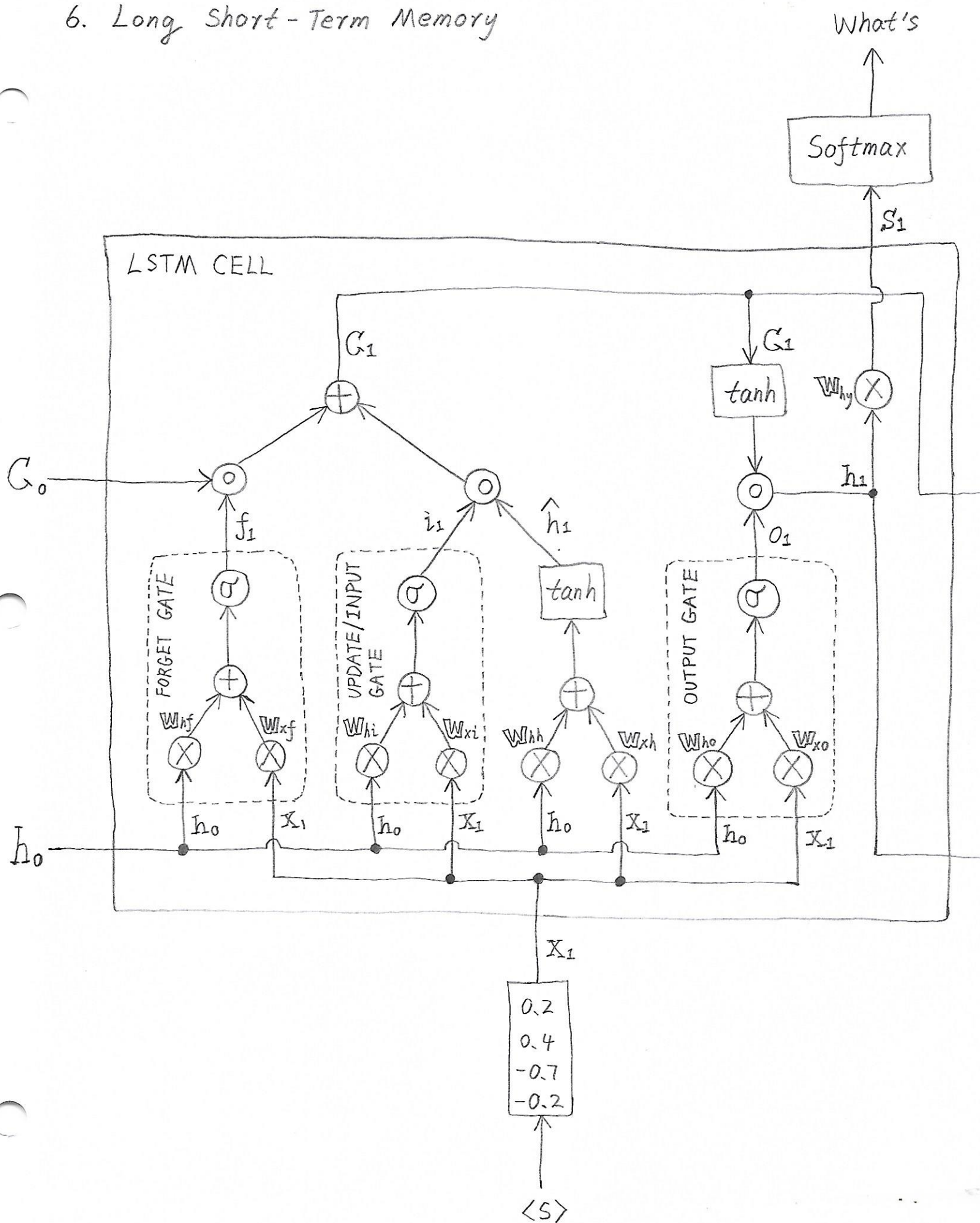
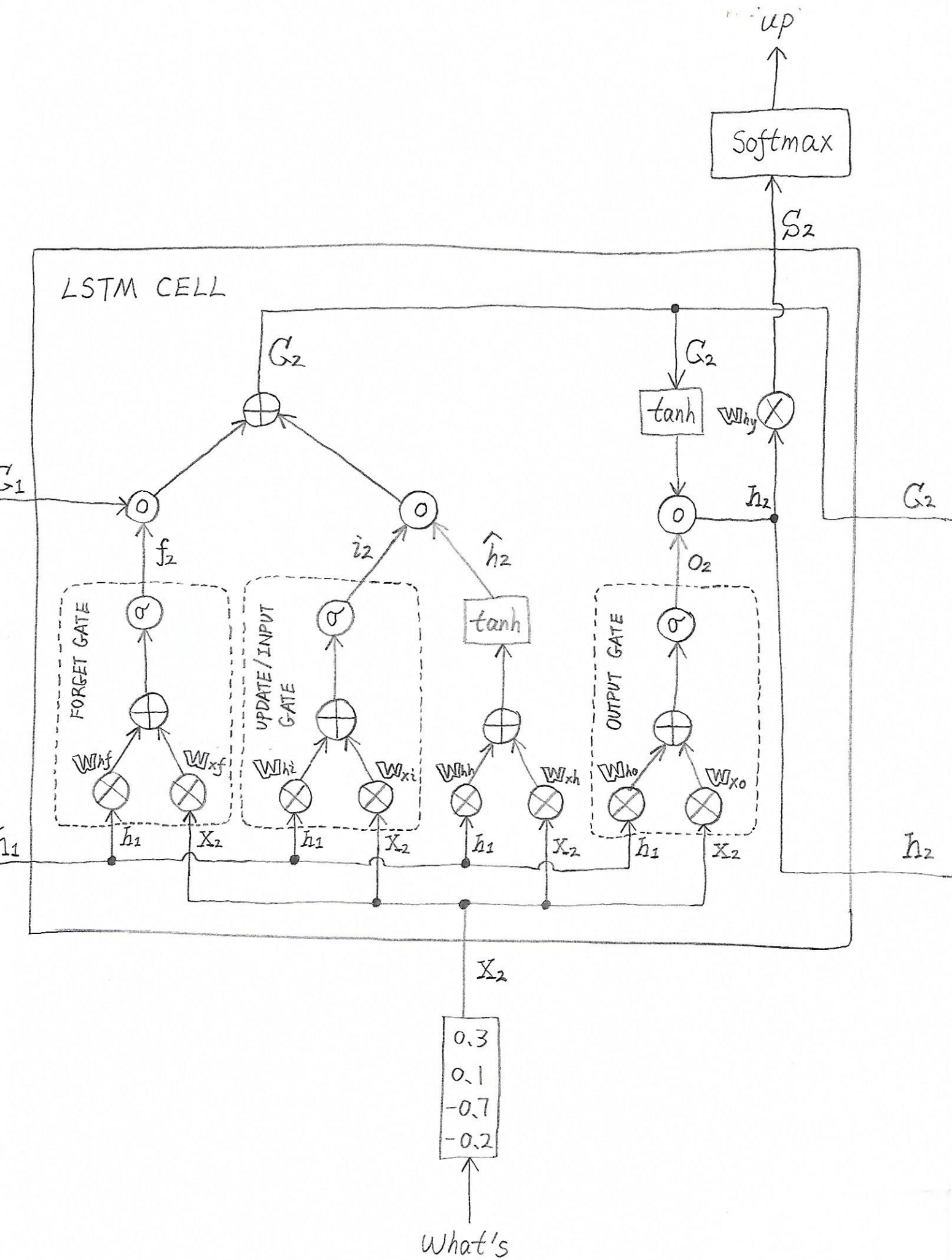
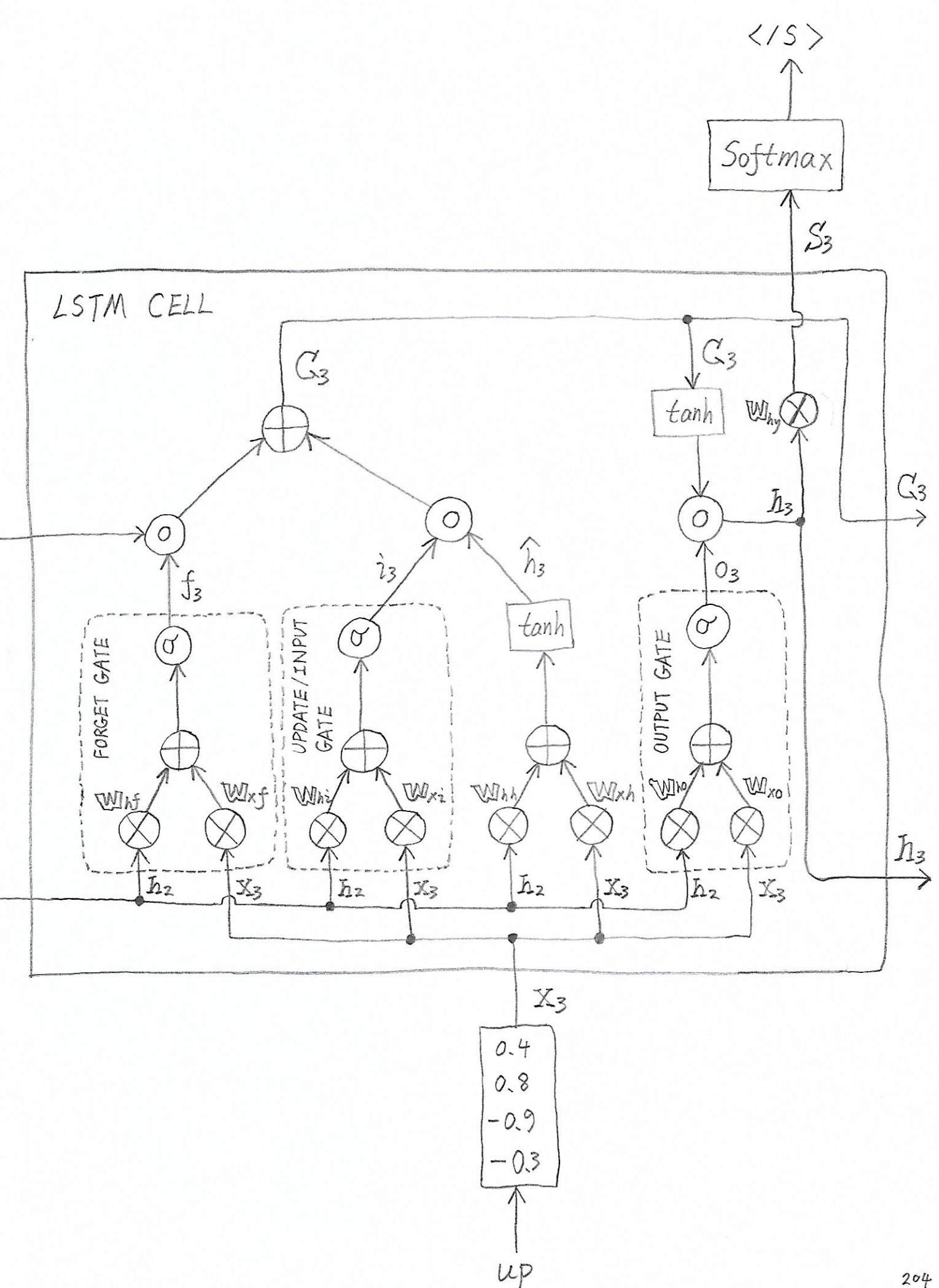


## 6. Long Short-Term Memory







## 7. Forward Pass

$$We = \{ \text{"hello"}: X_1, \text{"what"}: X_2, \dots, \text{"zaltan"}: X_{v\text{-size}} \}$$

$$X_t = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_{256} \end{bmatrix}, \text{ hidden-size} = 256$$

$$i_t = \sigma ( W_{xi} \cdot X_t + W_{hi} \cdot h_{t-1} ) \quad (\text{Update/Input Gate}) \quad (17)$$

$$f_t = \sigma ( W_{xf} \cdot X_t + W_{hf} \cdot h_{t-1} ) \quad (\text{Forget Gate}) \quad (18)$$

$$O_t = \sigma ( W_{xo} \cdot X_t + W_{ho} \cdot h_{t-1} ) \quad (\text{Output Gate}) \quad (19)$$

$$\hat{h}_t = \tanh ( W_{xh} \cdot X_t + W_{hh} \cdot h_{t-1} ) \quad (20)$$

$$C_t = f_t \circ C_{t-1} + i_t \circ \hat{h}_t \quad (20.1)$$

$$h_t = O_t \circ \tanh ( C_t ) \quad (20.2)$$

(17) - (20) can be written as

$$\begin{bmatrix} i_t \\ f_t \\ O_t \\ \hat{h}_t \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} \left( \begin{bmatrix} W_{xi} & W_{hi} \\ W_{xf} & W_{hf} \\ W_{xo} & W_{ho} \\ W_{xh} & W_{hh} \end{bmatrix} \cdot \begin{bmatrix} X_t \\ h_{t-1} \end{bmatrix} \right)$$

$$U_t = g ( T_{lstm} \cdot Z_t )$$

$$= g ( T_x \cdot X_t + T_h \cdot h_{t-1} )$$



$$S_t = W_{hy} \cdot h_t = \begin{bmatrix} S_t(y_1) \\ \vdots \\ S_t(y_m) \end{bmatrix} \quad (20.3)$$

$$P_t = \text{Softmax}(S_t) = \begin{bmatrix} \frac{e^{S_t(y_1)}}{\sum_{y' \in Y} e^{S_t(y')}} \\ \vdots \\ \frac{e^{S_t(y_m)}}{\sum_{y' \in Y} e^{S_t(y')}} \end{bmatrix} = \begin{bmatrix} P_t(y_1) \\ \vdots \\ P_t(y_m) \end{bmatrix} \quad (20.4)$$

$$\text{loss}_t = -\log P_t(y_t) = \log \sum_{y' \in Y} e^{S_t(y')} - S_t(y_t) \quad (\text{D.L. P19}) \quad (20.5)$$

8. Backward Pass.

$$dS_t = \frac{\partial \text{loss}_t}{\partial S_t} = P_t - 1_{y=y_t} \quad (\text{from } \textcircled{0})$$

$$dh_t = W_{hy}^T \cdot dS_t \quad (\text{from } \textcircled{9})$$

$$dW_{hy} = dS_t \cdot h_t^T \quad (\text{from } \textcircled{10})$$

$$dO_t = \tanh(G_t) \circ dh_t \quad (\text{by lemma 2}) \quad (21)$$

$$dG_t = \tanh'(G_t) \circ O_t \circ dh_t \quad (22)$$

Recall that

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\tanh'(x) = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

$$= 1 - \tanh^2(x)$$

$$df_t = C_{t-1} \circ dC_t \quad (\text{by corollary 2}) \quad (23)$$

$$dC_{t-1} = \frac{\partial l_t}{\partial C_{t-1}} = f_t \circ dC_t \quad (24)$$

$$di_t = \hat{h}_t \circ dC_t \quad (25)$$

$$d\hat{h}_t = i_t \circ dC_t \quad (26)$$

$$dU_t = \begin{bmatrix} di_t \\ df_t \\ dO_t \\ d\hat{h}_t \end{bmatrix} = [di_t; df_t; dO_t; d\hat{h}_t]$$

Recall that

$$U_t = g(T_{\text{estm}} \cdot Z_t)$$

$$= g(T_x \cdot X_t + T_h \cdot h_{t-1})$$

$$dT_x = (g'(T_{\text{estm}} \cdot Z_t) \circ dU_t) \cdot X_t^T \quad (\text{by lemma 1})$$

$$\begin{bmatrix} dW_{xi} \\ dW_{xf} \\ dW_{xo} \\ dW_{xh} \end{bmatrix} = \left( \begin{bmatrix} \sigma'(W_{xi} X_t + W_{hi} h_{t-1}) \\ \sigma'(W_{xf} X_t + W_{hf} h_{t-1}) \\ \sigma'(W_{xo} X_t + W_{ho} h_{t-1}) \\ \tanh'(W_{xh} X_t + W_{hh} h_{t-1}) \end{bmatrix} \circ \begin{bmatrix} di_t \\ df_t \\ dO_t \\ d\hat{h}_t \end{bmatrix} \right) \cdot X^T$$

$$d\mathbb{W}_{xi} = (\sigma'(\mathbb{W}_{xi} X_t + \mathbb{W}_{hi} h_{t-1}) \circ di_t) \cdot X^T \quad (27)$$

$$d\mathbb{W}_{xf} = (\sigma'(\mathbb{W}_{xf} X_t + \mathbb{W}_{hf} h_{t-1}) \circ df_t) \cdot X^T \quad (28)$$

$$d\mathbb{W}_{xo} = (\sigma'(\mathbb{W}_{xo} X_t + \mathbb{W}_{ho} h_{t-1}) \circ do_t) \cdot X^T \quad (29)$$

$$d\mathbb{W}_{xh} = (\tanh'(\mathbb{W}_{xh} X_t + \mathbb{W}_{hh} h_{t-1}) \circ d\hat{h}_t) \cdot X^T \quad (30)$$

Similarly,

$$dT_h = (g'(T_{lstm} \cdot Z_t) \circ dU_t) \cdot h_{t-1}^T \quad (\text{by lemma 1})$$

$$d\mathbb{W}_{hi} = (\sigma'(\mathbb{W}_{xi} X_t + \mathbb{W}_{hi} h_{t-1}) \circ di_t) \cdot h_{t-1}^T \quad (31)$$

$$d\mathbb{W}_{hf} = (\sigma'(\mathbb{W}_{xf} X_t + \mathbb{W}_{hf} h_{t-1}) \circ df_t) \cdot h_{t-1}^T \quad (32)$$

$$d\mathbb{W}_{ho} = (\sigma'(\mathbb{W}_{xo} X_t + \mathbb{W}_{ho} h_{t-1}) \circ do_t) \cdot h_{t-1}^T \quad (33)$$

$$d\mathbb{W}_{hh} = (\tanh'(\mathbb{W}_{xh} X_t + \mathbb{W}_{hh} h_{t-1}) \circ d\hat{h}_t) \cdot h_{t-1}^T \quad (34)$$

and

$$dX_t = T_x^T \cdot (g'(T_{lstm} \cdot Z_t) \circ dU_t) \quad (\text{by lemma 1})$$

$$= \begin{bmatrix} \mathbb{W}_{xi}^T & \mathbb{W}_{xf}^T & \mathbb{W}_{xo}^T & \mathbb{W}_{xh}^T \end{bmatrix} \cdot \begin{bmatrix} \sigma'(\mathbb{W}_{xi} X_t + \mathbb{W}_{hi} h_{t-1}) \circ di_t \\ \sigma'(\mathbb{W}_{xf} X_t + \mathbb{W}_{hf} h_{t-1}) \circ df_t \\ \sigma'(\mathbb{W}_{xo} X_t + \mathbb{W}_{ho} h_{t-1}) \circ do_t \\ \tanh'(\mathbb{W}_{xh} X_t + \mathbb{W}_{hh} h_{t-1}) \circ d\hat{h}_t \end{bmatrix}$$

$$\begin{aligned}
&= \mathbb{W}_{xi}^T \cdot ( \sigma'( \mathbb{W}_{xi} X_t + \mathbb{W}_{hi} h_{t-1} ) \circ di_t ) + \\
&\quad \mathbb{W}_{xf}^T \cdot ( \sigma'( \mathbb{W}_{xf} X_t + \mathbb{W}_{hf} h_{t-1} ) \circ df_t ) + \\
&\quad \mathbb{W}_{xo}^T \cdot ( \sigma'( \mathbb{W}_{xo} X_t + \mathbb{W}_{ho} h_{t-1} ) \circ do_t ) + \\
&\quad \mathbb{W}_{xh}^T \cdot ( \tanh'( \mathbb{W}_{xh} X_t + \mathbb{W}_{hh} h_{t-1} ) \circ d\hat{h}_t ) \quad (35)
\end{aligned}$$

Similarly,

$$\begin{aligned}
dh_{t-1} &= \frac{\partial \ell_t}{\partial h_{t-1}} = T_h^T \cdot ( g'( T_{lstm} \cdot z_t ) \circ dU_t ) \quad (\text{by lemma 1}) \\
&= \mathbb{W}_{hi}^T \cdot ( \sigma'( \mathbb{W}_{xi} X_t + \mathbb{W}_{hi} h_{t-1} ) \circ di_t ) + \\
&\quad \mathbb{W}_{hf}^T \cdot ( \sigma'( \mathbb{W}_{xf} X_t + \mathbb{W}_{hf} h_{t-1} ) \circ df_t ) + \\
&\quad \mathbb{W}_{ho}^T \cdot ( \sigma'( \mathbb{W}_{xo} X_t + \mathbb{W}_{ho} h_{t-1} ) \circ do_t ) + \\
&\quad \mathbb{W}_{hh}^T \cdot ( \tanh'( \mathbb{W}_{xh} X_t + \mathbb{W}_{hh} h_{t-1} ) \circ d\hat{h}_t ) \quad (36)
\end{aligned}$$



Backpropagation-Through-Time - LSTM  $(\{y_t\}_{t=1}^T, \{P_t\}_{t=1}^T, \{h_t\}_{t=1}^T, \{x_t\}_{t=1}^T,$

for  $t$  from  $T$  to  $1$ :  $\{W_{hy}, W_{xi}, W_{xf}, W_{xo}, W_{xh}, W_{hi}, W_{hf}, W_{ho}, W_{hh}\}^T$

$$\frac{\partial \ell_t}{\partial S_t} = P_t - 1 y = y_t \quad (20)$$

$$\frac{\partial \ell_t}{\partial W_{hy}} = \frac{\partial \ell_{t+1}}{\partial W_{hy}} + \frac{\partial \ell_t}{\partial S_t} \cdot h_t^T \quad (10)$$

$$\frac{\partial \ell_t}{\partial h_t} = \frac{\partial \ell_{t+1}}{\partial h_t} + W_{hy}^T \cdot \frac{\partial \ell_t}{\partial S_t} \quad (9)$$

$$\frac{\partial \ell_t}{\partial o_t} = \tanh(C_t) \circ \frac{\partial \ell_t}{\partial h_t} \quad (21)$$

$$\frac{\partial \ell_t}{\partial C_t} = \frac{\partial \ell_{t+1}}{\partial C_t} + \tanh'(C_t) \circ o_t \circ \frac{\partial \ell_t}{\partial h_t} \quad (22)$$

$$\frac{\partial \ell_t}{\partial f_t} = C_{t-1} \circ \frac{\partial \ell_t}{\partial C_t} \quad (23)$$

$$\frac{\partial \ell_t}{\partial i_t} = \hat{h}_t \circ \frac{\partial \ell_t}{\partial C_t} \quad (25)$$

$\partial C_t$  $\partial \hat{h}_t$ 

$$= \frac{\partial \ell_t}{\partial C_{t-1}} \circ f_t \circ \frac{\partial \ell_t}{\partial G_t} \quad (24)$$

$$\frac{\partial \ell_t}{\partial W_{xi}} = \frac{\partial \ell_{t+1}}{\partial W_{xi}} + (\sigma' (W_{xi} X_t + W_{hi} h_{t-1})) \circ \frac{\partial \ell_t}{\partial i_t} \cdot X_t^T \quad (27)$$

$$\frac{\partial \ell_t}{\partial W_{xf}} = \frac{\partial \ell_{t+1}}{\partial W_{xf}} + (\sigma' (W_{xf} X_t + W_{hf} h_{t-1})) \circ \frac{\partial \ell_t}{\partial f_t} \cdot X_t^T \quad (28)$$

$$\frac{\partial \ell_t}{\partial W_{xo}} = \frac{\partial \ell_{t+1}}{\partial W_{xo}} + (\sigma' (W_{xo} X_t + W_{ho} h_{t-1})) \circ \frac{\partial \ell_t}{\partial o_t} \cdot X_t^T \quad (29)$$

$$\frac{\partial \ell_t}{\partial W_{xh}} = \frac{\partial \ell_{t+1}}{\partial W_{xh}} + (\tanh' (W_{xh} X_t + W_{hh} h_{t-1})) \circ \frac{\partial \ell_t}{\partial \hat{h}_t} \cdot X_t^T \quad (30)$$

$$\frac{\partial \ell_t}{\partial W_{hi}} = \frac{\partial \ell_{t+1}}{\partial W_{hi}} + (\sigma' (W_{xi} X_t + W_{hi} h_{t-1})) \circ \frac{\partial \ell_t}{\partial i_t} \cdot h_{t-1}^T \quad (31)$$

$$\frac{\partial \ell_t}{\partial W_{hf}} = \frac{\partial \ell_{t+1}}{\partial W_{hf}} + (\sigma' (W_{xf} X_t + W_{hf} h_{t-1})) \circ \frac{\partial \ell_t}{\partial f_t} \cdot h_{t-1}^T \quad (32)$$

$$\frac{\partial \ell_t}{\partial W_{ho}} = \frac{\partial \ell_{t+1}}{\partial W_{ho}} + (\sigma' (W_{xo} X_t + W_{ho} h_{t-1})) \circ \frac{\partial \ell_t}{\partial o_t} \cdot h_{t-1}^T \quad (33)$$

GPU

$$= \frac{\partial \ell_t}{\partial w_{hh}} + \frac{\partial \ell_{t+1}}{\partial w_{hh}} + (\tanh' (w_{xh} x_t + w_{hh} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial \hat{h}_t} \quad (34)$$

$$= \frac{\partial \ell_t}{\partial x_t} = w_{xi}^T \cdot (\sigma' (w_{xi} x_t + w_{hi} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial i_t} + w_{xf}^T \cdot (\sigma' (w_{xf} x_t + w_{hf} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial f_t} + w_{xo}^T \cdot (\sigma' (w_{xo} x_t + w_{ho} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial o_t} + w_{xh}^T \cdot (\tanh' (w_{xh} x_t + w_{hh} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial \hat{h}_t} \quad (35)$$

$$= \frac{\partial \ell_t}{\partial h_{t-1}}$$

$$w_{hi}^T \cdot (\sigma' (w_{xi} x_t + w_{hi} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial i_t} + w_{hf}^T \cdot (\sigma' (w_{xf} x_t + w_{hf} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial f_t} + w_{ho}^T \cdot (\sigma' (w_{xo} x_t + w_{ho} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial o_t} + w_{hh}^T \cdot (\tanh' (w_{xh} x_t + w_{hh} h_{t-1})) \cdot \frac{\partial \ell_t}{\partial \hat{h}_t} \quad (36)$$

$$\text{return } \left\{ \frac{\partial \ell_{total}}{\partial w_{hy}}, \frac{\partial \ell_{total}}{\partial w_{xi}}, \frac{\partial \ell_{total}}{\partial w_{xf}}, \frac{\partial \ell_{total}}{\partial w_{xo}}, \frac{\partial \ell_{total}}{\partial w_{xh}}, \frac{\partial \ell_{total}}{\partial w_{hi}}, \frac{\partial \ell_{total}}{\partial w_{hf}}, \frac{\partial \ell_{total}}{\partial w_{ho}}, \frac{\partial \ell_{total}}{\partial w_{hh}} \right\}, \left\{ \frac{\partial \ell_t}{\partial x_t} \right\}_{t=}$$