

# LONG SHORT-TERM MEMORY NETWORKS

① You can view an RNN as a fancy way of representing a for-loop. For instance, we previously represented a loop that determines whether there are at least 3 Gs in a DNA sequence:

for  $x$  in  $X$ :

if  $x == "G"$  and not  $Q[1]$ :

$Q[1] = \text{True}$

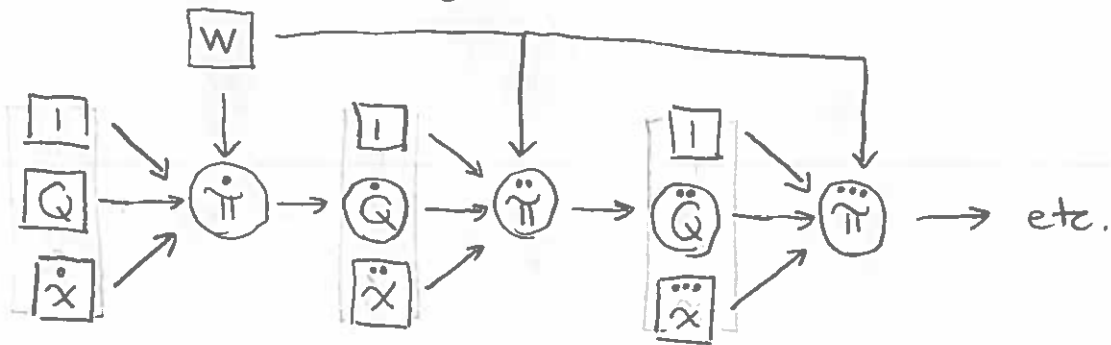
elif  $x == "G"$  and not  $Q[2]$ :

$Q[2] = \text{True}$

elif  $x == "G"$  and not  $Q[3]$ :

$Q[3] = \text{True}$

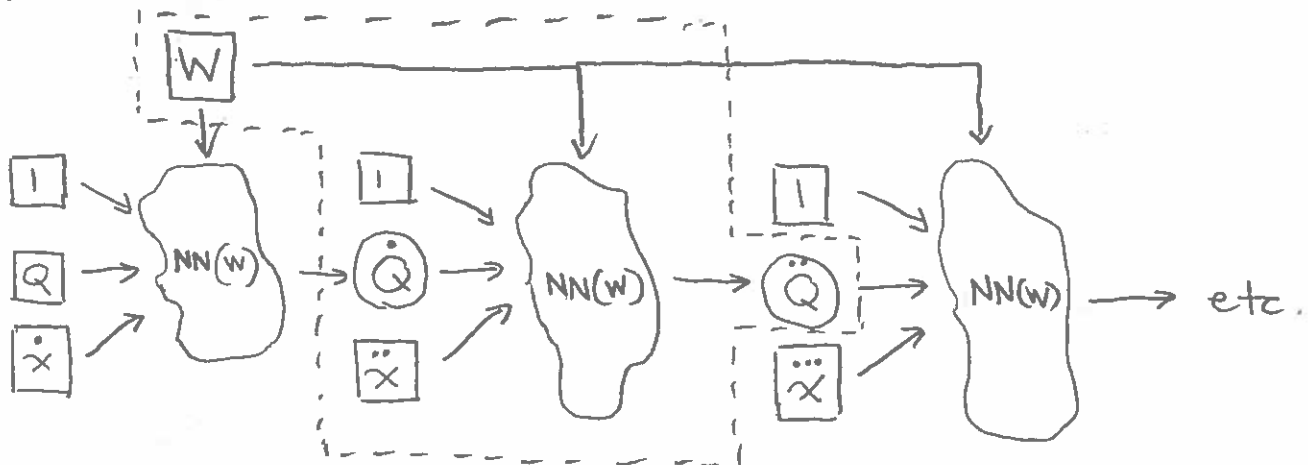
② We used the following neural net:



where:  $\pi^{(k)} = W^T \begin{bmatrix} 1 \\ Q^{(k)} \\ x^{(k)} \end{bmatrix}$  and  $Q^{(k)} = a(\pi^{(k)})$

# LONG SHORT-TERM MEMORY NETWORKS

- ③ For each iteration, we use a single fully connected layer, followed by an activation function. But there's nothing stopping us from using a more complicated neural network to compute the next state  $Q^{(k)}$  at each iteration:

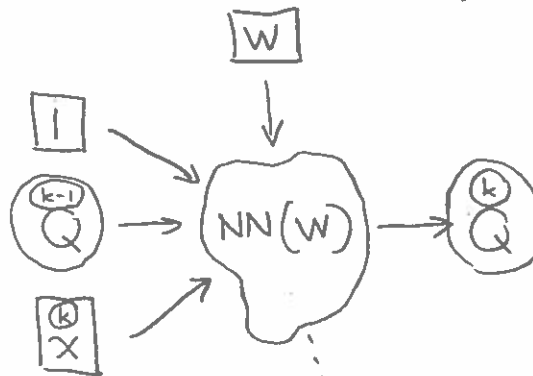


by defining different  $NN(W)$ , we can make different, more powerful RNNs.

- ④ Notice that we can structure  $NN(W)$  to enable more complex logic in each iteration of the loop:

for  $x$  in  $X$ :

complex logic goes here



complex logic goes here

## LONG SHORT-TERM MEMORY NETWORKS

---

- ⑤ Let's try encoding a somewhat more complicated function. Suppose we want to determine whether a bitstring contains at least four sequences of three consecutive ones, e.g.

0 1 1 1 1 0 1 1 1 0 1 1 1 ✓  
1 2 3 4

Note that overlapping sequences still count.

---

- ⑥ One way to write this as a function:

```
def f(X: bitstring):
```

```
    count = 0
```

```
    streak = 0
```

```
    for ch in X:
```

```
        if ch == 0:
```

```
            streak = 0
```

```
        else:
```

```
            streak += 1
```

```
            if streak >= 3:
```

```
                count += 1
```

```
    return count >= 4
```

## LONG SHORT-TERM MEMORY NETWORKS

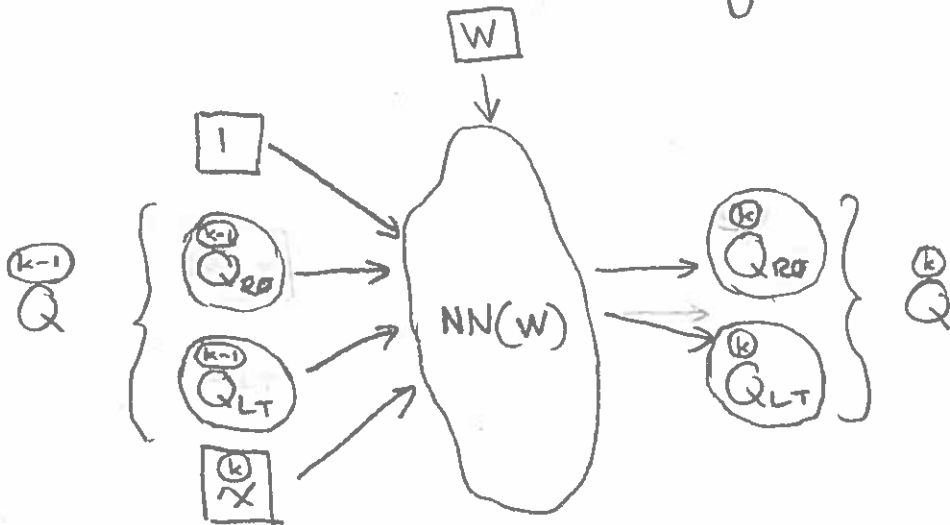
⑦ For our state variables  $Q$ , we'll use our variables count and streak.

$$Q = \begin{bmatrix} \text{count} \\ \text{streak} \end{bmatrix}$$

But we'll actually make two copies of each: a "long-term" version and a "read-only" version:

$$Q = \begin{bmatrix} \text{count}_{ro} \\ \text{streak}_{ro} \\ \text{count}_{LT} \\ \text{streak}_{LT} \end{bmatrix} \quad \left( \begin{array}{l} \text{or } Q_{ro} = \begin{bmatrix} \text{count}_{ro} \\ \text{streak}_{ro} \end{bmatrix} \\ Q_{LT} = \begin{bmatrix} \text{count}_{LT} \\ \text{streak}_{LT} \end{bmatrix} \end{array} \right)$$

So we'll be trying to find a good neural network architecture of the following form:



## LONG SHORT-TERM MEMORY NETWORKS

- ⑧ Let's trace out the desired functionality of  $NN(W)$  in terms of the long-term state variables. Suppose our bitstring is 0111011110 and we're starting the fifth iteration, i.e.

0111011110  
           $\wedge$

At this point, hopefully  $Q_{LT} = \begin{bmatrix} \text{count}_{LT} \\ \text{streak}_{LT} \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$

- ⑨ We'd like to go through the following steps:

$$\begin{array}{l} \text{count}_{LT} \\ \text{streak}_{LT} \end{array} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \xrightarrow[\text{(if applicable)}]{\begin{array}{l} \text{reset} \\ \text{streak} \end{array}} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \xrightarrow[\text{(if applicable)}]{\begin{array}{l} \text{increment} \\ \text{streak and} \\ \text{count} \end{array}} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

- ⑩ Similarly, for the next iteration:

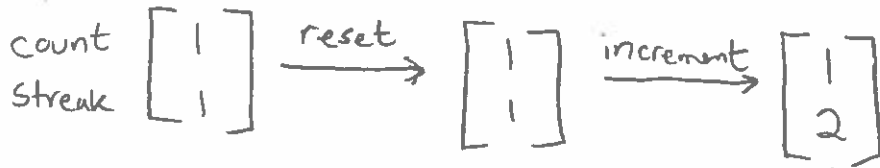
0111011110  
           $\wedge$

$$\begin{array}{l} \text{count} \\ \text{streak} \end{array} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \xrightarrow{\text{reset}} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \xrightarrow{\text{increment}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

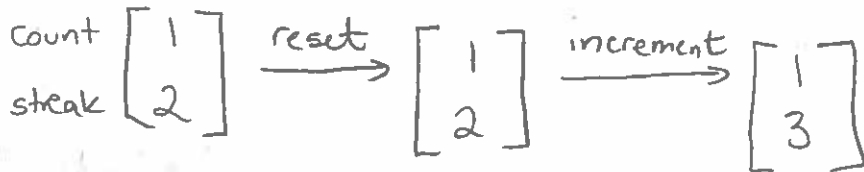
## LONG SHORT-TERM MEMORY NETWORKS

② Going through the next few iterations:

0 1 1 1 0 1 1 1 1 0  
                  ^

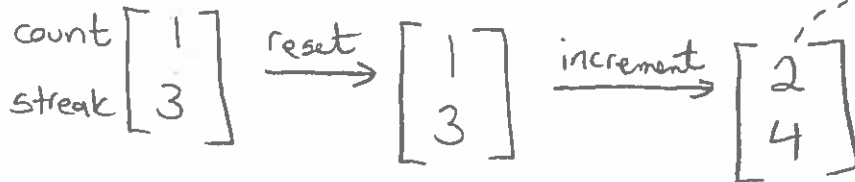


0 1 1 1 0 1 1 1 1 0  
                    ^



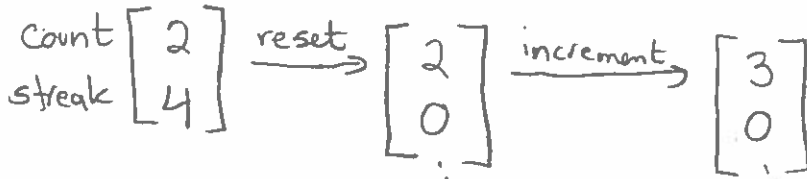
- streak gets incremented to 2 because we encountered another one in the bitstring

0 1 1 1 0 1 1 1 1 0  
                    ^



count gets incremented to 2 because the streak was  $\geq 3$  at the iteration's start

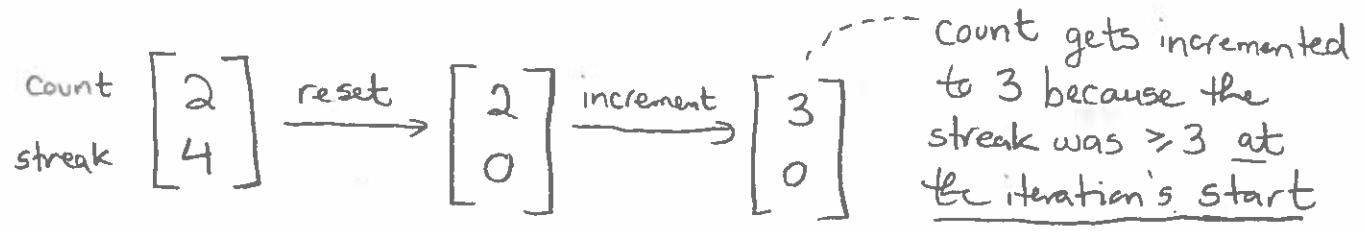
0 1 1 1 0 1 1 1 0  
                    ^



streak gets reset to 0  
because we encountered  
a zero in the bitstring

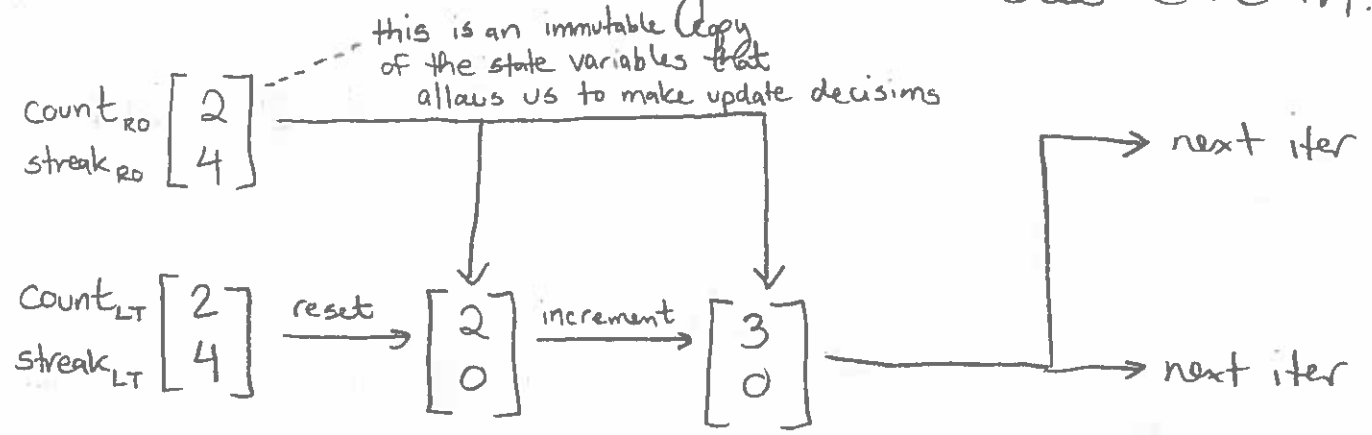
# LONG SHORT-TERM MEMORY NETWORKS

⑫ Notice particularly the last iteration:

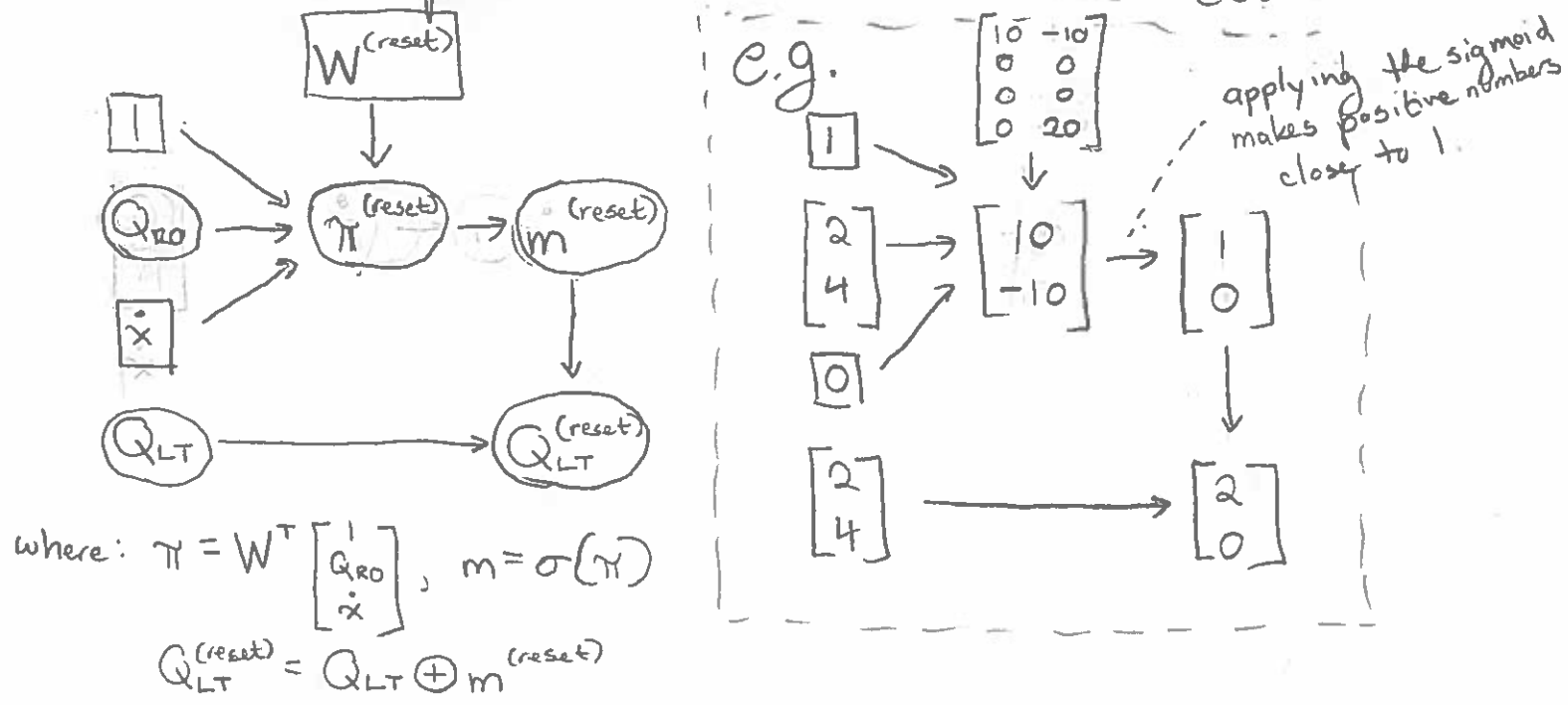


But we wiped out the streak during the reset, so how do we implement the increment of count?

⑬ This is where the read-only state variables come in:

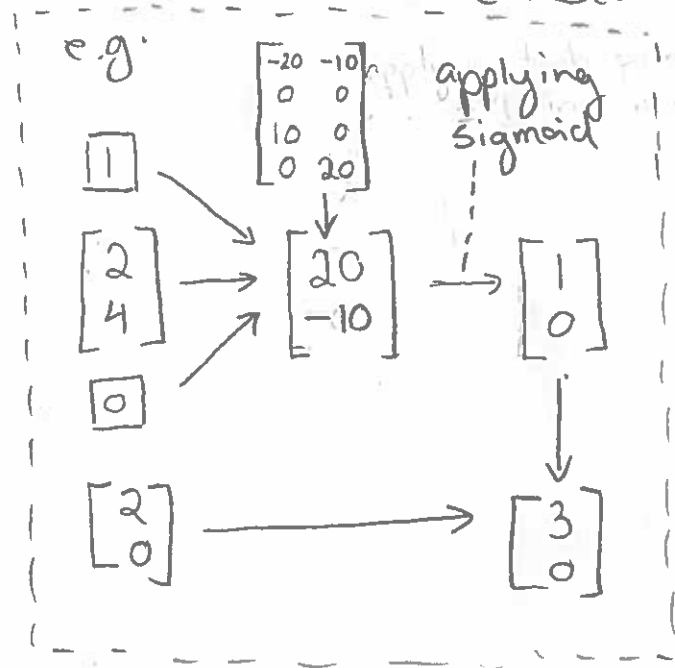
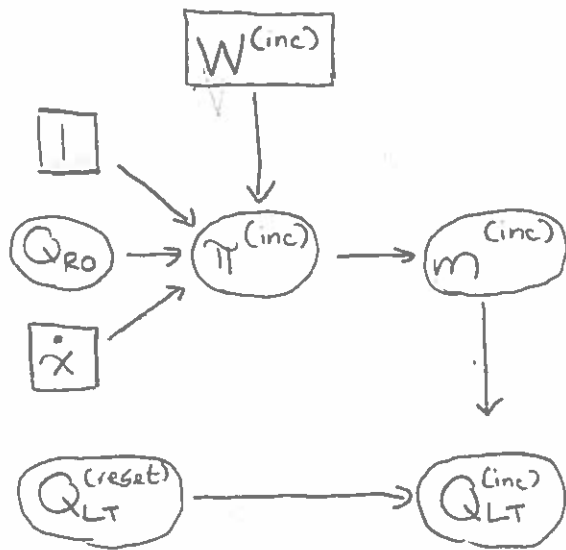


⑭ How do we implement "reset" as a neural network?



# LONG SHORT-TERM MEMORY NETWORKS

15) How do we implement "increment" as a neural network?



where:  $\pi = W^T \begin{bmatrix} 1 \\ Q_{RO} \\ \dot{x} \end{bmatrix}$

$m = \sigma(\pi)$

$Q_{LT}^{(inc)} = Q_{LT}^{(reset)} + m^{(inc)}$

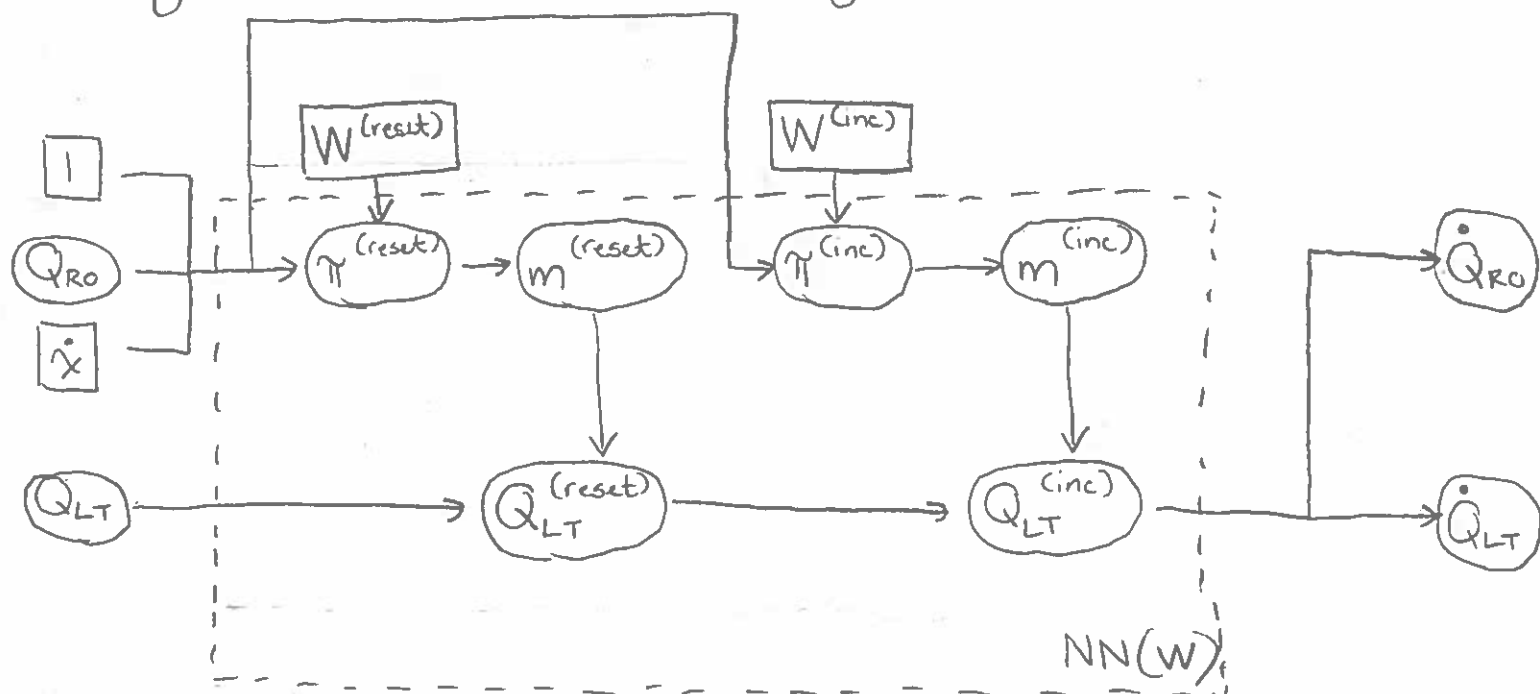
elementwise "plus", not Hadamard product





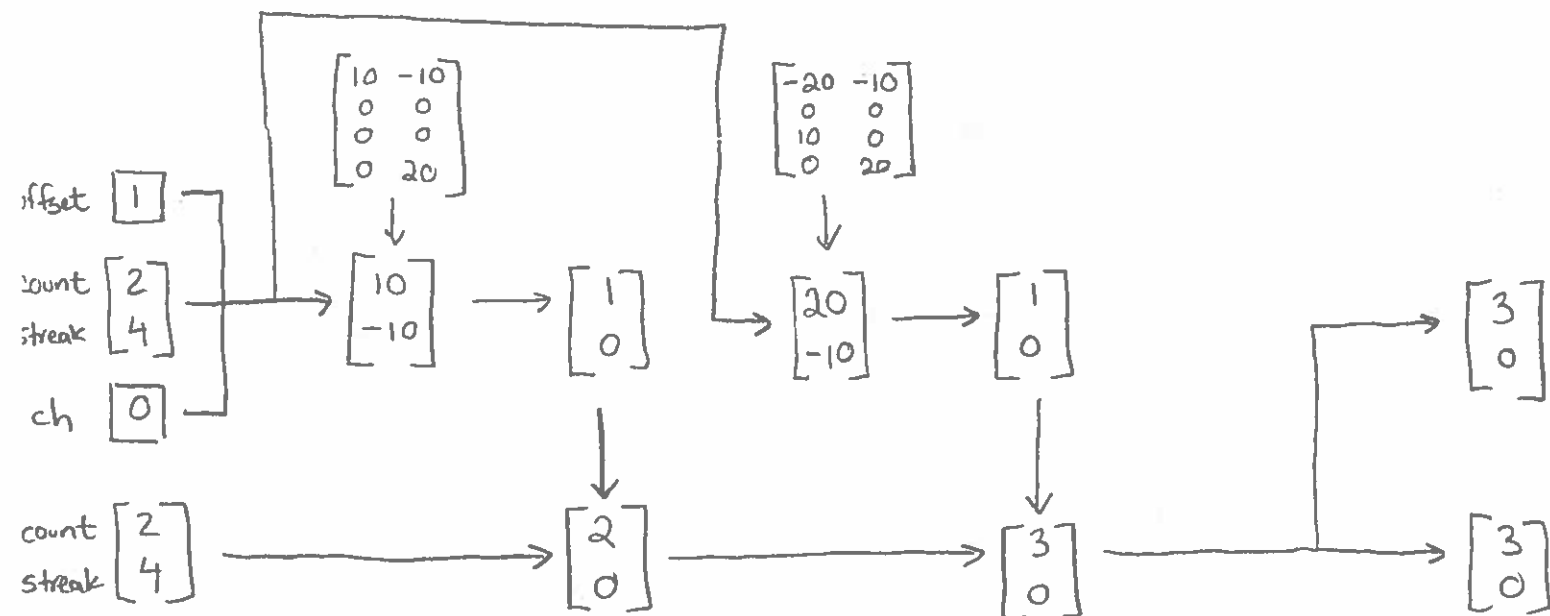
# LONG SHORT-TERM MEMORY NETWORKS

⑩ Putting "reset" and "increment" together:



(Compare this to ⑦)

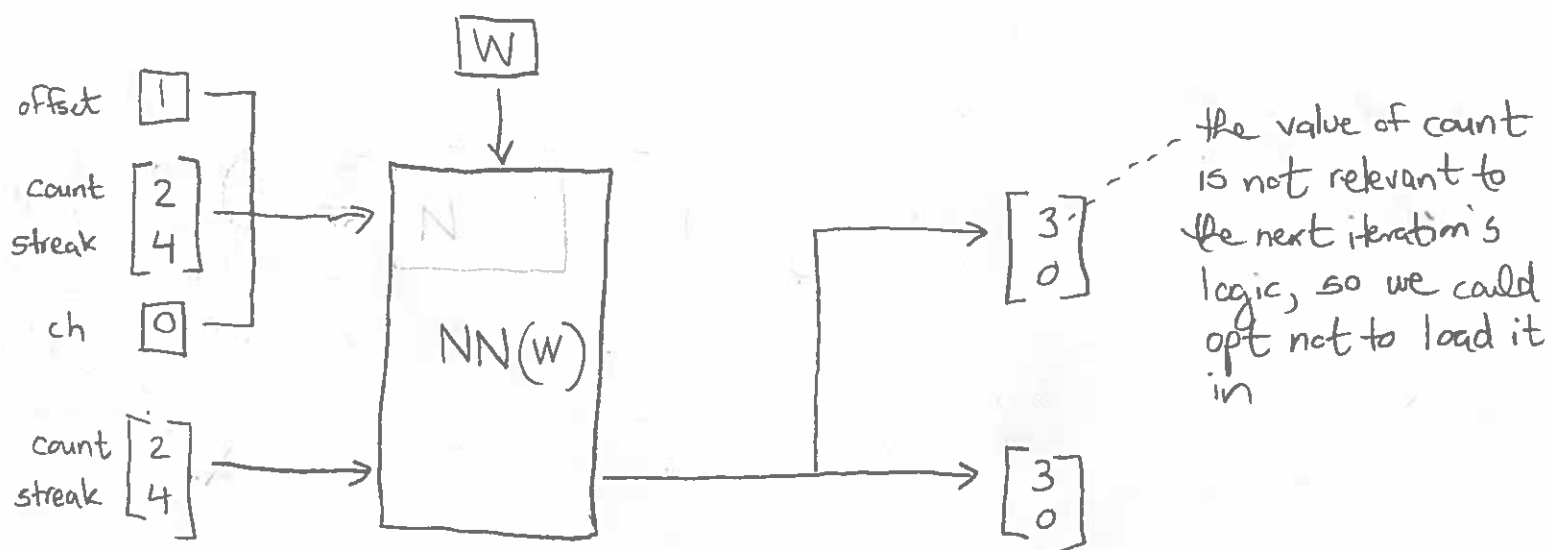
⑪ For example:



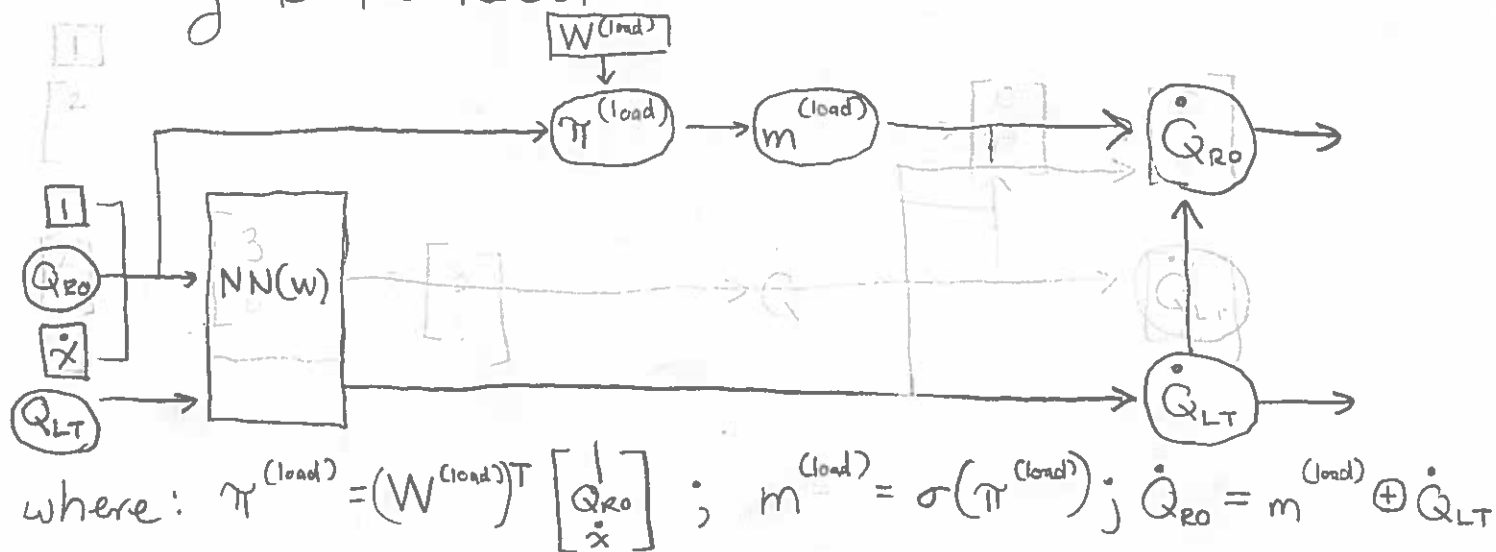
# LONG SHORT-TERM MEMORY NETWORKS

- ⑮ The RNN depicted in ⑮ is close to something called a Long Short-Term Memory network (LSTM). Full LSTMs have a couple additional extensions.

First, the LSTM does not have to load everything from the long-term memory back into the read-only buffer for the next iteration.

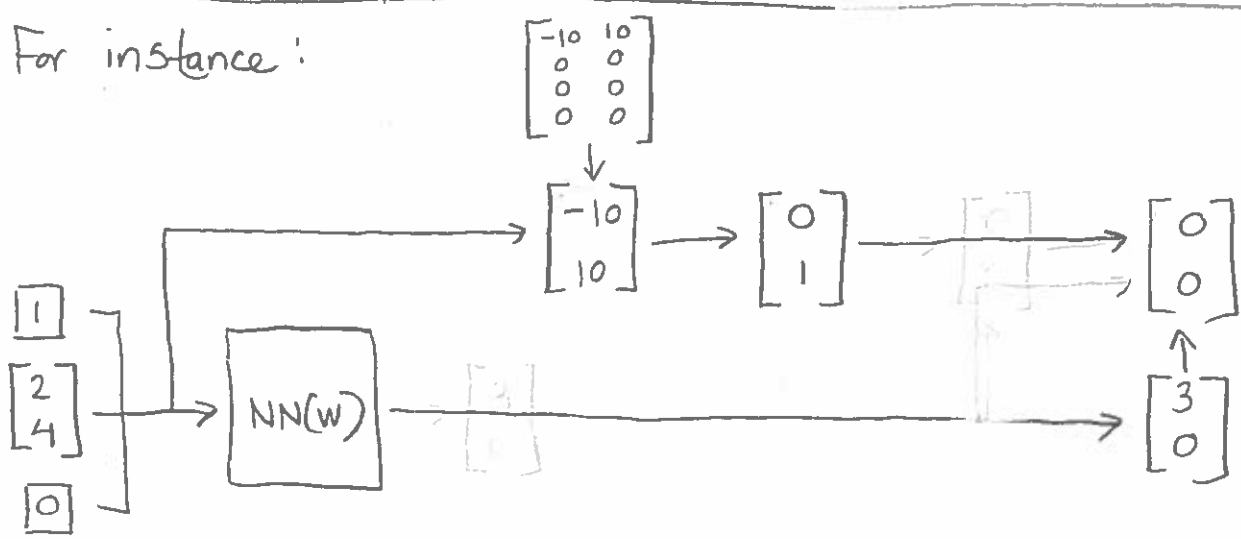


- ⑯ Zeroing out irrelevant state variables can be achieved the same way as the reset.



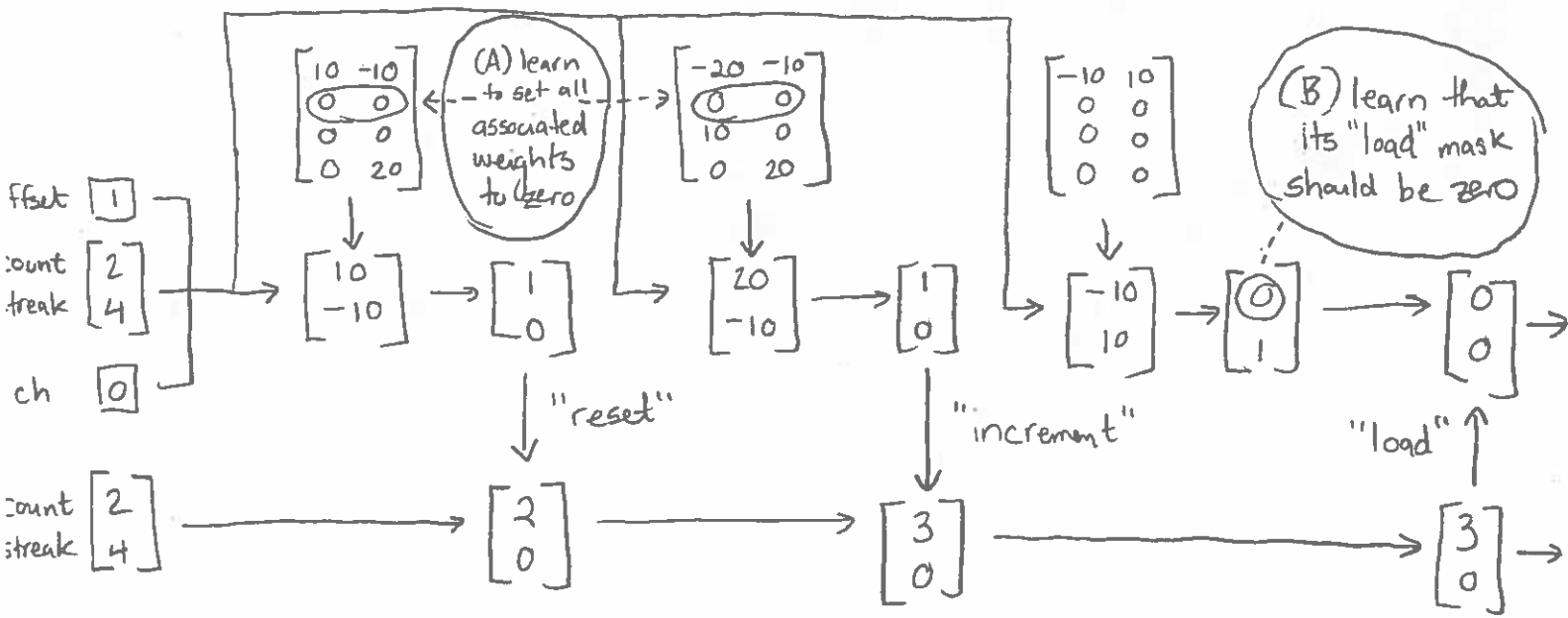
# LONG SHORT-TERM NETWORKS

20) For instance:



21) Why load selectively? It certainly doesn't harm things (from a representational perspective) to load the entire long-term memory into the read-only buffer every iteration.

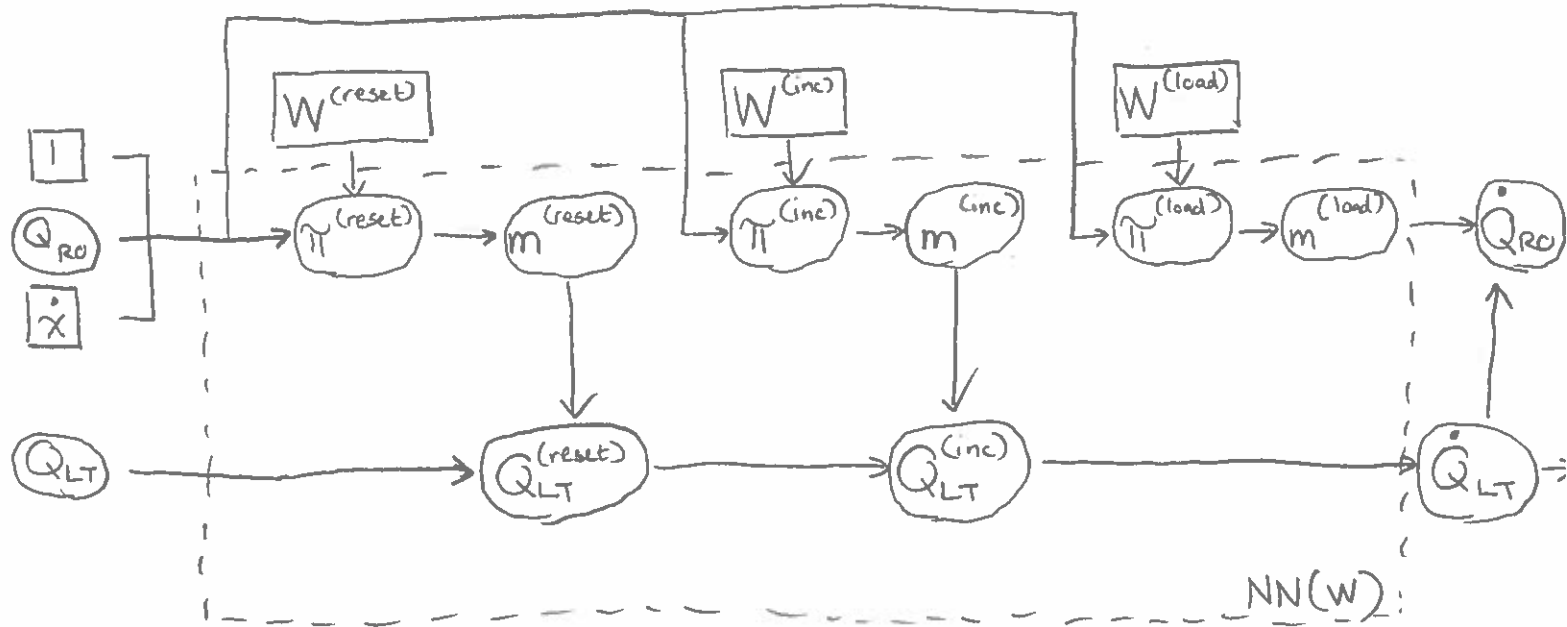
However, it may make it easier to learn the function. Consider <sup>the</sup> count variable, which is never needed in the read-only buffer. In order to ignore it, you can either:



Option B often requires less parameters to learn.

# LONG SHORT-TERM MEMORY NETWORKS

② Adding the "load" to our RNN gives us:



where:

$$\pi^{(reset)} = (W^{(reset)})^T \begin{bmatrix} 1 \\ Q_{RO} \\ \dot{x} \end{bmatrix}$$

$$m^{(reset)} = \sigma(\pi^{(reset)})$$

$$Q_{LT}^{(reset)} = m^{(reset)} \oplus Q_{LT}$$

$$\pi^{(inc)} = (W^{(inc)})^T \begin{bmatrix} 1 \\ Q_{RO} \\ \dot{x} \end{bmatrix}$$

$$m^{(inc)} = \sigma(\pi^{(inc)})$$

$$Q_{LT}^{(inc)} = m^{(inc)} + Q_{LT}^{(reset)}$$

$$\dot{Q}_{LT} = Q_{LT}^{(inc)}$$

$$\pi^{(load)} = (W^{(load)})^T \begin{bmatrix} 1 \\ Q_{RO} \\ \dot{x} \end{bmatrix}$$

$$m^{(load)} = \sigma(\pi^{(load)})$$

$$\dot{Q}_{RO} = m^{(load)} \oplus \dot{Q}_{LT}$$

not  $\oplus$ !

## LONG SHORT-TERM MEMORY NETWORKS

- 23) The final extension needed to turn this into a standard LSTM is the addition of a couple tanh activations in strategic places, e.g. to rescale the values in the read-only buffer to be between  $-1$  and  $1$ , i.e.

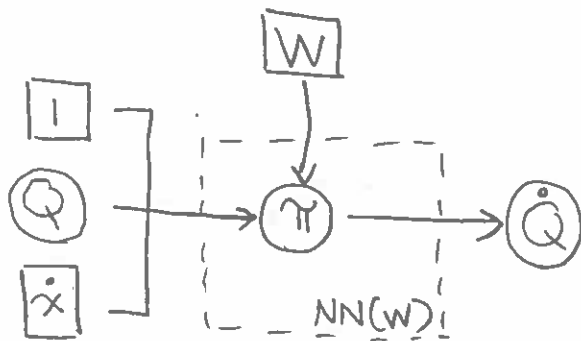


$$\dot{Q}_{RO} = \tanh(m^{(load)} \oplus \dot{Q}_{LT})$$

well, the rescaling can help gradient descent, by preventing gradients from exploding to  $\infty$



- 24) Anyways, even comparing our LSTM-lite with the vanilla RNN:



it's clear that it's much more complicated. However it does tend to work better, particularly for long sequential iterations. This is generally attributed to the long-term state, which allows the network to "remember" information that may only be relevant much later in the sequence.

## LONG SHORT-TERM MEMORY NETWORKS

---

②5 As an example of when these "long term relevancies" might occur, consider the task of coreference resolution in English.

"Suppose two integers sum to seven, and moreover, the product of the numbers is twelve. If so, what are the two numbers?"

If an RNN is processing this passage word by word, then it needs the capacity to store information about the "two integers" and retrieve it much later in the passage (both when it is stated that the product of the numbers is twelve, and when it is asked what the two numbers are).