

## Attention Mechanism: Mathematical Formulas

### Compute Attention Scores

Similarity

1

$$e_{t,i} = \text{score}(s_t, h_i) = s_t^T W_a h_i$$

Calculate similarity between decoder state  $s_t$  and each encoder hidden state  $h_i$  using learnable weight matrix  $W_a$

### Compute Attention Weights

Softmax

2

$$\alpha_{t,i} = \exp(e_{t,i}) / \sum_j \exp(e_{t,j})$$

Normalize scores using softmax to obtain attention weights that sum to 1. Higher scores → higher attention weights

### Compute Context Vector

Weighted Sum

3

$$c_t = \sum_i \alpha_{t,i} \times h_i$$

Weighted sum of encoder hidden states using attention weights. More relevant states contribute more to the context

### Generate Output

Prediction

4

$$\hat{y}_t = \text{softmax}(W_o [s_t; c_t] + b_o)$$

Concatenate decoder state and context vector, then pass through output layer to predict next token

## Key Components Summary

### Input Dimensions

$$s_t \in \mathbb{R}^d$$

$$h_i \in \mathbb{R}^d$$

$$W_a \in \mathbb{R}^{d \times d}$$

### Attention Properties

$$\sum_i \alpha_{t,i} = 1$$

$$\alpha_{t,i} \in [0, 1]$$

$$0 \leq \alpha_{t,i} \leq 1$$

### Output

$$c_t \in \mathbb{R}^d$$

Dynamic context

Per time step

**s<sub>t</sub>:** Decoder hidden state   **h<sub>i</sub>:** Encoder hidden state   **e<sub>t,i</sub>:** Attention score   **α<sub>t,i</sub>:** Attention weight   **c<sub>t</sub>:** Context vector  
**W<sub>a</sub>:** Alignment weight matrix



## Concrete Example with 3D Vectors



**Setup:** d = 3 (vector dimension), using 2 encoder hidden states (h<sub>1</sub>, h<sub>2</sub>).

1

### Given Input Data

Decoder hidden state (s<sub>t</sub>):

$$s_t = [1, 2, 1]^T$$

Encoder hidden states:

$$h_1 = [2, 0, 1]^T$$

$h_2 = [1, 1, 2]^T$

Weight matrix ( $W_a$ )  $\in \mathbb{R}^{3 \times 3}$ :

$W_a = [1 \ 0 \ 0] \ [0 \ 1 \ 0] \ [0 \ 0 \ 1]$

(Using Identity Matrix for simplicity)

## 2

### Compute Attention Scores

► Formula:  $e_{t,i} = s_t^T W_a h_i$

• Score for  $h_1$  :

$$\begin{aligned} e_{t,1} &= s_t^T W_a h_1 = [1, 2, 1] \cdot [2, 0, 1]^T \\ &= (1 \times 2) + (2 \times 0) + (1 \times 1) \\ &= 2 + 0 + 1 \\ &= 3 \end{aligned}$$

• Score for  $h_2$  :

$$\begin{aligned} e_{t,2} &= s_t^T W_a h_2 = [1, 2, 1] \cdot [1, 1, 2]^T \\ &= (1 \times 1) + (2 \times 1) + (1 \times 2) \\ &= 1 + 2 + 2 \\ &= 5 \end{aligned}$$



Attention Scores:

$$e_{t,1} = 3, e_{t,2} = 5$$

3

### Compute Attention Weights (Softmax)

► Formula:  $\alpha_{t,i} = \exp(e_{t,i}) / \sum_j \exp(e_{t,j})$

- Exponential calculation:

$$\exp(e_{t,1}) = \exp(3) \approx 20.09$$

$$\exp(e_{t,2}) = \exp(5) \approx 148.41$$

- Sum:

$$\sum \exp(e_{t,j}) = 20.09 + 148.41 = 168.50$$

- Attention Weights:

$$\alpha_{t,1} = 20.09 / 168.50 \approx 0.119$$

$$\alpha_{t,2} = 148.41 / 168.50 \approx 0.881$$

- Verification:

$$\alpha_{t,1} + \alpha_{t,2} = 0.119 + 0.881 = 1.000 \checkmark$$



Attention Weights:

$$\alpha_{t,1} \approx 0.119 \text{ (11.9\%)}, \alpha_{t,2} \approx 0.881 \text{ (88.1\%)}$$

 **Interpretation:**  $h_2$  receives about 88% attention because it's more similar to  $s_t$ !

## 4 Compute Context Vector

 **Formula:**  $c_t = \sum_i \alpha_{t,i} \times h_i$

- **Weighted sum calculation:**

$$c_t = \alpha_{t,1} \times h_1 + \alpha_{t,2} \times h_2$$

- **First term:**

$$\begin{aligned}\alpha_{t,1} \times h_1 &= 0.119 \times [2, 0, 1]^T \\ &= [0.238, 0.000, 0.119]^T\end{aligned}$$

- **Second term:**

$$\begin{aligned}\alpha_{t,2} \times h_2 &= 0.881 \times [1, 1, 2]^T \\ &= [0.881, 0.881, 1.762]^T\end{aligned}$$

- **Final Context Vector:**

$$\begin{aligned}c_t &= [0.238, 0.000, 0.119]^T + [0.881, 0.881, 1.762]^T \\ &= [1.119, 0.881, 1.881]^T\end{aligned}$$



**Final Context Vector:**

$$c_t = [1.119, 0.881, 1.881]^T$$

 **Result Interpretation:** The context vector is formed closer to  $h_2$  (88.1% weight). This means the current decoder state is more relevant to  $h_2$ .

## Key Takeaways

1. Attention Score measures similarity via dot product between two vectors
2. Softmax normalizes scores into probabilities between 0 and 1
3. Context Vector is a weighted average, focusing on important information
4. High similarity → high attention weight → greater influence