

Masking: We don't want to get affected by future tokens.

	Key1	Key2	Key3	KeyN
Query1	$s_{11}$	$-\infty$	$\dots$	$-\infty$
Query2	$s_{21}$	$s_{22}$	$\dots$	$-\infty$
Query3	$s_{31}$	$s_{32}$	$\dots$	$-\infty$
QueryN	$s_{N1}$	$s_{N2}$	$\dots$	$s_{NN}$

Masked  $QK^T$  with upper-triangular  $-\infty$

	Key1	Key2	Key3	KeyN
Query1	$s_{11}$	$-\infty$	$\dots$	$-\infty$
Query2	$s_{21}$	$\alpha_{22}$	$\dots$	$-\infty$
Query3	$\vdots$	$\vdots$	$\ddots$	$\vdots$
QueryN	$s_{N1}$	$s_{N2}$	$\dots$	$s_{NN}$

Softmax applied per query row (row highlighted)

	Key1	Key2	Key3	KeyN
Query1	$\alpha_{11}$	0	$\dots$	0
Query2	$\alpha_{21}$	$\alpha_{22}$	$\dots$	0
Query3	$\alpha_{31}$	$\alpha_{32}$	$\dots$	0
QueryN	$\alpha_{N1}$	$\alpha_{N2}$	$\dots$	$\alpha_{NN}$

Attention Weight(After softmax): masked positions ( $-\infty$ ) become 0