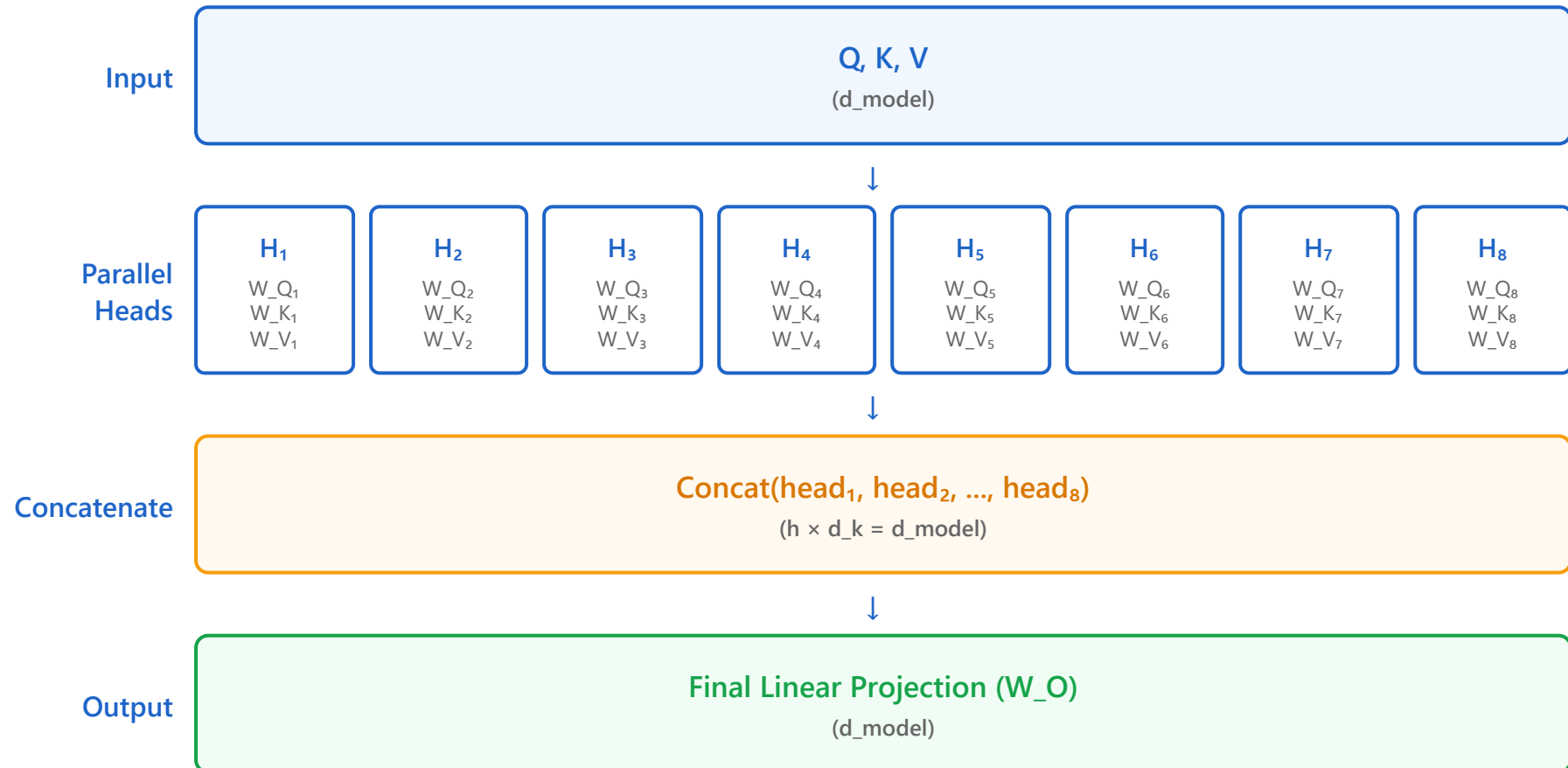


## Multi-Head Attention Architecture



### Configuration

Typically  $h = 8$  heads  
Head dimension:  $d_k = d_{\text{model}} / h$



### Parameters

Each head: separate  $W_Q$ ,  $W_K$ ,  $W_V$   
Total params  $\approx$  single large head



### Formula

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W_O$$