# Deep Learning

## A course about theory & practice

## Attention and Transformers
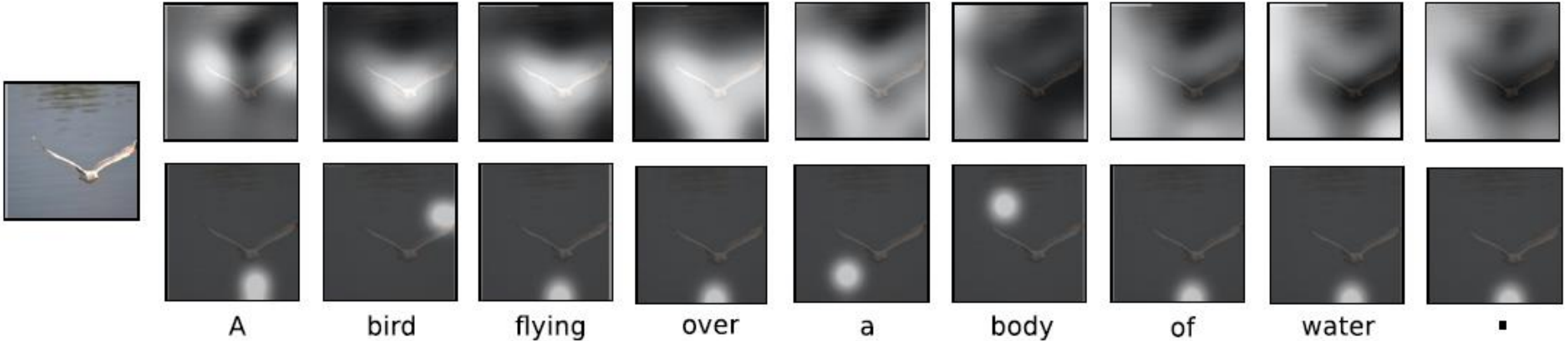
Marco Piastra

# Attention is what we need? (intuition)

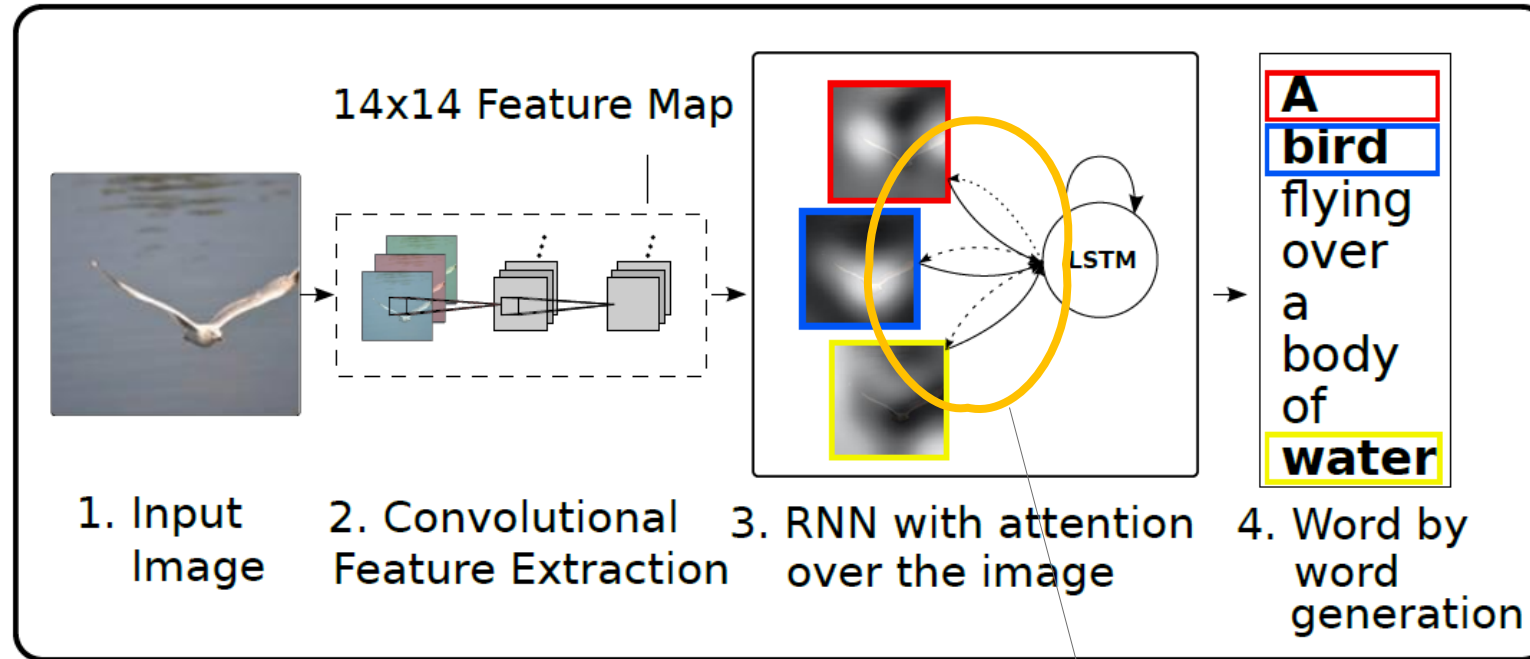# Generating Text Captions from Images

- **DCNN + RNN**

  [*Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*, Xu et al., 2015]

# Generating Text Captions from Images

- ## DCNN + RNN

  [*Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*, Xu et al., 2015]



*The 'trick' is here: when generating each word the LSTM focuses on a specific region in the image*
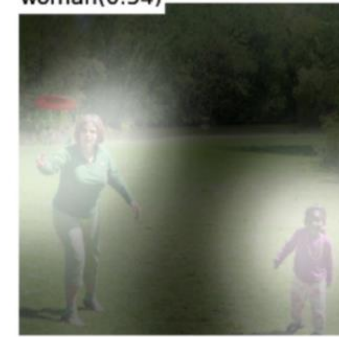
# Generating Text Captions from Images

- **DCNN + RNN**

  [*Show, Attend and Tell: Neural Image Caption Generation with Visual Attention,* Xu et al., 2015]

# Generating Text Captions from Images

- **DCNN + RNN**

  [*Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*, Xu et al., 2015]
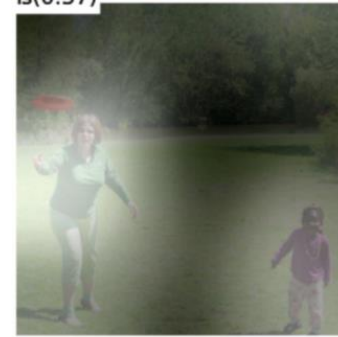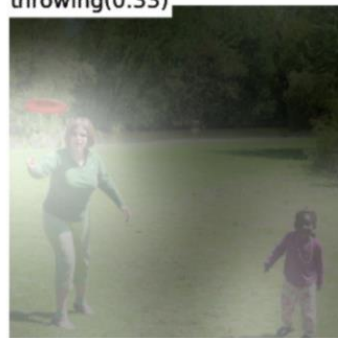


A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

*Look at this: <u>attention</u> focuses on regions that are <u>far apart</u> in the image*

# Natural Language requires Attention

- **Encoder / Decoder with attention**

  [*Long Short-Term Memory-Networks for Machine Reading*, Cheng, Dong and Lapata, 2016]



Current word being read

The machine learns a hidden representation, for *sentiment analysis*,
by focusing on different previous words while reading a sentence

# Attention as a Kernel
# (yet another view of convolution)

# Attention as Kernel

- **Attention Pooling**

  Consider an input-output relation and a dataset

  $$(\boldsymbol{x}, \boldsymbol{y}), \;\; \boldsymbol{x} \in \mathbb{R}^m, \boldsymbol{y} \in \mathbb{R}^n \qquad\qquad D := \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^N$$

  In general *Attention Pooling* is defined as a function on data items

  $$\tilde{\boldsymbol{y}} := \sum_{i=1}^N \alpha(\boldsymbol{x}, \boldsymbol{x}^{(i)}) \, \boldsymbol{y}^{(i)},$$

  Example:

  $$f^*(x) = 2\sin(x_i) + x_i^{0.8}$$

  $$D := \{(x^{(i)}, y^{(i)})\}_{i=1}^N \qquad \textcolor{red}{\textit{Dataset is noisy}}$$

  $$y^{(i)} = f^*(x^{(i)}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.5)$$



[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention as Kernel

- **Attention Pooling**

  Consider an input-output relation

  $$(\boldsymbol{x}, \boldsymbol{y}), \ \ \boldsymbol{x} \in \mathbb{R}^m, \boldsymbol{y} \in \mathbb{R}^n \qquad\qquad D := \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^N$$

  *Attention Pooling* is defined as a function on each input component

  $$\tilde{\boldsymbol{y}} := \sum_{i=1}^{N} \alpha(\boldsymbol{x}, \boldsymbol{x}^{(i)})\, \boldsymbol{y}^{(i)},$$

  Example:

  Global average (i.e., *no attention*)

  $$\alpha(x, x_i) = \frac{1}{N}$$

[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention as Kernel

- **Attention Pooling**

    Consider an input-output relation

    $$(\boldsymbol{x}, \boldsymbol{y}), \ \ \boldsymbol{x} \in \mathbb{R}^m, \boldsymbol{y} \in \mathbb{R}^n \qquad\qquad D := \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^N$$

    *Attention Pooling* is defined as a function on each input component

    $$\tilde{\boldsymbol{y}} := \sum_{i=1}^N \alpha(\boldsymbol{x}, \boldsymbol{x}^{(i)}) \, \boldsymbol{y}^{(i)},$$

    Example:

    Gaussian Kernel [Nadaraya & Watson, 1964]

$$\alpha(x, x^{(i)}) = \frac{K(x - x^{(i)})}{\sum_{j=1}^N K(x - x^{(j)})} \qquad K(u) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{u^2}{2})$$



[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention as Kernel

- **Gaussian Kernel and Softmax**

$$\alpha(x, x^{(i)}) = \frac{K(x - x^{(i)})}{\sum_{j=1}^{N} K(x - x^{(j)})} \qquad K(u) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{u^2}{2})$$

*'Attention Map'*



$$\alpha(x, x^{(i)}) = \frac{\exp\left(-\frac{1}{2}(x - x^{(i)})^2\right)}{\sum_{j=1}^{N} \exp\left(-\frac{1}{2}(x - x^{(j)})^2\right)}$$

$$= \text{softmax}\left(-\frac{1}{2}(x - x^{(i)})^2\right)$$

Gaussian kernel regression converges to the optimal solution, as the dataset increases

*Note that Gaussian Kernel is non-parametric: it is a pure pooling operation*

[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention as Kernel

- **(Simple) Parametric Attention Pooling**

$$\alpha(x, x_i) = \frac{\exp\left(-\frac{1}{2}(x - x_i)^2 \, w\right)}{\sum_{j=1}^{N} \exp\left(-\frac{1}{2}(x - x_j)^2 \, w\right)}$$

$$= \mathrm{softmax}\left(-\frac{1}{2}(x - x_i)^2 w\right)$$

$$K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2 w}{2}\right)$$

$$\Rightarrow \ \sigma^2 = \frac{1}{w}$$

This requires training of the (unique) parameter $w$

Consider an MSE *loss* function:

$$L(D) = \frac{1}{N} \sum_{i=1}^{N} (f(x^{(i)}) - y^{(i)})^2$$

and perform *gradient descent*

# Attention as Kernel

- **(Simple) Parametric Attention Pooling**

$$\alpha(x, x_i) = \frac{\exp\left(-\frac{1}{2}(x - x_i)^2\, w\right)}{\sum_{j=1}^{N} \exp\left(-\frac{1}{2}(x - x_j)^2\, w\right)}$$

$$= \mathrm{softmax}\left(-\frac{1}{2}(x - x_i)^2 w\right)$$

This requires training of the (unique) parameter $w$

Consider an MSE *loss* function:

$$L(D) = \frac{1}{N} \sum_{i=1}^{N} (f(x^{(i)}) - y^{(i)})^2$$

and perform *gradient descent*

<span style="color:red">The (Gaussian) attention field becomes 'sharper'</span>



[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention as Kernel

- **Terminology**

    In the following:

    - *data* items will be referred to as *keys*

    - *input* items will be referred to as *queries*

    *(This is field-specific jargon)*



queries

keys

[image from http://d2l.ai/chapter_attention-mechanisms/]

# Attention: Queries, Key and Values

# Is Attention a Kernel?

*The phenomenon of inversion in language translation*



[image from https://arxiv.org/pdf/1409.0473]

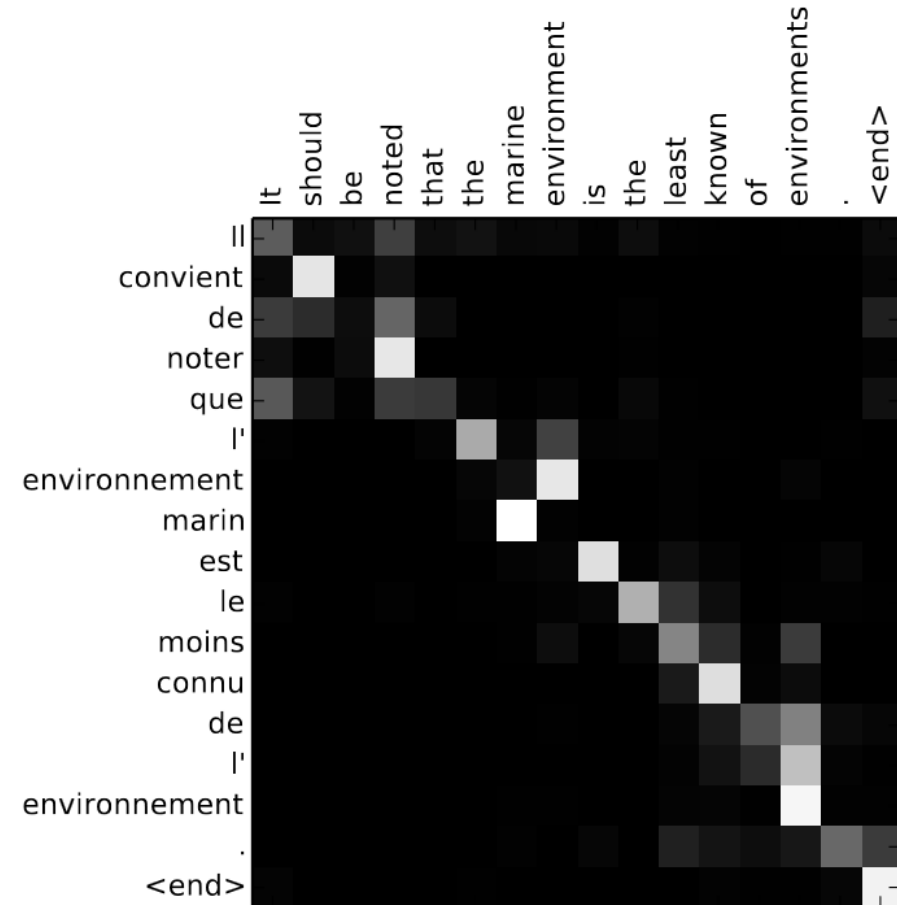# Attention Pooling: Queries, Keys and Values

- **Generalized model**

$$\text{Attention}(\boldsymbol{q}, \boldsymbol{k}, \boldsymbol{v}) := \sum_{i=1}^{m} \alpha(\boldsymbol{q}, \boldsymbol{k}_i) \boldsymbol{v}_i$$

$$\boldsymbol{q} \in \mathbb{R}^q, \ \boldsymbol{k} \in \mathbb{R}^k, \ \boldsymbol{v} \in \mathbb{R}^v$$

In such Attention Pooling:

- *queries* and *keys* could come from different spaces
- the *attention map* $\alpha$ is <u>normalized</u>: it describes how attention is *distributed*
- *values* are specific contributions (in general, sizes are equal $k = v$ )
- $m$ is the width of the receptive field (=*how many keys are in it*)

# Attention: Queries, Keys and Values

- **Generalized model**

$$\text{Attention}(\boldsymbol{q}, \boldsymbol{k}, \boldsymbol{v}) := \sum_{i=1}^{m} \alpha(\boldsymbol{q}, \boldsymbol{k}_i)\boldsymbol{v}_i$$

$$\boldsymbol{q} \in \mathbb{R}^q, \ \boldsymbol{k} \in \mathbb{R}^k, \ \boldsymbol{v} \in \mathbb{R}^v$$

The *attention map* is defined as:

$$\alpha(\boldsymbol{q}, \boldsymbol{k}_i) = \text{softmax}(a(\boldsymbol{q}, \boldsymbol{k}_i)) = \frac{\exp(a(\boldsymbol{q}, \boldsymbol{k}_i))}{\sum_{j=1}^{m} \exp(a(\boldsymbol{q}, \boldsymbol{k}_j))}$$

where $a$ is the *attention scoring function* of choice

# Attention Scoring Function

- **Scaled Dot-Product Attention**

    Assume that both queries and keys encoded as vectors of size $d$

    The *attention scoring function* is defined as:

    $$a(\boldsymbol{q}, \boldsymbol{k}) := \frac{\boldsymbol{q} \cdot \boldsymbol{k}}{\sqrt{d}}, \;\; \boldsymbol{q}, \boldsymbol{k} \in \mathbb{R}^d$$

    In the line of principle, $\boldsymbol{q}$ and $\boldsymbol{k}$ could be anything, including the output of other *layers*

    *The normalizing term $\sqrt{d}$ comes from the assumptions that each component of the encodings is an independent random variable with zero mean and unit standard deviation*
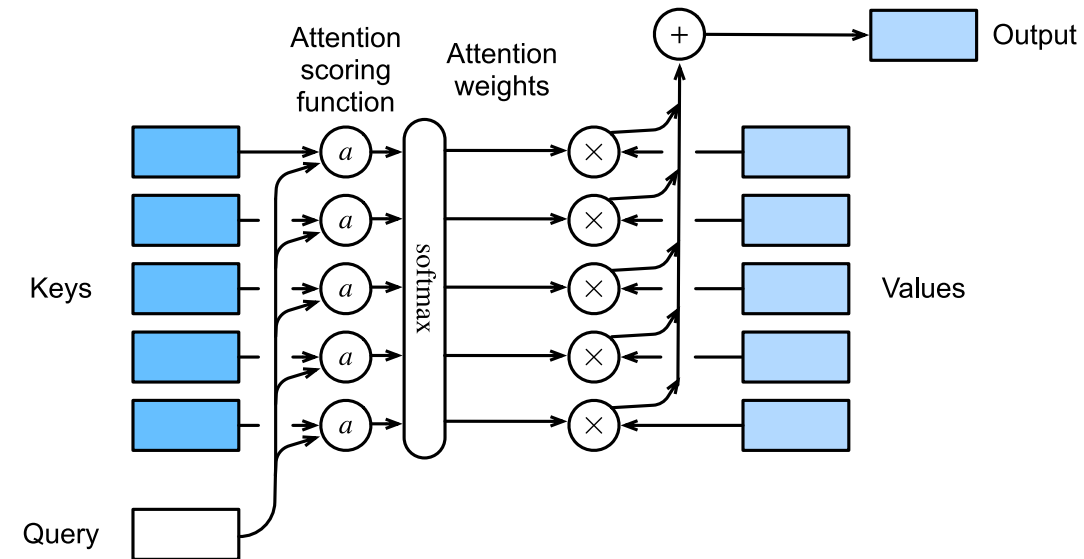
# Attention Scoring Function

- ### Scaled Dot-Product Attention

Using a tensorial representation, assume there are $m$ keys, $n$ queries and $v$ values:

$$\boldsymbol{Q} \in \mathbb{R}^{n \times d}, \ \ \boldsymbol{K} \in \mathbb{R}^{m \times d}, \ \ \boldsymbol{V} \in \mathbb{R}^{m \times v}$$

Attention Pooling becomes:

$$\text{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) := \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d}}\right)\boldsymbol{V} \ \ \in \mathbb{R}^{n \times v}$$

# Attention Map

- **Example: a square matrix**

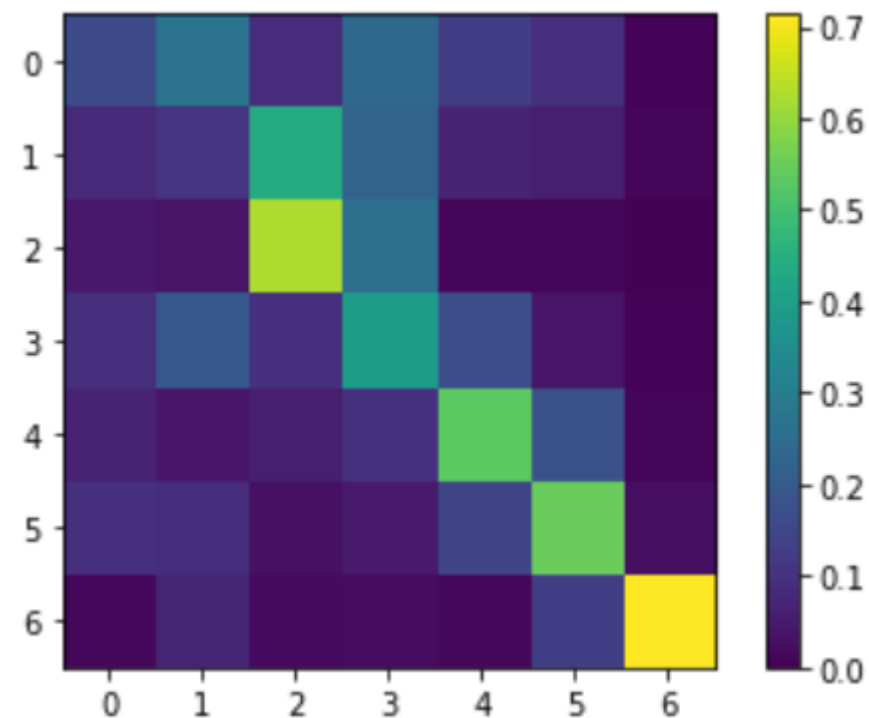  When both *queries* and *keys* come from <u>the same source</u> (i.e., *self-attention*)
  Namely:

  $$\boldsymbol{Q} \in \mathbb{R}^{n \times d}, \ \ \boldsymbol{K} \in \mathbb{R}^{n \times d}$$

  Then:

  $$\alpha(\boldsymbol{Q}, \boldsymbol{K}) := \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d}}\right) \ \ \in \mathbb{R}^{n \times n}$$

# Scaled Dot-Product Attention

$$\boldsymbol{Q} \in \mathbb{R}^{n \times d}, \ \ \boldsymbol{K} \in \mathbb{R}^{n \times d}, \ \ \boldsymbol{V} \in \mathbb{R}^{n \times v}$$

$$\text{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) := \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d}}\right) \boldsymbol{V} \quad \in \mathbb{R}^{n \times v}$$

*This is the basic building block*



Scaled Dot-Product Attention

[image from https://arxiv.org/pdf/1706.03762.pdf]

# Scaled Dot-Product Self-Attention

$$\boldsymbol{Q} \in \mathbb{R}^{n \times d}, \ \boldsymbol{K} \in \mathbb{R}^{n \times d}, \ \boldsymbol{V} \in \mathbb{R}^{n \times n}$$

$$\mathrm{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) := \mathrm{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d}}\right)\boldsymbol{V} \quad \in \mathbb{R}^{n \times n}$$

In the case of *self-attention*, further flexibility can be gained by adding a *linear transformation*:
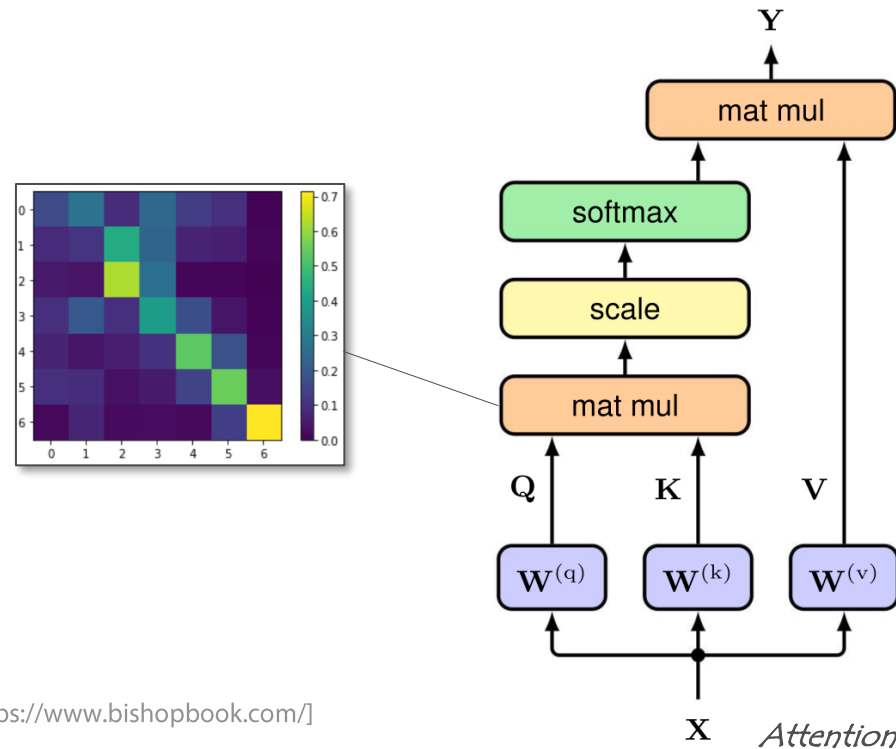
$$\boldsymbol{Q} = \boldsymbol{X}\boldsymbol{W}_q^T$$
$$\boldsymbol{K} = \boldsymbol{X}\boldsymbol{W}_k^T$$
$$\boldsymbol{V} = \boldsymbol{X}\boldsymbol{W}_v^T$$

where:

$$\boldsymbol{W}_q, \ \boldsymbol{W}_k, \ \boldsymbol{W}_v \in \mathbb{R}^{n \times d}$$

are parameters to be trained



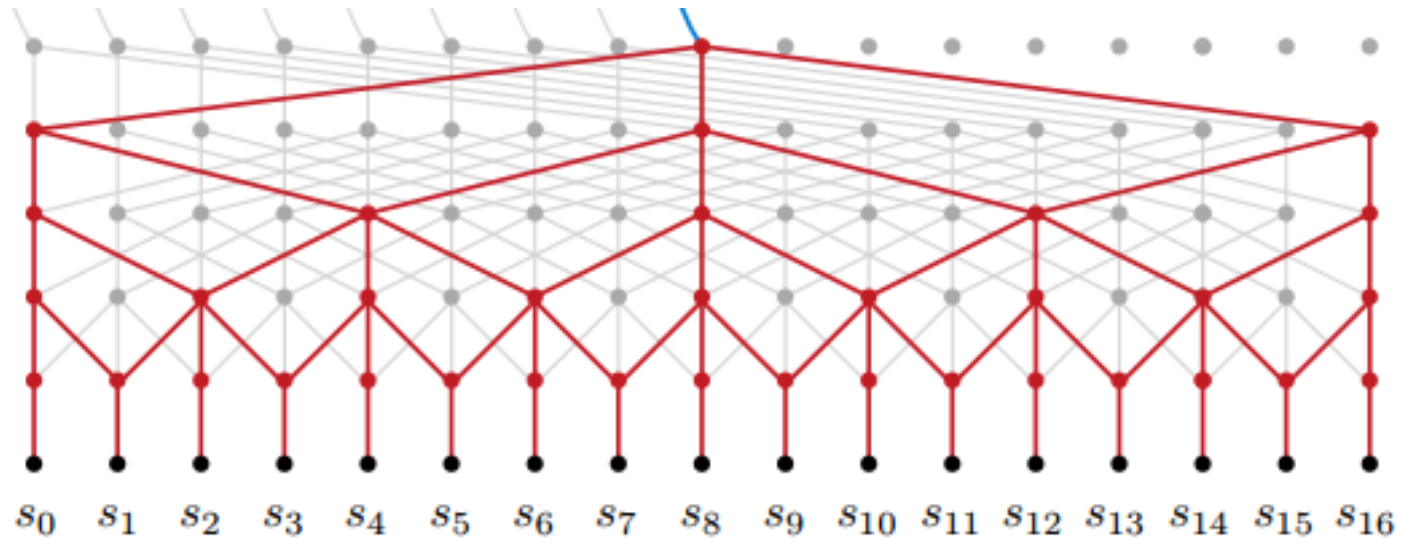[image from https://www.bishopbook.com/]

*Attention and Transformers* [24]

# Attention vs. Convolution: why should it be better?

# Attention vs Convolution

- **Progressively widening receptive field**

  Consider 1D convolution, size 3:
  the receptive field of each filter grows progressively



  $s_0 \quad s_1 \quad s_2 \quad s_3 \quad s_4 \quad s_5 \quad s_6 \quad s_7 \quad s_8 \quad s_9 \quad s_{10} \quad s_{11} \quad s_{12} \quad s_{13} \quad s_{14} \quad s_{15} \quad s_{16}$

*Problem: four layers are required in this case to have a receptive field of 16*

# Attention vs Convolution vs RNN



CNN

RNN

Self-attention

Convolution

Fully Connected layer

Global attention

Local attention

[image from http://d2l.ai/]

[image from https://theaisummer.com/attention/]

# Positional Encoding

# Positional Encoding

- **Using sine and cosine**

  Assume we want to sum a *positional encoding* vector $p$ to a data vector $x$ ( $x, p \in \mathbb{R}^d$ )

  $$x + p$$

  To do so, we can use a *sine − cosine* representation:

  $$p_{i,2j} = \sin\left(\frac{i}{s^{2j/d}}\right)$$

  $$p_{i,2j+1} = \cos\left(\frac{i}{s^{2j/d}}\right)$$

  Where $i$ is the *position index* of data vector $x$ (*query* or *key*), $j$ is one of the $d$ components

  of vector $p$ and $s$ is a suitable scale constant (in the original paper $s = 1000$ )

# Positional Encoding

- **Using sine and cosine**

  Assume we want to sum a *positional encoding* vector $p$ to a data vector $x$ ( $x, p \in \mathbb{R}^d$ )

  $$x + p$$

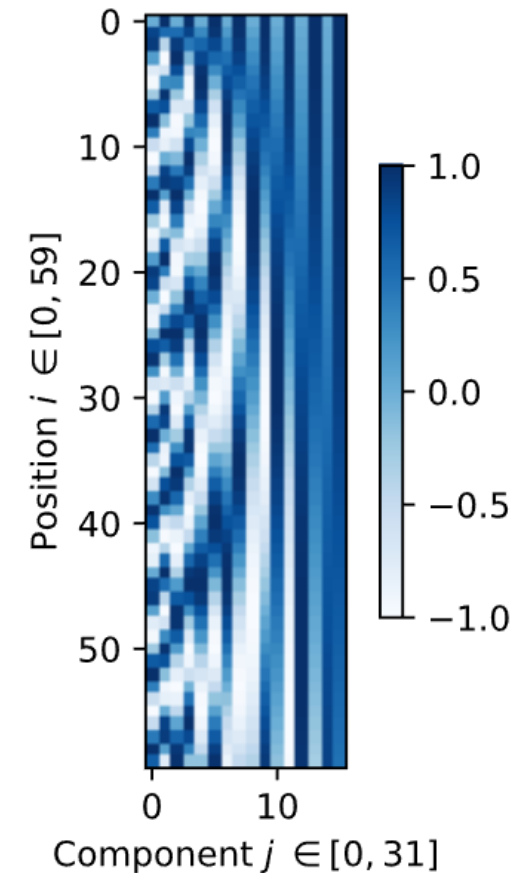  To do so, we can use a *sine – cosine* representation:

  $$p_{i,2j} = \sin\left(\frac{i}{s^{2j/d}}\right)$$

  $$p_{i,2j+1} = \cos\left(\frac{i}{s^{2j/d}}\right)$$

  *Position can be either absolute or relative,*
  *to the position of each <u>query</u> or <u>key</u>*

  <span style="color:red">*In high dimensions,*
  *two randomly chosen vectors tend to be <u>orthogonal</u>*</span>

# Positional Encoding

- **Relative displacements**

$$
\begin{bmatrix} \cos(\delta\omega_j) & \sin(\delta\omega_j) \\ -\sin(\delta\omega_j) & \cos(\delta\omega_j) \end{bmatrix} \begin{bmatrix} p_{i,2j} \\ p_{i,2j+1} \end{bmatrix}
$$

$$
= \begin{bmatrix} \cos(\delta\omega_j)\sin(i\omega_j) + \sin(\delta\omega_j)\cos(i\omega_j) \\ -\sin(\delta\omega_j)\sin(i\omega_j) + \cos(\delta\omega_j)\cos(i\omega_j) \end{bmatrix}
$$

$$
= \begin{bmatrix} \sin\left((i+\delta)\omega_j\right) \\ \cos\left((i+\delta)\omega_j\right) \end{bmatrix}
$$

$$
= \begin{bmatrix} p_{i+\delta,2j} \\ p_{i+\delta,2j+1} \end{bmatrix},
$$

Displacements can be represented via a linear transformation.
This means that relative positions can be _learnt_

# Transformers:
# a network architecture

# Multiple Attention Heads

Multi-head attention consists of four parts:

1. Linear layers *(i.e., fully connected, no activation function)*

2. Scaled dot-product attention

3. Output concatenation

4. Final linear layer

*Each input combination of Q (query), K (key), V (values)*
*is passed to each a separate linear layer hence to an attention head*

*The output of multiple attention heads is the concatenated*
*and fed to a final linear layer*

[image from https://arxiv.org/pdf/1706.03762.pdf]

# Encoder Layer

Each encoder layer includes:

1. Multi-head attention

2. Addition *(ResNet style)*

3. Normalization *(per each input)*

4. Feed-forward network
   *(one hidden layer with ReLU plus one linear layer)*

5. Addition

6. Normalization

*There could be many encoder layers in the overall architecture*



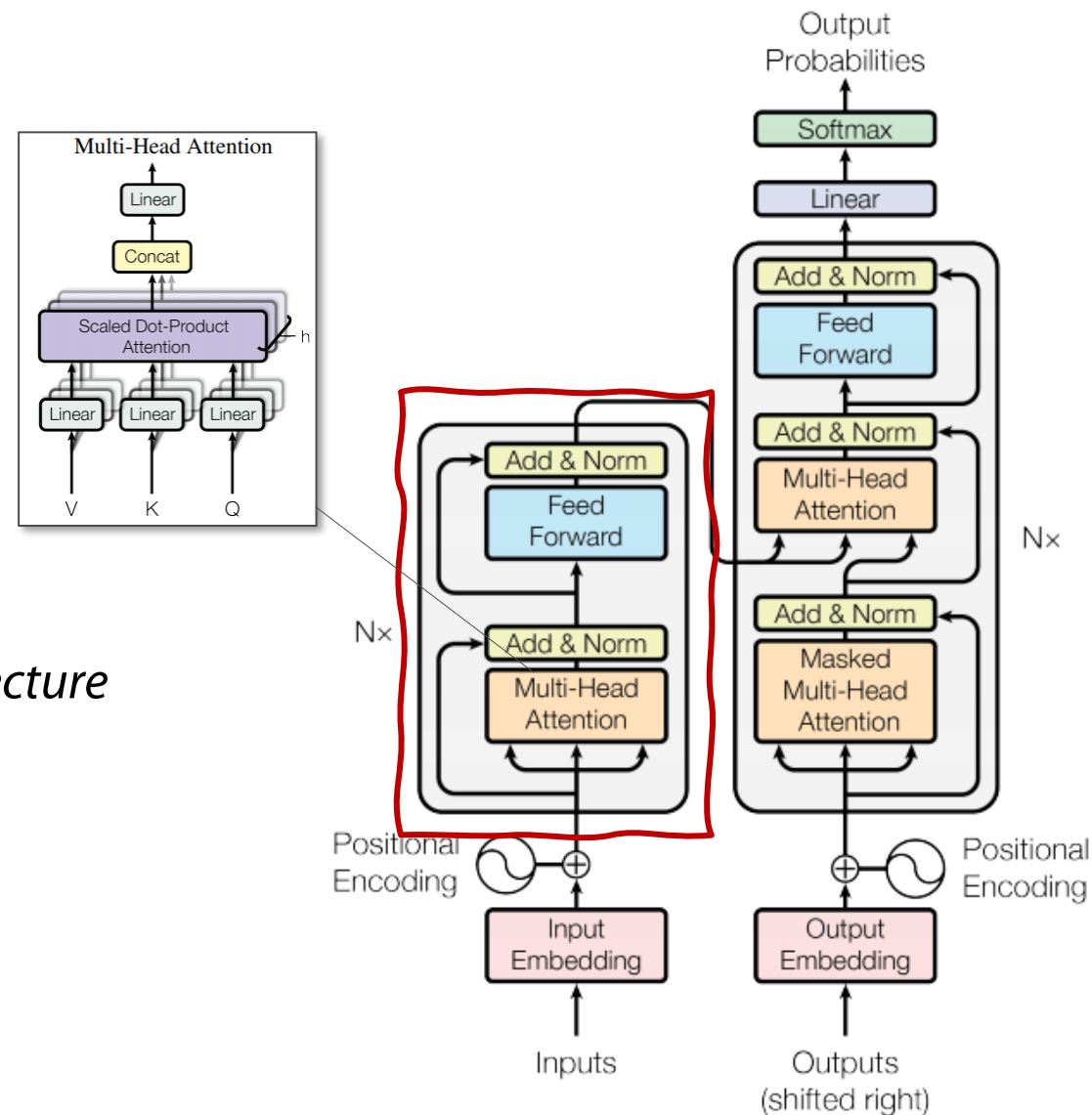[image from https://arxiv.org/pdf/1706.03762.pdf]
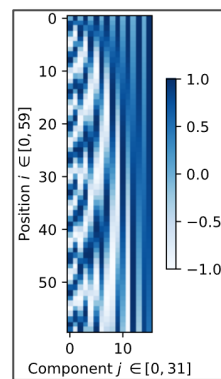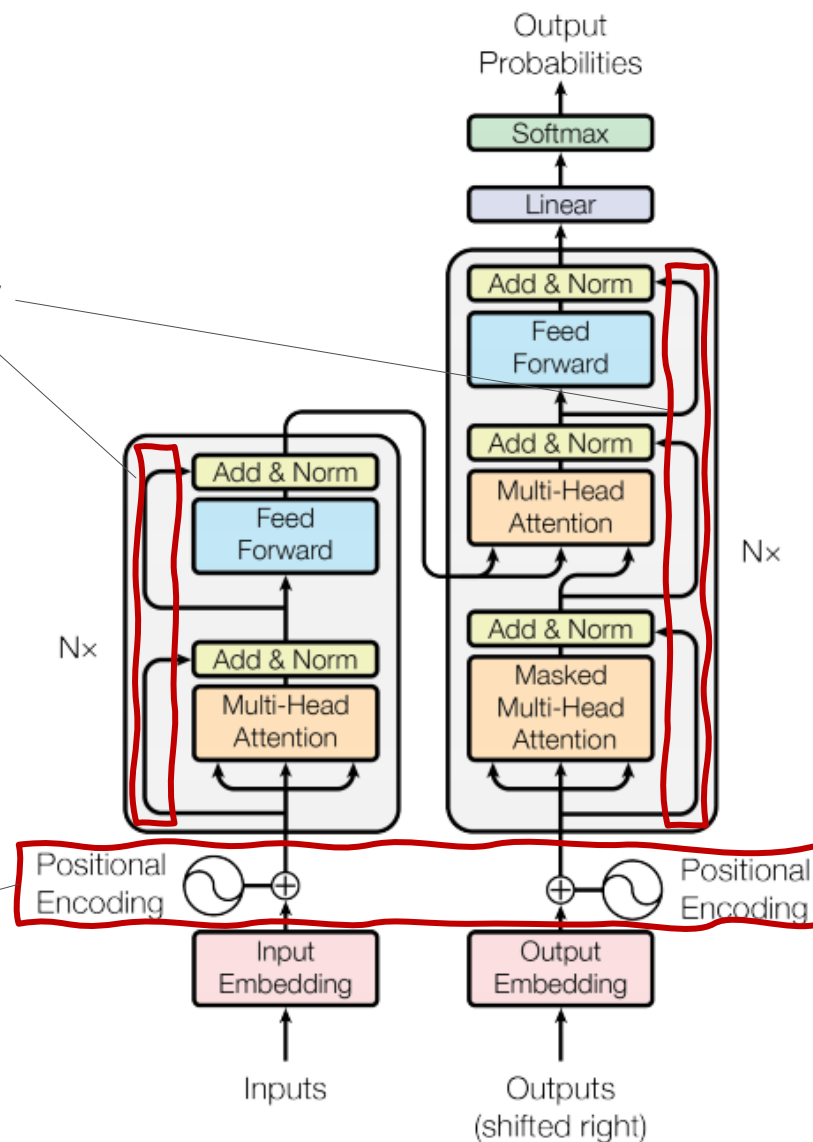
# Encoder Layer

Each encoder layer includes:

1. Multi-head attention

2. Addition *(ResNet style)*

3. Normalization *(per each input)*

4. Feed-forward network
   *(one hidden layer with ReLU plus one linear layer)*

5. Addition

6. Normalization

*There could be many encoder layers in the overall architecture*

residual connections



[image from https://arxiv.org/pdf/1706.03762.pdf]

# Decoder Layer

Each decode layer includes:

1. Multi-head attention

2. Addition

3. Normalization

4. Multi-head attention
   *values* and *keys* come from the encoder output
   while *queries* come from the previous decoder layer

5. Addition

6. Normalization

7. Feed-forward network
   *(one hidden layer with ReLU plus one linear layer)*

8. Addition

9. Normalization

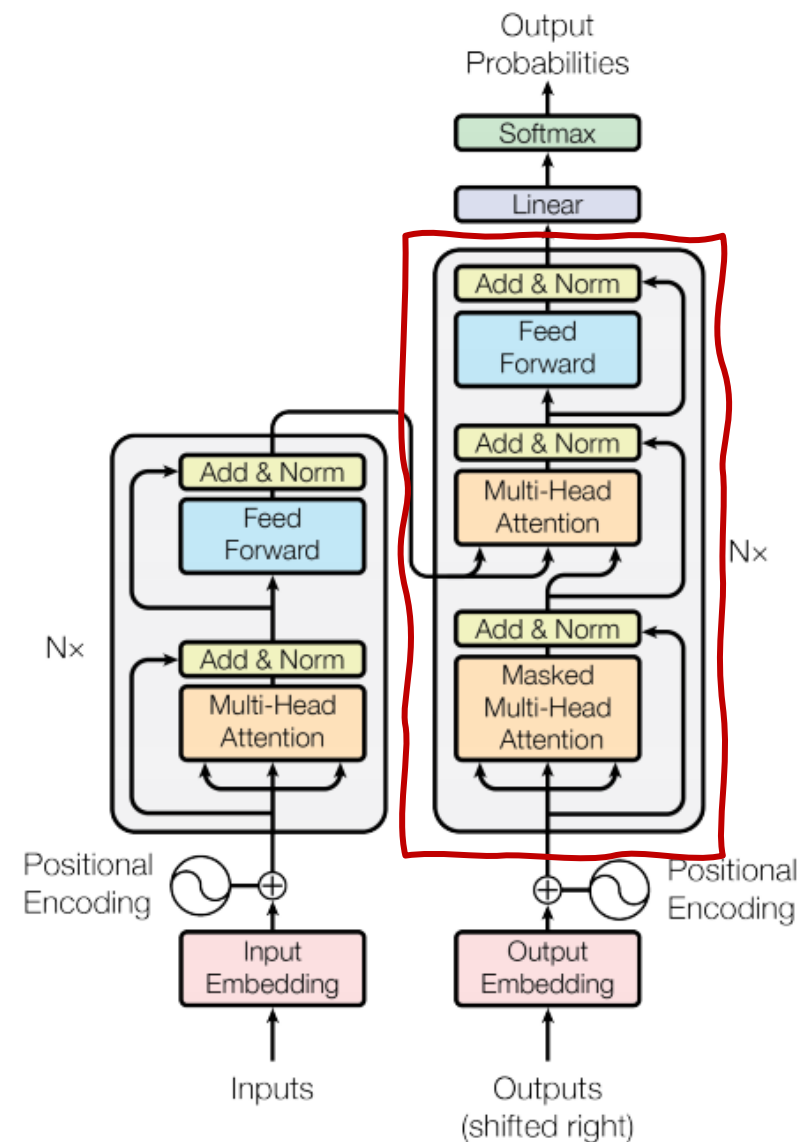*There could be many decoder layers in the overall architecture*

[image from https://arxiv.org/pdf/1706.03762.pdf]

# Decoder Layer

*Why <u>masked</u> multi-head attention in the decoder layer?*

The production of the output is incremental: one word at time

The output embedding 'input' can only see what has been generated thus far

Masks are not trained, they are 'superimposed' as the generation process advances

# Encoder

The encoder block includes:

1. Input embedding *(word2vec style)*
2. Positional encoding
3. Addition
4. N encoder layers



[image from https://arxiv.org/pdf/1706.03762.pdf]

# Decoder

The decoder block includes:

1.  Output embedding *(word2vec style)*
    *It encodes the output produced so far*

2.  Positional encoding

3.  Addition

4.  N decoder layers
    Each connected to a corresponding encoder layer

5.  Linear layer

6.  Softmax layer
    *It predicts the next token in the sequence*



[image from https://arxiv.org/pdf/1706.03762.pdf]

# Translator (Encoder-Decoder)

this is the first book i've ever done.

Feedback



Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

este é o primeiro livro que eu fiz.

this is the first book i've ever done

[image from https://arxiv.org/pdf/1706.03762.pdf]

# Translator (Encoder-Decoder)

**During Training**

**`this is the first book i've ever done.`**

*This is the predicted output*

*Loss is computed on the softmax*



Softmax

Linear

Add & Norm
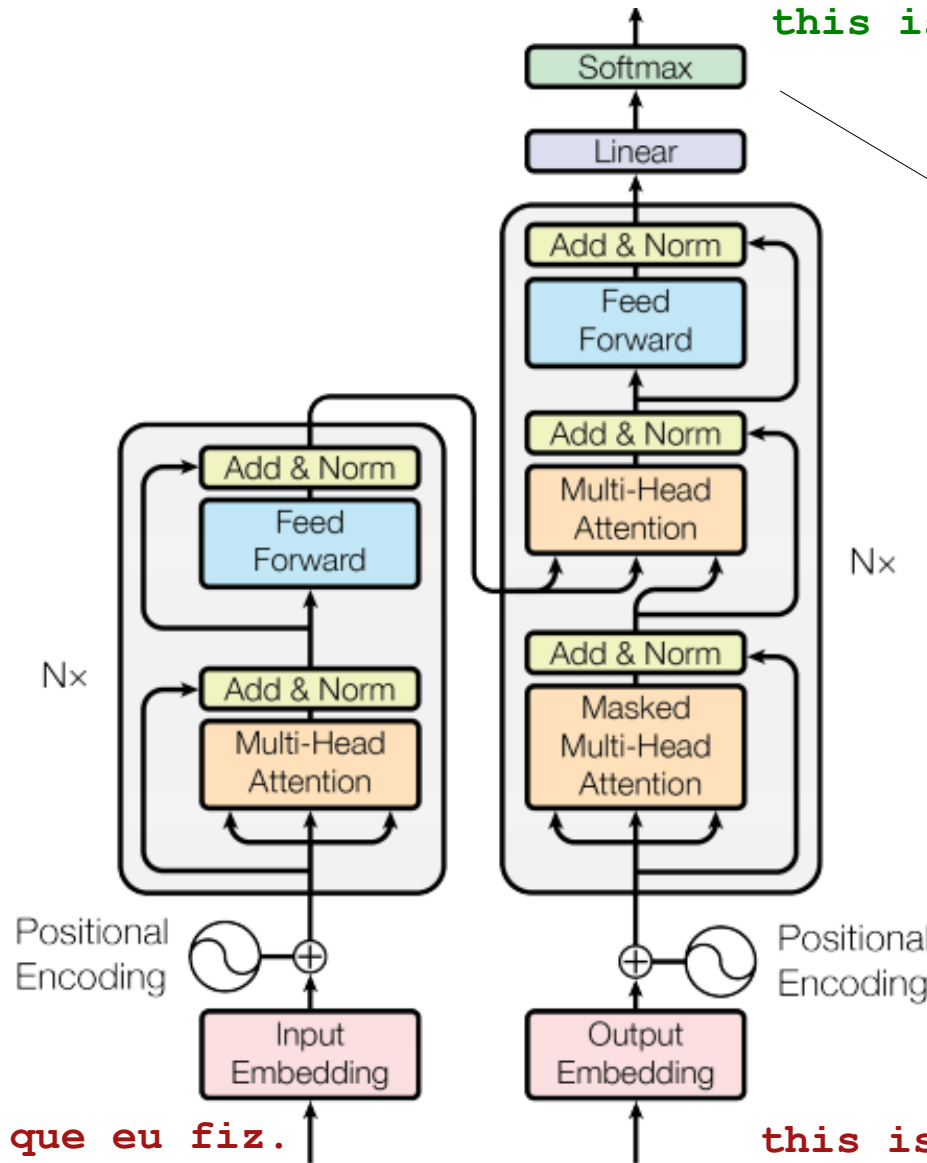
Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

*This one input
(in one data item)*

*This the true output
(teacher forcing)*

**`este é o primeiro livro que eu fiz.`**

**`this is the first book i've ever done`**

[image from https://arxiv.org/pdf/1706.03762.pdf]

# Attention Maps

Multiple heads, mixing layer ("Mixer"), topmost decoder block

# Attention in Blocks

1. Encoder (all layers)
   *input tokens over input tokens* (self-attention)


2. Decoder (lower layers)
   *output tokens over output tokens* (self-attention)


3. Decoder (top layer, "Mixer")
   *input tokens over output tokens* (cross-attention)
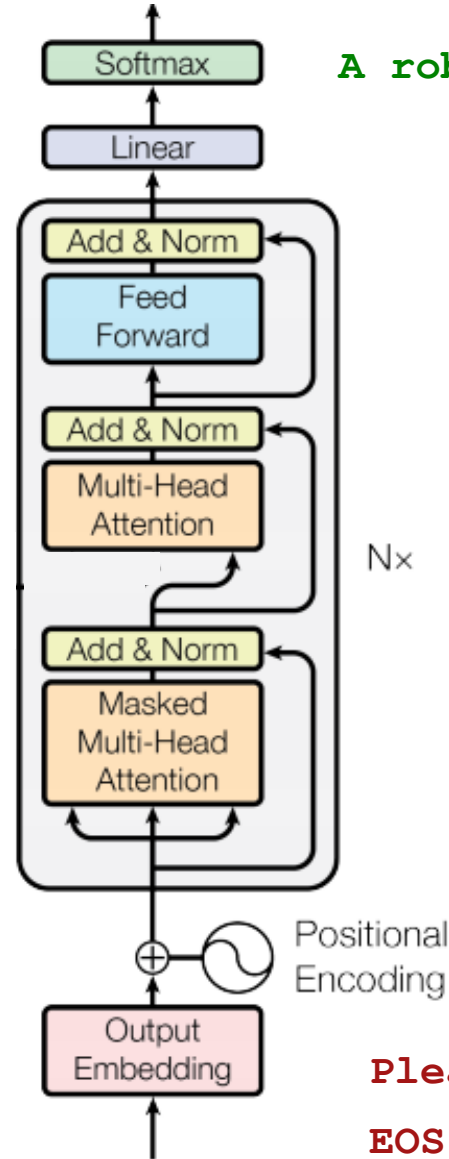
# Chat (Decoder-only)



**A robot may not injure a human being**

When in chat mode, the transformer network only performs
<u>text completions</u>

**Please recite the first law of robotics**

**EOS A robot may not injure a human**

[image from https://arxiv.org/pdf/1706.03762.pdf]

# GPT Secret Sauce: Multi-Task Learning (MTL)

**Incorporating input and task description into the prompt**

(translate from language A to B, <text in language A>, <text in language B>)

task           input           label

(choose an answer to a question, <context>, <question>, <answer$_1$>, <answer$_2$>, ..., <answer$_n$>, <correct answer>)

task                           input                   label

"I'm not the cleverest man in the world, but like they say in French: Je ne suis pas un imbecile [I'm not a fool].

phrase in English          language indicator          phrase in French

**Training the system to perform multiple different tasks**

language indicator

If this sounds like a bit of a stretch, consider this question in French:
As-tu aller au cinéma?, or Did you go to the movies?, which literally translates as Have-you to go to movies/theater?

phrase in French          phrase in English

# ChatGPT 4o: So clever even when hallucinates..

> Who was the first elephant to swim across the Strait of Messina?

The first elephant to swim across the Strait of Messina was a famous circus elephant named Tusko. This remarkable feat occurred in 1979. Tusko, weighing over three tons, made the swim from the Italian mainland to the island of Sicily, covering a distance of about three kilometers. This event was widely publicized and showcased the impressive abilities of elephants, demonstrating their strength and endurance.

> Could you provide more details?

I apologize for the confusion earlier. It seems there was a misunderstanding. There isn't a widely known or documented case of an elephant named Tusko or any other elephant swimming across the Strait of Messina.

In historical context, the notion of elephants in relation to the Strait of Messina often brings to mind Hannibal's famous crossing of the Alps with elephants during the Second Punic War. However, these elephants did not swim across the Strait of Messina.

[Conversation made by myself with ChatGPT 4o, on May 17, 2024]

# Links

https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
Language Models are Unsupervised Multitask Learners

https://arxiv.org/pdf/1904.02679
Visualizing Attention in Transformer-Based Language Representation Models

https://arxiv.org/pdf/2203.02155
Training language models to follow instructions with human feedback

https://medium.com/@row3no6/why-chatgpt-uses-decoder-only-eaf0223143e6
Why ChatGPT Uses Decoder-Only

https://cameronrwolfe.substack.com/p/decoder-only-transformers-the-workhorse
Decoder-Only Transformers: The Workhorse of Generative LLMs