# A Zoo of Models
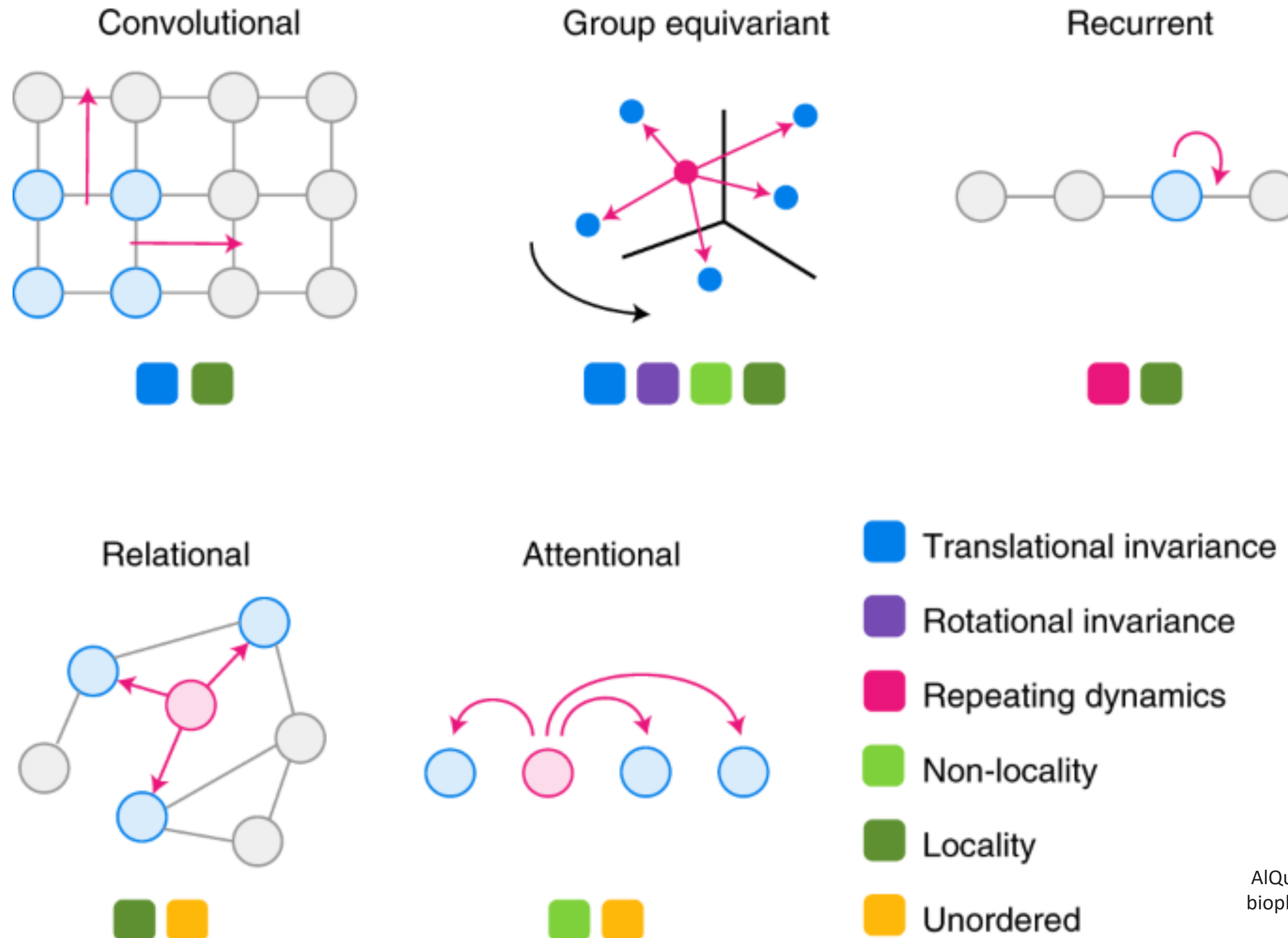
L3, Structural Bioinformatics

**WiSe 2023/24, Heidelberg Universitys**

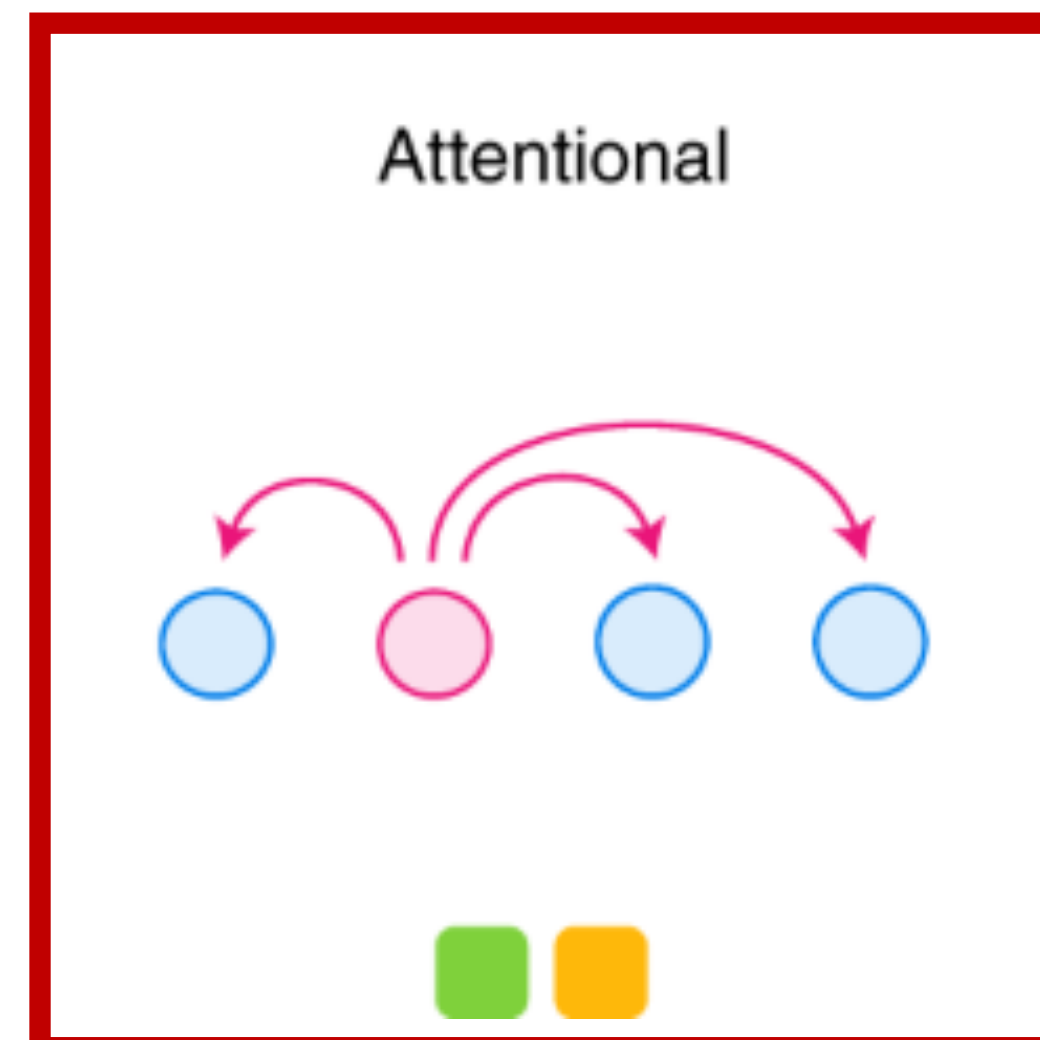# How to make sense of all these models?

## Find the inductive biases they instill in the network



AlQuraishi, M., Sorger, P.K. Differentiable biology: using deep learning for biophysics-based and data-driven modeling of molecular mechanisms. *Nat Methods* **18**, 1169–1180 (2021).
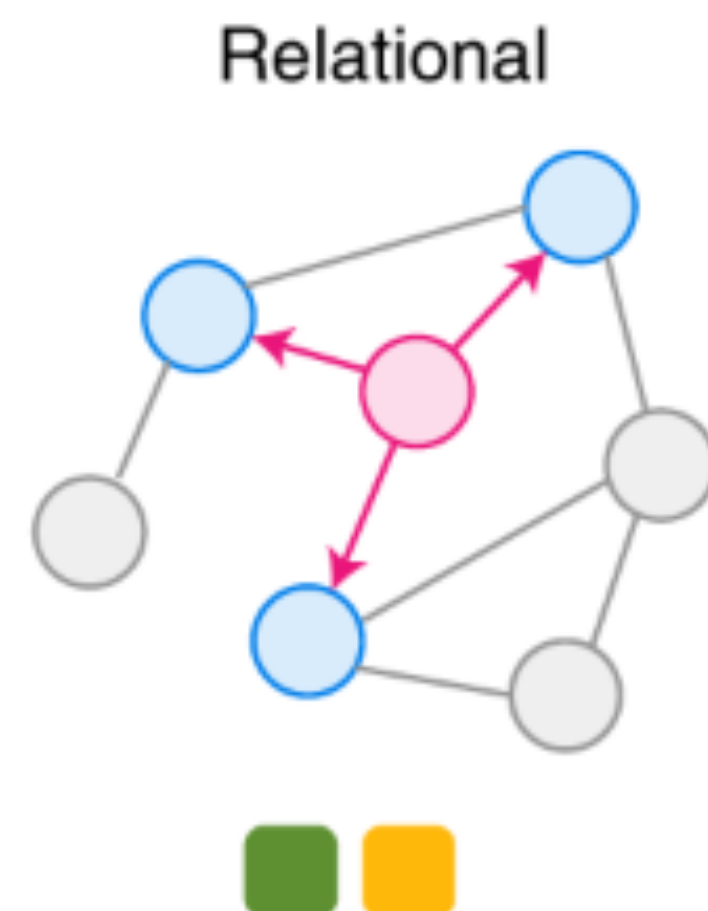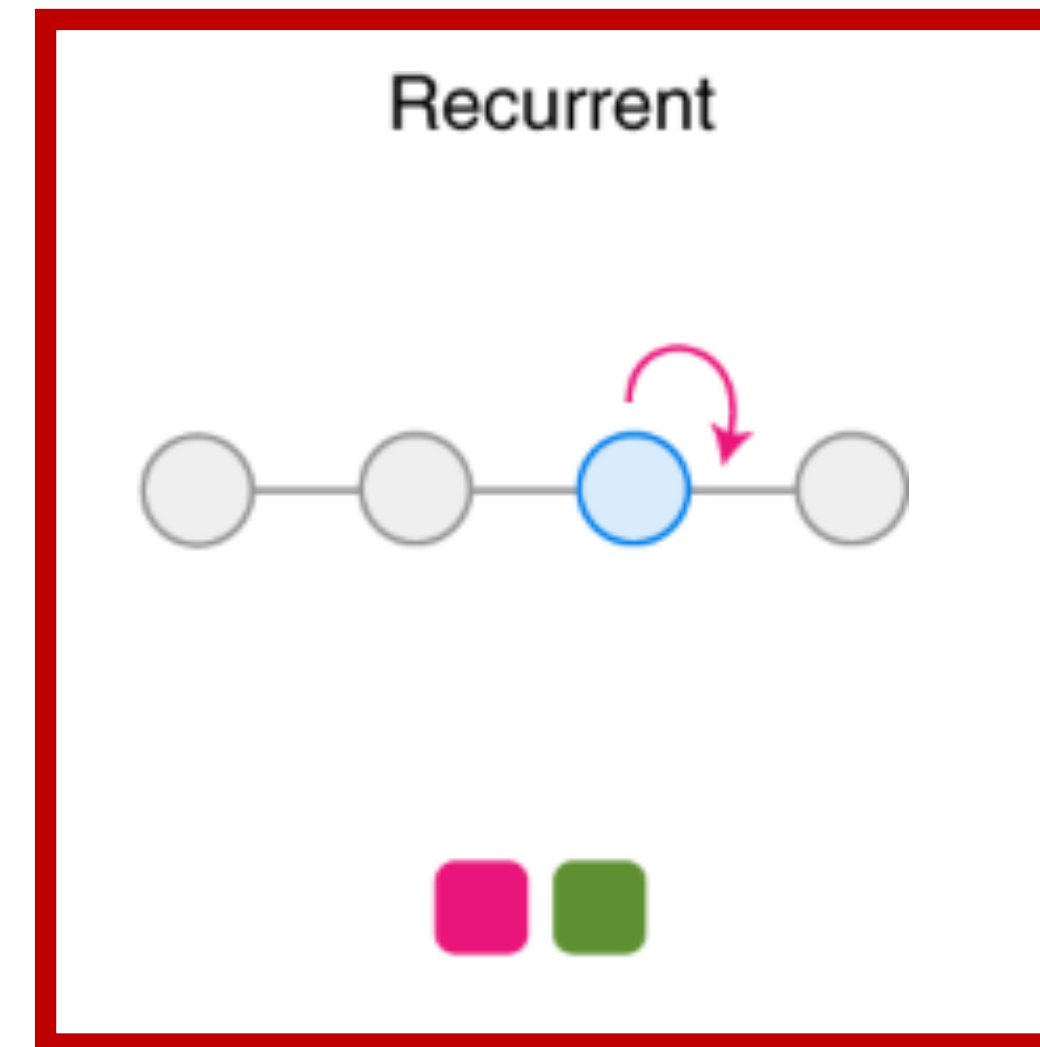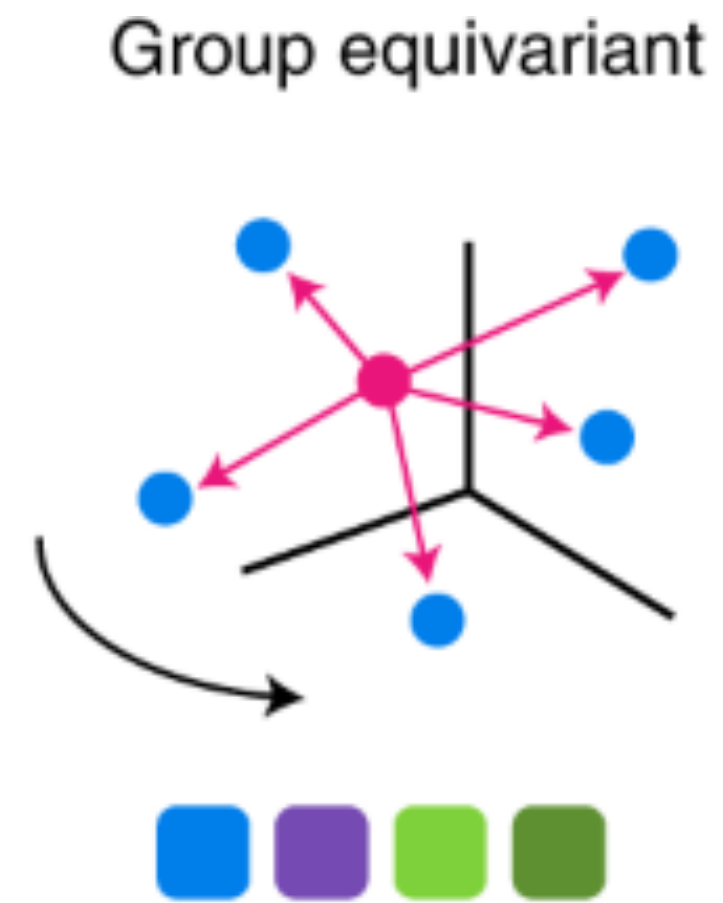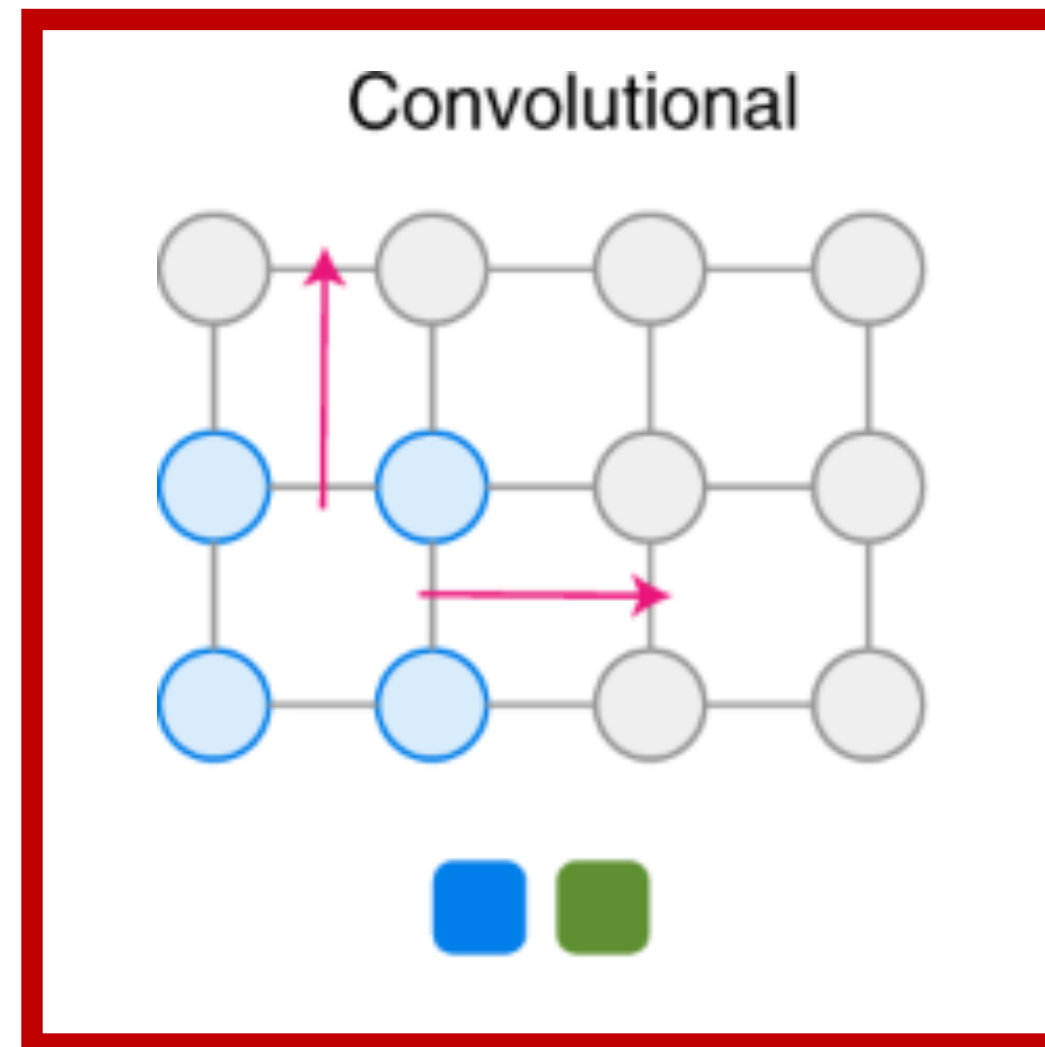
# How to make sense of all these models?

## Find the inductive biases they instill in the network



**This week**

# How to make sense of all these models?

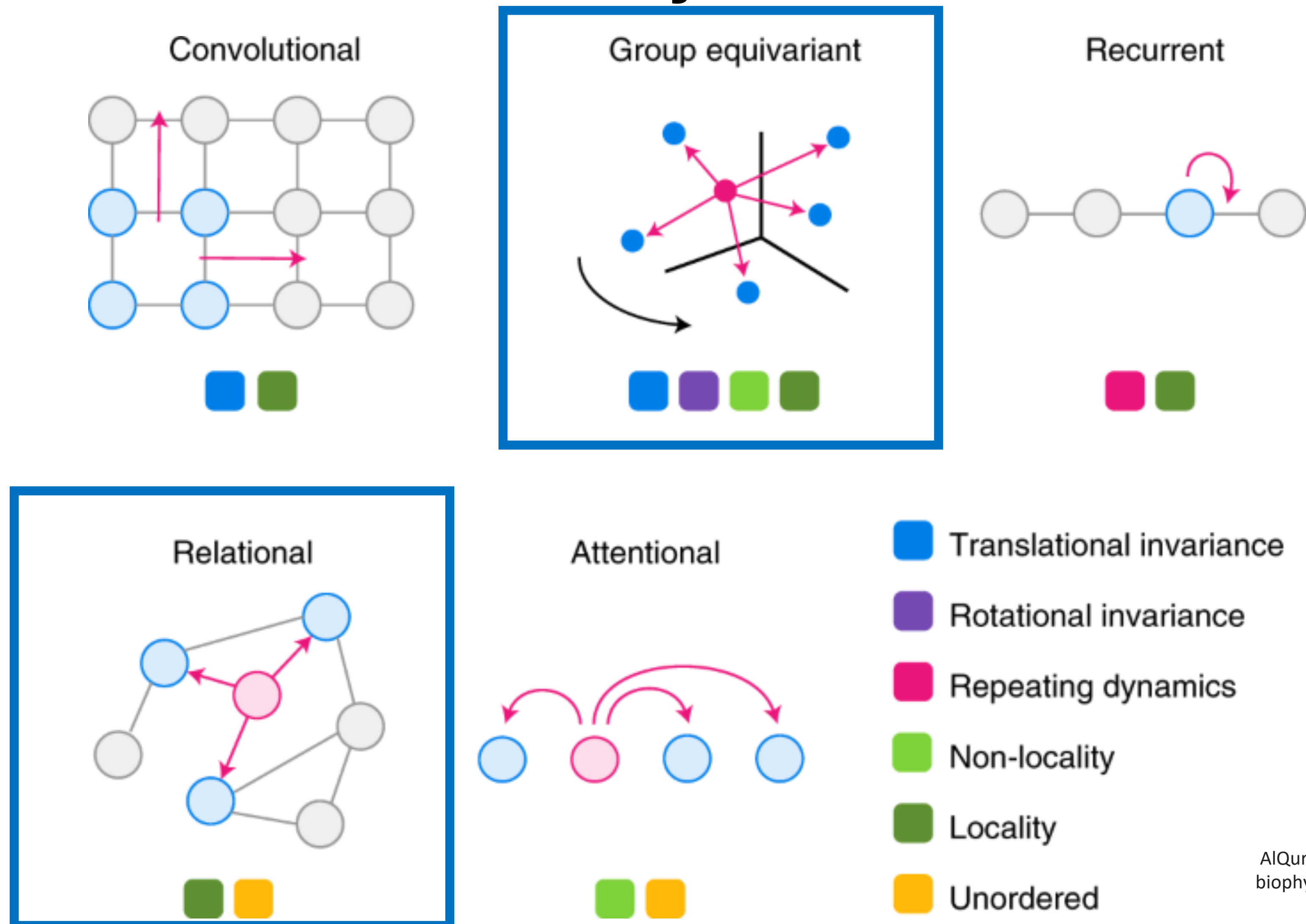## Find the inductive biases they instill in the network



Convolutional

Group equivariant

Recurrent

**Next week**

Relational

Attentional

- ■ Translational invariance
- ■ Rotational invariance
- ■ Repeating dynamics
- ■ Non-locality
- ■ Locality
- ■ Unordered

AlQuraishi, M., Sorger, P.K. Differentiable biology: using deep learning for biophysics-based and data-driven modeling of molecular mechanisms. *Nat Methods* **18**, 1169–1180 (2021).

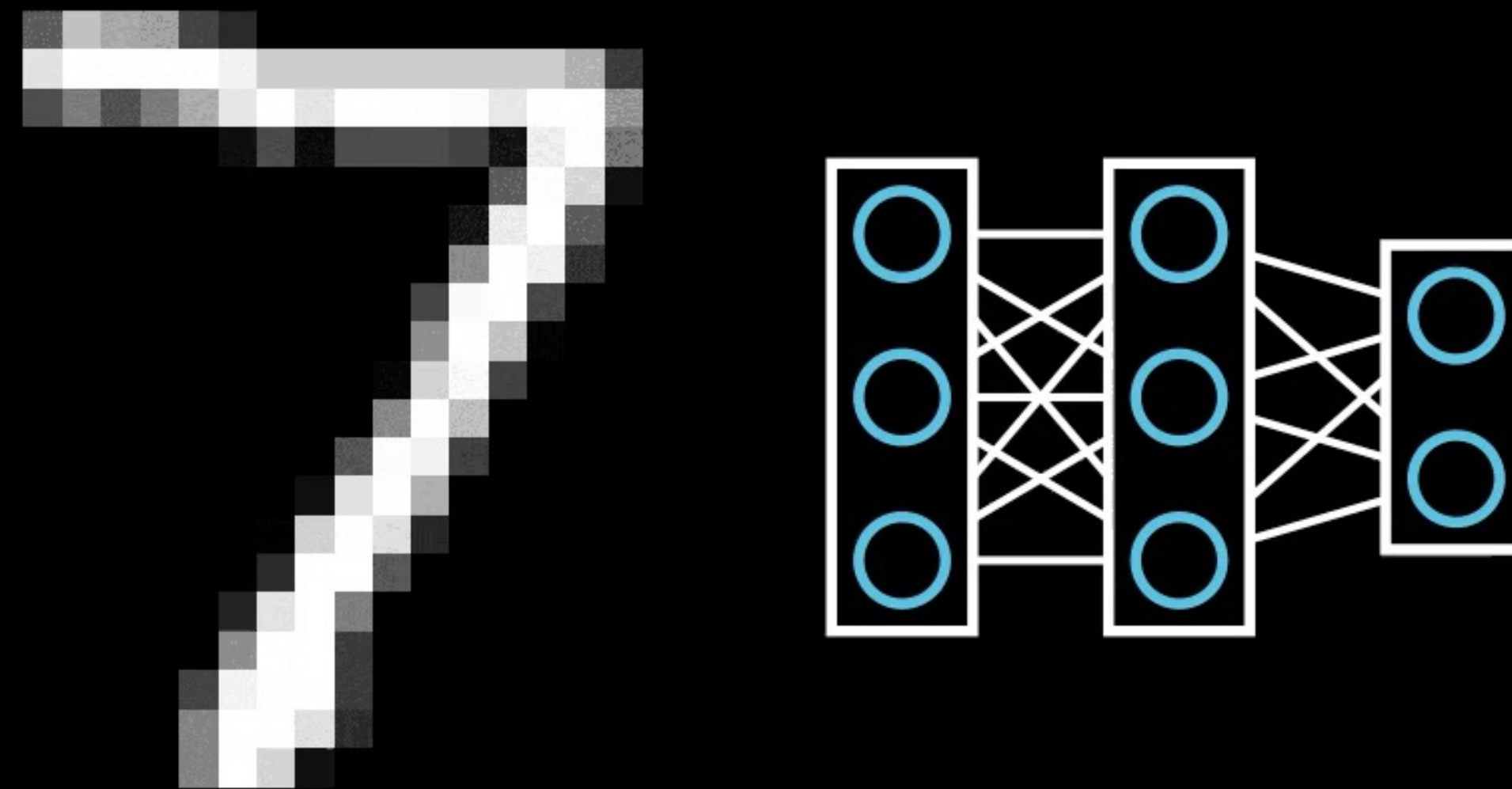# Overview

1. Images: Convolutional Neural Networks

2. Sequences: RNNs

3. Transformers

4. Current developments

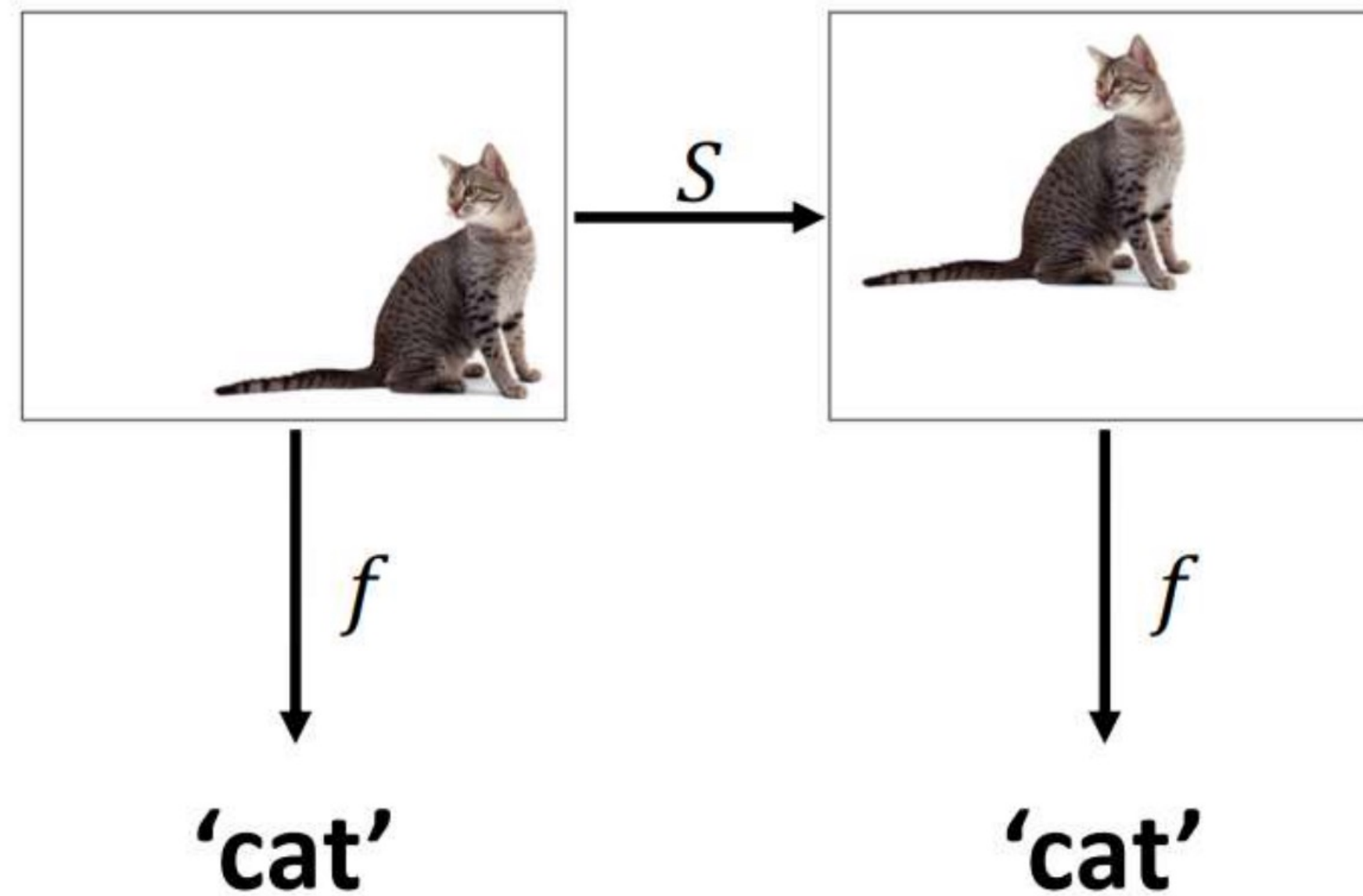# 1. Convolutional Neural Networks

# How to deal with images

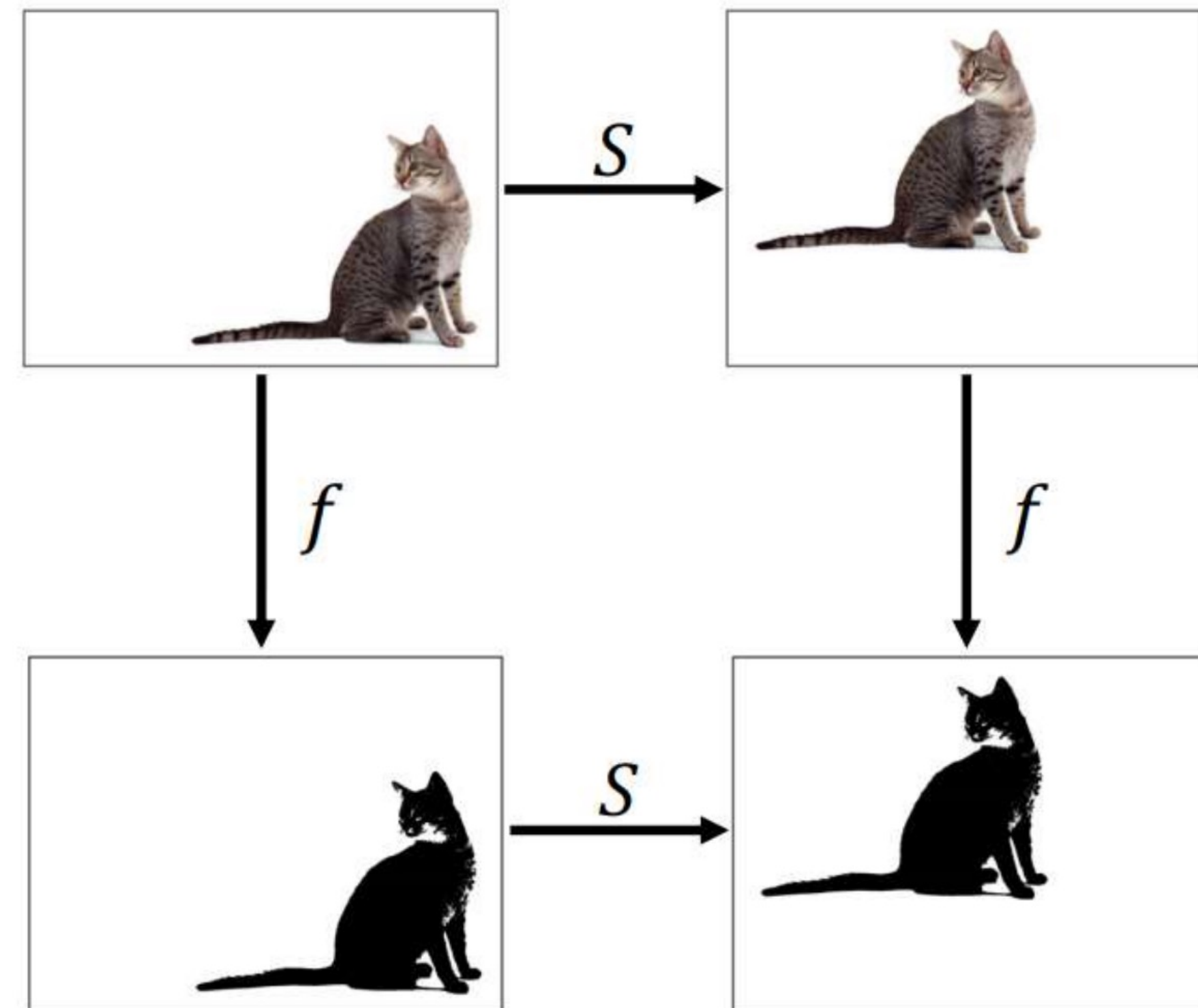**Naive approach: unroll them and passt them into an MLP**

# Inductive Bias: Translational In-/Equivariance

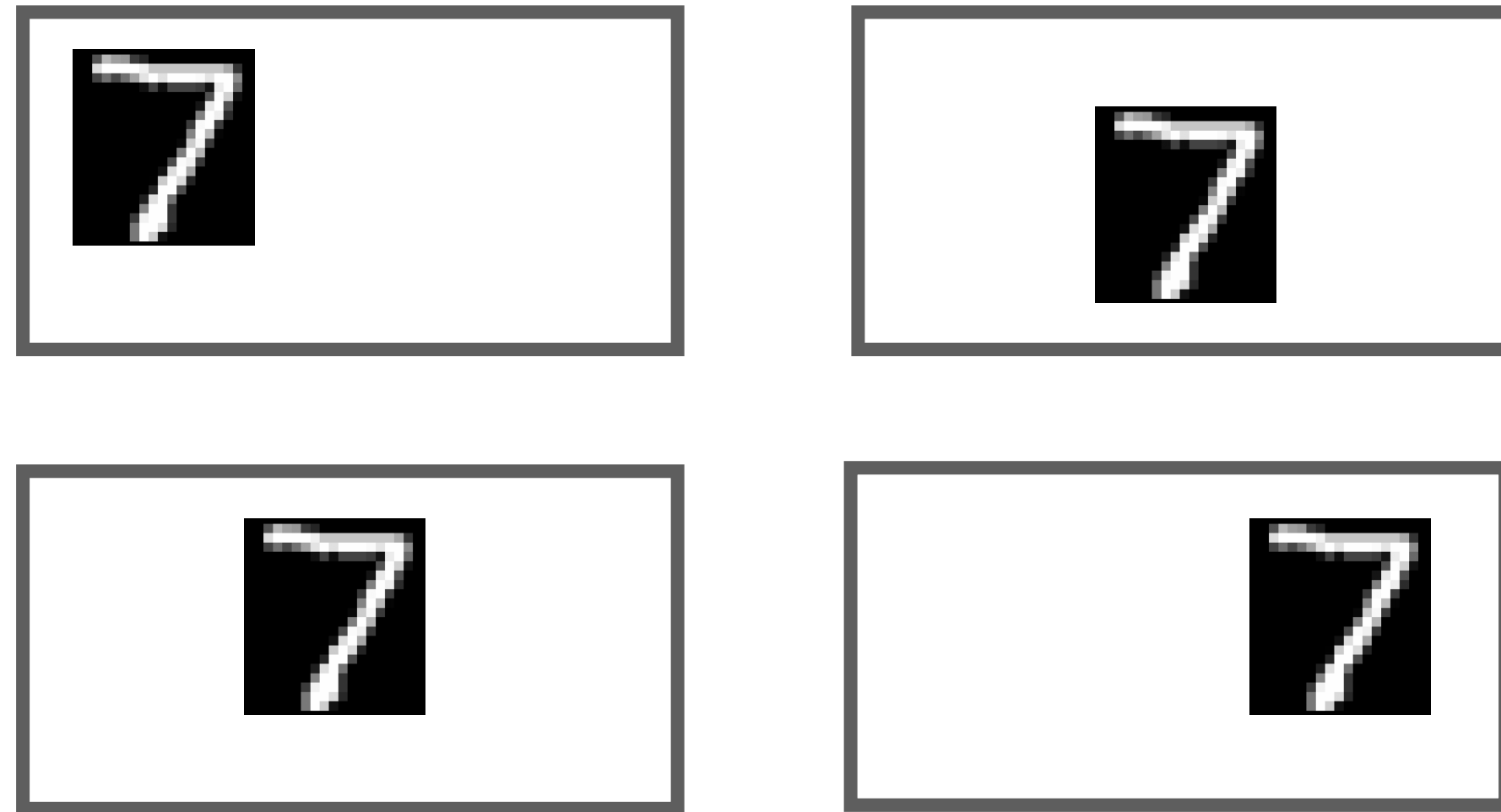## Leverage the symmetry of your data

# Why leverage symmetries?

**We need more data = our network is more efficient!**
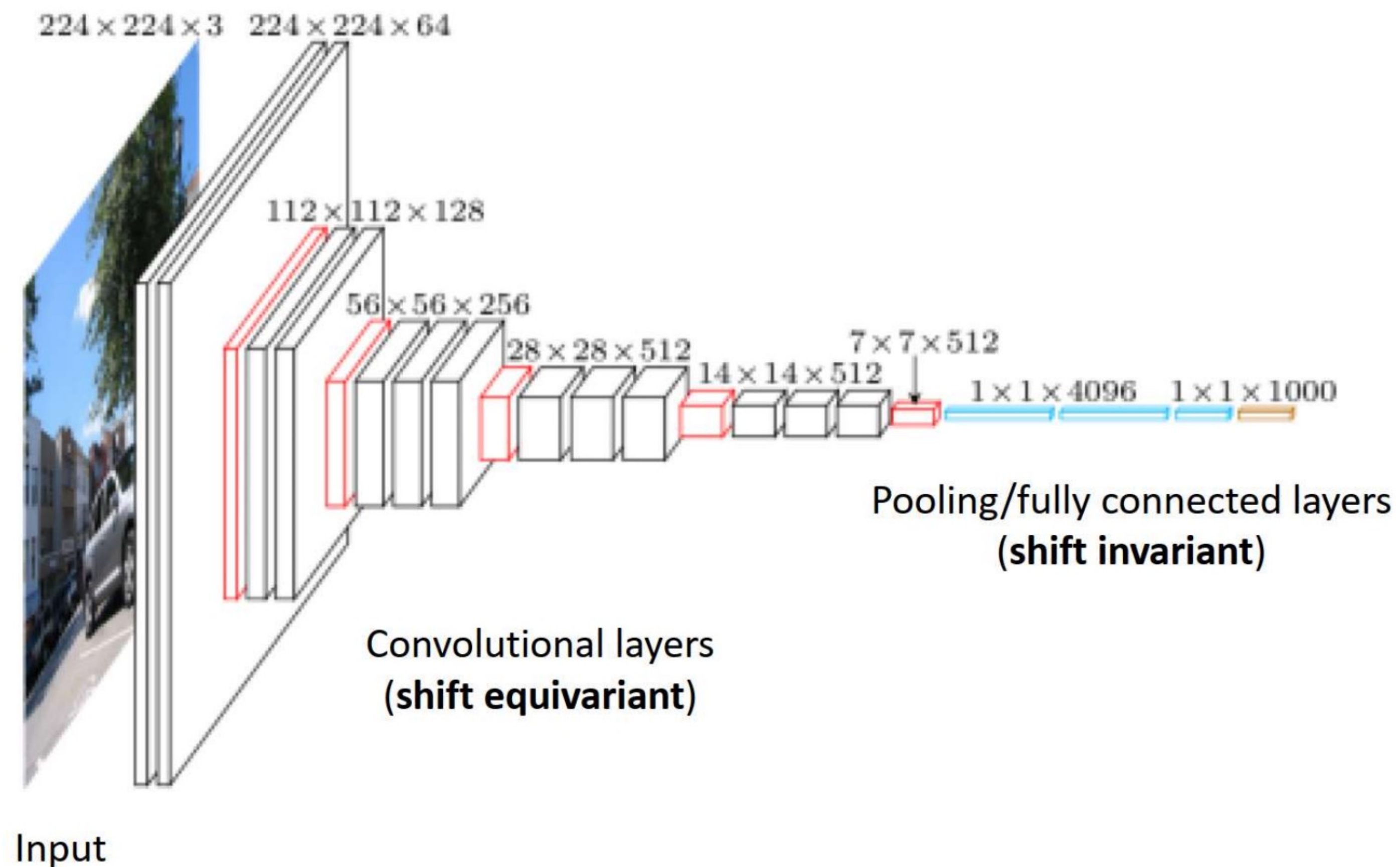
**Training without translational symmetry**



**Training with translational symmetry**

# How do we do this in practice?
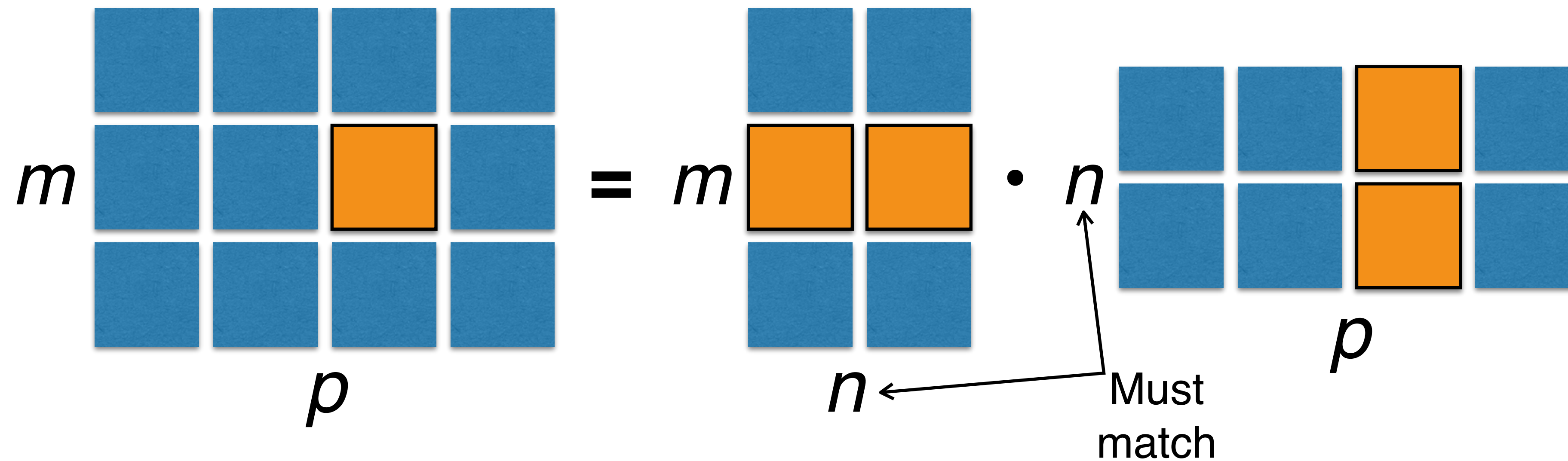
**Implement neural network layers that respect these symmetries**

# Convolutional Layers
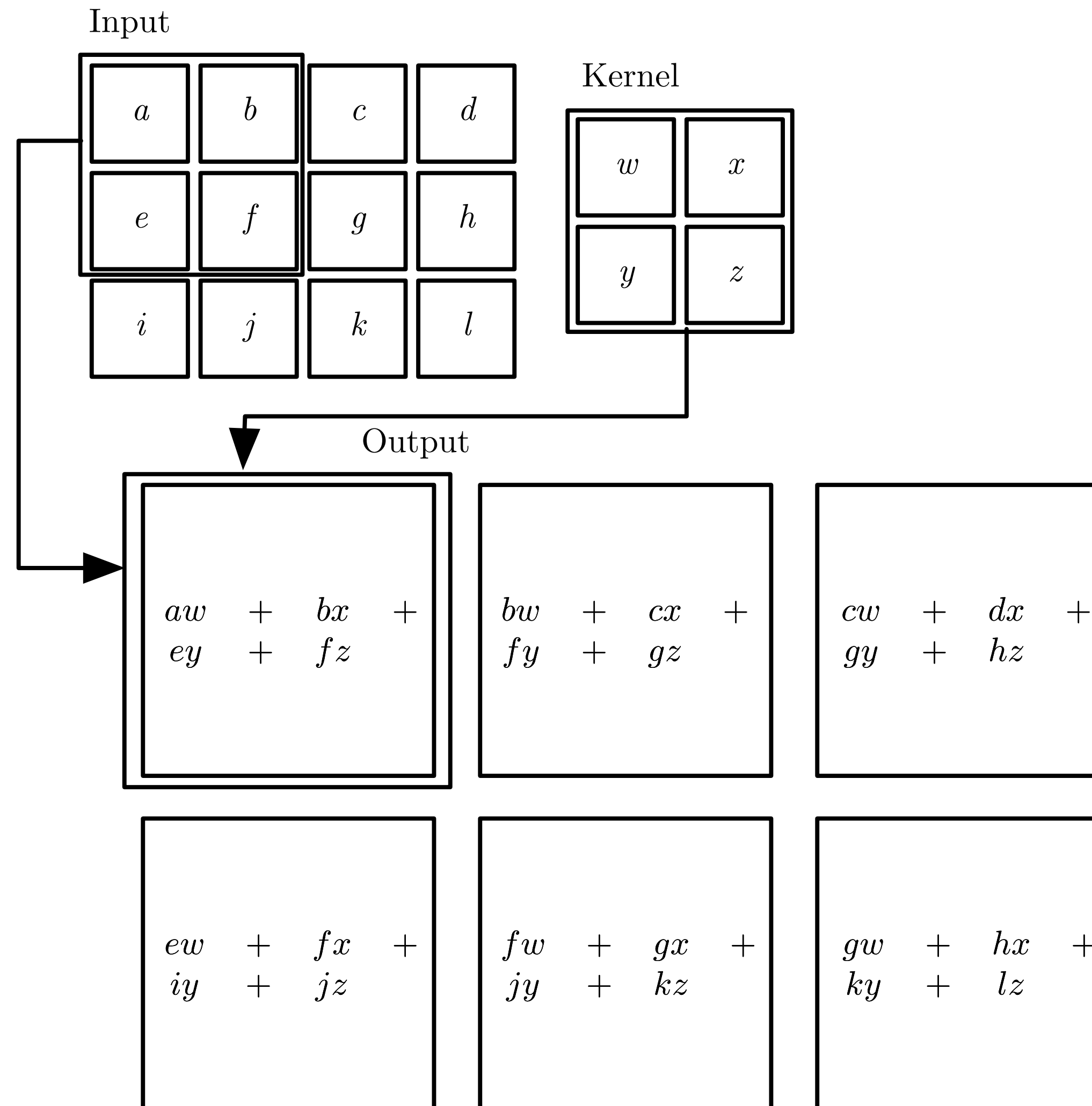
## Reminder: Matrix multiplication

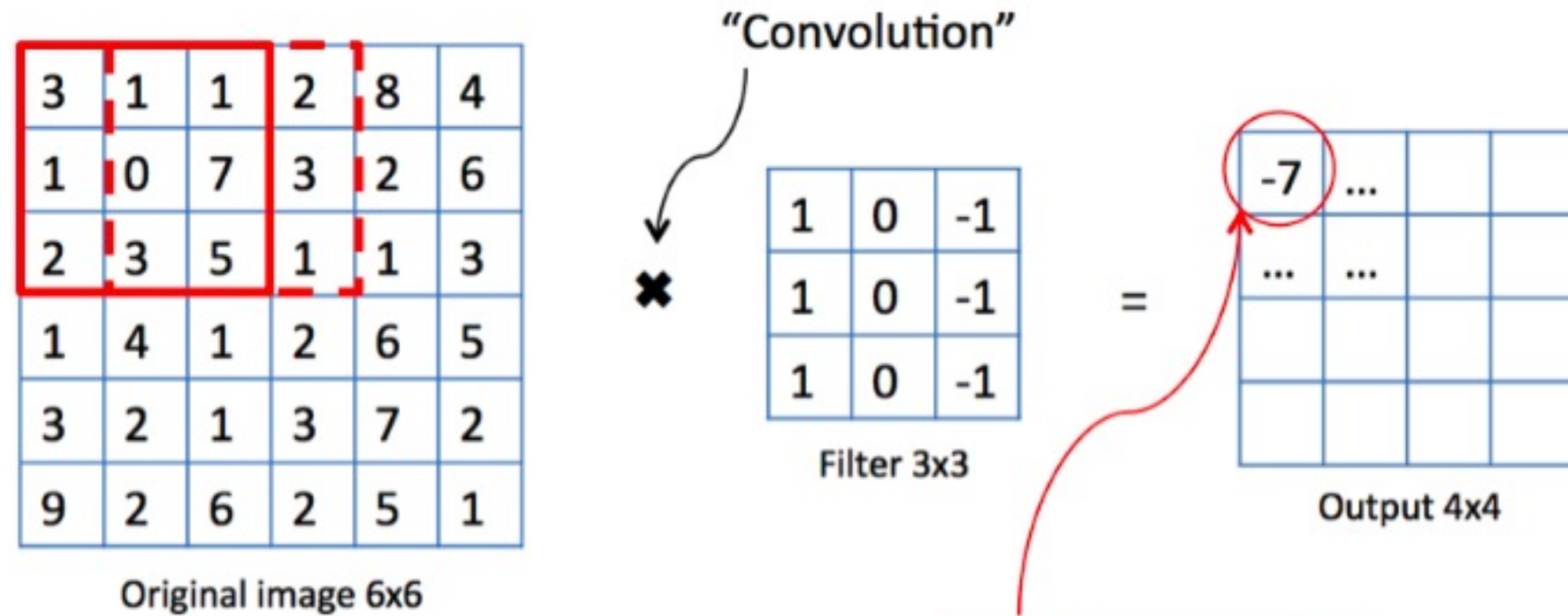$$C = AB. \qquad\qquad (2.4)$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}. \qquad\qquad (2.5)$$

# Convolutional Layers

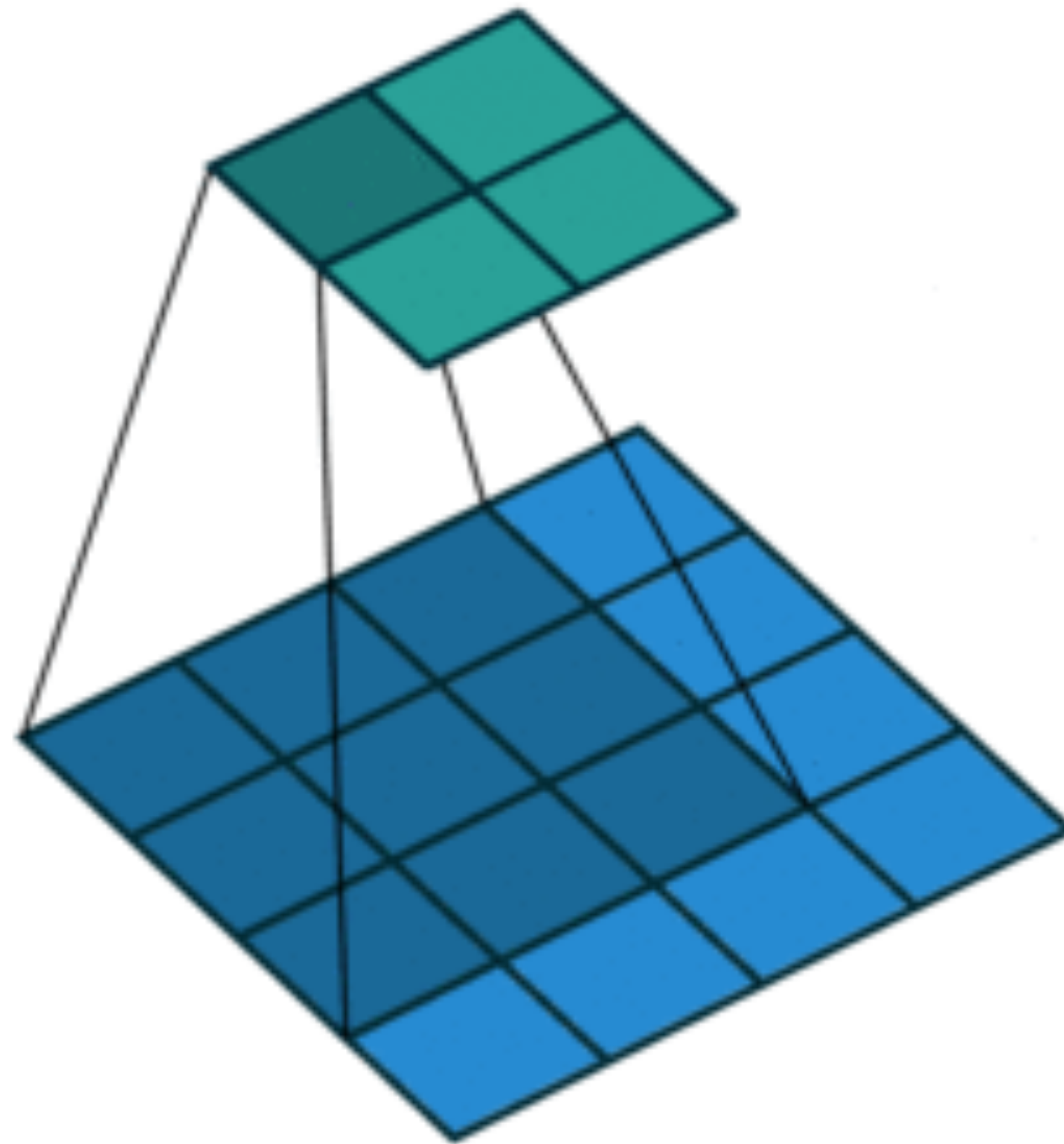## The weights are in the kernel



Input

Kernel

Output

# Convolutional Layers

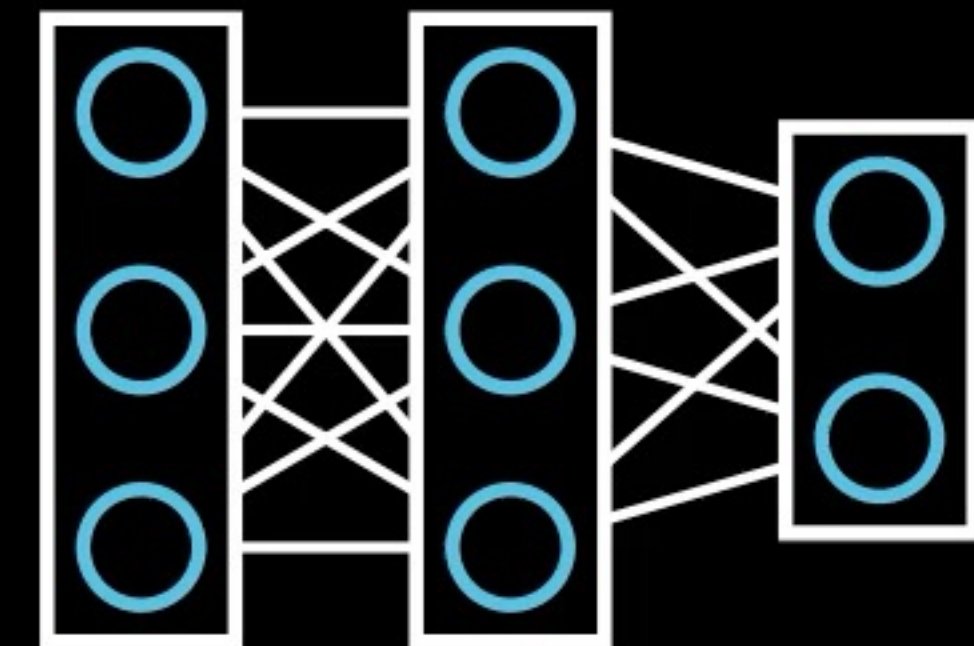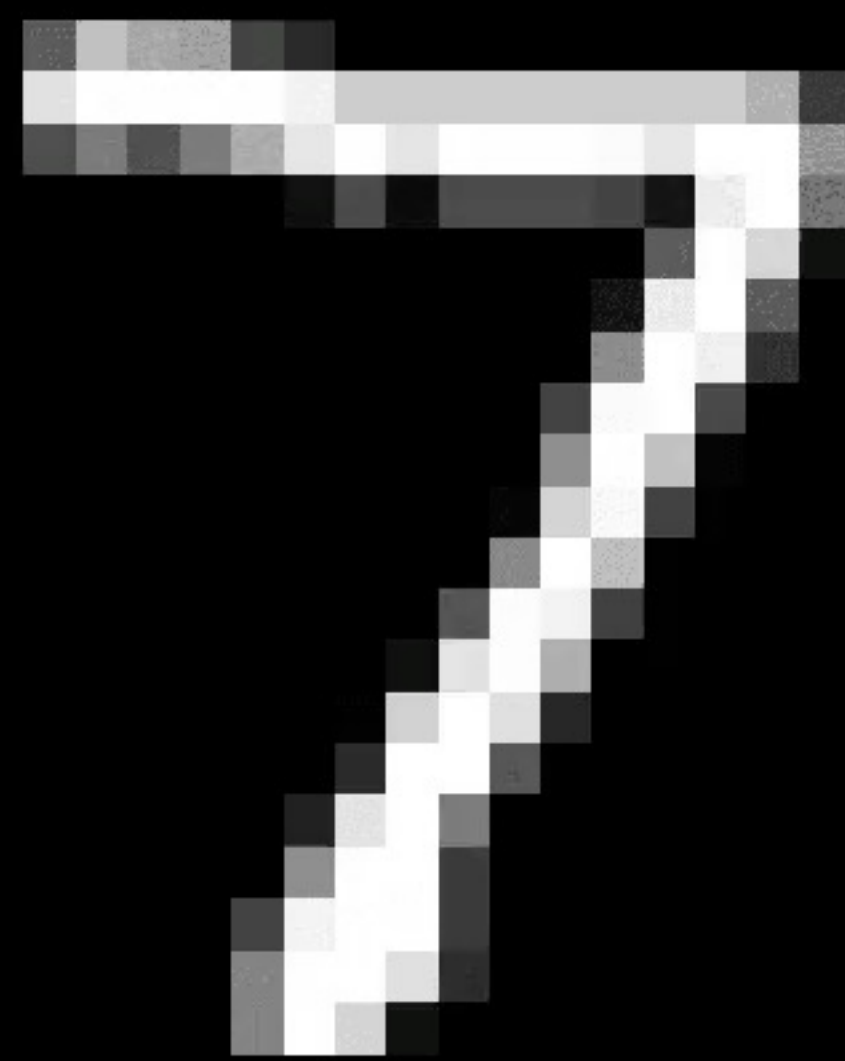**Convolution = Repeated Matrix Multiplication**

# How can I imagine that?

## Sliding the kernel over the image
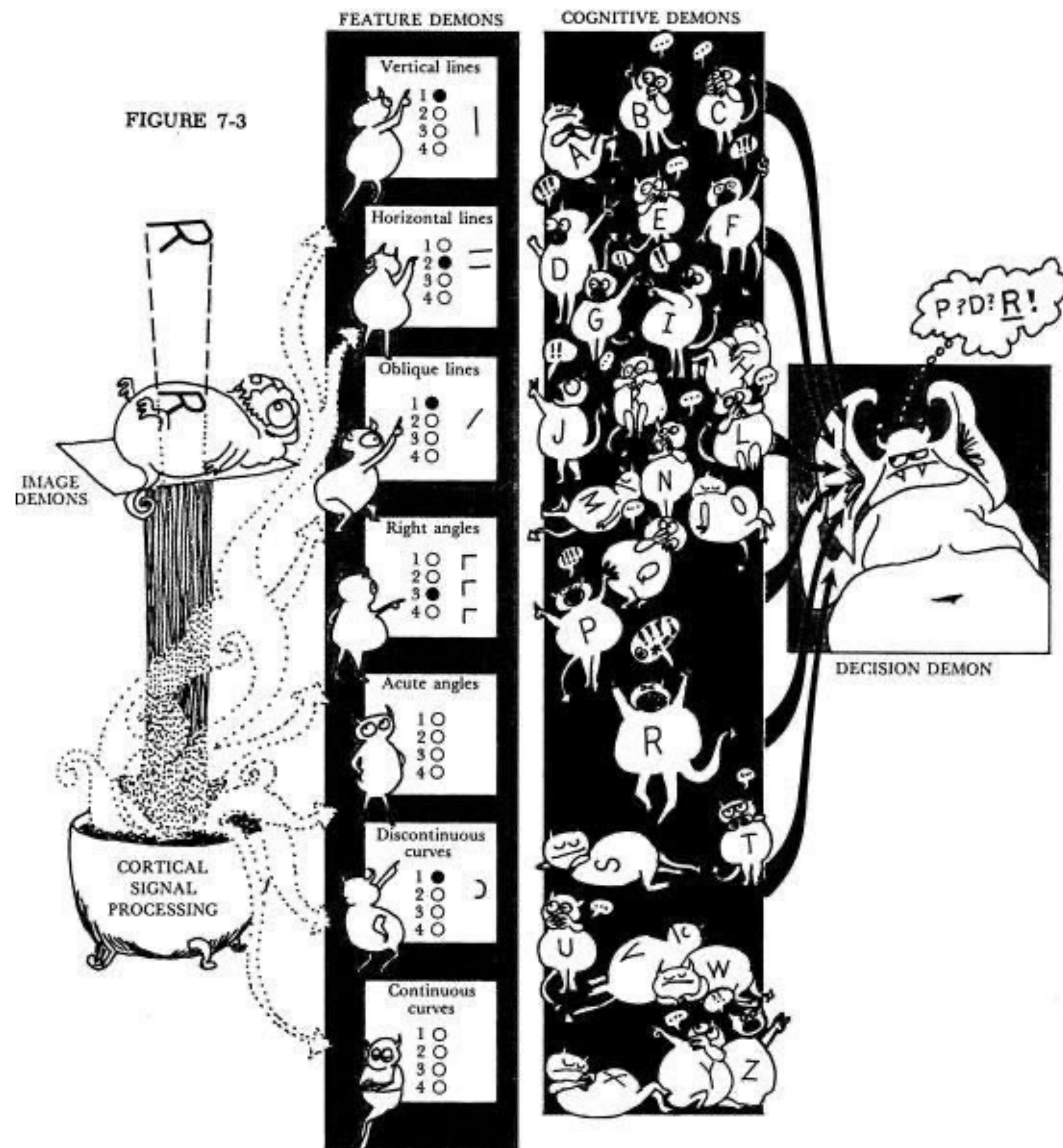
# How can I imagine that?

**Multiple Kernels allow detecting multiple features**

# Pattern Recognition all over again
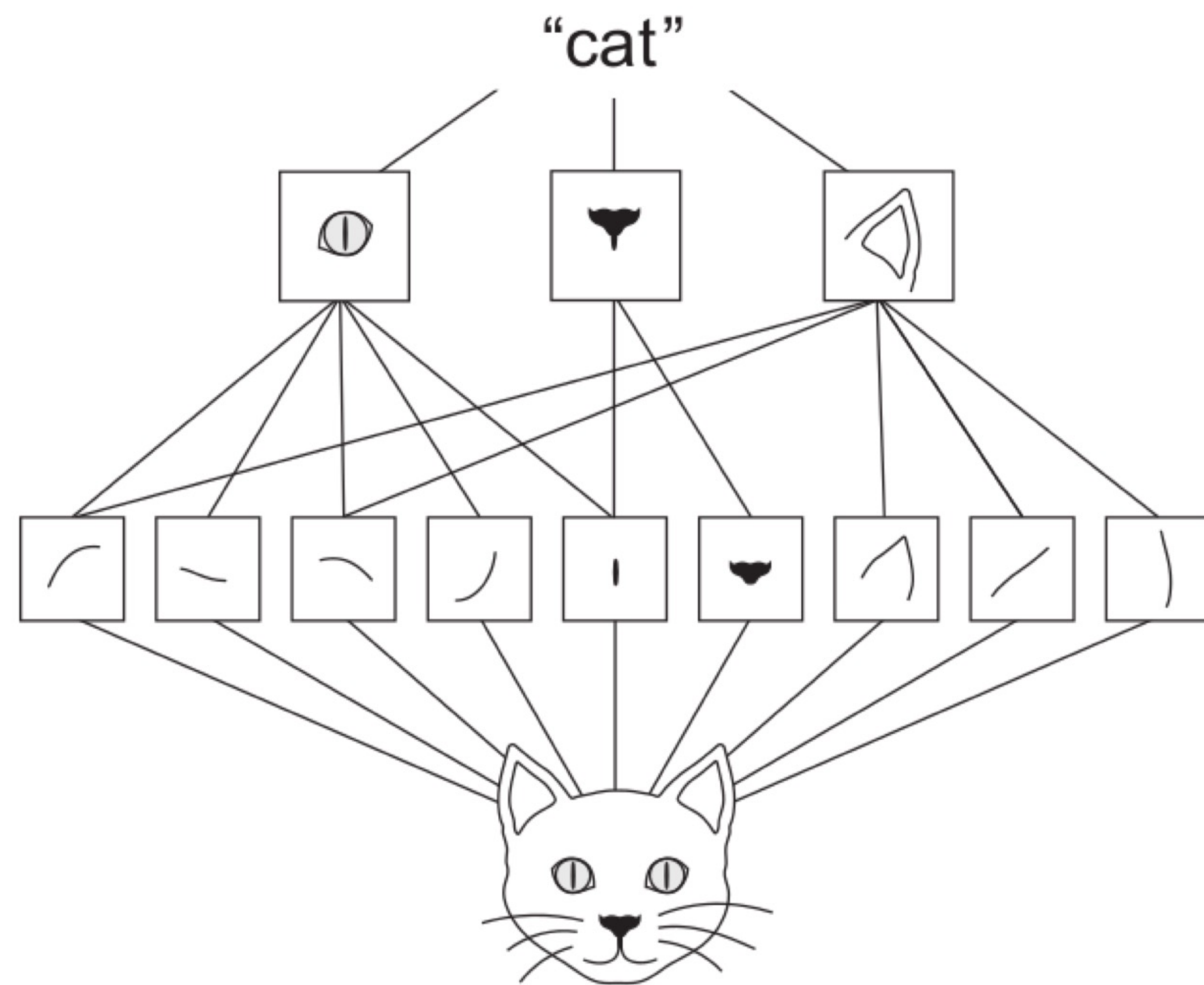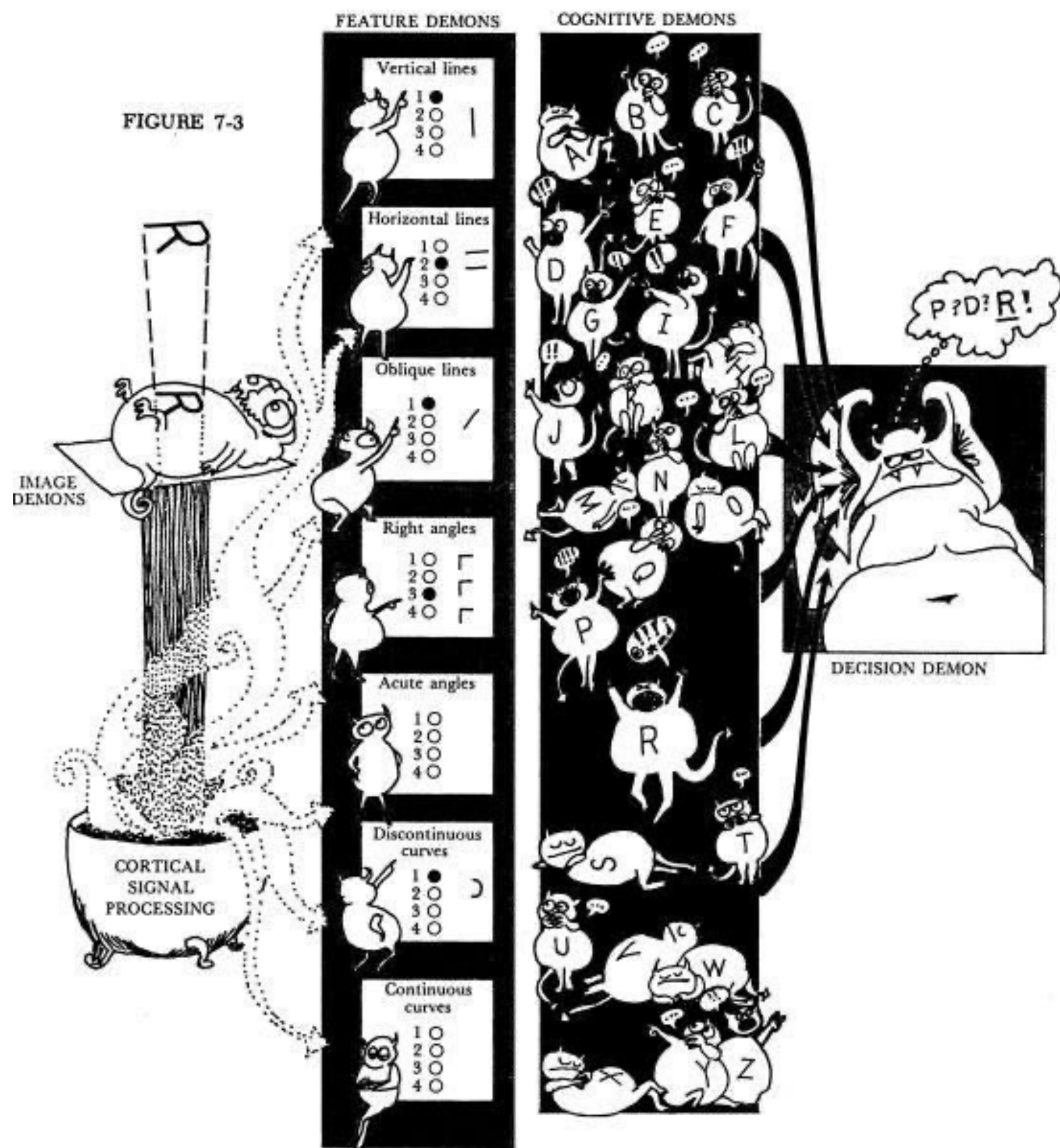
## This time adjusted to the image case

# Pattern Recognition all over again

## This time adjusted to the image case

# Pattern Recognition all over again

## Look at it yourself!

[Google Brain: Feature Visualisation](#)                    [OpenAI: Microscope](#)



**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.
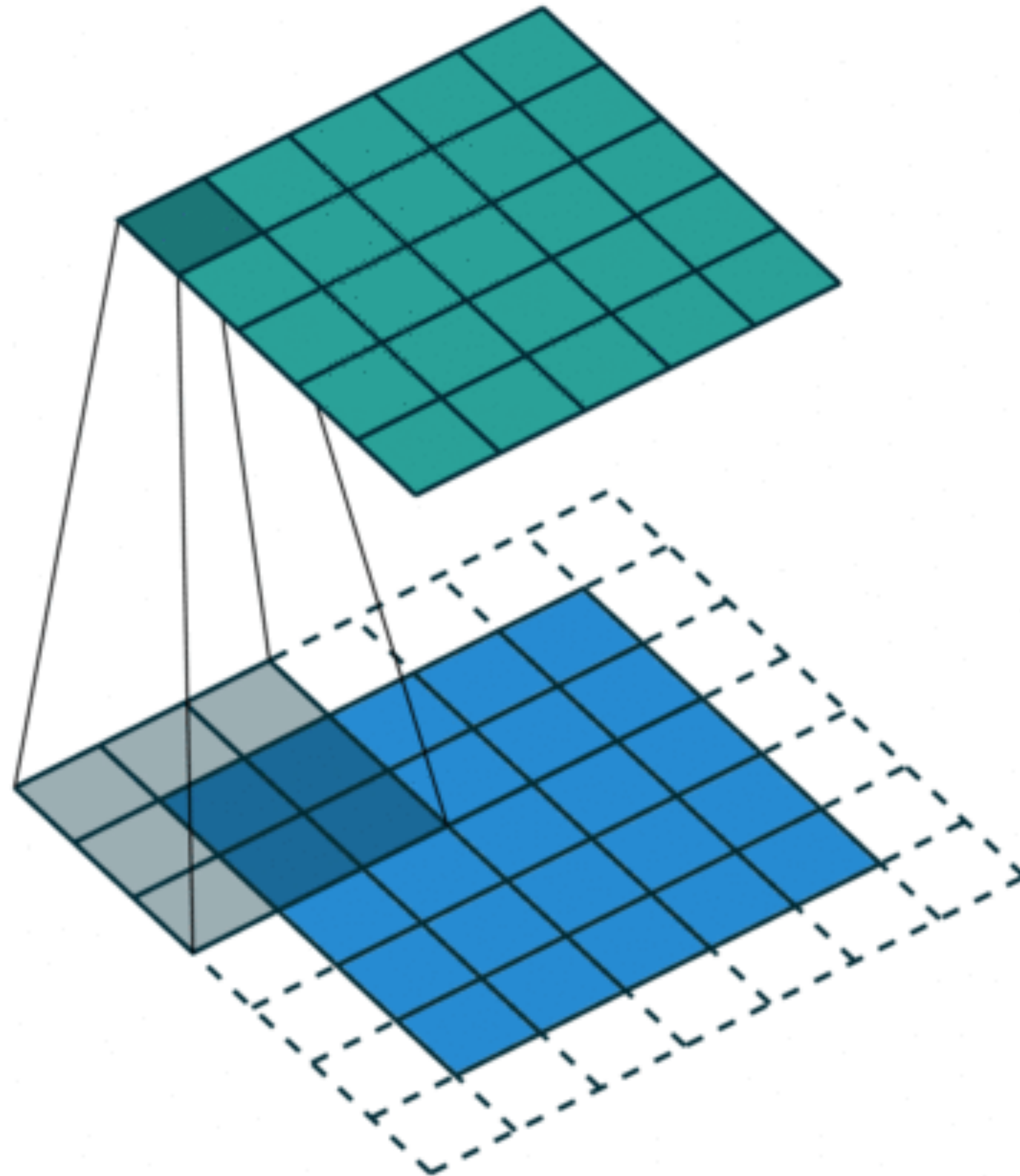
Baseball—or stripes?
*mixed4a, Unit 6*

# Avoid reducing size with padding

**Different ways to pad (zero-pad, mean-pad, …)**



Vincent Dumoulin, Francesco Visin - A guide to convolution arithmetic for deep learning

# Make bigger jumps with strides

Vincent Dumoulin, Francesco Visin - A guide to convolution arithmetic for deep learning

# Pooling: Shift-invariant operation

## Reduce size, but no learning involved
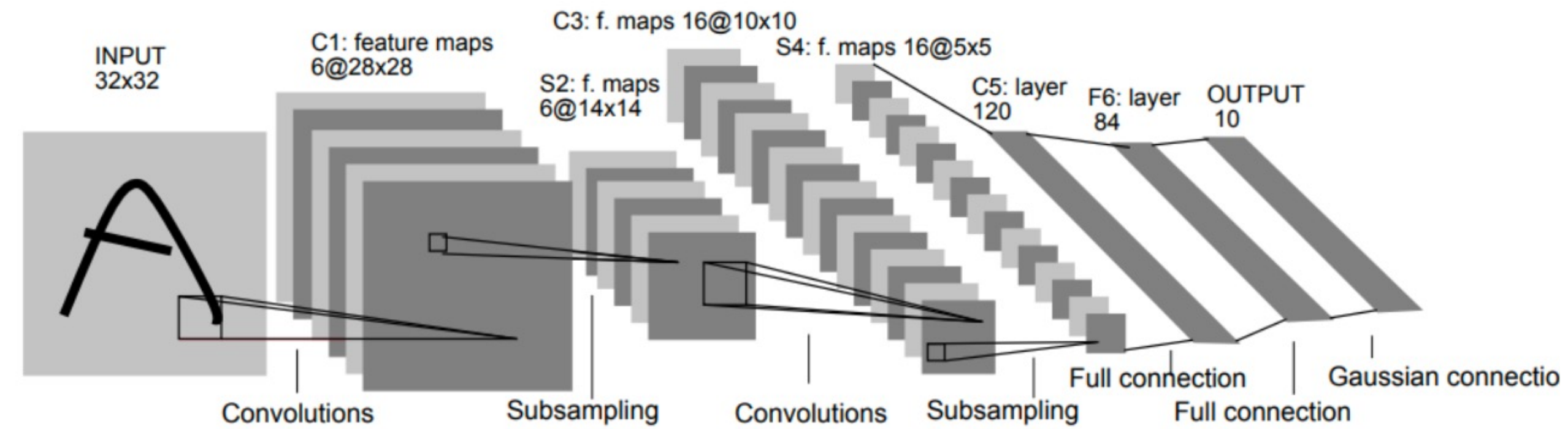


2x2 Max Pooling

# Putting things together: A full CNN
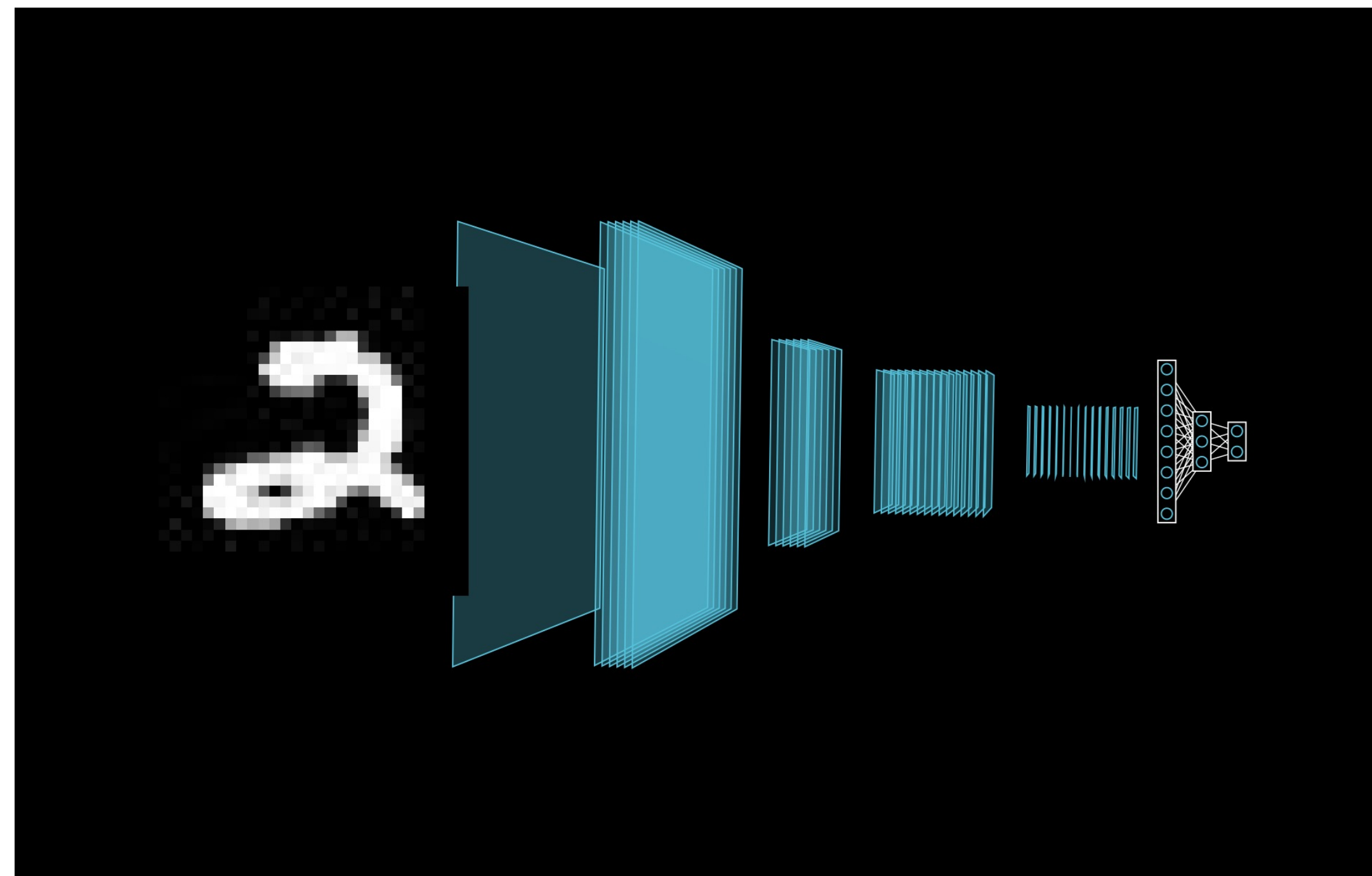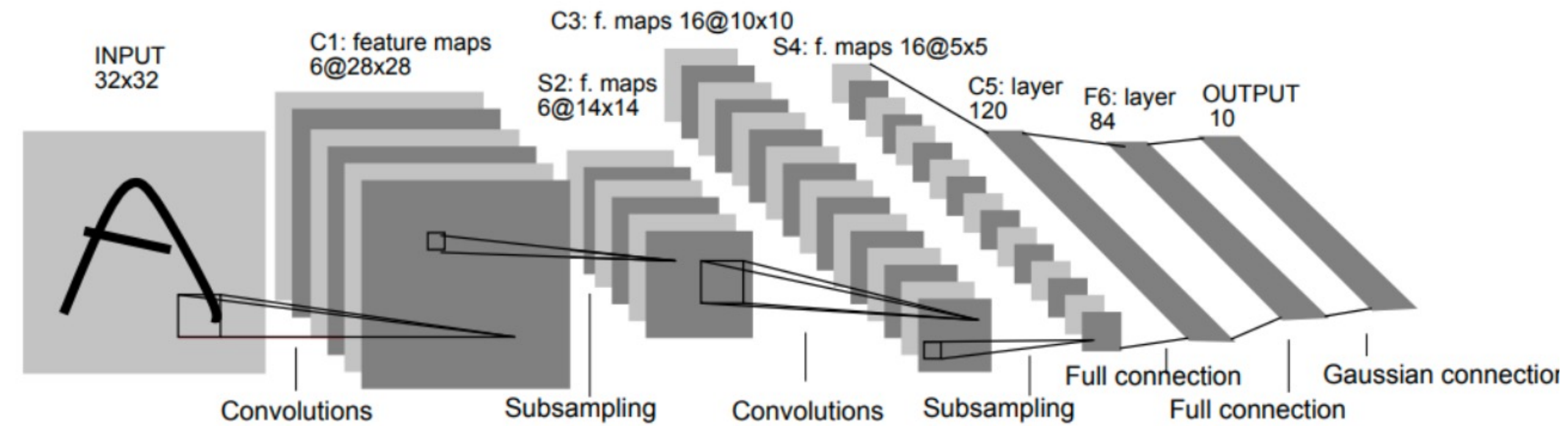
**Conv. Layers -> Pooling -> FC Layers**

# LeNet (1998): CNNs become a thing

**Exactly what we discussed, just bigger**

# LeNet (1998): CNNs become a thing

**Exactly what we discussed, just bigger**

# 2. RNNs

# The classic landscape

## One architecture per community

Lucas Beyer, Transformer Talk

# The transformer's takeover

## One community at a time



Computer Vision

Natural Lang. Proc.

Reinf. Learning

Speech

Translation

Graphs/Science

transformer image source: "Attention Is All You Need" paper

# How to deal with sequential data?

## You can only look into the past, not into the future

# How to deal with sequential data?

## You can only look into the past, not into the future

# How to deal with sequential data?

## You can only look into the past, not into the future



Global Average Temperature Chan

# How to deal with sequential data?

## You can only look into the past, not into the future



Global Average Temperature Change

# RNN: Recurrent Neural Networks

**Making predictions with respect to time**

# RNN: Recurrent Neural Networks

**Making predictions with respect to time**

# Different tasks, different architectures

**Making predictions with respect to time**



one to one        one to many        many to one        many to many        many to many

# RNNs have problems

## Vanishing Gradients cause short context lengths



**Vanishing Gradient:** where the contribution from the earlier steps becomes insignificant in the gradient for the vanilla RNN unit.

# RNNs have problems

## Vanishing Gradients cause short context lengths



**Vanishing Gradient:** where the contribution from the earlier steps becomes insignificant in the gradient for the vanilla RNN unit.

# RNNs have problems

## Vanishing Gradients cause short context lengths



**Vanishing Gradient:** where the contribution from the earlier steps becomes insignificant in the gradient for the vanilla RNN unit.

# RNN variants tackle vanishing gradients

## Still, the problem of limited context length remains

# RNN variants tackle vanishing gradients

**Still, the problem of limited context length remains**

# RNNs have problems

**Vanishing Gradients cause short context lengths**



# Visualizing memorization in RNNs

Inspecting gradient magnitudes in context can be a powerful tool to see when recurrent units use short-term or long-term contextual understanding.

context the formal study of grammar is an important part of education — Nested LSTM

context the formal study of grammar is an important part of education — LSTM

context the formal study of grammar is an important part of education — GRU

# RNNs have other problems, too

**No parallelisation possible**



many to many

# 3. Transformers

# Transformers to the rescue

**Parallel instead of sequential encoding with <u>attention</u>**



RNN based Encoder

RNN → RNN → RNN → RNN →

The Cat Is Black

Transformer's Encoder

Encoder
(Transformer)

The Cat Is Black

# What is Attention?

**Allowing every word to be influenced by any other word**



Bahdanau et al, 2014, *arxiv*

# What is Attention?

## Apparently it is all you need

---

### Attention Is All You Need

---

**Ashish Vaswani*** 
Google Brain 
avaswani@google.com

**Noam Shazeer*** 
Google Brain 
noam@google.com

**Niki Parmar*** 
Google Research 
nikip@google.com

**Jakob Uszkoreit*** 
Google Research 
usz@google.com

**Llion Jones*** 
Google Research 
llion@google.com

**Aidan N. Gomez*** [†] 
University of Toronto 
aidan@cs.toronto.edu

**Łukasz Kaiser*** 
Google Brain 
lukaszkaiser@google.com

**Illia Polosukhin*** [‡] 
illia.polosukhin@gmail.com

# The Transformer

## Not as scary as it looks like

# The Transformer

## Not as scary as it looks like

# Input Embedding

**Our computer does not understand English**

Vocabulary

One-hot vectors

# Input Embedding

**From one-hot encodings to word embeddings**

One-hot vectors

Word embeddings

# Input Embedding

**Play with a few word embeddings yourself**

https://lamyiowce.github.io/word2viz/

https://ronxin.github.io/wevi/

http://projector.tensorflow.org/

# Positional Embedding

**We must tell our computer what comes first and what later**



EMBEDDINGS    x₁ [ ][ ][ ][ ]        x₂ [ ][ ][ ][ ]        x₃ [ ][ ][ ][ ]

INPUT              Je                    suis                   étudiant

# Positional Embedding

**We must tell our computer what comes first and what later**

# Positional Embedding

## We must tell our computer what comes first and what later



$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

# Attention

**Looking at everyone around you to determine your update**

- Input: sequence of tensors

  $$x_1, x_2, \ldots x_t$$

# Attention

**Looking at everyone around you to determine your update**

- Input: sequence of tensors

$$x_1, x_2, \dots x_t$$

- Output: sequence of tensors, each one a weighted sum of the input sequence

$$y_1, y_2, \dots, y_t$$

$$y_i = \sum_j w_{ij} x_j$$

# Attention

**Looking at everyone around you to determine your update**

- Input: sequence of tensors

$$x_1, x_2, \dots x_t$$

- Output: sequence of tensors, each one a weighted sum of the input sequence

$$y_1, y_2, \dots, y_t$$

$$y_i = \sum_j w_{ij} x_j$$

  – weight is just a dot product $\quad w'_{ij} = x_i^{\mathsf{T}} x_j$

# Attention

## Looking at everyone around you to determine your update

- Input: sequence of tensors
    $$x_1, x_2, \ldots x_t$$

- Output: sequence of tensors, each one a weighted sum of the input sequence
    $$y_1, y_2, \ldots, y_t$$

    $$y_i = \sum_j w_{ij} x_j$$

  – weight is just a dot product $\quad w'_{ij} = x_i{}^\top x_j$

  – make it sum to 1 $\qquad w_{ij} = \dfrac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$

# Attention

## Looking at everyone around you to determine your update

# Attention

## Learning the weights

### Query, Key, Value

- Every input vector x_i is used in 3 ways:

  – Query

  – Key

  – Value

# Attention

## Learning the weights

### Query, Key, Value

- Every input vector x_i is used in 3 ways:

  – Query  **What am I looking for?**

  – Key  **What do I have?**

  – Value  **What do I reveal/give to others?**

# Attention

## Learning the weights

- We can process each input vector to fulfill the three roles with matrix multiplication

- Learning the matrices → learning attention

**What am I looking for?**     **What do I have?**     **What do I reveal/give to others?**

$$q_i = W_q x_i \qquad k_i = W_k x_i \qquad v_i = W_v x_i$$

$$w'_{ij} = q_i^\top k_j$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

$$y_i = \sum_j w_{ij} v_j \,.$$

# Imagine you are in a library

## How do you answer a question you have?



- Query — The question you have

- Key — The titles books have on their spines

- Value — Information the book contains

# Multi-head attention

## Looking at everyone around you to determine your update

- Multiple "heads" of attention just means learning different sets of W_q, W_k, and W_v matrices simultaneously.

- Implemented as just a single matrix...

# Multi-head attention

## Looking at everyone around you to determine your update

- Multiple "heads" of attention just means learning different sets of W_q, W_k, and W_v matrices simultaneously.

- Implemented as just a single matrix...



$x_i$

$W_q^1$    $q_i^1$

$W_q^2$    $q_i^2$

$W_q^3$    $q_i^3$

# Multi-head attention

**Different heads attend to different parts in a sentence**
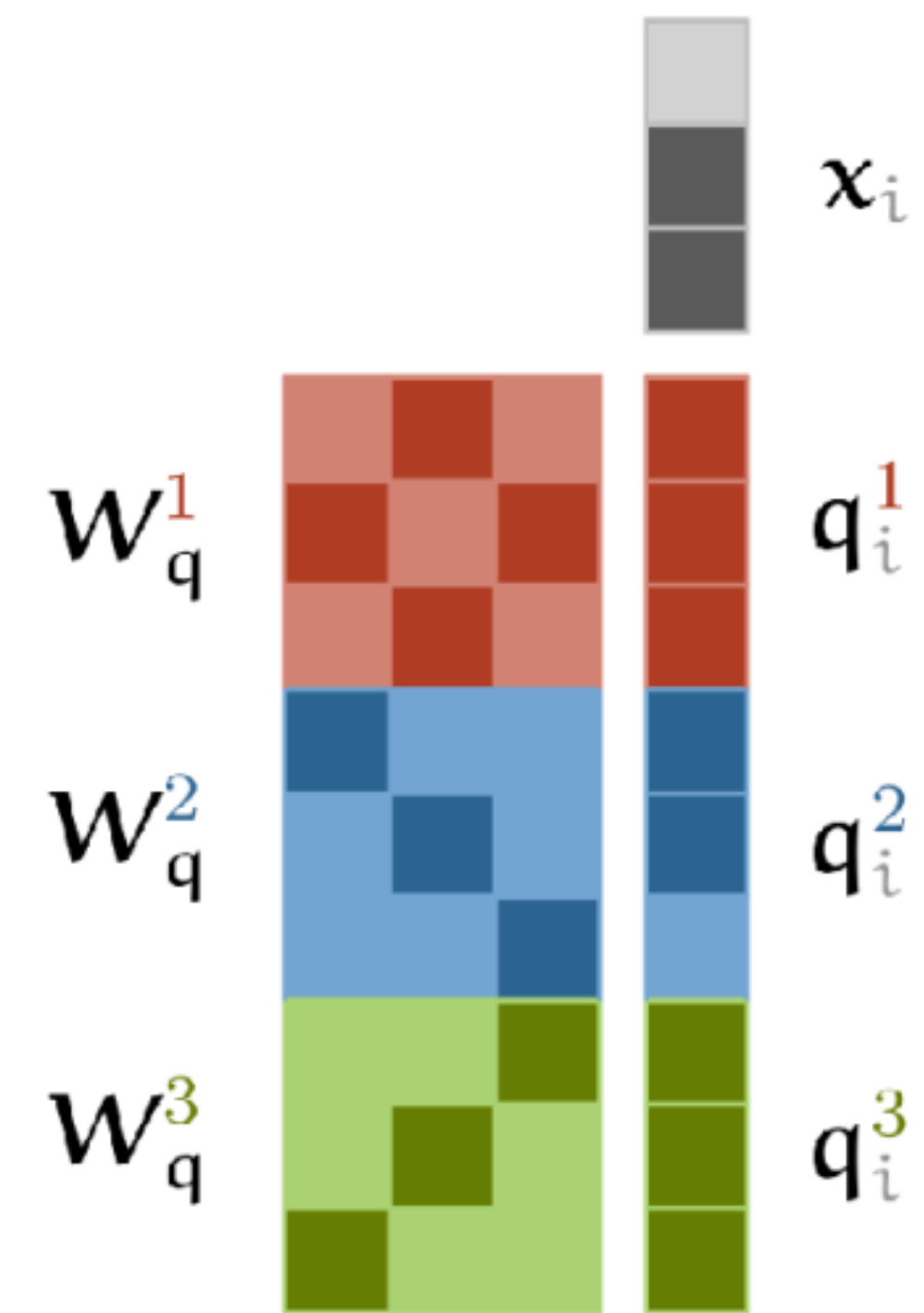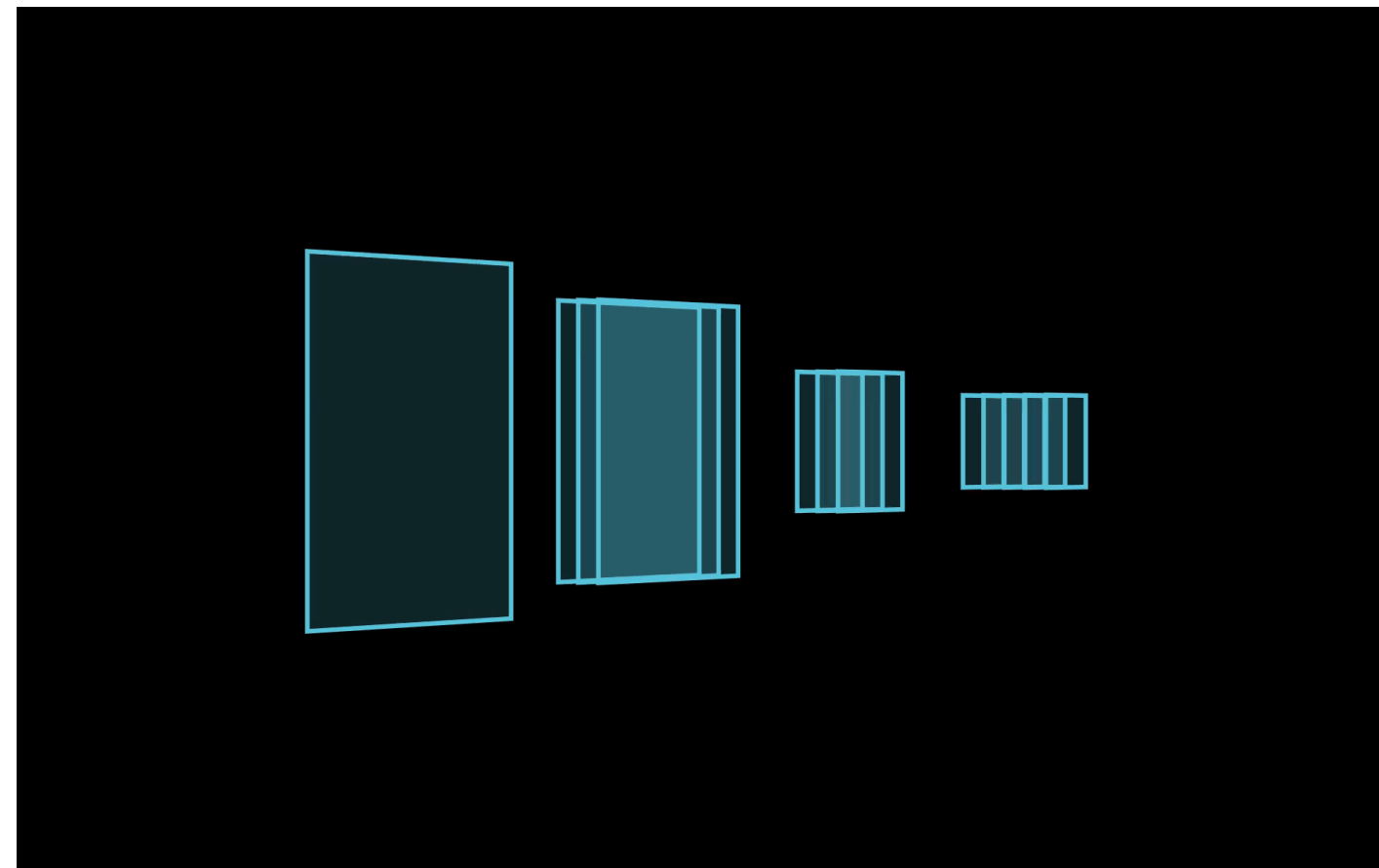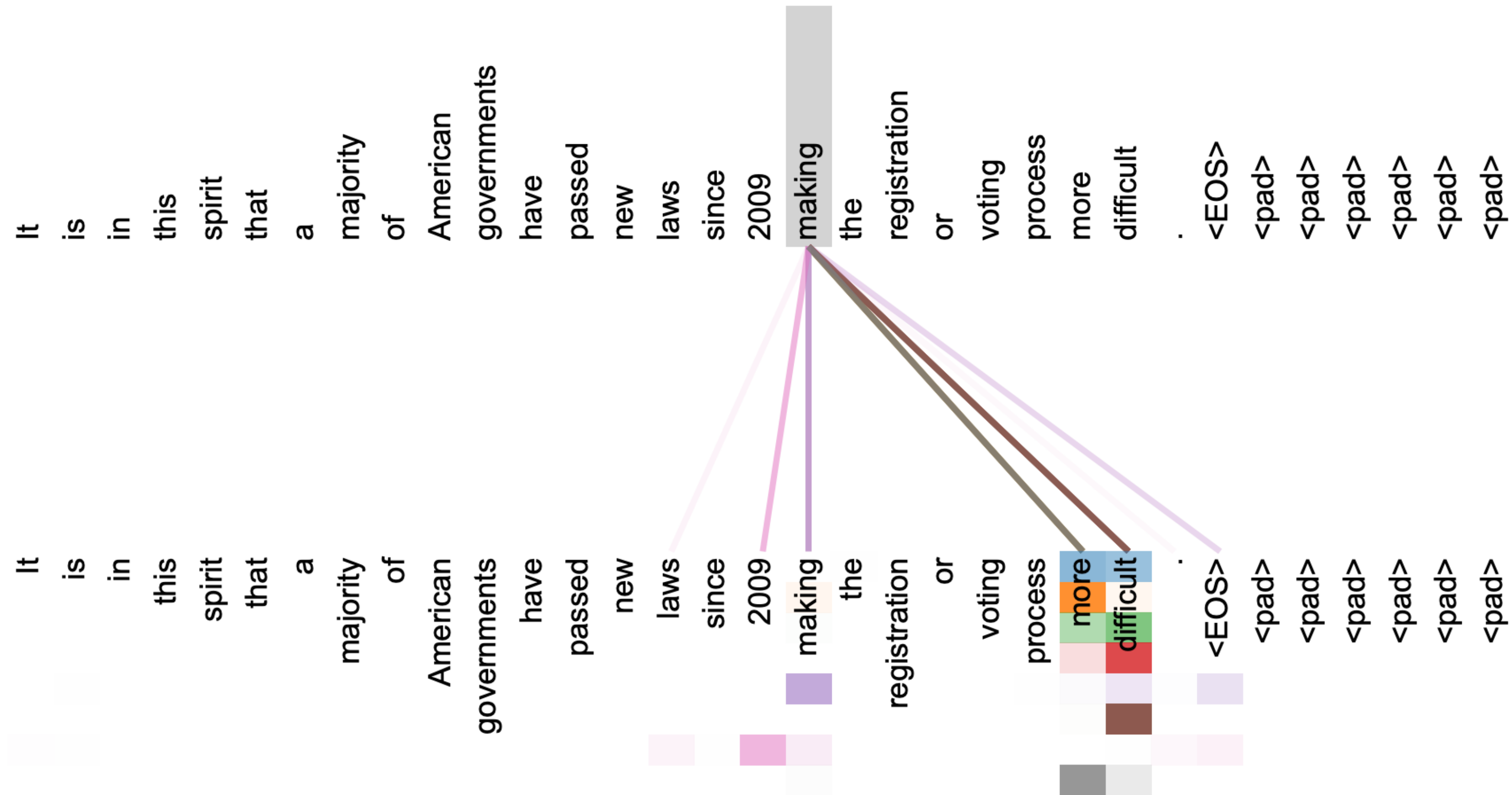


Attention Visualizations

# Multi-head attention

## The same applies for proteins



(a) Attention in head 12-4, which targets amino acid pairs that are close in physical space (see inset subsequence 117D-157I) but lie apart in the sequence. Example is a *de novo* designed TIM-barrel (5BVL) with characteristic symmetry.

(b) Attention in head 7-1, which targets binding sites, a key functional component of proteins. Example is HIV-1 protease (7HVP). The primary location receiving attention is 27G, a binding site for protease inhibitor small-molecule drugs.

Vig et al, ILCR 2021

# Layer Normalization

**Standardize means and stds of input vectors**

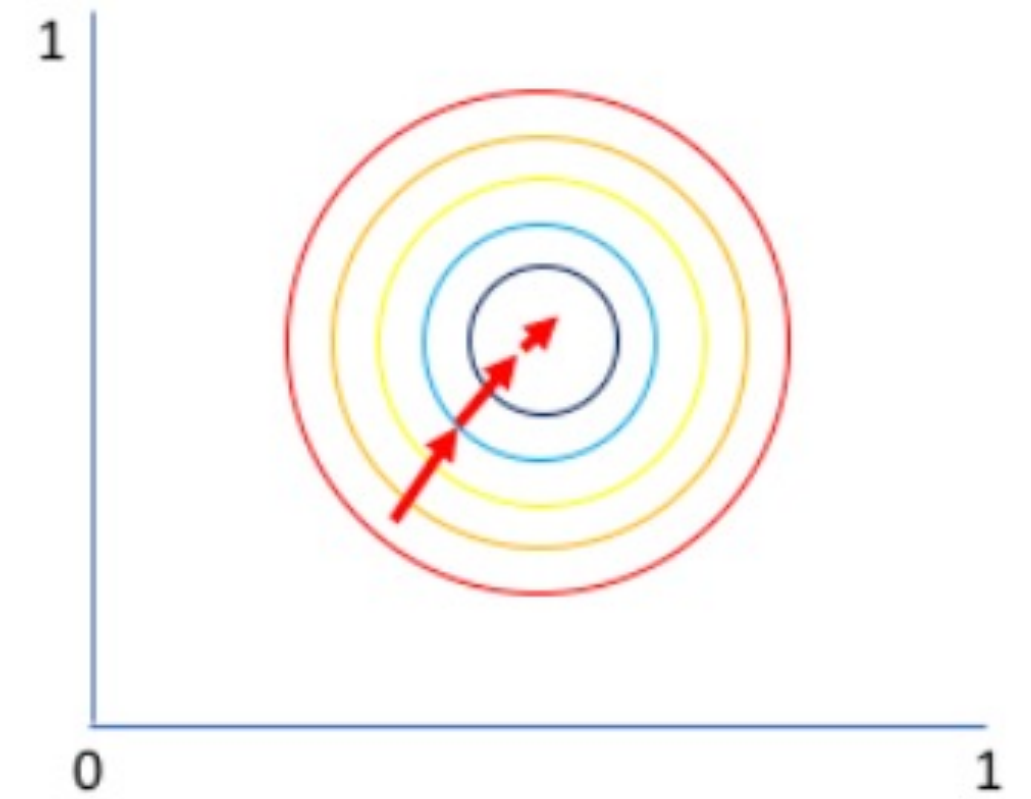- Neural net layers work best when input vectors have uniform mean and std in each dimension

- As inputs flow through the network, means and std's get blown out.

- Layer Normalization is a hack to reset things to where we want them in between layers.



Both parameters can be updated in equal proportions



Gradient of larger parameter dominates the update

# The Transformer

## Not as scary as it looks like

# Many good blogs about Transformers
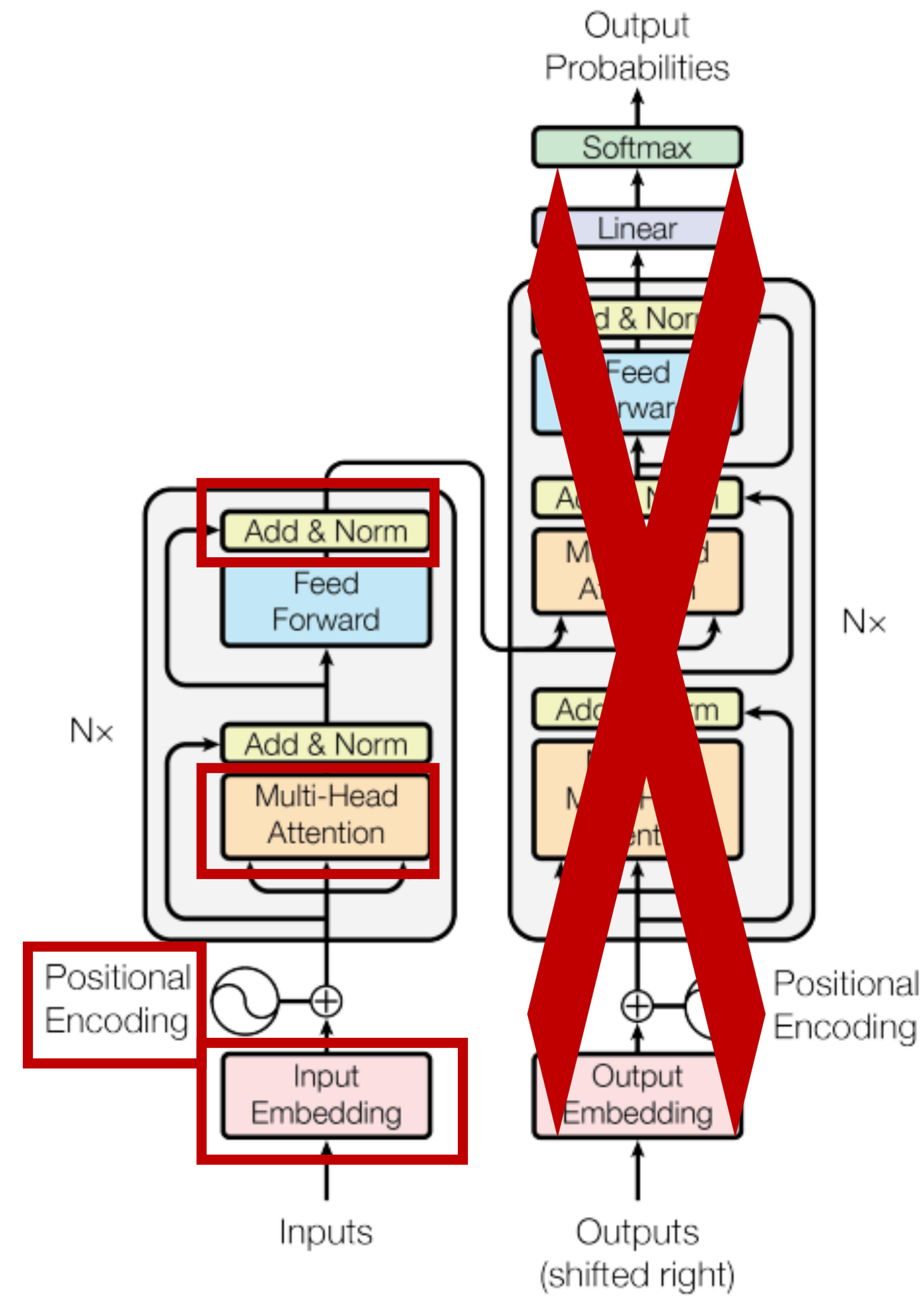
**I leave it to you to choose the ones you like best**

1. [The Illustrated Transformer ](#)(Pictures)
2. [The Annotated Transformer ](#)(Code)
3. [Transformers from Scratch ](#)(Code)
4. [Transformers from Scratch ](#)(Again, this time long detailed deep dive)
5. [An Intuitive Introduction to Transformers ](#)(Pictures)
6. [The Transformer – Attention is All You Need ](#)()
7. [Primers – Transformer ](#)(Long, detailed Deep Dive)
8. [Some Intuition on Attention and the Transformer ](#)(Short insights)
9. [Transformer Math ](#)(If you want to implement a big one in practice)

🥡 **Takeaway** 🥡

**Model architectures** are influenced by the **inductive bias of the data.** is present, while **Respecting symmetry** and **making models scale well** are two popular approaches these days.