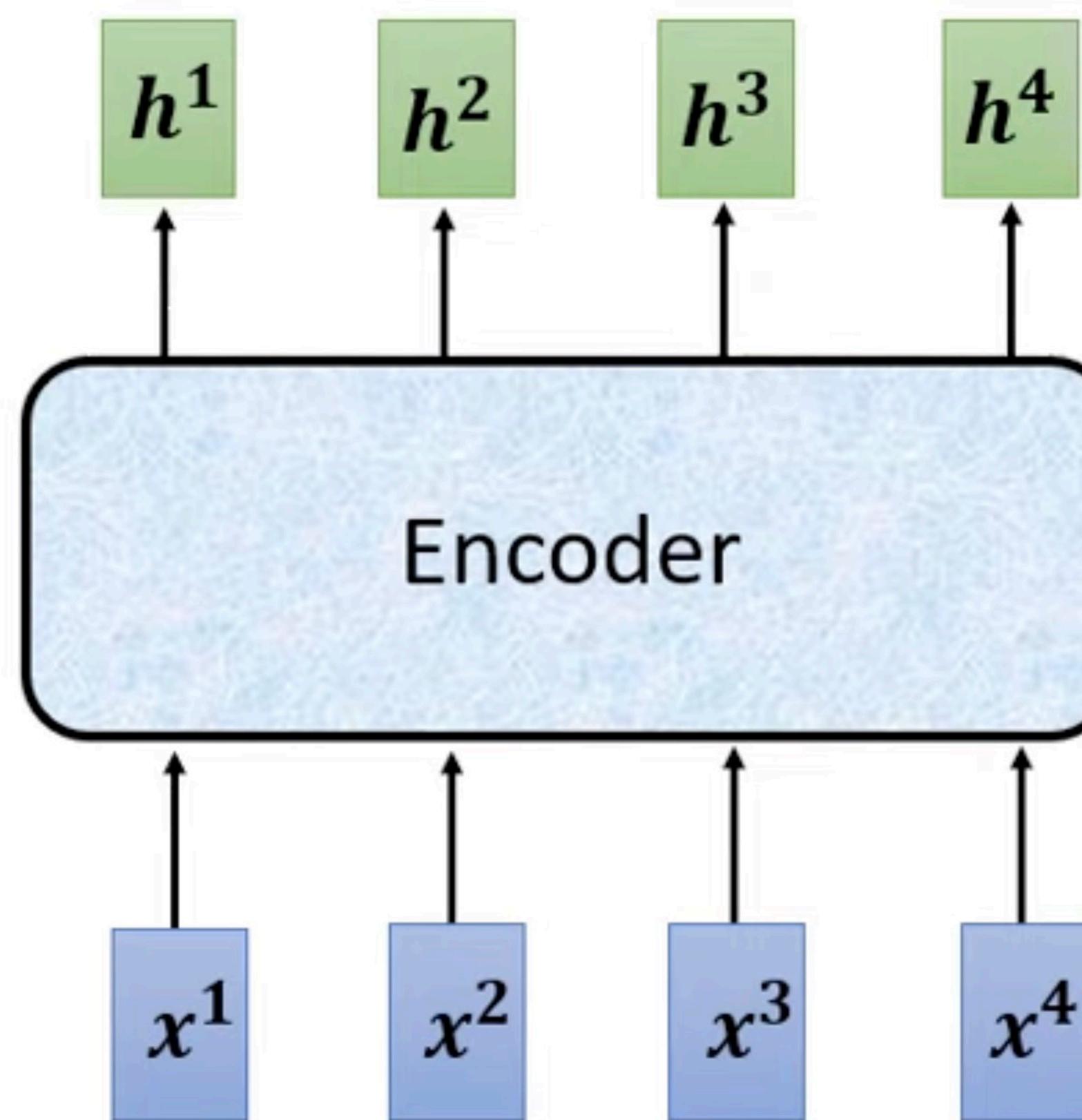
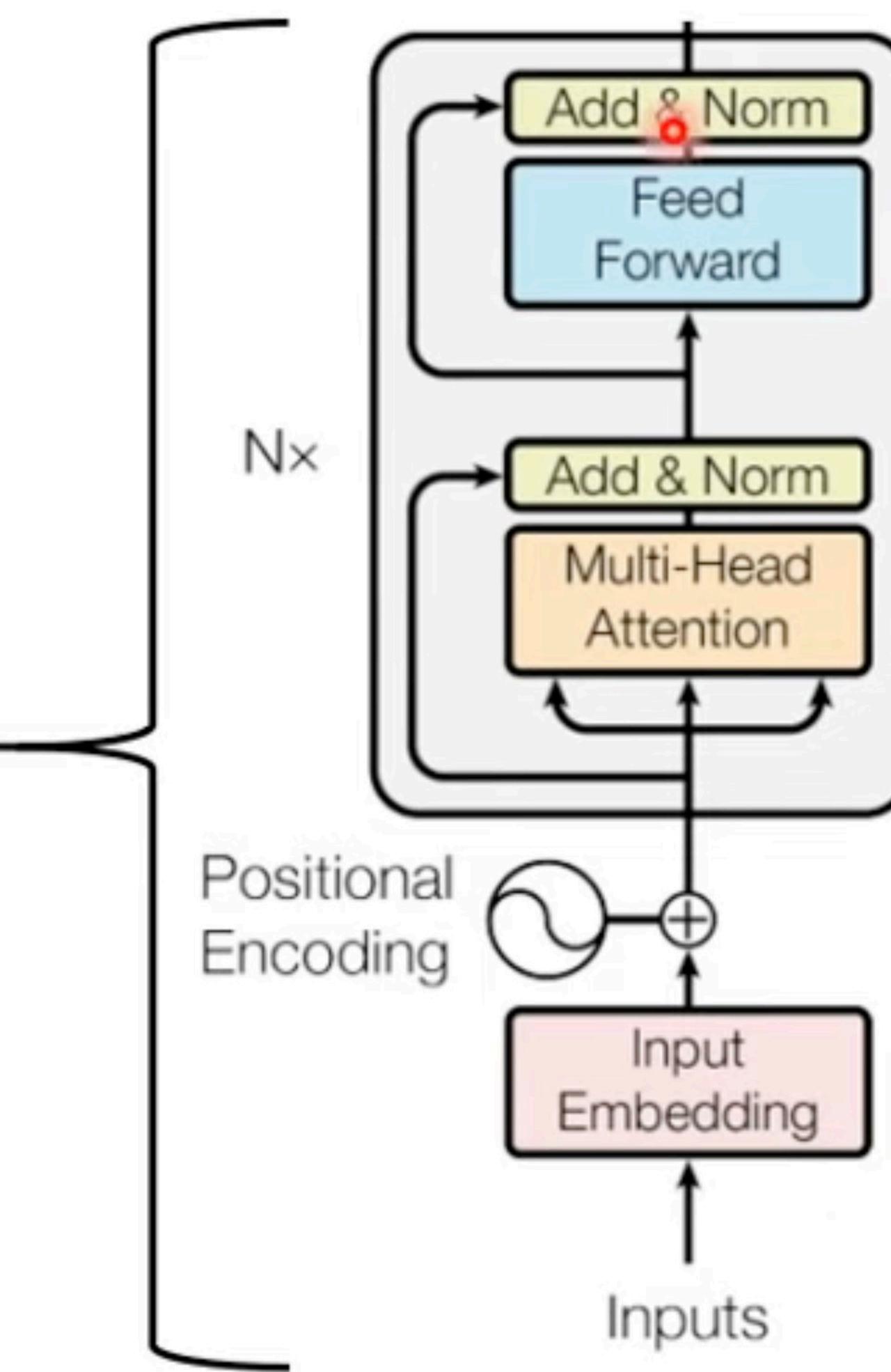


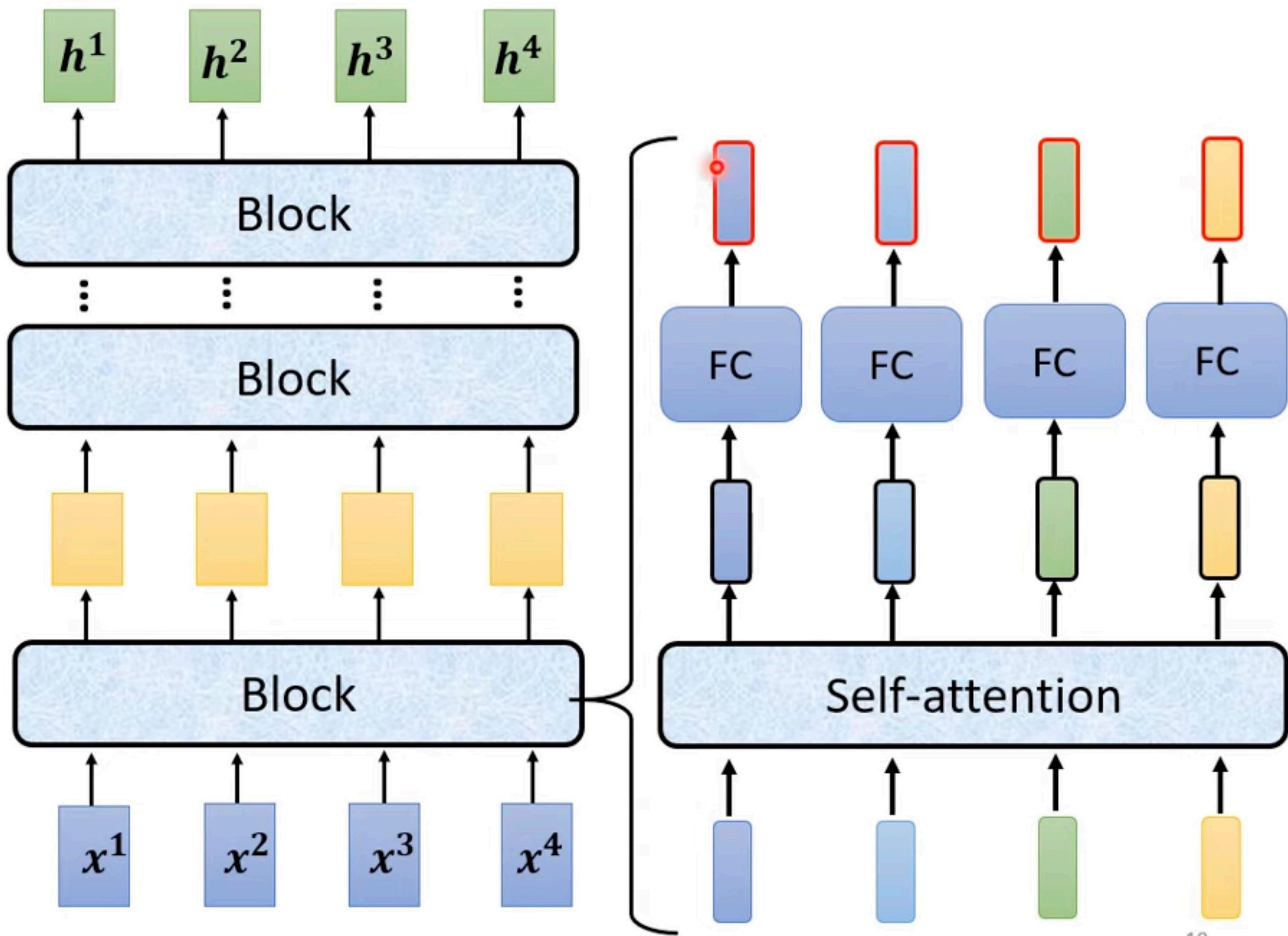
Encoder

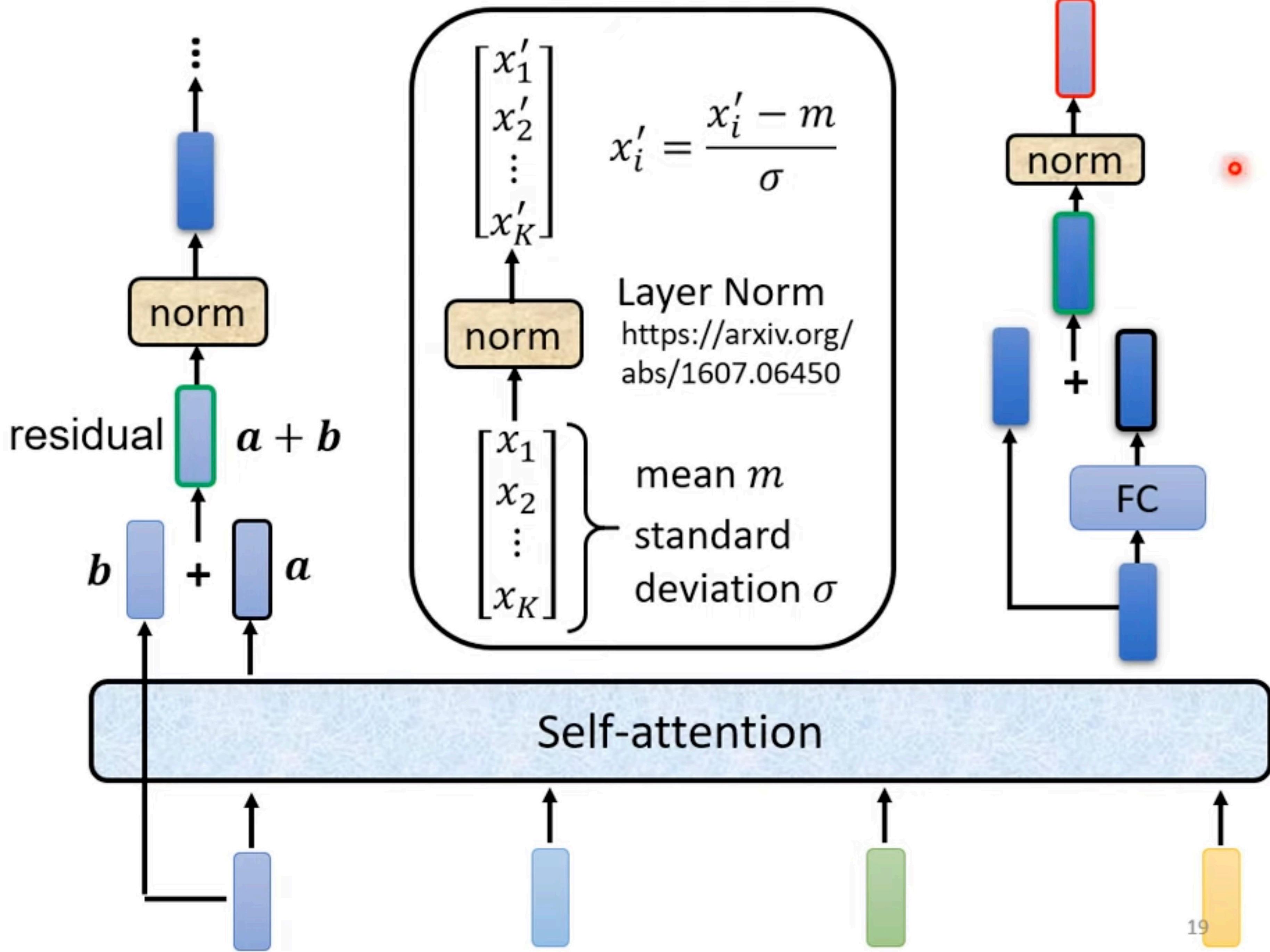
You can use **RNN** or **CNN**.



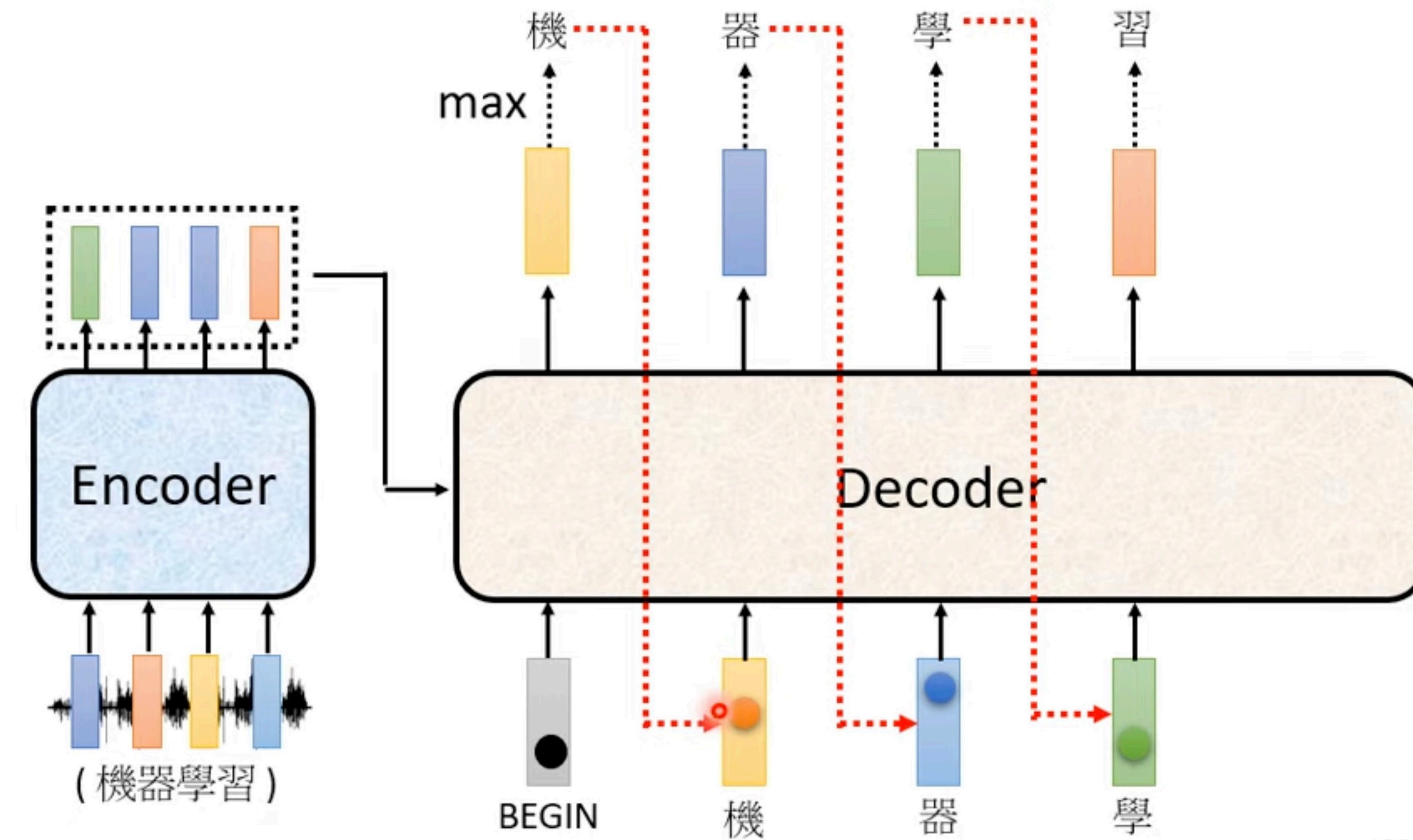
Transformer's Encoder



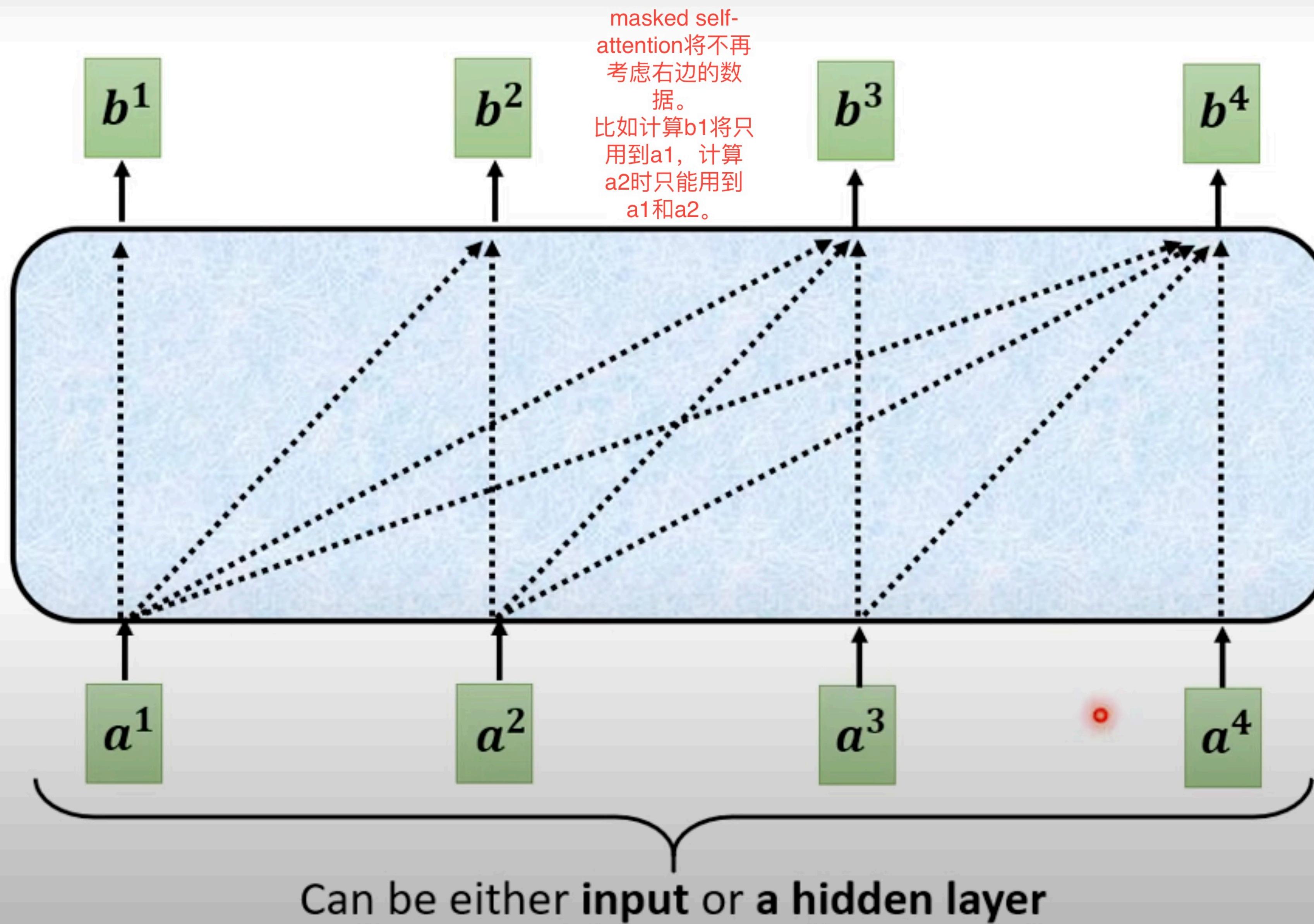




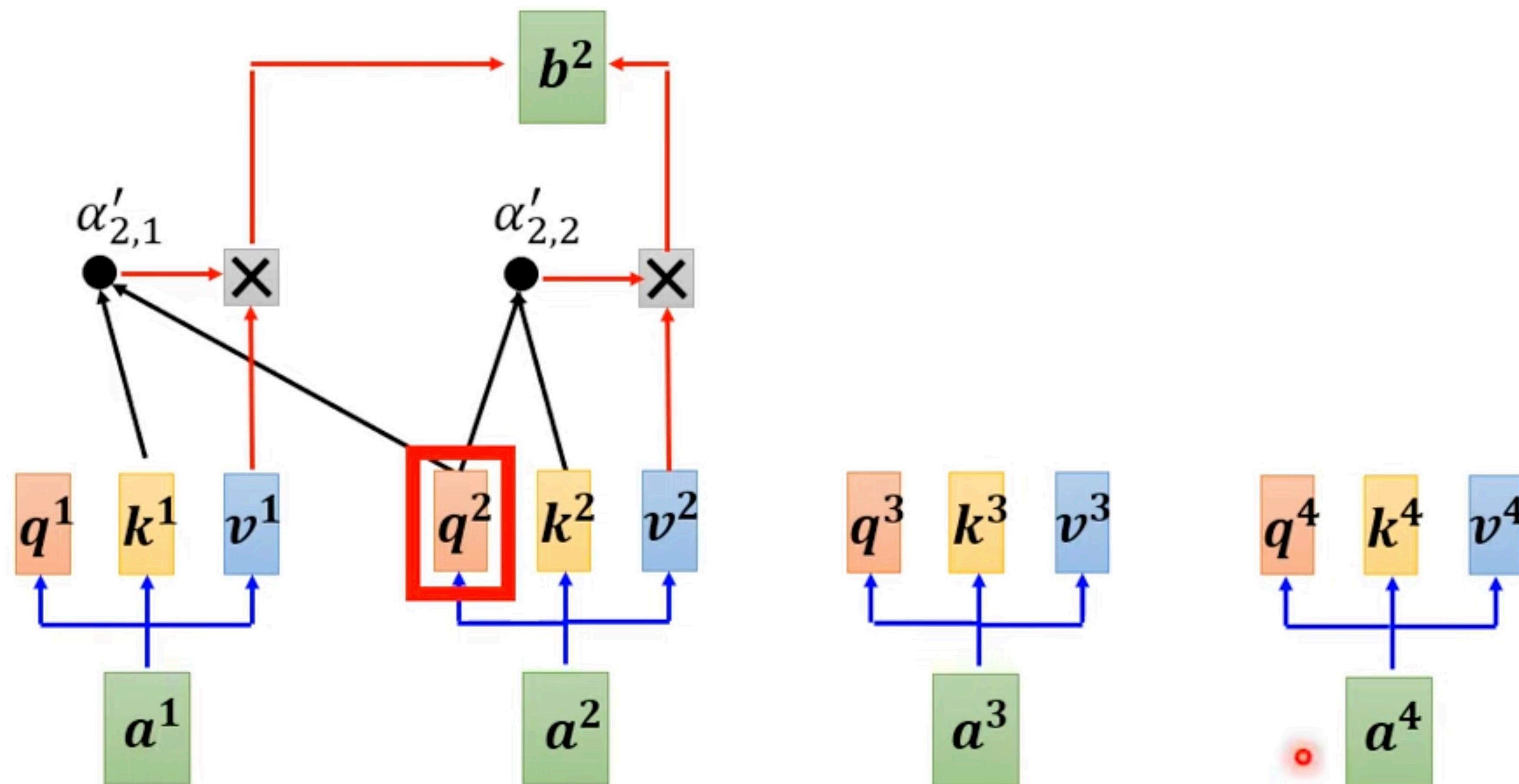
Autoregressive



Self-attention → Masked Self-attention

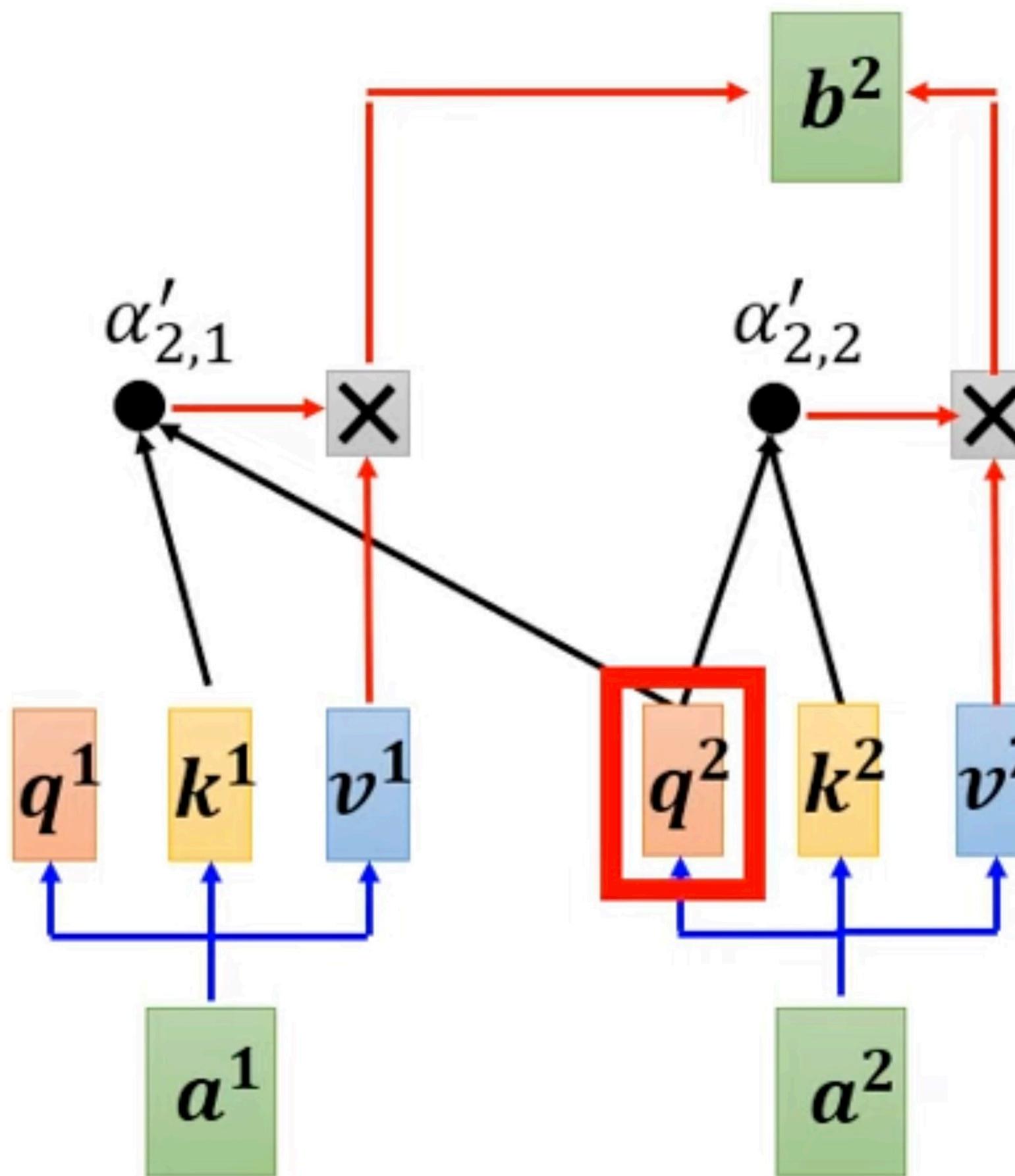


Self-attention → Masked Self-attention



Why masked?

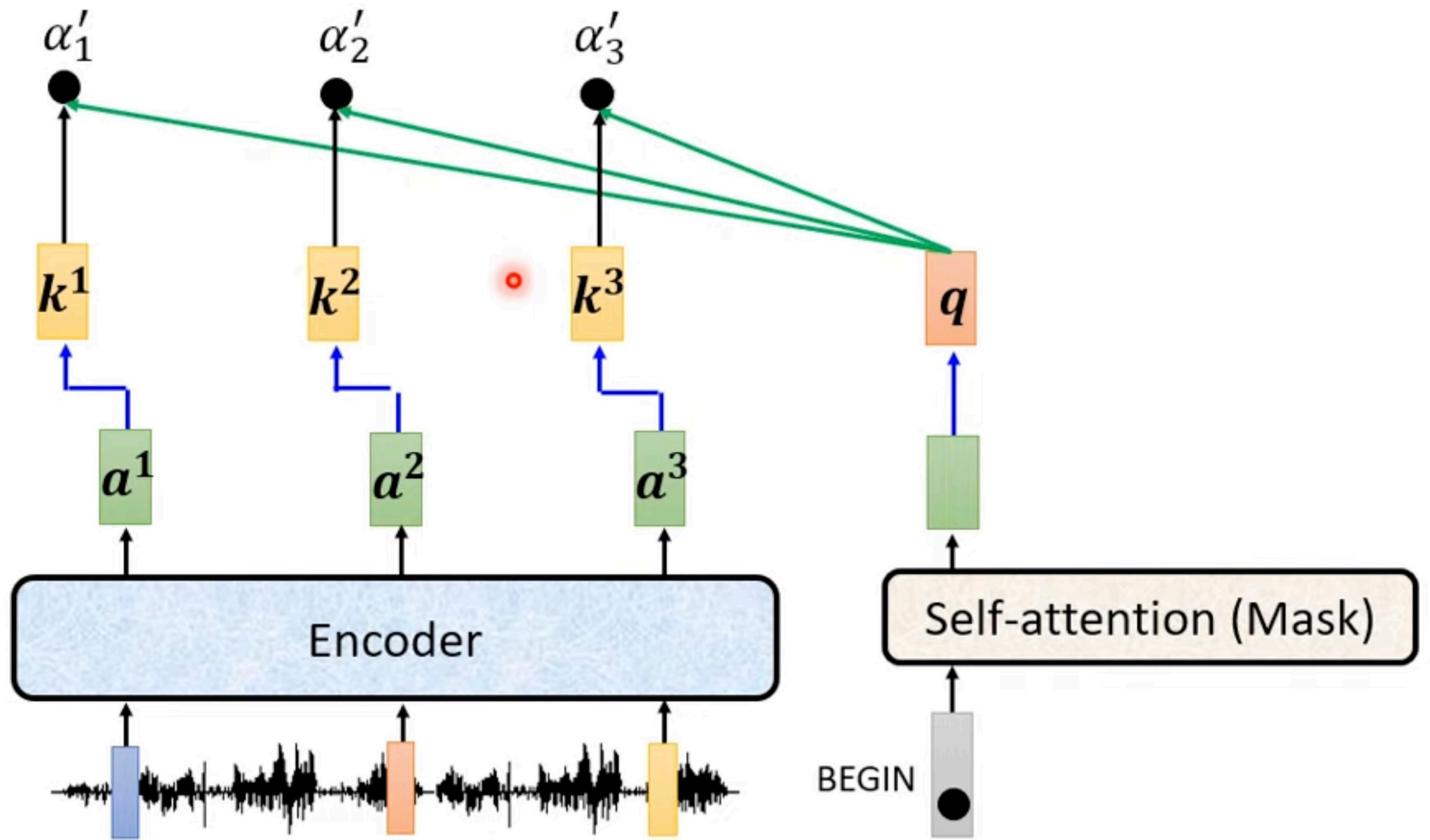
Self-attention → Masked Self-attention



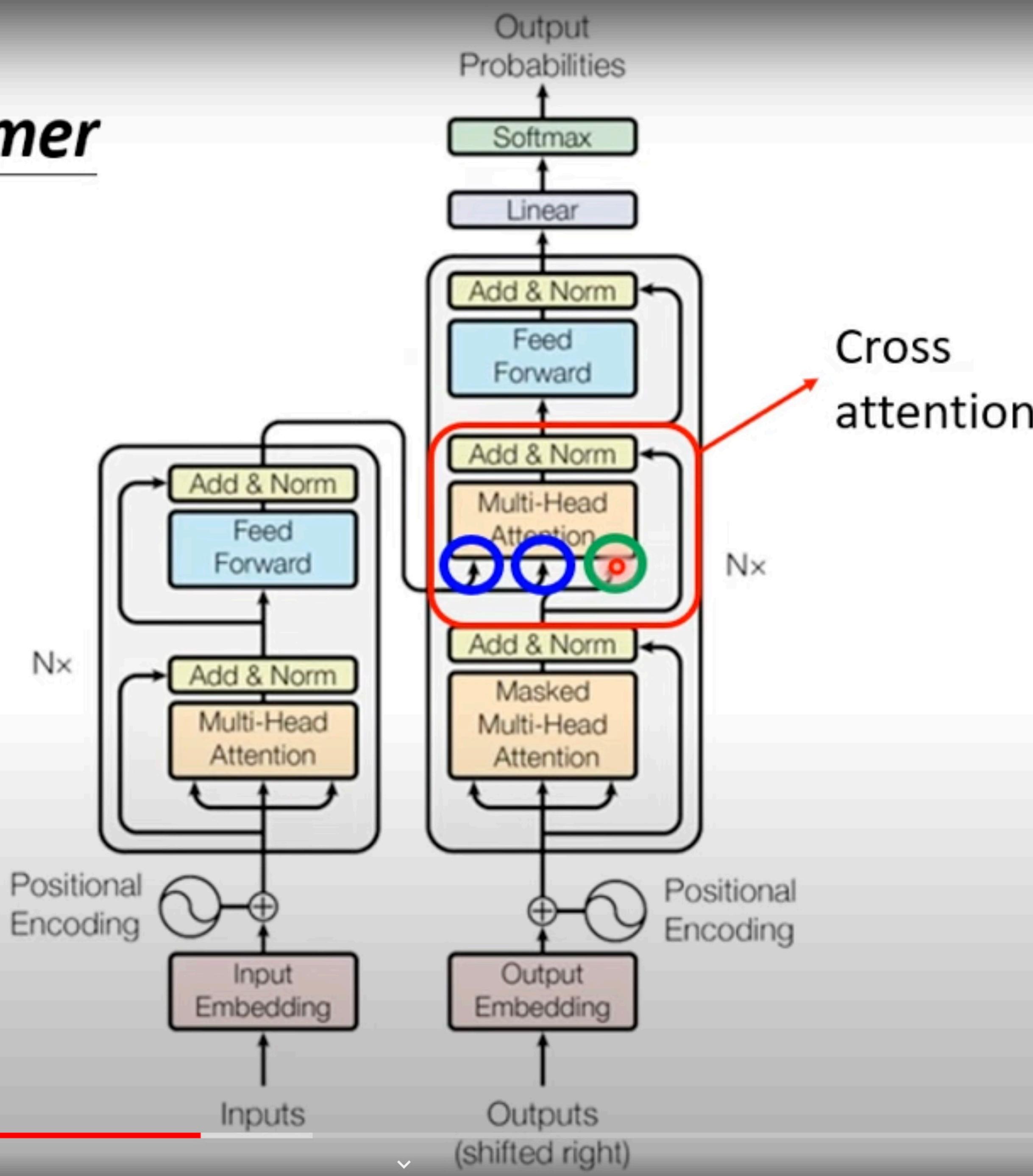
Why masked? :
由于decoder的输入是一个一个产生，因此在计算b2时也没有办法得到a3、a4。

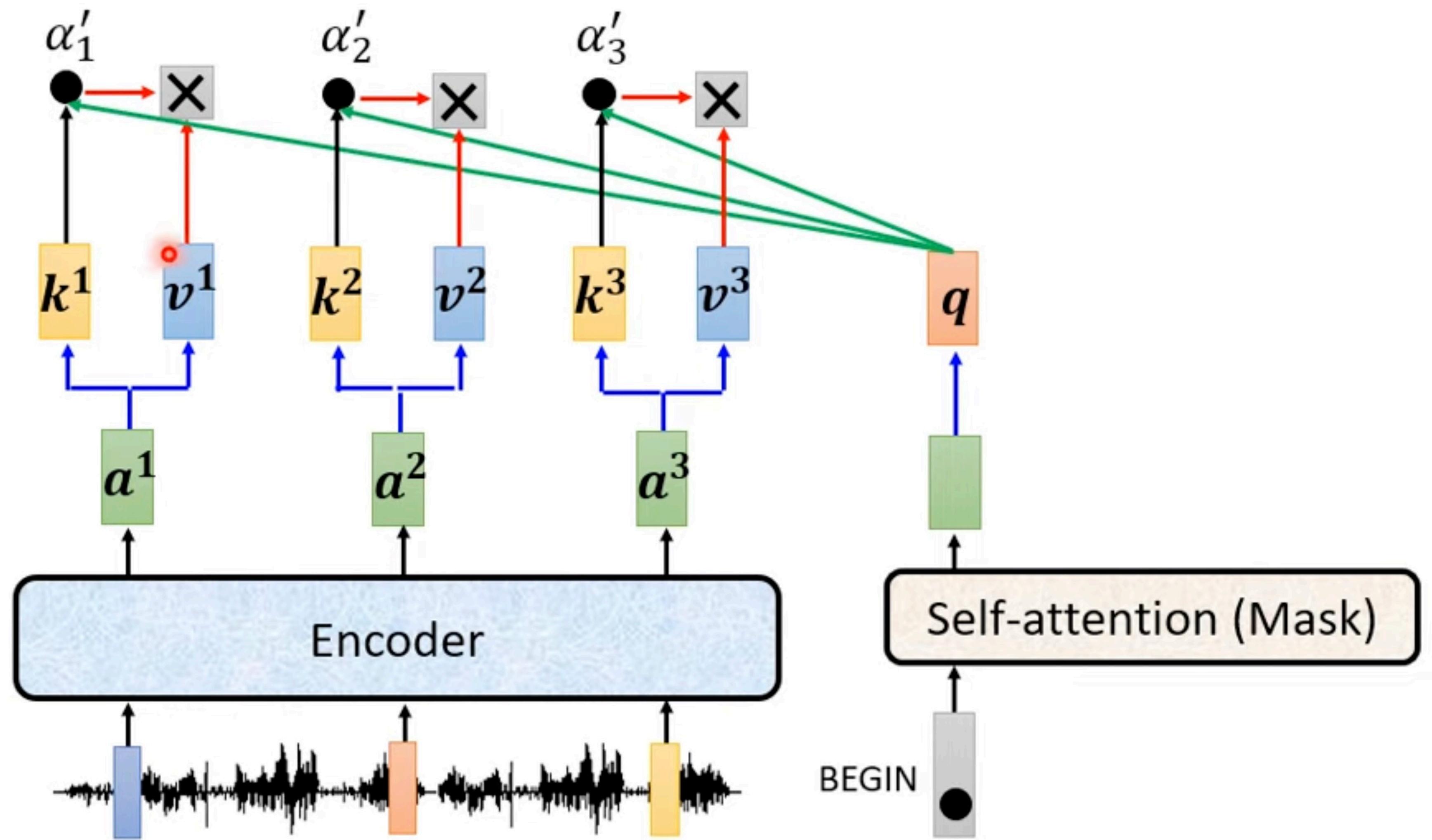
◦

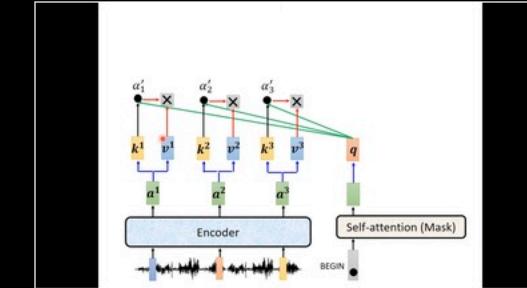
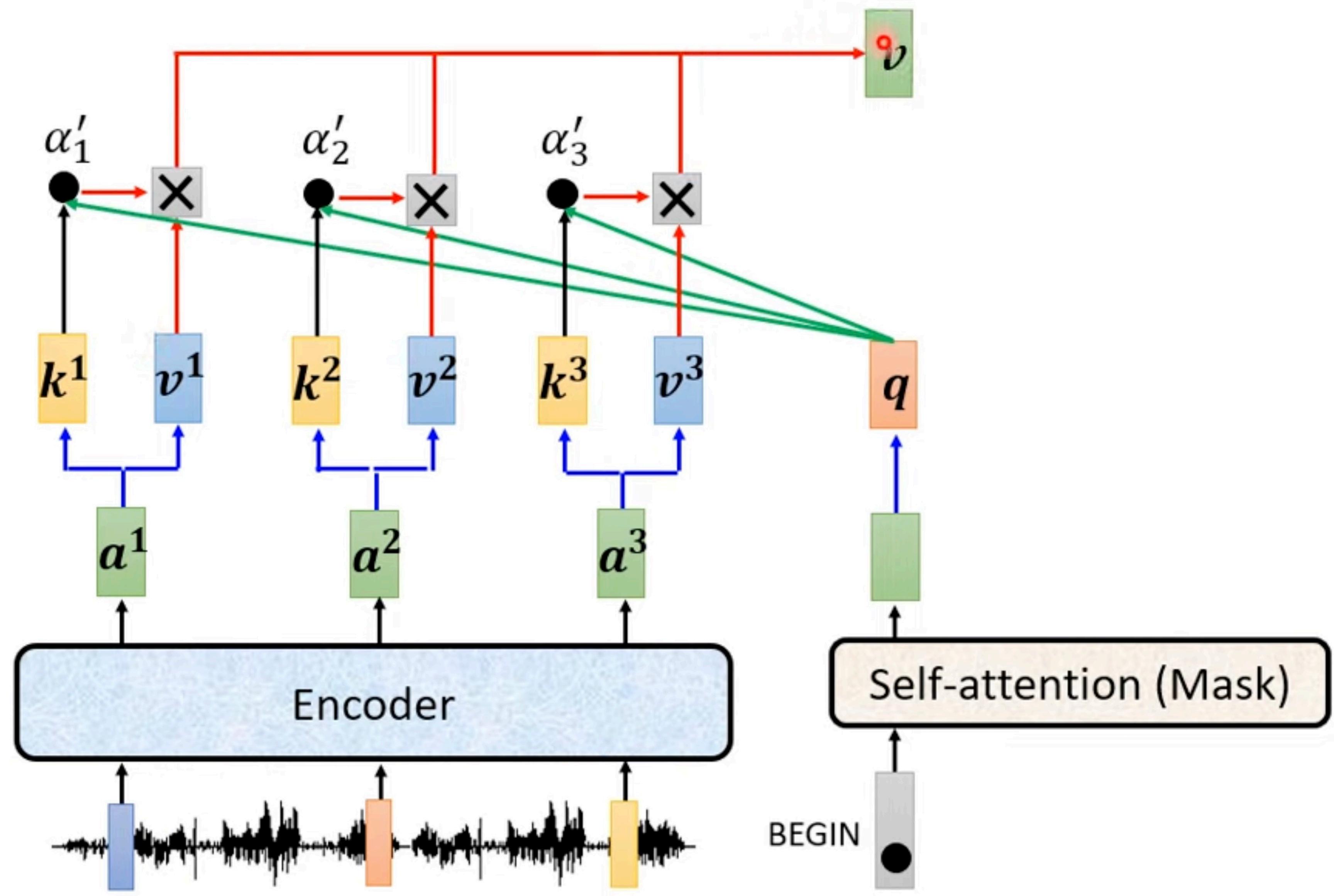
Why masked? Consider how does decoder work

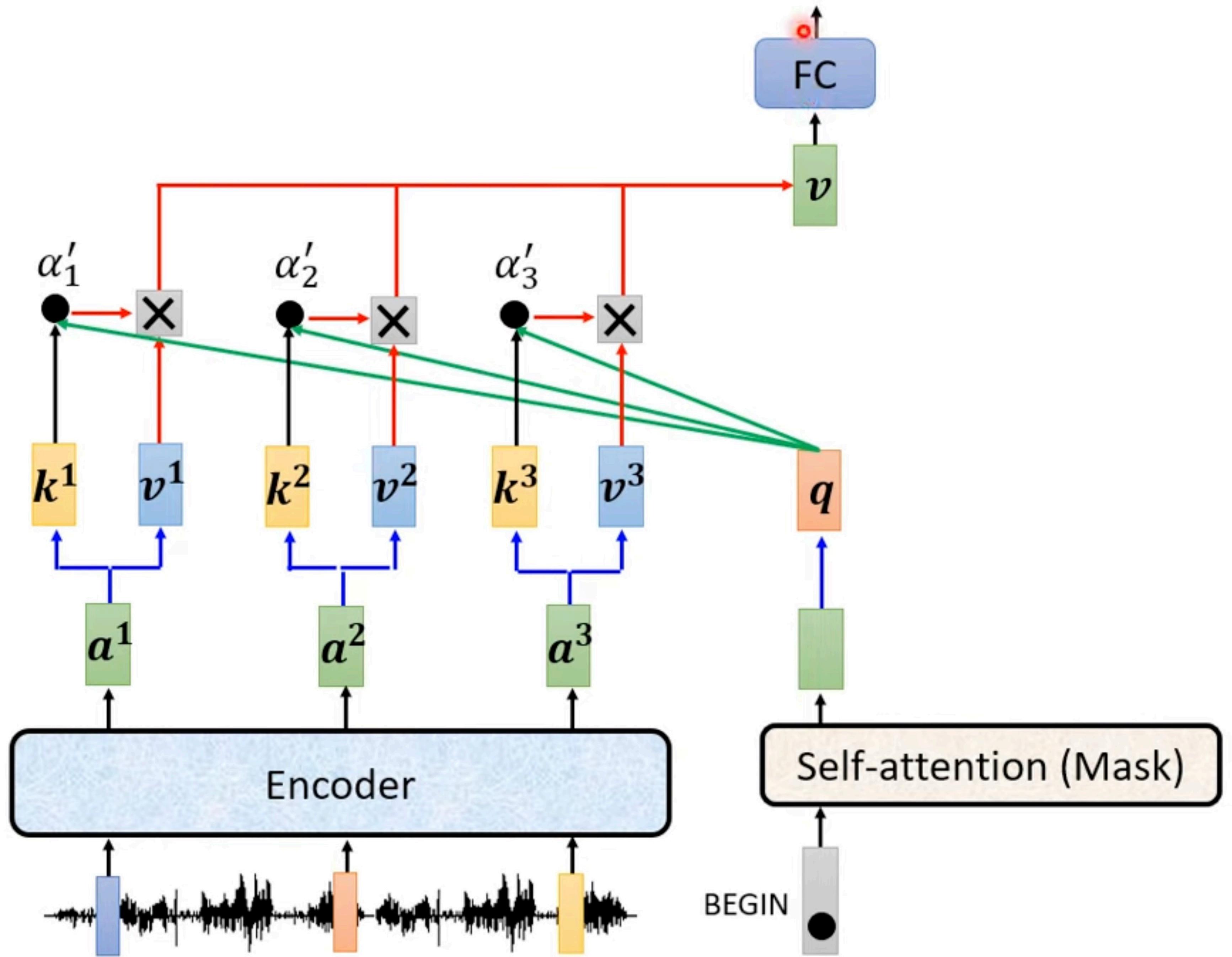


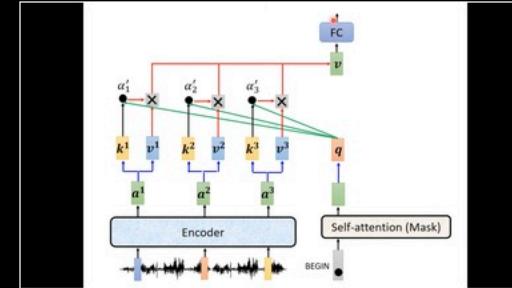
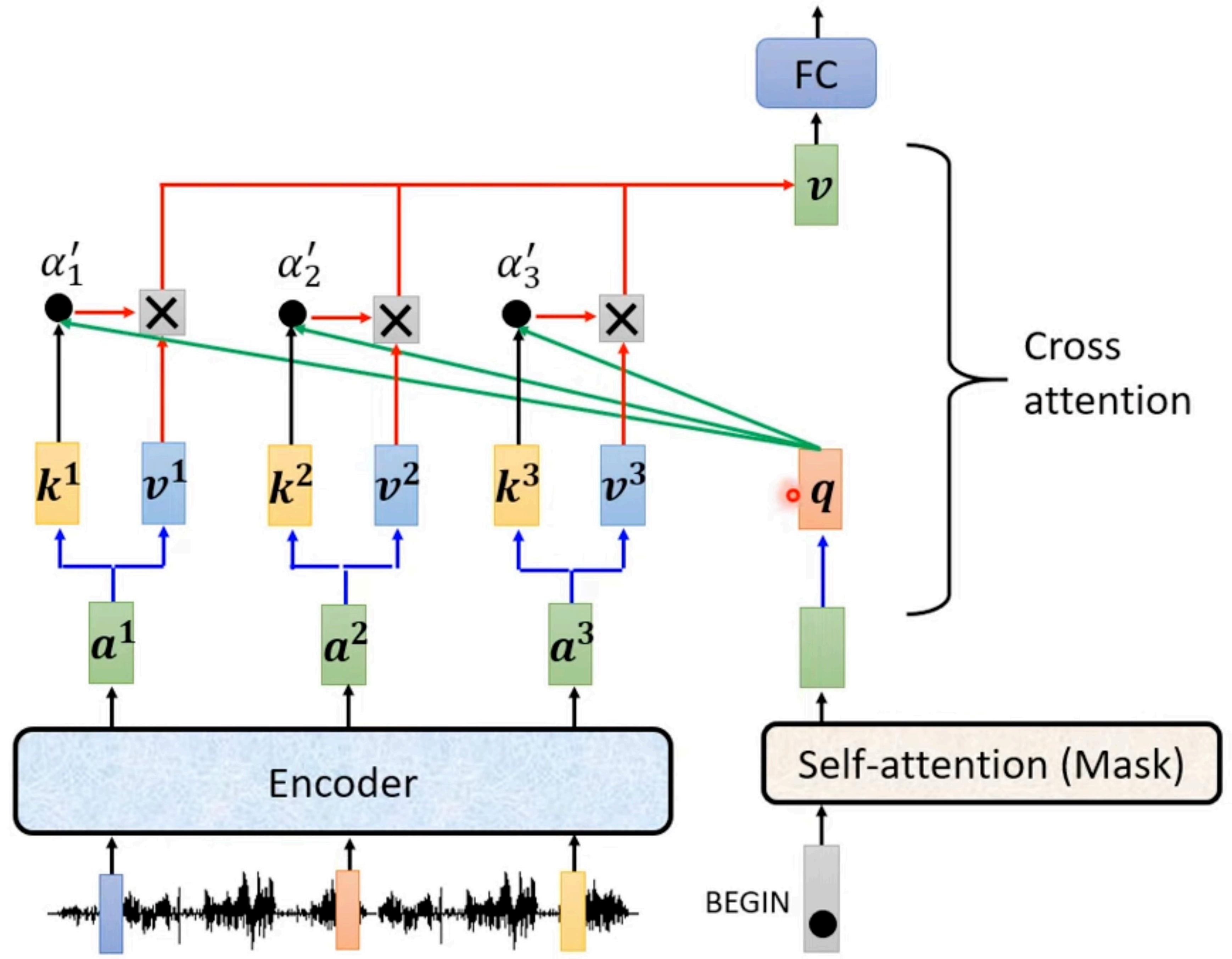
Transformer

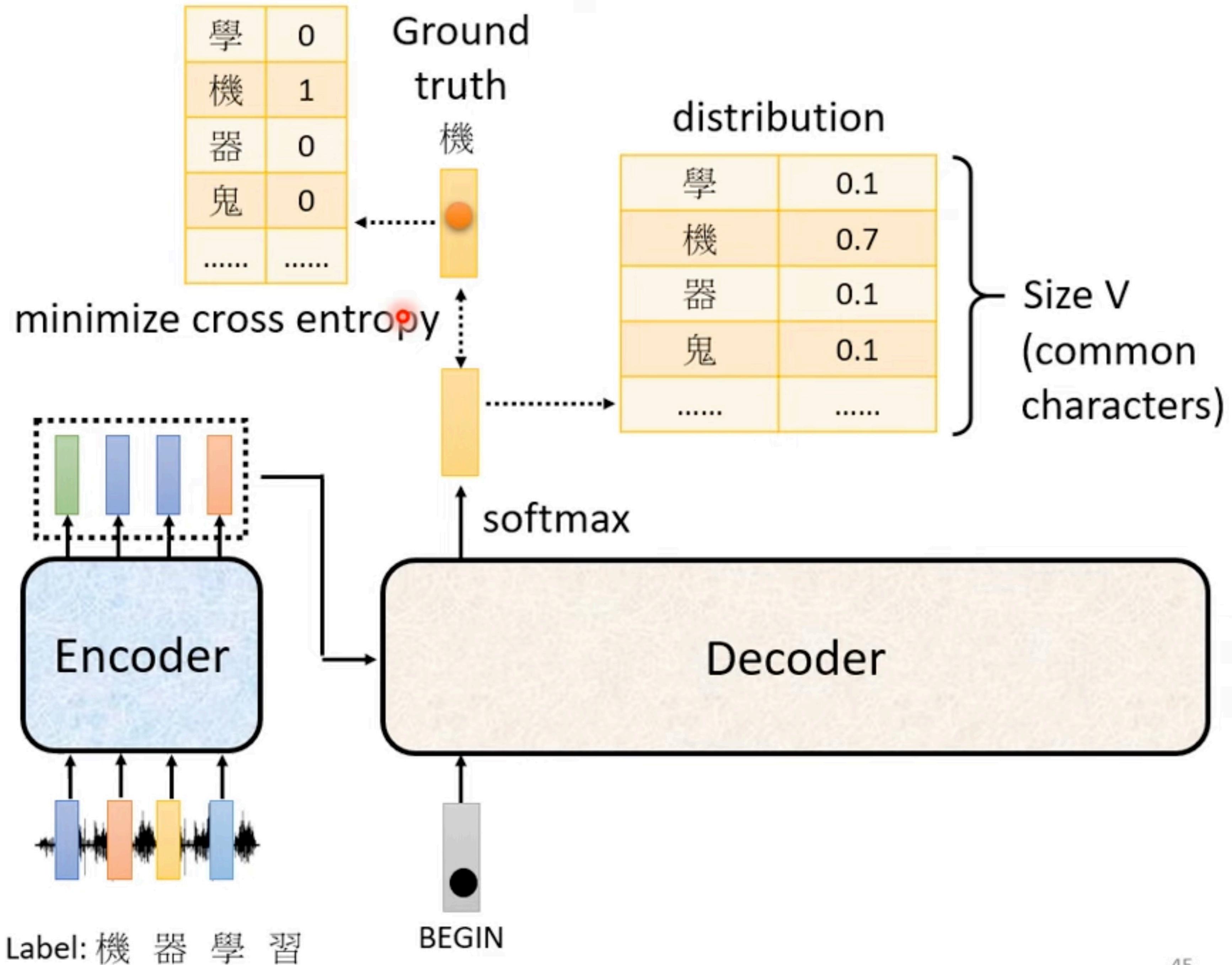


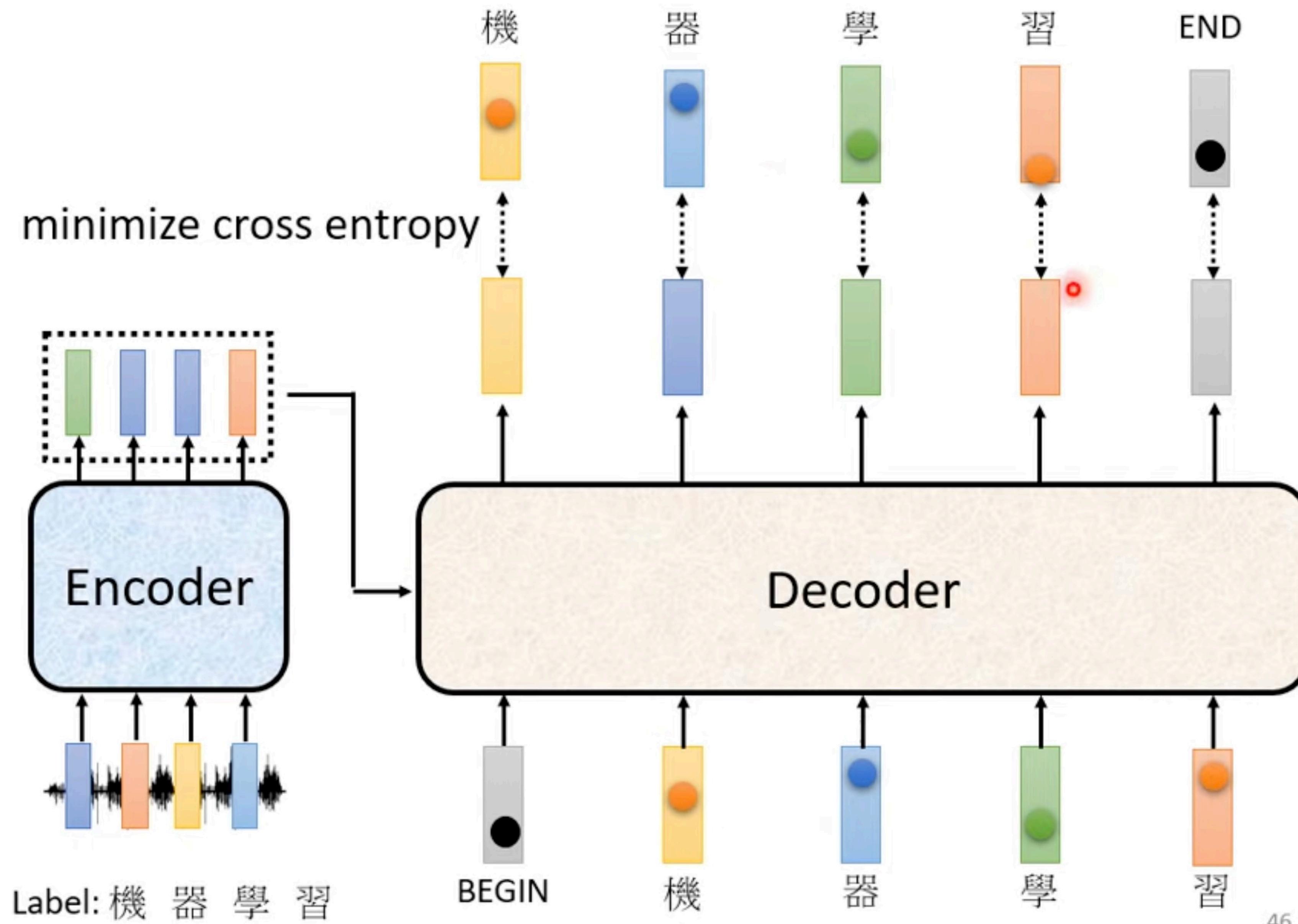




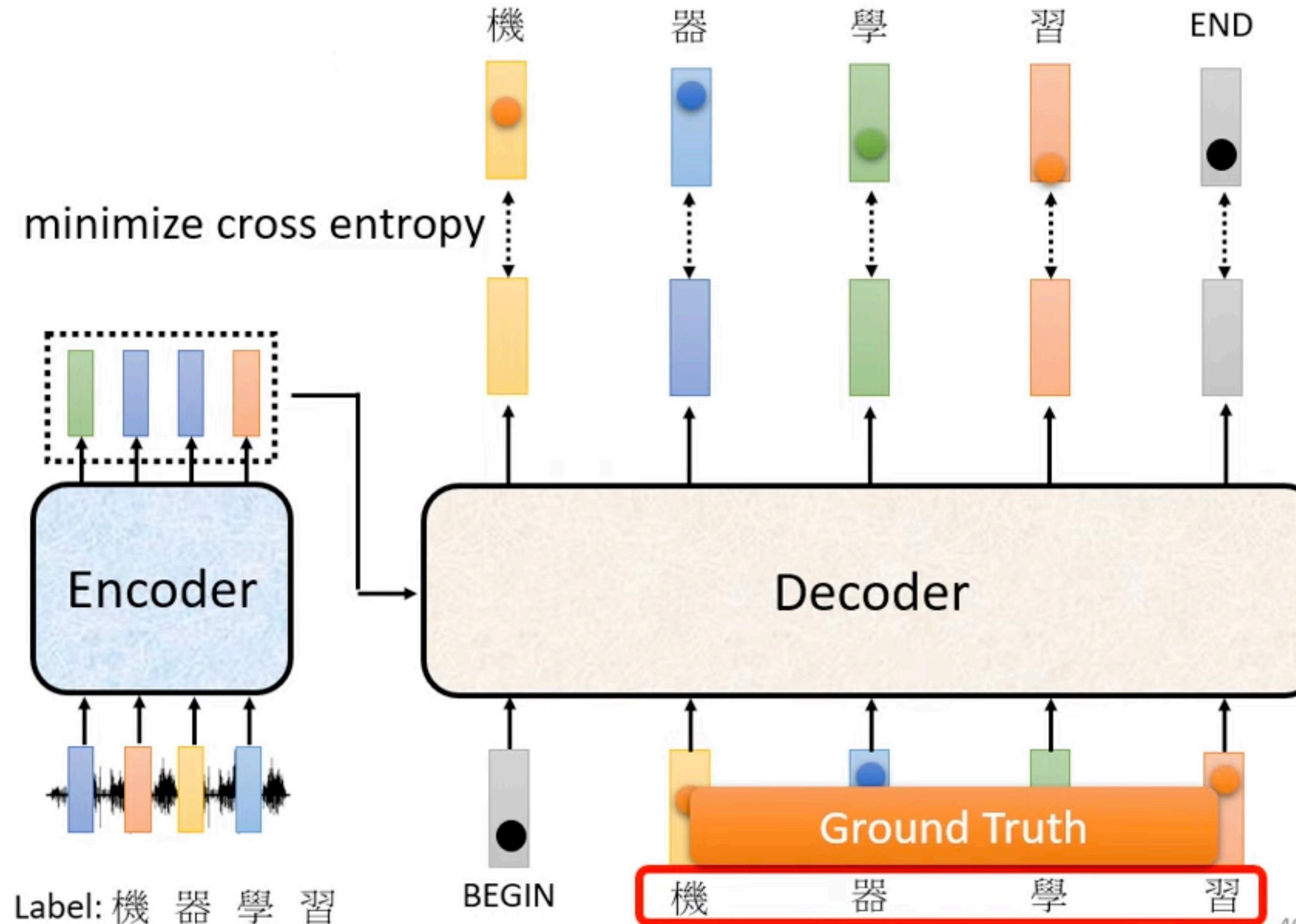








Teacher Forcing: using the ground truth as input.



Optimizing Evaluation Metrics?

