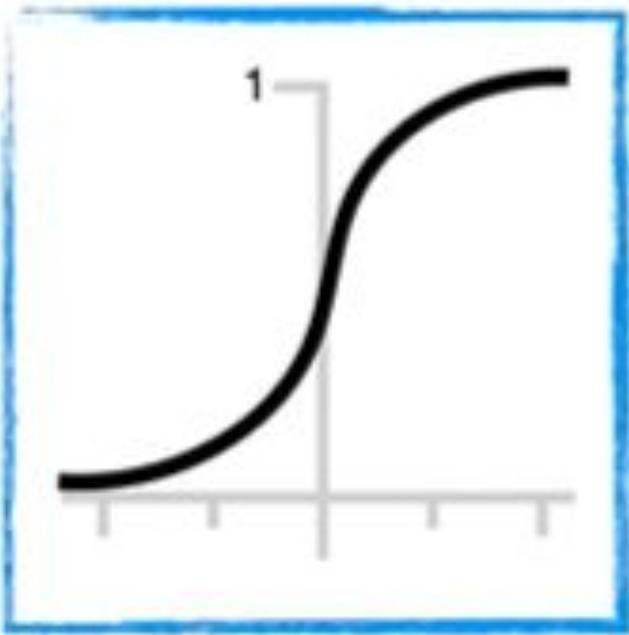
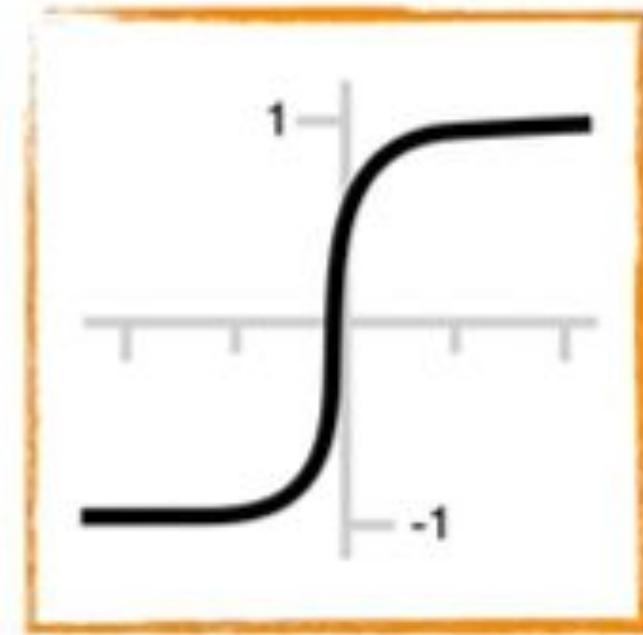


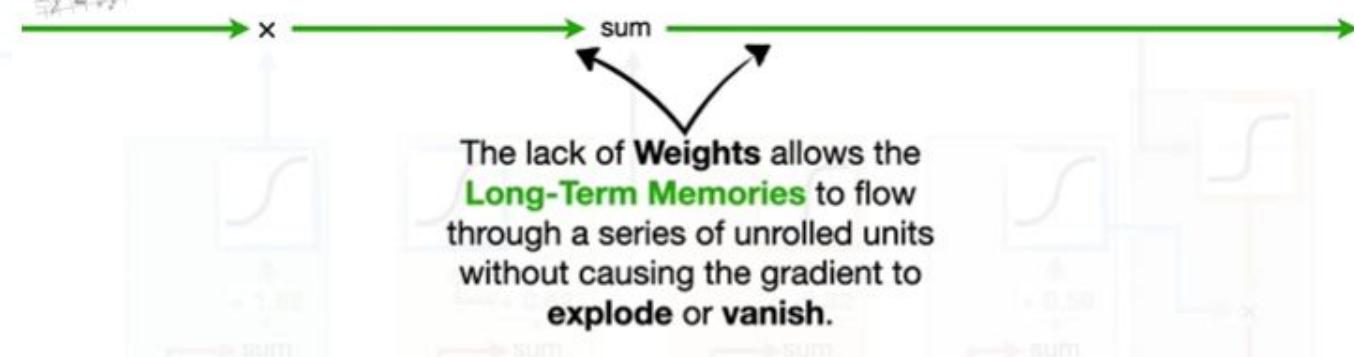
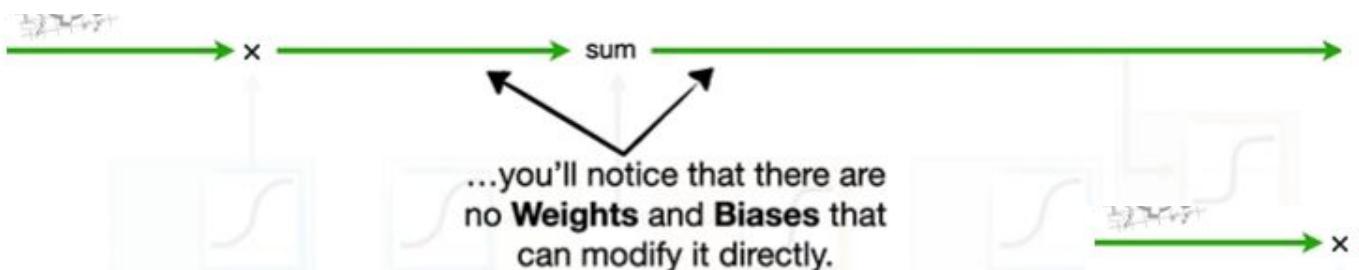
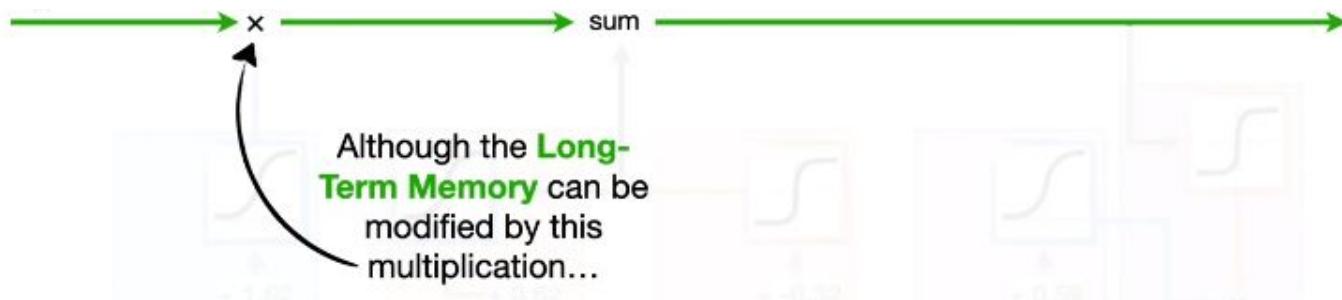
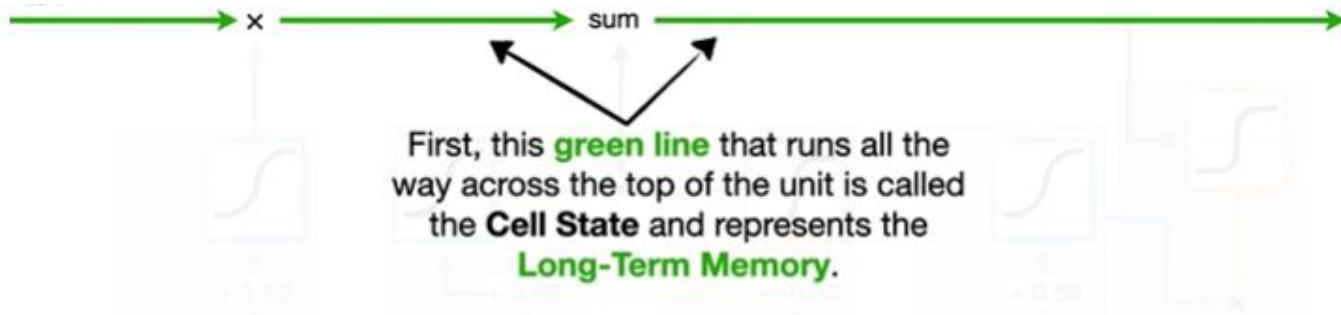
LSTM

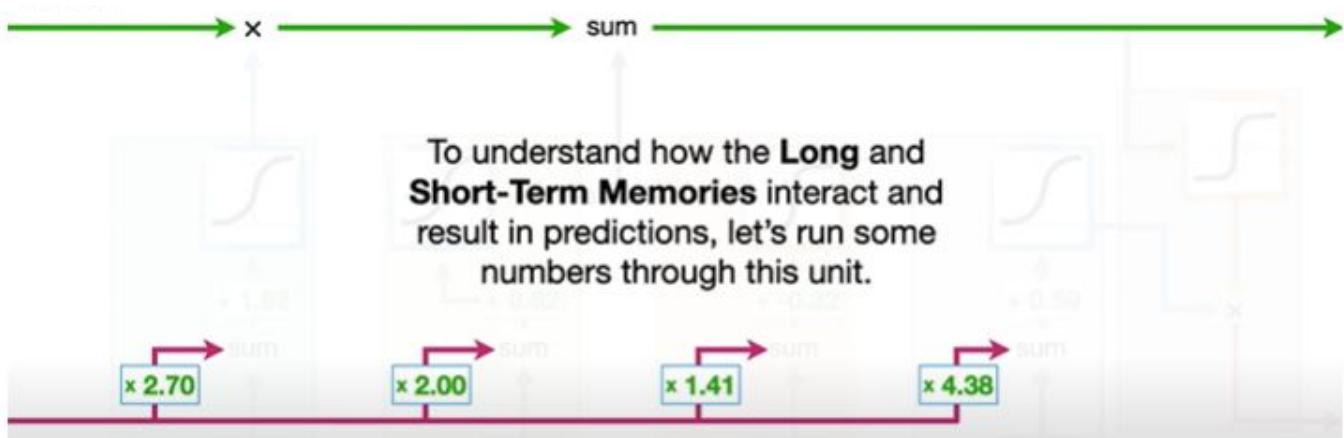
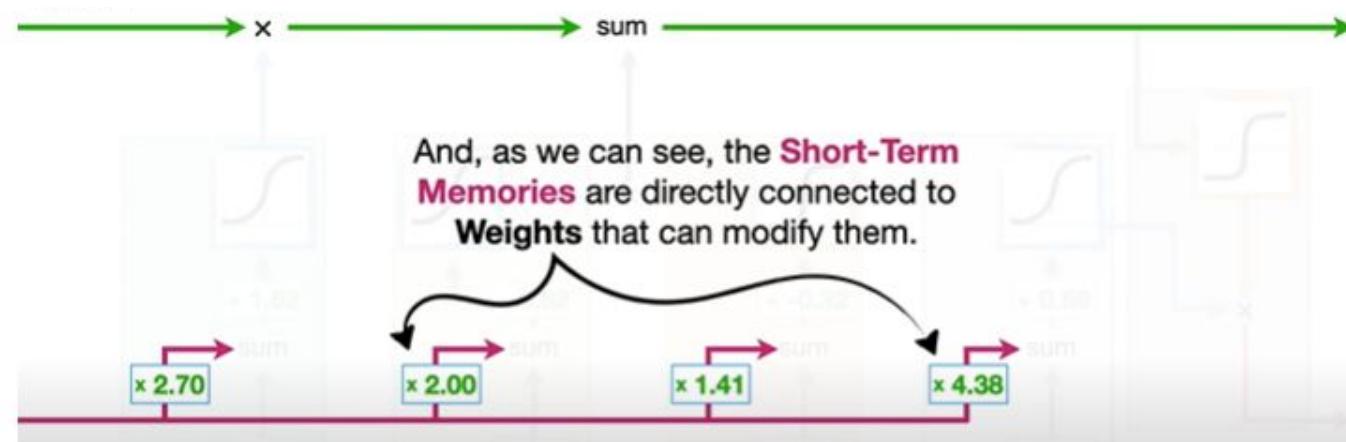
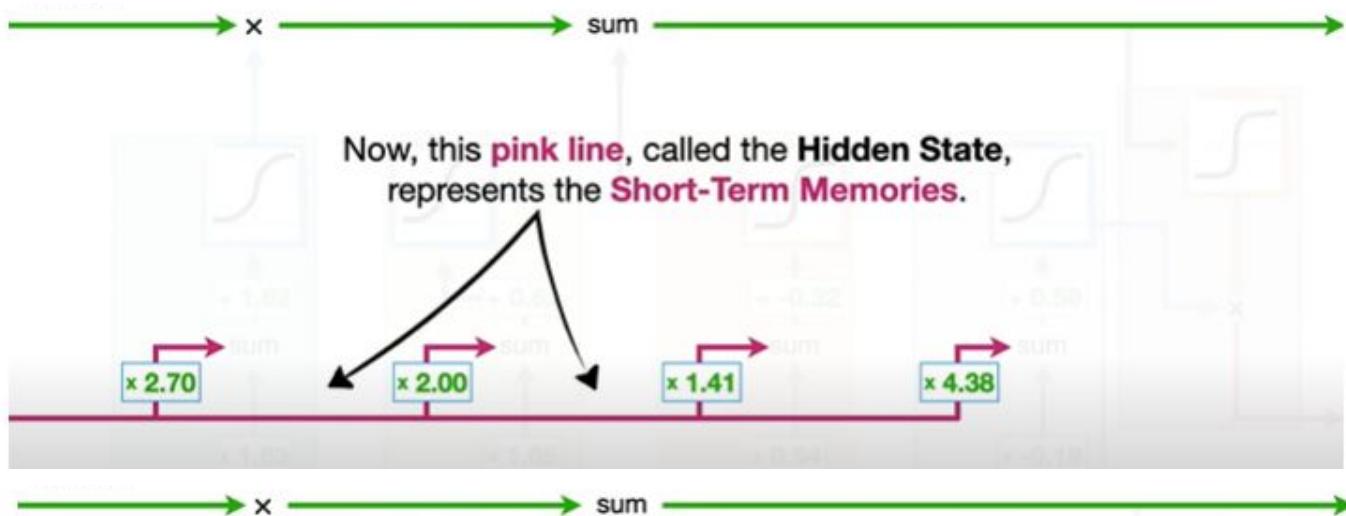
So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...

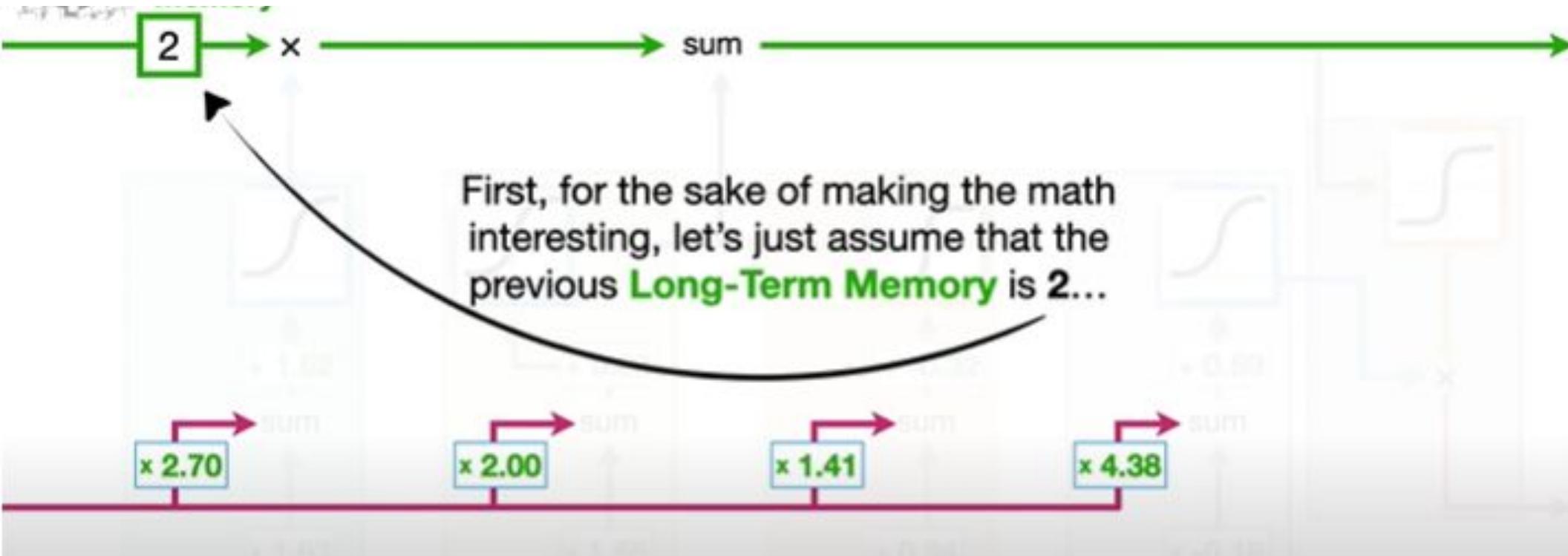


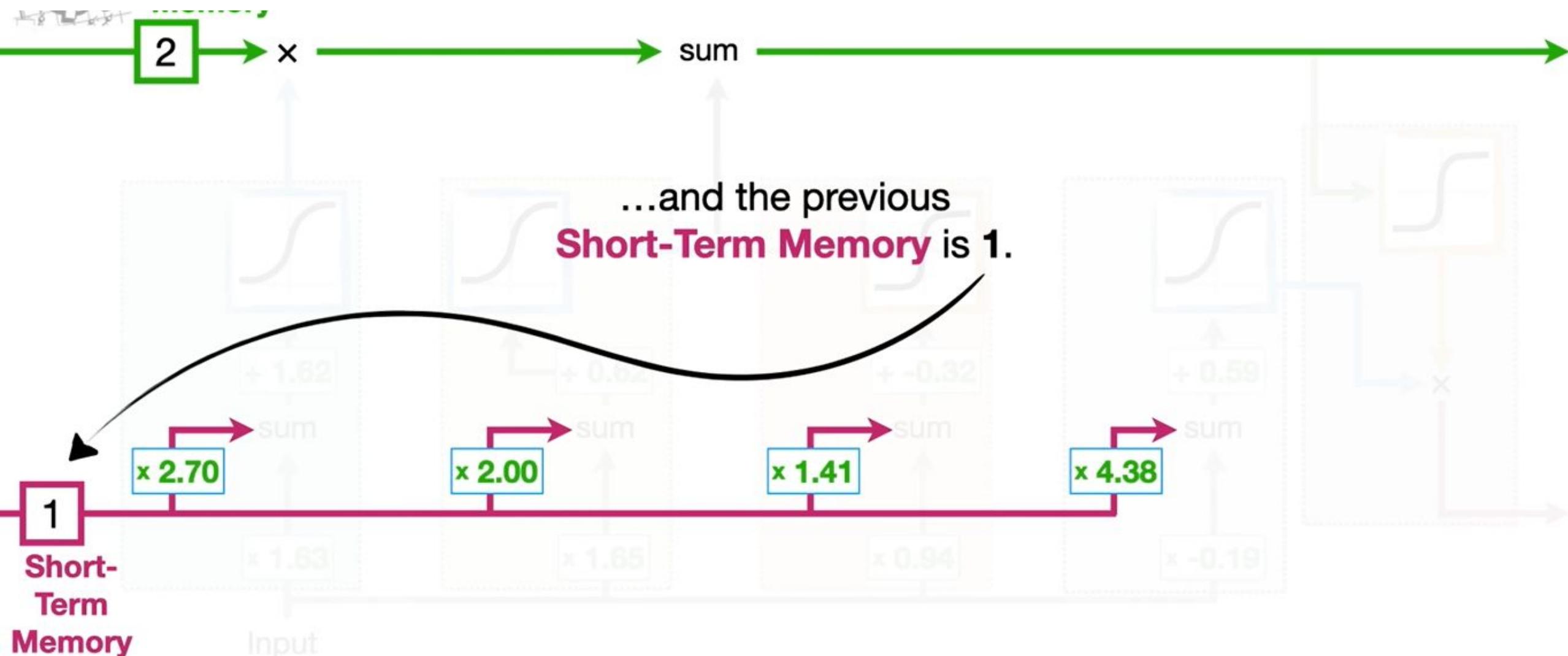
...and the **Tanh Activation Function** turns any input into a number between **-1** and **1**...





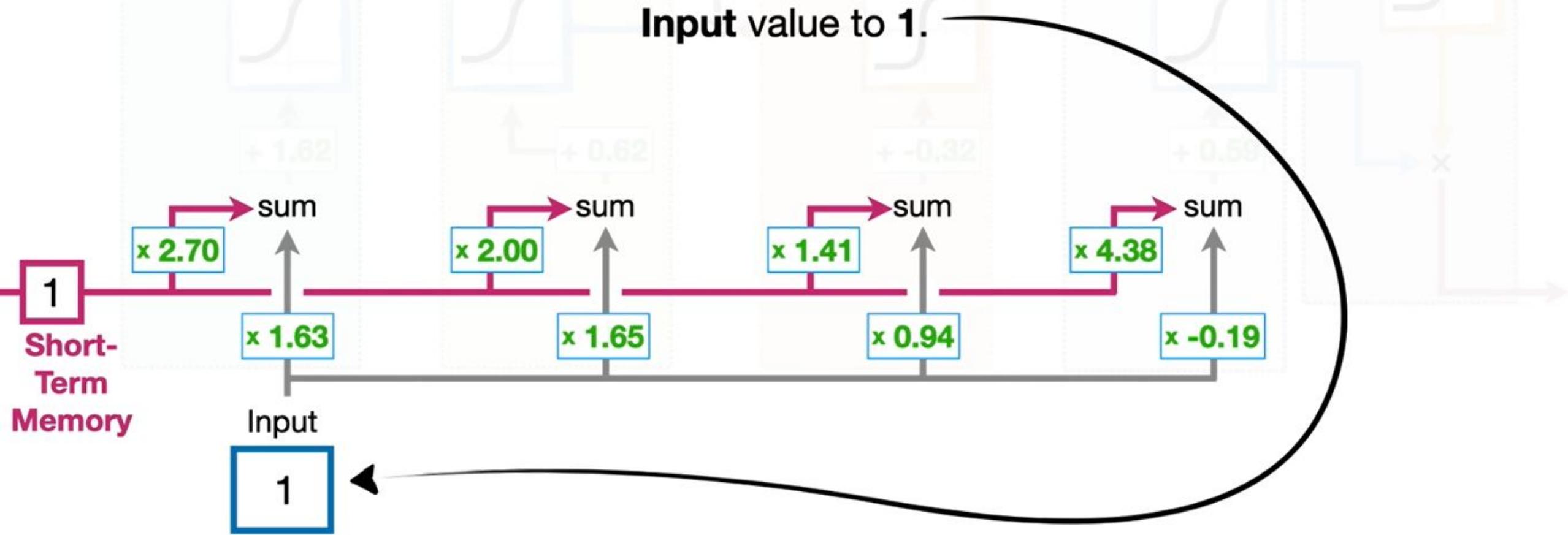






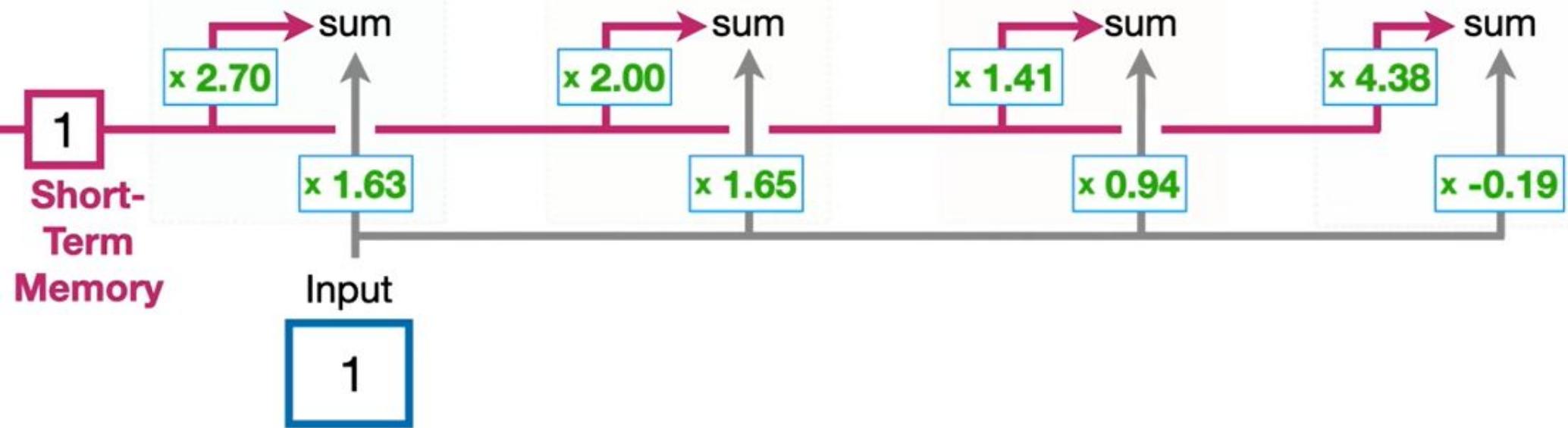


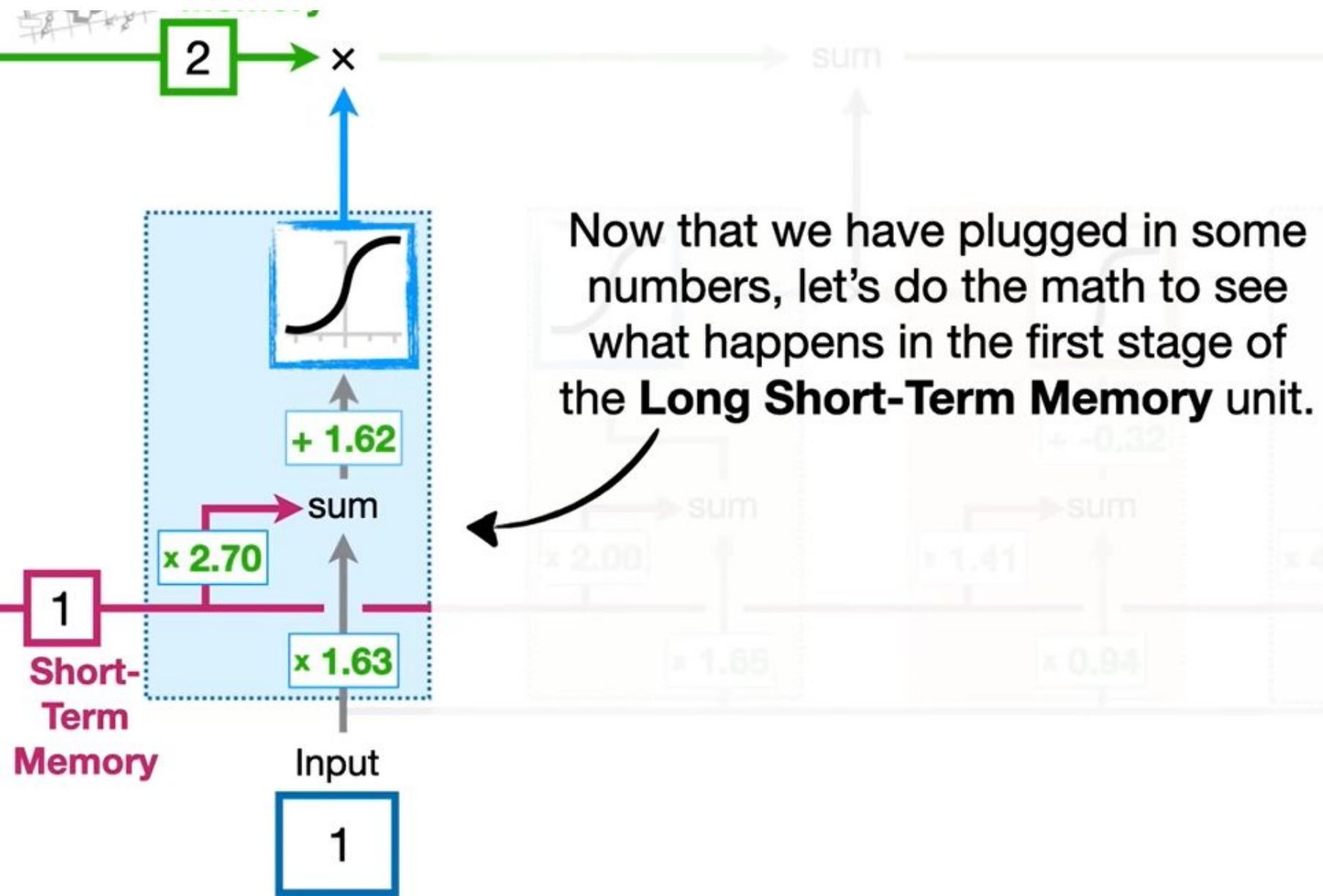
And let's also set the
Input value to 1.



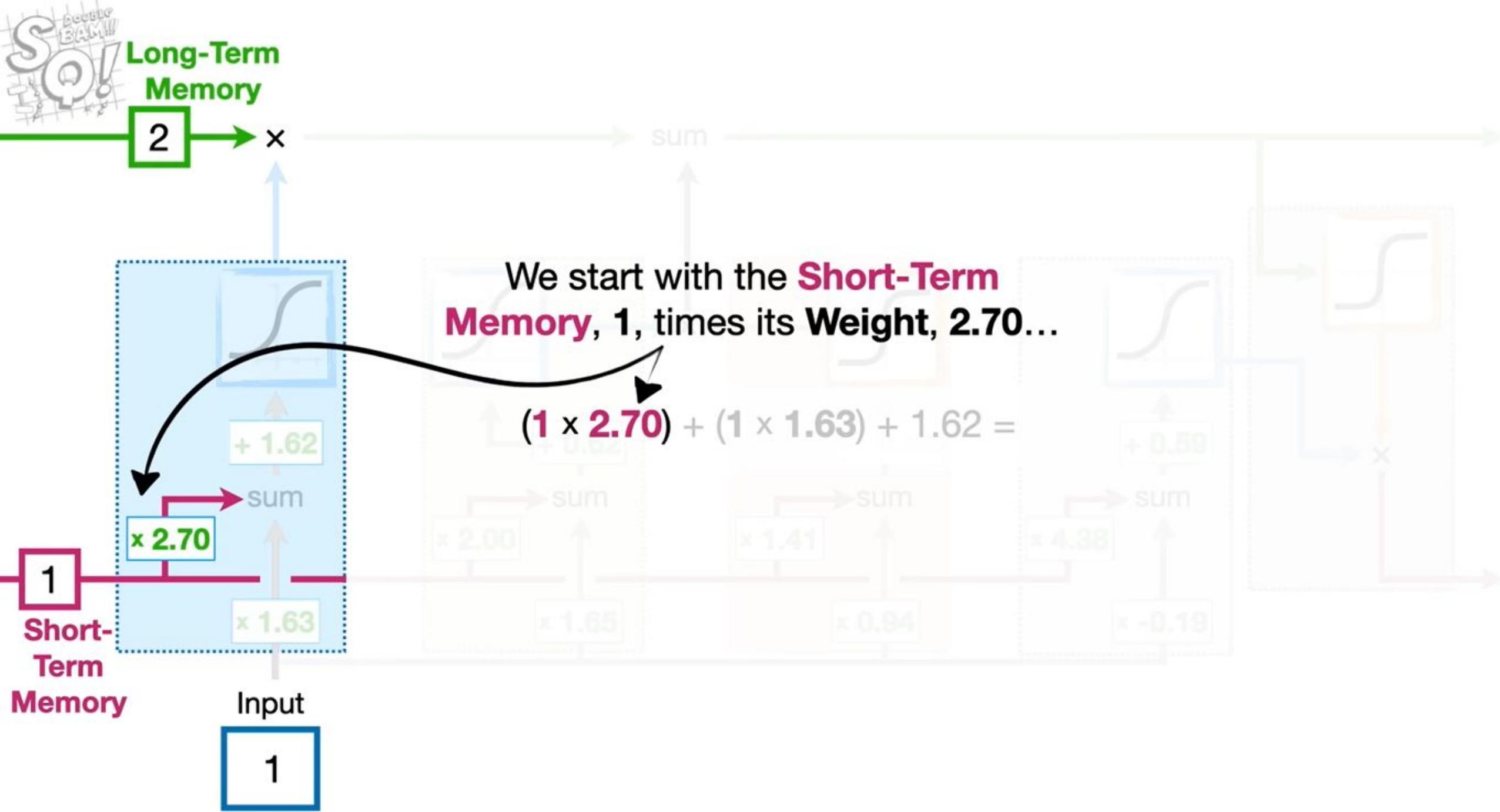


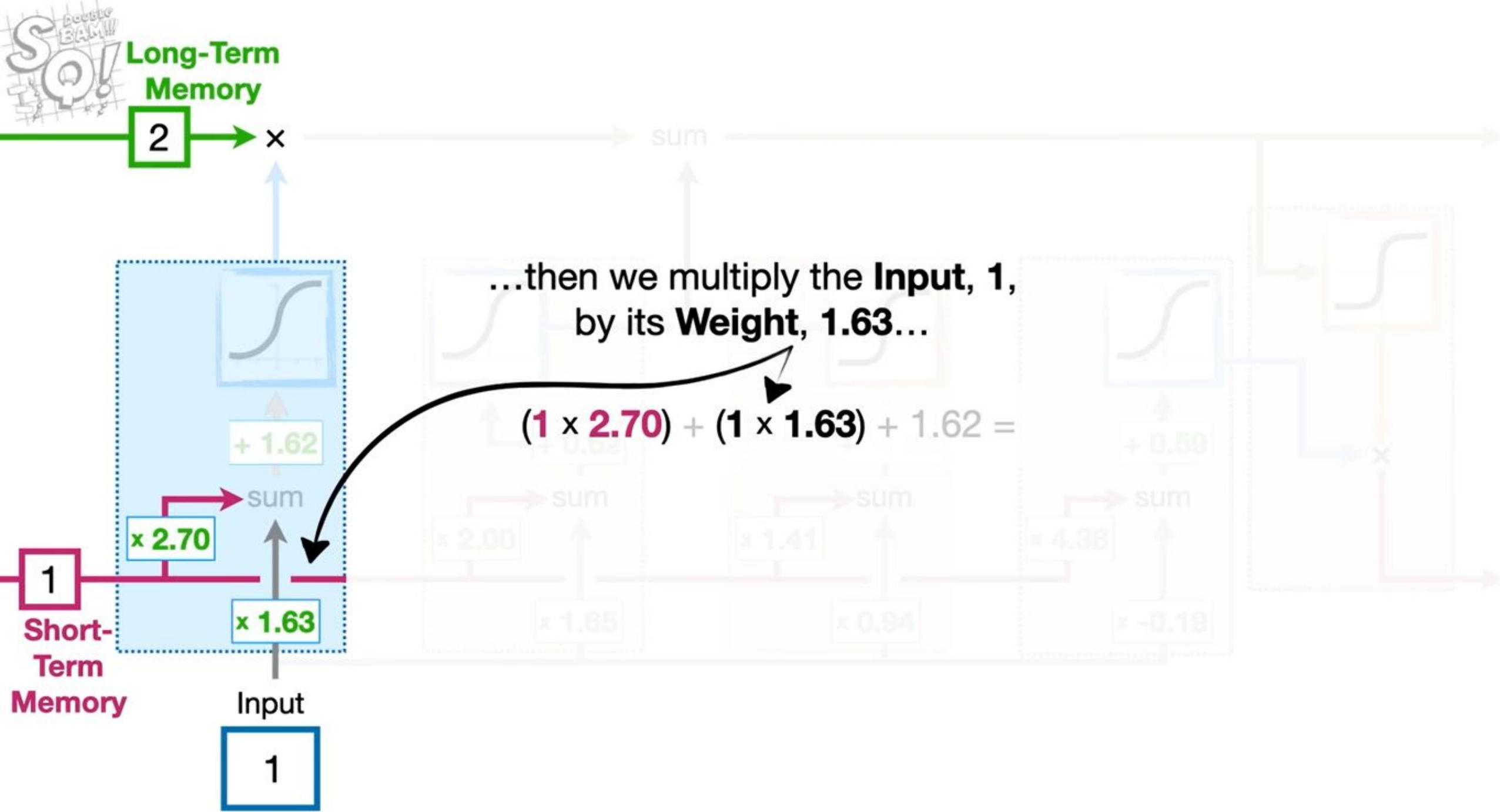
Now that we have plugged in some numbers, let's do the math to see what happens in the first stage of the **Long Short-Term Memory** unit.

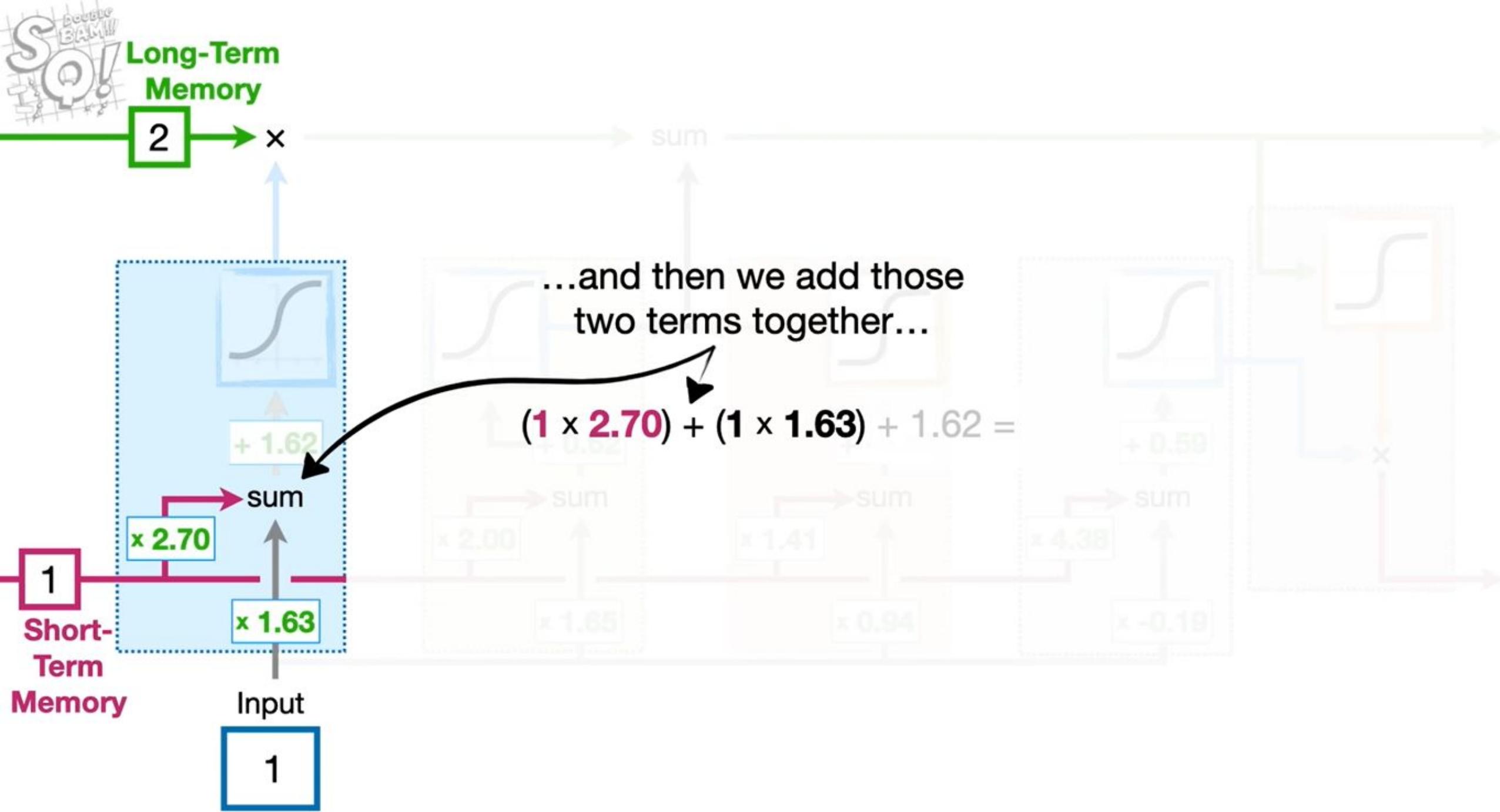


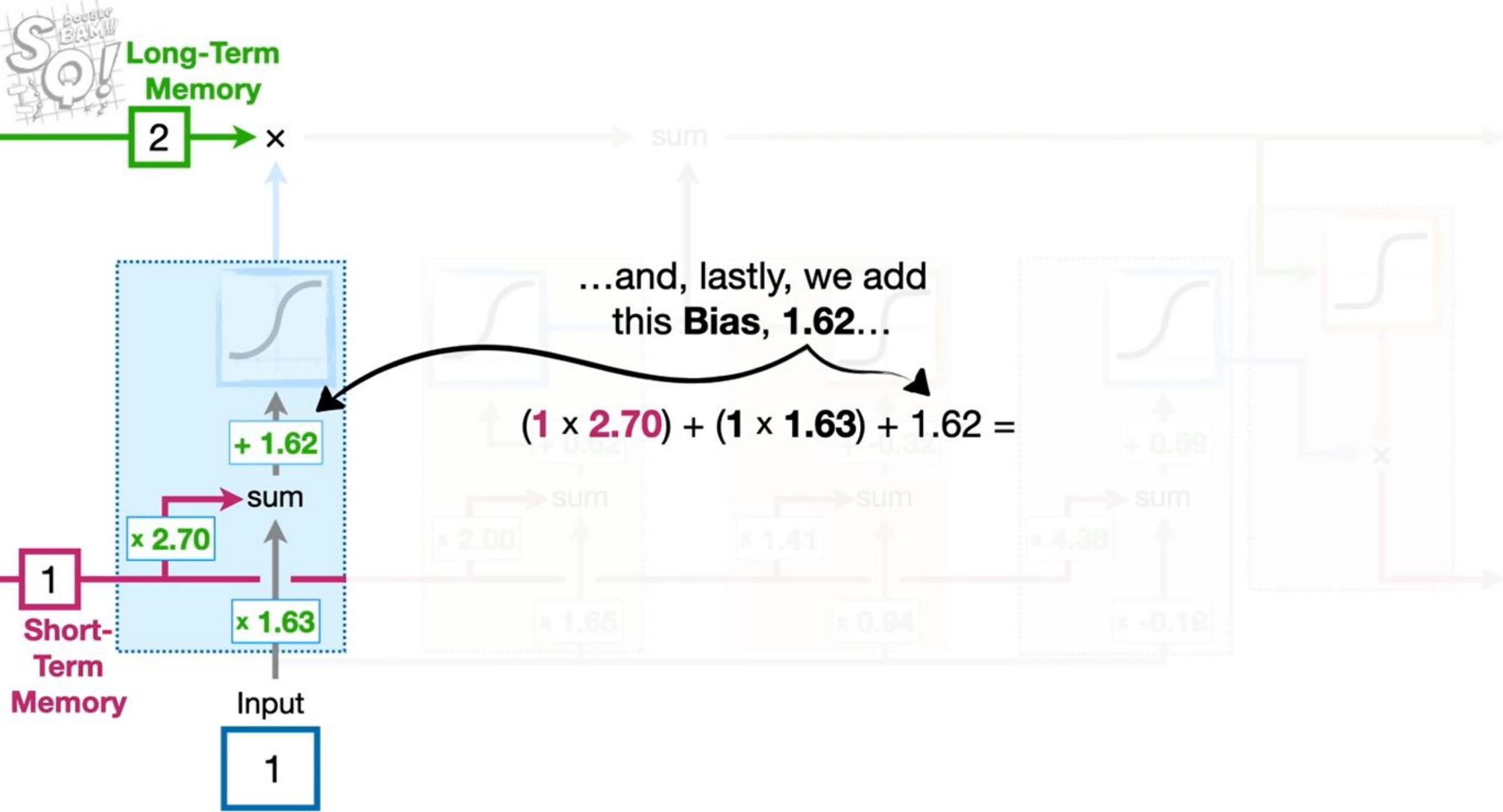


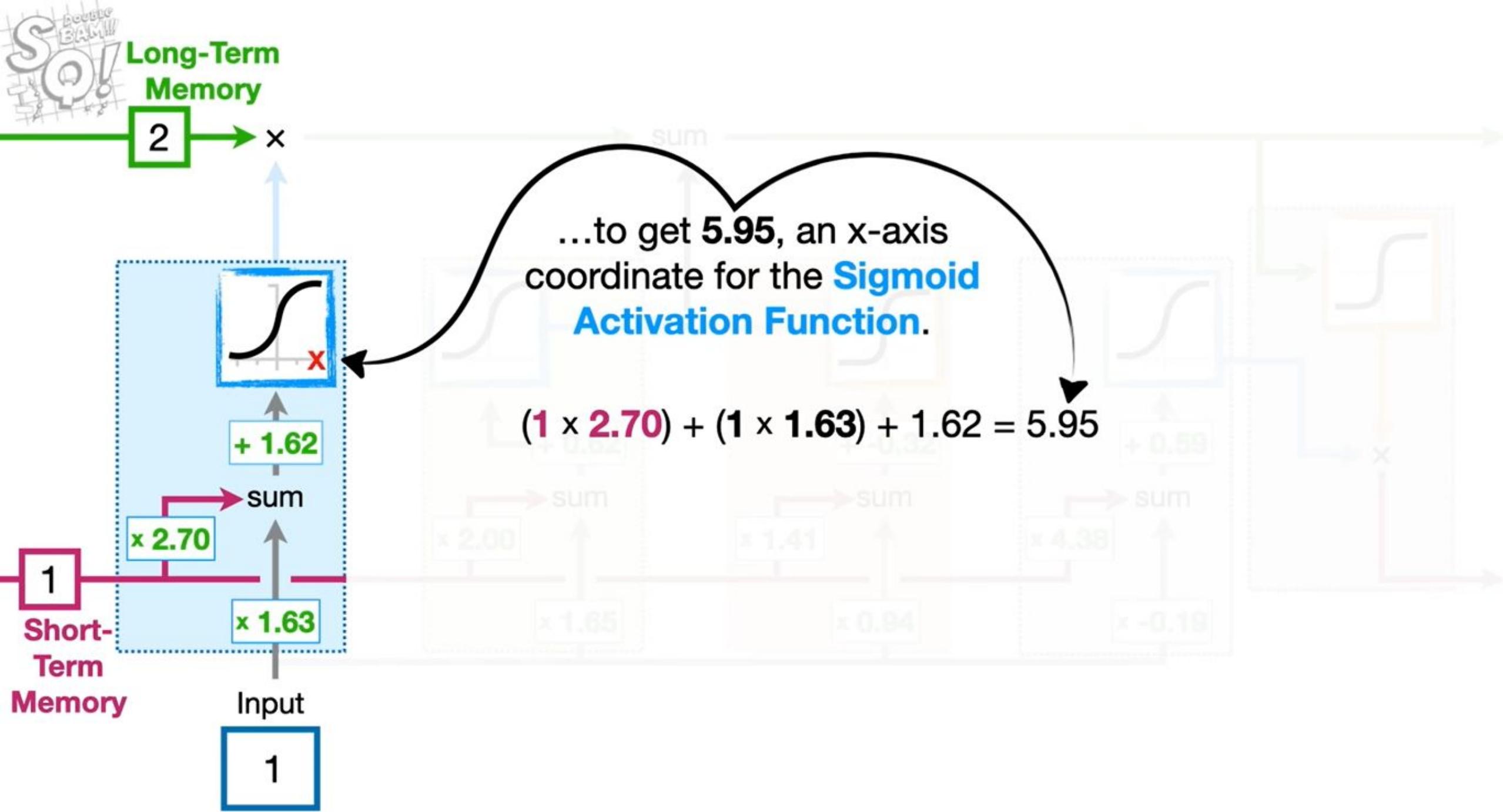
Now that we have plugged in some numbers, let's do the math to see what happens in the first stage of the **Long Short-Term Memory** unit.

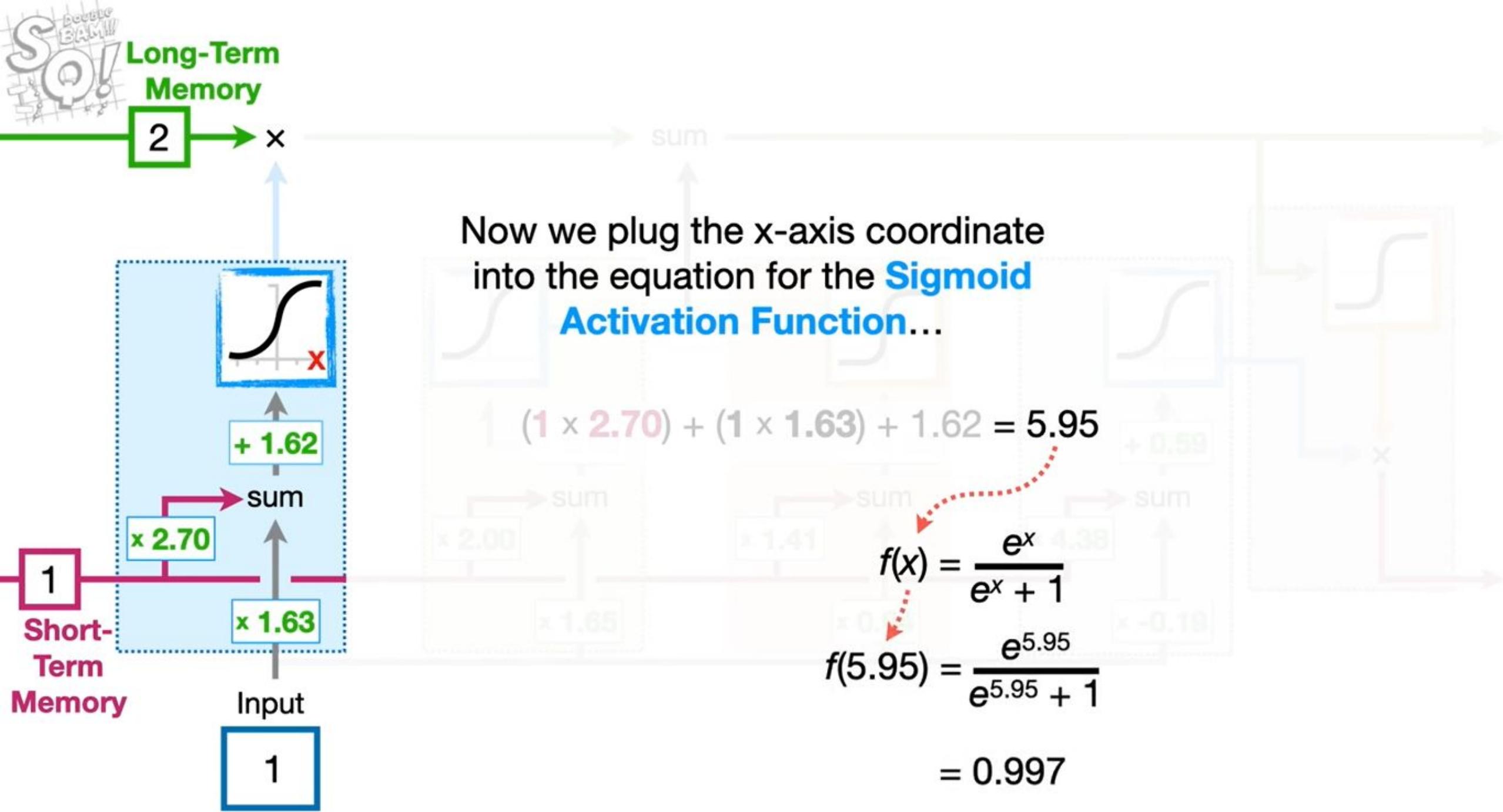


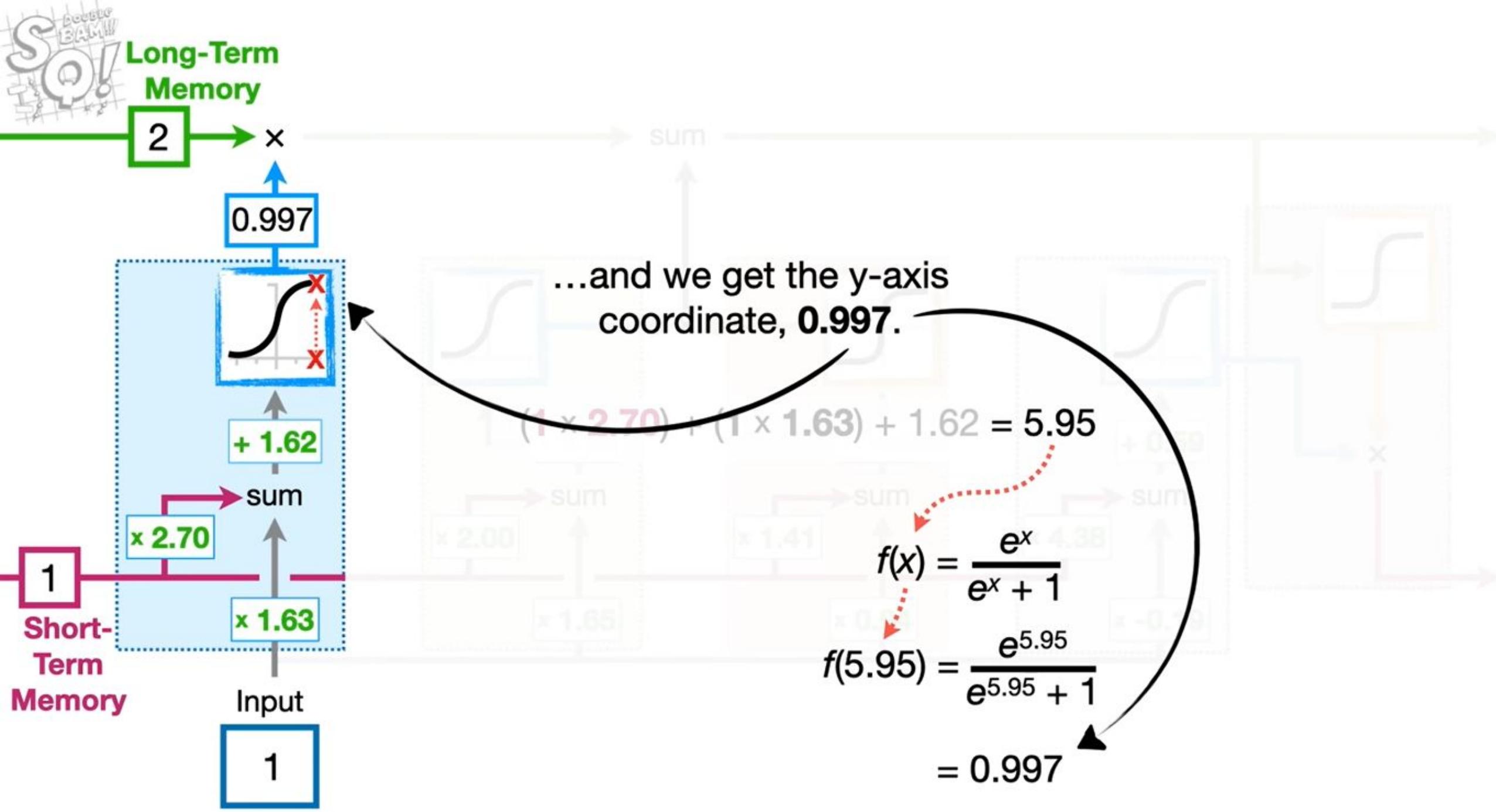


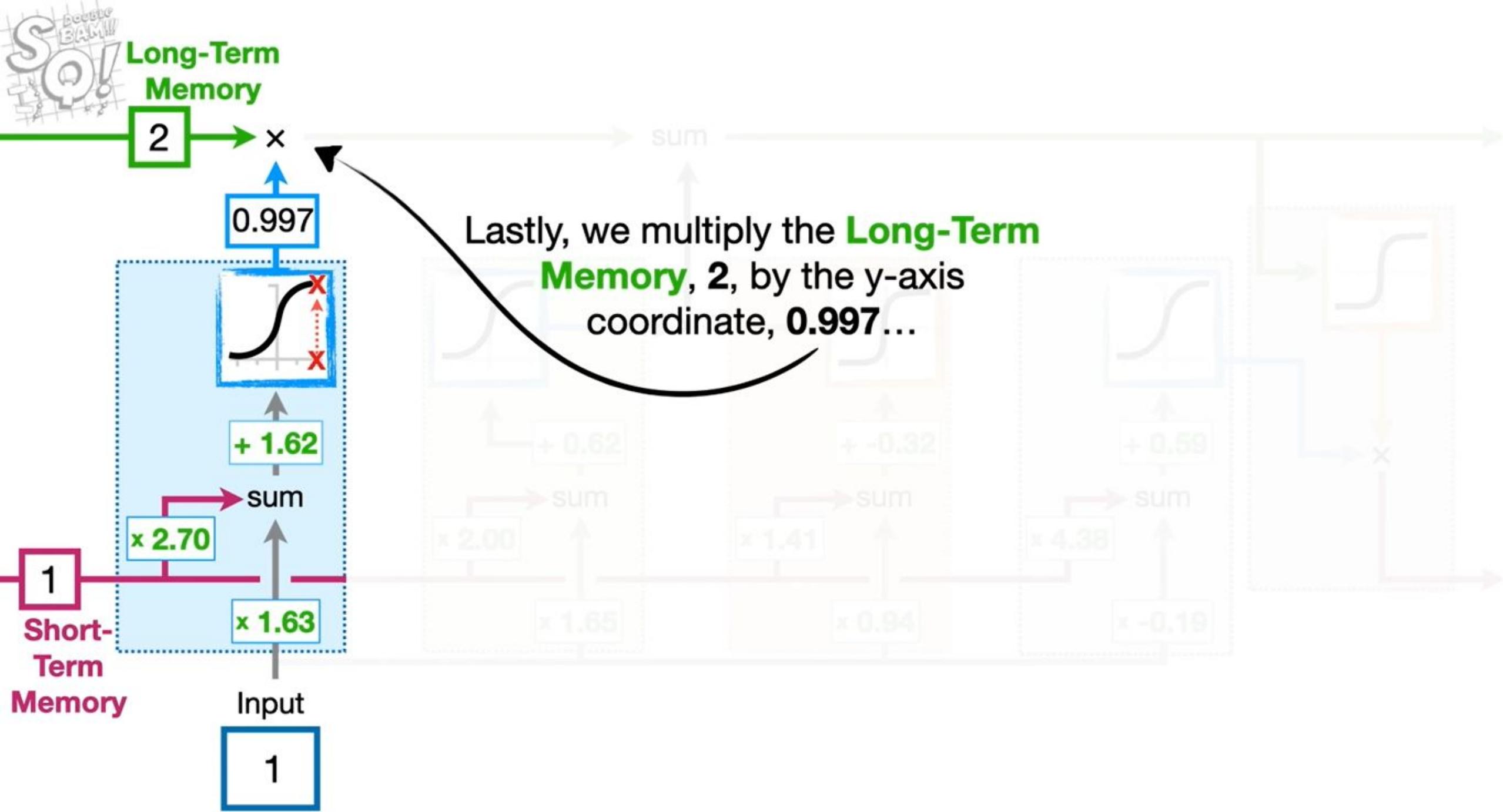


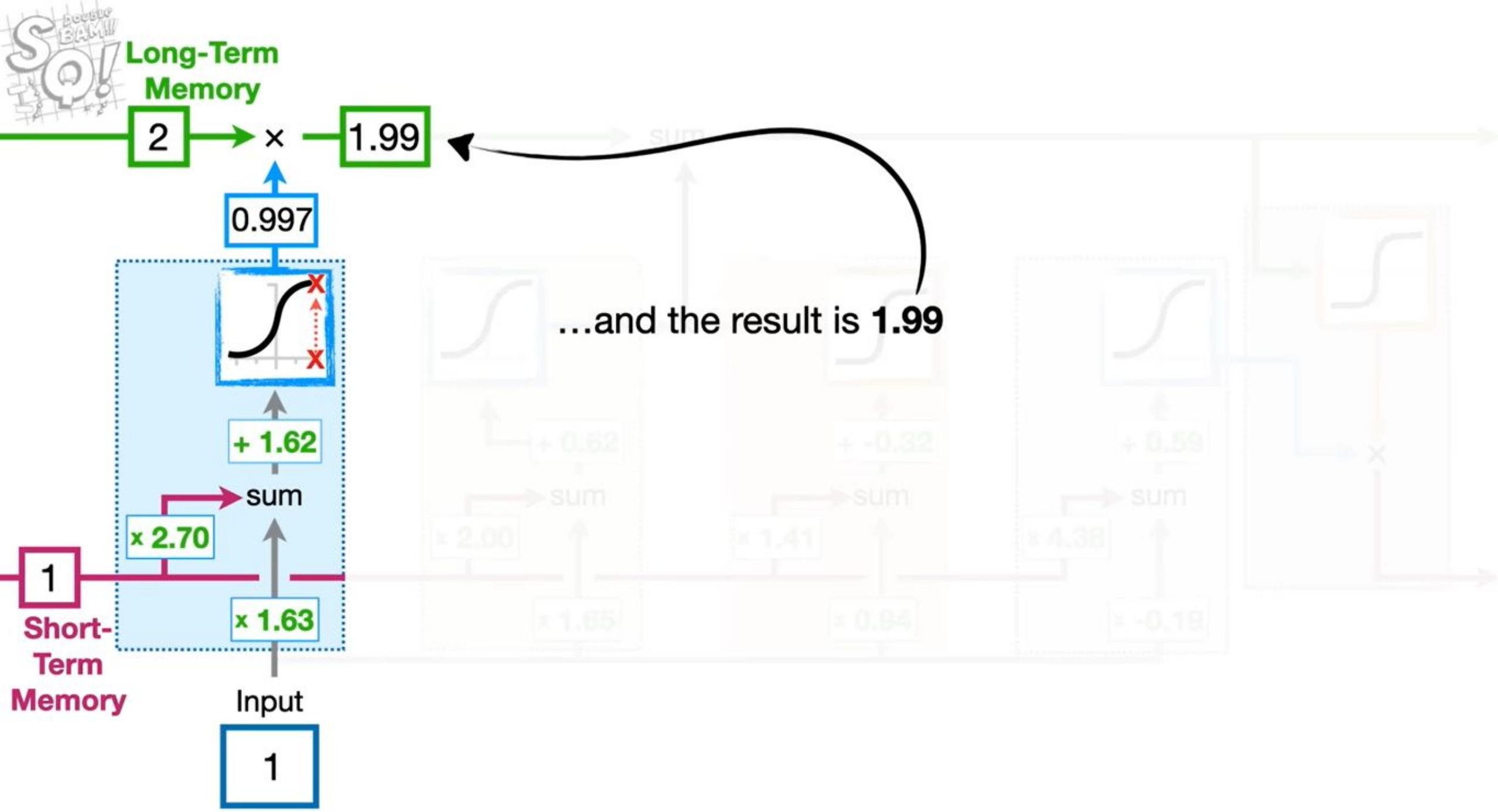


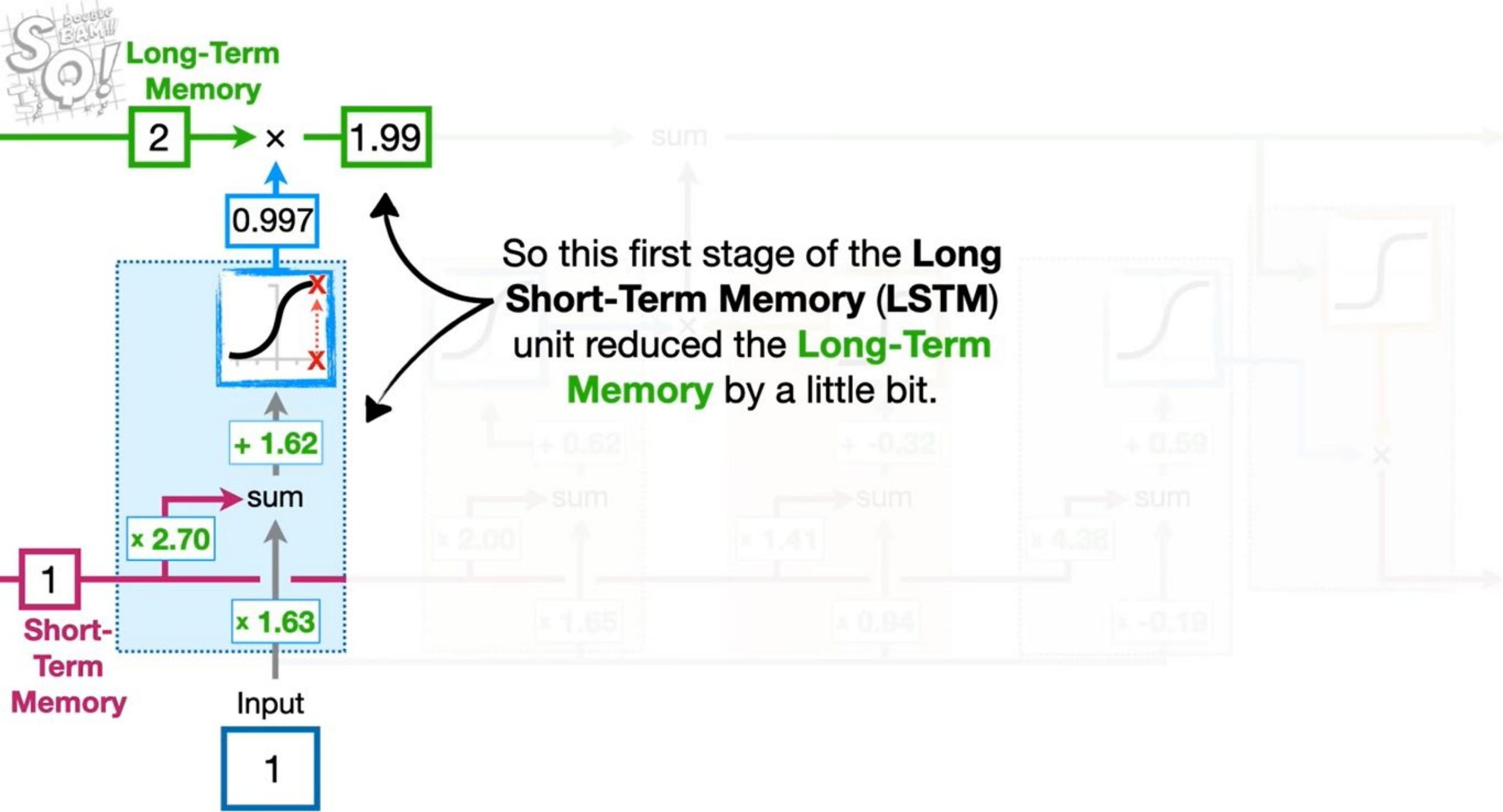


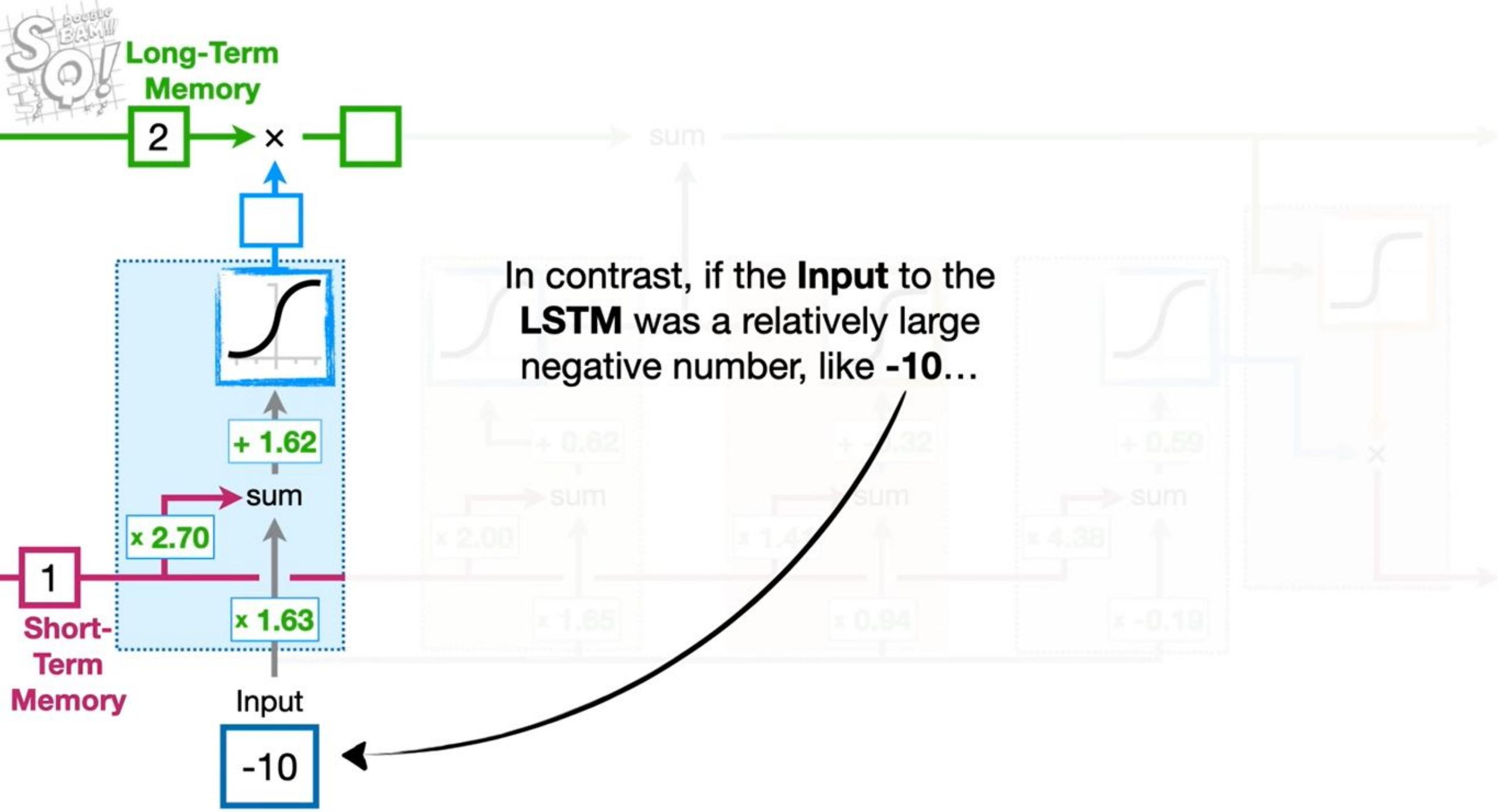


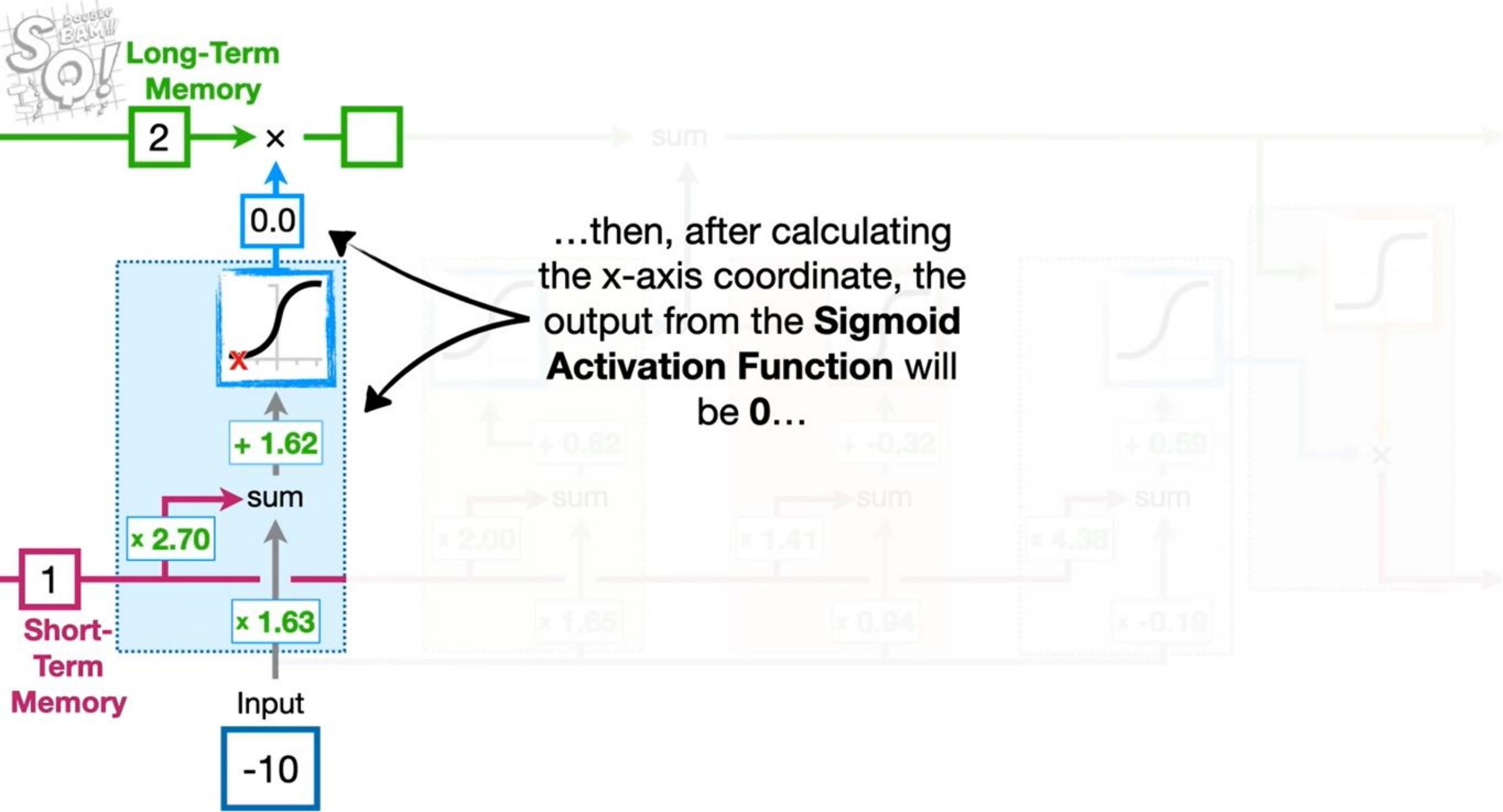


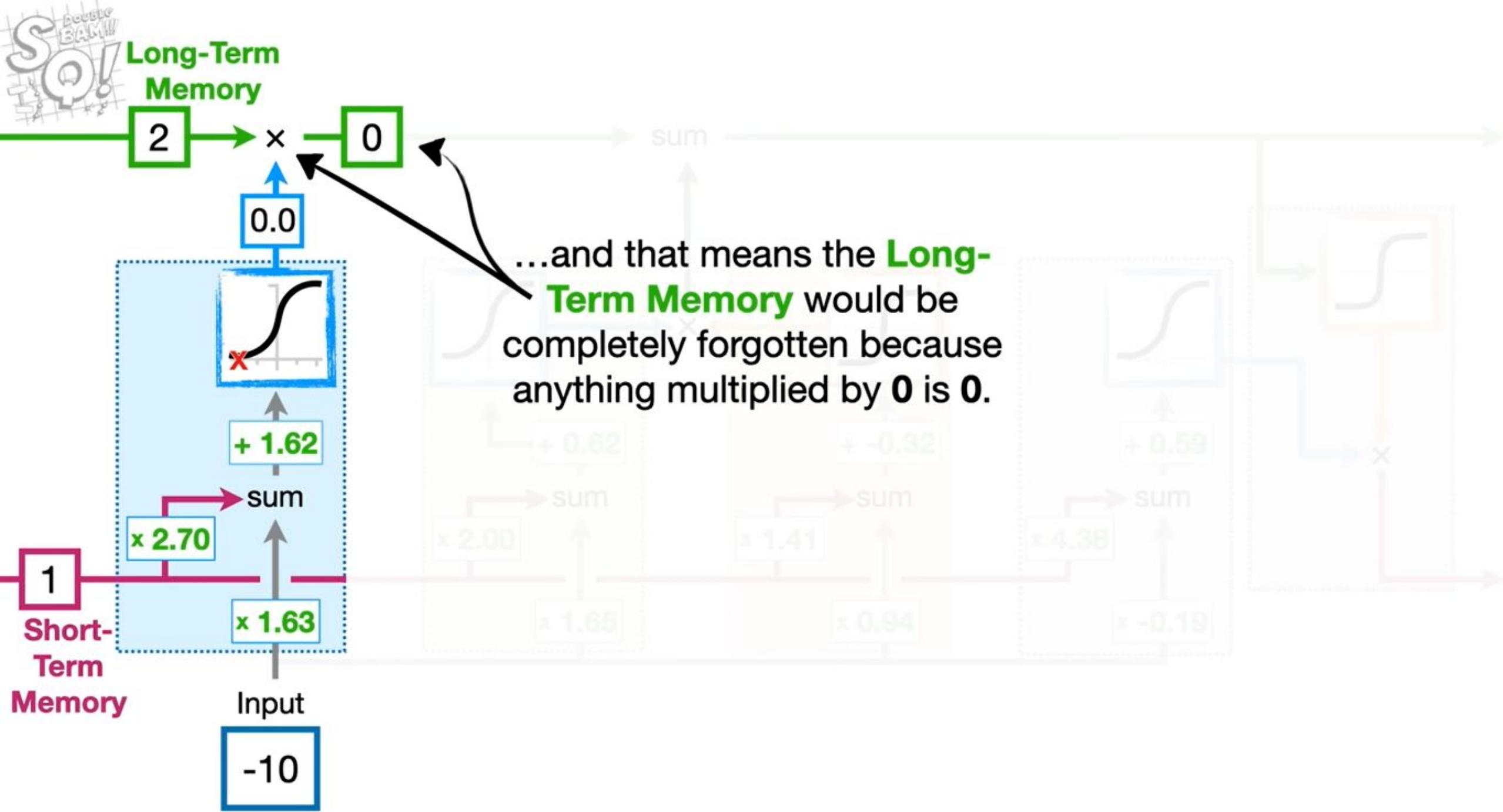


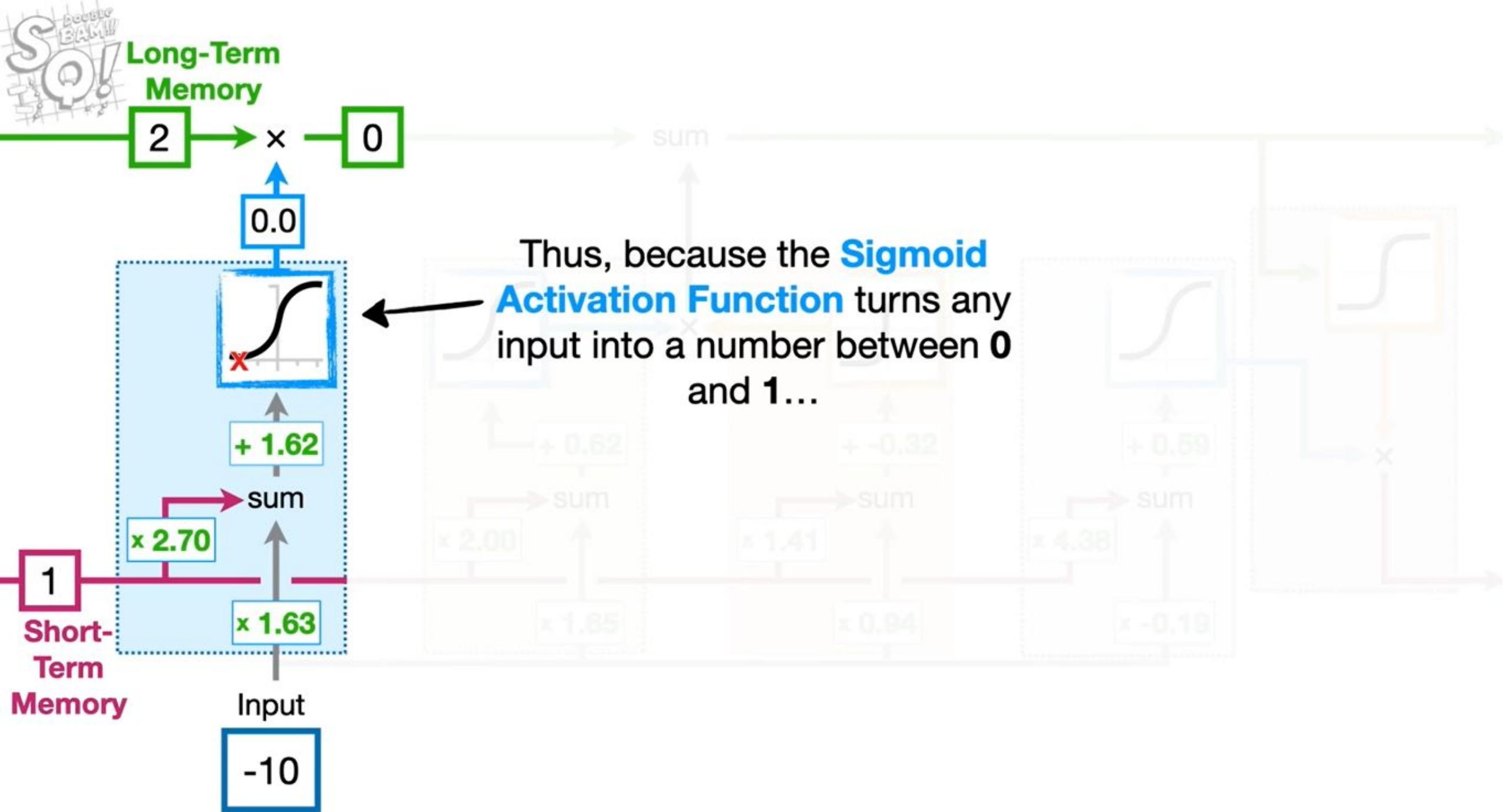


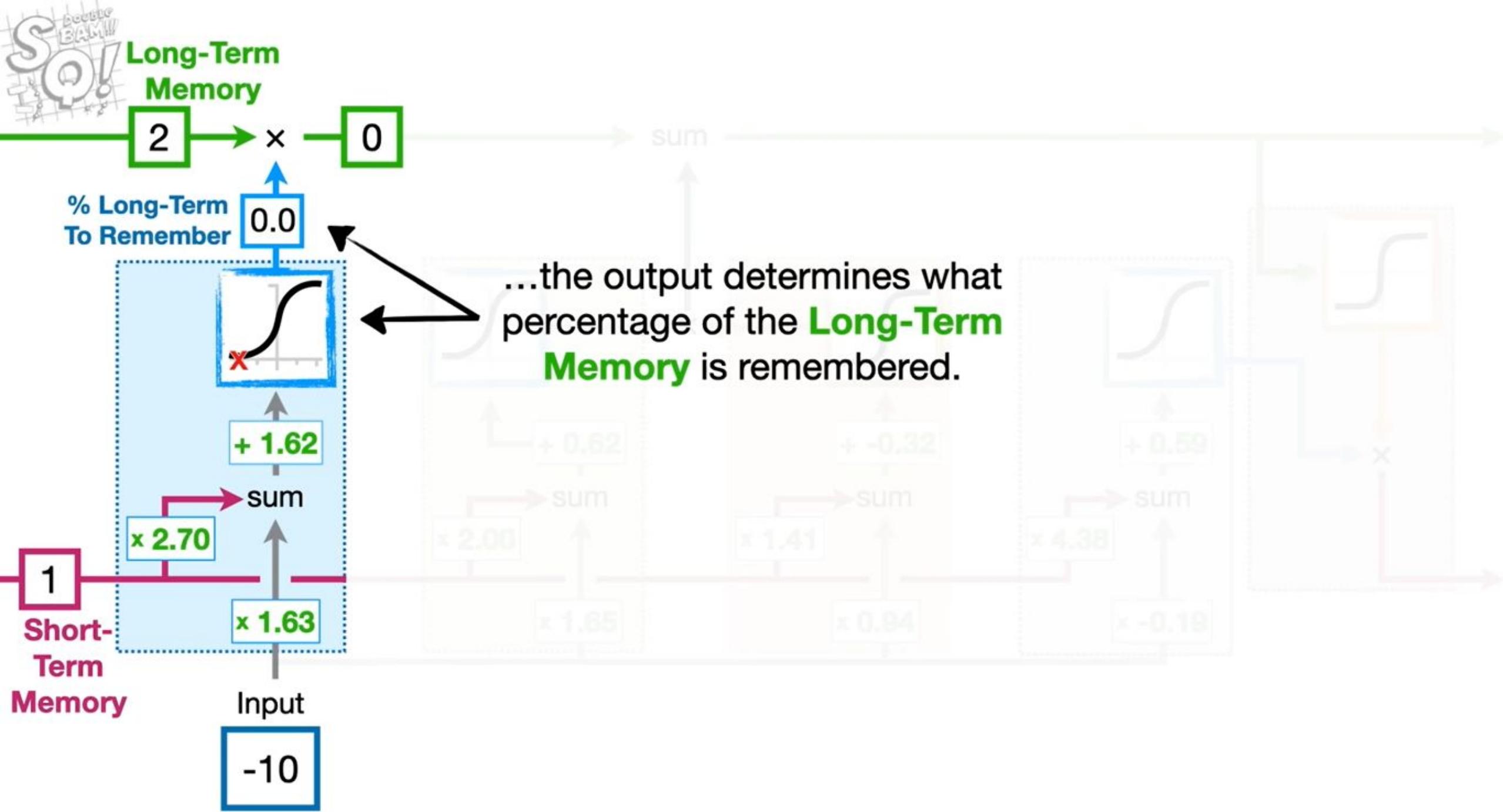


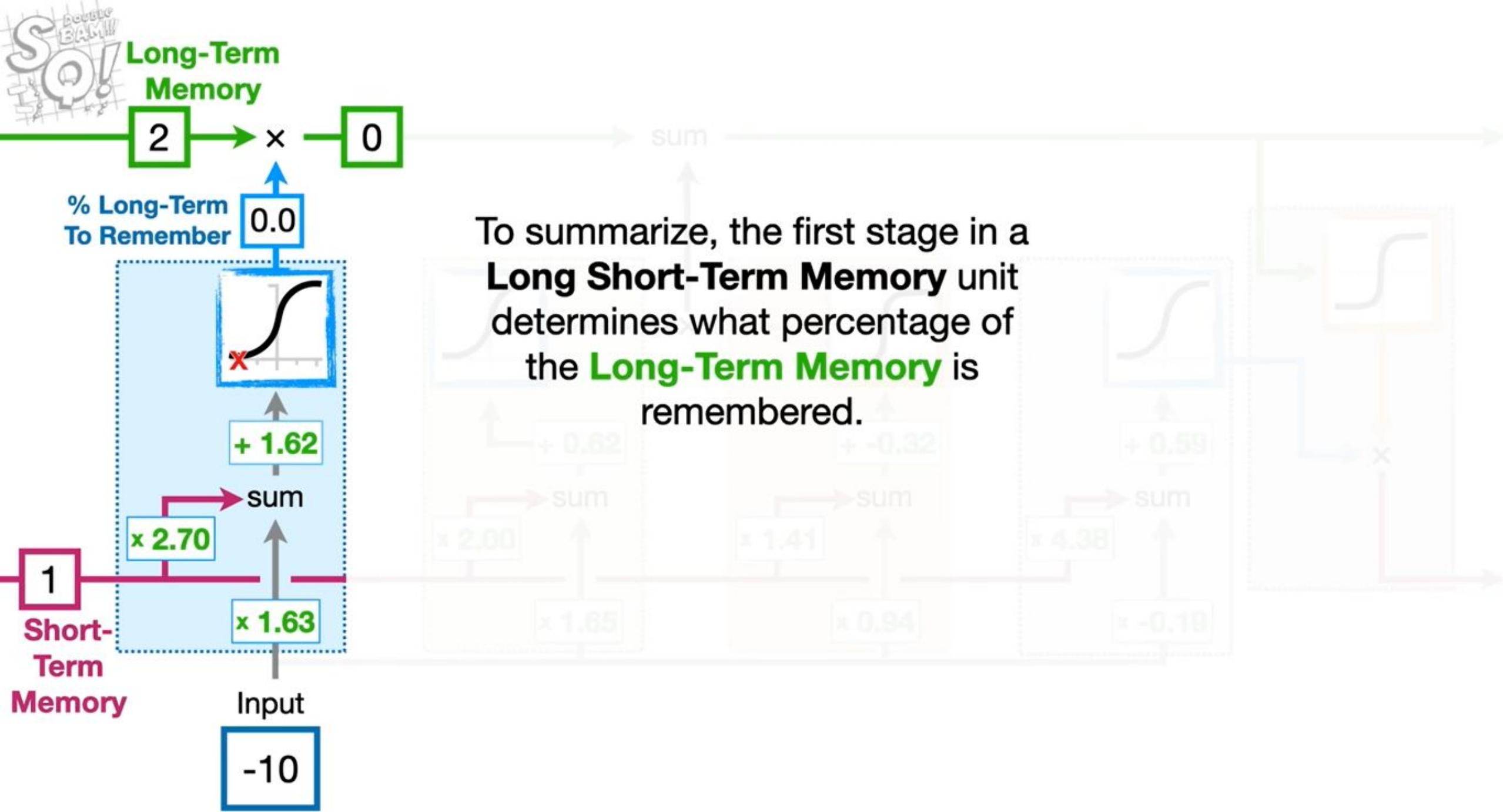


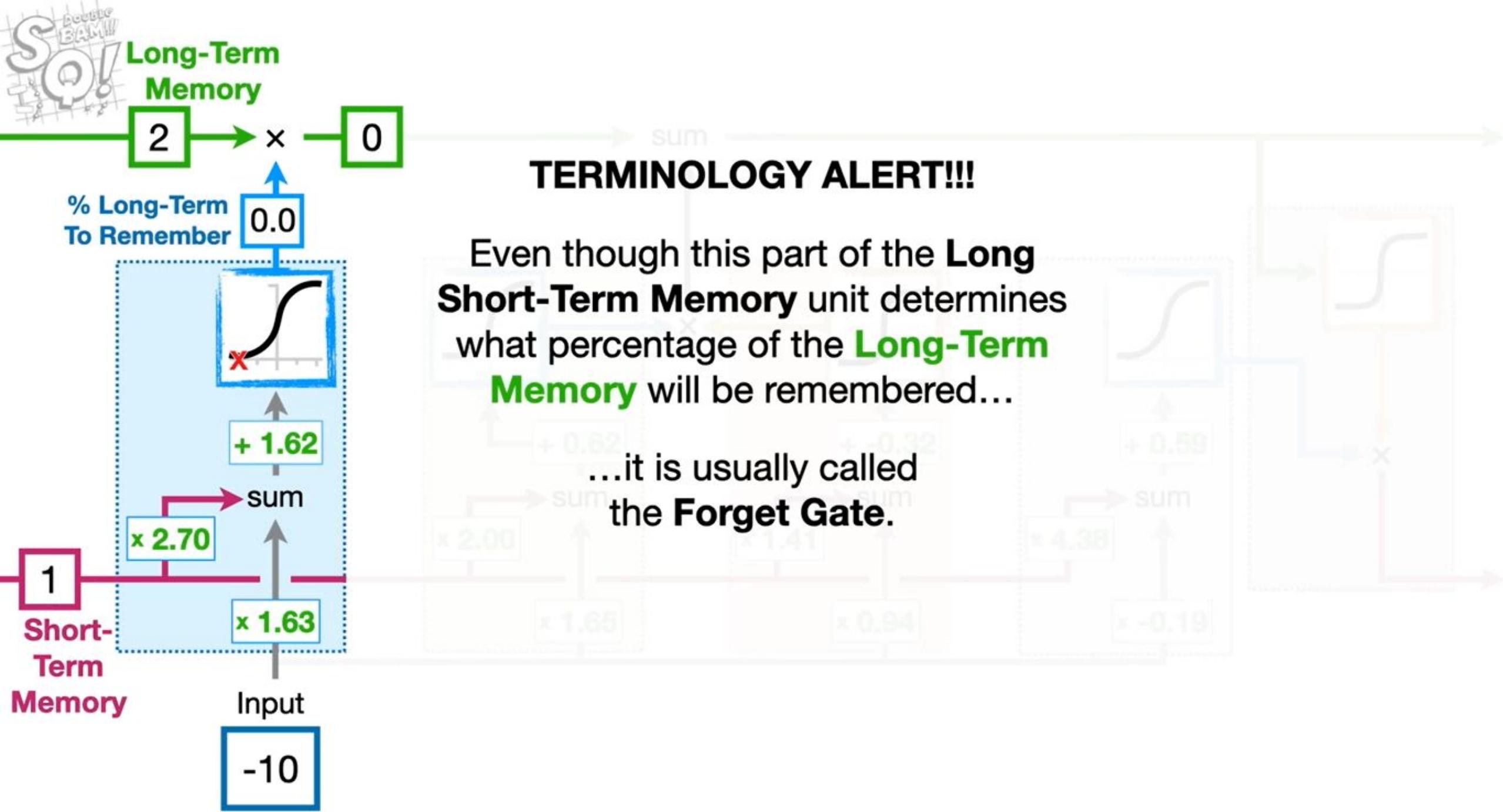


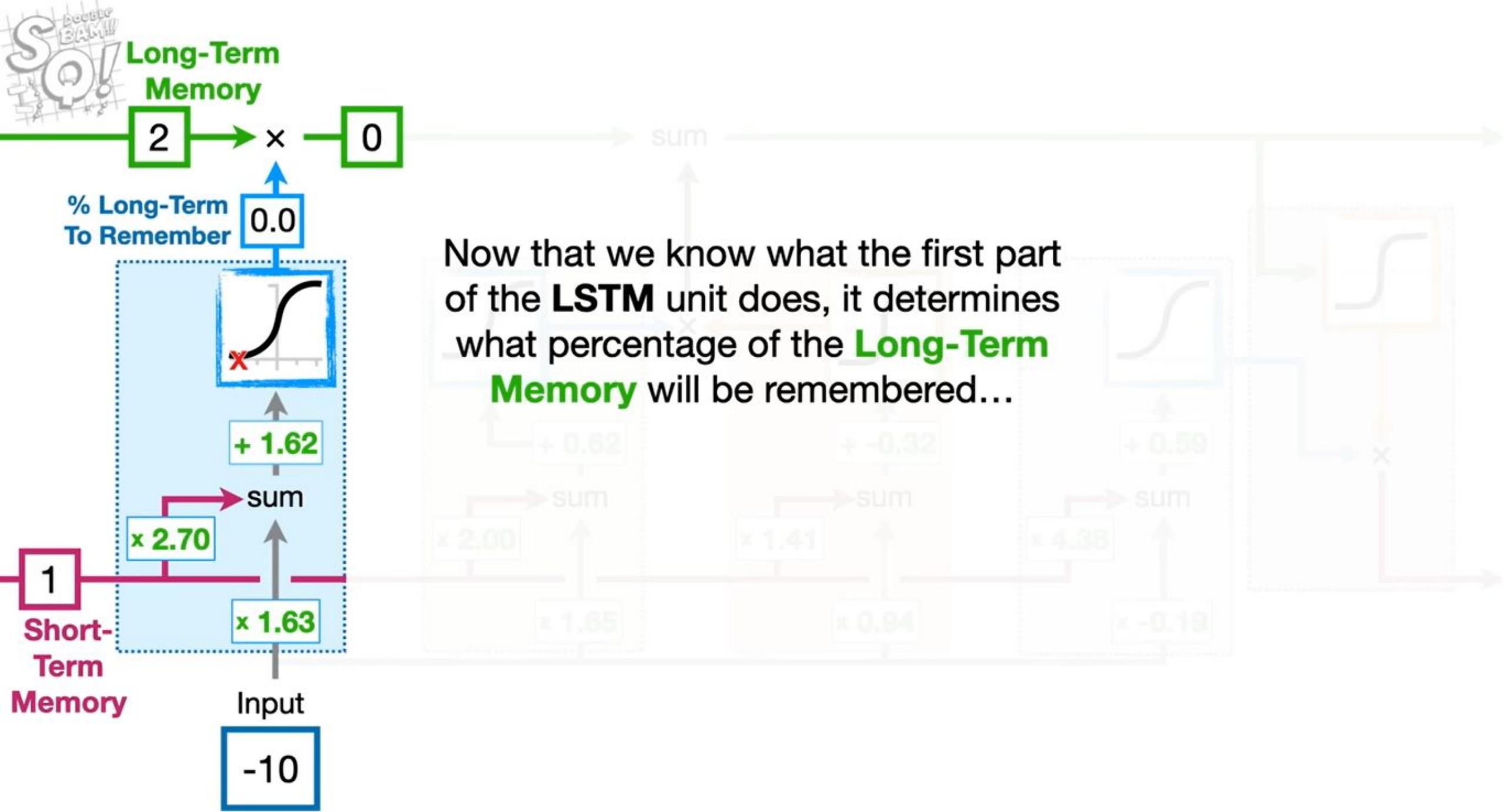


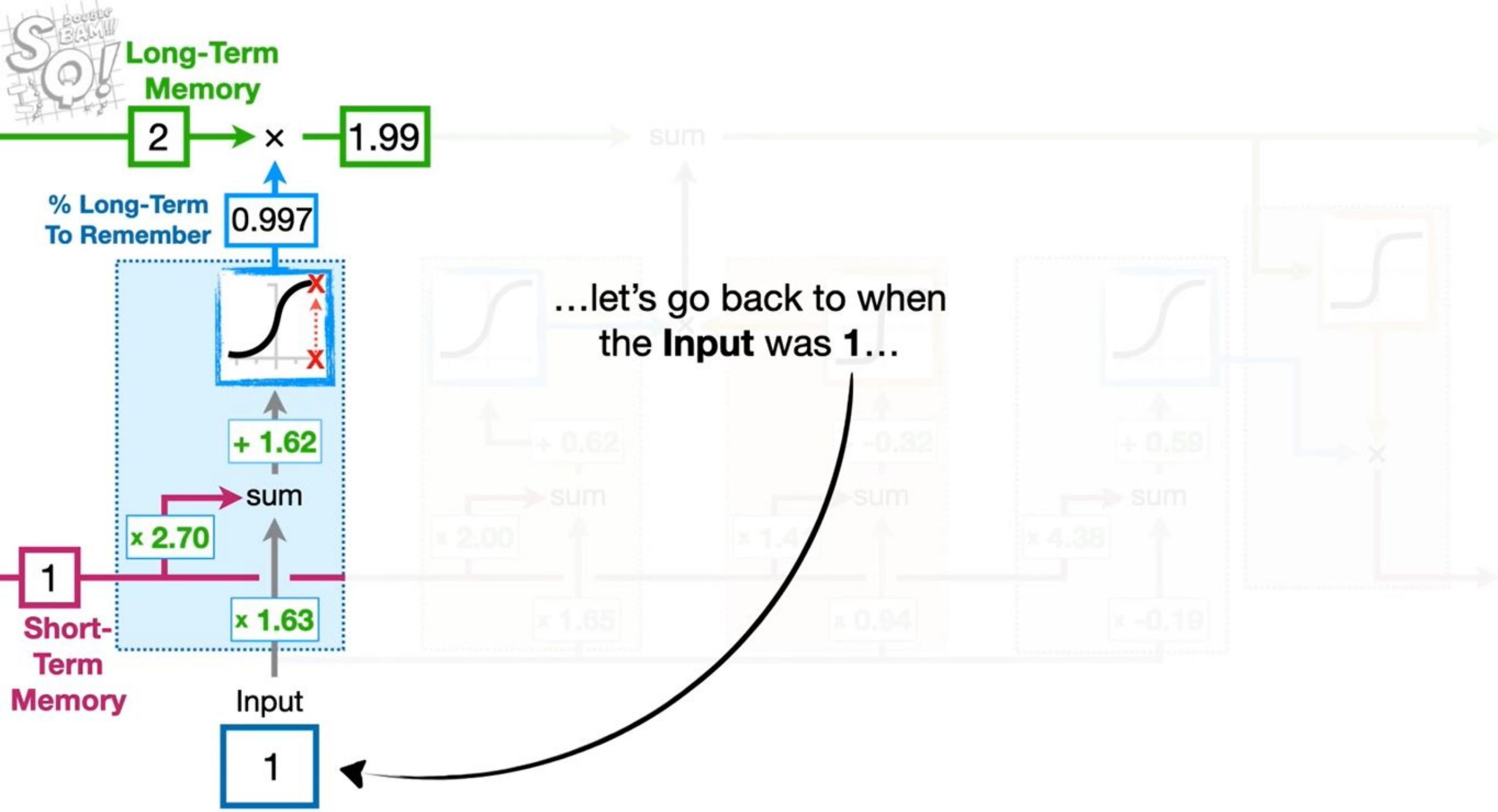


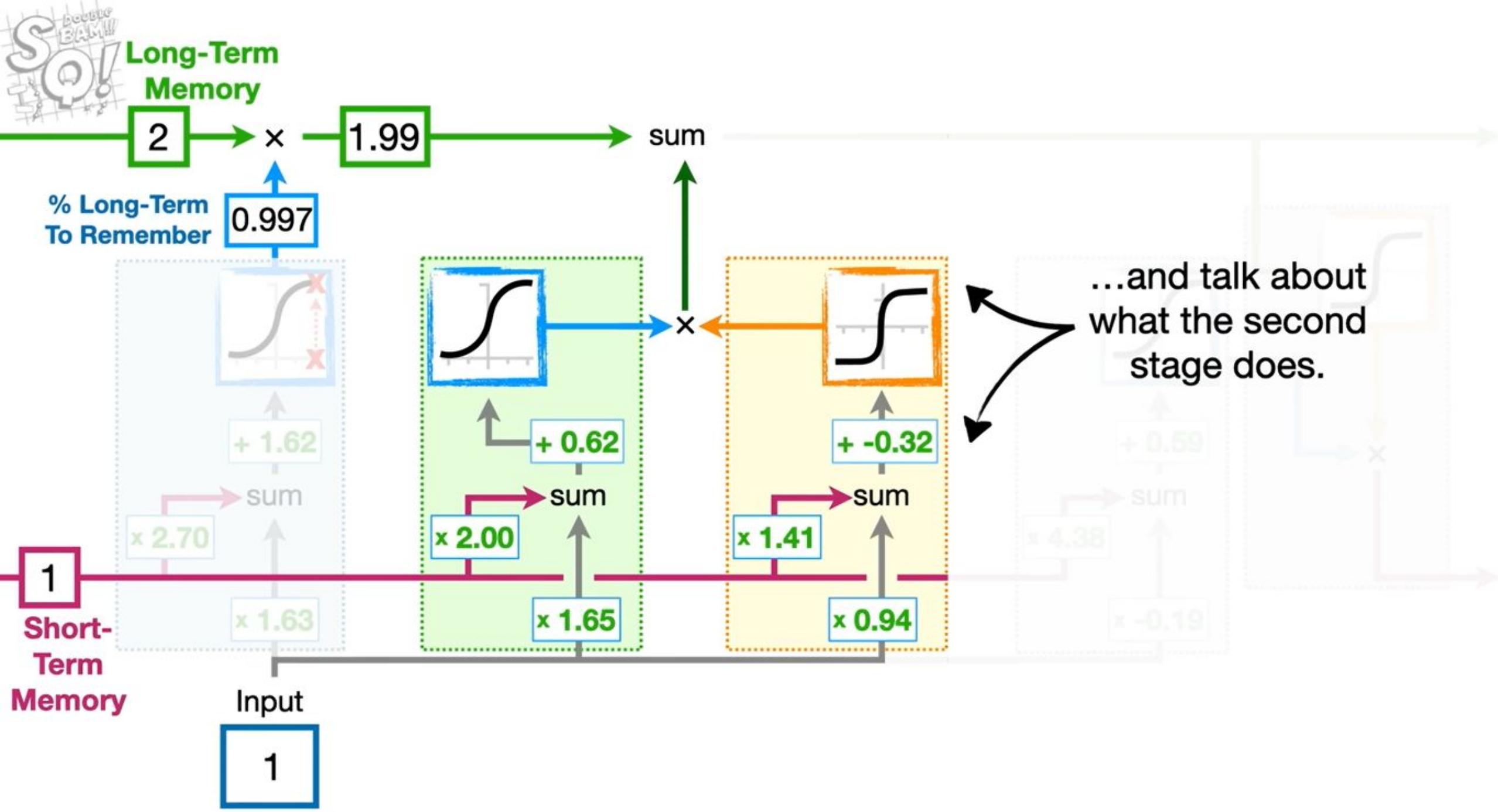


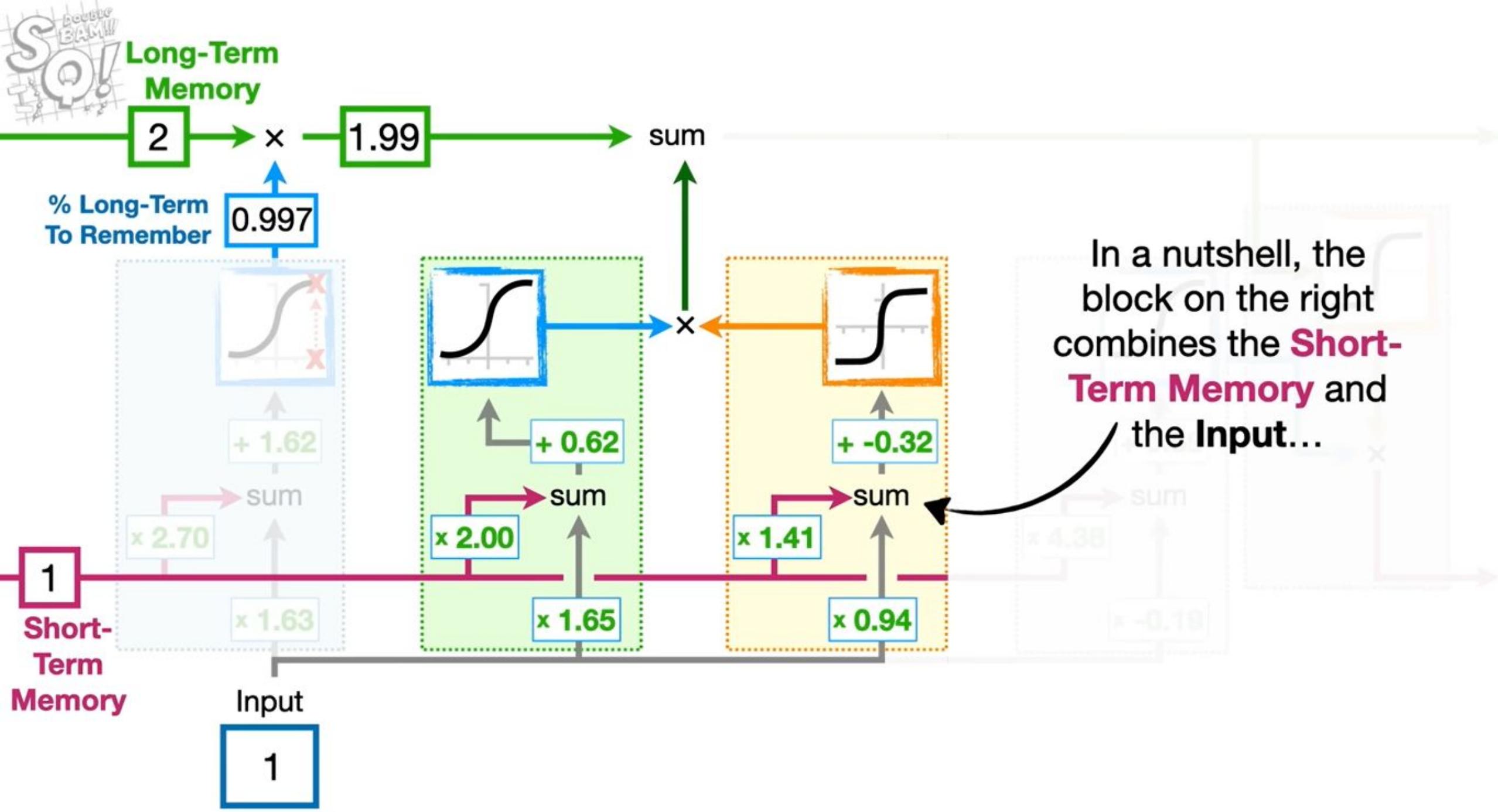


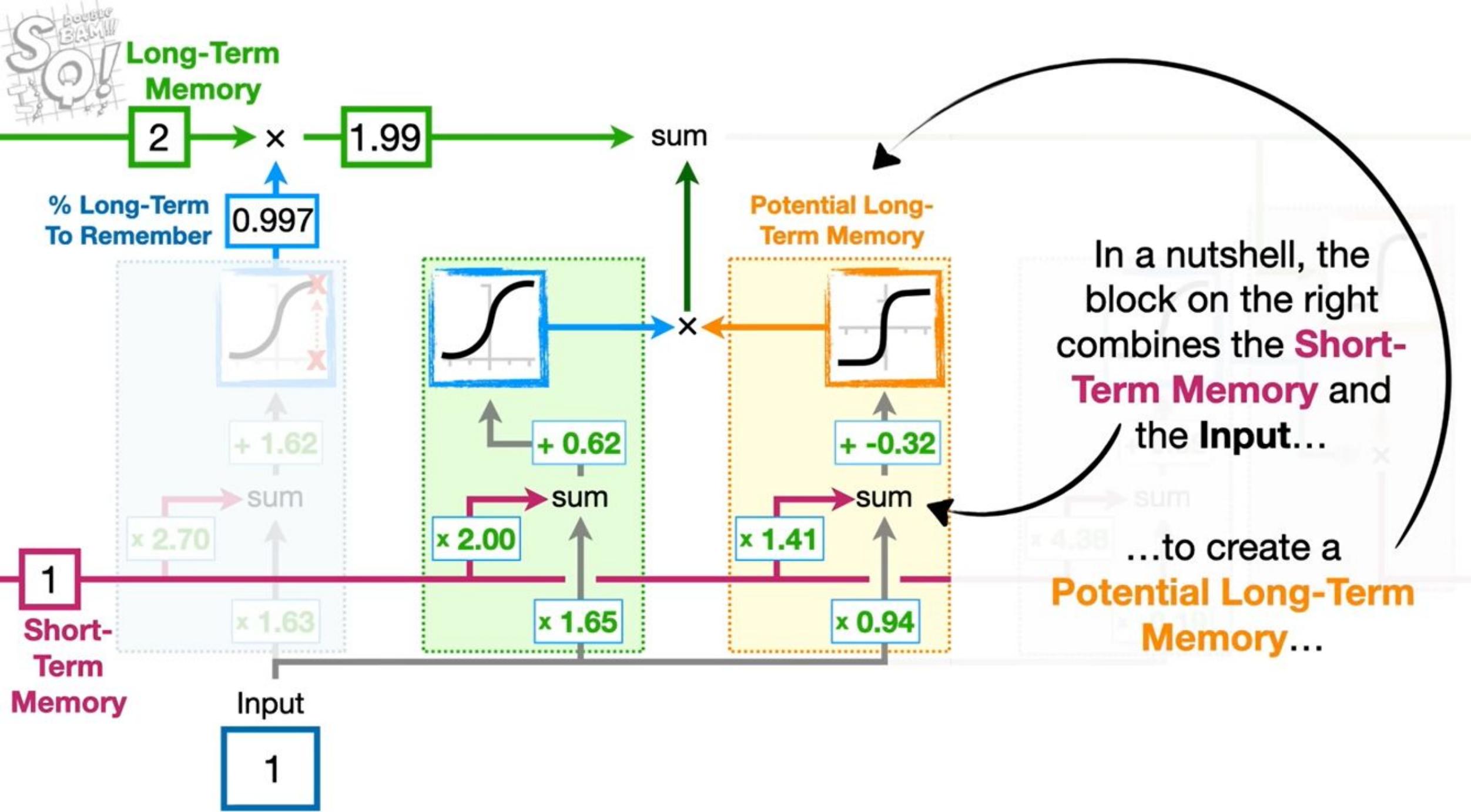


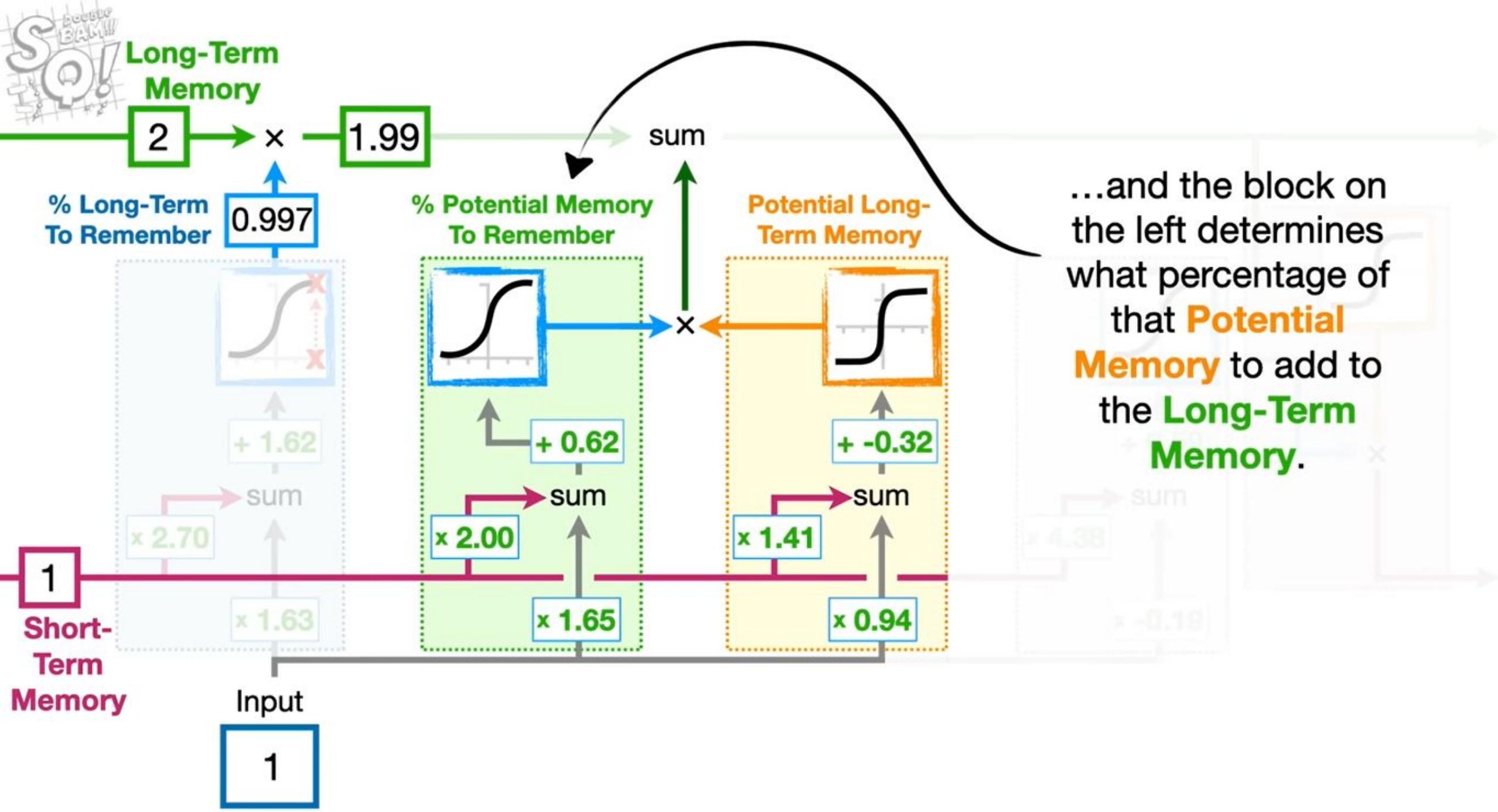


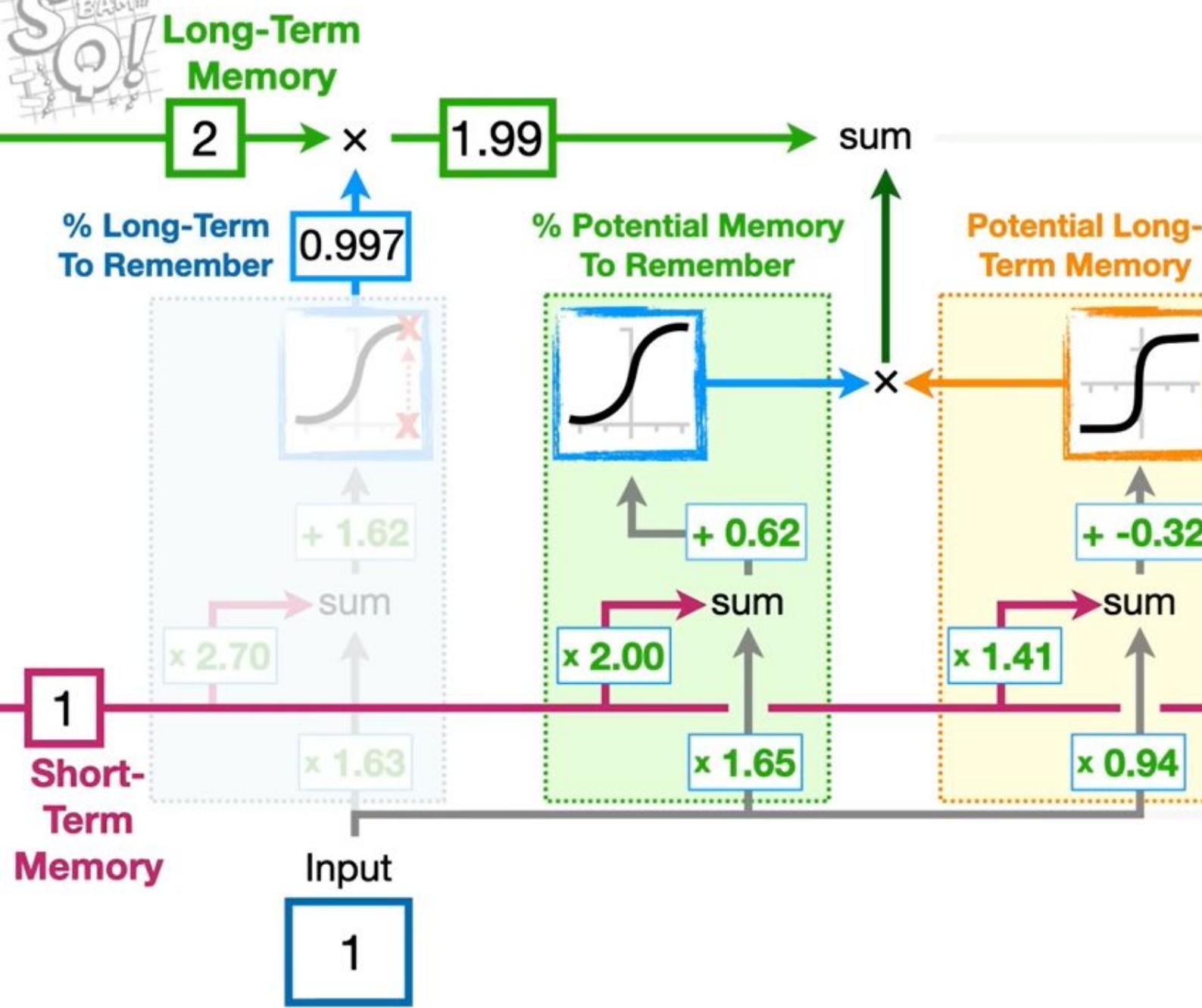




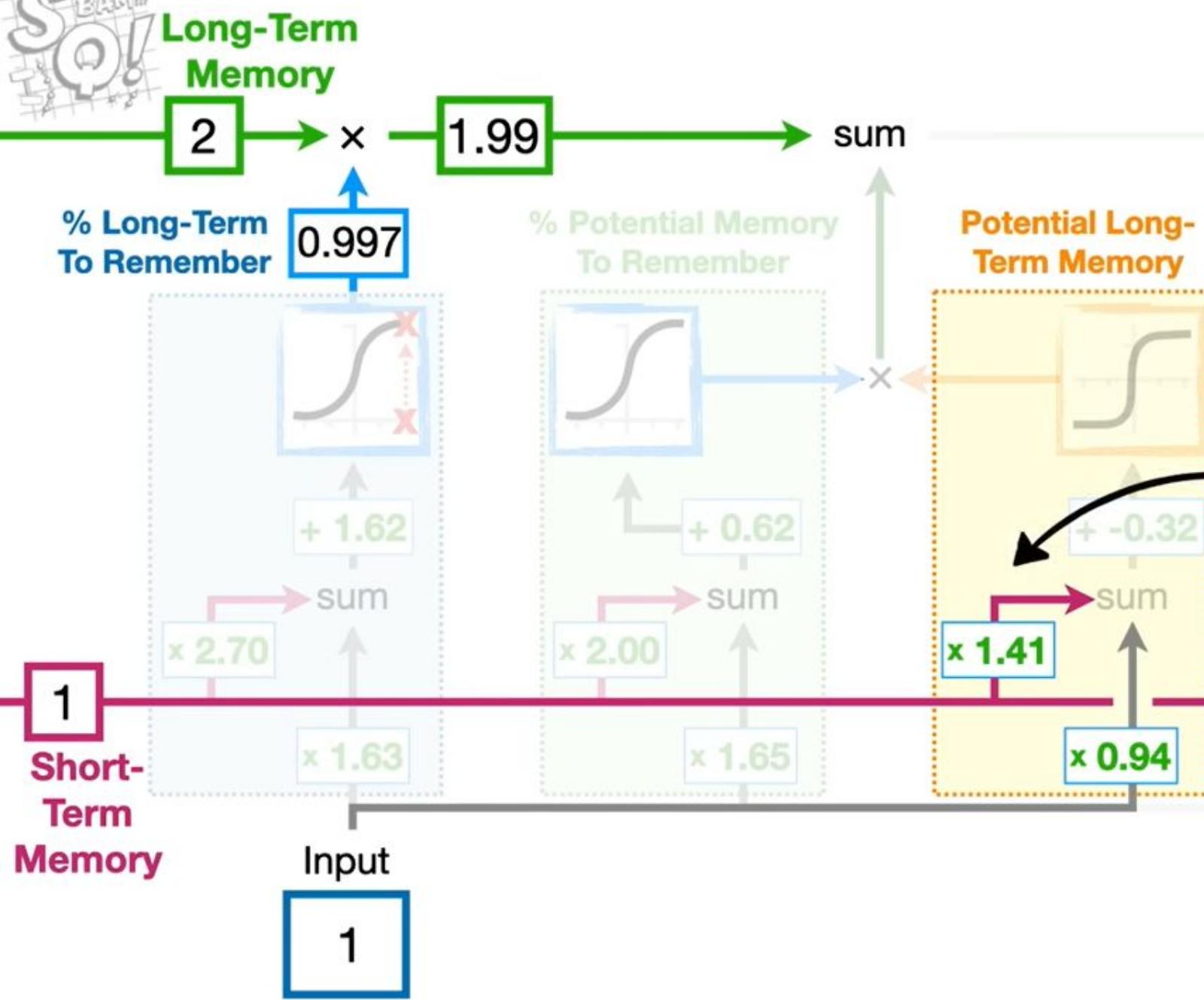








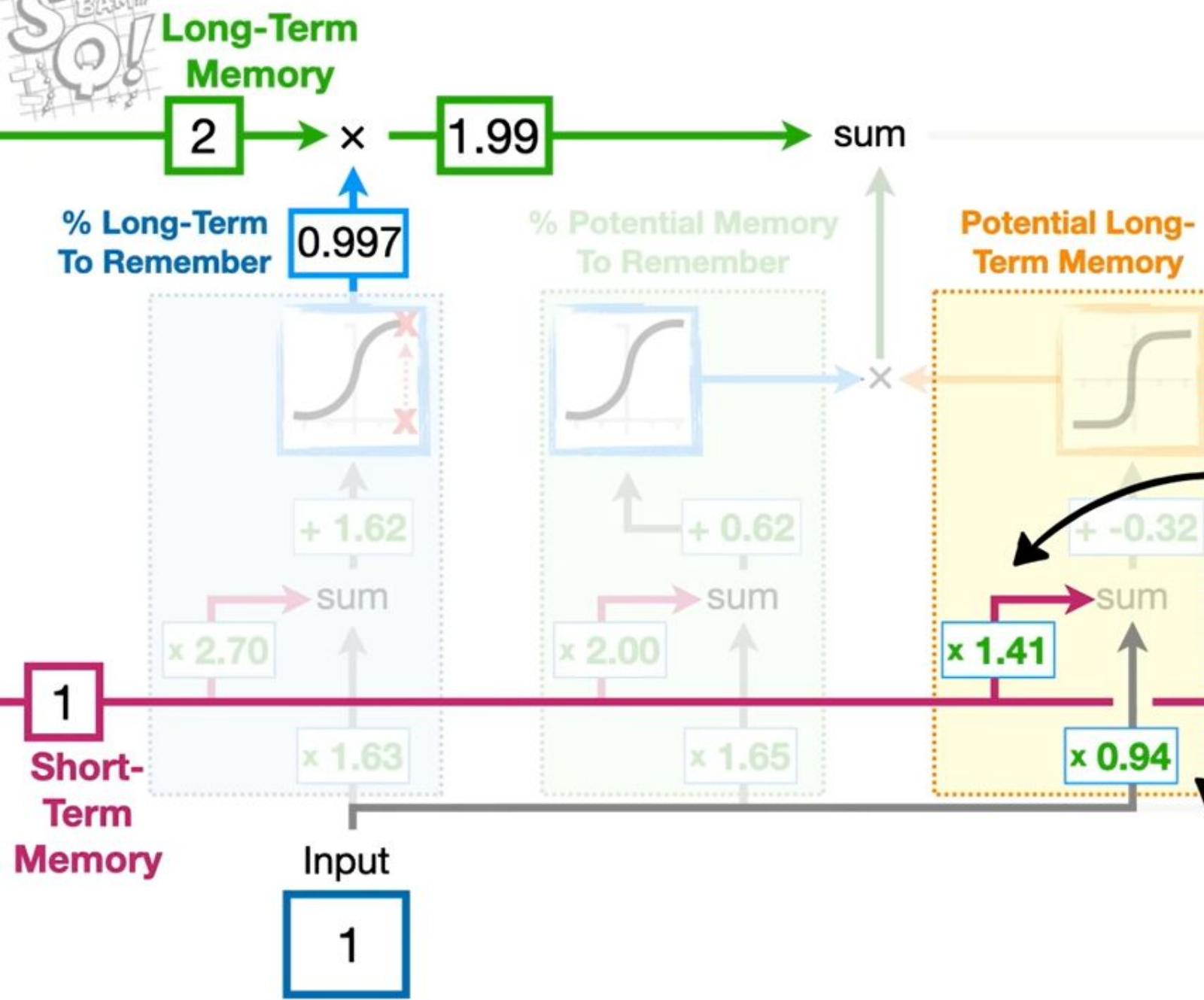
So let's plug the numbers in and do the math to see how a **Potential Memory** is created and how much of it is added to the **Long-Term Memory**.



Starting with the block furthest to the right, we multiply the **Short-Term Memory** and **Input** by their respective **Weights**.

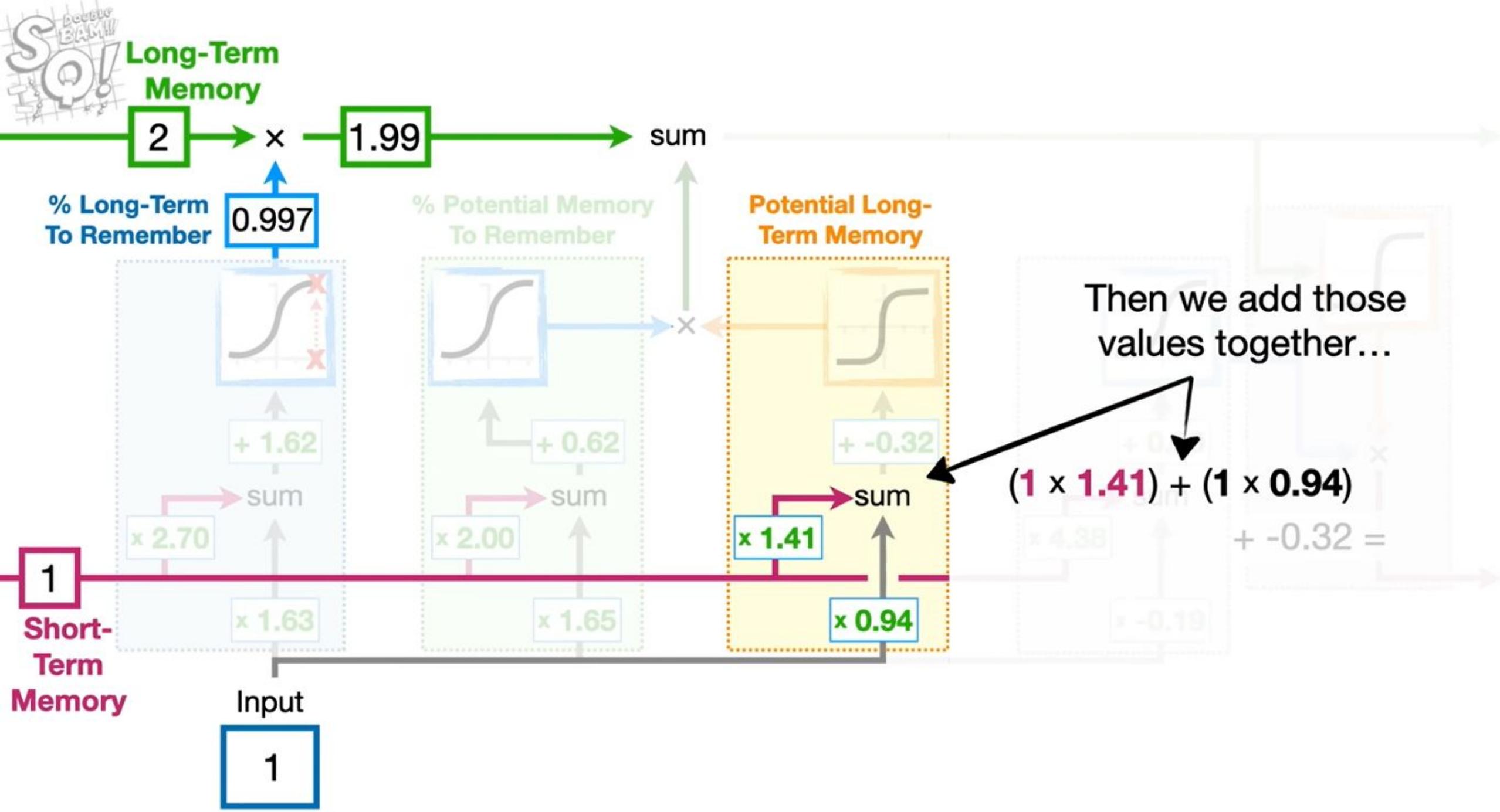
$$(1 \times 1.41) + (1 \times 0.94)$$

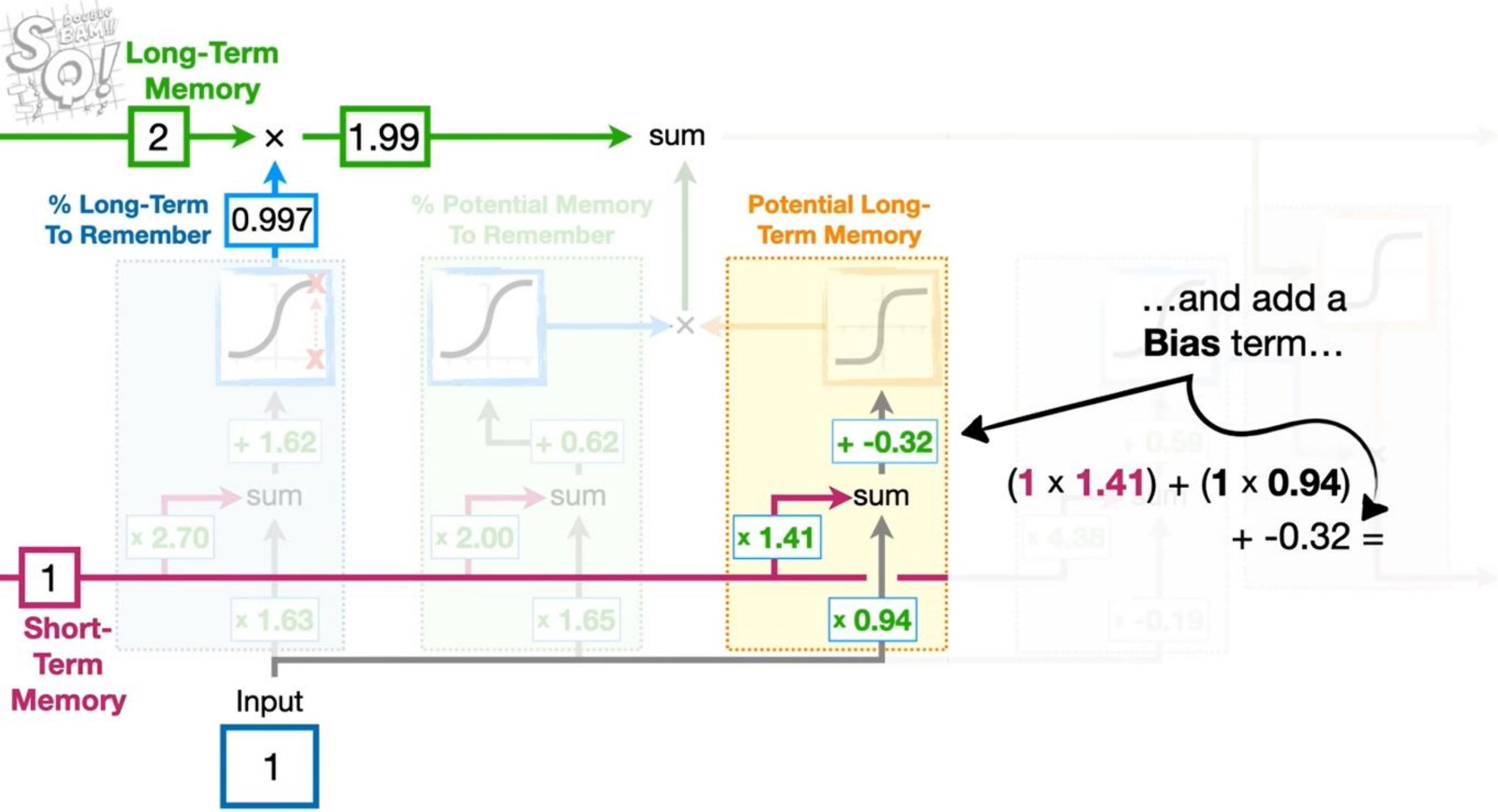
$$+ -0.32 =$$

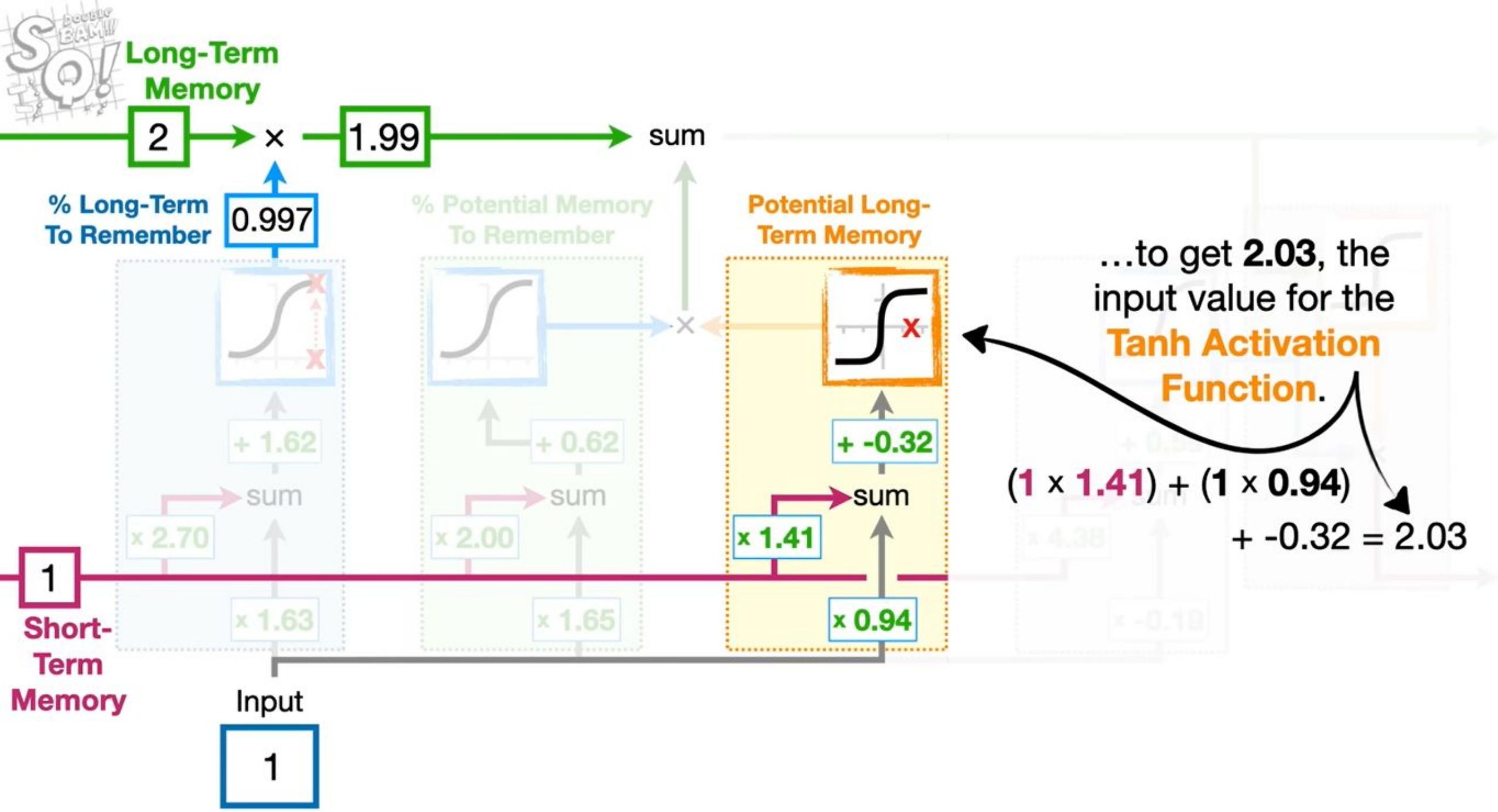


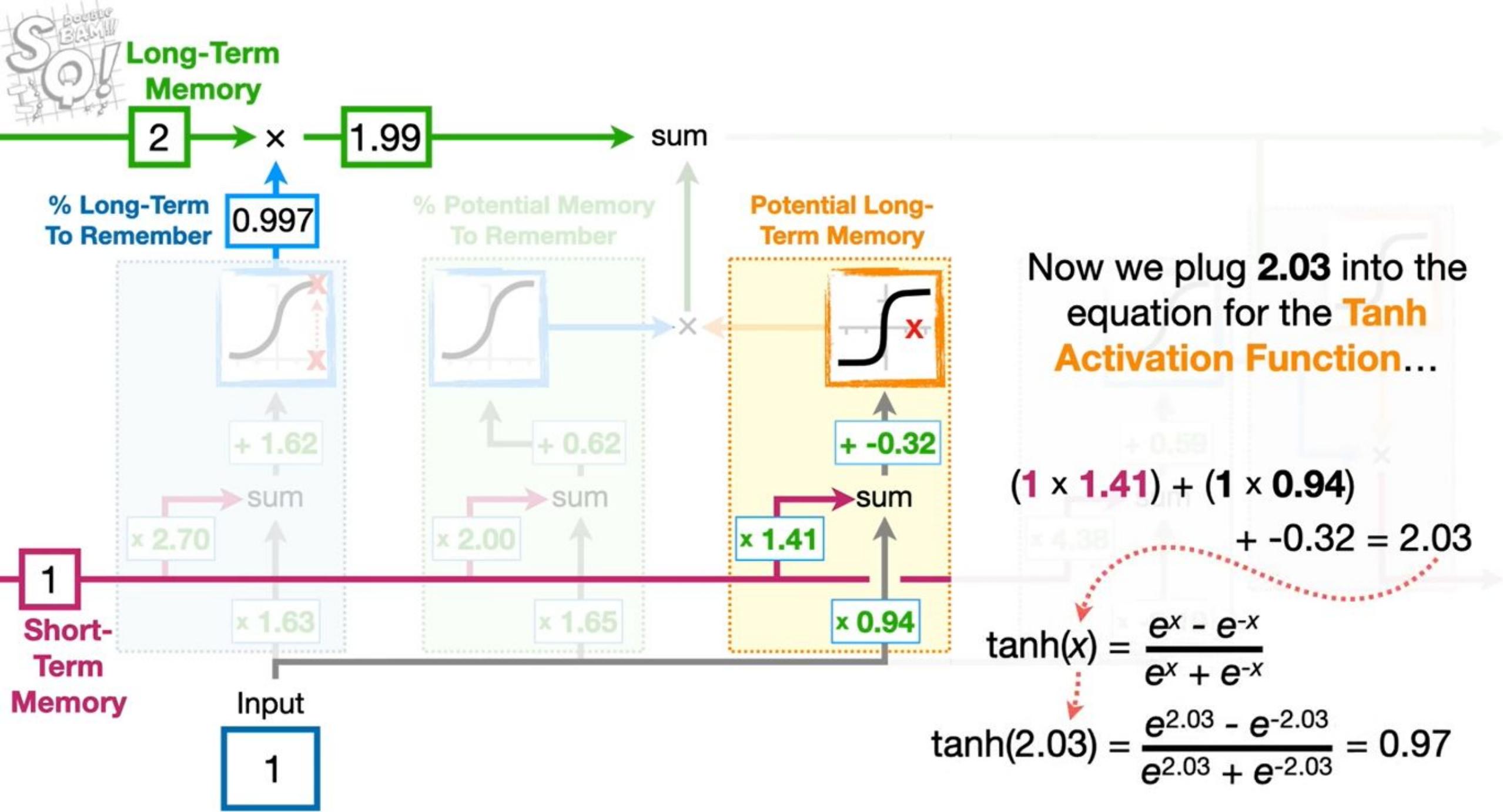
Starting with the block furthest to the right, we multiply the **Short-Term Memory** and **Input** by their respective **Weights**.

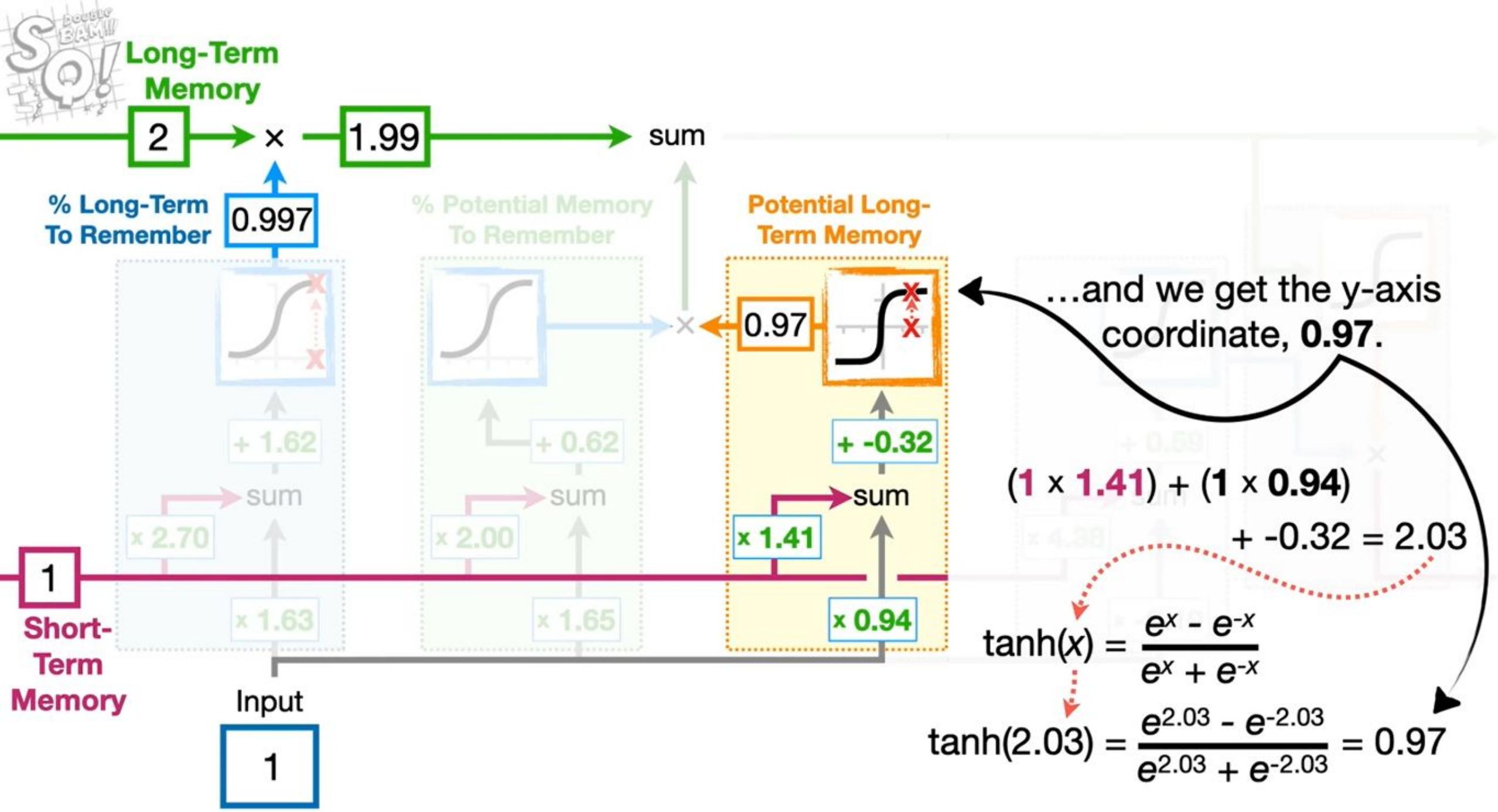
$$(1 \times 1.41) + (1 \times 0.94) + 0.32 =$$

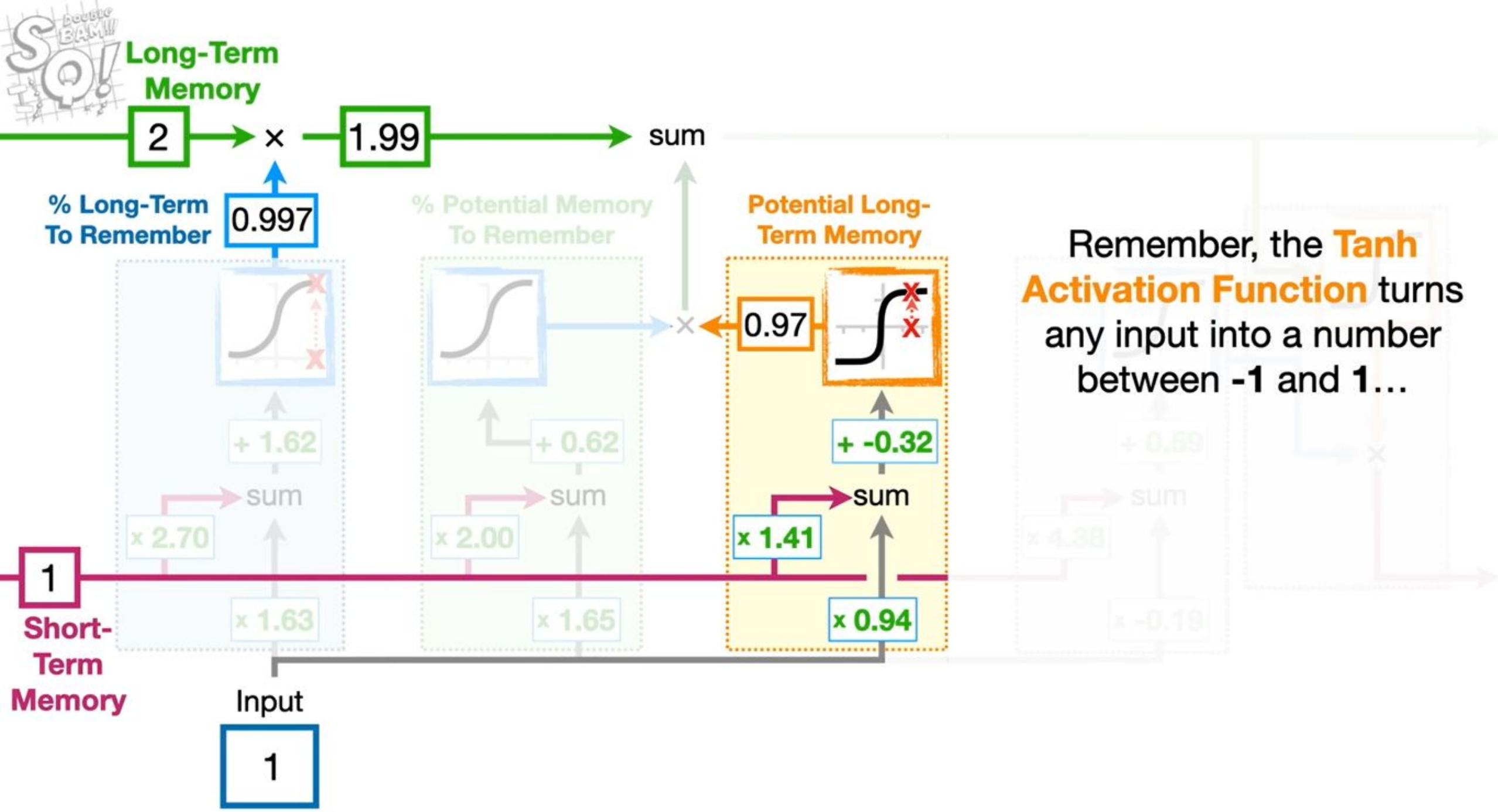


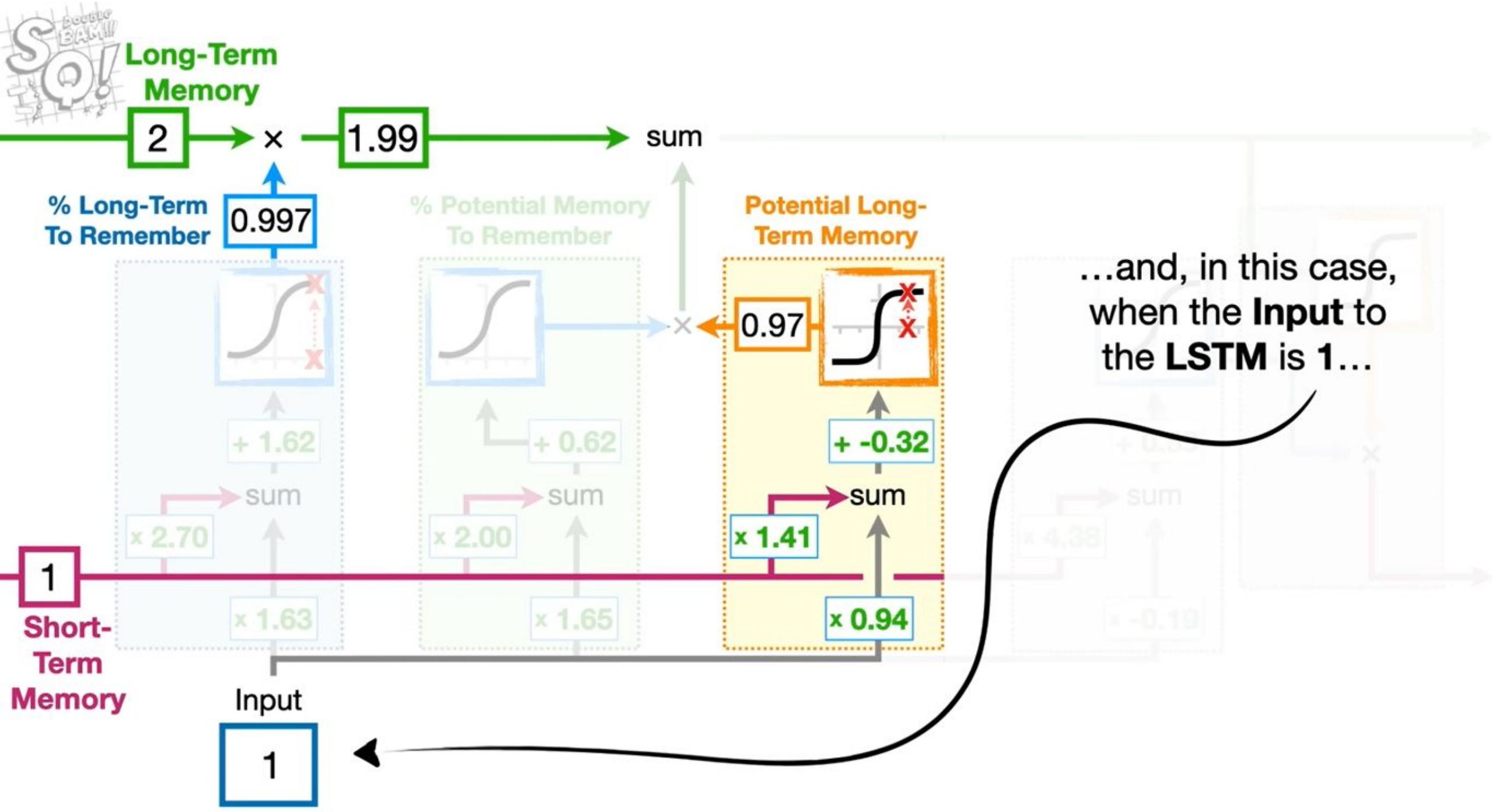


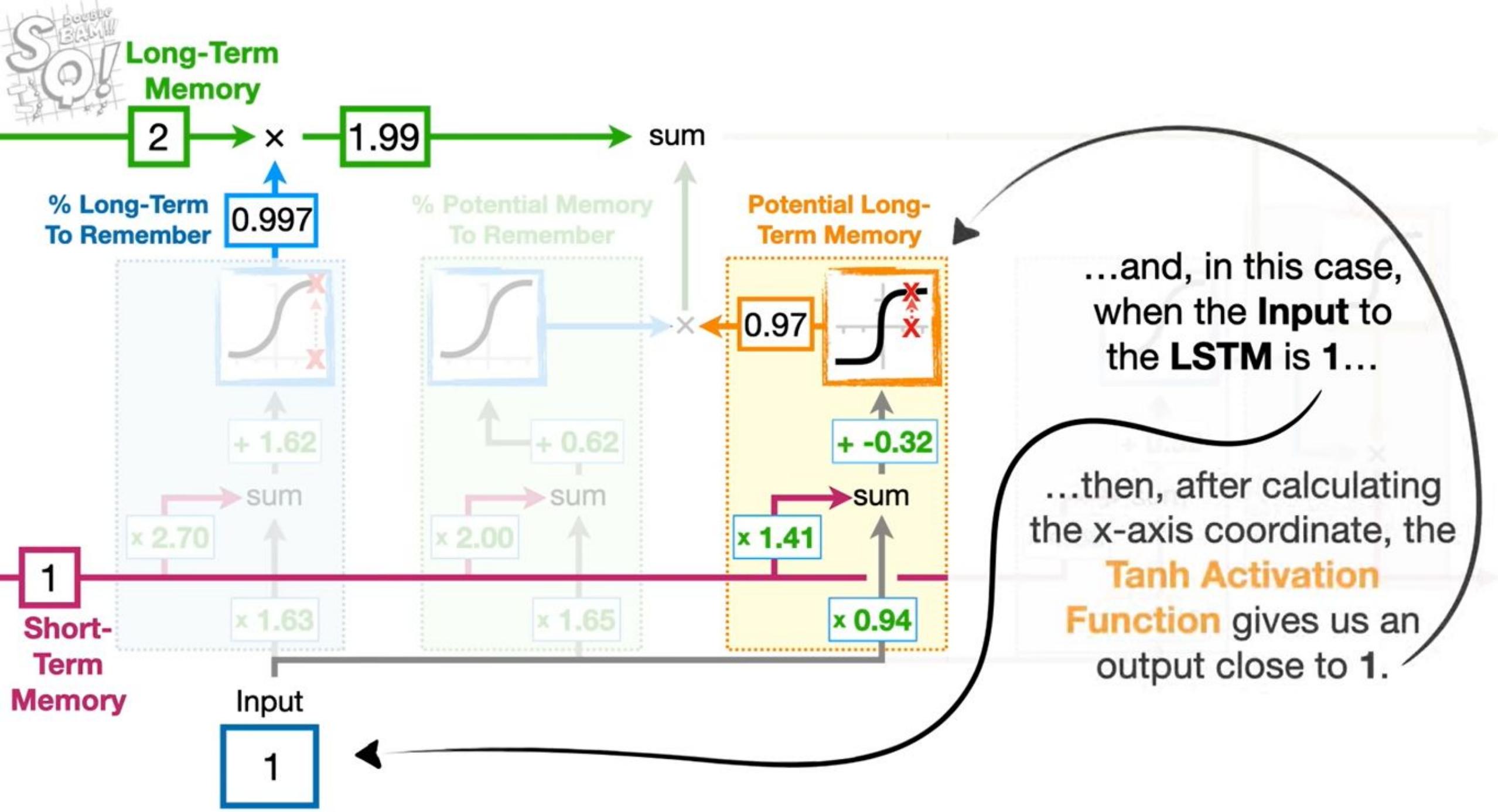


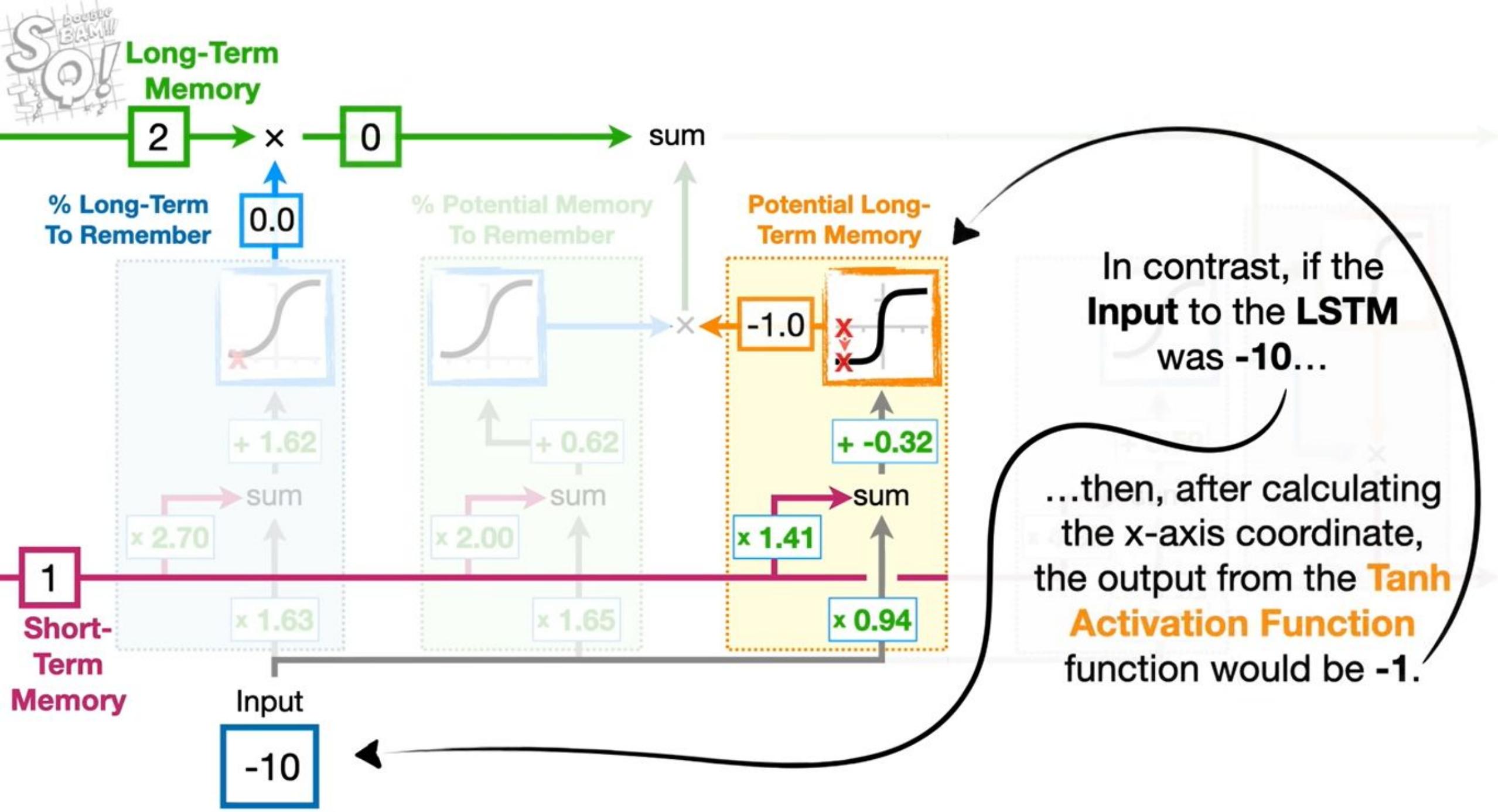


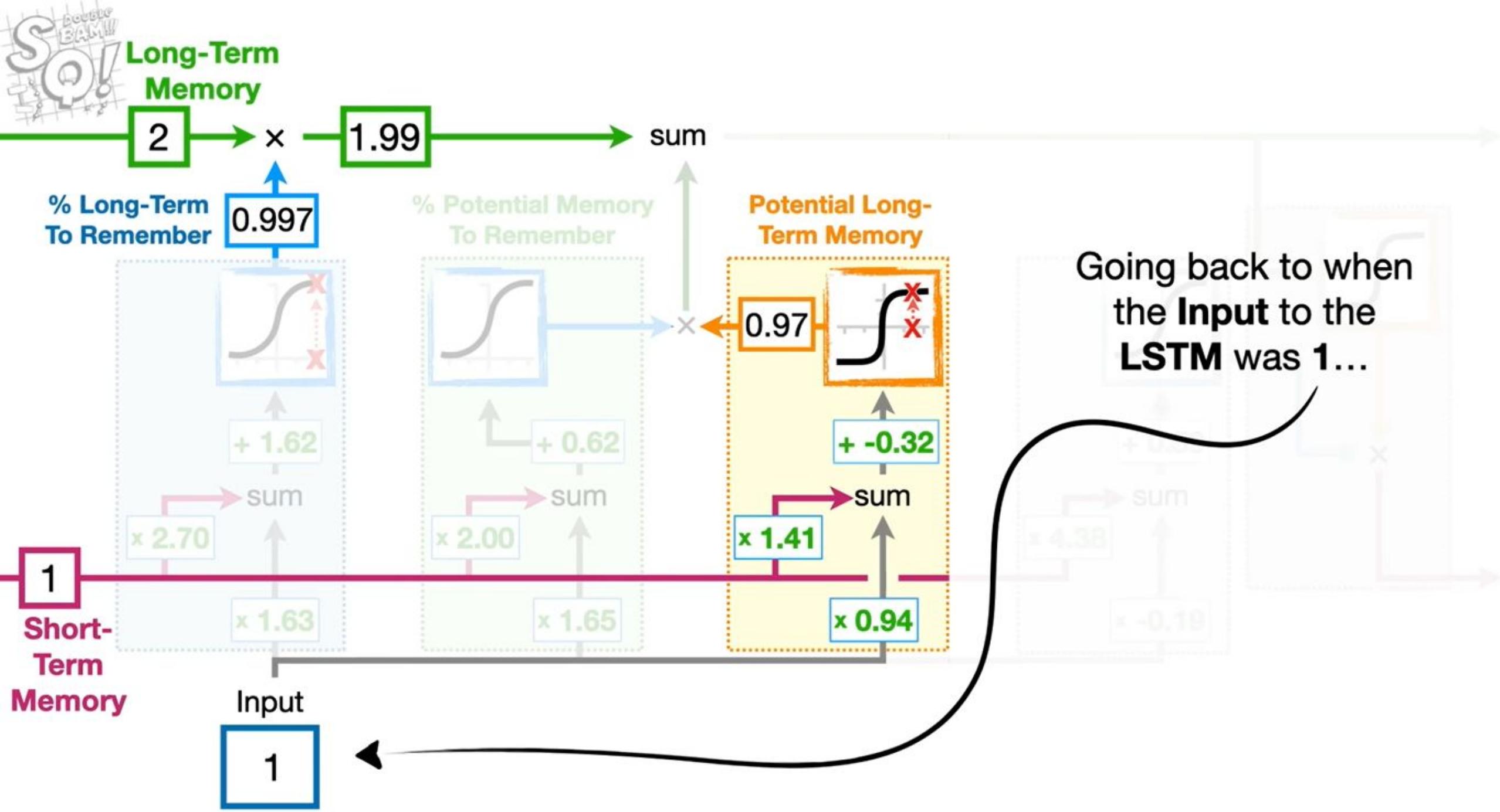




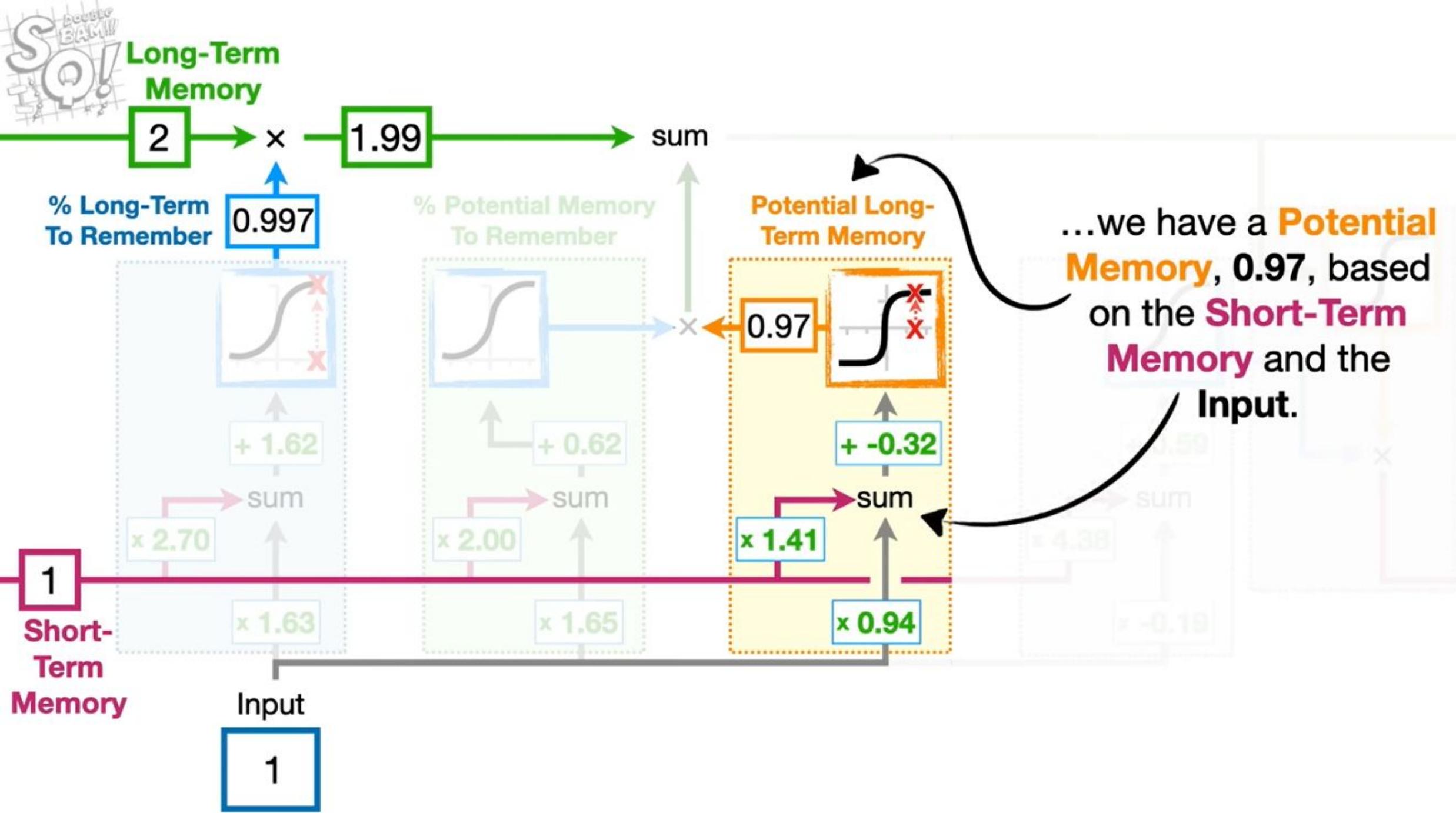


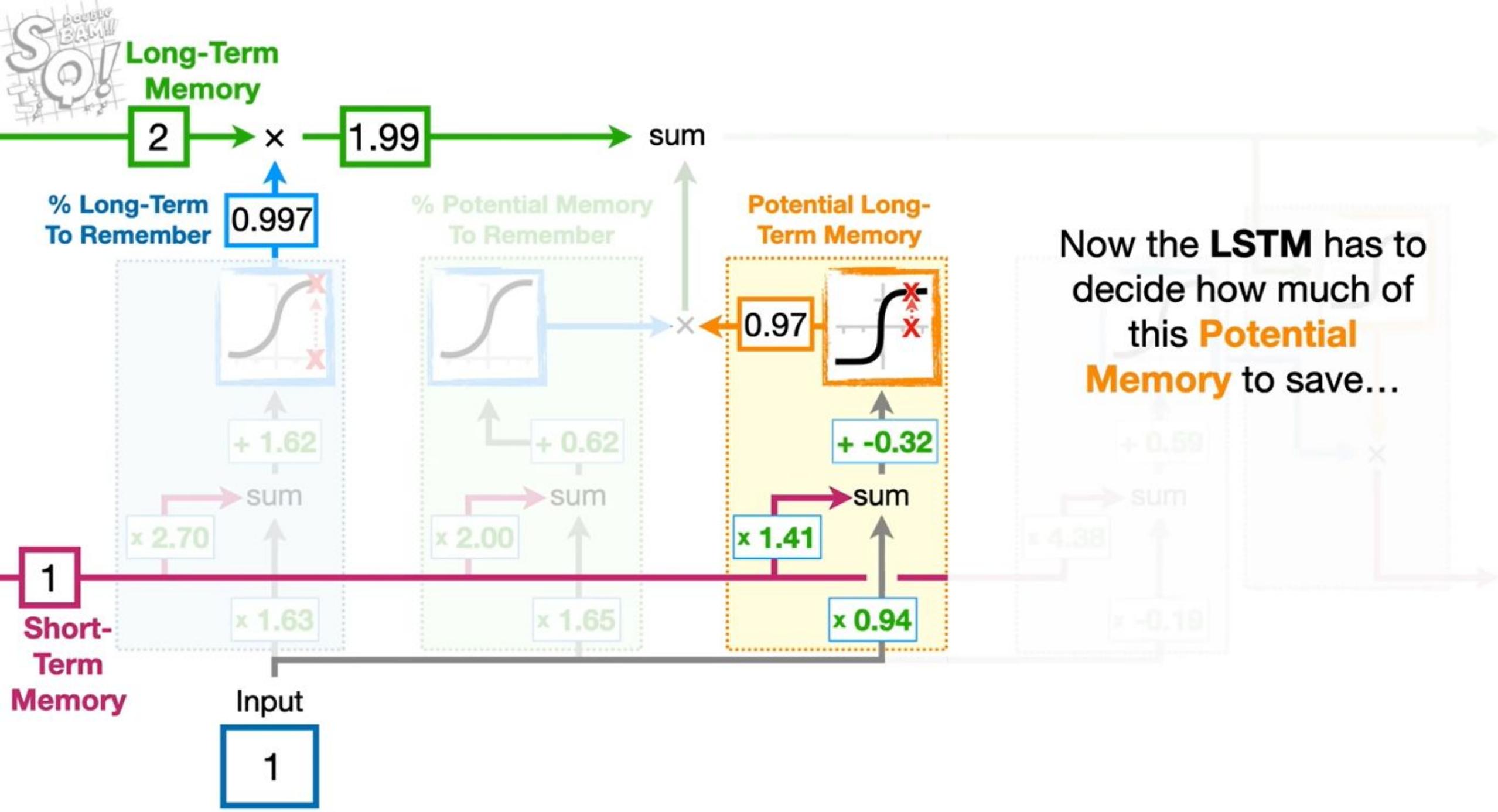


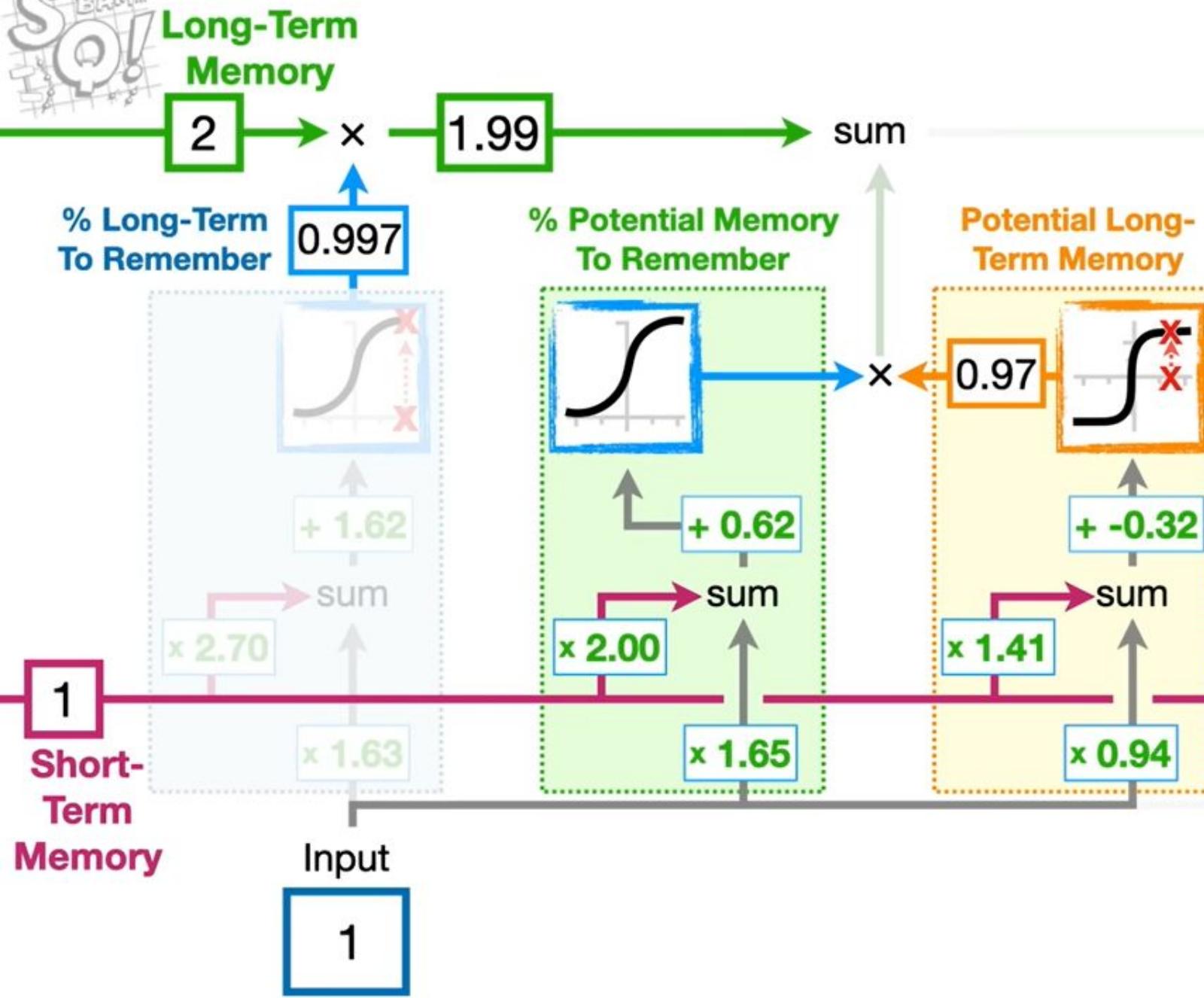




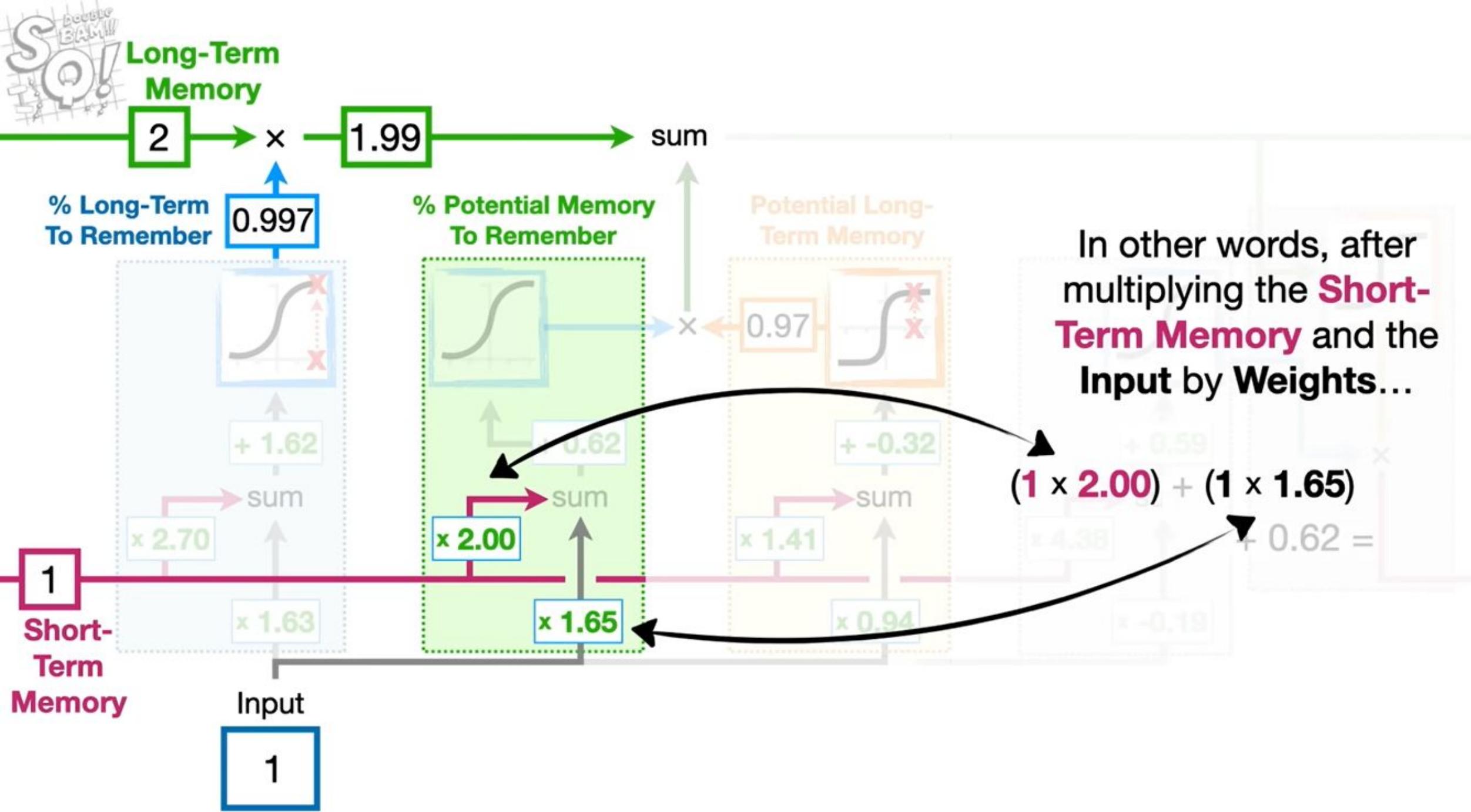
Going back to when
the **Input** to the
LSTM was 1...

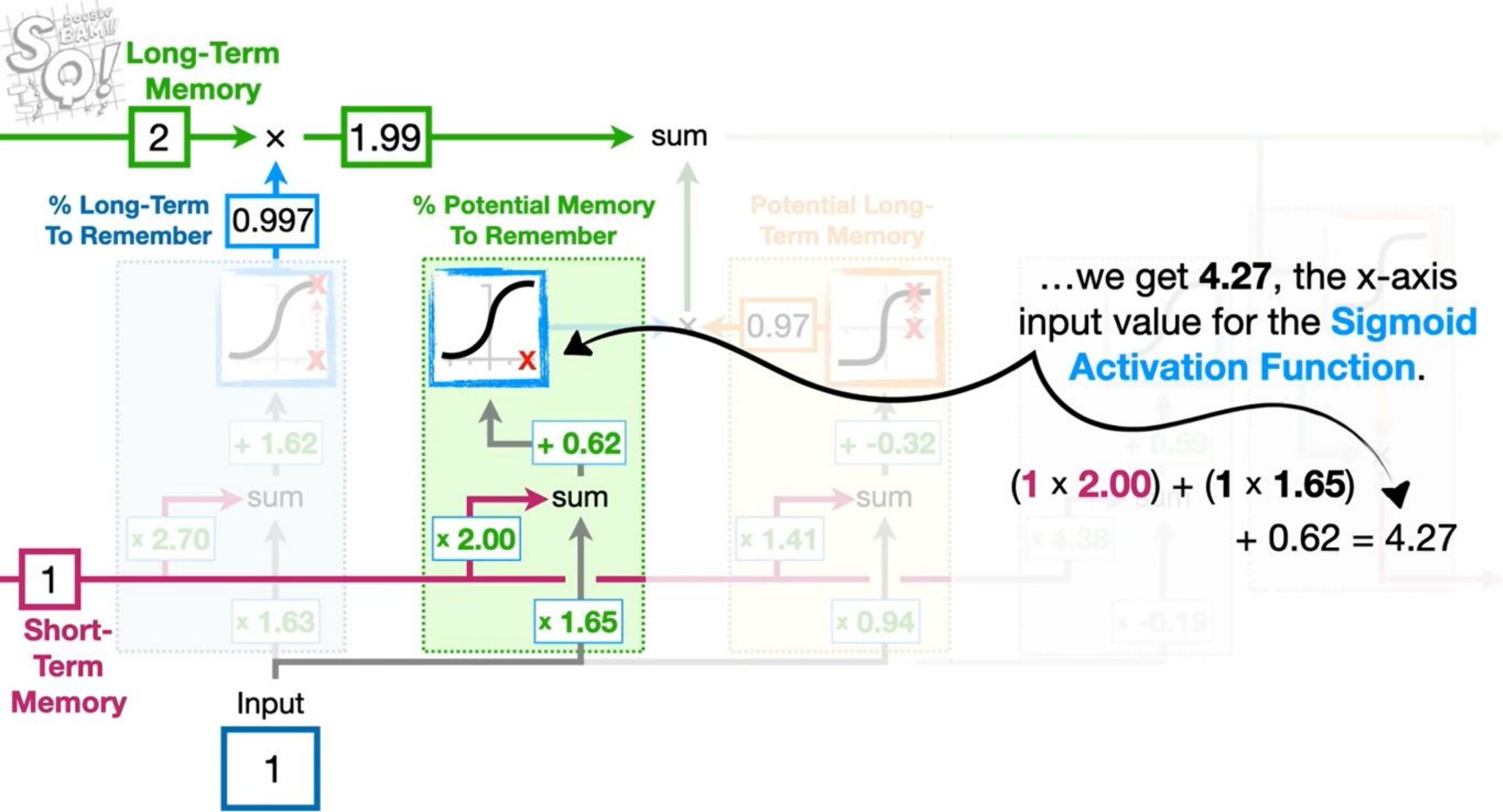


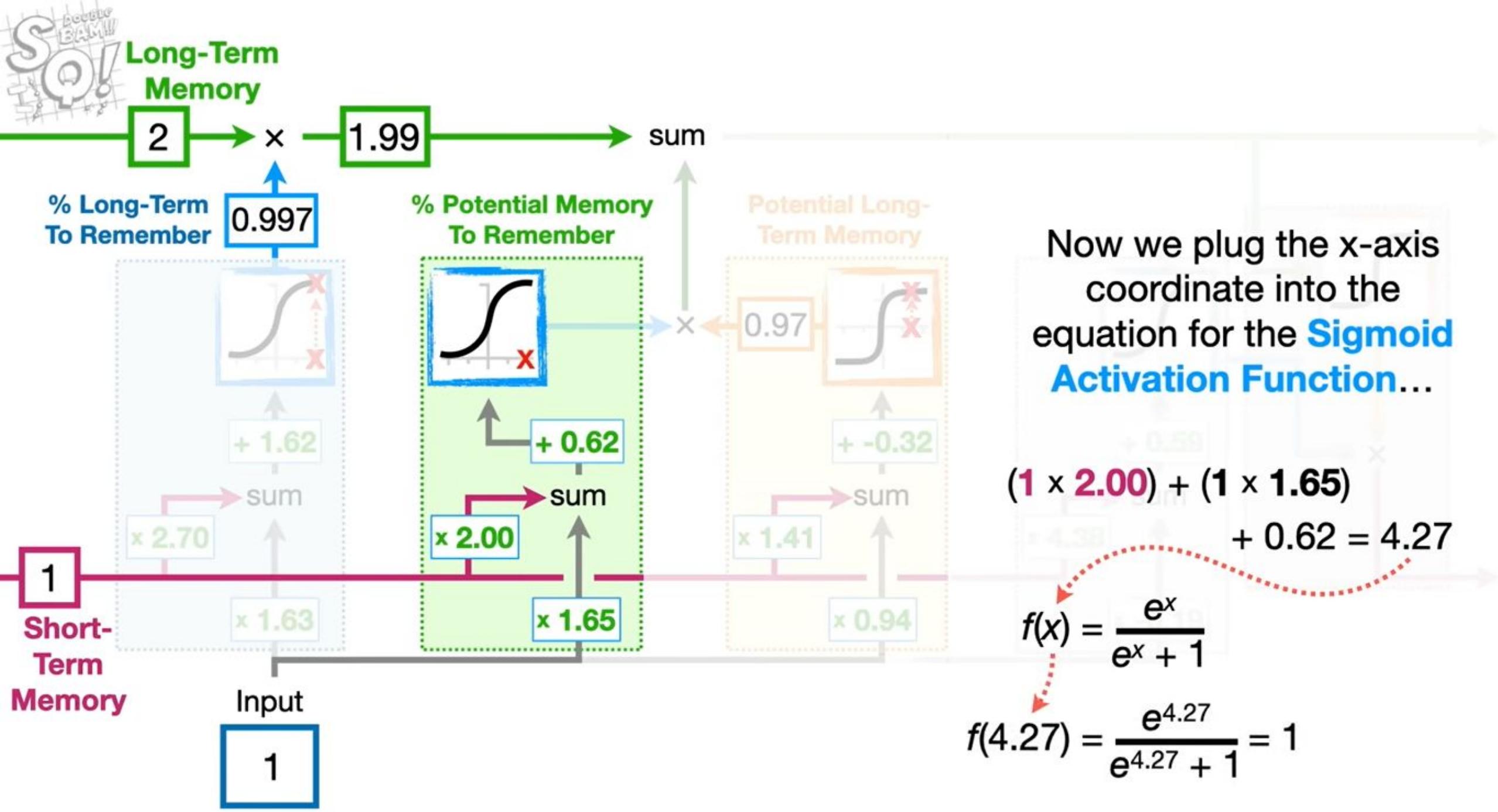


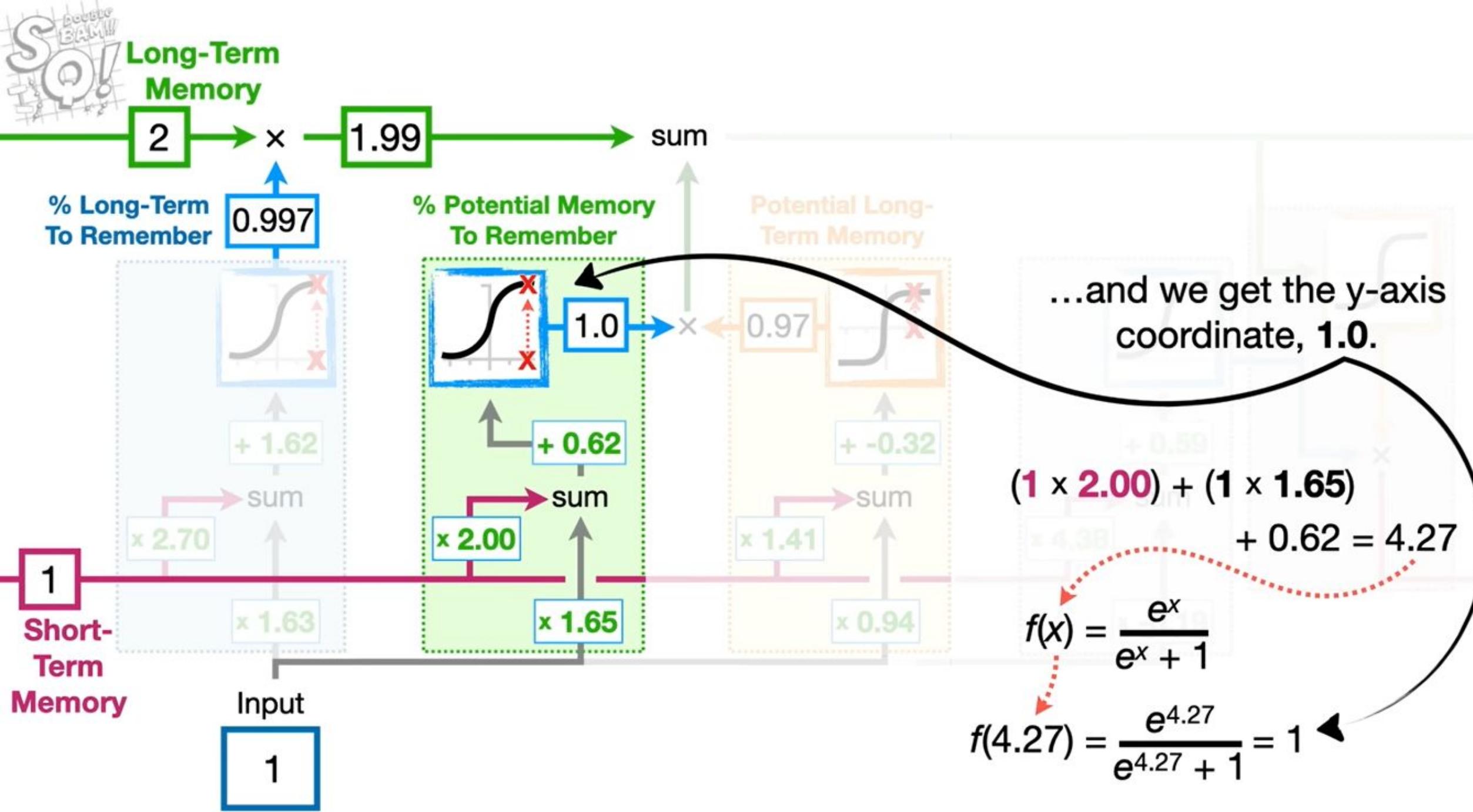


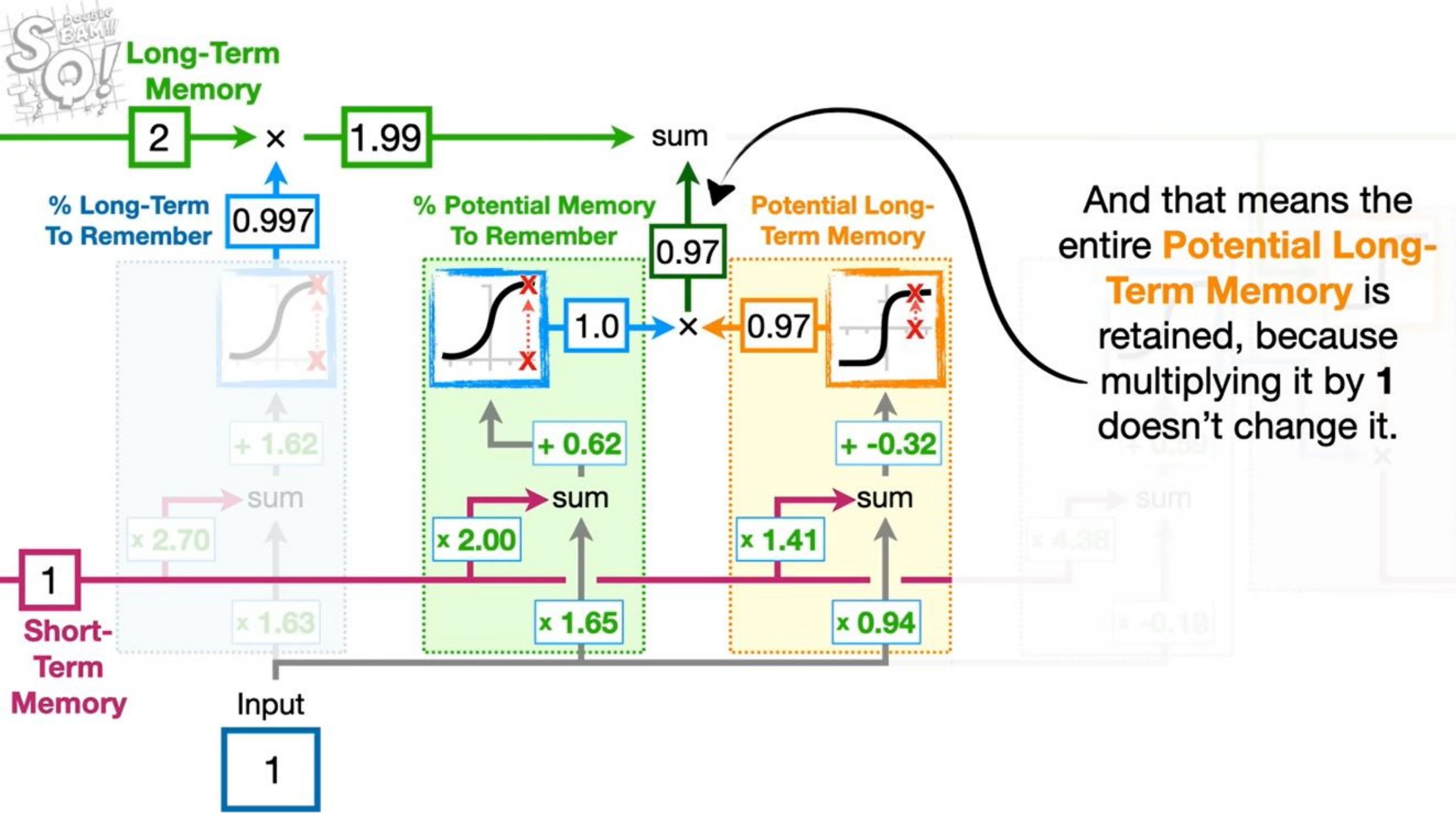
...and this is done using the exact same method we used earlier, when we determined what percentage of the **Long-Term Memory** to remember.

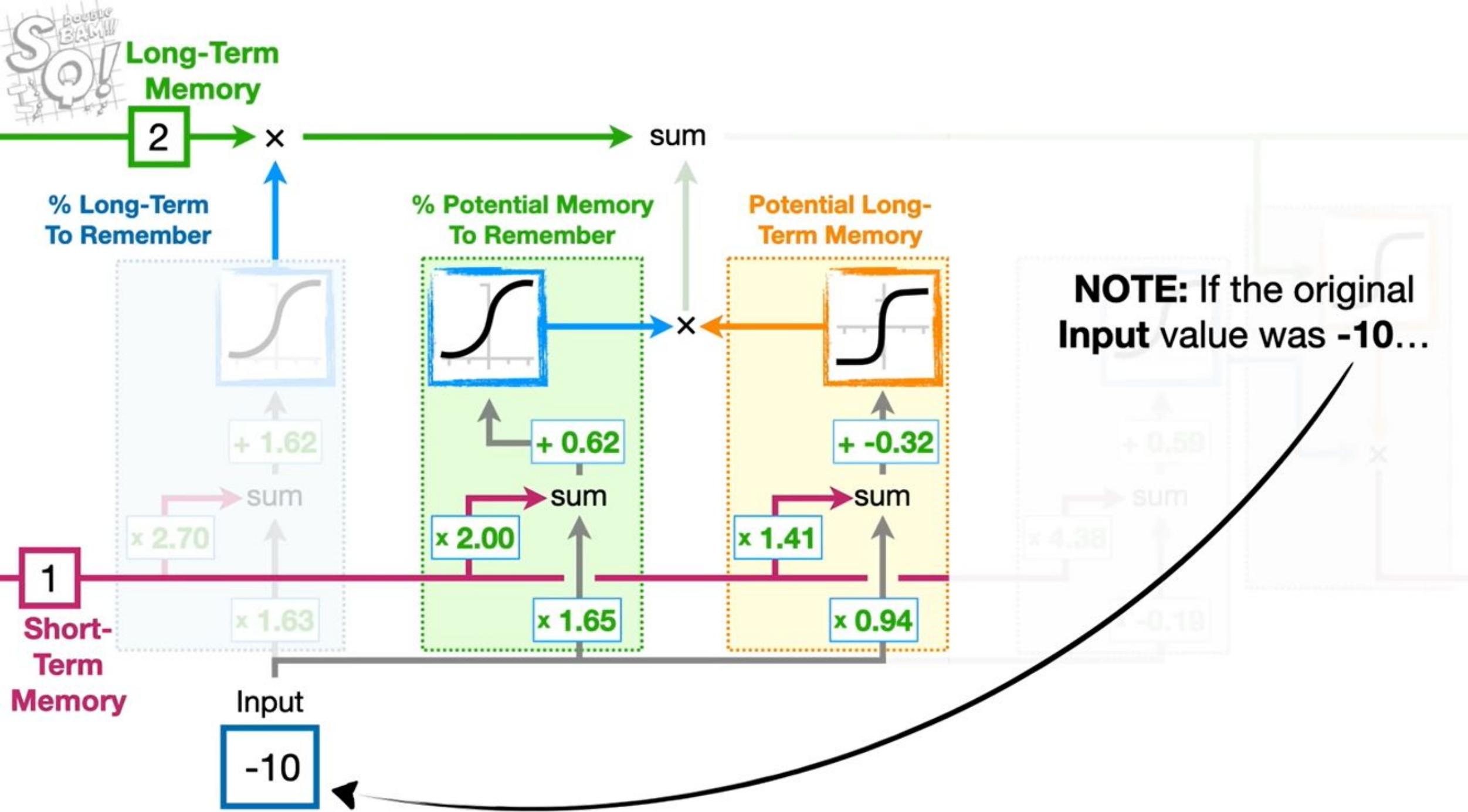


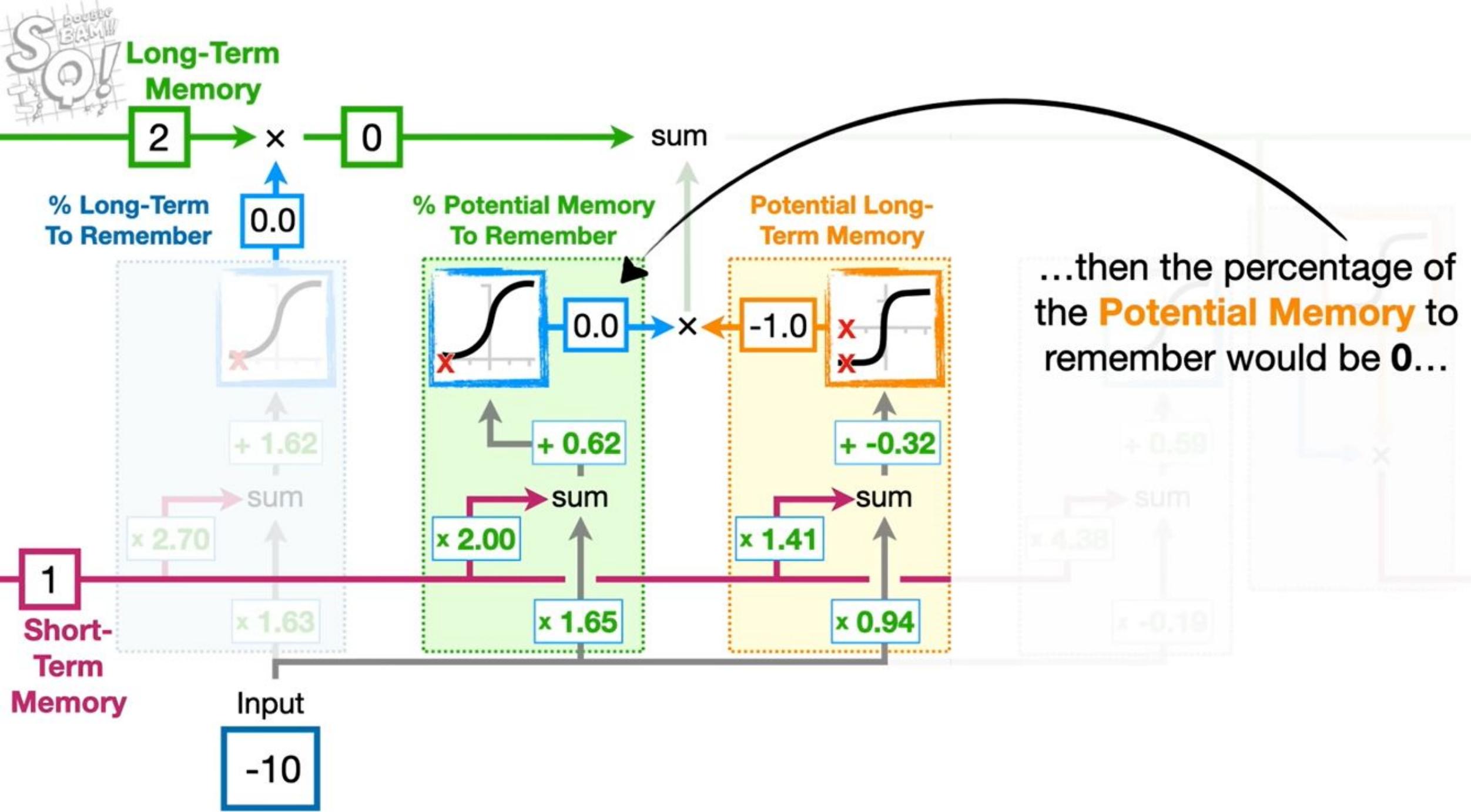






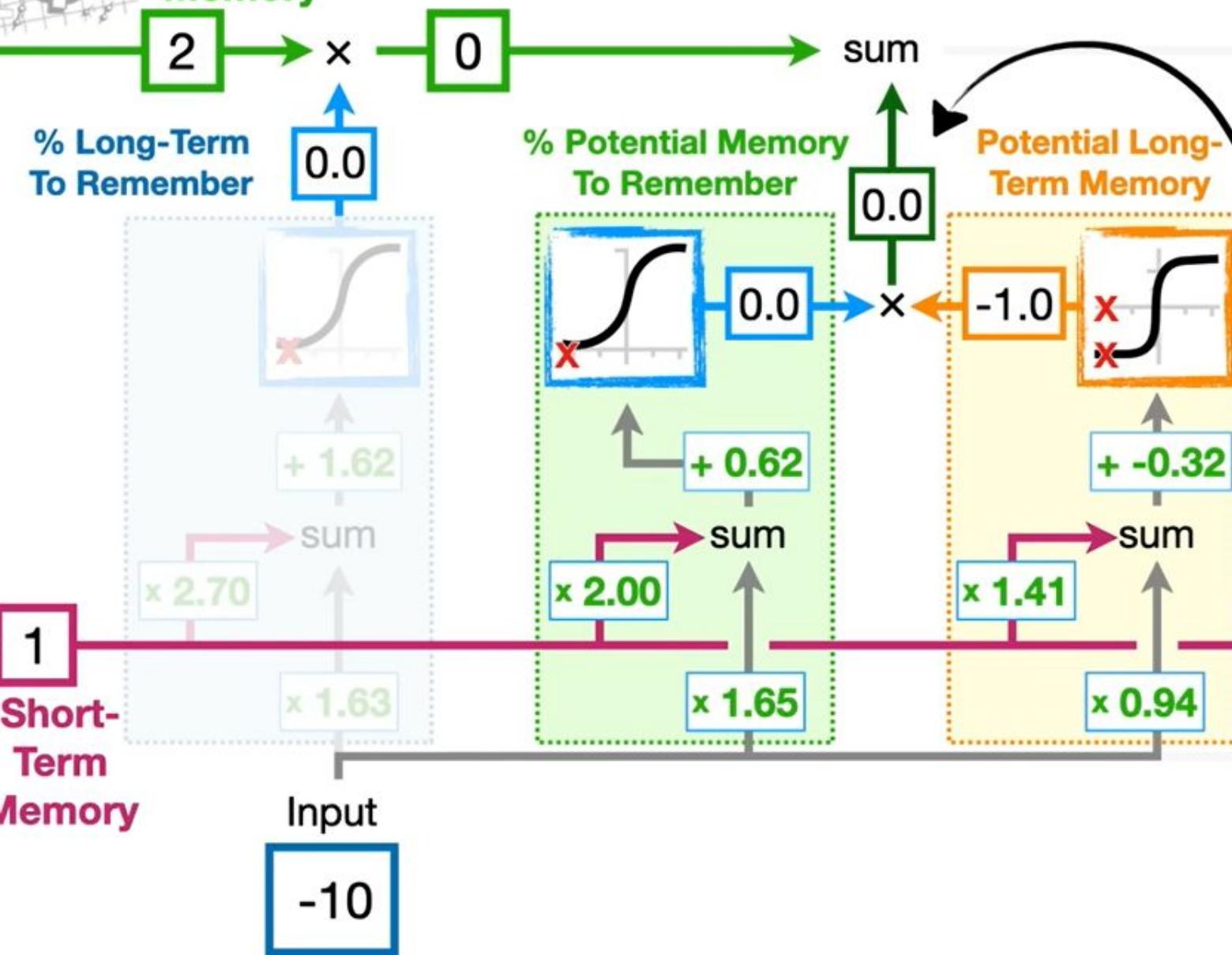




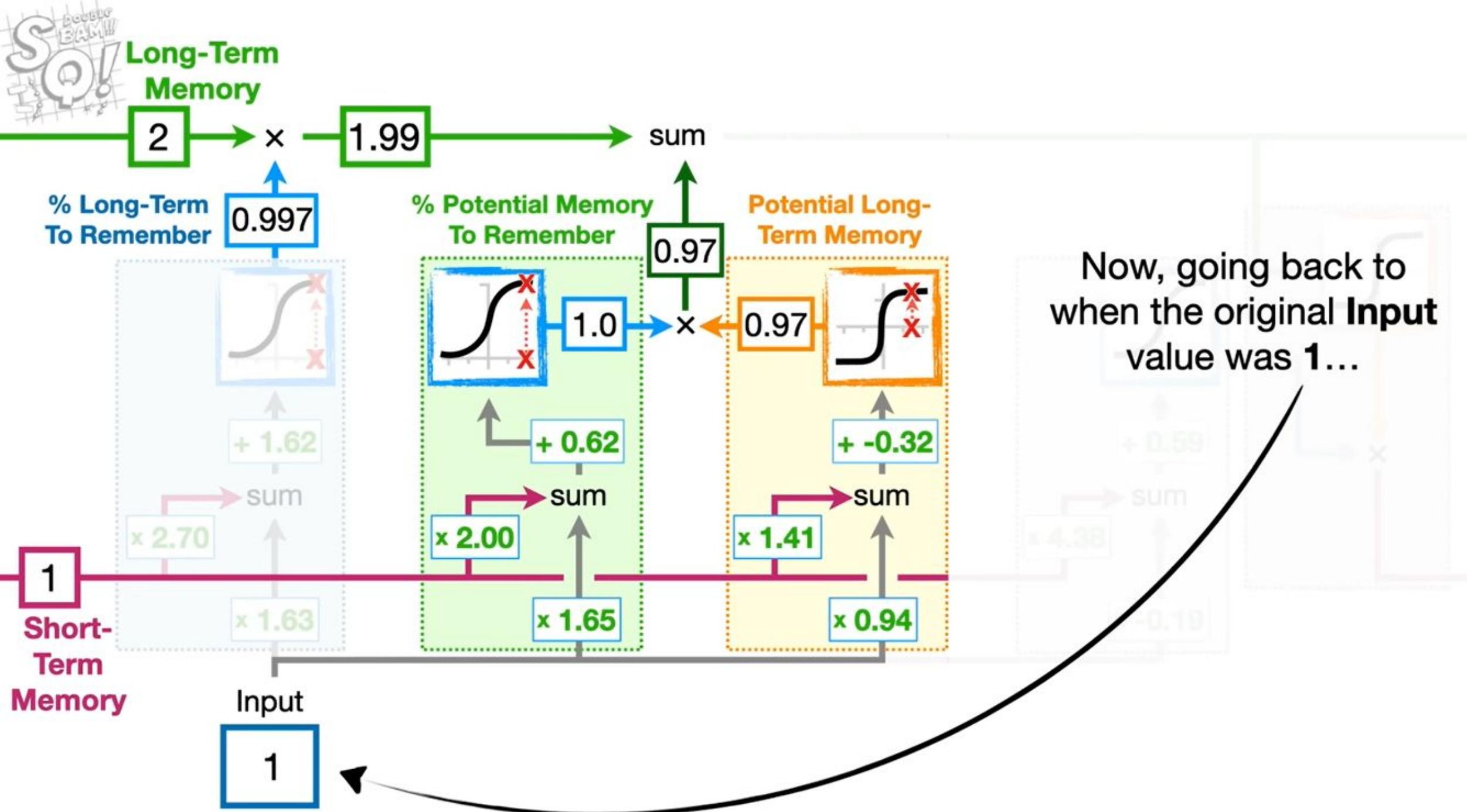


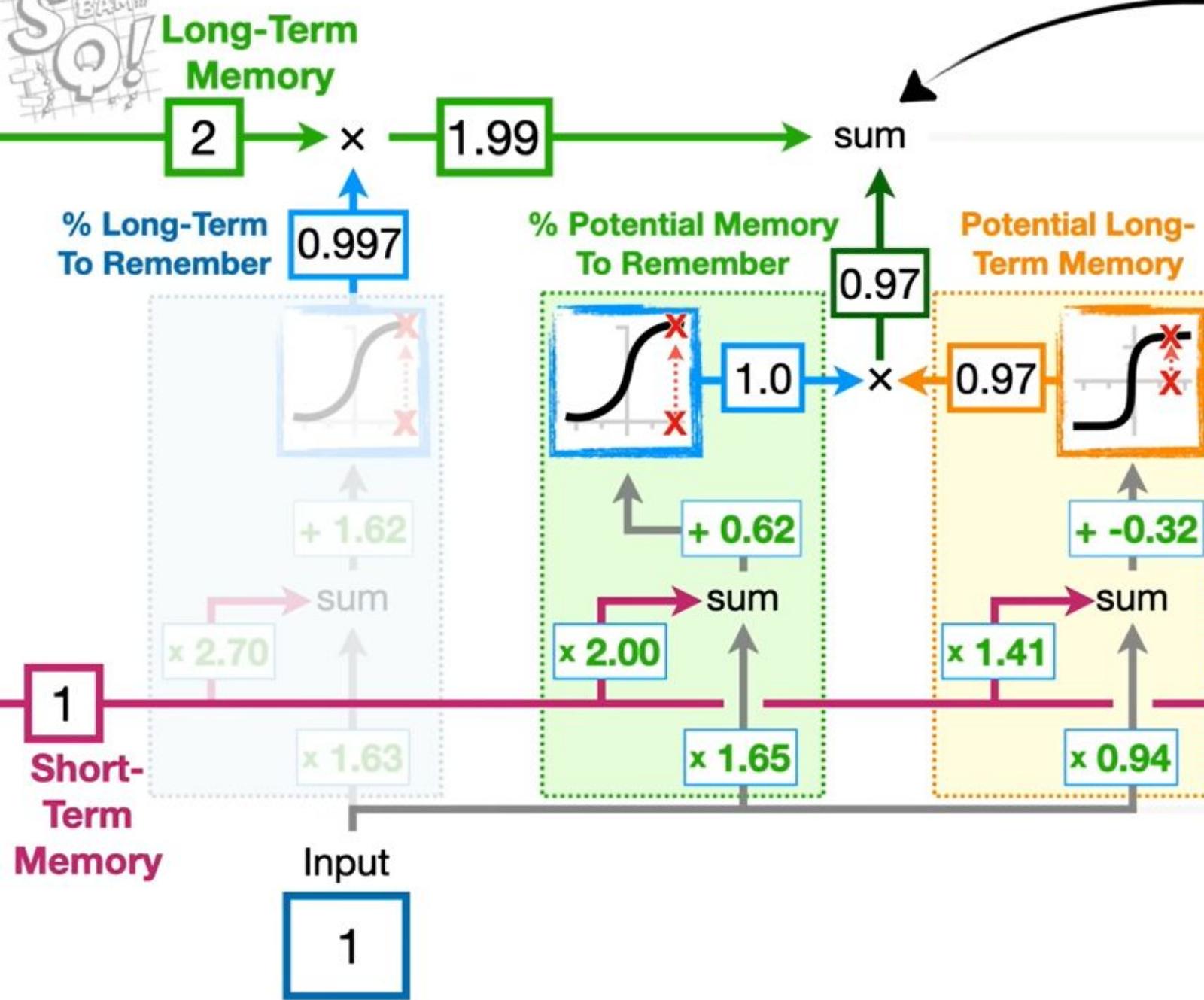


Long-Term Memory

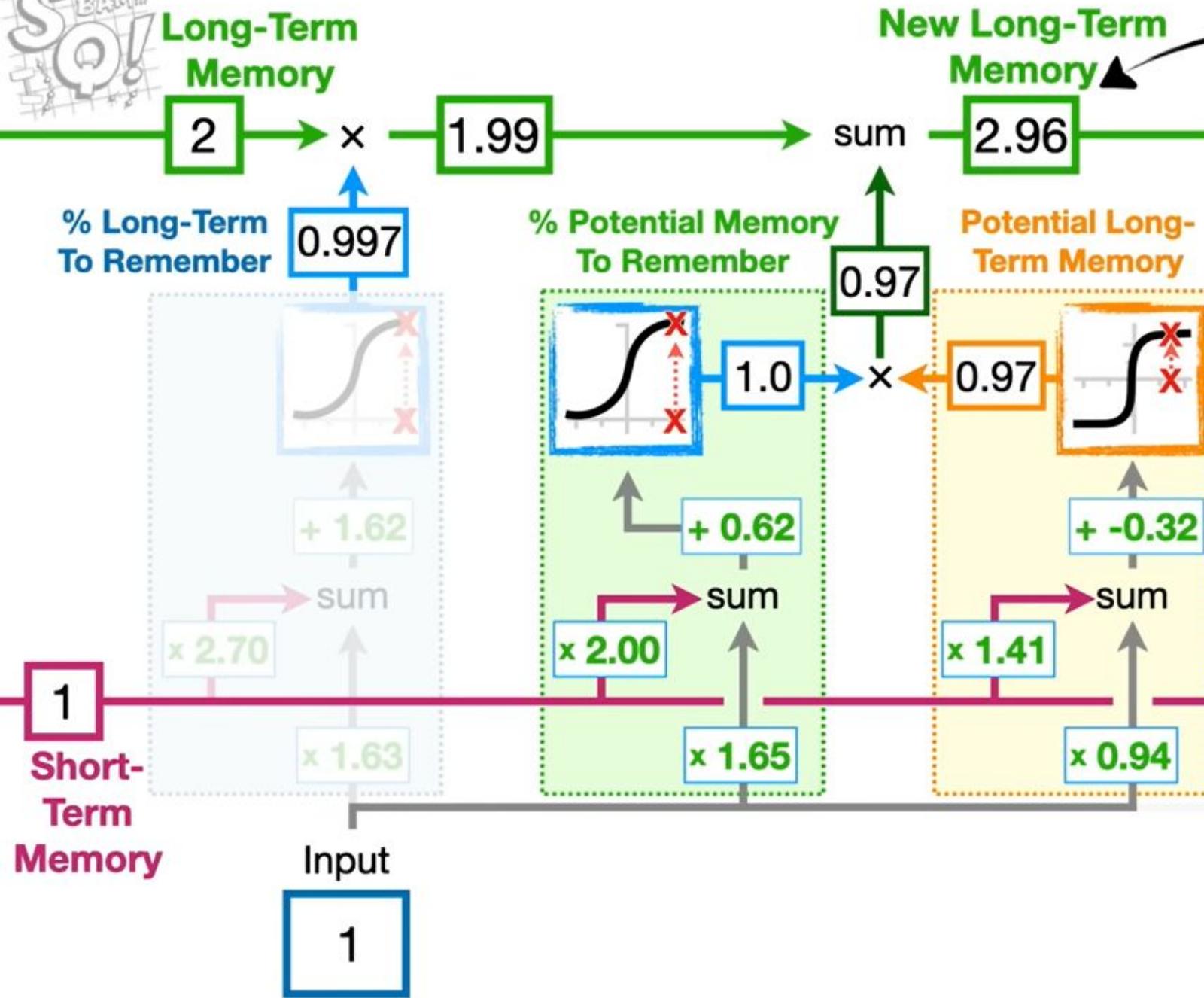


...and so we would not add anything to the **Long-Term Memory**.

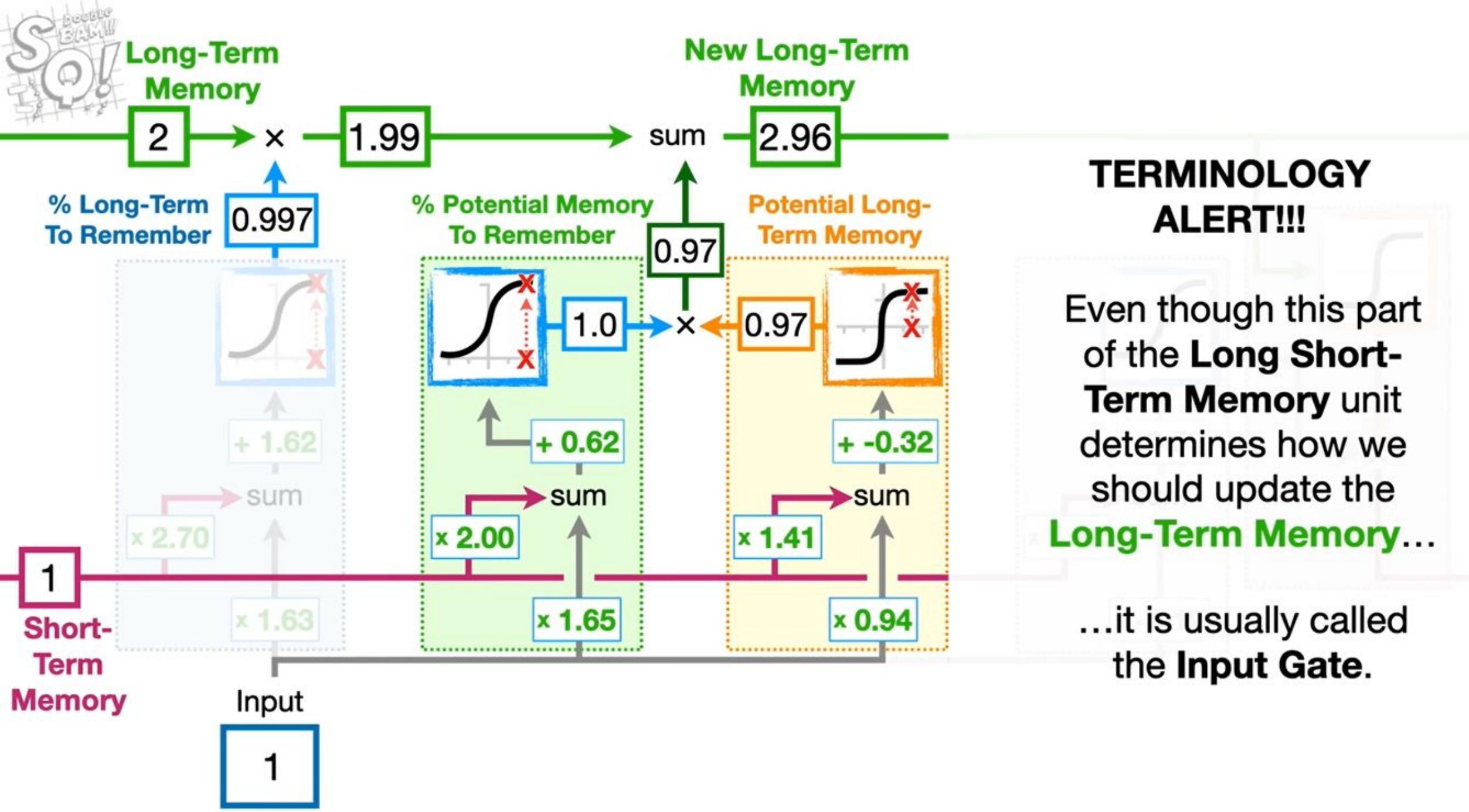




...we add 0.97 to the existing **Long-Term Memory**...



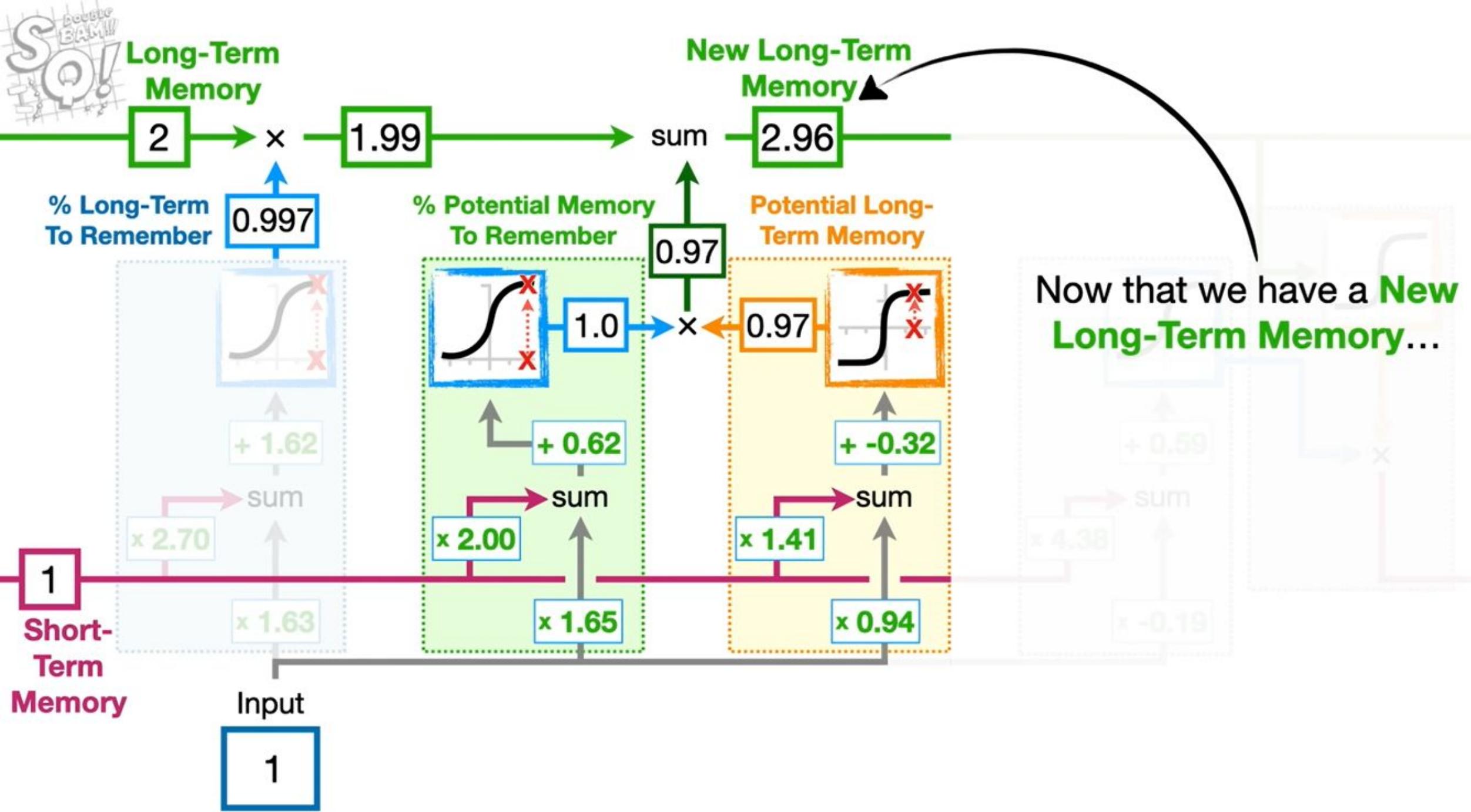
...and we get a new
Long-Term Memory,
2.96.

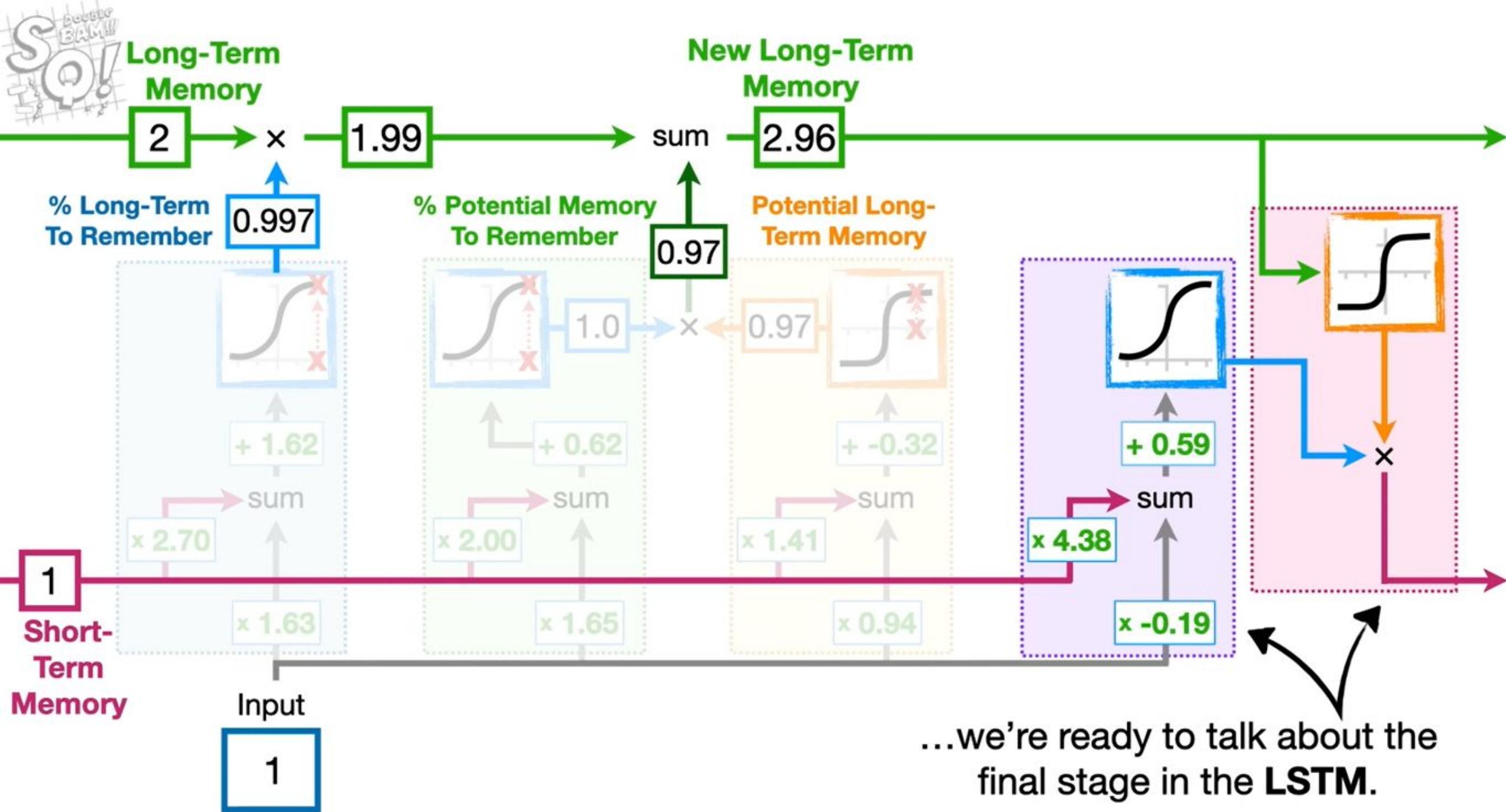


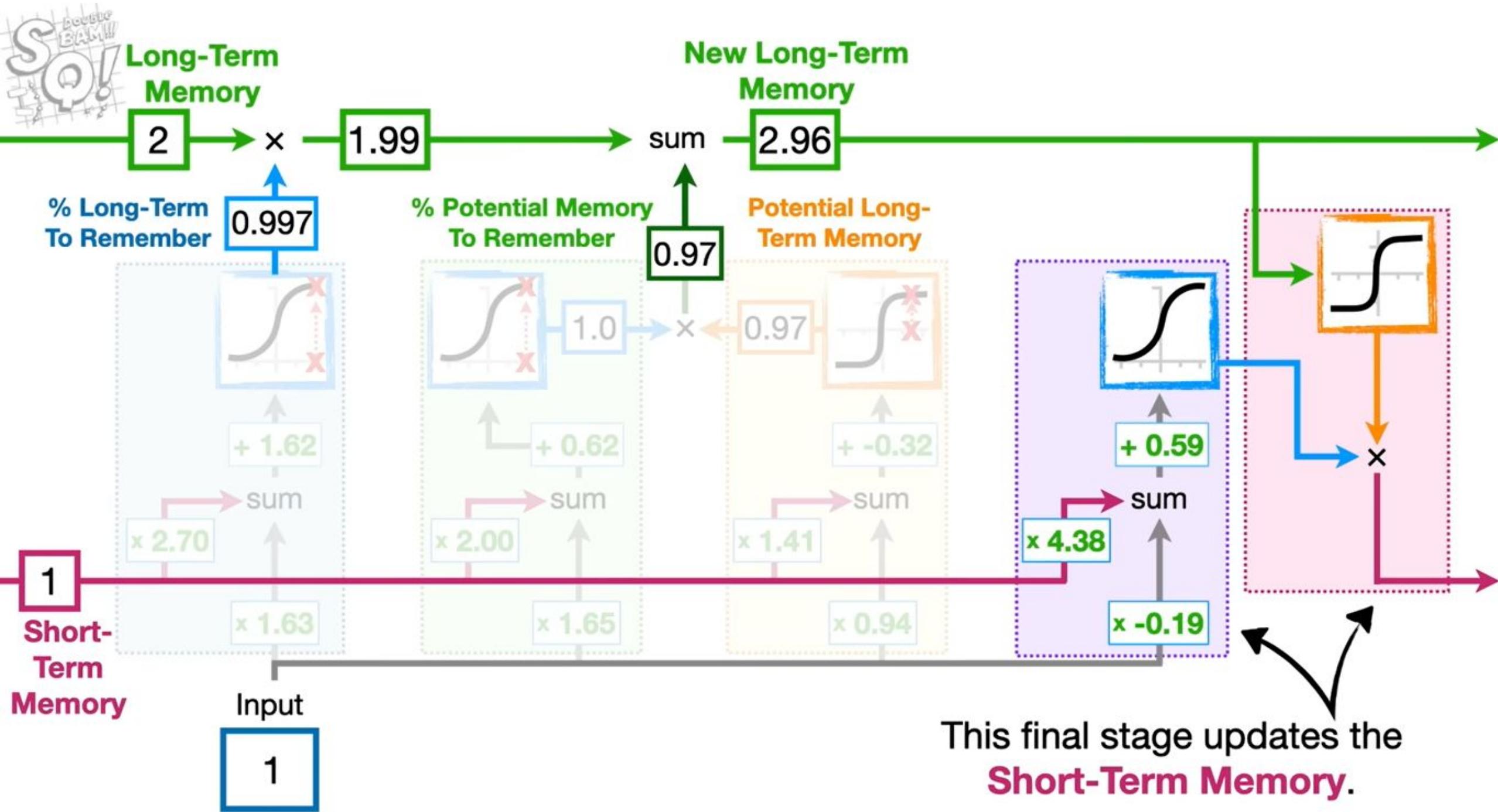
TERMINOLOGY ALERT!!!

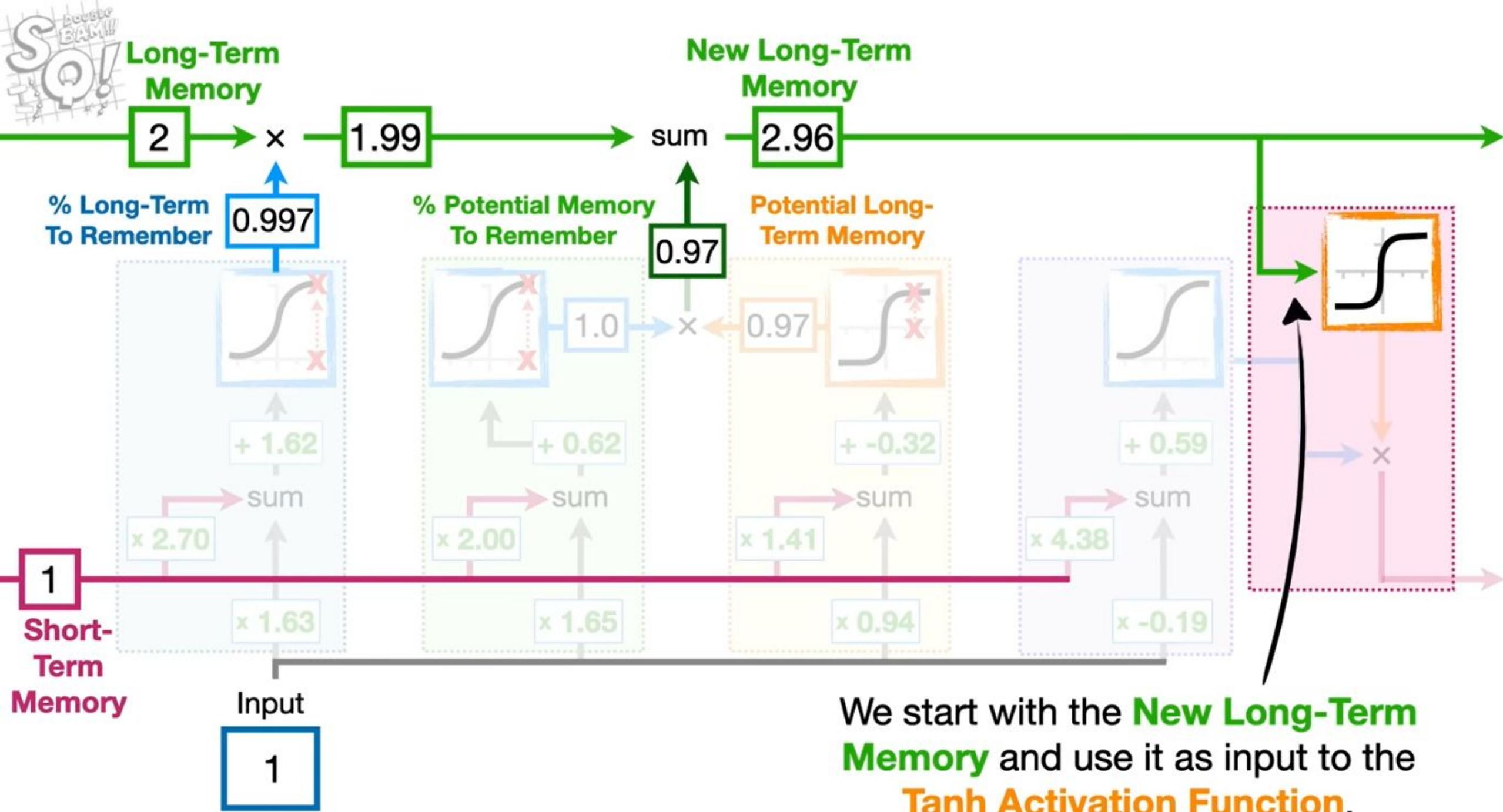
Even though this part of the **Long Short-Term Memory** unit determines how we should update the **Long-Term Memory**...

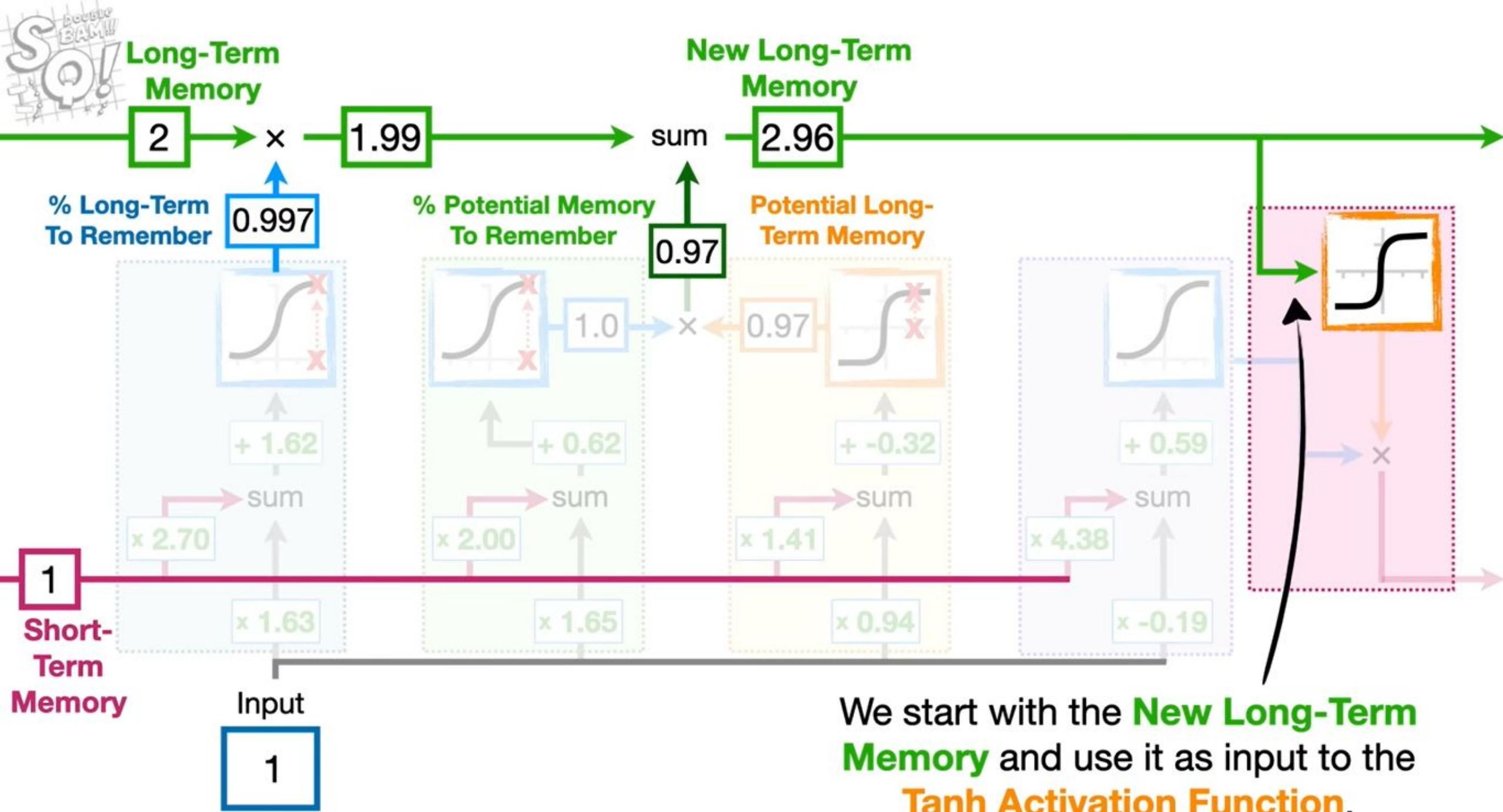
...it is usually called the **Input Gate**.

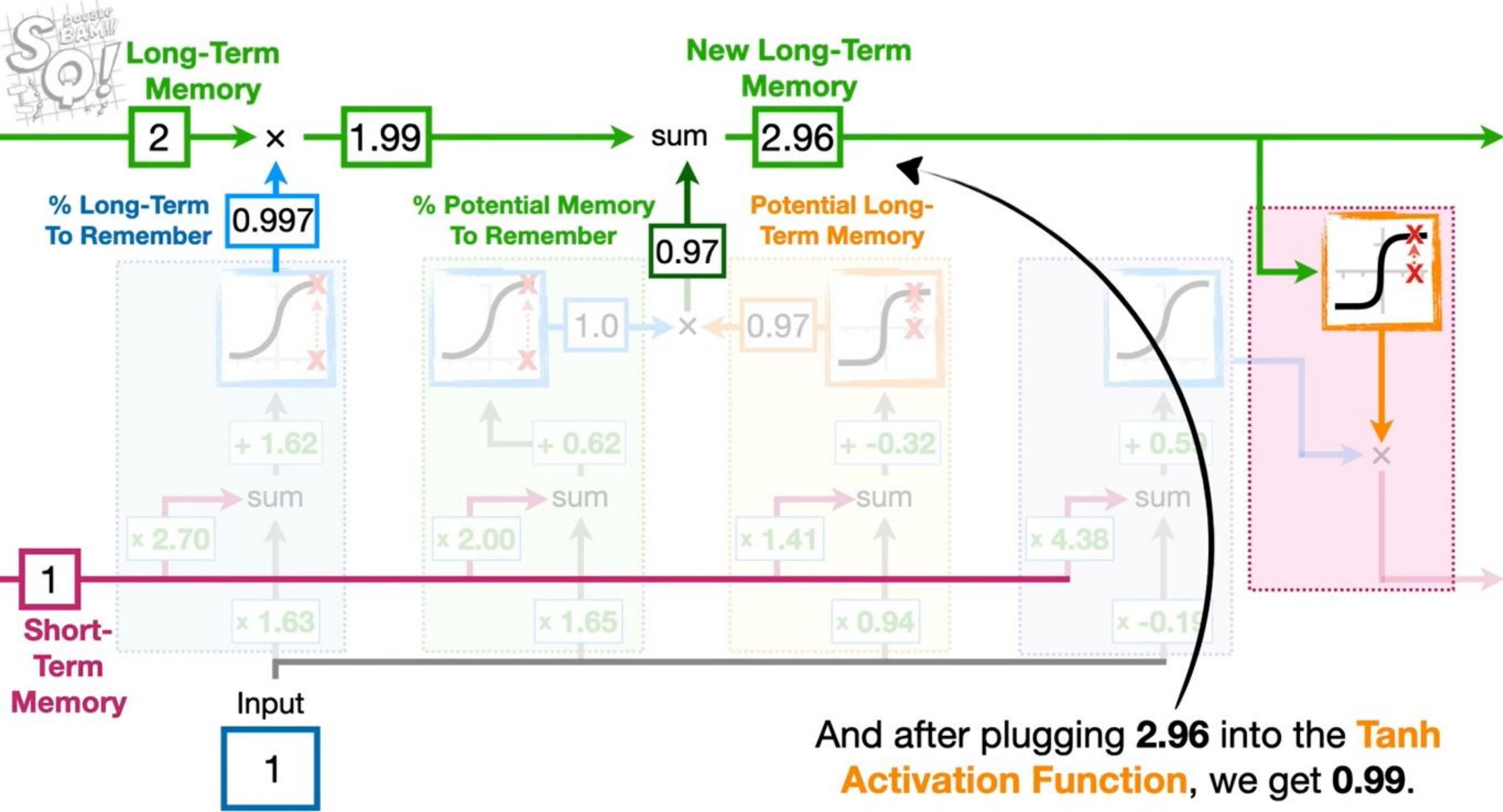


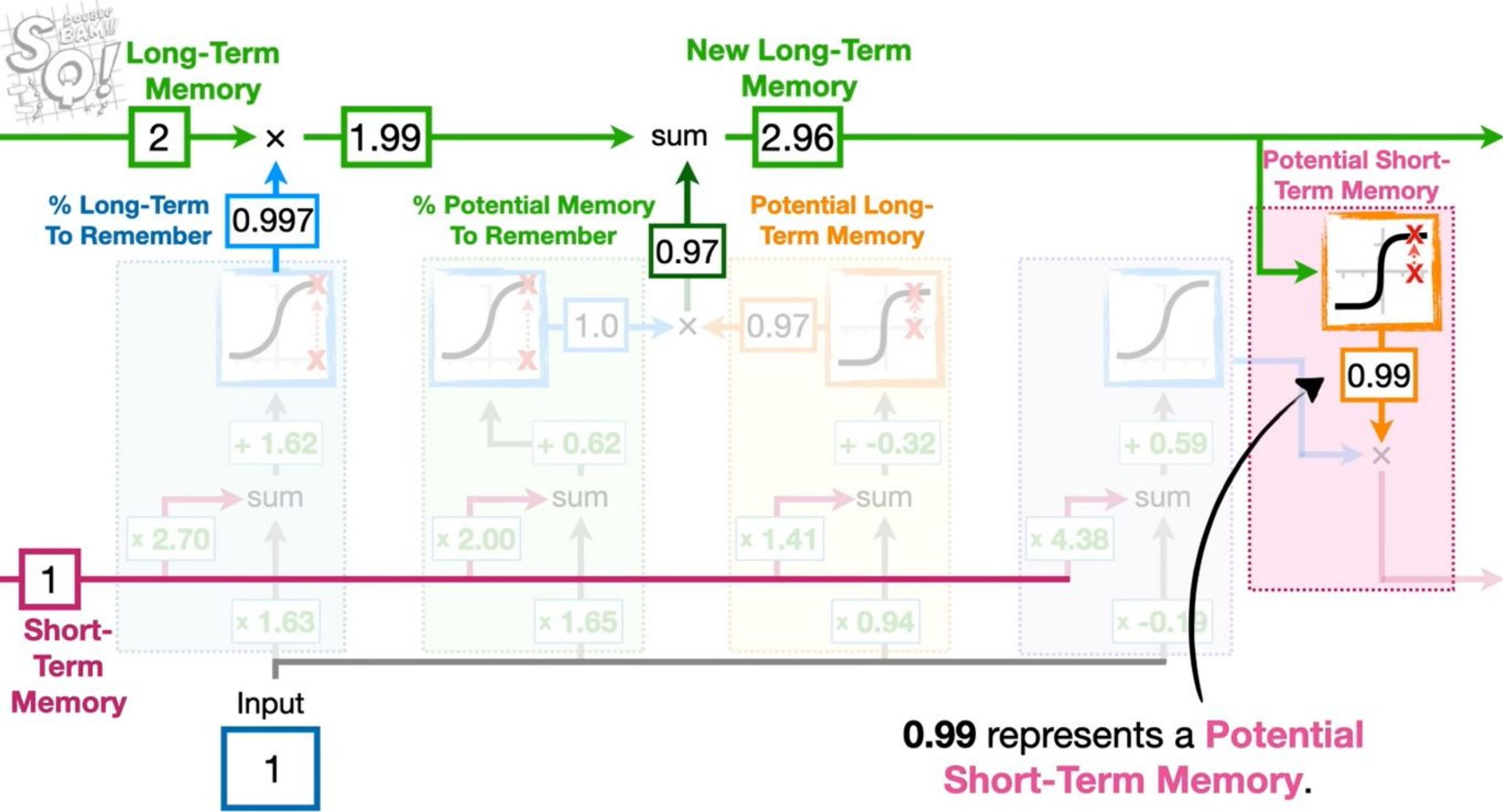


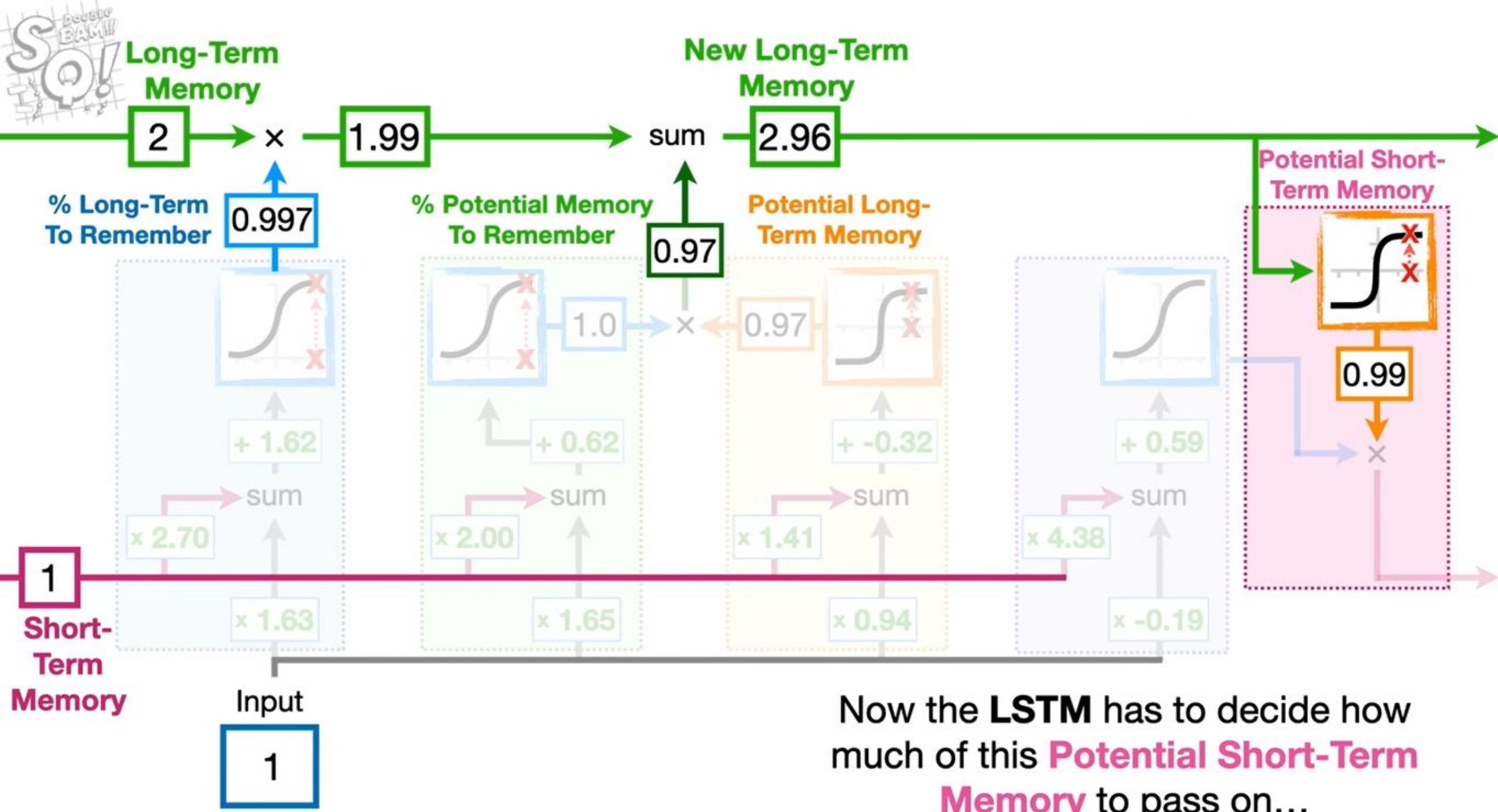




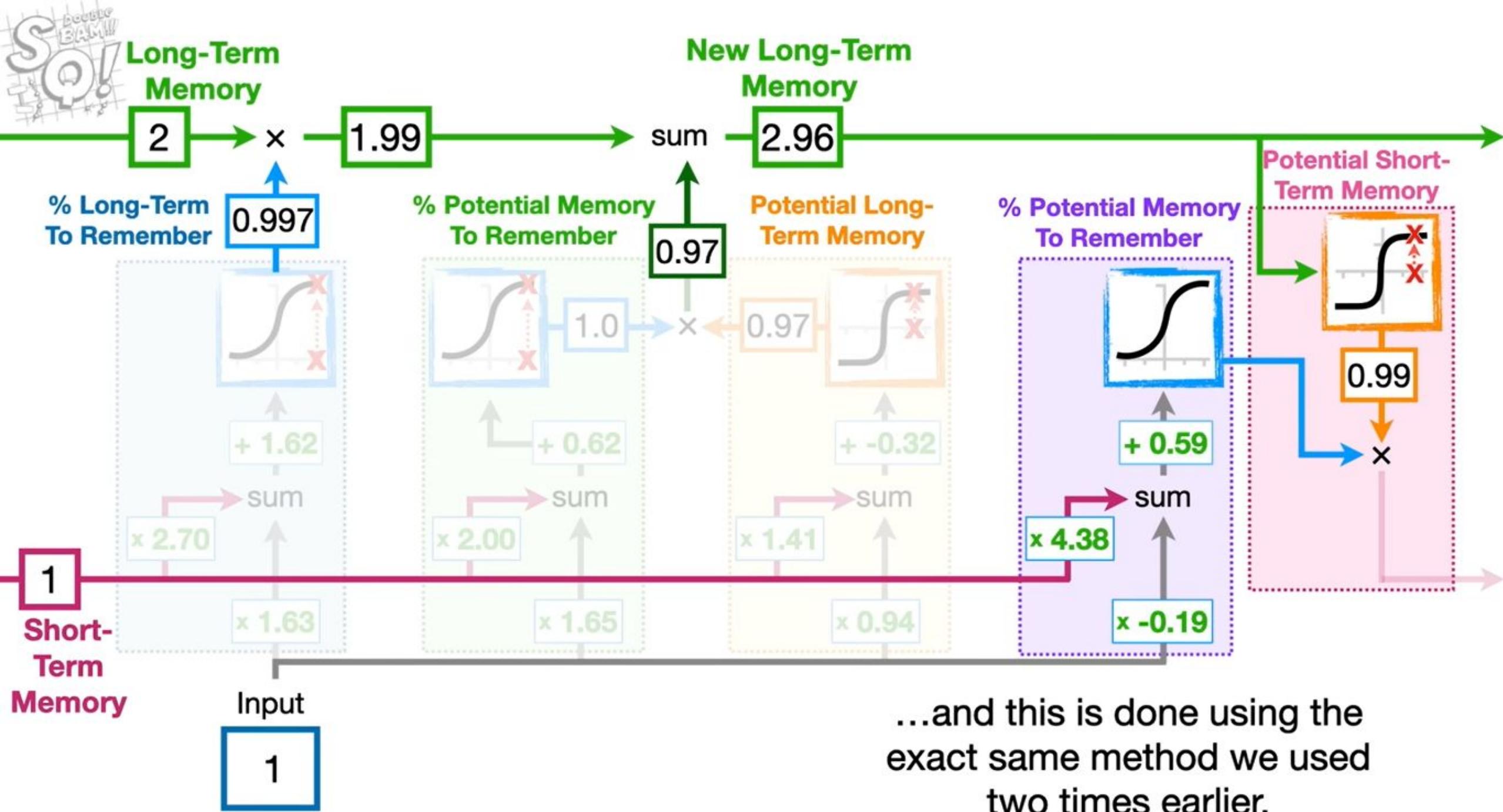


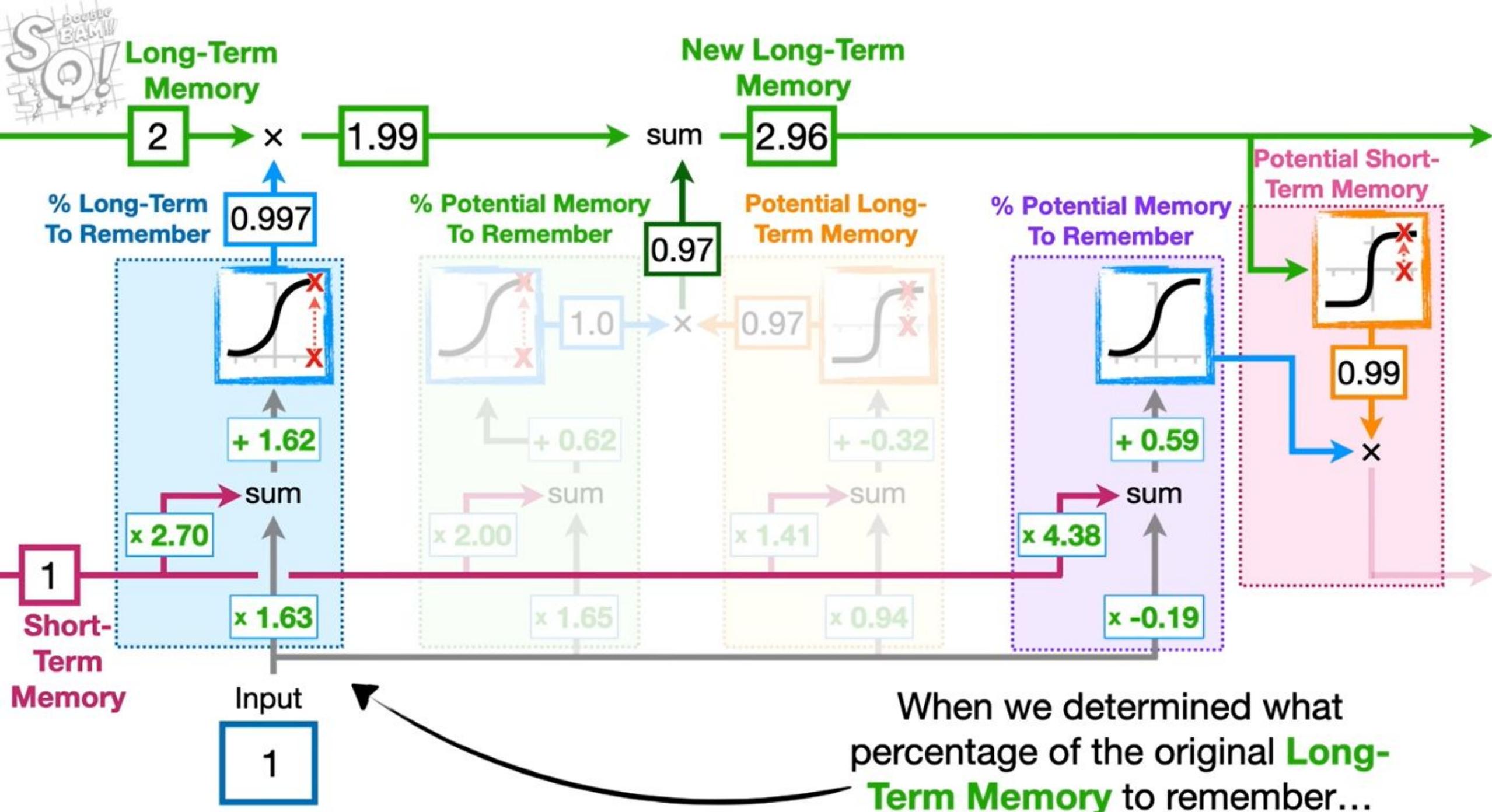


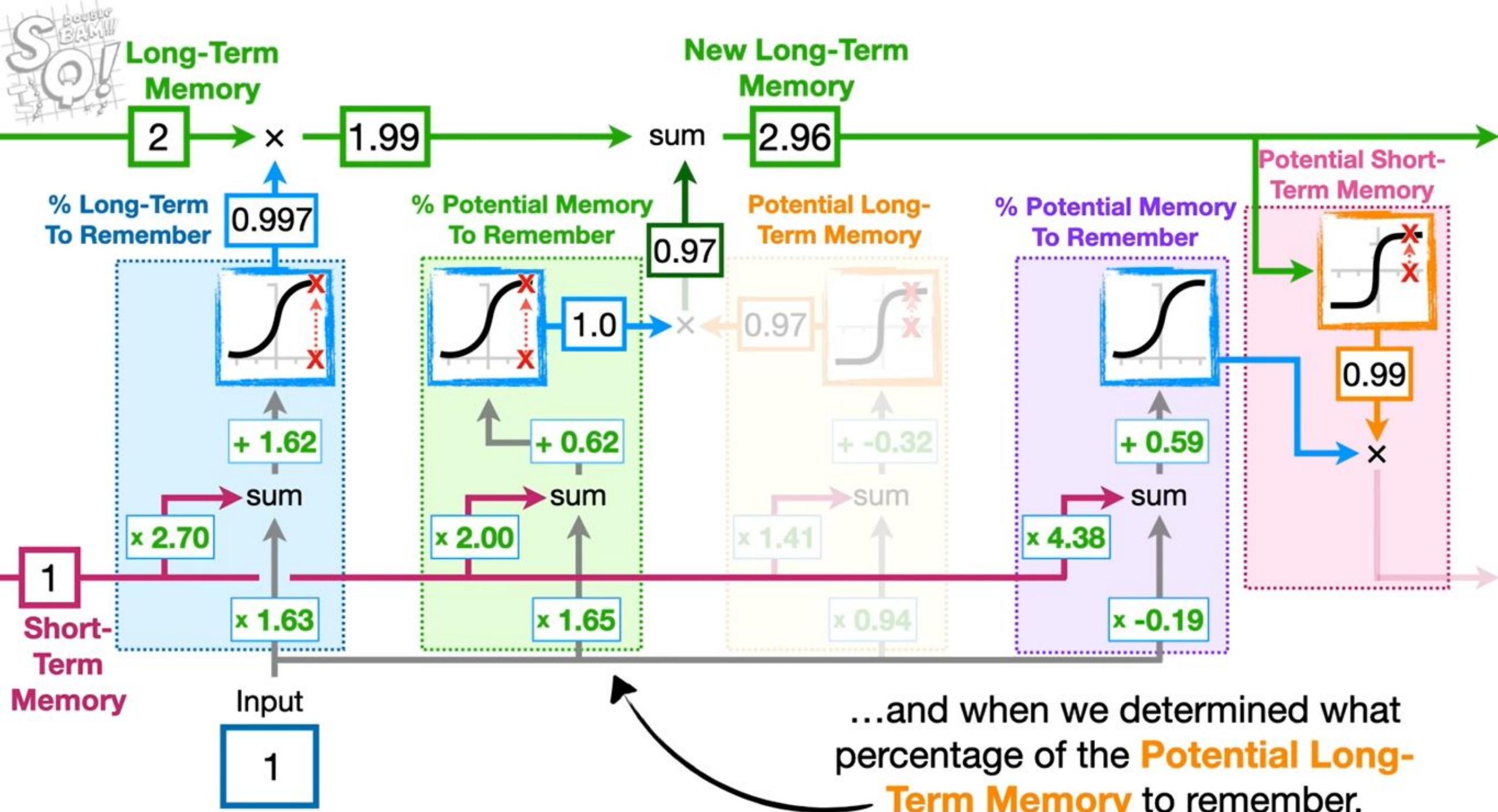


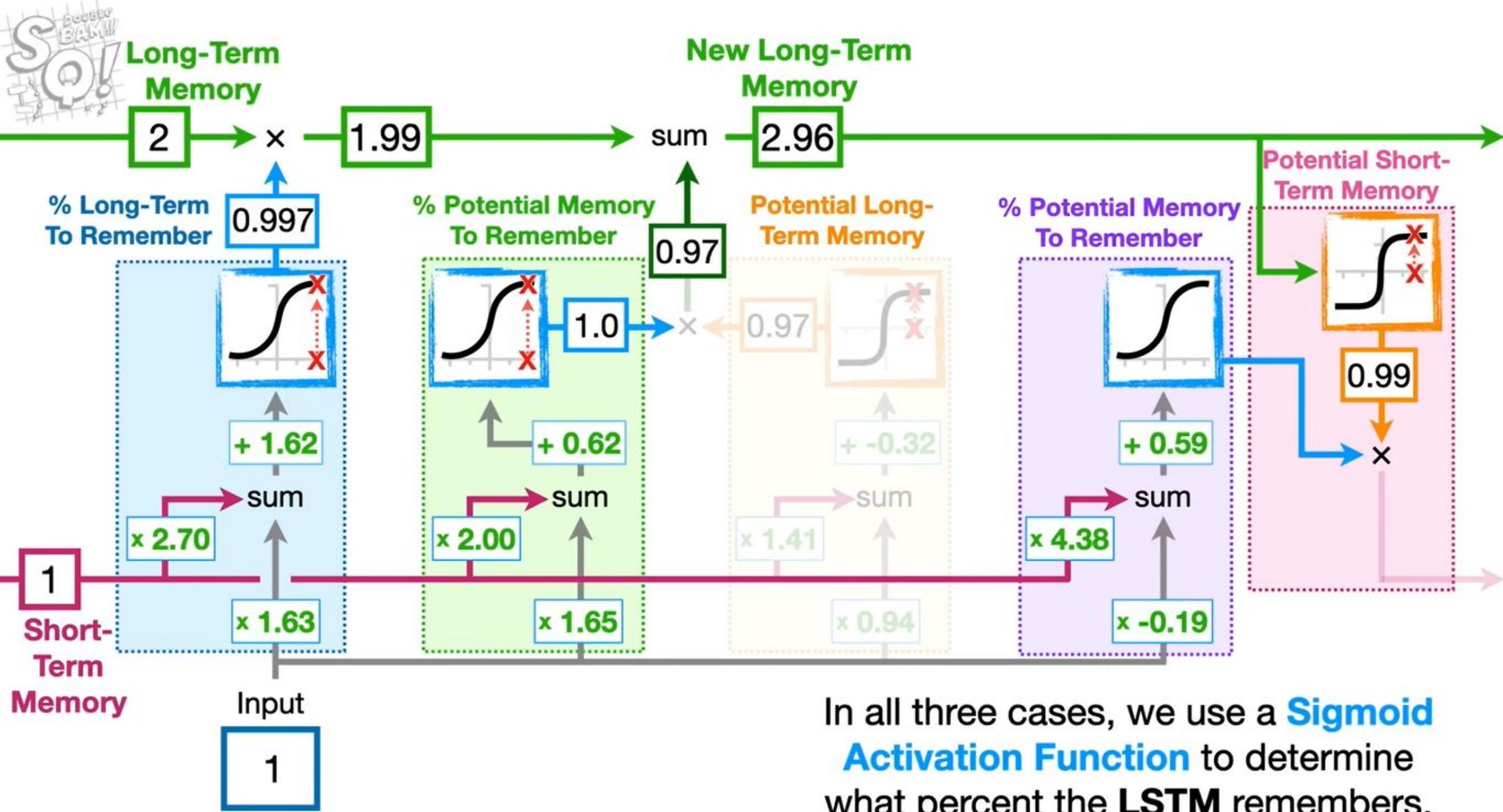


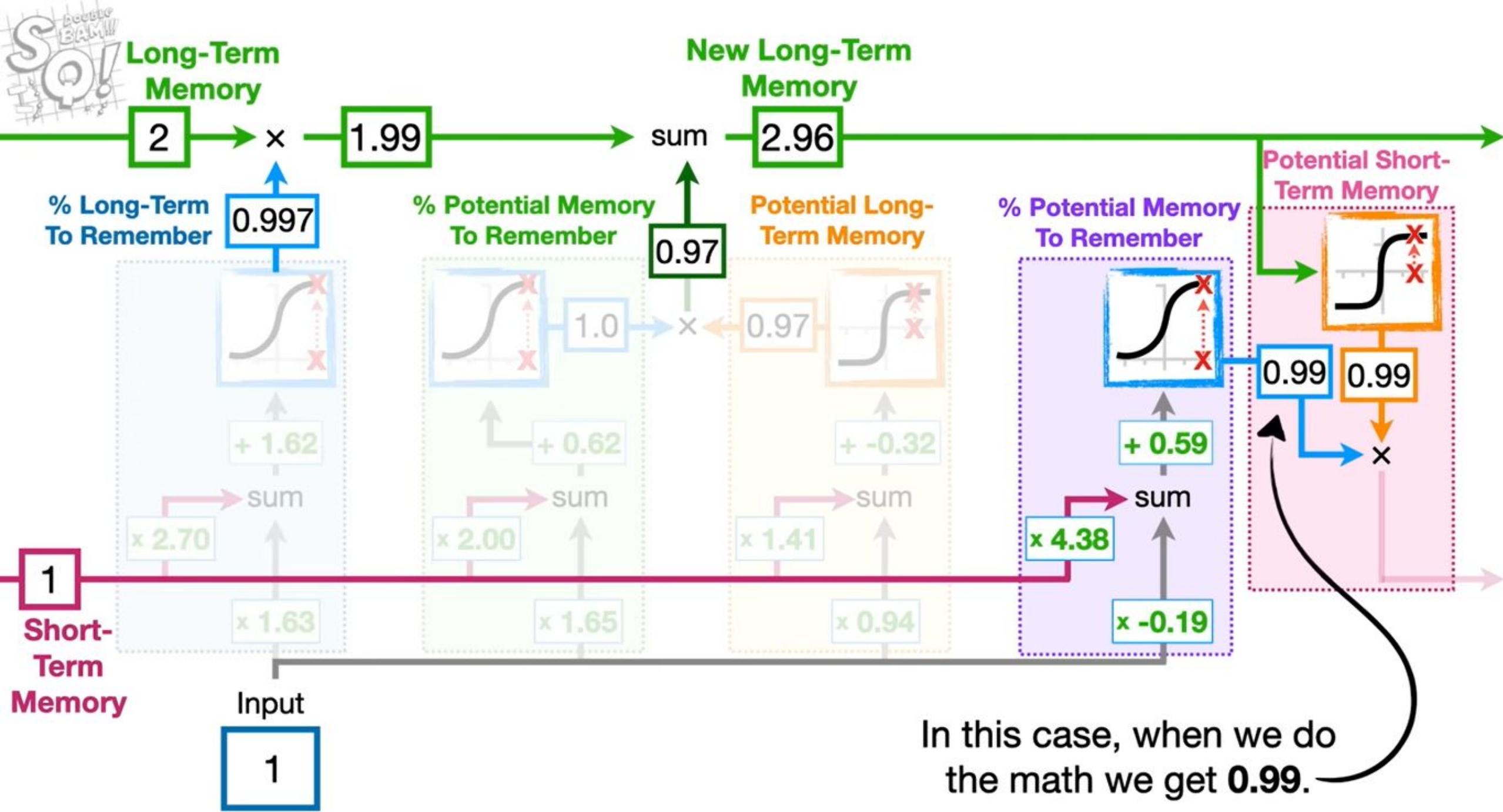
Now the **LSTM** has to decide how much of this **Potential Short-Term Memory** to pass on...

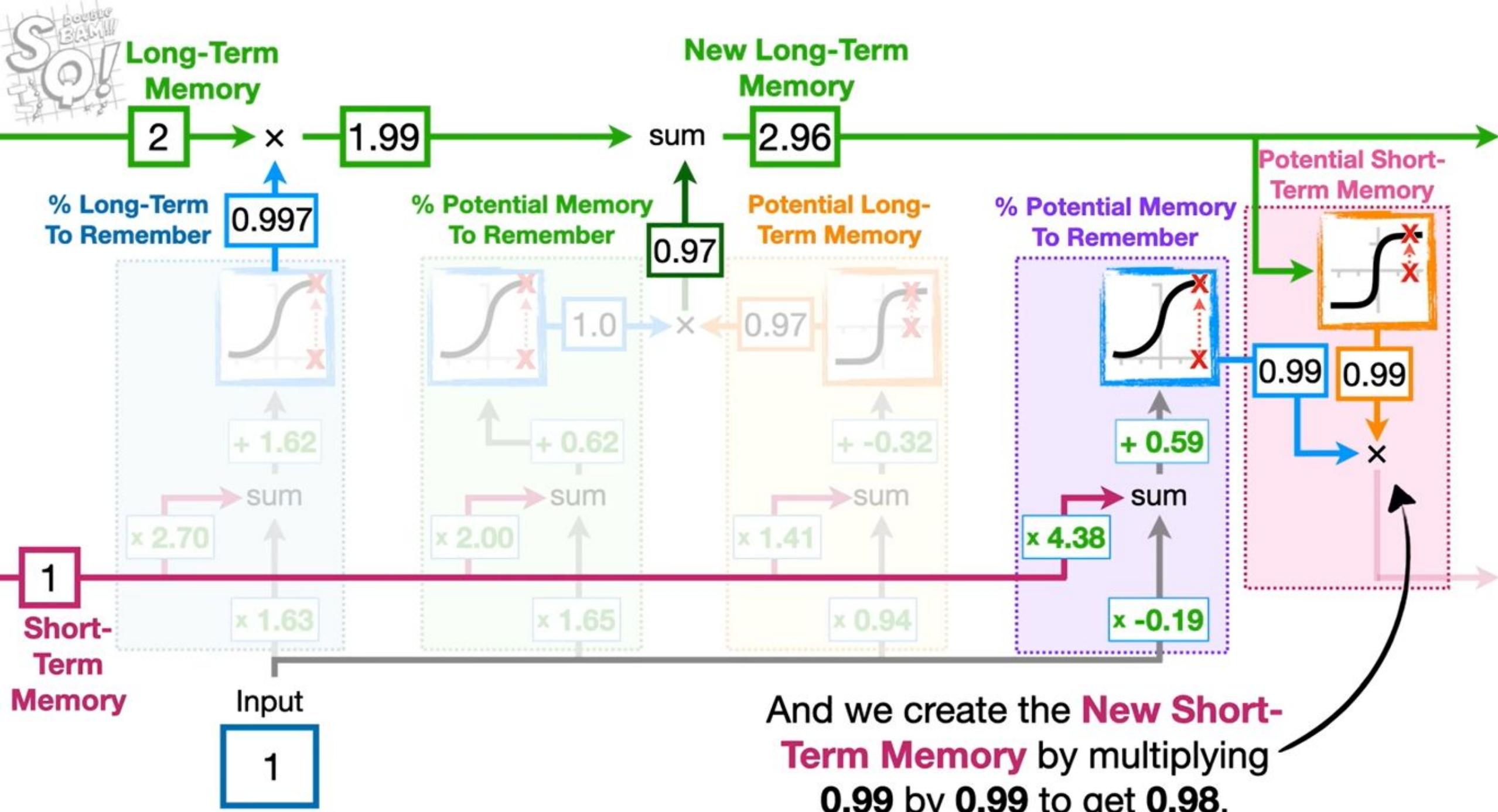


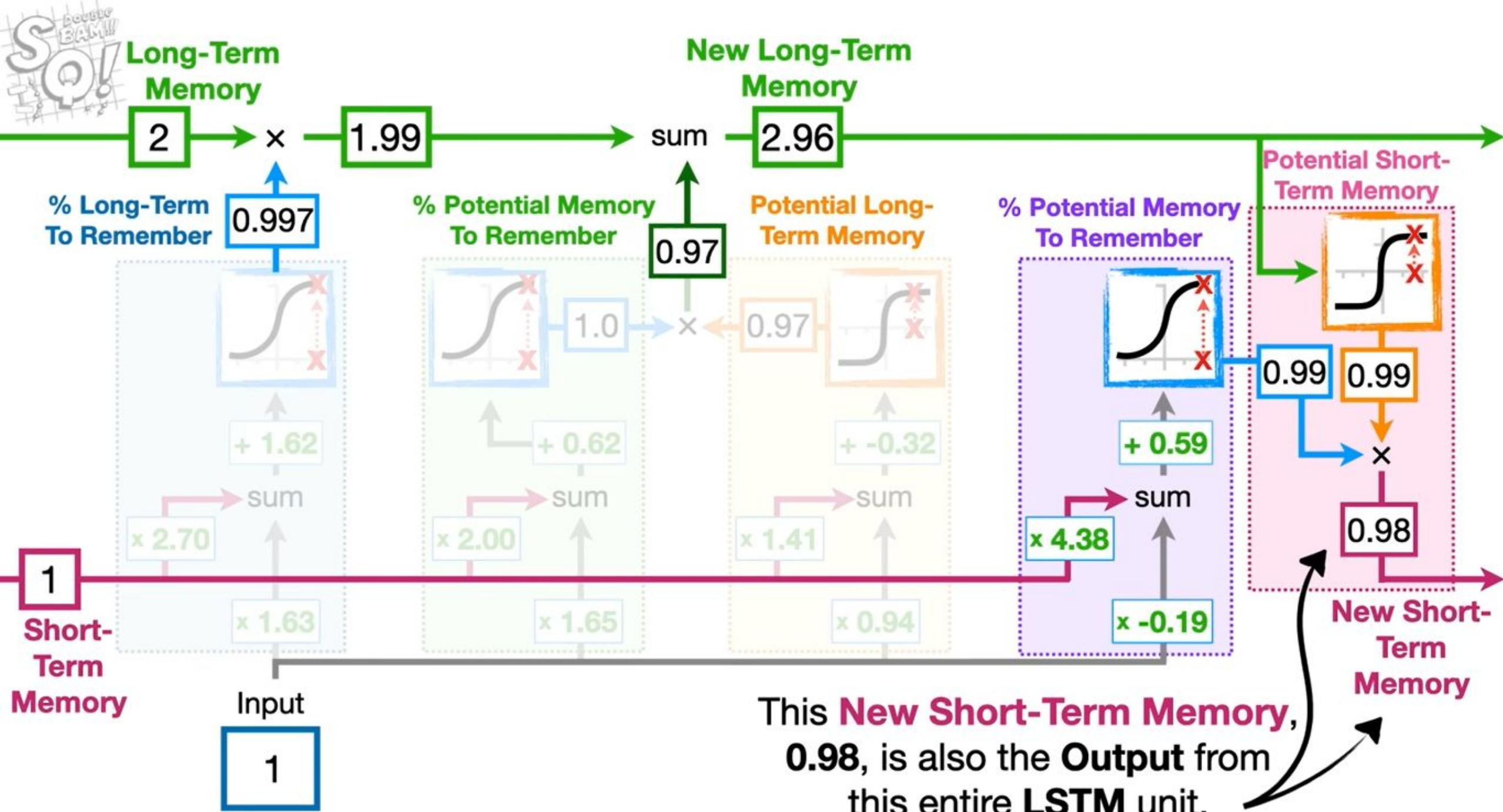


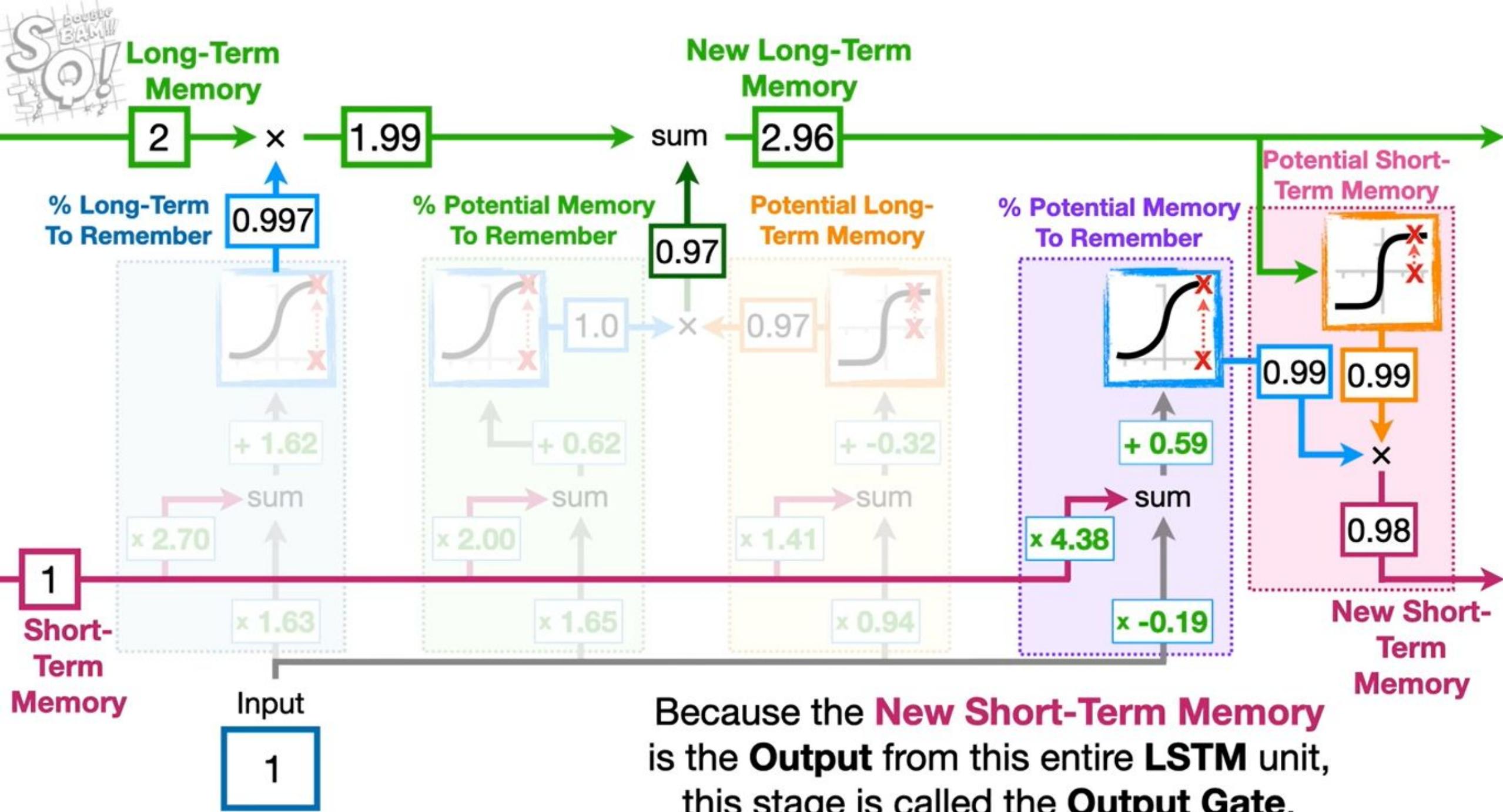


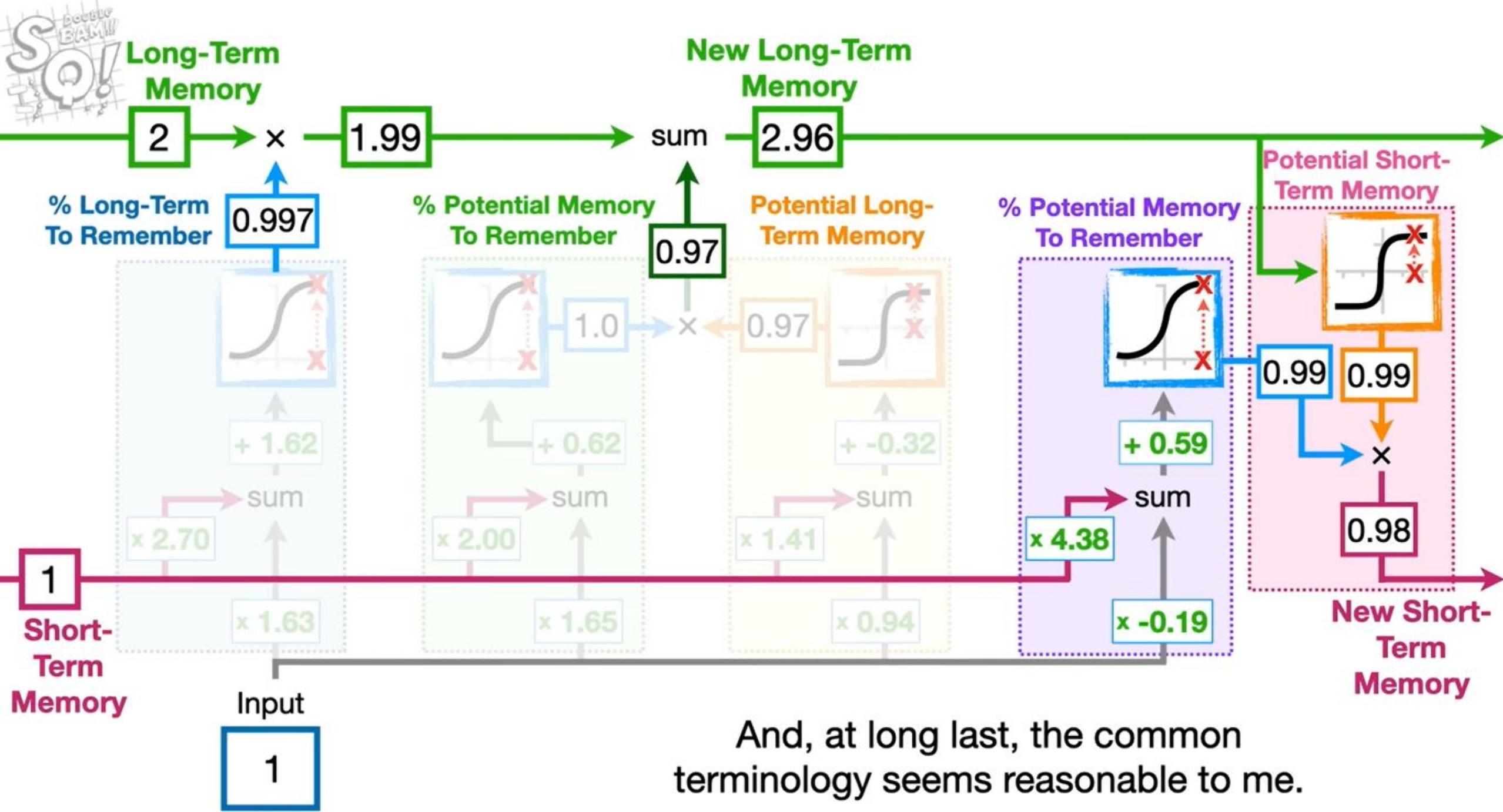




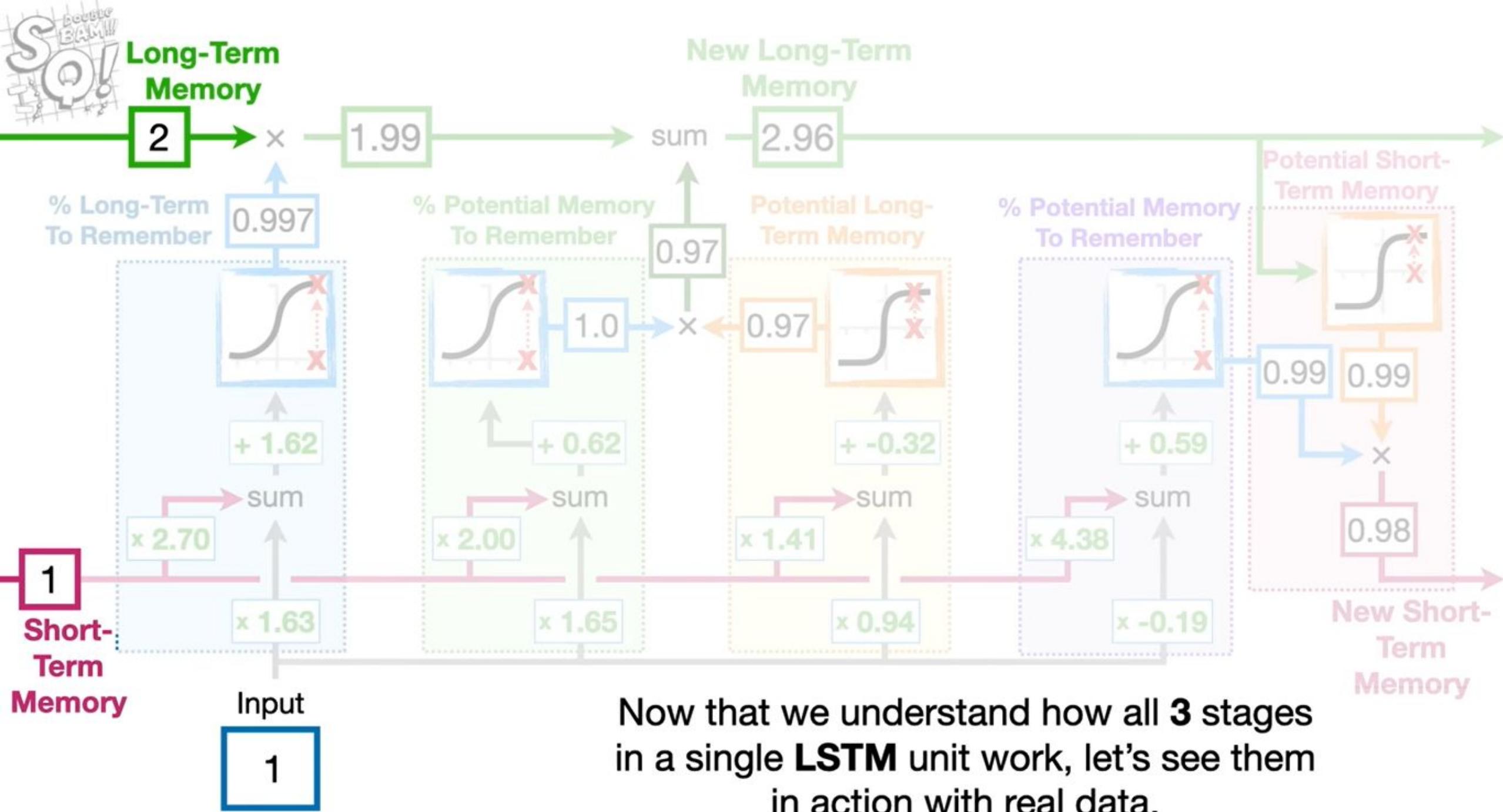


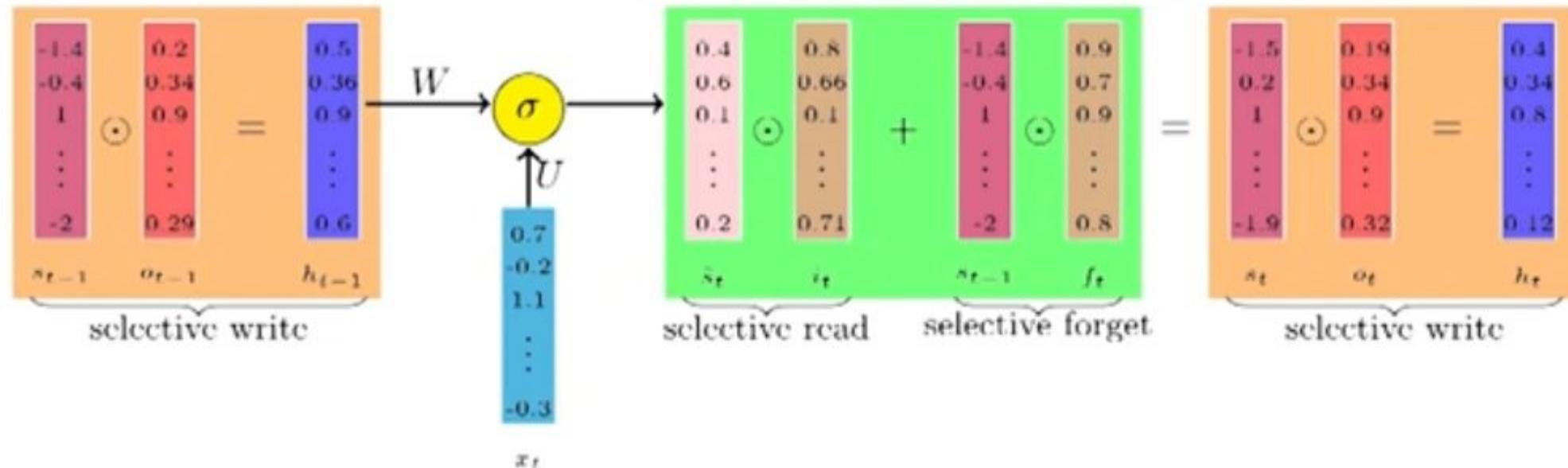






And, at long last, the common terminology seems reasonable to me.





- We now have the full set of equations for LSTMs
- The green box together with the selective write operations following it, show all the computations which happen at timestep t

Gates:

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

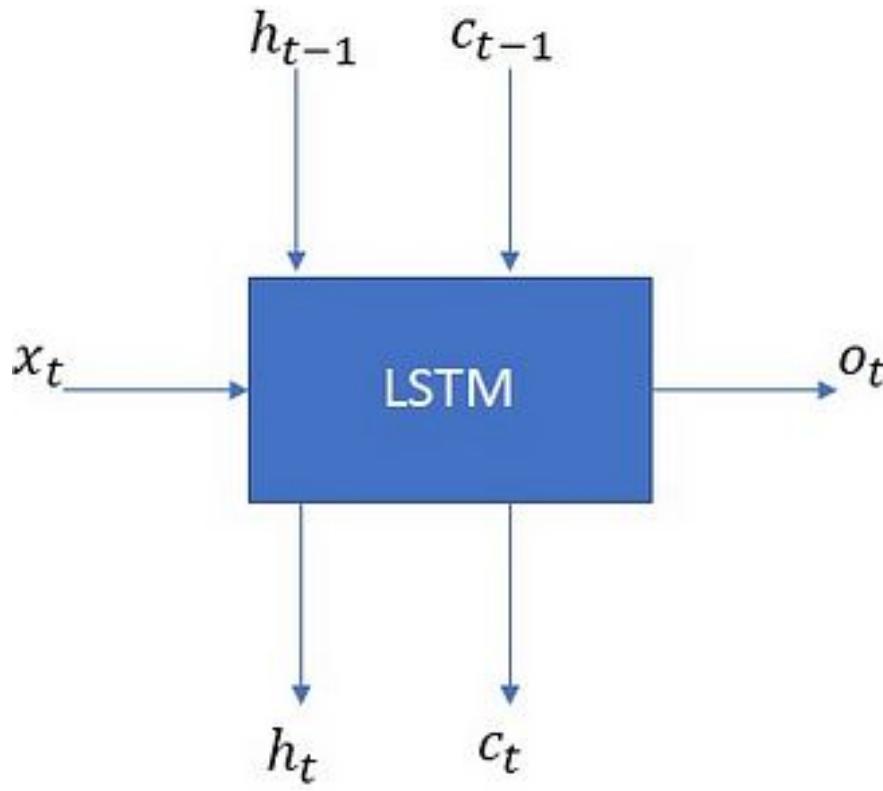
States:

$$\tilde{s}_t = \sigma(W h_{t-1} + U x_t + b)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

$$h_t = o_t \odot \sigma(s_t)$$

LSTM Equations



σ_g : sigmoid

σ_c : tanh

. : element wise multiplication

$$f_t = \sigma_g (W_f \times x_t + U_f \times h_{t-1} + b_f)$$

$$i_t = \sigma_g (W_i \times x_t + U_i \times h_{t-1} + b_i)$$

$$o_t = \sigma_g (W_o \times x_t + U_o \times h_{t-1} + b_o)$$

$$c'_t = \sigma_c (W_c \times x_t + U_c \times h_{t-1} + b_c)$$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot c'_t$$

$$h_t = o_t \cdot \sigma_c(c_t)$$

f_t is the forget gate

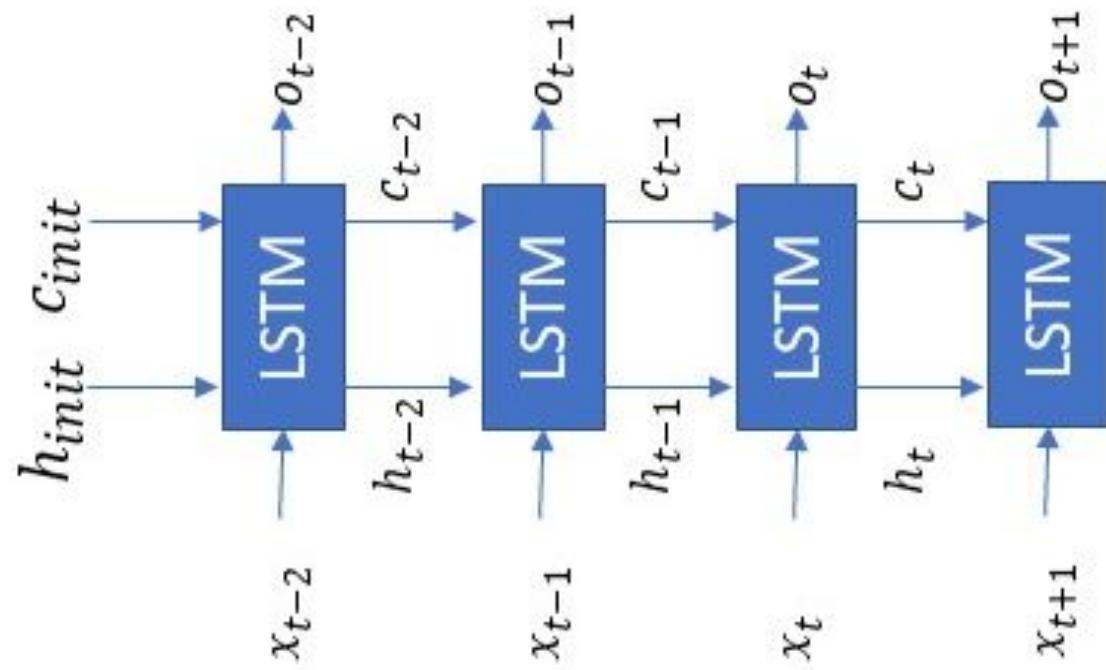
i_t is the input gate

o_t is the output gate

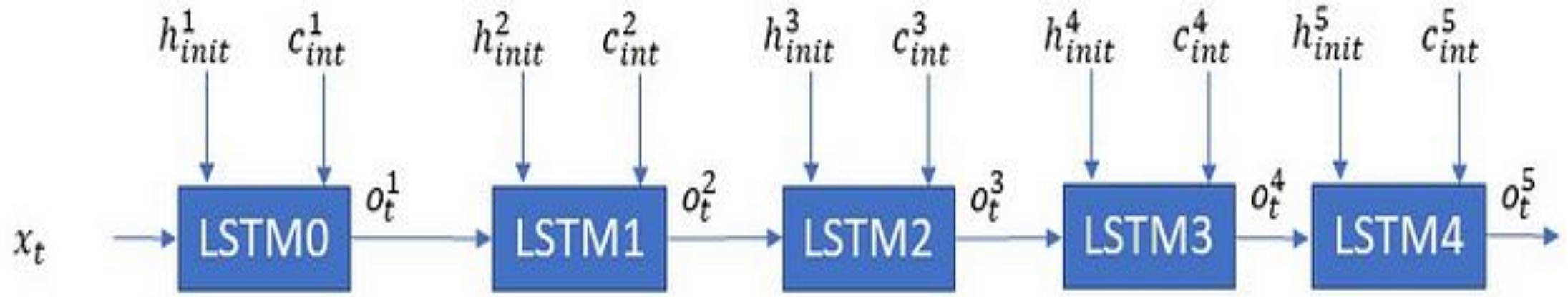
c_t is the cell state

h_t is the hidden state

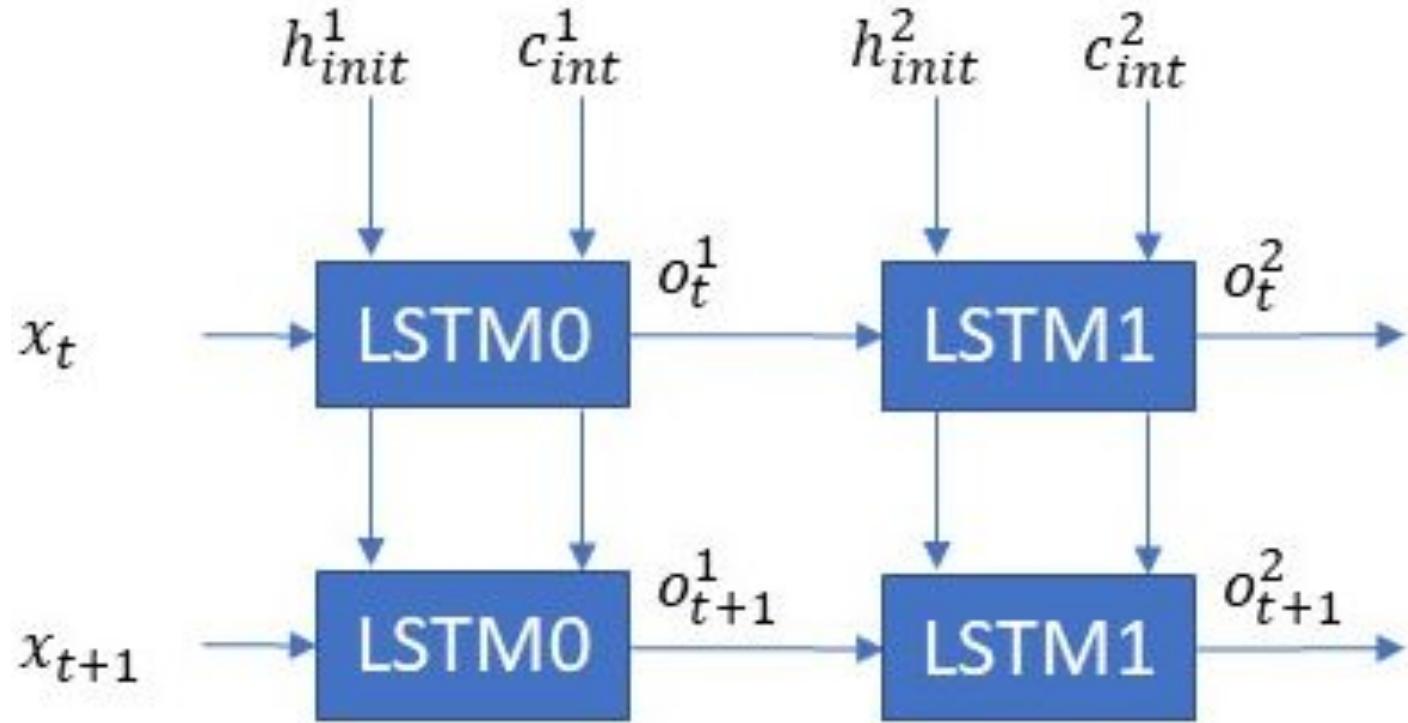
Network with 1 layer of LSTM with 4 timestamps



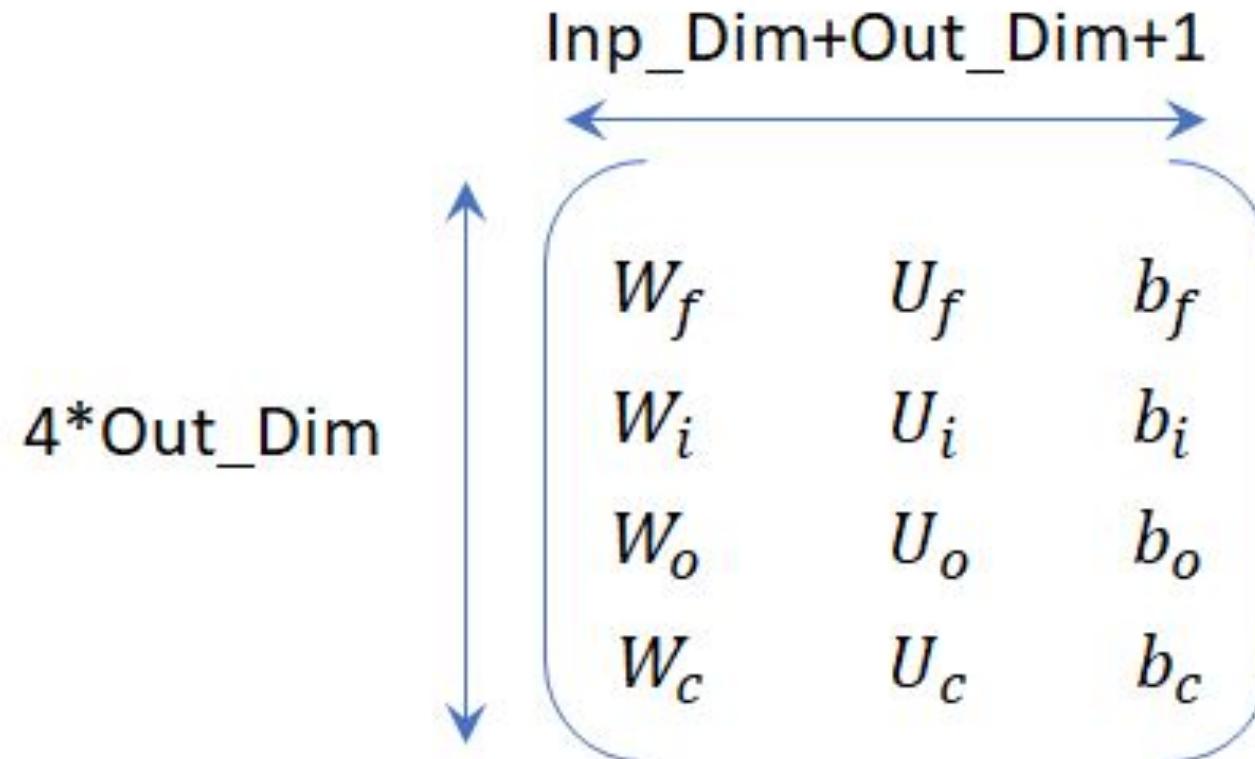
Network with 5 layers of LSTM with time stamp =1



Network with 2 layers of LSTM with time stamp =2

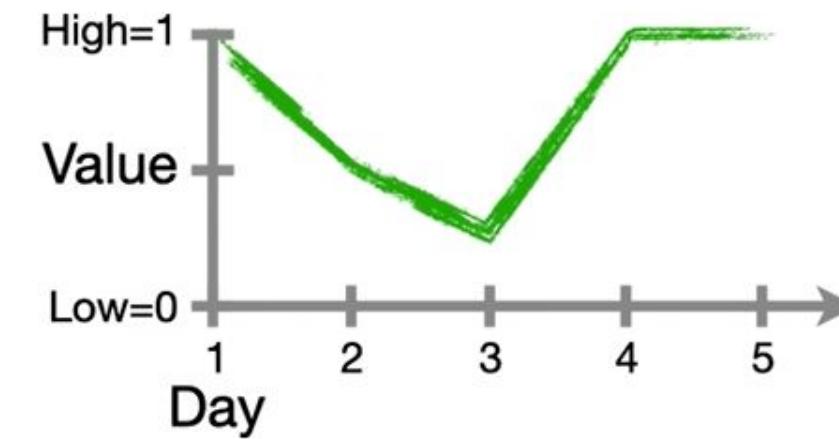
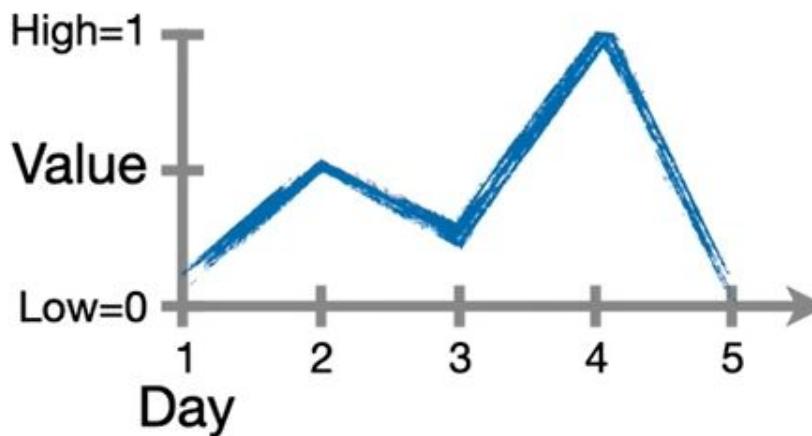


LSTM Weights





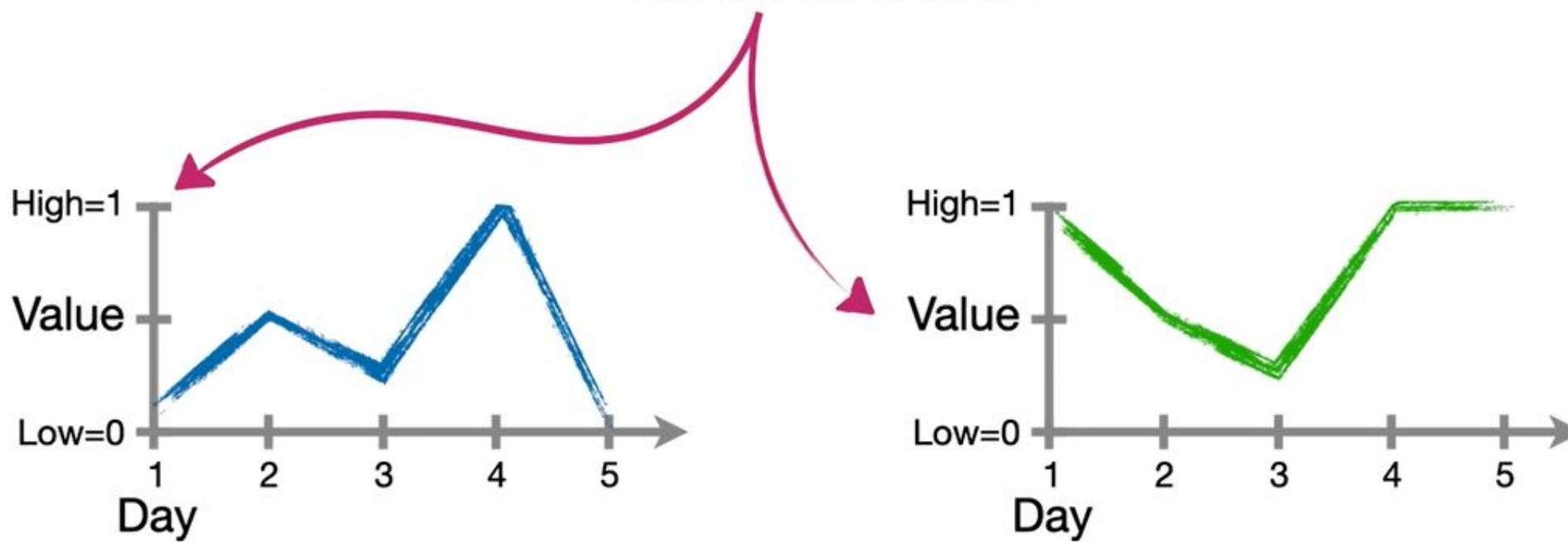
Here we have stock prices for two companies, **Company A** and **Company B**.



Company A = 
Company B = 



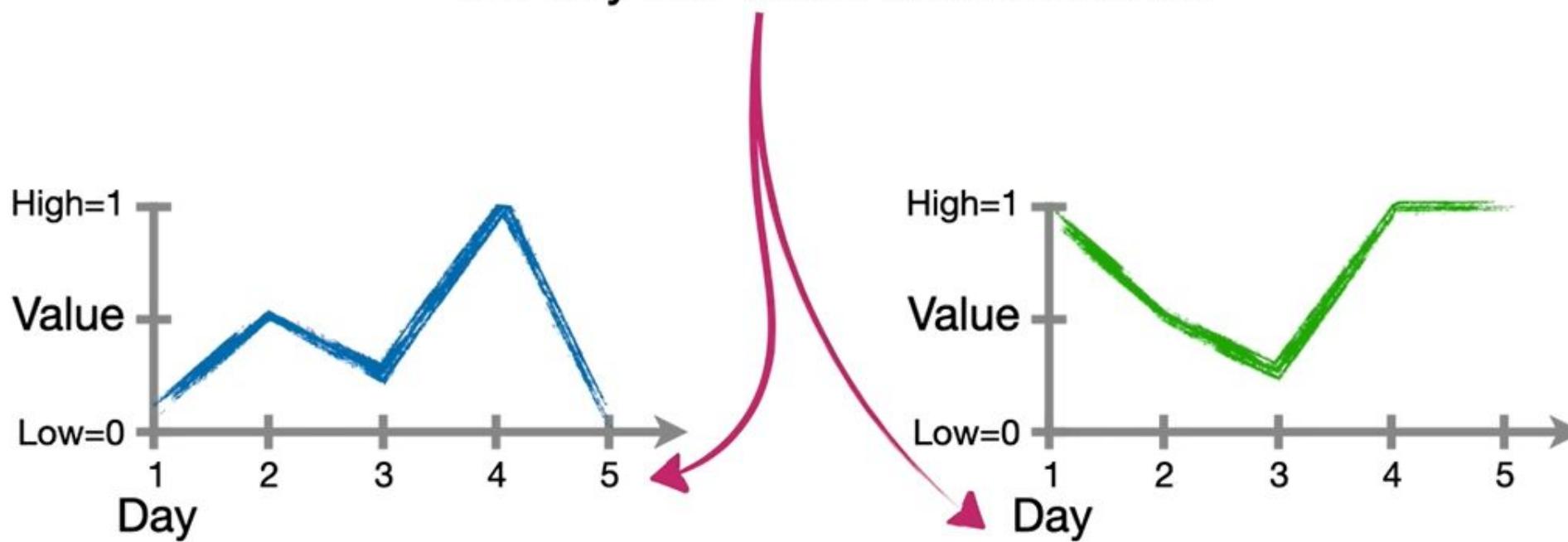
On the y-axis, we have
the stock value...



Company A =
Company B =



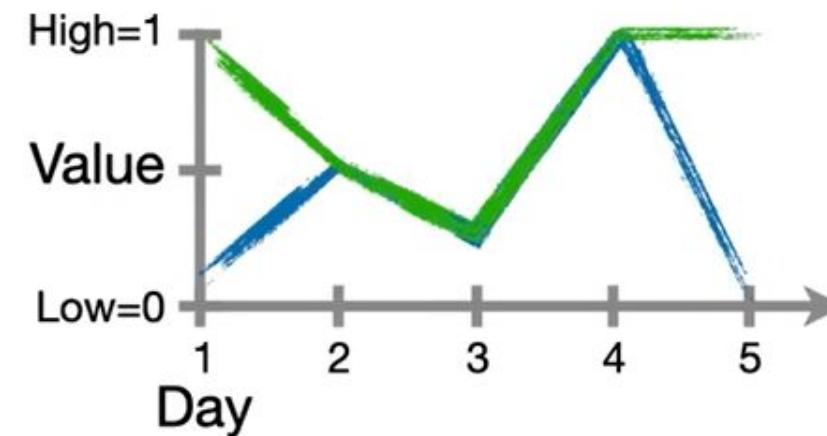
...and on the x-axis, we have
the day the value was recorded.



Company A = 
Company B = 



NOTE: If we overlap the data from the two companies...



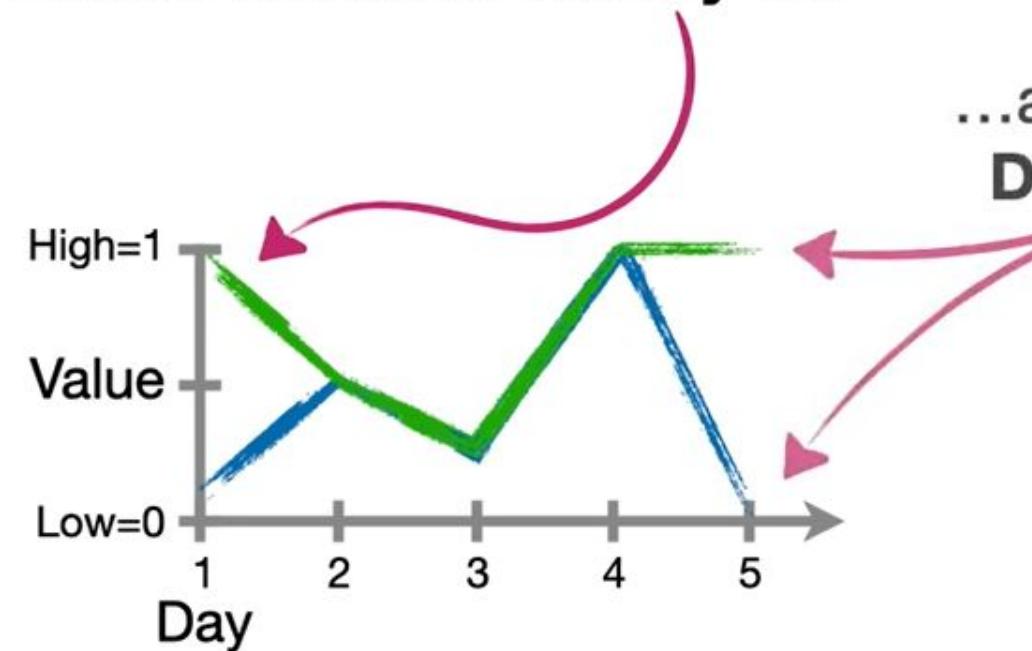
Company A =

Company B =



...we see that the only differences occur on **Day 1**...

...and on
Day 5.

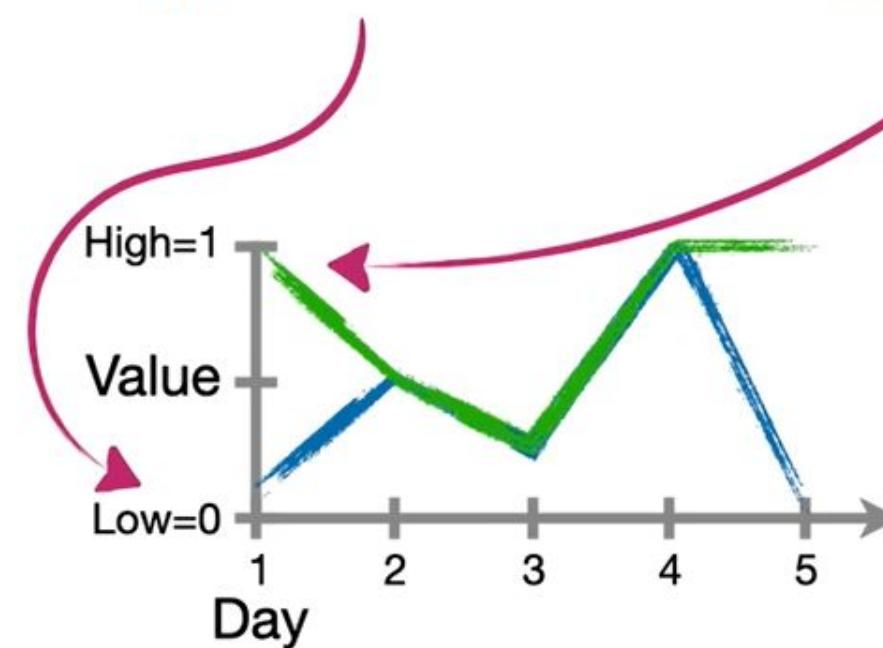


Company A = /

Company B = /



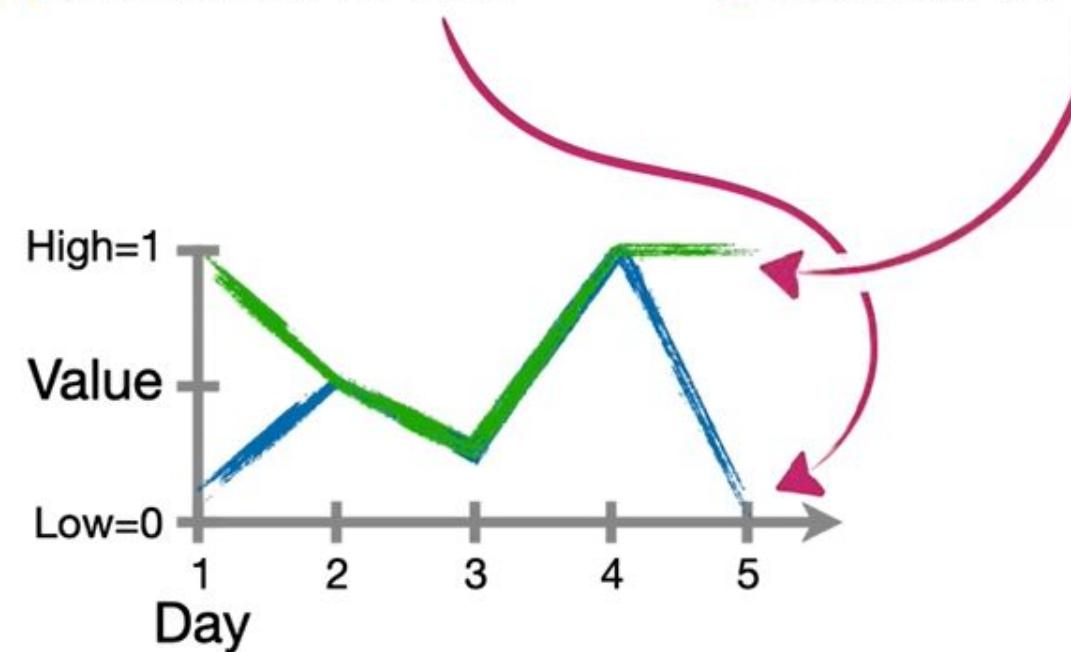
On Day 1, Company
A is at 0...
...and Company
B is at 1.



Company A =
Company B =



And on **Day 5**, **Company A** is returns to 0... ...and **Company B** returns to 1.

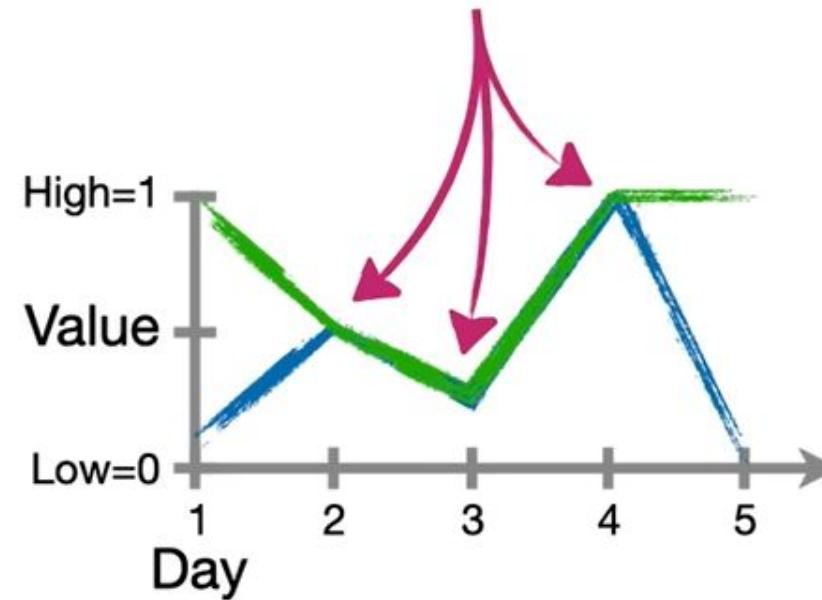


Company A =

Company B =



On all of the other days, **Days 2, 3 and 4**, both companies have the exact same values.

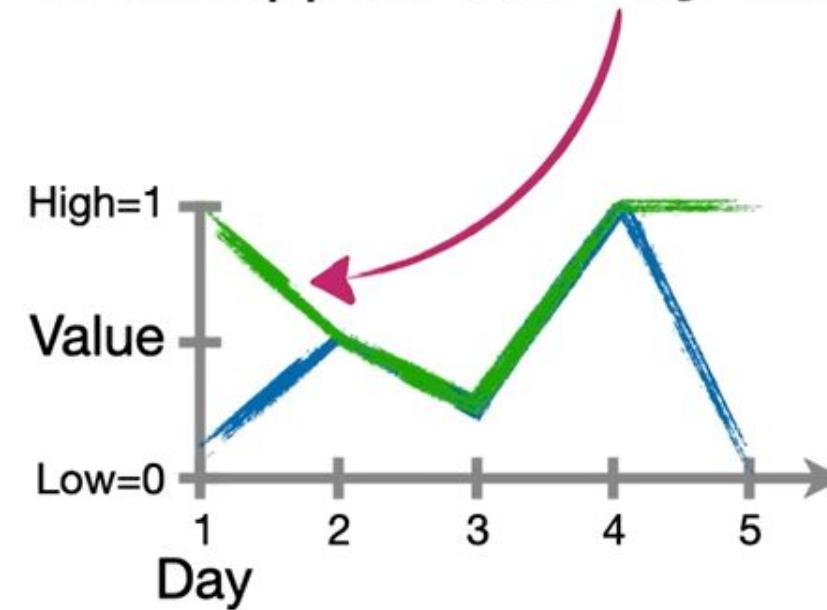


Company A = /

Company B = /



Given this sequential data, we want the **LSTM** to remember what happened on **Day 1...**

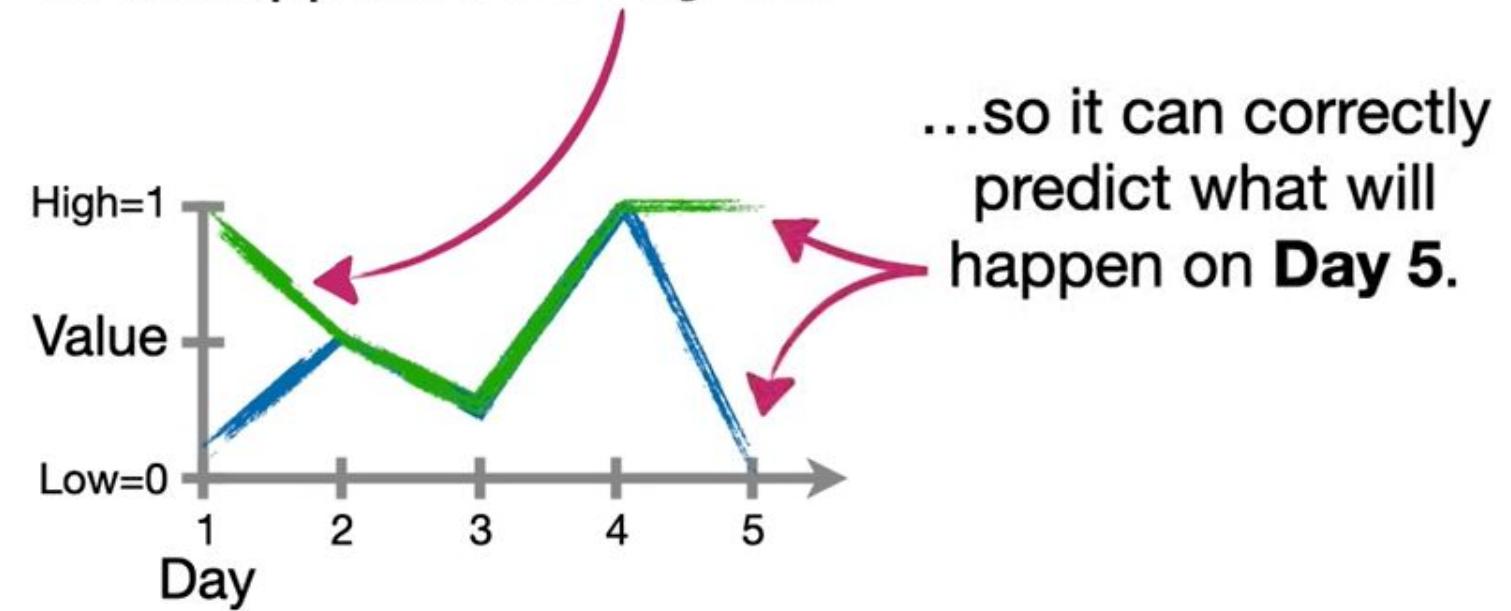


Company A =

Company B =



Given this sequential data, we want the **LSTM** to remember what happened on **Day 1**...



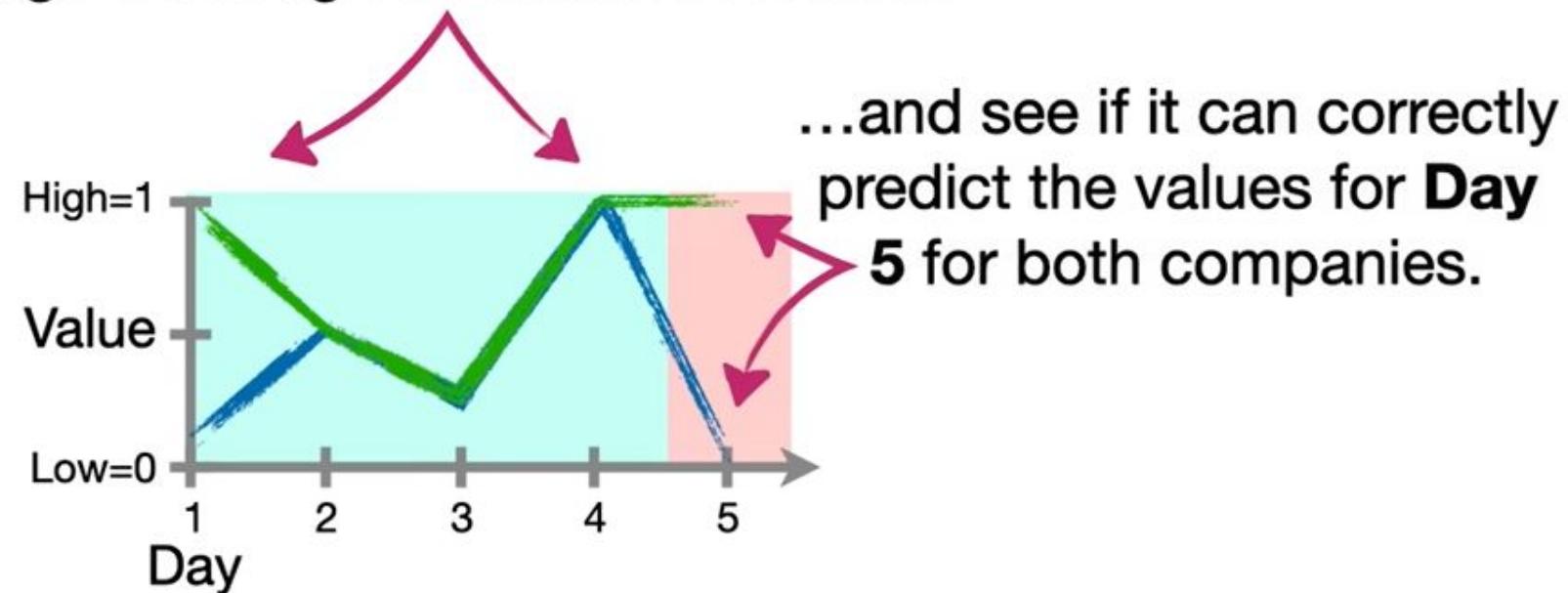
...so it can correctly predict what will happen on **Day 5**.

Company A = /

Company B = /



In other words, we're going to sequentially run the data from **Days 1** through **4** through an unrolled **LSTM**...

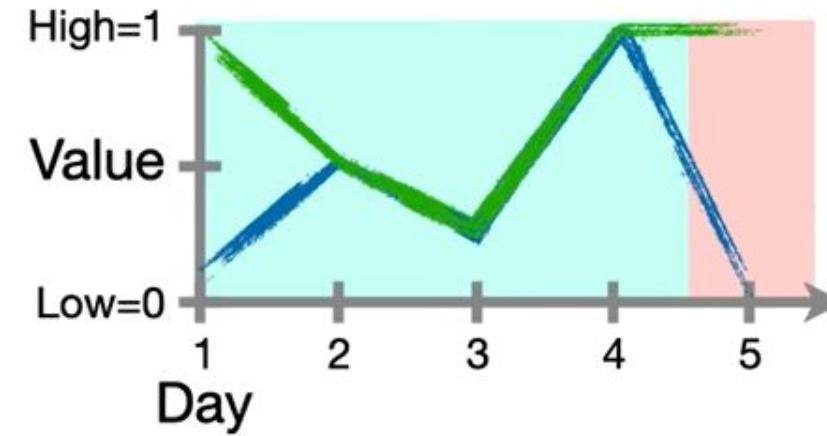


Company A =

Company B =

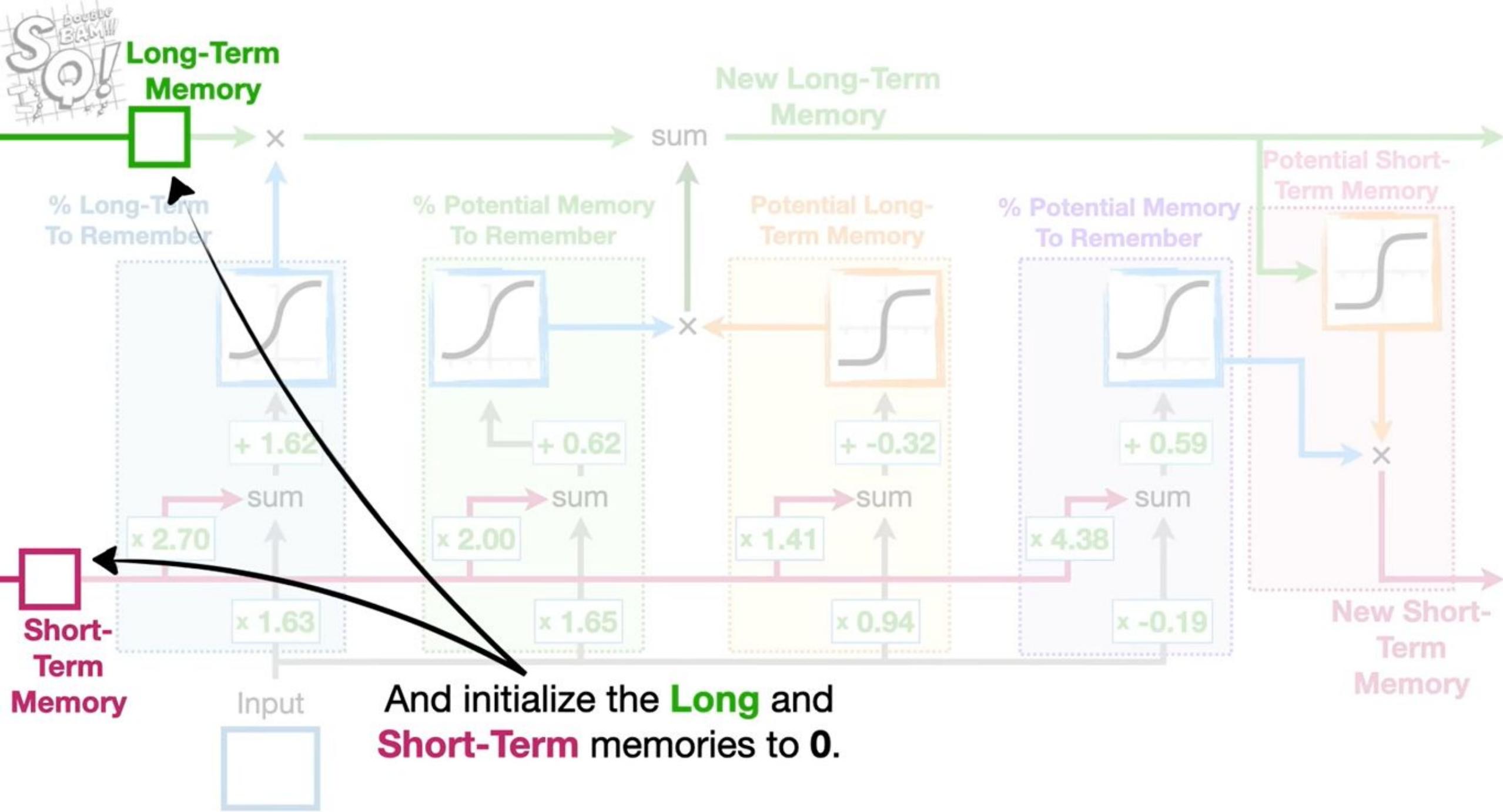


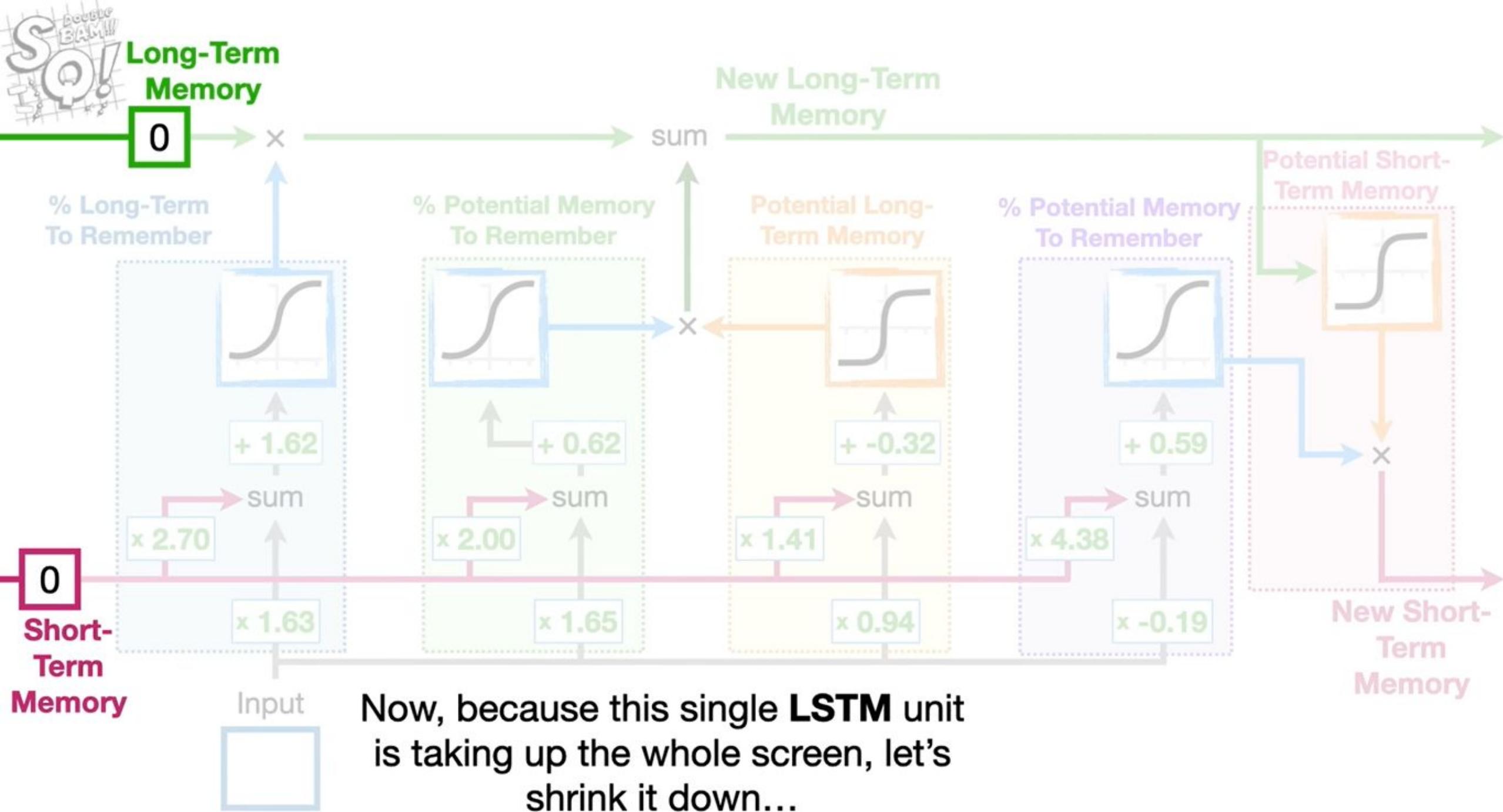
So let's go back to the **LSTM**...

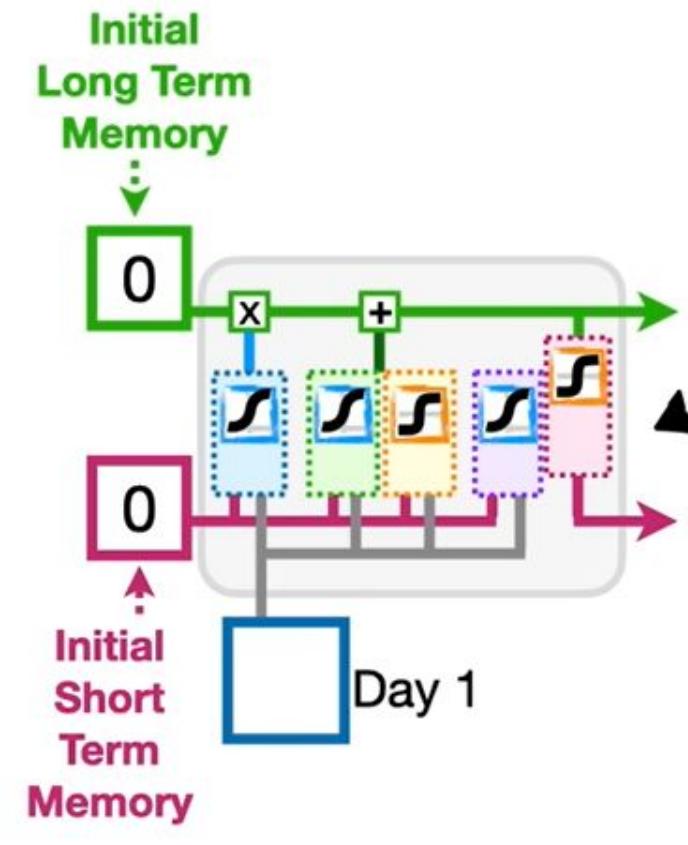


Company A =

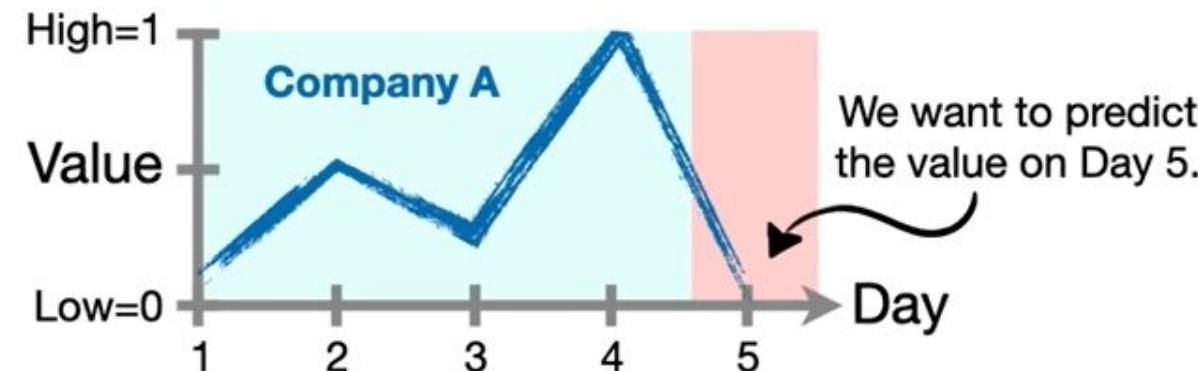
Company B =



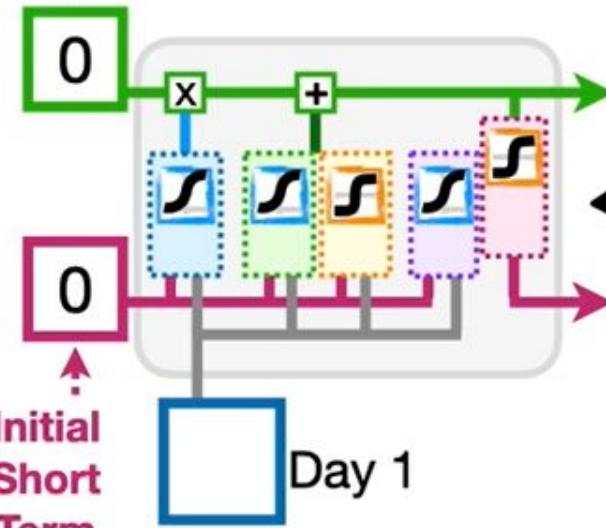




...to this, smaller,
simplified diagram.



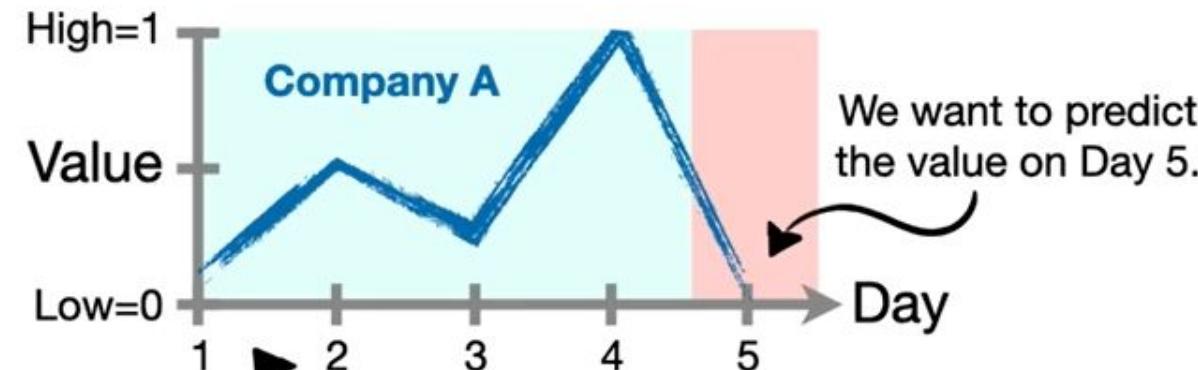
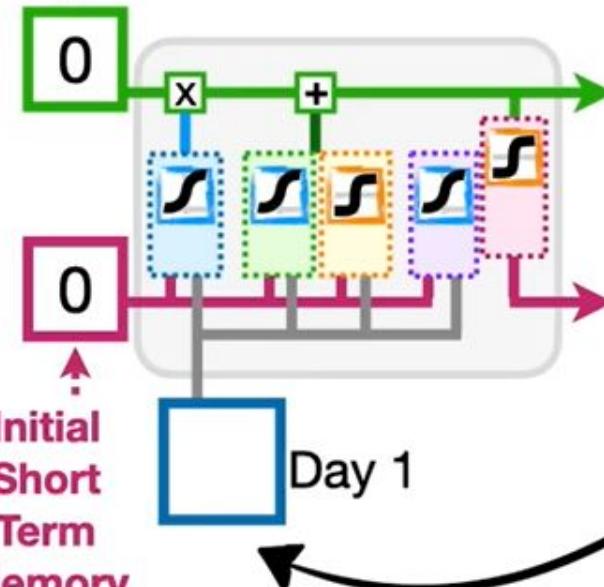
Initial
Long Term
Memory



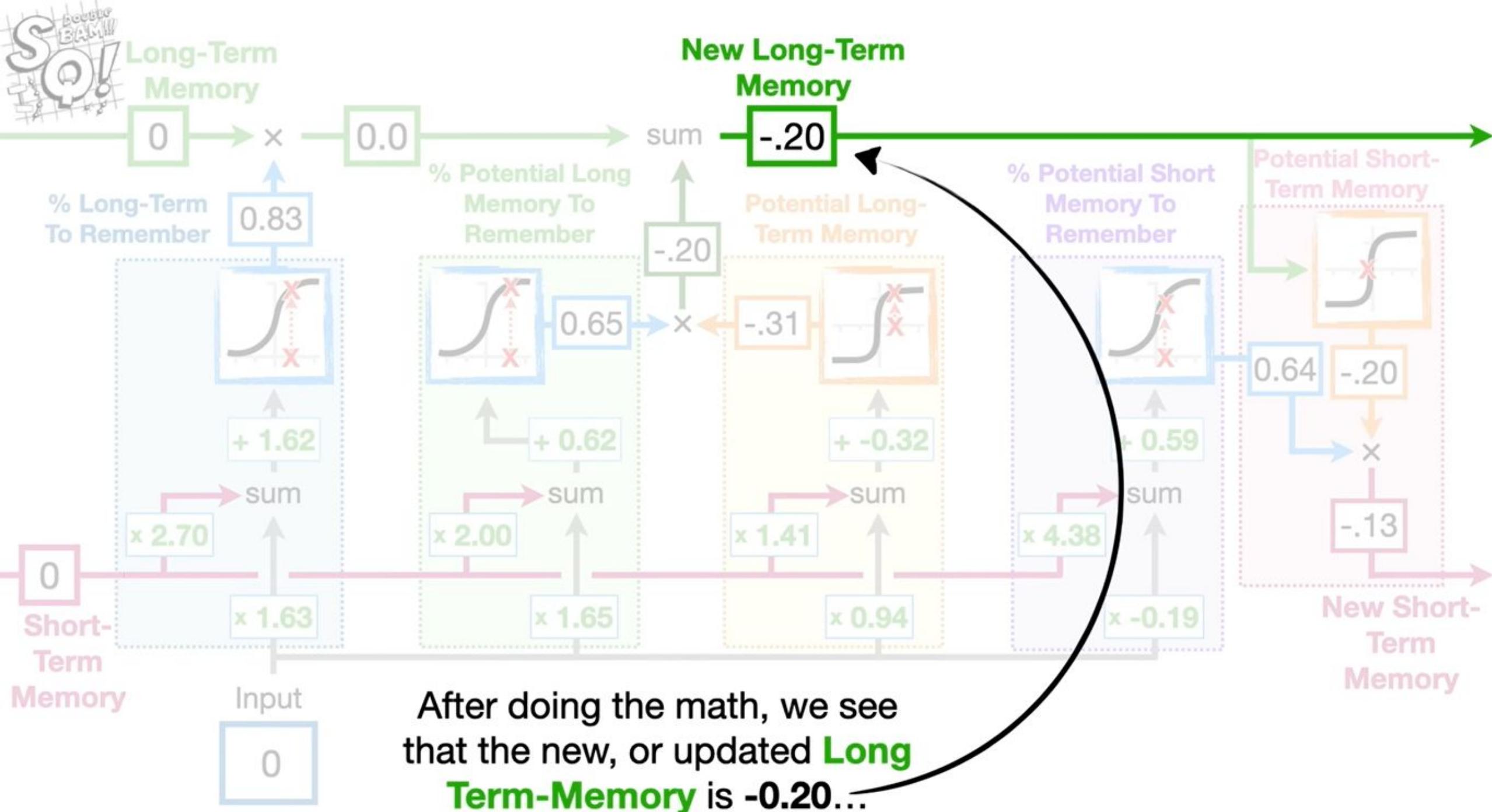
Now, if we want to sequentially run
Company A's values from **Days 1**
through **4** through this **LSTM**...

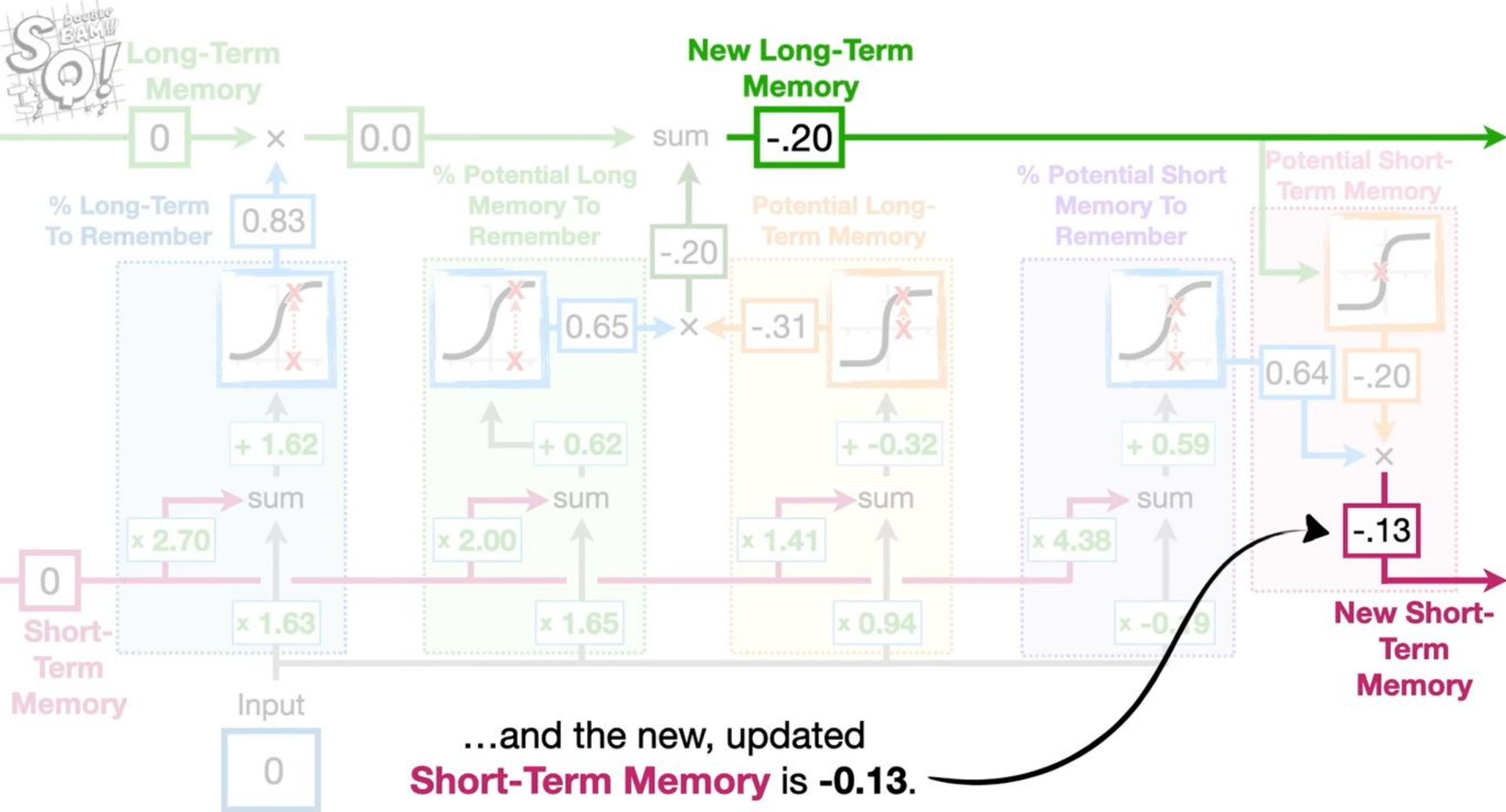


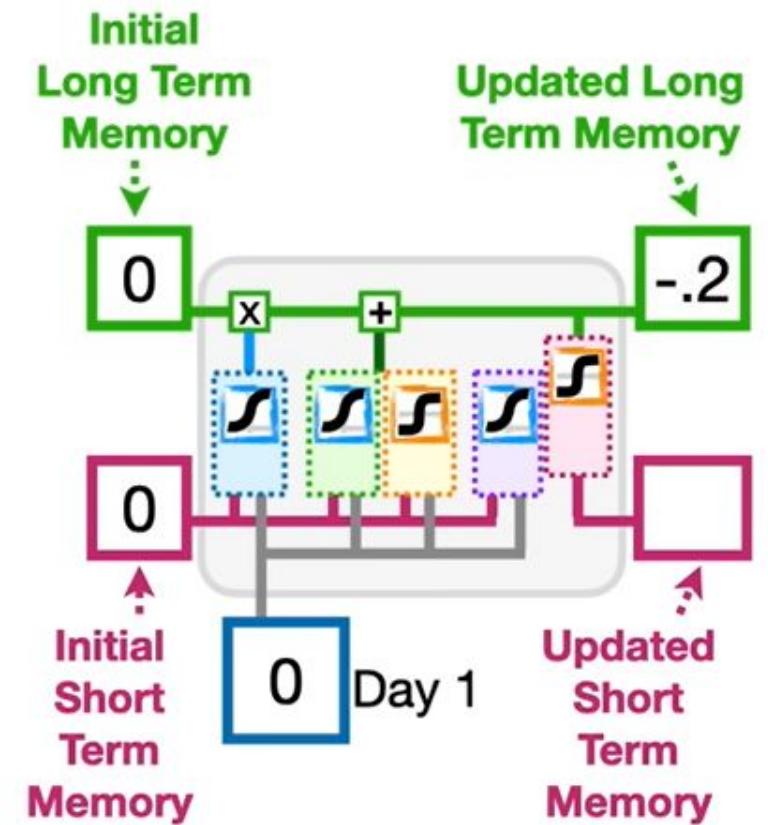
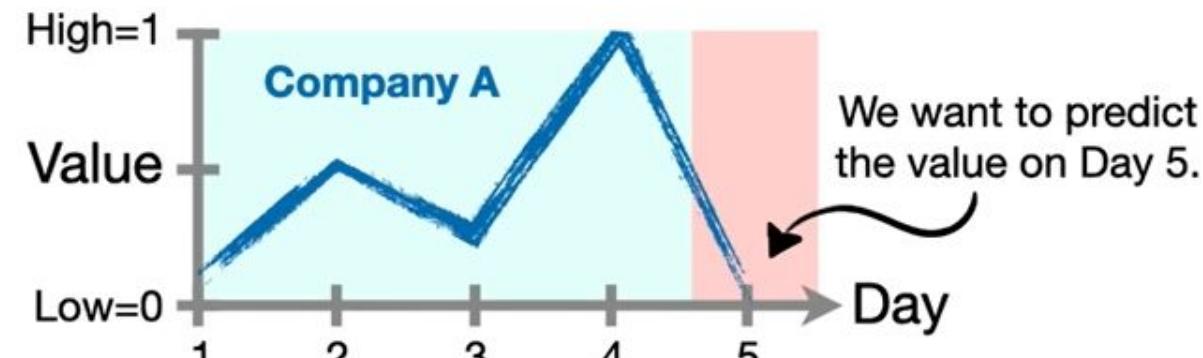
Initial
Long Term
Memory



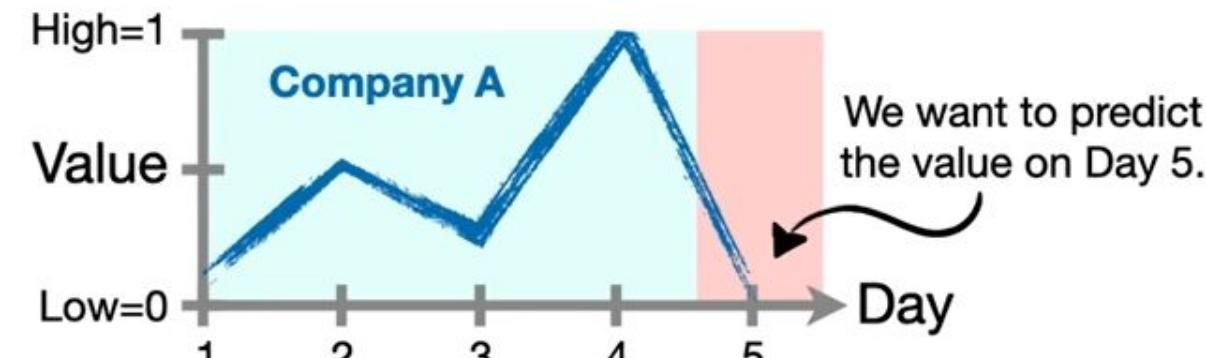
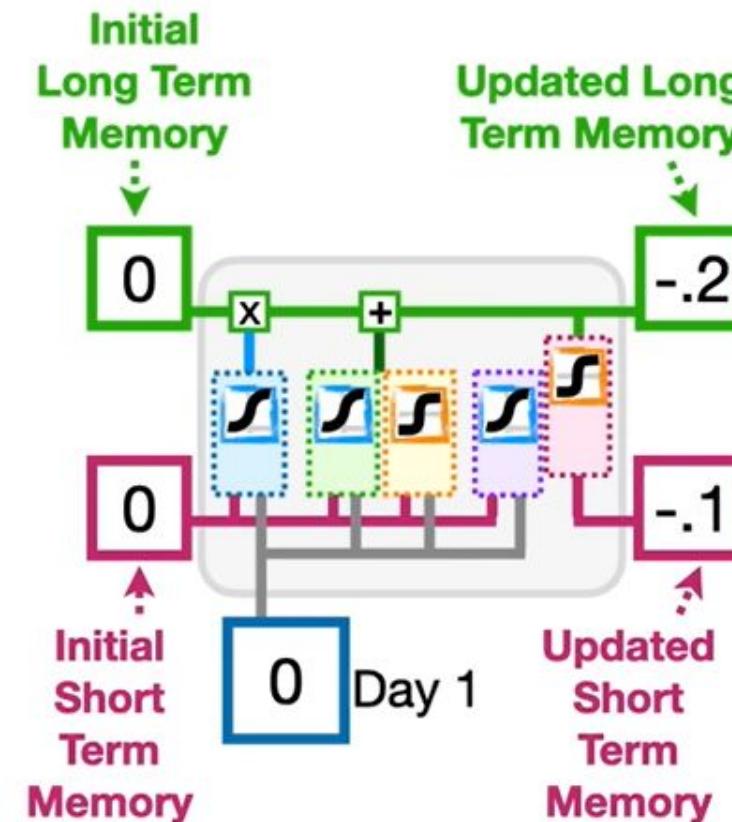
...then we'll start by plugging the value for **Day 1**, which is **0**, into the **Input**.



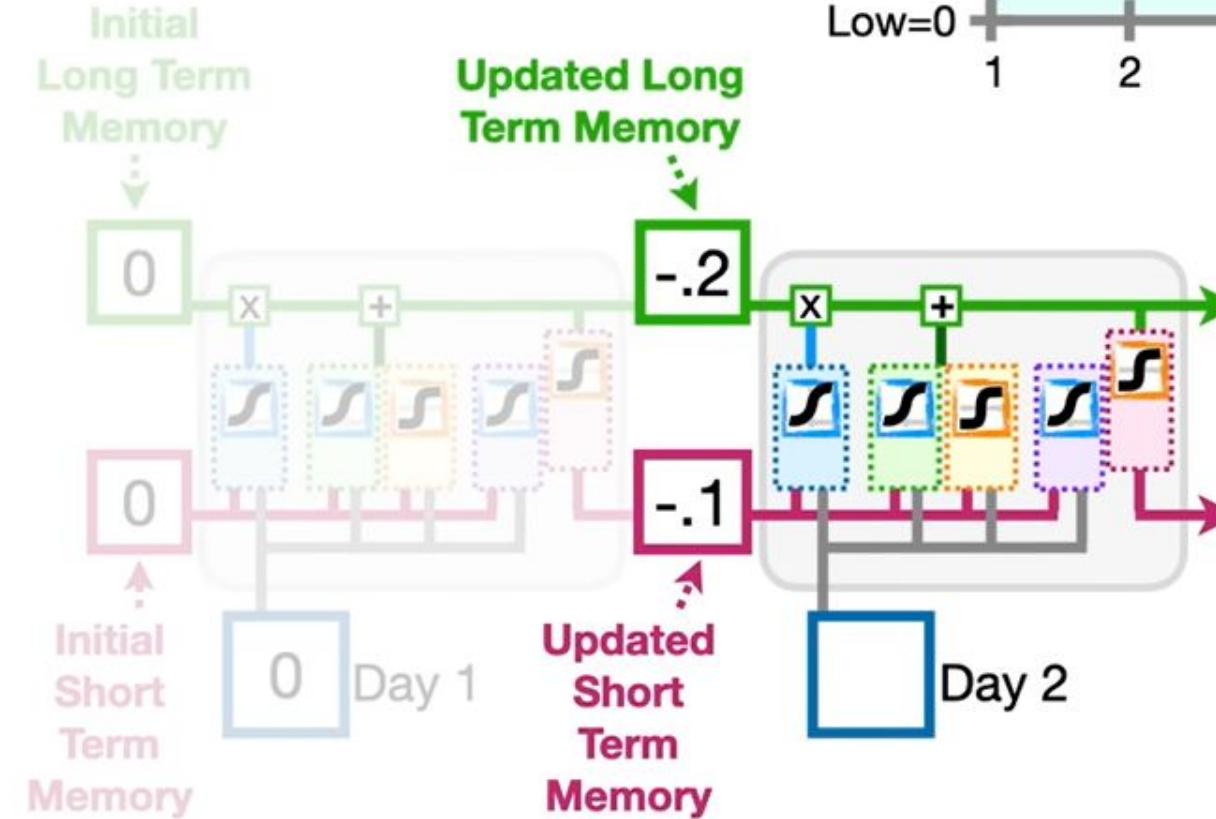
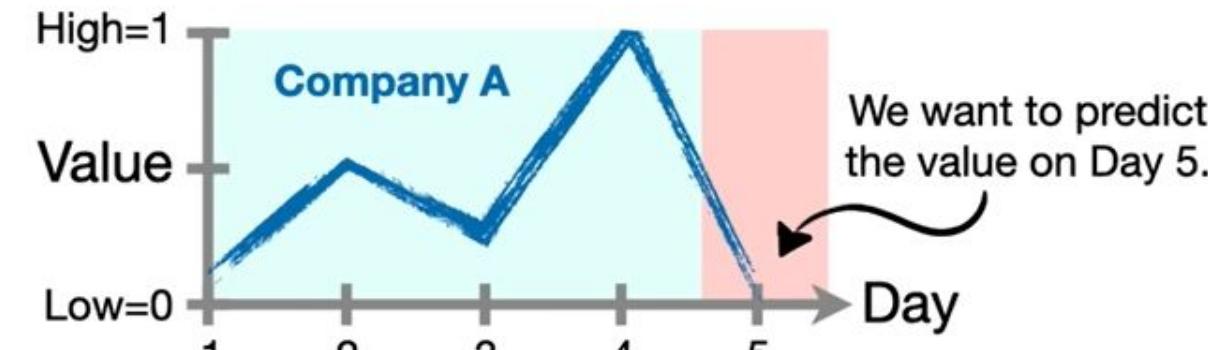




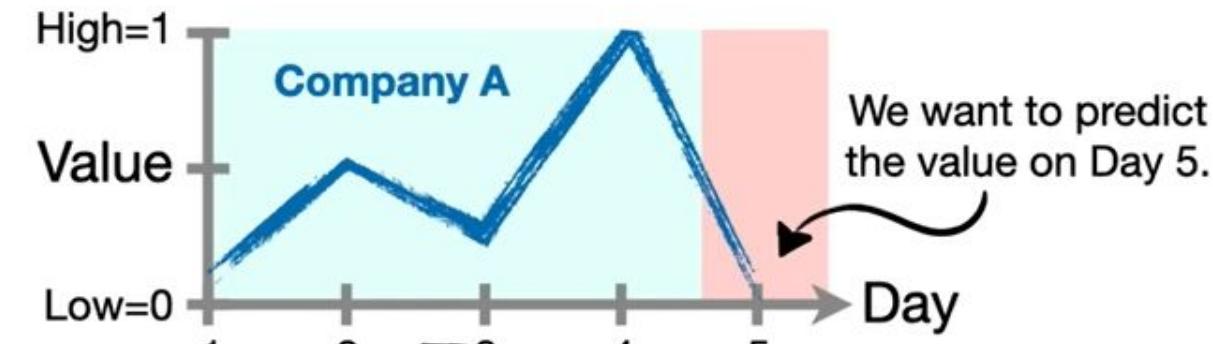
...so we plug in **-0.2** for the updated **Long-Term Memory**...



...and **-0.1** (rounded) for the
updated **Short-Term**
Memory.

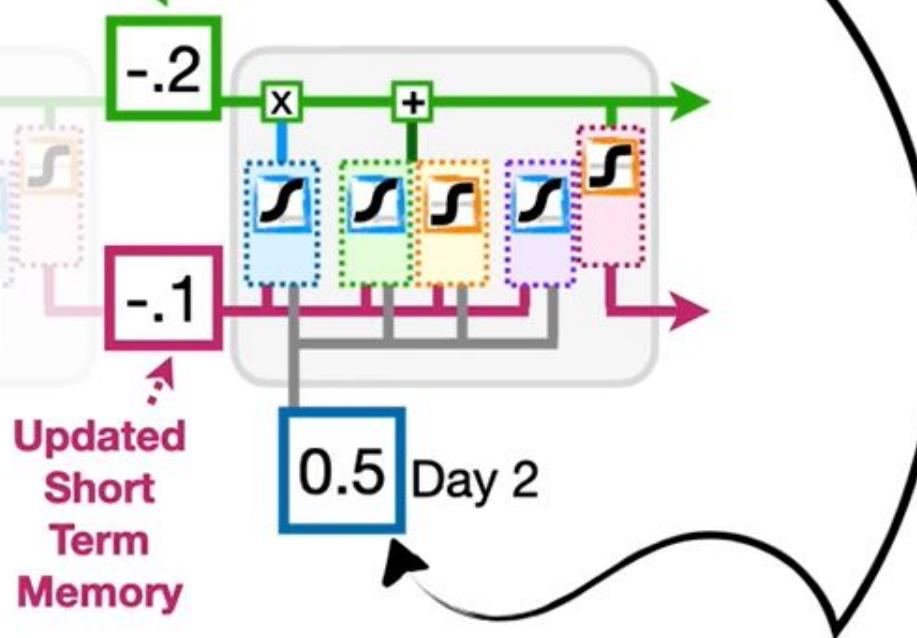
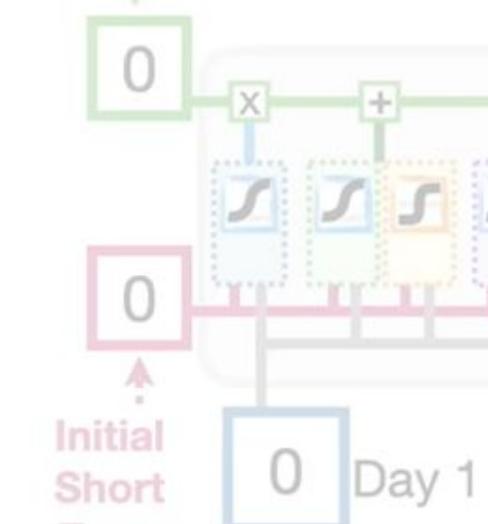


Now we unroll the **LSTM**, using
the updated memories...

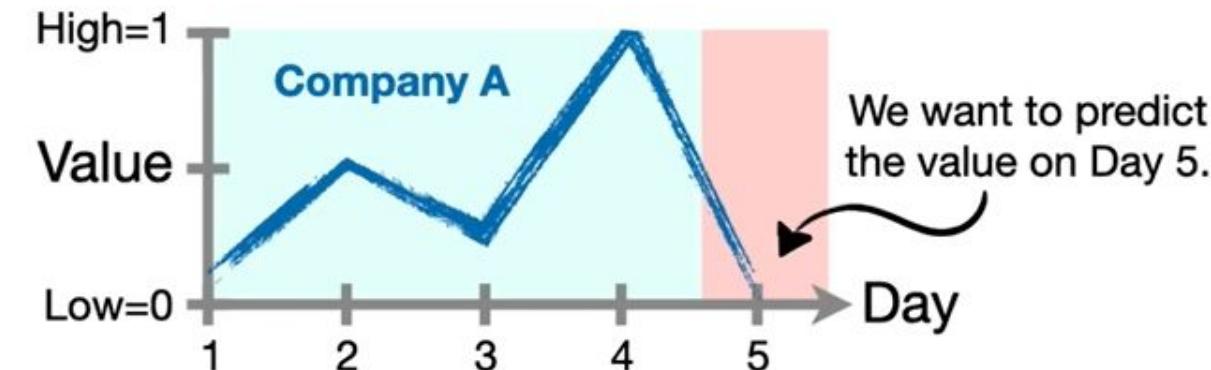


Initial
Long Term
Memory

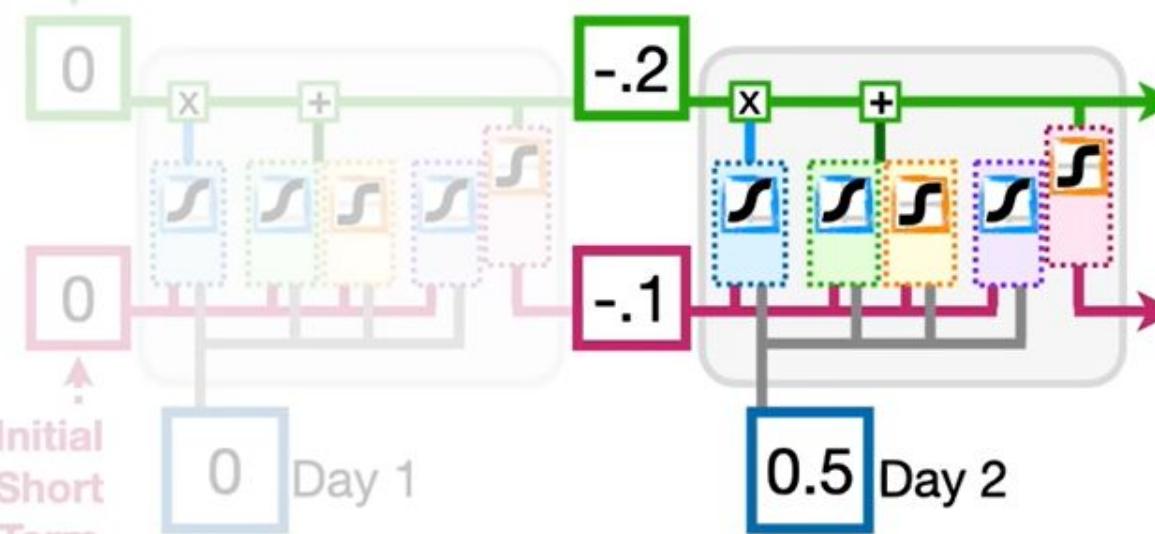
Updated Long
Term Memory



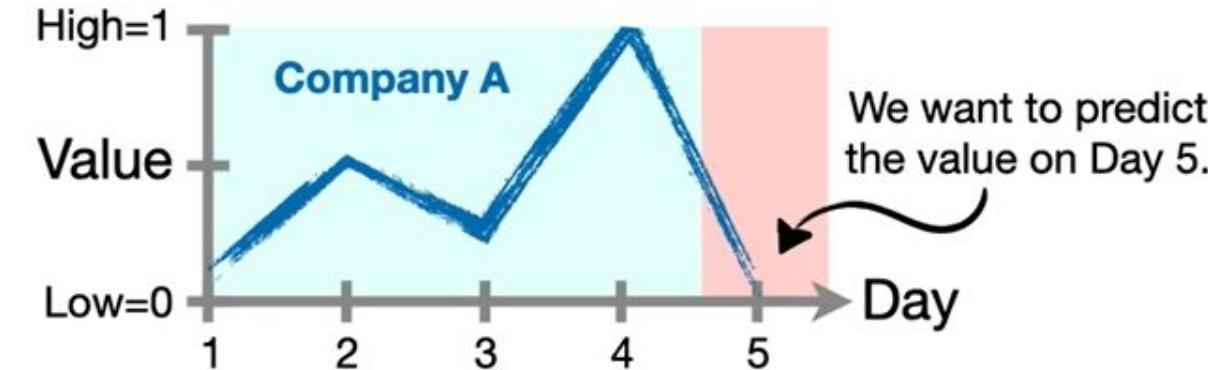
...and plug the value from
Day 2, 0.5, into the Input.



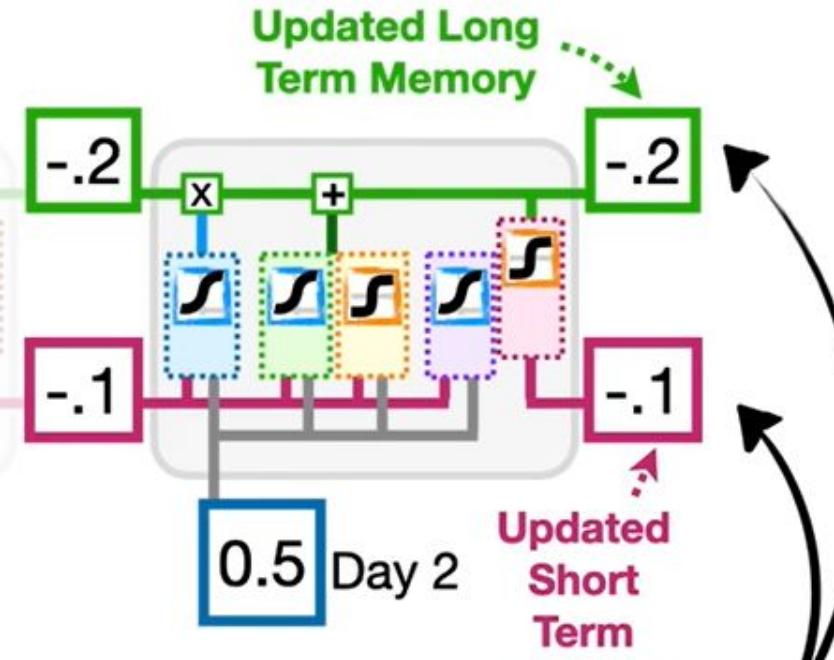
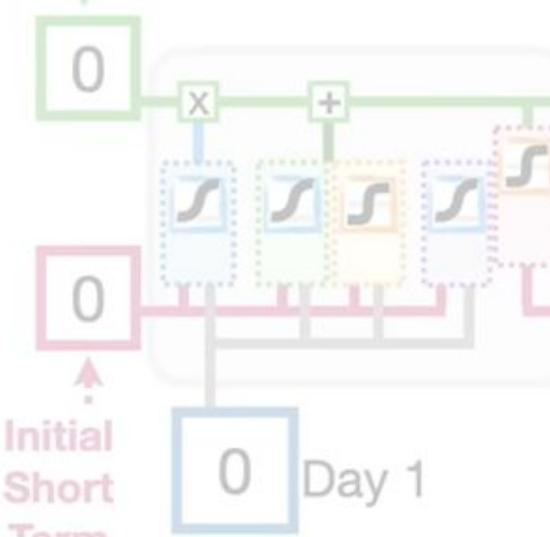
Initial
Long Term
Memory



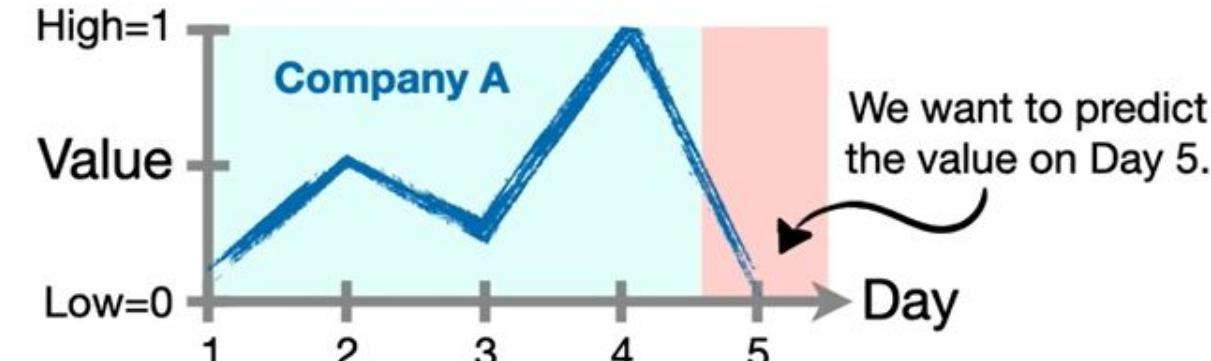
Then the **LSTM** does it's math,
using the exact same **Weights**
and **Biases** as before...



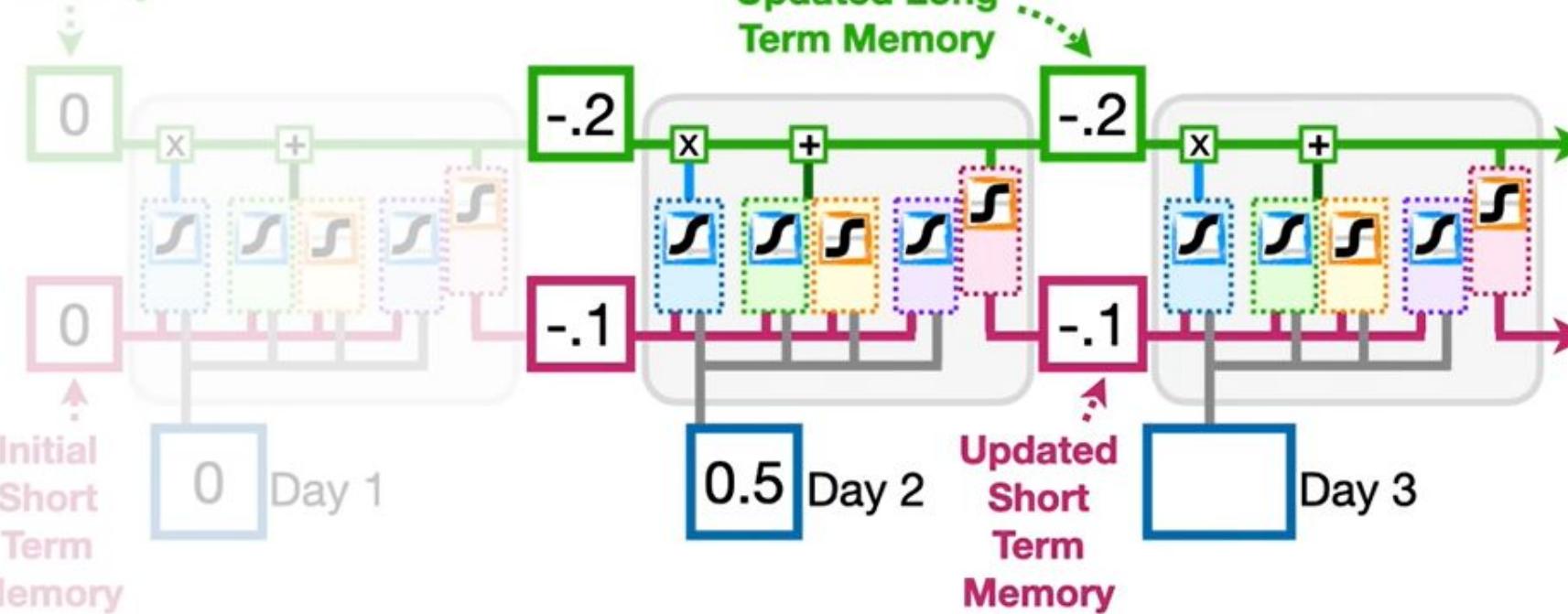
Initial
Long Term
Memory



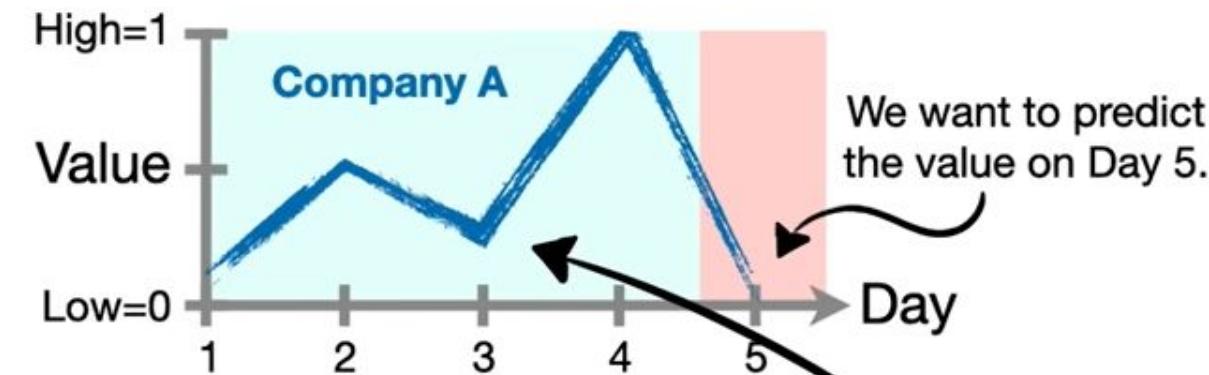
...and we end up with these updated
Long and Short-Term Memories.



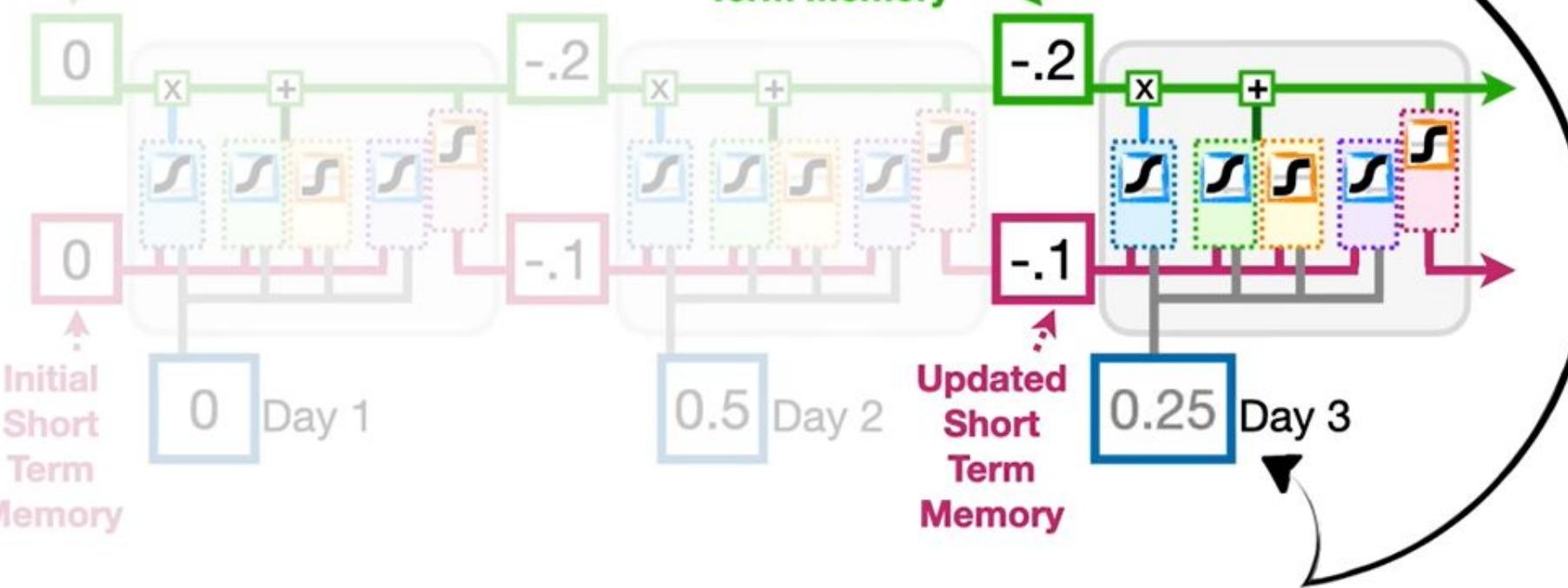
Initial
Long Term
Memory



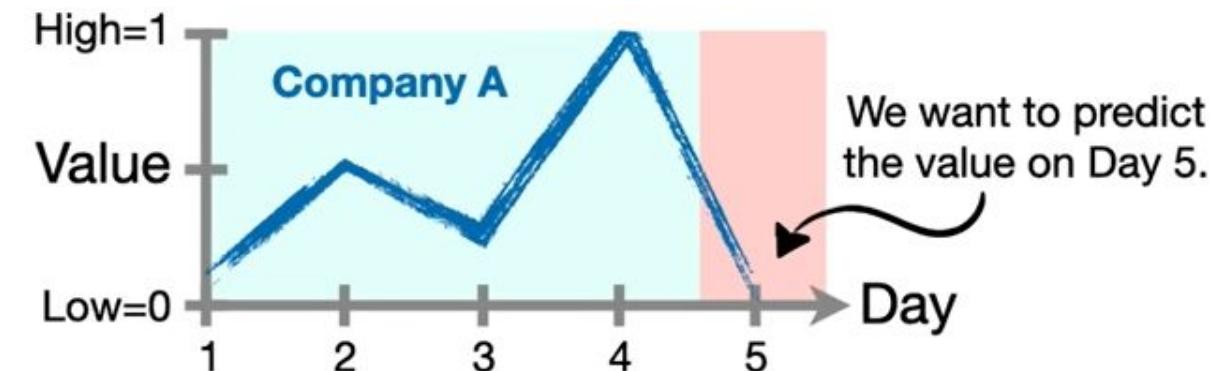
Anyway, we unroll the **LSTM** again...



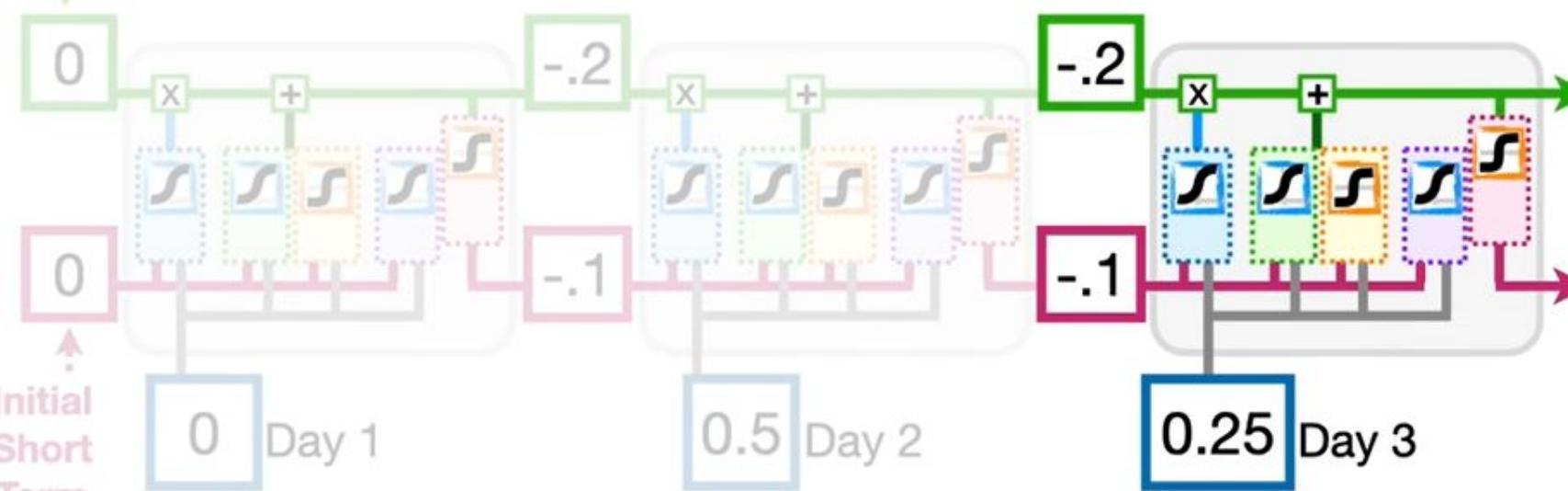
Initial
Long Term
Memory



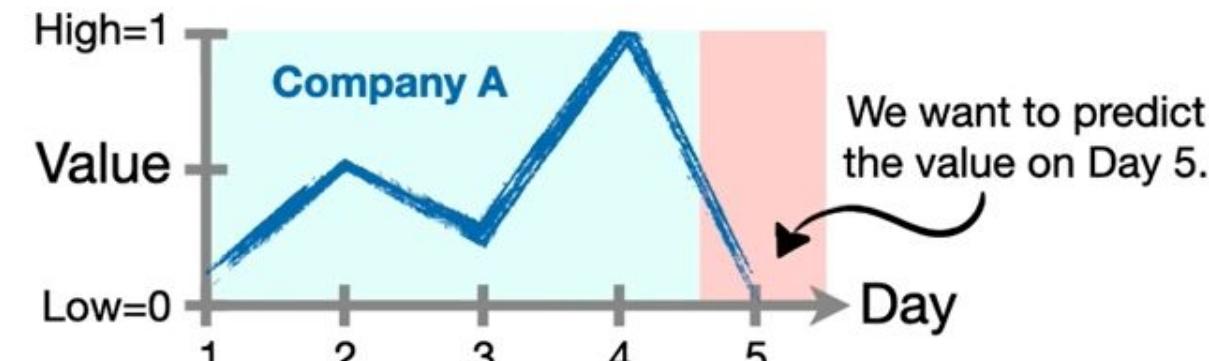
...and plug in the value for **Day 3**.



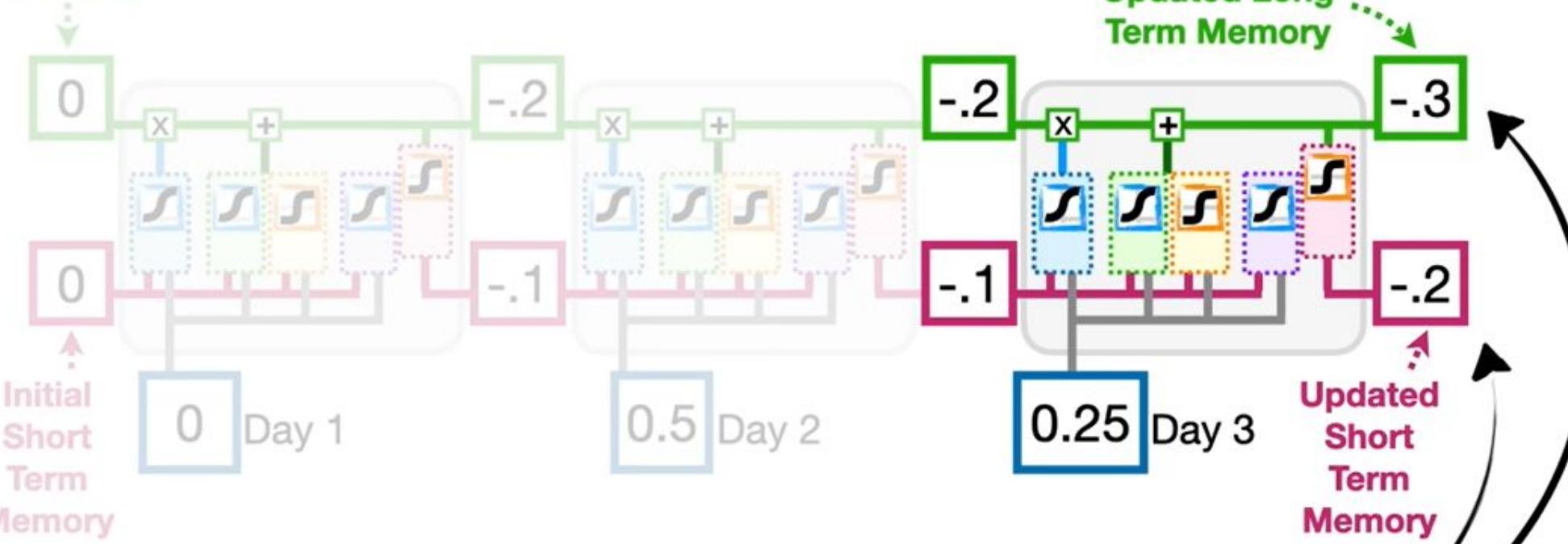
Initial
Long Term
Memory



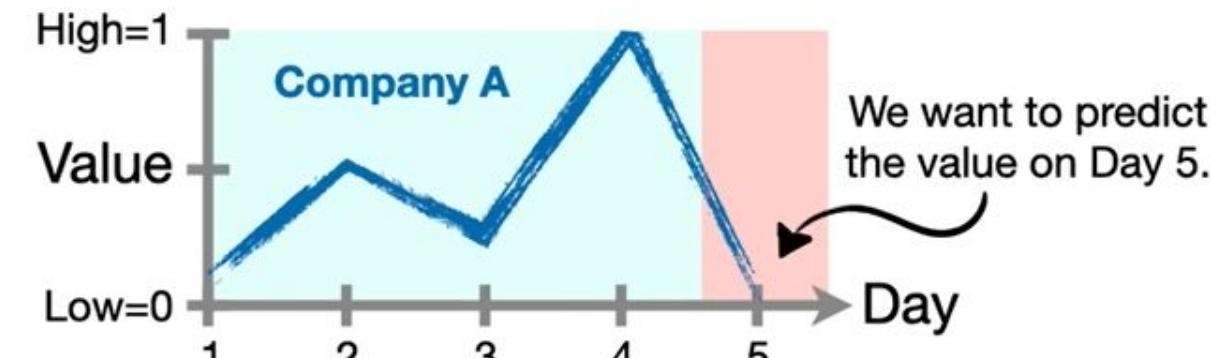
And then the **LSTM** does the
math, again, using the exact
same **Weights and Biases...**



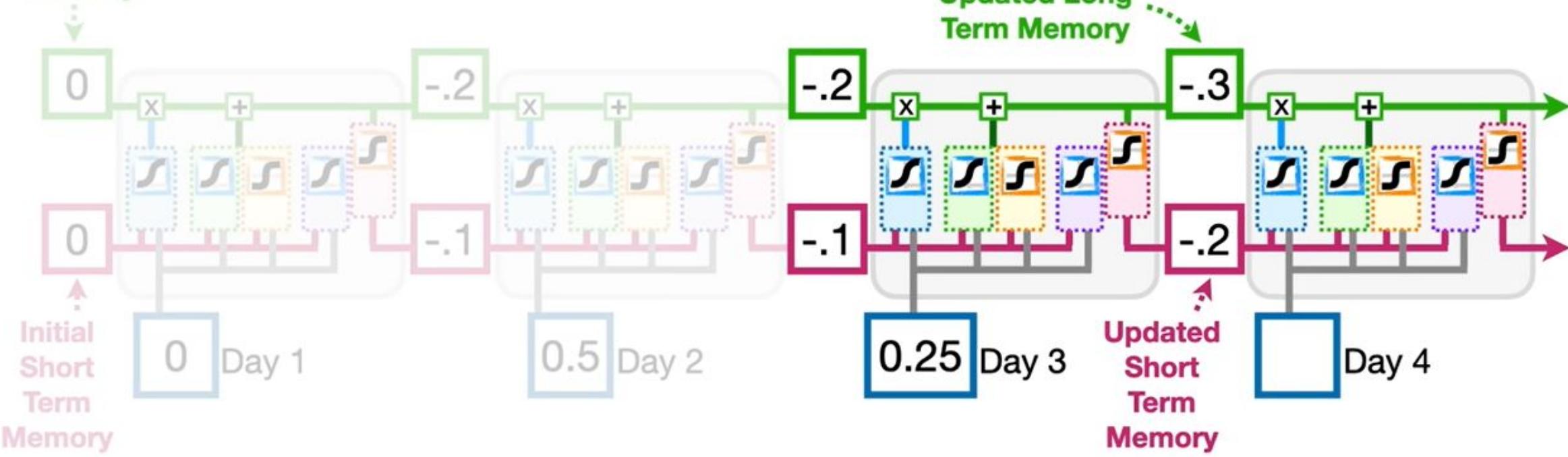
Initial
Long Term
Memory



...and gives us these
updated memories.



Initial
Long Term
Memory



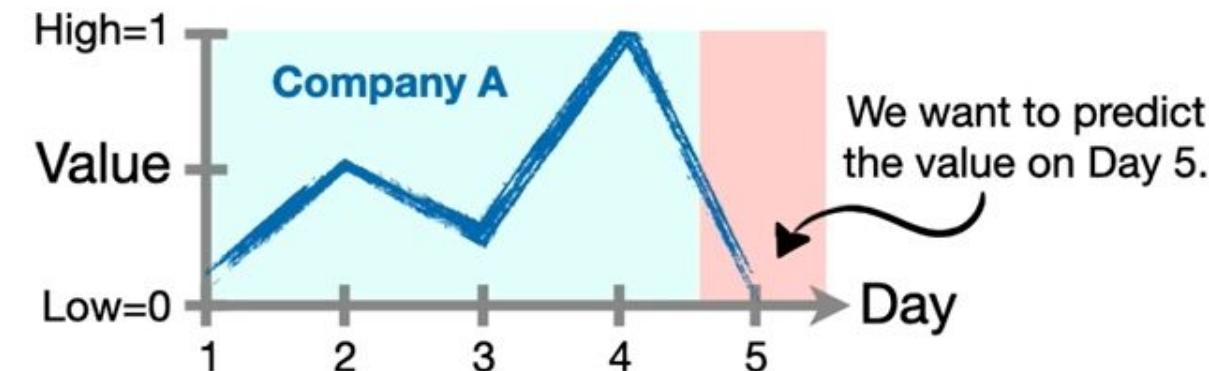
Then we unroll the **LSTM**
one last time...



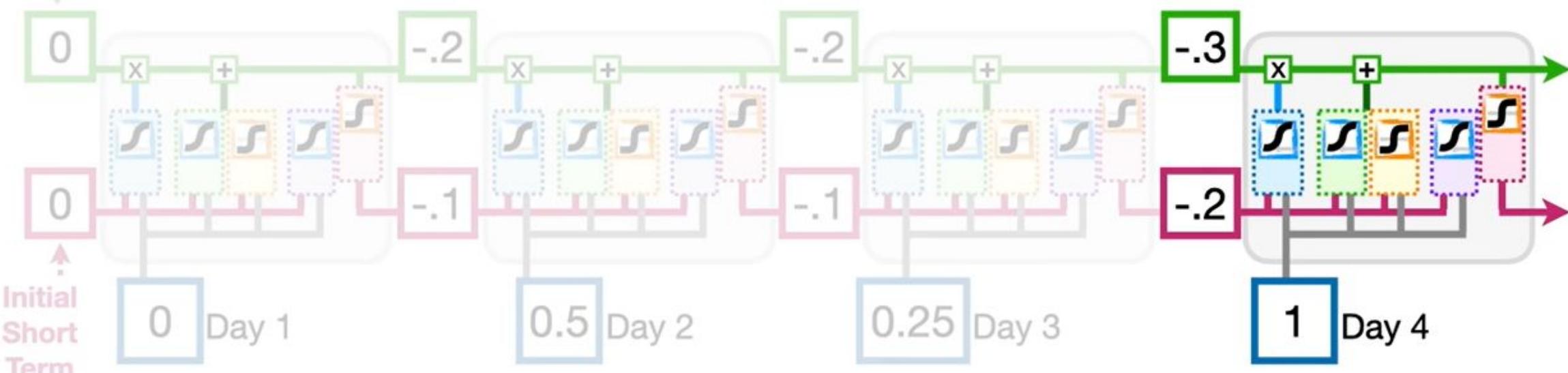
Initial
Long Term
Memory



...and plug in the
value for **Day 4**.



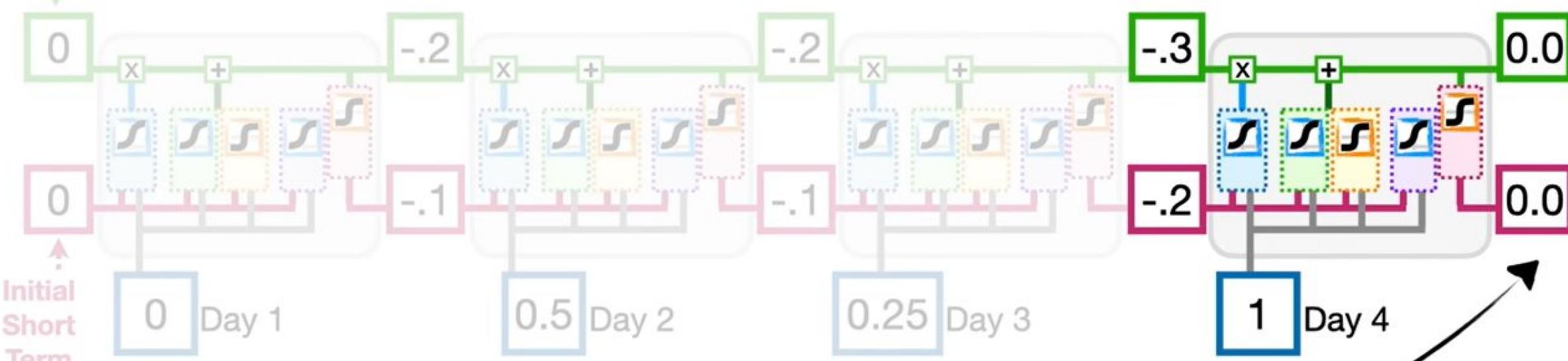
Initial
Long Term
Memory



And the **LSTM** does the math,
again, using the exact same
Weights and Biases...



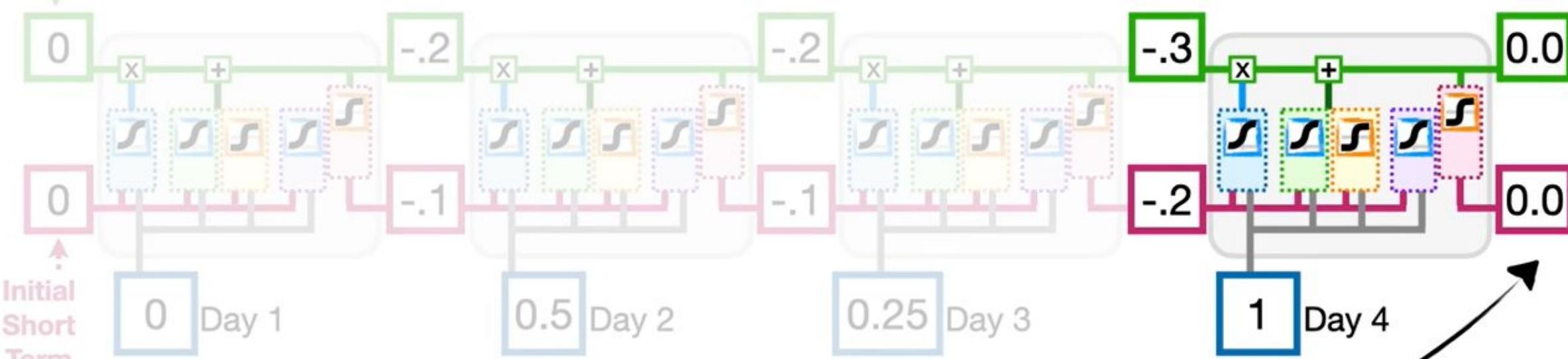
Initial
Long Term
Memory



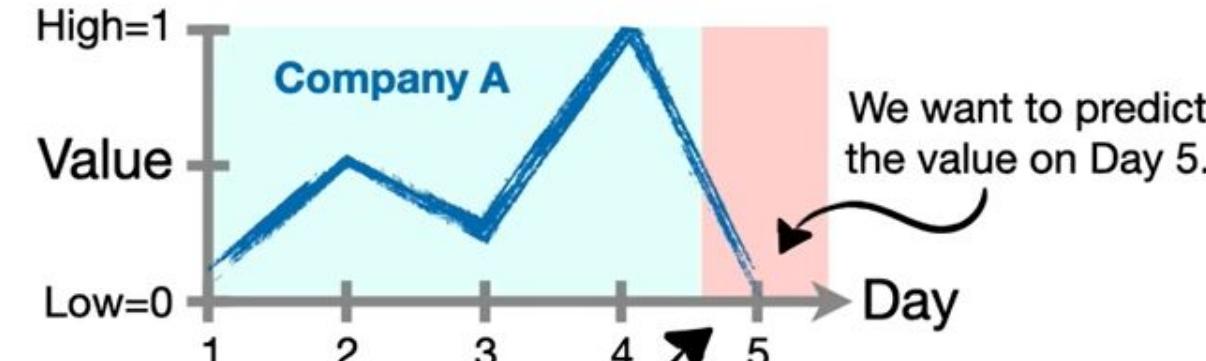
...and gives us the
final memories.



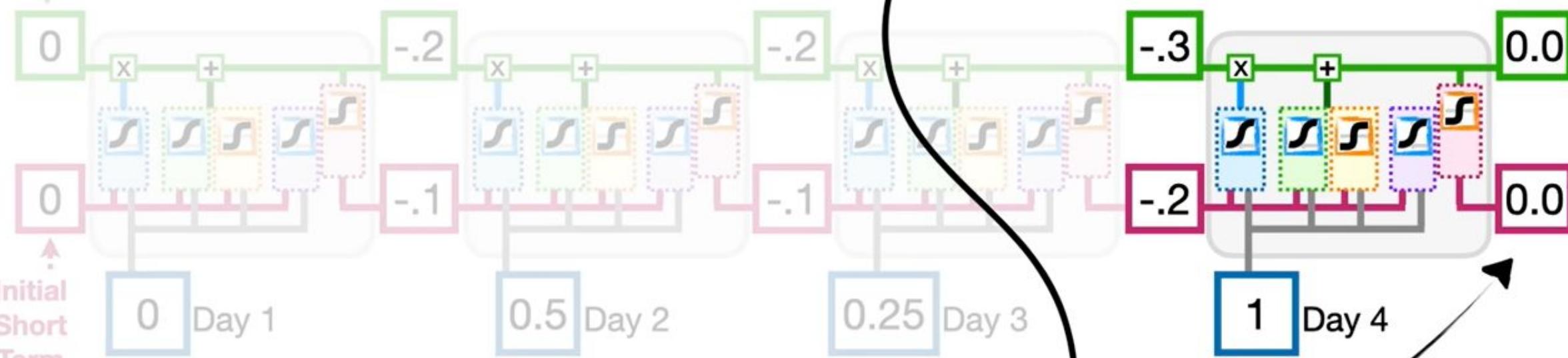
Initial
Long Term
Memory



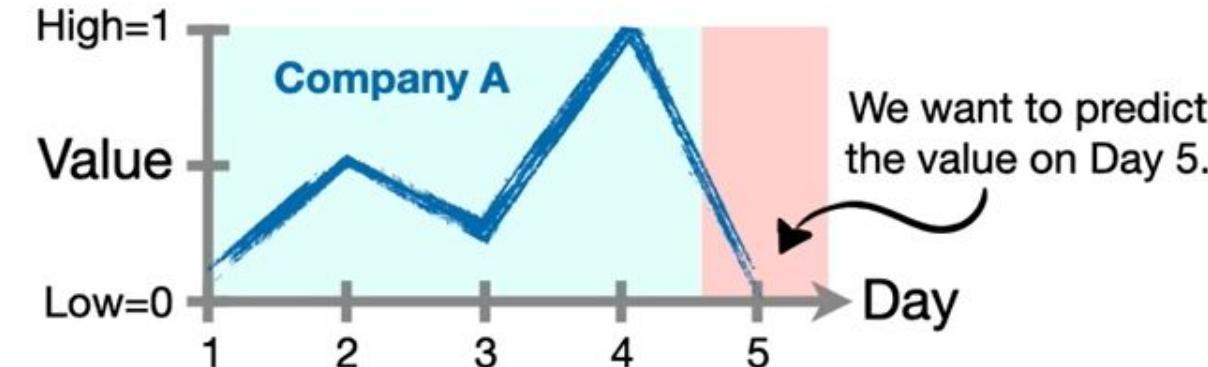
And the final **Short-Term
Memory**, **0.0**, is the **Output**
from the unrolled **LSTM**.



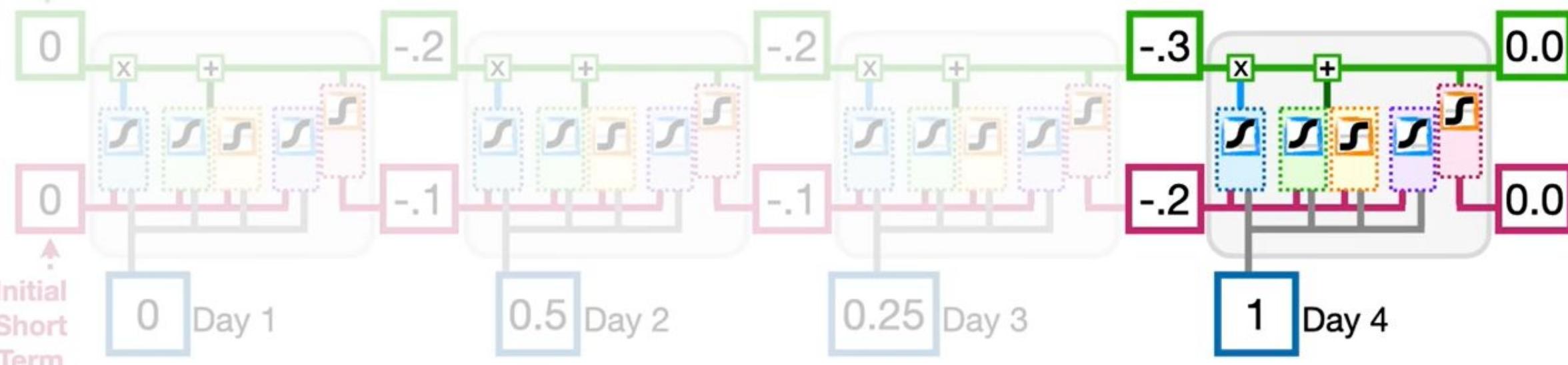
Initial
Long Term
Memory



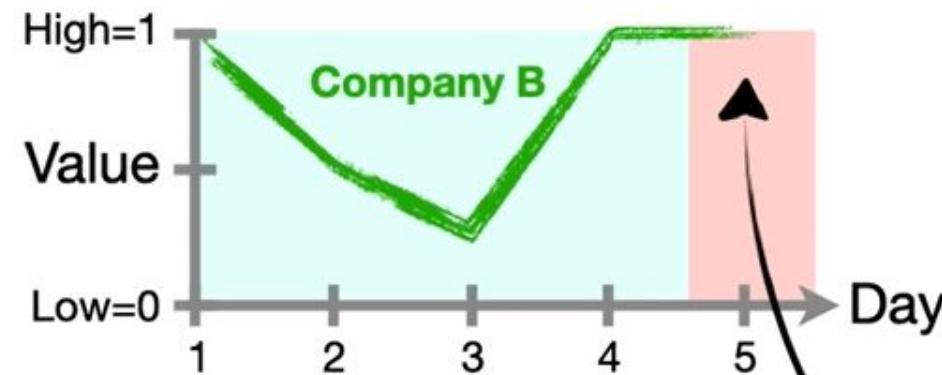
And that means that the **Output**
from the **LSTM** correctly predicts
Company A's value for Day 5.



Initial
Long Term
Memory



Now that we have shown that the
LSTM can correctly predict the value
on **Day 5** for **Company A**...

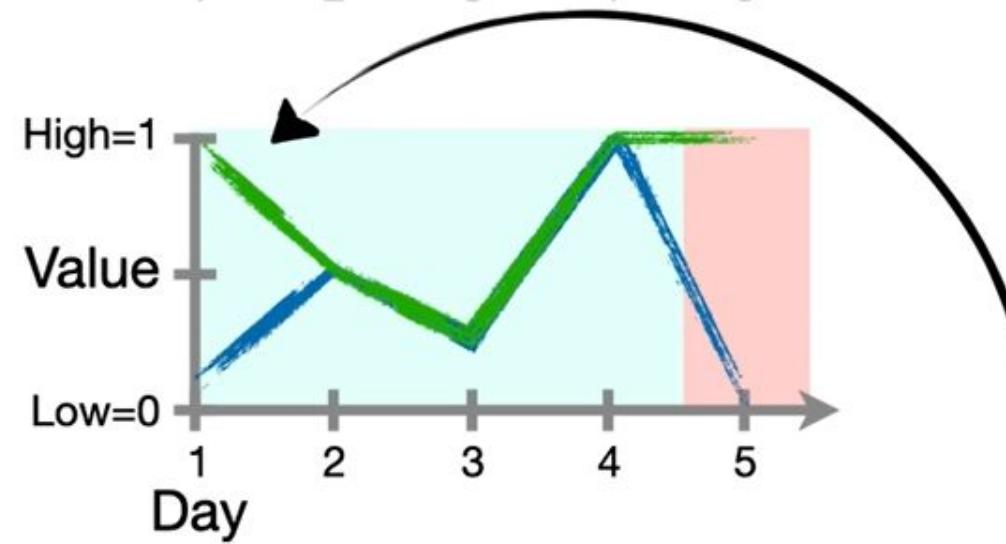


...let's show how the same **LSTM**, with the same **Weights** and **Biases**, can correctly predict the value on **Day 5** for **Company B**.



Company A =

Company B =

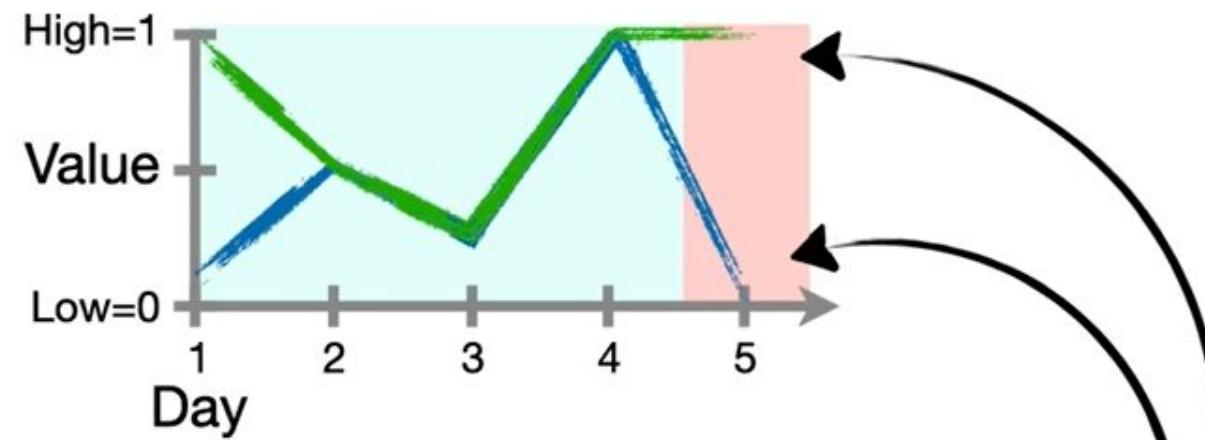


NOTE: Remember, on **Days 1** through **4**, the only difference between the companies occurs on **Day 1**...

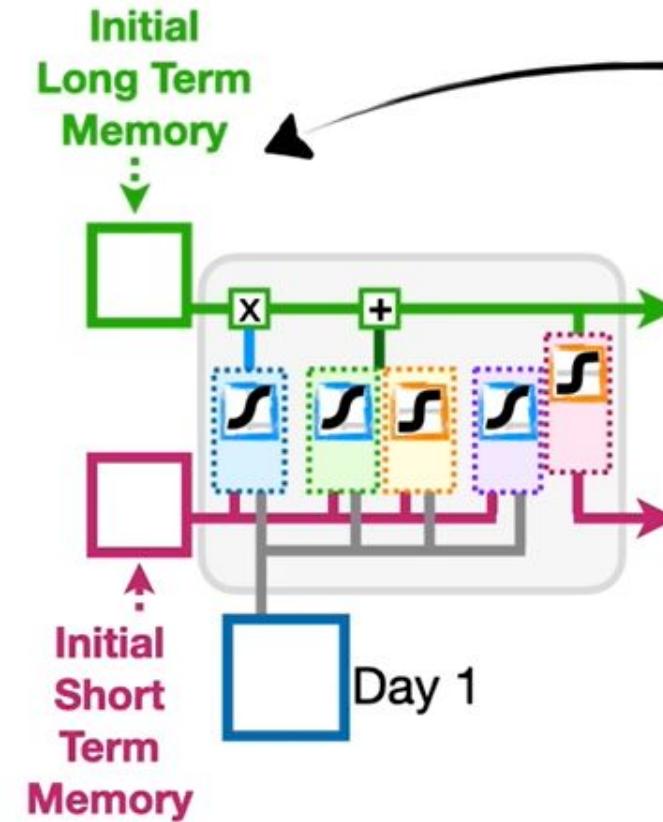


Company A =

Company B =



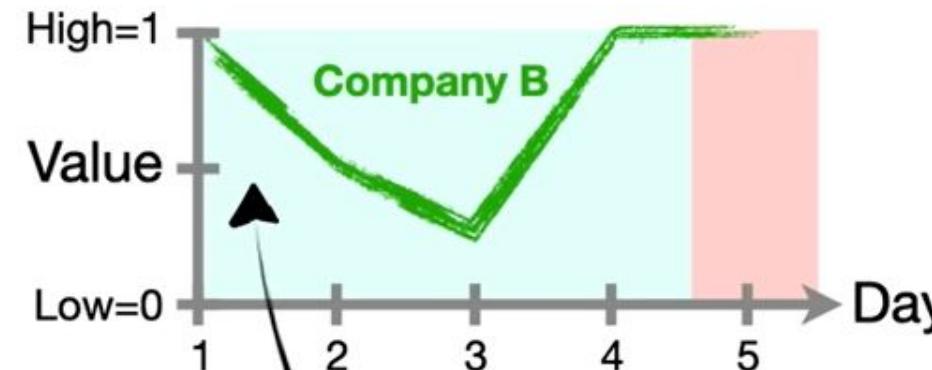
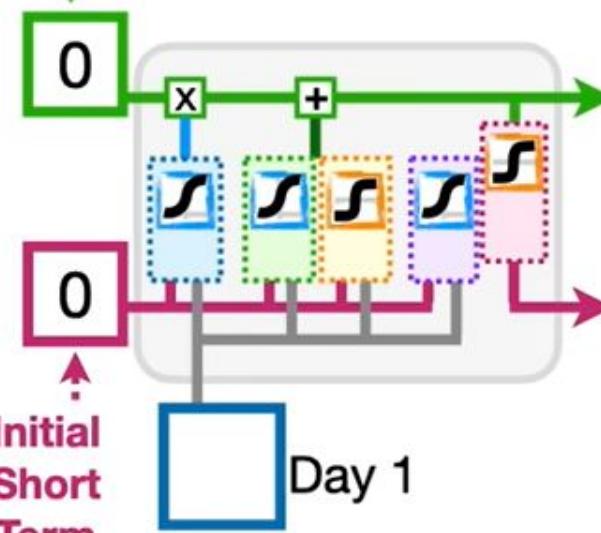
...and that means the **LSTM** has to remember
what happened on **Day 1** in order to correctly
predict the different output values on **Day 5**.



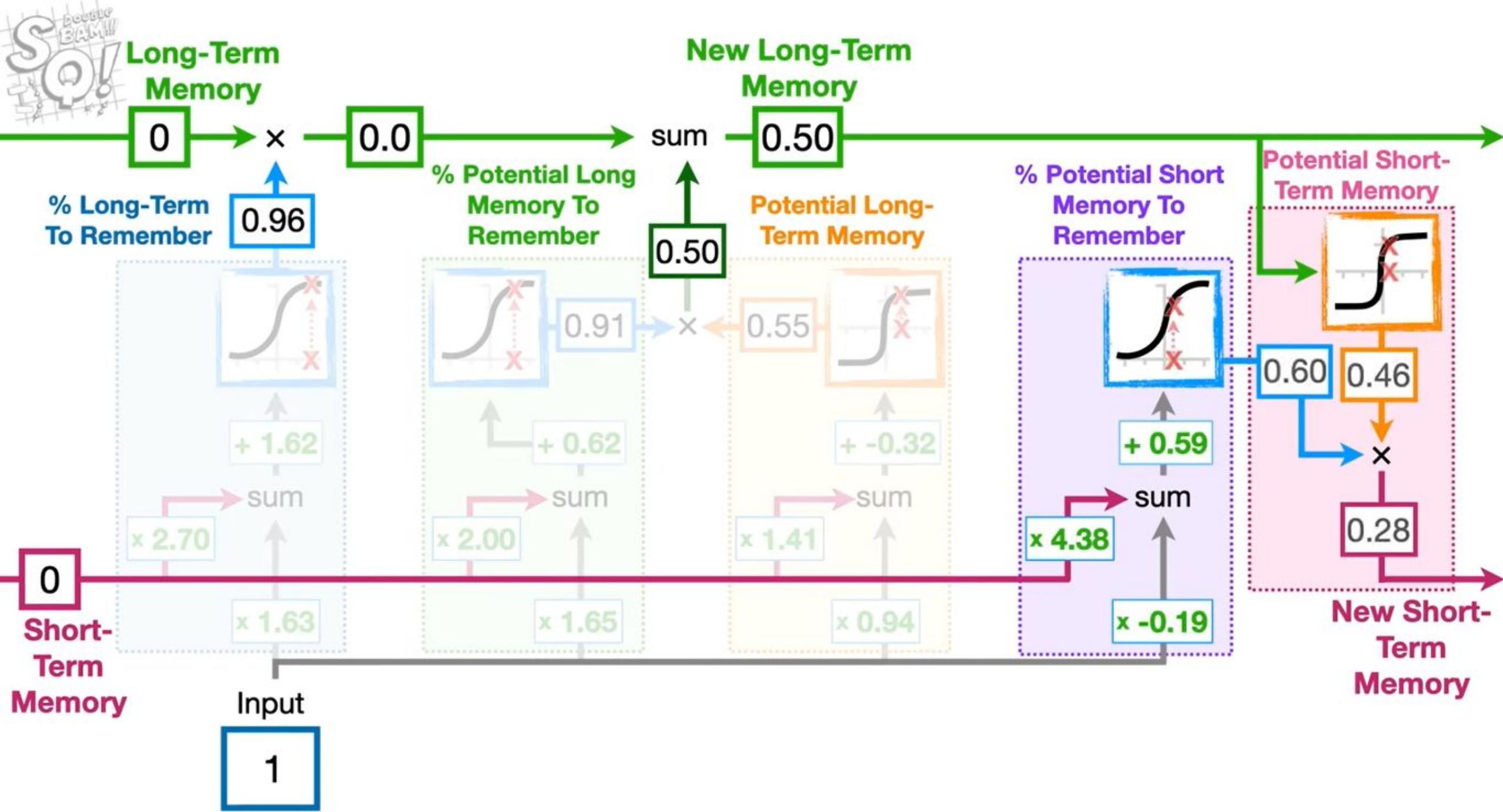
So let's start by initializing the **Long** and **Short-Term Memories** to 0.

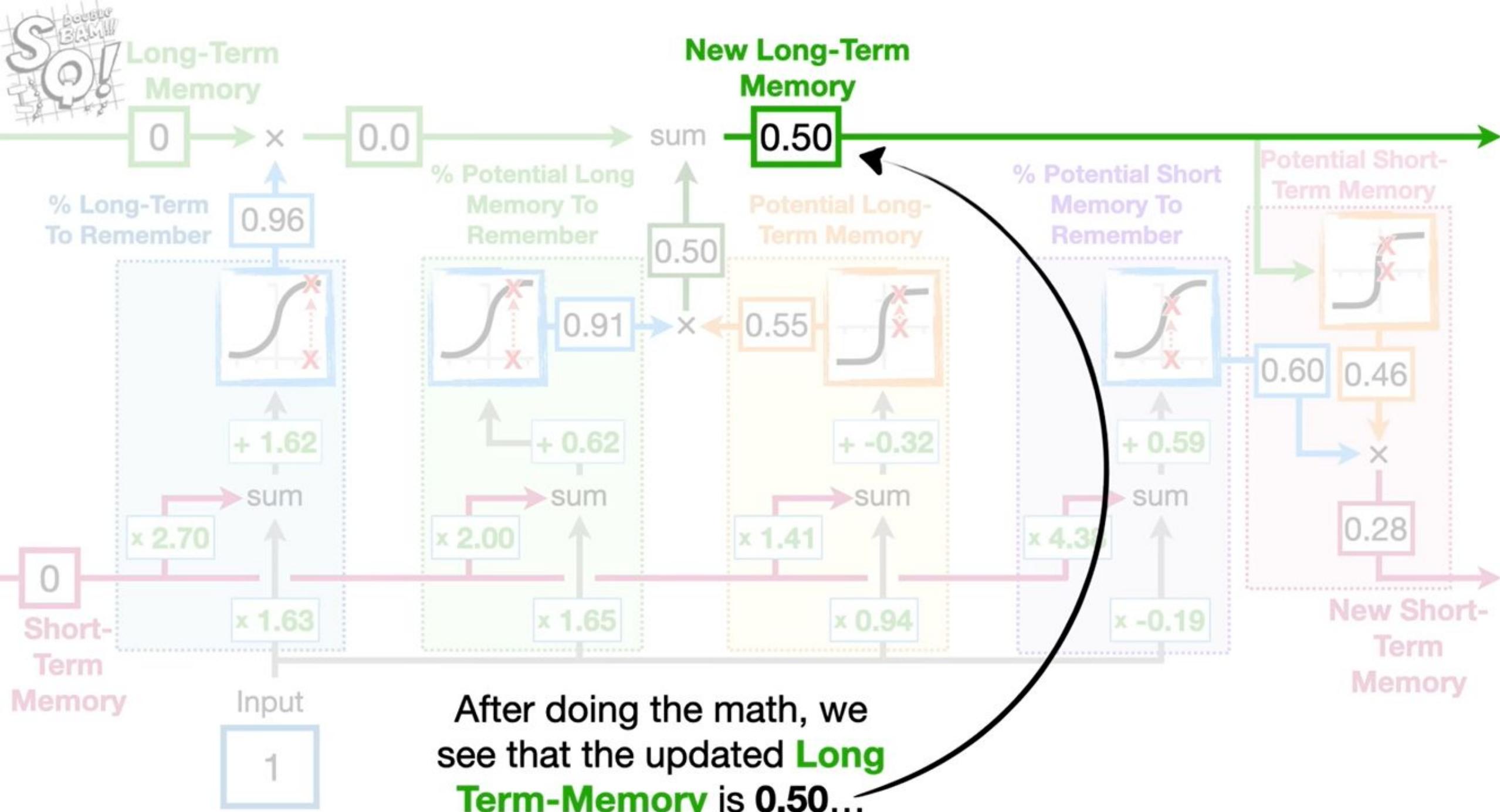


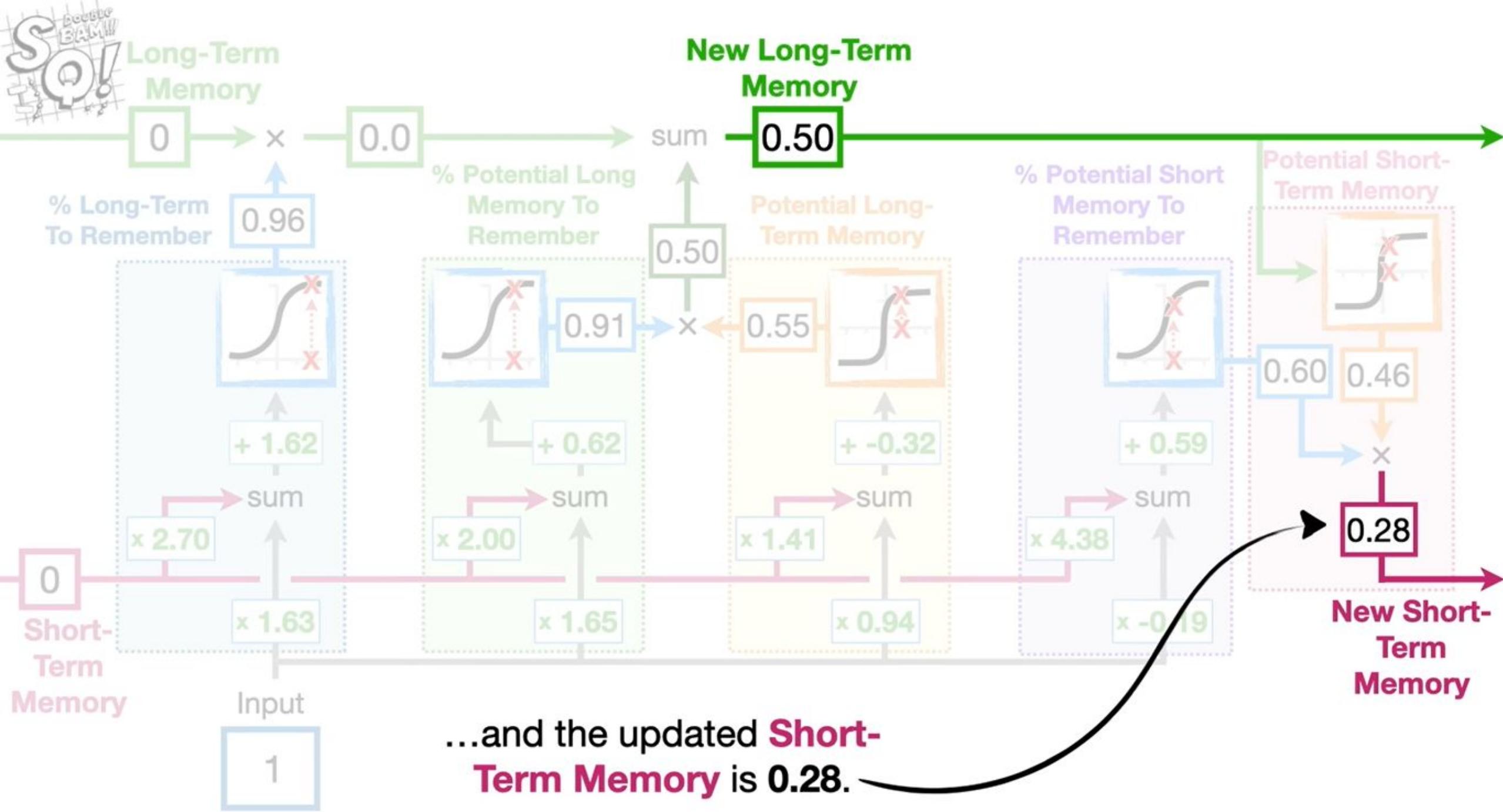
Initial
Long Term
Memory

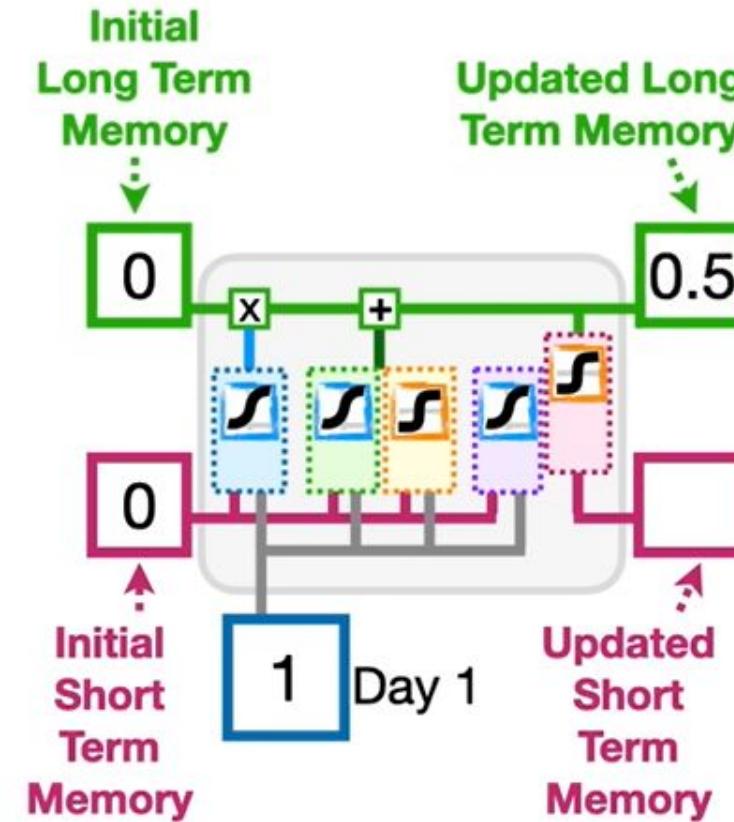
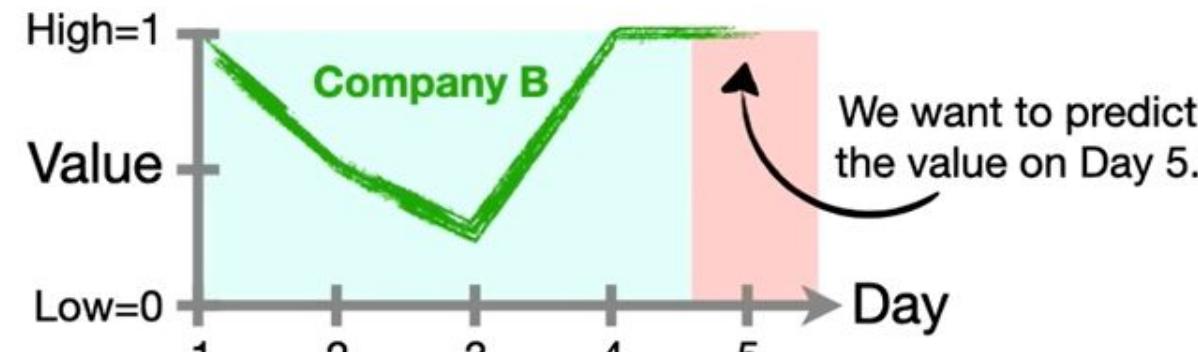


Now let's plug in the
value for **Day 1** from
Company B, 1...

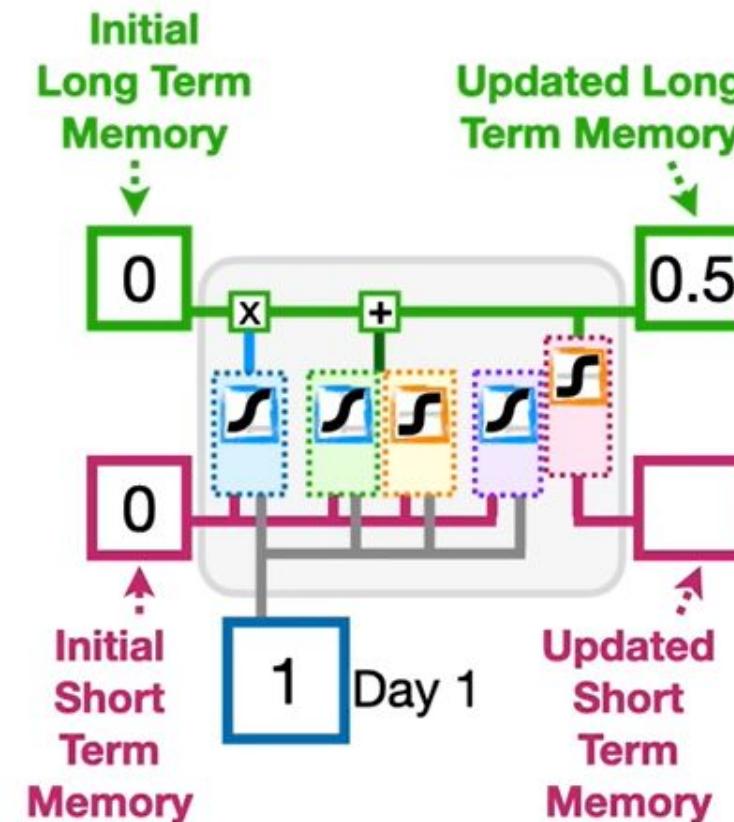




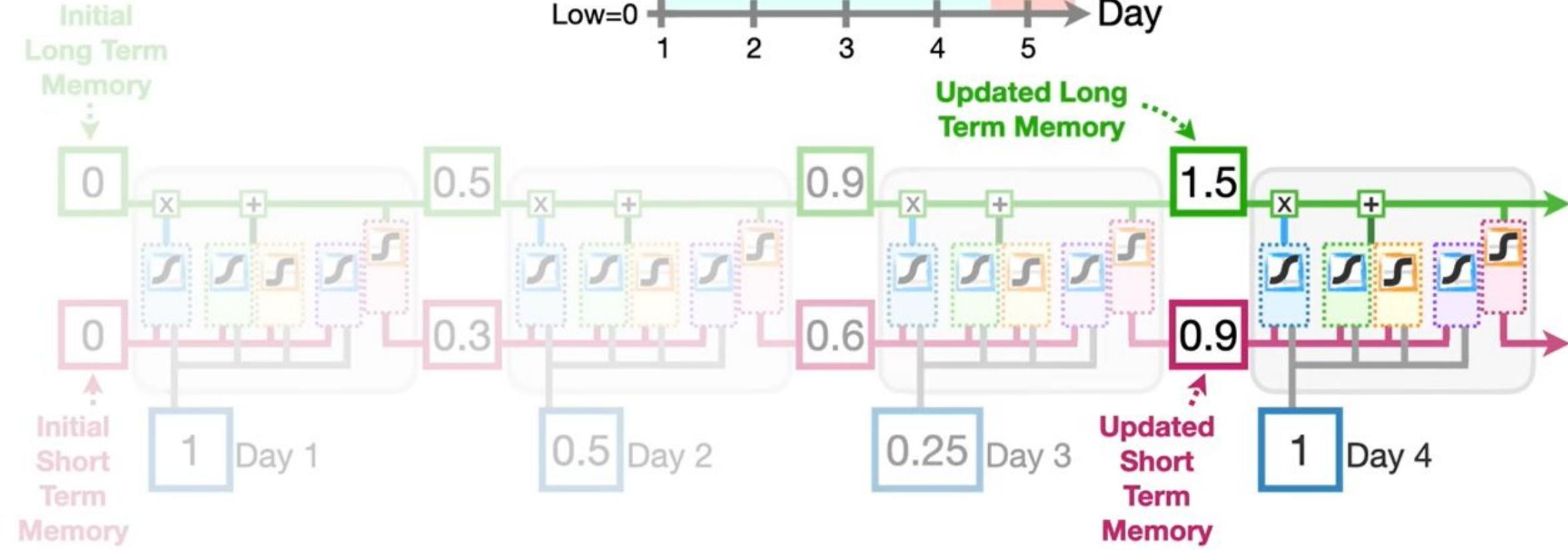
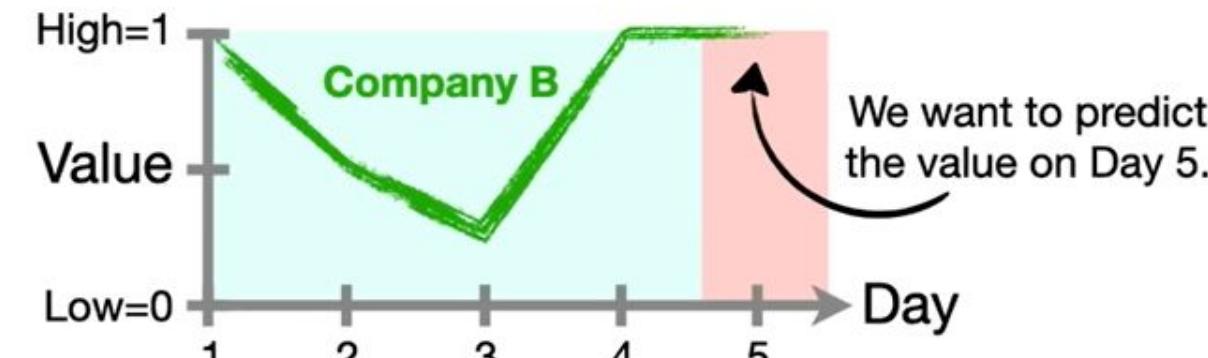




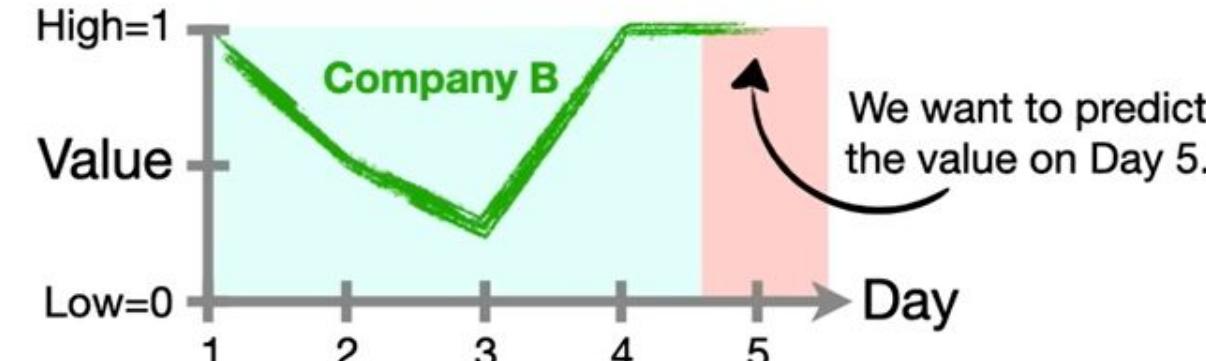
...so we plug in **0.5** for the updated **Long-Term Memory**...



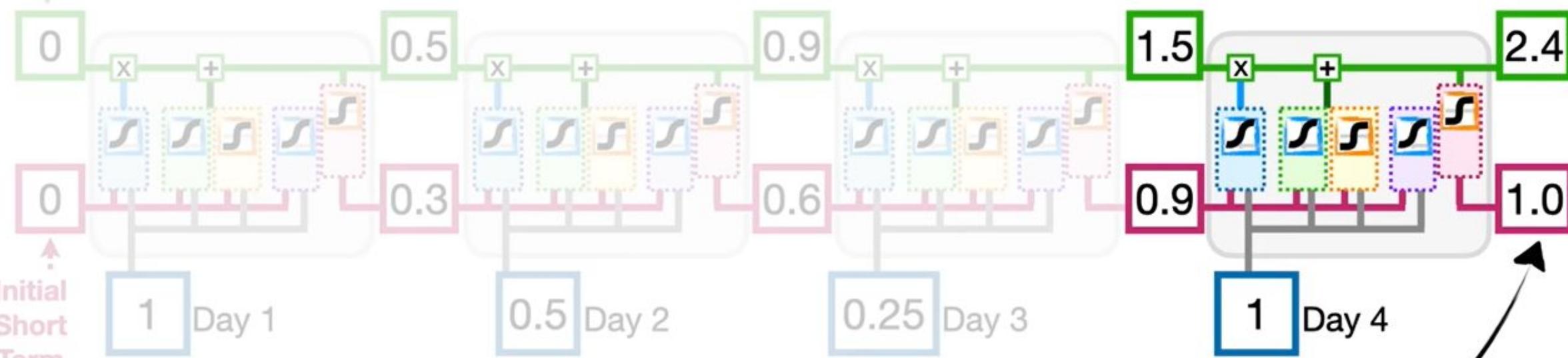
...and **0.3** (rounded) for the updated **Short-Term Memory**.



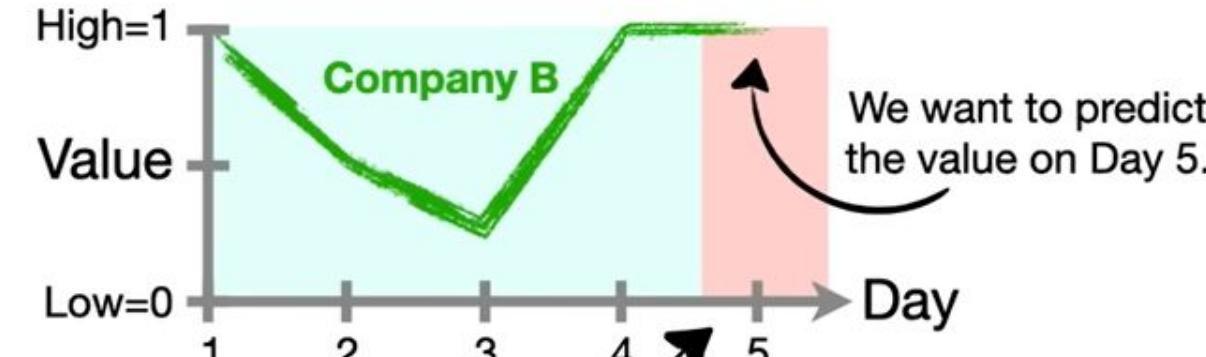
Now we unroll the **LSTM** and do the math with the remaining input values...



Initial
Long Term
Memory



...and the final **Short-Term Memory**, 1.0, is the **Output** from the **unrolled LSTM**.



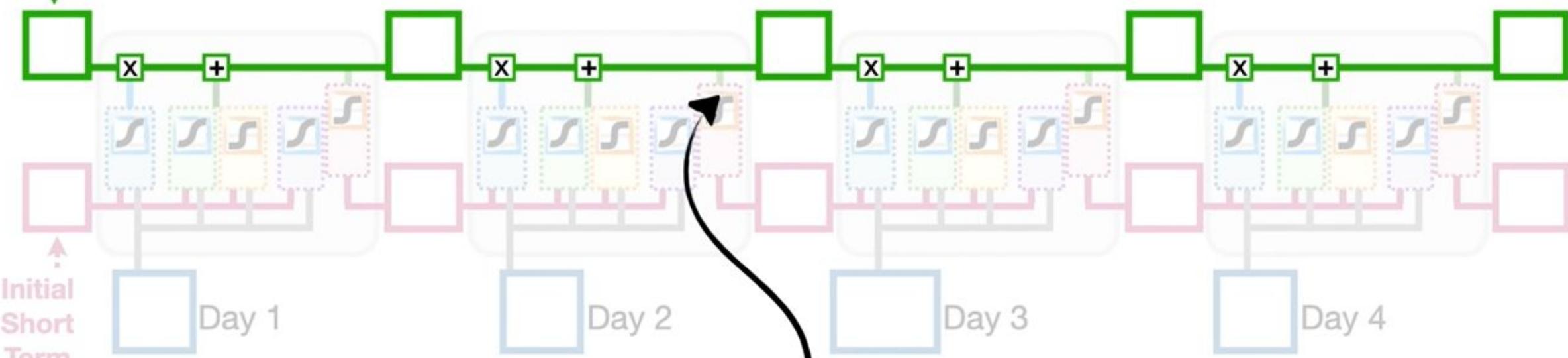
Initial
Long Term
Memory



And that means that the **Output**
from the **LSTM** correctly predicts
Company B's value for Day 5.



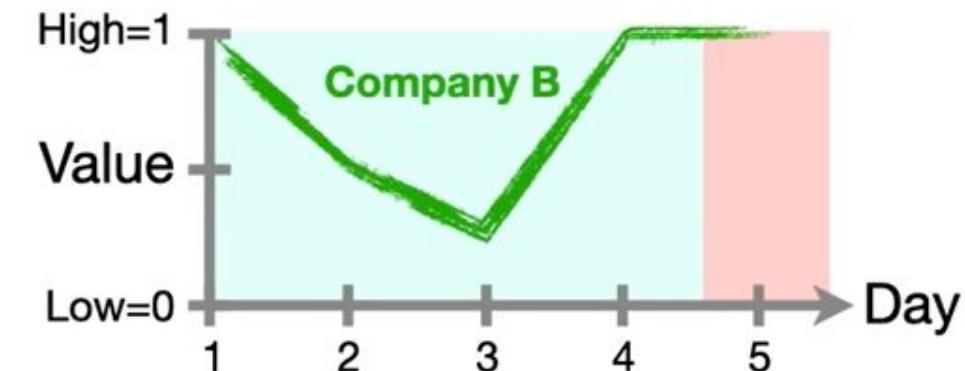
Initial
Long Term
Memory



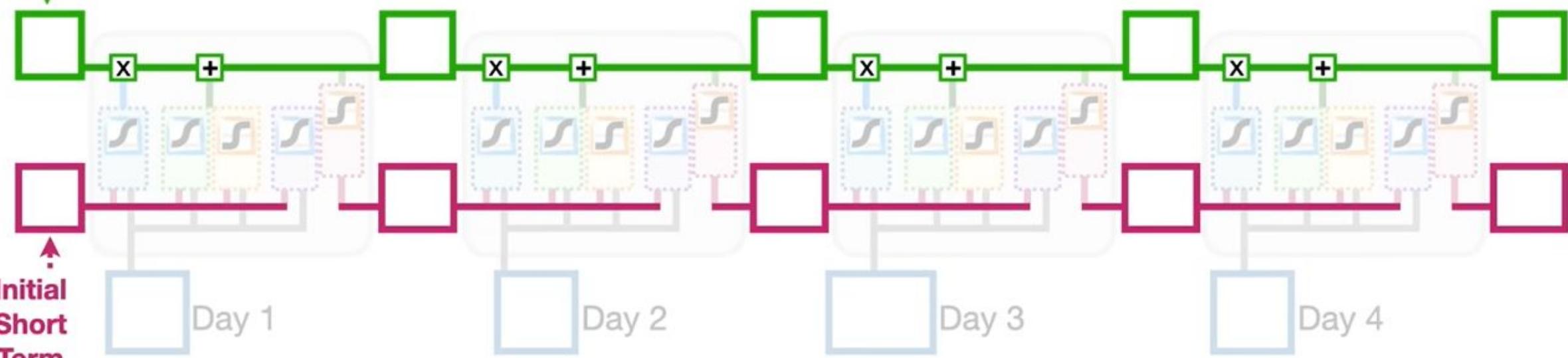
In summary, using
separate paths for **Long-
Term Memories...**



Initial
Long Term
Memory



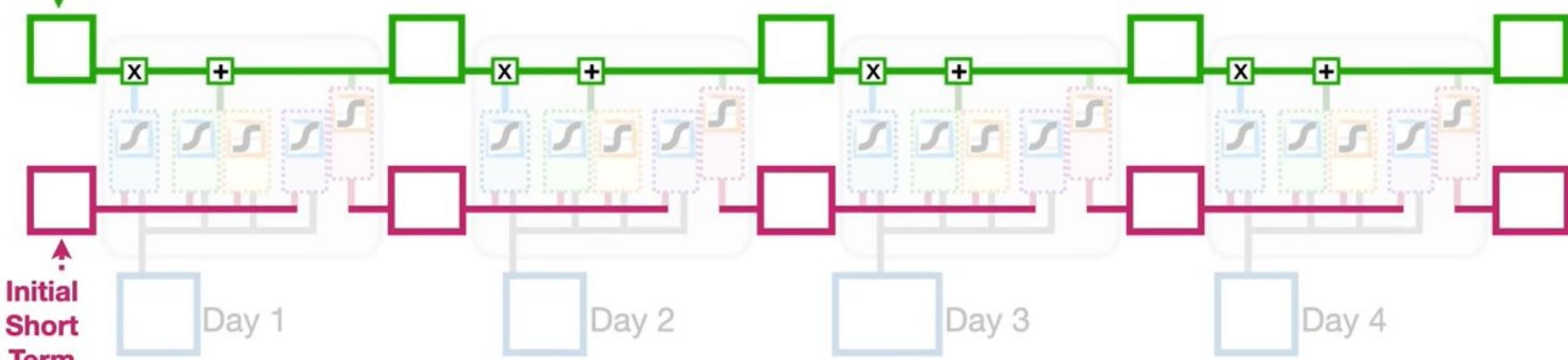
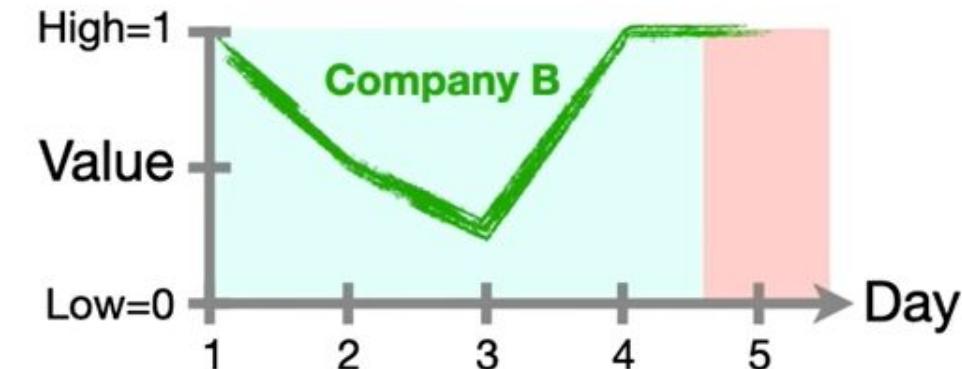
Initial
Short Term
Memory



...Long Short-Term Memory
Networks avoid the exploding/
vanishing gradient problem...



Initial
Long Term
Memory



...and that means we can unroll them more times
to accommodate longer sequences of input data
than a vanilla **Recurrent Neural Network**.

<https://towardsdatascience.com/tutorial-on-lstm-a-computational-perspective-f3417442c2cd>

<https://www.youtube.com/watch?v=YCzL96nL7j0&t=659s>