

Frontiers in Generative AI for Medical Imaging and Healthcare

Session 3
Retrieval-Augmented Generation & Clinical NLP
Sep 20 2025

Recurrent Neural Networks

output distribution

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

hidden states

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

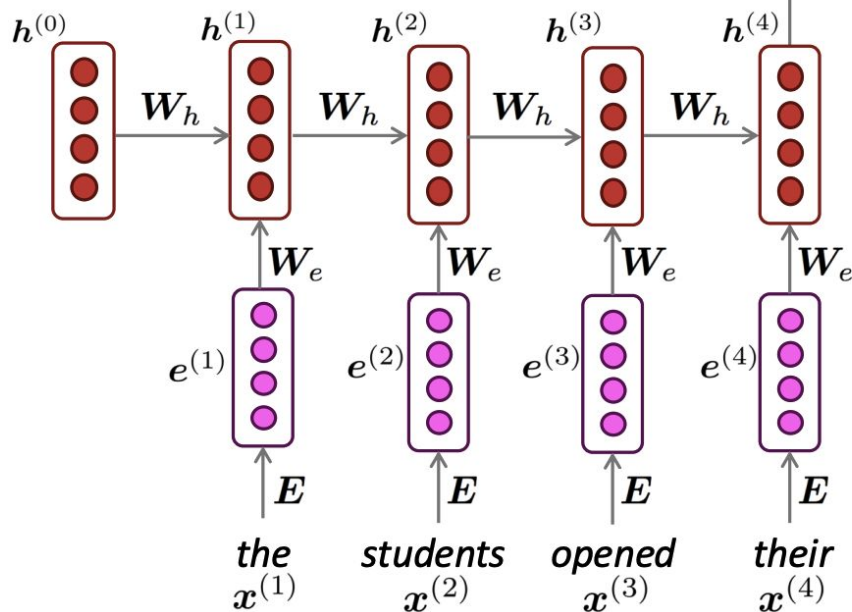
$h^{(0)}$ is the initial hidden state

word embeddings

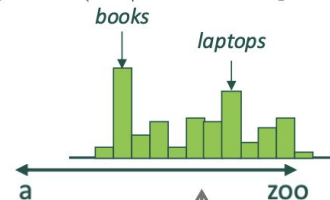
$$e^{(t)} = Ex^{(t)}$$

words / one-hot vectors

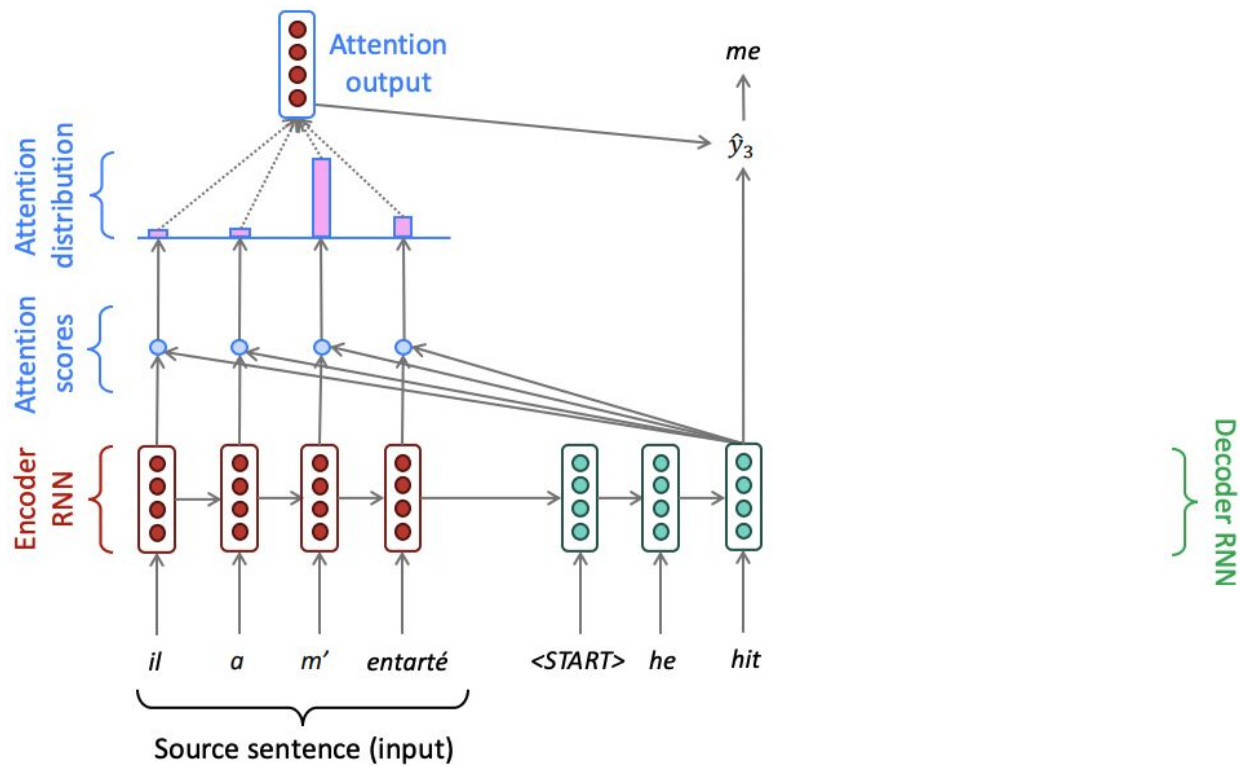
$$x^{(t)} \in \mathbb{R}^{|V|}$$



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



Attention in RNNs



Self-attention: Keys, Queries, and Values

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

For each w_i , let $x_i = Ew_i$, 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 $\mathbb{R}^{d \times d}$

$$\textcolor{brown}{q}_i = Qx_i \text{ (queries)} \quad \textcolor{teal}{k}_i = Kx_i \text{ (keys)} \quad \textcolor{teal}{v}_i = Vx_i \text{ (values)}$$

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

$$e_{ij} = \textcolor{brown}{q}_i^\top \textcolor{teal}{k}_j \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

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

$$o_i = \sum_j \alpha_{ij} \textcolor{teal}{v}_j$$

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

$\mathbf{p}_i \in \mathbb{R}^d$, for $i \in \{1, 2, \dots, n\}$ are position vectors

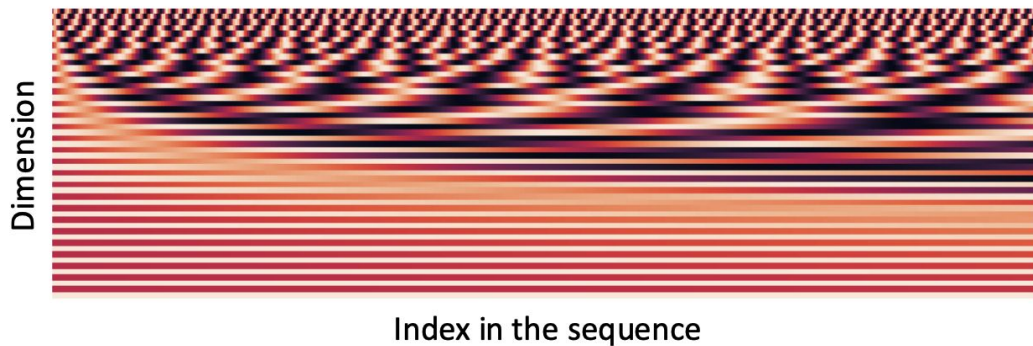
$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \mathbf{p}_i$$

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

Sinusoidal 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!

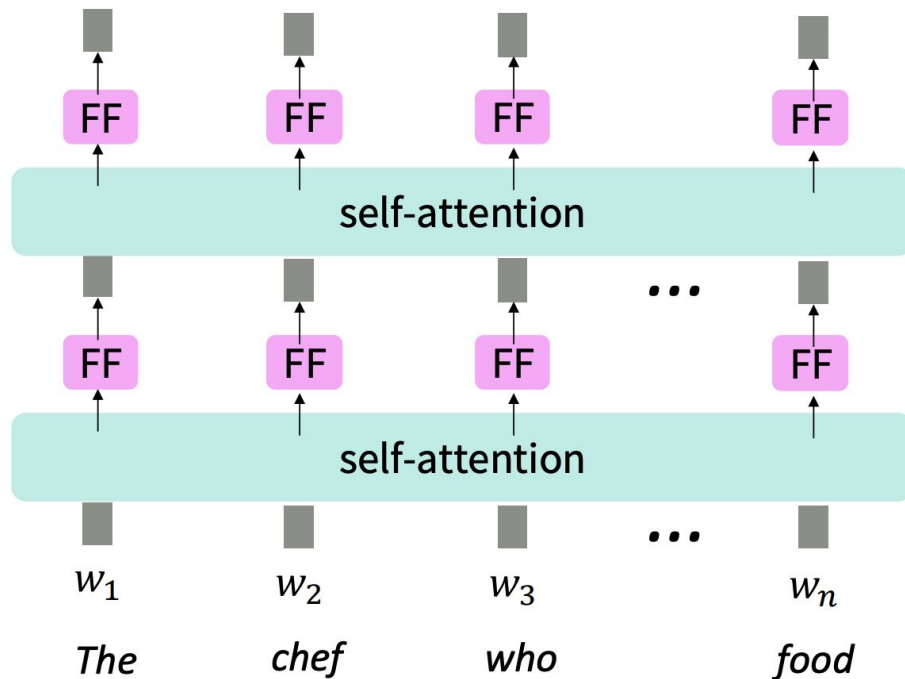
$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*d/2/d}) \\ \cos(i/10000^{2*d/2/d}) \end{pmatrix}$$



Non-Linearity in Self-Attention

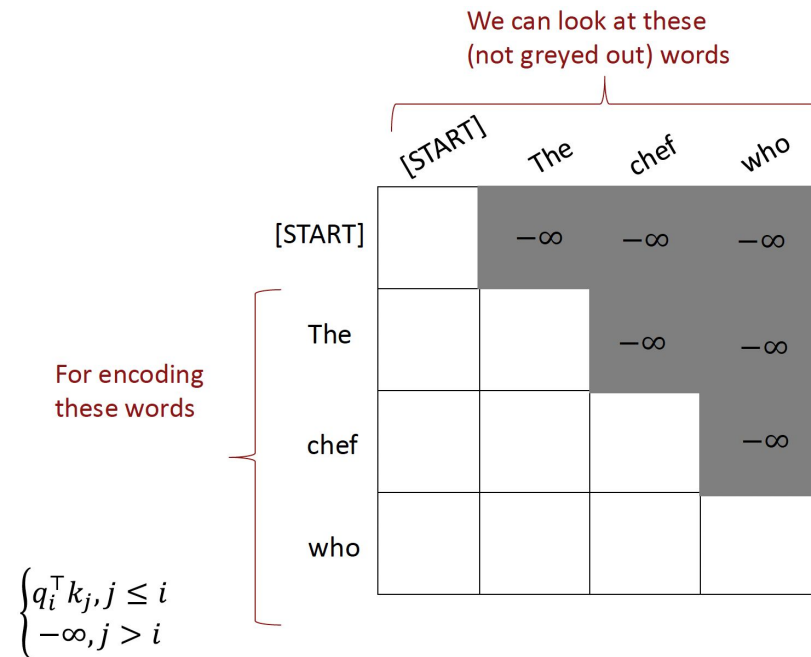
- Easy fix: add a **feed-forward network** to post-process each output vector.

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



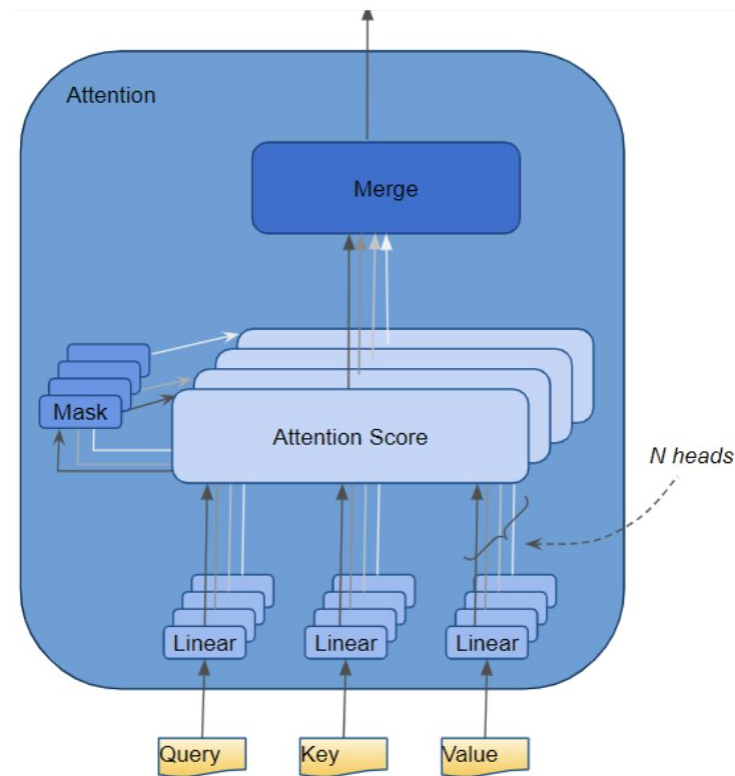
Causal Masking in Self-Attention

- 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!**



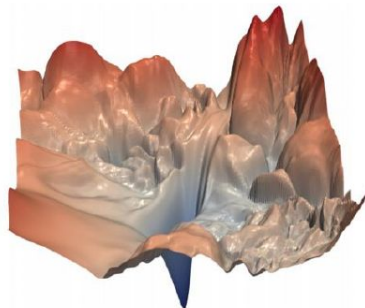
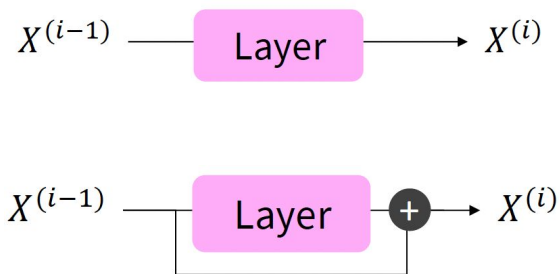
Multi-Head Self-Attention Layer

- 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.

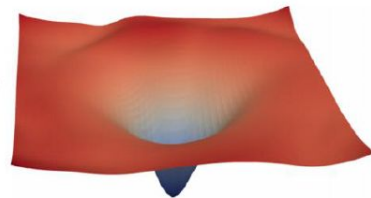


Residual Connections

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



[no residuals]



[residuals]

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

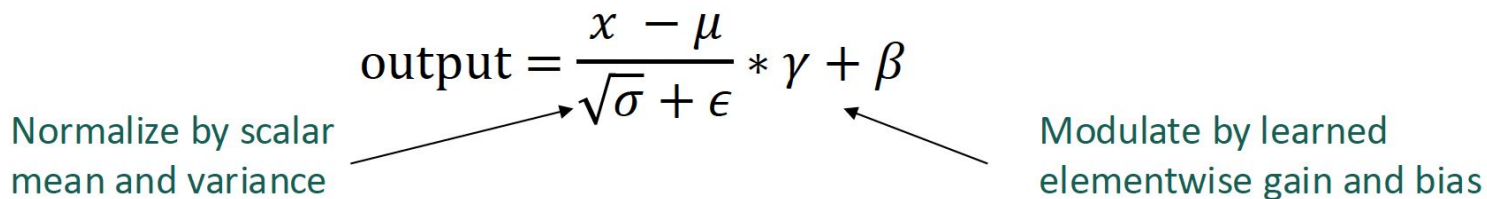
Layer Normalization

- 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.

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

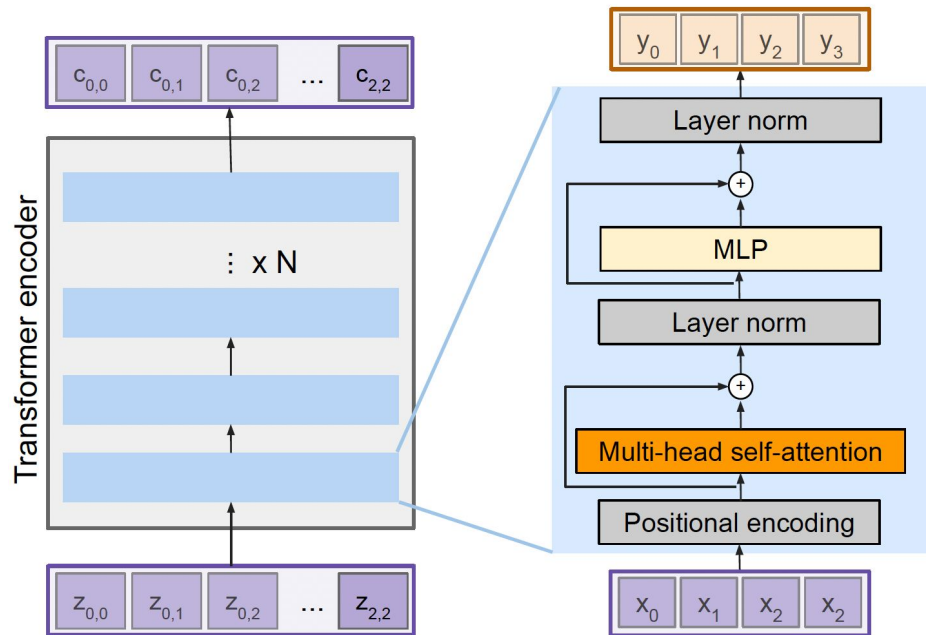
Normalize by scalar mean and variance

Modulate by learned elementwise gain and bias

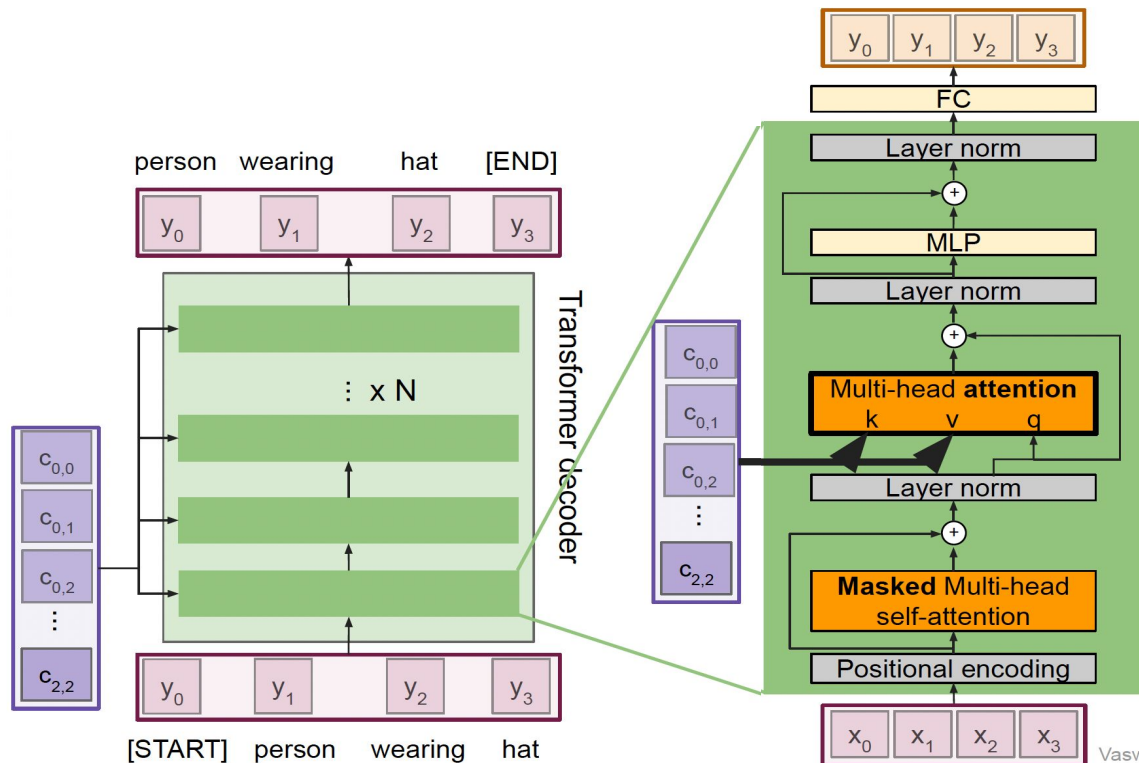
The diagram shows the Layer Normalization formula: output = (x - μ) / (√σ + ε) * γ + β. Two arrows point from descriptive text to parts of the formula. One arrow points from 'Normalize by scalar mean and variance' to the denominator (√σ + ε). Another arrow points from 'Modulate by learned elementwise gain and bias' to the multiplication by γ and the addition of β.

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



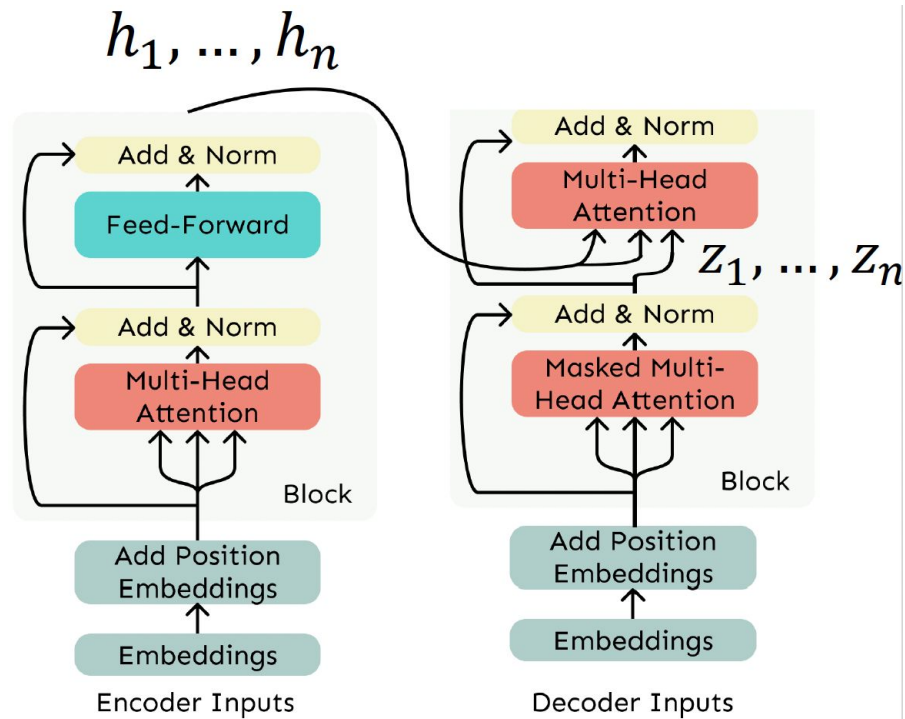
Multi-head attention block attends over the transformer encoder outputs.

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