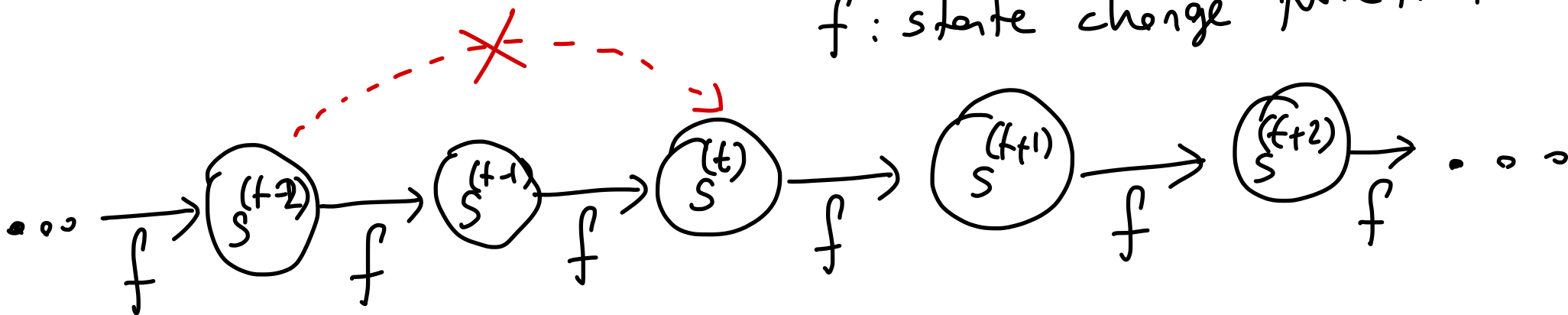


# Recurrent Neural Networks (RNNs)

→ Dynamic Systems

$s^{(t)}$ : state vector at time  $t$ .  
 $f$ : state change function



$$f(s^{(t)}) = \underbrace{W}_{D \times D} \cdot \underbrace{s^{(t)}}_{D \times 1} + \underbrace{b}_{D \times 1}$$

$s^{(t+1)} \in \mathbb{R}^{D \times 1}$

parameters:  $[W, b]$

$$\underbrace{s^{(t+1)}}_{\substack{\uparrow \\ \text{next} \\ \text{state}}} = f\left(\underbrace{s^{(t)}}_{\substack{\uparrow \\ \text{current} \\ \text{state}}}\right) \quad \forall t.$$

We are going to learn the parameters of  $f$  function.

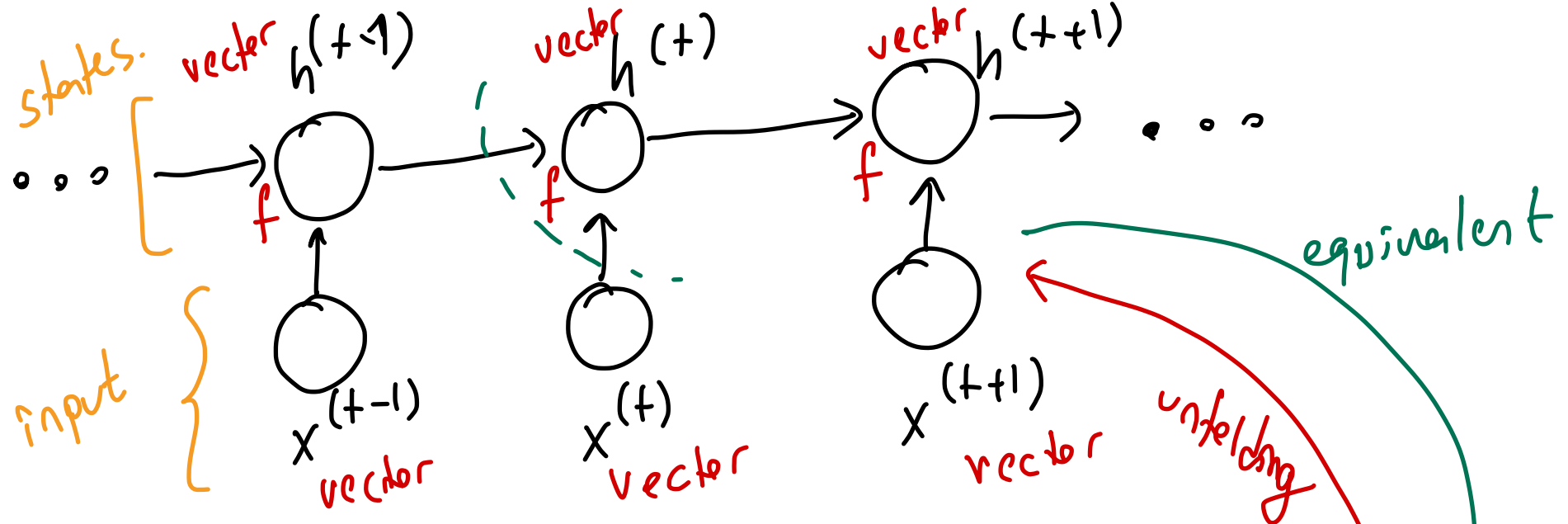
$$s^{(2)} = f(s^{(1)}) \Rightarrow s^{(2)} = W \cdot s^{(1)} + b$$

$$s^{(1)} = f(s^{(0)}) \Rightarrow s^{(1)} = \boxed{W \cdot s^{(0)} + b}$$

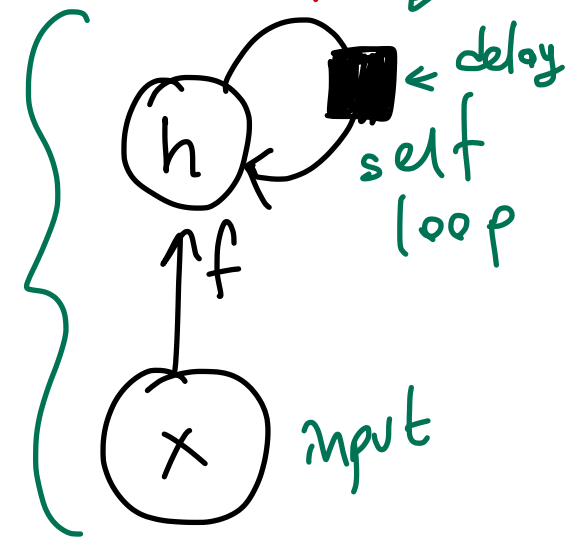
$$s^{(2)} = W [W s^{(0)} + b] + b$$

$$= \underbrace{W \cdot W}_{\tilde{W}} \cdot s^{(0)} + \underbrace{W \cdot b + b}_{\tilde{b}}$$

$$= \tilde{W} \cdot s^{(0)} + \tilde{b}$$

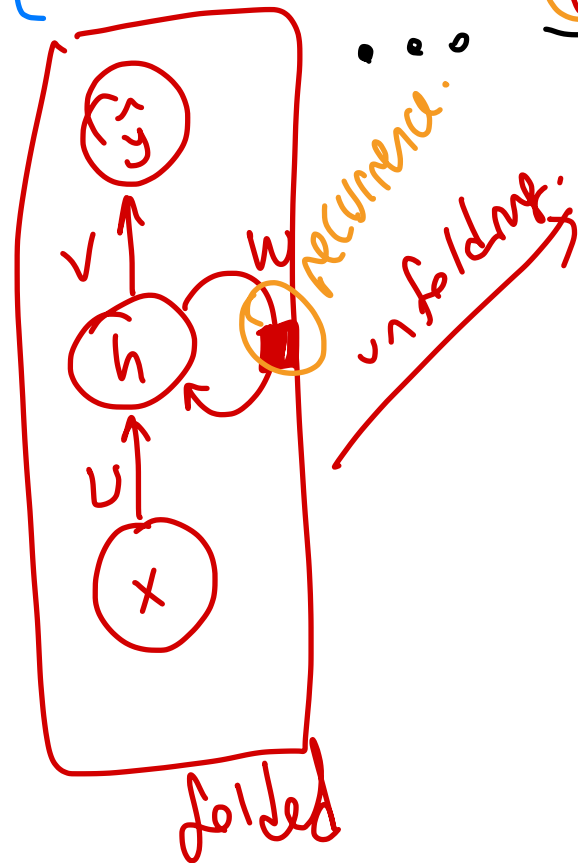
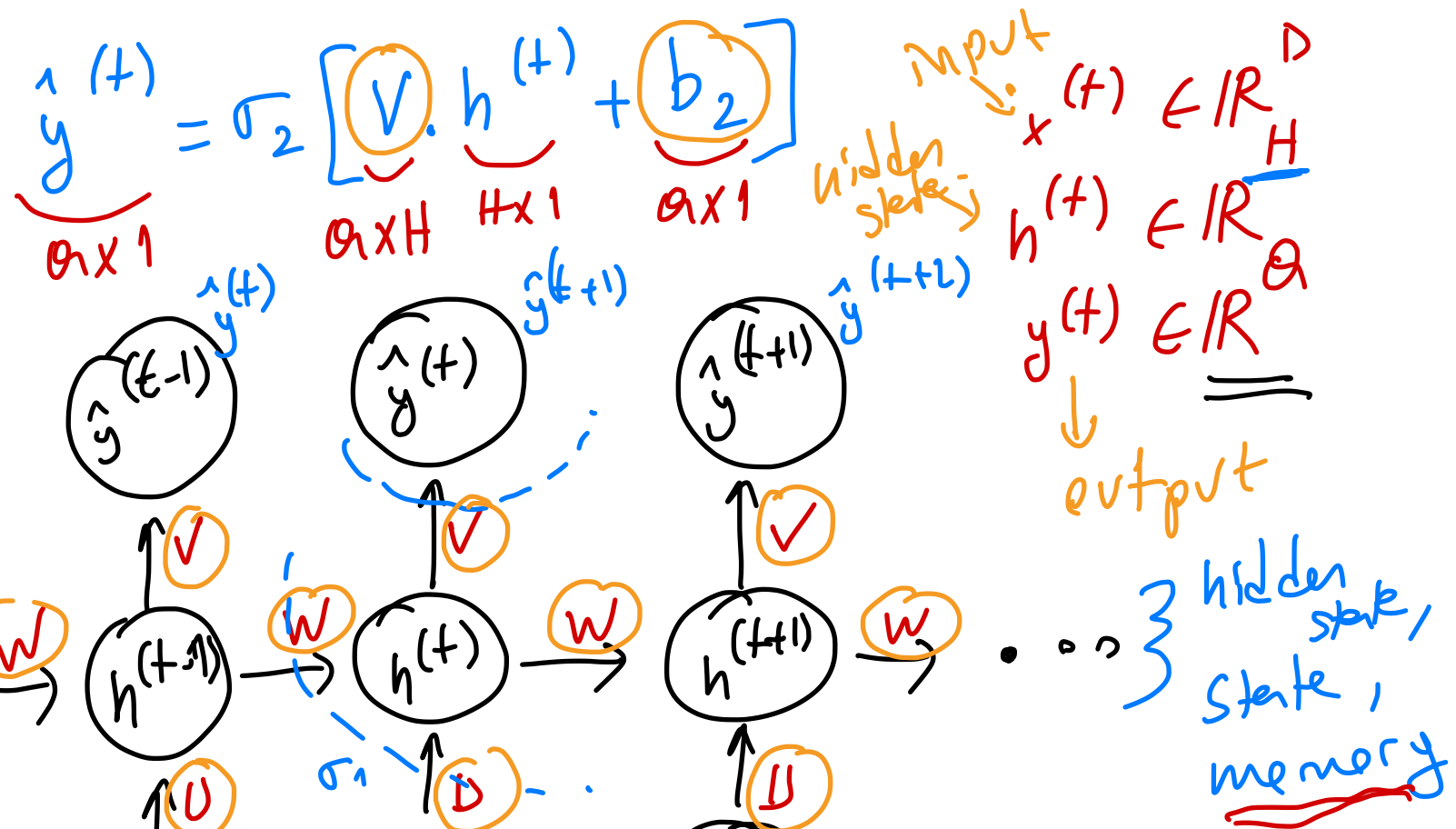


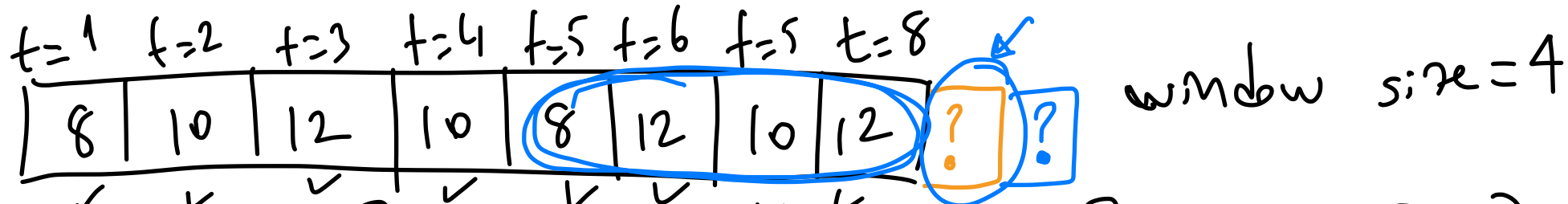
$$\begin{aligned}
 \underline{h^{(t)}} &= f(\text{previous state, current input}) \\
 &= f(\underline{h^{(t-1)}}, \underline{x^{(t)}})
 \end{aligned}$$



$$\left[ y^{(t)} - \hat{y}^{(t)} \right]^2$$

$$\left[ y^{(t+1)} - \hat{y}^{(t+1)} \right]^2$$





$$x^{(1)} = \begin{bmatrix} 8 \\ 10 \\ 12 \\ 10 \end{bmatrix}$$

$$x^{(2)} = \begin{bmatrix} 10 \\ 12 \\ 10 \\ 8 \end{bmatrix}$$

$$x^{(3)} = \begin{bmatrix} 12 \\ 10 \\ 8 \\ 12 \end{bmatrix}$$

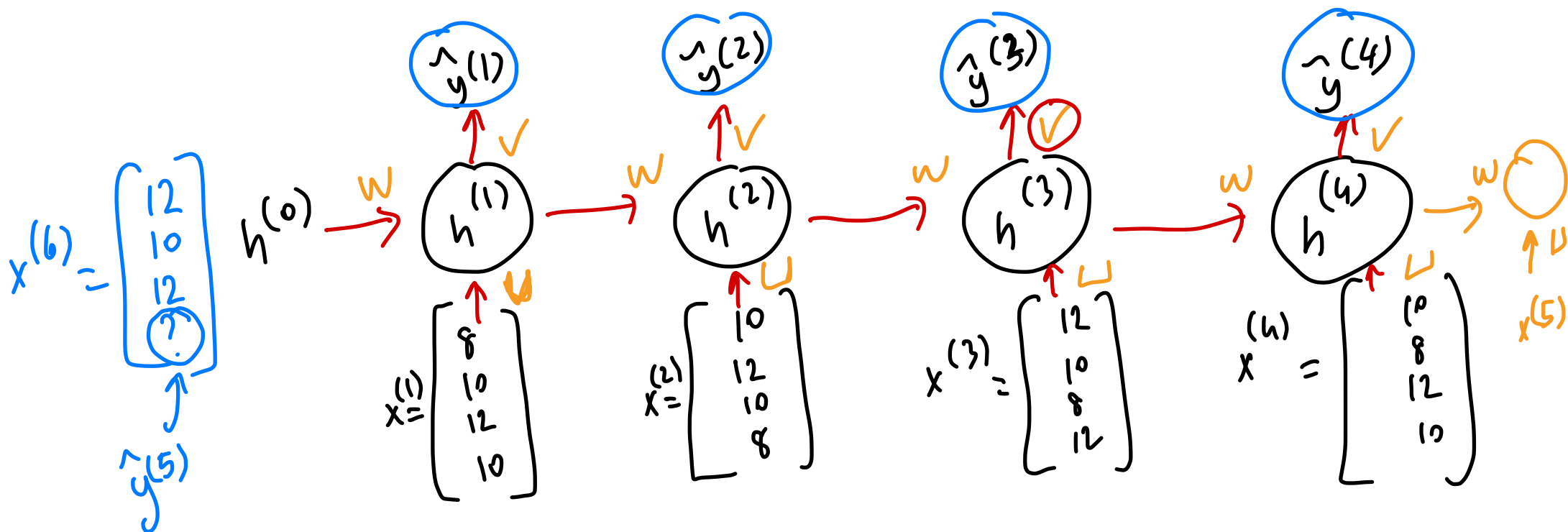
$$x^{(4)} = \begin{bmatrix} 10 \\ 8 \\ 12 \\ 10 \end{bmatrix} \quad x^{(5)} = \begin{bmatrix} 8 \\ 12 \\ 10 \\ 12 \end{bmatrix}$$

$$\hat{y}^{(1)} = [8]$$

$$\hat{y}^{(2)} = [12]$$

$$\hat{y}^{(3)} = [10]$$

$$\hat{y}^{(4)} = [12] \quad y^{(5)} = ?$$



$$\text{Loss [MSE]} = \frac{1}{4} \left[ (y^{(1)} - \hat{y}^{(1)})^2 + \dots + (y^{(4)} - \hat{y}^{(4)})^2 \right]$$

$$\checkmark h^{(5)} = \sigma_1 \left[ \underbrace{\checkmark W}_{20 \times 20} \cdot \underbrace{\checkmark h^{(4)}}_{20 \times 1} + \underbrace{\checkmark U}_{20 \times 4} \cdot \underbrace{\checkmark x^{(5)}}_{4 \times 1} + \checkmark b_1 \right]$$

$$\underbrace{\checkmark \hat{y}^{(5)}}_{1 \times 1} = \sigma_2 \left[ \checkmark V \cdot \underbrace{\checkmark h^{(5)}}_{20 \times 1} + \checkmark b_2 \right]$$

assume.

$$\boxed{H=20}$$

$$\left[ \begin{array}{c} \underline{(H \times H)} \\ \underline{(H \times \text{window size})} \end{array} \right]$$

$$\hat{y}^{(6)} = ?$$

521

parameters

$$\left\{ \begin{array}{l} W \in \mathbb{R}^{20 \times 20} \\ U \in \mathbb{R}^{20 \times 4} \\ V \in \mathbb{R}^{1 \times 20} \\ b_1 \in \mathbb{R}^{20 \times 1} \\ b_2 \in \mathbb{R}^1 \end{array} \right.$$

$$\underline{(a \times H)} \leftarrow$$

$$\underline{(H)}$$

$$\underline{(a)} \leftarrow$$

8, 10, 12, 10, 8, 12, 10, 12, ? - - -

$$x^{(1)} = \begin{bmatrix} 8 \\ 10 \\ 12 \\ 10 \end{bmatrix}$$

$$x^{(2)} = \begin{bmatrix} 10 \\ 12 \\ 10 \\ 8 \end{bmatrix}$$

$$x^{(3)} = \begin{bmatrix} 12 \\ 10 \\ 8 \\ 12 \end{bmatrix}$$

$$y^{(1)} = \begin{bmatrix} 8 \\ 12 \end{bmatrix}$$

$$y^{(2)} = \begin{bmatrix} 12 \\ 10 \end{bmatrix}$$

$$y^{(3)} = \begin{bmatrix} 10 \\ 12 \end{bmatrix}$$

~~$$x^{(4)} = \begin{bmatrix} 10 \\ 8 \\ 12 \\ 10 \end{bmatrix}$$~~

~~$$y^{(4)} = \begin{bmatrix} 12 \\ ? \\ ? \end{bmatrix}$$~~

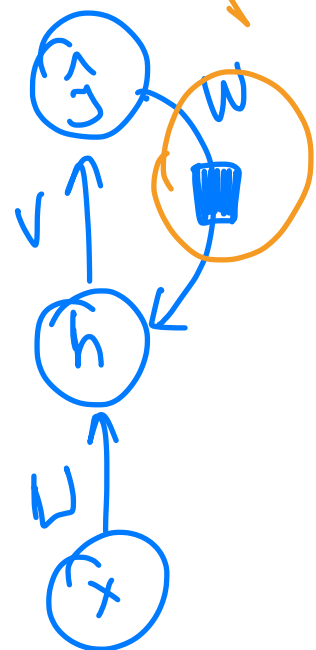
$\frac{20 \times 20}{20 \times 4}$   
 $\frac{20 \times 20}{20 \times 1}$   
 $\frac{20 \times 1}{2 \times 1}$

542 parameters

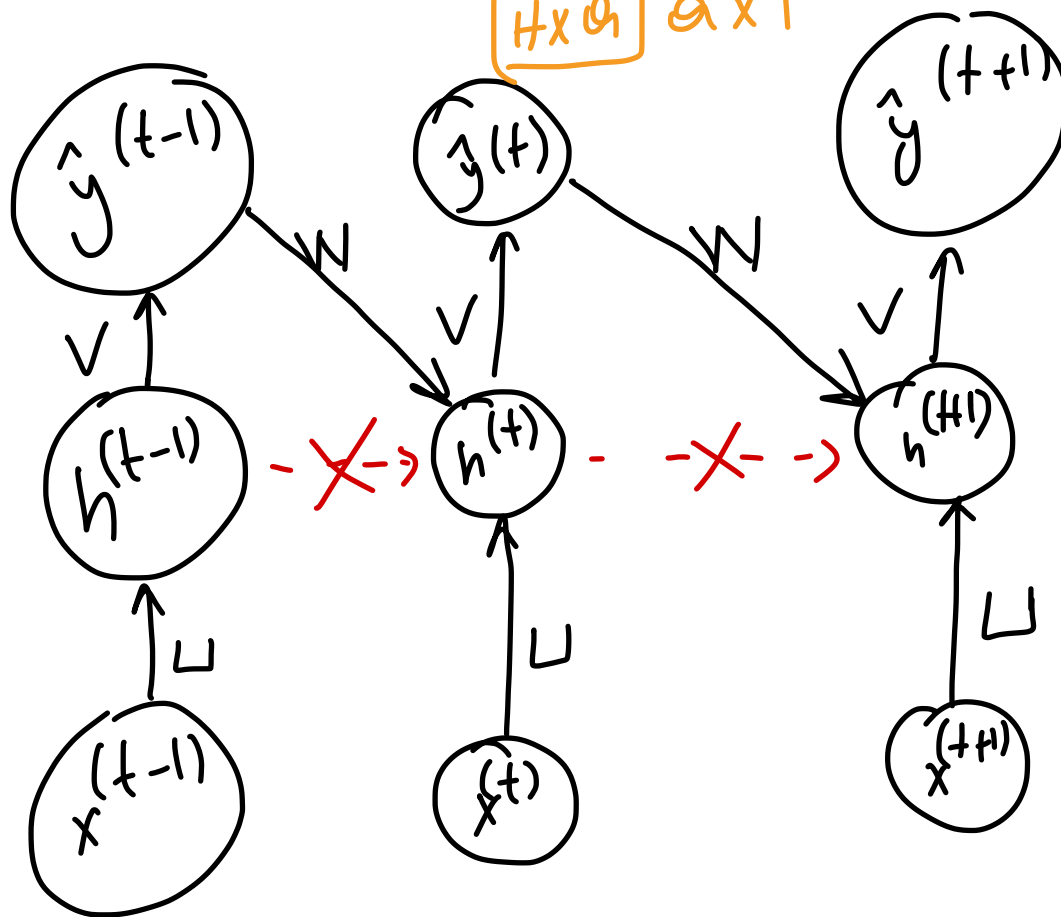
folded  
version

$$h^{(t)} = \sigma_1 \left[ \underbrace{W \cdot \hat{y}^{(t-1)}}_{\text{different } H \times 1 \text{ ax 1}} + \underbrace{U \cdot x^{(t)}}_{\text{some}} + \underbrace{b_1}_{\text{some}} \right]$$

$W$ 's size  
changes from  
 $H \times H$  to  $H \times 1$ .



...



...

$$\hat{y}^{(t)} = \sigma_2 \left[ V \cdot h^{(t)} + b_2 \right]$$



# Sequence Input / Single Output Models

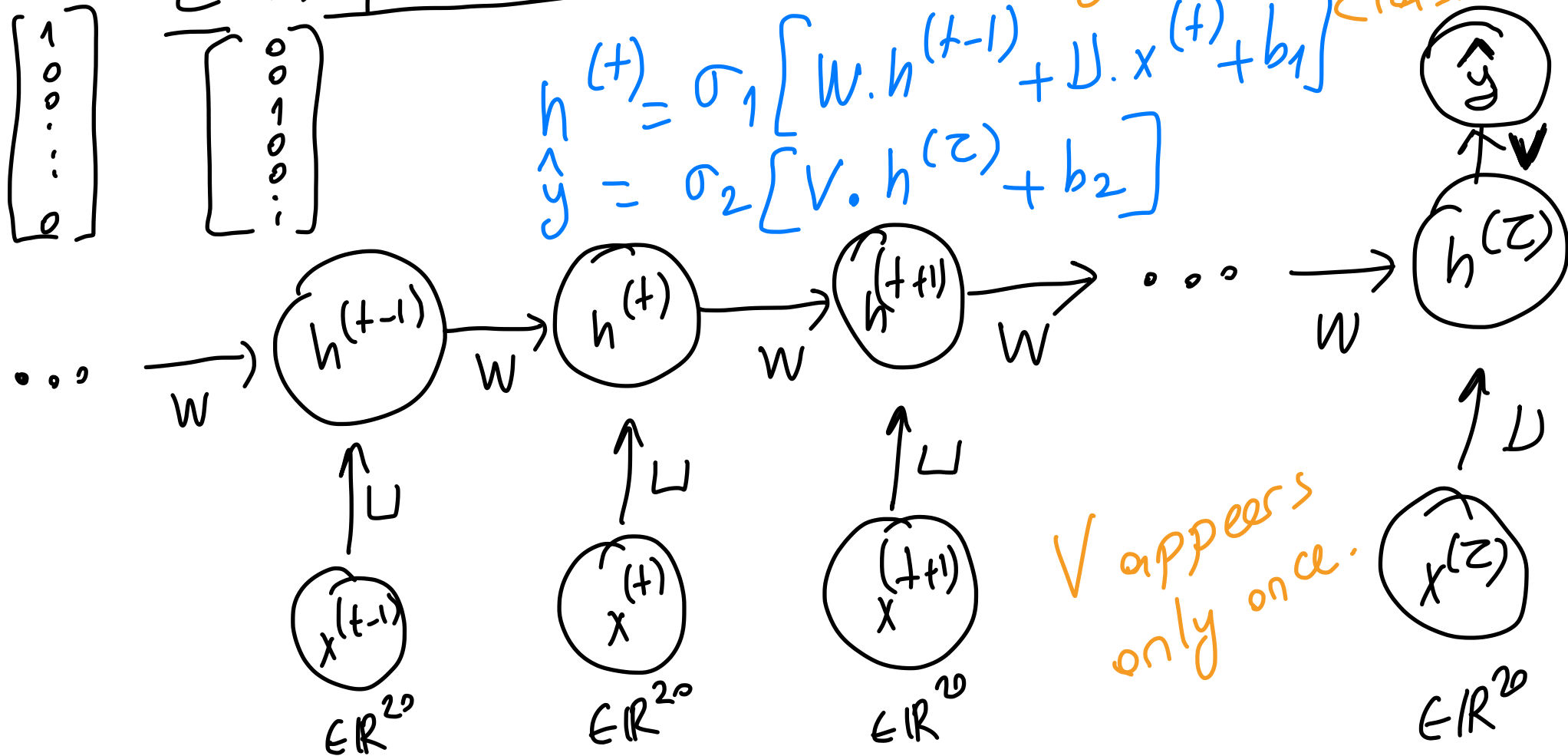
↳ Sentiment Analysis

↳ signal (time series) classification

$Z$  : sequence length.

$$h^{(t)} = \sigma_1 [W \cdot h^{(t-1)} + U \cdot x^{(t)} + b_1]$$

$$\hat{y} = \sigma_2 [V \cdot h^{(Z)} + b_2]$$



$V$  appears only once.

P1: A K L V 0 . . . . L 0

$$Z_1 = 100$$

$$y_1 = +$$

P2: K L V 0 A K K K L 0

$$Z_2 = 10$$

$$y_2 = +$$

⋮

P<sub>1000</sub>: L 0 0 A A K L . . . . . K L 0

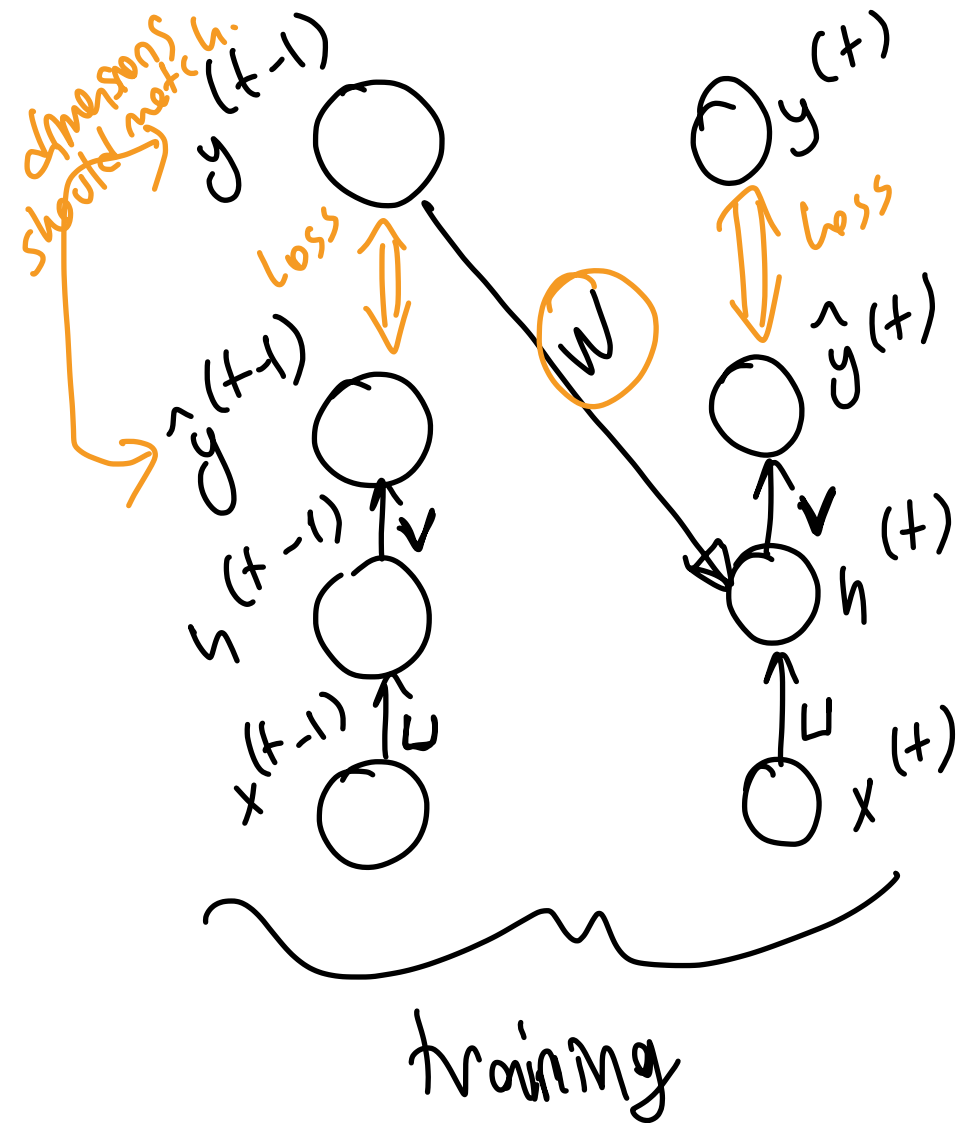
$$y_{1000} = -$$

$$Z_{1000} = 428$$

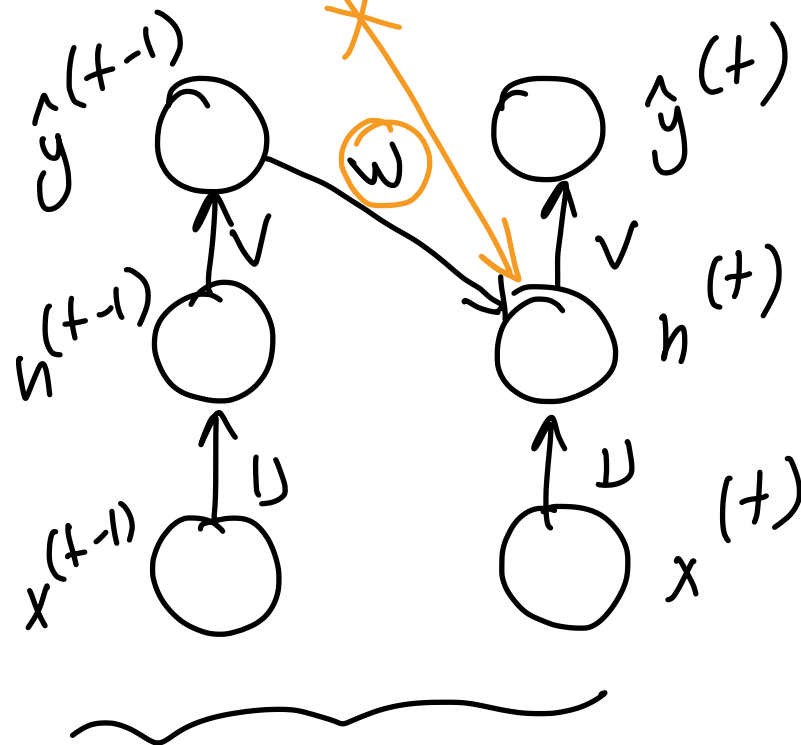
$$x_1 = \begin{bmatrix} A \\ C \\ D \\ \vdots \\ 0 \end{bmatrix}_{20 \times 1}$$

P<sub>1</sub> ⇒ 20 × 100 dimensional  
P<sub>2</sub> ⇒ 20 × 10 dimensional  
P<sub>1000</sub> ⇒ 20 × 428 dimensional.

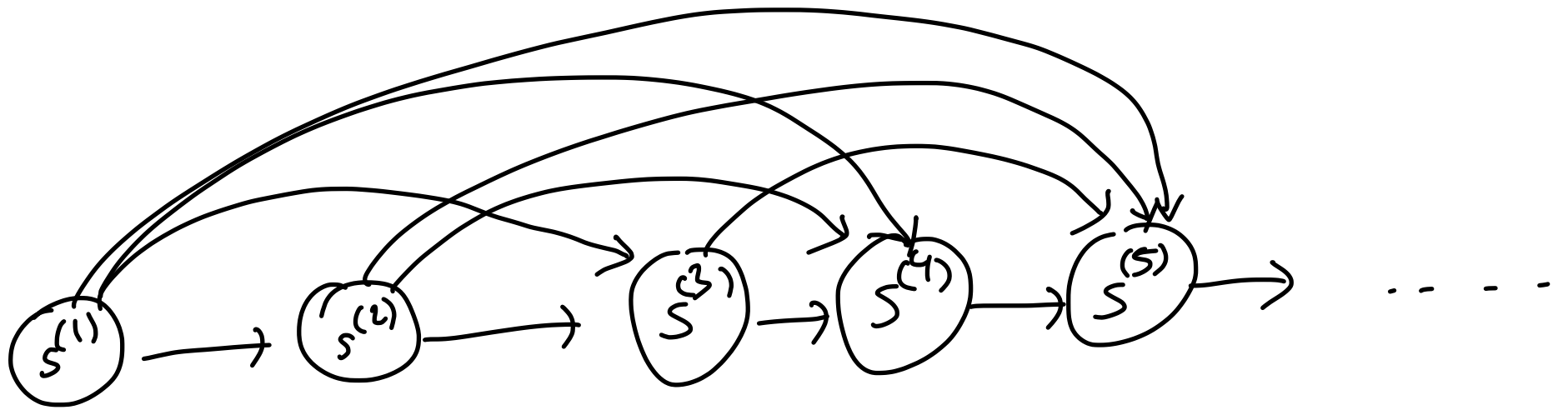
# Teacher Forcing



$$h^{(t)} = \sigma_1 \left[ \underbrace{W y^{(t-1)}}_{y^{(t-1)}} + U \cdot x^{(t)} + b_1 \right]$$



testing.

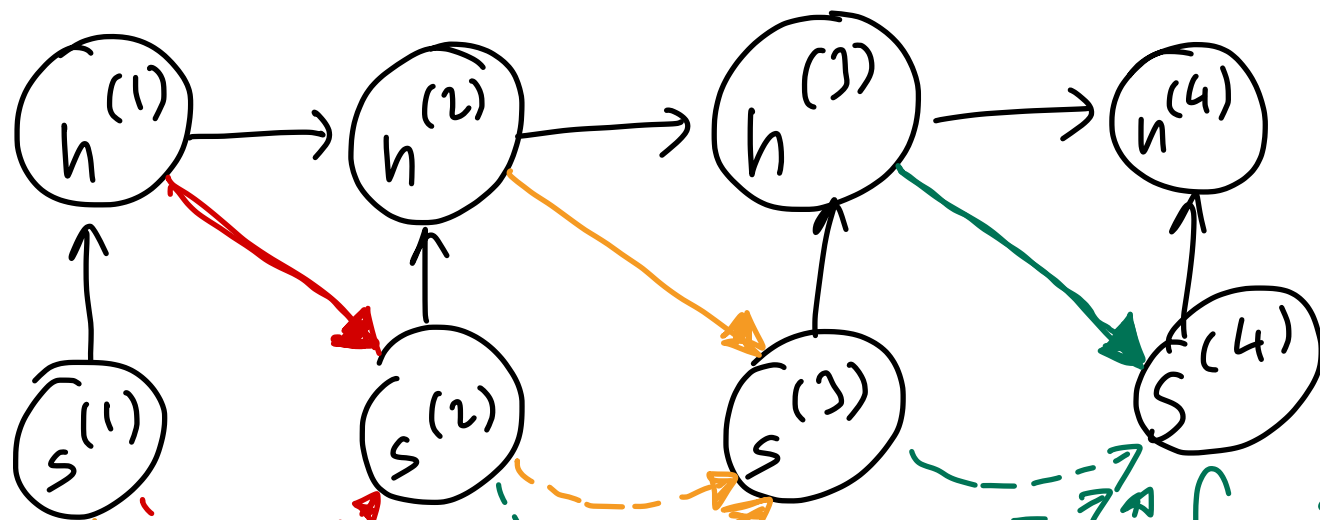


$$s^{(2)} = f(s^{(1)})$$

$$s^{(3)} = f(s^{(2)}, s^{(1)})$$

$$s^{(4)} = f(s^{(3)}, s^{(2)}, s^{(1)})$$

$$s^{(5)} = f(s^{(4)}, s^{(3)}, s^{(2)}, s^{(1)})$$



$h^{(1)}$  is a function of  $\underline{s^{(1)}}$ .  
 $h^{(2)}$  is a function of  $\underline{h^{(1)}}$  and  $\underline{s^{(2)}}$ .

$s^{(3)}$  is a function of  $h^{(2)}$ .

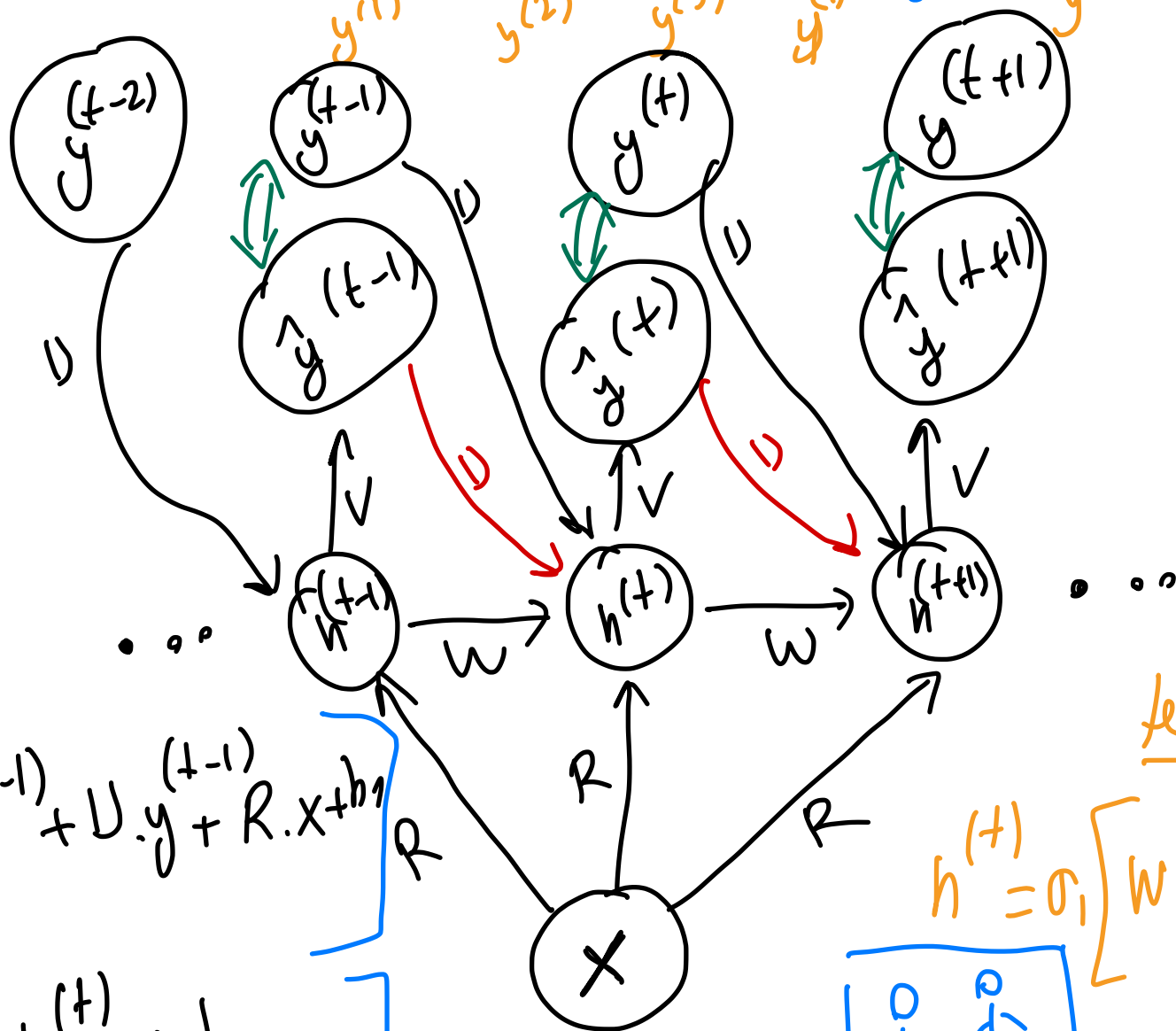
$s^{(4)}$  is a function of  $h^{(3)}$ .

$h^{(3)}$  is a function of  $\underline{h^{(2)}}$  and  $s^{(3)}$ .

# Vector to Sequence

image captioning.

Two kids are playing basketball



only for testing.

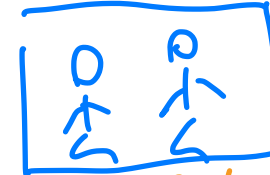
$$h^{(t)} = \sigma_1 \left[ W \cdot h^{(t-1)} + U \cdot y^{(t-1)} + R \cdot x + b_1 \right]$$

$$\hat{y}^{(t)} = \sigma_2 \left[ V \cdot h^{(t)} + b_2 \right]$$

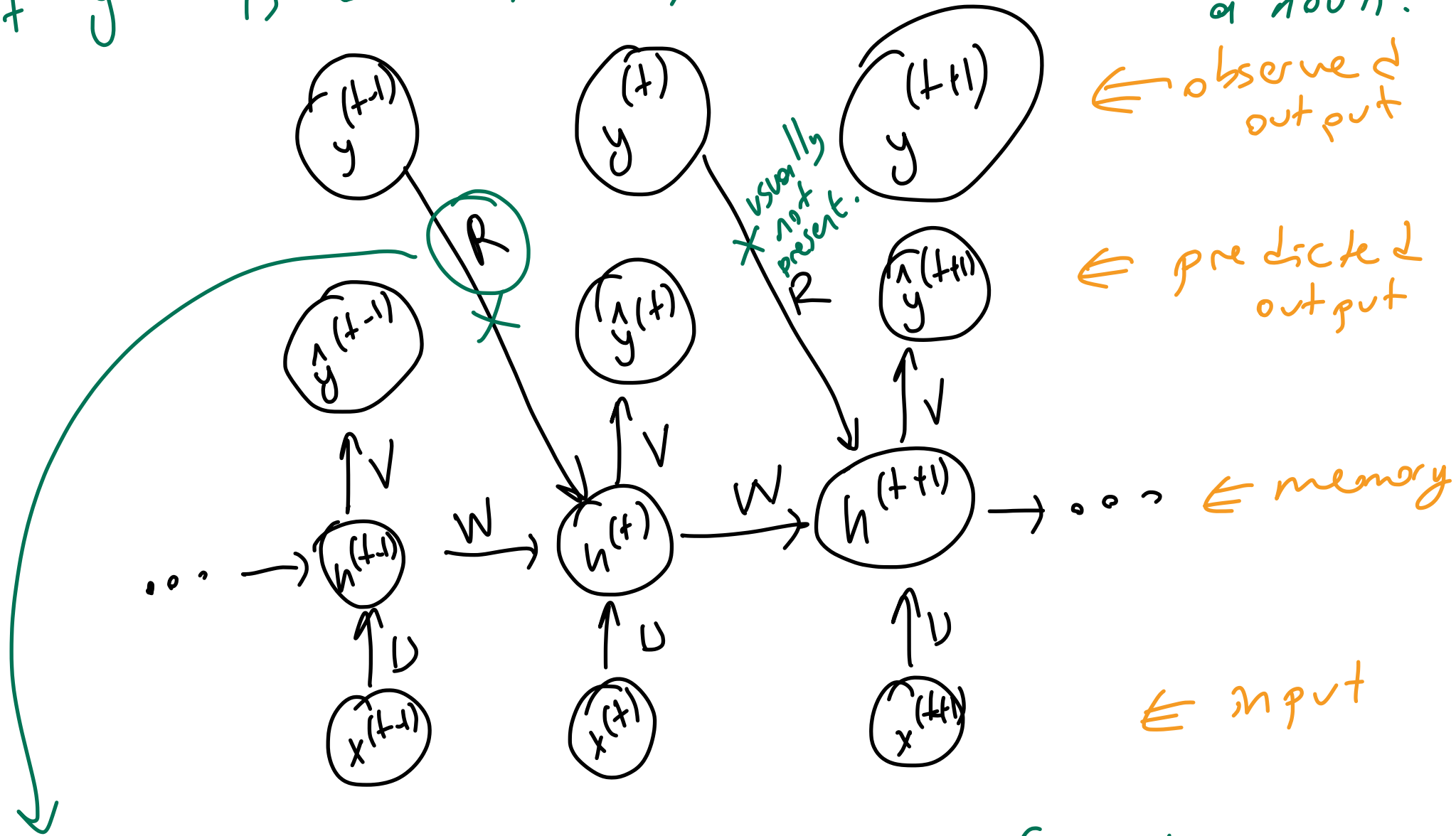
testing.

$$h^{(t)} = \sigma_1 \left[ W \cdot h^{(t-1)} + U \cdot \hat{y}^{(t-1)} + R \cdot x + b_1 \right]$$

$$\hat{y}^{(t)} = \sigma_2 \left[ V \cdot h^{(t)} + b_2 \right]$$

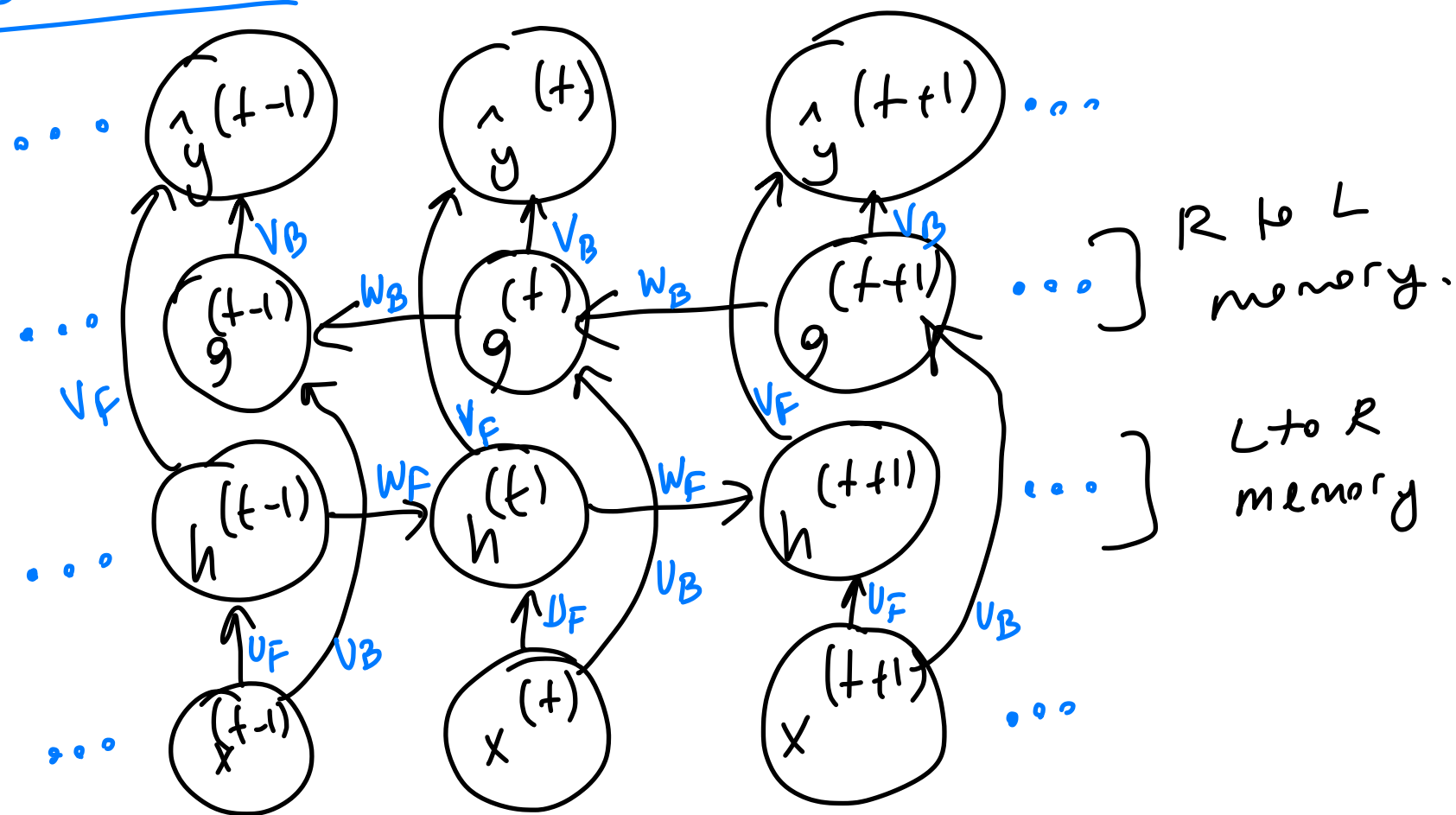


if  $y^{(t-1)}$  is an adjective, the model will force  $\hat{y}^{(t)}$  to be a noun.



∴ going to learn grammatical rules from training data.

# Bidirectional RNNs



$$h^{(t)} = \sigma_1 \left[ W_F \cdot h^{(t-1)} + V_F x^{(t)} + b_F \right]$$

$$g^{(t)} = \sigma_2 \left[ W_B \cdot g^{(t+1)} + V_B x^{(t)} + b_B \right]$$

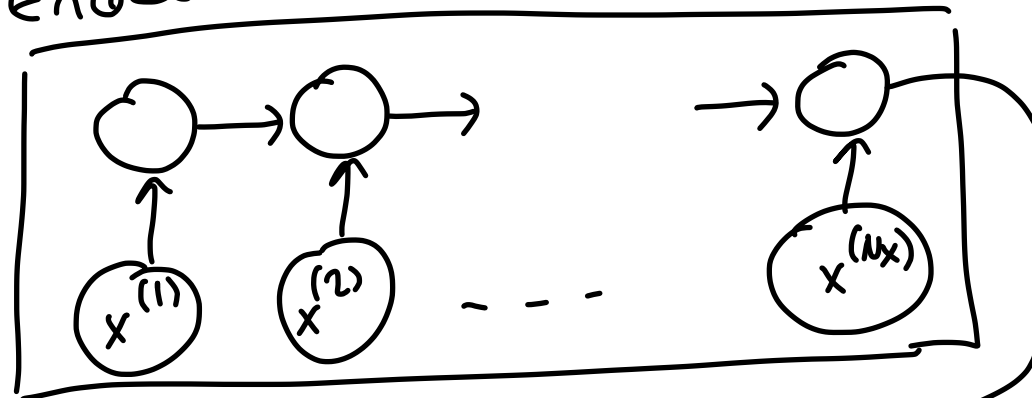
$$\hat{y}^{(t)} = \sigma_3 \left[ V_F \cdot h^{(t)} + V_B \cdot g^{(t)} + b \right]$$



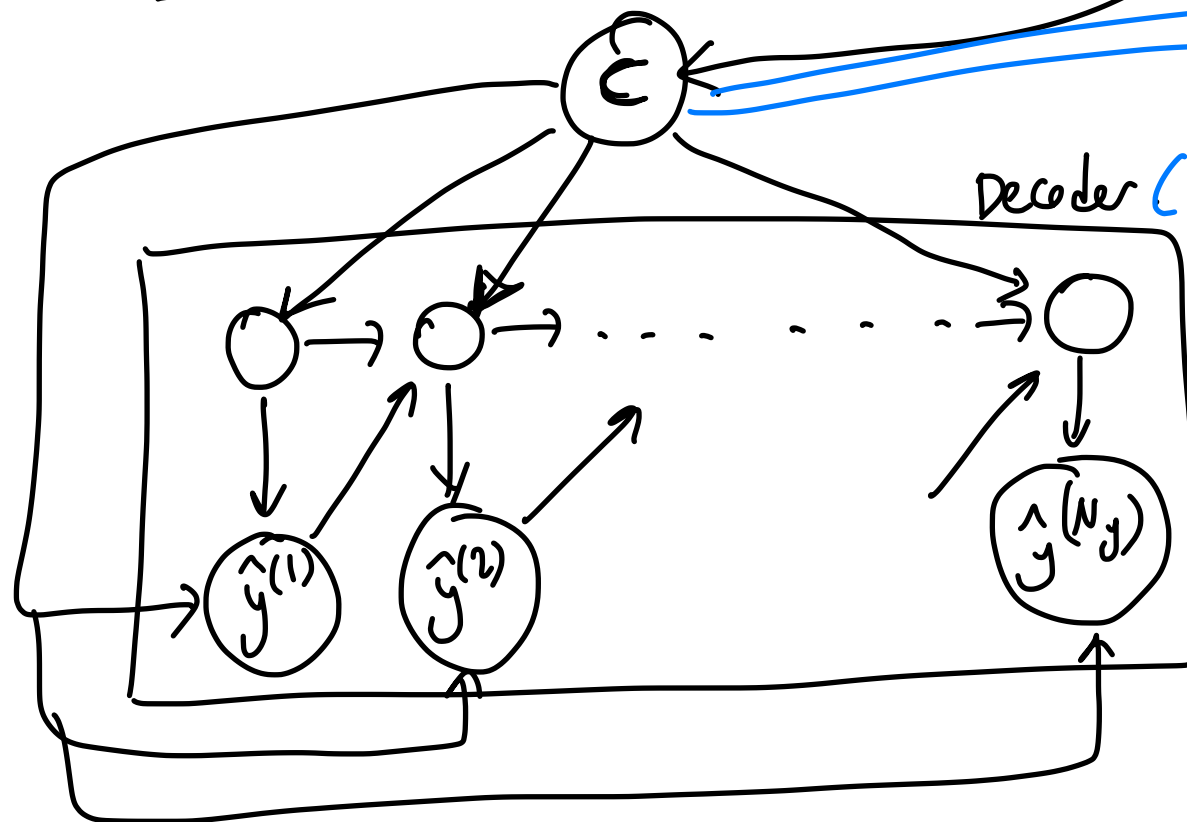
# (Sequence to Sequence) Models.

machine translation

Encoder.



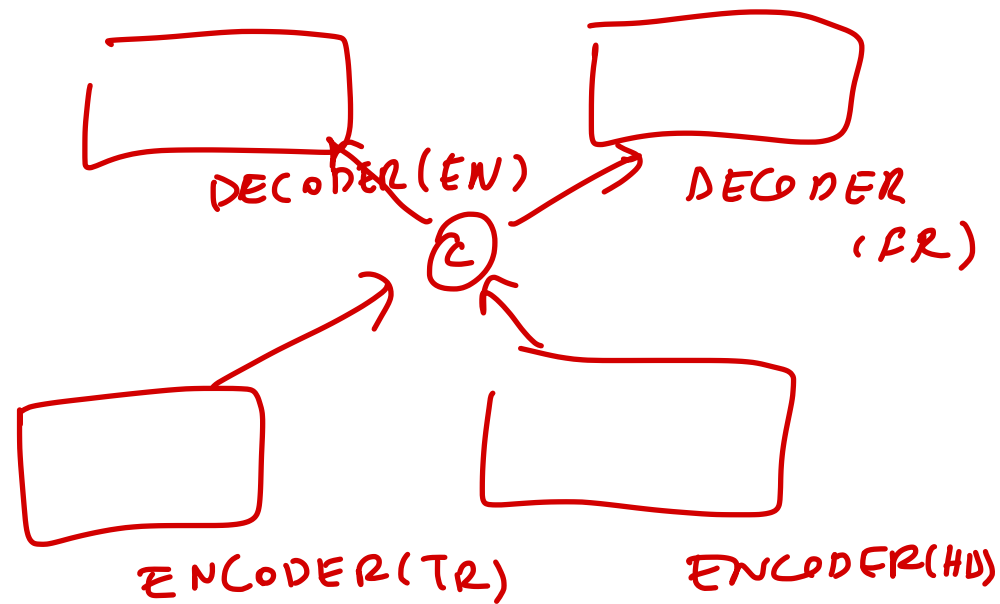
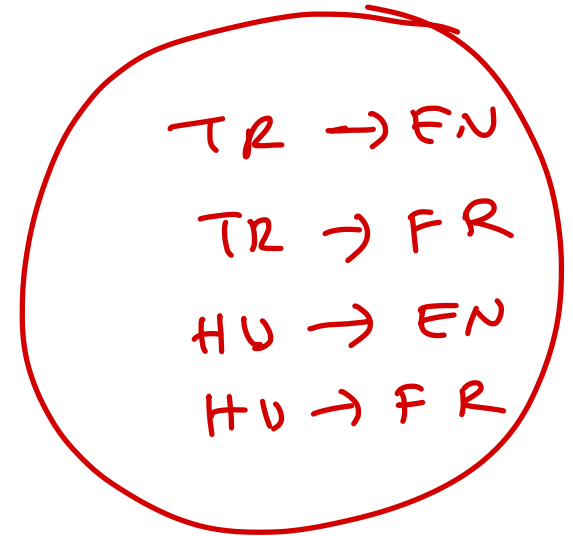
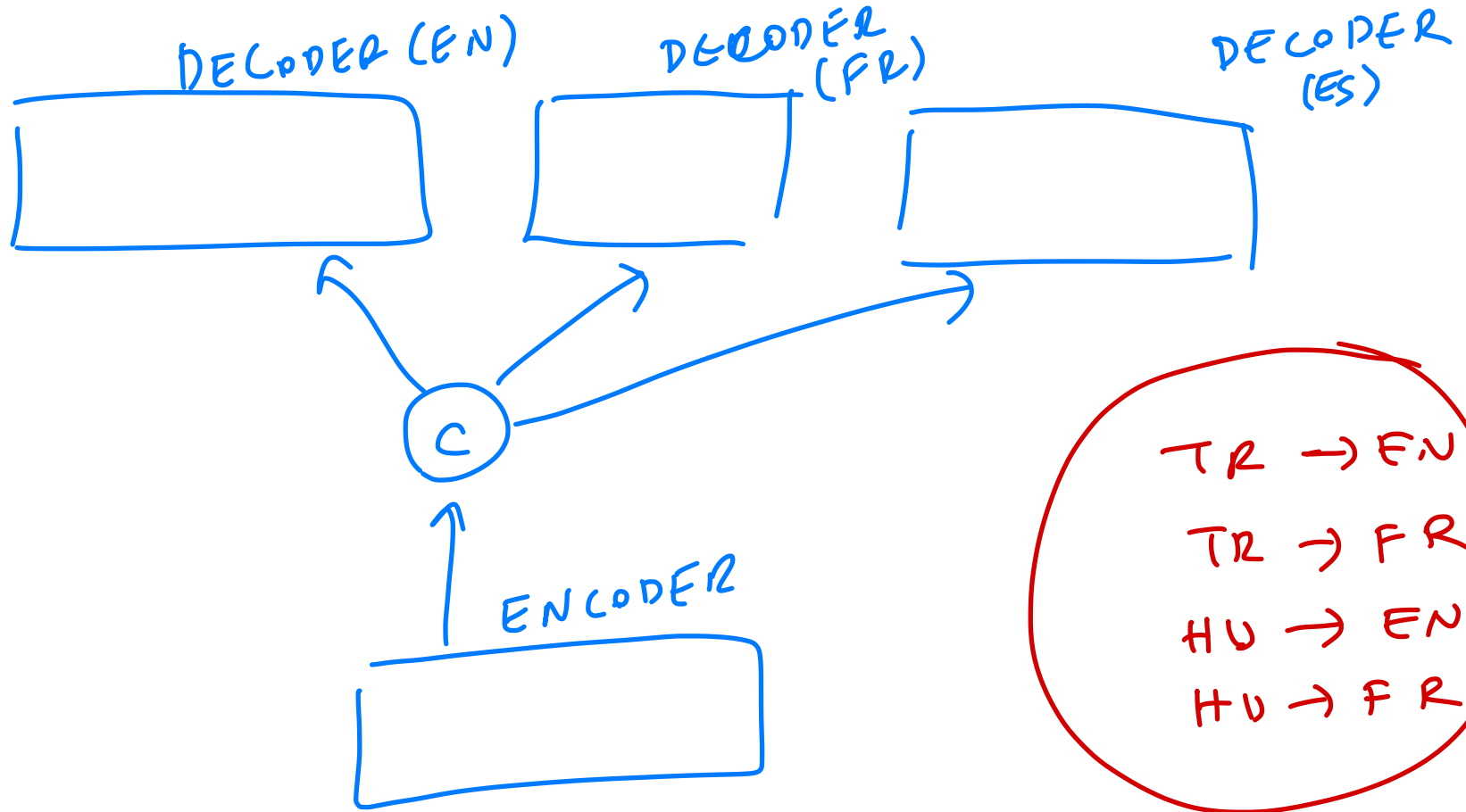
500000 TR  $\rightarrow$  EN  
500000 TR  $\rightarrow$  FR  
500000 TR  $\rightarrow$  ES



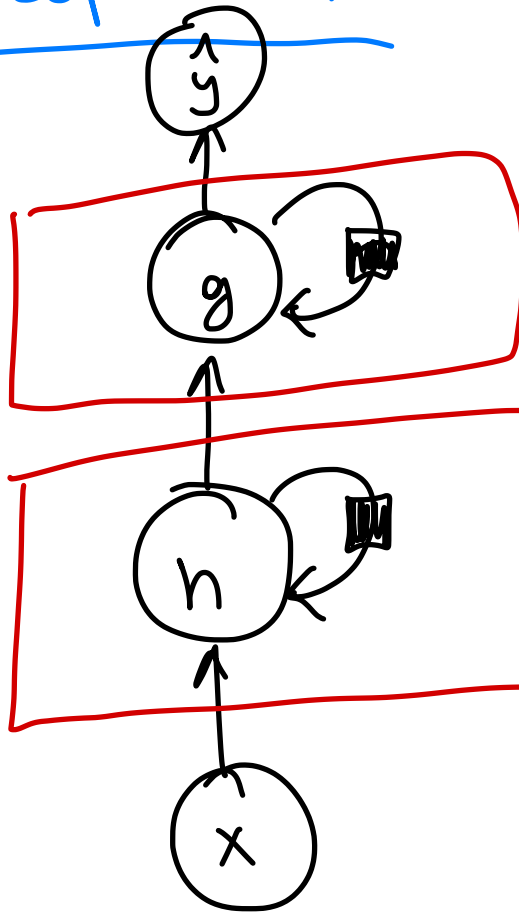
Decoder (EN)

DECODER(FR)

DECODER(ES)



# Deep RNNs



recurrent  
layer.

recurrent  
layer

↓  
{ GRU  
LSTM  
RNN ← ↗