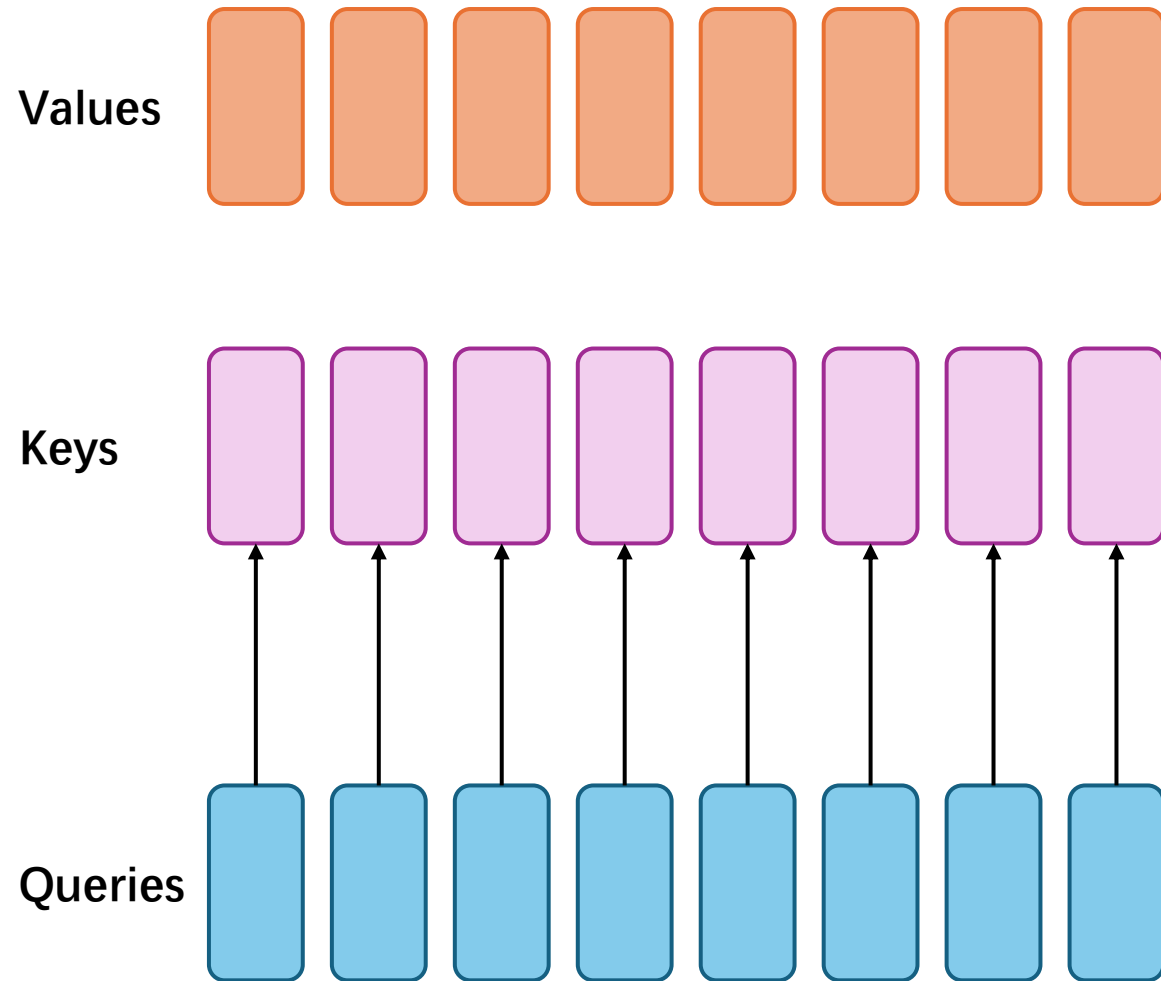


Multi-head Attention



Grouped-query Attention

