

LECTURE 4: ATTENTION IS ALL YOU NEED

CSC5991: Introduction to LLMs

OUTLINE



Recap

Seq2Seq
CNN



Image Captioning with Visual Attention



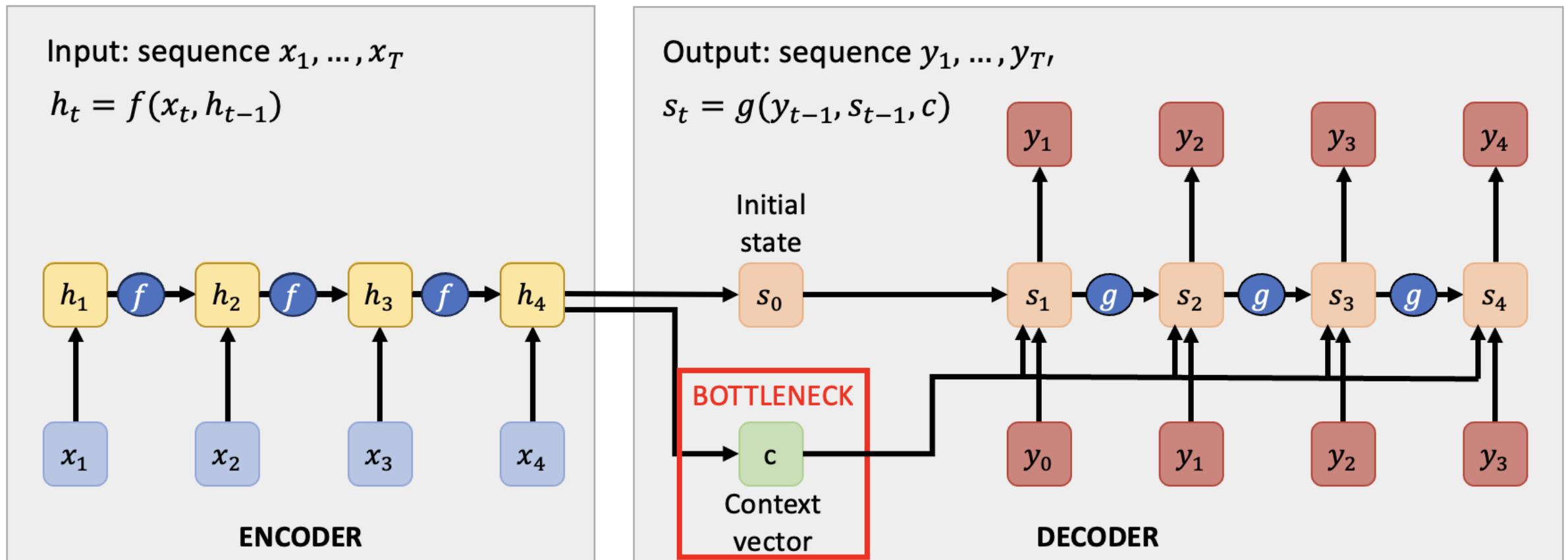
Attention Is All You Need:

Motivation
Attention decoupled from the RNN
Transformer architecture

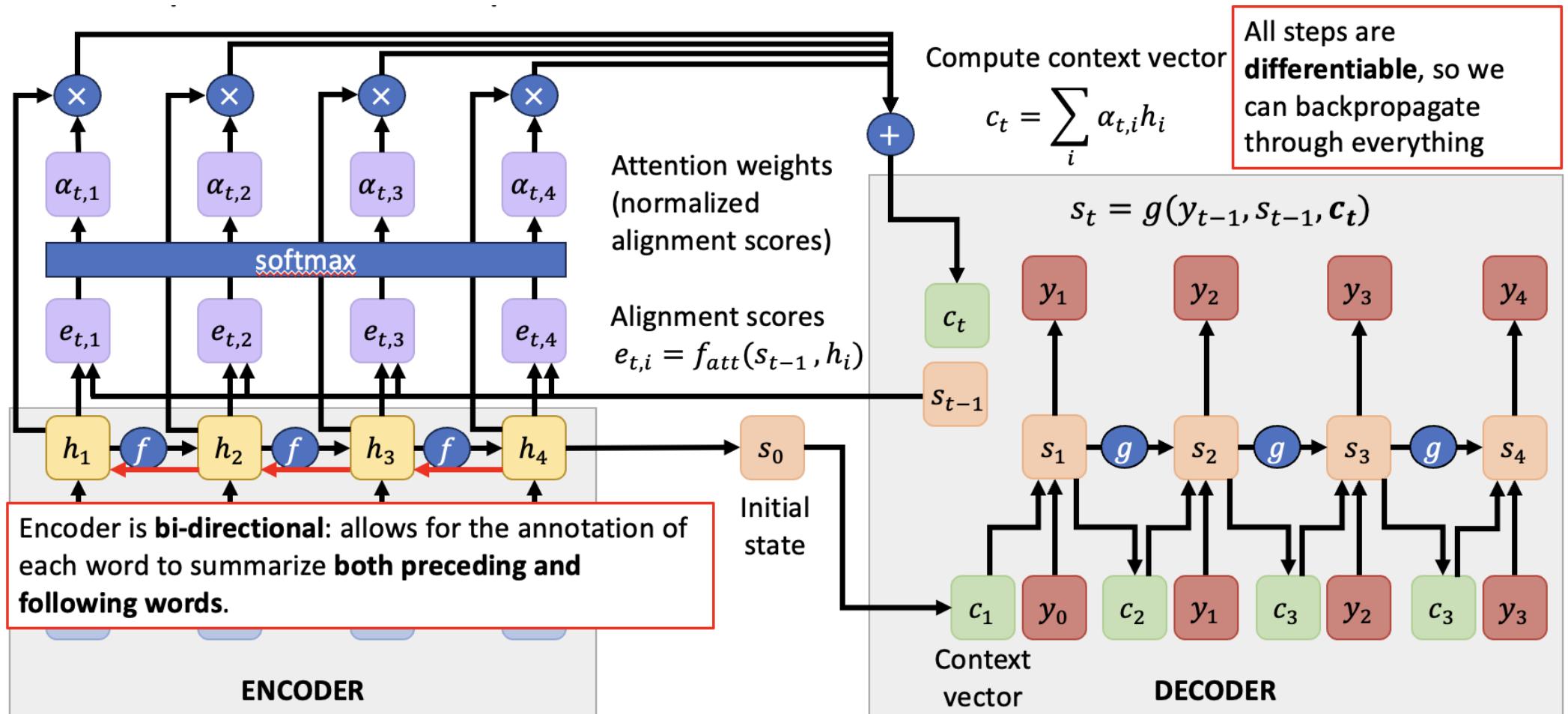
ATTENTION

- What we now call “attention” in DL
 - The idea of paying “**attention**” to the most relevant parts of I/P
 - Ideally, we’d like to **learn** this!
- There are many forms of attention mechanisms
 - **Additive**
 - **Dot-product**
- We have names to distinguish attention based on what is attended to
 - **Self-attention** (intra attention)
 - **Cross-attention** (encoder-decoder attention/inter attention)

BOTTLENECK

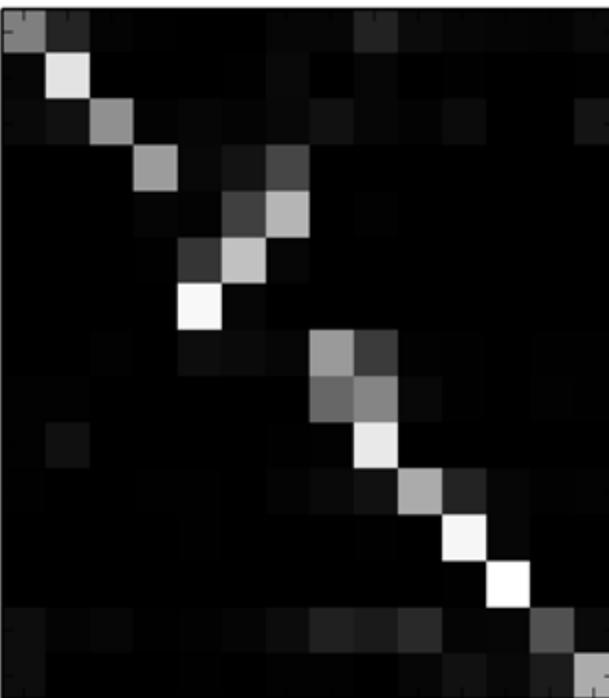


SEQ2SEQ: RNNs + ATTENTION



ATTENTION MAP

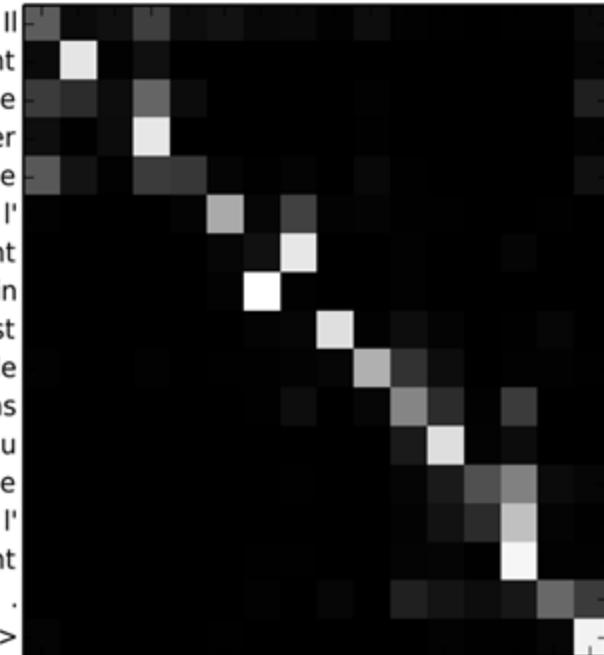
The
agreement
on
the
European
Economic
Area
was
signed
in
August
1992
. <end>



This attention map visualizes the importance of words in the English sentence. The words are arranged in a 12x12 grid. The most prominent words are 'The' (top-left), 'agreement' (second row), 'on' (third row), 'European' (fourth row), 'Economic' (fifth row), 'Area' (sixth row), 'was' (seventh row), 'signed' (eighth row), 'in' (ninth row), 'August' (tenth row), '1992' (eleventh row), and the punctuation ' .' and '<end>' (bottom). The background is black, and the grid cells are shaded in white, light gray, medium gray, and dark gray.

L'
accord
sur
la
zone
économique
européenne
a
été
signé
en
août
1992
. <end>

It
should
be
noted
that
the
marine
environment
is
the
least
known
of
environments
. <end>



This attention map visualizes the importance of words in the French sentence. The words are arranged in a 12x12 grid. The most prominent words are 'L'' (top-left), 'accord' (second row), 'sur' (third row), 'la' (fourth row), 'zone' (fifth row), 'économique' (sixth row), 'européenne' (seventh row), 'a' (eighth row), 'été' (ninth row), 'signé' (tenth row), 'en' (eleventh row), and the punctuation ' .' and '<end>' (bottom). The background is black, and the grid cells are shaded in white, light gray, medium gray, and dark gray.

IMAGE CAPTIONING

- The task of the task of describing the content of an image in words.

What is Image Captioning?



"Man in black shirt is playing guitar."



"Construction worker in orange safety vest is working on road."



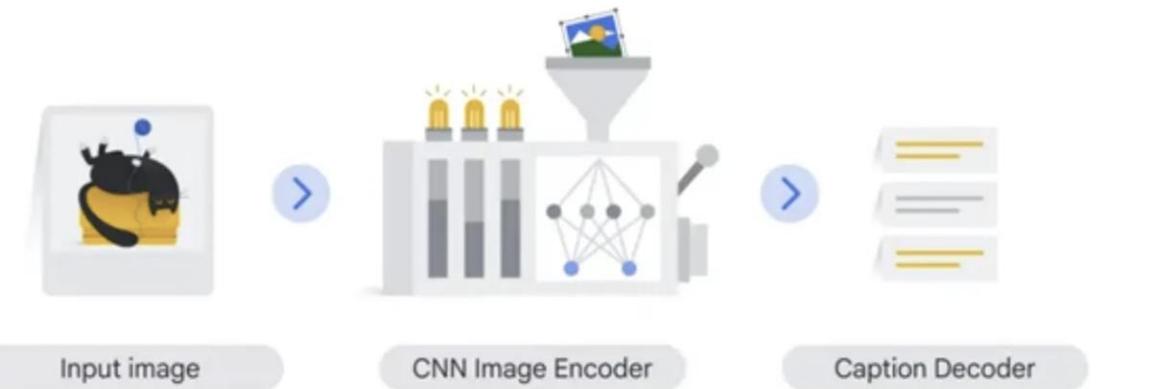
"Two young girls are playing with lego toy."

ARCHITECTURE

Goal: Build and train a model that can generate text captions based on images.

- **Encoder-decoder model**
- **Encoder:** Images are passed to the encoder
 - Extracts information from the images and creates feature vectors.
 - The encoder can use any image backbone like CNNs, ResNets, ViTs.
- **Decoder:** Builds captions by generating words.

Encoder-Decoder based Image Captioning Architecture



Reference: <https://arxiv.org/abs/1502.03044>



ENCODER: CNN-BASED ARCHITECTURE

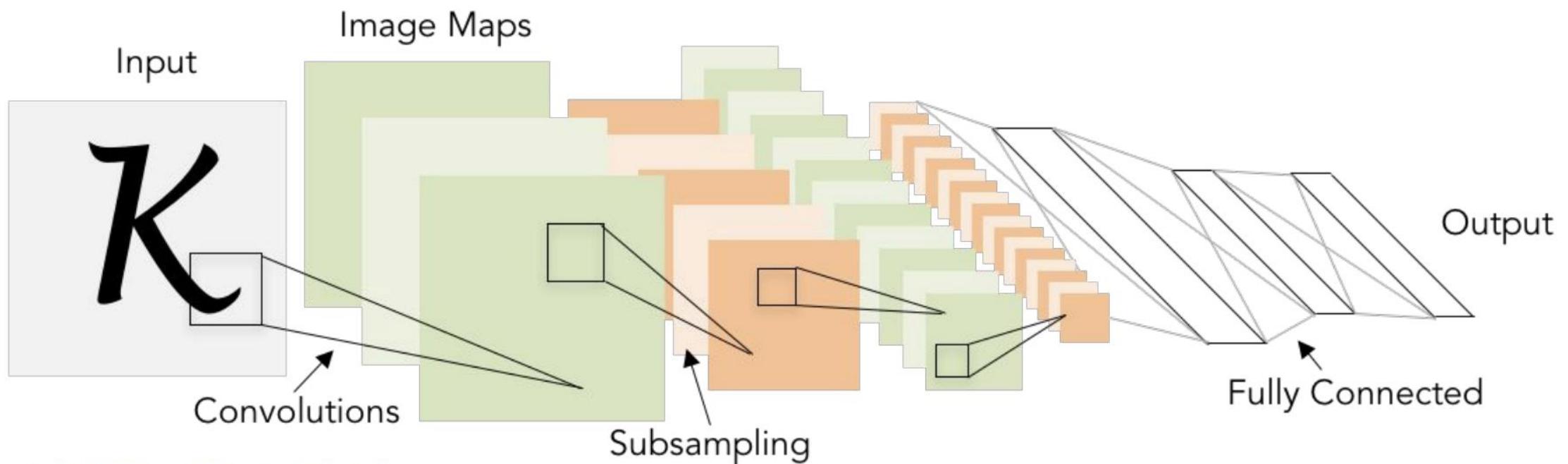
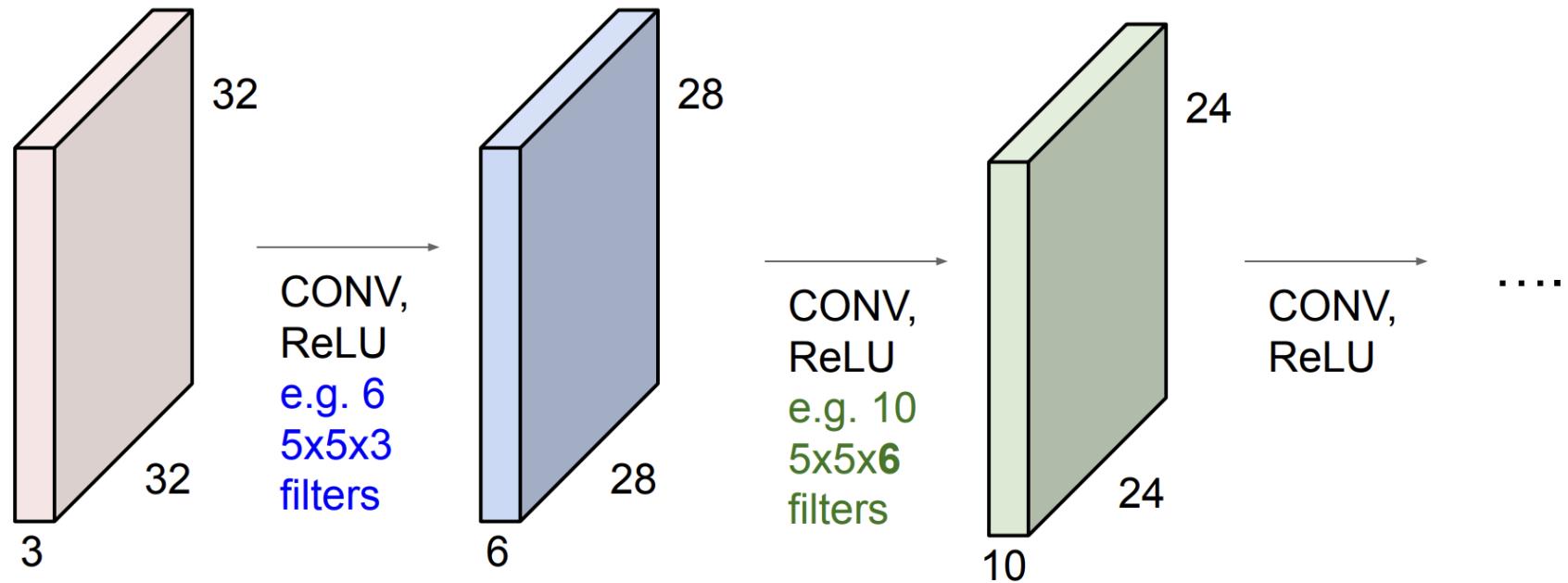


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

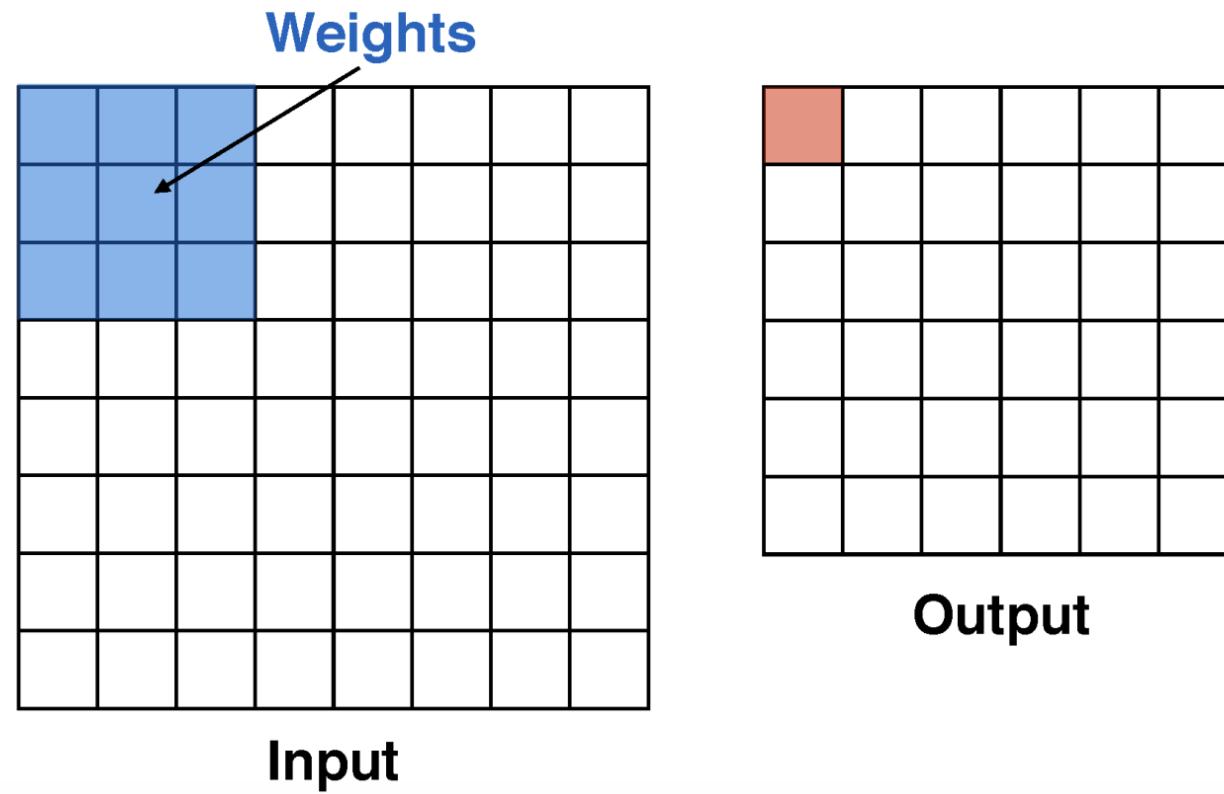
CNN: SEQUENCE OF CONVOLUTION LAYERS

Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



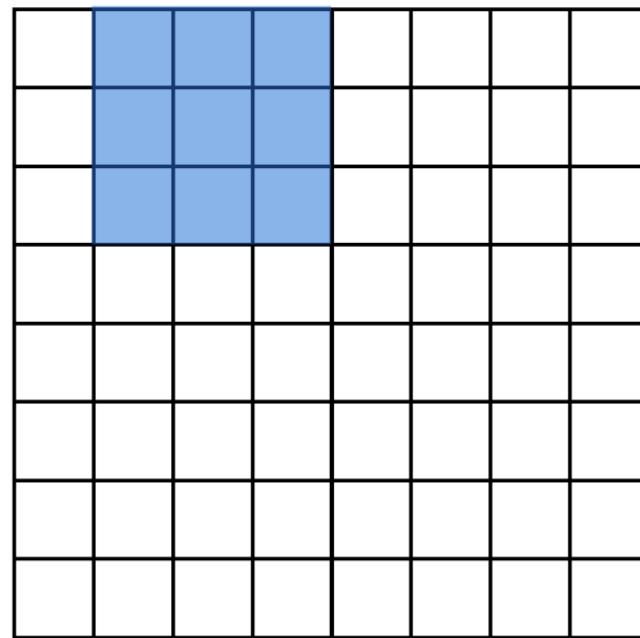
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

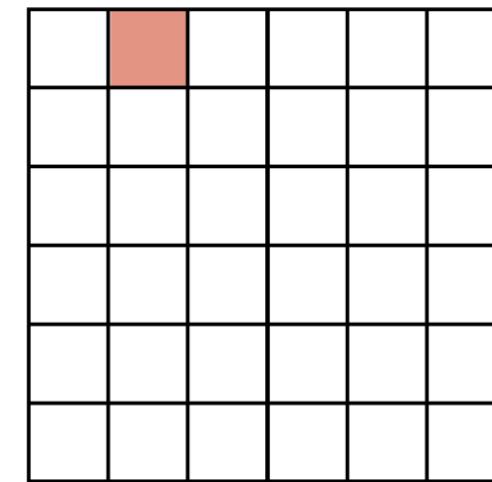


Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



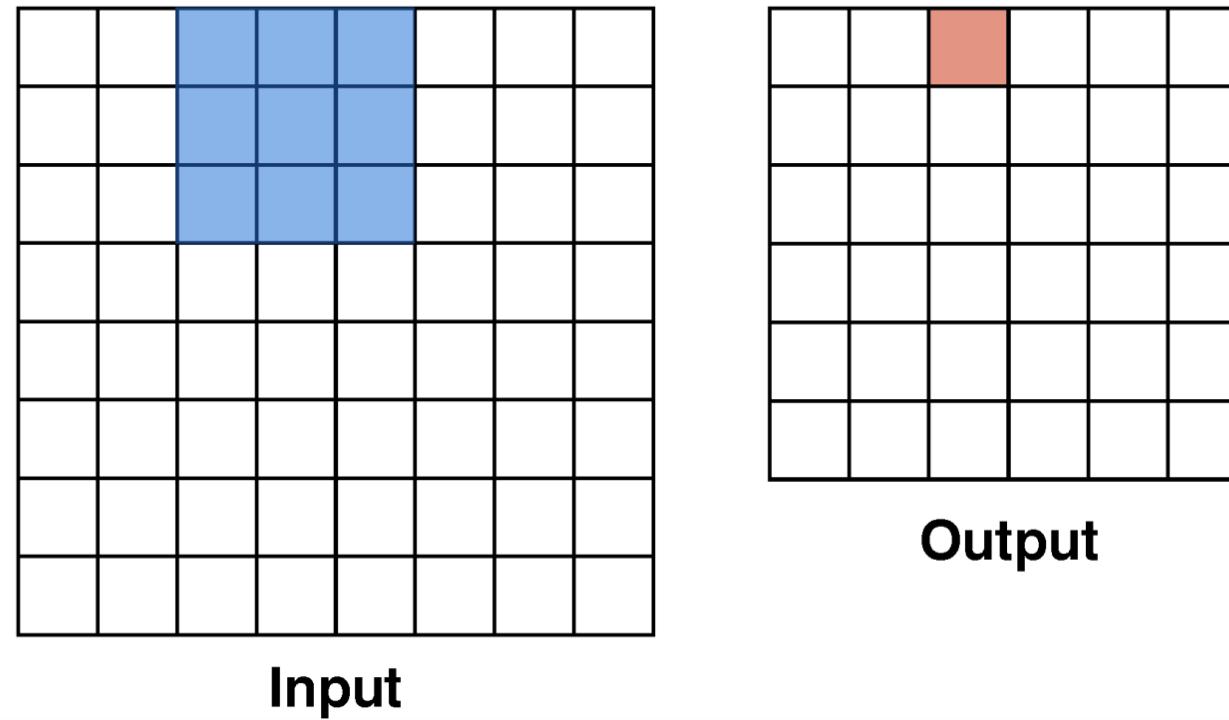
Input



Output

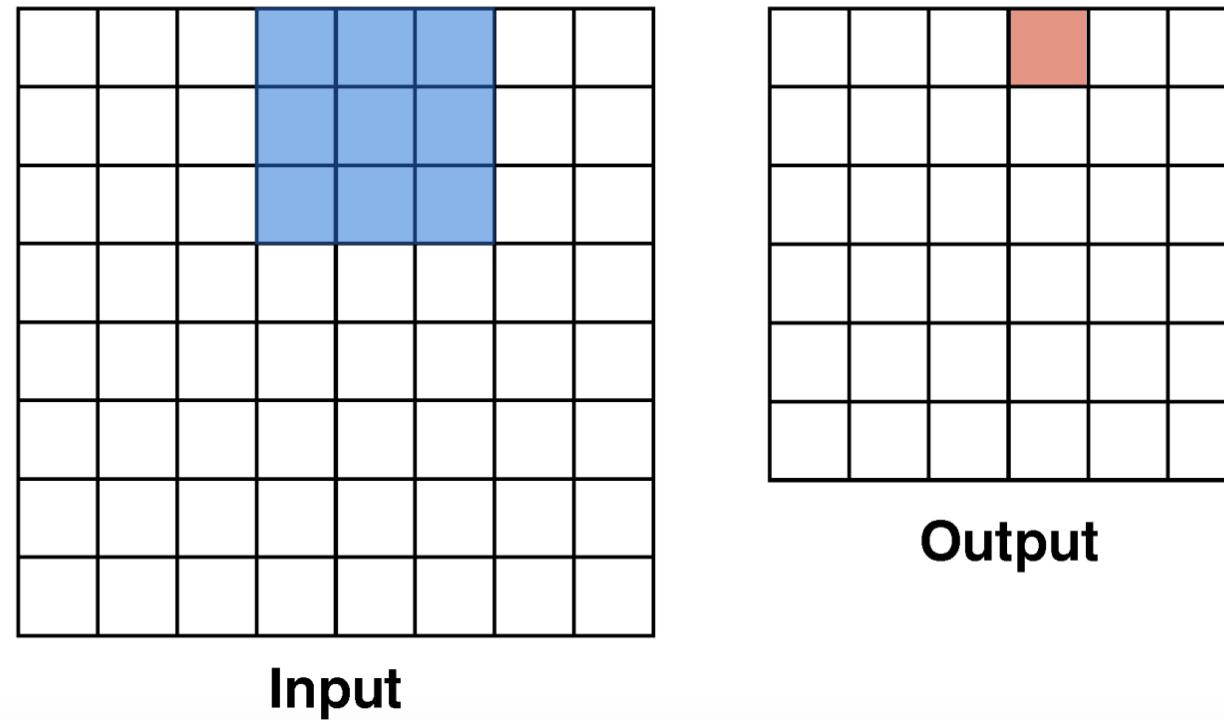
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



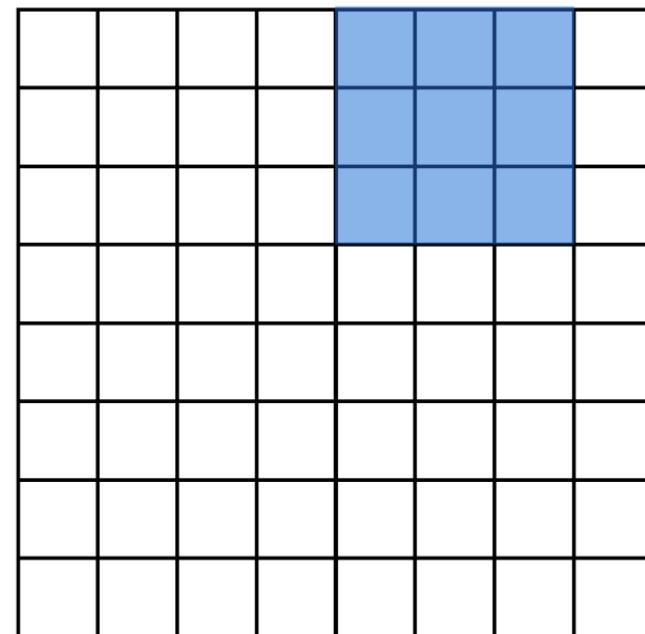
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

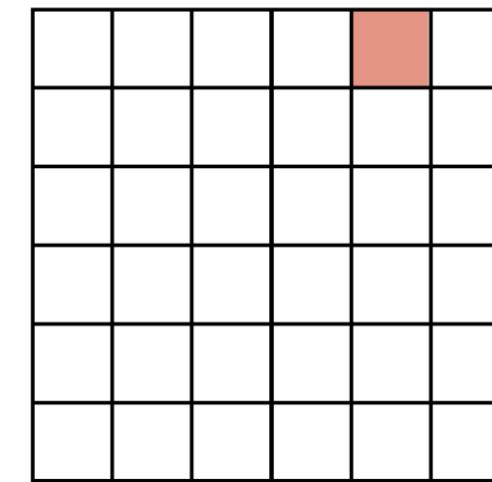


Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



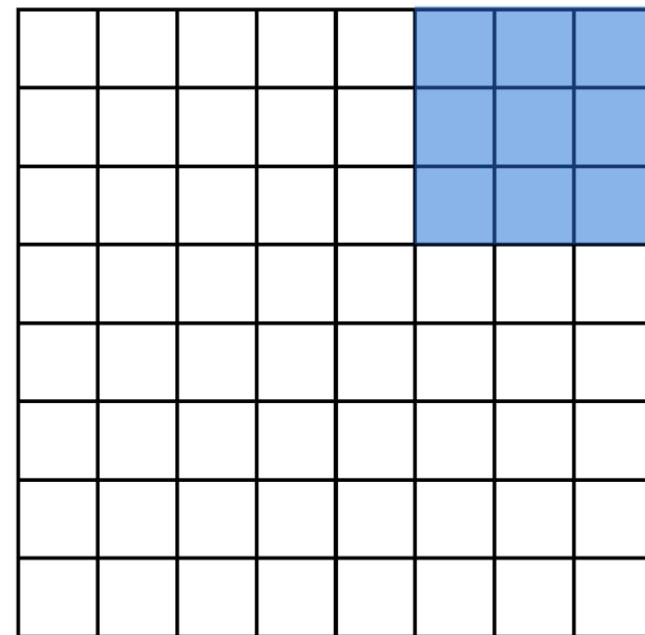
Input



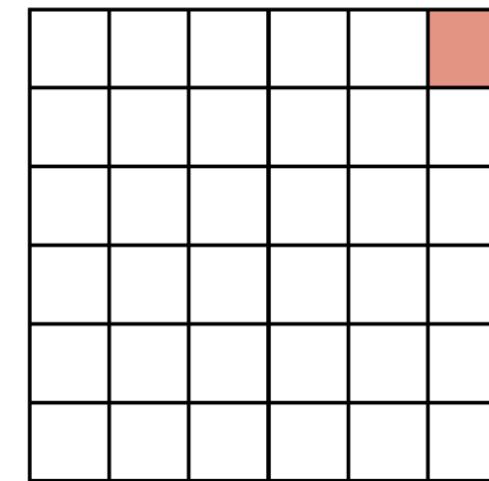
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



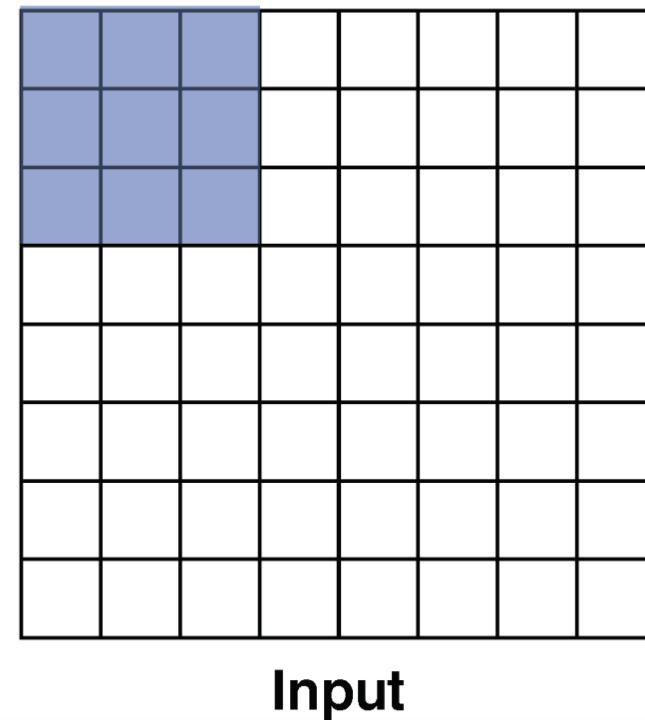
Input



Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



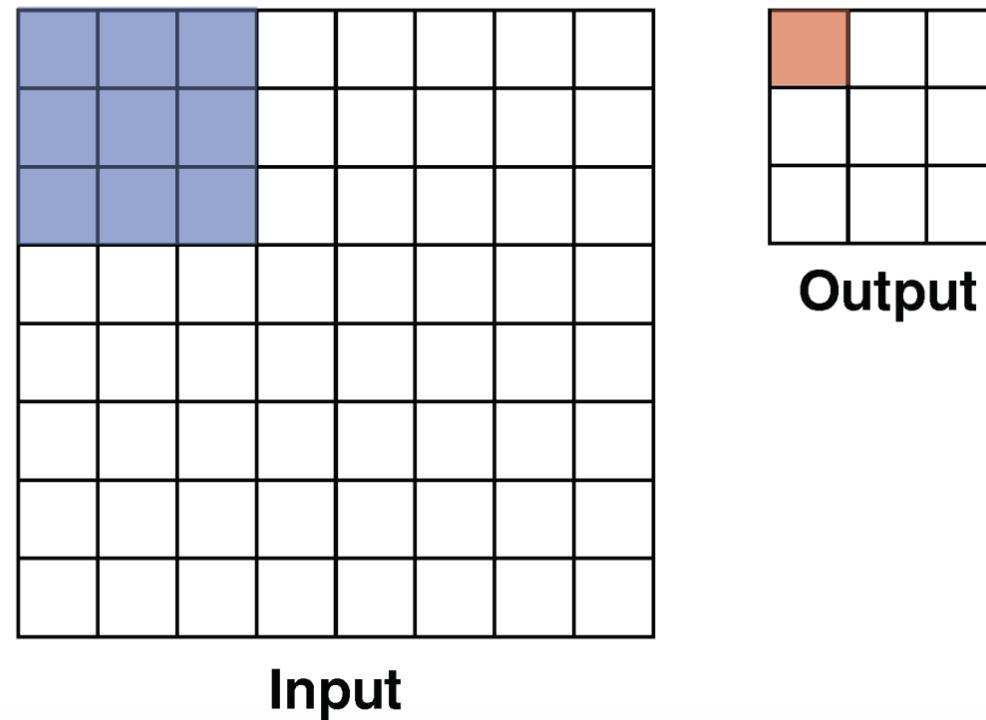
Recall that at each position we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

(channel, row, column)

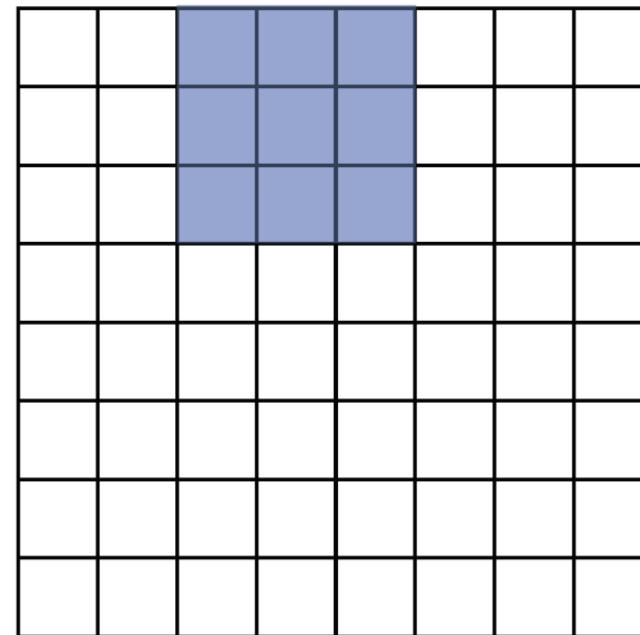
Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2

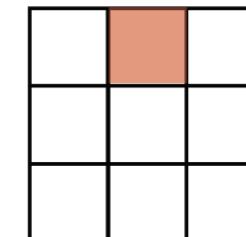


Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



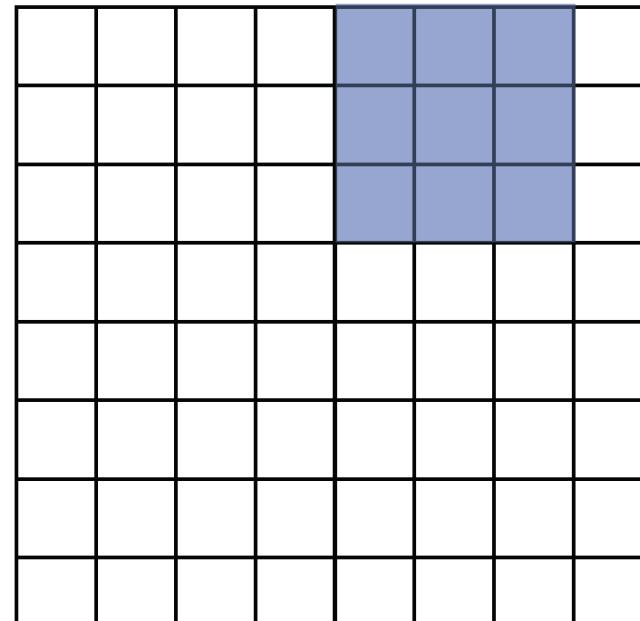
Input



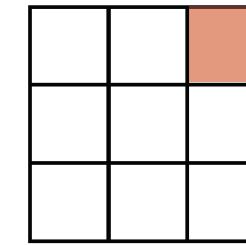
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



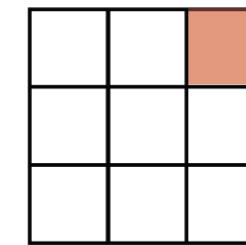
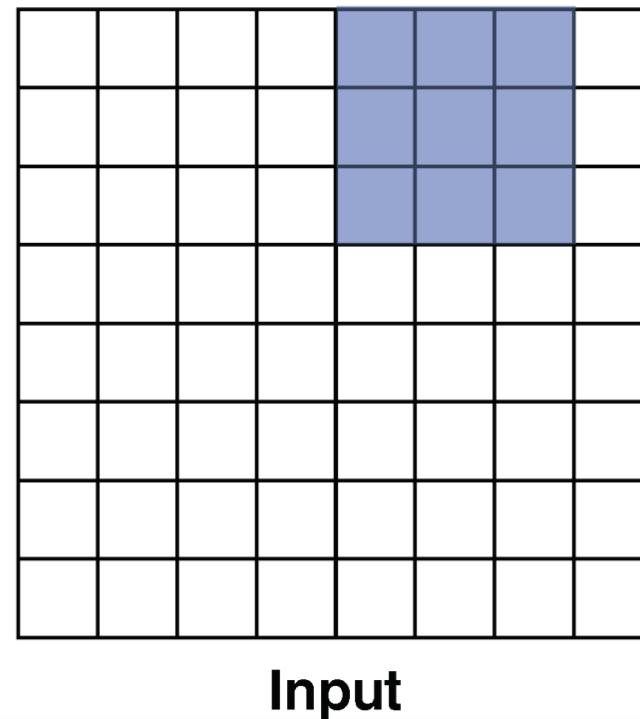
Input



Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0	0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0	0	0	0	0	0	0	0	0	0

Input

0			
	0		
		0	
			0

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

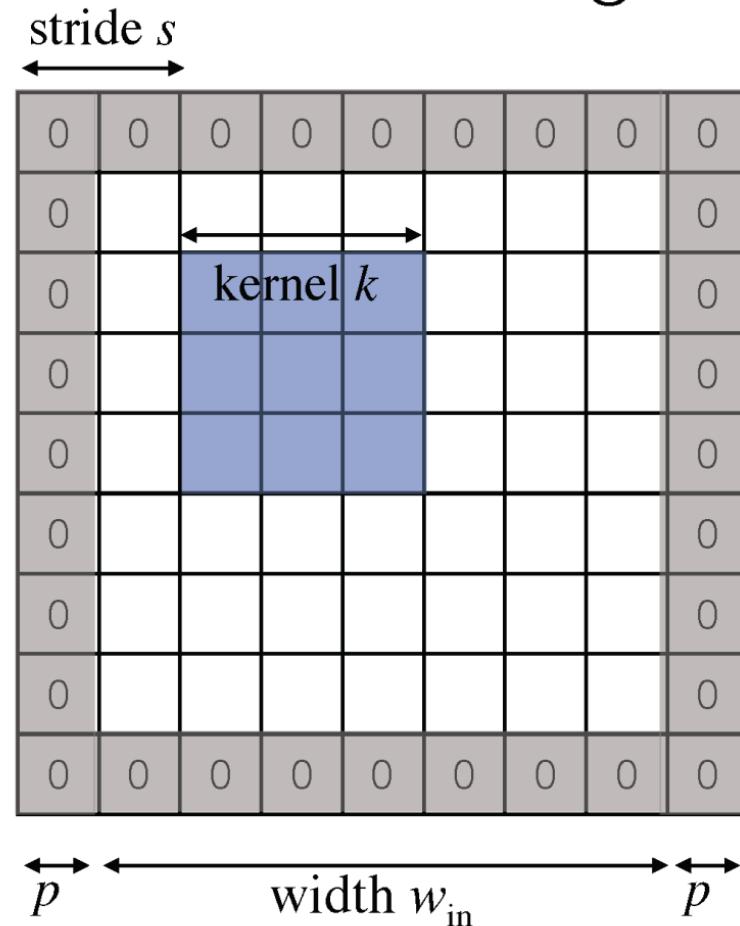
Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: How big is the output?

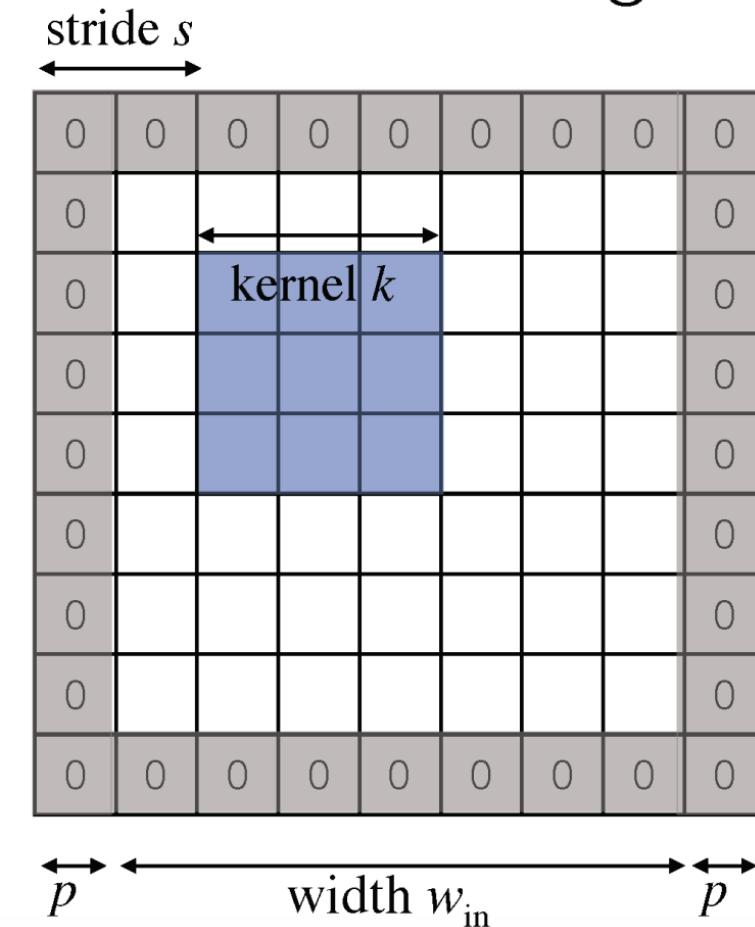


In general, the output has size:

$$w_{\text{out}} = \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1$$

Convolution:

How big is the output?



Example: $k=3$, $s=1$, $p=1$

$$\begin{aligned}
 w_{\text{out}} &= \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1 \\
 &= \left\lfloor \frac{w_{\text{in}} + 2 - 3}{1} \right\rfloor + 1 \\
 &= w_{\text{in}}
 \end{aligned}$$

VGGNet [Simonyan 2014]
uses filters of this shape

Pooling

For most ConvNets, **convolution** is often followed by **pooling**:

- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common
- Why might “avg” be a poor choice?

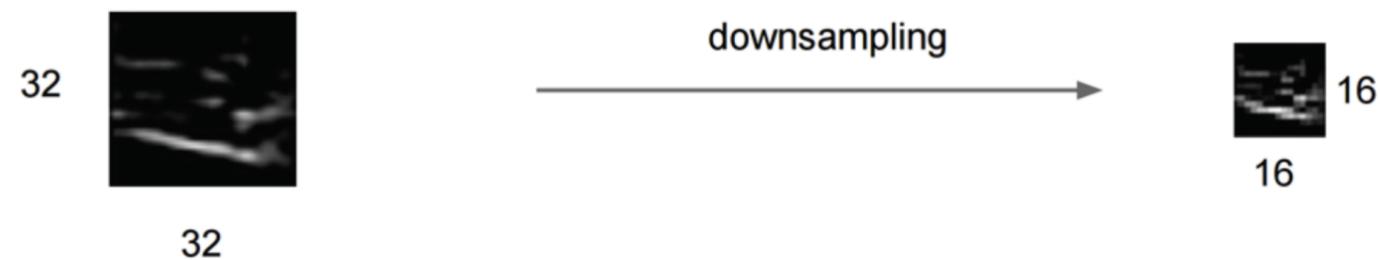
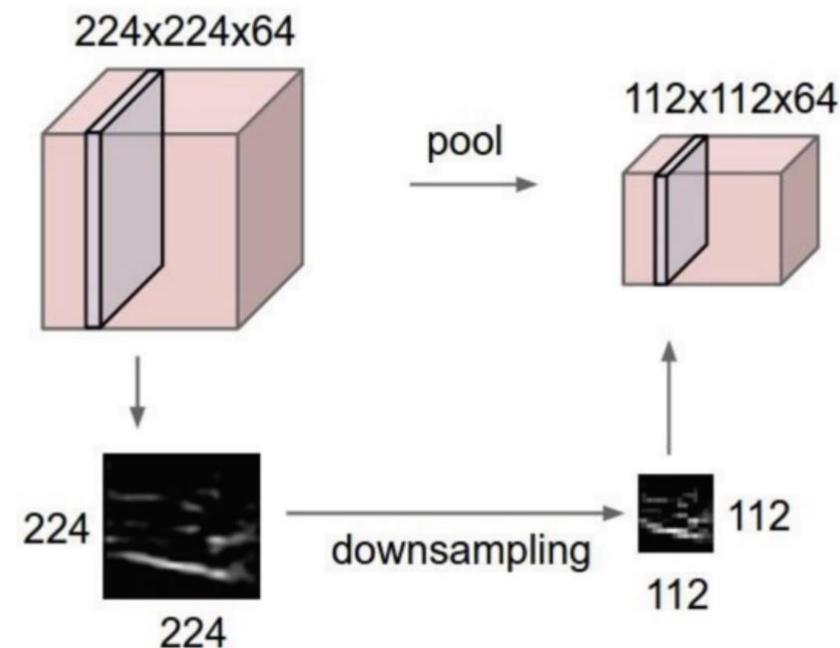


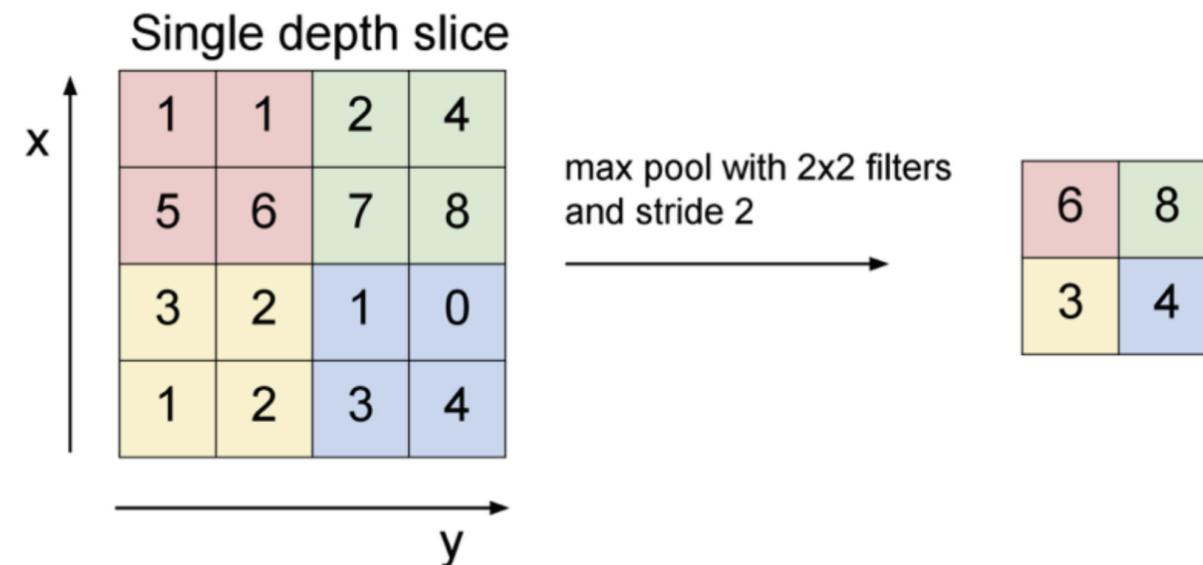
Figure: Andrej Karpathy

Pooling

- makes the representations smaller and more manageable
- operates over each activation map independently:



Max Pooling



What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

Figure: Andrej Karpathy

Example ConvNet

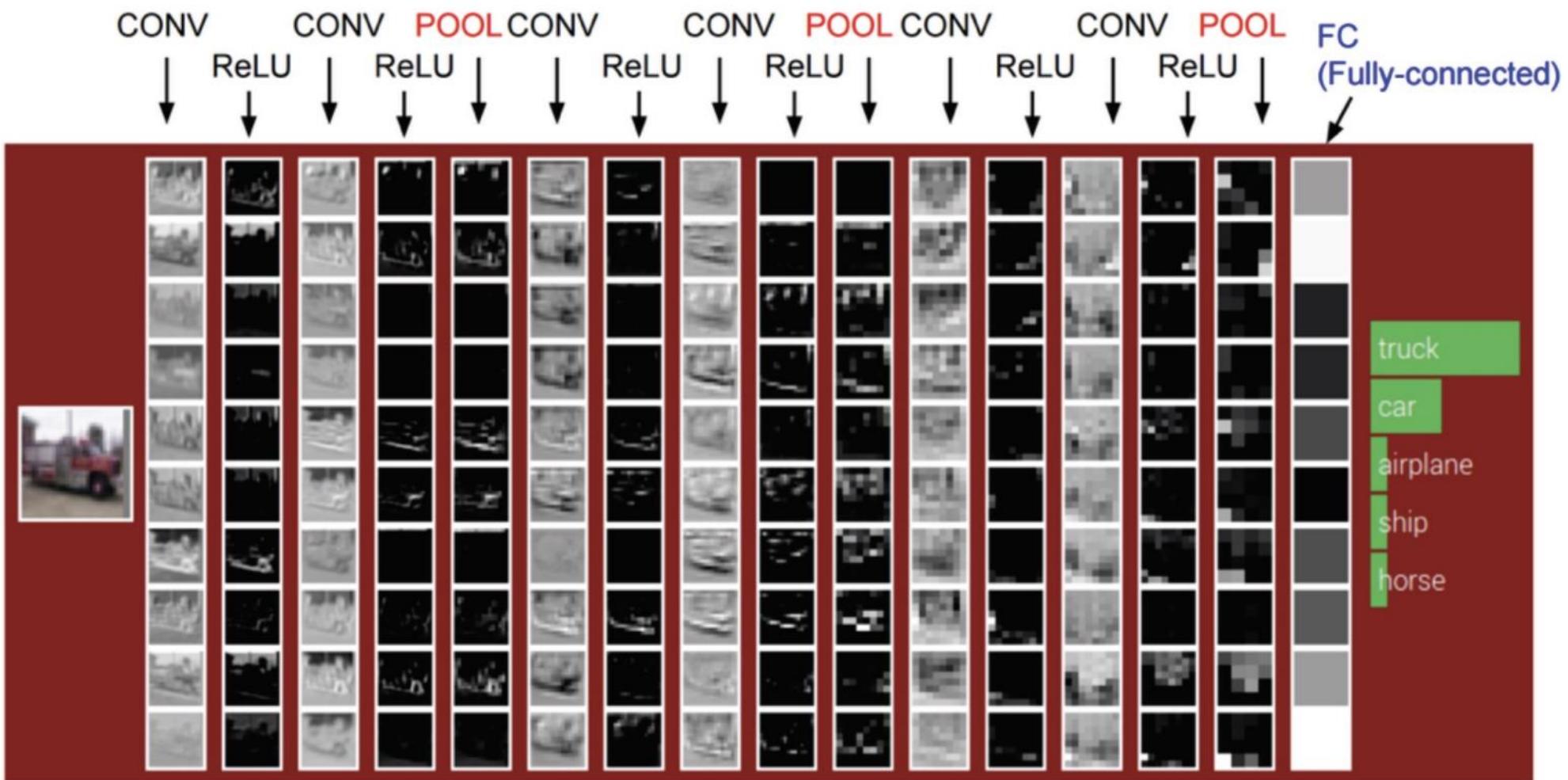


Figure: Andrej Karpathy

RESOURCES

- **CIFAR-10**
- <https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>
- **MNIST**
- <https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

IMAGE CAPTIONING WITH VISUAL ATTENTION

- Similarly to Seq2Seq we can use attention for image captioning (image → text)
- Builds directly on previous work

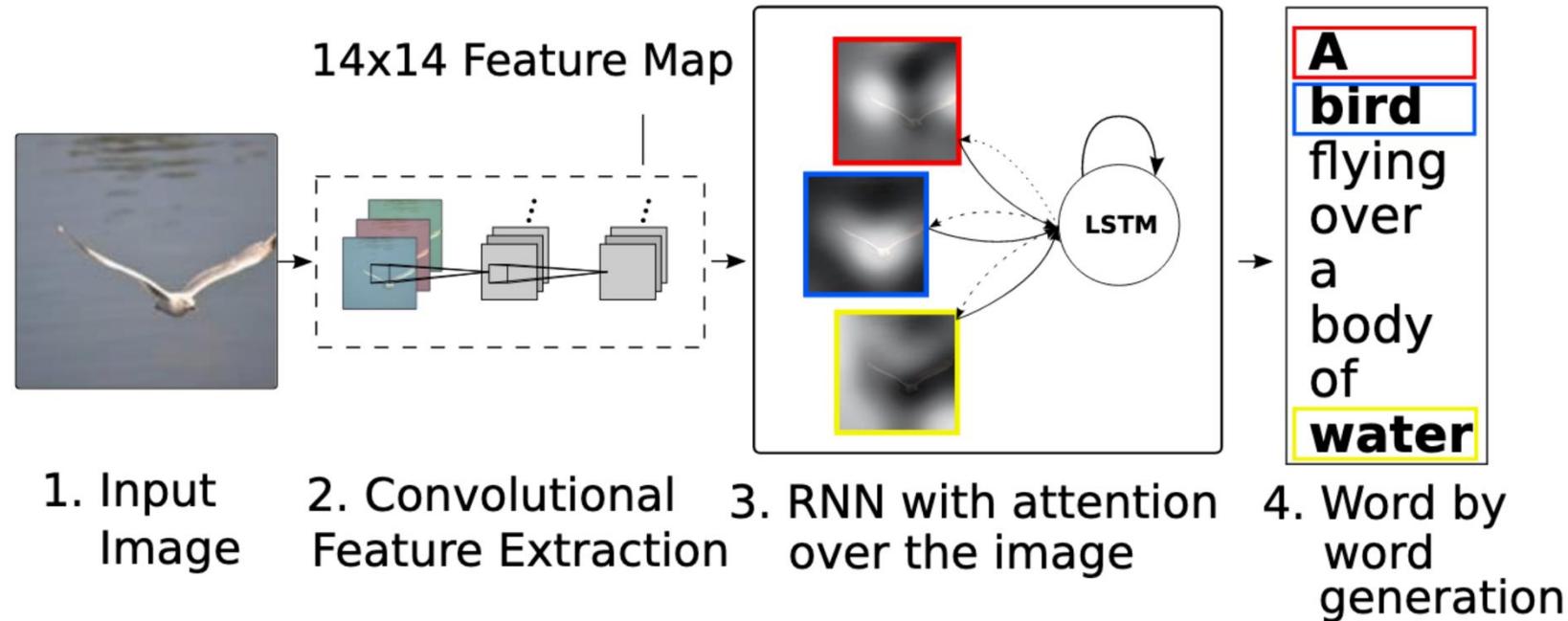


IMAGE CAPTIONING WITH VISUAL ATTENTION

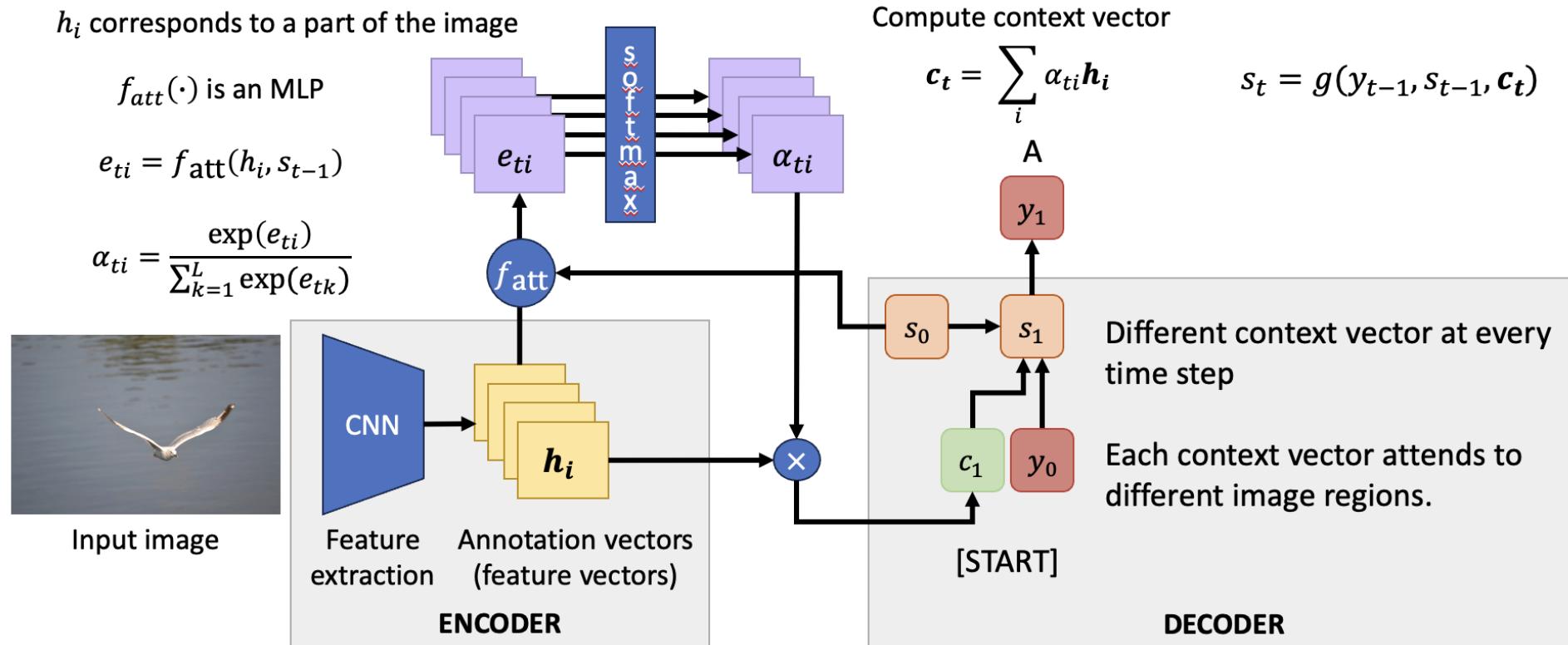


IMAGE CAPTIONING WITH VISUAL ATTENTION

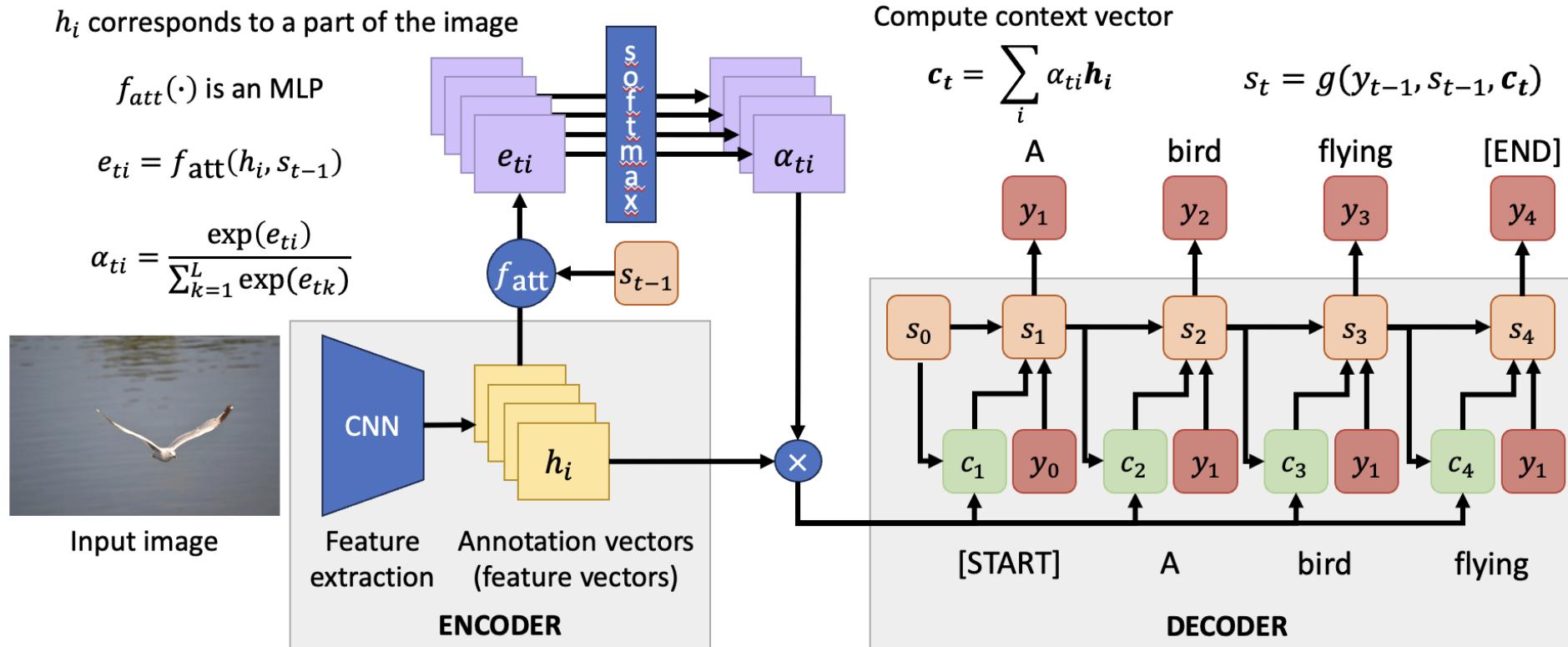


IMAGE CAPTIONING WITH VISUAL ATTENTION

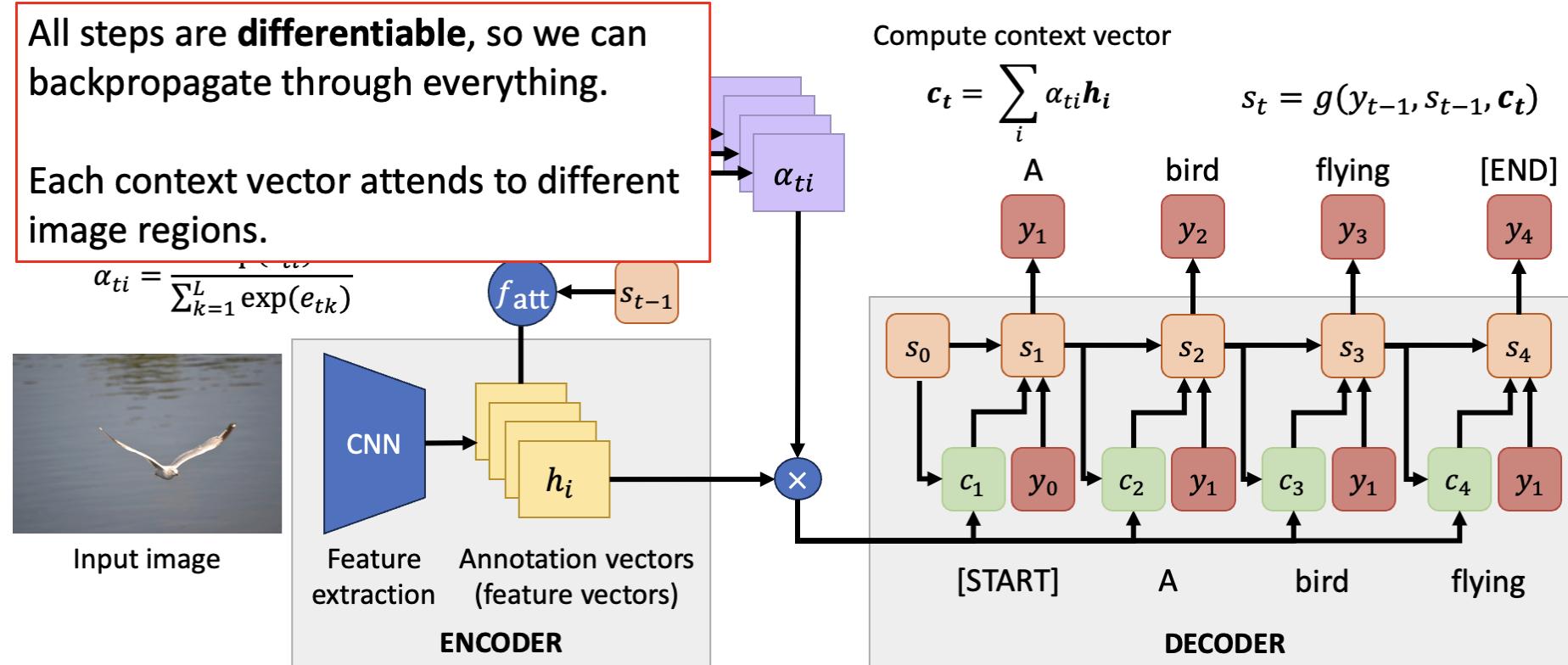


IMAGE CAPTIONING WITH VISUAL ATTENTION

- Visualization of the attention for each generated word
 - Gives insight to “where” and “what” the attention focused on when generating each word

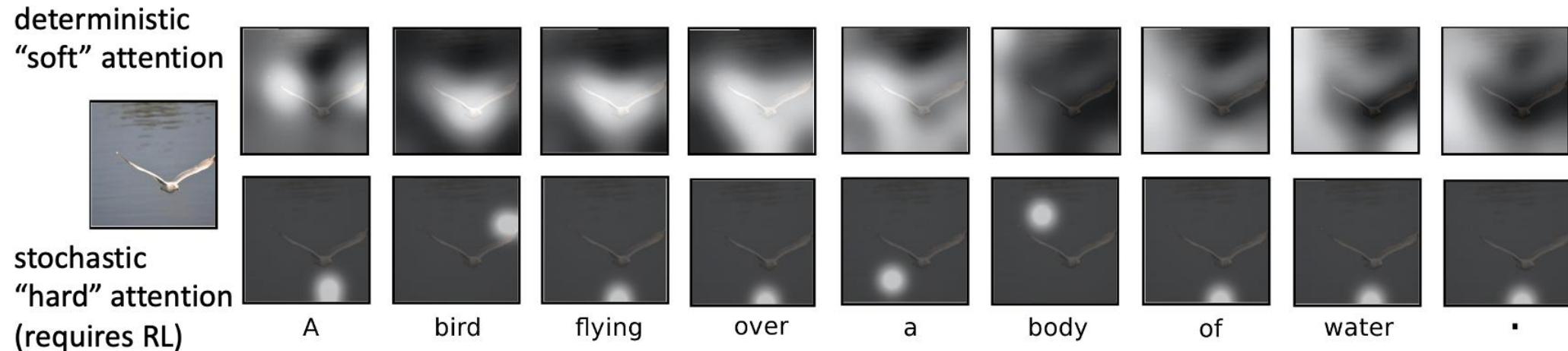


IMAGE CAPTIONING WITH VISUAL ATTENTION



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

RESOURCES

- <https://medium.com/@deepeshrishu09/automatic-image-captioning-with-pytorch-cf576c98d319>
- <https://www.kaggle.com/code/mdteach/image-captioning-with-attention-pytorch>
- Tensorflow: https://www.tensorflow.org/text/tutorials/image_captioning

ATTENTION IS ALL YOU NEED

- **Key Idea:**

- Decouple attention from RNNs
- Use self-attention to make this efficient

- **Contributions:**

- Multi-head attention
- Transformer architecture

- Highly impactful (as we'll touch on later)

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

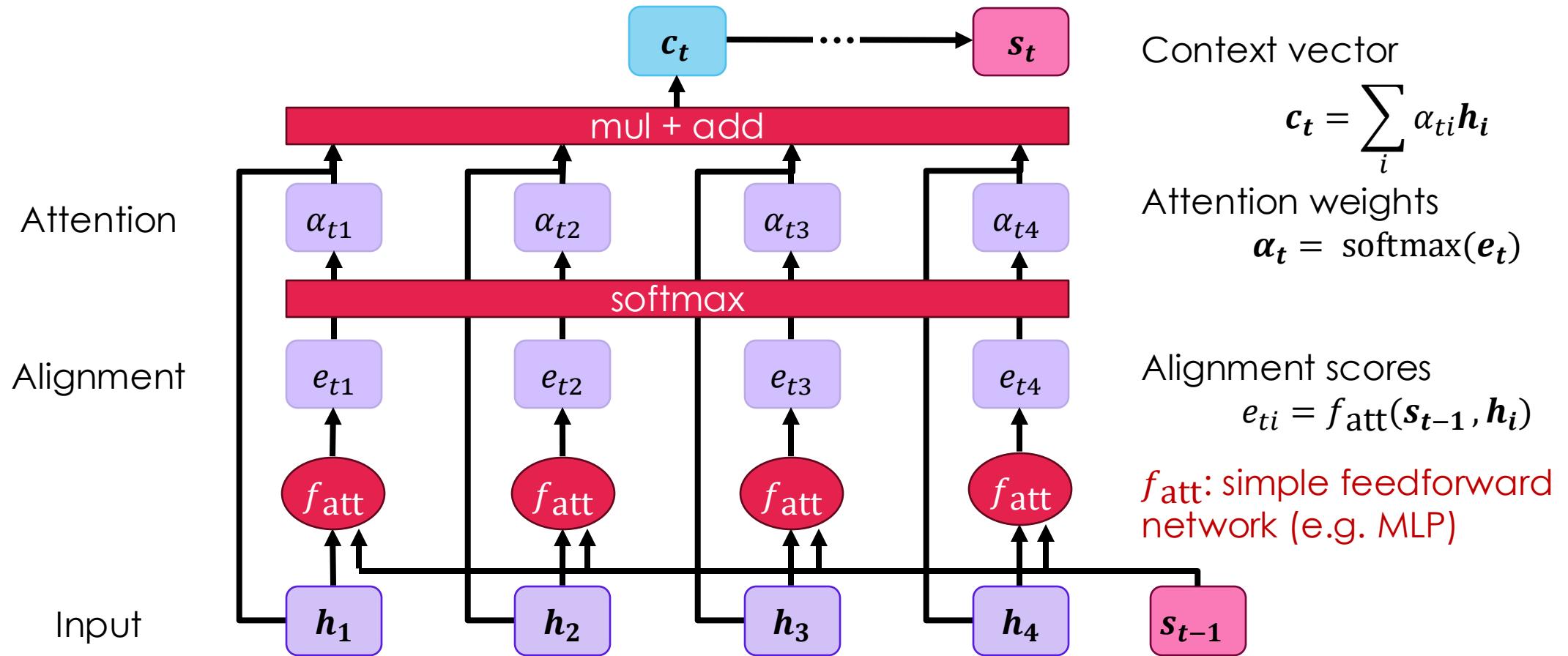
Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

ATTENTION WE'VE SEEN SO FAR

- Known as “**additive**” recurrent attention (type of encoder-decoder attention)



ISSUES WITH RECURRENT ATTENTION

- **Scalability issues:** Performance degrades as the distance between words increases
- **Parallelization limitations:** Recurrent processes lacks ability to be parallelized
- **Memory constraints**
 - RNNs have limited memory and struggle with long-range dependencies
 - Diluted impact of earlier elements on output as sequence progresses
- **Potential solution:** Decouple attention from RNNs
 - How? Separate the attention mechanism into smaller, self-contained components

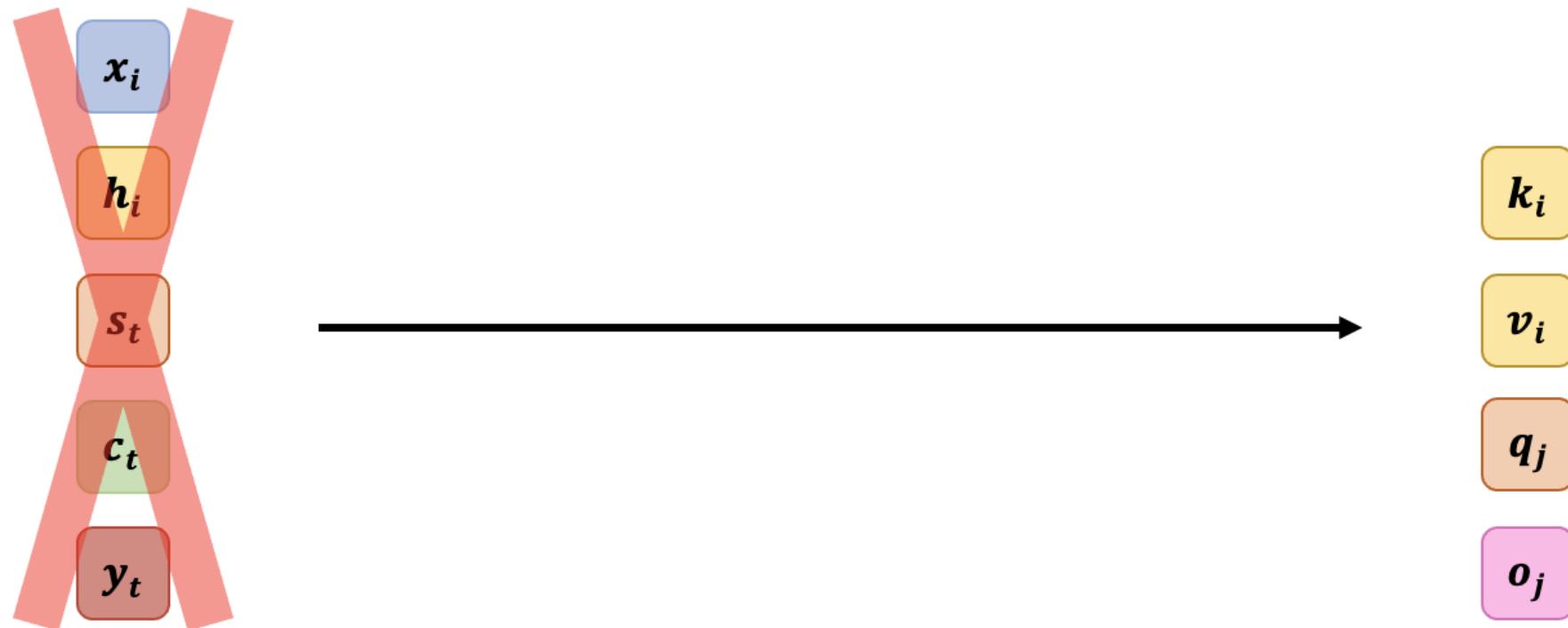
DECOUPLING FROM RNNs

- **Recall:** Attention determines the **importance of elements** to be passed forward in the model.
 - These weights let the model pay **attention** to the **most significant parts**
- **Objective:** A more general attention mechanism not confined to RNNs
- We need a modified procedure to:
 1. Determine weights based on context that indicate the elements to attend to
 2. Apply these weights to enhance attended features

RNN NOTATIONS

x_i	Input for position i in source sequence
h_i	Hidden states for position i in source sequence
s_t	Hidden states for position t in target sequence
c_t	Context vector for position t in target sequence
y_t	Output for position t in target sequence

NEW NOTATION



NEW NOTATION

 k_i

Key vector for position i in an arbitrary sequence

 v_i

Value vector for position i in an arbitrary sequence

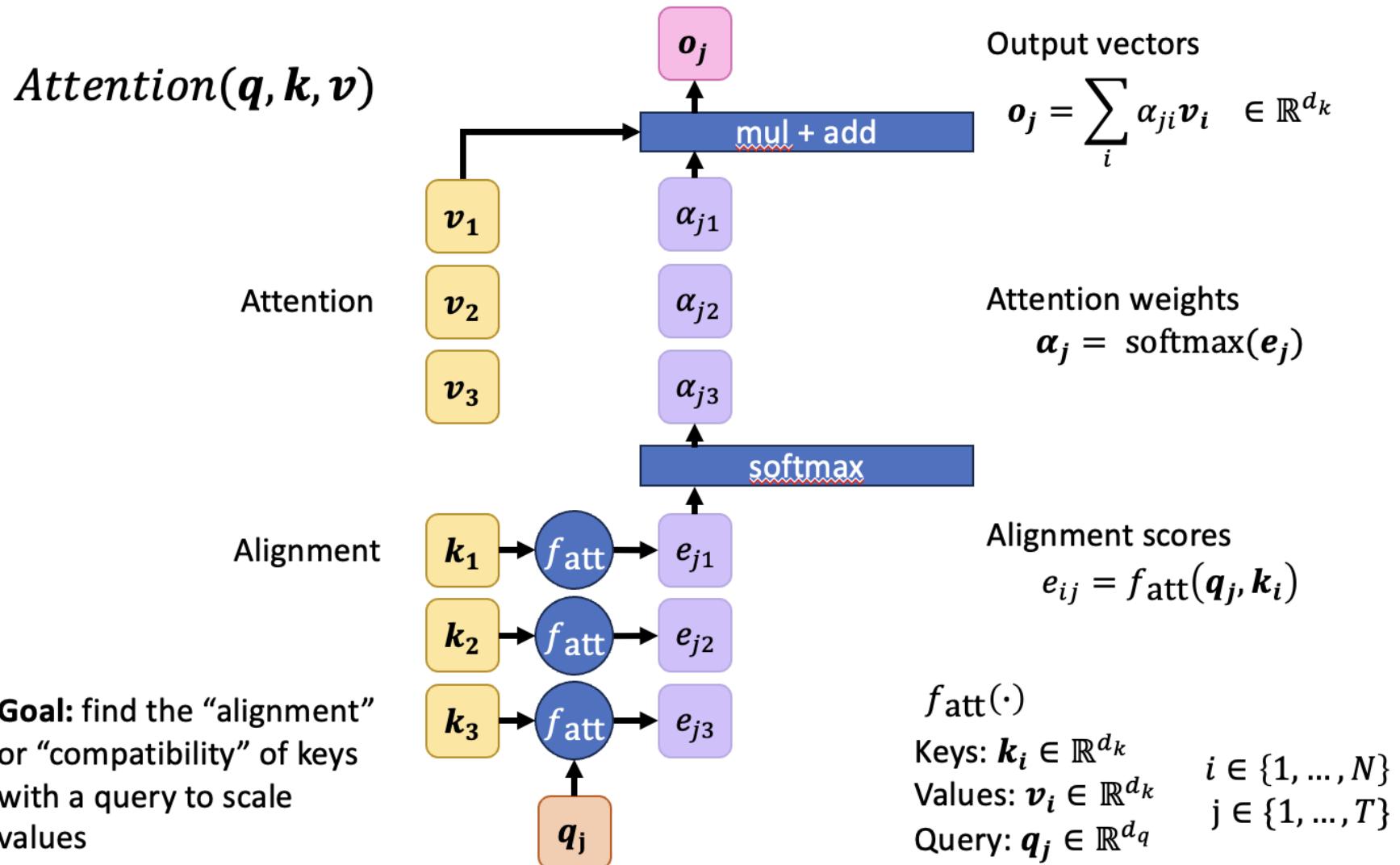
 q_j

Query vector for position j in a (same/different) arbitrary sequence

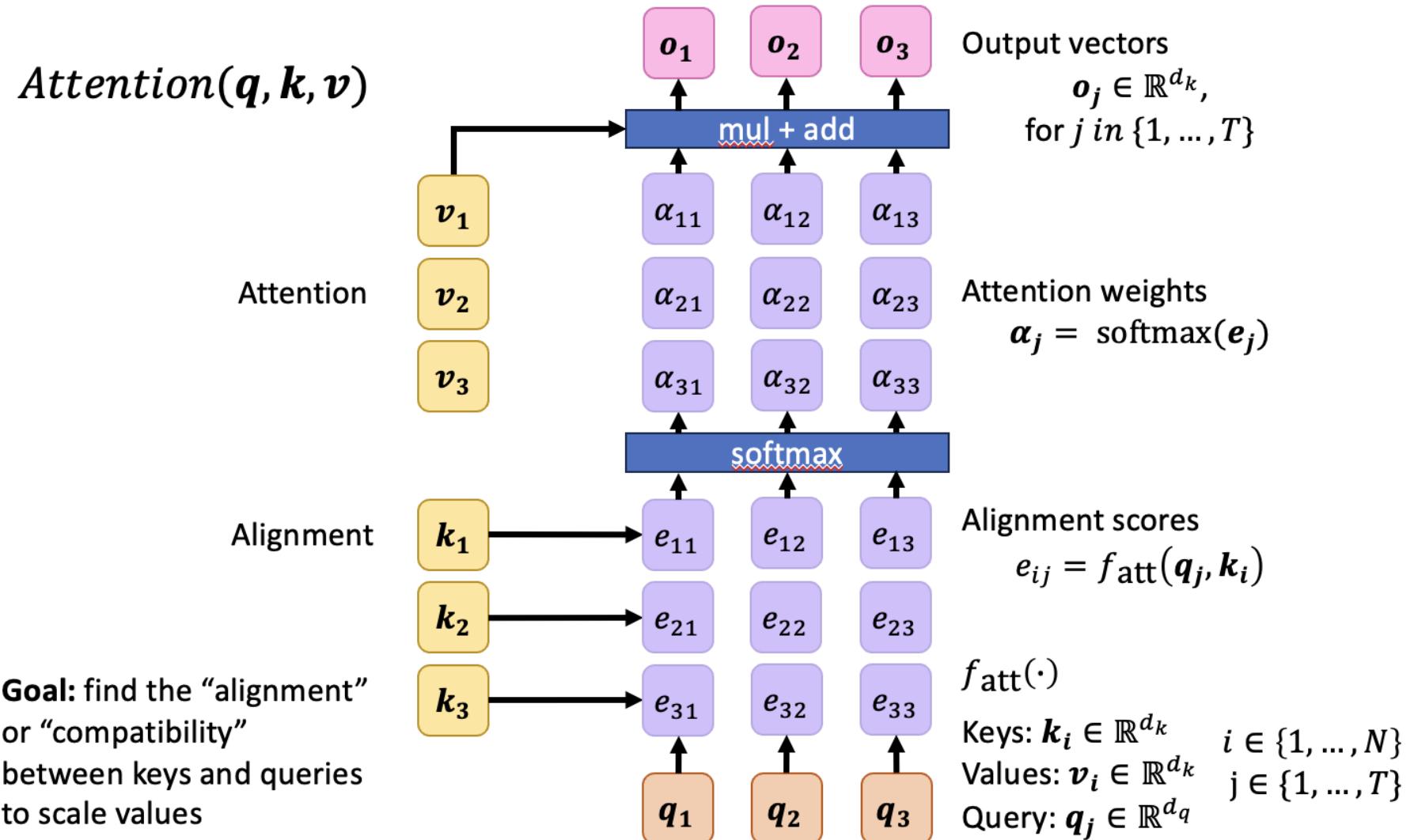
 o_j

Output vector corresponding to position j

A GENERAL ATTENTION MECHANISM



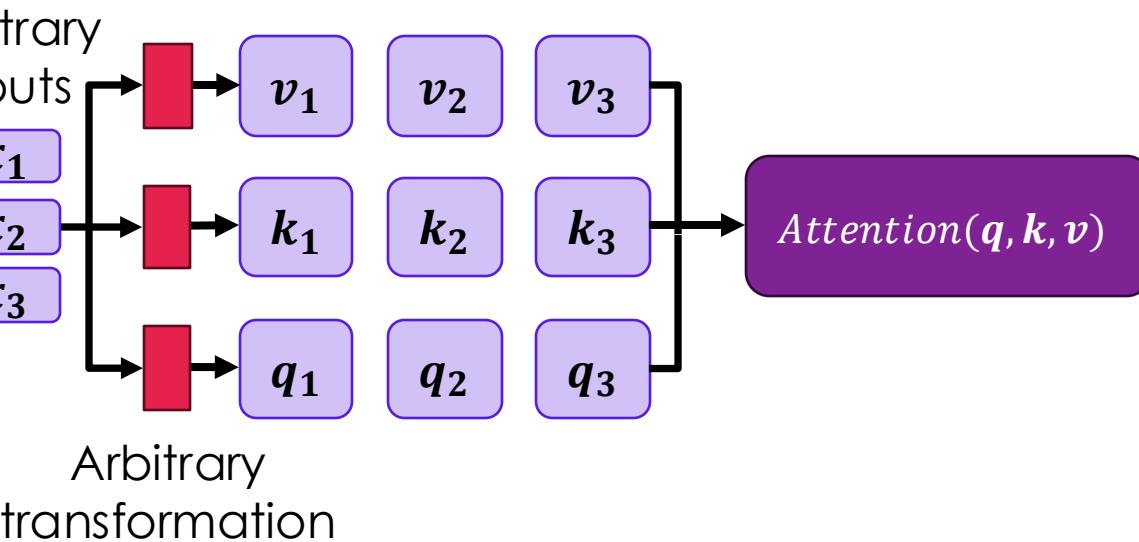
A GENERAL ATTENTION MECHANISM



APPLYING THE ATTENTION MECHANISM

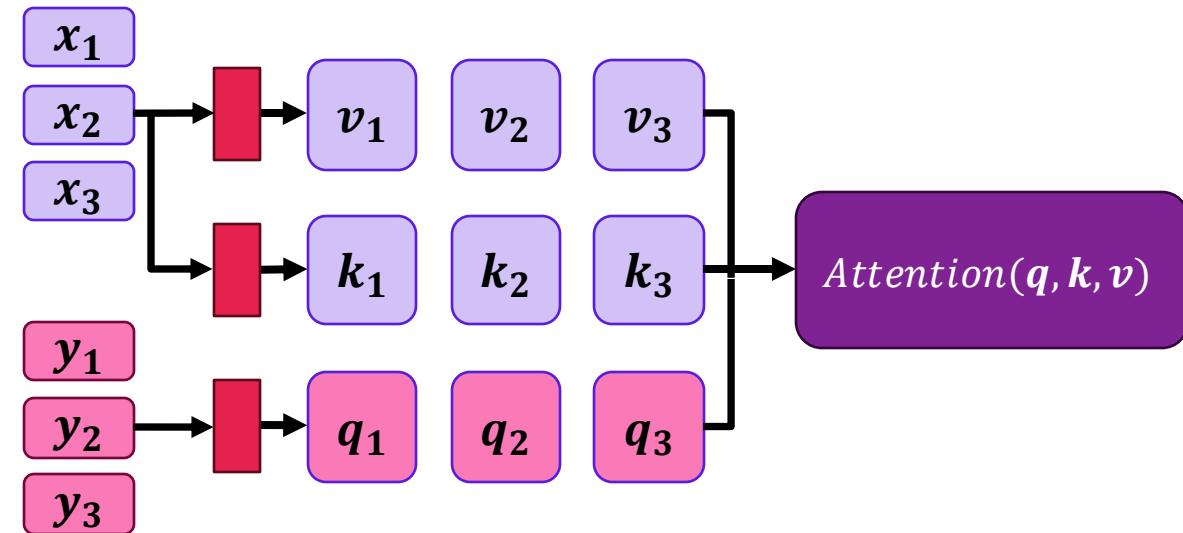
- **Self-Attention**

- Keys, values, and queries are all derived from the same source.



- **Cross-Attention**

- Keys-values and queries are derived from separate sources.



ATTENTION MECHANISM

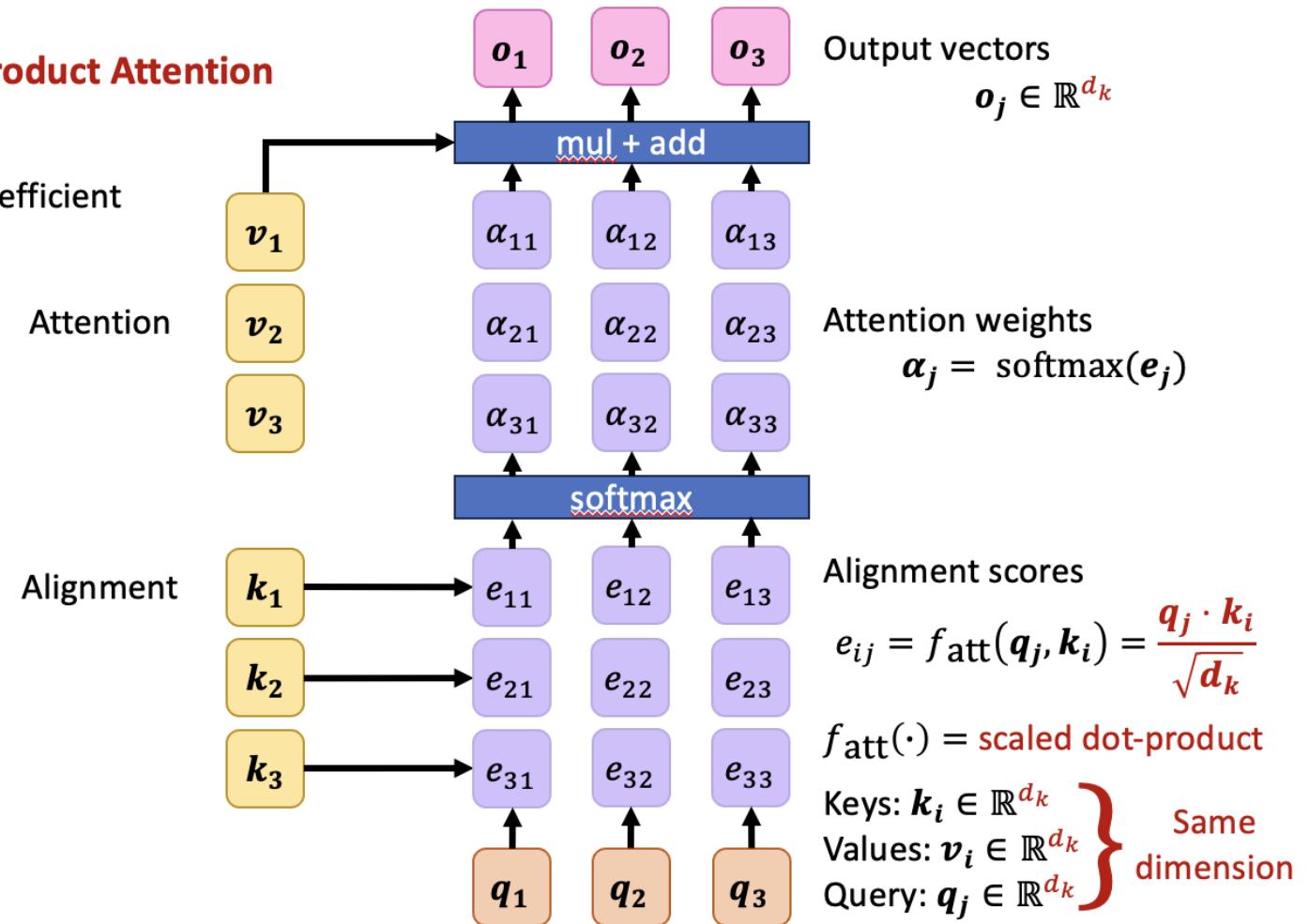
To use a **decoupled attention mechanism**, it is implemented with properties:

1. $f_{\text{att}}(\cdot)$ = **scaled dot-product attention**
 - Good representation of compatibility
 - Fast and interpretable computation
 - Parallelizable evaluation across all queries (can leverage GPUs)
 - Scaled dot-products for stable softmax gradients in high dimensions (prevents large magnitudes)
2. Imposed a **common dimension** for keys, values, and queries
 - Requirement for dot-product
 - Simplifies architecture with predictable attention output shape
 - Provides consistent hidden state dimensions for easier model analysis

ATTENTION MECHANISM

Scaled Dot-Product Attention

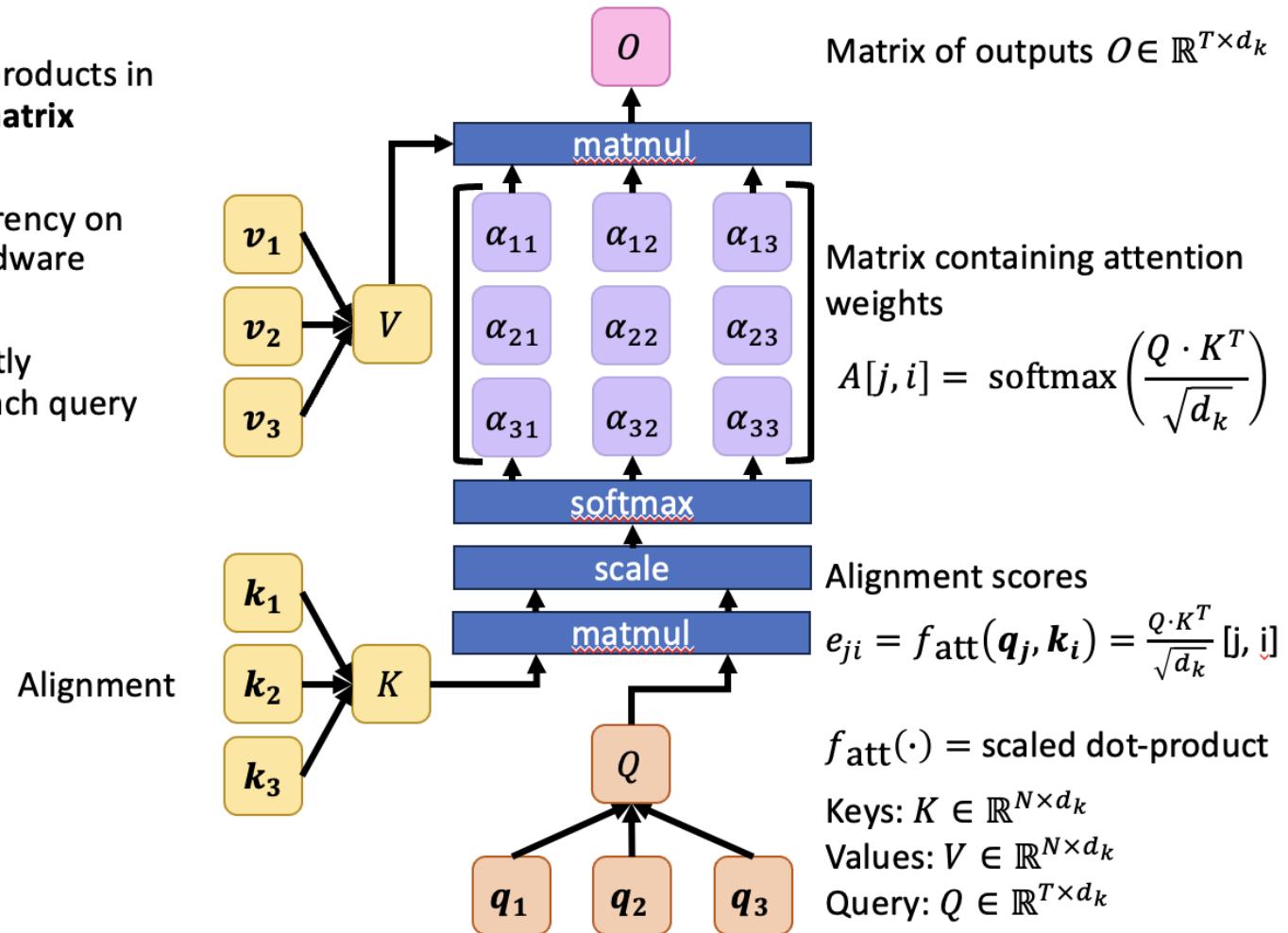
- Faster
- More space-efficient



ATTENTION MECHANISM

Calculate dot-products in **parallel with matrix multiplication**

- High concurrency on modern hardware (GPUs)
- Independently calculates each query



TRANSFORMERS

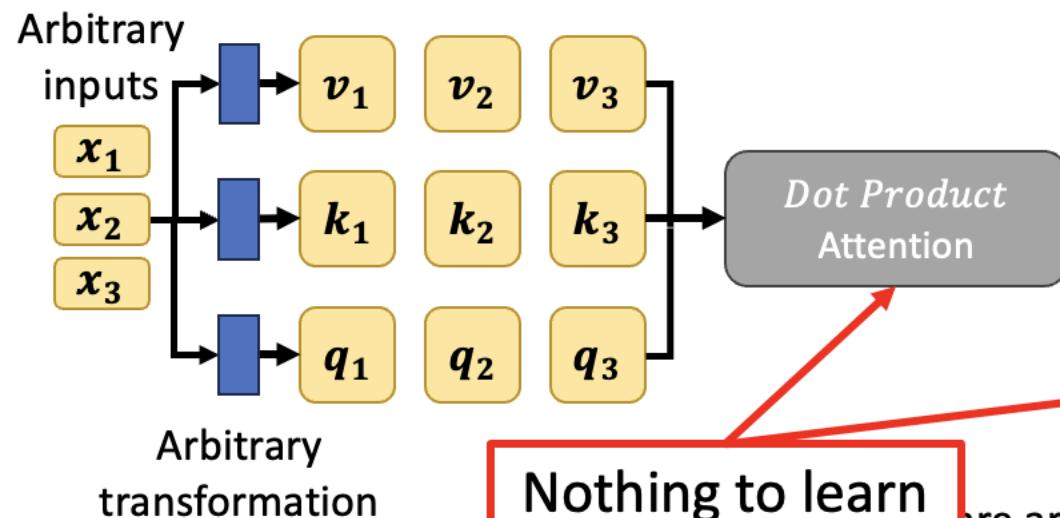
- How do we get Q, K, and V?
- What are we learning?
- Is this parametric or non-parametric?

$$\max \left(\frac{\begin{array}{c} \text{Q} \quad \text{K}^T \\ \text{---} \times \text{---} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \quad \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{array}}{\sqrt{d_k}} \right)$$
$$\begin{array}{c} \text{Z} \\ \text{---} \\ \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{array}$$

TRANSFORMER ATTENTION

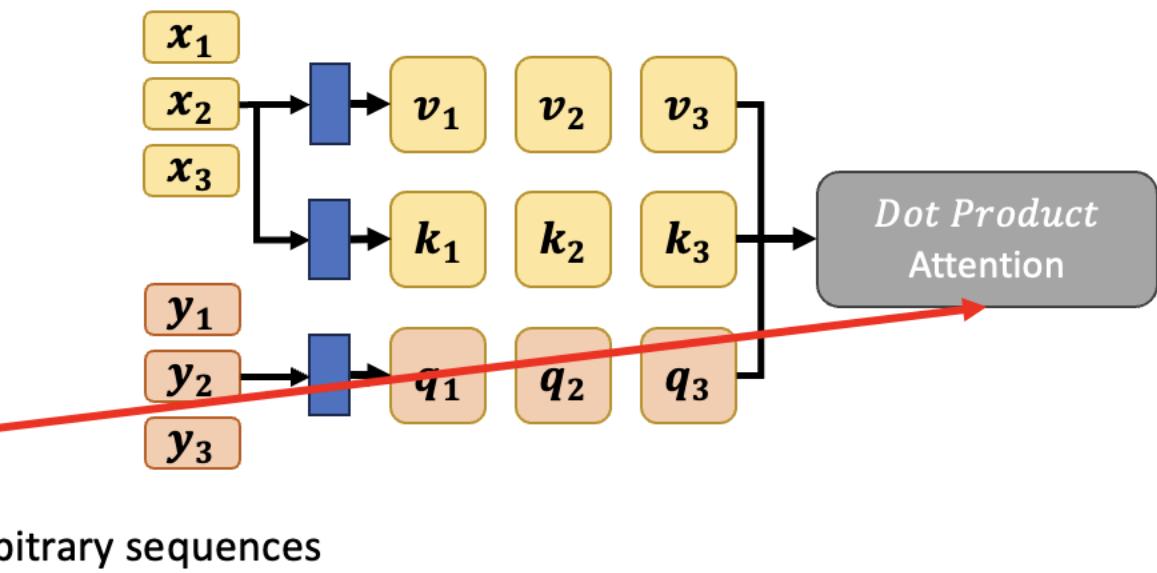
- **Self-Attention**

- Keys, values, and queries are all derived from the same source.



- **Cross-Attention**

- Keys-values and queries are derived from separate sources.



Nothing to learn
inside of this

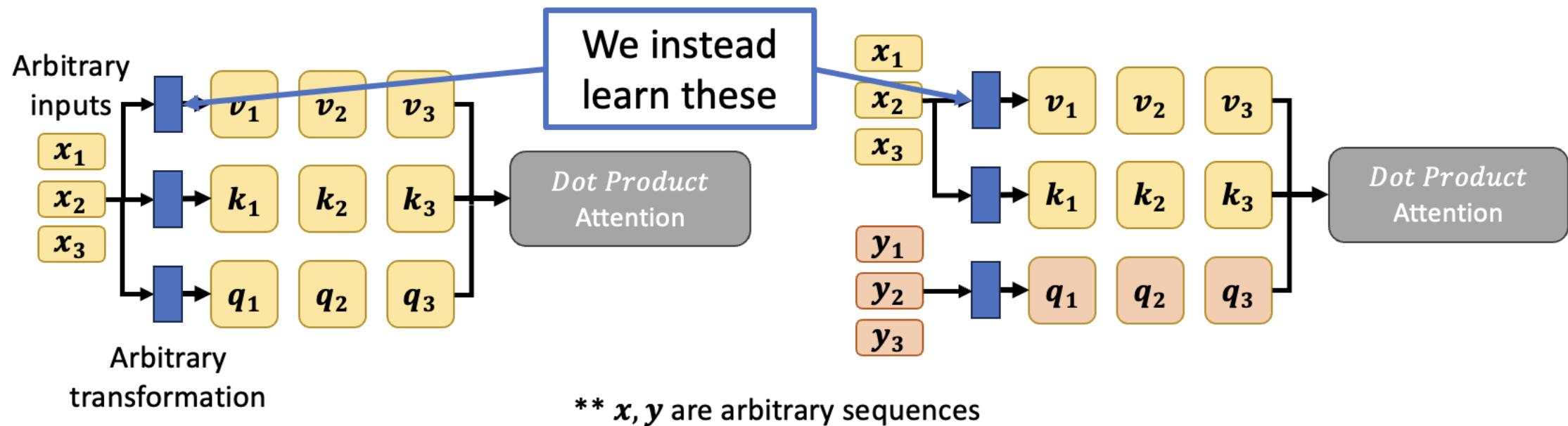
TRANSFORMER ATTENTION

- **Self-Attention**

- Keys, values, and queries are all derived from the same source.

- **Cross-Attention**

- Keys-values and queries are derived from separate sources.



LEARNING TRANSFORMER ATTENTION

Self-Attention

$$\begin{matrix} X \\ \text{[purple grid]} \end{matrix} \times \begin{matrix} WQ \\ \text{[purple grid]} \end{matrix} = \begin{matrix} Q \\ \text{[yellow grid]} \end{matrix}$$

$$\begin{matrix} X \\ \text{[purple grid]} \end{matrix} \times \begin{matrix} WK \\ \text{[orange grid]} \end{matrix} = \begin{matrix} K \\ \text{[yellow grid]} \end{matrix}$$

$$\begin{matrix} X \\ \text{[purple grid]} \end{matrix} \times \begin{matrix} WV \\ \text{[blue grid]} \end{matrix} = \begin{matrix} V \\ \text{[purple grid]} \end{matrix}$$

Cross-Attention

$$\begin{matrix} Y \\ \text{[brown grid]} \end{matrix} \times \begin{matrix} WQ \\ \text{[purple grid]} \end{matrix} = \begin{matrix} Q \\ \text{[brown grid]} \end{matrix}$$

$$\begin{matrix} X \\ \text{[purple grid]} \end{matrix} \times \begin{matrix} WK \\ \text{[orange grid]} \end{matrix} = \begin{matrix} K \\ \text{[yellow grid]} \end{matrix}$$

$$\begin{matrix} X \\ \text{[purple grid]} \end{matrix} \times \begin{matrix} WV \\ \text{[blue grid]} \end{matrix} = \begin{matrix} V \\ \text{[purple grid]} \end{matrix}$$

RESOURCES

- <https://pytorch.org/docs/stable/generated/torch.nn.Transformer.html>
- <https://www.datacamp.com/tutorial/building-a-transformer-with-py-torch>