

* LSTM Backprop \rightarrow

$a < t >$

$$\frac{\partial a_{re} < t >}{\partial \Gamma_{o,ij} < t >} = \tanh(C_{re} < t >) \delta_{k,i} \delta_{l,j} ; \frac{\partial a_{re} < t >}{\partial G_{ij} < t >} = \Gamma_{re} < t > (1 - \tanh^2(C_{re} < t >)) \delta_{k,i} \delta_{l,j}$$

$$\frac{\partial C_{re} < t >}{\partial \Gamma_{f,ij} < t >} = C_{re} < t-1 > \delta_{k,i} \delta_{l,j} ; \frac{\partial C_{re} < t >}{\partial \Gamma_{g,ij} < t >} = \tilde{C} < t > \delta_{k,i} \delta_{l,j}$$

$$\frac{\partial C_{re} < t >}{\partial C_{ij} < t-1 >} = \Gamma_{re} < t > \delta_{k,i} \delta_{l,j} ; \frac{\partial C_{re} < t >}{\partial \tilde{C} < t >}} = \Gamma_{g,ij} < t > \delta_{k,i} \delta_{l,j}$$

$$\begin{aligned} \frac{\partial h}{\partial \Gamma_{o,ij} < t >}} &= \sum_{k=1}^{N_k} \sum_{l=1}^M \frac{\partial h}{\partial a_{re} < t >}} \cdot \frac{\partial a_{re} < t >}{\partial \Gamma_{o,ij} < t >}} = \sum_{k=1}^{N_k} \sum_{l=1}^M \frac{\partial h}{\partial a_{re} < t >}} \delta_{l,j} \cdot \tanh(C_{re} < t >) \cdot \delta_{k,i} \\ &= \frac{\partial h}{\partial a_{ij} < t >}} \tanh(G_{ij} < t >) \end{aligned}$$

$$\begin{aligned} \frac{\partial h}{\partial G_{ij} < t >}} &= \left(\sum_{k=1}^{N_k} \sum_{l=1}^M \frac{\partial h}{\partial a_{re} < t >}} \cdot \frac{\partial a_{re} < t >}{\partial G_{ij} < t >}} \right) + \frac{\partial h}{\partial G_{ij} < t >}} = \\ &= \frac{\partial h}{\partial a_{ij} < t >}} \cdot \Gamma_{o,ij} < t > (1 - \tanh^2(G_{ij} < t >)) + \frac{\partial h}{\partial G_{ij} < t >}} \end{aligned}$$

$$\begin{aligned} \frac{\partial h}{\partial \Gamma_{f,ij} < t >}} &= \sum_{k=1}^{N_k} \sum_{l=1}^M \frac{\partial h}{\partial C_{re} < t >}} \cdot \frac{\partial C_{re} < t >}{\partial \Gamma_{f,ij} < t >}} = \sum_{k=1}^{N_k} \sum_{l=1}^M \frac{\partial h}{\partial C_{re} < t >}} \cdot C_{re} < t-1 > \delta_{k,i} \delta_{l,j} \\ &= \frac{\partial h}{\partial C_{ij} < t >}} \cdot C_{ij} < t-1 > = \\ &\left(\frac{\partial h}{\partial a_{ij} < t >}} * \Gamma_{o,ij} < t > (1 - \tanh^2(G_{ij} < t >)) + \frac{\partial h}{\partial G_{ij} < t >}} \right) * C_{ij} < t-1 > \end{aligned}$$

$$\frac{\partial h}{\partial \tilde{C}_{ij}^{<+>}} = \sum_{k=1}^N \sum_{l=1}^M \frac{\partial h}{\partial C_{kl}^{<+>}} \cdot \Gamma_{kl}^{<+>} \cdot \delta_{ki} \delta_{lj} =$$

$$= \left(\frac{\partial h}{\partial a_{ij}^{<+>}} * \Gamma_{ij}^{<+>} * (1 - \tanh^2(a_{ij}^{<+>})) + \frac{\partial h}{\partial b_{ij}^{<+>}} \right) \Gamma_{ij}^{<+>}$$

$$\frac{\partial h}{\partial \tilde{\mu}_{ij}^{<+>}} = \left(\frac{\partial h}{\partial a_{ij}^{<+>}} * \Gamma_{ij}^{<+>} * (1 - \tanh^2(a_{ij}^{<+>})) + \frac{\partial h}{\partial b_{ij}^{<+>}} \right) * \tilde{C}_{ij}^{<+>}$$

$$\frac{\partial h}{\partial \tilde{C}_{ij}^{<+>}} = \left(\frac{\partial h}{\partial a_{ij}^{<+>}} * \Gamma_{ij}^{<+>} * (1 - \tanh^2(a_{ij}^{<+>})) + \frac{\partial h}{\partial b_{ij}^{<+>}} \right) * \Gamma_{ij}^{<+>}$$

↔

$$\Gamma = \sigma(z) ; z = WX \rightarrow \frac{\partial \Gamma_{kl}}{\partial z_{ij}} = \Gamma_{kl}(1 - \Gamma_{kl}) \delta_{ki} \delta_{lj}$$

$$\frac{\partial \Gamma_{kl}}{\partial w_{ij}} = \delta_{ki} x_{jl} \rightarrow \frac{\partial \Gamma_{kl}}{\partial w_{ij}} = \sum_{r=1}^N \sum_{s=1}^M \frac{\partial \Gamma_{kl}}{\partial z_{rs}} \cdot \frac{\partial z_{rs}}{\partial w_{ij}}$$

$$\frac{\partial \Gamma_{kl}}{\partial w_{ij}} = \sum_{r=1}^N \sum_{s=1}^M \Gamma_{kl}(1 - \Gamma_{kl}) \delta_{kr} \delta_{ls} \cdot \delta_{ri} x_{js} = \delta_{ki} \Gamma_{kl}(1 - \Gamma_{kl}) x_{jl}$$

↔

$$\frac{\partial h}{\partial w_{ij}} = \sum_{k=1}^N \sum_{l=1}^M \frac{\partial h}{\partial \Gamma_{kl}} \frac{\partial \Gamma_{kl}}{\partial w_{ij}} = \sum_{k=1}^N \sum_{l=1}^M \frac{\partial h}{\partial \Gamma_{kl}} \cdot \Gamma_{kl}(1 - \Gamma_{kl}) x_{jl} \delta_{ki}$$

$$= \sum_{l=1}^M \frac{\partial h}{\partial \Gamma_{il}} \Gamma_{il}(1 - \Gamma_{il}) x_{jl} , \rightarrow x = [a^{<+>}, x^{<+>}]$$

$$\frac{\partial h}{\partial w_{ij}} = \sum_{l=1}^M \frac{\partial h}{\partial \Gamma_{il}} \cdot \Gamma_{il}(1 - \Gamma_{il}) [a^{<+>}, x^{<+>}]_j$$

$$\frac{\partial h}{\partial w_{ij}} = \sum_{l=1}^m \frac{\partial h}{\partial \Gamma_{ij}^{l,t}} \cdot \Gamma_{ij}^{l,t} (1 - \Gamma_{ij}^{l,t}) [a^{l,t-1}, x^{l,t}]_{je}$$

$$\frac{\partial h}{\partial w_{ij}} = \sum_{l=1}^m \frac{\partial h}{\partial \hat{C}_{ie}^{l,t}} (1 - \tanh^2(w_{ie} [a^{l,t-1}, x^{l,t}] + b_c)_{ie}) [a^{l,t-1}, x^{l,t}]_{je}$$

↔

$$z = wx \rightarrow \frac{\partial z_{ke}}{\partial x_{ij}} = \delta_{e,j} w_{k,i} \rightarrow \frac{\partial \Gamma_{ke}}{\partial x_{ij}} = \delta_{e,j} \Gamma_{ke} (1 - \Gamma_{ke}) w_{k,i}$$

↔

$$\frac{\partial h}{\partial [a^{l,t-1}, x^{l,t}]_{ij}} = \sum_{k=1}^{n_c} \sum_{l=1}^m \left(\frac{\partial h}{\partial \Gamma_{ke}} \frac{\partial \Gamma_{ke}}{\partial [a^{l,t-1}, x^{l,t}]_{ij}} + \frac{\partial h}{\partial \Gamma_{ke}^{l,t}} \frac{\partial \Gamma_{ke}^{l,t}}{\partial [a^{l,t-1}, x^{l,t}]_{ij}} \right)$$

↑
concatenated, for

the gradients,
just update
neural network
element.

$$+ \frac{\partial h}{\partial \Gamma_{ke}^{l,t}} \frac{\partial \Gamma_{ke}^{l,t}}{\partial [a^{l,t-1}, x^{l,t}]_{ij}} + \frac{\partial h}{\partial \hat{C}_{ke}^{l,t}} \frac{\partial \hat{C}_{ke}^{l,t}}{\partial [a^{l,t-1}, x^{l,t}]_{ij}} \Bigg)$$

$$= \sum_{k=1}^{n_c} \sum_{l=1}^m \delta_{e,j} \left[\frac{\partial h}{\partial \Gamma_{ke}} \Gamma_{ke} (1 - \Gamma_{ke}) w_{k,i} + \frac{\partial h}{\partial \Gamma_{ke}^{l,t}} \Gamma_{ke}^{l,t} (1 - \Gamma_{ke}^{l,t}) w_{k,i} \right.$$

$$\left. + \frac{\partial h}{\partial \hat{C}_{ke}^{l,t}} (1 - \tanh^2(w_{ie} [a^{l,t-1}, x^{l,t}]_{je} + b_c)) w_{k,i} \right]$$

$$+ \frac{\partial h}{\partial \hat{C}_{ke}^{l,t}} (1 - \tanh^2(w_{ie} [a^{l,t-1}, x^{l,t}]_{je} + b_c)) w_{k,i}$$

$$= \sum_{k=1}^{n_c} \left[\frac{\partial h}{\partial \Gamma_{kj}} \Gamma_{kj} (1 - \Gamma_{kj}) w_{k,i} + \frac{\partial h}{\partial \Gamma_{kj}^{l,t}} \Gamma_{kj}^{l,t} (1 - \Gamma_{kj}^{l,t}) w_{k,i} \right.$$

$$\left. + \frac{\partial h}{\partial \hat{C}_{kj}^{l,t}} (1 - \tanh^2(w_{ie} [a^{l,t-1}, x^{l,t}]_{je} + b_c)) w_{k,i} \right]$$

- Factorized representation.

$$d\Gamma_0<t> = da_{next} * \tanh(C_{next})$$

$$d\Gamma_f<t> = (da_{next} * \Gamma_0<t> * (1 - \tanh^2(C_{next})) + dC_{next}) * C_{prev}^{t-1}$$

$$d\Gamma_\mu<t> = (da_{next} * \Gamma_0<t> * (1 - \tanh^2(C_{next})) + dC_{next}) * \hat{C}^2<t>$$

$$d\hat{C}^2<t> = (da_{next} * \Gamma_0<t> * (1 - \tanh^2(C_{next})) + dC_{next}) * \Gamma_\mu<t>$$

$$dW_f = [d\Gamma_f<t> * \Gamma_f<t> * (1 - \Gamma_f<t>)] \cdot \begin{pmatrix} a_{prev} \\ x<t> \end{pmatrix}^T$$

$$dW_0 = [d\Gamma_0<t> * \Gamma_0<t> * (1 - \Gamma_0<t>)] \cdot \begin{pmatrix} a_{prev} \\ x<t> \end{pmatrix}^T$$

$$dW_\mu = [d\Gamma_\mu<t> * \Gamma_\mu<t> * (1 - \Gamma_\mu<t>)] \cdot \begin{pmatrix} a_{prev} \\ x<t> \end{pmatrix}^T$$

$$dW_c = [d\hat{C}^2<t> * (1 - \tanh^2(W_c[a<t-1>, x<t>] + b_c))] \cdot \begin{pmatrix} a_{prev} \\ x<t> \end{pmatrix}^T$$

$$d\begin{bmatrix} a<t-1> \\ x<t> \end{bmatrix} = W_f^T [d\Gamma_f<t> * \Gamma_f<t> * (1 - \Gamma_f<t>)]$$

$$+ W_0^T [d\Gamma_0<t> * \Gamma_0<t> * (1 - \Gamma_0<t>)]$$

$$+ W_\mu^T [d\Gamma_\mu<t> * \Gamma_\mu<t> * (1 - \Gamma_\mu<t>)]$$

$$+ W_c^T [d\hat{C}^2<t> * (1 - \tanh^2(W_c[a<t-1>, x<t>] + b_c))]$$

$$dC_{prev} = (da_{next} * \Gamma_0<t> * (1 - \tanh^2(C_{next})) + dC_{next}) * \Gamma_f<t>$$