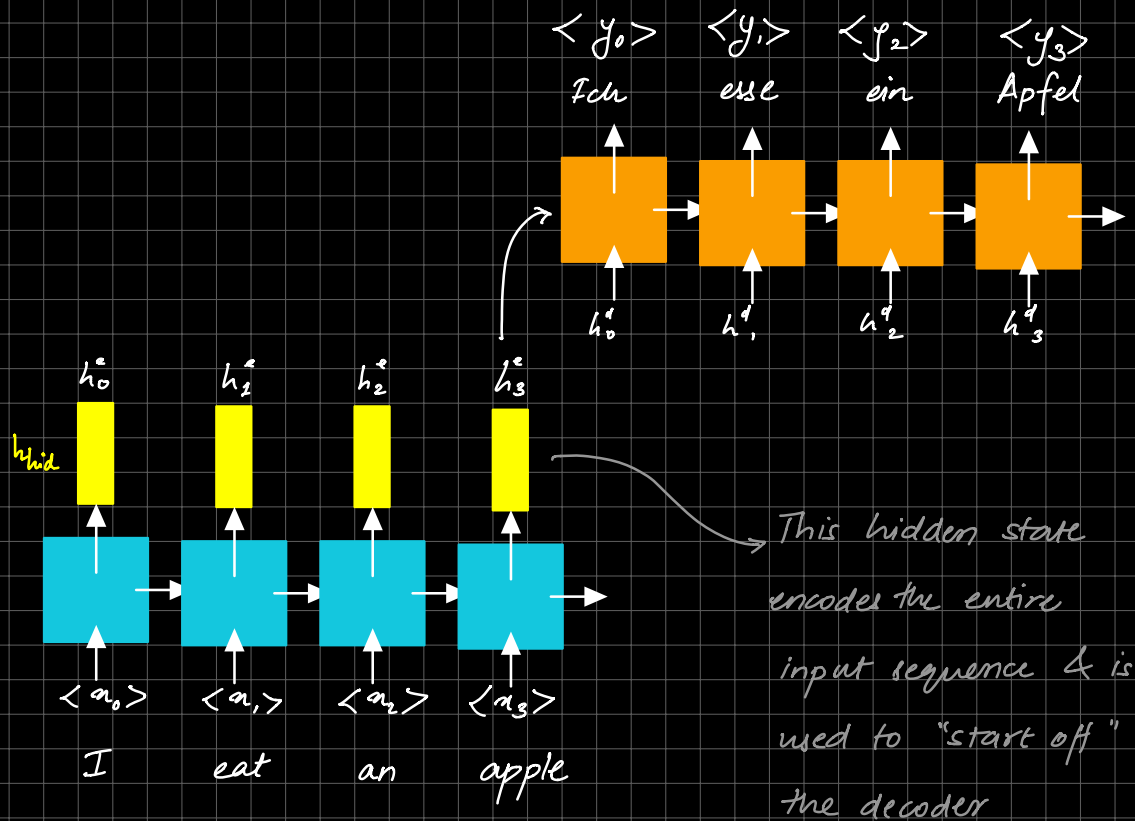
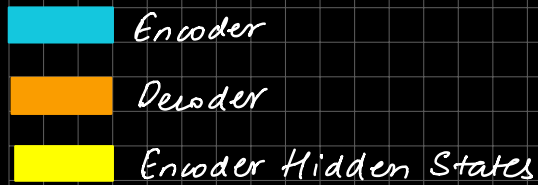
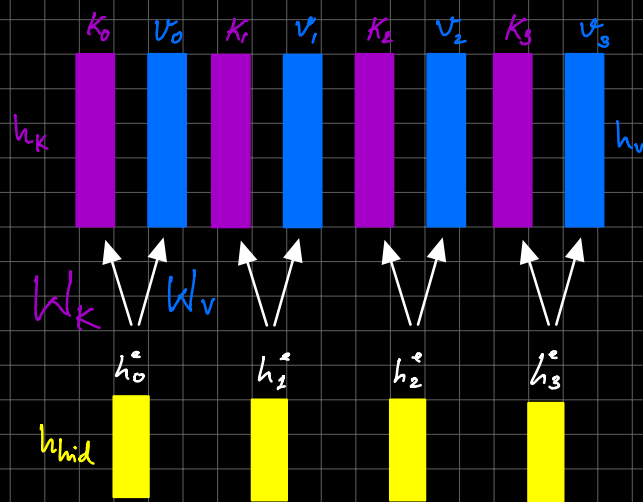


\* Recitation 8 : Attention :



## "Construct Keys & Values from Encoder Hidden States"



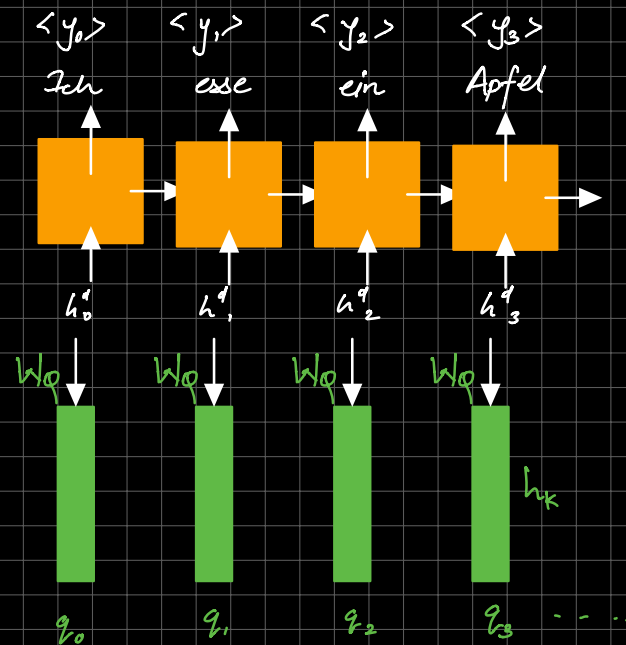
- Using matrices  $W_k$  and  $W_v$ , we construct **keys** & **values** for all the Encoder Hidden States

- Matrix Dimensions:

$$W_k = h_{hid} \times h_k$$

$$W_v = h_{hid} \times h_v$$

## "Construct Queries from Decoder Hidden States"

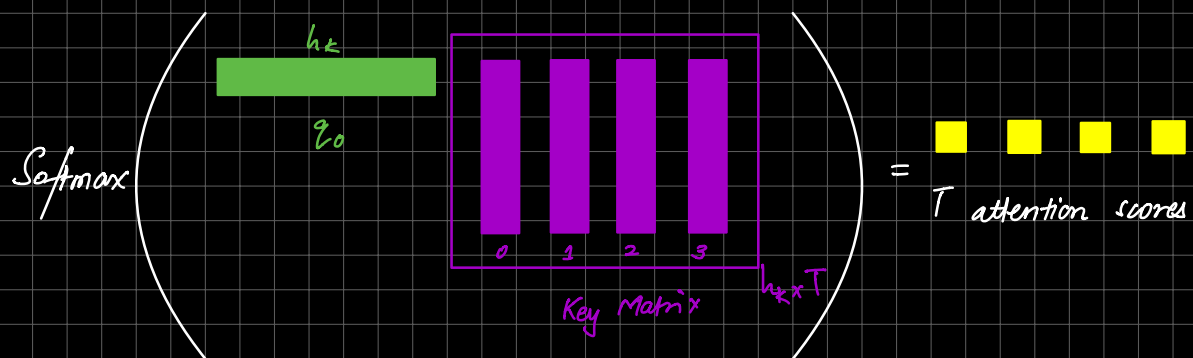


- Using matrix  $W_Q$ , we construct a query from the Decoder hidden State.

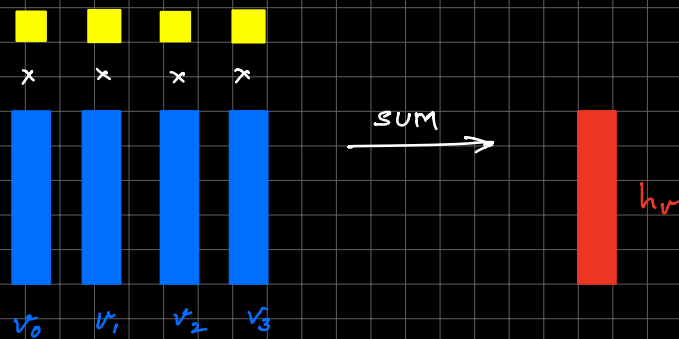
$$W_Q = h_{hid} \times h_k$$

- NOTE: Queries & Keys are of the same dimension in this example but this is not true for all types of attention.

"Compute Attention Scores"

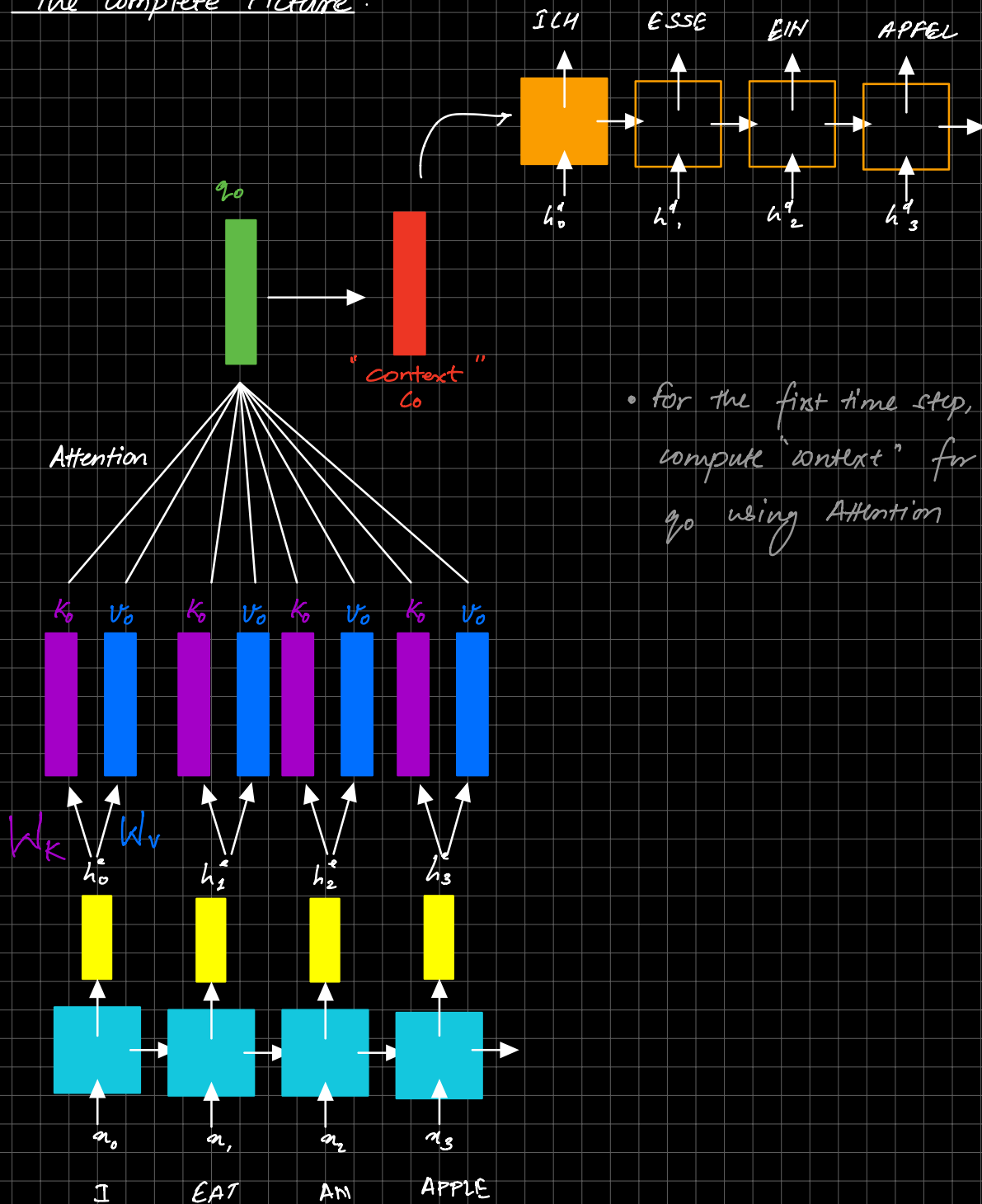


"Scale values by attention scores"

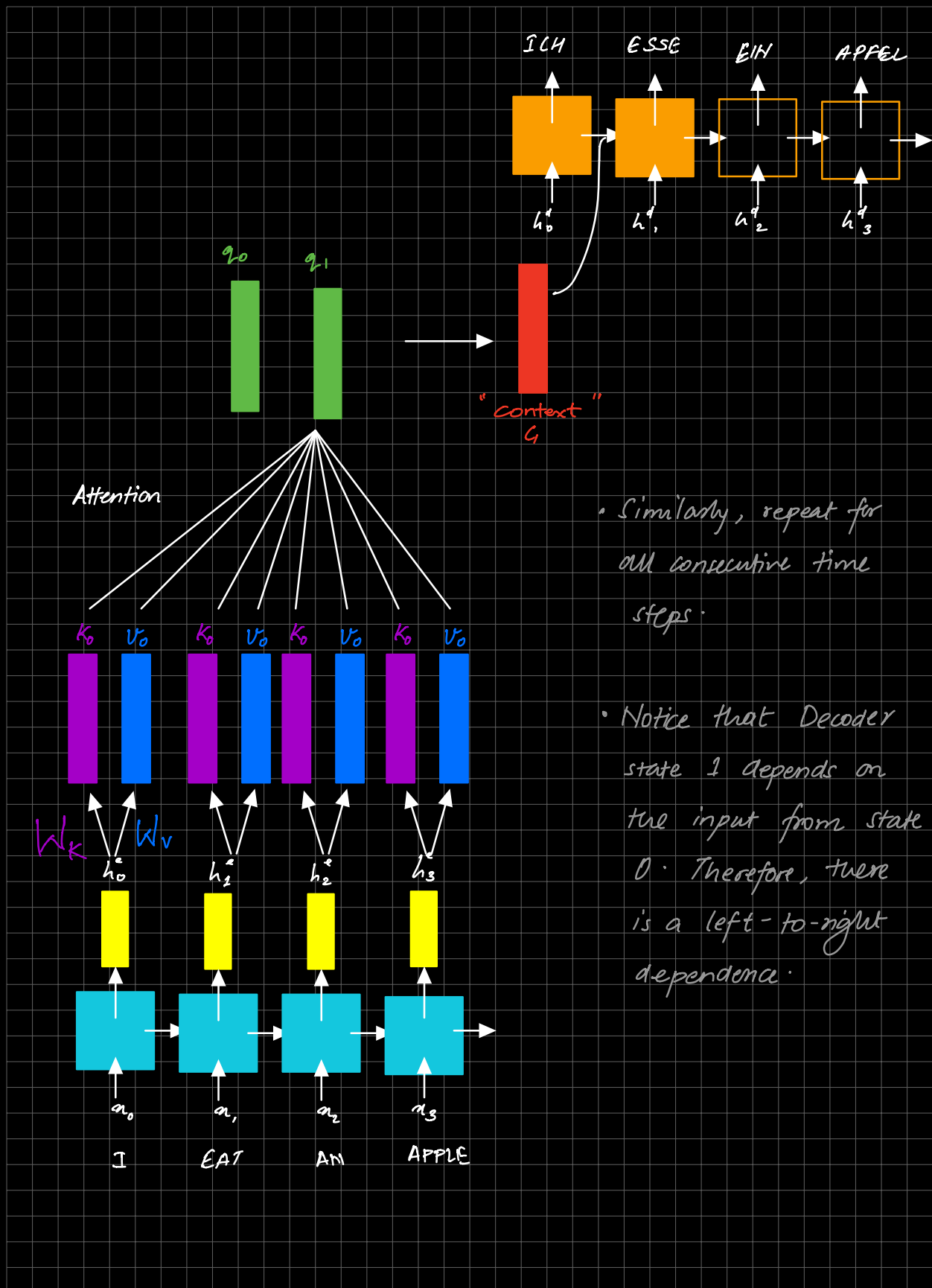


FINAL REPRESENTATION  
FOR ONE DECODER STATE

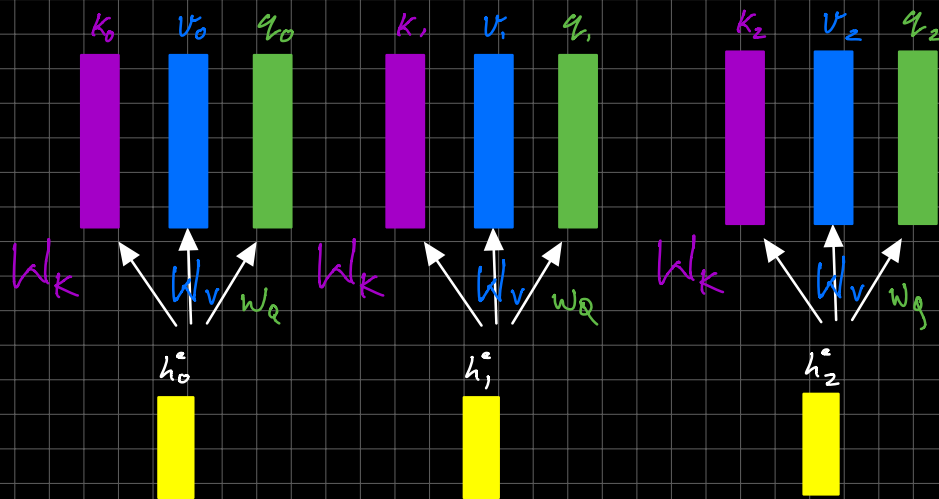
# The Complete Picture:



- for the first time step, compute 'context' for  $q_0$  using Attention



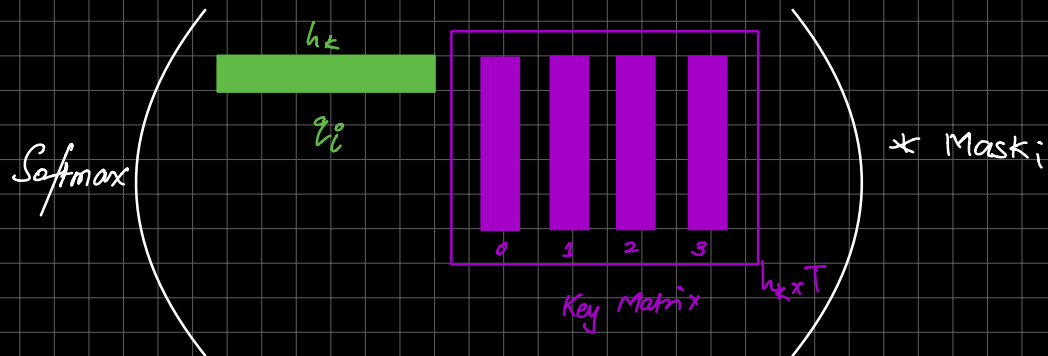
## "Self Attention"



- Keys, Values, and Queries are all computed from Encoder Hidden States.
- Compute Attention between queries, Keys & values exactly like before.
- BUT, no left-to-right dependence, so can compute all context vectors in parallel.

## "Masking in Transformer Decoder"

- want to condition only on previous tokens' context.  $\therefore$  "Mask" out the keys from consecutive tokens while computing attention.



- Mask for  $q_0 = [1 \ 0 \ 0 \ 0]$

- Mask for  $q_1 = [1 \ 1 \ 0 \ 0]$

- Mask for  $q_2 = [1 \ 1 \ 1 \ 0]$

- Mask for  $q_3 = [1 \ 1 \ 1 \ 1]$

\* To compute all the attention scores at once :

