

# RECURRENT NEURAL NETWORKS

---

- ① Let's reconsider the task of classifying <sup>DNA</sup> sequences. Suppose that a DNA sequence is in a particular family if it contains 3 or more "G" bases, e.g.

A G G T A A C G

G A T A G G G A T A

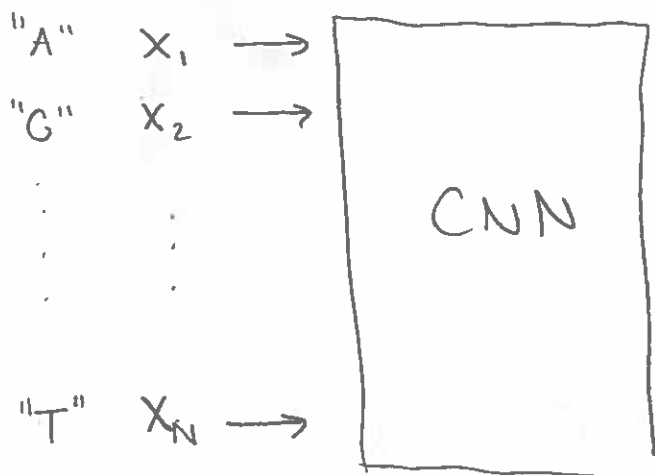
belong to the family, but:

T A T A C C A G A T T A

G A T T A C A G A T T A C A

do not.

- 
- ② One could try using a CNN, but maybe we don't know upfront the maximum length of a DNA sequence. That's problematic, because we need the input layer to be an upper bound on any sequence length



# RECURRENT NEURAL NETWORKS

- ③ Alternatively, we could try to mimic the process of identifying positive sequences programmatically:

```
def f(X: string)
    Q = [False, False, False]
    for x in X:
        if x == "G" and not Q[1]:
            Q[1] = True
        elif x == "G" and not Q[2]:
            Q[2] = True
        elif x == "G" and not Q[3]:
            Q[3] = True
    return Q[3]
```



this looks like Python,  
but let's count from  
1 for simplicity

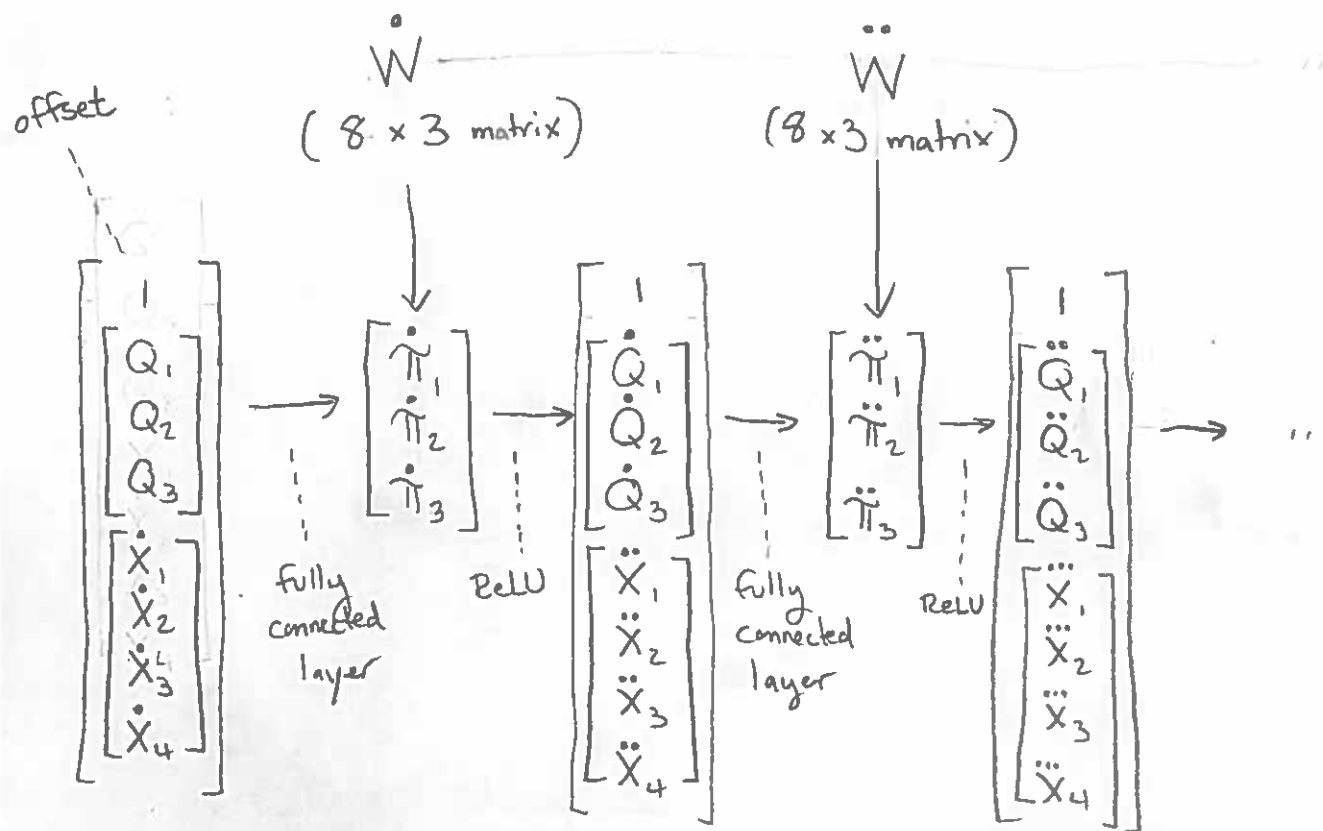
- ④ We could try encoding this as a neural network, in which rather than giving the entire input sequence as input to the first layer, we give the sequence one character at a time to each layer.

First, let's represent an input sequence A G G T G as one-hot vectors:

$\overset{\cdot}{X}$	$\overset{\cdot\cdot}{X}$	$\overset{\cdot\cdot\cdot}{X}$	$\overset{\cdot\cdot\cdot\cdot}{X}$	$\overset{\textcircled{5}}{X}$
"A"	"G"	"G"	"T"	"G"
$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

# RECURRENT NEURAL NETWORKS

⑤ The neural network representation of our Python function could look as follows:



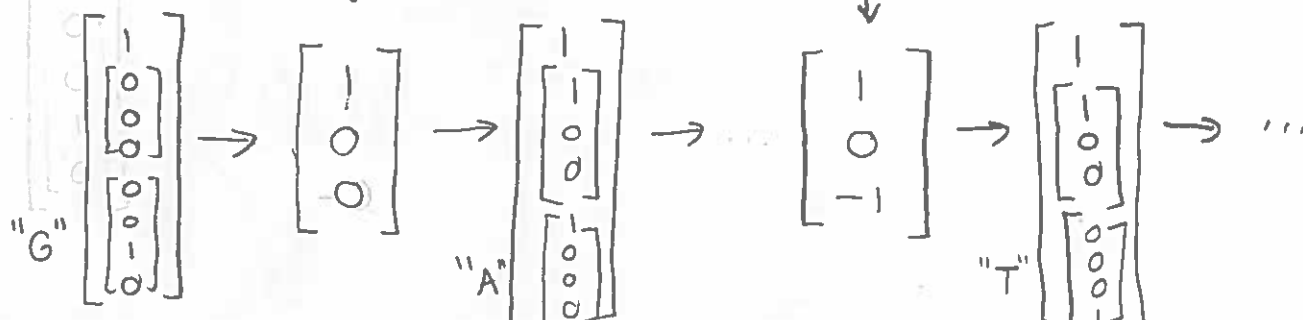
⑥ We can encode our logic (the if-statement) as the following

fully connected layer

$$\begin{bmatrix} 0 & -1 & -1 \\ 1 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & -1 & -1 \\ 1 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

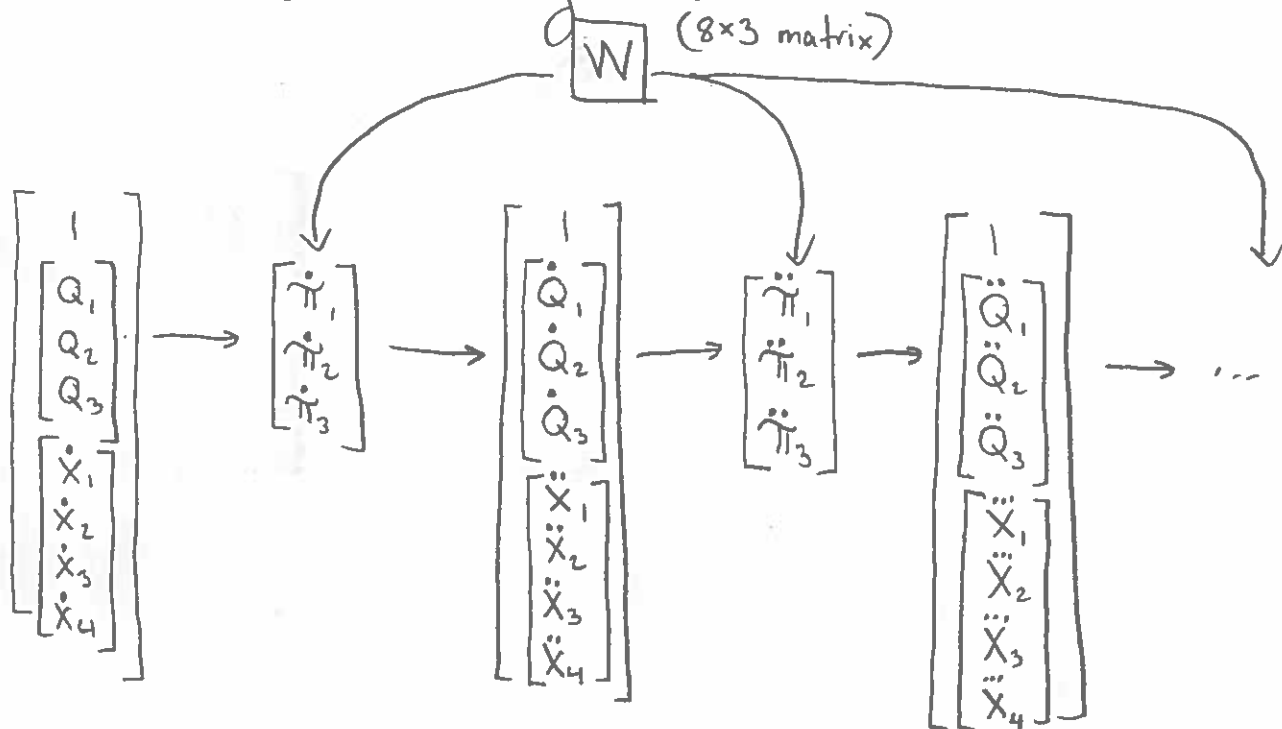
$\begin{bmatrix} Q_3 \text{ becomes positive} \\ \text{if } Q_2 \text{ was positive} \\ \text{and } X_3 = 1 \text{ (i.e. "G")} \end{bmatrix}$   
OR  
 $\begin{bmatrix} Q_3 \text{ becomes positive} \\ \text{if } Q_3 \text{ was already positive} \end{bmatrix}$



# RECURRENT NEURAL NETWORKS

⑦ But just as, in a for-loop, the logic doesn't change from iteration to iteration, there's no need to change our weight matrix (which encodes the logic) from iteration to iteration.

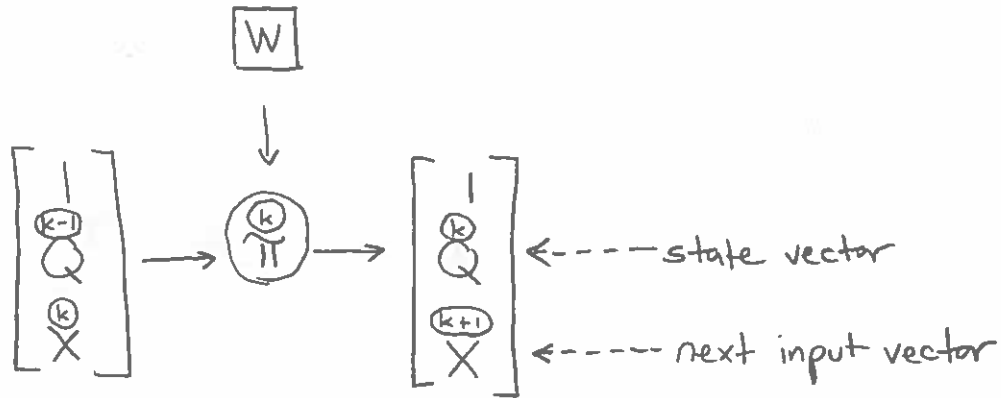
If we share these parameters between all layers, we obtain the following network:



This is called a recurrent neural network (RNN).

# RECURRENT NEURAL NETWORKS

③ Observe that you can specify an RNN by showing how a single layer looks:

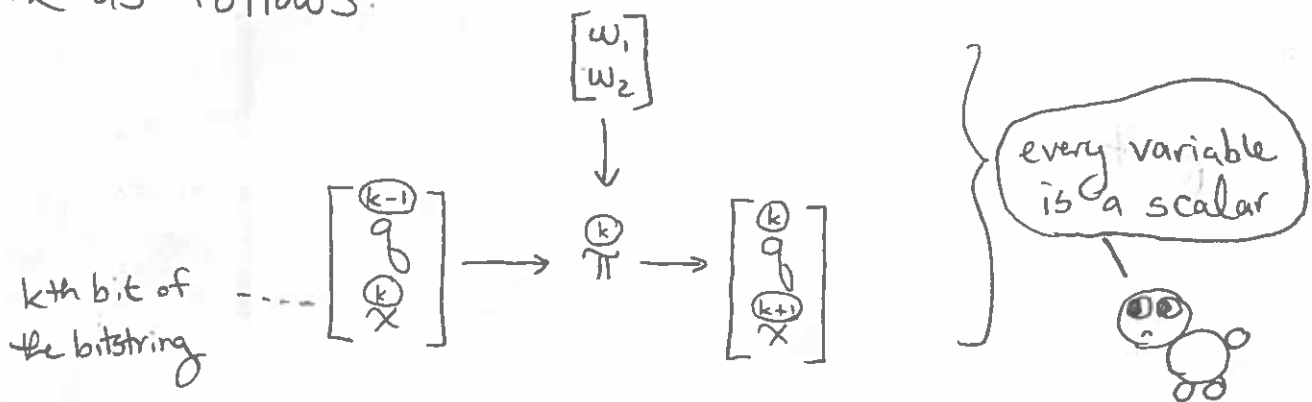


where  $\pi^{(k)} = W^T \begin{bmatrix} 1 \\ Q^{(k-1)} \\ X^{(k)} \end{bmatrix}$ ,  $Q^{(k)} = a(\pi^{(k)})$

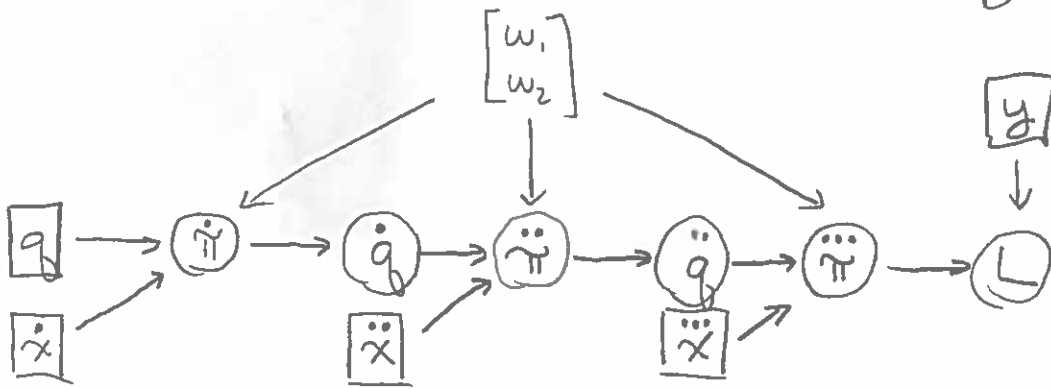
# RECURRENT NEURAL NETWORKS

⑨ Although we had weight sharing in CNNs, this is the first time we've encountered weight sharing between layers, so it's worth a look at how gradients are computed.

Consider an RNN for bitstrings, where the state vectors  $Q$  each have size 1. For simplicity, we won't use an offset. Thus, the RNN layers look as follows:



⑩ Suppose we run this RNN on a bitstring of length 3:

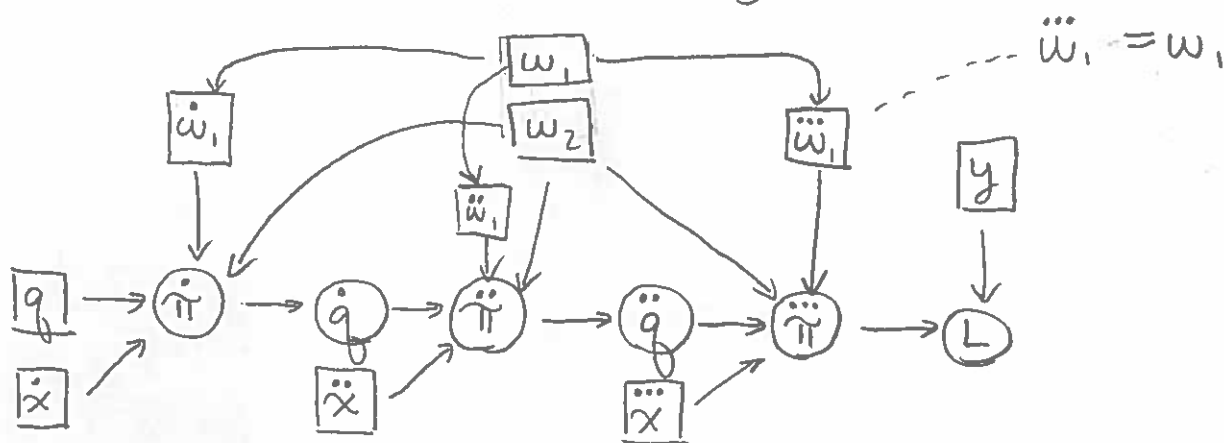


# RECURRENT NEURAL NETWORKS

⑪ Exercise: Compute  $\frac{\partial L}{\partial w_1}$  for datum  $(x, y) = ("101", 15)$

at the point in the weight space  $(w_1, w_2) = (4, 1)$   
if  $L(\tilde{\pi}, y) = (y - \tilde{\pi})^2$

It's more straightforward to apply the Chain Rule if we create copies of  $w_1$  for each layer, i.e.



⑫ Now, we can use  $\{\dot{w}_1, \ddot{w}_1, \ddot{\pi}\}$  as a separating set between  $w_1$  and  $\ddot{\pi}$ :

$$\begin{aligned}
 \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial \ddot{\pi}} \frac{\partial \ddot{\pi}}{\partial w_1} \\
 &= \frac{\partial L}{\partial \ddot{\pi}} \left( \frac{\partial \ddot{\pi}}{\partial \dot{w}_1} \cdot \frac{\partial \dot{w}_1}{\partial w_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} \cdot \frac{\partial \ddot{w}_1}{\partial w_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} \cdot \frac{\partial \ddot{\pi}}{\partial w_1} \right) \\
 &= \frac{\partial L}{\partial \ddot{\pi}} \left( \frac{\partial \ddot{\pi}}{\partial \dot{w}_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} \right) \\
 &= -2(y - \ddot{\pi}) \left( \frac{\partial \ddot{\pi}}{\partial \dot{w}_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} + \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} \right)
 \end{aligned}$$

all of these equal 1

## RECURRENT NEURAL NETWORKS

⑬ The remaining derivatives can be computed efficiently through backpropagation:

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} = \ddot{q}$$

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} = \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} \frac{\partial \ddot{\pi}}{\partial \ddot{w}_1} = \ddot{q} \cdot \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}}$$

$$\frac{\partial \ddot{\pi}}{\partial \dot{w}_1} = \frac{\partial \ddot{\pi}}{\partial \dot{\pi}} \frac{\partial \dot{\pi}}{\partial \dot{w}_1} = \dot{q} \cdot \frac{\partial \ddot{\pi}}{\partial \dot{\pi}}$$

where:

$$\frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} = \frac{\partial \ddot{\pi}}{\partial \ddot{q}} \frac{\partial \ddot{q}}{\partial \ddot{\pi}} = \ddot{w}_1 \cdot a'(\ddot{\pi}) = w_1 \cdot a'(\ddot{\pi})$$

$$\begin{aligned} \frac{\partial \ddot{\pi}}{\partial \dot{\pi}} &= \frac{\partial \ddot{\pi}}{\partial \dot{q}} \cdot \frac{\partial \dot{q}}{\partial \dot{\pi}} \\ &= \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}} \cdot \frac{\partial \ddot{\pi}}{\partial \dot{q}} \cdot \frac{\partial \dot{q}}{\partial \dot{\pi}} = \left( w_1 \cdot a'(\ddot{\pi}) \right) \ddot{w}_1 \cdot a'(\dot{\pi}) = w_1^2 \cdot a'(\ddot{\pi}) \cdot a'(\dot{\pi}) \end{aligned}$$

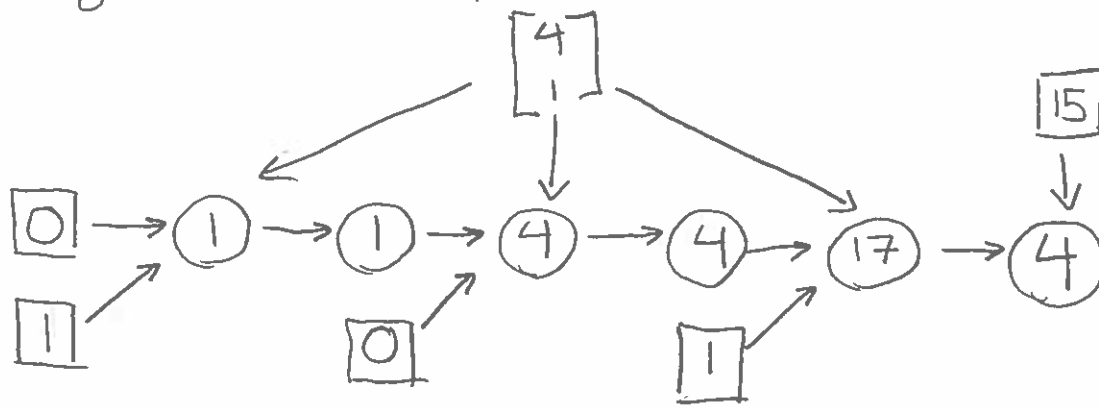
⑭ So then:

$$\frac{\partial L}{\partial w_1} = -2(y - \ddot{\pi}) \left( \ddot{q} + \dot{q} \cdot w_1 \cdot a'(\ddot{\pi}) + q \cdot w_1^2 \cdot a'(\ddot{\pi}) \cdot a'(\dot{\pi}) \right)$$



# RECURRENT NEURAL NETWORKS

⑮ Doing the forward pass on the specific data, we get:



⑯ So our partial derivative evaluates to:

$$\begin{aligned}\frac{\partial L}{\partial w_1} &= -2(15-17) \left( 4 + 1 \cdot 4 \cdot 1 + 0 \cdot 4^2 \cdot 1 \cdot 1 \right) \\ &= -2(-2)(8) \\ &= 32\end{aligned}$$