

Transformers

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

June 11, 2025

1. Motivation
2. Learning Outcomes
3. General Attention Layer
4. Self-Attention Layer
5. Positional Encoding
6. Masked Self-Attention Layer
7. Multi-Head Attention
8. General Attention Versus Self-Attention
9. CNN with Self-Attention
10. Transformer Architecture
11. Summary

► Limited Context:

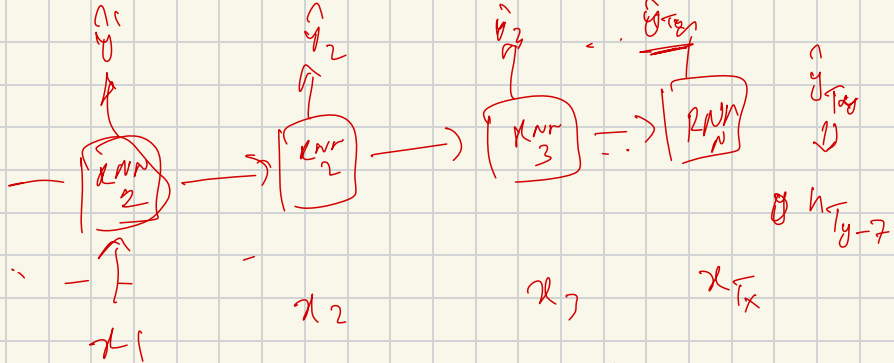
- RNNs process sequences step-by-step → Difficult to capture long-range dependencies.
- CNNs are effective for local patterns → Struggle with global context in sequences.

► Sequential Bottlenecks:

- RNNs process data sequentially.
- Prevents efficient parallelization.
- Leads to slow training and inference.

► Slow Training:

- RNNs and CNNs are computationally intensive.
- Especially challenging for long sequences or large images.
- Results in increased training times.



Vanilla RNNs \rightarrow 7 cells

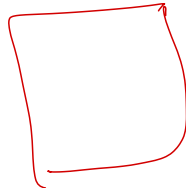
The¹ boy², who³ was⁴ playing⁵ football,⁶
 also⁷ loved⁸ to⁹ watch¹⁰ cricket.

There is a need for models that:

- ▶ Capture comprehensive, global context efficiently.
- ▶ Allow for parallel computation.
- ▶ Accelerate training and inference.

Transformers address these needs by:

- ▶ Leveraging self-attention mechanisms.
- ▶ Modeling dependencies across entire sequences.
- ▶ Enabling highly parallelizable computation.

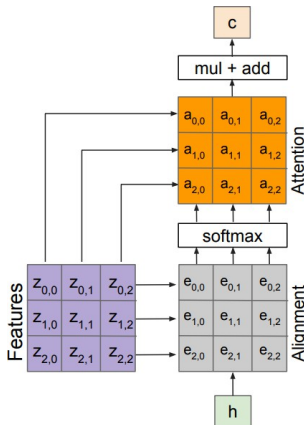


After this lecture, you will be able to:

- ▶ Explain the motivation behind the development of Transformers.
- ▶ Describe the architecture and key components of Transformer models.
- ▶ Understand and implement attention and self-attention mechanisms.
- ▶ Explain the role of positional encoding in Transformers.
- ▶ Distinguish between general attention and self-attention.
- ▶ Understand the concept and benefits of multi-head attention.
- ▶ Compare CNNs with and without self-attention.
- ▶ Summarize the advantages of Transformers over traditional sequence models.

Attention we just saw in image captioning

- Previously, we saw how attention helps image captioning models focus on important regions, generating more relevant captions.
- However, we noticed limitations—attention sometimes misses subtle details or struggles with complex scenes.



Outputs:

context vector: **c** (shape: D)

Operations:

Alignment: $e_{i,j} = f_{\text{att}}(h, z_{i,j})$

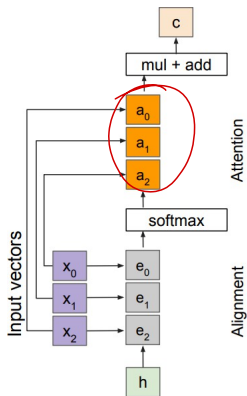
Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$

Output: $\mathbf{c} = \sum_{i,j} a_{i,j} z_{i,j}$

Inputs:

Features: **z** (shape: H x W x D)

Query: **h** (shape: D)

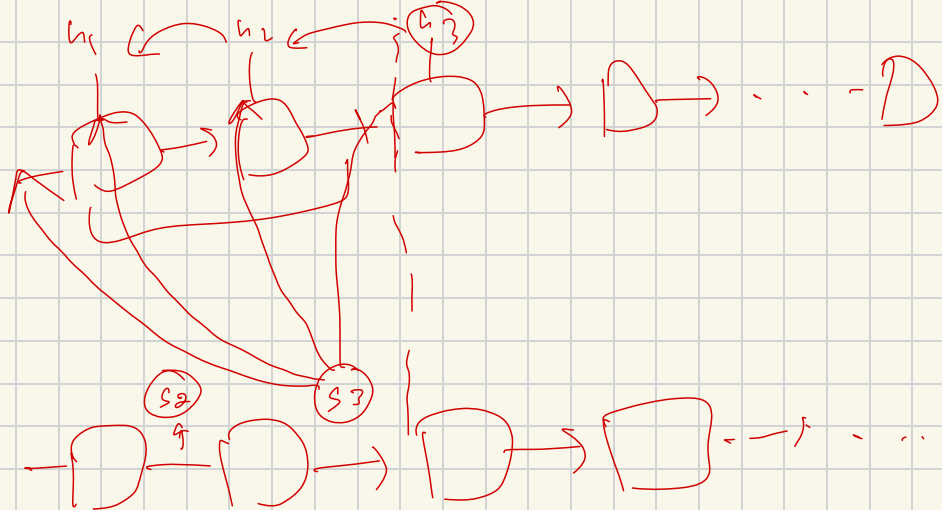


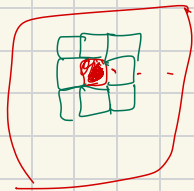
Outputs:
context vector: c (shape: D)

Operations:
Alignment: $e_i = f_{\text{att}}(h, x_i)$
Attention: $a = \text{softmax}(e)$
Output: $c = \sum_i a_i x_i$

Inputs:
Input vectors: x (shape: $N \times D$)
Query: h (shape: D)

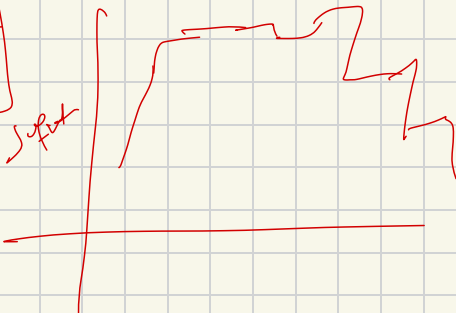
- The **attention** operation is permutation invariant.
- It does not depend on the ordering of features.
- Reshape the input of size $H \times W = N$ into N feature vectors.



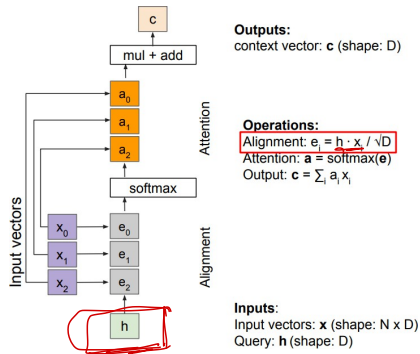


Attention \rightarrow high

0.9
0.001
0.001

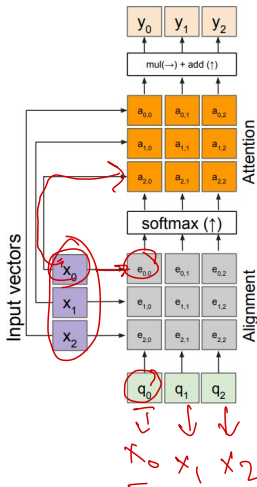


General Attention Layer (cont.)



- ▶ The original attention mechanism $f_{att}(\cdot)$ uses a simple dot product to compute similarity between queries and keys.
- ▶ As the feature dimension D increases \rightarrow the dot product sum involves more terms \rightarrow larger variance in the resulting logits.
- ▶ Larger magnitude vectors produce higher logits \rightarrow causing the softmax output to become more peaked (lower entropy, assuming logits are IID).
- ▶ This means the **model may focus too narrowly**, assigning very **high attention to a few positions** and **almost none to others**.
- ▶ To counteract this \rightarrow we scale the dot product by dividing by \sqrt{D} .
- ▶ This normalization keeps the variance of the logits more consistent, resulting in a softer, more balanced attention distribution.

General Attention Layer (cont.)



Outputs:

context vectors: \mathbf{y} (shape: D)

Operations:

Alignment: $e_{i,j} = q_i \cdot x_j / \sqrt{D}$

Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$

Output: $y_j = \sum_i a_{i,j} x_i$

Inputs:

Input vectors: \mathbf{x} (shape: $N \times D$)

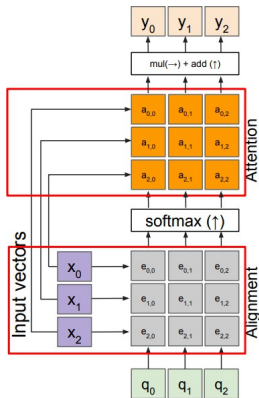
Queries: \mathbf{q} (shape: $M \times D$)



- We can use multiple **query vectors** in the attention mechanism.
- Each query attends to the input independently, producing its own output context vector.
- This allows the model to extract different types of information from the same input, enabling richer and more flexible representations.

$$\begin{aligned} \rightarrow e_1 &= (q_0 \cdot x_0) / \sqrt{D} \\ \rightarrow e_2 &= (q_1 \cdot x_1) / \sqrt{D} \end{aligned}$$

General Attention Layer (cont.)

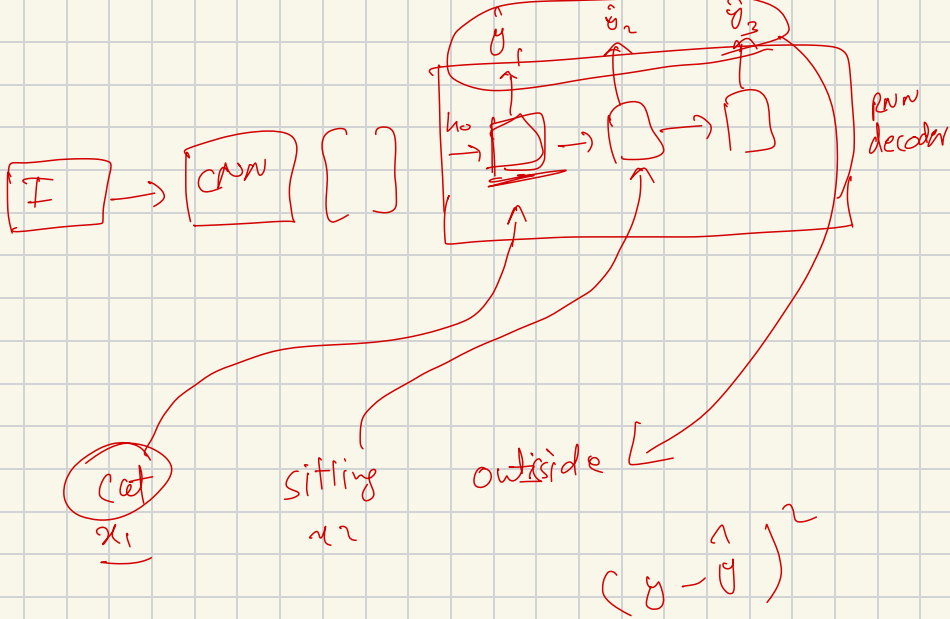


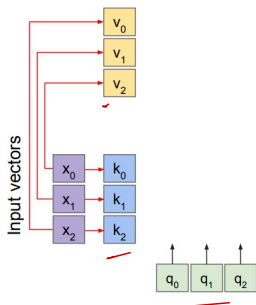
Outputs:
context vectors: \mathbf{y} (shape: D)

Operations:
Alignment: $e_{i,j} = \mathbf{q}_j \cdot \mathbf{x}_i / \sqrt{D}$
Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$
Output: $y_j = \sum_i a_{i,j} x_i$

Inputs:
Input vectors: \mathbf{x} (shape: N x D)
Queries: \mathbf{q} (shape: M x D)

- Observe that the same input vectors are used for both computing the alignment scores (queries and keys) and for generating the attention-weighted output (values).
- This dual use lets the model efficiently use the same features for both attending and aggregating information.





Operations:

Key vectors: $\mathbf{k} = \mathbf{x}\mathbf{W}_k$

Value vectors: $\mathbf{v} = \mathbf{x}\mathbf{W}_v$

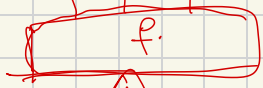
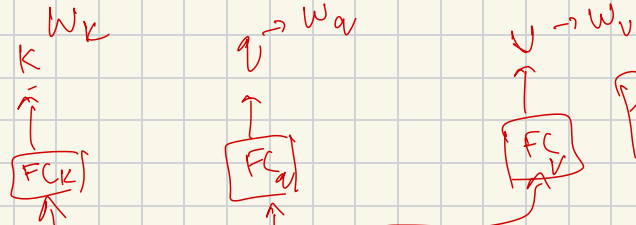
Inputs:

Input vectors: \mathbf{x} (shape: $N \times D$)

Queries: \mathbf{q} (shape: $M \times D_k$)

- ▶ Notice that the same input vectors are used for both computing alignment scores (as queries and keys) and for generating the attention-weighted output (as values).
- ▶ To increase the expressiveness of the attention layer, we can introduce separate fully connected (FC) layers before each step \rightarrow one for queries, one for keys, and one for values.
- ▶ This allows the model to learn different transformations for each role, enabling richer and more flexible representations.

Projections



1000
 $f_i = \text{pool}_{1000}(\text{protraine}(v_{qgr}) \cdot f_{in})$

128 $\leftarrow \textcircled{K} = K = v_{qgr} \cdot FC_1(f_i)$
 $\underline{K} = FC(f_i)$

CNN

FCV

nn.Linear (1000, 512)

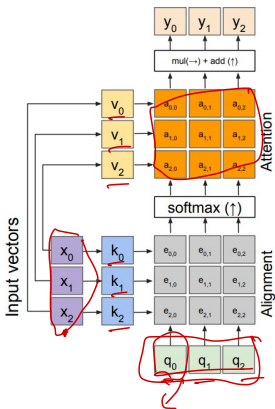
.nn - - (512, 256)

.nn.lin - (256, 128)

FC2

nn.Linear (1000, 128)

General Attention Layer (cont.)



Outputs:

context vectors: y (shape: D_v)

Operations:

Key vectors: $k = xW_k$

Value vectors: $v = xW_v$

Alignment: $e_{ij} = q_i \cdot k_j / \sqrt{D}$

Attention: $a = \text{softmax}(e)$

Output: $y_j = \sum_i a_{ij} v_i$

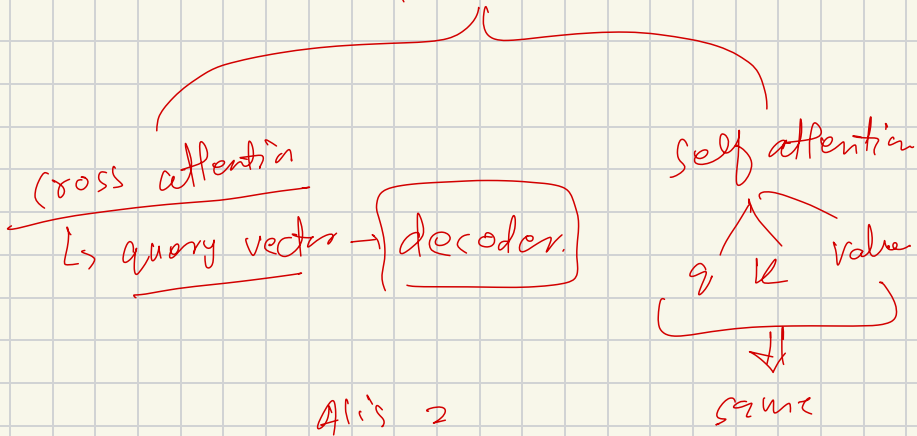
Inputs:

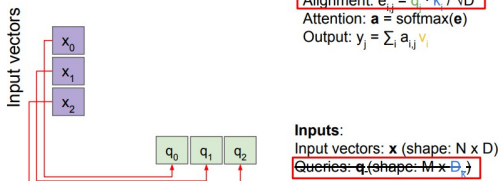
Input vectors: x (shape: $N \times D$)

Queries: q (shape: $M \times D_k$)

- ▶ The input vectors are used for both computing alignment scores (as queries and keys) and for generating the attention-weighted output (as values).
- ▶ To enhance the expressiveness of the attention layer, we introduce separate fully connected (FC) layers for queries, keys, and values.
- ▶ Each FC layer can learn a different transformation, allowing the model to capture more complex relationships.
- ▶ With this setup, the input and output dimensions can differ, depending on the transformations applied by the key and value FC layers.

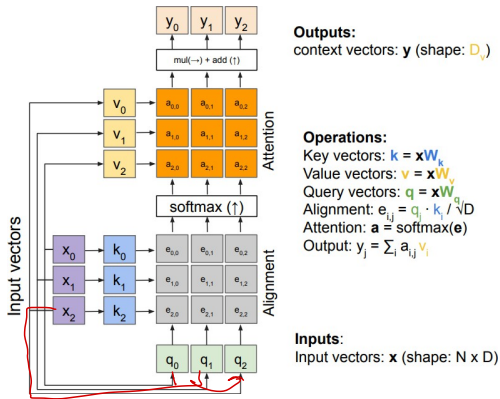
Attention





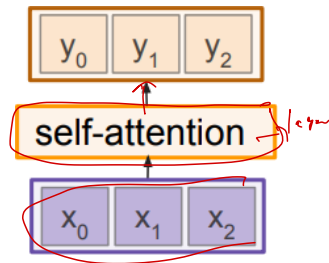
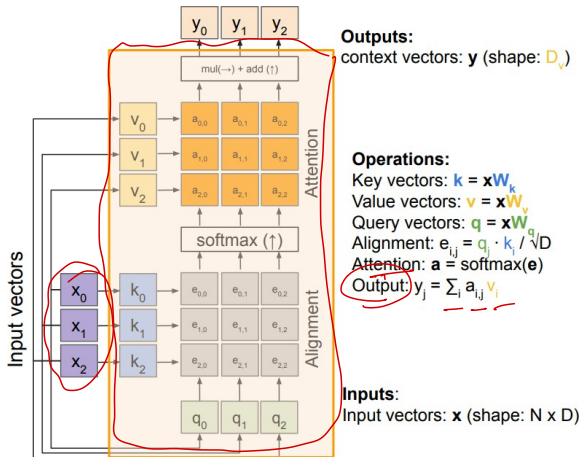
- ▶ Recall: the query vector is derived from the input vectors.
- ▶ In a self-attention layer, query, key, and value vectors are all computed from the same input.
- ▶ There are no separate input query vectors; instead, they are generated internally.
- ▶ Typically, fully connected (FC) layers are used to compute the query, key, and value vectors from the input.
- ▶ This allows each position in the input to attend to all other positions, enabling the model to capture contextual relationships.

Self Attention Layer (cont.)



- Recall: the query vector is derived from the input vectors.
- In a self-attention layer, query, key, and value vectors are all computed from the same input.
- There are no separate input query vectors; instead, they are generated internally.
- Typically, fully connected (FC) layers are used to compute the query, key, and value vectors from the input.
- This allows each position in the input to attend to all other positions, enabling the model to capture contextual relationships.

Self Attention Layer (cont.)



$$Q_w = K_w = V_w = \text{self.FC}(\underline{P_1})$$

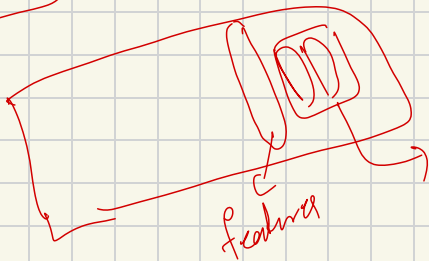
projection

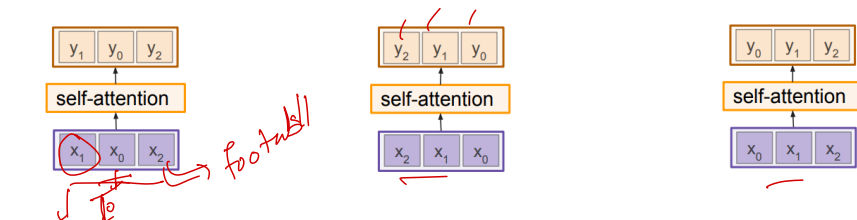
$$\text{Alignment} = Q \cdot K^T / \sqrt{D}$$

$$\text{Attention} = \text{Softmax}(\text{Alignment})$$

$$\boxed{\text{Output}} = (\text{Attention} \cdot V_w)$$

class Embed



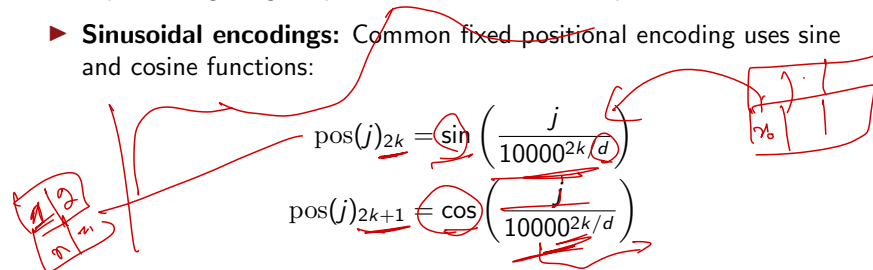


- ▶ The self-attention layer is **permutation equivariant**: it produces the same output regardless of the order of the input elements.
- ▶ This means the model does not inherently capture the order of the sequence.
- ▶ **Challenge:** For tasks involving ordered data, such as language or images, we need a way to encode positional information.
- ▶ *How can we enable the model to distinguish between different positions in a sequence?*

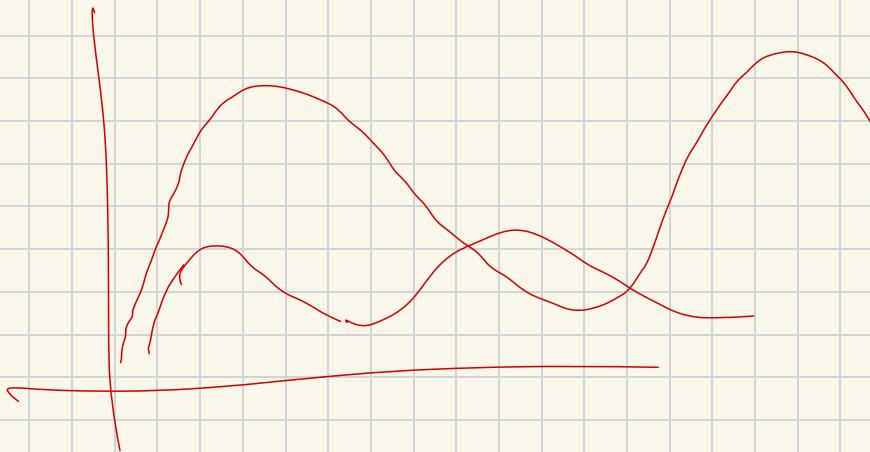
- Love
Football



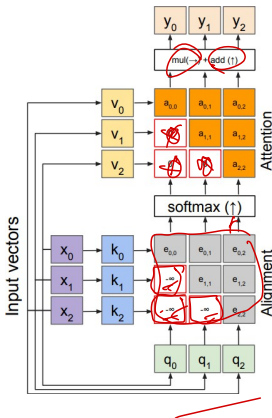
- ▶ The positional encoding is added to the input vectors before computing the attention scores.
- ▶ This allows the model to incorporate information about the position of each vector in the sequence.
- ▶ The positional encoding can be learned or fixed, depending on the implementation.
- ▶ The choice of positional encoding can affect the model's ability to capture long-range dependencies and relationships in the data.
- ▶ **Sinusoidal encodings:** Common fixed positional encoding uses sine and cosine functions:


$$\text{pos}(j)_{2k} = \sin\left(\frac{j}{10000^{2k/d}}\right)$$

$$\text{pos}(j)_{2k+1} = \cos\left(\frac{j}{10000^{2k/d}}\right)$$



Masked Self-Attention Layer



Outputs:
context vectors: y (shape: D_v)

Operations:
Key vectors: $k = xW_k$
Value vectors: $v = xW_v$
Query vectors: $q = xW_q$
Alignment: $e_{ij} = q_i \cdot k_j / \sqrt{D}$
Attention: $a = \text{softmax}(e)$
Output: $y_i = \sum_j a_{ij} v_j$

Inputs:
Input vectors: x (shape: $N \times D$)

$$\text{mask} = \begin{bmatrix} 1 & & & \\ -\infty & 1 & & \\ -\infty & -\infty & 1 & \\ -\infty & -\infty & -\infty & 1 \end{bmatrix}$$

- Prevents each position from attending to subsequent (future) positions.
- Achieved by setting (manually) alignment scores of future tokens to $-\infty$ before softmax.
- Ensures predictions for position i depend only on positions $\leq i$.

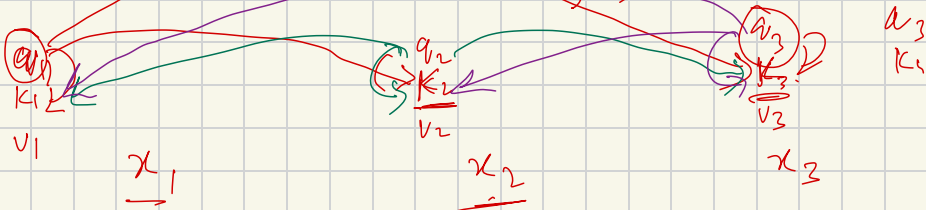


$I \rightarrow \text{love}$

$$a_1 \cdot k_1 (I \leftrightarrow I) =$$

$$\boxed{a_2 \cdot k_2} (\bar{I} \leftrightarrow \text{love}) = 1.00$$

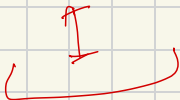
$$a_3 \cdot k_3 (I \leftrightarrow \text{football}) =$$



x_1

x_2

x_3



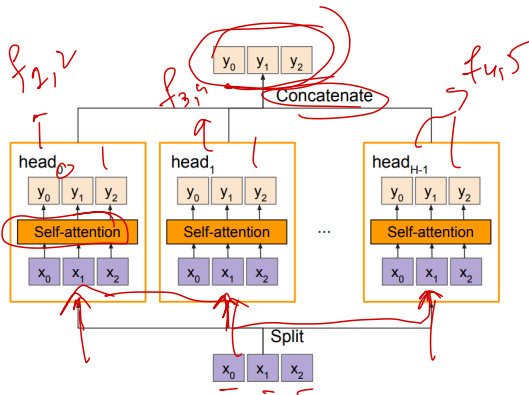
love

football

	$a_1 \cdot k_1$	$a_2 \cdot k_1$	$a_3 \cdot k_1$
\times	$a_1 \cdot k_2$	$a_2 \cdot k_2$	$a_3 \cdot k_2$
	$a_1 \cdot k_3$	$a_2 \cdot k_3$	$a_3 \cdot k_3$
\times			

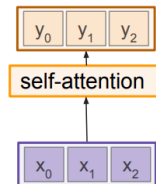
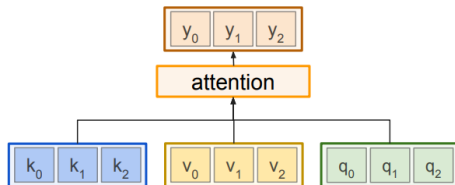
$a_3 k_5$
 $a_2 k_5$

Multi-Head Self-Attention Layer



- Instead of a single attention mechanism, use multiple “heads” operating in parallel.
- Each head attends to different parts or aspects of the input, capturing diverse relationships (e.g., subject-verb, coreferences).
- The outputs from all heads are concatenated and projected, resulting in a richer and more expressive context representation.
- This modular approach allows the model to learn various types of dependencies simultaneously, improving overall performance.

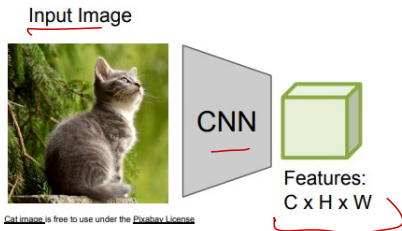
General-Attention Vs Self-Attention



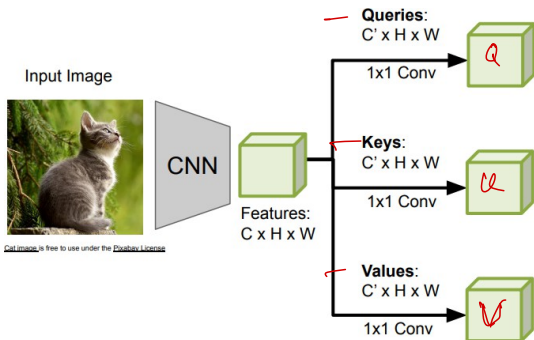
Comparison	General Attention	Self-Attention
Q, K, V origins	From separate source & target	From same input sequence
Use case	Encoder \rightarrow Decoder cross-attention	Encoder/Decoder internal relation info
Information flow	<u>Across representations</u>	Within single representation <u> </u>

Table 1: General Attention vs. Self-Attention

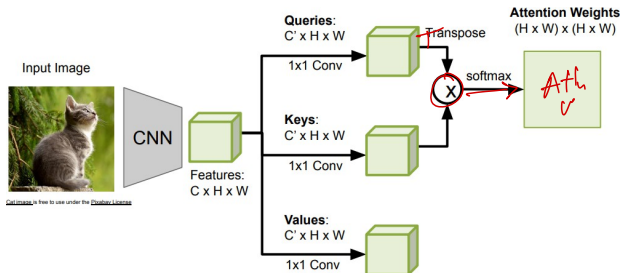
Example: CNN with Self-Attention



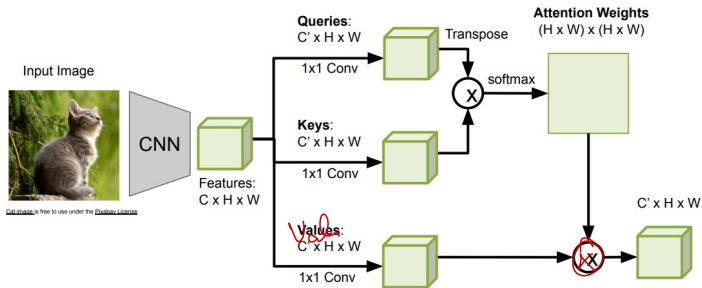
Example: CNN with Self-Attention (cont.)



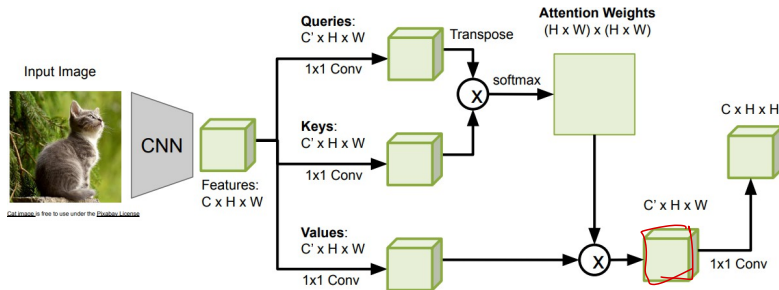
Example: CNN with Self-Attention (cont.)



Example: CNN with Self-Attention (cont.)

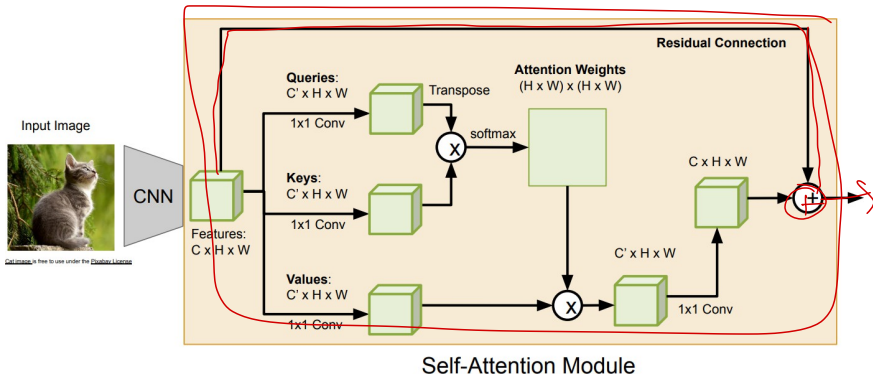


Example: CNN with Self-Attention (cont.)



Example: CNN with Self-Attention (cont.)

mul + Add



Attention is all you need

Vaswani et al, NeurIPS 2017

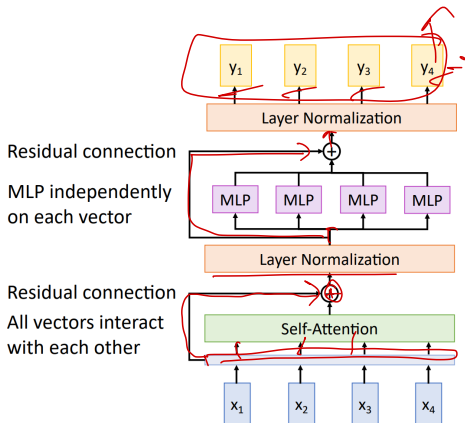


Figure 2: Detailed view of a single Transformer block.

- **Input:** A set of vectors \mathbf{x} , one per token.
- **Output:** A set of vectors \mathbf{y} , one per token.
- **Self-Attention:** Allows each token to attend to all others, capturing contextual relationships.
- **Layer Normalization & MLP:** Applied independently to each token, ensuring stability and expressiveness.
- **Key Properties:** Highly scalable and parallelizable due to independent operations across tokens.

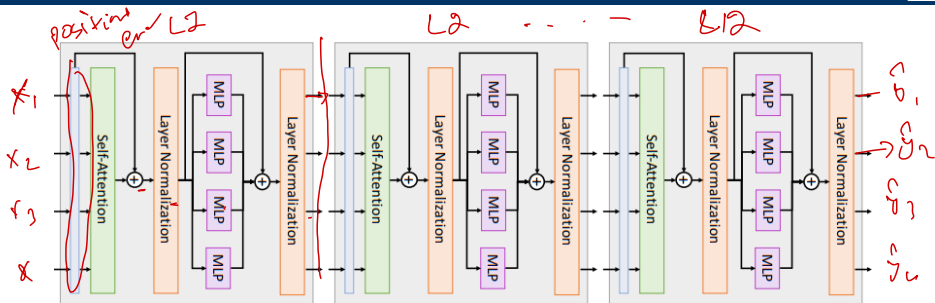


Figure 3: Stacking multiple Transformer blocks.

- ▶ A Transformer consists of a sequence of identical blocks, each refining the token representations.
- ▶ In the original architecture (Vaswani et al., 2017): 12 blocks, hidden size $D_Q = 512$ and 6 attention heads were used.
- ▶ The modular design enables deep stacking and efficient parallel computation.

- ▶ Transformers use self-attention to capture relationships within sequences.
- ▶ Self-attention allows each token to attend to all other tokens, enabling context-aware representations.
- ▶ Positional encoding is crucial for maintaining the order of tokens in the sequence.
- ▶ Masked self-attention prevents future information leakage during training.
- ▶ Multi-head attention enhances model capacity by allowing multiple attention mechanisms to learn different aspects of the data.

- [1] Fei-Fei Li, Yunzhu Li, and Ruohan Gao. *Stanford CS231n: Deep Learning for Computer Vision*.
<http://cs231n.stanford.edu/index.html>

- [2] Assaf Shocher, Shai Bagon, Meirav Galun, and Tali Dekel. *WAIC DL4CV Deep Learning for Computer Vision: Fundamentals and Applications*. <https://dl4cv.github.io/index.html>

- [3] Justin Johnson. *UMich EECS 498.008/598.008: Deep Learning for Computer Vision*. <https://web.eecs.umich.edu/~justincj/teaching/eecs498/WI2022/>

- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. *Attention Is All You Need*. Advances in Neural Information Processing Systems, 2017. <https://arxiv.org/abs/1706.03762>

- [5] Stanford CS231n. *Lecture 11: Attention & Transformers.*
<http://cs231n.stanford.edu/slides/2023/lecture11.pdf>

- [6] Carnegie Mellon University. *Deep Learning Lecture: Vision Transformers and Transformers.*
<https://deeplearning.cs.cmu.edu/F23/index.html>
<https://www.youtube.com/watch?v=1gq2y1E5g6M>

- [7] MIT 6.S191. *Introduction to Deep Learning Resources.*
<https://introtodeeplearning.com/>

- [8] Satya Mallick. *LearnOpenCV: A Comprehensive Guide to Attention Mechanisms in Deep Learning.* <https://learnopencv.com/attention-mechanism-in-neural-networks/>

Credits

Dr. Prashant Aparajeya

Computer Vision Scientist — Director(AISimply Ltd)

p.aparajeya@aisimply.uk

This project benefited from external collaboration, and we acknowledge their contribution with gratitude.