

# KV-Cache

$$\text{softmax}\left(\frac{Q \times k^T}{\sqrt{d_{\text{head}}}} + \text{MASK}\right) = Q$$

	1	LOVE	PEPPERONI
1	1.0	0	0
LOVE	0.6	0.4	0
PEPPERONI	0.2	0.4	0.4

	1	LOVE	PEPPERONI
1	1.0	0	0
LOVE	0.6	0.4	0
PEPPERONI	0.2	0.4	0.4

(3, 3)

$$\times \begin{bmatrix} [1 \dots 128] \\ [1 \dots 128] \\ [1 \dots 128] \end{bmatrix}$$

(3, 128)

$$= \begin{bmatrix} [1 \dots 128] \\ [1 \dots 128] \\ [1 \dots 128] \end{bmatrix}$$

(3, 128)

CONTEXTUALIZED EMBEDDINGS