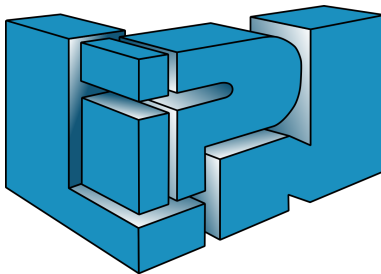


# Unveiling the Power of Attention: A Deep Dive into Attention Mechanisms

Bilal FAYE, Nicolas FLOQUET



# Table of contents

- 1 Introduction
- 2 Broad Understanding of Attention Mechanisms
- 3 Categories of Attention Mechanisms

# Introduction To Attention Mechanism

An attention mechanism in the context of machine learning refers to a computational model that allows a system to selectively focus on certain parts of input data while ignoring others.

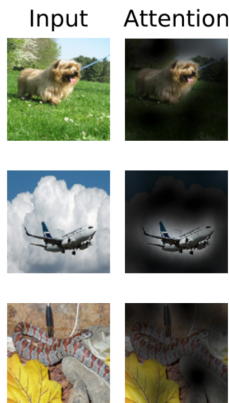


Figure – Attention Visualization on Image Classification<sup>1</sup>

---

1. Dosovitskiy et al., An Image is Worth 16x16 Words :  
Transformers for Image Recognition, 2020

Some key points highlighting the powerful aspects of attention mechanisms :

- Resource Allocation Efficiency
- Wide Range Applications
- Interpretable Neural Architectures
- Enhanced Performance in Specific Tasks
- Intuitive Explanations for Neural Behaviors

# Broad Understanding of Attention Mechanisms

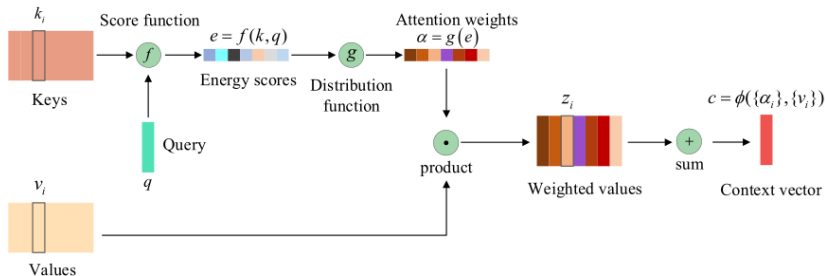


Figure – Unified Attention Model<sup>2</sup>

2. Zhaoyang et al., A review on the attention mechanism of deep learning, 2021



- **Keys** : This refers to the input data representations on which the model relies to identify relevant information or patterns in the sequence.
- **Query** : This represents what the model aims to search for or extract from the input data.
- **Values** : This corresponds to the actual information associated with each part of the input data (sequence).

The score function  $\mathbf{e} = \mathbf{f}(\mathbf{k}, \mathbf{q})$  determines the matching or combination of keys and queries.

Additive attention<sup>3</sup> combines keys and queries through a summation operation :

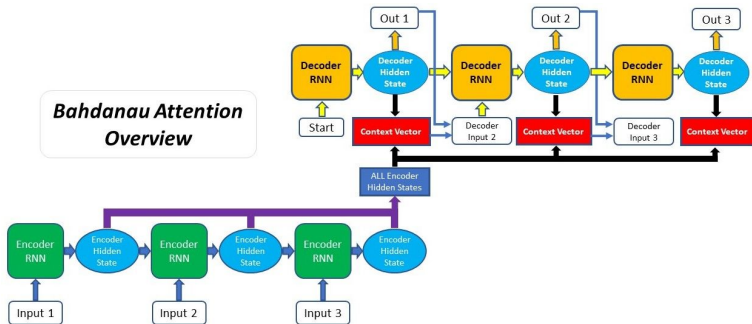
Additive score function

$$f(\mathbf{q}, \mathbf{k}) = \mathbf{v}^T \text{activation}(\mathbf{W}_1 \mathbf{k} + \mathbf{W}_2 \mathbf{q})$$

where  $\mathbf{v}$ ,  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are learnable parameters.

---

3. Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate, 2014

**Bahdanau Attention  
Overview**

# Broad Understanding of Attention Mechanisms

Multiplicative (dot-product) attention<sup>4</sup> computes the relevance between keys and queries by taking their dot product :

## Multiplicative score function

$$f(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$$

In WMT'15 English  $\rightarrow$  German task<sup>5</sup>, authors found that parameterized additive attention slightly outperformed multiplicative attention.

---

4. Luong et al., Effective Approaches to Attention-based Neural Machine Translation, 2015

5. Britz et al., Massive exploration of neural machine translation architectures, 2017

# Broad Understanding of Attention Mechanisms

Scaled multiplicative (dot-product) attention<sup>6</sup> computes the relevance between keys and queries by taking their dot product :

Scaled multiplicative score function

$$f(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{d_k}}$$

where  $d_k$  is the dimension of keys.

For small values of  $d_k$ , both mechanisms perform similarly, but additive attention outperforms multiplicative attention without scaling for larger  $d_k$ .

---

6. Vaswani et al., Attention is all you need, 2017

General attention<sup>7</sup> extends the concept of multiplicative attention by introducing a learnable matrix parameter  $\mathbf{W}$  :

### General score function

$$f(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{W} \mathbf{k}$$

where  $\mathbf{W}$  is a learnable parameter.

This approach is applicable to keys and queries with distinct representations.

---

7. Luong et al., Effective Approaches to Attention-based Neural Machine Translation, 2015

Concat attention<sup>8</sup> aims to derive the joint representation of the keys and queries instead of comparing them :

#### Concat score function

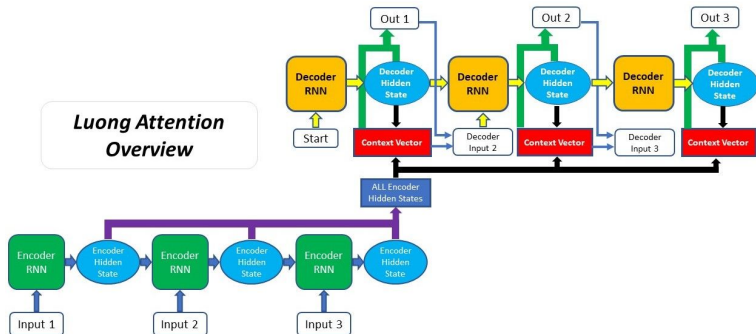
$$\mathbf{f}(\mathbf{q}, \mathbf{k}) = \mathbf{v}^T \text{activation}(\mathbf{W}[\mathbf{k}; \mathbf{q}])$$

where  $\mathbf{W}$  is a learnable parameter.

---

8. Luong et al., Effective Approaches to Attention-based Neural Machine Translation, 2015



***Luong Attention Overview***

Location-based attention<sup>9</sup> are solely computed from the target hidden state :

### Location-based score function

$$f(\mathbf{q}, \mathbf{k}) = f(\mathbf{q})$$

Energy scores ( $f$ ) depend solely on  $\mathbf{q}$  rather than  $\mathbf{K}$ . Conversely, self-attention is calculated solely based on  $\mathbf{K}$ , without requiring  $\mathbf{q}$ .

---

9. Luong et al., Effective Approaches to Attention-based Neural Machine Translation, 2015

Similarity attention<sup>10</sup> compares the similarity between  $\mathbf{K}$  and  $\mathbf{q}$ , which relied on cosine similarity. :

### Similarity score function

$$f(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q} \cdot \mathbf{k}}{\|\mathbf{q}\| \cdot \|\mathbf{k}\|}$$

Similarity attention is crucial in :

- semantic similarity assessments (in NLP)
- feature-based comparisons (in CV)

---

10. Graves et al., Neural Turing machines, 2014

The distribution function  $\mathbf{g}$  corresponds to the softmax, logistic or sigmoid, which normalize all the energy scores to a probability distribution.

After calculating attention weights and values, the context vector  $\mathbf{c}$  is computed as follows :

### Context vector

$$\mathbf{c} = \phi(\{\alpha_i\}, \{\mathbf{v}_i\}),$$

where  $\phi$  is a function that returns a single vector given the set of values and their corresponding weights.

$$\mathbf{z}_i = \alpha_i \mathbf{v}_i,$$

and

$$\mathbf{c} = \sum_{i=1}^n \mathbf{z}_i,$$

where  $\mathbf{z}_i$  is a weighted representation of an element in values and  $n$  is the dimension of  $\mathbf{Z}$ .

# Categories of Attention Mechanisms

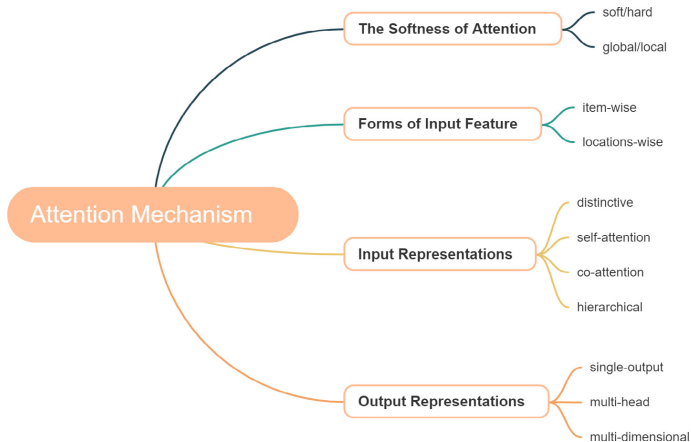


Figure – Unified Attention Model<sup>11</sup>

11. Zhaoyang et al., A review on the attention mechanism of deep learning, 2021

## *Categorie 1 : The Softness of Attention*



- **Soft Attention :**

- **Definition :** Soft (deterministic) attention calculates a context vector through a weighted average of all keys, facilitating differentiability with respect to inputs, thus enabling training via standard backpropagation methods.

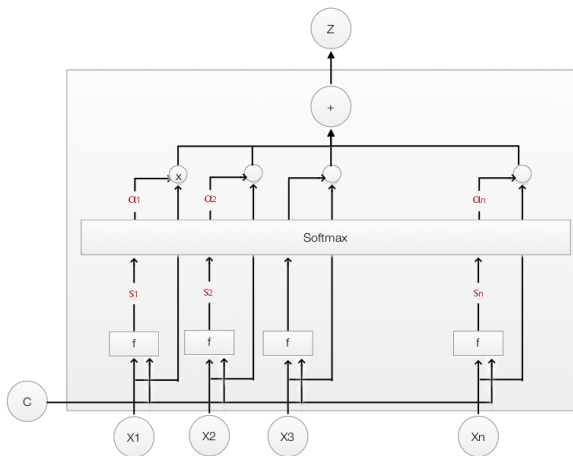


Figure – An instance demonstrating the application of soft attention.

- **Hard Attention :**

- **Definition :** Hard (stochastic) attention is an attention mechanism that makes discrete decisions regarding which parts of the input sequence to focus on, resulting in non-differentiability with respect to inputs, thereby complicating training using standard backpropagation methods.

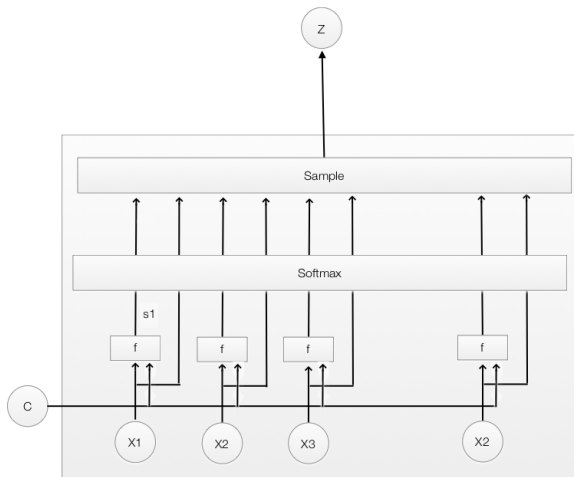
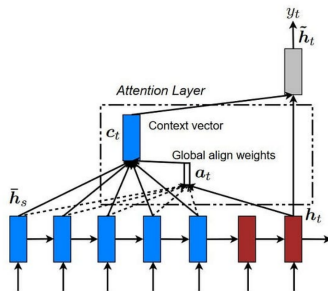


Figure – An instance demonstrating the application of hard attention.

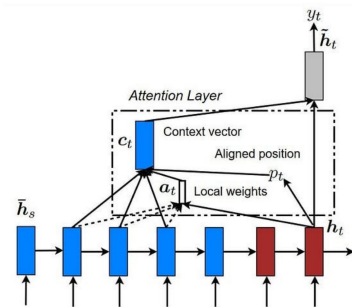
- **Global Attention :**

- **Definition :** Global attention is similar to soft attention, all source words are considered at a time.



- **Local Attention :**

- **Definition :** In Local attention, only a subset of source words are considered at a time.



## *Categorie 2 : Form of Input Feature*

- **Item-wise attention**

- **Definition :** The item-wise attention requires that the input is either explicit items or an additional preprocessing (ex. word embeddings) step is added to generate a sequence of items (vectors) from the source data.
- item-wise soft attention calculates a weight for each item, and then makes a linear combination of them.
- Instead of a linear combination of all items, the item-wise hard attention stochastically picks one or some items based on their probabilities.



- **Location-wise attention**

- **Definition :** location-wise attention is aimed at tasks that are difficult to obtain distinct input items (visual tasks).
- The location-wise soft attention accepts an entire feature map as input and generates a transformed version through the attention module.
- The location-wise hard attention stochastically picks a sub-region as input and the location of the sub-region to be picked is calculated by the attention module.

## *Categorie 3 : Input Representations*

- **Distinctive Attention :**

- **Definition :** Distinctive attention, as defined in the mentioned context, involves attention models with a single input and corresponding output sequence, where keys and queries are derived from two independent sequences.

- **Self-Attention :**

- **Definition :** Self-attention is an attention mechanism in which each element in a sequence attends to all other elements in the same sequence.

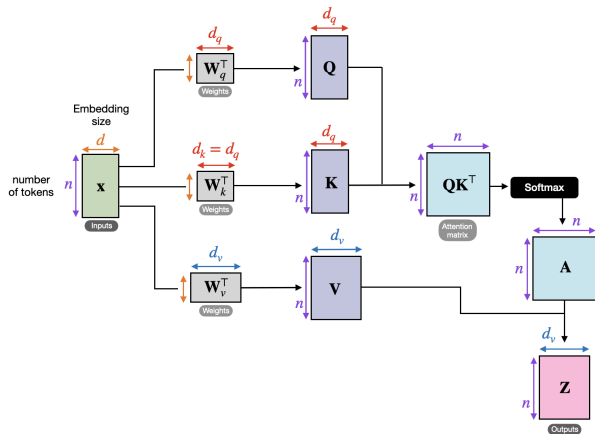


Figure – Application of self-attention

- **Cross-Attention :**

- **Definition :** Cross-attention refers to scenarios where attention is applied between different parts of the input and output sequences.

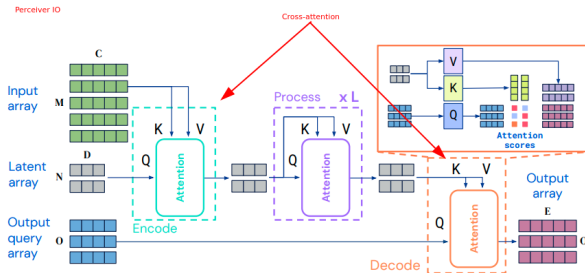


Figure – Application of cross-attention

- **Co-Attention :**

- **Definition :** Co-attention refers to an attention mechanism that simultaneously considers and aligns information from multiple input sequences or modalities. co-attention can be coarse-grained<sup>12</sup> or fine-grained<sup>13</sup>.

---

12. Coarse-grained attention computes attention on each input, using an embedding of the other input as a query

13. Fine-grained attention evaluates how each element of an input affects each element of the other input



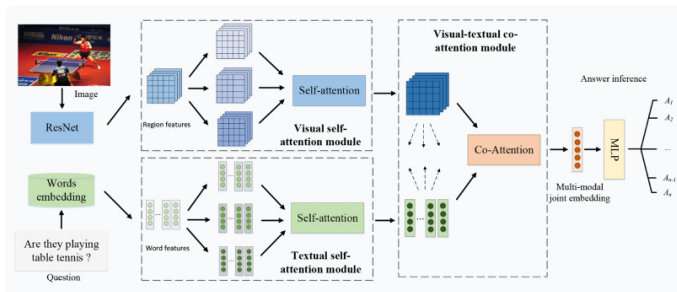


Figure – Application of co-attention

- **Hierarchical Attention :**

- **Definition :** Hierarchical attention allows the computation of attention weights not only from the original input sequence but also from different abstraction levels (ex. document classification).

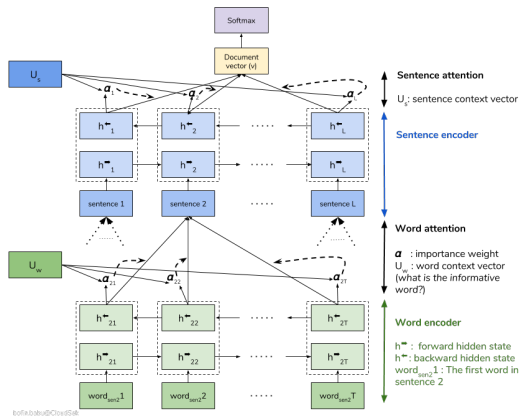


Figure – Application of hierarchical attention

## *Categorie 4 : Output Representations*

- **Single Output Attention :**

- **Definition :** The energy scores are represented by one and only one vector at each time step.

- **Multi-Head Attention :**

- **Definition :** An attention mechanism that employs multiple attention heads to capture diverse features and relationships in parallel.

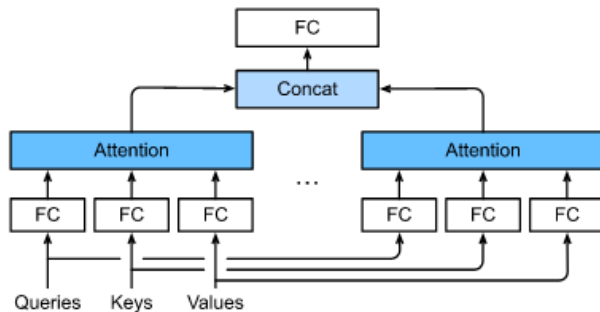


Figure – Application of multi-head attention

- **Multi-Dimensional Attention :**
  - **Definition :** An approach that computes a feature-wise score vector for keys by replacing weight scores vector with a matrix. In this way, the neural network can calculate multiple attention distributions for the same data.



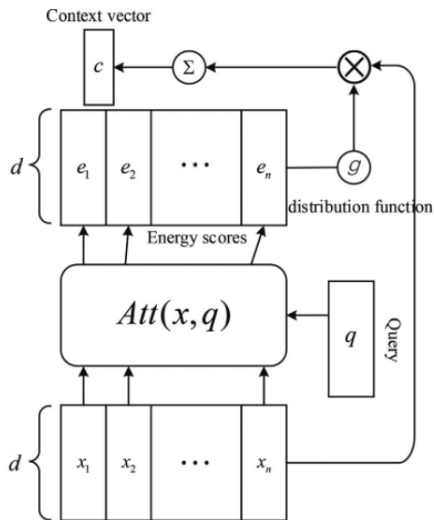


Figure – Application of multi-dimensional attention