

Attention (p.4) (LLM notes)

- from before to compare key Queries in Attention blocks we use dot products here is a visualization:

★ Attention 8x8 Pattern
 • = ignore (cardinal)
 ! has normalized values

	a ↓ E ₁ ↓ W _a Q ₁	Fluffy ↓ E ₂ ↓ W _a Q ₂	blue ↓ E ₃ ↓ W _a Q ₃	creature ↓ E ₄ ↓ W _a Q ₄	roamed ↓ E ₅ ↓ W _a Q ₅	the ↓ E ₆ ↓ W _a Q ₆	verdant ↓ E ₇ ↓ W _a Q ₇	Forest ↓ E ₈ ↓ W _a Q ₈
a → E ₁ → K ₁	K ₁ .Q ₁	K ₁ .Q ₂	K ₁ .Q ₃	K ₁ .Q ₄	K ₁ .Q ₅	K ₁ .Q ₆	K ₁ .Q ₇	K ₁ .Q ₈
Fluffy → E ₂ → K ₂	K ₂ .Q ₁	K ₂ .Q ₂	K ₂ .Q ₃	K ₂ .Q ₄	K ₂ .Q ₅	K ₂ .Q ₆	K ₂ .Q ₇	K ₂ .Q ₈
blue → E ₃ → K ₃	K ₃ .Q ₁	K ₃ .Q ₂	K ₃ .Q ₃	K ₃ .Q ₄	K ₃ .Q ₅	K ₃ .Q ₆	K ₃ .Q ₇	K ₃ .Q ₈
creature → E ₄ → K ₄	K ₄ .Q ₁	K ₄ .Q ₂	K ₄ .Q ₃	K ₄ .Q ₄	K ₄ .Q ₅	K ₄ .Q ₆	K ₄ .Q ₇	K ₄ .Q ₈
roamed → E ₅ → K ₅	K ₅ .Q ₁	K ₅ .Q ₂	K ₅ .Q ₃	K ₅ .Q ₄	K ₅ .Q ₅	K ₅ .Q ₆	K ₅ .Q ₇	K ₅ .Q ₈
the → E ₆ → K ₆	K ₆ .Q ₁	K ₆ .Q ₂	K ₆ .Q ₃	K ₆ .Q ₄	K ₆ .Q ₅	K ₆ .Q ₆	K ₆ .Q ₇	K ₆ .Q ₈
verdant → E ₇ → K ₇	K ₇ .Q ₁	K ₇ .Q ₂	K ₇ .Q ₃	K ₇ .Q ₄	K ₇ .Q ₅	K ₇ .Q ₆	K ₇ .Q ₇	K ₇ .Q ₈
forest → E ₈ → K ₈	K ₈ .Q ₁	K ₈ .Q ₂	K ₈ .Q ₃	K ₈ .Q ₄	K ₈ .Q ₅	K ₈ .Q ₆	K ₈ .Q ₇	K ₈ .Q ₈

→ to measure the dot product we compute it between each key Query pair. the larger Dot products mean the key and Queries align and the dot product would be some large num. meaning the embedding of keys (Fluffy, blue) attend to the embedding of Query (creature), and the neg dot products like "the" and "creature" represents that these are unrelated to each other.

- so we have numbers from $-\infty \rightarrow \infty$ that give us a score for how relevant each word is to updating the meaning of every other word

Attends to

	Q ₁ a	Q ₂ Fluffy	Q ₃ blue	Q ₄ creature...
a K ₁	0.7	-0.93	-0.24	-5
Fluffy K ₂	-0.73	2.9	-0.5	+93
blue K ₃	-0.53	-0.57	1.9	+94
creature K ₄	-0.4	-0.29	-0.86	4.9

• = size of Dot product range $(-\infty, \infty)$

- the way we will use this score is to take a weighted sum across each column, weighted by relevance (size •) so we need normalization specifically range $(-\infty, \infty) \rightarrow \text{range}(0,1)$ sum(vals) = 1

- here we use softmax for each col

Ex $K_1 \begin{bmatrix} 0.7 \\ -0.93 \\ -0.24 \\ -5 \end{bmatrix} \xrightarrow{\text{softmax}} K_1 \begin{bmatrix} 0.02 \\ 0.92 \\ 0.89 \\ 0 \end{bmatrix} \rightarrow \text{sum} = 1$

$\text{sum(col)} = 1 \rightarrow \text{range}(0,1)$