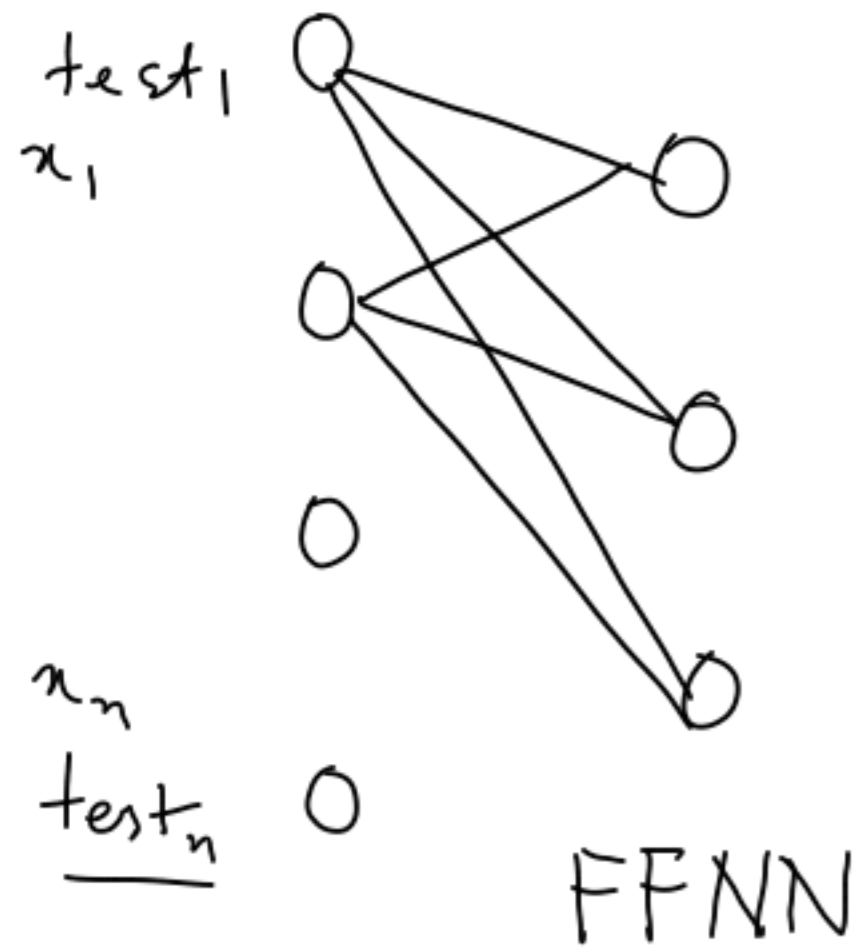


# Lec 13: Recurrent Neural Networks



$0 \leftarrow dis 1$

$0$

$0 \leftarrow dis m$

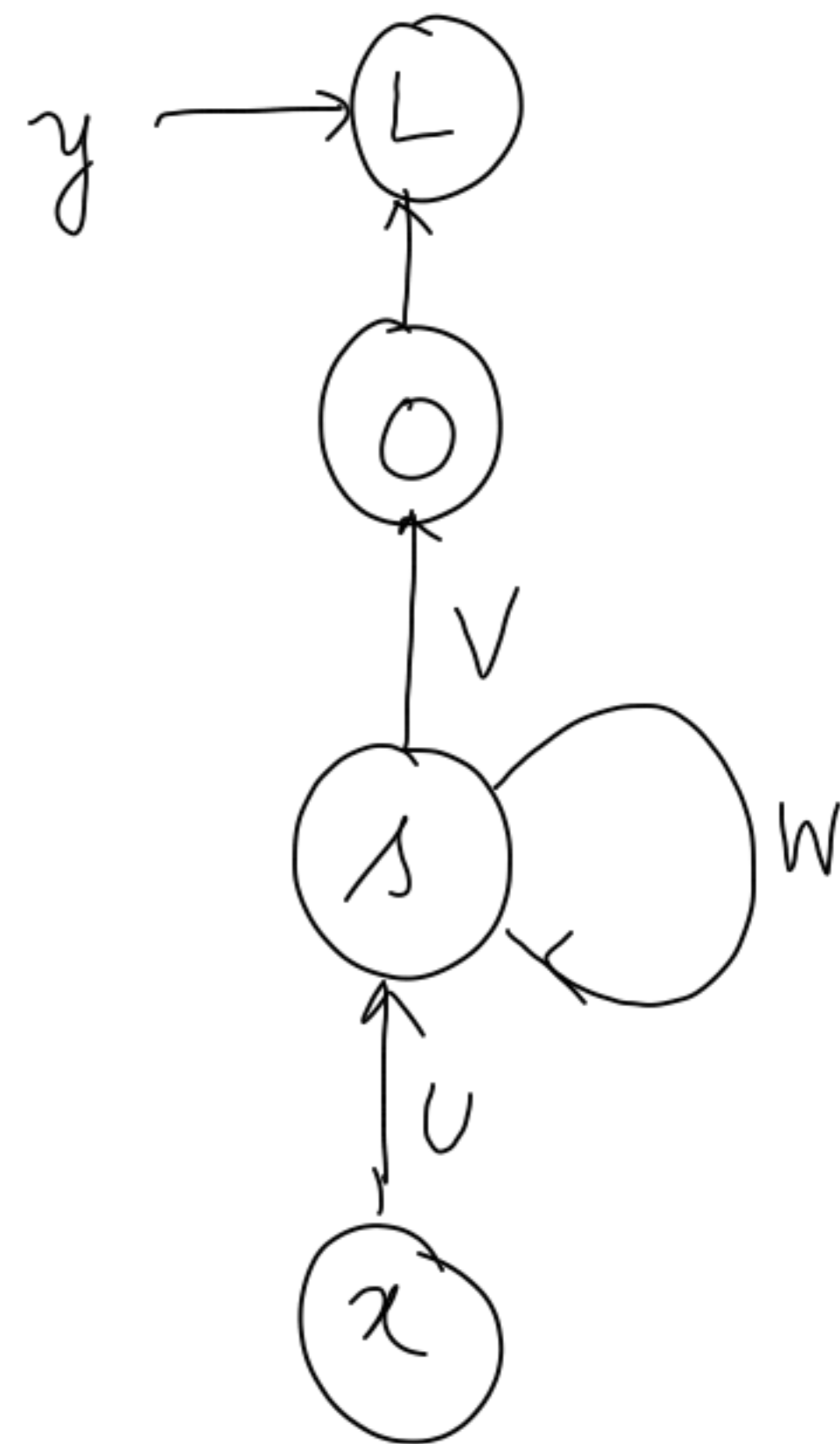
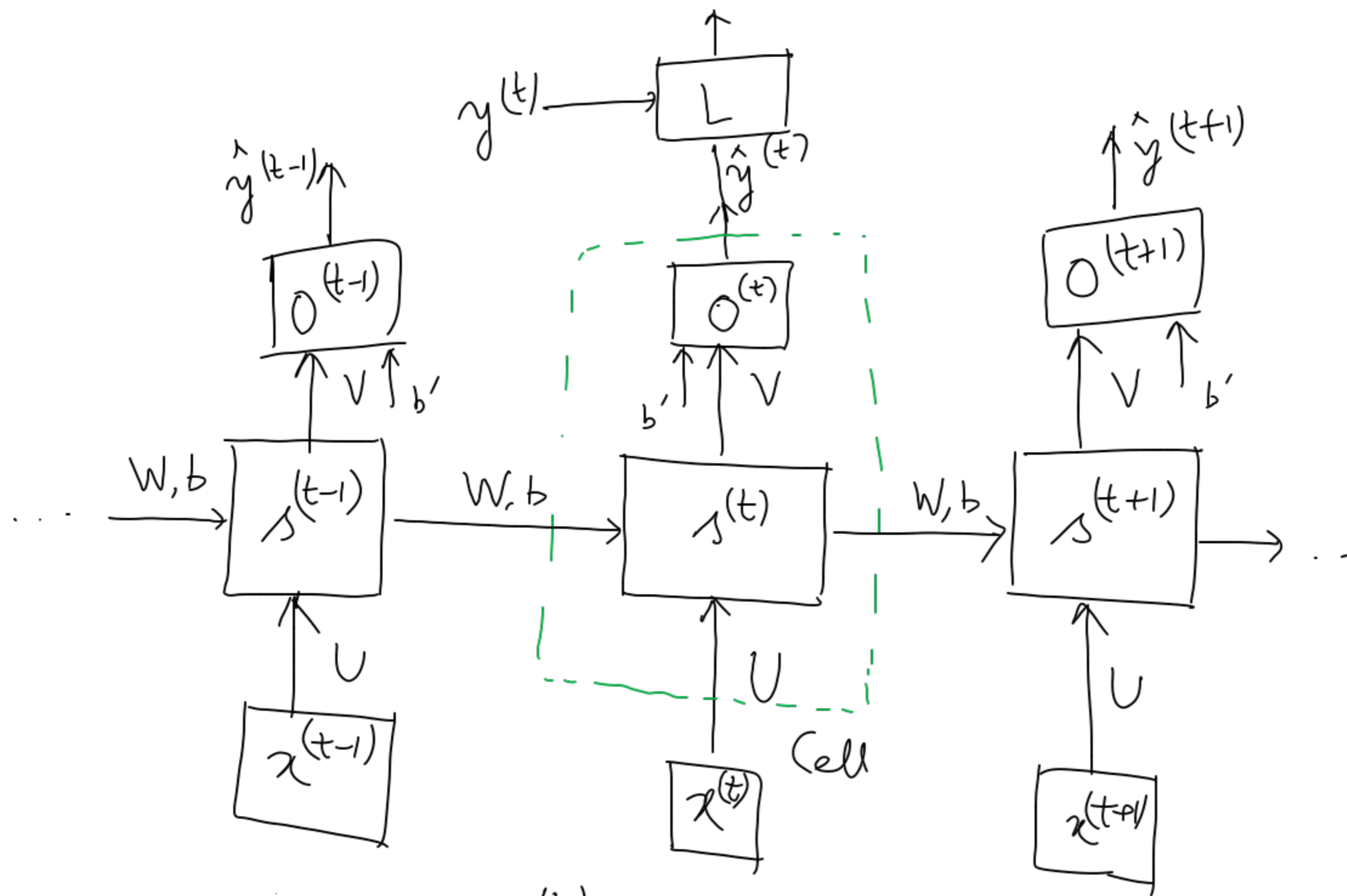
RNN is a way to get around this  $\rightarrow$  changes the architecture.

Data is sequential.  $\leftarrow$  Ordering matters

$\leftarrow$  Length is variable (input/output)

Sentence:

Examples: auto complete, grammar error, virtual assistants, language translation.



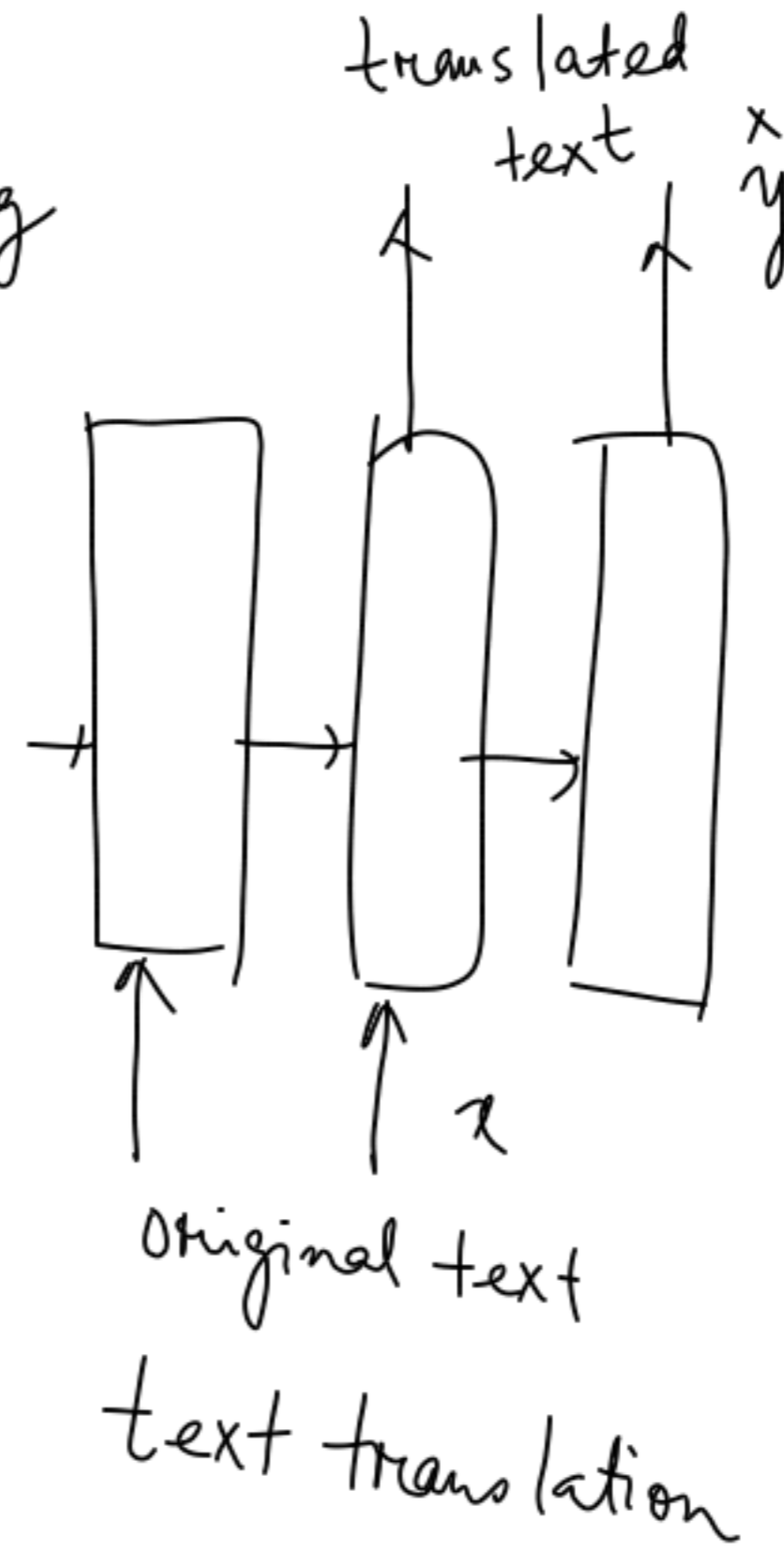
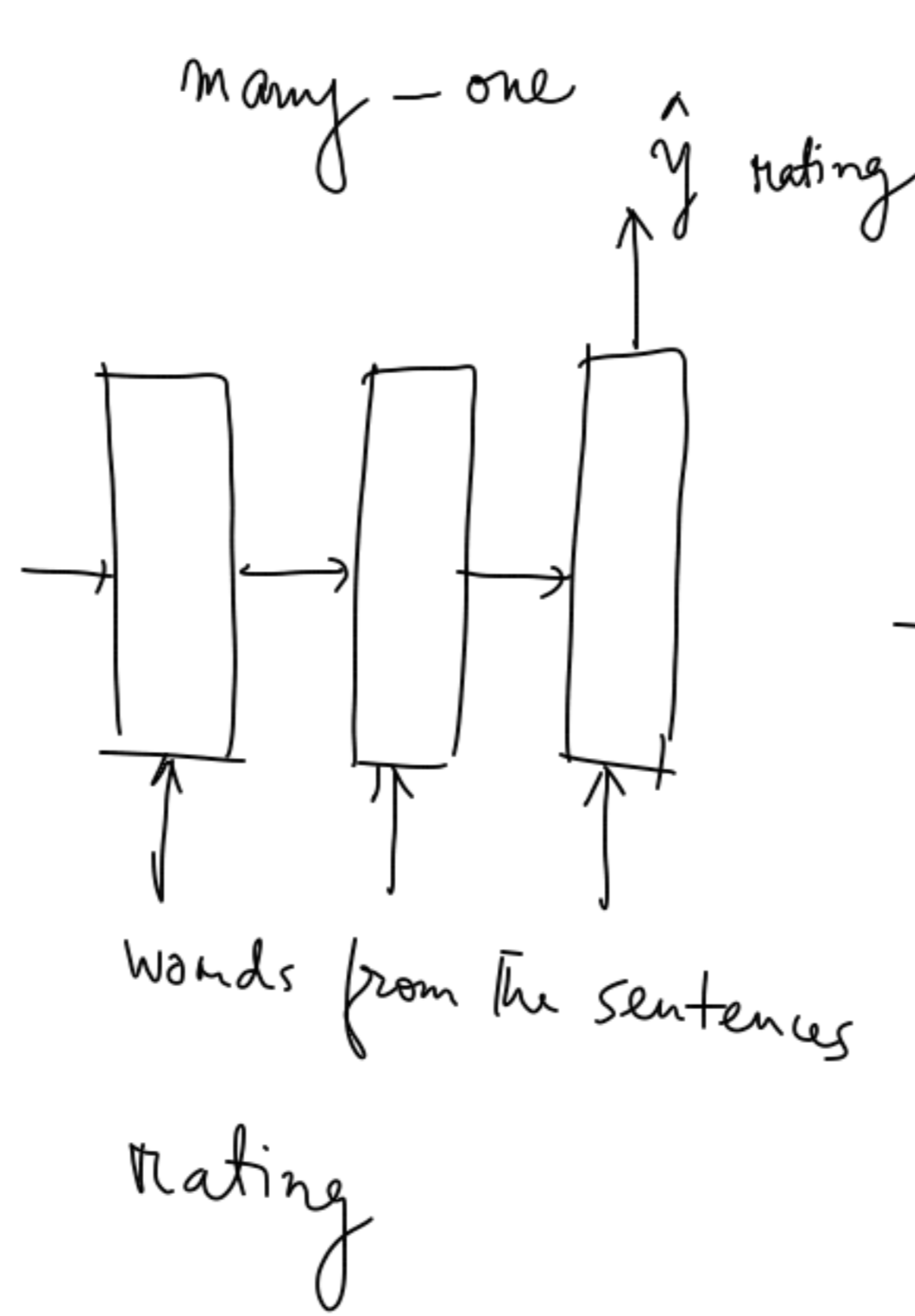
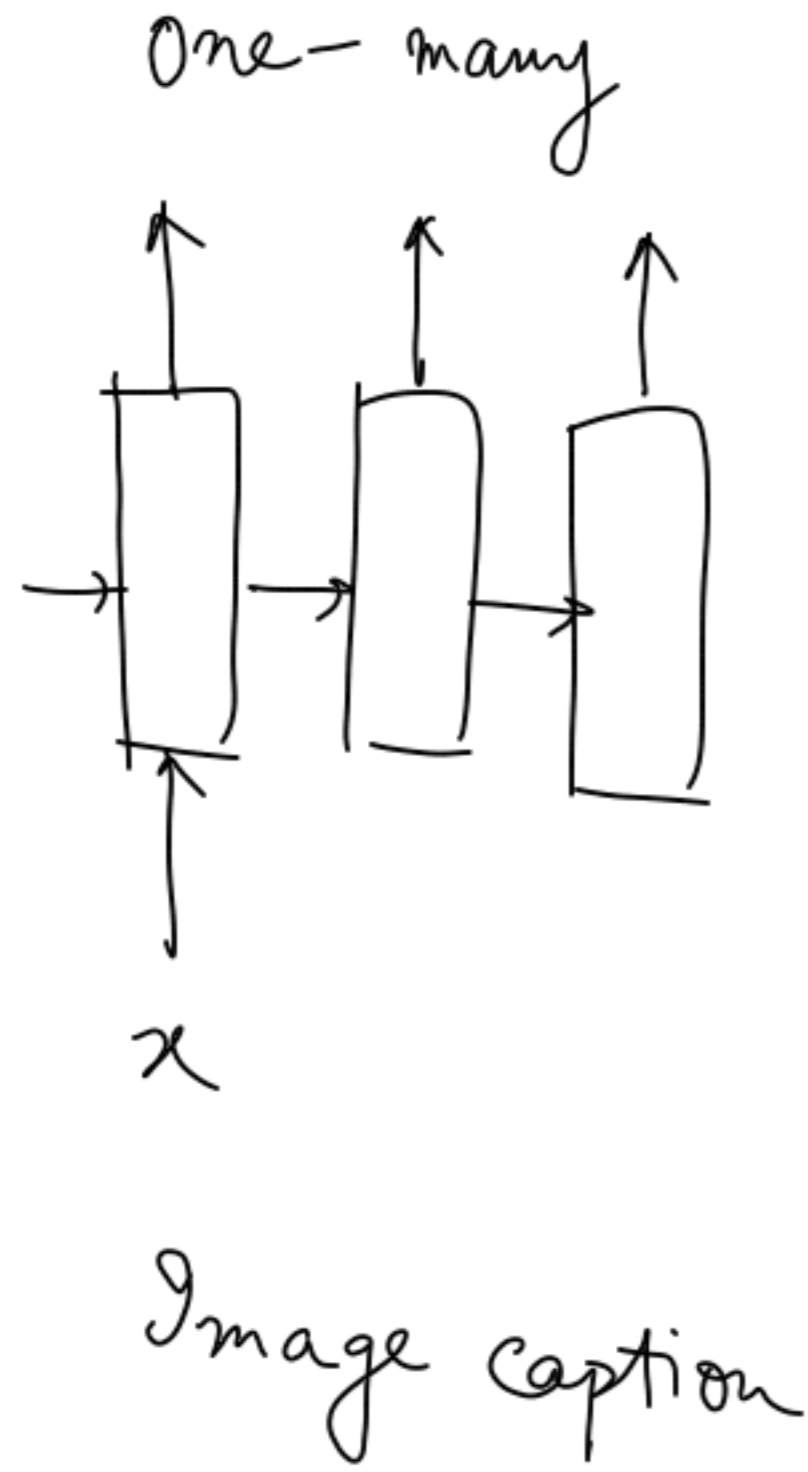
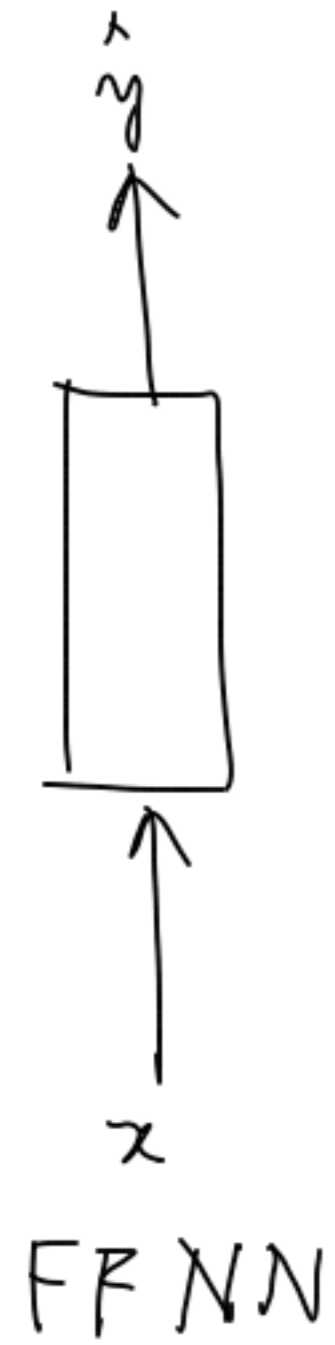
$$s^{(t)} = g(W s^{(t-1)} + U x^{(t)} + b), \quad g \equiv \sigma, \tanh, \text{ReLU}$$

$$o^{(t)} = V s^{(t)} + b'$$

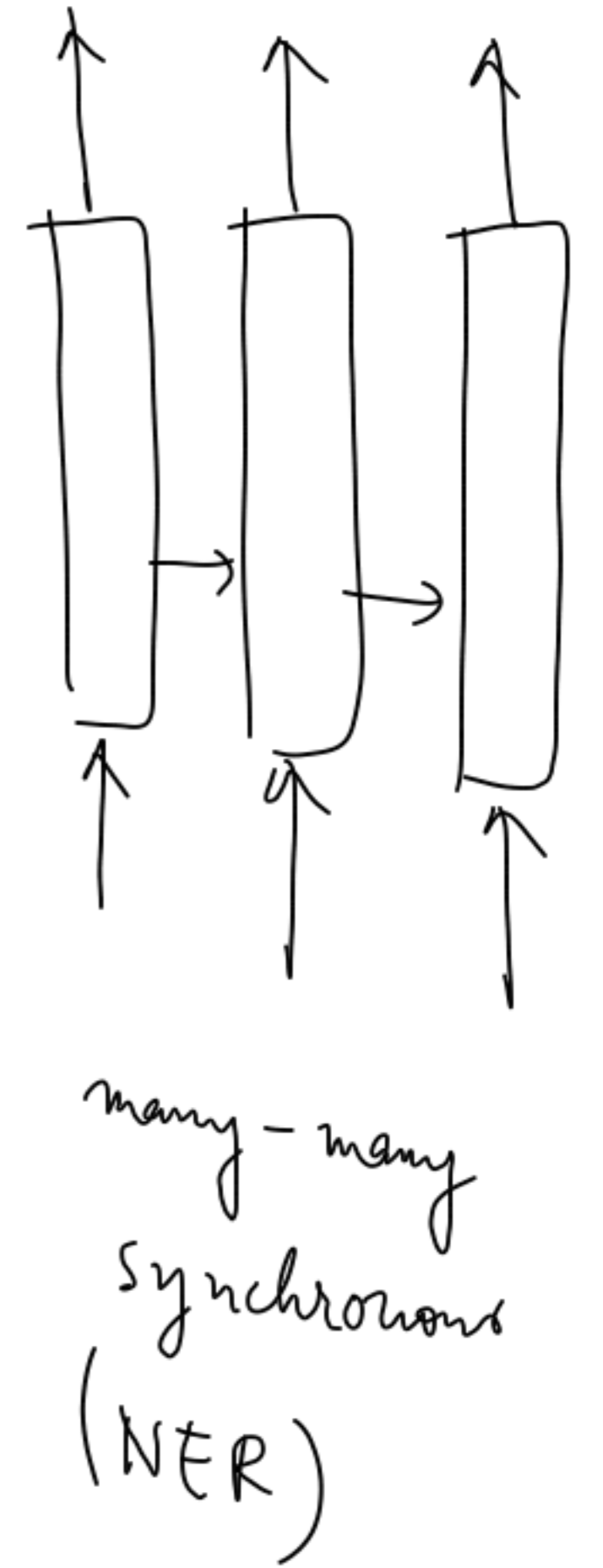
$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$



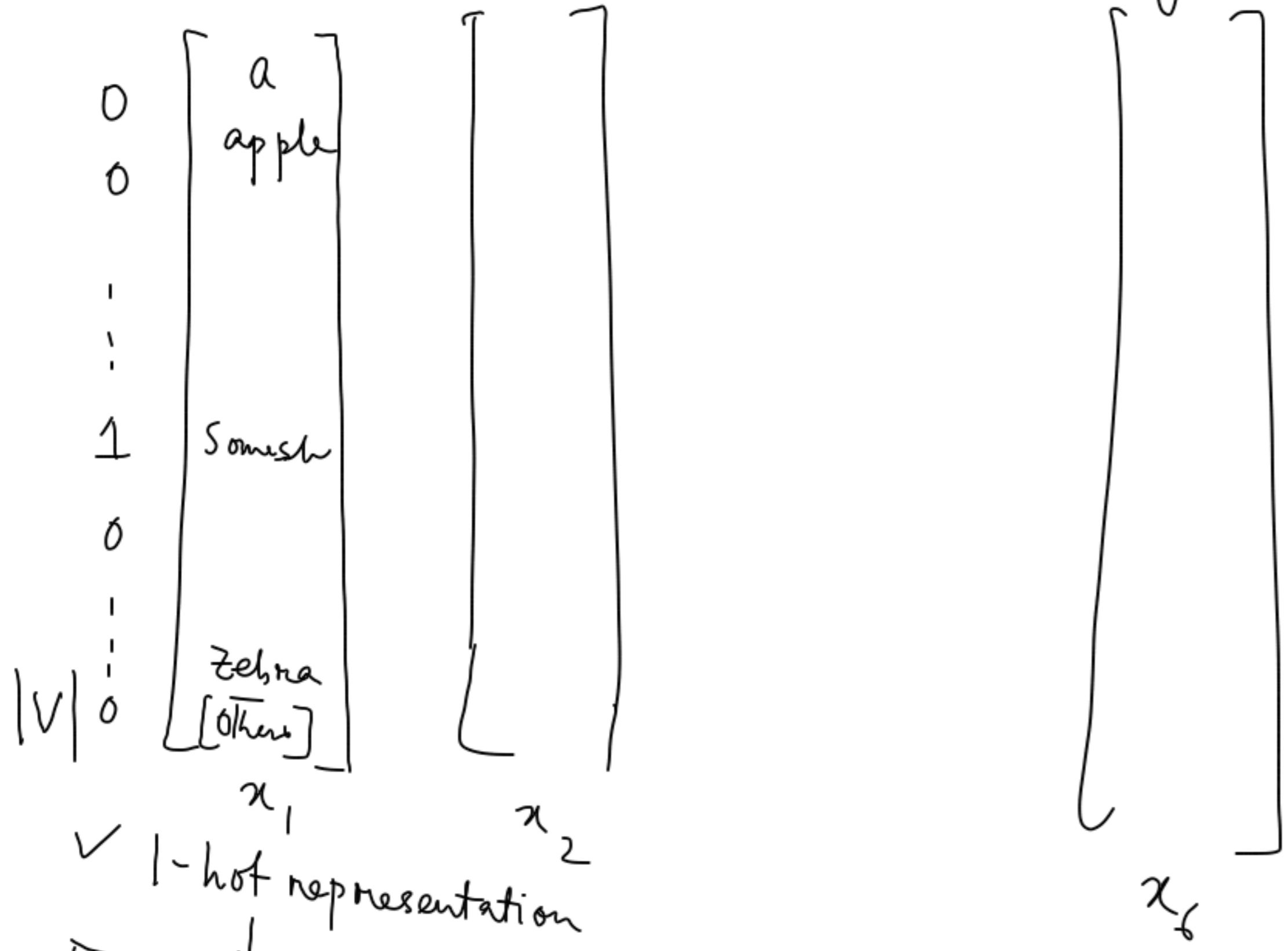
# RNN possible structures



many-many  
staggered  
 $\text{len}(x) \neq \text{len}(\hat{y})$

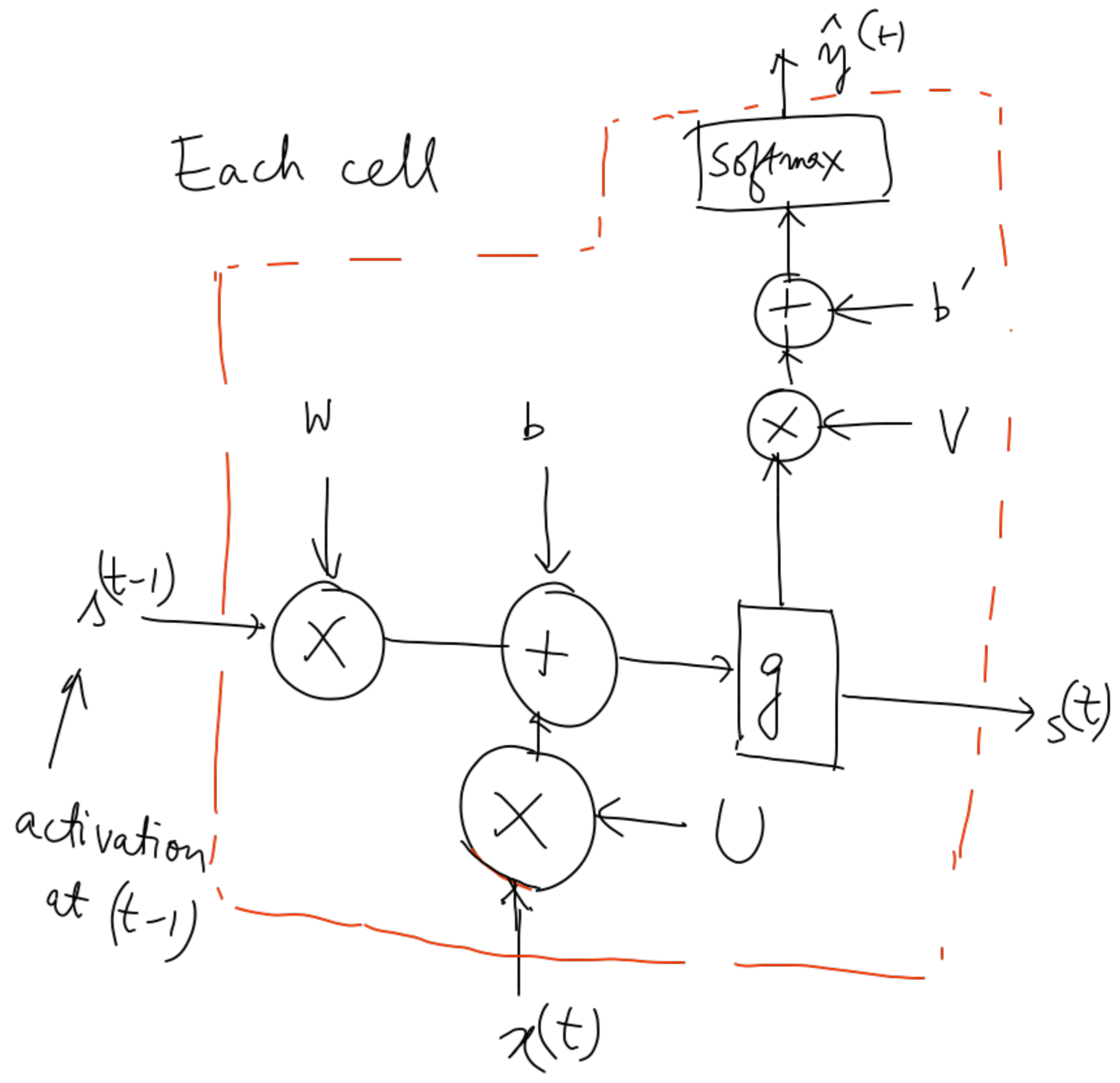


Some sh worked at Cisco in Bengaluru



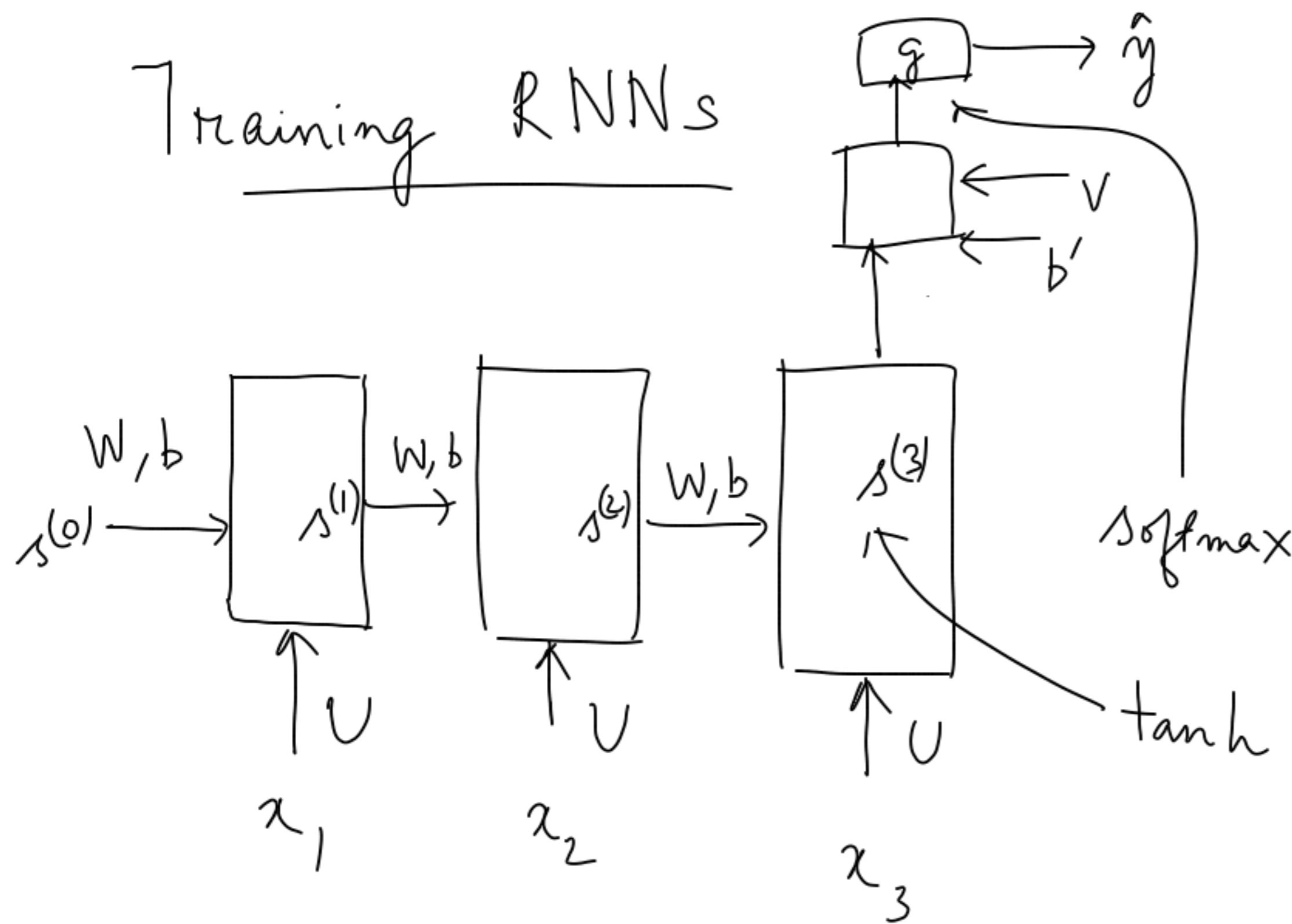
Embedding

Each cell





# Training RNNs



$$\text{Loss}(y, \hat{y}) = -\sum y_k \log \hat{y}_k = L$$

↓ Categorical cross entropy fn.

$$= -\sum_{k=1}^K I\{y=k\} \log \left( \frac{\exp\{w^T \dots\}}{\sum \exp\{\dots\}} \right)$$

BPTT: Backprop through time

Forward pass: compute  $s^{(1)}, s^{(2)}, \dots$

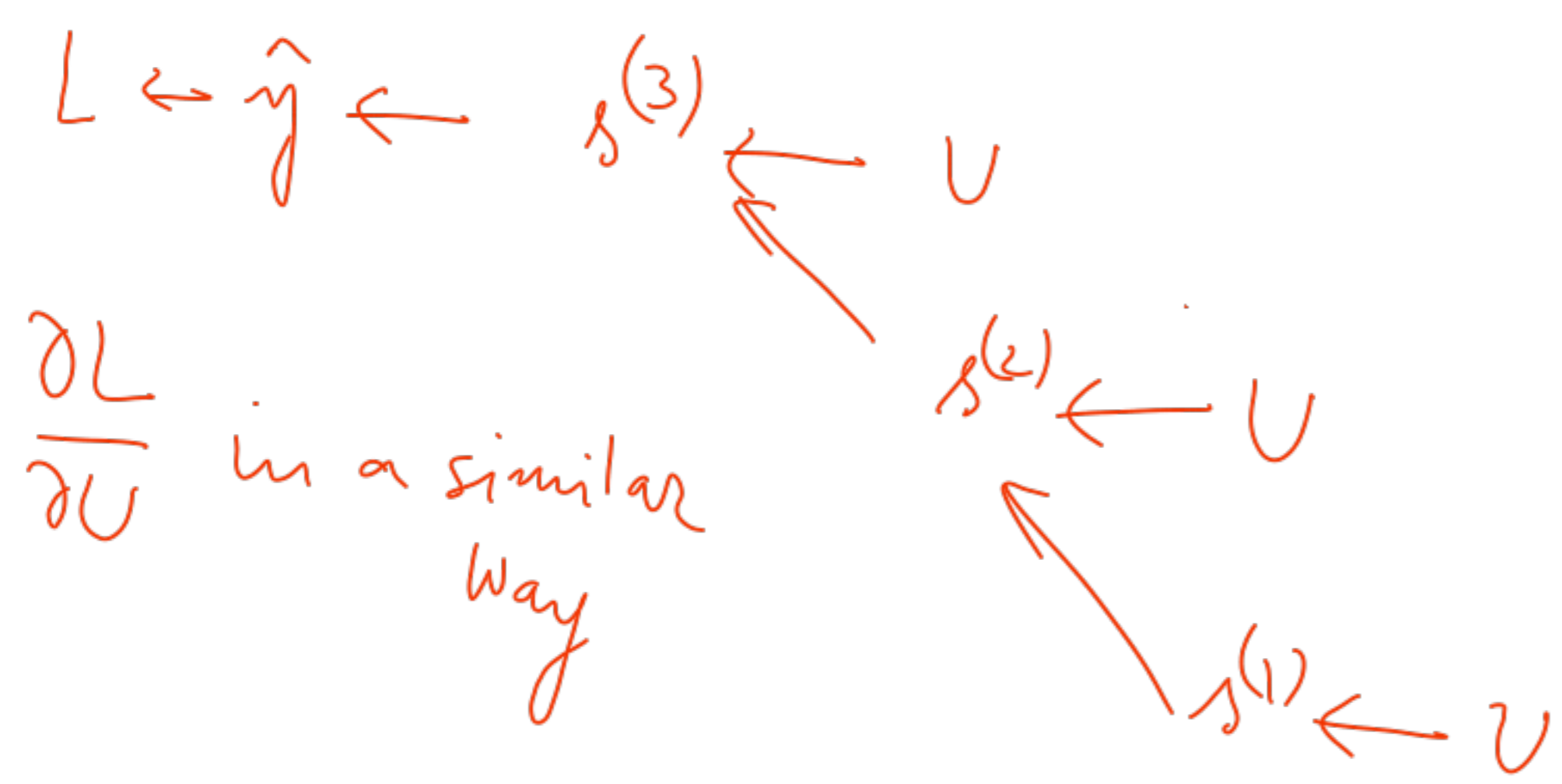
$$s^{(t)} = \tanh \left( W s^{(t-1)} + U x^{(t)} + b \right)$$

$$\hat{y} = \text{softmax} \left( V s^{(T)} + b' \right)$$

kind  $\leftarrow \frac{\partial L}{\partial \hat{y}_k}, \frac{\partial L}{\partial \hat{y}_k}, \frac{\partial L}{\partial \hat{y}_k}, \frac{\partial L}{\partial \hat{y}_k}, \frac{\partial L}{\partial \hat{y}_k}, \frac{\partial L}{\partial \hat{y}_k}$

$$\frac{\partial L}{\partial \hat{y}_k} = \frac{\partial L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial \hat{y}_k}; \quad \frac{\partial L}{\partial \hat{y}_k} =$$

$$\frac{\partial L}{\partial \hat{y}_k} = \frac{\partial}{\partial \hat{y}_k} \left( -\sum_k \hat{y}_k \log \hat{y}_k \right)$$



$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial w}$$

$$= \frac{\partial L}{\partial \hat{y}_k} \left( \frac{\partial \hat{y}_k}{\partial s^{(3)}} \right) \frac{\partial s^{(3)}}{\partial w}$$

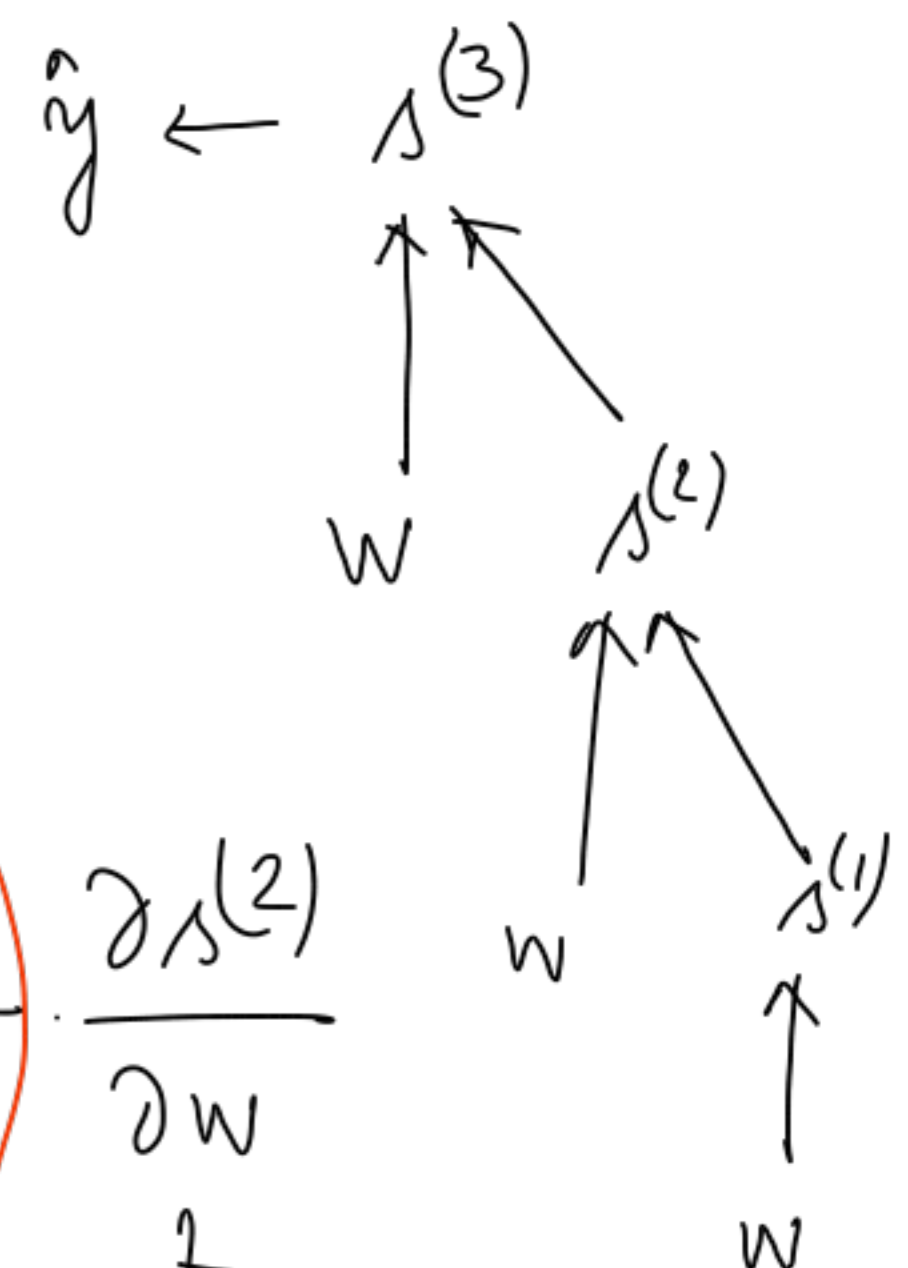
$$+ \frac{\partial L}{\partial \hat{y}_k} \left( \frac{\partial \hat{y}_k}{\partial s^{(3)}} \cdot \frac{\partial s^{(3)}}{\partial s^{(2)}} \right) \frac{\partial s^{(2)}}{\partial w}$$

$$+ \frac{\partial L}{\partial \hat{y}_k} \left( \frac{\partial \hat{y}_k}{\partial s^{(3)}} \cdot \frac{\partial s^{(3)}}{\partial s^{(1)}} \right) \times$$

$$\frac{\partial s^{(1)}}{\partial w} = \frac{\partial s^{(2)}}{\partial s^{(1)}} \cdot \frac{\partial s^{(1)}}{\partial w}$$

$$= \sum_{i=1}^I \frac{\partial L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial s^{(i)}} \cdot \frac{\partial s^{(i)}}{\partial w}$$

in a similar way



$$f(g_1(x), g_2(x))$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g_1} \cdot \frac{\partial g_1}{\partial x} + \frac{\partial f}{\partial g_2} \cdot \frac{\partial g_2}{\partial x}$$

Disadvantages: Vanishing gradient problem

Isha went for a vacation, . . . . . He was overjoyed  
↑

→ In traditional RNN, grammar checker won't be able to detect this.

→ To mitigate, a modification to the cell structure.

