

Self-Attention Heads

$\text{concat}(\mathbf{head}_1, \dots, \mathbf{head}_h)$

\mathbf{head}_1

\dots

\mathbf{head}_h

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

Weights \mathbf{W}^O

\times

w_{11}^O	\dots	$w_{1d_{model}}^O$
\vdots	\ddots	\vdots
$w_{(hd_v)1}^O$	\dots	$w_{(hd_v)d_{model}}^O$



Multi-Head Self-Attention

$\text{concat}(\mathbf{head}_1, \dots, \mathbf{head}_h)\mathbf{W}^O$

MHA_{11}	\dots	$MHA_{1d_{model}}$
\vdots	\ddots	\vdots
MHA_{n1}	\dots	$MHA_{nd_{model}}$