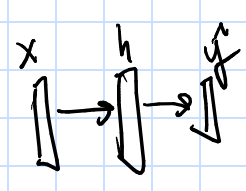


Our previous work, focused on training a feedforward neural net. That is given a point  $x \in \mathbb{R}^D$ , the sample was passed through the net to get  $\hat{y}$ . In this context, we are classifying a single point. However, what if we have a sequence?  $x_1, x_2, \dots, x_t$

- In some settings we have a history  $y$



### Example: Language Model

Consider the task of making a prediction on the next word given the current word  $w(t)$  and a history  $h(t-1)$ . To do this prediction task, we need to know some context. So given  $w(t)$  and  $h(t-1)$ , what is

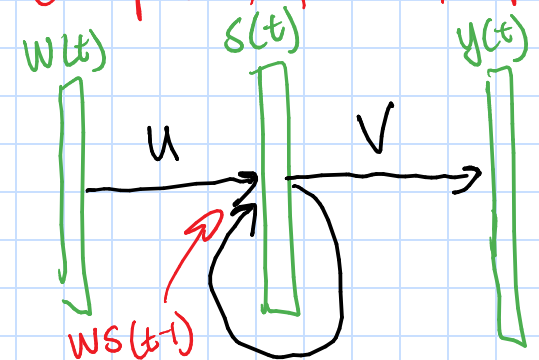
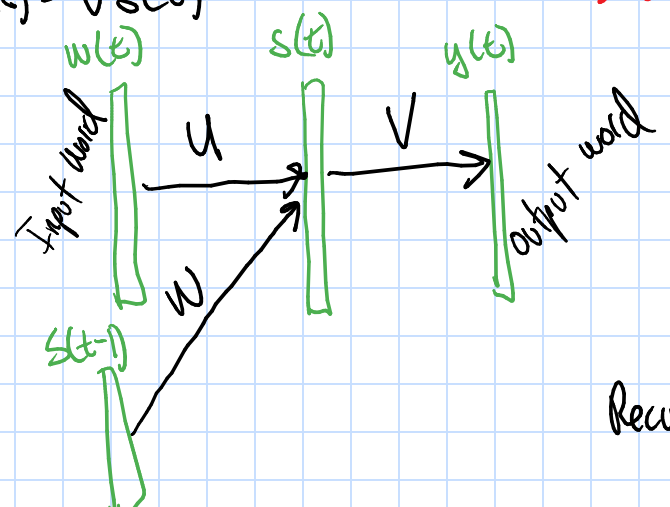
$$\Pr \{w(t+1) | w(t), h(t-1)\}$$

Can you pick up milk at the " " " " ?  
 $\downarrow$   $\downarrow$   
 $w(t)$   $w(t+1)$

$$s(t) = U w(t) + W s(t-1)$$

$$y(t) = V s(t)$$

$h(t-1)$   $\rightarrow$  or some representation of this part



$U \in \mathbb{R}^{10k \times 5k}$   
 $W \in \mathbb{R}^{5k \times 5k}$   
 $V \in \mathbb{R}^{5k \times 10k}$

Recurrent Neural Net  $W$

$D = 10k$   
 $H = 5k$

How do we encode  $w(t)$  if it is a word?

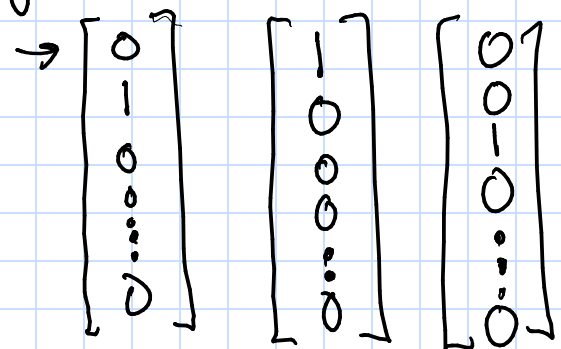
- "One hot encoding"

- Word embeddings (word2vec)

Dictionary has  $D$  words

Mikolov (2013)

BERT



one

not

encoding

$n$ -grams

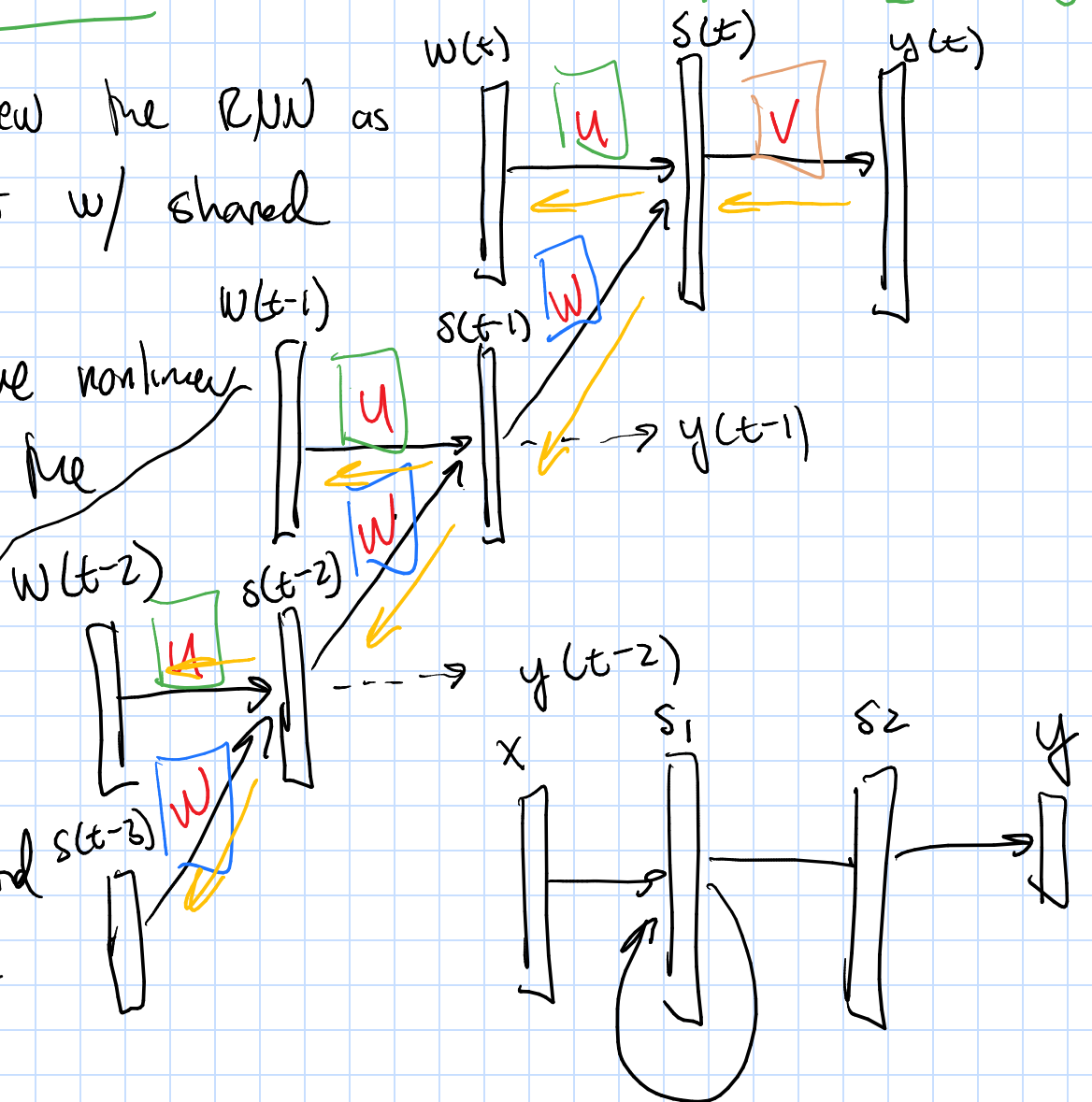
$\{a, g, c, t\}$

## Unraveling the RNN

- We can view the RNN as a deep net w/ shared weights

- We still have nonlinear functions @ the hidden and outputs

- Errors need to be propagated back through time



Def:

$$e_o(t) = d(t) - y(t) \quad \{\text{error}\}$$

$\alpha \rightarrow$  learning rate

$\beta \rightarrow$  L2 regularization

$$e_n(t) = d_n(e_o^T v)$$

$$d_{n,i}(x,t) = x \cdot s_j(t) (1 - s_j(t))$$

$$V(t+1) = V(t) + \alpha s(t) e_o(t)^T - \beta V(t)$$

$$U(t+1) = U(t) + \alpha W(t) \underline{e_n(t)}^T - \beta U(t)$$

$$W(t+1) = W(t) + \alpha s(t-1) e_n(t)^T - \beta W(t)$$

$$s(t) = f(U_w(t) + W_s(t-1))$$

$$y(t) = g(V_o(t))$$