labels=text.labels, mask=text.mask)
```

`vit_encoder`: vision transformer based encoder; `lm_transformer`: language-model transformers.

# Semantic segmentation

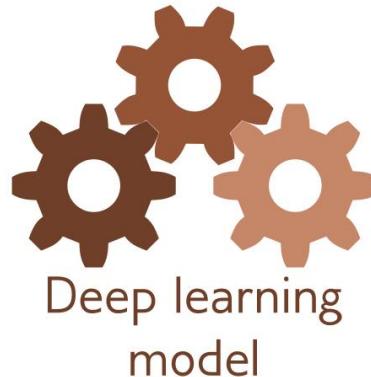
Real world input



Model  
input

$$\begin{bmatrix} 183 \\ 204 \\ 231 \\ 185 \\ 204 \\ 232 \\ \vdots \end{bmatrix}$$

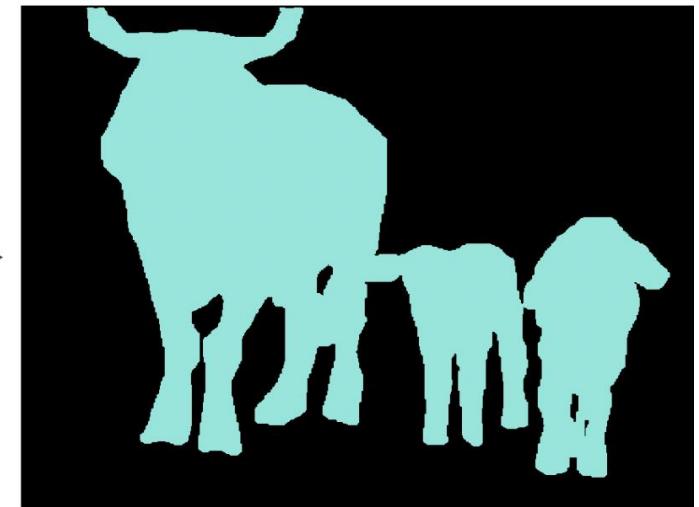
Model



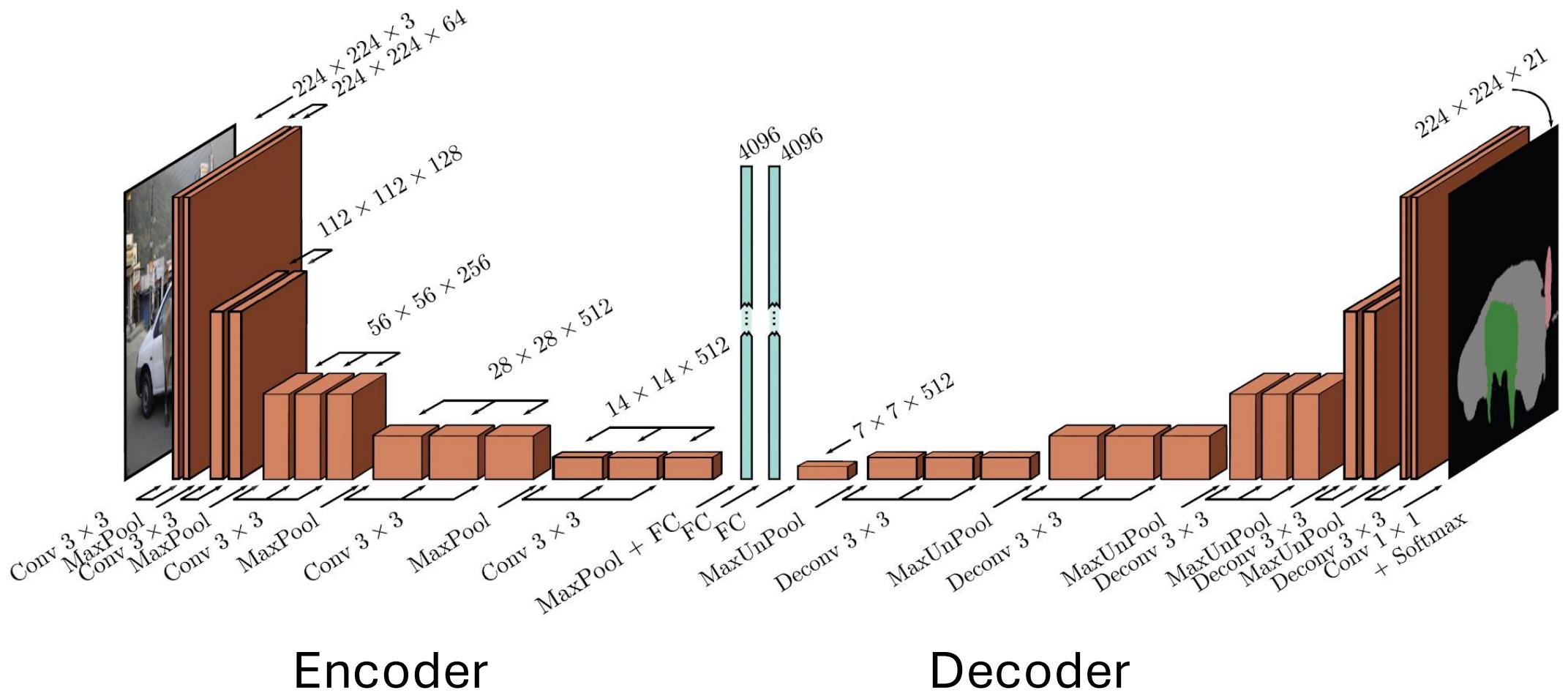
Model  
output

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

Real world output



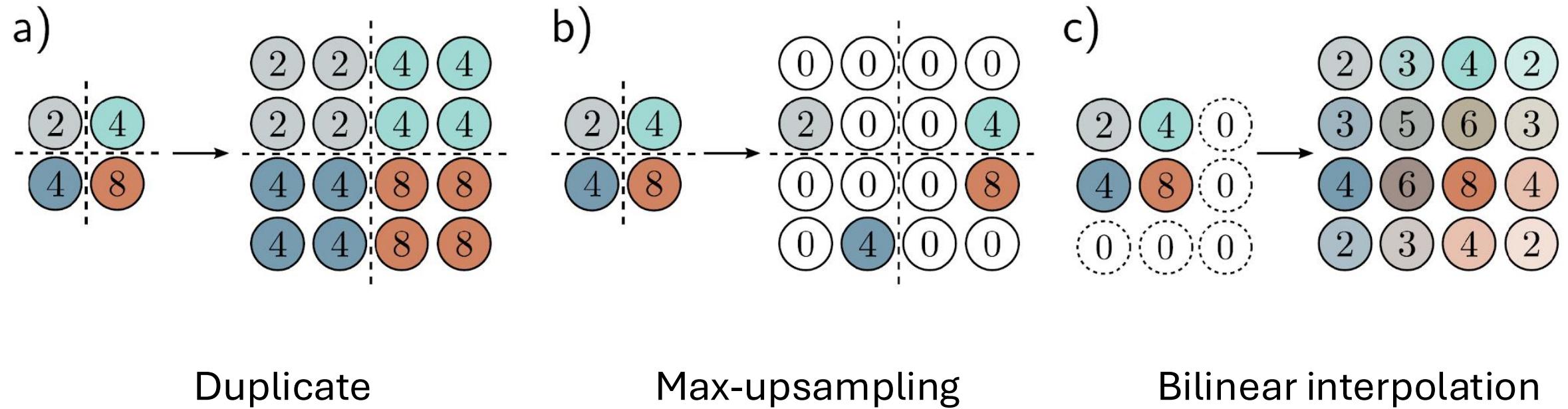
# Semantic segmentation – how it started



Encoder

Decoder

# Upsampling



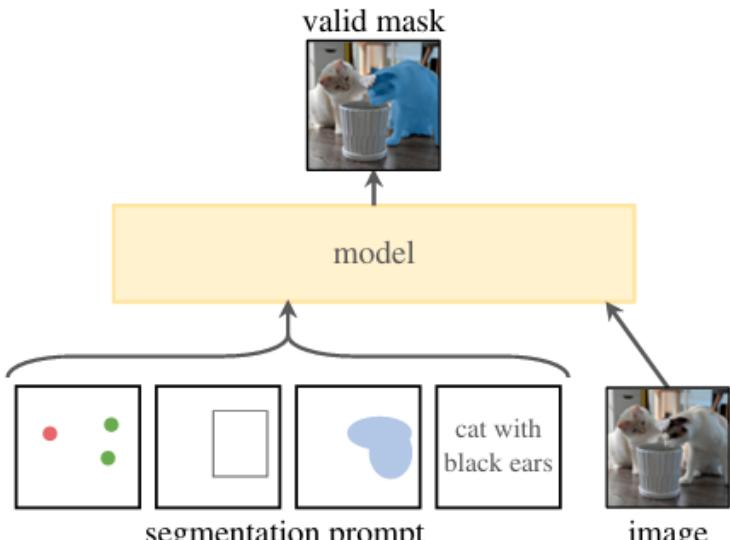
Duplicate

Max-upsampling

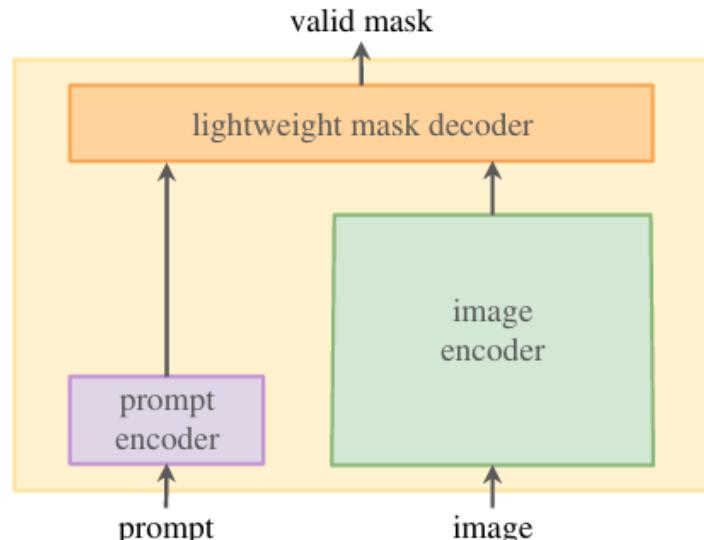
Bilinear interpolation

# Semantic segmentation – how it's going

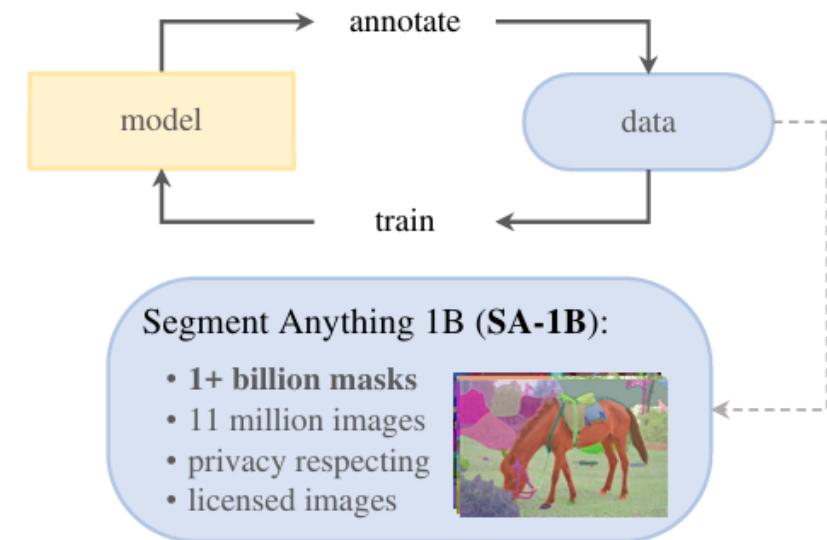
What task will enable zero-shot generalization?



(a) **Task:** promptable segmentation



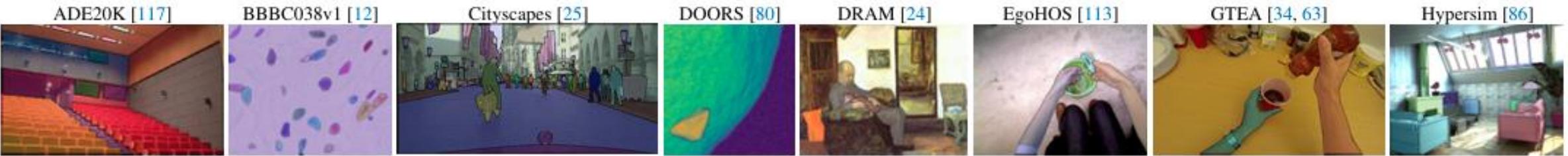
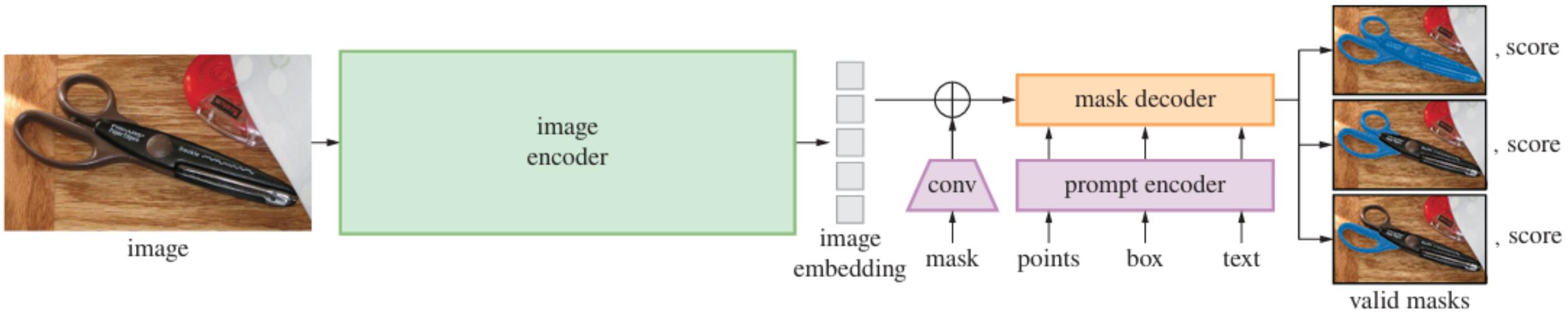
(b) **Model:** Segment Anything Model (SAM)



(c) **Data:** data engine (top) & dataset (bottom)

# Segment Anything

Zero-shot capabilities



**EIT** Oxford

## A Deep Dive into the Vision Transformer (ViT)

Google Colab [link](#)

# We'd love to hear your feedback!

---



**Scan the QR code or use the link**

<https://forms.office.com/e/J9y561KuDV>