# Attention (Query, Key, Value)

• What it is?

A mechanism that lets each token decide which other tokens are important.

• Goal:

To compute a context - aware representation for each token.

• Math :

- Queries, keys, and values are linear projections:

$$Q = X W_Q, \quad K = X W_K, \quad V = X W_V$$

- Attention scores :

$$scores = \frac{Q K^T}{\sqrt{d_K}}$$

- Attention weighs :

$$A = softmax(scores)$$

- Output :

$$Output = AV$$

- Conclusion:

Attention allows each token to dynamically focus on relevant parts of the sequence.