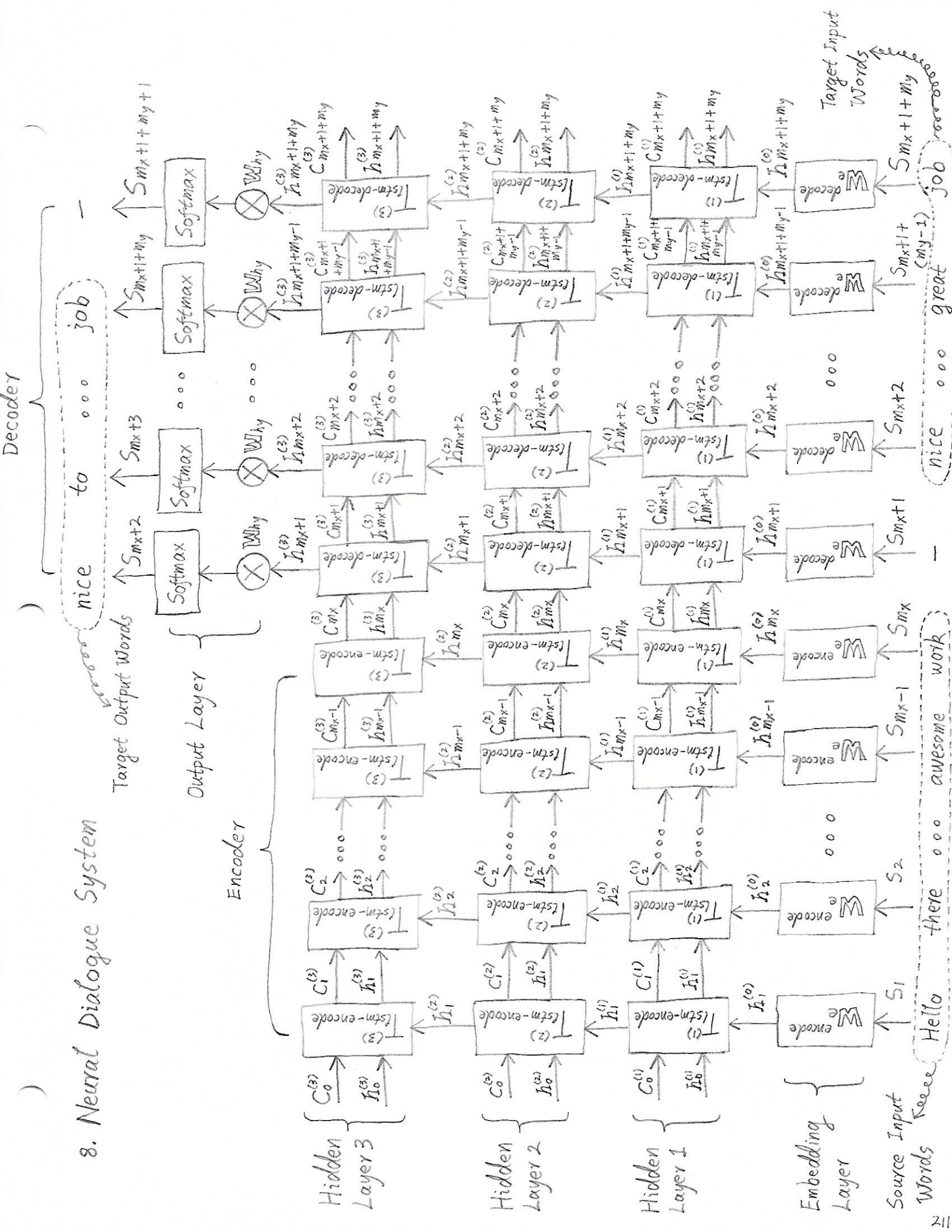# 8. Neural Dialogue System

9.

$\text{Neural\_Dialogue\_System\_Forward\_Pass}\ (\ \{X_i\}_{i=1}^{m_x},\ \{Y_i\}_{i=1}^{m_y},$

$We^{encode},\ We^{decode},\ \{T_{lstm-encode}^{(l)}\}_{l=1}^{L},\ \{T_{lstm-decode}^{(l)}\}_{l=1}^{L},\ W_{hy}\ )\ :$

$$S = [\ X_1, X_2, \cdots, X_{m_x}, -, Y_1, Y_2, \cdots, Y_{m_y}, -\ ] = \{S_t\}_{t=1}^{m_x+1+m_y+1}$$

for $l$ from 1 to $L$:

    $h_0^{(l)} = 0$   # Reset hidden state for the new training sequence

    $C_0^{(l)} = 0$   # so that the model doesn't depend on the previous
                   # training sequence. But who said it's not a good idea

    $T_{lstm}^{(l)} = T_{lstm-encode}^{(l)}$ # for an intelligent chat bot?

$We = We^{encode}$

for $t$ from 1 to $m_x + 1 + m_y$:

    if $\quad t == m_x + 1$:

        for $l$ from 1 to $L$:

            $T_{lstm}^{(l)} = T_{lstm-decode}^{(l)}$

        $We = We^{decode}$

    $h_t^{(0)} = \text{Embedding\_Look\_Up}\ (\ S_t,\ We\ )$

    for $l$ from 1 to $L$:

        $h_t^{(l)},\ C_t^{(l)} = \text{LSTM\_Forward\_Pass}\ (\ h_{t-1}^{(l)},\ C_{t-1}^{(l)},$
                                $h_t^{(l-1)},\ T_{lstm}^{(l)}\ )$

        # By D.L. P205, ⑰ − ⑳.2, replace $X_t$ by $h_t^{(l-1)}$

    if $\quad t >= m_x + 1$:

        $l_t,\ P_t = \text{Predict}\ (\ S_{t+1},\ h_t^{(L)},\ W_{hy}\ )$

        # By D.L. P206, ⑳.3 − ⑳.5, replace $y_t$ by $S_{t+1}$

return $\{l_t,\ P_t\}_{t=m_x+1}^{m_x+1+m_y+1},\ \{h_t^{(l)},\ C_t^{(l)}\}_{t=1,2,\cdots,m_x+1+m_y,\ l=0,1,\cdots,L}$

10. Neural_Dialogue_System_Backward_Pass ( $\{x_i\}_{i=1}^{m_x}$, $\{y_i\}_{i=1}^{m_y}$, $\{T_{lstm-encode}^{(\ell)}\}_{t=1}^{L}$,

$\{T_{lstm-decode}^{(\ell)}\}_{t=1}^{L}$, $\{h_t^{(\ell)}, C_t^{(\ell)}\}_{t=1,2,\cdots,m_x+1+m_y}$, $\ell = 0, 1, \cdots, L$, $We^{encode}$, $We^{decode}$, $Why$):

$$S = [x_1, x_2, \cdots, x_{m_x}, \underline{\ \ }, y_1, y_2, \cdots, y_{m_y}, \underline{\ \ }] = \{S_t\}_{t=1}^{m_x+1+m_y+1}$$

for $t$ from $m_x+1+m_y$ to $1$:

$\quad$ for $\ell$ from $L$ to $1$:

$$\frac{\partial L_t}{\partial h_t^{(\ell)}}, \quad \frac{\partial L_t}{\partial C_t^{(\ell)}}, \quad \frac{\partial L_t}{\partial T_{lstm}^{(\ell)}} = 0$$

$$\frac{\partial L_t}{\partial Why} = 0$$

$$\frac{\partial L_t}{\partial We} = 0$$

for $t$ from $m_x+1+m_y$ to $1$:

$\quad$ if $t == m_x$:

$\quad\quad$ for $\ell$ from $1$ to $L$:

$$\frac{\partial L}{\partial T_{lstm-decode}^{(\ell)}} = \frac{\partial L_{t+1}}{\partial T_{lstm}^{(\ell)}}$$

$$\frac{\partial L_{t+1}}{\partial T_{lstm}^{(\ell)}} = 0$$

$$\frac{\partial L}{\partial We^{decode}} = \frac{\partial L_{t+1}}{\partial We}$$

$$\frac{\partial L_{t+1}}{\partial We} = 0$$

$$\text{if } t >= m_x + 1 :$$

$$\frac{\partial L_t}{\partial h_t^{(L)}}, \frac{\partial L_t}{\partial W_{hy}} = \text{Output\_Grad}\left(S_{t+1}, P_t, h_t^{(L)}, \frac{\partial L_{t+1}}{\partial h_t^{(L)}}, \frac{\partial L_{t+1}}{\partial W_{hy}}, W_{hy}\right)$$

\# Refactored from Backpropagation\_Through\_Time\_LSTM() —

\# ⓪, ⑨ and ⑩, and replace $y_t$ by $S_{t+1}$

$$\text{for } t \text{ from } L \text{ to } 1 :$$

$$\text{if } t >= m_x + 1 :$$

$$T_{lstm}^{(\ell)} = T_{lstm-decode}^{(\ell)}$$

$$\text{else} :$$

$$T_{lstm}^{(\ell)} = T_{lstm-encode}^{(\ell)}$$

$$\frac{\partial L_t}{\partial C_{t-1}^{(\ell)}}, \frac{\partial L_t}{\partial X}, \frac{\partial L_t}{\partial T_{lstm}^{(\ell)}}, \frac{\partial L_{t+1}}{\partial T_{lstm}^{(\ell)}} = \text{LSTM\_Grad}($$

$$\left(C_{t-1}^{(\ell)}, h_{t-1}^{(\ell)}, h_t^{(\ell-1)}, h_t^{(\ell)}, T_{lstm}^{(\ell)}, \frac{\partial L_{t+1}}{\partial C_t^{(\ell)}}, C_t^{(\ell)}, \frac{\partial L_t}{\partial h_t^{(\ell)}}\right)$$

\# Refactored from Backpropagation\_Through\_Time\_LSTM() —

\# ㉑ — ㊱, and replace $X_t$ by $h_t^{(\ell-1)}$

$$\frac{\partial L_t}{\partial h_t^{(\ell-1)}} = \frac{\partial L_{t+1}}{\partial h_t^{(\ell-1)}} + \frac{\partial L_t}{\partial X} *$$

$$*$$

$$\text{if } t >= m_x + 1:$$

$$\quad W_e = W_e^{decode}$$

$$\text{else :}$$

$$\quad W_e = W_e^{encode}$$

$$\frac{\partial L_t}{\partial W_e} = Embedding\text{-}Grad\left( S_t, \; \boxed{\frac{\partial L_t}{\partial h_t^{(o)}}}, \; W_e, \; \frac{\partial L_{t+1}}{\partial W_e} \right)$$

$$\text{for } l \text{ from } 1 \text{ to } L:$$

$$\frac{\partial L}{\partial T_{lstm-encode}^{(l)}} = \boxed{\frac{\partial L_{t=1}}{\partial T_{lstm}^{(l)}}} \qquad \frac{\partial L_{t=1}}{\partial W_e}$$

$$\frac{\partial L}{\partial W_e^{encode}}$$

$$return \quad \frac{\partial L_{t=m_x+1}}{\partial W_{hy}}, \; \frac{\partial L}{\partial W_e^{encode}}, \; \frac{\partial L}{\partial W_e^{decode}}, \; \left\{ \frac{\partial L}{\partial T_{lstm-encode}^{(l)}} \right\}_{l=1}^L, \; \left\{ \frac{\partial L}{\partial T_{lstm-decode}^{(l)}} \right\}_{l=1}^L$$

11. Stochastic_Gradient_Descent ( $\mathbb{X} = \{ X^{(1)}, X^{(2)}, \cdots, X^{(m)} \}$,

$\quad \mathbb{Y} = \{ Y^{(1)}, Y^{(2)}, \cdots, Y^{(m)} \}, \varepsilon, \delta )$:

$$\theta' = Random()$$

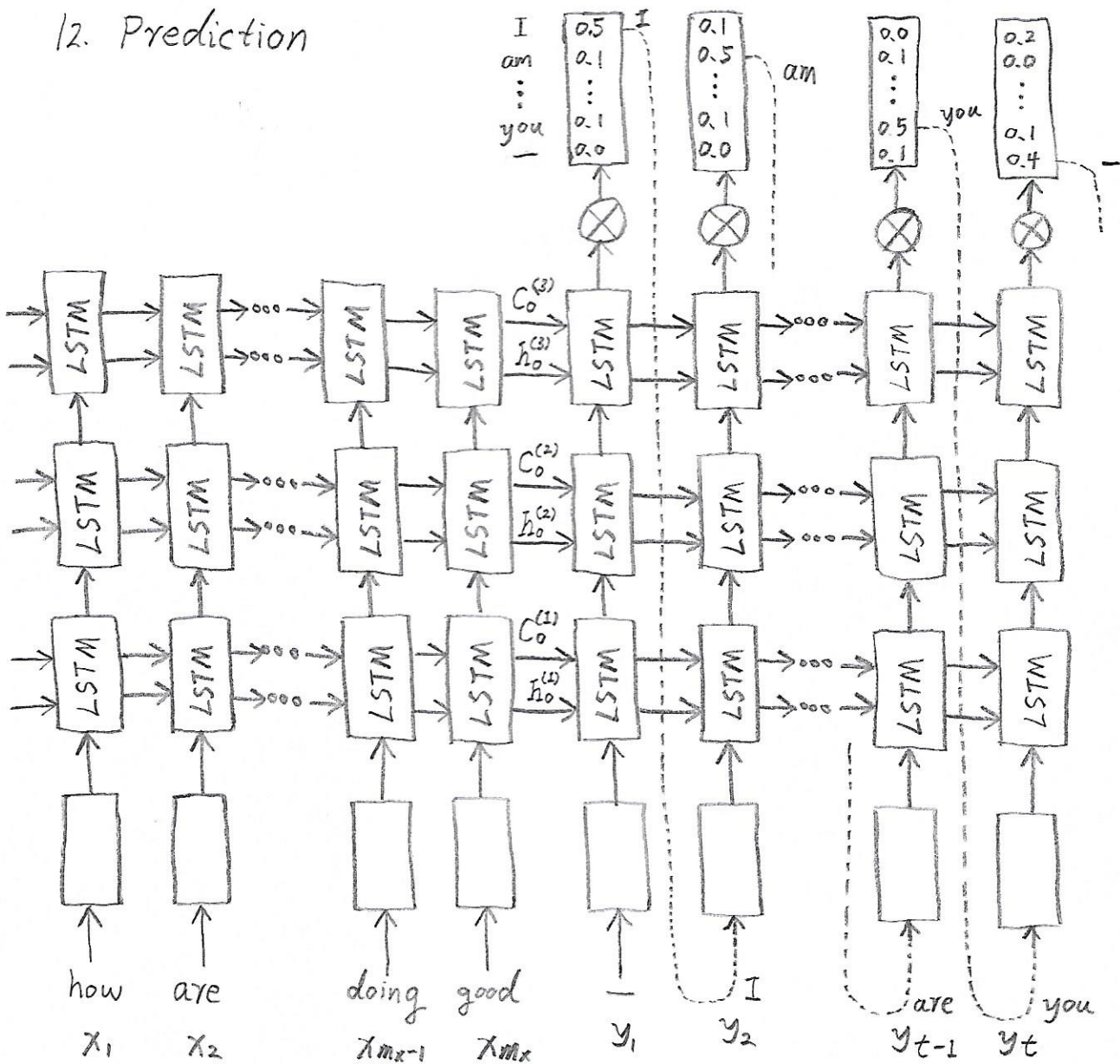$$while \; ( \; \| \frac{1}{m} \cdot \sum_{i=1}^{m} L(X^{(i)}, Y^{(i)}, \theta) \big|_{\theta=\theta'} \|_2 > \delta \; ):$$

$$( \mathbb{X}' = \{ X^{(i)} \}_{i=1}^{m'}, \; \mathbb{Y}' = \{ Y^{(i)} \}_{i=1}^{m'} )$$

$$= Sample\_Mini\_Batch ( \mathbb{X}, \mathbb{Y} )$$

$$\theta' = \theta' - \varepsilon \cdot ( \frac{1}{m'} \cdot \sum_{i=1}^{m'} \nabla_\theta L(X^{(i)}, Y^{(i)}, \theta) \big|_{\theta=\theta'} )$$

$$return \; \theta'$$

12. Prediction

$\text{Neural\_Dialogue\_System\_Greedy\_Predict} \left( \{X_i\}_{i=1}^{m_x}, \alpha, \right.$

$\left\{ T_{lstm\text{-}encode}^{(\ell)} \right\}_{\ell=1}^{L}, \left\{ T_{lstm\text{-}decode}^{(\ell)} \right\}_{\ell=1}^{L}, W_e^{encode}, W_e^{decode}, W_{hy} \left. \right):$

$\quad \left\{ h_0^{(\ell)} \right\}_{\ell=1}^{L}, \left\{ C_0^{(\ell)} \right\}_{\ell=1}^{L} = \text{Encode} \left( \{X_i\}_{i=1}^{m_x}, \left\{ T_{lstm\text{-}encode}^{(\ell)} \right\}_{\ell=1}^{L}, W_e^{encode} \right.$

$\quad t = 1$

$\quad y_1 = \text{``\_''}$

$\quad \text{while} \quad t <= \alpha \cdot m_x :$

$\qquad \left\{ h_t^{(\ell)}, C_t^{(\ell)} \right\}_{\ell=1}^{L} = \text{LSTM} \left( \left\{ h_{t-1}^{(\ell)}, C_{t-1}^{(\ell)} \right\}_{\ell=1}^{L}, \left\{ T_{lstm\text{-}decode}^{(\ell)} \right\}_{\ell=1}^{L}, \right.$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad W_e^{decode}, y_t \left. \right)$

$\qquad P_t = \text{Softmax} \left( h_t^{(L)}, W_{hy} \right) = \left[ p_t(y_1), \cdots, p_t(y_m) \right]^T$

$\qquad y_{t+1} = \underset{y_i}{\arg\max} \; p_t(y_i)$

$\qquad \text{if} \quad y_{t+1} == \text{``\_''} :$

$\qquad\qquad \text{break}$

$\qquad t = t + 1$

$\quad \text{return} \quad \{ y_2, y_3, \cdots, y_t \}$

# 13. Attention Mechanism

$$p\left(y_t \mid y_{<t}, \{x_i\}_{i=1}^{m_x}\right)$$

Softmax

$\otimes$ $W_s$

$\tilde{h}_t \rightsquigarrow$ Attention Vector

tanh

$\otimes$ $W_c$

Context Vector

$C_t$

$C_t$

$h_t$

Attention Layer

$\alpha_{t1}$ $\alpha_{t2}$ $\alpha_{ts-1}$ $\alpha_{ts}$ Attention Weights

$\bar{h}_1$ $\bar{h}_2$ $\bar{h}_{s-1}$ $\bar{h}_s$

$\oplus$

$\bar{h}_1$ $\bar{h}_2$ $\bar{h}_{s-1}$ $\bar{h}_s$ $h_t$

align

. . .

Encoder

Decoder

216

$$\alpha_{ts} = \text{align}(h_t, \bar{h}_s)$$

$$= \frac{e^{\text{score}(h_t, \bar{h}_s)}}{\sum_{s'} e^{\text{score}(h_t, \bar{h}_{s'})}}$$

$$\text{score}(h_t, \bar{h}_s) = \begin{cases} h_t^T \cdot \bar{h}_s \\ h_t^T \cdot W_a \cdot \bar{h}_s \\ v_a^T \cdot \tanh(W_a \cdot [h_t; \bar{h}_s]) \end{cases}$$

$$C_t = \sum_s \alpha_{ts} \cdot \bar{h}_s$$

$$\tilde{h}_t = \tanh(W_c \cdot [C_t; h_t])$$

$$p(y_t \mid y_{<t}, \{x_i\}_{i=1}^{m_x}) = \text{Softmax}(W_s \cdot \tilde{h}_t)$$