# Training RNNs

In RNN: Back Propagation Through Time (BPTT).

$$\overline{S}_t = \phi \left( \overline{x}_t \cdot W_x + \overline{S}_{t-1} \cdot W_S \right)$$

e.g., $\quad \overline{S}_t = \tanh \left( \overline{x}_t \cdot W_x + \overline{S}_{t-1} \cdot W_S \right)$

$$\overline{y}_t = \sigma \left( \overline{S}_t \cdot W_y \right)$$

$\downarrow$
Softmax

$$E_t = \left( \overline{d}_t - \overline{y}_t \right)^2 \qquad \text{error function is MSE}$$

How does BPTT work?

$t = 3$

$$E_3 = (\overline{d}_3 - \overline{y}_3)^2$$

$\downarrow \qquad \llcorner_\circ$ Calculated output

desired
output

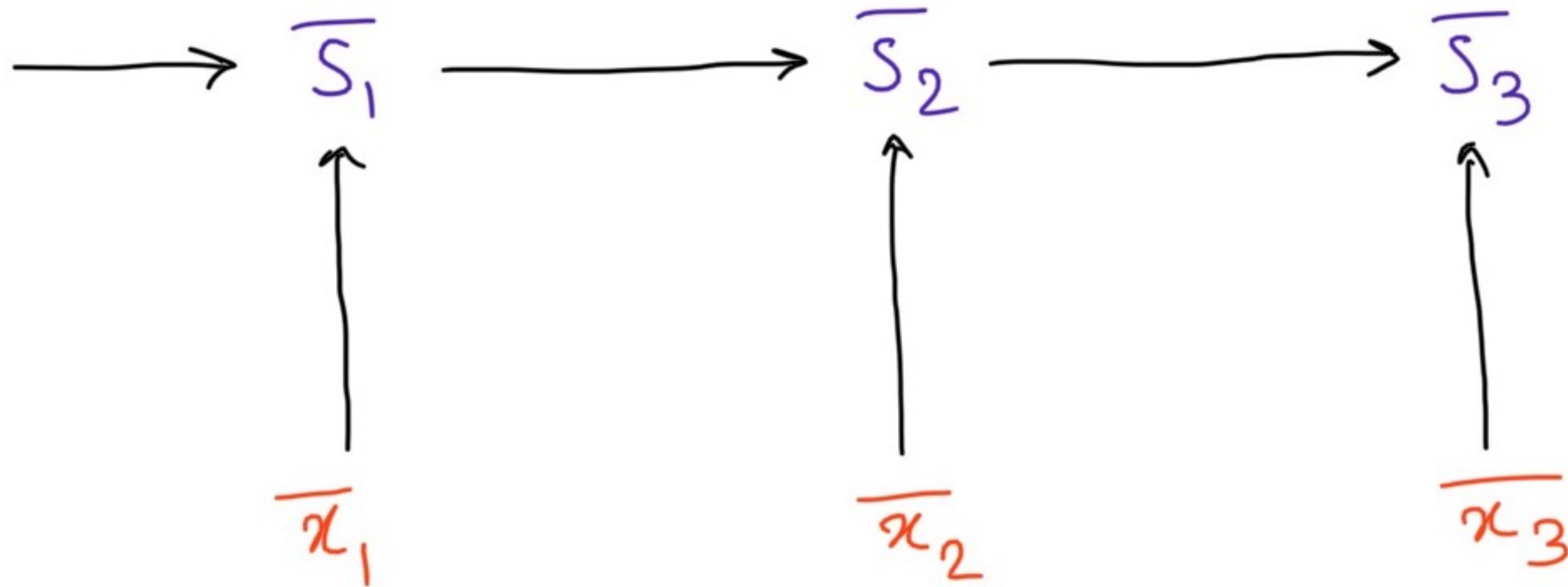but we also need info from $t=2$ and $t=1$.
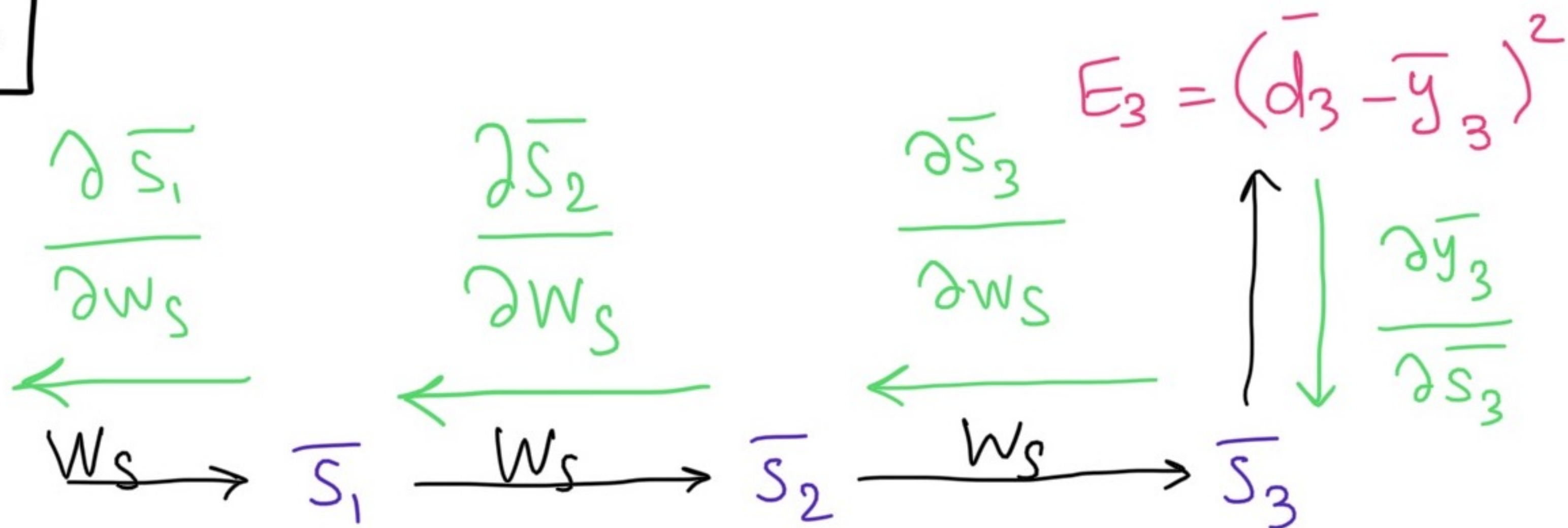
$$\boxed{W_y}$$

$$\boxed{\dfrac{\partial E_3}{\partial w_y} = \dfrac{\partial E_3}{\partial \bar{y}_3} \cdot \dfrac{\partial \bar{y}_3}{\partial w_y}}$$

$$E_3 = (\bar{d}_3 - \bar{y}_3)^2$$

$$w_y \uparrow \quad \downarrow \dfrac{\partial \bar{y}_3}{\partial w_y}$$

$$\longrightarrow \bar{S}_1 \longrightarrow \bar{S}_2 \longrightarrow \bar{S}_3$$

$$\uparrow \bar{x}_1 \qquad \uparrow \bar{x}_2 \qquad \uparrow \bar{x}_3$$

$$\boxed{W_S}$$

$$E_3 = (\bar{d}_3 - \bar{y}_3)^2$$

$$\dfrac{\partial \bar{S}_1}{\partial w_S} \qquad \dfrac{\partial \bar{S}_2}{\partial w_S} \qquad \dfrac{\partial \bar{S}_3}{\partial w_S} \qquad \uparrow \downarrow \dfrac{\partial \bar{y}_3}{\partial \bar{S}_3}$$

$$\xleftarrow{} \qquad \xleftarrow{} \qquad \xleftarrow{}$$

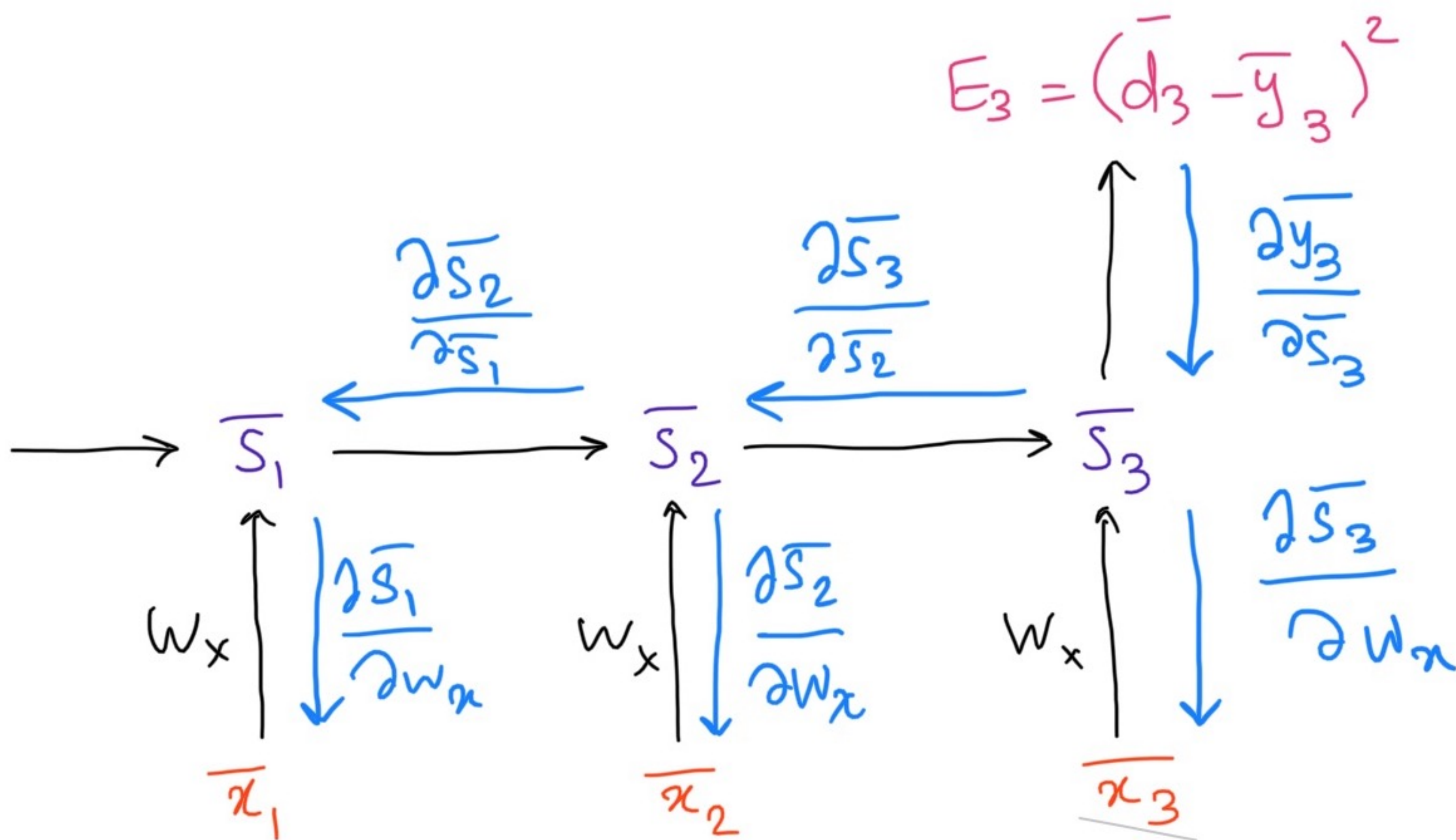$$\xrightarrow{W_S} \bar{S}_1 \xrightarrow{W_S} \bar{S}_2 \xrightarrow{W_S} \bar{S}_3$$

$$\dfrac{\partial E_3}{\partial W_S} = \dfrac{\partial E_3}{\partial \bar{y}_3} \cdot \dfrac{\partial \bar{y}_3}{\partial \bar{S}_3} \cdot \dfrac{\partial \bar{S}_3}{\partial w_S}$$

$$+ \dfrac{\partial E_3}{\partial \bar{y}_3} \cdot \dfrac{\partial \bar{y}_3}{\partial \bar{S}_3} \cdot \dfrac{\partial \bar{S}_3}{\partial \bar{S}_2} \cdot \dfrac{\partial \bar{S}_2}{\partial w_S}$$

$$+ \dfrac{\partial E_3}{\partial \bar{y}_3} \cdot \dfrac{\partial \bar{y}_3}{\partial \bar{S}_3} \cdot \dfrac{\partial \bar{S}_3}{\partial \bar{S}_2} \cdot \dfrac{\partial \bar{S}_2}{\partial \bar{S}_1} \cdot \dfrac{\partial \bar{S}_1}{\partial w_S}$$

$\boxed{W_x}$



$$E_3 = \left(\bar{d}_3 - \bar{y}_3\right)^2$$

$$\frac{\partial E_3}{\partial W_x} = \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial W_x}$$

$$+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial W_x}$$

$$+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \cdot \frac{\partial \bar{s}_1}{\partial W_x}$$

More generally: $\dfrac{\partial E_N}{\partial W_x} = \sum\limits_{i=1}^{N} \dfrac{\partial E_N}{\partial \bar{y}_N} \cdot \dfrac{\partial \bar{y}_N}{\partial \bar{s}_i} \cdot \dfrac{\partial \bar{s}_i}{\partial W_x}$

- for training: updating weights every $N$ steps
  $$\underbrace{\phantom{\text{every N steps}}}_{\text{mini-batch}}$$

- what happens if we have too many steps?
  up to 8-10 steps, this BPTT works,
  but beyond that, the vanishing gradient
  problem happens. $\longrightarrow$ LSTM is born!

- Gradient clipping: → avoiding Exploding Gradient Problem.

At eatch timestep $t$:

$$\delta = \frac{\partial y}{\partial W_{ij}} > \text{threshold ?}$$

If so, normalize the gradient.