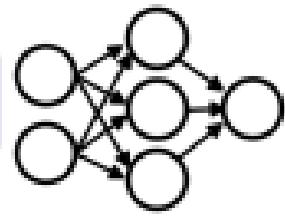
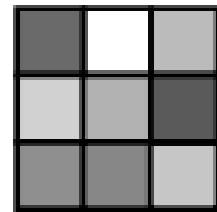


$[B, \text{seq\_len}, n_{\text{embd}}]$ 

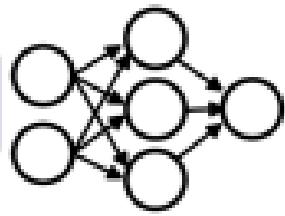
MLP



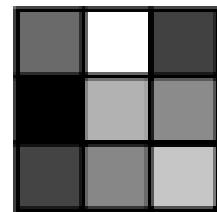
self-attention



MLP



self-attention

 $[B, \text{seq\_len}, n_{\text{embd}}]$