

UNIVERSITÄT
HEIDELBERG



A Zoo of Models

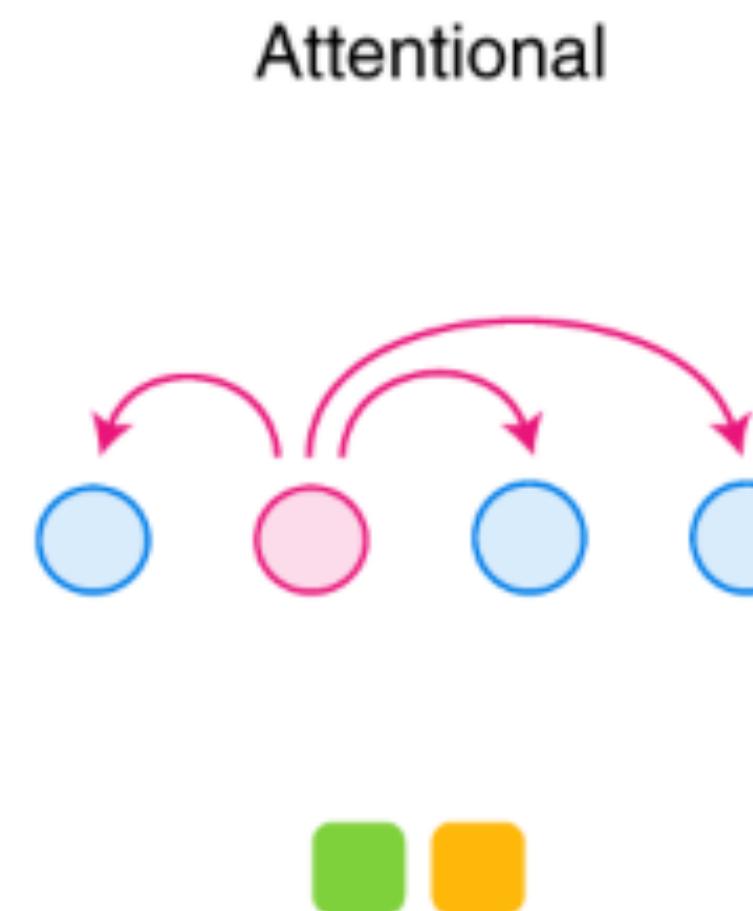
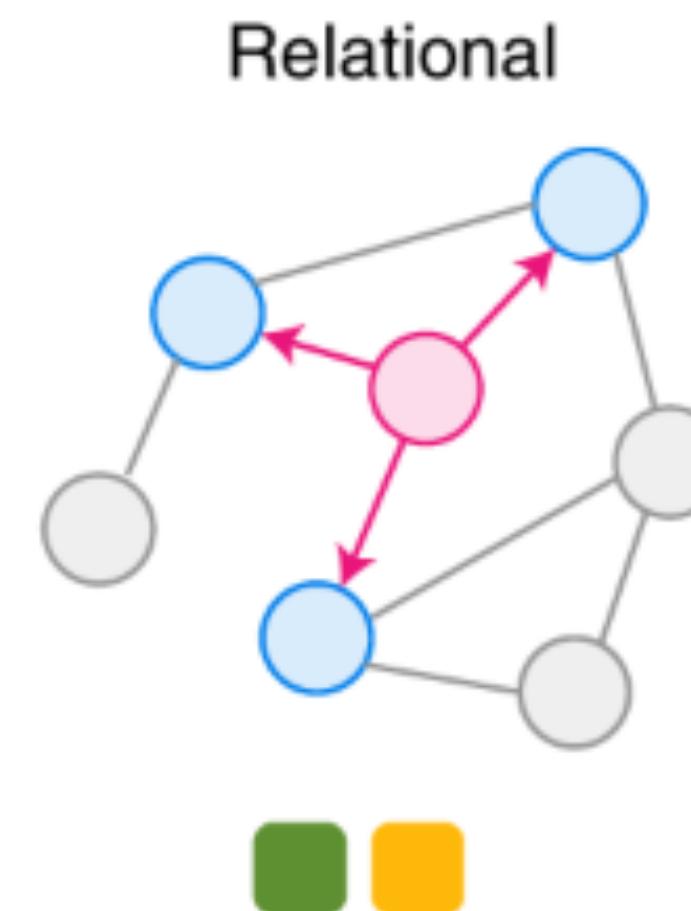
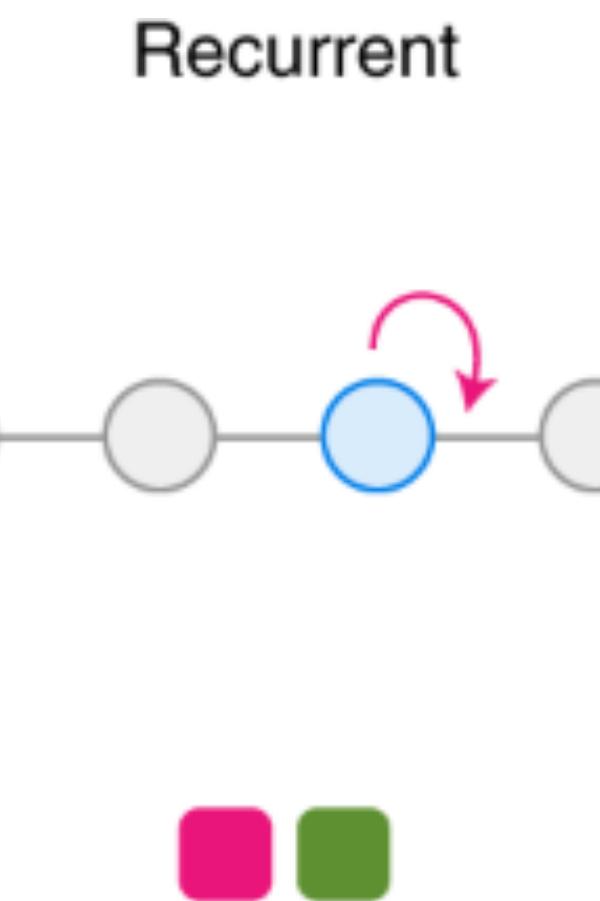
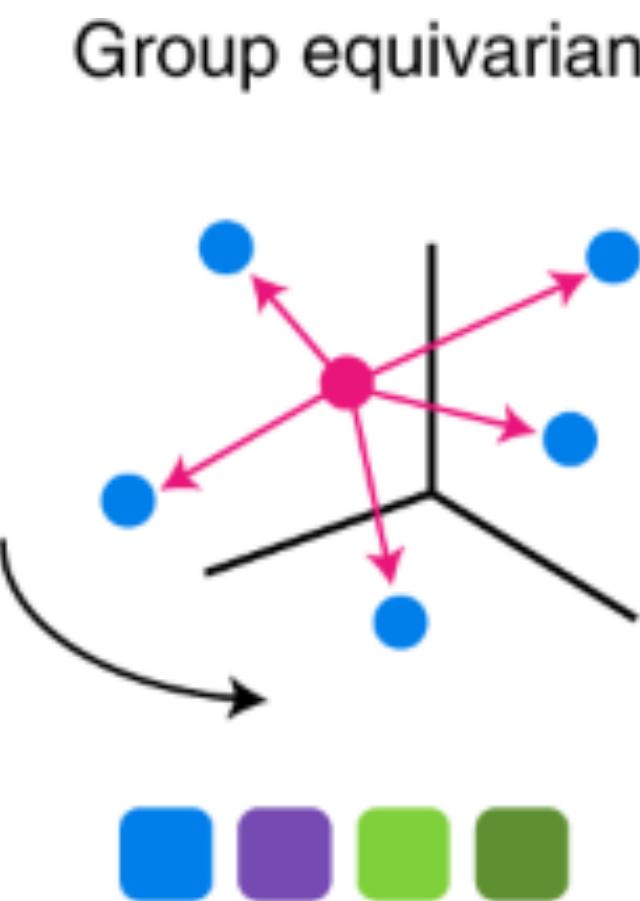
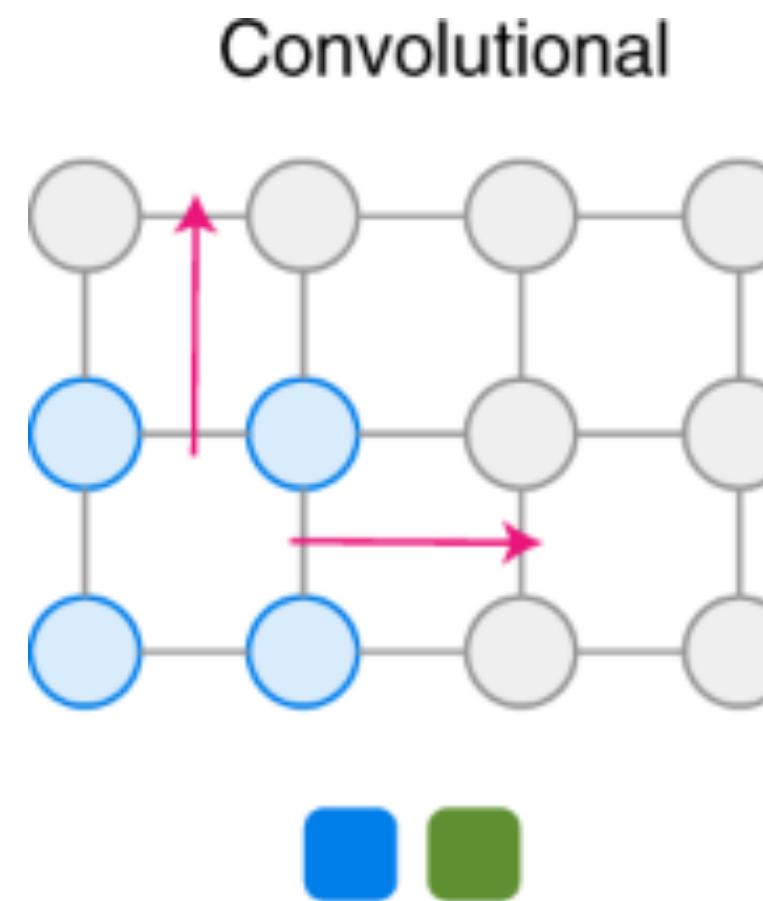
L3, Structural Bioinformatics

WiSe 2023/24, Heidelberg University

Kieran Didi

How to make sense of all these models?

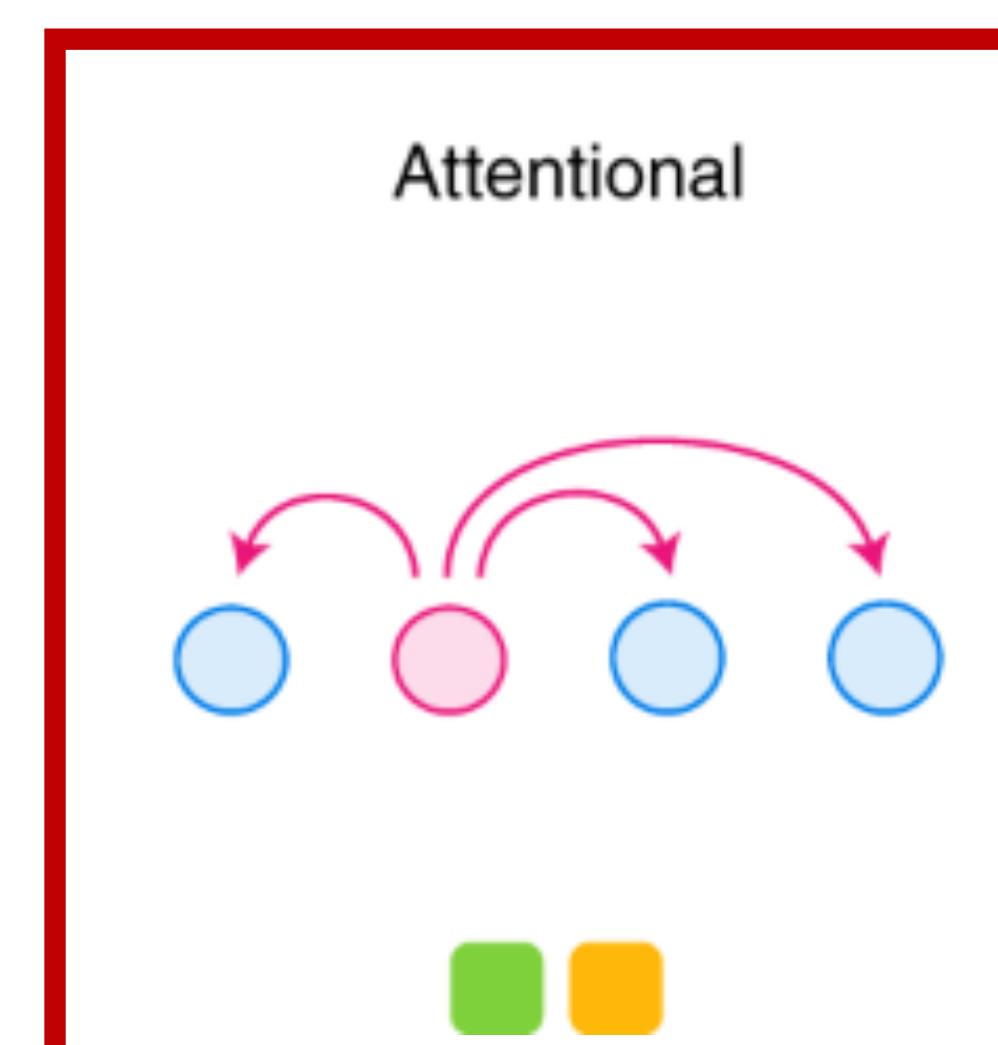
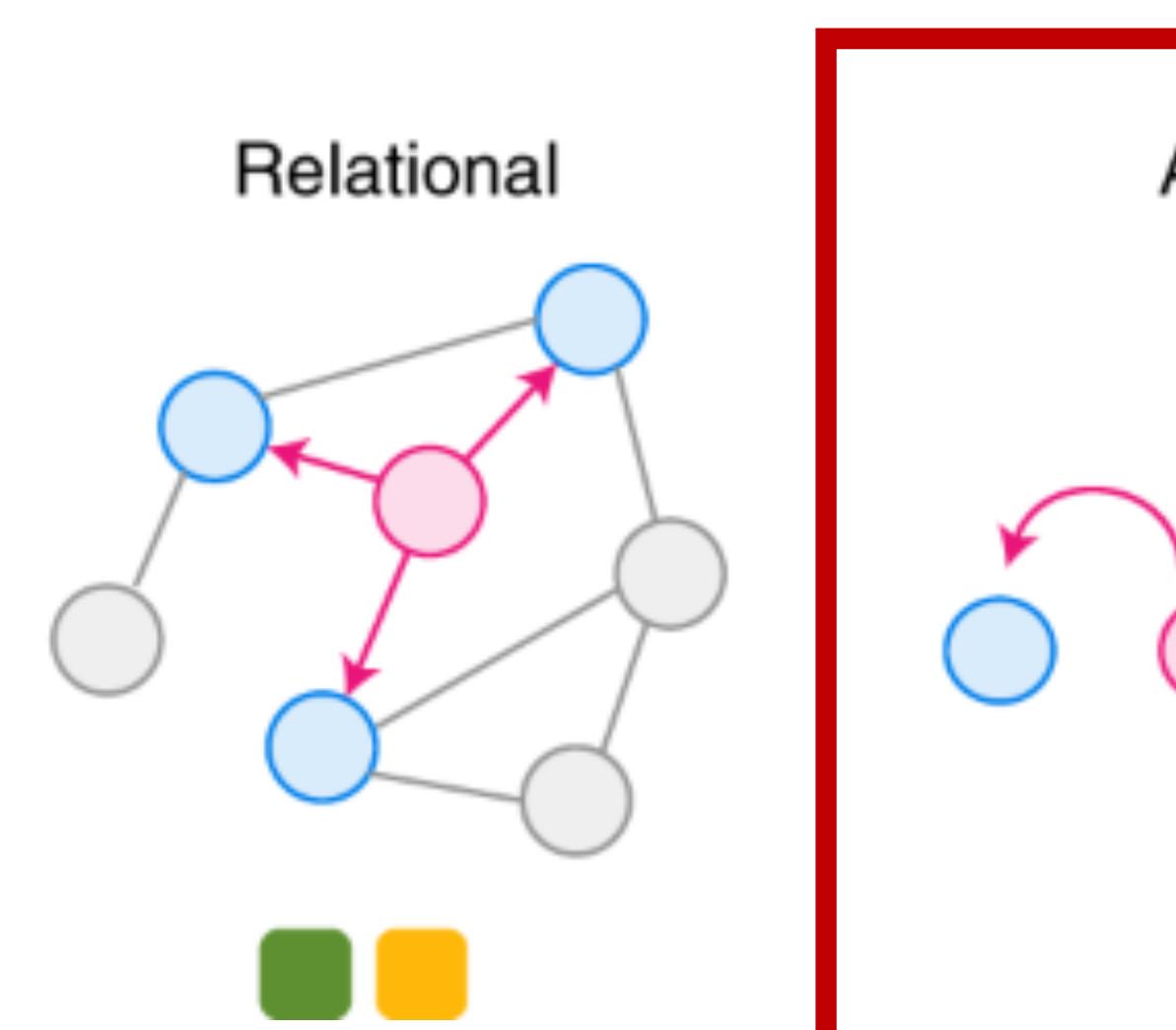
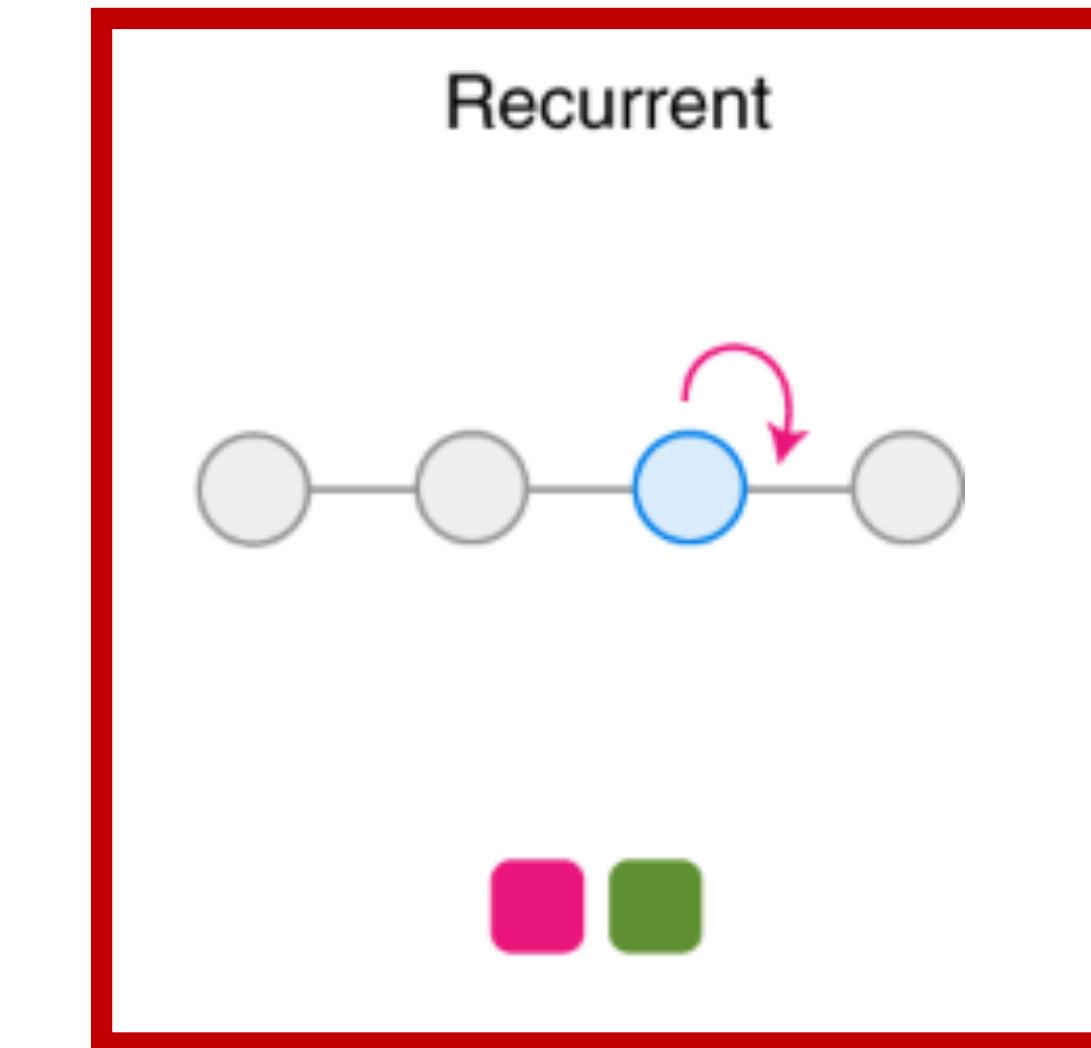
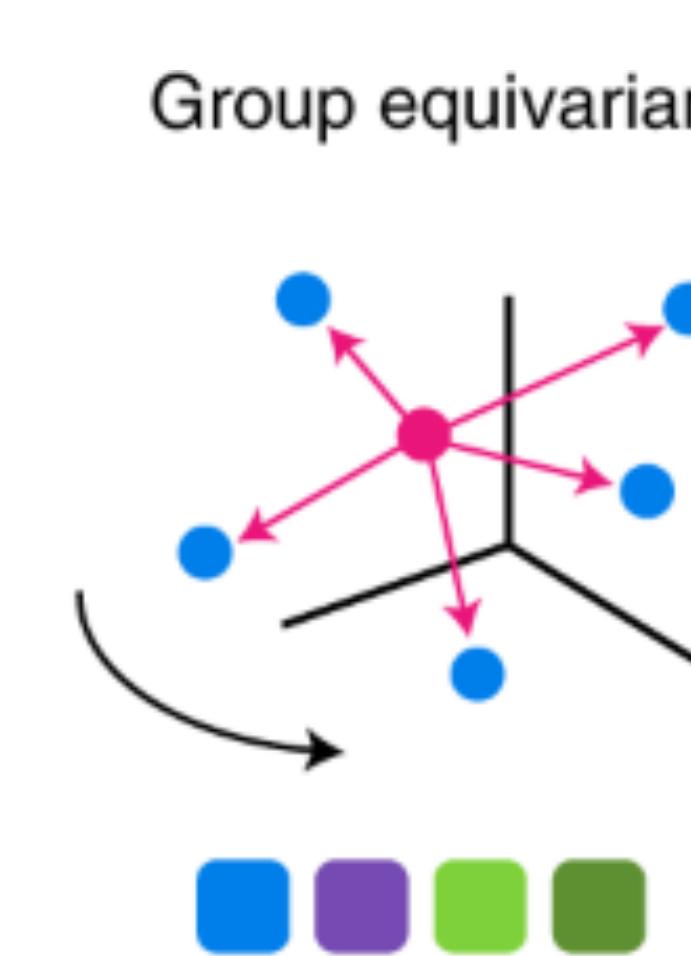
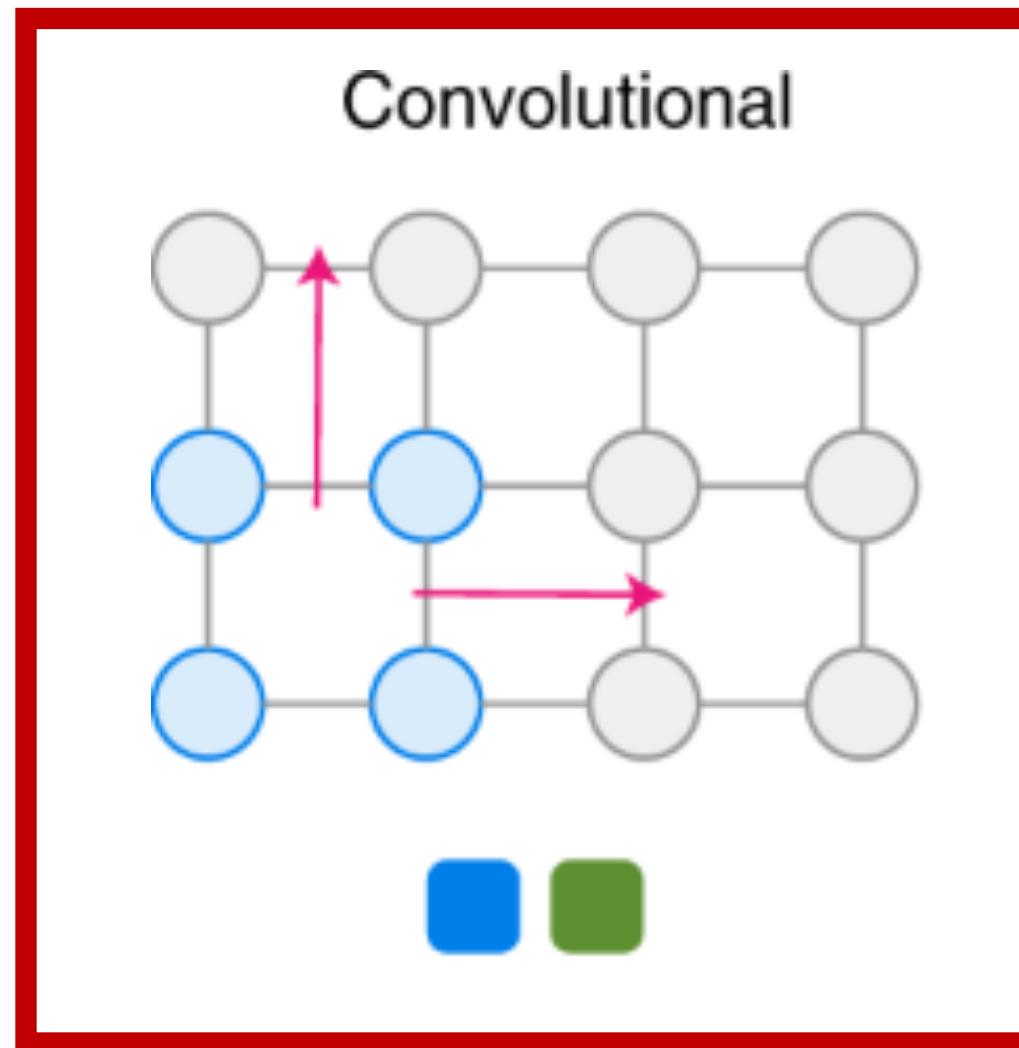
Find the inductive biases they instill in the network



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

How to make sense of all these models?

Find the inductive biases they instill in the network

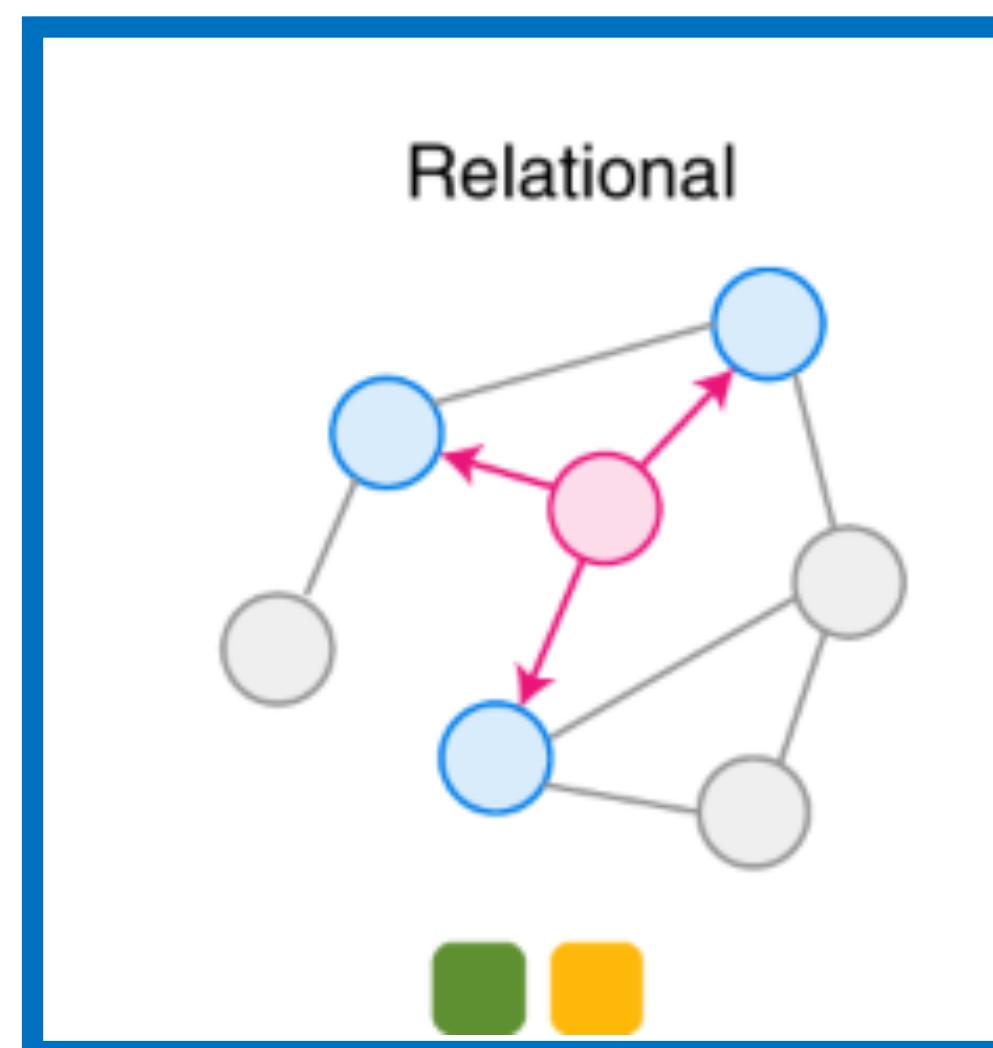
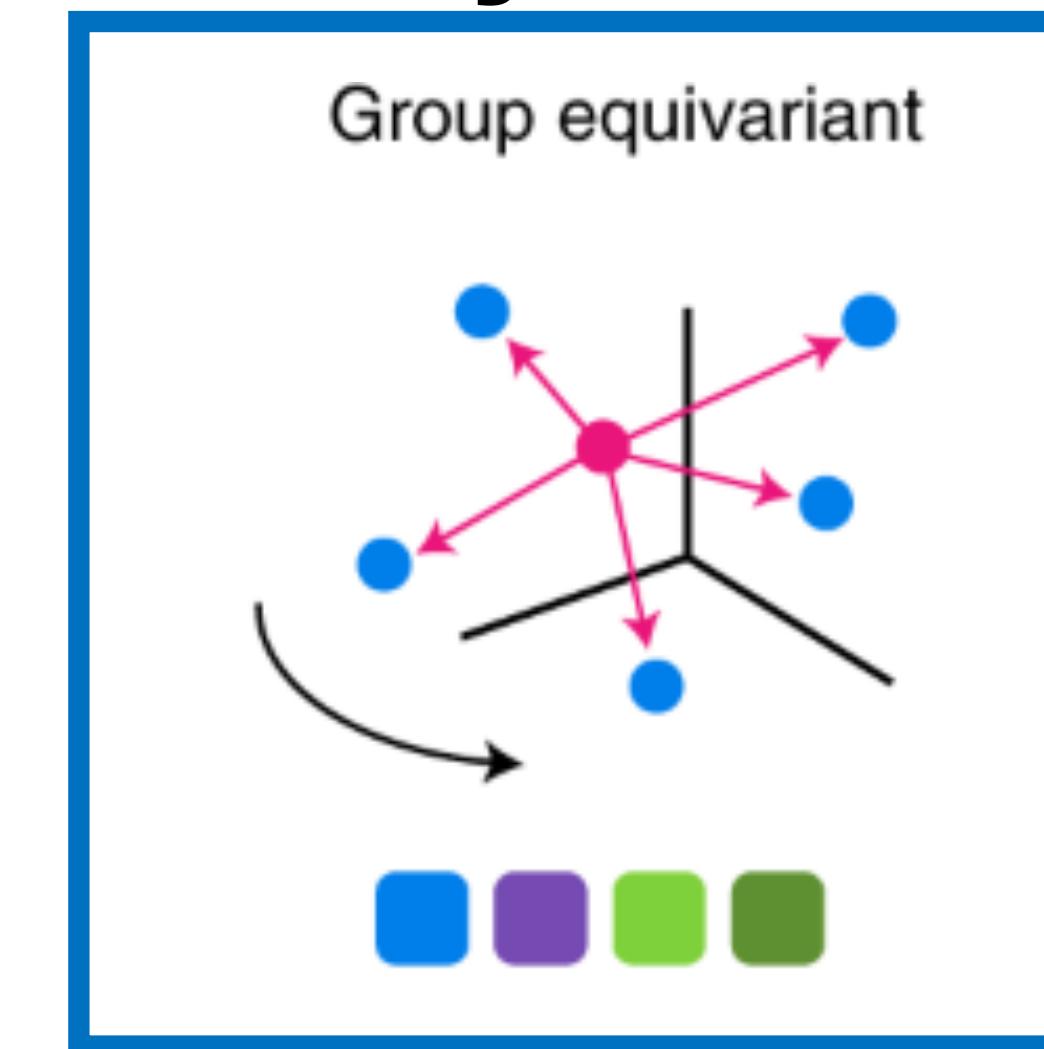
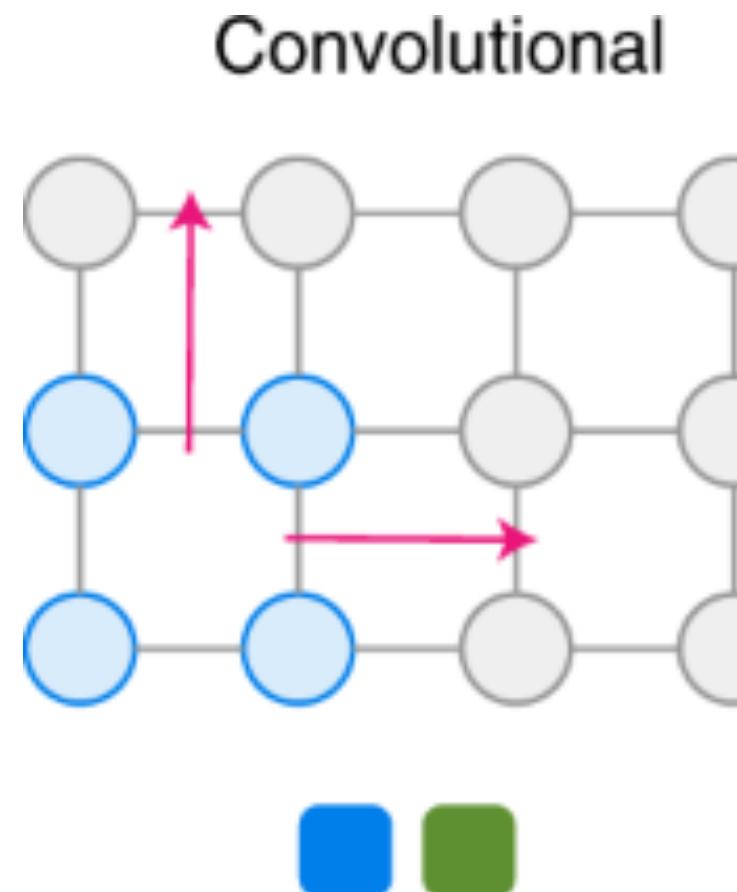


- █ Translational invariance
- █ Rotational invariance
- █ Repeating dynamics
- █ Non-locality
- █ Locality
- █ Unordered

This week

How to make sense of all these models?

Find the inductive biases they instill in the network



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

Next week

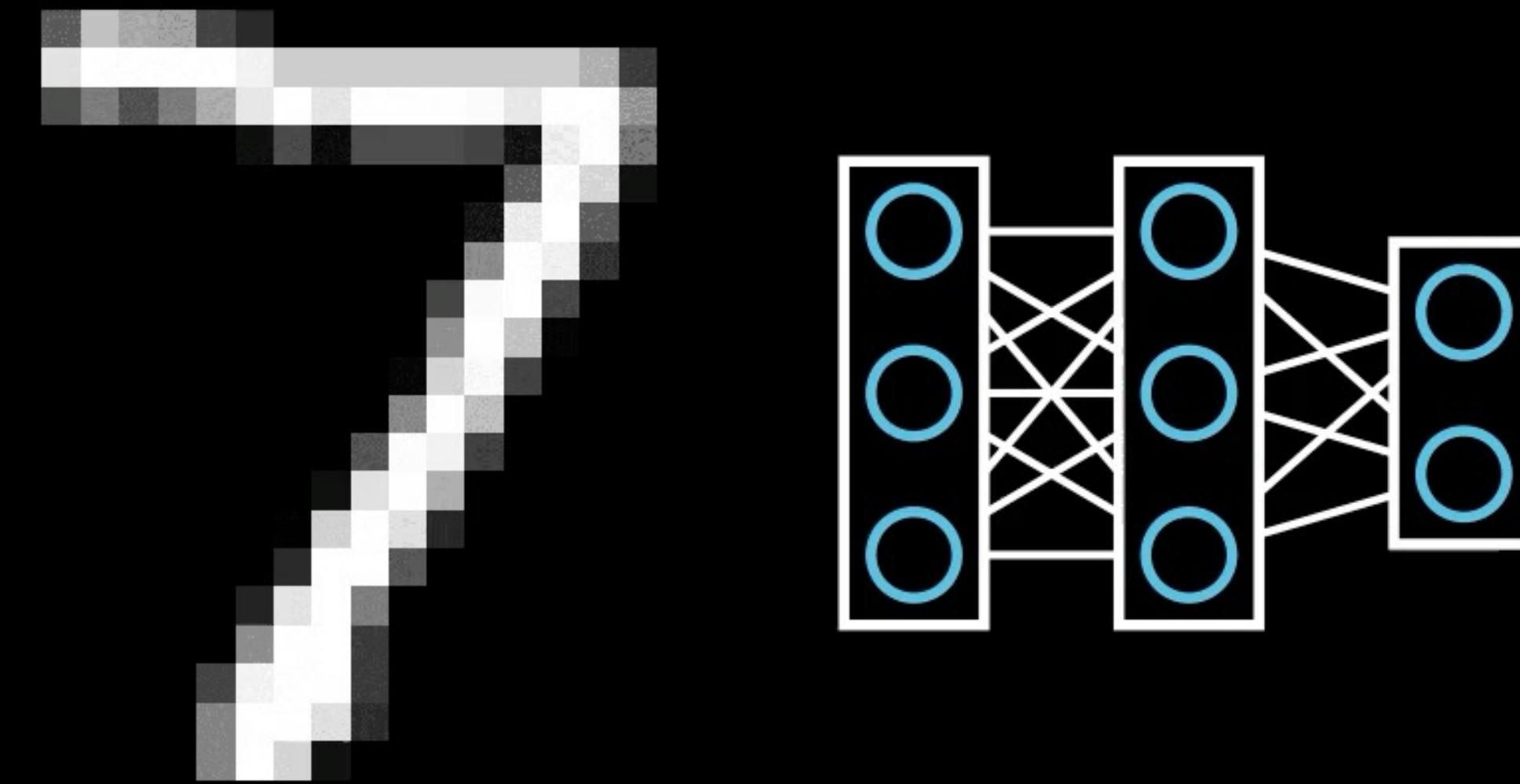
Overview

- 1. Images: Convolutional Neural Networks**
- 2. Sequences: RNNs**
- 3. Transformers**
- 4. Current developments**

1. Convolutional Neural Networks

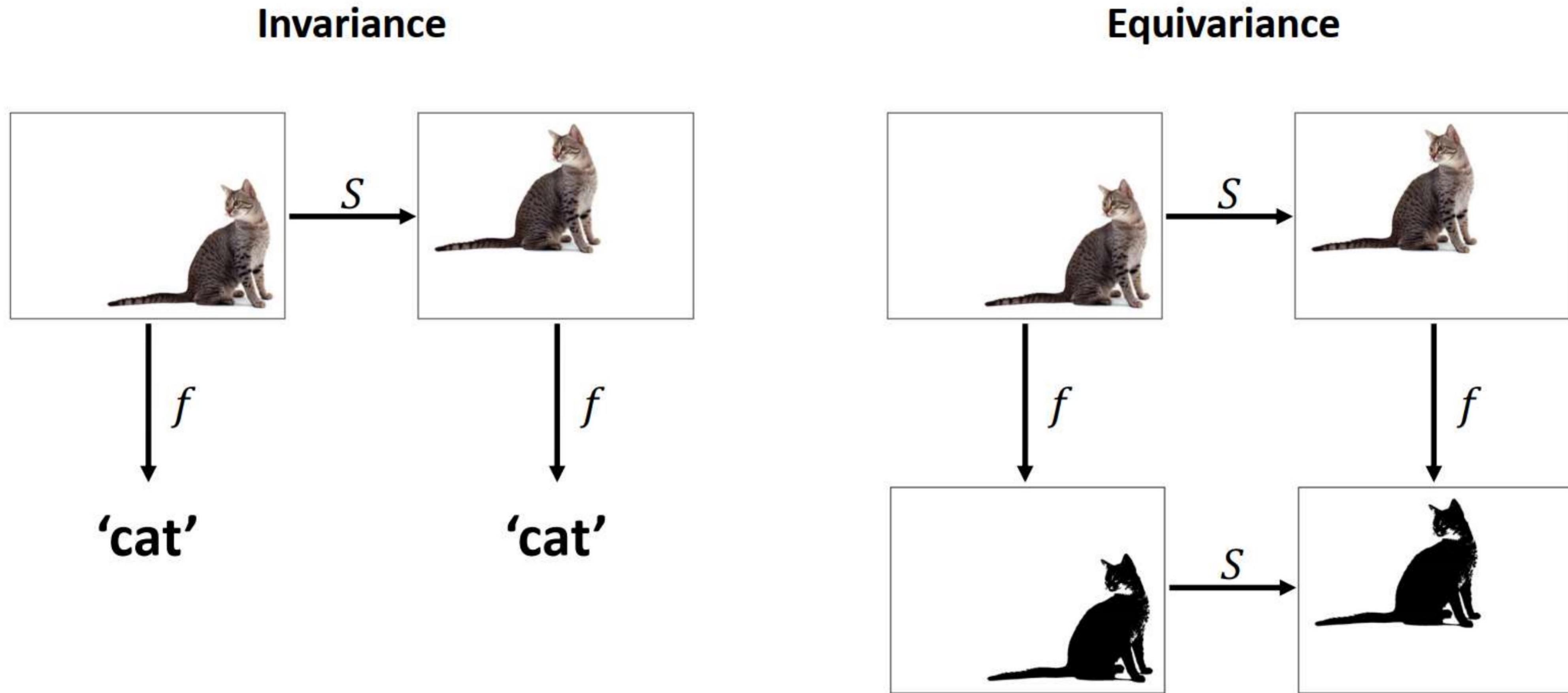
How to deal with images

Naive approach: unroll them and pass them into an MLP



Inductive Bias: Translational In-/Equivariance

Leverage the symmetry of your data



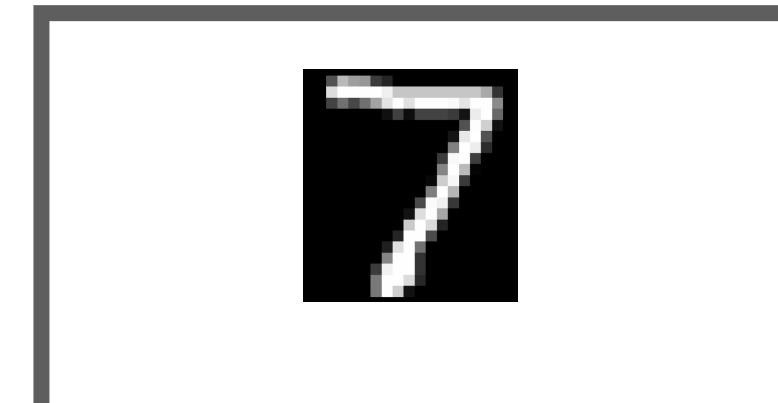
Why leverage symmetries?

We need more data = our network is more efficient!

Training without translational symmetry

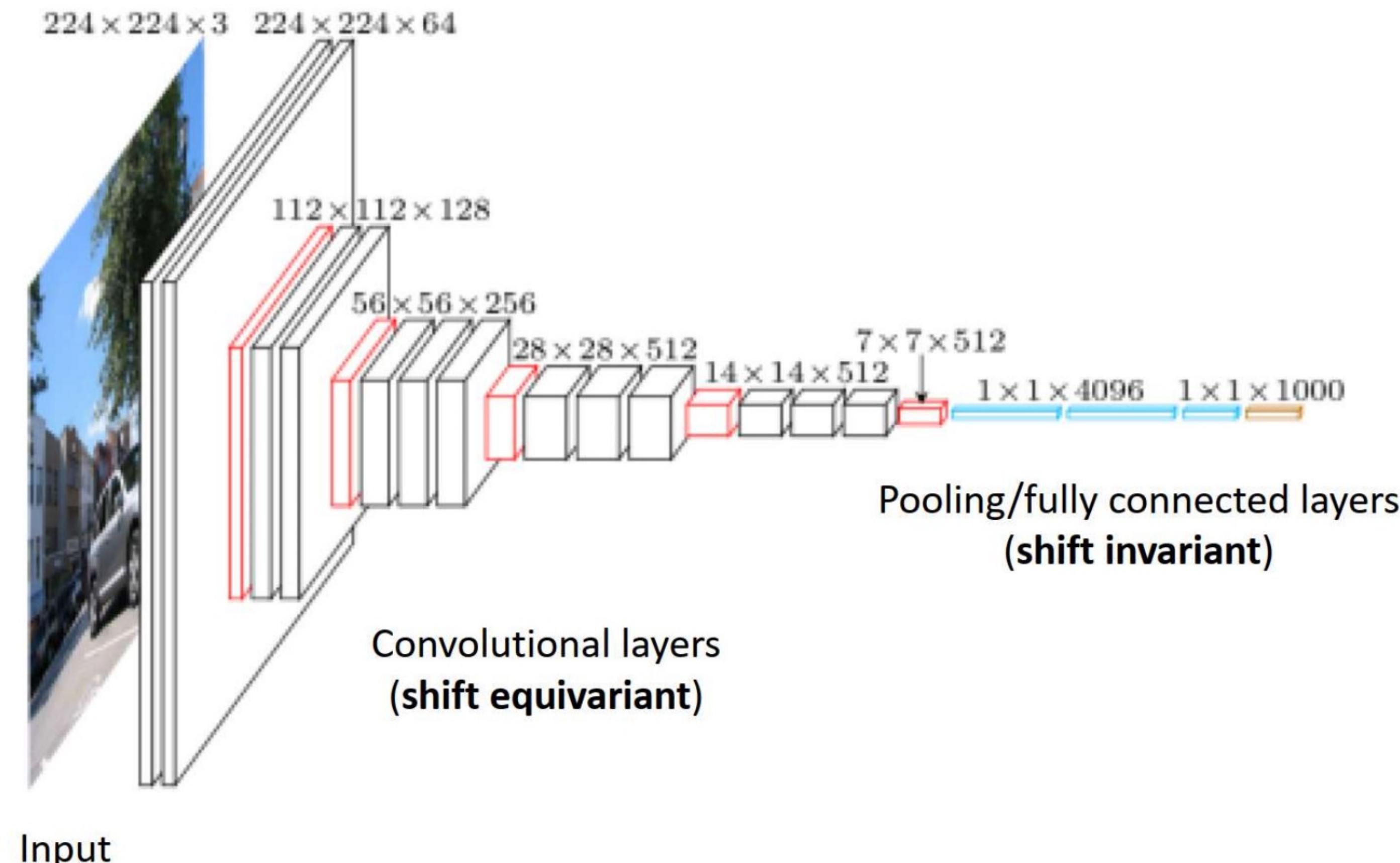


Training with translational symmetry



How do we do this in practice?

Implement neural network layers that respect these symmetries

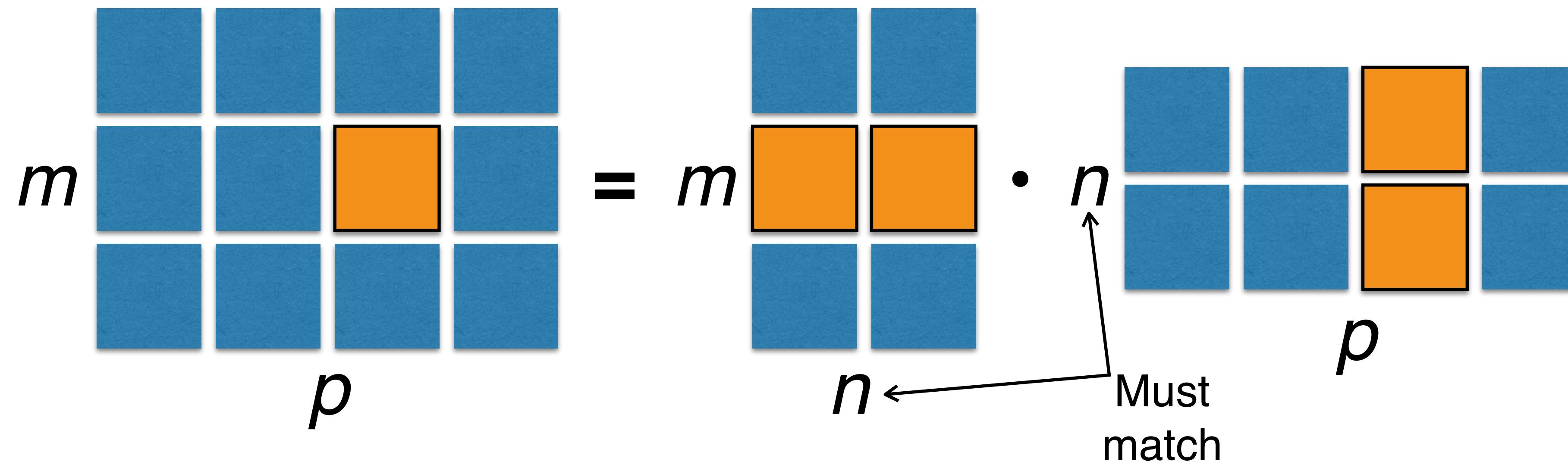


Convolutional Layers

Reminder: Matrix multiplication

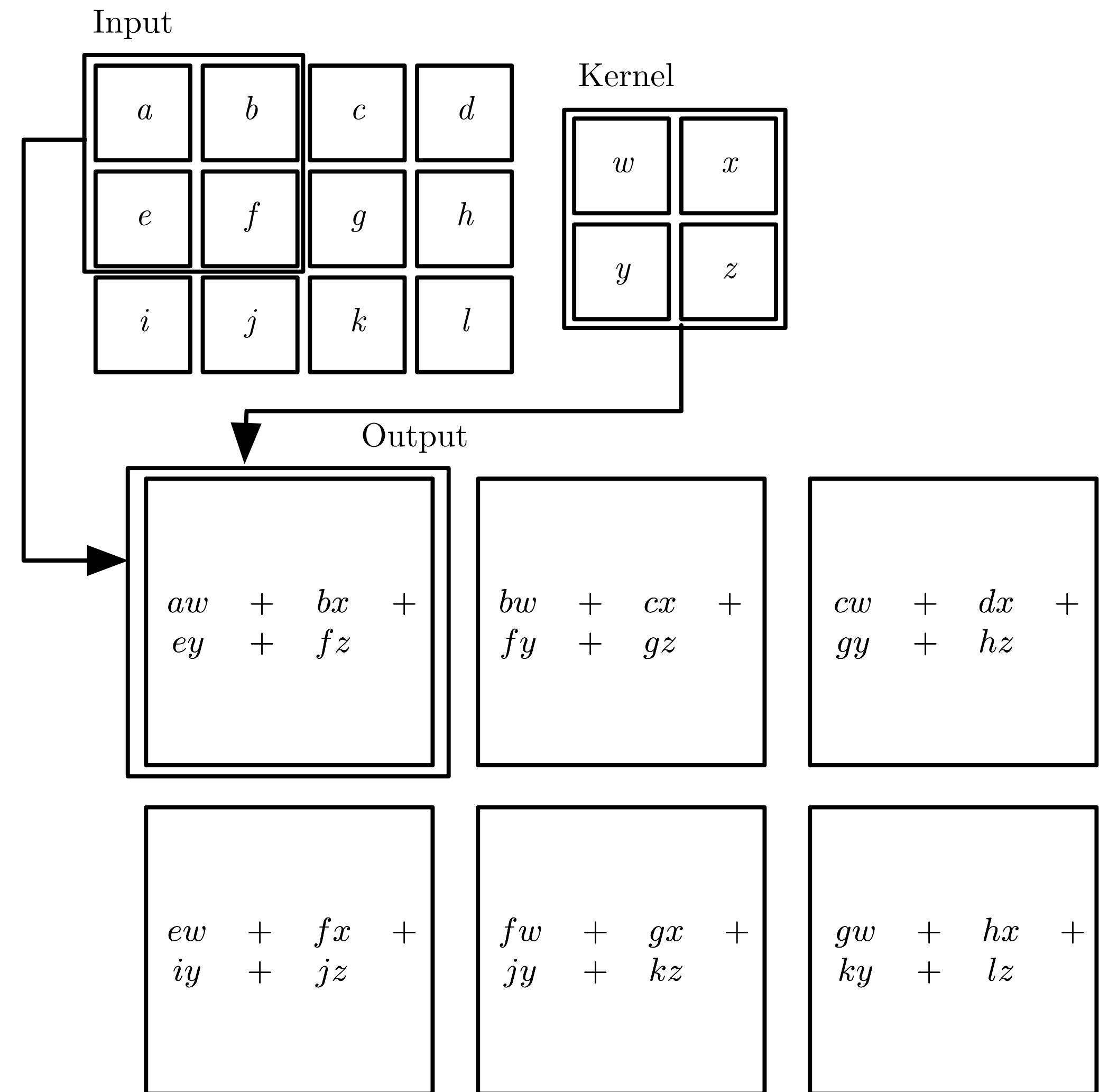
$$C = AB. \quad (2.4)$$

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



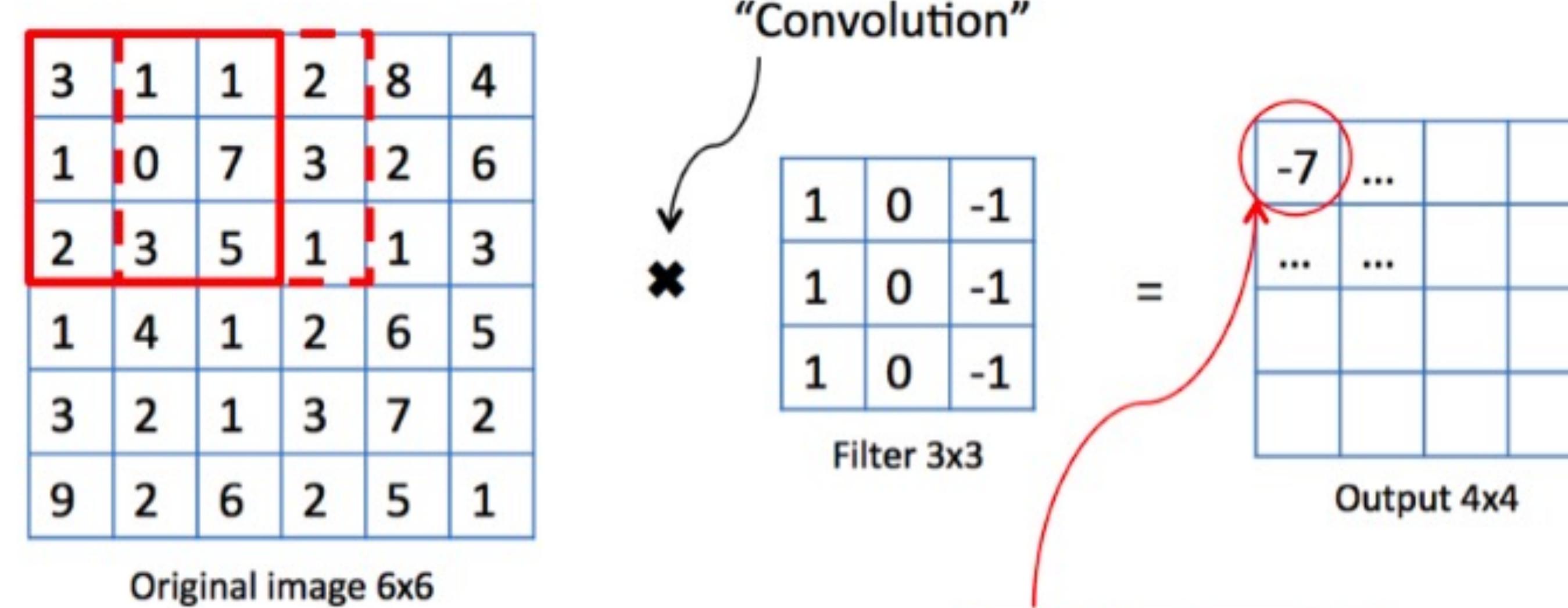
Convolutional Layers

The weights are in the kernel



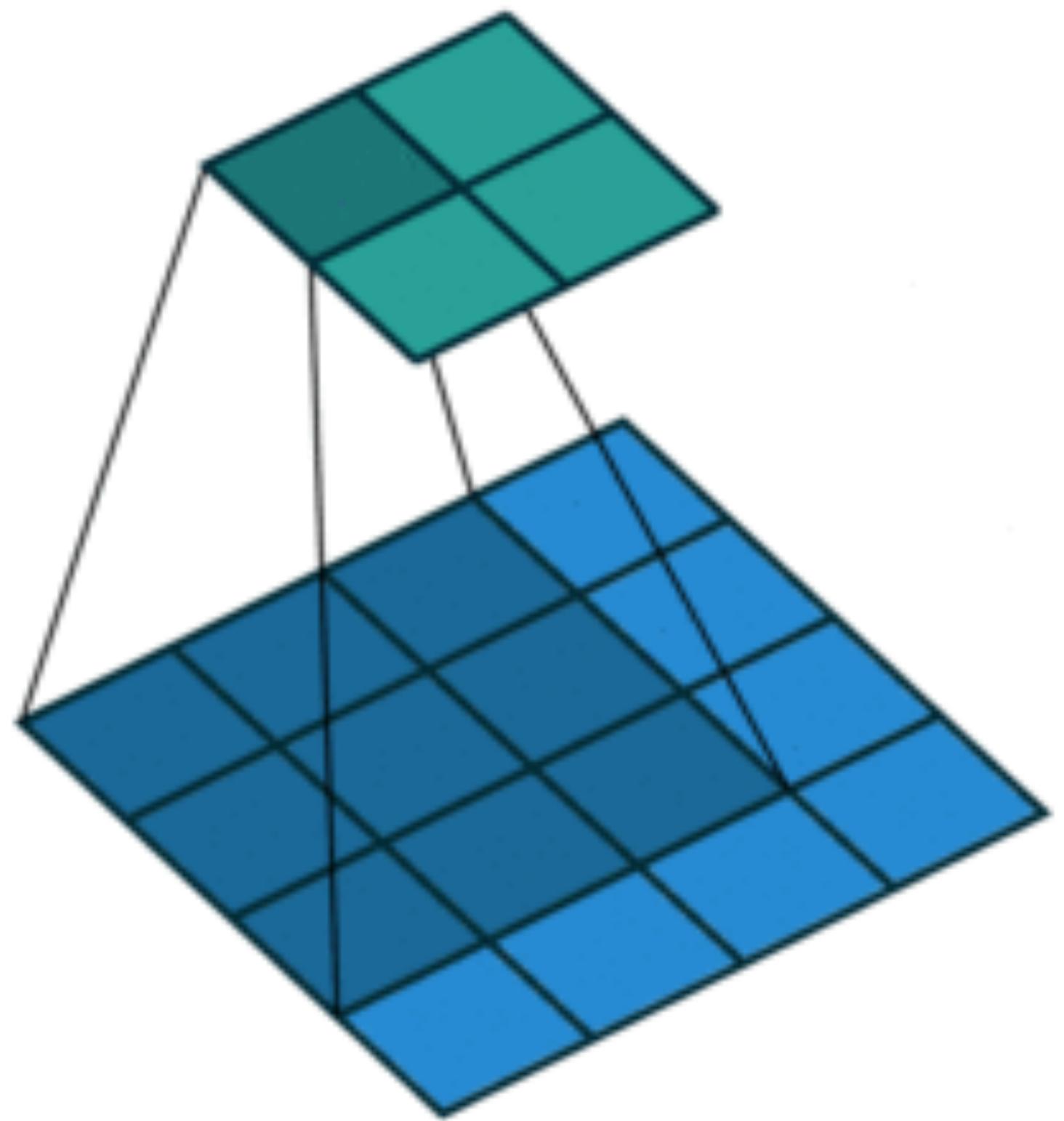
Convolutional Layers

Convolution = Repeated Matrix Multiplication



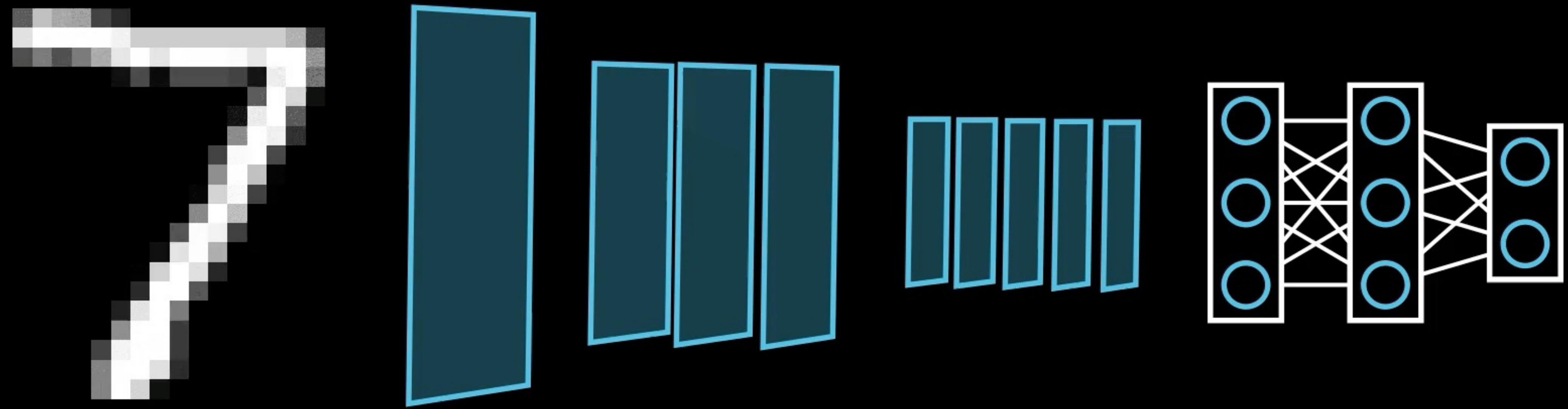
How can I imagine that?

Sliding the kernel over the image



How can I imagine that?

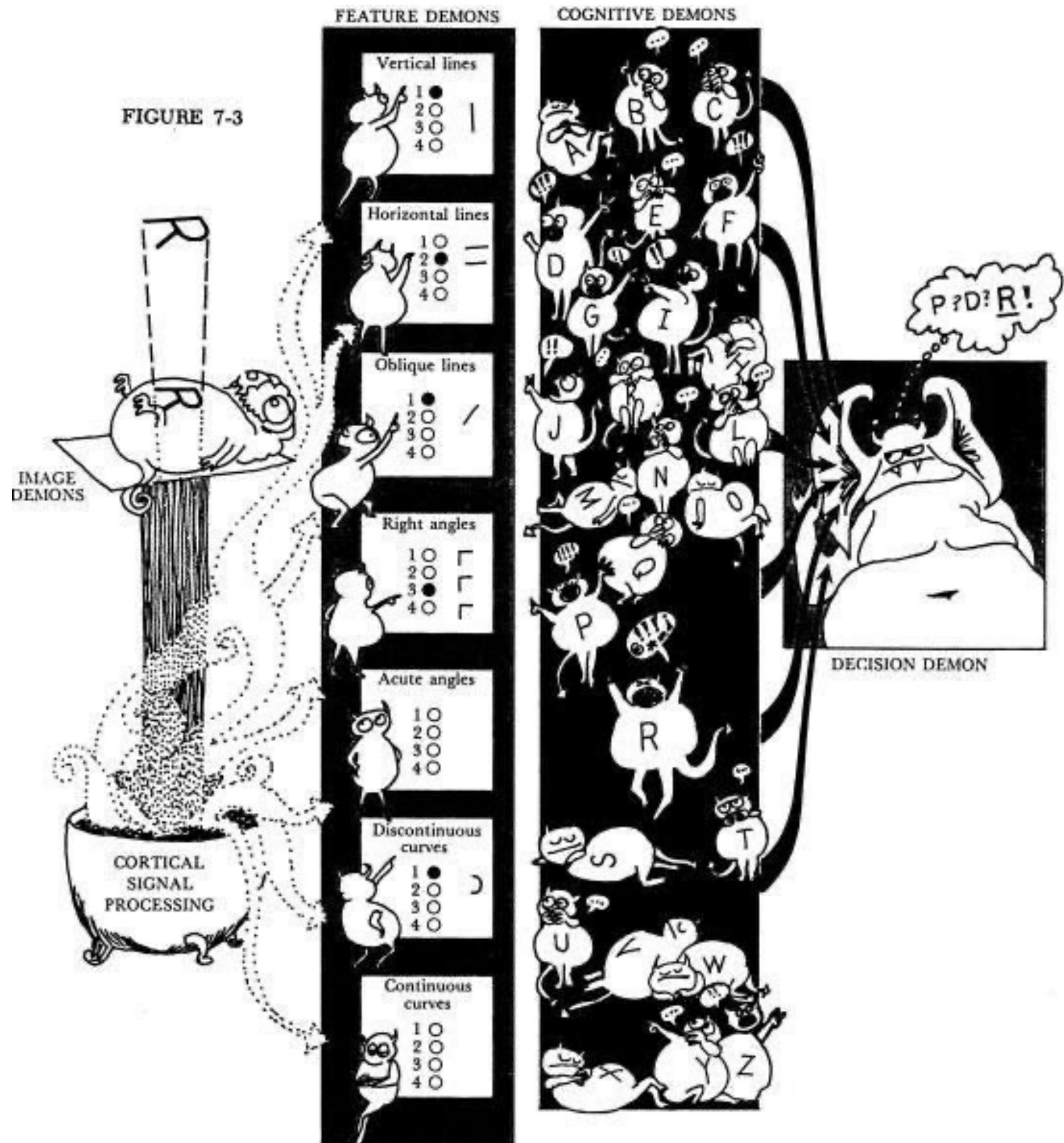
Multiple Kernels allow detecting multiple features



Pattern Recognition all over again

This time adjusted to the image case

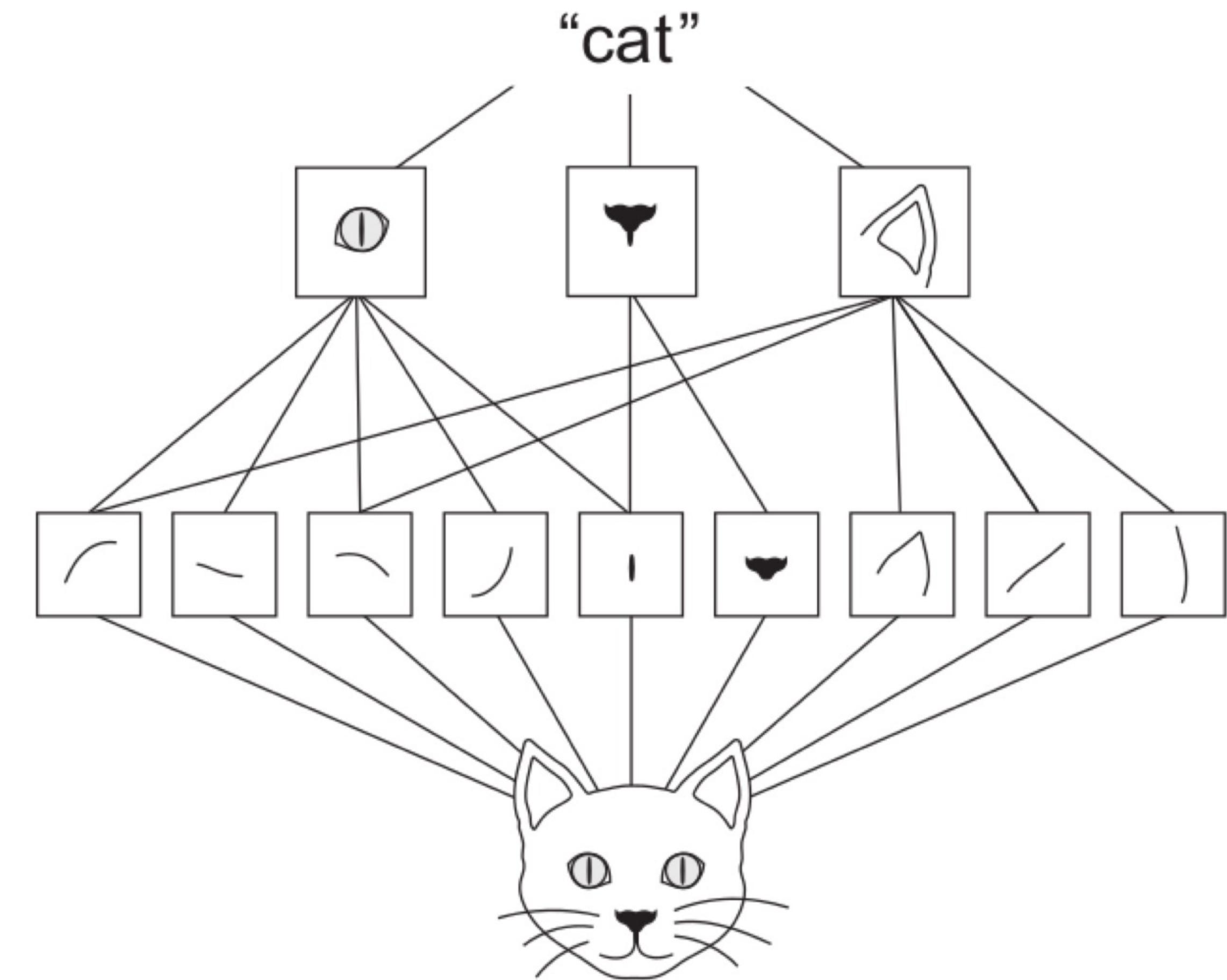
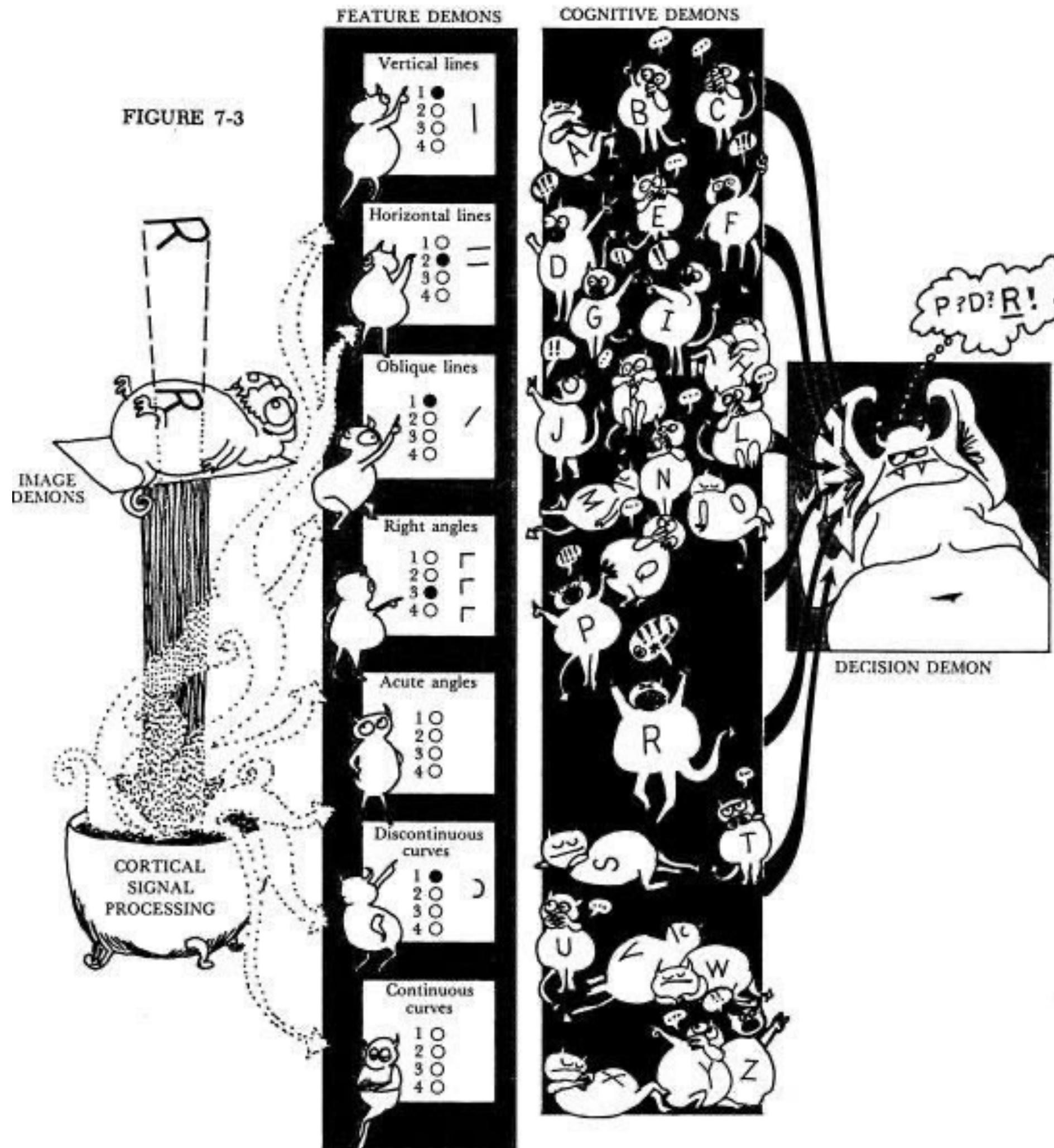
266 7. Pattern recognition and attention



Pattern Recognition all over again

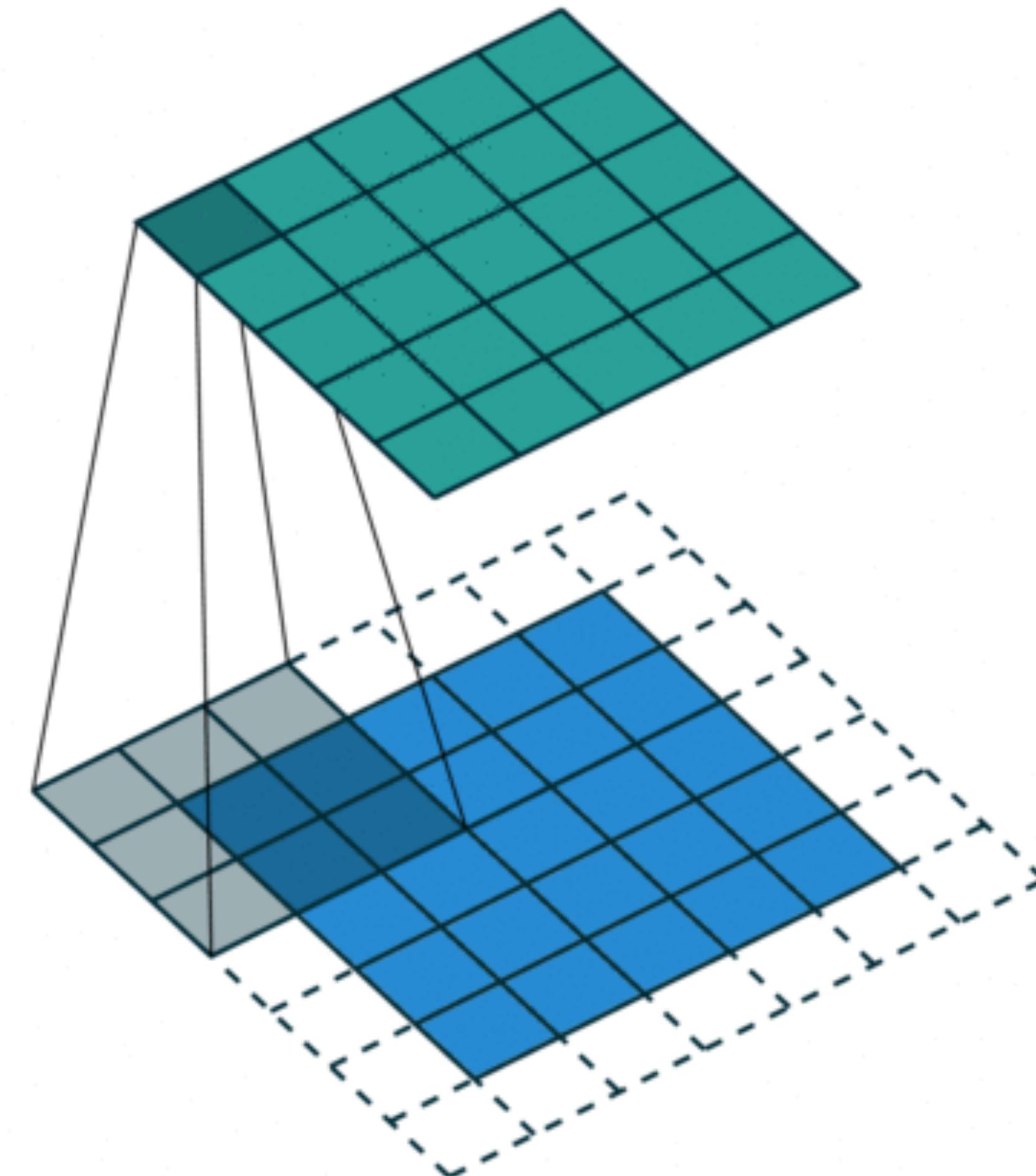
This time adjusted to the image case

266 7. Pattern recognition and attention

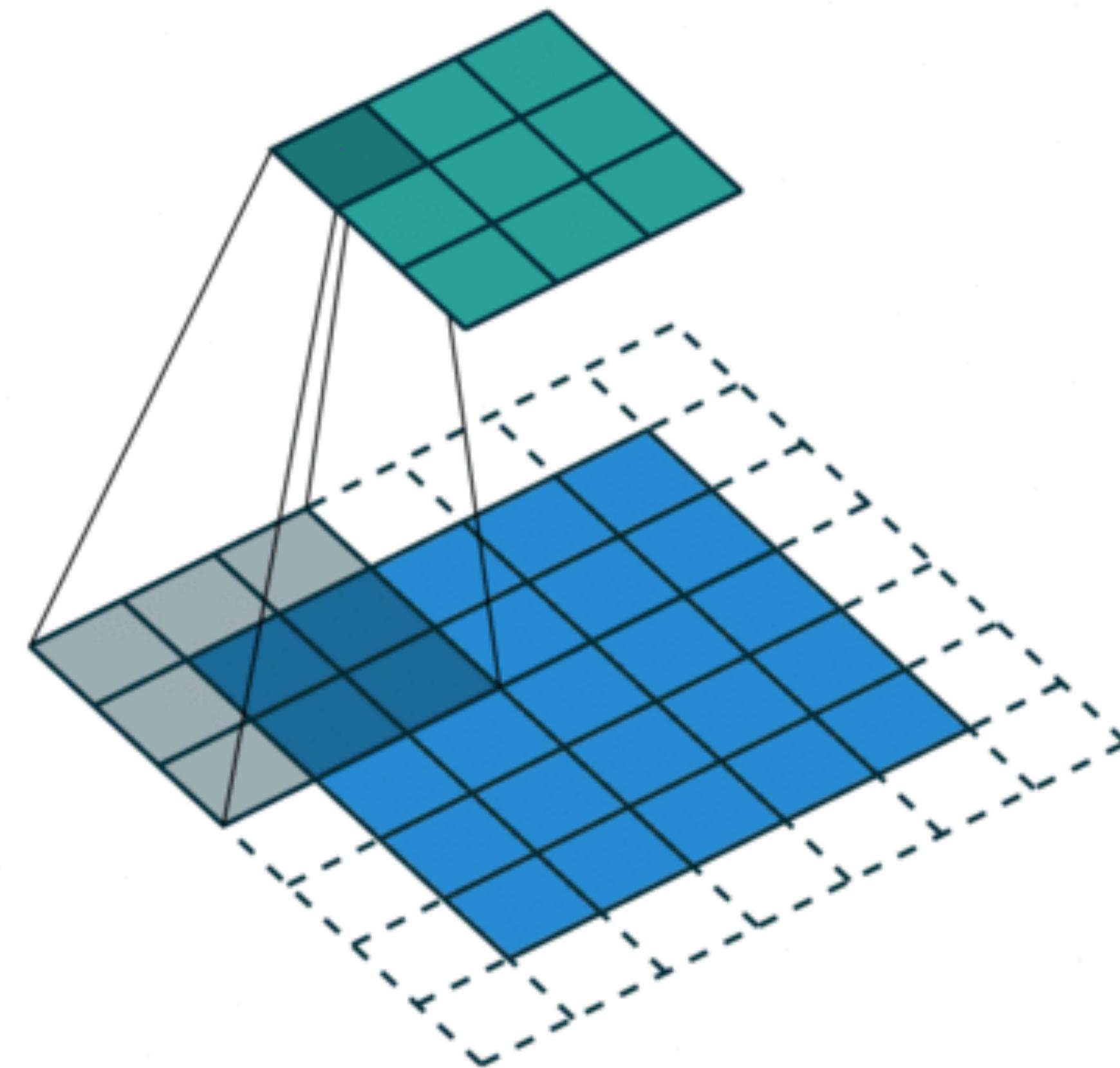


Avoid reducing size with padding

Different ways to pad (zero-pad, mean-pad, ...)

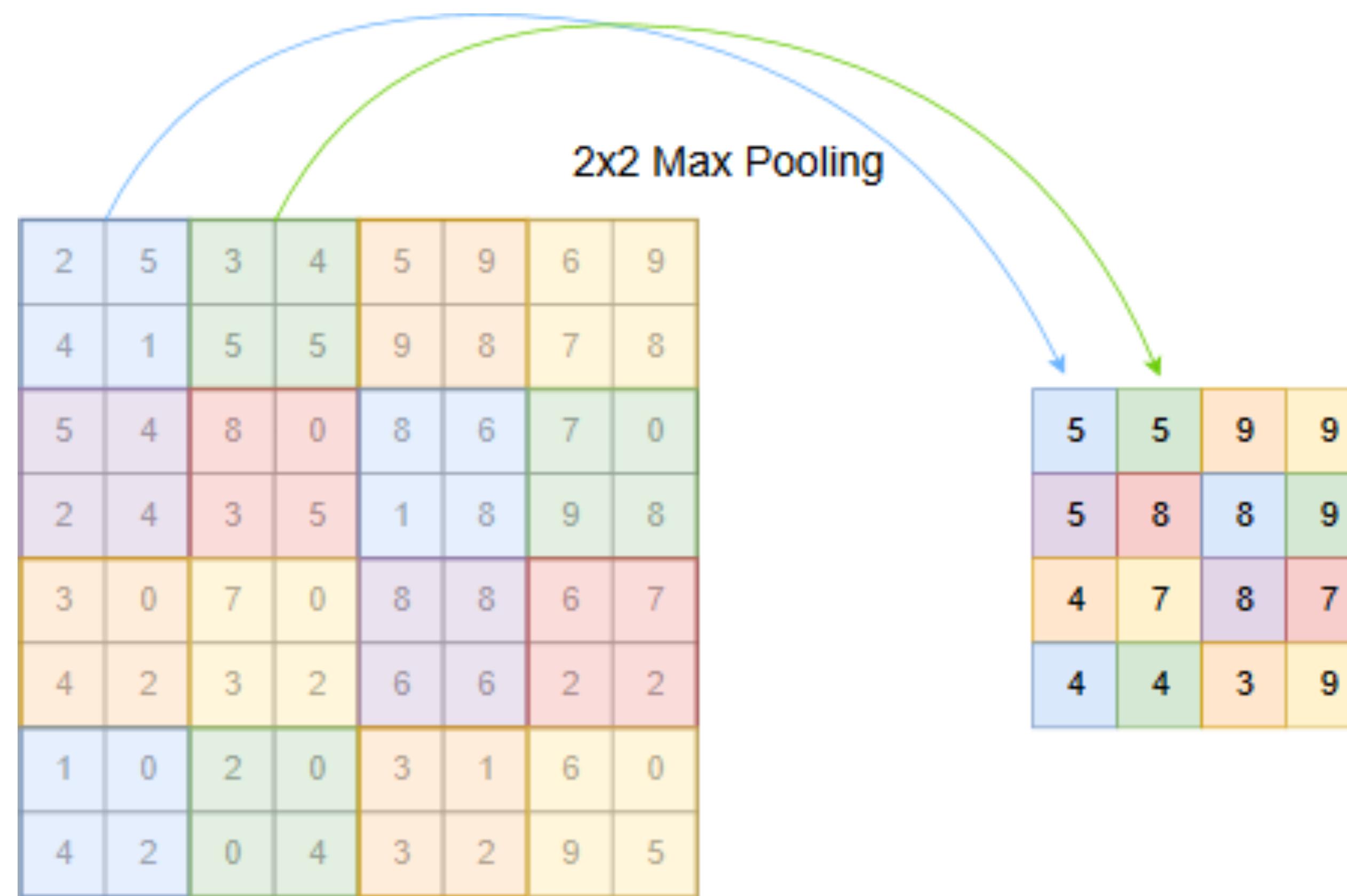


Make bigger jumps with strides



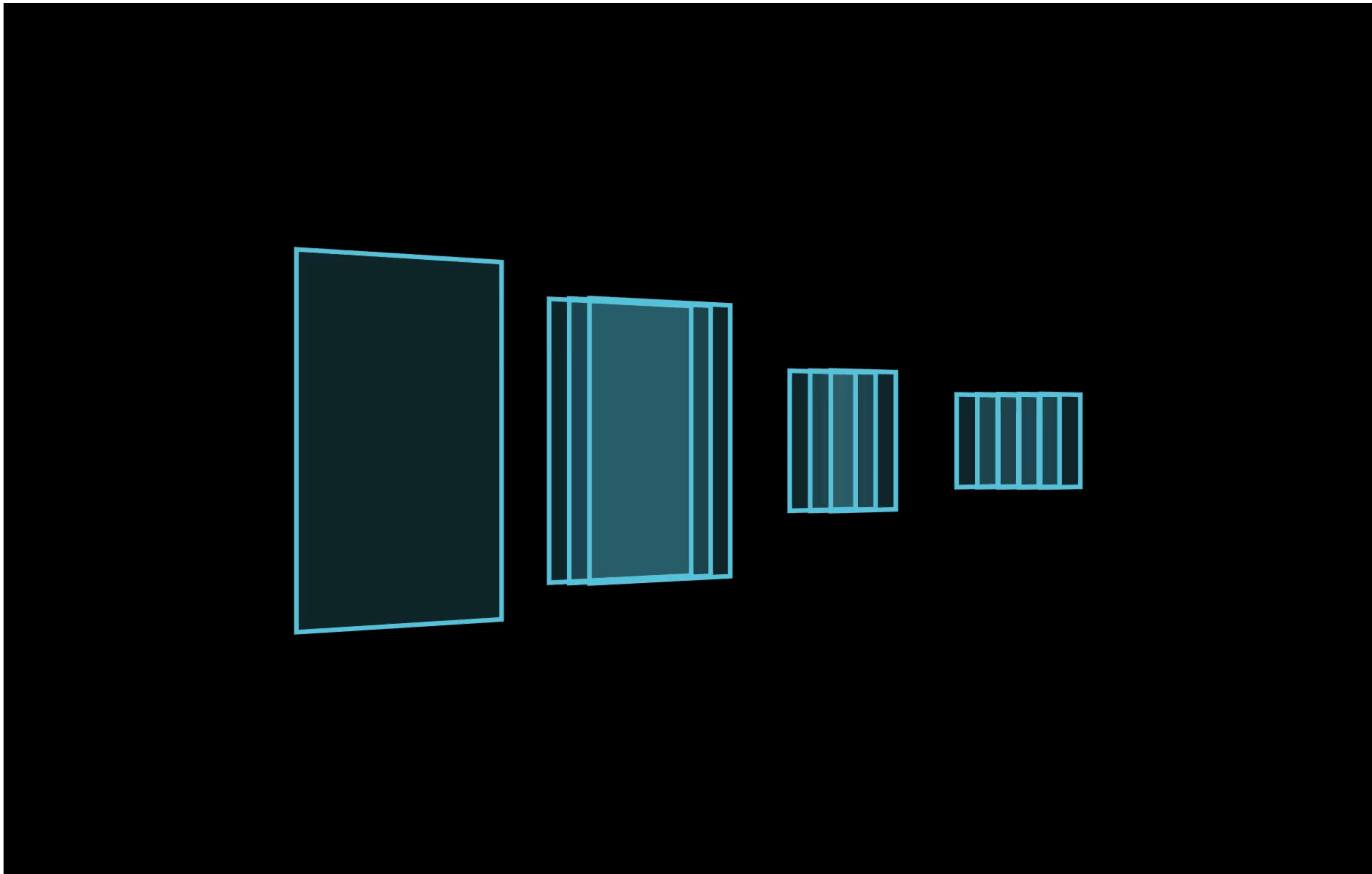
Pooling: Shift-invariant operation

Reduce size, but no learning involved



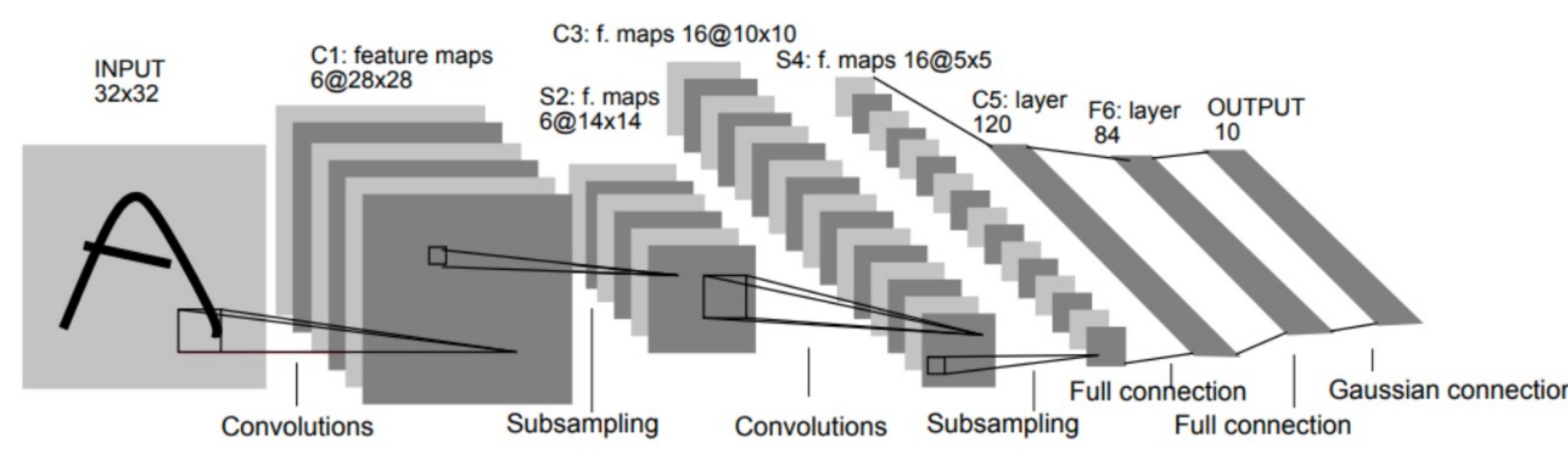
Putting things together: A full CNN

Conv. Layers -> Pooling -> FC Layers



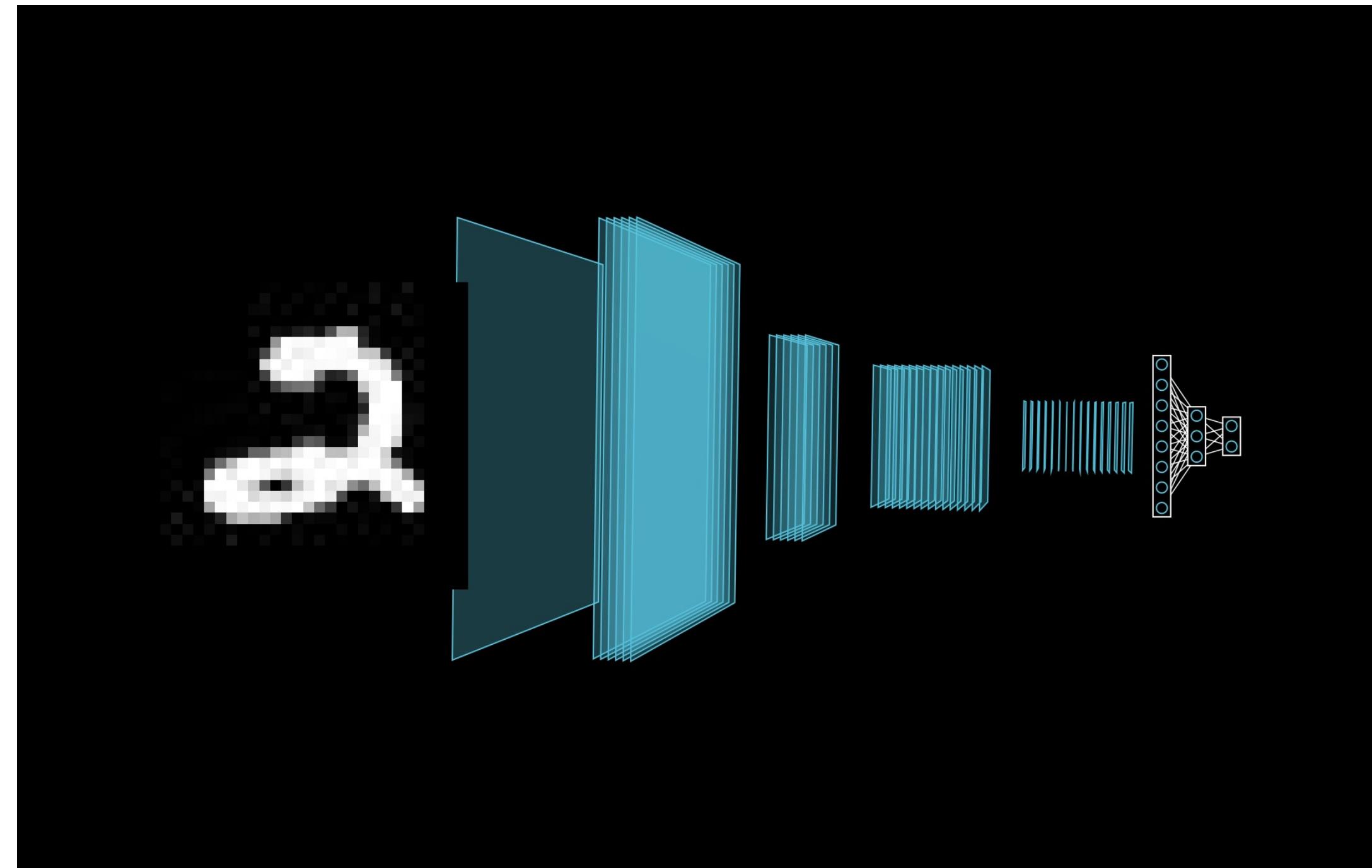
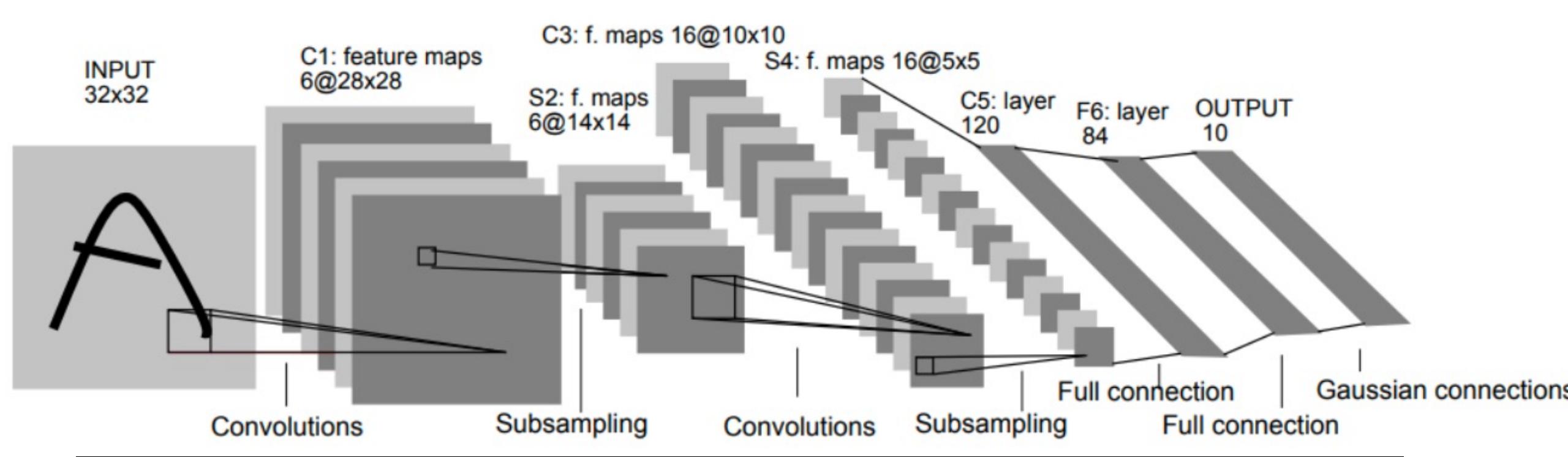
LeNet (1998): CNNs become a thing

Exactly what we discussed, just bigger



LeNet (1998): CNNs become a thing

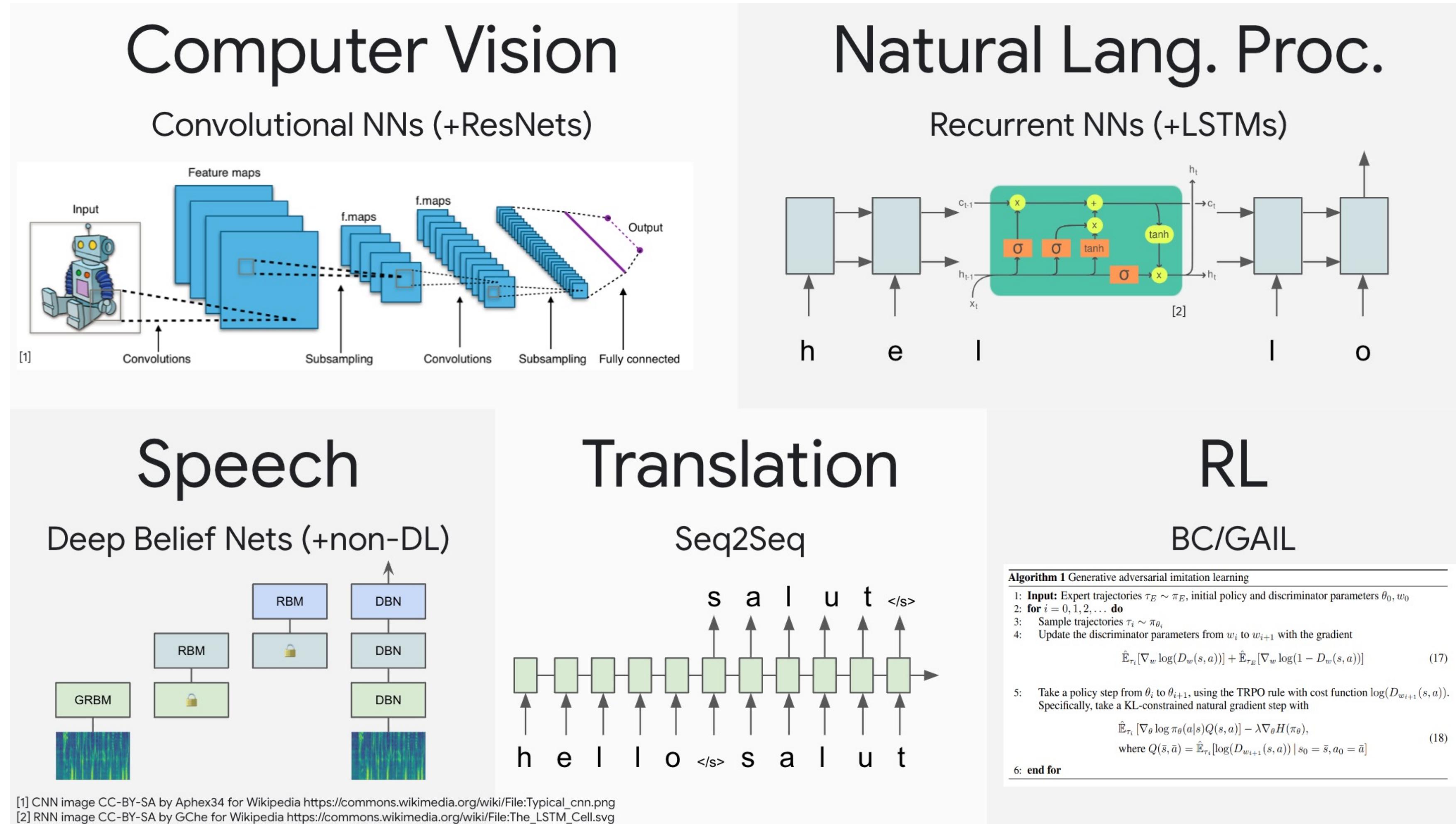
Exactly what we discussed, just bigger



2. RNNs

The classic landscape

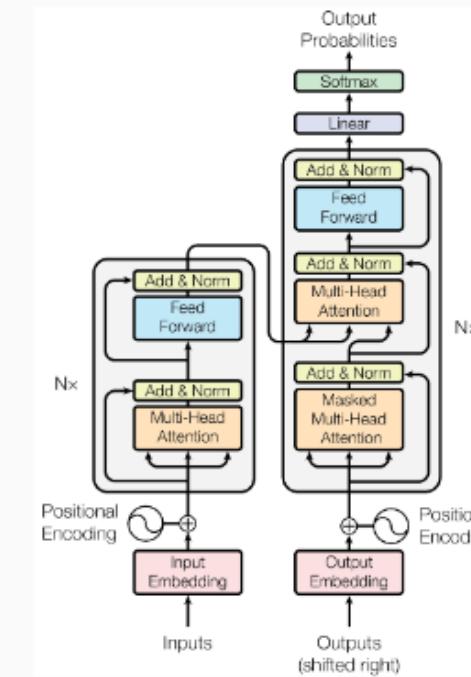
One architecture per community



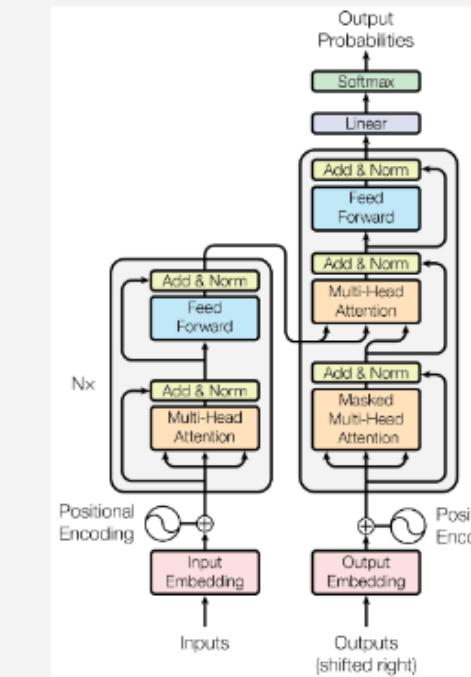
The transformer's takeover

One community at a time

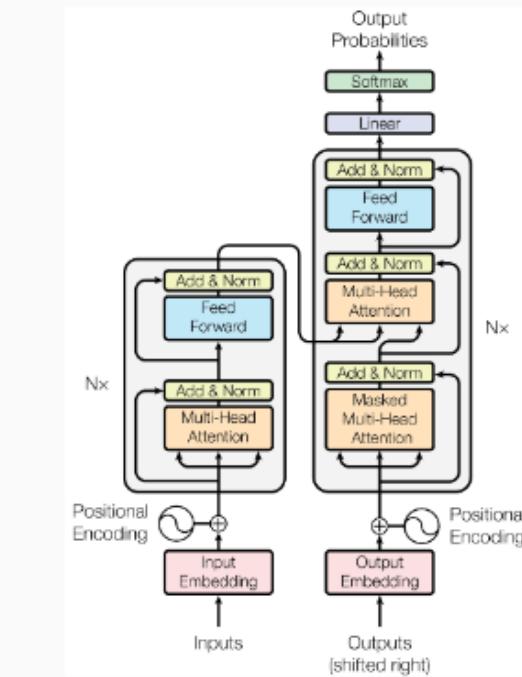
Computer Vision



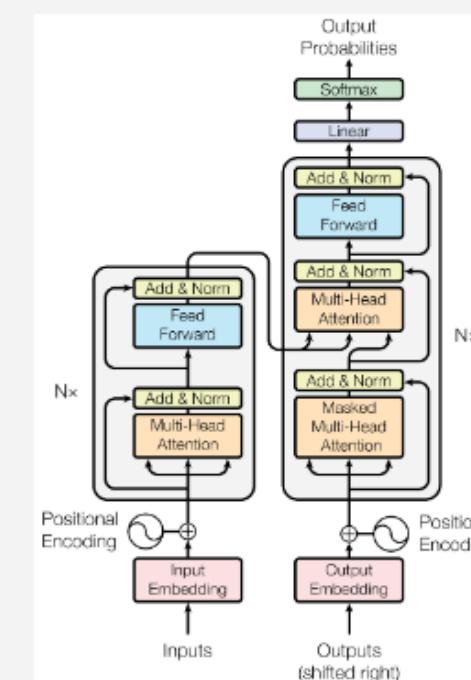
Natural Lang. Proc.



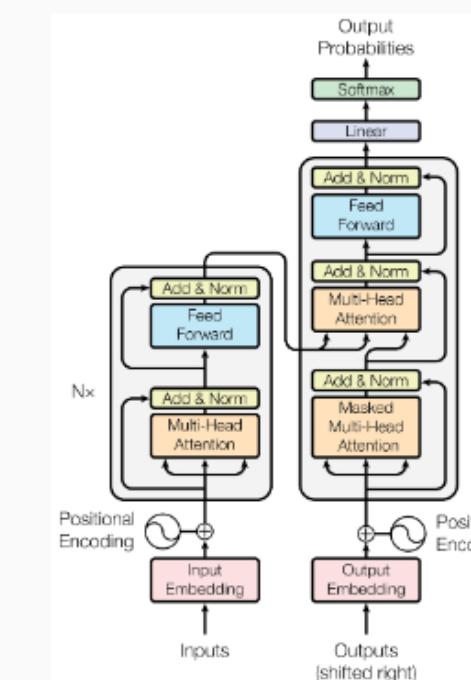
Reinf. Learning



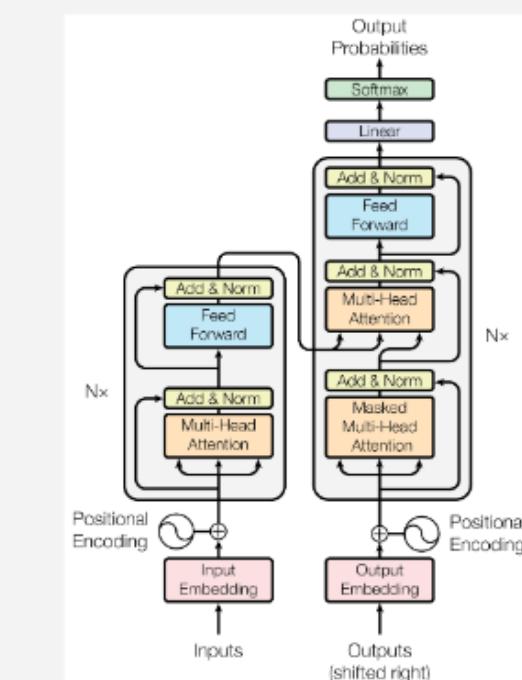
Speech



Translation

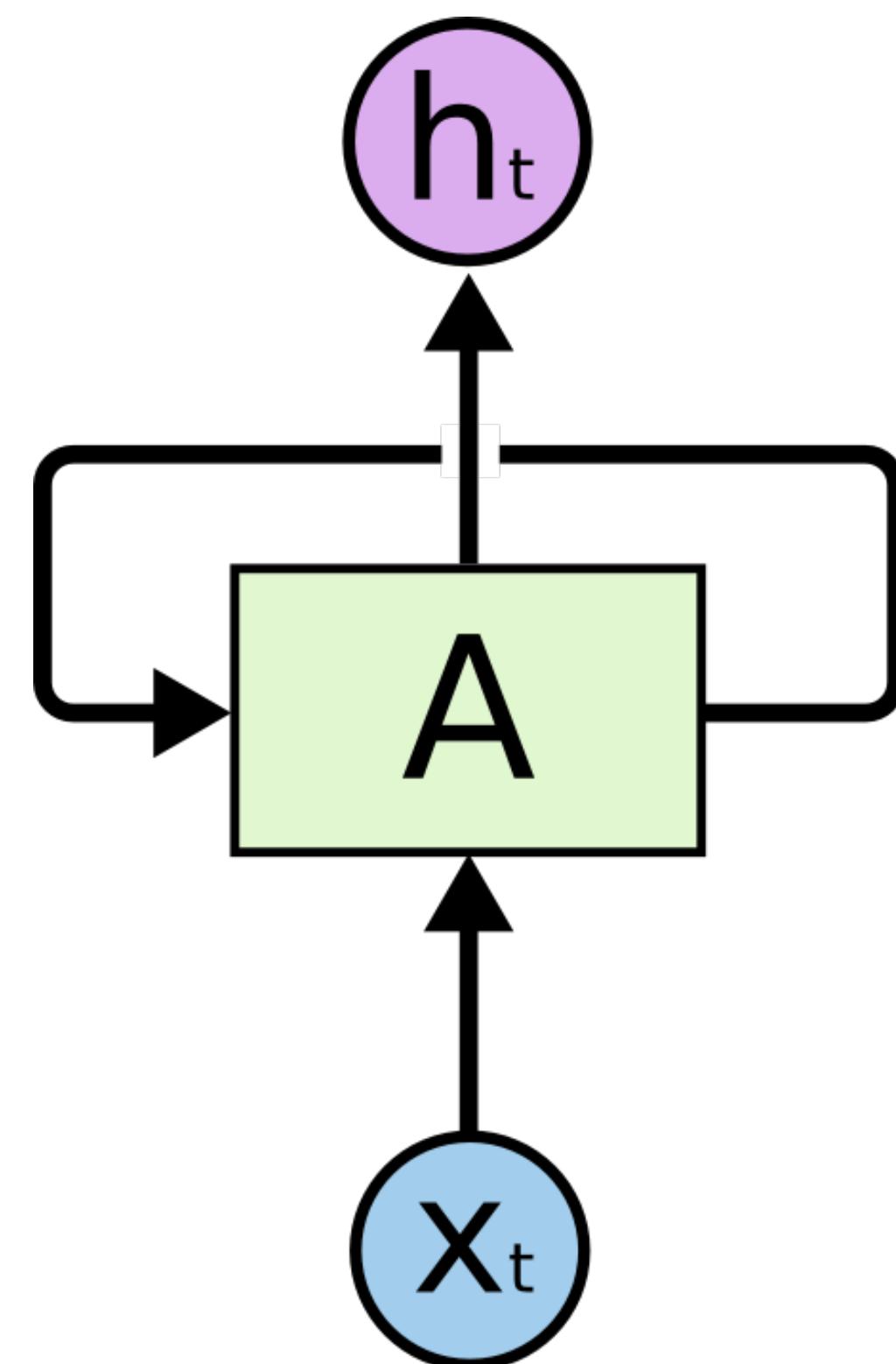


Graphs/Science



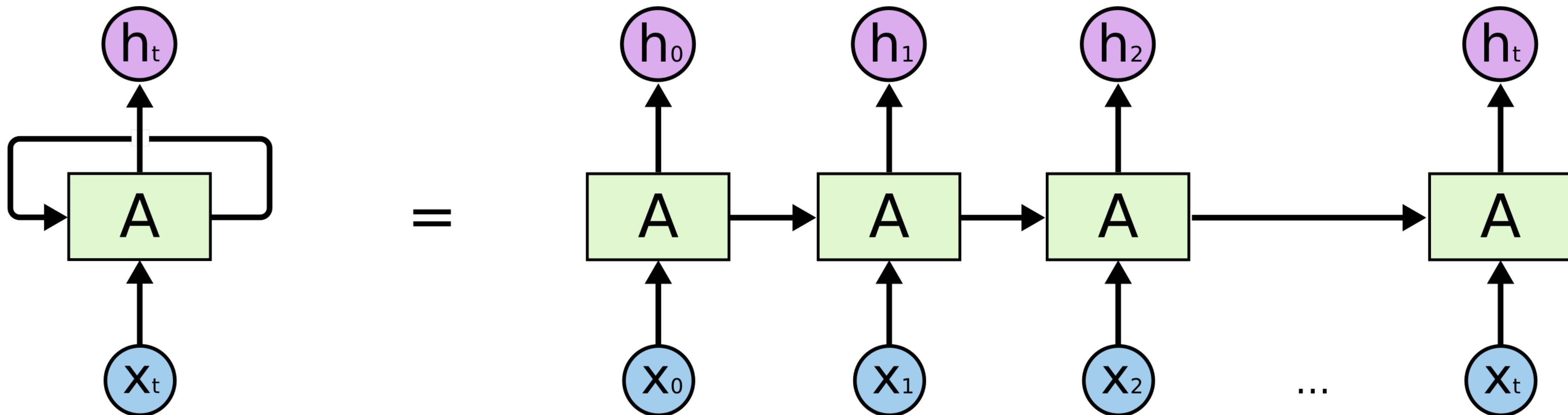
RNN: Recurrent Neural Networks

Making predictions with respect to time



RNN: Recurrent Neural Networks

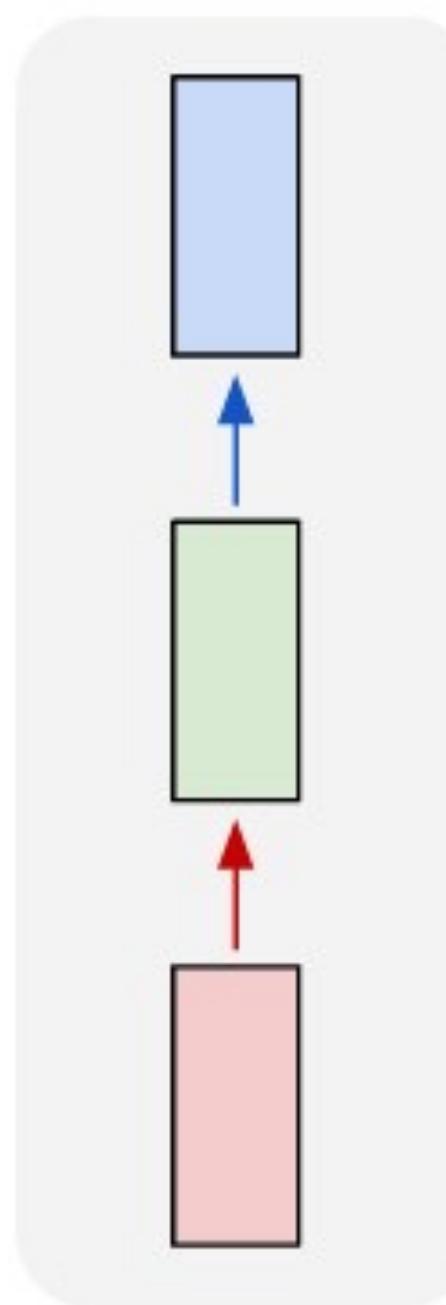
Making predictions with respect to time



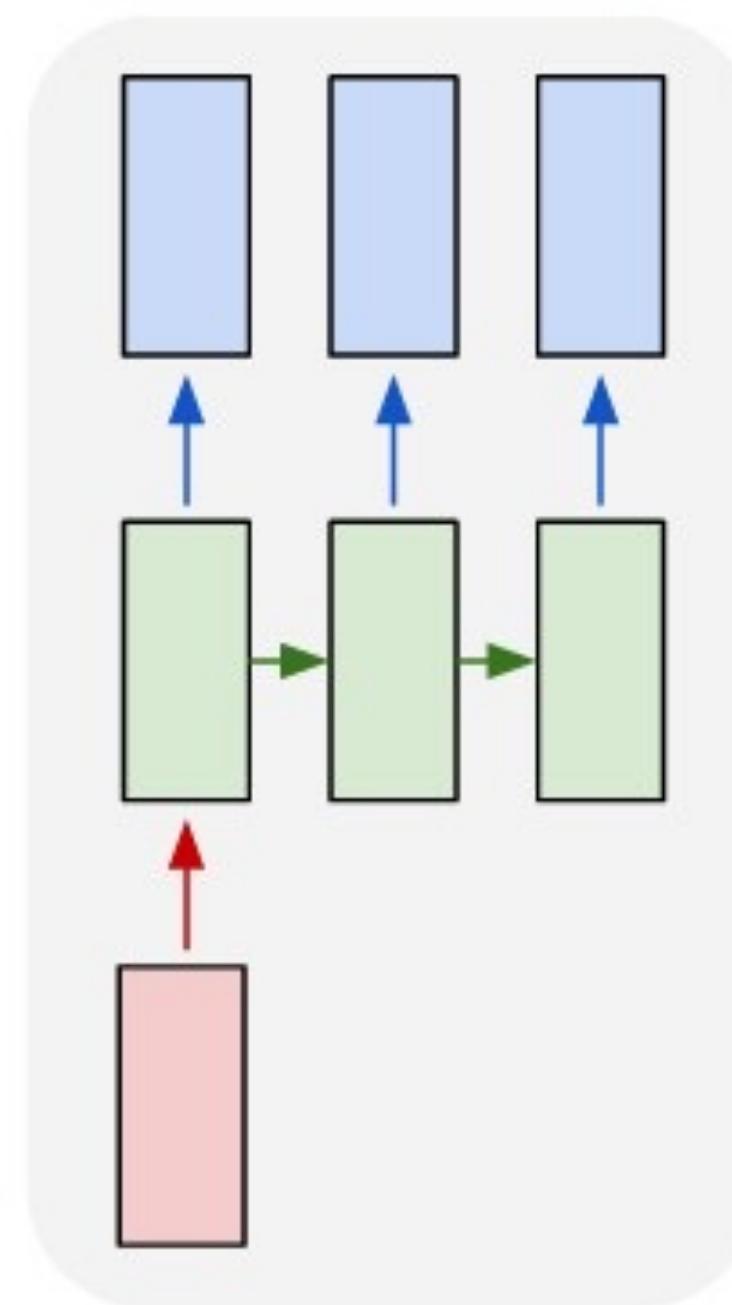
Different tasks, different architectures

Making predictions with respect to time

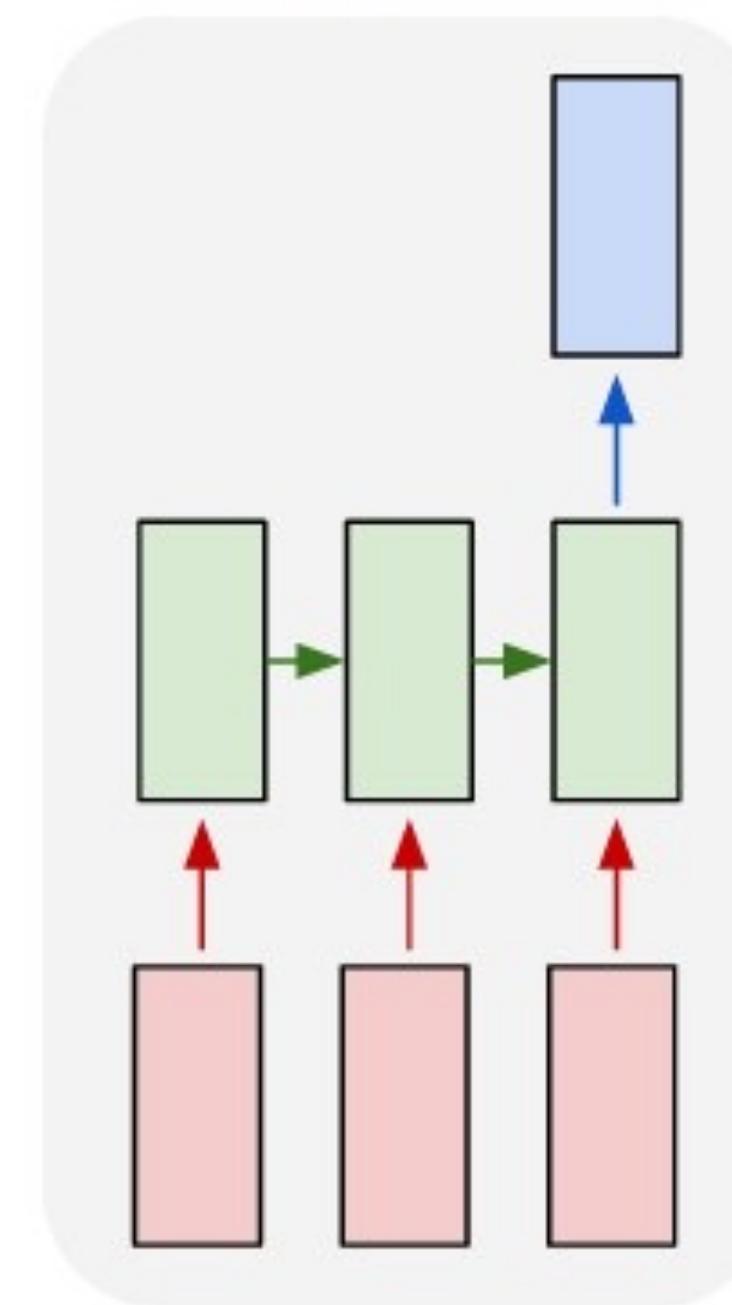
one to one



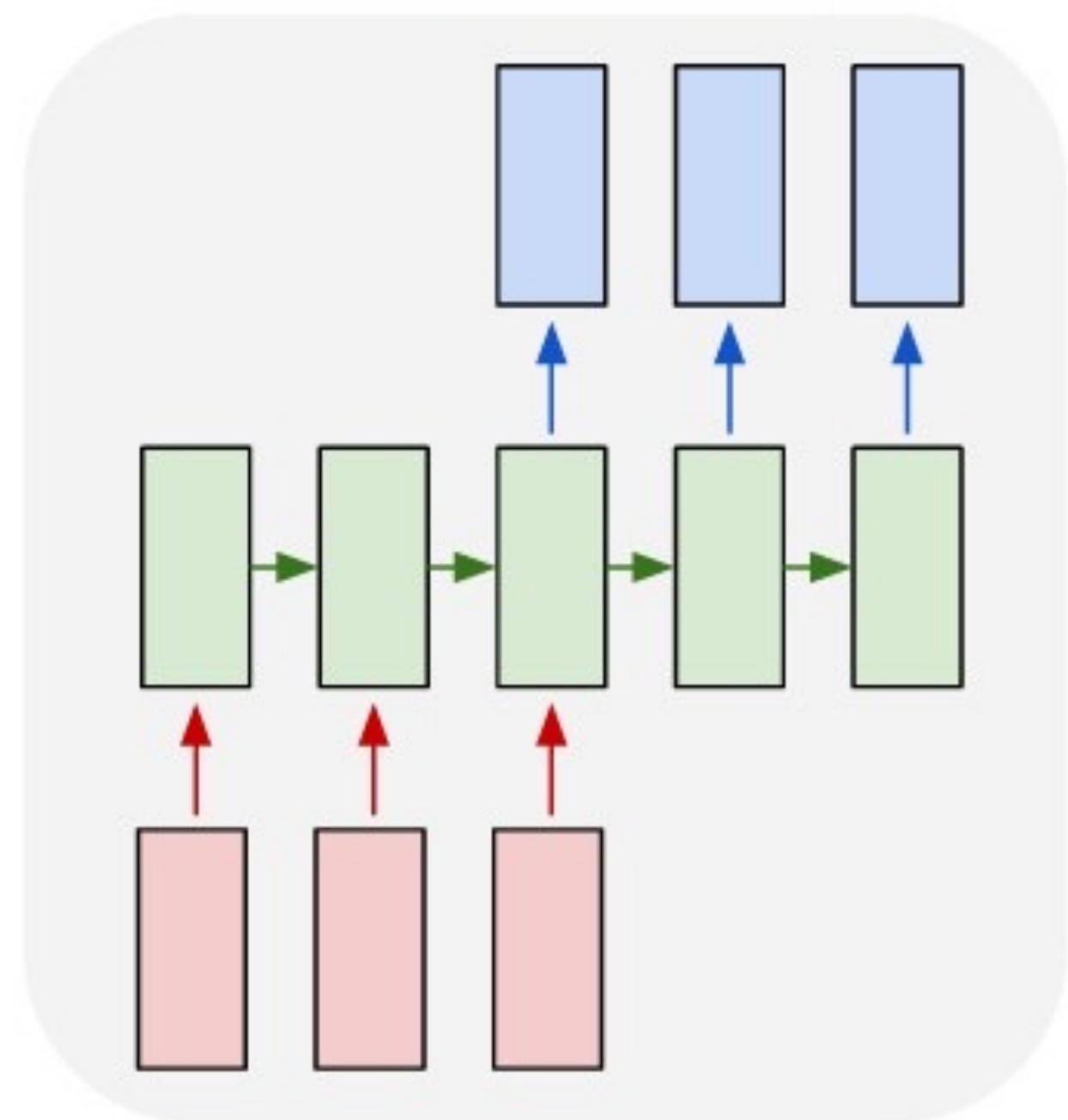
one to many



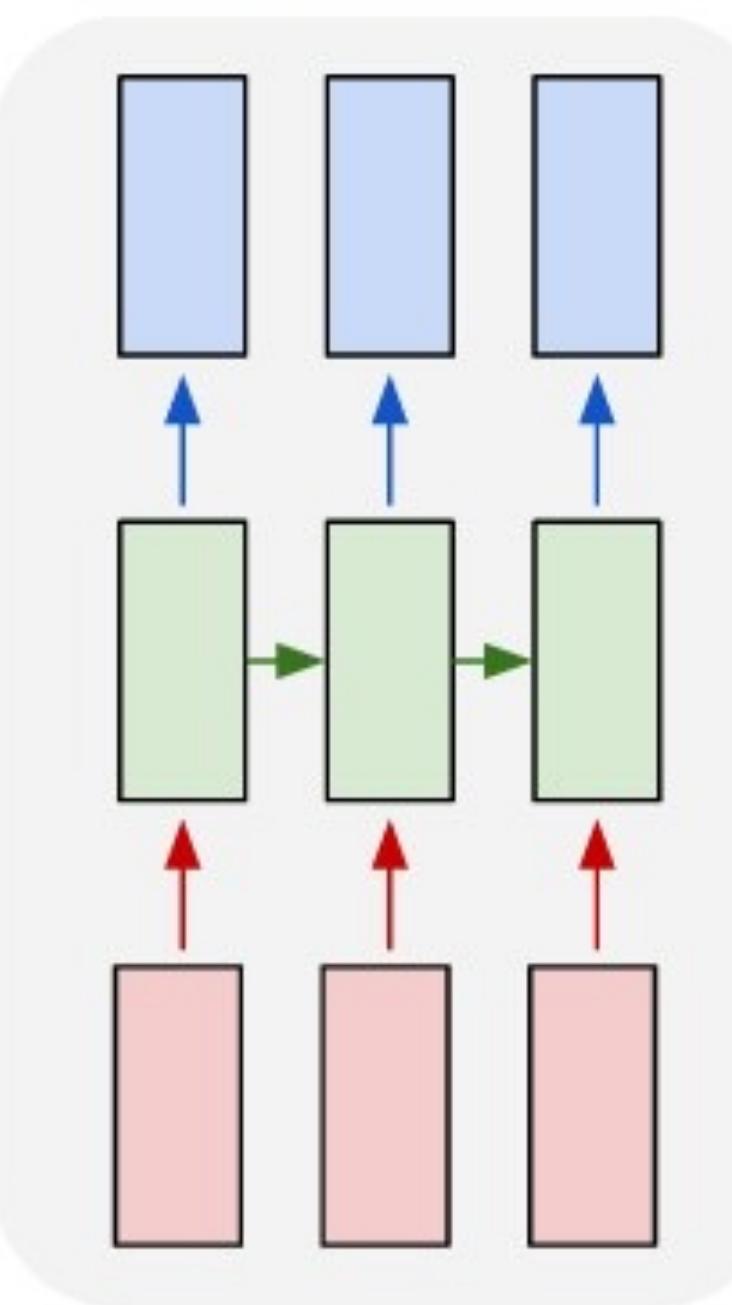
many to one



many to many

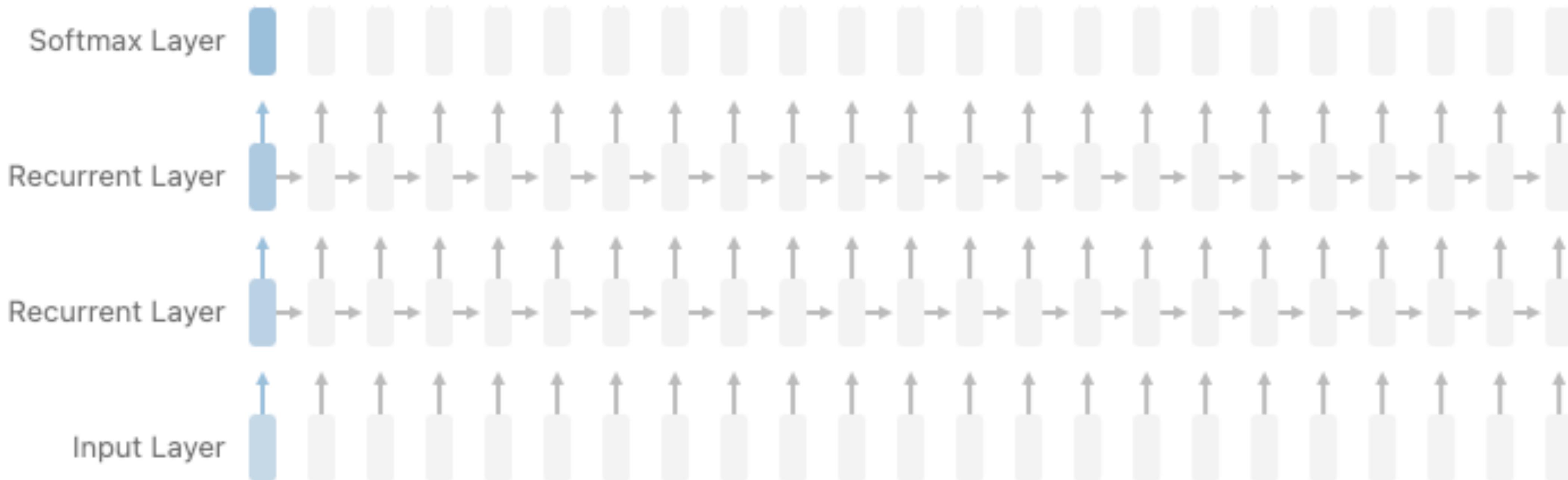


many to many



RNNs have problems

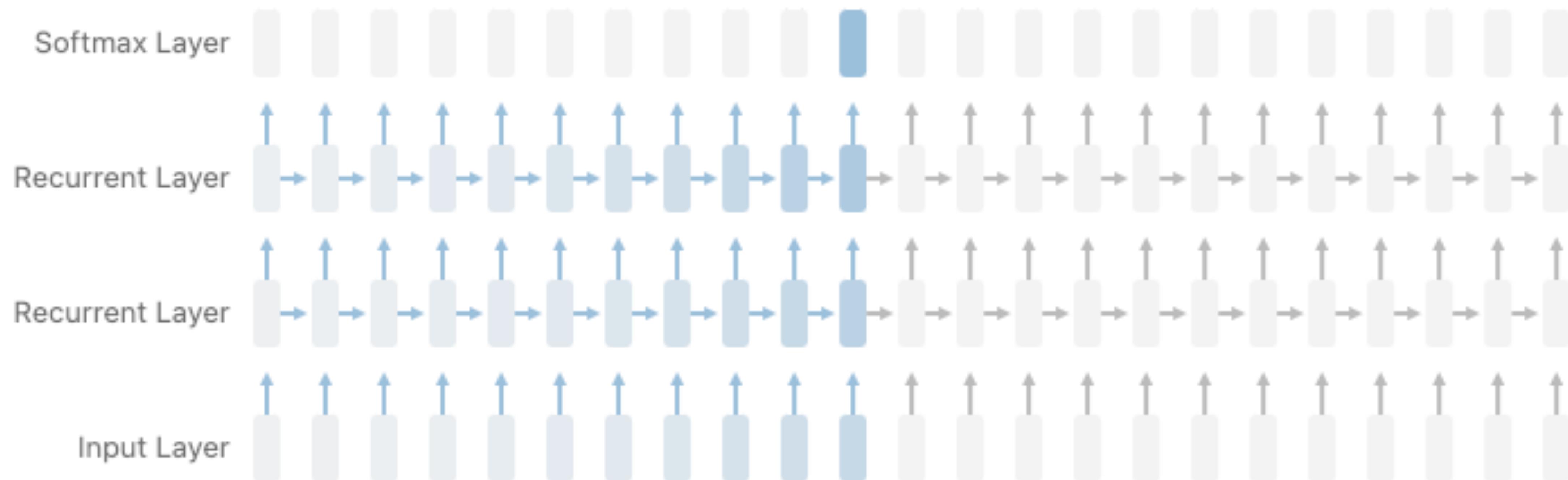
Vanishing Gradients cause short context lengths



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

RNNs have problems

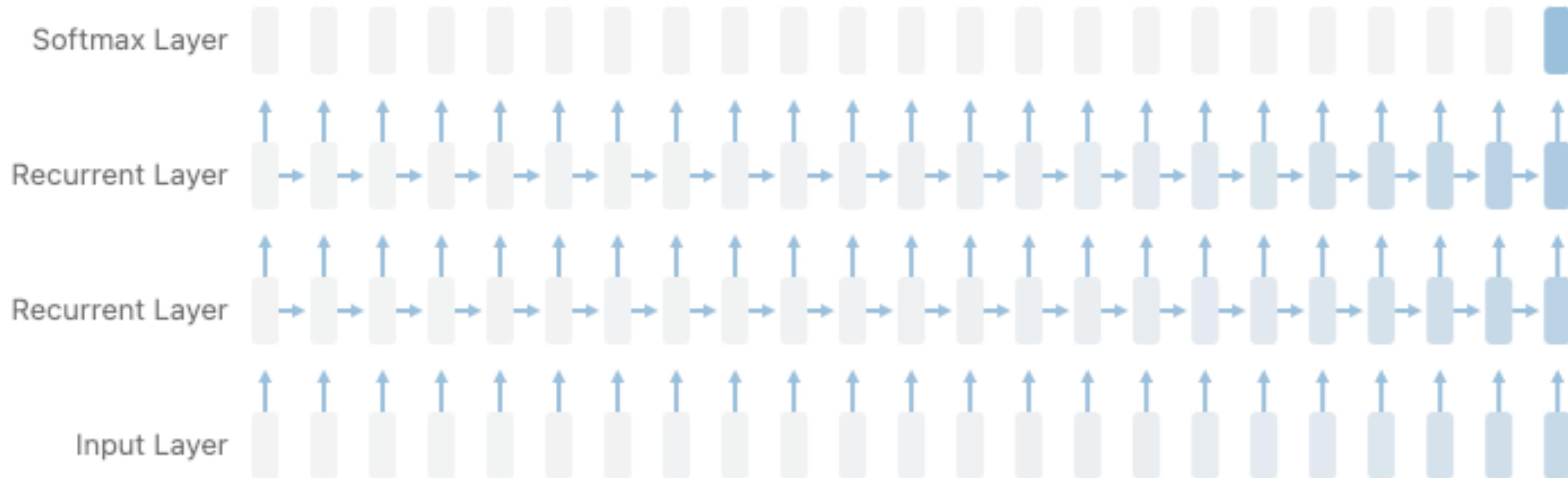
Vanishing Gradients cause short context lengths



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

RNNs have problems

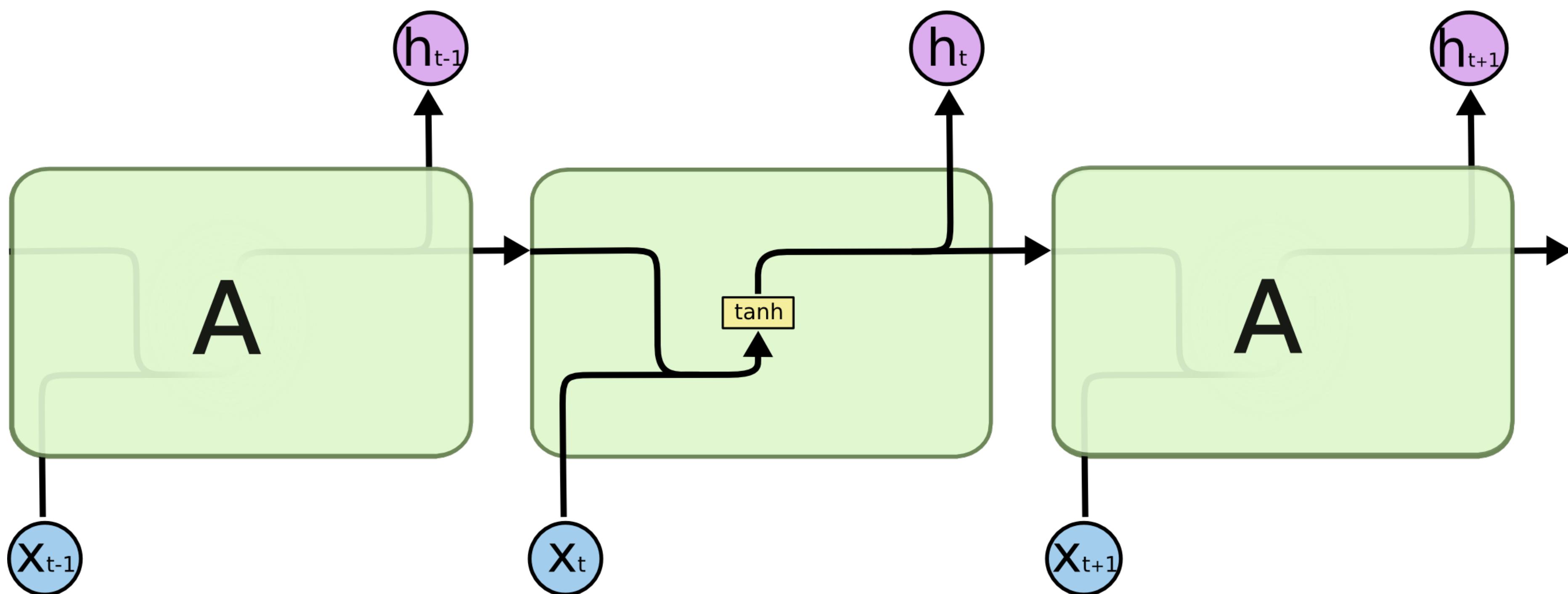
Vanishing Gradients cause short context lengths



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

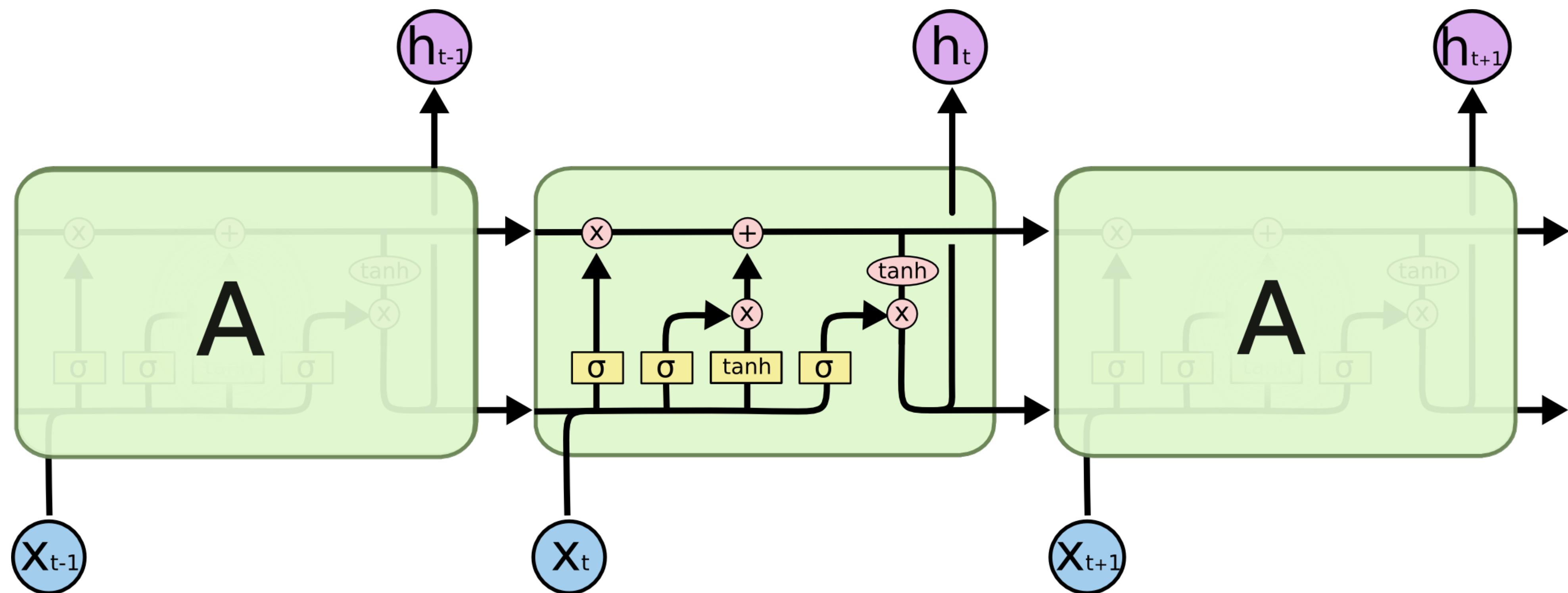
RNN variants tackle vanishing gradients

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

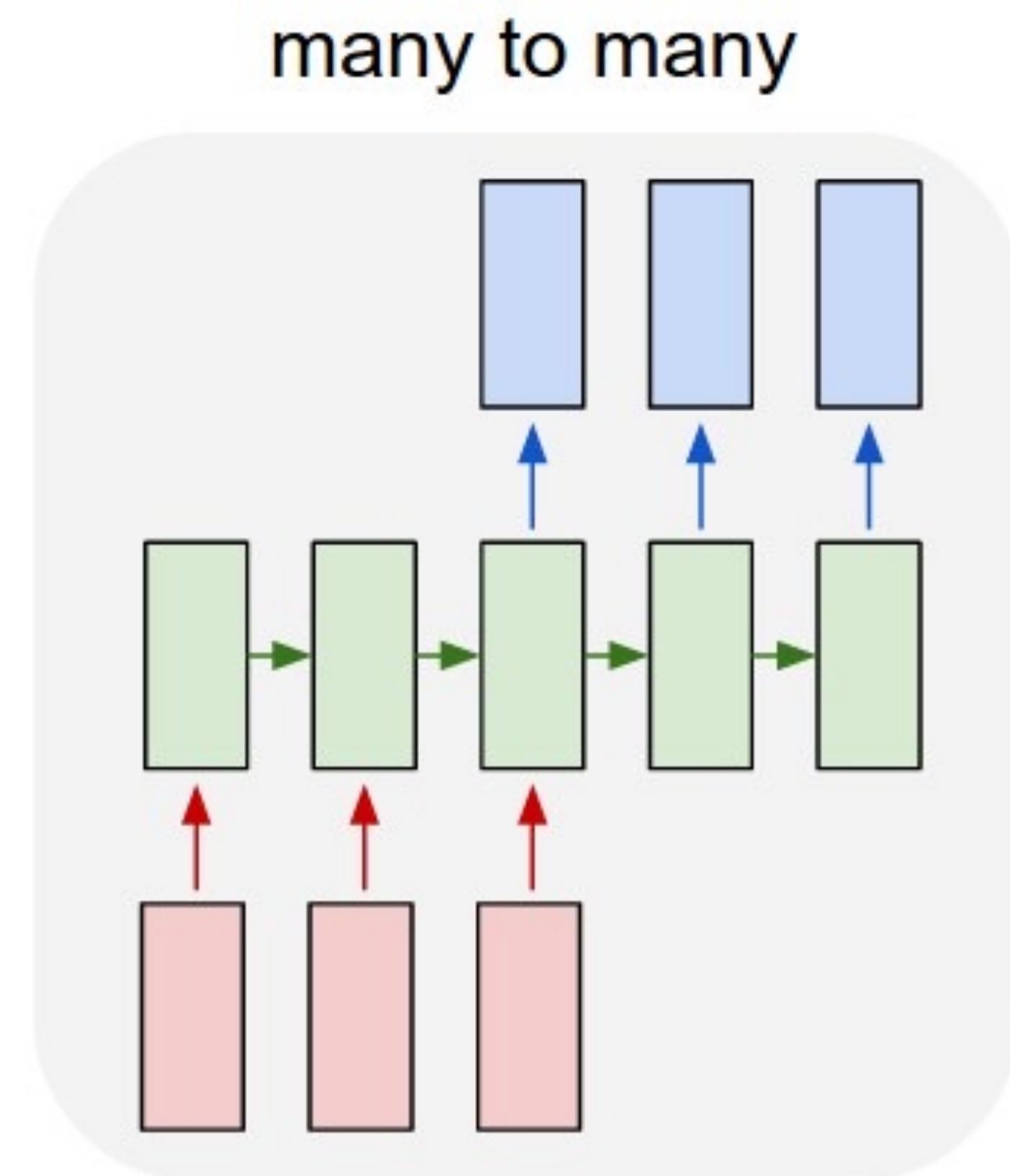
LSTM

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

GRU

RNNs have other problems, too

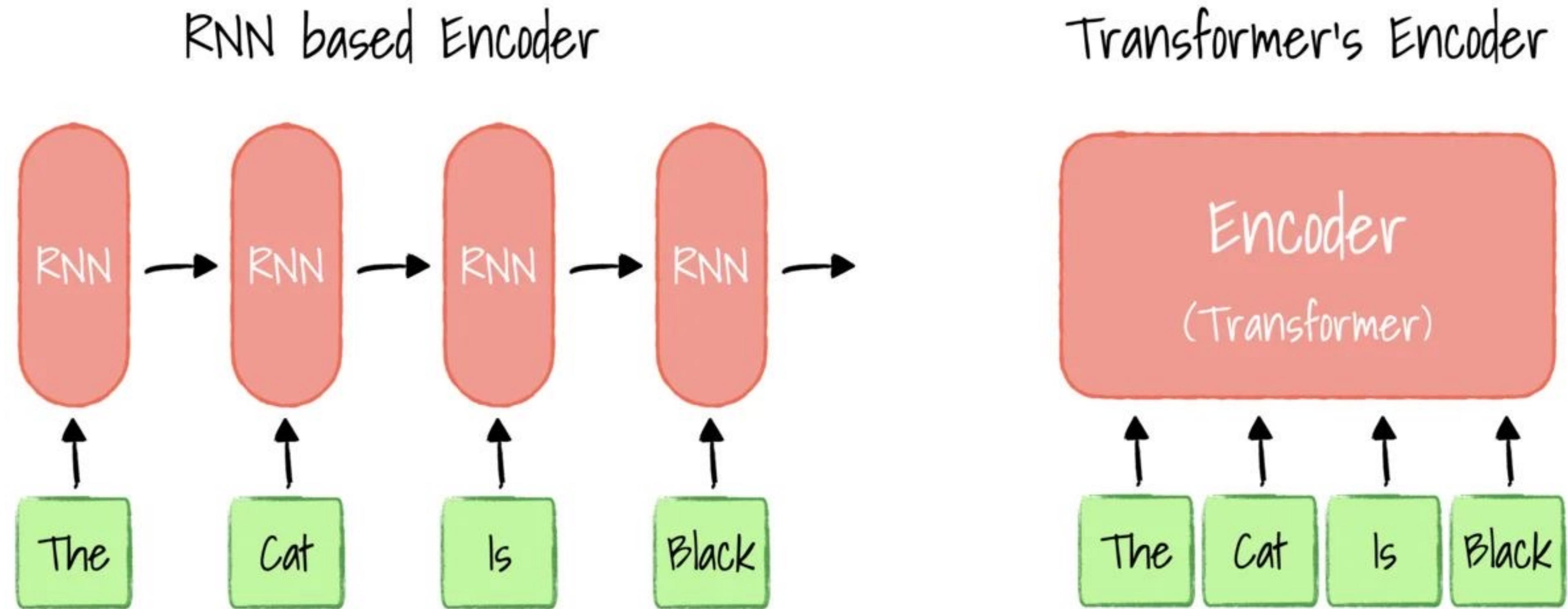
No parallelisation possible



3. Transformers

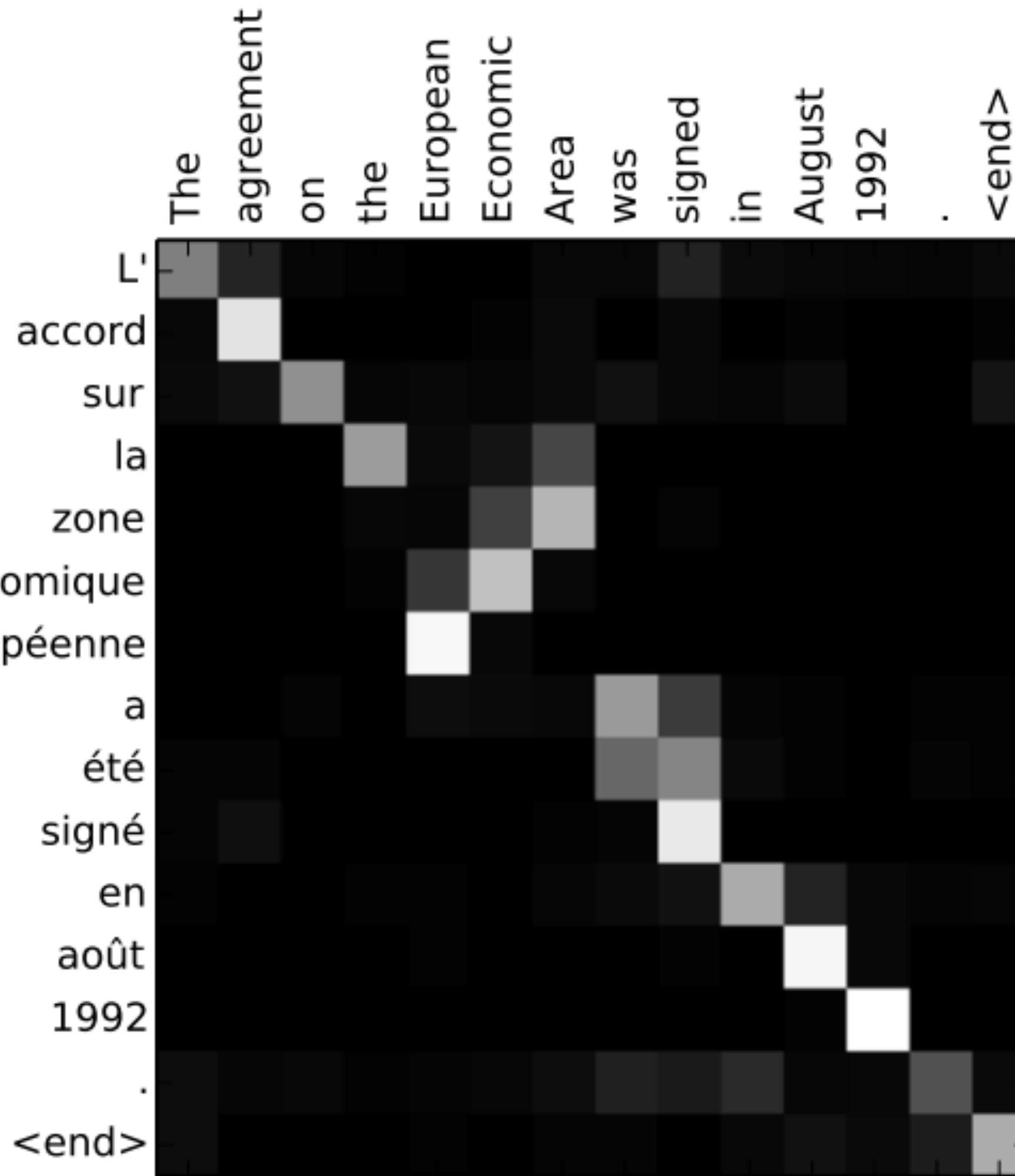
Transformers to the rescue

Parallel instead of sequential encoding with attention



What is Attention?

Allowing every word to be influenced by any other word



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

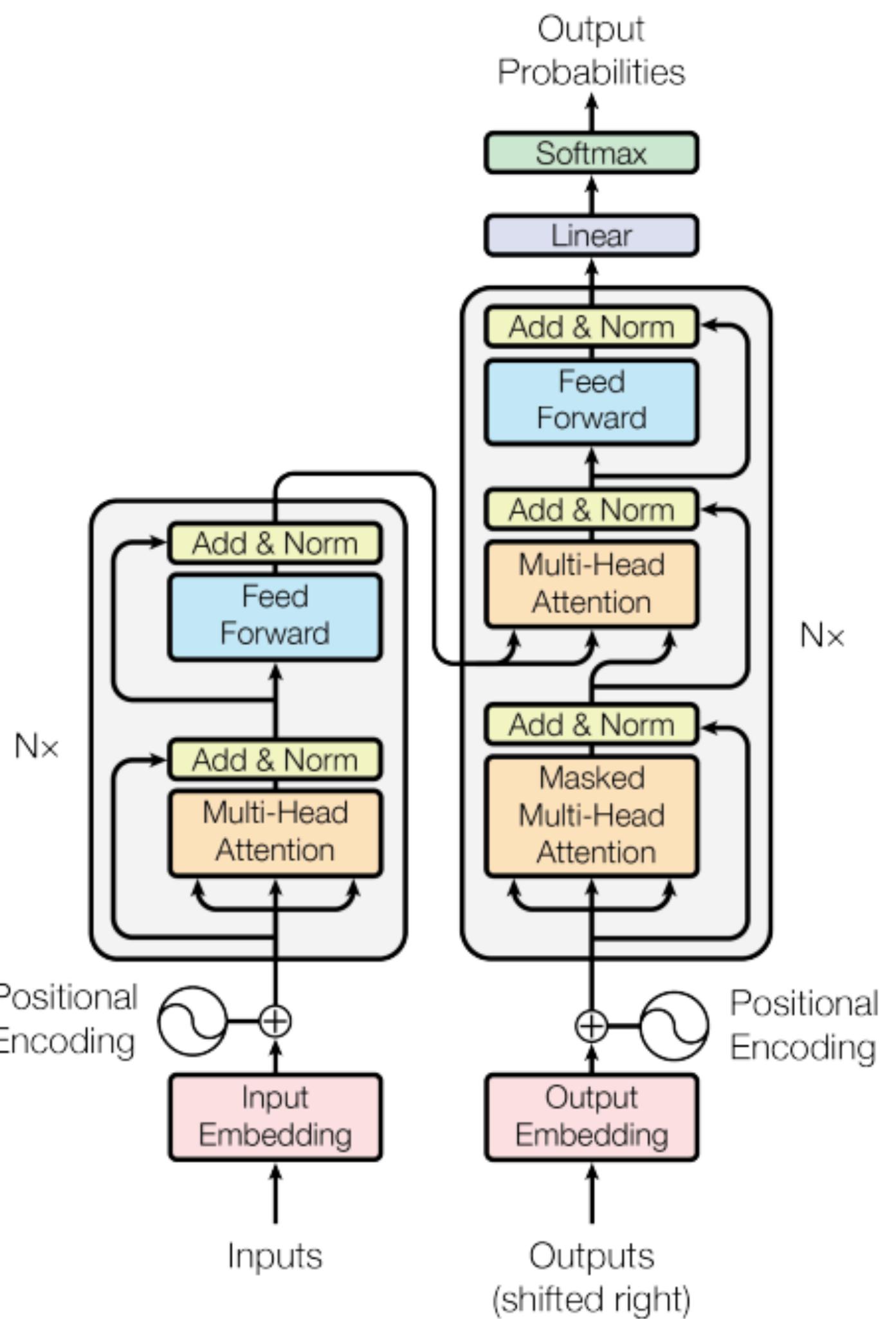
lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

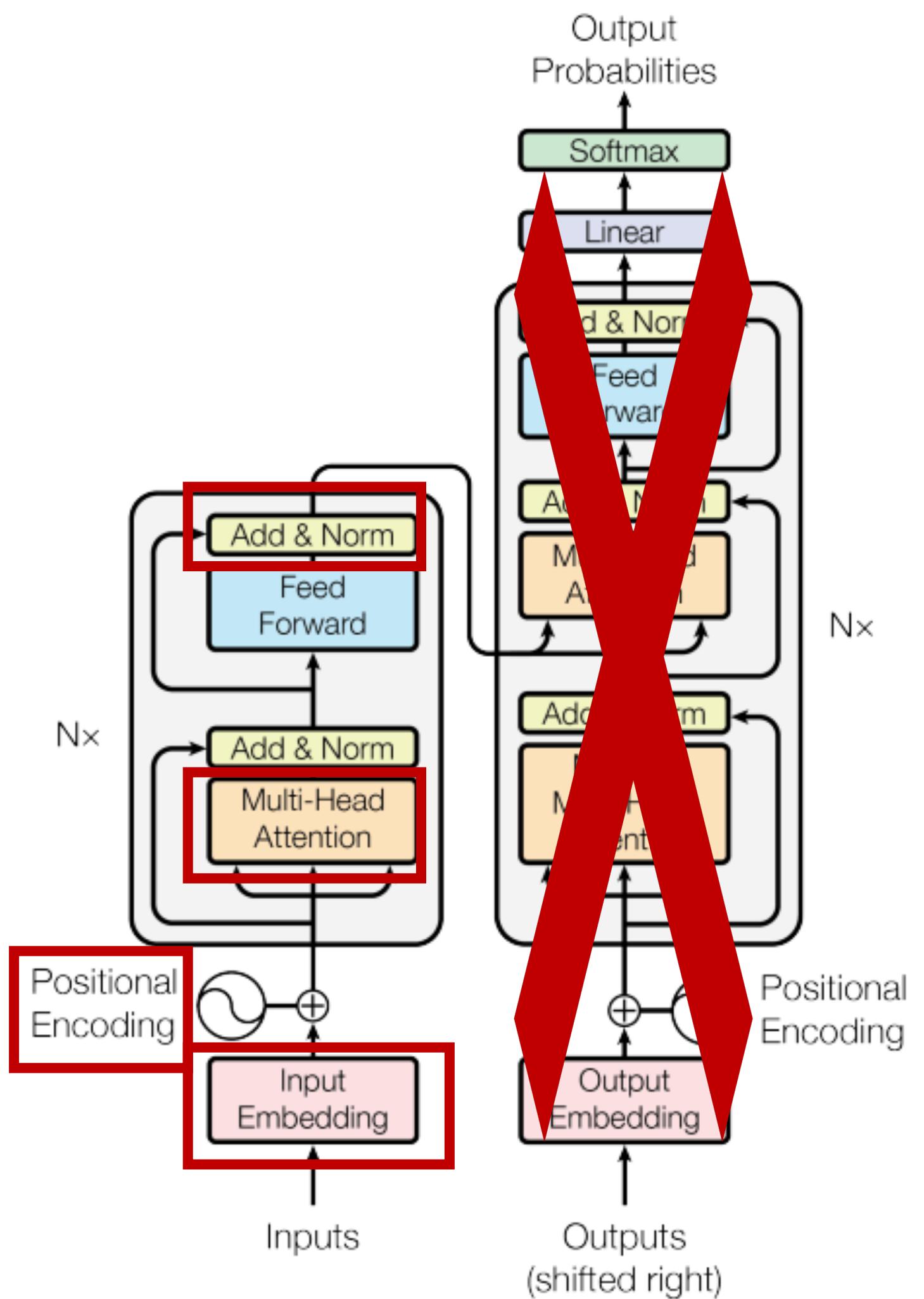
The Transformer

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/>

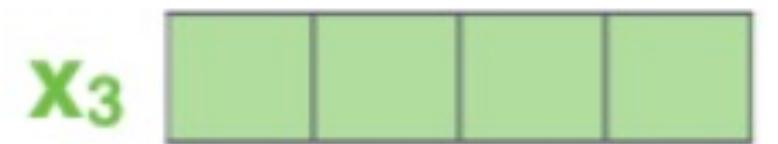
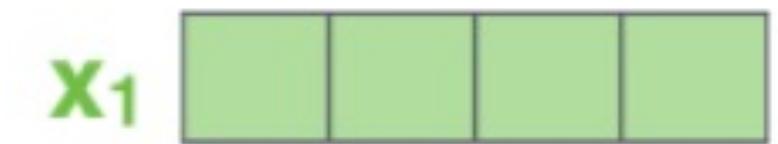
<https://ronxin.github.io/wevi/>

<http://projector.tensorflow.org/>

Positional Embedding

We must tell our computer what comes first and what later

EMBEDDINGS



INPUT

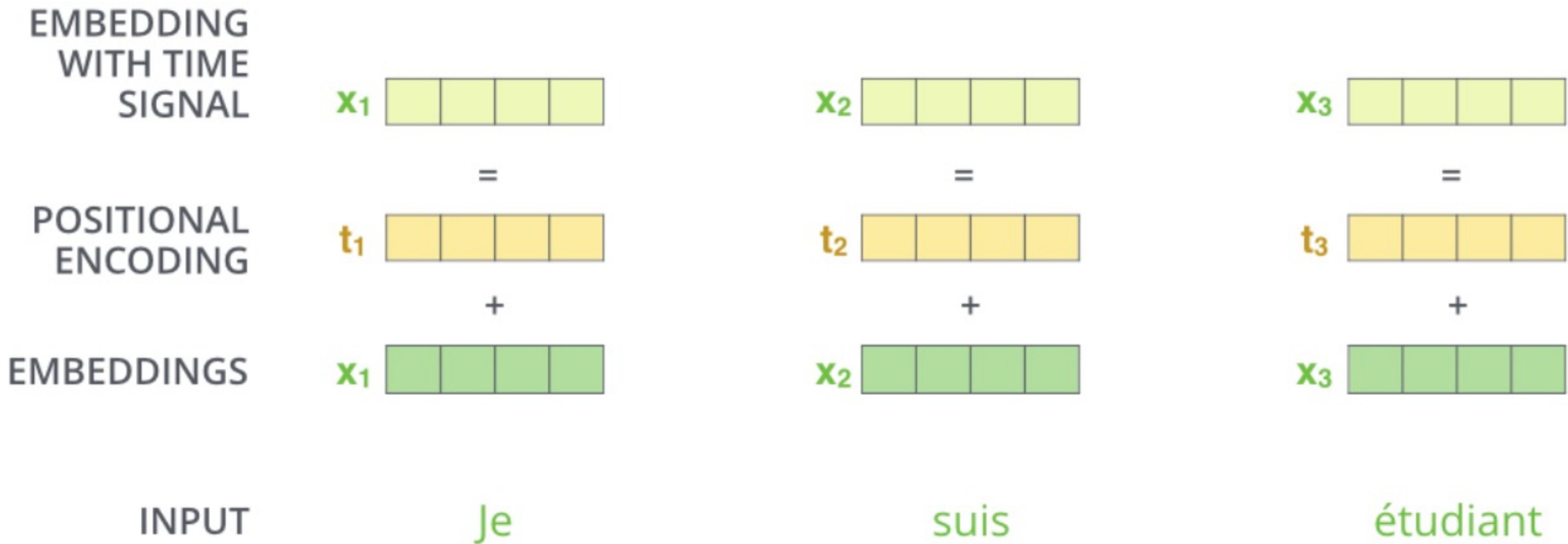
Je

suis

étudiant

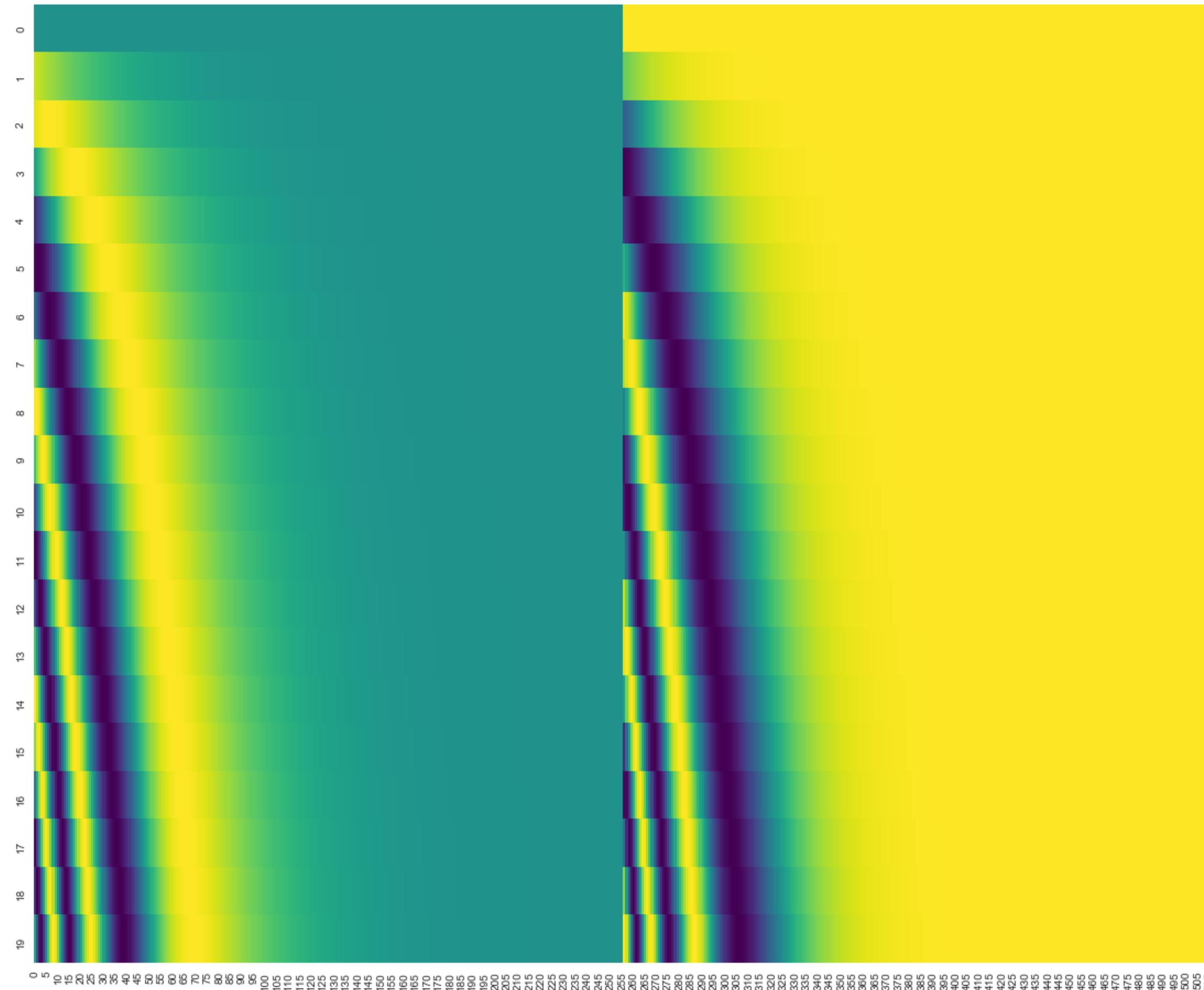
Positional Embedding

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
 $x_1, x_2, \dots x_t$

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

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$$

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

$$\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

- weight is just a dot product $w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$

Attention

Looking at everyone around you to determine your update

- Input: sequence of tensors

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$$

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

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

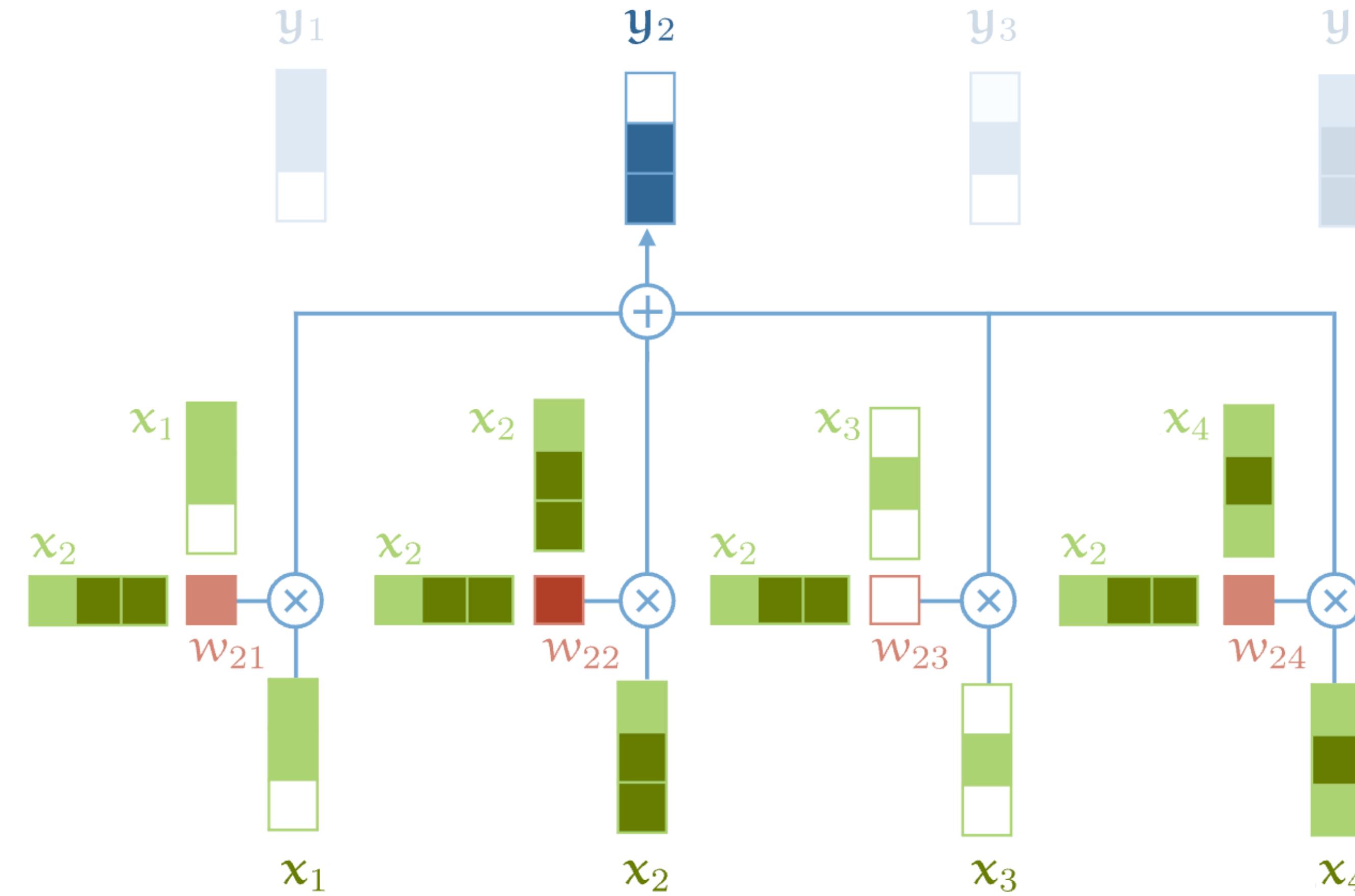
- weight is just a dot product $w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$

- make it sum to 1

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

Attention

Looking at everyone around you to determine your update

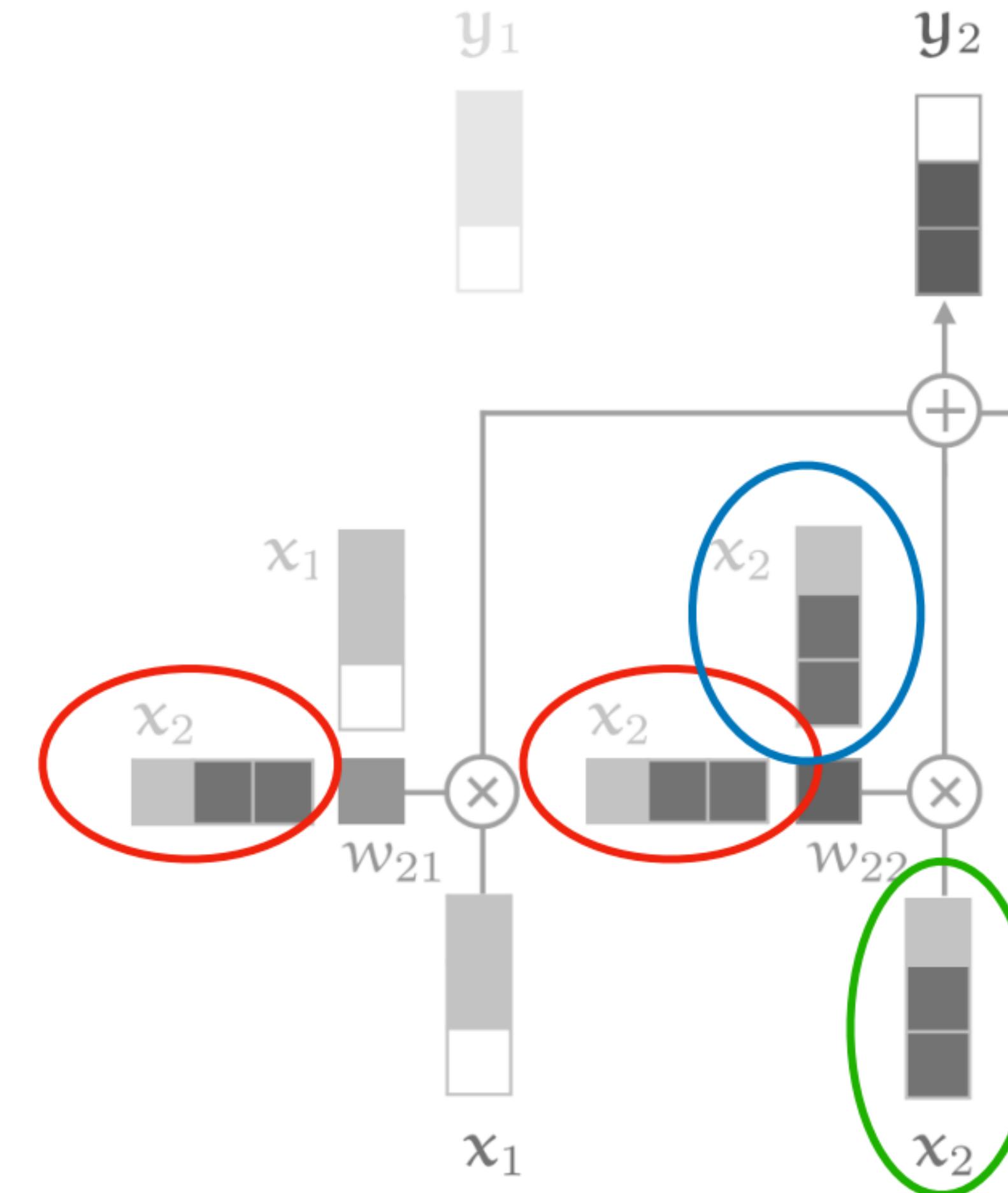


Attention

Learning the weights

Query, Key, Value

- Every input vector x_i is used in 3 ways:
 - Query
 - Key
 - Value

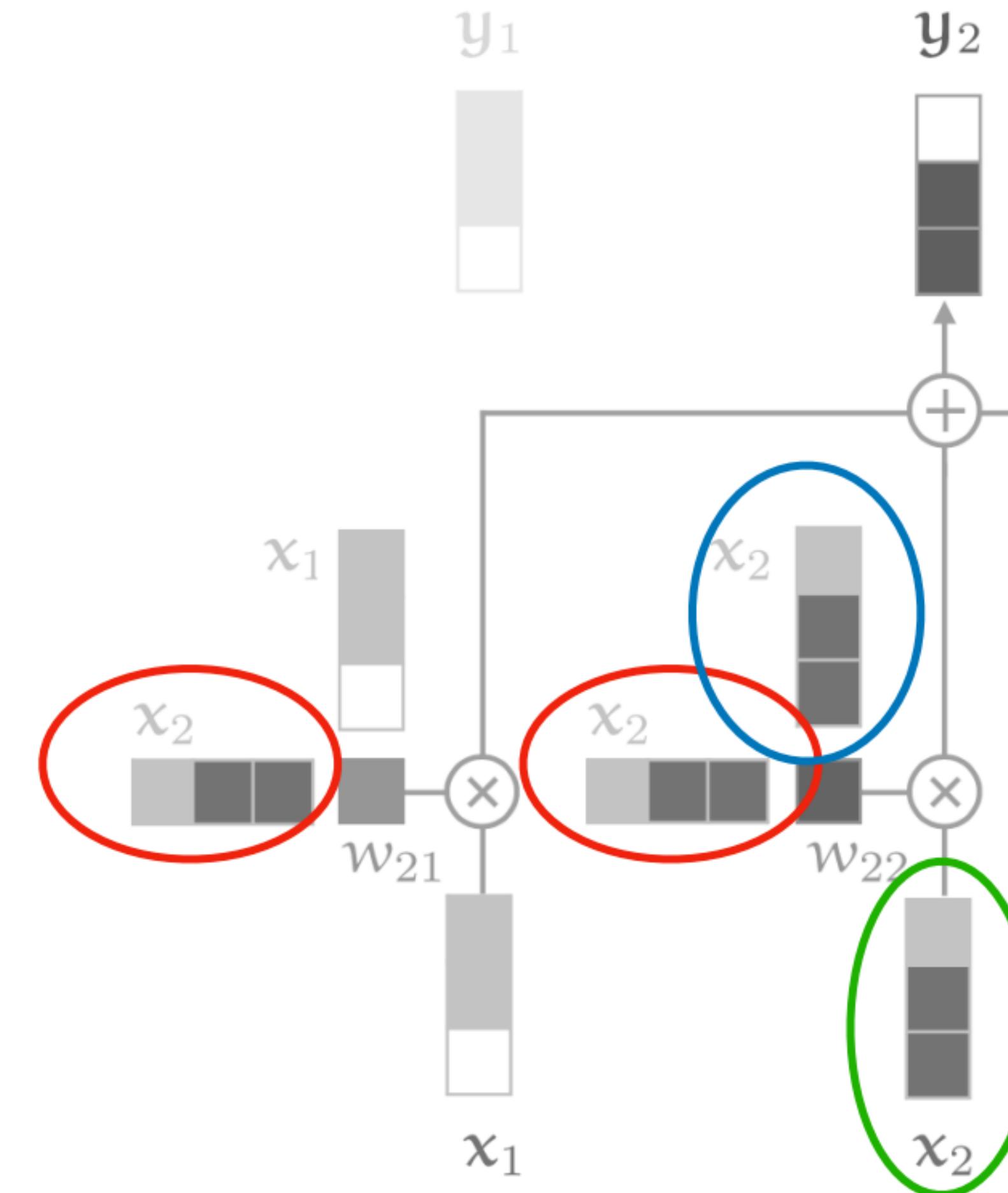


Attention

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?

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$$

What do I have?

$$\mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i$$

What do I reveal/give to others?

$$\mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

$$w'_{ij} = \mathbf{q}_i^\top \mathbf{k}_j$$

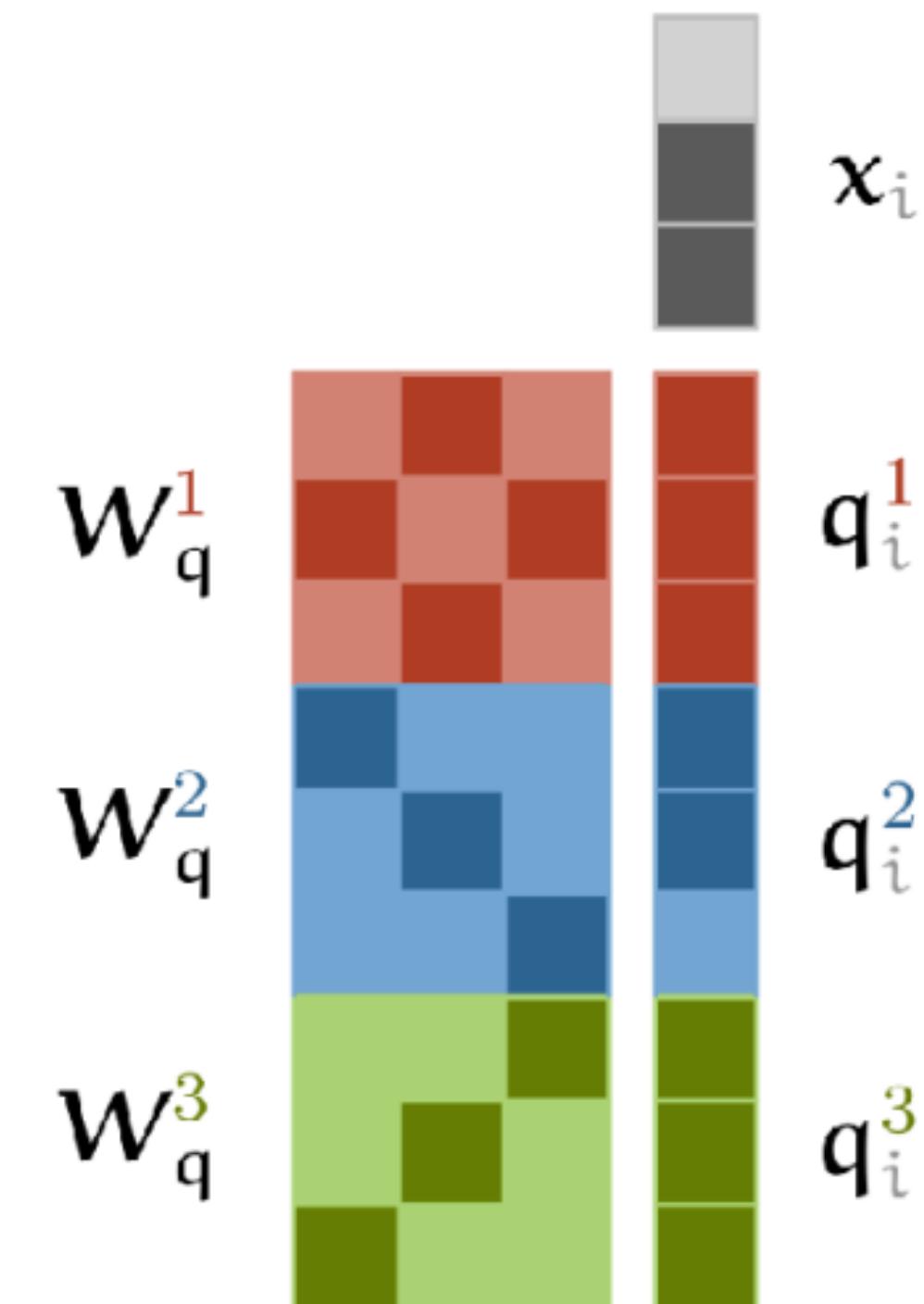
$$w_{ij} = \text{softmax}(w'_{ij})$$

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{v}_j .$$

Multi-head attention

Looking at everyone around you to determine your update

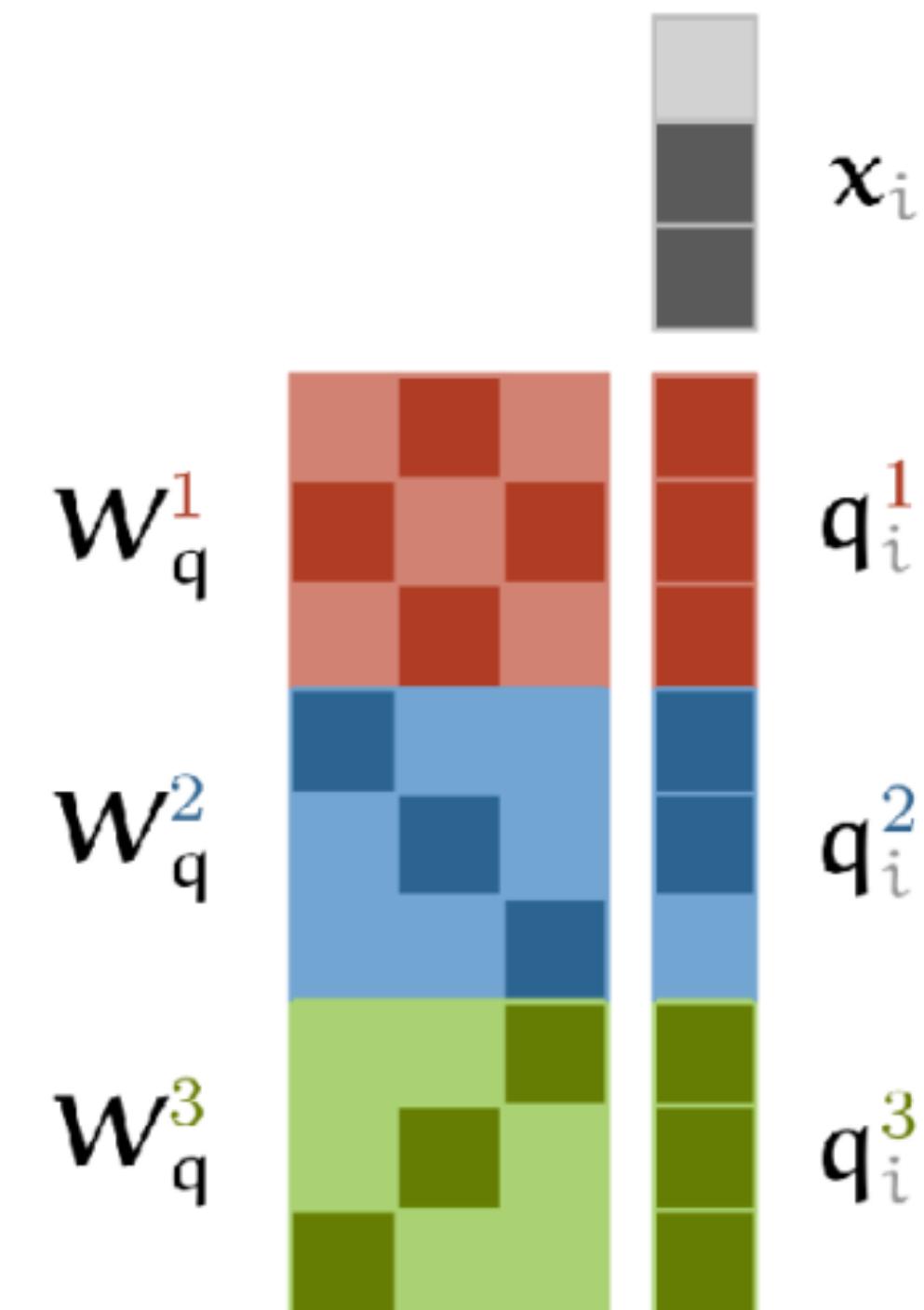
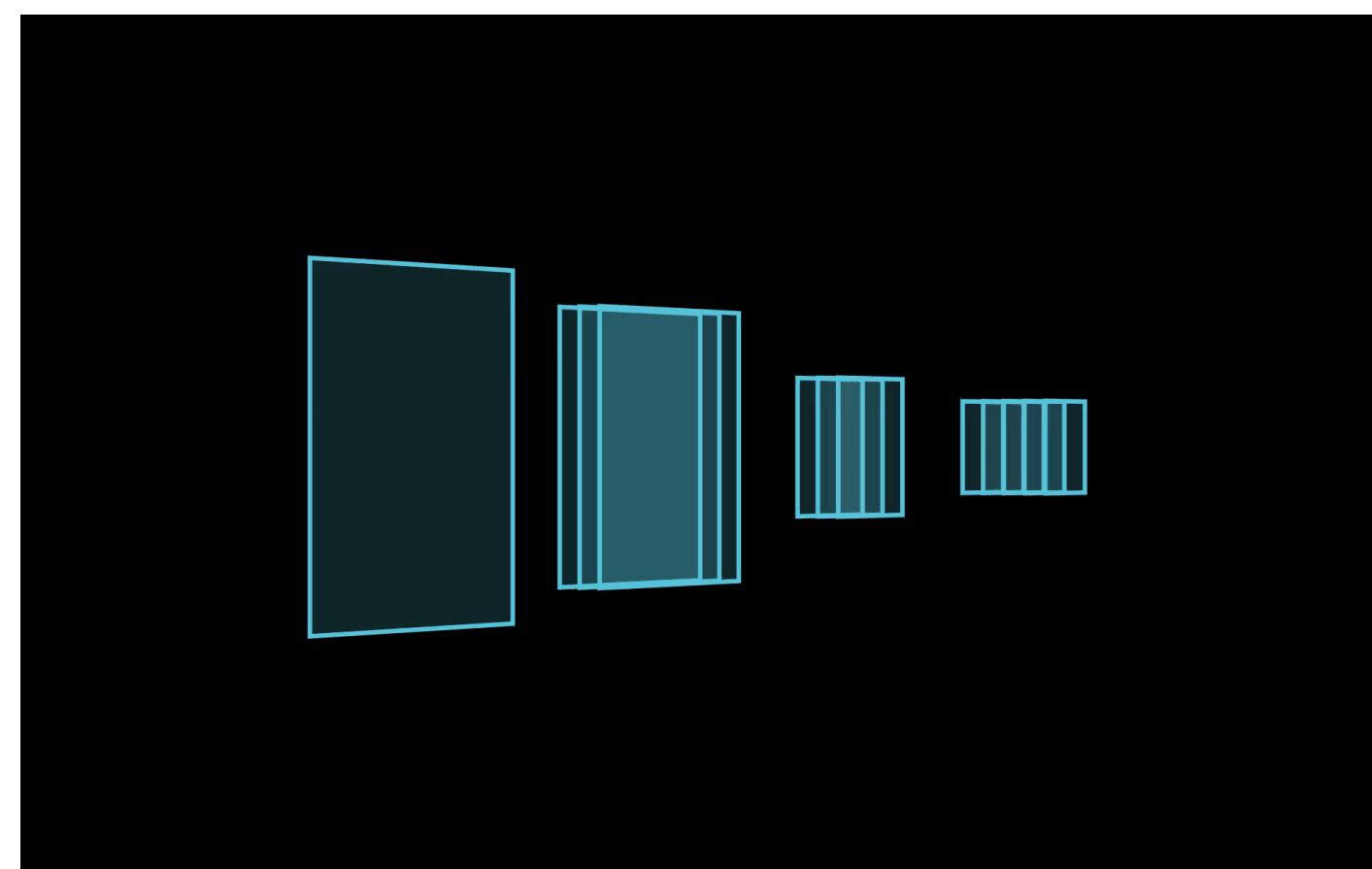
- Multiple "heads" of attention just means learning different sets of W_q , W_k , and W_v matrices simultaneously.
- Implemented as just a single matrix...



Multi-head attention

Looking at everyone around you to determine your update

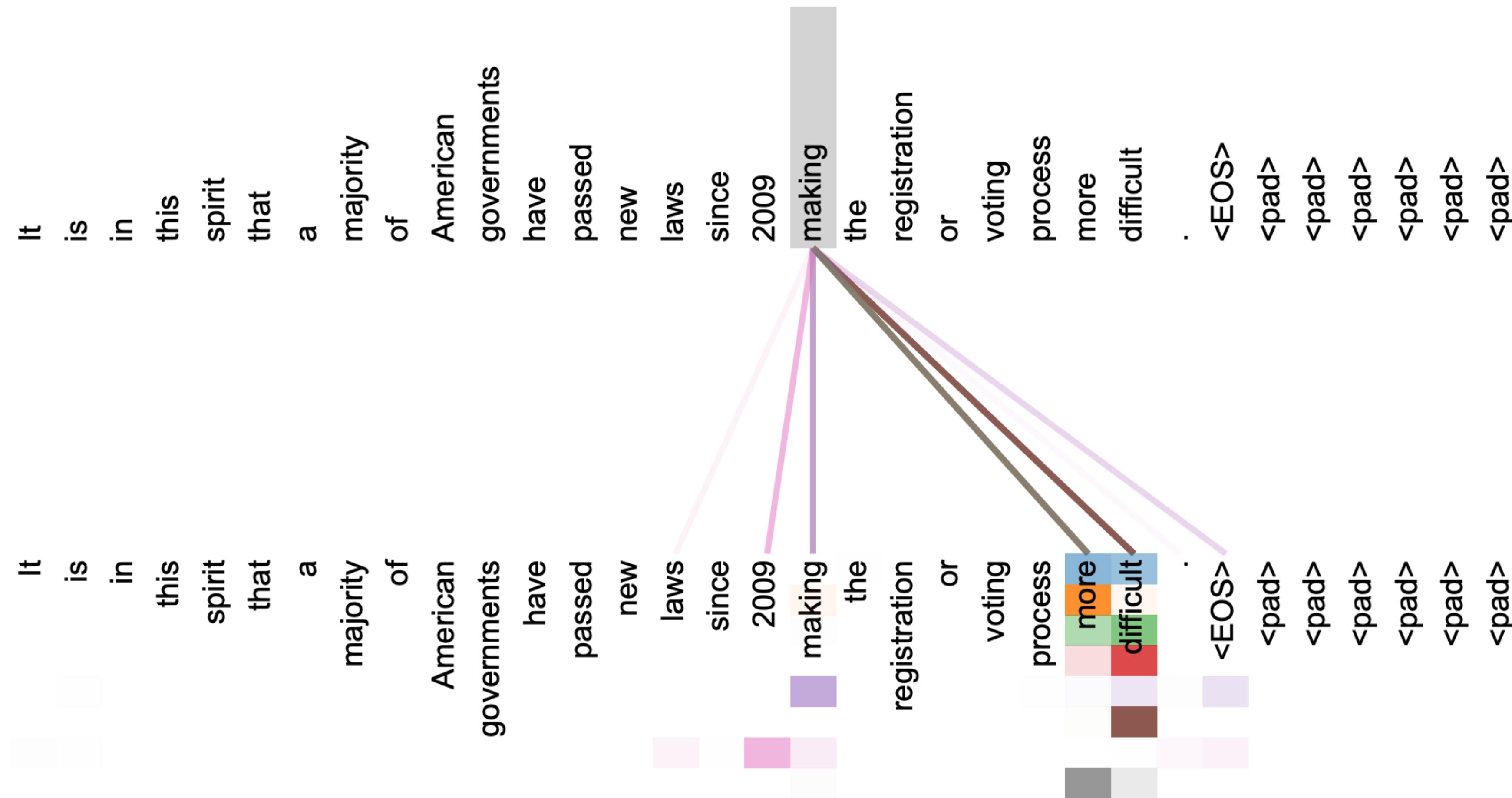
- Multiple "heads" of attention just means learning different sets of W_q , W_k , and W_v matrices simultaneously.
- Implemented as just a single matrix...



Multi-head attention

Different heads attend to different parts in a sentence

Attention Visualizations

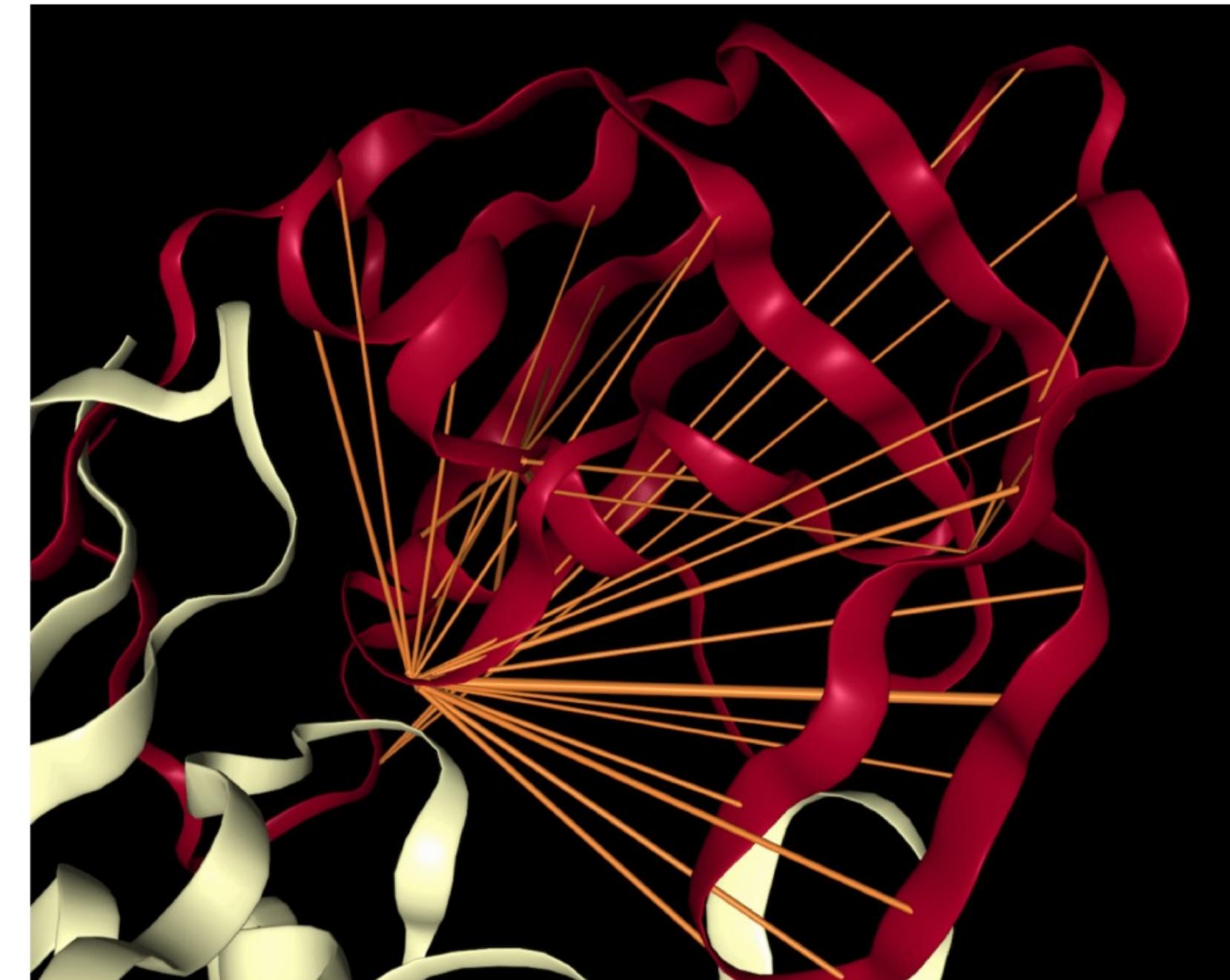


Multi-head attention

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.

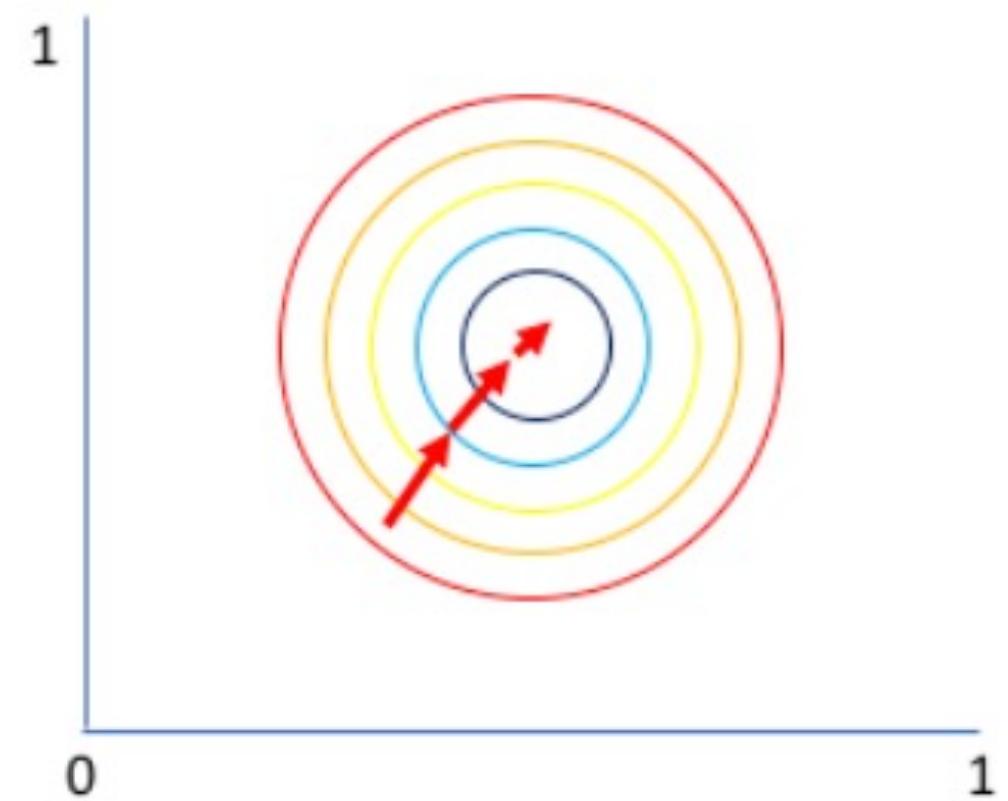


(b) Attention in head 7-1, which targets binding sites, a key functional component of proteins. Example is HIV-1 protease (7HVP). The primary location receiving attention is 27G, a binding site for protease inhibitor small-molecule drugs.

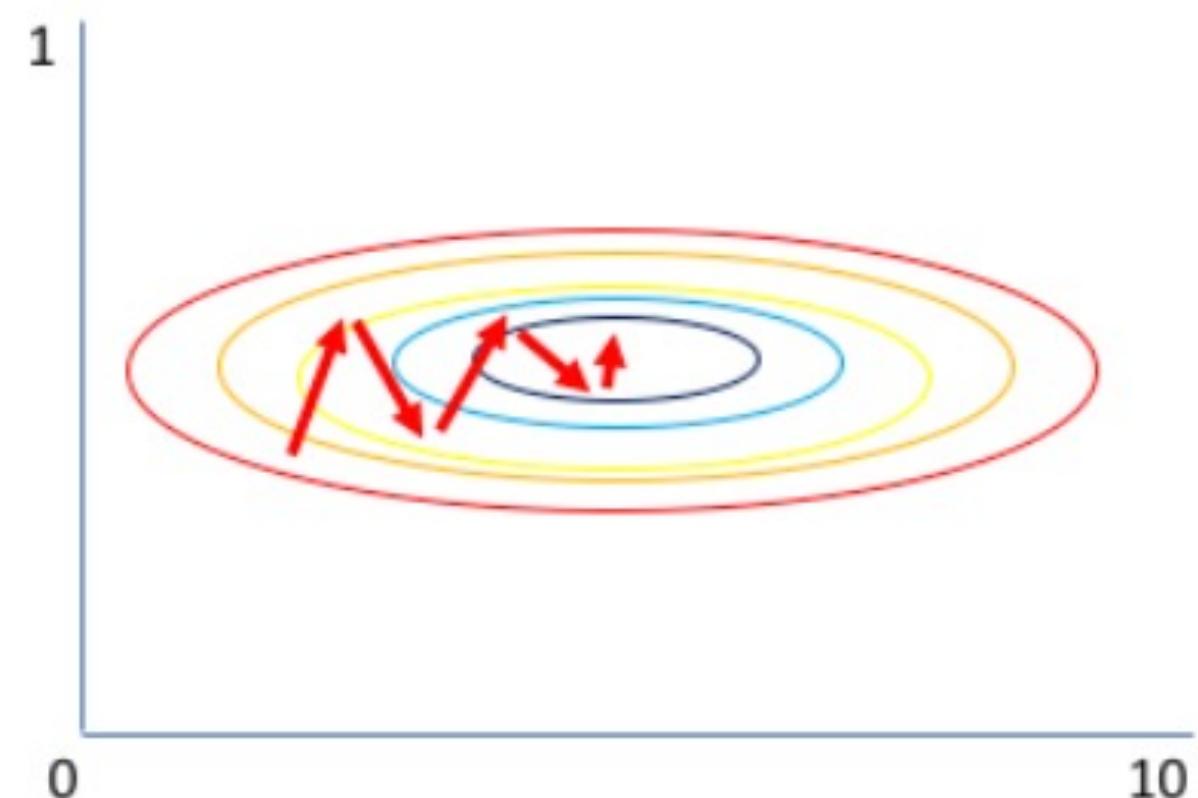
Layer Normalization

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.



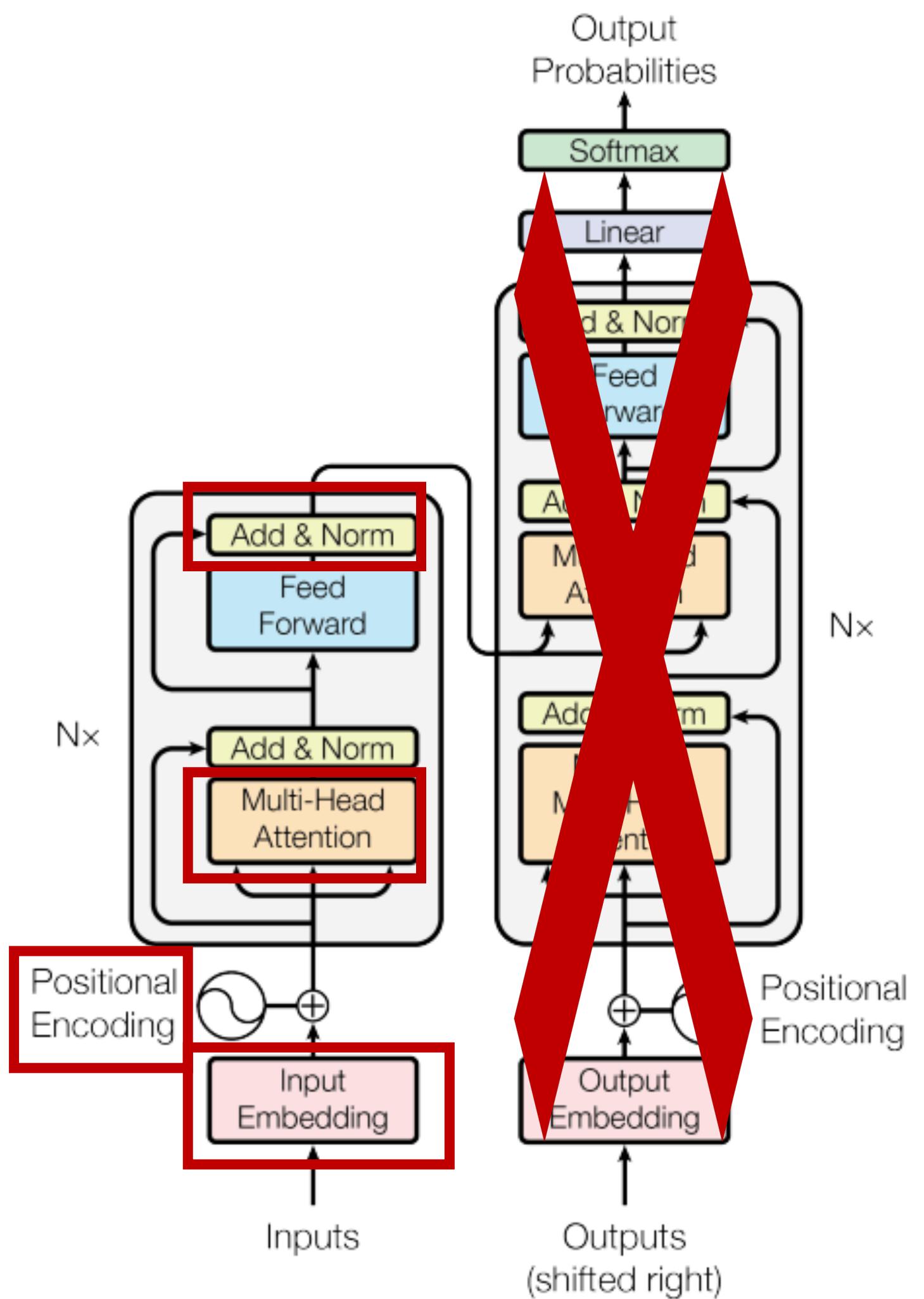
Both parameters can be updated in equal proportions



Gradient of larger parameter dominates the update

The Transformer

Not as scary as it looks like



Many good blogs about Transformers

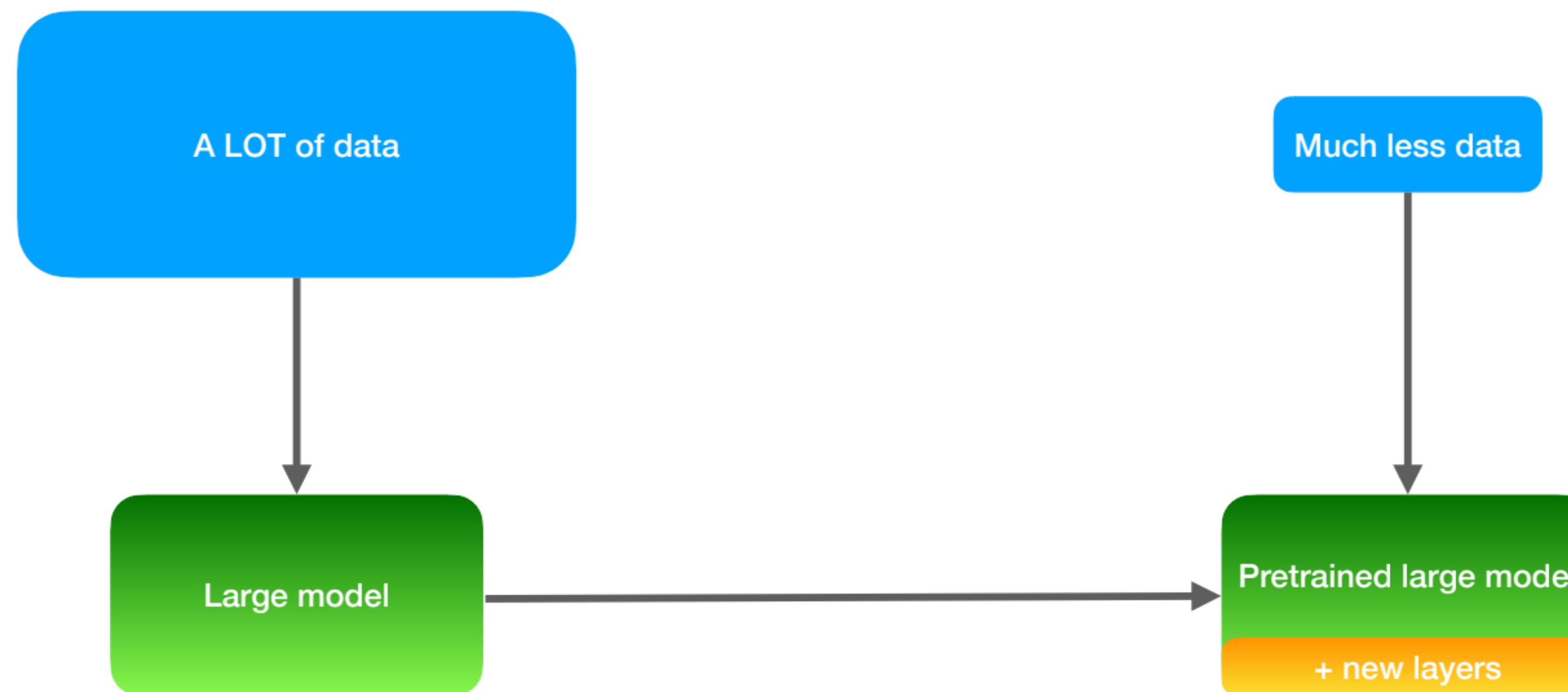
I leave it to you to choose the ones you like best

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. [Transformer Math](#) (If you want to implement a big one in practice)

4. Current Developments: Chat GPT and Proteins

Transformers are good at Transfer Learning

Use unlabeled data to get better on specific tasks



Traditional Machine Learning:
slow training on a lot of data

Transfer learning:
fast training on a little data

Language Models

Bigger = Better ?

Use the output of the masked word's position to predict the masked word

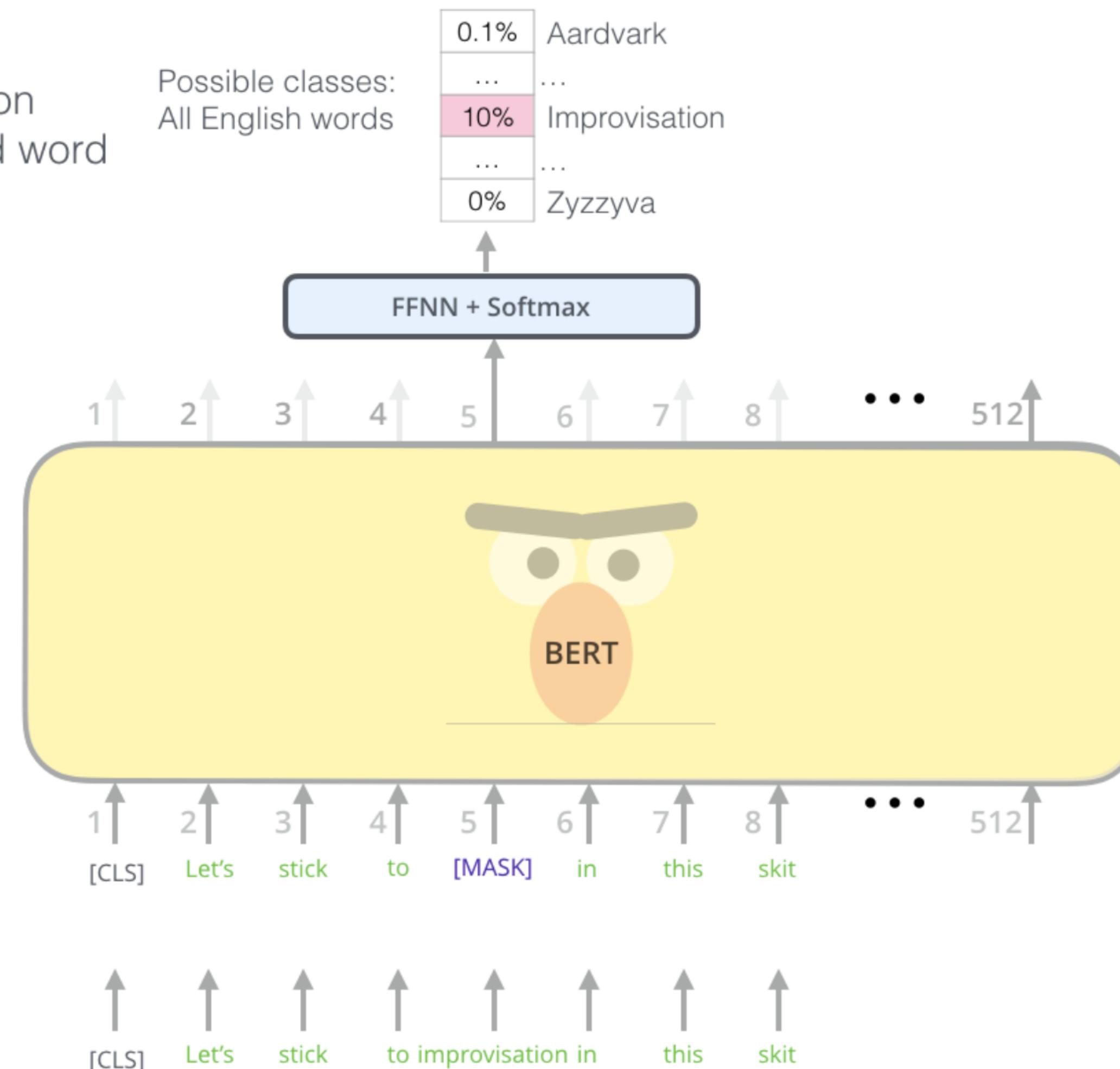
Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax

Randomly mask
15% of tokens

Input



BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

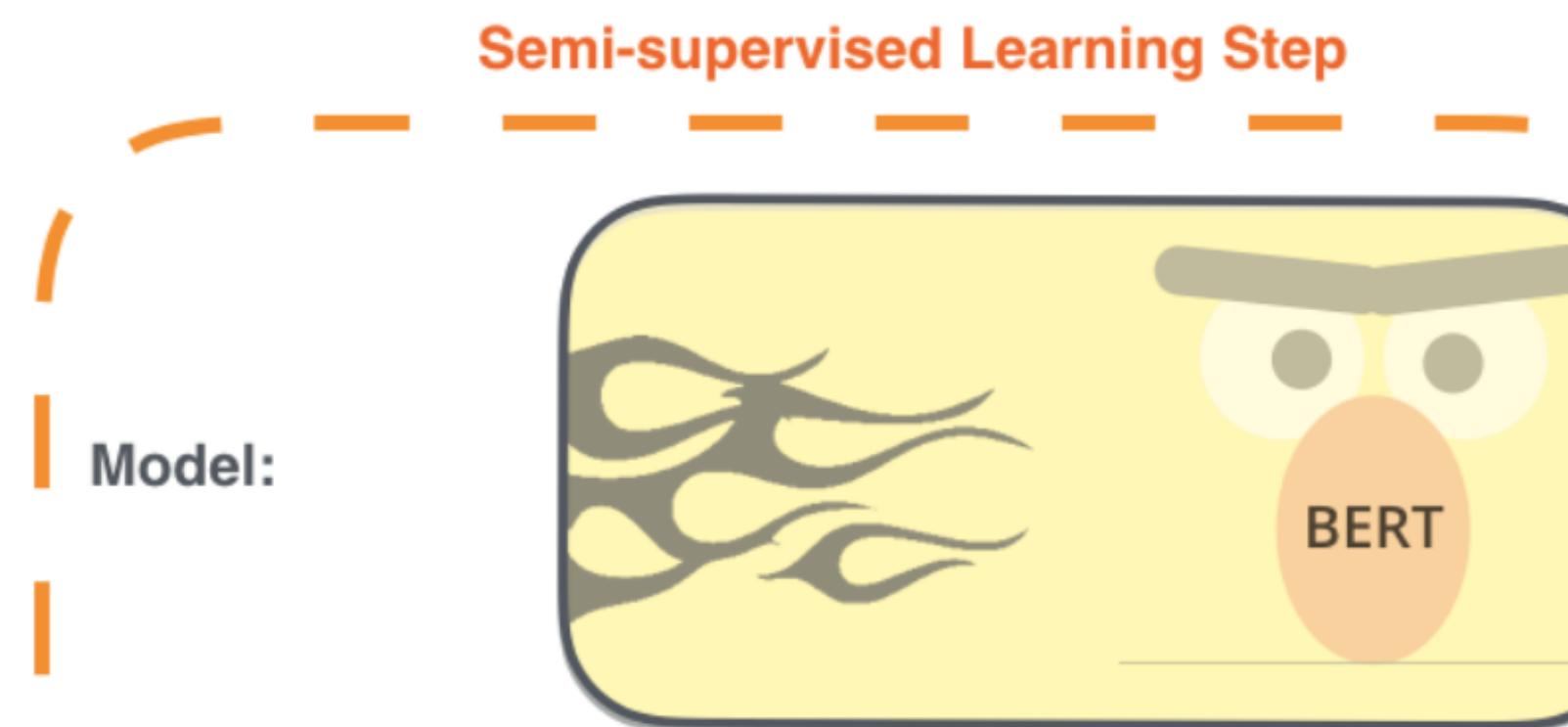
Jay Allamar – Illustrated BERT

Language Models

Bigger = Better ?

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



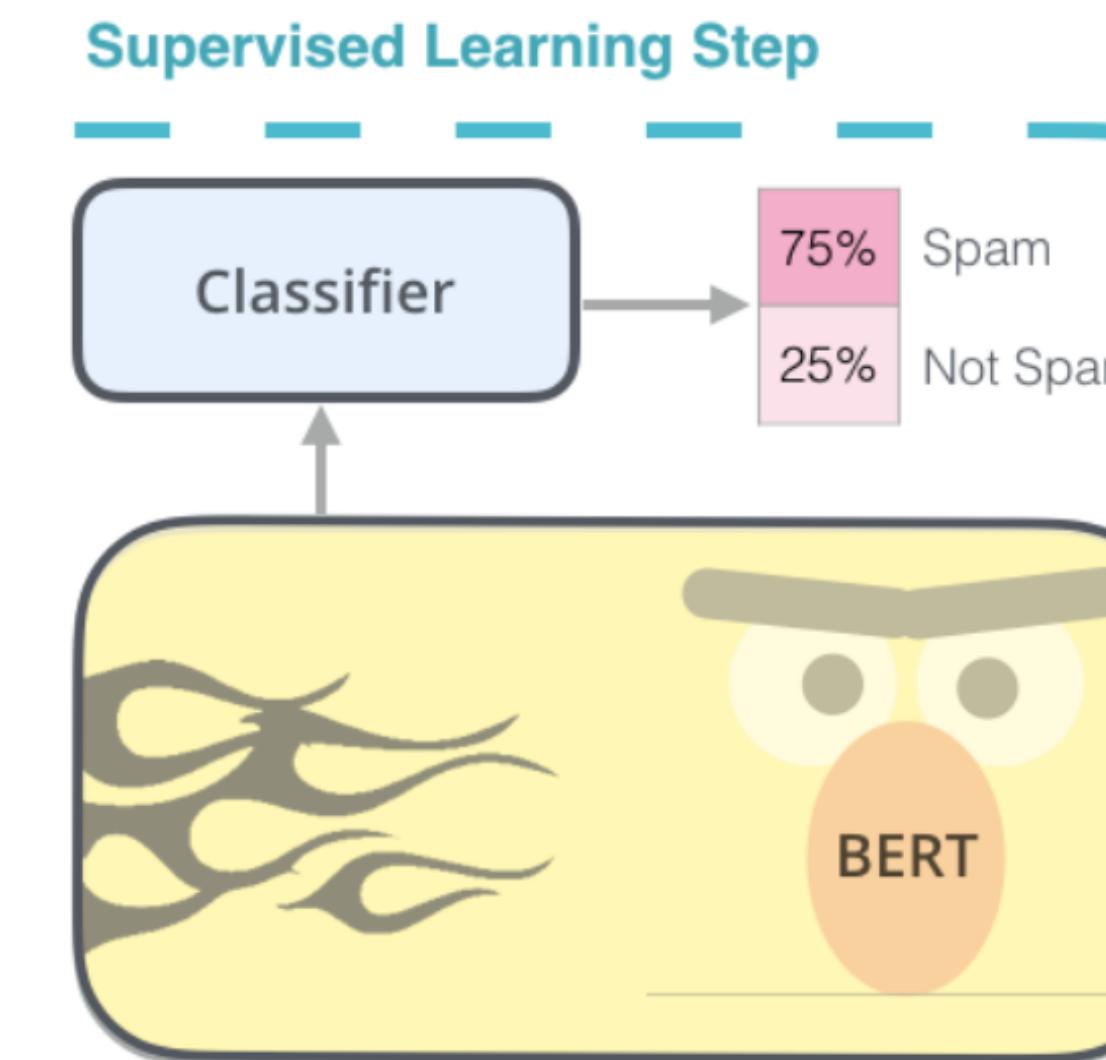
Model:



WIKIPEDIA
Die freie Enzyklopädie

Predict the masked word
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.



Model:
(pre-trained
in step #1)

Dataset:

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

Objective:

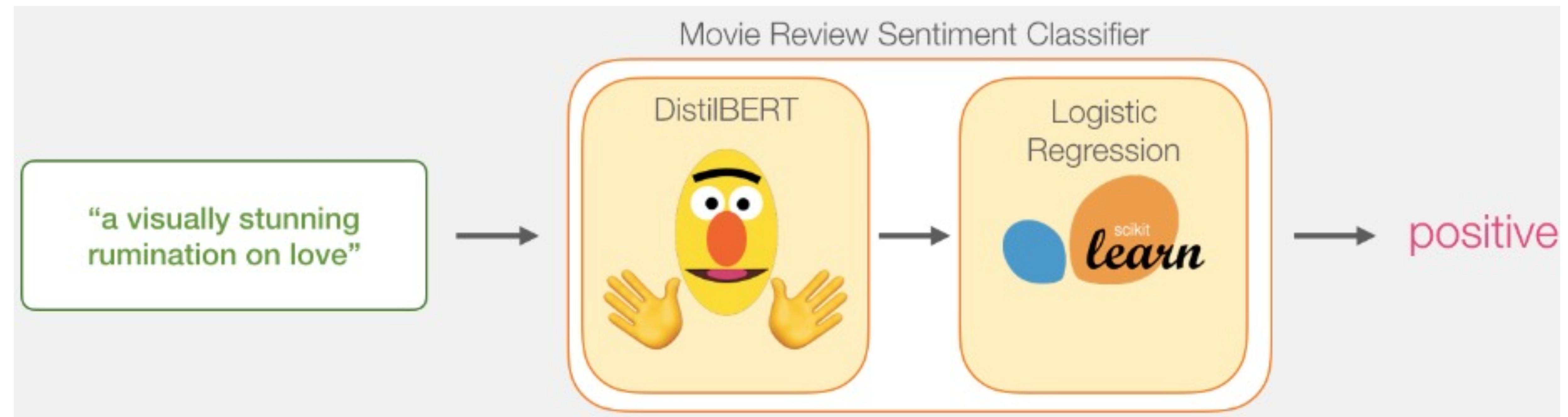
Transformers are good at Transfer Learning

Pre-Training improves downstream performance



Transformers are good at Transfer Learning

Pre-Training improves downstream performance

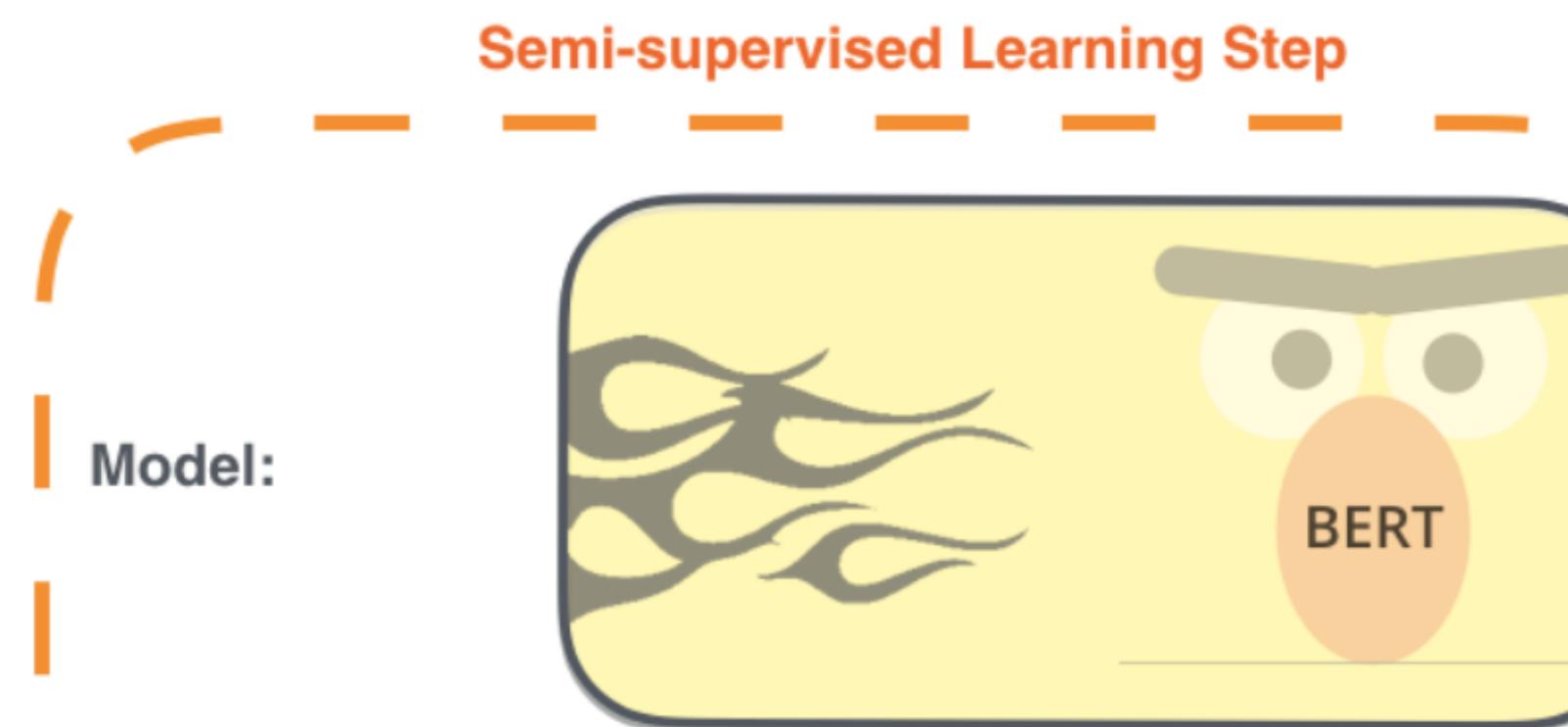


Language Models

Bigger = Better ?

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



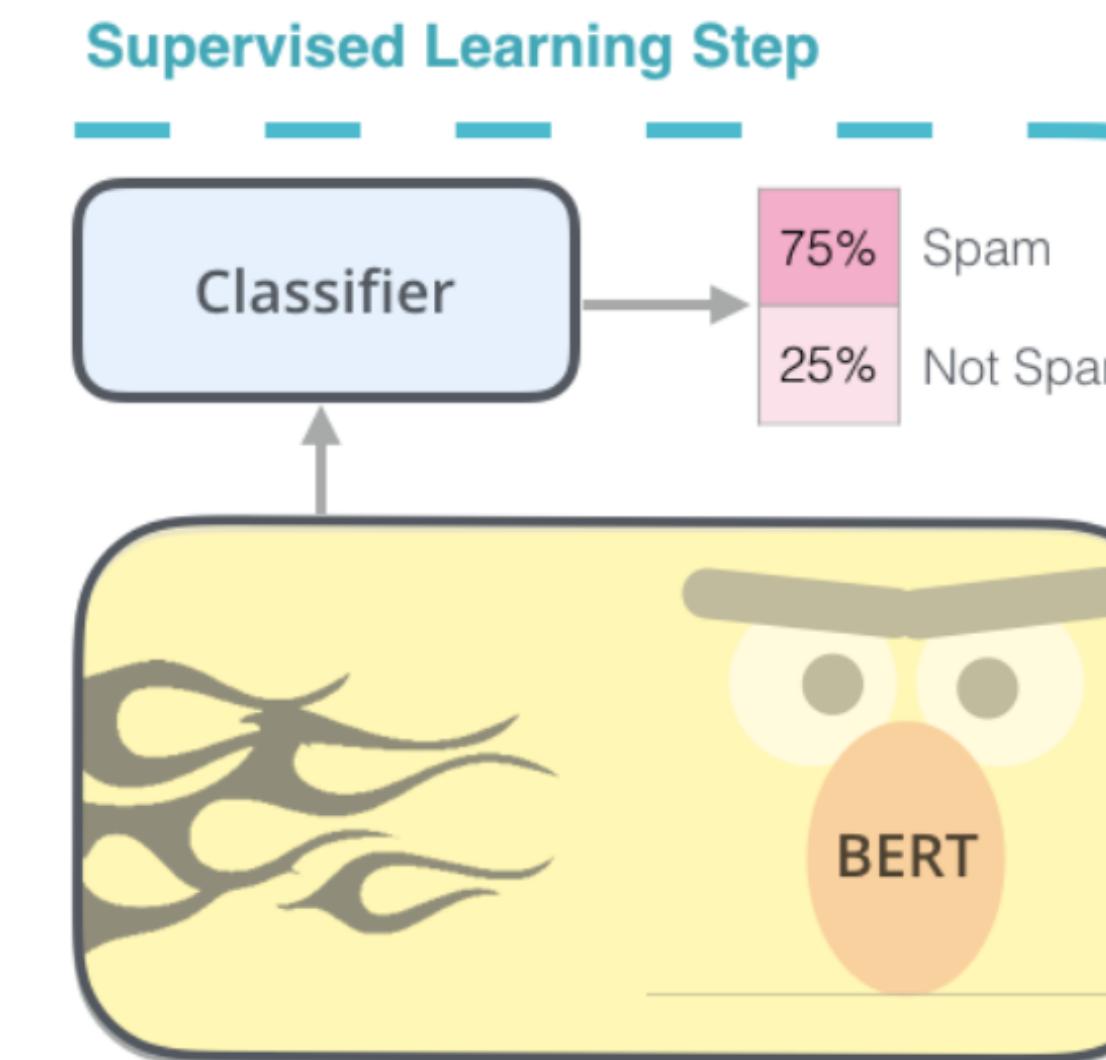
Model:



WIKIPEDIA
Die freie Enzyklopädie

Predict the masked word
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.



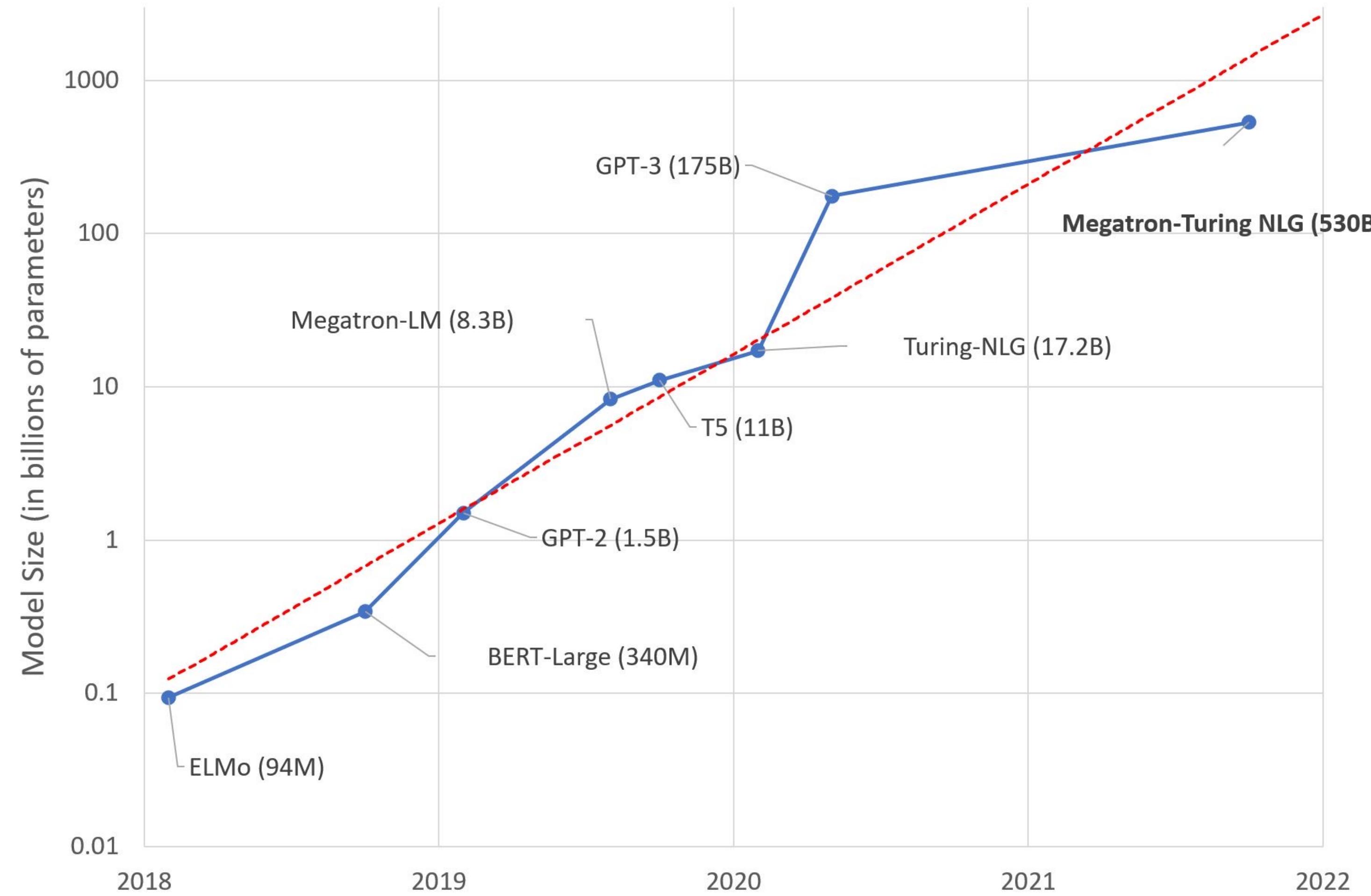
Model:
(pre-trained
in step #1)

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

Dataset:

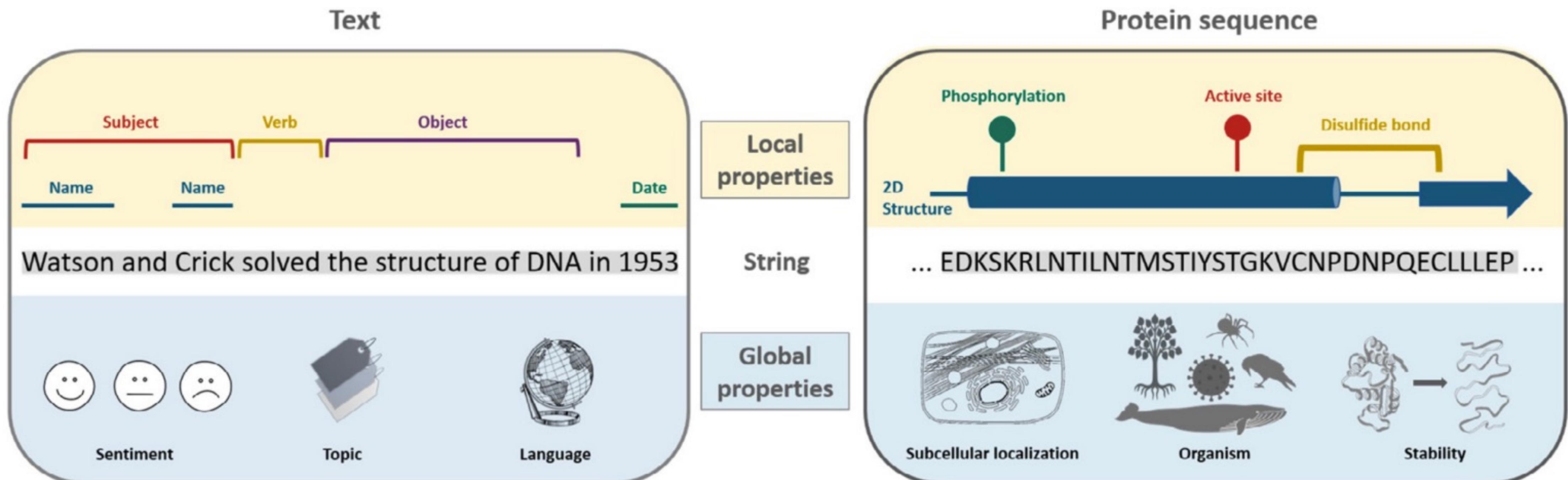
Language Models

Bigger = Better ?



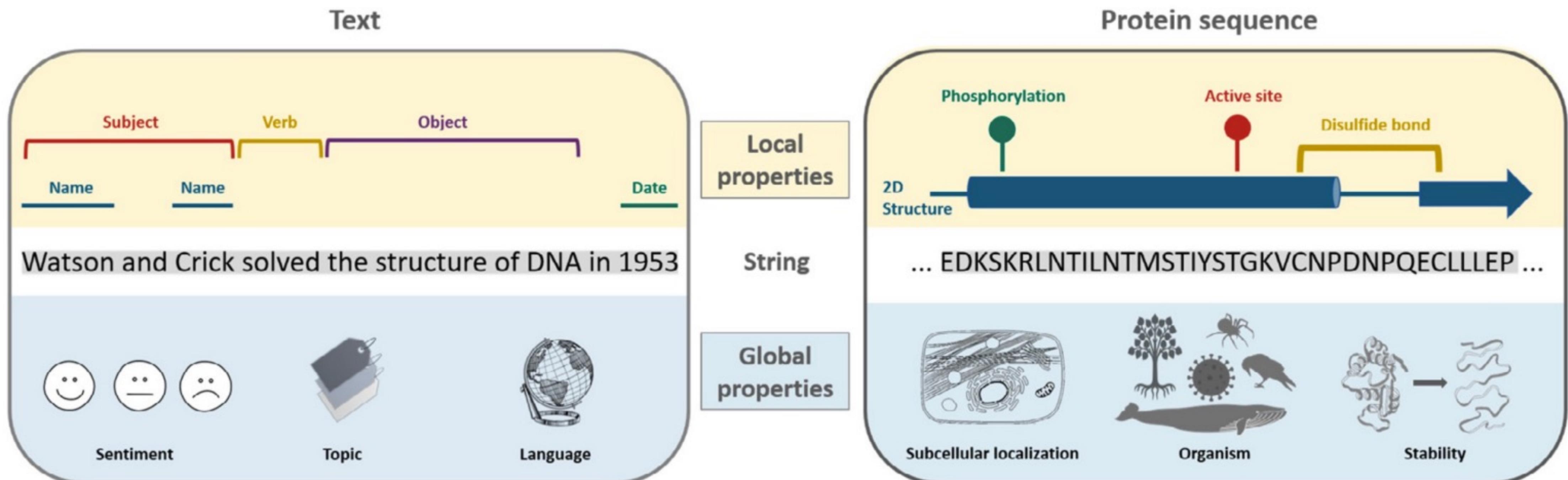
Proteins vs Sentences: The same?

Similar, but also important differences



Proteins vs Sentences: The same?

Similar, but also important differences



Can you *read* a protein?

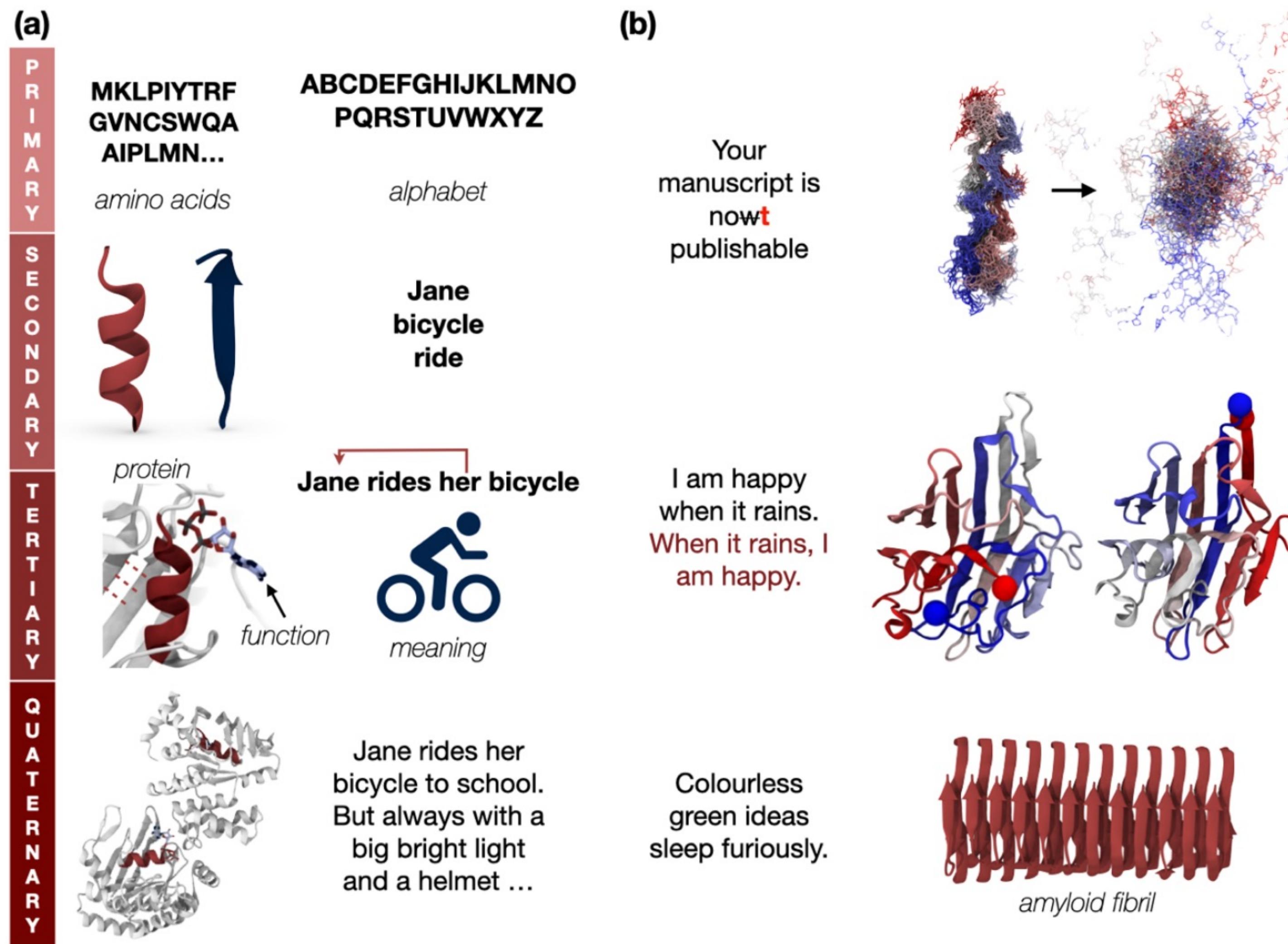
Past and Future tense?

Distant interactions in
protein structures

Bias in
Sequencing/Research

Proteins vs Sentences: The same?

Similar, but also important differences



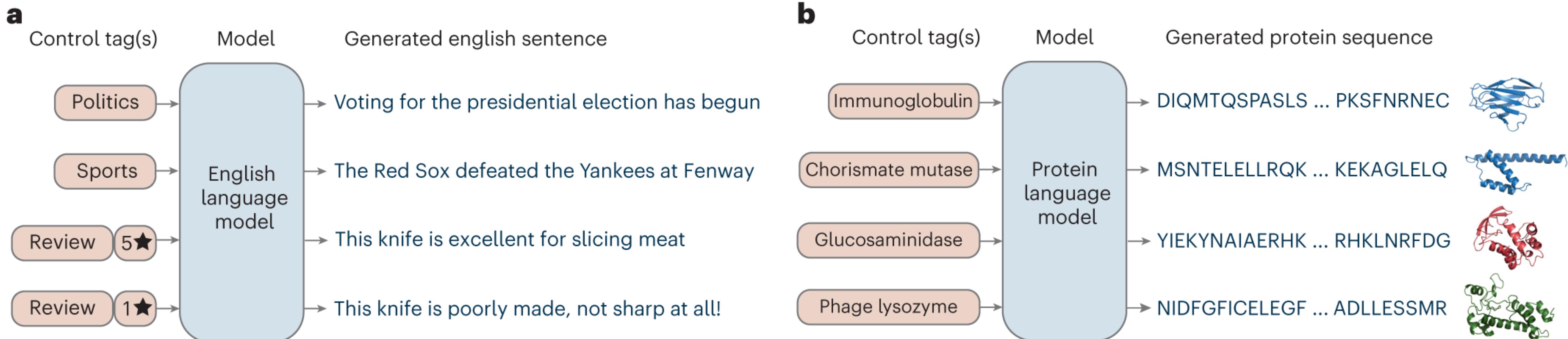
The linguistic hypothesis

Did evolution force proteins to develop a “language”?

- The space of naturally occurring proteins occupies a learnable manifold.
- This manifold emerges from evolutionary pressures that heavily encourage the reuse of components at many scales: from short motifs of secondary structure, to entire globular domains.

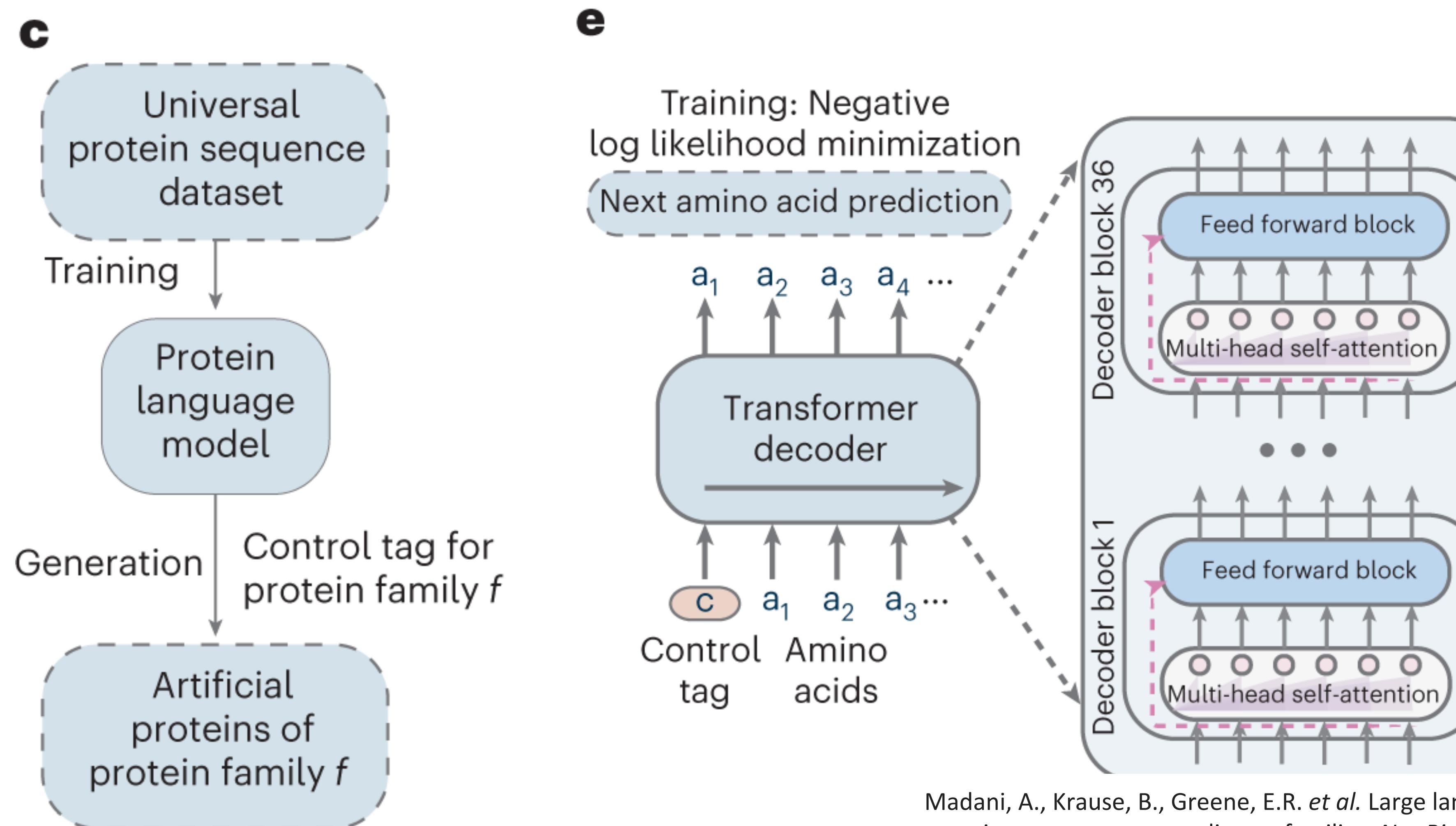
Protein Language Models

Train a model to understand the language of proteins



Protein Language Models

Train a model to understand the language of proteins





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.