

LSTM Network gradient
through time

$$\begin{aligned}
 \frac{\partial E}{\partial w_f} = & \left[\frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial w_f} \right] + \left[\frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial w_f} + \right. \\
 & + \left. \frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial w_f} + \frac{\partial E_1}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial w_f} \right] + \\
 & + \left[\frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial c_{C0}} \cdot \frac{\partial c_{C0}}{\partial w_f} \right. \\
 & + \left. \frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial c_{C0}} \cdot \frac{\partial c_{C0}}{\partial w_f} \right. \\
 & + \left. \frac{\partial E_2}{\partial s_{C2}} \cdot \frac{\partial s_{C2}}{\partial c_{C2}} \cdot \frac{\partial c_{C2}}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial w_f} \cdot \frac{\partial c_{C0}}{\partial w_f} + \right. \\
 & + \left. \frac{\partial E_1}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial c_{C0}} \cdot \frac{\partial c_{C0}}{\partial w_f} + \frac{\partial E_1}{\partial s_{C1}} \cdot \frac{\partial s_{C1}}{\partial c_{C1}} \cdot \frac{\partial c_{C1}}{\partial c_{C0}} \cdot \frac{\partial c_{C0}}{\partial w_f} \right. \\
 & + \left. \frac{\partial E_0}{\partial s_{C0}} \cdot \frac{\partial s_{C0}}{\partial c_{C0}} \cdot \frac{\partial c_{C0}}{\partial w_f} \right]
 \end{aligned}$$

$$\begin{aligned}
& - \frac{\partial F_2}{\partial S_{C2>}} \cdot \frac{\partial S_{C2>}}{\partial C_{C2>}} \cdot \frac{\partial C_{C2>}}{\partial W_f} + \\
& + \left[\frac{\partial F_2}{\partial S_{C2>}} \cdot \frac{\partial S_{C2>}}{\partial C_{C2>}} \left(\frac{\partial C_{C2>}}{\partial C_{C1>}} + \frac{\partial C_{C2>}}{\partial S_{C1>}} \cdot \frac{\partial S_{C1>}}{\partial C_{C1>}} \right) \frac{\partial C_{C1>}}{\partial W_f} + \frac{\partial F_1}{\partial S_{C1>}} \frac{\partial S_{C1>}}{\partial C_{C1>}} \cdot \frac{\partial C_{C1>}}{\partial W_f} \right] + \\
& + \left[\frac{\partial F_2}{\partial S_{C2>}} \cdot \frac{\partial S_{C2>}}{\partial C_{C2>}} \left(\frac{\partial C_{C2>}}{\partial C_{C1>}} + \frac{\partial C_{C2>}}{\partial S_{C1>}} \cdot \frac{\partial S_{C1>}}{\partial C_{C1>}} \right) \left(\frac{\partial C_{C1>}}{\partial C_{C0>}} + \frac{\partial C_{C1>}}{\partial S_{C0>}} \cdot \frac{\partial S_{C0>}}{\partial C_{C0>}} \right) \frac{\partial C_{C0>}}{\partial W_f} \right. \\
& \left. + \frac{\partial F_1}{\partial S_{C1>}} \cdot \frac{\partial S_{C1>}}{\partial C_{C1>}} \left(\frac{\partial C_{C1>}}{\partial C_{C0>}} + \frac{\partial C_{C1>}}{\partial S_{C0>}} \cdot \frac{\partial S_{C0>}}{\partial C_{C0>}} \right) \frac{\partial C_{C0>}}{\partial W_f} + \right. \\
& \left. + \frac{\partial F_1}{\partial S_{C0>}} \cdot \frac{\partial S_{C0>}}{\partial C_{C0>}} \cdot \frac{\partial C_{C0>}}{\partial W_f} \right] \Rightarrow
\end{aligned}$$

$$\frac{\partial F}{\partial W_f} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial F_k}{\partial S_{Ck>}} \cdot \frac{\partial S_{Ck>}}{\partial C_{Ck>}} \left[\prod_{j=l+1}^K \left(\frac{\partial C_{Cj>}}{\partial C_{Cj-1>}} + \frac{\partial C_{Cj>}}{\partial S_{Cj-1>}} \cdot \frac{\partial S_{Cj-1>}}{\partial C_{Cj-1>}} \right) \right] \frac{\partial C_{C0>}}{\partial W_f}$$

$$\frac{\partial F}{\partial W_h} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial F_k}{\partial S_{Ch>}} \cdot \frac{\partial S_{Ch>}}{\partial C_{Ch>}} \left[\prod_{j=l+1}^K \left(\frac{\partial C_{Ch>}}{\partial C_{Ch-1>}} + \frac{\partial C_{Ch>}}{\partial S_{Ch-1>}} \cdot \frac{\partial S_{Ch-1>}}{\partial C_{Ch-1>}} \right) \right] \frac{\partial C_{Ch>}}{\partial W_h}$$

$$\frac{\partial F}{\partial W_c} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial F_k}{\partial S_{Cc>}} \cdot \frac{\partial S_{Cc>}}{\partial C_{Cc>}} \left[\prod_{j=l+1}^K \left(\frac{\partial C_{Cc>}}{\partial C_{Cc-1>}} + \frac{\partial C_{Cc>}}{\partial S_{Cc-1>}} \cdot \frac{\partial S_{Cc-1>}}{\partial C_{Cc-1>}} \right) \right] \frac{\partial C_{Cc>}}{\partial W_c}$$

$$\frac{\partial F}{\partial W_b} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial F_k}{\partial S_{Cb>}} \cdot \prod_{j=l+1}^K \left(\frac{\partial S_{Cb>}}{\partial S_{Cb-1>}} \right) \frac{\partial S_{Cb>}}{\partial W_b}$$

$$\frac{\partial C_{Cj>}}{\partial C_{Cj-1>}} = P_{Cj>} \rightarrow \prod \frac{\partial C_{Cj>}}{\partial C_{Cj-1>}}$$

will not result in update if
the $\frac{\partial C_{Cj>}}{\partial C_{Cj-1>}}$ is on the
path

\rightarrow Value $C_{Cj>}$ dependency of
it had importance \rightarrow
If enabled and fit, enabled
update.

For this gradient, it won't vanish because

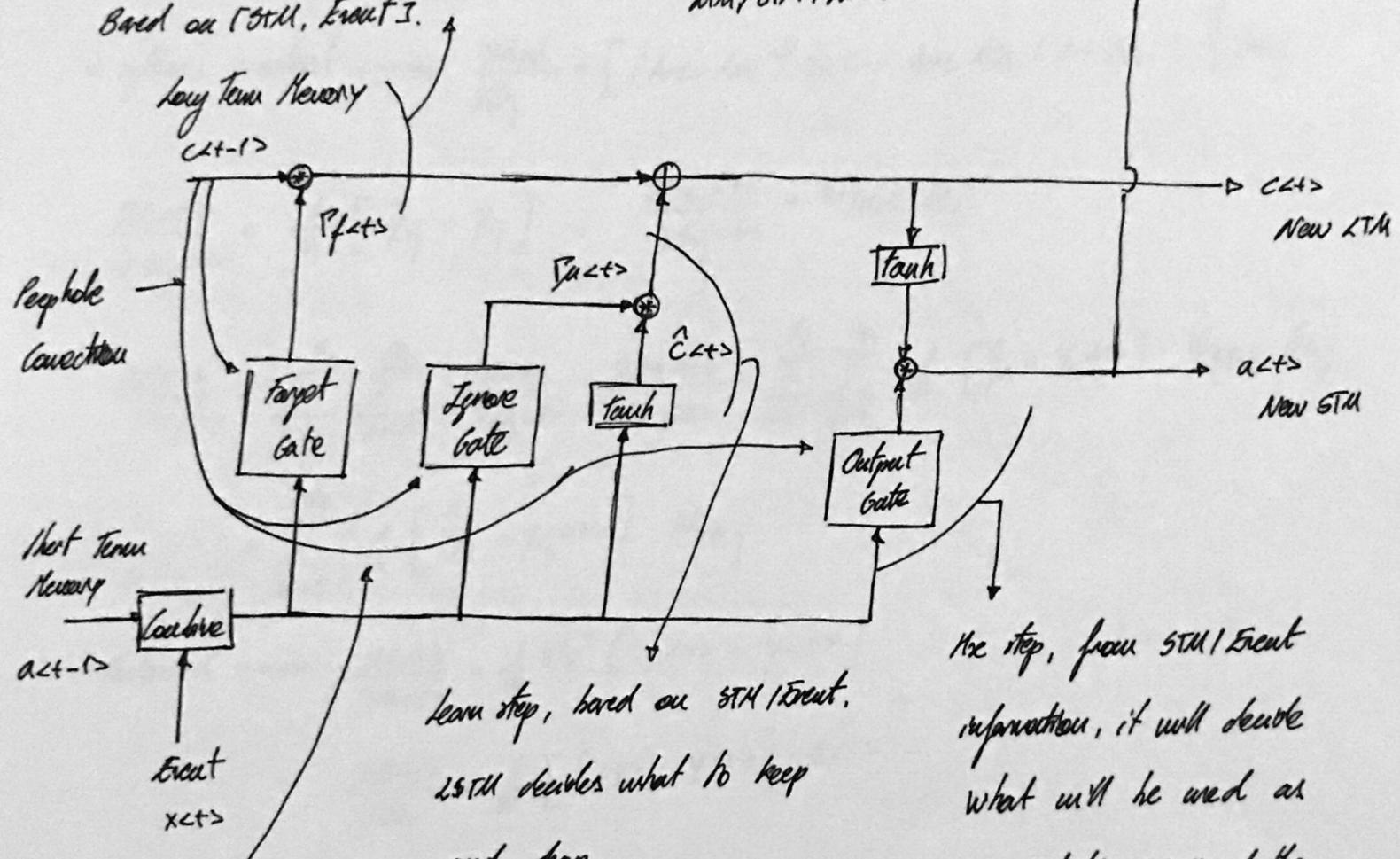
of the path until previous block $\frac{\partial C_{Ct+1>}}{\partial C_{Ct>}} \cdot \frac{\partial C_{Ct>}}{\partial W_b}$

$$\frac{\partial \tilde{F}}{\partial x_{\leq l>}} = \sum_{k=0}^{Ty} \frac{\partial F_k}{\partial s_{\leq k>}} \cdot \frac{\partial s_{\leq k>}}{\partial c_{\leq k>}} \left[\prod_{j=l+1}^K \left(\frac{\partial c_{\leq j>}}{\partial c_{\leq j-1>}} + \frac{\partial c_{\leq j>}}{\partial s_{\leq j-1>}} \cdot \frac{\partial s_{\leq j-1>}}{\partial c_{\leq j-1>}} \right) \right] \frac{\partial c_{\leq l>}}{\partial x_{\leq l>}}$$

day 10: Long Term Memory

Forget step, in which LSTM
tags LSTM information not useful
Based on STM, Forget I.

Remember step, combines
useful information from
STM/LSTM / Input



Peephole connections:

LSTM variant, includes
LSTM in gate decisions.

Learn step, based on STM/Event.
LSTM decides what to keep
and drop.
Usefull information of STM/Event.

The step, from STM/Event
information, it will decide
what will be used as
a prediction/STM of the
combination LSTM/STM/Event.

$$P_f^{t+>} = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

$$P_i^{t+>} = \sigma(W_i[h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_{t+>} = \tanh(W_c[h_{t-1}, x_t] + b_c)$$

$$C_t = C_{t-1} * P_f^{t+>} + P_i^{t+>} * \tilde{C}_{t+>}$$

$$P_o^{t+>} = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$h_t = P_o^{t+>} * \tanh(C_t)$$

all backprop →

$$\hat{y}^{<t>} = \text{softmax}(z^{<t>}) ; z^{<t>} = W y^{<t>} + b$$

* cross-entropy gradient → $\frac{\partial E}{\partial \hat{y}_{ij}} = -\frac{1}{m} \cdot \frac{y_{ij}}{\hat{y}_{ij}}$

* softmax gradient → $\frac{\partial \text{Ave}}{\partial z_{ij}} = [(Ave - Ave)^2 \delta_{K,i} - Ave A_{i,c} (1 - \delta_{K,i})] \delta_{c,j}$

$$\frac{\partial E^{<t>}}{\partial z_{ij}^{<t>}} = \frac{1}{m} [\hat{y}_{ij} - y_{ij}] ; \quad \frac{\partial \text{Ave}^{<t>}}{\partial a_{ij}^{<t>}} = W_{K,i} \delta_{c,j}$$

$$\begin{aligned} \frac{\partial E^{<t>}}{\partial a_{ij}^{<t>}} &= \sum_{K=1}^{N_p} \sum_{S=1}^m \frac{\partial E^{<t>}}{\partial z_{KS}^{<t>}} \cdot \frac{\partial z_{KS}^{<t>}}{\partial a_{ij}^{<t>}} = \sum_{K=1}^{N_p} \sum_{S=1}^m \frac{1}{m} [\hat{y}_{KS} - y_{KS}^{<t>}] \cdot W_{K,i} \delta_{S,j} \\ &= \sum_{K=1}^{N_p} \frac{1}{m} [\hat{y}_{rj} - y_{rj}^{<t>}] \cdot W_{r,i} \end{aligned}$$

Vectorized → $\frac{\partial E^{<t>}}{\partial a^{<t>}} = \frac{1}{m} W y^T [\hat{y}^{<t>} - y^{<t>}]$

$$\frac{\partial E^{<t>}}{\partial W_y} = \frac{1}{m} [\hat{y}^{<t>} - y^{<t>}] \cdot a^{<t>}$$

$$\frac{\partial C_{RE}^{<t>}}{\partial G_{ij}^{<t-1>}} = (\ell + \Gamma_{Ave}^{<t>}) \delta_{K,i} \delta_{G,j} ; \quad \frac{\partial \text{Ave}^{<t>}}{\partial G_{ij}^{<t>}} = \Gamma_{ave}^{<t>} (1 - tan^2(C_{ave}^{<t>})) \delta_{K,i} \delta_{c,j}$$

$$\begin{aligned} \frac{\partial E_{\text{softmax}}}{\partial C_{ij}^{<t>}} &= \frac{\partial E^{<t+1:T_y>}}{\partial C_{ij}^{<t>}} + \sum_{R=1}^{N_p} \sum_{S=1}^m \left(\frac{\partial E^{<t+1:T_y>}}{\partial a_{RS}^{<t>}} + \frac{\partial E^{<t>}}{\partial a_{RS}^{<t>}} \right) \Gamma_{ave}^{<t>} (1 - tan^2(C_{ave}^{<t>})) \delta_{R,i} \delta_{c,j} \\ &= \frac{\partial E^{<t+1:T_y>}}{\partial C_{ij}^{<t>}} + \left[\frac{\partial E^{<t+1:T_y>}}{\partial a_{ij}^{<t>}} + \frac{\partial E^{<t>}}{\partial a_{ij}^{<t>}} \right] \Gamma_{ave}^{<t>} (1 - tan^2(C_{ave}^{<t>})) \end{aligned}$$

Vectorized → $\frac{\partial E}{\partial C^{<t>}} = \frac{\partial E^{<t+1:T_y>}}{\partial C^{<t>}} + \left[\frac{\partial E^{<t+1:T_y>}}{\partial a^{<t>}} + \frac{\partial E^{<t>}}{\partial a^{<t>}} \right] \Gamma_{ave}^{<t>} (1 - tan^2(C_{ave}^{<t>}))$

$$\frac{\partial E}{\partial C_{ij}^{<+1>}} = \sum_{k=1}^{n_c} \sum_{l=1}^m \frac{\partial E}{\partial G_{kl}^{<+1>}} \cdot \underbrace{P_{f_{kl}}^{<+1>} \delta_{ki} \delta_{lj}}_{\partial G_{kl}^{<+1>} / \partial C_{ij}^{<+1>}} = \frac{\partial E}{\partial C_{ij}^{<+1>}} \cdot P_{f_{ij}}^{<+1>}$$

$$\text{Vectorized} \rightarrow \frac{\partial E}{\partial C^{<+1>}} = \frac{\partial E}{\partial C^{<+1>}} * P_f^{<+1>}$$

$$\frac{\partial E}{\partial P_{fij}^{<+1>}} = \frac{\partial E}{\partial C_{ij}^{<+1>}} \cdot C_{ij}^{<+1>} ; \quad \frac{\partial E}{\partial P_{\mu}(f)} = \frac{\partial E}{\partial C_{ij}^{<+1>}} * \hat{C}_{ij}^{<+1>} ; \quad \frac{\partial E}{\partial \hat{C}_{ij}^{<+1>}} = \frac{\partial E}{\partial C_{ij}^{<+1>}} + P_{\mu_{ij}}^{<+1>}$$

$$\text{Vectorized} \rightarrow \frac{\partial E}{\partial P_f^{<+1>}} = \frac{\partial E}{\partial C^{<+1>}} * C^{<+1>} ; \quad \frac{\partial E}{\partial P_{\mu}^{<+1>}} = \frac{\partial E}{\partial C^{<+1>}} * \hat{C}^{<+1>} ; \quad \frac{\partial E}{\partial \hat{C}^{<+1>}} = \frac{\partial E}{\partial C^{<+1>}} * P_{\mu}^{<+1>}$$

$$\frac{\partial E}{\partial P_{0ij}^{<+1>}} = \sum_{r=1}^{n_c} \sum_{s=1}^m \left[\frac{\partial E^{<+1>}}{\partial a_{rs}^{<+1>}} + \frac{\partial E^{<+1>} : T_y}{\partial a_{rs}^{<+1>}} \right] \cdot \tan(C_{ij}^{<+1>}) \delta_{ri} \delta_{sj}$$

$$\text{Vectorized} \rightarrow \frac{\partial E}{\partial P_0^{<+1>}} = \frac{\partial E}{\partial a^{<+1>}} * \tan(C^{<+1>}) ; \quad \frac{\partial E}{\partial w_0} = \left[\frac{\partial E}{\partial a^{<+1>}} + \tan(C^{<+1>}) \right] [a^{<-1>} : x^{<+1>}]^T$$

$$\frac{\partial E}{\partial [a^{<-1>} : x^{<+1>}]} = \sum_{r=1}^{n_c} \sum_{s=1}^m \left[\frac{\partial E}{\partial P_{fri}^{<+1>}} \cdot w_{fri} s_{sj} + \frac{\partial E}{\partial P_{\mu rs}^{<+1>}} \cdot w_{\mu rs} s_{sj} + \frac{\partial E}{\partial C^{<+1>}} w_{ki} s_{sj} \right]$$

$$+ \frac{\partial E}{\partial P_{0ij}^{<+1>}} \cdot w_{ki} \cdot s_{ej} \Big] = \sum_{r=1}^{n_c} \left[\frac{\partial E}{\partial P_{fri}^{<+1>}} \cdot w_{fri} + \frac{\partial E}{\partial P_{\mu rj}^{<+1>}} \cdot w_{\mu rj} + \right.$$

$$+ \frac{\partial E}{\partial P_{0ij}^{<+1>}} w_{ki} + \frac{\partial E}{\partial C^{<+1>}} w_{ki} \Big]$$

Vectorized

$$\frac{\partial E}{\partial [a^{<-1>} : x^{<+1>}]} = \overbrace{\frac{\partial E}{\partial P_f^{<+1>}} \cdot w_f + w_{\mu}^T \frac{\partial E}{\partial P_{\mu}^{<+1>}} + w_c^T \frac{\partial E}{\partial C^{<+1>}} + w_0^T \frac{\partial E}{\partial P_0^{<+1>}} * P_0^{<+1>} * (1 - P_0^{<+1>})}^{* P_f(1-P_f)} + \overbrace{\underbrace{w_{\mu}^T \frac{\partial E}{\partial P_{\mu}^{<+1>}} * P_{\mu}(1-P_{\mu})}_{\text{hybrids}} + \underbrace{w_c^T \frac{\partial E}{\partial C^{<+1>}} * (1-C^{<+1>})}_{(1-C^{<+1>})}}_{(1-P_0^{<+1>})}$$

$$[a^{<-1>} : x^{<+1>}] \equiv \text{distr}(n_a + n_x, w)$$

$$\hookrightarrow \frac{\partial E}{\partial a^{<-1>}} = \frac{\partial E}{\partial [a^{<-1>} : x^{<+1>}]} \quad [: n_a, m] ; \quad \frac{\partial E}{\partial x^{<+1>}} = \frac{\partial E}{\partial [a^{<-1>} : x^{<+1>}]} \quad [n_a : n_x, m]$$

Kalman \rightarrow

$$\frac{\partial E}{\partial w_f} = \left[\frac{\partial E}{\partial P_{f,t+>}} * P_{f,t+>} * (1 - P_{f,t+>}) \right] \cdot [a^{t-1}, x^{t+>}]^T$$

$$\frac{\partial E}{\partial w_u} = \left[\frac{\partial E}{\partial P_{u,t+>}} * P_{u,t+>} * (1 - P_{u,t+>}) \right] \cdot [a^{t-1}, x^{t+>}]^T$$

$$\frac{\partial E}{\partial w_b} = \left[\frac{\partial E}{\partial C_{t+>}} * (1 - C_{t+>}^2) \right] \cdot [a^{t-1}, x^{t+>}]^T$$

