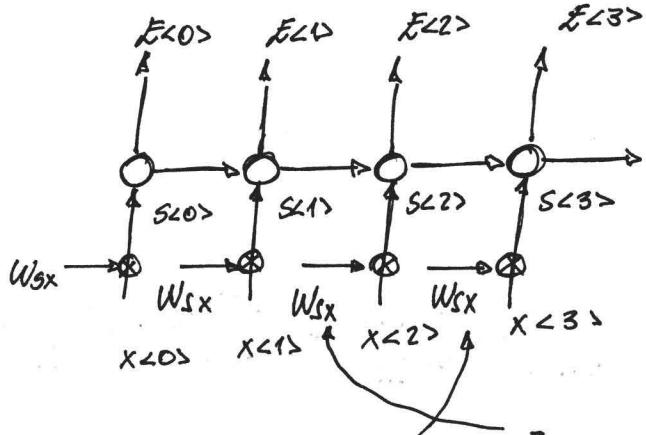


RNN Network Backprop through time



$$\frac{\partial E}{\partial W_{sx}} = \frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial W_{sx}} + \left[\frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} + \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} \right]$$

$$+ \left[\frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} + \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} + \right.$$

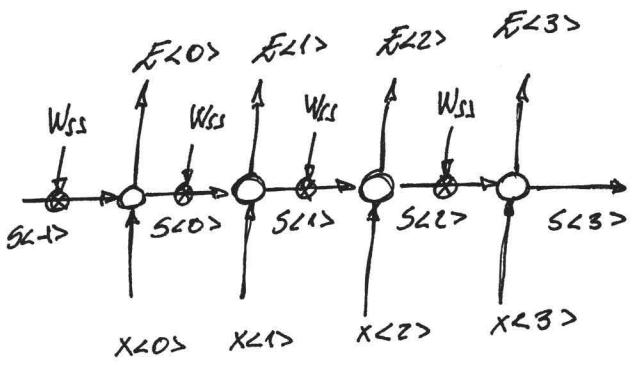
$$+ \left. \frac{\partial E^{<1>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial W_{sx}} \right] + \left[\frac{\partial E^{<3>}}{\partial s^{<3>}} \cdot \frac{\partial s^{<3>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} + \right.$$

$$+ \frac{\partial E^{<2>}}{\partial s^{<2>}} \cdot \frac{\partial s^{<2>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} + \frac{\partial E^{<1>}}{\partial s^{<1>}} \cdot \frac{\partial s^{<1>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} +$$

$$+ \left. \frac{\partial E^{<0>}}{\partial s^{<0>}} \cdot \frac{\partial s^{<0>}}{\partial W_{sx}} \right] = \sum_{k=3}^3 \frac{\partial E^{<k>}}{\partial s^{<k>}} \cdot \prod_{j=3+1}^K \left(\frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<3>}}{\partial W_{sx}}$$

$$+ \sum_{k=2}^3 \frac{\partial E^{<k>}}{\partial s^{<k>}} \cdot \prod_{j=3+1}^K \left(\frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<2>}}{\partial W_{sx}} + \dots =$$

$$\frac{\partial E}{\partial W_{sx}} = \sum_{t=0}^{T_x} \sum_{k=c}^{T_y} \frac{\partial E^{<k>}}{\partial s^{<k>}} \prod_{j=l+1}^K \left(\frac{\partial s^{<j>}}{\partial s^{<j-1>}} \right) \cdot \frac{\partial s^{<l>}}{\partial W_{sx}}$$



$$\begin{aligned}
 \frac{\partial E}{\partial w_{ss}} &= \left[\frac{\partial E_{<3>}}{\partial s_{<3>}} \frac{\partial s_{<3>}}{\partial w_{ss}} \right] + \left[\frac{\partial E_{<3>}}{\partial s_{<3>}} \cdot \frac{\partial s_{<3>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial w_{ss}} + \frac{\partial E_{<2>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial w_{ss}} \right] + \\
 &+ \left[\frac{\partial E_{<3>}}{\partial s_{<3>}} \cdot \frac{\partial s_{<3>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial w_{ss}} + \frac{\partial E_{<2>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial w_{ss}} + \right. \\
 &+ \left. \frac{\partial E_{<1>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial w_{ss}} \right] + \left[\frac{\partial E_{<0>}}{\partial s_{<0>}} \cdot \frac{\partial s_{<0>}}{\partial w_{ss}} \right] + \frac{\partial E_{<3>}}{\partial s_{<3>}} \cdot \frac{\partial s_{<3>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial s_{<0>}} \Rightarrow \\
 &\frac{\partial E_{<0>}}{\partial w_{ss}} + \frac{\partial E_{<2>}}{\partial s_{<2>}} \cdot \frac{\partial s_{<2>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial s_{<0>}} \cdot \frac{\partial s_{<0>}}{\partial w_{ss}} + \frac{\partial E_{<1>}}{\partial s_{<1>}} \cdot \frac{\partial s_{<1>}}{\partial s_{<0>}} \cdot \frac{\partial s_{<0>}}{\partial w_{ss}}
 \end{aligned}$$

$$\frac{\partial E}{\partial w_{ss}} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_k}{\partial s_{<k>}} \prod_{j=l+1}^K \left(\frac{\partial s_{<j>}}{\partial s_{<j-1>}} \right) \frac{\partial s_{<l>}}{\partial w_{ss}}$$

$$\frac{\partial E}{\partial x_{<l>}} = \sum_{k=c}^{T_y} \frac{\partial E_k}{\partial s_{<k>}} \prod_{j=l+1}^K \left(\frac{\partial s_{<j>}}{\partial s_{<j-1>}} \right) \frac{\partial s_{<l>}}{\partial x_{<l>}}$$

* These consecutive multiplications may cause vanishing gradient or exploding gradient:

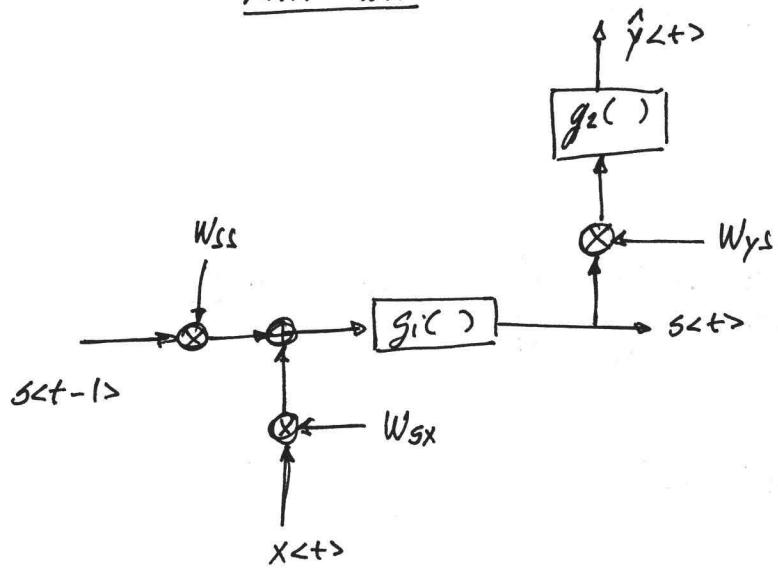
- Exploding gradient $\rightarrow \frac{\partial E}{\partial \theta} \rightarrow \frac{\partial E}{\partial \theta} > \text{Threshold} \rightarrow \frac{\partial E}{\partial \theta} = \frac{\text{Thres}}{\frac{\partial E}{\partial \theta}} \cdot \frac{\partial E}{\partial \theta}$

- Vanishing gradient $\rightarrow \text{LSTM, GRU}$

$$\frac{\partial s_{<l>}}{\partial s_{<l>}} = \tanh'(\tau_{<l>}) \cdot w_{sk} \delta_j$$

$0 \cdot x \leftarrow \begin{cases} z = -1 & \\ z = 1 & \end{cases} \quad | \quad \begin{cases} < 1 : 0 \cdot x^T \rightarrow \text{Vanishing} \\ > 1 : 1 \cdot x^T \rightarrow \text{Exploding} \end{cases}$

RNN Cell



* Forward \rightarrow

$$s^{t+>} = g_1(z^{t+>})$$

$$z^{t+>} = W_{ss} s^{t-1>} +$$

$$W_{sx} x^{t+>} + b_s$$

$$\hat{y}^{t+>} = g_2(p^{t+>})$$

$$p^{t+>} = W_{ys} s^{t+>} + b_y$$

$$\text{Dimensions} \rightarrow x^{t+>} \in \mathbb{R}^{(n_x, m)}; s^{t+>} \in \mathbb{R}^{(n_s, m)}; y^{t+>} \in \mathbb{R}^{(n_y, m)}$$

$$W_{ss} \in \mathbb{R}^{(n_s, n_s)}; W_{sx} \in \mathbb{R}^{(n_s, n_x)}; W_{ys} \in \mathbb{R}^{(n_y, n_s)}$$

$$\frac{\partial s_{ke}^{t+>}}{\partial z_{ij}^{t+>}} = g_1'(z_{ke}^{t+>}) \delta_{ki} \delta_{ej}; \quad \frac{\partial \hat{y}_{ke}^{t+>}}{\partial p_{ij}^{t+>}} = g_2'(s_{ke}^{t+>}) \delta_{ki} \delta_{ej}$$

$$\frac{\partial z_{ke}^{t+>}}{\partial W_{ss,ij}} = \delta_{ki} s_{je}^{t+>}; \quad \frac{\partial z_{ke}^{t+>}}{\partial W_{sx,ij}} = \delta_{ki} x_{je}^{t+>}; \quad \frac{\partial z_{ke}^{t+>}}{\partial s_{ij}^{t-1>}} = W_{ss,ki} \delta_{ej}$$

$$\frac{\partial z_{ke}^{t+>}}{\partial x_{ij}^{t+>}} = W_{sx,ki} \delta_{ej}; \quad \frac{\partial p_{ke}^{t+>}}{\partial W_{ys,ij}} = \delta_{ki} s_{je}^{t+>}; \quad \frac{\partial p_{ke}^{t+>}}{\partial s_{ij}^{t+>}} = W_{ys,ki} \delta_{ej}$$

$$\frac{\partial s_{ke}^{t+>}}{\partial x_{ij}^{t+>}} = \sum_{r=1}^{n_a} \sum_{s=1}^m \frac{\partial s_{ke}^{t+>}}{\partial z_{rs}^{t+>}} \frac{\partial z_{rs}^{t+>}}{\partial x_{ij}^{t+>}} = \sum_{r=1}^{n_a} \sum_{s=1}^m g_1'(z_{ke}^{t+>}) \delta_{ki} \delta_{sj}$$

$$\cdot W_{sx,ki} \delta_{sj} = g_1'(z_{ke}^{t+>}) \cdot W_{sx,ki} \delta_{ej}; \quad \frac{\partial s_{ke}^{t+>}}{\partial s_{ij}^{t-1>}} = g_1'(z_{ke}^{t+>}) \cdot W_{ss,ki} \delta_{ej}$$

$$\begin{aligned} \frac{\partial s_{ke}^{t+>}}{\partial W_{ss,ij}} &= \sum_{r=1}^{n_a} \sum_{s=1}^m \frac{\partial s_{ke}^{t+>}}{\partial z_{rs}^{t+>}} \cdot \frac{\partial z_{rs}^{t+>}}{\partial W_{ss,ij}} = \sum_{r=1}^{n_a} \sum_{s=1}^m g_1'(z_{ke}^{t+>}) \delta_{kr} \delta_{is} \delta_{ri} s_{js}^{t-1>} \\ &= g_1'(z_{ke}^{t+>}) \delta_{je} \delta_{ki} \rightarrow s_{je}^{t-1>} \end{aligned}$$

$$\frac{\partial \hat{y}_{Kc}^{<t>}}{\partial \cancel{w_{x,i,j}}} = g'_1(z_{Kc}^{<t>}) \cdot \cancel{g_2'(p_{Kc}^{<t>})} \cdot \cancel{x_{je}^{<t>}} \delta_{K,i}$$

$$\frac{\partial \hat{y}_{Kc}^{<t>}}{\partial s_{i,j}^{<t>}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{\partial \hat{y}_{Kc}^{<t>}}{\partial p_{rs}^{<t>}} \cdot \frac{\partial p_{rs}^{<t>}}{\partial s_{i,j}^{<t>}} = \sum_{r=1}^{n_y} \sum_{s=1}^m g'_2(p_{Kc}^{<t>}) \delta_{K,r} \cancel{\delta_{i,s}}$$

$$w_{rc} \delta_{s,j} = g'_2(p_{Kc}^{<t>}) \cdot w_{s,Kr} \delta_{e,j}$$

$$\frac{\partial \hat{y}_{Kc}^{<t>}}{\partial w_{x,i,j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{\partial \hat{y}_{Kc}^{<t>}}{\partial p_{rs}^{<t>}} \frac{\partial p_{rs}^{<t>}}{\partial w_{x,i,j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m g'_2(p_{Kc}^{<t>}) \delta_{K,r} \delta_{e,s}$$

$$\cdot s_{j,e}^{<t>} \delta_{r,i} = g'_2(p_{Kc}^{<t>}) \cdot s_j \cancel{s_e^{<t>}} \delta_{K,i}$$

$$\frac{\partial \hat{y}_{Kc}^{<t>}}{\partial w_{x,i,j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{\partial \hat{y}_{Kc}^{<t>}}{\partial s_{rs}^{<t>}} \frac{\partial s_{rs}^{<t>}}{\partial w_{x,i,j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m g'_1(z_{Kc}^{<t>}) \cdot w_{s,Kr} \delta_{e,s}$$

$$\cdot g'_1(z_{Kc}^{<t>}) \cdot x_{js}^{<t>} \delta_{r,i} = g'_2(p_{Kc}^{<t>}) \cdot w_{s,Kr} g'(z_{ie}^{<t>}) x_{je}^{<t>}$$

$$\frac{\partial \hat{y}_{Kc}^{<t>}}{\partial s_{i,j}^{<t-1>}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{\partial \hat{y}_{Kc}^{<t>}}{\partial g_1' z_{rs}^{<t>}} \cdot \frac{\partial g_1' z_{rs}^{<t>}}{\partial s_{i,j}^{<t-1>}} = \sum_{r=1}^{n_y} \sum_{s=1}^m g'_2(p_{Kc}^{<t>}) \cdot w_{s,Kr} \delta_{e,s}$$

$$\cdot g'_1(z_{rs}^{<t>}) \cdot w_{s,r,i} \cdot \delta_{s,j} =$$

$$= \sum_{r=1}^{n_y} g'_2(p_{Kc}^{<t>}) \cdot w_{s,Kr} \cdot g'_1(z_{rj}^{<t>}) w_{s,r,i} \delta_{e,j}$$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{ij} \langle t-1 \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot \frac{\partial S_{kl} \langle t \rangle}{\partial S_{ij} \langle t-1 \rangle} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot g'_k(z_{kl} \langle t \rangle) \cdot w_{sski} \cdot s_{lj}$$

$$= \sum_{k=1}^{n_s} \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kj} \langle t \rangle} \cdot g'_k(z_{kj} \langle t \rangle) \cdot w_{sski} \rightarrow$$

vektorielle
Implementation $\rightarrow \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{ij} \langle t-1 \rangle} = W_{ss}^T \left(\frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_i(z \langle t \rangle) \right)$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial X_{ij} \langle t \rangle} = W_{sx}^T \left(\frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_i(z \langle t \rangle) \right)$$

$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{ss,ij}} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot \frac{\partial S_{kl} \langle t \rangle}{\partial w_{ss,ij}} = \sum_{k=1}^{n_s} \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{kl} \langle t \rangle} \cdot g'_k(z_{kl} \langle t \rangle) \cdot s_{jl} \langle t-1 \rangle \cdot s_{ki}$$

$$= \sum_{l=1}^m \frac{\partial E \langle t+1 : T_y \rangle}{\partial S_{il} \langle t \rangle} \cdot g'_i(z_{il} \langle t \rangle) \cdot s_{je} \langle t-1 \rangle$$

vektorielle
Implementation $\rightarrow \frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{ss}} = \left(\frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_i(z \langle t \rangle) \right) \cdot S^T \langle t-1 \rangle$

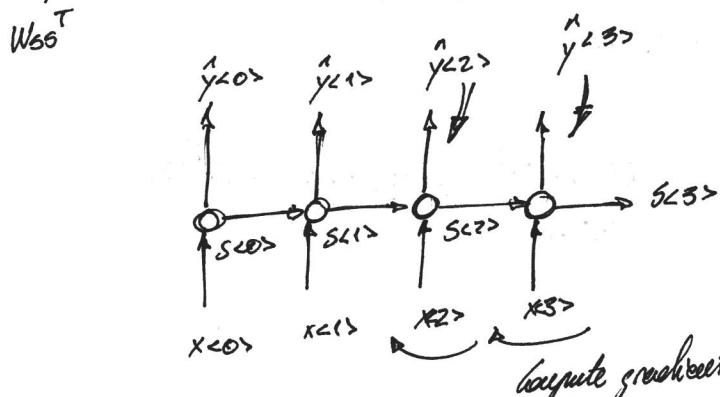
$$\frac{\partial E \langle t+1 : T_y \rangle}{\partial w_{sx}} = \left(\frac{\partial E \langle t+1 : T_y \rangle}{\partial S \langle t \rangle} * g'_i(z \langle t \rangle) \right) \cdot X^T \langle t \rangle$$

$$\begin{aligned}
 \frac{\partial E_{t+>}}{\partial \delta_{ij}^{<t+>}} &= \sum_{k=1}^{n_y} \sum_{l=1}^m \frac{\partial E_{t+>}}{\partial y_{kl}^{<t+>}} \cdot \frac{\partial y_{kl}^{<t+>}}{\partial \delta_{ij}^{<t+>}} = \sum_{k=1}^{n_y} \sum_{l=1}^m \frac{\partial E_{t+>}}{\partial y_{kl}^{<t+>}} \cdot g_2'(\rho_{kl}^{<t+>}) \cdot w_{y_{kl}} \cdot \delta_{ij} \\
 &= \sum_{k=1}^{n_y} \frac{\partial E_{t+>}}{\partial y_{kj}^{<t+>}} \cdot g_2'(\rho_{kj}^{<t+>}) \cdot w_{y_{kj}} - \\
 \frac{\partial E_{t+>}}{\partial w_{ys_{ij}}} &= \sum_{k=1}^{n_y} \sum_{l=1}^m \frac{\partial E_{t+>}}{\partial y_{kl}^{<t+>}} \cdot \frac{\partial y_{kl}^{<t+>}}{\partial w_{ys_{ij}}} = \sum_{k=1}^{n_y} \sum_{l=1}^m \frac{\partial E_{t+>}}{\partial y_{kl}^{<t+>}} \cdot g_2'(\rho_{kl}^{<t+>}) \cdot \delta_{je}^{<t+>} \delta_{ki} \\
 &= \sum_{l=1}^m \frac{\partial E_{t+>}}{\partial y_{le}^{<t+>}} \cdot g_2'(\rho_{le}^{<t+>}) \cdot \delta_{je}^{<t+>}
 \end{aligned}$$

Kettenregel
Implementation $\rightarrow \frac{\partial E_{t+>}}{\partial s_{<t+>}} = W_{yc}^T \left(\frac{\partial E_{t+>}}{\partial \hat{y}_{<t+>}} * g_2'(\rho_{<t+>}) \right)$

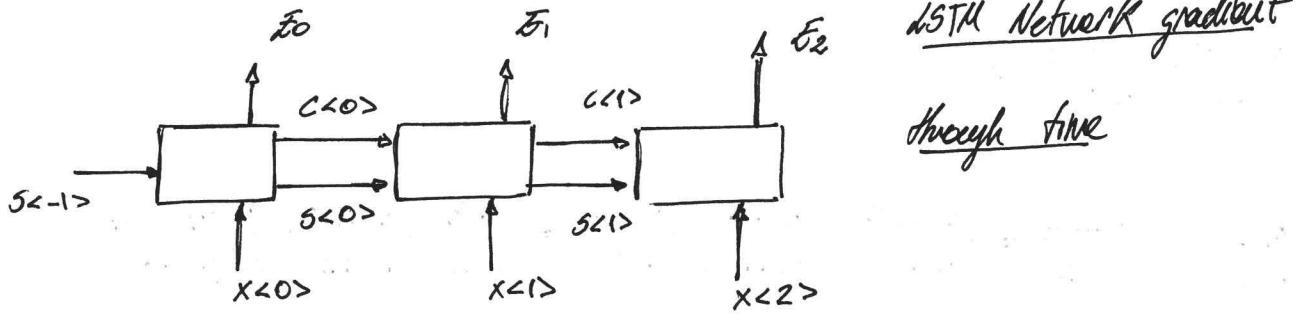
$$\frac{\partial E_{t+>}}{\partial w_{ys}} = \left(\frac{\partial E_{t+>}}{\partial \hat{y}_{<t+>}} * g_2'(\rho_{<t+>}) \right) \cdot \delta_{e,t+>}^T$$

$$\frac{\partial E}{\partial s_{<t+>}} = \left[W_{yc}^T \left(\frac{\partial E_{t+>}}{\partial \hat{y}_{<t+>}} * g_2'(\rho_{<t+>}) \right) * j_1'(z_{t+>}) \right] + W_{ss}^T \left(\frac{\partial E_{t+1:T_X}}{\partial s_{<t+>}} * g_1'(z_{t+>}) \right)$$



Implementation.

Compute gradients, iterate across tree & accumulate grad.



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_{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_{C0}} \cdot \frac{\partial S_{C0}}{\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 S_{C0}} \cdot \frac{\partial S_{C0}}{\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 S_{C0}} \cdot \frac{\partial S_{C0}}{\partial W_f} \cdot \frac{\partial S_{C0}}{\partial C_{C0}} \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 W_f} + \frac{\partial E_1}{\partial S_{C1}} \cdot \frac{\partial S_{C1}}{\partial C_{C1}} \cdot \frac{\partial C_{C1}}{\partial S_{C0}} \cdot \frac{\partial S_{C0}}{\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}$$

$$\begin{aligned}
&= \frac{\partial E_2}{\partial S<2>} \cdot \frac{\partial S<2>}{\partial C<2>} \cdot \frac{\partial C<2>}{\partial W_f} + \\
&+ \left[\frac{\partial E_2}{\partial S<2>} \cdot \frac{\partial S<2>}{\partial C<2>} \left(\frac{\partial C<2>}{\partial C<1>} + \frac{\partial C<2>}{\partial S<1>} \cdot \frac{\partial S<1>}{\partial C<1>} \right) \frac{\partial C<1>}{\partial W_f} + \frac{\partial E_1}{\partial S<1>} \frac{\partial S<1>}{\partial C<1>} \cdot \frac{\partial C<1>}{\partial W_f} \right] + \\
&+ \left[\frac{\partial E_2}{\partial S<2>} \cdot \frac{\partial S<2>}{\partial C<2>} \left(\frac{\partial C<2>}{\partial C<1>} + \frac{\partial C<2>}{\partial S<1>} \cdot \frac{\partial S<1>}{\partial C<1>} \right) \left(\frac{\partial C<1>}{\partial C<0>} + \frac{\partial C<1>}{\partial S<0>} \cdot \frac{\partial S<0>}{\partial C<0>} \right) \frac{\partial C<0>}{\partial W_f} \right. \\
&\quad \left. + \frac{\partial E_1}{\partial S<1>} \cdot \frac{\partial S<1>}{\partial C<1>} \left(\frac{\partial C<1>}{\partial C<0>} + \frac{\partial C<1>}{\partial S<0>} \cdot \frac{\partial S<0>}{\partial C<0>} \right) \frac{\partial C<0>}{\partial W_f} \right. \\
&\quad \left. + \frac{\partial E_1}{\partial S<0>} \cdot \frac{\partial S<0>}{\partial C<0>} \cdot \frac{\partial C<0>}{\partial W_f} \right] \Rightarrow
\end{aligned}$$

$$\begin{aligned}
\frac{\partial F}{\partial W_f} &= \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_K}{\partial S<k>} \cdot \frac{\partial S<k>}{\partial C<k>} \left[\prod_{j=l+1}^K \left(\frac{\partial C<j>}{\partial C<j-1>} + \frac{\partial C<j>}{\partial S<j-1>} \cdot \frac{\partial S<j-1>}{\partial C<j-1>} \right) \right] \frac{\partial C<l>}{\partial W_f} \\
\frac{\partial F}{\partial W_h} &= \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_K}{\partial S<k>} \cdot \frac{\partial S<k>}{\partial C<k>} \left[\prod_{j=l+1}^K \left(\frac{\partial C<j>}{\partial C<j-1>} + \frac{\partial C<j>}{\partial S<j-1>} \cdot \frac{\partial S<j-1>}{\partial C<j-1>} \right) \right] \frac{\partial C<l>}{\partial W_h} \\
\frac{\partial F}{\partial W_c} &= \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_K}{\partial S<k>} \cdot \frac{\partial S<k>}{\partial C<k>} \left[\prod_{j=l+1}^K \left(\frac{\partial C<j>}{\partial C<j-1>} + \frac{\partial C<j>}{\partial S<j-1>} \cdot \frac{\partial S<j-1>}{\partial C<j-1>} \right) \right] \frac{\partial C<l>}{\partial W_c} \\
\frac{\partial F}{\partial W_o} &= \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_K}{\partial S<k>} \cdot \prod_{j=l+1}^K \left(\frac{\partial S<j>}{\partial S<j-1>} \right) \frac{\partial S<l>}{\partial W_o}
\end{aligned}$$

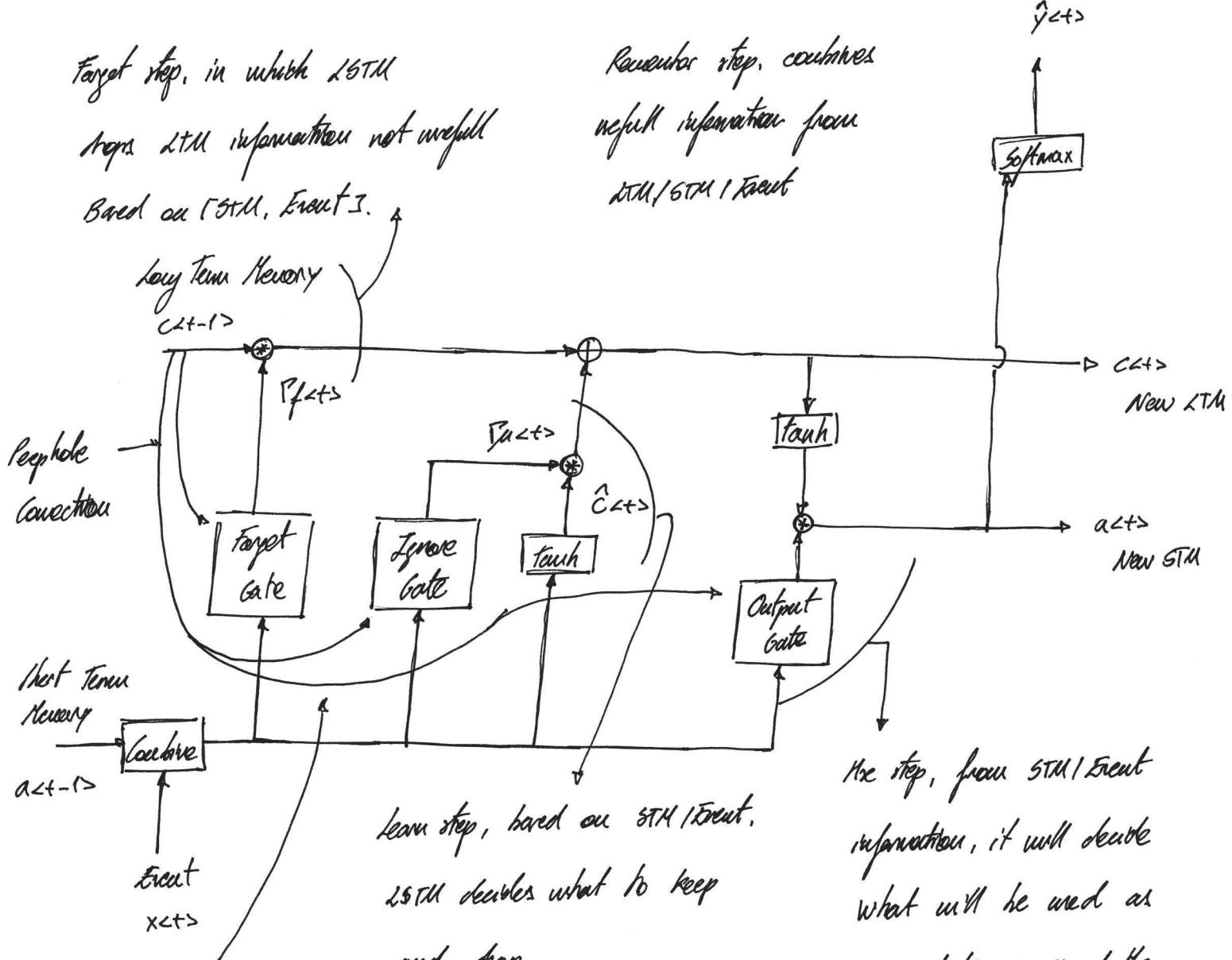
$$\frac{\partial C<j>}{\partial C<j-1>} = P_{f<+>} \rightarrow \prod \frac{\partial C<j>}{\partial C<j-1>}$$

with next round or explore if
the size of $P_{f<+>}$ is on the limit
part \rightarrow value (0,1) depending if
it had importance \rightarrow
If causaled and f^+ , enabled
update.

For this gradient, it won't vanish because
of the path until previous block $\frac{\partial C<t+1>}{\partial C<t>} \cdot \frac{\partial C<t>}{\partial W_o}$

$$\frac{\partial \tilde{E}}{\partial x^{<\ell>}} = \sum_{K=\ell}^{Ty} \frac{\partial E_K}{\partial s^{<K>}} \cdot \frac{\partial s^{<K>}}{\partial c^{<K>}} \left[\prod_{j=\ell+1}^K \left(\frac{\partial c^{<j>}}{\partial c^{<j-1>}} + \frac{\partial c^{<j>}}{\partial s^{<j-1>}} \cdot \frac{\partial s^{<j-1>}}{\partial c^{<j-1>}} \right) \right] \frac{\partial c^{<\ell>}}{\partial x^{<\ell>}}$$

day 9: Long Term Memory



Peephole connections:

LSTM variant, includes STM in gate decisions.

Learn step, based on STM/Event.

LSTM decides what to keep and drop.

Useful information of STM/Event.

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

$$f_t = \sigma(W_f[a_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i[a_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_c[a_{t-1}, x_t] + b_c)$$

$$c_t = c_{t-1} * f_t + i_t * \tilde{c}_t$$

$$o_t = \sigma(W_o[a_{t-1}, x_t] + b_o)$$

$$a_t = o_t * \tanh(c_t)$$

Cell backprop →

$$\hat{y}^{<t>} = \text{softmax}(z^{<t>}) ; z^{<t>} = W_y a^{<t>} + b_y$$

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

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

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

$$\begin{aligned} \frac{\partial E^{<t>}}{\partial a_{ij}^{<t>}} &= \sum_{k=1}^{n_y} \sum_{s=1}^m \frac{\partial E^{<t>}}{\partial a_{ks}^{<t>}} \cdot \frac{\partial a_{ks}^{<t>}}{\partial a_{ij}^{<t>}} = \sum_{k=1}^{n_y} \sum_{s=1}^m \frac{1}{m} [\hat{y}_{rs} - y_{rs}^{<t>}] \cdot W_{y_{k,i}} \delta_{s,j} \\ &= \sum_{k=1}^{n_y} \frac{1}{m} [\hat{y}_{rj} - y_{rj}^{<t>}] \cdot W_{y_{k,i}} \end{aligned}$$

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

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

$$\frac{\partial C_{RE}^{<t>}}{\partial c_{ij}^{<t-1>}} = (\epsilon + \Gamma_{A_{RE}}^{<t>}) \delta_{k,i} \delta_{e,j} ; \frac{\partial A_{RE}^{<t>}}{\partial c_{ij}^{<t>}} = \Gamma_{O_{RE}}^{<t>} (1 - \tanh^2(C_{RE}^{<t>})) \delta_{k,i} \delta_{e,j}$$

$$\begin{aligned} \frac{\partial E_{\text{absolute}}}{\partial c_{ij}^{<t>}} &= \frac{\partial E^{<t+1:T_y>}}{\partial c_{ij}^{<t>}} + \sum_{r=1}^{n_c} \sum_{s=1}^m \left(\frac{\partial E^{<t+1:T_y>}}{\partial a_{rs}^{<t>}} + \frac{\partial E^{<t>}}{\partial a_{rs}^{<t>}} \right) \Gamma_{O_{rs}}^{<t>} (1 - \tanh^2(C_{ij}^{<t>})) \delta_{r,i} \delta_{s,j} \\ &= \frac{\partial E^{<t+1:T_y>}}{\partial c_{ij}^{<t>}} + \left[\frac{\partial E^{<t+1:T_y>}}{\partial a_{ij}^{<t>}} + \frac{\partial E^{<t>}}{\partial a_{ij}^{<t>}} \right] \Gamma_{O_{ij}}^{<t>} (1 - \tanh^2(C_{ij}^{<t>})) \end{aligned}$$

Vectorized → $\frac{\partial E}{\partial c^{<t>}} = \frac{\partial E^{<t+1:T_y>}}{\partial c^{<t>}} + \left[\frac{\partial E^{<t+1:T_y>}}{\partial a^{<t>}} + \frac{\partial E^{<t>}}{\partial a^{<t>}} \right] * \Gamma_{O^{<t>}} * (1 - \tanh^2(c^{<t>}))$

$$\frac{\partial \delta}{\partial G_{ij}^{<+1>}} = \sum_{K=1}^{N_c} \sum_{l=1}^m \frac{\partial \delta}{\partial G_{kl}^{<+>}} \cdot \underbrace{P_{f_{kl}}^{<+>} \delta_{kl} \delta_{ij}}_{\partial G_{kl}^{<+>} / \partial G_{ij}^{<+1>}} = \frac{\partial \delta}{\partial G_{kj}^{<+>}} \cdot P_{f_{ij}}^{<+>}$$

$$\text{Varied} \rightarrow \frac{\partial E}{\partial C^{<+1>}} = \frac{\partial E}{\partial C^{<+3>}} + V_f^{<+3>}$$

$$\frac{\partial F}{\partial P_{ij}^{<+>}} = \frac{\partial F}{\partial C_{ij}^{<+>}} \cdot C_{ij}^{<+ -1>} ; \quad \frac{\partial F}{\partial P_{ij}(t)} = \frac{\partial F}{\partial C_{ij}^{<+>}} * C_{ij}^{<+>} ; \quad \frac{\partial F}{\partial E_{ij}^{<+>}} = \frac{\partial F}{\partial C_{ij}^{<+>}} + C_{ij}^{<+>}$$

$$\text{Vektorielles } \rightarrow \frac{\partial E}{\partial P_{\mu}^{L+}} = \frac{\partial E}{\partial C^{L+}} * C^{L+} ; \frac{\partial E}{\partial P_{\mu}^{L+}} = \frac{\partial E}{\partial \tilde{C}^{L+}} * \tilde{C}^{L+} ; \frac{\partial E}{\partial \tilde{C}^{L+}} = \frac{\partial E}{\partial C^{L+}} * P_{\mu}^{L+}$$

$$\frac{\partial f}{\partial P_{ij}^{<+>}} = \sum_{r=1}^{n_r} \sum_{s=1}^{m_s} \left[\frac{\partial f^{<+>}}{\partial a_{rs}^{<+>}} + \frac{\partial f^{<+1:T_y>}}{\partial a_{rs}^{<+>}} \right] \cdot \tan(C_{ij}^{<+>}) \quad \text{for } i \in S, j \in T$$

$$\text{Vektorisiert} \rightarrow \frac{\partial E}{\partial \Gamma_0^{C \leftrightarrow S}} = \frac{\partial E}{\partial \alpha^{C \leftrightarrow S}} * \text{ham}(C \leftrightarrow S) ; \quad \frac{\partial E}{\partial W_0} = \left[\frac{\partial E}{\partial \alpha^{C \leftrightarrow S}} * \text{ham}(C \leftrightarrow S) \right] [\alpha^{C \leftrightarrow S}, x^{C \leftrightarrow S}]^T$$

$$\frac{\partial F}{\partial [\alpha_{\leftarrow \rightarrow}, x_{\leftarrow \rightarrow}]} = \sum_{r=1}^{n_c} \sum_{s=1}^m \left[\frac{\partial F}{\partial p_{rs}^{\leftarrow \rightarrow}} \cdot w_{fri\; ssj} + \frac{\partial F}{\partial p_{rs}^{\rightarrow \leftarrow}} \cdot w_{\mu_{rs}^s ssj} + \frac{\partial F}{\partial c_{\leftarrow \rightarrow}} w_{ksi\; ssj} \right]$$

$$+ \frac{\partial \Delta}{\partial P_{fri}^{ij}} \cdot W_{fri} \cdot S_{fri} \Big] = \sum_{r=1}^R \left[\frac{\partial \Delta}{\partial P_{fri}^{ij}} \cdot W_{fri} + \frac{\partial \Delta}{\partial P_{Mrij}} \cdot W_{Mrij} + \right]$$

$$+ \frac{\partial E}{\partial T_{\text{env}} \leftarrow \rightarrow} W_{K,i} + \frac{\partial E}{\partial C \leftarrow \rightarrow} W_{C,i}]$$

Vektorrech

$$\frac{\partial E}{\partial [a_{ct+1}, x_{t+1}]} = \underbrace{\frac{\partial E}{\partial p_{ct+1}} \cdot w_f^T}_{\leftarrow * P_f(1-P_f)} + w_n^T \underbrace{\frac{\partial E}{\partial p_{ct+1}}}_{\leftarrow * P_n(1-P_n)} + w_c^T \underbrace{\frac{\partial E}{\partial c_{t+1}}}_{\substack{\leftarrow (1-\tan^2 c) \\ \text{symbol}}} + w_o^T \underbrace{\frac{\partial E}{\partial o_{t+1}}}_{\substack{\leftarrow (1-\bar{c}_{t+1})}} * P_o^{ct+1} * (1-P_o^{ct+1})$$

$$[a<+,-1>, x<+>] \equiv \text{distr}(Na + Nx, u)$$

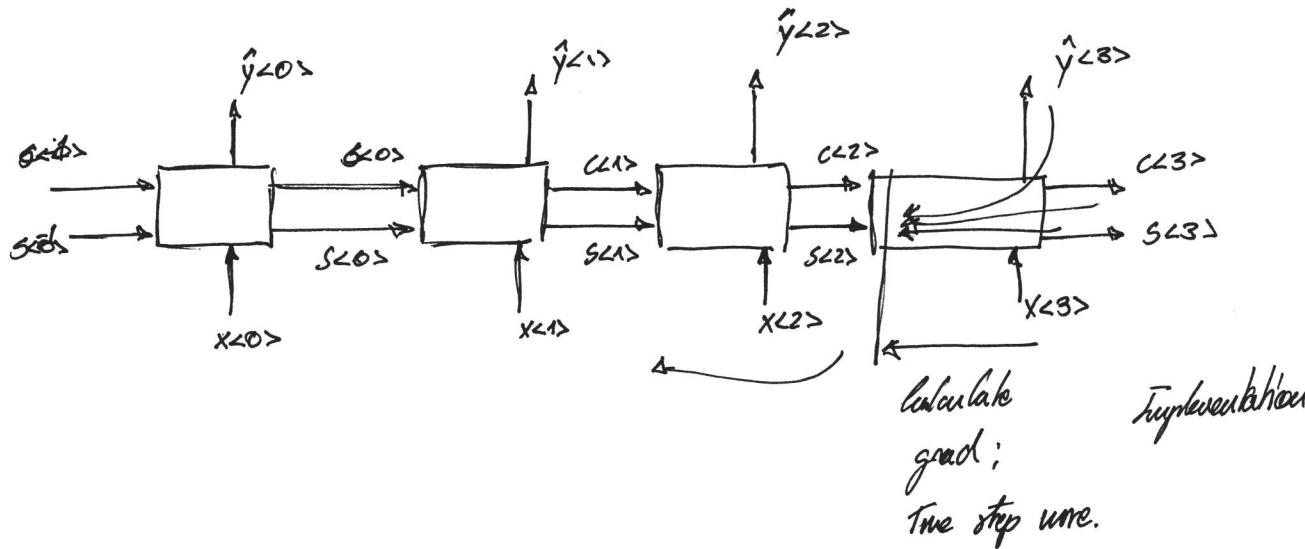
$$\rightarrow \frac{\partial F}{\partial x^{t+1}} = \frac{\partial F}{\partial x^{t+1}, \alpha^{t+1}} [: n_a, m] ; \frac{\partial E}{\partial x^t} = \frac{\partial E}{\dots} [n_a : n_a + n_k, m]$$

Vektorisiert \rightarrow

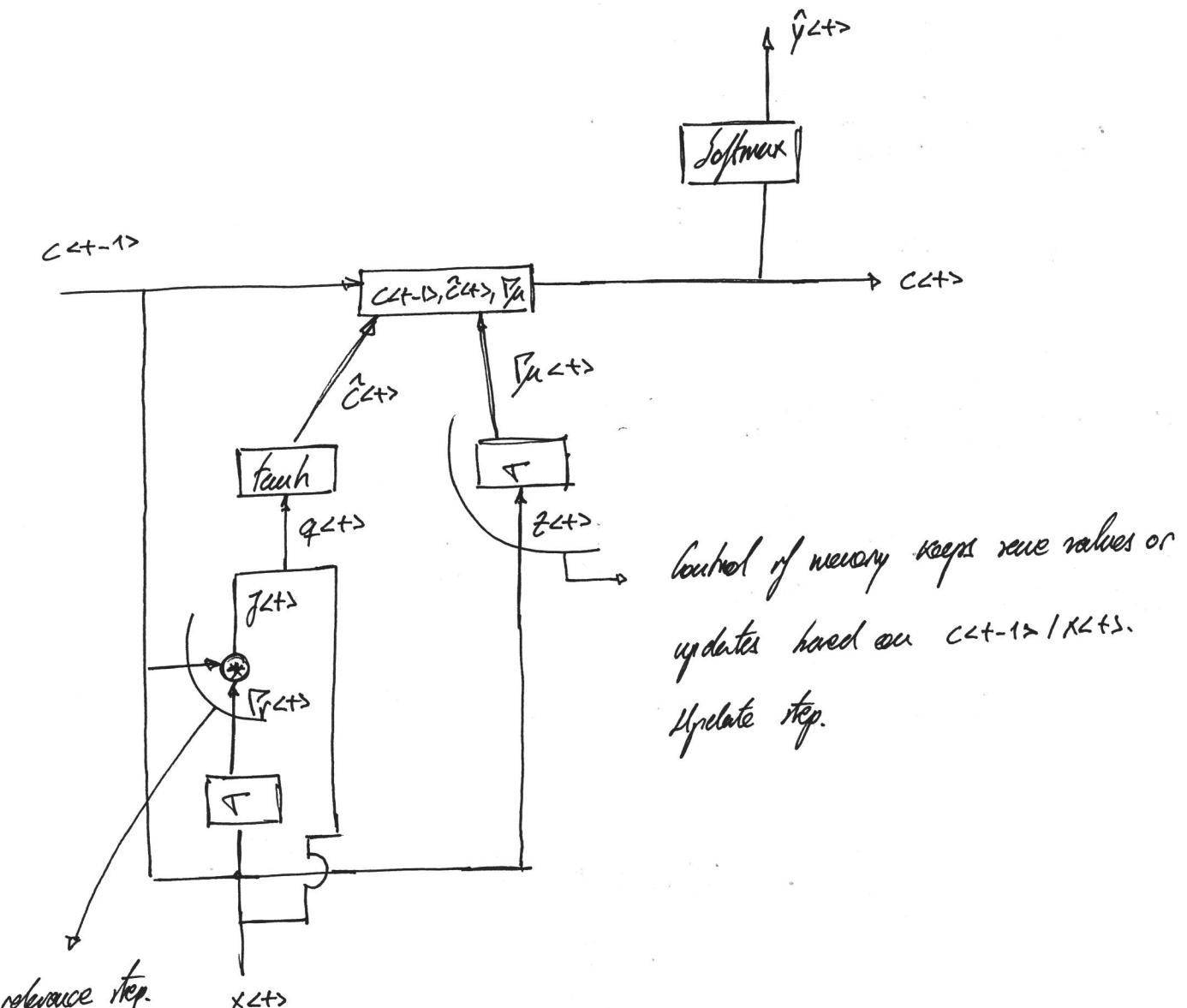
$$\frac{\partial E}{\partial w_f} = \left[\frac{\partial E}{\partial p_f^{<t>}} * p_f^{<t>} * (1 - p_f^{<t>}) \right] \cdot [a^{<t-1>}, x^{<t>}]^T$$

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

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



Gated Recurrent Unit



control of memory keeps same values or updates based on C^{t+1-1} / x^t .
Update step.

Reset, relevance step.

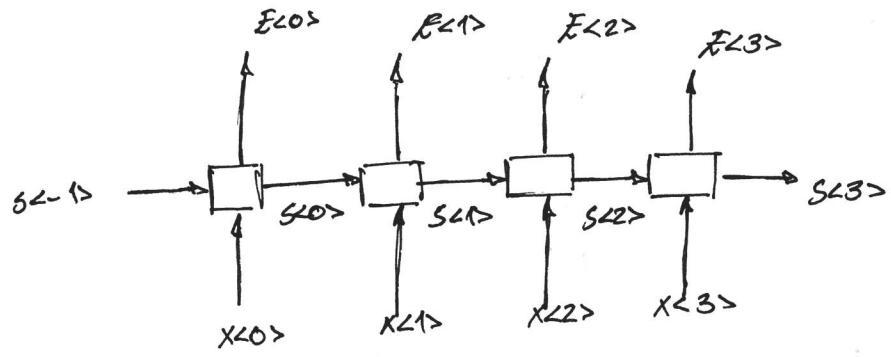
Decides which features to of memory we want to keep.

$$\tilde{f}_r^{t+1} = \sigma(W_r[C^{t+1-1}, x^{t+1}] + b_r)$$

$$\tilde{c}^{t+1} = \text{tanh}(W_r[\tilde{f}_r^{t+1} * C^{t+1-1}, x^{t+1}] + b_c)$$

$$\tilde{\beta}_u^{t+1} = \sigma(W_u[C^{t+1-1}, x^{t+1}] + b_u)$$

$$C^{t+1} = \tilde{\beta}_u^{t+1} * \tilde{c}^{t+1} + (1 - \tilde{\beta}_u^{t+1}) * C^{t+1-1}$$



$$\frac{\partial E}{\partial w_\mu} = \sum_{l=0}^{T_x} \sum_{k=l}^{T_y} \frac{\partial E_k}{\partial s^{(k)}} \prod_{j=l+1}^K \left(\frac{\partial s^{(j)}}{\partial s^{(j-1)}} \right) \frac{\partial s^{(l)}}{\partial w_\mu}$$

$$\frac{\partial E}{\partial x^{(l)}} = \sum_{k=l}^{T_y} \frac{\partial E^{(k)}}{\partial s^{(k)}} \prod_{j=l+1}^K \left(\frac{\partial s^{(j)}}{\partial s^{(j-1)}} \right) \frac{\partial s^{(l)}}{\partial x^{(l)}}$$

due to RNN (one cell architecture), but $\frac{\partial s^{(j)}}{\partial s^{(j-1)}}$ has a linear path to prevent vanishing gradient \rightarrow

$$(1 - r_{\mu, l+1}) \cancel{\text{linear path}} \rightarrow \cancel{\text{vanishing gradient}}$$

If enable or disable (0,1)

If enabled and $E \uparrow$, path to any update. If disable an $E \uparrow$, w_μ will change to minuscule \rightarrow carry gradient when enabled.

GRU Cell Backprop →

$$\hat{y}_{t+j} = \text{softmax}; z_{t+j} = w_y \alpha_{t+j} + b_y$$

$$\frac{\partial E_{t+j}}{\partial a_{ij|t+j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{\partial E_{t+j}}{\partial z_{rs|t+j}} \cdot \frac{\partial z_{rs|t+j}}{\partial a_{ij|t+j}} = \sum_{r=1}^{n_y} \sum_{s=1}^m \frac{1}{m} [\hat{y}_{rs|t+j} - y_{rs|t+j}] \cdot w_{yni} s_{s,j}$$

$$= \sum_{k=1}^{n_y} \frac{1}{m} [\hat{y}_{kj|t+j} - y_{kj|t+j}] \cdot w_{ykj}$$

Vectorized → $\frac{\partial E_{t+j}}{\partial \alpha_{t+j}} = \frac{1}{m} w_y^T [\hat{y}_{t+j} - y_{t+j}]$

$$\frac{\partial E_{t+j}}{\partial w_y} = \frac{1}{m} [\hat{y}_{t+j} - y_{t+j}] \cdot \alpha_{t+j}^T$$

$$\frac{\partial E}{\partial C_{t+j}} = \frac{\partial E_{t+j}}{\partial C_{t+j}} + \frac{\partial E_{t+1:T_j}}{\partial C_{t+j}}$$

$$\frac{\partial C_{k|t+j}}{\partial P_{\mu|i|j|t+j}} = \hat{C}_{k|t+j} s_{ki} s_{lj}; \frac{\partial C_{k|t+j}}{\partial \hat{C}_{i|j|t+j}} = P_{\mu|i|j|t+j} s_{ki} s_{lj}$$

$$\frac{\partial E}{\partial P_\mu} = \frac{\partial E}{\partial C_{t+j}} * \hat{C}_{t+j}; \frac{\partial E}{\partial \hat{C}} = \frac{\partial E}{\partial C_{t+j}} * P_\mu \leftarrow \text{Vectorized.}$$

$$\frac{\partial E}{\partial z_{ij|t+j}} = \sum_{r=1}^{n_c} \sum_{s=1}^m \left(\frac{\partial E}{\partial P_\mu} * P_{\mu|rs|t+j} \cdot (1 - P_{\mu|rs|t+j}) \right) \cdot w_{\mu|ni} s_{s,j}$$

$$\frac{\partial E}{\partial z_{t+j}} = w_\mu^T \left[\frac{\partial E}{\partial P_\mu} * P_\mu * (1 - P_\mu) \right] \leftarrow \text{Vectorized.}$$

$$\frac{\partial E}{\partial q_{t+j}} = w_c^T \left[\frac{\partial E}{\partial \hat{C}} * (1 - \hat{C}_{t+j}^2) \right]$$

$J^{<+>} = q^{<+>} [: n_c, :] \rightarrow$ Just elements corresponding to $C^{<+>-1}$

$$\frac{\partial E}{\partial J^{<+>}} = \frac{\partial E}{\partial q^{<+>}} [: n_c, :]$$

$$J^{<+>} = P_r^{<+>} * C^{<+>-1} \rightarrow \frac{\partial f^{<+>}_{ke}}{\partial P_{rij}^{<+>}} = C_{ke}^{<+>-1} \delta_{ki} \delta_{je}; \frac{\partial f^{<+>}_{ke}}{\partial G_j^{<+>1}} = V_{rme}^{<+>-1} \delta_{ki} \delta_{ej}$$

$$\frac{\partial E}{\partial p^{<+>}} = (\text{def}) \quad W_r^T \left[\frac{\partial E}{\partial J^{<+>}} * C_k^{<+>} * P_r^{<+>} * (1 - P_r^{<+>}) \right]$$

$$\frac{\partial E}{\partial x^{<+>}} = \frac{\partial E}{\partial p^{<+>}} [n_c :, :] + \frac{\partial E}{\partial z^{<+>}} [n_c :, :] + \frac{\partial E}{\partial q^{<+>}} [n_c :, :]$$

$$\frac{\partial E}{\partial C^{<+>-1}} = \frac{\partial E}{\partial C^{<+>}} * (1 - P_u^{<+>}) + \frac{\partial E}{\partial J^{<+>}} * P_r^{<+>} + \frac{\partial E}{\partial z^{<+>}} [: n_c, :]$$

←
Vektorschreibweise.

Sequence Models

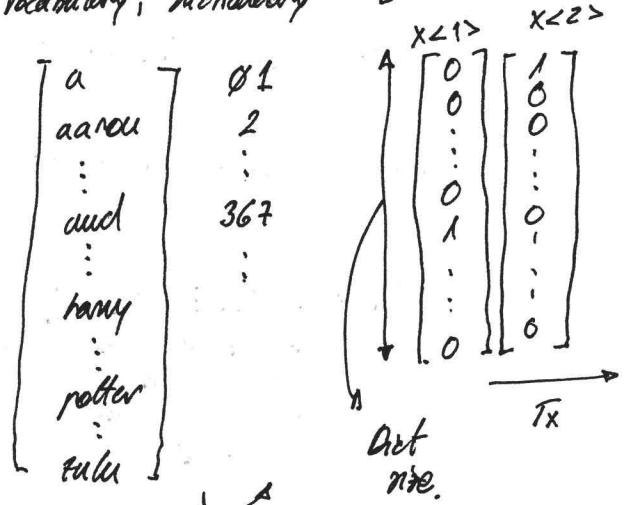
* Notations →

$x^{<t>} = \text{sample } t \text{ on the temporal space.}$

$T_x = \text{length of the sequence.}$

$x_{(i)}^{<t>} = \text{sample } t \text{ on the i-th sample set.}$

Vocabulary, dictionary →

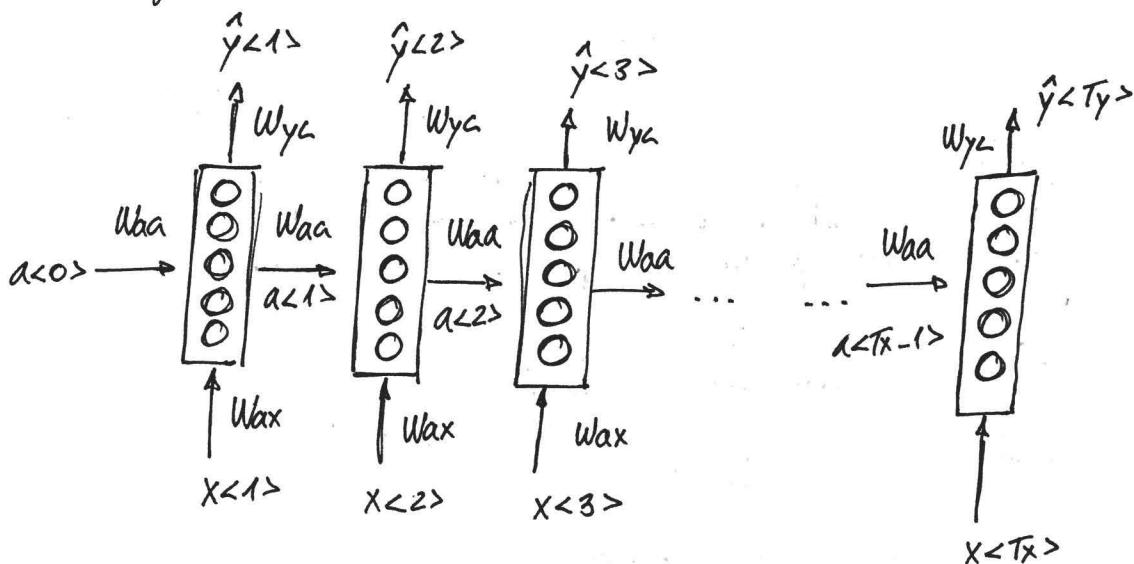


* RNN model →

Reason not to use DNN →

one hot representation

- Input or output can have different length excepted.
- They don't share features learned across different portions of the sequence.
- If you want to use the whole set of words layer become huge.
- huge.



$$\alpha^{<t>} = g_1(W_{aa} \alpha^{<t-1>} + W_{ax} x^{<t>} + b_a) \quad \leftarrow \text{Tanh / ReLU}$$

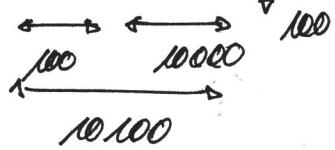
$$\hat{y}^{<t>} = g_2(W_{ya} \alpha^{<t>} + b_y) \quad \leftarrow \text{softmax} \rightarrow \text{usually, not necessarily.}$$

- Simplifying equations →

$$a_{t-1} = g_1(W_a a_{t-1} + W_x x_t + b_a)$$

$$a_t = g_1(W_a [a_{t-1}, x_t] + b_a)$$

$W_a = [W_{aa} : W_{ax}]$ → stacked horizontally.

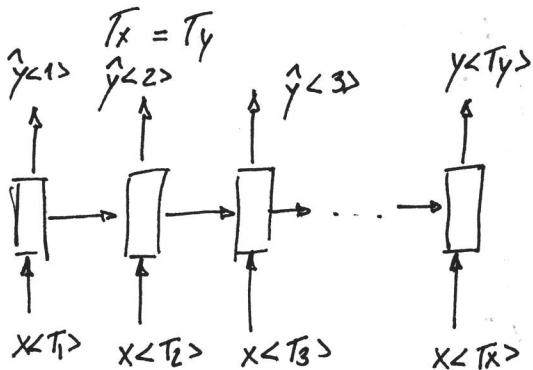


$$[a_{t-1}, x_t] = \begin{bmatrix} a_{t-1} \\ x_t \end{bmatrix} \begin{bmatrix} 100 & 100 \\ 10000 & 10,100 \end{bmatrix}$$

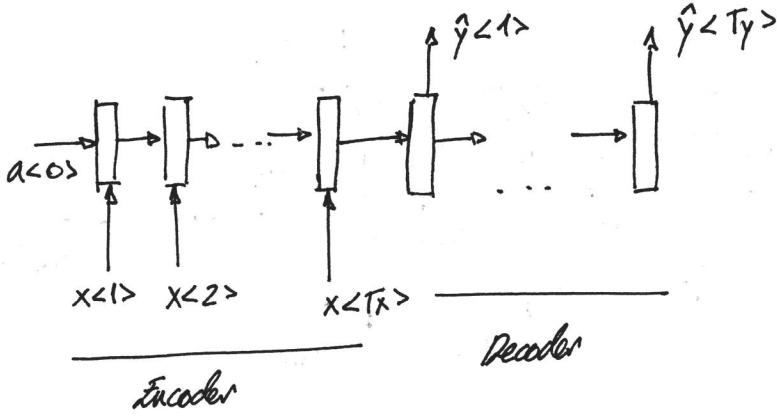
$$[W_{aa} : W_{ax}] \begin{bmatrix} a_{t-1} \\ x_t \end{bmatrix} = W_{aa} a_{t-1} + W_{ax} x_t$$

* Different RNNs →

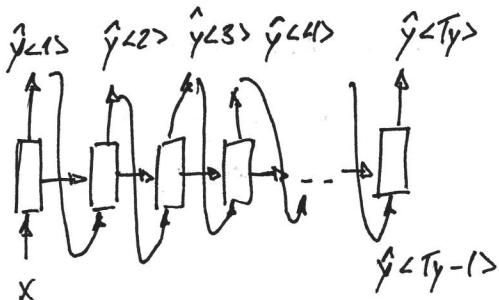
Many-to-Many



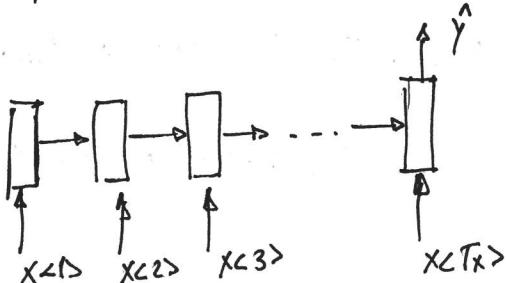
Many-to-Many



One-to-Many



Many-to-one



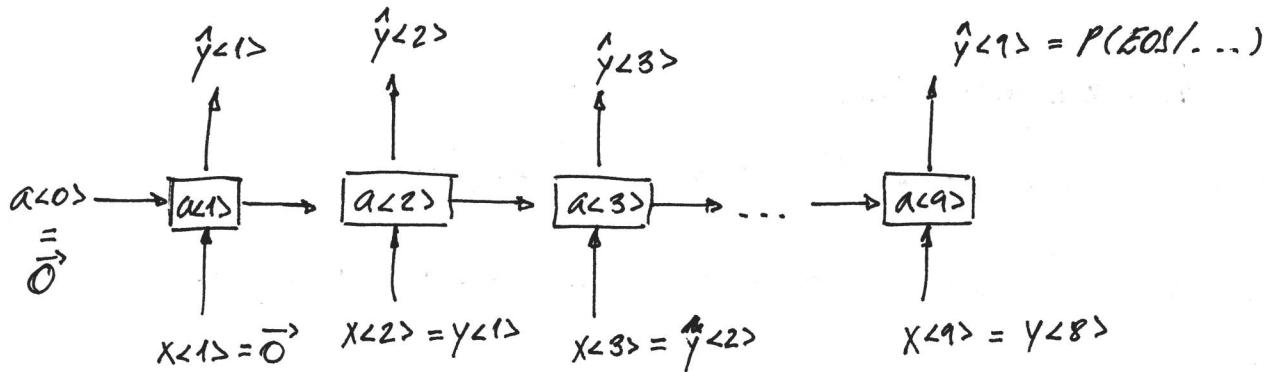
* language and sentence generation →

- Does probability assignment to word sequence → sentence.

$$P(y_{<1>}^{}, y_{<2>}^{}, y_{<3>}^{}, \dots, y_{<T_y>}^{})$$

"The Egyptian rat is a breed of a cat" → $\langle \text{EOS} \rangle$
Not in dict, LMKs.

- RNN model →



$\hat{y}_{<1>}^*$ = softmax prediction of a word, $x_{<1>}^*$. → dim of dictionary.

$$\hat{y}_{<2>}^* = P(\text{Egyptian} / y_{<1>}^* = \text{the})$$

- Cost function →

$$L_{<t>}(\hat{y}_{<t>}^*, y_{<t>}^*) = - \sum_i y_{i<t>}^* \log_i \hat{y}_{i<t>}^*$$

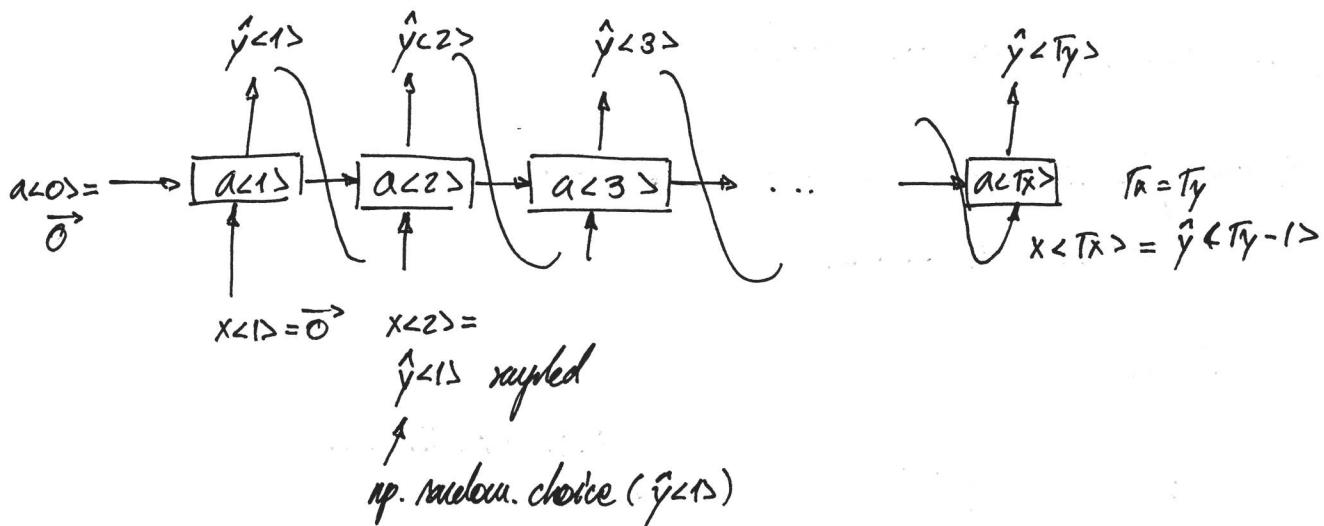
$$L(\hat{y}, y) = \sum_t L_{<t>}(\hat{y}_{<t>}^*, y_{<t>}^*)$$

- Given new sentence; you can determine the probability →

$$P(y_{<1>}^{}, y_{<2>}^{}, y_{<3>}^{}) = P(y_{<1>}^{}) P(y_{<2>}^{} / y_{<1>}^{}) P(y_{<3>}^{} / y_{<2>}^{} y_{<1>}^{})$$

* Sampling novel sequences \rightarrow

The models represent the chance of having sequence words.



- * Randomly sample across the softmax distribution $\hat{y}_{<1>}$

- * You can sample until the sentence has x words or when you hit EOS.

* Character level. \rightarrow

vocabulary = [a, b, c, ..., z, ., , , ;, 0, 1, ..., 9, A, ..., Z]

- Same idea as before $\rightarrow \hat{y}_{<1>} = \text{char} ; \text{E.g. } 'a'$.

- This setup is able to assign P=0 to unknown words such as Name.

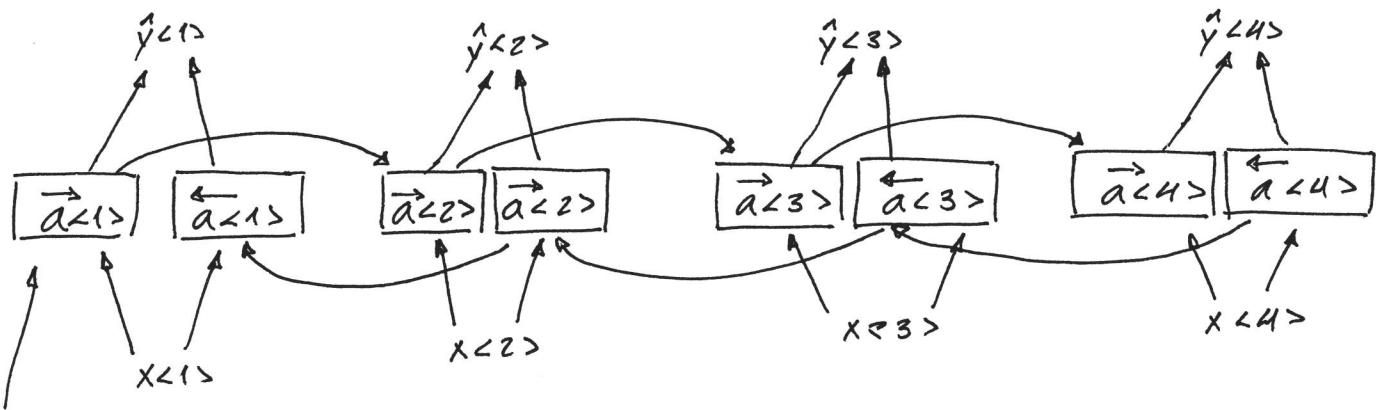
* Vanishing gradient \rightarrow

- Long sequences \rightarrow Problem with the gradients as in DNNs.

- Exploding gradients can also happen \rightarrow

Name; gradient clipping \rightarrow rescale gradient.

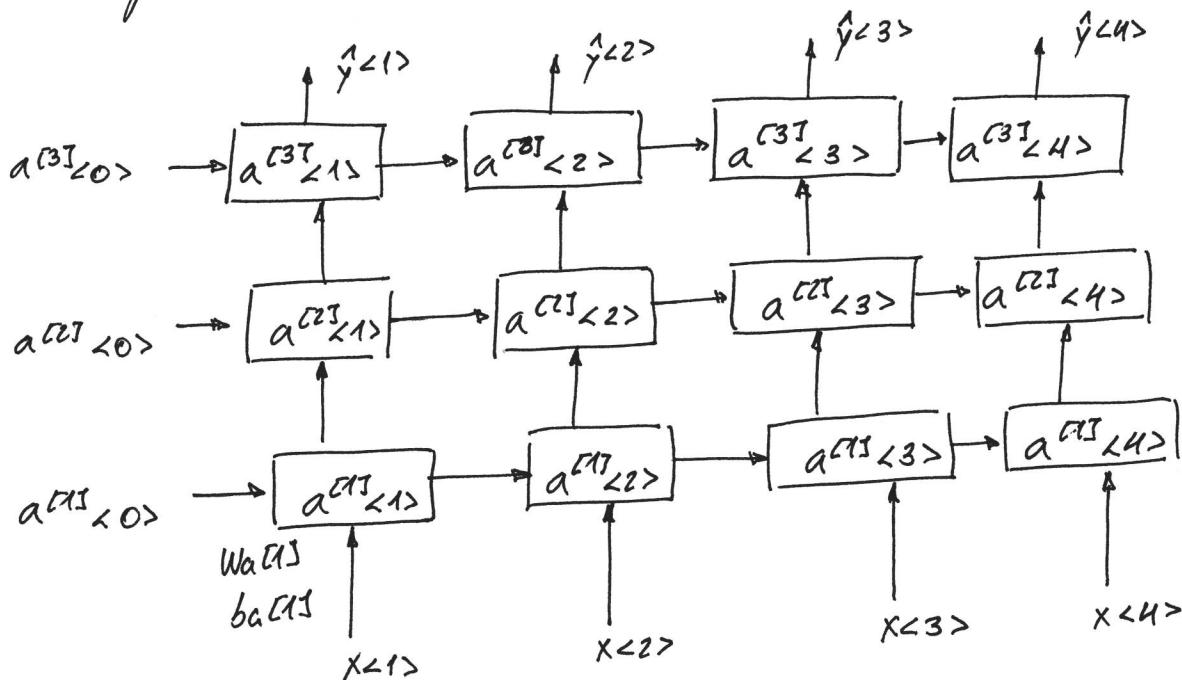
* Bidirectional RNNs →



RNN; GRU; LSTM.

$$\hat{y}^{<t>} = g(W_y [\vec{a}^{<t>}, \overleftarrow{a}^{<t>}] + b_y)$$

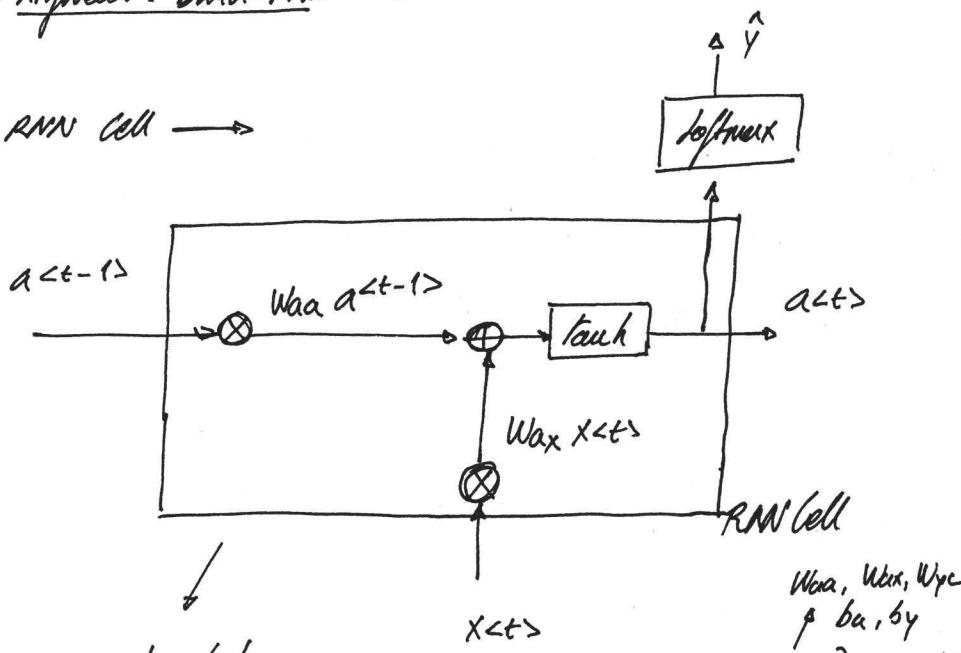
* Deep RNN →



$$a^{(2)}_{<3>} = g(W^{(2)} [a^{(2)}_{<2>}, a^{(2)}_{<3>}] + b^{(2)})$$

- Algorithm: Build RNN →

* RNN Cell →



$$a^t = \tanh(W_{aa} a^{t-1} + W_{ya} x^t + b_a)$$

$$\hat{y}^t = \text{softmax}(W_{ya} a^t + b_y)$$

$x^t \rightarrow n_x$ vector.

$a^t \rightarrow n_a$ vector

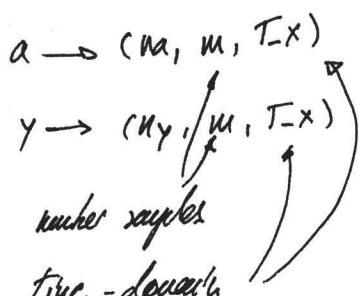
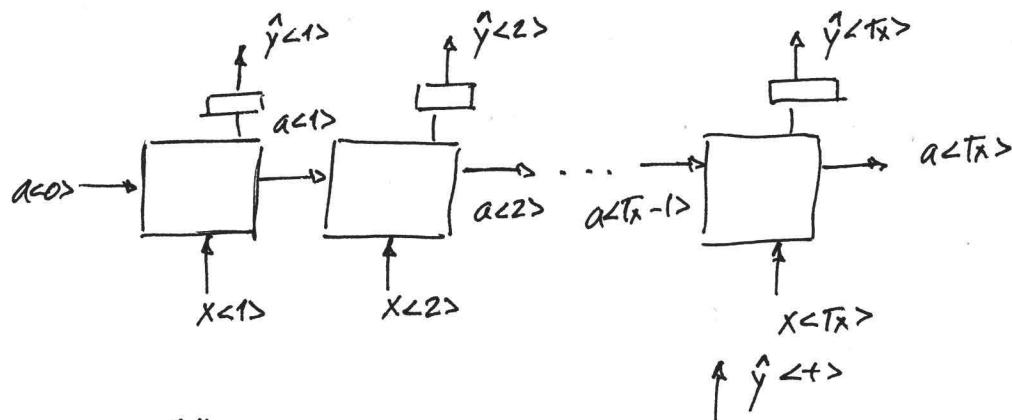
$\hat{y}^t \rightarrow n_y$ vector.

$$W_{ya} \rightarrow (n_a, n_x) \quad W_{ya} \rightarrow (n_y, n_a)$$

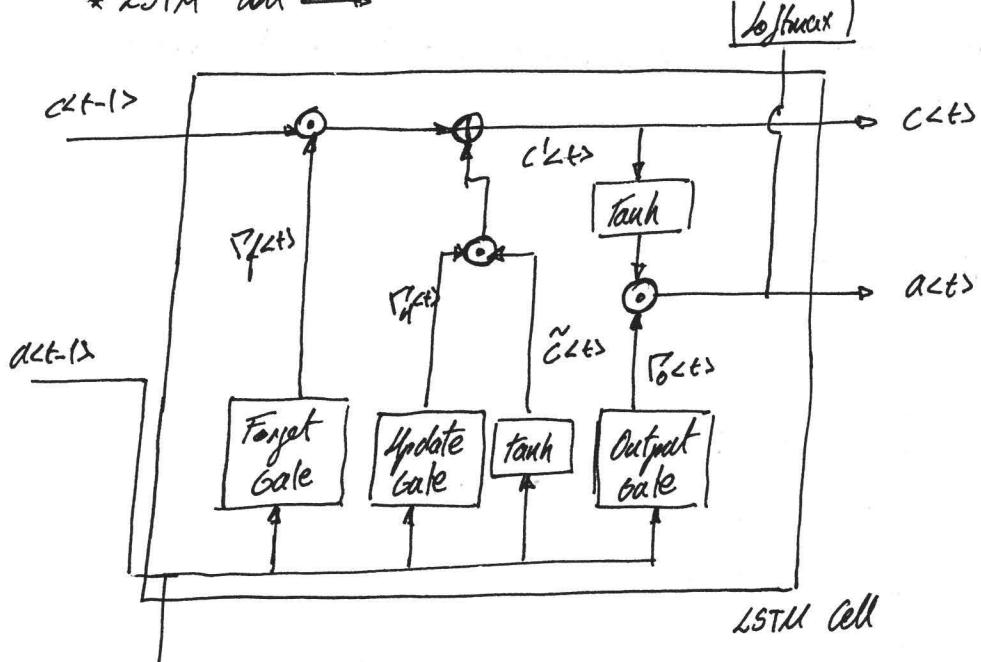
$$W_{aa}, W_{ya}, b_a, b_y$$

cache backprop → (a^t , a^{t-1} , x^t , parent)

* Forward pass →



* LSTM Cell →



$$f^t = \sigma(W_f [a^{t-1}, x^t] + b_f)$$

$$i^t = \sigma(W_i [a^{t-1}, x^t] + b_i)$$

$$\tilde{c}^t = \tanh(W_c [a^{t-1}, x^t] + b_c)$$

$$c^t = c^{t-1} * f^t + i^t * \tilde{c}^t$$

$$o^t = \sigma(W_o [a^{t-1}, x^t] + b_o)$$

$$a^t = o^t * \tanh(c^t)$$

x^t

- Cache → (a^t , a^{t-1} , c^t , c^{t-1} , f^t , i^t , \tilde{c}^t , O^t , x^t),
parameters.

* Concepts about gates →

- Forget gate →

* E.g.: piece of text, keep track of grammatical structures, such as subject is singular or plural. If it changes from singular to plural need a way to get rid of previous value.

- Update gate →

* Once we forget the subject, we need to update the value so it reflects the new plural.

- Updating the cell →

* To update the new subject, we need to create a new vector that we can add to our previous cell state.

- Output gate →

* To decide which of the outputs we will use.

* Dimension →

$$W \in (\text{Na}, \text{N}_x + \text{Na})$$

$$x \in (\text{N}_x, M)$$

$$\rightarrow [a_{\leq t-1}, a_{\leq t}] \in (\text{N}_x + \text{Na}, M)$$

$$c \in (\text{Na}, M)$$

$$a \in (\text{Na}, M)$$

$$y \in (\text{Ny}, M)$$

$$\tau \Rightarrow \tau(W[a_{\leq t-1}, x_{\leq t}] + b)$$

$$\left(\begin{array}{cccc} w_{11} & \dots & w_{1N_x+Na} \\ \vdots & & \vdots \\ w_{Na1} & \dots & w_{NaN_x+Na} \end{array} \right) \left(\begin{array}{c} a_1 \dots a_{1M} \\ \vdots \\ a_{Na1} \dots a_{NaM} \\ \hline x_{11} \dots x_{1M} \\ \vdots \\ x_{N_x1} \dots x_{N_xM} \end{array} \right) \rightarrow (\text{Na}, M)$$

* LSTM = long Short Term Memory \rightarrow

$$\tilde{C}^{<t>} = \tanh(W_c [a^{<t-1>}, x^{<t>}] + b_c)$$

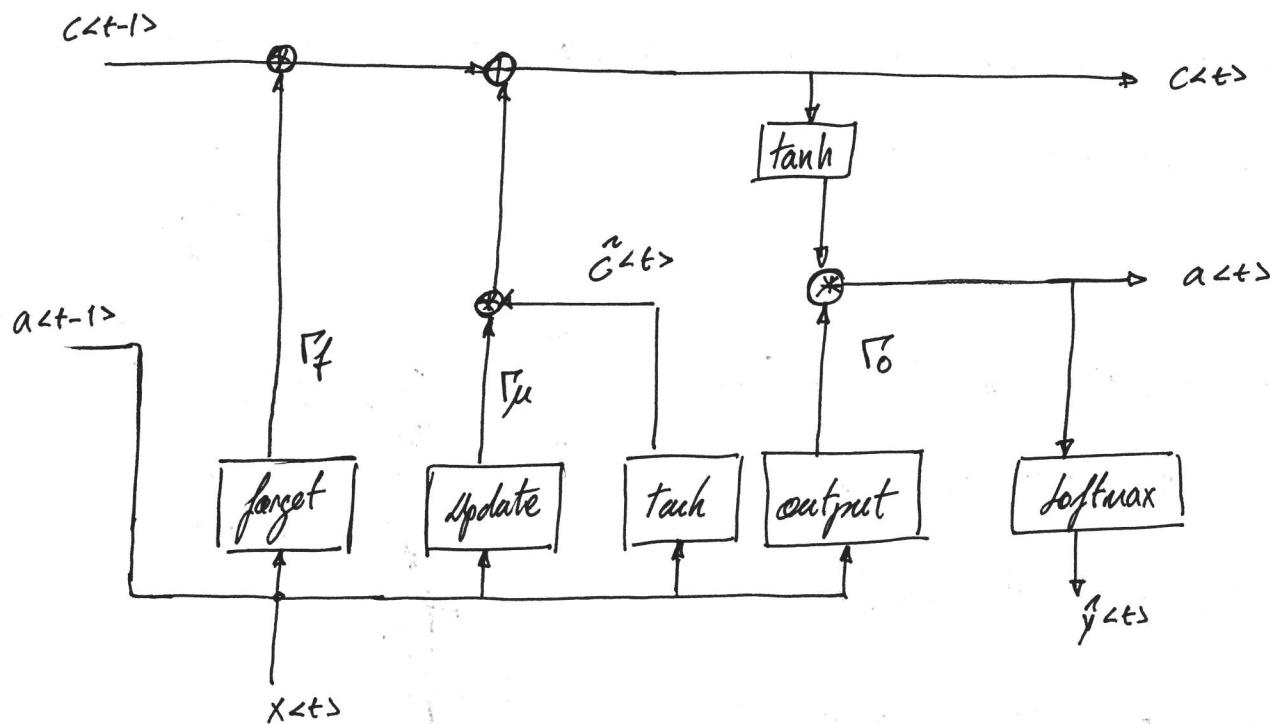
$$\Gamma_u = \sigma(W_u [a^{<t-1>}, x^{<t>}] + b_u) ; \text{ update.}$$

$$\Gamma_f = \sigma(W_f [a^{<t-1>}, x^{<t>}] + b_f) ; \text{ forget.}$$

$$\Gamma_o = \sigma(W_o [a^{<t-1>}, x^{<t>}] + b_o) ; \text{ output.}$$

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + \Gamma_f * C^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \text{tanh}(C^{<t>})$$



- Possible variant \rightarrow Peephole connection.

* Gate values depend on previous memory values.

$$\Gamma_o, \Gamma_f, \Gamma_u = \sigma(W [a^{<t-1>}, x^{<t>}] + b) \quad C^{<t-1>}$$

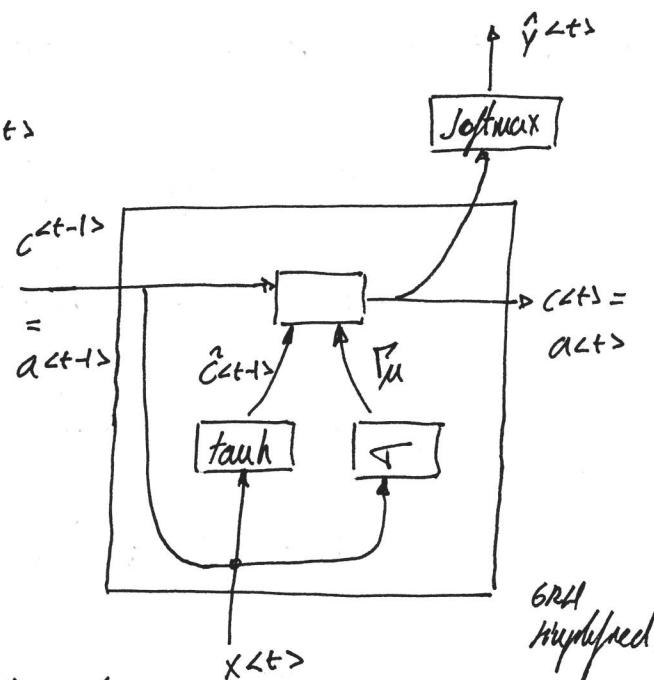
* GRU = Gated Recurrent Unit →

- Simplified → $C \equiv \text{Memory Cell} ; C^{<t>} = a^{<t>}$

$$\tilde{C}^{<t>} = \tanh(W_C [C^{<t-1>} , x^{<t>}] + b_C)$$

$$\Gamma_u = \sigma(W_u [C^{<t-1>} , x^{<t>}] + b_u)$$

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>}$$



GRU Simplified.

* $C^{<t>} , \tilde{C}^{<t>} , \Gamma_u \equiv \text{have same dimension.}$

* $\Gamma_u \equiv \text{Update gate, the sigmoid can get values close to 1, 0 for high, low values.}$

* Γ_u will control if $C^{<t>}$ will update its value depending on new value $x^{<t>}$ or previous $C^{<t-1>}$.

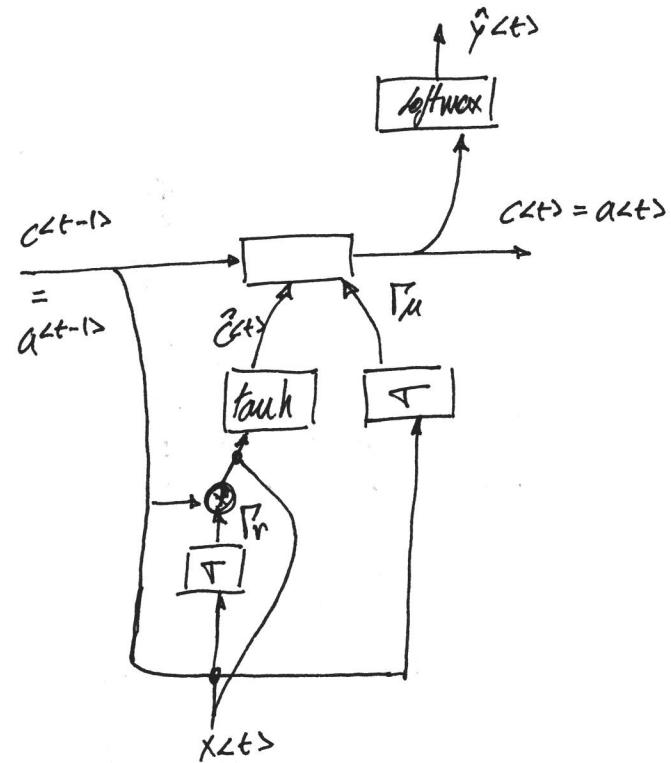
- Full GRU →

$$\tilde{C}^{<t>} = \tanh(W_C [\Gamma_r * C^{<t-1>} , x^{<t>}] + b_C)$$

$$\Gamma_u = \sigma(W_u [C^{<t-1>} , x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r [C^{<t-1>} , x^{<t>}] + b_r)$$

$$C^{<t>} = \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>}$$



* Introduced relevance gate. $\equiv \Gamma_r$

NLP & Word Embedding →

* Word representation →

- One-hot representation for word has limitations, tell nothing about meaning or relationship between them.

Features	Man 5391	Woman 9853	King 4914	Queen 7157	Apple 456	Orange 6257
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.01
Food	0.04	0.01	0.02	0.01	0.95	0.97

↓
300
features.
↓
C⁵³⁹¹

- Visualisation →

key:
 Man . . Woman fish dog
 King . . Queen cat
 tree . Fair Apple Orange
 one two Grape

t-SNE

3000 → 2D

non-linear

- Only word embedding makes the RNN infer meaning based on features →

If RNN was trained with "Sally Johnson is an orange farmer".
 Since "Robert Lin is a dumpling cultivator", dumpling features will
 tell that it is a fruit and cultivator similar to farmer.

* Transfer learning →

1. Learn word embeddings from large dataset (1-100B words) (or download pre-trained model).
2. Train embeddings to a new task with smaller training set → 100K words.
3. Continue to finetune the word embeddings with new data. (Optional).

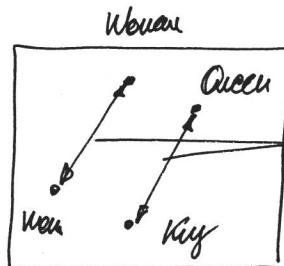
* Properties →

- Analogy →

$$e^{\text{man}} - e^{\text{woman}} = \begin{bmatrix} -2 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$e^{\text{key}} - e^{\text{queen}} = \begin{bmatrix} -2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

new
difference
is zero.



vector that represents difference
in gender

t-SNE might not show this difference, since it
is a non-linear mapping

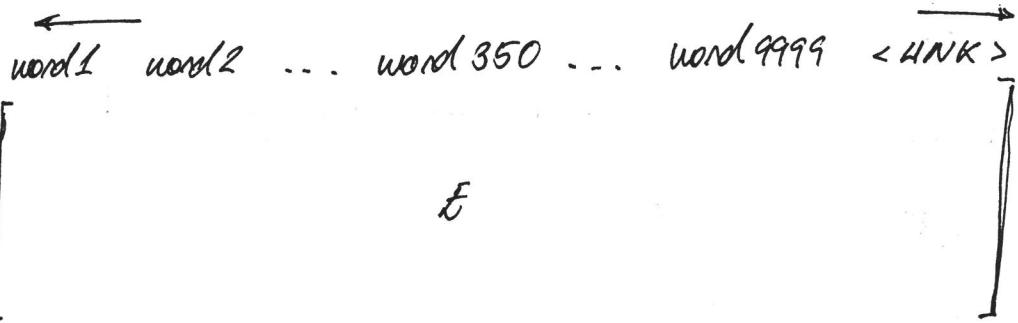
* Similarity function →

$$\text{sim}(e_w, e_{w'} - e_{w''} + e_{w'''})$$

$$\text{sim}(U, V) = \frac{U^T V}{\|U\|_2 \|V\|_2} ; \|U\|_2 = \sqrt{\sum_{i=1}^n \mu_i^2}$$

- You could also use the square distance, it will work out
but come works better.

* Embedding matrix \rightarrow



orange vector

e_{6257}

$$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \text{6257}$$

$$E \cdot e_{6257} = e_{6257}$$

$$(300, 10000) (10000, 1) = (300, 1)$$

$$E \cdot e_j = e_j$$

Embedding for word j

- Objective \rightarrow learn vectors e_j

- Initialize randomly, do gradient descent on it.

- Keras \rightarrow embedding layer (Multiplication is different).

* learning word embeddings \rightarrow

"I want a glass of orange" \rightarrow Target.

04343 09665 01 03852 06163 06257

- What context?

E E E E E E
↓ ↓ ↓ ↓ ↓ ↓
e4343 e9665 e1 e3852 e6163 e6257

- Necessity to find a context-target.

0 0 0 0 0 0 1 \leftarrow word, best

softmax \leftarrow word, best
prediction target.

* Word2Vec \rightarrow

- Skip-gram \rightarrow

"I want a glass of orange juice to go along with my cereal."

<u>Context</u>	<u>Target</u>
orange	juice
orange	glass
orange	any
orange	4

Pick randomly Based on target, choose randomly.

- Model \rightarrow

$$\text{Vocabulary size} = 10.000$$

orange \rightarrow target

$$o_c \rightarrow f \rightarrow e_c \Rightarrow \text{softmax} \rightarrow \hat{y}$$

$$c = E \cdot o_c$$

This is really
expensive.

For each one the words of
the dictionary.

$$\rightarrow p(\text{target} | \text{context}) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

$$L(q, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$$

$$y_i = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{one hot vector.}$$

* Negative raybox →

<u>Context</u>	<u>Word</u>	<u>Target</u>	
Orange	juice	1	
Orange	Kiwi	0	
Orange	hook	0	→ K random samples from dictionary.
Orange	the	0	small data sets → K = 5/20
Orange	of	0	large data sets → K = 275
<hr/>			
Input.		Output	

- Now it becomes a logistic regression problem →

$$P(y=1 / \text{context}, \text{target}) = \sigma(\theta_c^T \mathbf{x}_c)$$

- It is now a 10 000 binary logistic regression classifier.

- Instead we only train K+1

- Instead of a 10K softmax → 10K binary classification for K+1.

- Selecting negative examples →

* sample according to the freq words, but you want to take the ones most freq or do it uniformly → synonyme.

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10k} f(w_j)^{3/4}}$$

* GloVe = Global Word Vector →

"I want a glass of ..."

$x_{ij} \equiv \# \text{ times that word } i \text{ appears in context of word } j$ "

$j \equiv \text{context word}; i \equiv \text{target word}$.

$$\text{minimize} \left(\sum_{i=1}^{10000} \sum_{j=1}^{10000} [\theta_i^T e_j + b_i + b_j - \log x_{ij}]^2 \right)$$

↑
Relationship between 2 words

$$f(x_{ij}) \Rightarrow =0 \text{ if } x_{ij}=0 \text{ (An eval log 0)}$$

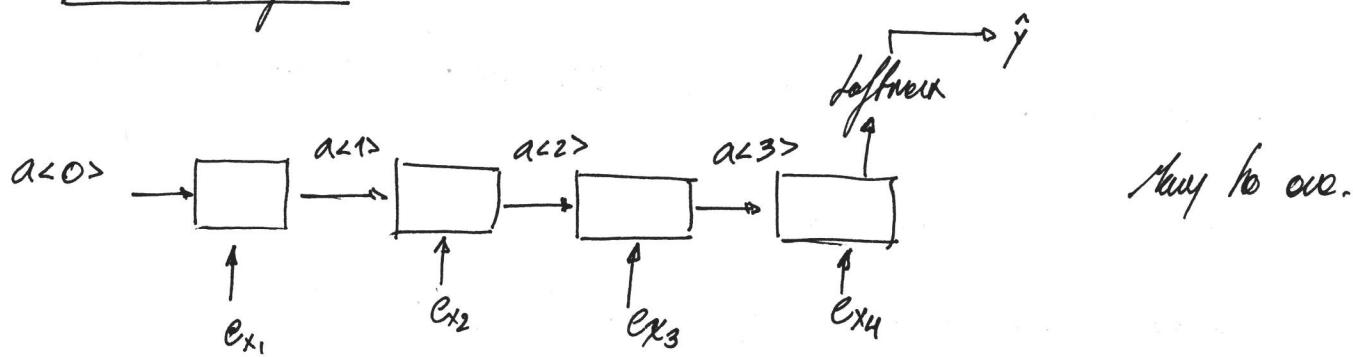
It can also be handle to treat freq words such as

the, of, a, ...

- θ_i, e_j are symmetric; initialize random uniform distribution.

- θ_i, e_j gradient descent.

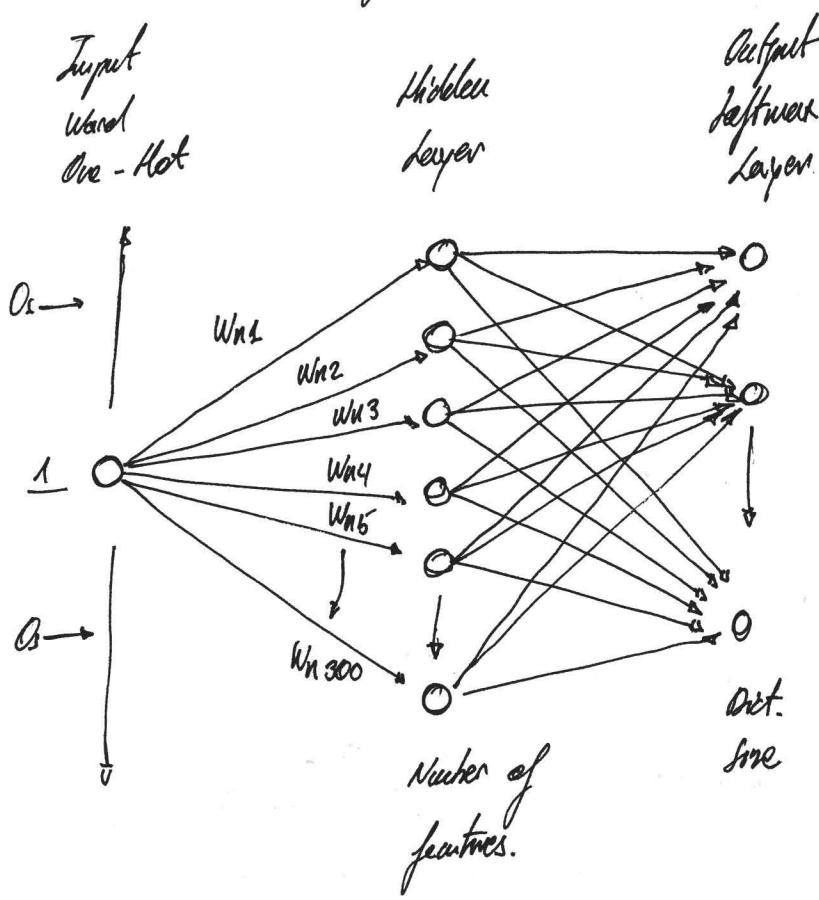
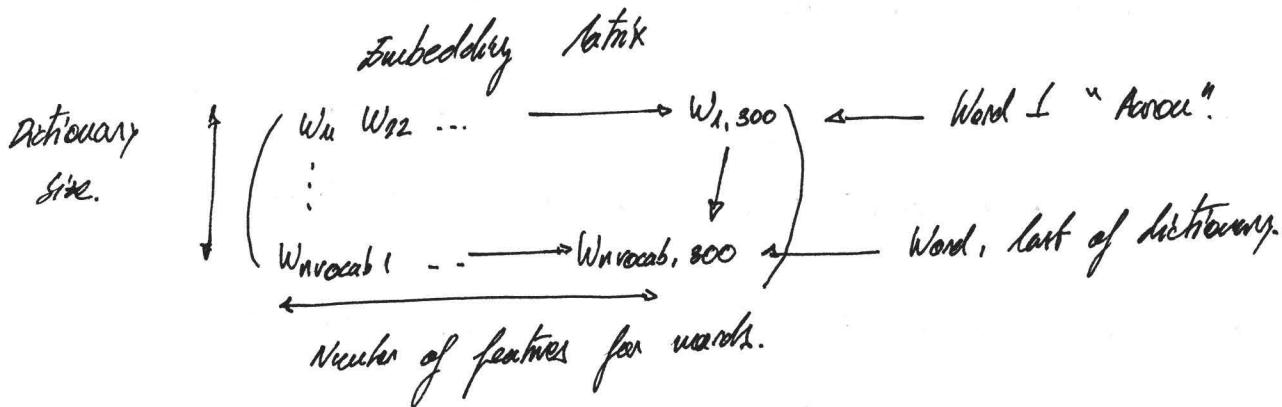
- softmax →



Word Embedding

* skip-gram representation →

- like a MLP with a softmax classifier
- build semantic relationships between words based on context.



$$P(\text{Target} | \text{Context}) = \hat{y}_i =$$

$$= \frac{e^{\theta_i^T \cdot c_c}}{\sum_{j=1}^{10K} e^{\theta_j^T \cdot c_c}}$$

$$\text{Loss} = - \sum_{i=1}^{10K} y_i \log \hat{y}_i$$

* Per input only one row of the embedding matrix will be enabled →
No multiplication, just indexing.

* calculating softmax and loss is really expensive. needs improvement.

ignorants are top →

* Word pairs & phrase →

- "Bottom Globe" as a word, not "Bottom" and "Globe".
- Count number of times two words appear together in the training set →
 - Equation decides to make phrase out of those pairs based on the phrase count and the individual occurrences.
 - Favors phrases from infrequent words, to avoid always like "This" or "and the".

* Infrequently frequent words →

- "the" example →

1. Pairs like "fox", "the" doesn't tell much about the meaning of fox.
2. Too many examples of "the" to learn.

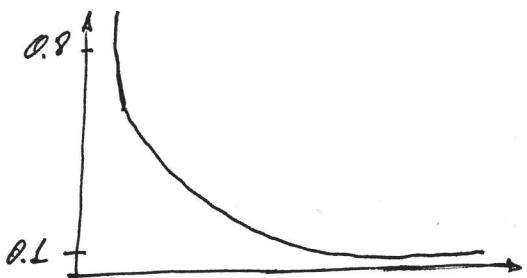
- Supply rate →

* Probability to keep a word in the text. $\equiv P(w_i)$

* $f(w_i) \equiv$ frequency of a word in the text. $= \frac{\text{counts}}{\text{total words text}}$

* Parameter $b \equiv$ parameter, who used to maintain probability between 0,1.

$$P(w_i) = \left(\sqrt{\frac{f(w_i)}{t}} + 1 \right) \frac{t}{f(w_i)} ; \quad 1 - P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$



$$P(w_i) = 1 \longrightarrow z(w_i) \leq 0.0026$$

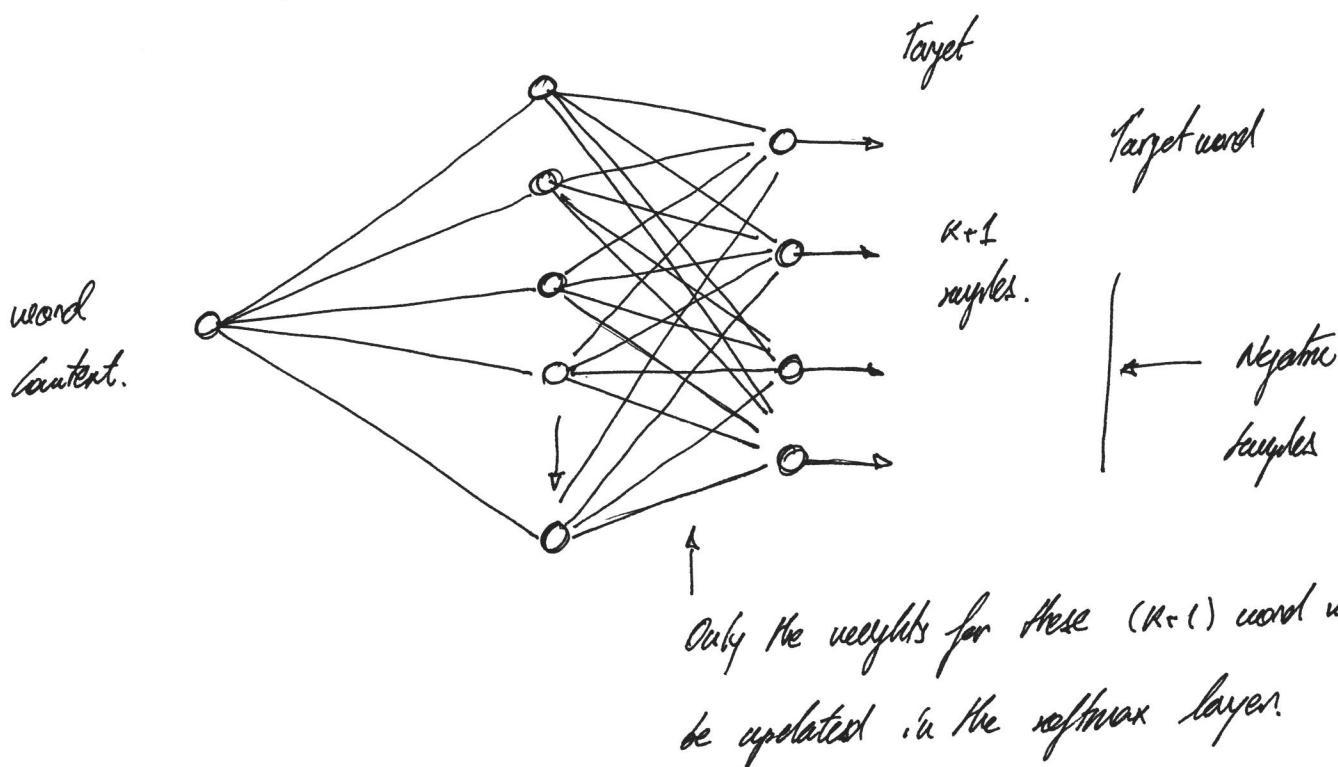
Only words that will be 0.26% of the text will be underlined.

$$P(w_i) = 0.5 \longrightarrow z(w_i) = 0.00746$$

- * Part → lower down the frequency difference for some really common words, like "the", "and", "it".
- * No single word should be a large part of the corpus.

- Negative sampling →

- * softmax step is really expensive → $\sum_{i=1}^{10K} e^{-\theta w_i}$; few billions of words.
- * One approach →



* Selecting negative samples \rightarrow

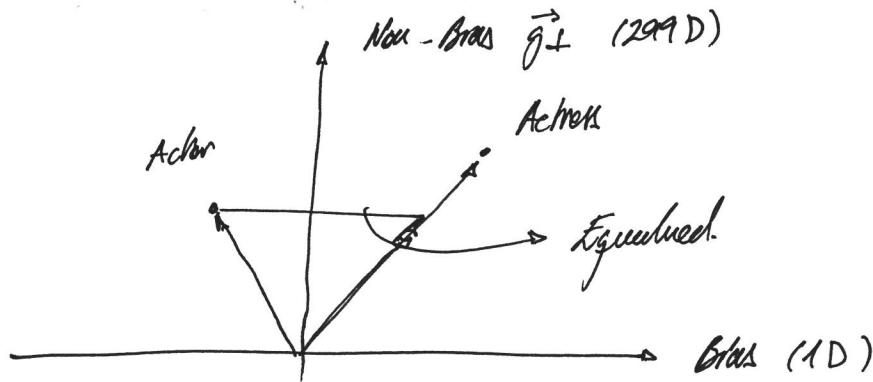
- Choose any an "inverse distribution"
- Probability of selecting a negative word is proportional to its frequency

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n (f(w_j)^{3/4})} ; \quad f(w_i) = \text{word frequency}$$

$3/4 \equiv \text{Empirical factor.}$

* Debiasing word embeddings →

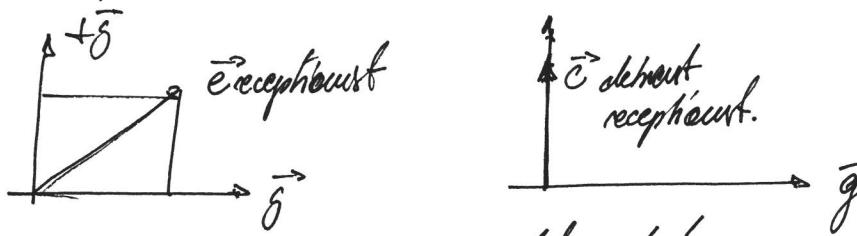
- Word embeddings can reflect gender, ethnicity, age, sexual orientation and other biases of the text used for training.



1. Identify bias vector:

2. Neutralize bias →

For every word that is not deontic



$$e_{bias-component} = \frac{e \cdot g}{\|g\|^2} \circ g$$

dot product.

$$e_{deontic} = e - e_{bias-component}.$$

3. Equalize pairs → vector word only difference is bias

$$\mu = \frac{\mu_W + \mu_B}{2}; \quad \mu_B = \frac{\mu - bias-axis}{\|bias-axis\|^2} * bias-axis; \quad \mu_{\perp} = \mu - \mu_B$$

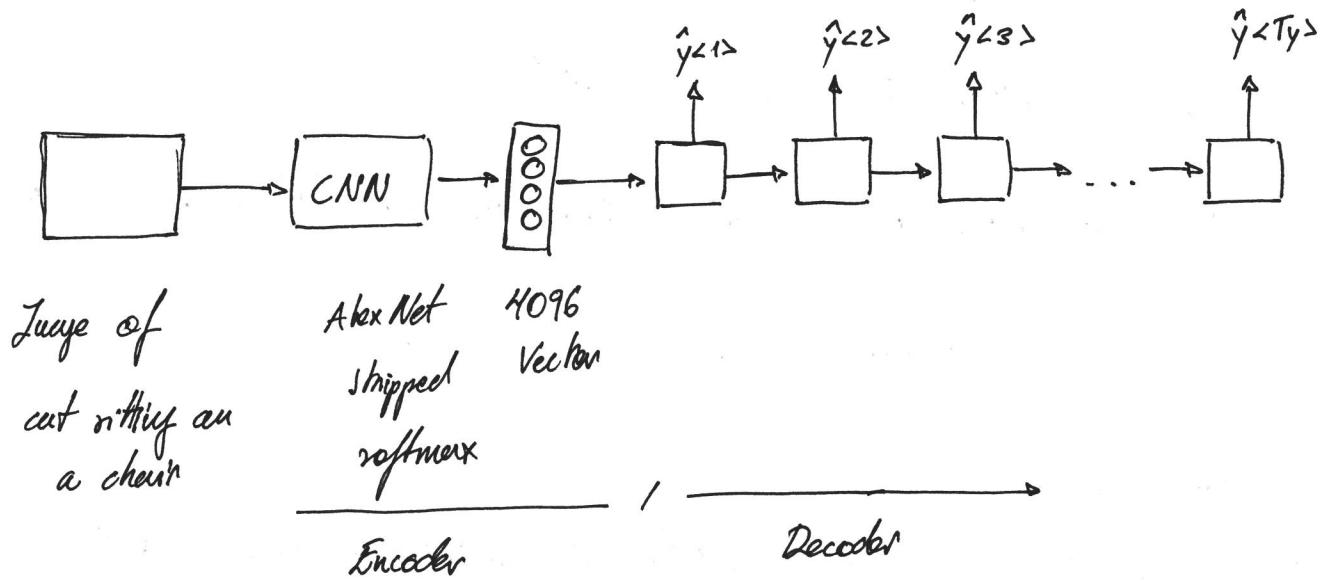
$$\mu_{WIB} = \frac{\mu_W + bias-axis}{\|bias-axis\|^2} * bias-axis; \quad \mu_{WB} = \frac{\mu_B + bias-axis}{\|bias-axis\|^2} * bias-axis$$

$$e_{W1B}^{\text{corrected}} = \sqrt{|1 - |\mu_4 u_2^2|} * \frac{e_{W1B} - \mu_B}{|e_{W1} - \mu_4 - \mu_B|}; e_{W2B}^{\text{corrected}} = \sqrt{|1 - |\mu_4 u_2^2|} * \frac{e_{W2B} - \mu_B}{|e_{W2} - \mu_4 - \mu_B|}$$

$$e_1 = e_{W1B}^{\text{corrected}} + \mu_4; e_2 = e_{W2B}^{\text{corrected}} + \mu_4$$

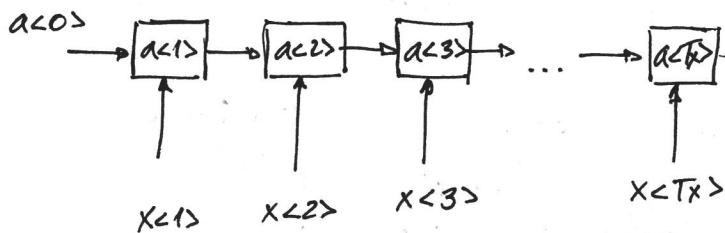
Sequence to sequence \rightarrow

- Large captioning \rightarrow



- Machine translation \rightarrow

Jane is writing Africa in depthener.



Jane unité l'Afrique en deptenre.

Decoder = language model.

Encoder

- The language model allows to estimate the probability of a sentence

$$P(\hat{y}_{<1>} \dots \hat{y}_{<Ty>} | x_{<1>} \dots x_{<Tx>})$$

Probability of the english sentence given the french.

- * We don't sample in the language model, randomly.

- * We don't either do greedy search because →

$$P(\hat{y}^{<1>} \dots, \hat{y}^{<T_y>} | x) \neq P(y^{<1>}) \cdot P(y^{<2>}) \dots P(y^{<T_y>})$$

"Jane is visiting Africa in September"

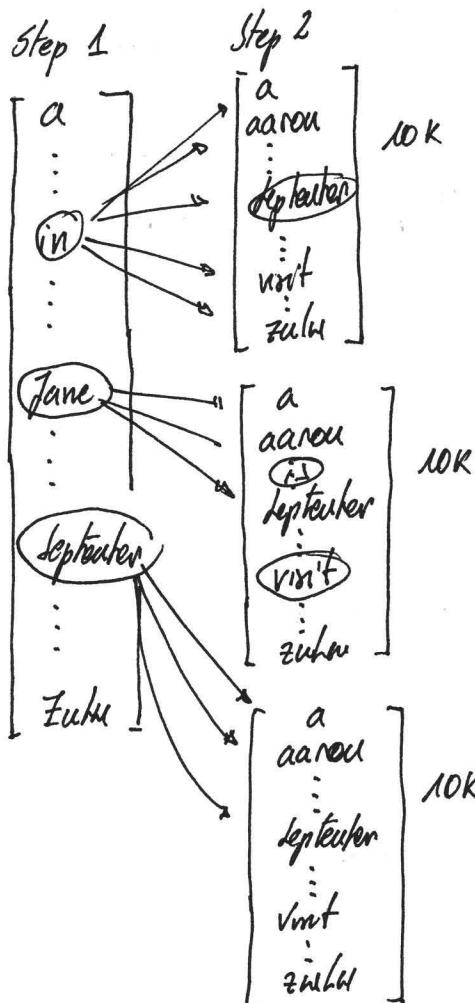
"Jane is going to be visiting Africa in September" → Worse solution.

$$P(\text{Jane is going } | x) > P(\text{Jane is visiting } | x)$$

- Beam search →

- * B = Beam width parameter; algorithm will consider 3 seed at a time

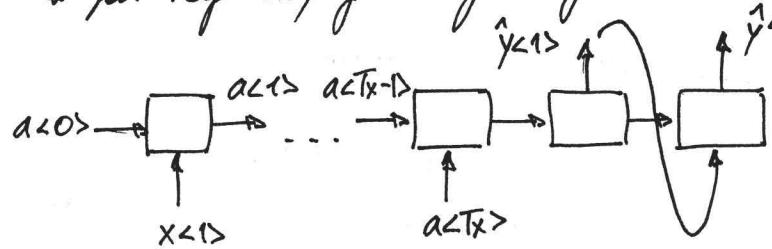
if $B=3$. Objective → $\arg \max_{y^{<1>} \dots, y^{<T_y>}} P(y^{<1>} \dots, y^{<T_y>} | x)$



- * For all $30K$ combinations, pick top B . $B=3$

$$P(y^{<1>} \dots, y^{<T_y>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

- * You keep carrying B copies of the network →



$$\hat{y}^{<1>} = \text{in}; \hat{y}^{<2>} = \text{September}$$

$$\hat{y}^{<1>} = \text{Jane}; \hat{y}^{<2>} = \text{in}$$

$$\hat{y}^{<1>} = \text{September}; \hat{y}^{<2>} = \text{visit}$$

Jane

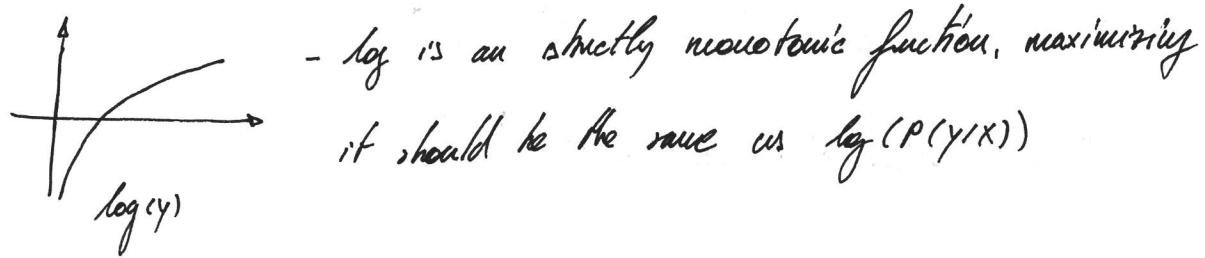
- * For each step you carry B copied.

* length normalization \rightarrow

$$\arg \max_y \prod_{t=1}^{T_y} P(y_{\leq t} | x, y_{\leq 1}, y_{\leq 2}, \dots, y_{\leq t-1})$$

↓ products of 0.x, resulting in really small numbers and round up at that low value can result in accuracy loss.

$$\arg \max_y \sum_{t=1}^{T_y} \log P(y_{\leq t} | x, y_{\leq 1}, y_{\leq 2}, \dots, y_{\leq t-1})$$



With longer sentences the Σ tends to be more negative, so it will tend to make shorter sentences.

$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y_{\leq t} | x, y_{\leq 1}, \dots, y_{\leq t-1}) = \text{Score.}$$

↑ Normalized by number of words in sentence, averaging them.

κ = Normalization factor, hyperparameter.

- check different $T_y = 1, 2, 3, \dots, 30. \rightarrow$

$\delta = 3 \rightarrow$ keep 3 possibilities per T_y

- look at score for all output and keep highest

* Trade-off between B values →

- Hyper → Better, slower
- Smaller → Worse, faster.

* Error Analysis →

- Erratic algorithm, doesn't guarantee best option.

~~Human reference~~ "Jave visita l'Afrique au septembre"

Human reference "Jave visits Africa in September" y^*

Algorithm "Jave visited Africa in September" \hat{y}

- RNN computes $P(y|x)$ (Encoder-decoder pair)

- Compute $P(y^*|x)$ and $P(\hat{y}|x) \rightarrow$

* $P(y^*|x) > P(\hat{y}|x) \rightarrow$ Beam search, not performing well.

* $P(y^*|x) < P(\hat{y}|x) \rightarrow y^*$ is a better translation; RNN.

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	Label
Sentence 1	T1	2×10^{-11}	1×10^{-11}	B
2	T2	↓	↓	R
3	T3	↓	↓	B
4	T4	↓	↓	B

% B, % R

- Blue score →

* Evaluates machine translation performance.

French "le chat est sur le mat"

Reference 1 "the cat is on the mat"

Reference 2. "there is a cat on the mat"

MT output: "The cat the cat on the mat"

	Count	Clipped Count	Bi-gram
1. 3 the cat	2	1	→ # of times it appears on reference 1,2.
2 cat the	1	0	
4 cat on	1	1	
5 on the	1	1	
6 the mat	1	1	

$$\text{Pn} = \frac{\sum_{\text{bigram}} \text{CountClip(n-gram)}}{\sum_{\text{bigram}} \text{Count(n-gram)}}$$

$$Pn = \frac{\sum_{n\text{-grams}} \text{CountClip}(n\text{-gram})}{\sum_{n\text{-grams}} \text{Count}(n\text{-gram})}$$

$$P = 4/6$$

; $P_1, P_2, \dots, P_n = 1.0 \rightarrow$ When MT is equal to reference.

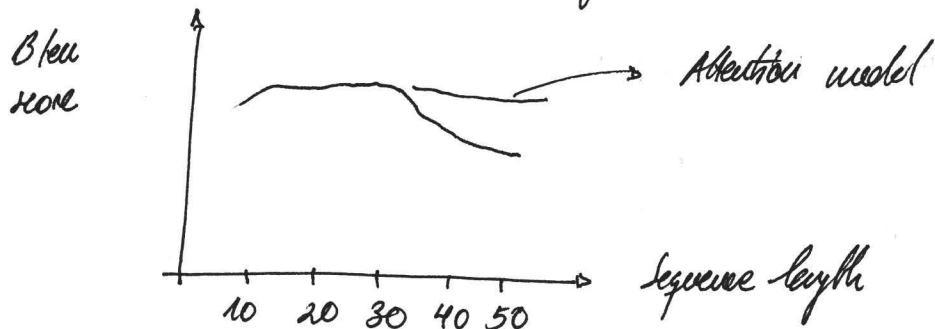
Corrected Blue Score → BP · exp $(\frac{1}{n} \sum_{n=1}^N P_n)$

BP = Brevity Penalty, easier to get better translations on short versions.

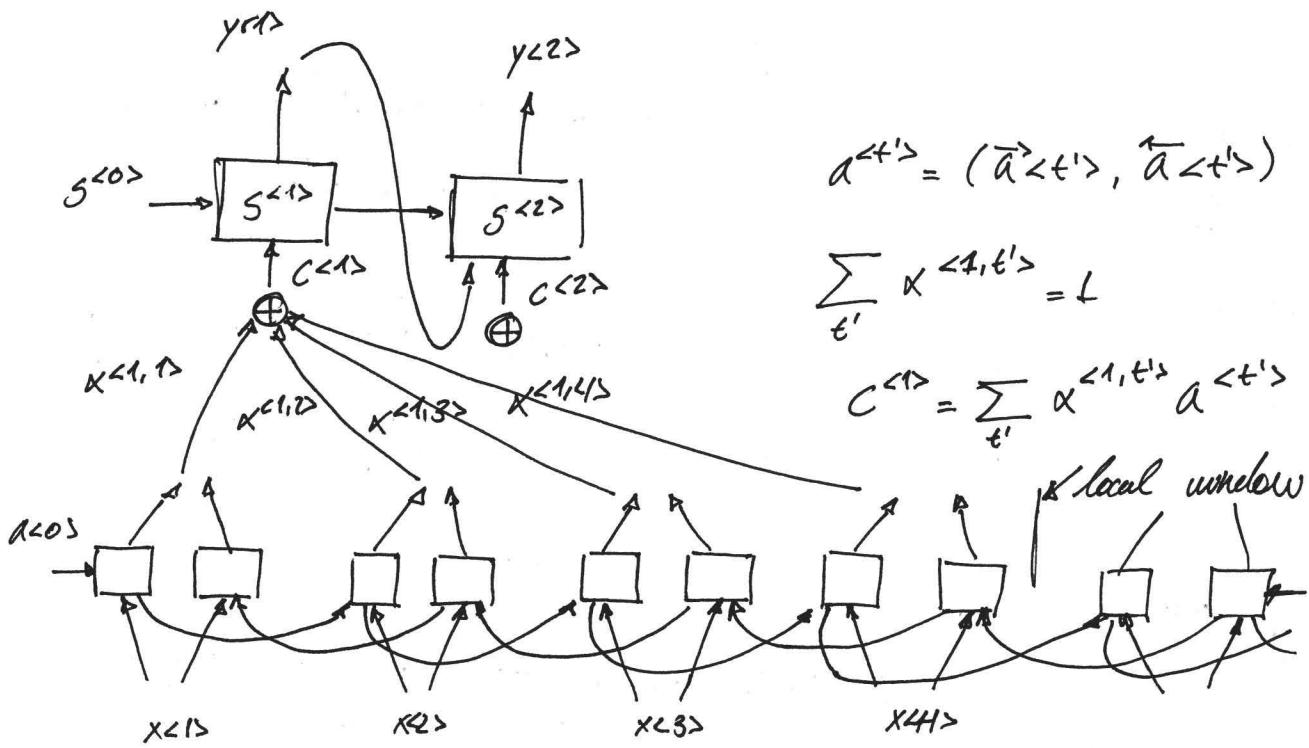
$$BP = \begin{cases} \exp(1 - \text{MT-length}/\text{Ref-length}) & \text{Otherwise} \\ 0 & \text{MT-output-length} > \text{reference-length.} \end{cases}$$

- Attention model →

* Problem with long sequence of sentences.



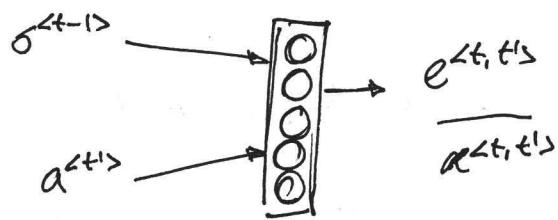
* As a human wouldn't try to translate after memorizing 10-20 sentences, etc., piece by piece.



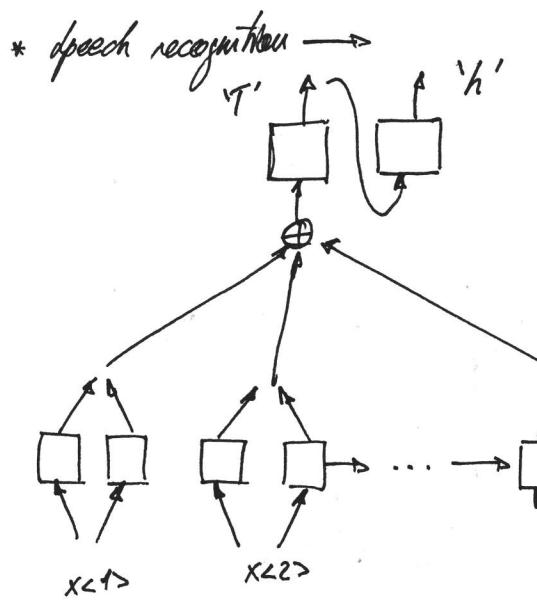
$\alpha^{<t, t'>} = \text{Amount of attention that } y^{<t>} \text{ should pay to } a^{<t'>}$

$$\alpha^{<t, t'>} = \frac{\exp(e^{<t, t'>})}{\sum_{t'=1}^T \exp(e^{<t, t'>})}$$

definition



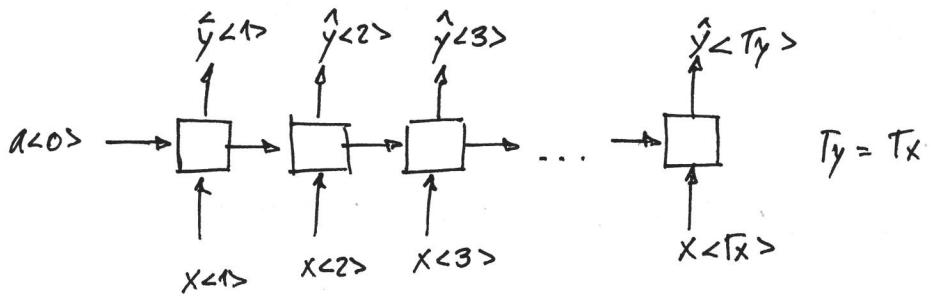
function to be learned by the network.



* Attention model

* like T_x time frames, e.g.: $T_x = 1000$

* CTC cost for speech recognition →



- Representation with regular RNN → Needs to be a bidirectional LSTM/GRU

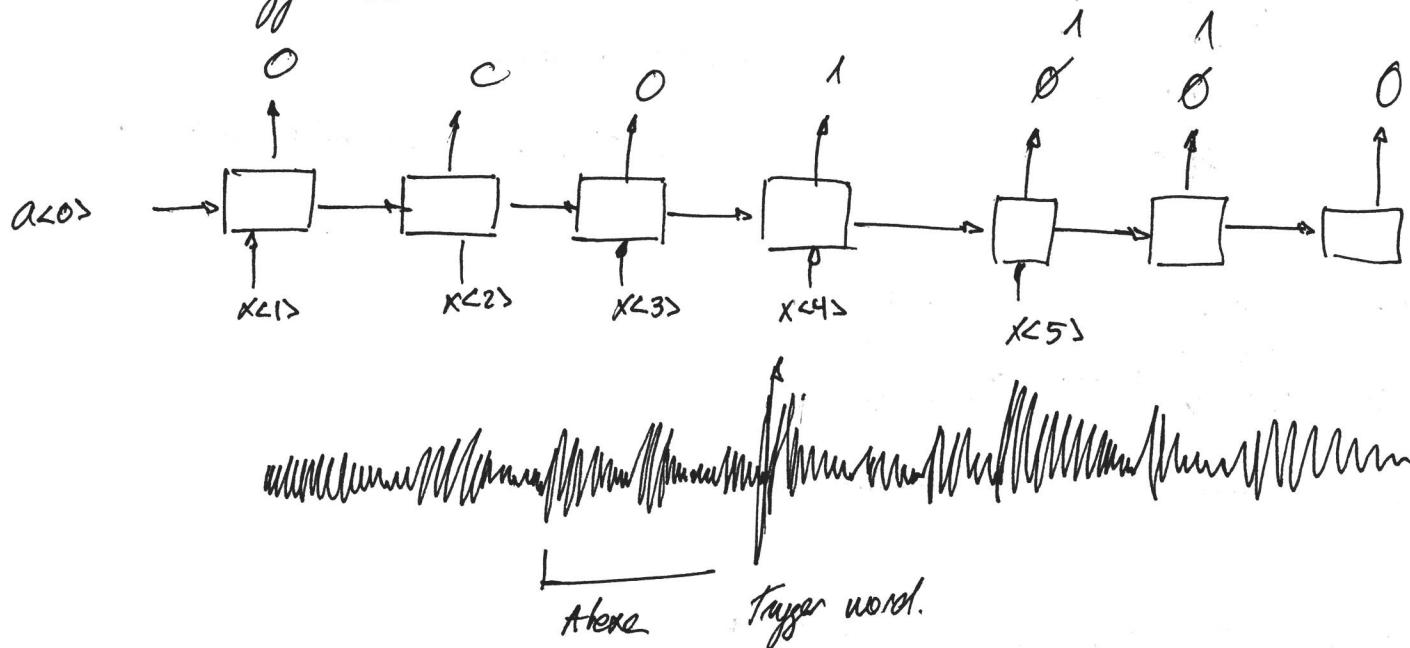
- For 10sec, at $T_x = 10000$, more input than choices →

ttt-h-eee---lll---qq-qllll-I--ccccccc-KK

"The quick" special character blank

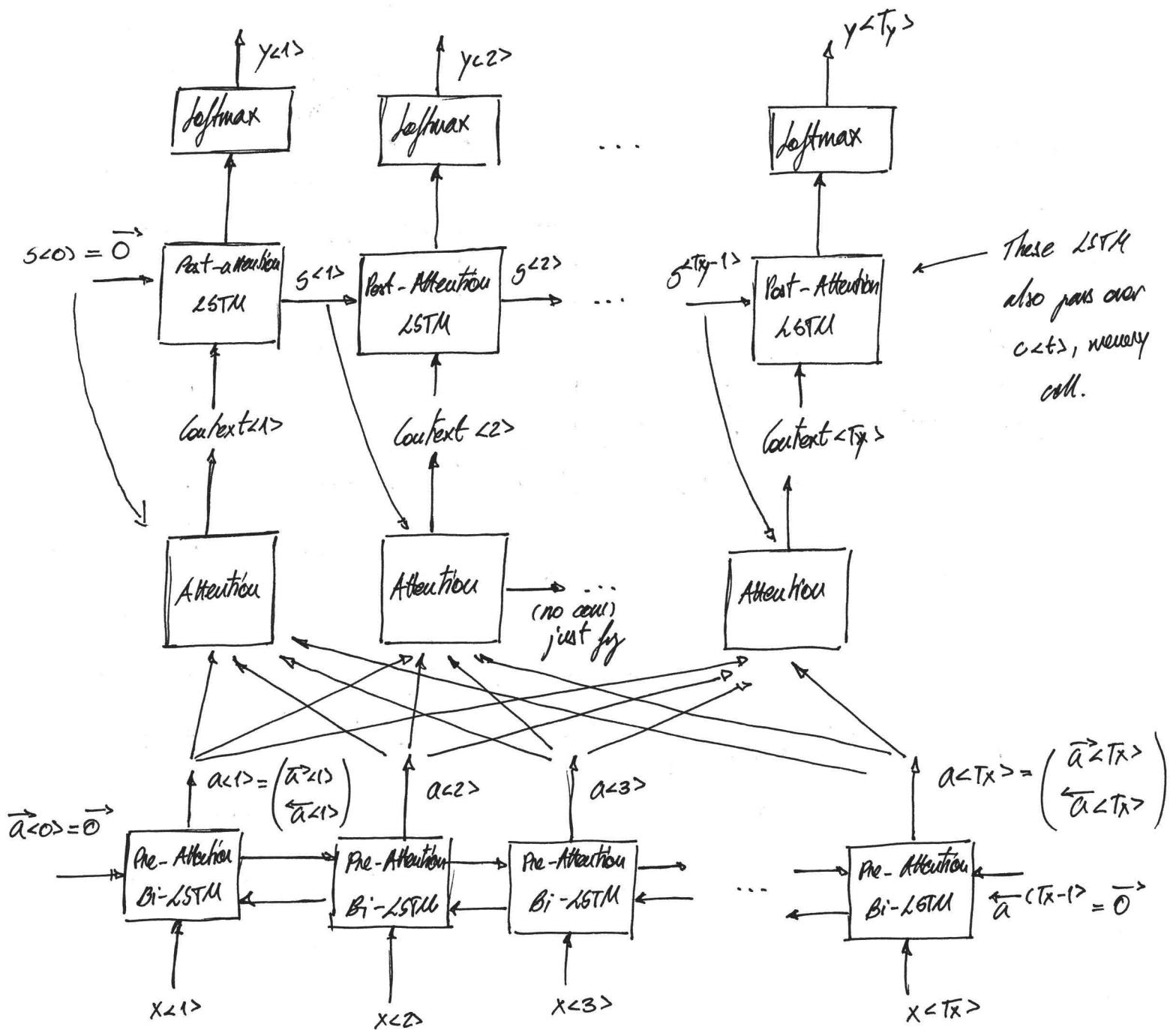
- Unrolled correct output, collapse characters not separated by blank.

* Trigger word detection →



- This sequence will work but the timing set is unbalanced a lot of times.
- You can force the output to 1 after a few cycles of the trigger word.

Attention Model



+ *Context* →

- Post-attention LSTM also passes through the memory cell $c^{<i>}$, handles $s^{<i>}$
- Unlike for other purposes (language modeling), in here $y^{<i-1>}$ is not seen as an input to the next time step. There no high correlation between previous output and next.
- $\alpha^{<i>} = [\vec{\alpha}^{<i>}, \overleftarrow{\alpha}^{<i>}]$, concatenation of bi-dir activations.

These LSTM
also pass over
 $c^{<t>}$, memory
cell.

