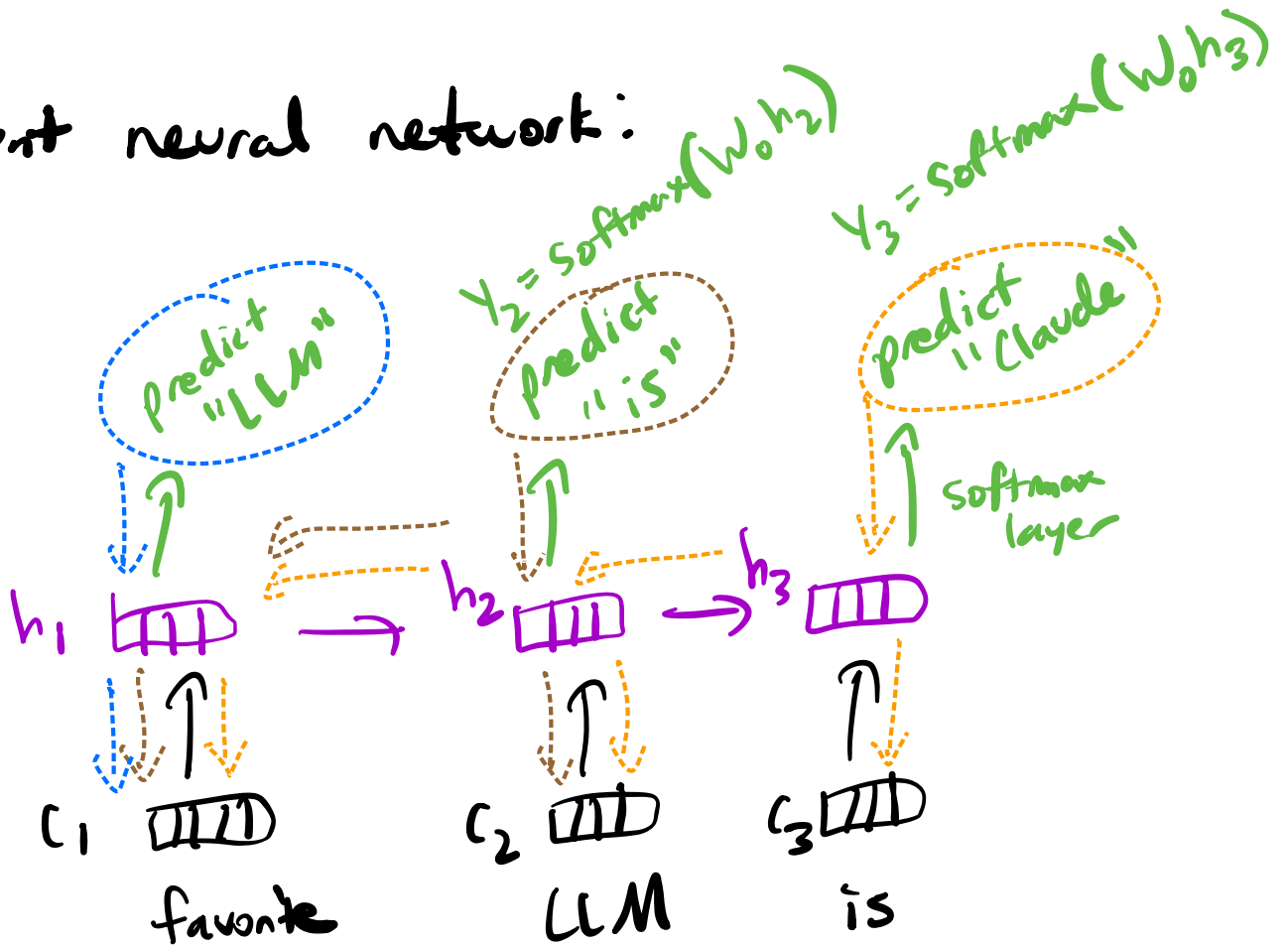


Recurrent neural network:



$$L_3 = -\log p(\text{claudio} \mid \text{"favorite LLM is"})$$

$$L_2 = -\log p(\text{is} \mid \text{favorite LLM})$$

$$L_1 = -\log p(\text{LLM} \mid \text{favorite})$$

$$L = \frac{L_3 + L_2 + L_1}{3} \quad \left. \vphantom{\frac{L_3 + L_2 + L_1}{3}} \right\} \begin{array}{l} \text{ave. token-level} \\ \text{neg. log likelihood} \end{array}$$

batch:

<SOS> my favorite LLM is claudio <EOS>

<SOS> i am a zebra <EOS>

i

issues w/ RNNs:

↳ vanishing / exploding gradient

$$\begin{aligned} h_2 &= W_h h_1 + W_c c_2 \\ h_3 &= W_h h_2 + W_c c_3 \\ &\vdots \end{aligned} \quad \left. \vphantom{\begin{aligned} h_2 &= W_h h_1 + W_c c_2 \\ h_3 &= W_h h_2 + W_c c_3 \\ &\vdots \end{aligned}} \right\} \text{linear RNN}$$

$$h_3 = W_h W_h h_1$$

$$h_n = W_h^{n-1} h_1$$

↳ bottleneck

↳ $h_t \Rightarrow$ current hidden state

↳ entire prefix is represented with a single vector

↳ cannot parallelize computation of h_t

↳ h_t depends on h_{t-1}

attention mechanism:

↳ originated as a way to alleviate bottleneck in RNNs

↳ Bahdanau, Cho et al 2014

↳ Transformer, Vaswani et al 2017

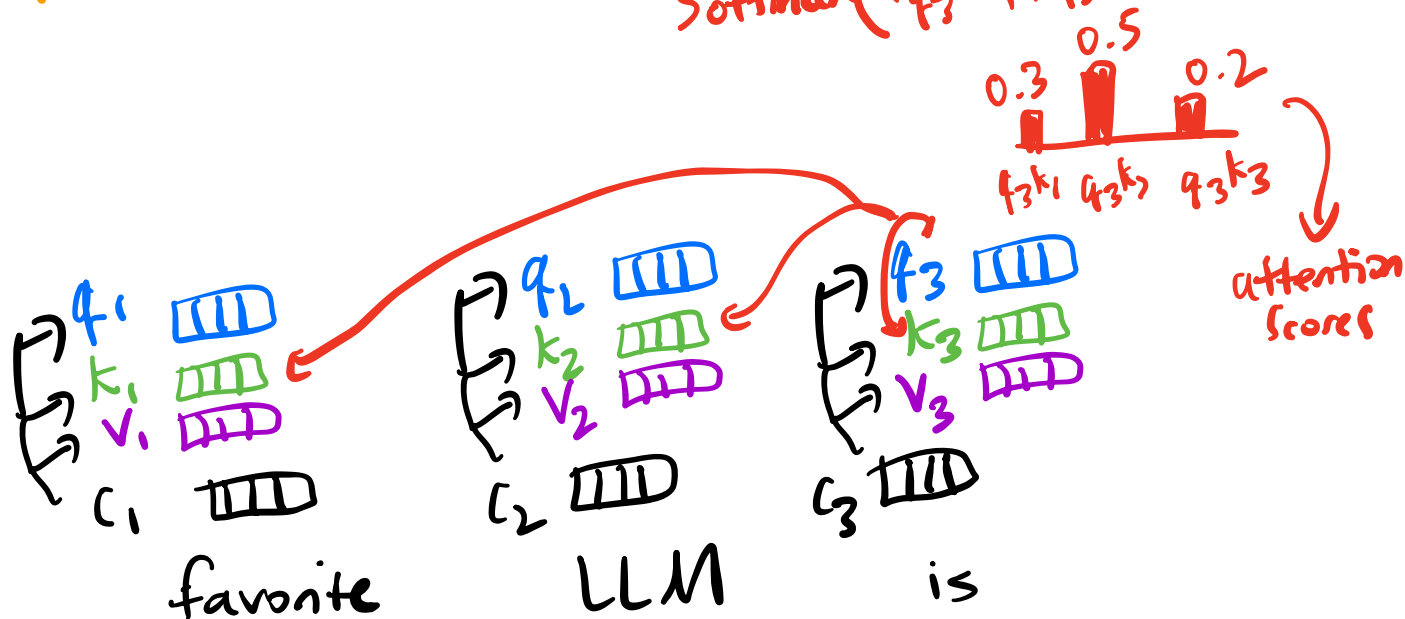
↳ drop recurrence

↳ "self-attention"

self-attention

computation of h_3

$h_3 = \text{[vector]} \rightarrow \text{Softmax layer} \rightarrow \text{predict claude}$
 $h_3 = 0.3v_1 + 0.5v_2 + 0.2v_3$
 $\text{Softmax}(\langle q_3 \cdot k_1, q_3 \cdot k_2, q_3 \cdot k_3 \rangle)$



query $q_1 = f(W_q c_1)$
key $k_1 = f(W_k c_1)$
value $v_1 = f(W_v c_1)$

} linear / FF layers
 $f = \text{ReLU}$

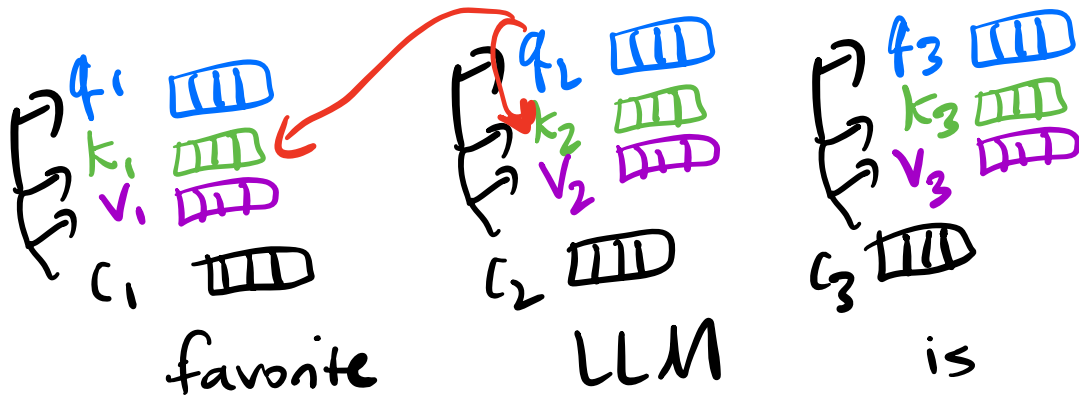
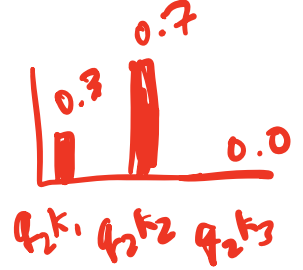
$q_2 = f(W_q c_2)$

Computation of h_2

softmax
layer predict "is"

$$h_2 = 0.3 \cdot v_1 + 0.7 v_2 \rightarrow$$

$$\text{softmax}(\langle q_2 \cdot k_1, q_2 \cdot k_2 \rangle)$$



↳ in self-attn, no dependency between h_t and h_{t-1}

parallelizing self-attention:

q_1  k_1 

q_2  k_2 

q_3  k_3 

attn vectors:

$$a_1 : \langle q_1 \cdot k_1 \rangle$$

$$a_2 : \langle q_2 \cdot k_1, q_2 \cdot k_2 \rangle$$

$$a_3 : \langle q_3 \cdot k_1, q_3 \cdot k_2, q_3 \cdot k_3 \rangle$$




step 1:

q_1
 q_2
 q_3

 \times

k_1	k_2	k_3

 $=$

	k_1	k_2	k_3
q_1	.		
q_2	.	.	
q_3	.	.	.

these cells contain info about the future, need to mask

step 2:

Softmax

	k_1	k_2	k_3
q_1			
q_2			
q_3			

	1	0	0
	1	1	0
	1	1	1

\rightarrow PyTorch
mask-fill
(0, -10)

↓

Softmax

	k_1	k_2	k_3
q_1	.	-10	-10
q_2	.	.	-10
q_3	.	.	.

step 3

	k_1	k_2	k_3
q_1	1.0	0	0
q_2	0.3	0.7	0
q_3	0.3	0.5	0.2

 \times

v_1	
v_2	
v_3	

$$=$$

h_1	
h_2	
h_3	