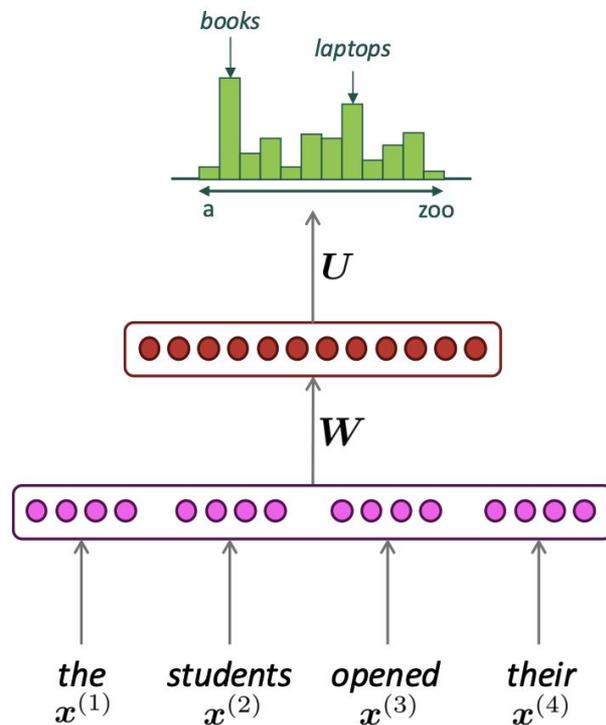# Deep Learning & Generative AI in Healthcare

Session 08

# Neural Network Model of Language

- A neural probabilistic language model (Y. Bengio, et al.)
- Fixed window is small
- No window is large enough

We need a neural architecture that can process *any length input*



books

laptops

a                                    zoo

$U$

$W$

the          students        opened         their
$x^{(1)}$    $x^{(2)}$       $x^{(3)}$      $x^{(4)}$

# Recurrent Neural Networks

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

*books*

*laptops*

**output distribution**

$$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

a                     zoo

$U$

$h^{(0)}$      $h^{(1)}$      $h^{(2)}$      $h^{(3)}$      $h^{(4)}$

**hidden states**

$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

$W_h$      $W_h$      $W_h$      $W_h$

$W_e$      $W_e$      $W_e$      $W_e$

**word embeddings**

$$e^{(t)} = Ex^{(t)}$$

$e^{(1)}$      $e^{(2)}$      $e^{(3)}$      $e^{(4)}$

$E$      $E$      $E$      $E$

**words / one-hot vectors**

$$x^{(t)} \in \mathbb{R}^{|V|}$$

*the*      *students*      *opened*      *their*

$x^{(1)}$      $x^{(2)}$      $x^{(3)}$      $x^{(4)}$

# Recurrent Neural Networks

- Advantages
  - The process any length!
  - Can use information from previous steps
  - Model size does not increase for longer input context
  - Same weights applied on every timestep
- Disadvantages
  - Slow
  - Difficult to access information from many steps back

$$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$$

*books*

*laptops*

a                           zoo

$U$

$h^{(0)}$        $h^{(1)}$        $h^{(2)}$        $h^{(3)}$        $h^{(4)}$

$W_h$        $W_h$        $W_h$        $W_h$

$W_e$        $W_e$        $W_e$        $W_e$

$e^{(1)}$        $e^{(2)}$        $e^{(3)}$        $e^{(4)}$

$E$        $E$        $E$        $E$

*the*        *students*        *opened*        *their*
$x^{(1)}$        $x^{(2)}$        $x^{(3)}$        $x^{(4)}$

# Attention is a solution!

- Attention provides a solution to the bottleneck problem!
- Core idea: on each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sequence!
- In attention, the query matches all keys softly, to a weight between 0 and 1. The key's values are multiplied by the weights and summed!
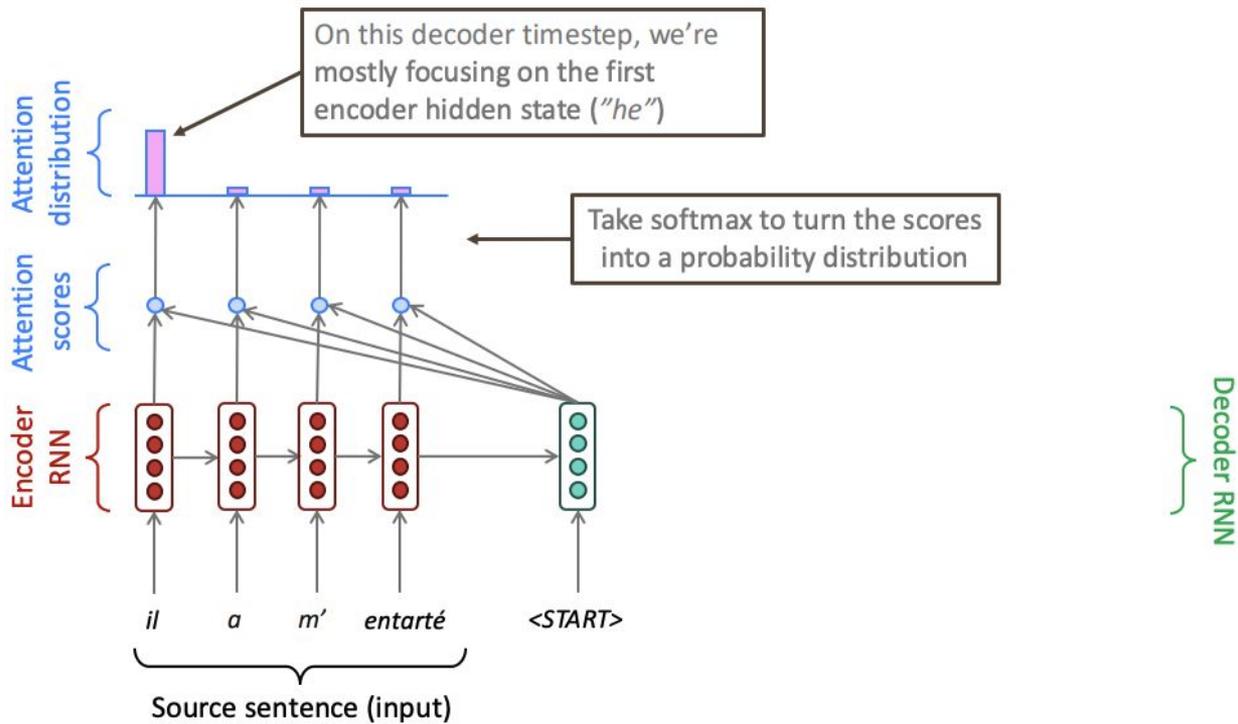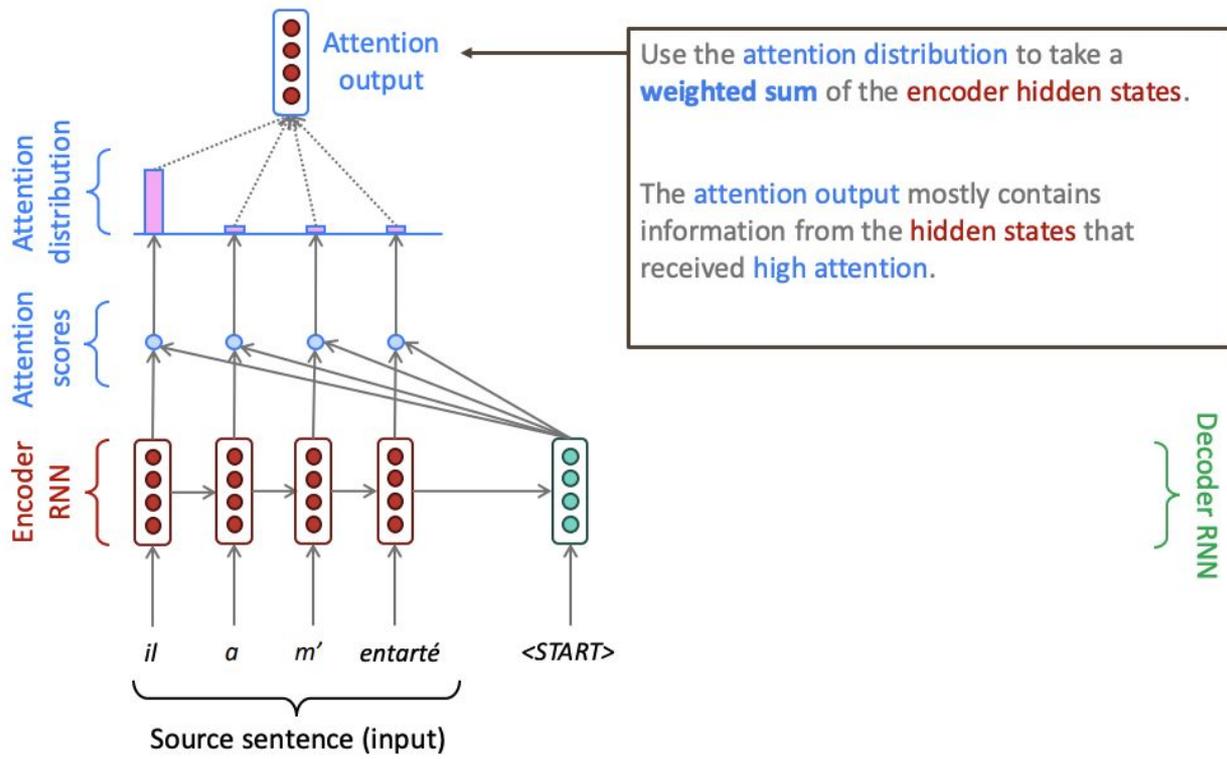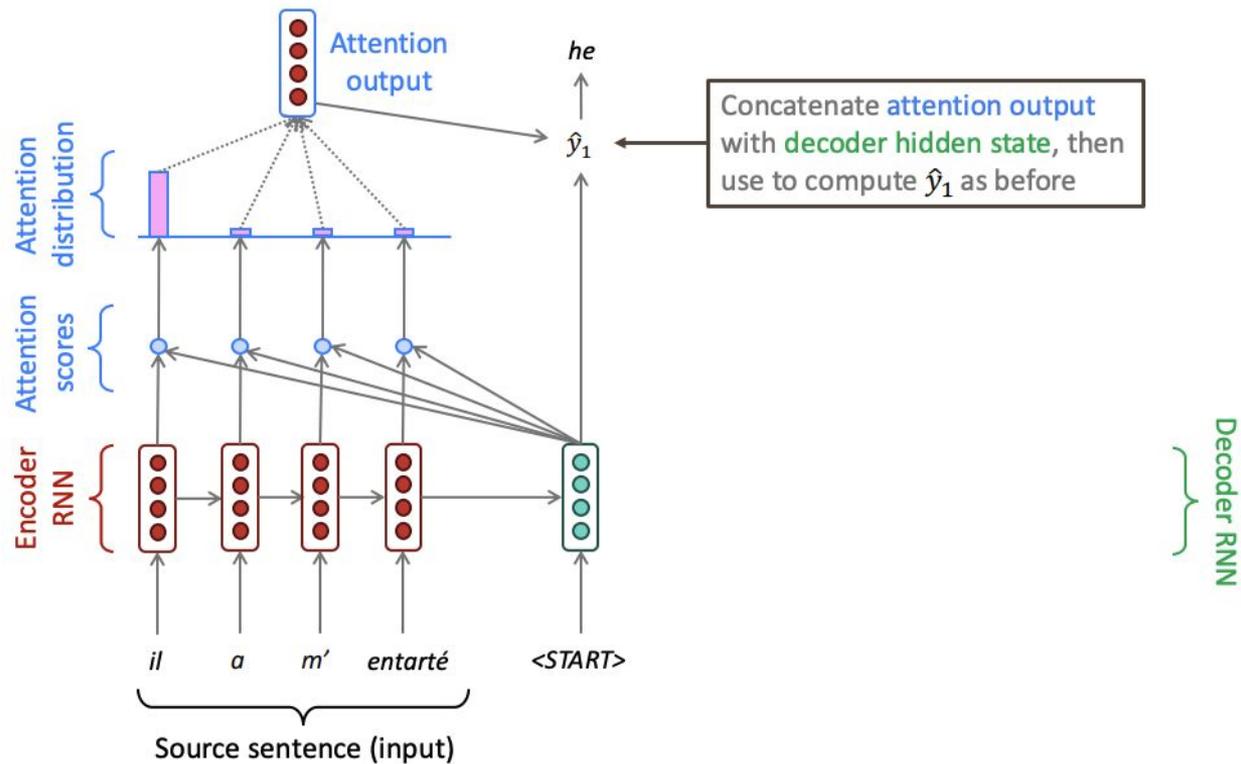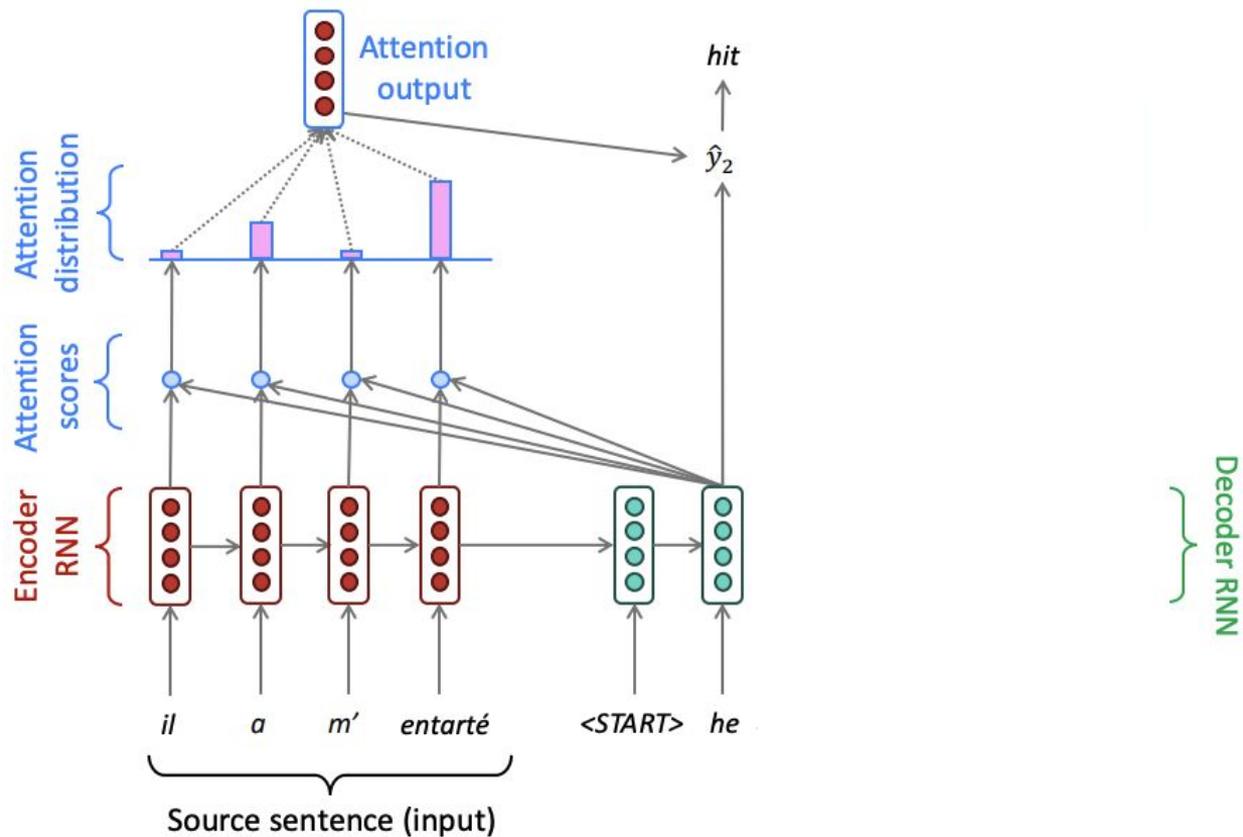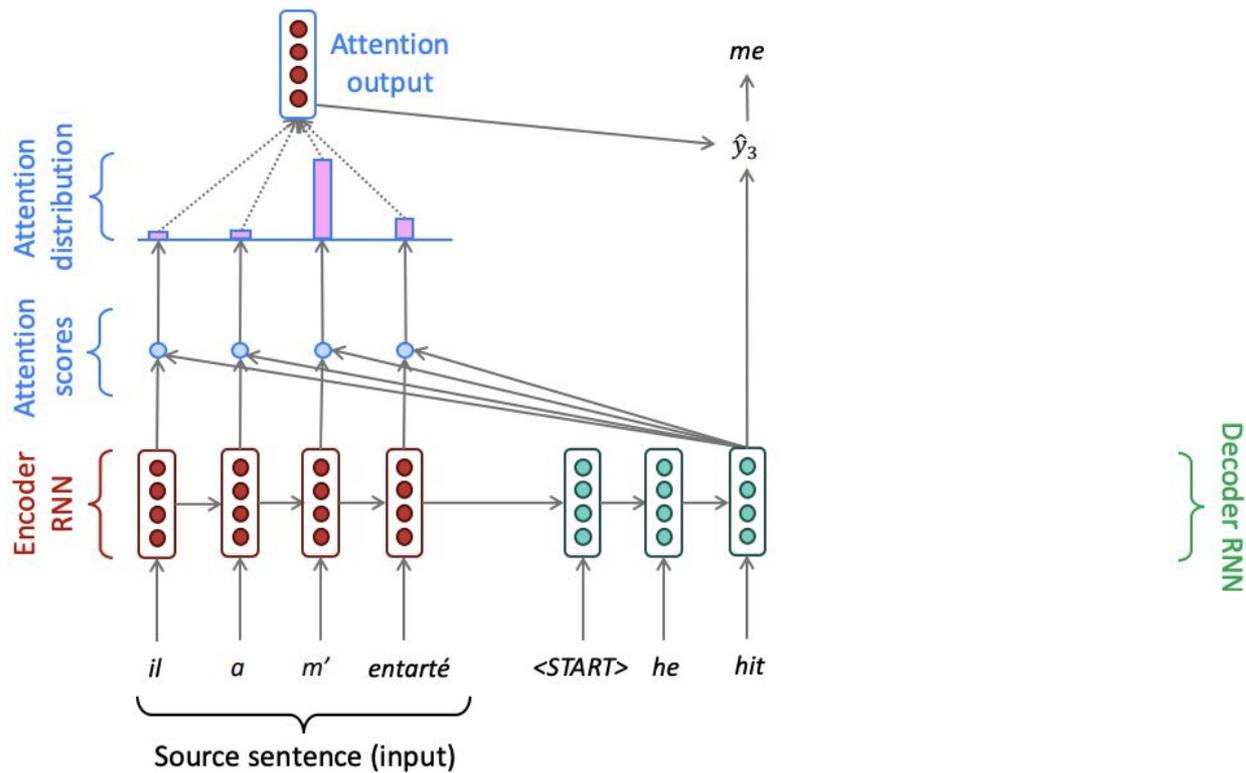
# Attention in RNNs

- On each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sequence.

# Attention in RNNs

# Attention in RNNs



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

# Attention in RNNs

# Attention in RNNs

# Attention in RNNs

# Is Recurrent Necessary at All?

- Abstractly: Attention is a way to pass information from a sequence ($x$) to a neural network input. ($h_t$)
- This is also exactly what RNNs are used for – to pass information!
- Can we just get rid of the RNN entirely? Maybe attention is just a better way to pass information!
- The building block we need is self Attention!
- So far we saw cross-attention!

# Attention in RNNs

# Self-attention: Keys, Queries, and Values

Let $\boldsymbol{w}_{1:n}$ be a sequence of words in vocabulary $V$, like *Zuko made his uncle tea*.

For each $\boldsymbol{w}_i$, let $\boldsymbol{x}_i = E\boldsymbol{w_i}$, where $E \in \mathbb{R}^{d \times |V|}$ is an embedding matrix.

1. Transform each word embedding with weight matrices $Q, K, V$, each in $\mathbb{R}^{d \times d}$

$$\boldsymbol{q}_i = Q\boldsymbol{x_i} \text{ (queries)} \qquad \boldsymbol{k}_i = K\boldsymbol{x_i} \text{ (keys)} \qquad \boldsymbol{v}_i = V\boldsymbol{x_i} \text{ (values)}$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

$$\boldsymbol{e}_{ij} = \boldsymbol{q}_i^\top \boldsymbol{k_j} \qquad \boldsymbol{\alpha}_{ij} = \frac{\exp(\boldsymbol{e}_{ij})}{\sum_{j'} \exp(\boldsymbol{e}_{ij'})}$$

3. Compute output for each word as weighted sum of values

$$\boldsymbol{o}_i = \sum_j \boldsymbol{\alpha}_{ij} \, \boldsymbol{v}_i$$

# Positional Embedding in Self-Attention

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values
- Consider representing each sequence index as a vector

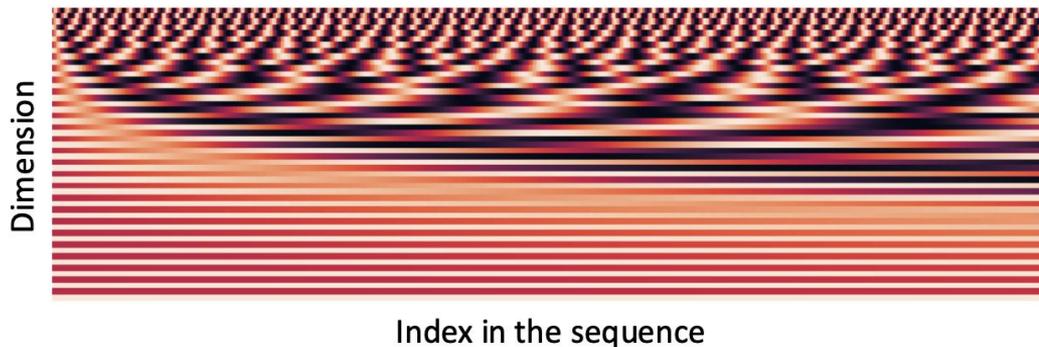$$p_i \in \mathbb{R}^d, \text{ for } i \in \{1, 2, \ldots, n\} \text{ are position vectors}$$

$$\widetilde{x}_i = x_i + p_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...

# Sinusoida Positional Embedding

- Sinusoidal position representations: concatenate sinusoidal functions of varying periods
- Periodicity indicates that maybe "absolute position" isn't as important
- It can extrapolate to longer sequences as periods restart!

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



Dimension

Index in the sequence

# Non-Linearity in Self-Attention

- Easy fix: add a feed-forward network to post-process each output vector.

$$m_i = MLP(\text{output}_i)$$
$$= W_2 * \text{ReLU}(W_1 \text{ output}_i + b_1) + b_2$$

# Causal Masking in Self-Attention

- For causality, we need to ensure not to peek at the future.
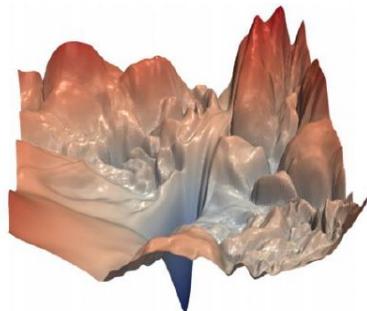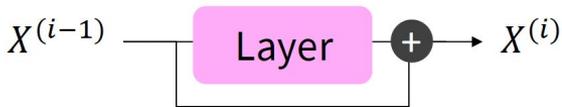- At each timestep, we could change the set of keys and queries to only include past words!

We can look at these (not greyed out) words



$$\begin{cases} q_i^\top k_j, j \le i \\ -\infty, j > i \end{cases}$$

For encoding these words

# Multi-Head Self-Attention Layer

- The Attention module splits its Query, Key, and Value parameters N-ways and passes each split independently through a separate head.
- Calculations are combined together to produce a final attention score.
- Greater power to encode multiple relationships and nuances for each word.

# Residual Connections

- A trick to help models learn better!
- Gradient is 1 through residual connection
- Bias toward identity function.

$X^{(i-1)}$ —— [ Layer ] —— $X^{(i)}$

$X^{(i-1)}$ —— [ Layer ] + —— $X^{(i)}$

[no residuals]   [residuals]

[Loss landscape visualization, Li et al., 2018, on a ResNet]

# Layer Normalization

- A trick to help models train faster.
- Cut down on uninformative variation in hidden vectors by normalizing to unit mean and standard deviation within each layer.

$$\text{output} = \frac{x - \mu}{\sqrt{\sigma} + \epsilon} * \gamma + \beta$$

Normalize by scalar mean and variance

Modulate by learned elementwise gain and bias

# Transformer Encoder

- Position representation
  - Specify the sequence order, since self-attention is an unordered function of its inputs.
- Nonlinearities
  - Frequently implemented as a simple feedforward network.
- Masking
  - Keep information about the future from "leaking" to the past.
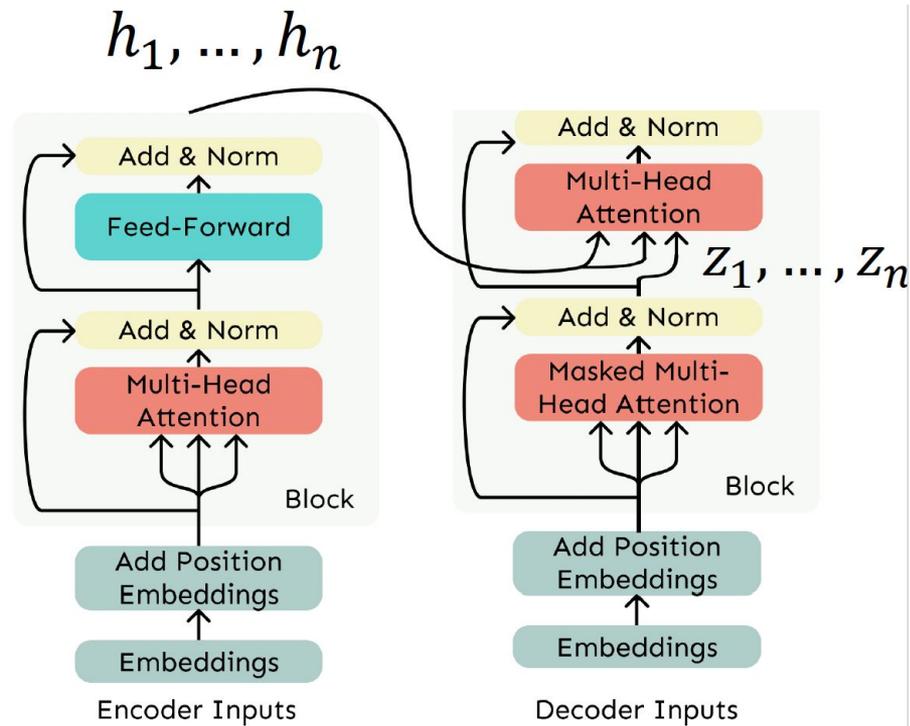
# Transformer Decoder



person wearing hat [END]

$y_0$ $y_1$ $y_2$ $y_3$

x N

Transformer decoder

$c_{0,0}$ $c_{0,1}$ $c_{0,2}$ $c_{2,2}$

$y_0$ $y_1$ $y_2$ $y_3$

[START] person wearing hat

$y_0$ $y_1$ $y_2$ $y_3$

FC

Layer norm

+

MLP

Layer norm

+

Multi-head **attention**
k v q

$c_{0,0}$ $c_{0,1}$ $c_{0,2}$ $c_{2,2}$

Layer norm

+

**Masked** Multi-head self-attention

Positional encoding

$x_0$ $x_1$ $x_2$ $x_3$

Multi-head attention block attends over the transformer encoder outputs.

For image captions, this is how we inject image features into the decoder.

Vaswani et al, "Attention is all you need", NeurIPS 2017

# Cross Attention

- Self-attention:
  - Keys, queries, and values from same source
- Cross-attention
  - The keys and values are from encoder (like a memory)
  - The queries are from the decoder

Inside an LLM

# Vision-Language Models

## Prompt

You are an expert radiologist. Please succinctly describe the findings for the above chest x-ray.



### MedGemma

The chest X-ray shows a large left pneumothorax with a chest tube in place. There is also a possible nodule in the right upper lobe. The mediastinum is midline. The heart size is normal. There is no evidence of pulmonary edema.

**Reference Information: Radiologist Impression**
Moderately large left hydropneumothorax (predominantly air-containing) with incompletely expanded left lung. If persistent, consider ex vacuo pneumothorax/trapped lung condition. Multiple right lung nodules, highly suspicious for metastases.

## Prompt

Describe this image and provide the most likely condition. Keep your answer brief



### MedGemma

The chest X-ray shows a significant opacity in the left lower lung field, obscuring the costophrenic angle. This suggests a possible pleural effusion or consolidation. The right lung appears relatively clear. The heart size is within normal limits. The mediastinum is unremarkable.

**Reference Information: Radiologist Impression**
Lower left lung consolidation. Small-medium left pleural effusion. Similar though lesser findings right side.