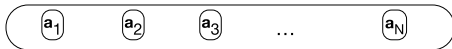


$[N \times d]$



Project from
 hd_v to d

w^O $[hd_v \times d]$

Concatenate
Outputs
 $[N \times hd_v]$

head1 output val
 $[N \times d_v]$

head2 output val
 $[N \times d_v]$

head3 output val
 $[N \times d_v]$

head4 output val
 $[N \times d_v]$

Multihead
Attention Layer
with $h=4$ heads

w^Q_1, w^K_1, w^V_1 Head 1

w^Q_2, w^K_2, w^V_2 Head 2

w^Q_3, w^K_3, w^V_3 Head 3

w^Q_4, w^K_4, w^V_4 Head 4

$[N \times d]$

