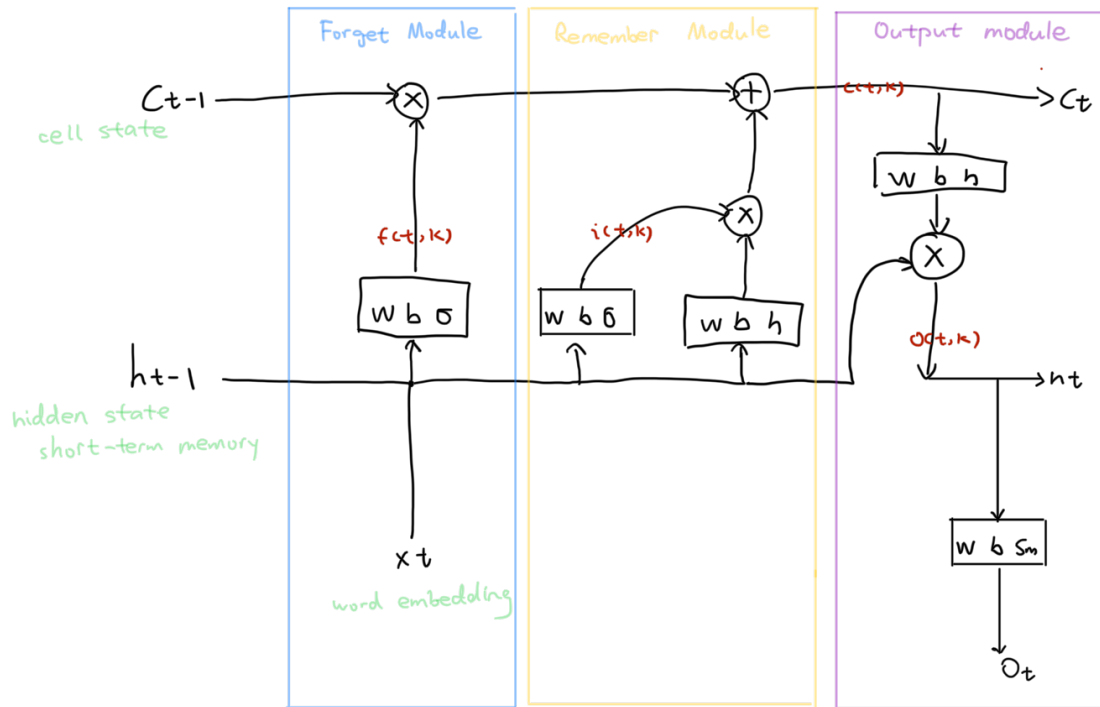


# LSTM

- Plain LSTM Cell



- Input

one word at a time in the order

① document delimiter query

delimiter: |||

② query delimiter document

each document query pair is passed as a long sequence

- Architecture

input  $x(t)$

output  $y(t)$

$$y(t) = y'(t, 1) || \dots || y'(t, k)$$

output for entire LSTM is the concatenation of outputs from each layer

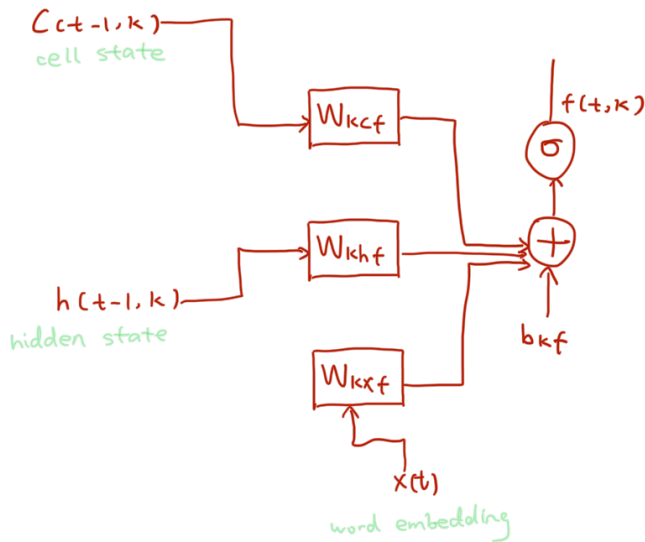
in our two layer deep LSTM Reader, output of the entire LSTM reader is the concatenation of output from each layer

$$x'(t, k) = x(t) || y'(t, k-1)$$

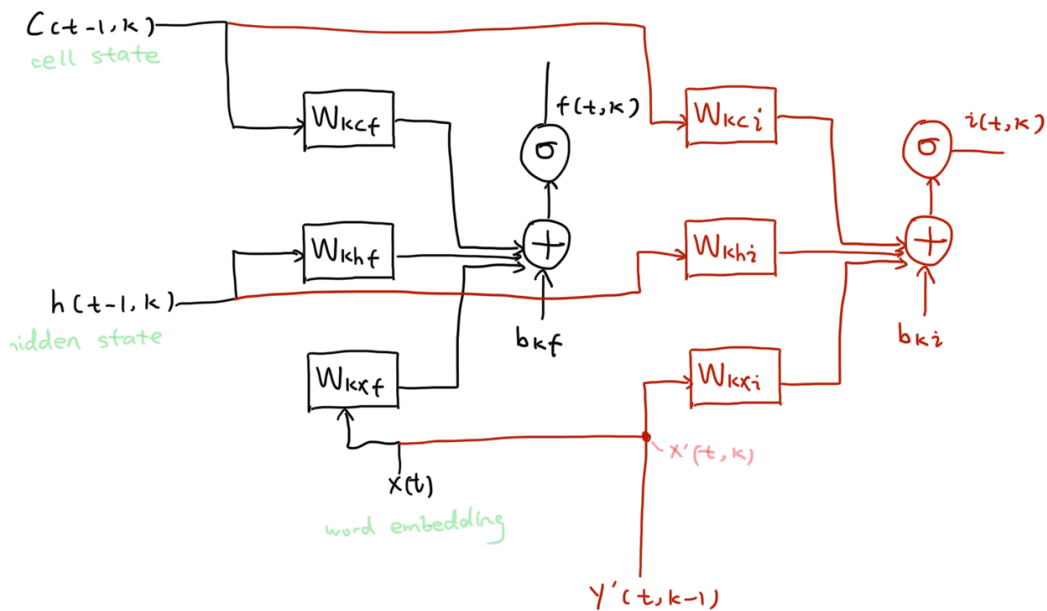
input for each hidden layer of LSTM Reader is the concatenation of original input to the LSTM Reader (embedded document query pair) and the output from previous layer

the output from previous layer

$$f(t, k) = \sigma(W_{kxf} x(t) + W_{khf} h(t-1, k) + W_{kcf} c(t-1, k) + b_{kf})$$

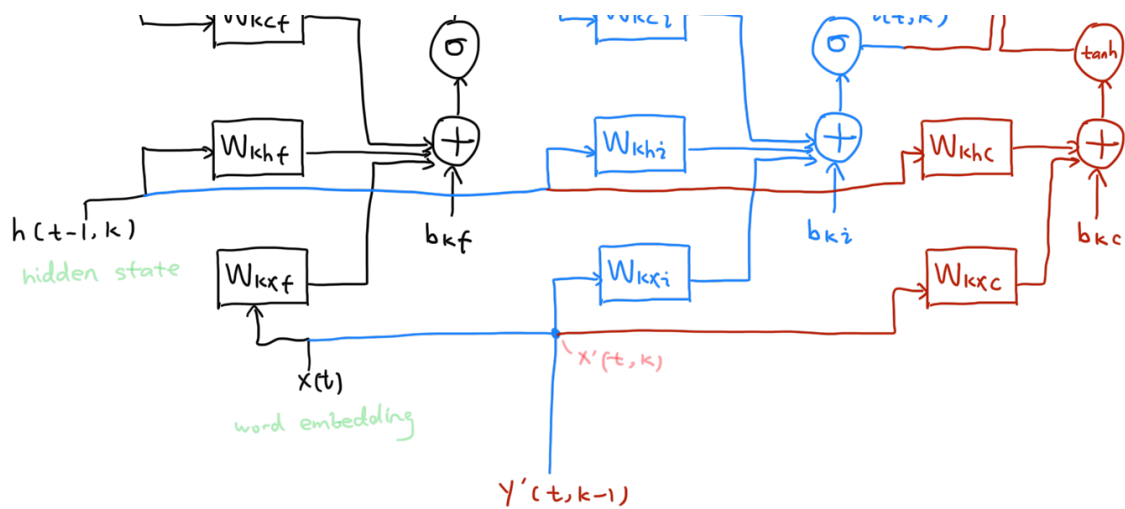


$$i(t, k) = \sigma(W_{kxi} x'(t, k) + W_{khi} h(t-1, k) + W_{kci} c(t-1, k) + b_{ki})$$

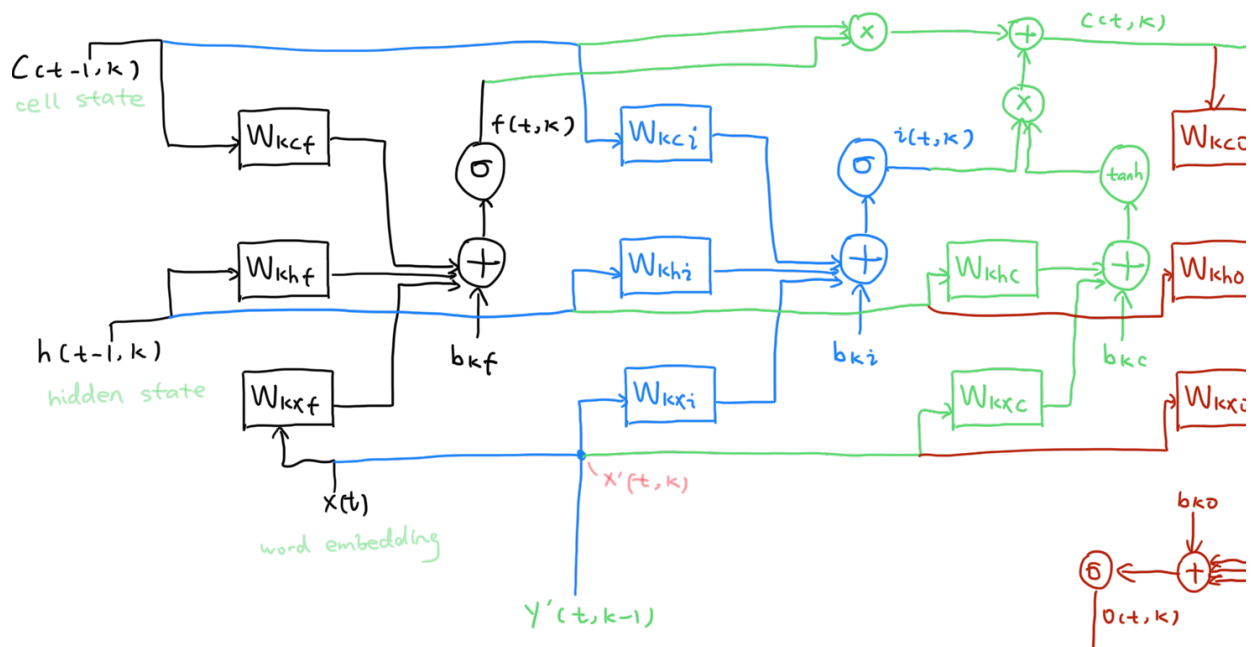


$$c(t, k) = f(t, k) c(t-1, k) + i(t, k) \tanh(W_{kxc} x'(t, k) + W_{khc} h(t-1, k) + b_{kc})$$

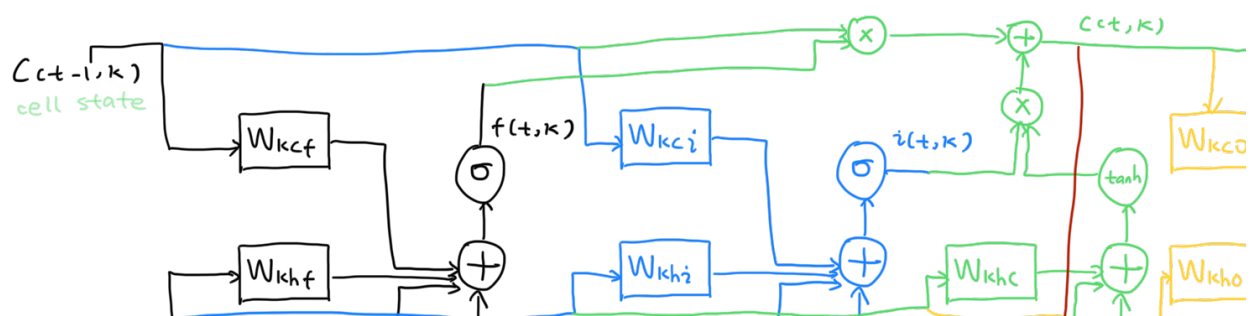


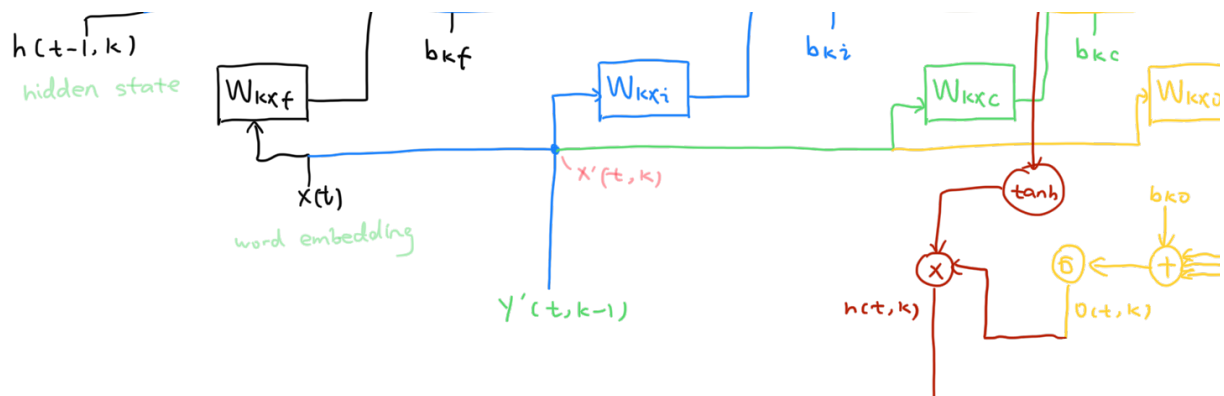


$$O(t, k) = \sigma(W_{kxo} x'(t, k) + W_{kxo} h(t-1, k) + W_{kxo} c(t, k) + b_{ko})$$

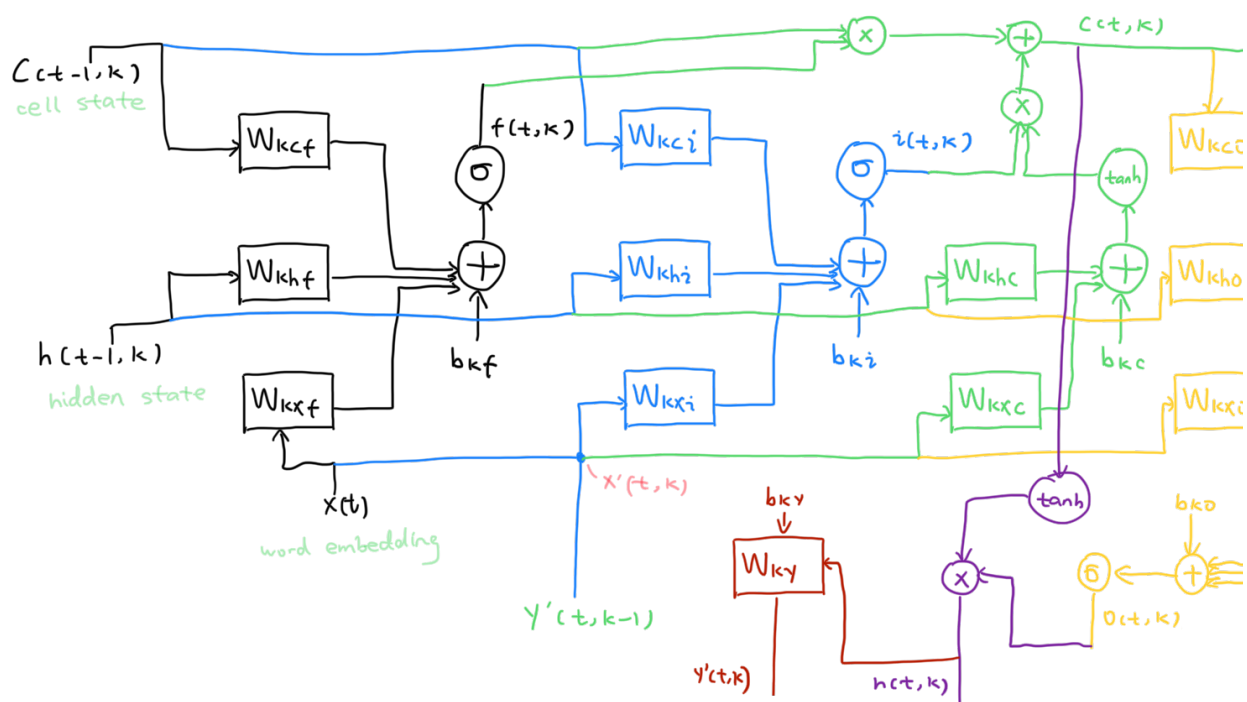


$$h(t, k) = O(t, k) \tanh(c(t, k))$$





$$y'(t, k) = W_{ky} h(t, k) + b_{ky}$$



$$g^{LSTM}(d, q) = y(|d| + |q|)$$

$$\text{Recall } y(t) = y'(t, 1) || \dots || y'(t, k)$$

### Model Hyperparameters

• indicates optimal configuration

Hidden Size	Learning Rate	Batch Size	Dropout	depth # layers
64	1-3	16	0.0	
128	$5 \times 10^{-4}$		0.1	
256	$10^{-4}$	32	0.2	
	$5 \times 10^{-5}$			

uri

Input Form

cqa (context document, delimiter, question)

qca (question, delimiter, context document)

RMS Drop momentum 0.9 decay 0.95

• Model Summary

