

Attention Computation: Step-by-Step Flow

STEP 1:

Query from "mat" meets all Keys



STEP 2:

Calculate $Q \cdot K$ (dot products)

0.1

0.3

0.4

0.8

0.6

Higher score
= more relevant

STEP 3:

Apply Softmax (convert to percentages)

14%

17%

19%

28%

23%

Sum = 100%

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

STEP 4:

Multiply weights with Values

V("The")

V("cat")

V("sat")

V("on")

V("the")

$\times 14\%$

$\times 17\%$

$\times 19\%$

$\times 28\%$

$\times 23\%$

Each Value
contributes
proportionally

STEP 5:

Sum all weighted values

$$14\% \times V("The") + 17\% \times V("cat") + 19\% \times V("sat") + \dots$$

OUTPUT: Context-aware representation of "mat"

Knows it follows "on the" (location pattern)

Key Insight: Attention learns to identify and combine relevant information automatically