

Attention Mechanism Deep Dive

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للعلوم والتقنية

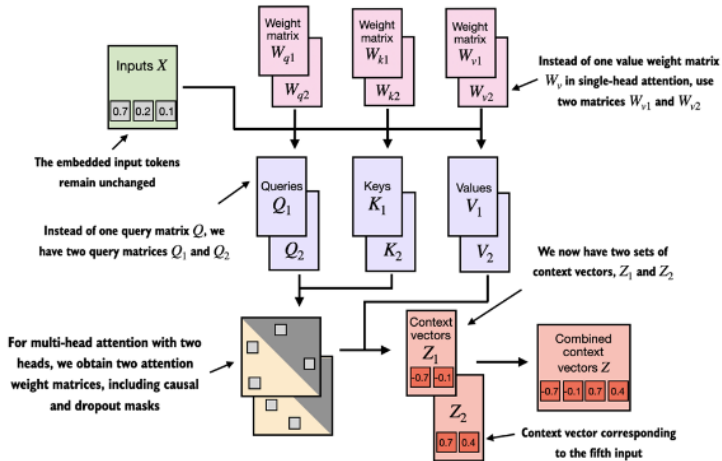
King Abdullah University of
Science and Technology



LMH

Lady Margaret Hall

July 26, 2025



1. Motivation for Attention Mechanisms
2. Learning Outcomes
3. Alignment Function
4. Query, Key, Value (QKV)
5. Attention-based Decoding Sequences
6. Summary
7. References

- ▶ Sequence-to-sequence models struggle with **long-range dependencies**, making it difficult to capture context over lengthy inputs.
- ▶ Attention mechanisms address this by allowing the model to **selectively focus** on relevant parts of the input sequence during processing.
- ▶ They form the **backbone** of modern architectures such as **Transformer models** (e.g., BERT, GPT, ViT).
- ▶ Attention is crucial in various tasks including **translation, summarization, and image captioning**.

Example: In machine translation, rather than encoding an entire sentence into a single fixed-size vector, attention enables the model to dynamically focus on relevant source words when generating each target word.

After this session, you should be able to:

- ▶ Understand how **alignment functions** operate.
- ▶ Comprehend **QKV (Query, Key, Value)** mechanism in attention.
- ▶ Learn how attention is used in **decoding sequences**.
- ▶ Recognize the **strengths and limitations** of attention mechanisms.
- ▶ Get familiar with **modern attention research**.

- ▶ The alignment function is a crucial component of attention mechanisms, determining how much focus to place on different parts of the input sequence.
- ▶ It computes a score for each input token relative to the current output token being generated.
- ▶ The scores are then normalized (often using softmax) to create a probability distribution, indicating the relative importance of each input token.
- ▶ Common alignment functions include:
 - **Dot Product:** Computes the dot product between the query and key vectors.
 - **Scaled Dot Product:** Similar to dot product but scaled by the square root of the dimension of the key vectors to prevent large values that can saturate the softmax function.

- **Additive Attention:** Combines the query and key vectors using a feed-forward neural network to compute the alignment scores.
- ▶ The choice of alignment function can significantly impact the model's performance, especially in tasks with long sequences or complex dependencies.

Example: In a translation task, the alignment function helps the model determine which words in the source language should be emphasized when generating each word in the target language. For instance, when translating "The cat sat on the mat," the model might focus more on "cat" when generating "gato" in Spanish, while still considering "sat" and "on" to maintain the context.

► **Dot Product:**

$$\text{score} = q^{\top} k$$

► **Scaled Dot Product:**

$$\text{score} = \frac{q^{\top} k}{\sqrt{d_k}}$$

(used in Transformers to prevent softmax saturation)

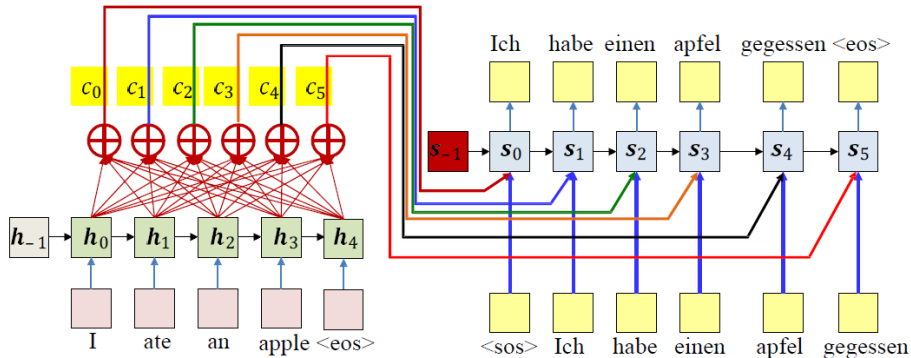
► **Additive (Bahdanau) Attention:**

$$\text{score} = v^{\top} \tanh(W_1 h_i + W_2 s_t)$$

► **General:**

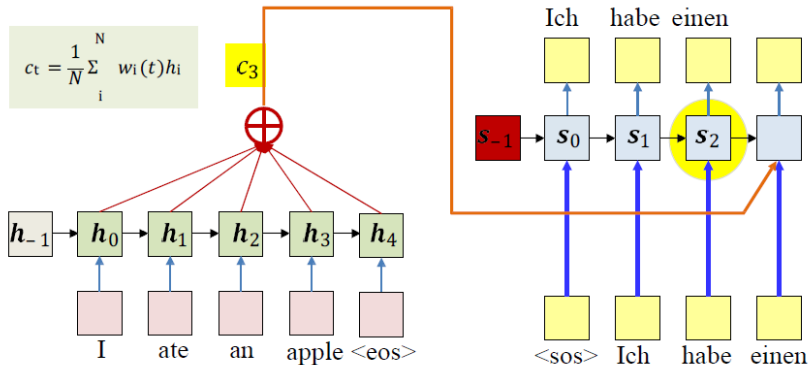
$$\text{score} = q^{\top} W k$$

Each function balances efficiency vs expressiveness.



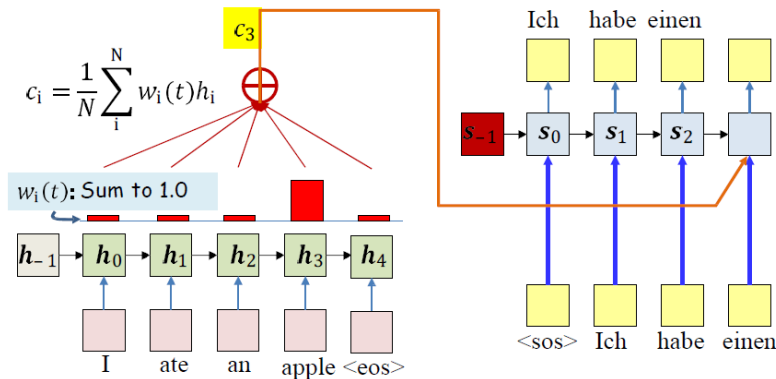
$$c_t = \frac{1}{N} \sum_i^N w_i(t) h_i$$

- **Attention weights:** The weights $w_i(t)$ are dynamically computed as functions of decoder state
 - Expectation: if the model is well-trained, this will automatically “highlight” the correct input
- But how are these computed?



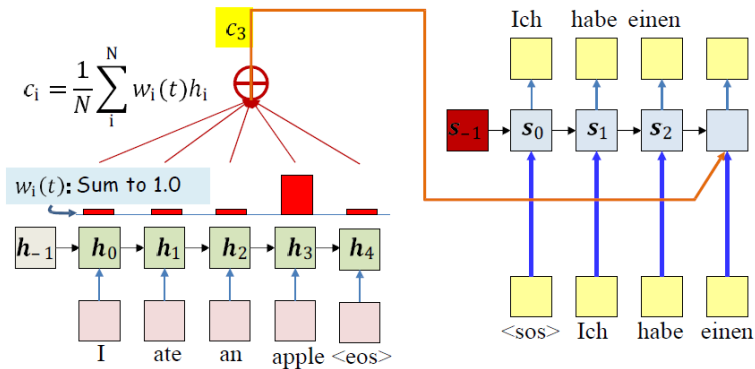
- The “attention” weights $w_i(t)$ at time t must be computed from available information at time t
- The primary information is s_{t-1} (the state at time $t - 1$)
 - Also, the input word at time t , but generally not used for simplicity

Requirement on attention weights



- The weights $w_i(t)$ must be positive and sum to 1.0
 - I.e. be a distribution
 - Ideally, they must be high for the most relevant inputs for the i th output and low elsewhere

Requirement on attention weights (cont.)



- The weights $w_i(t)$ must be positive and sum to 1.0
 - I.e. be a distribution
 - Ideally, they must be high for the most relevant inputs for the i th output and low elsewhere
- Solution: A two step weight computation
 - First compute *raw* weights (which could be +ve or -ve)
 - Then softmax them to convert them to a distribution

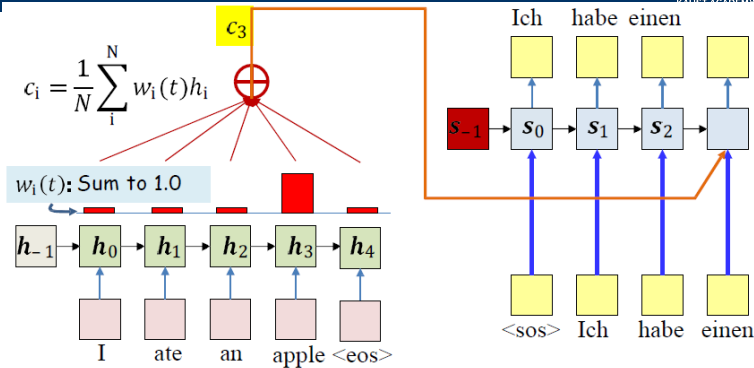
$$e_i(t) = g(h_i, s_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

Attention weights



Attention weights (cont.)



- Typical options for $g()$ (**variables in red must be learned**)

$$g(h_i, s_{t-1}) = h_i^T s_{t-1}$$

$$g(h_i, s_{t-1}) = h_i^T \mathbf{w}_g s_{t-1}$$

$$g(h_i, s_{t-1}) = \mathbf{v}_g^T \tanh\left(\mathbf{w}_g \begin{bmatrix} h_i \\ s_{t-1} \end{bmatrix}\right)$$

$$g(h_i, s_{t-1}) = \text{MLP}([h_i, s_{t-1}])$$

$$e_i(t) = g(h_i, s_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

- ▶ **Query:** What you want (the search)
- ▶ **Key:** Index of content (where to look)
- ▶ **Value:** Actual content (what you get)

Intuition

Think of QKV like a search engine:

- ▶ **Query (Q):** Comes from the current decoder state (or token)
- ▶ **Key (K), Value (V):** Come from the encoder (or input sequence)

Given:

- ▶ Input sequence \mathbf{X}

Linear projections:

$$\mathbf{Q} = \mathbf{XW}_Q$$

$$\mathbf{K} = \mathbf{XW}_K$$

$$\mathbf{V} = \mathbf{XW}_V$$

Attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right) \mathbf{V}$$

Attention-Based Decoding – Overview

- ▶ The decoder uses attention to select relevant encoder hidden states at each decoding step.
- ▶ **Combines:**
 - Previous decoder state
 - Encoder hidden states
 - Attention context vector
- ▶ The output at time t :

$$y_t = \text{Decoder}(y_{t-1}, \text{context}_t)$$

1. **Compute query:** Generate a query vector from the current decoder hidden state.
2. **Align with encoder outputs:** Compare the query with encoder outputs (keys) to measure relevance.
3. **Get attention weights:** Apply softmax to the alignment scores to obtain attention weights.
4. **Compute context:** Calculate the context vector as the weighted sum of encoder outputs using the attention weights.
5. **Feed into decoder:** Use the context vector and previous outputs as input for the next decoding step.

Example: English to French Translation

Source Sentence (English)

The cat sat on the mat.

Target Sentence (French)

Le chat s'est assis sur le tapis.

- ▶ At each decoding step, the decoder generates a French word by:
 1. Using the current decoder hidden state.
 2. Attending to relevant English words via attention weights.
 3. Computing a context vector as a weighted sum of encoder outputs.
- ▶ For example, when generating “chat”, the attention focuses on “cat”.
- ▶ The process repeats for each word, dynamically shifting attention.

- ▶ **Quadratic Complexity:** Attention computation scales as $\mathcal{O}(n^2)$ with sequence length, making it expensive for long sequences.
- ▶ **Shallow Reasoning:** Standard attention mechanisms typically consider only one step of interaction, limiting deep reasoning capabilities.
- ▶ **Noisy Attention Maps:** Attention distributions can be diffuse or hard to interpret, making model decisions less transparent.
- ▶ **Limited Long-Term Memory:** Attention alone may struggle to capture dependencies over very long contexts.

► Efficient Attention Variants:

- *Linformer, Performer, Longformer, FlashAttention* – Reduce memory and computation for long sequences.

► Structured Attention:

- *Sparsemax, Routing Attention* – Impose structure or sparsity for interpretability and efficiency.

► Cross-modal Attention:

- *Vision + Language (e.g., Flamingo, Gato)* – Integrate information across different modalities.

► Learnable Routing:

- *Mixture-of-Experts, Dynamic Attention* – Dynamically select computation paths for scalability.

Concept	Key Takeaway
Alignment	Scores relevance between states
QKV	Core mechanism to compute attention
Decoding	Uses attention to generate outputs
Limitation	Costly, hard to scale to long seqs
Future	Efficient, structured, multimodal attention

References

- ▶ Bahdanau, D., Cho, K., & Bengio, Y. (2015). *Neural Machine Translation by Jointly Learning to Align and Translate*.
- ▶ Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). *Attention Is All You Need*.
- ▶ Raffel, C., Shazeer, N., Roberts, A., et al. (2020). *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*.
- ▶ Choromanski, K., Likhoshesterov, V., Dohan, D., et al. (2021). *Rethinking Attention with Performers*.
- ▶ Tay, Y., Dehghani, M., Bahri, D., & Metzler, D. (2023). *Efficient Transformers: A Survey*.
- ▶ Google Research Slides: *The Annotated Transformer*.

Further Reading

- ▶ Attention Mechanisms in NLP: A Comprehensive Survey
- ▶ Advances in Transformer Architectures
- ▶ Efficient Attention Mechanisms for Long Sequences

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.