

Attention Weights

	\mathbf{k}_1	\dots	\mathbf{k}_n
\mathbf{q}_1	$Weight_{q_1 k_1}$	\dots	$Weight_{q_1 k_n}$
\vdots	\vdots	\ddots	\vdots
\mathbf{q}_n	$Weight_{q_n k_1}$	\dots	$Weight_{q_n k_n}$

\times

Values V

\mathbf{v}_1	v_{11}	\dots	v_{1d_v}
\vdots	\vdots	\ddots	\vdots
\mathbf{v}_n	v_{n1}	\dots	v_{nd_v}



Self-Attention Attention(Q, K, V)

Att_{11}	\dots	Att_{1d_v}
\vdots	\ddots	\vdots
Att_{n1}	\dots	Att_{nd_v}