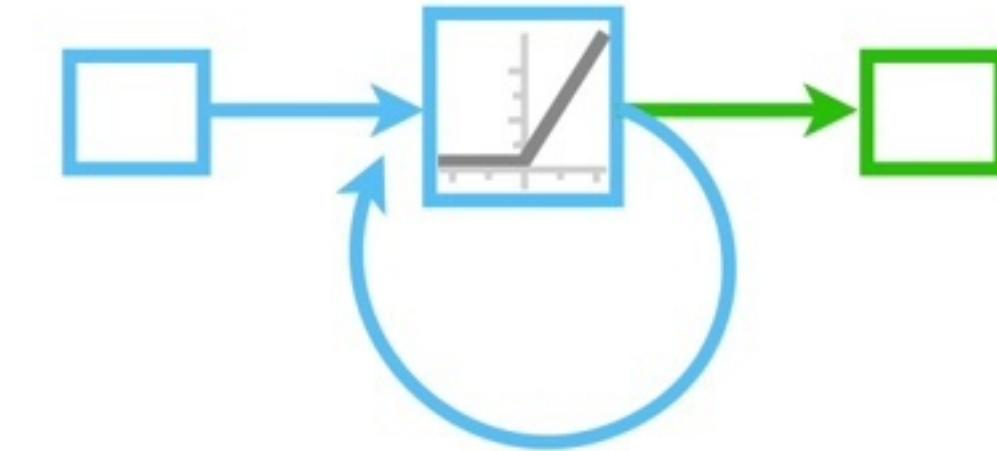




...to ***unroll*** a network that
works well with *different*
amounts of sequential data.



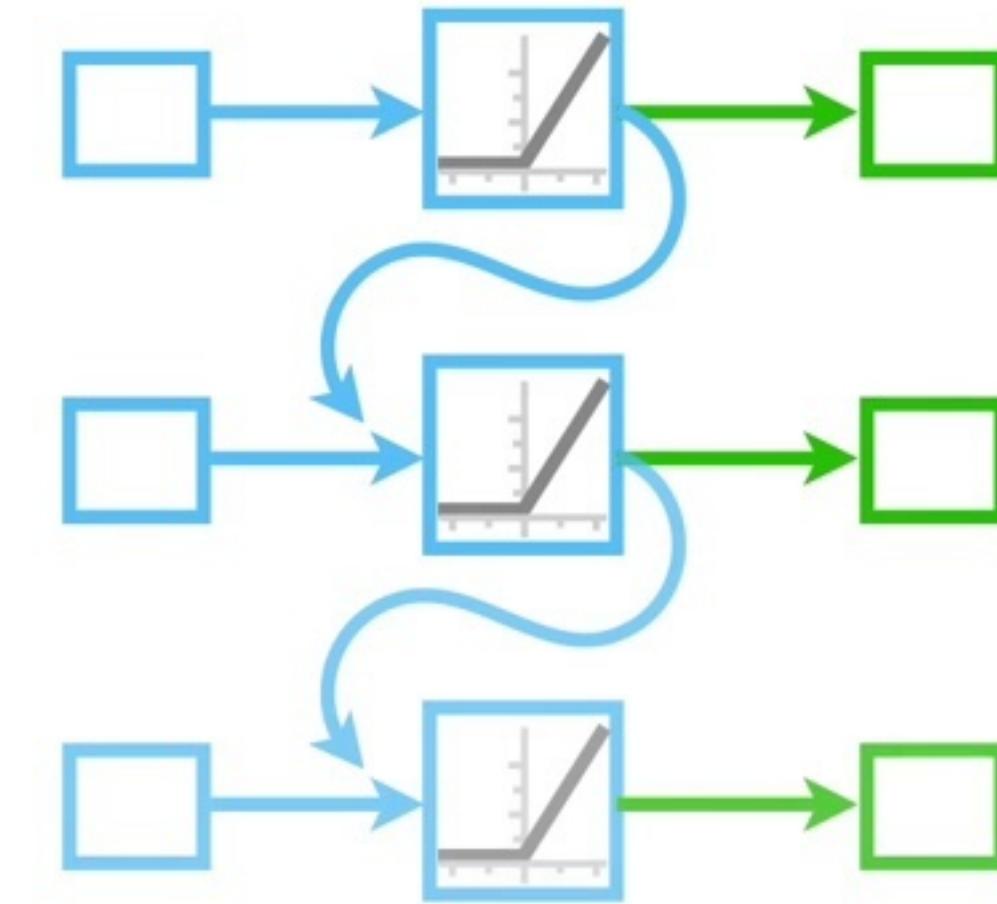


...to ***unroll*** a network that
works well with *different*
amounts of sequential data.



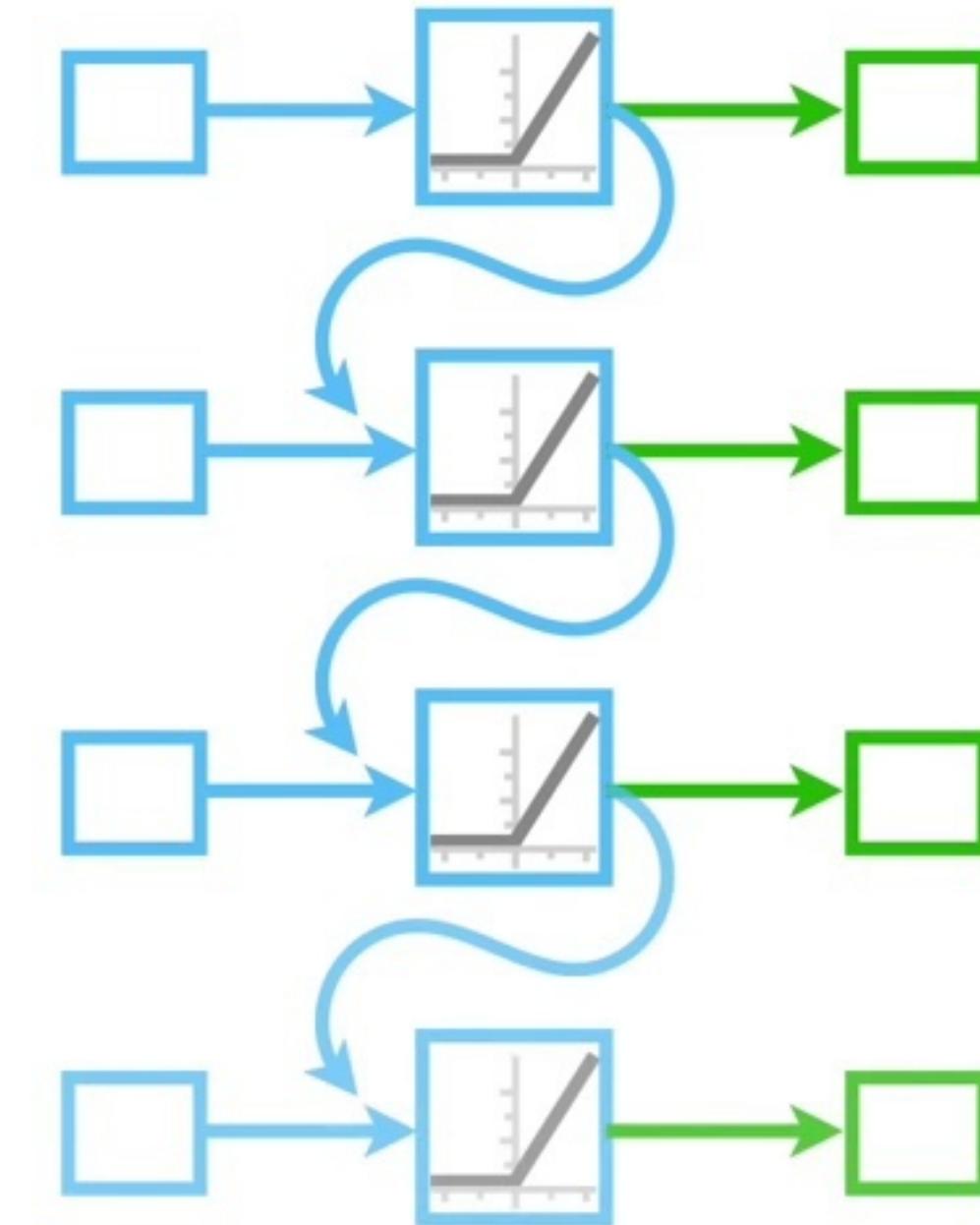


...to ***unroll*** a network that
works well with *different*
amounts of sequential data.



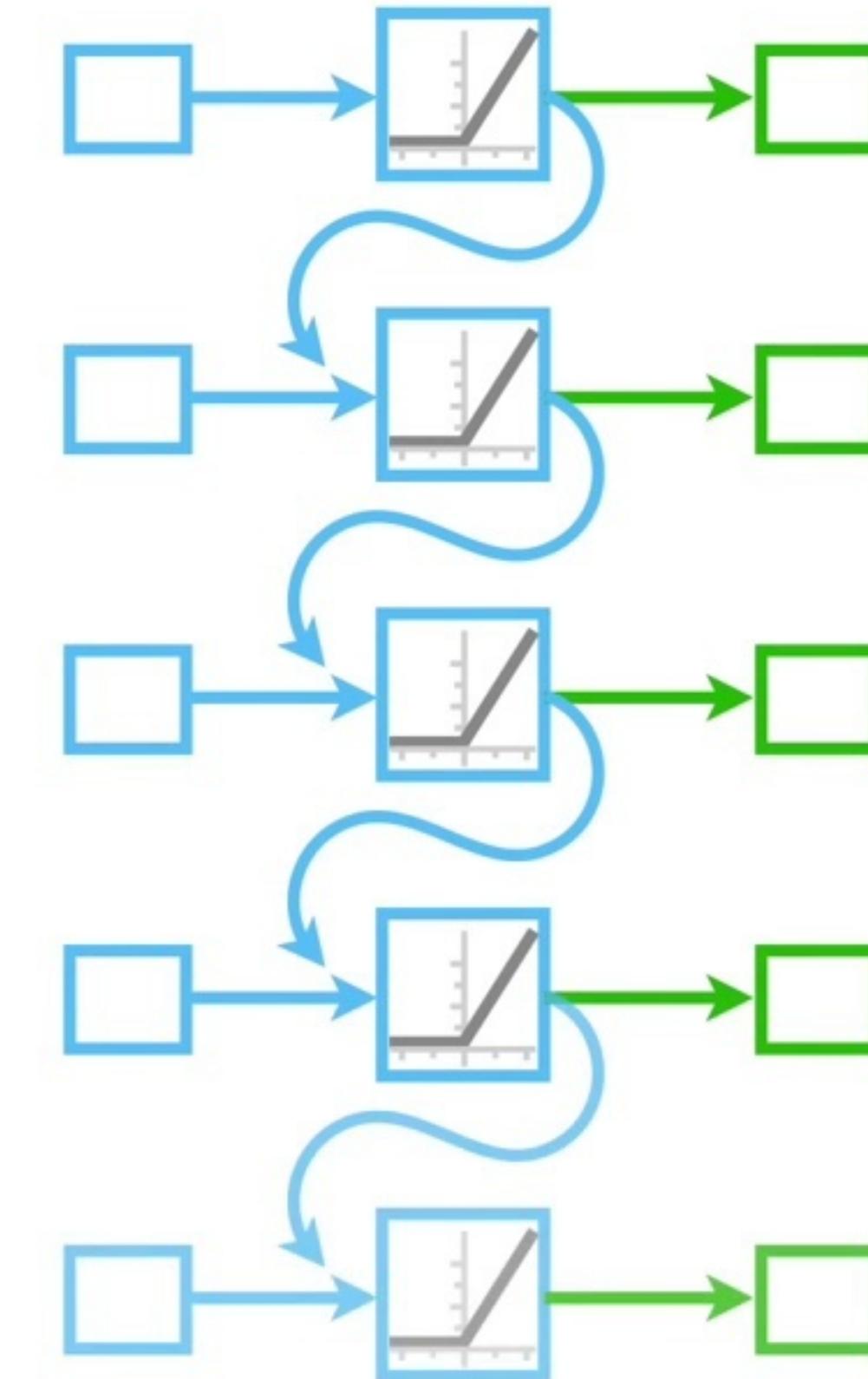


...to ***unroll*** a network that
works well with *different*
amounts of sequential data.



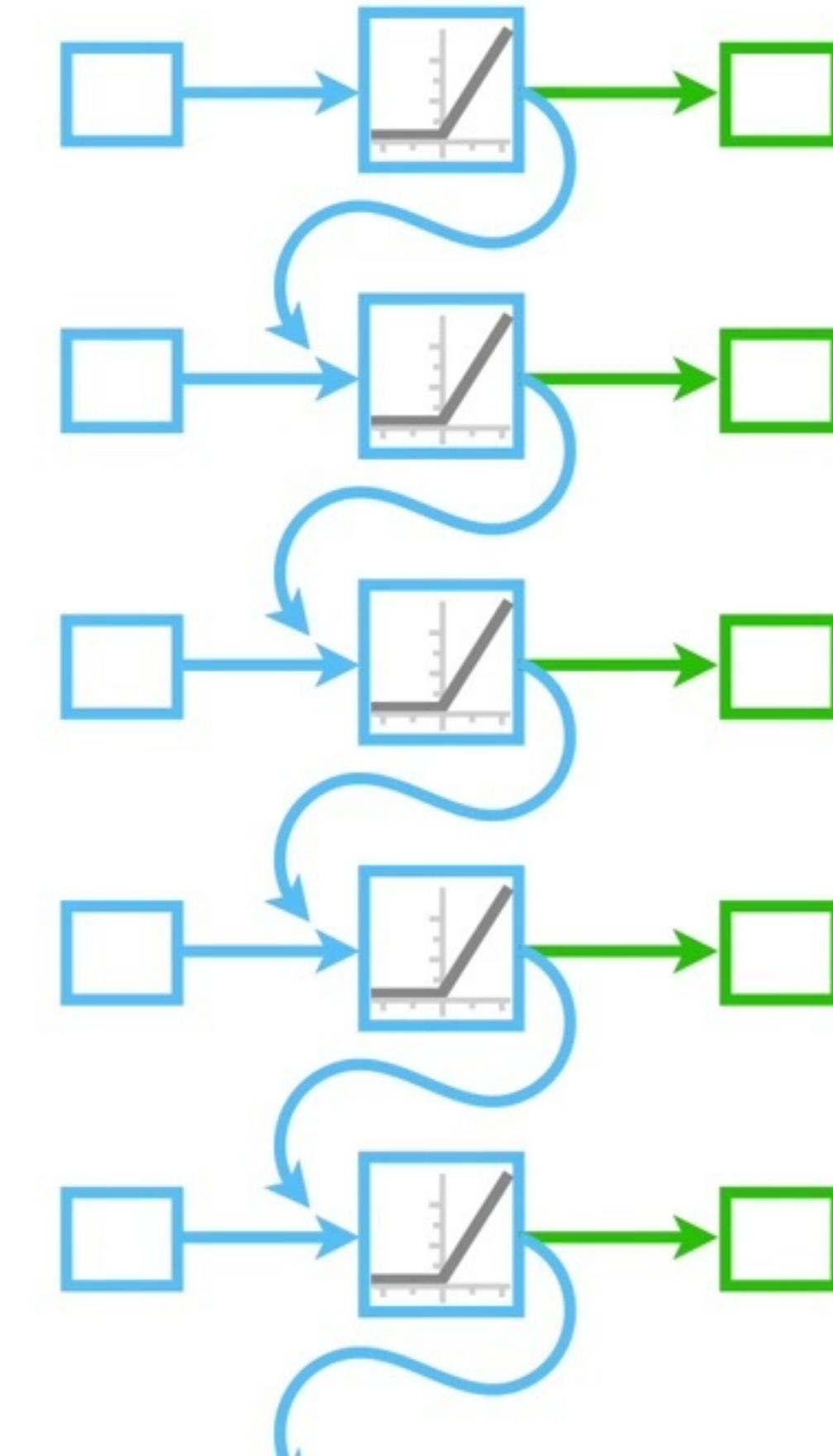


...to ***unroll*** a network that
works well with *different*
amounts of sequential data.



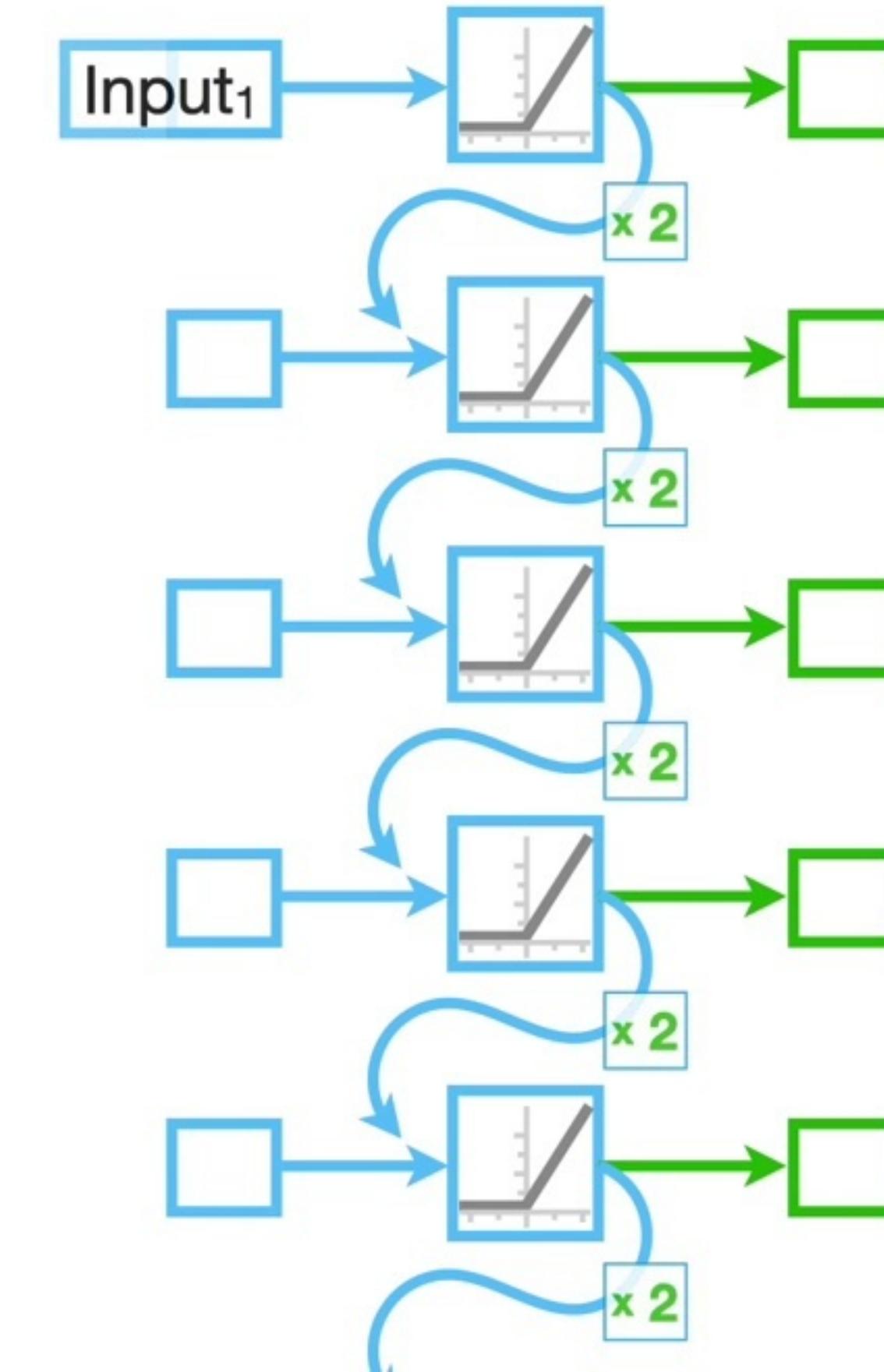


However, we also saw that
when we plug in the
numbers...





However, we also saw that when we plug in the numbers...





...when the **Weight** on the feedback loop is greater than 1, and in this example the **Weight is 2...**





...and **2** to the **50th** power is a
HUGE NUMBER...

$$= \text{Input}_1 \times 2^{50}$$

$$= \text{Input}_1 \times \text{A HUGE NUMBER}$$





...and this **HUGE NUMBER**
would cause the gradient,
which we need for **Gradient
Descent**...

$$= \text{Input}_1 \times 2^{50}$$

$$= \text{Input}_1 \times \text{A HUGE NUMBER}$$





...to **EXPLODE!!!**

$$= \text{Input}_1 \times 2^{50}$$

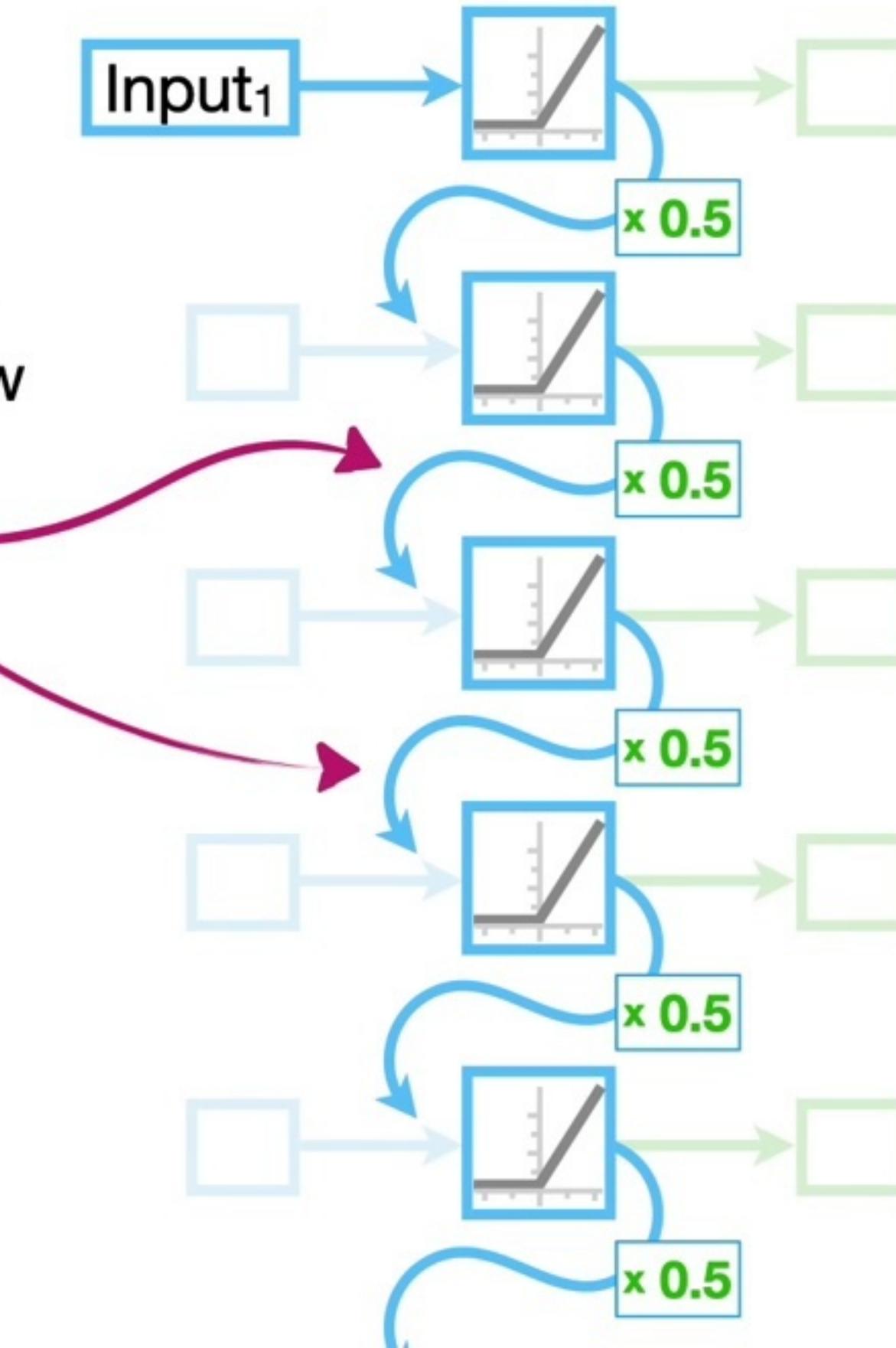
$$= \text{Input}_1 \times \text{A HUGE NUMBER}$$





Alternatively, we saw that if the **Weight** on the feedback loop was less than 1, and now we have it set to **0.5**...

$$= \text{Input}_1 \times 0.5^{50}$$





...then we'll end up
multiplying the **Input** value
by **0.5** raised to the **50th**
power...

$$= \text{Input}_1 \times 0.5^{50}$$





...and **0.5** raised to the **50th** power is a number super close to **0**...

$$= \text{Input}_1 \times 0.5^{50}$$

$$= \text{Input}_1 \times \text{a number super close to } 0$$

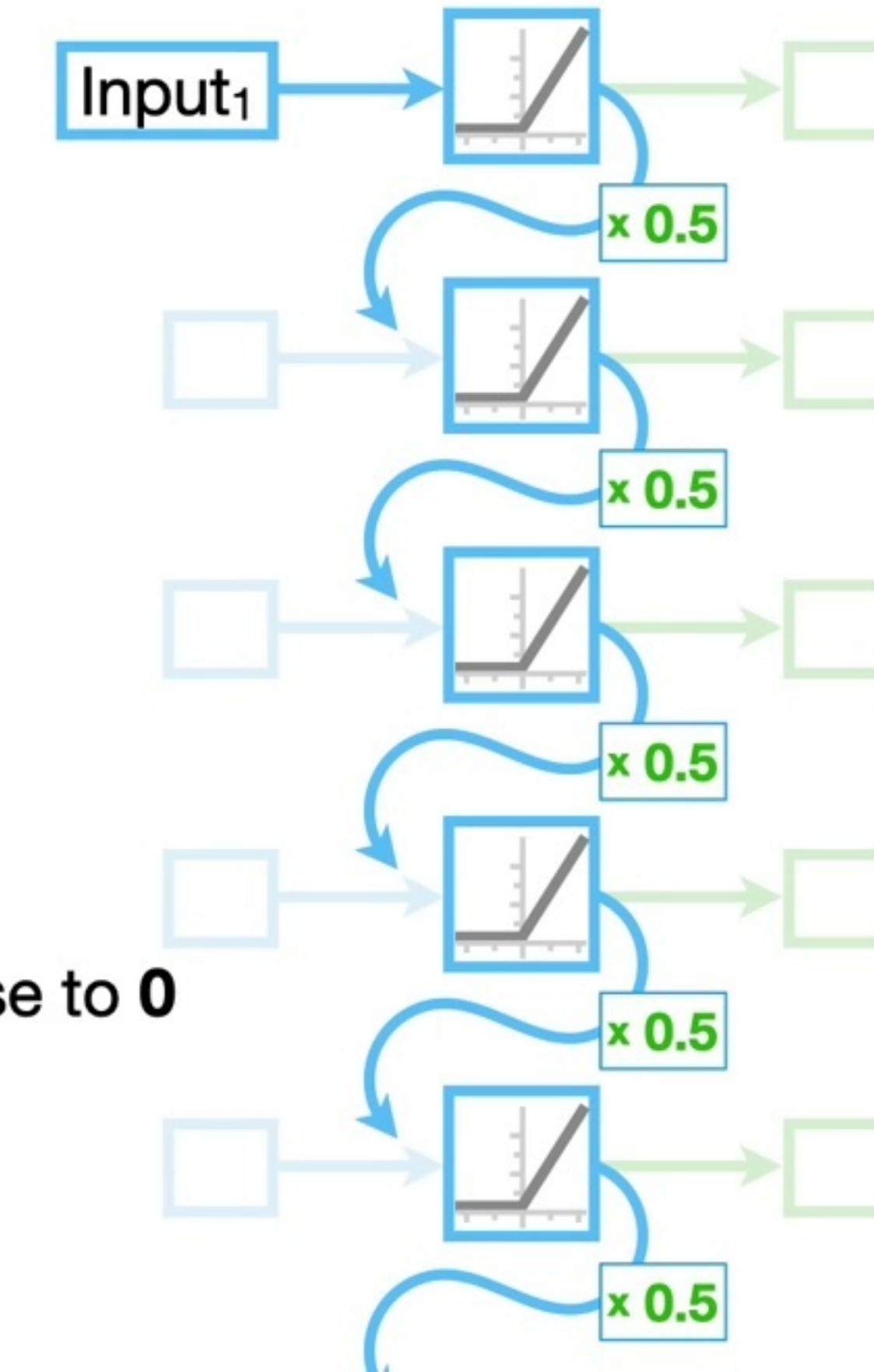




...and this number super close to **0** would cause the gradient, which we need for **Gradient Descent**...

$$= \text{Input}_1 \times 0.5^{50}$$

$$= \text{Input}_1 \times \text{a number super close to } 0$$

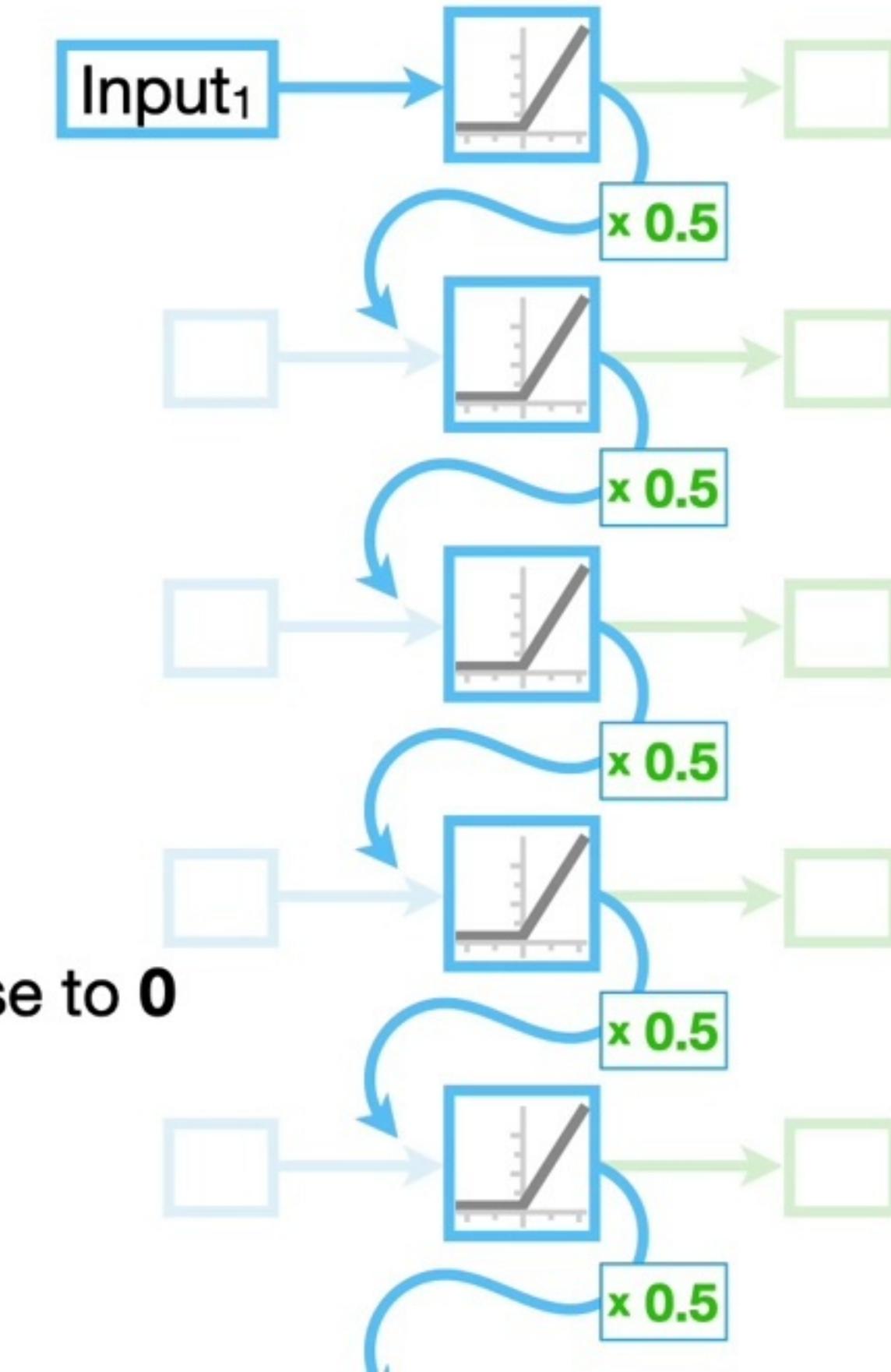




...to vanish.

$$= \text{Input}_1 \times 0.5^{50}$$

$$= \text{Input}_1 \times \text{a number super close to } 0$$





In summary, basic, vanilla
Recurrent Neural Networks
are hard to train because the
gradients can **explode**, or
vanish.

$$= \text{Input}_1 \times \mathbf{W}_2^{\text{Num. Unroll}}$$





The good news is that it doesn't take much to extend the basic, vanilla **Recurrent Neural Network** so that we can avoid this problem.



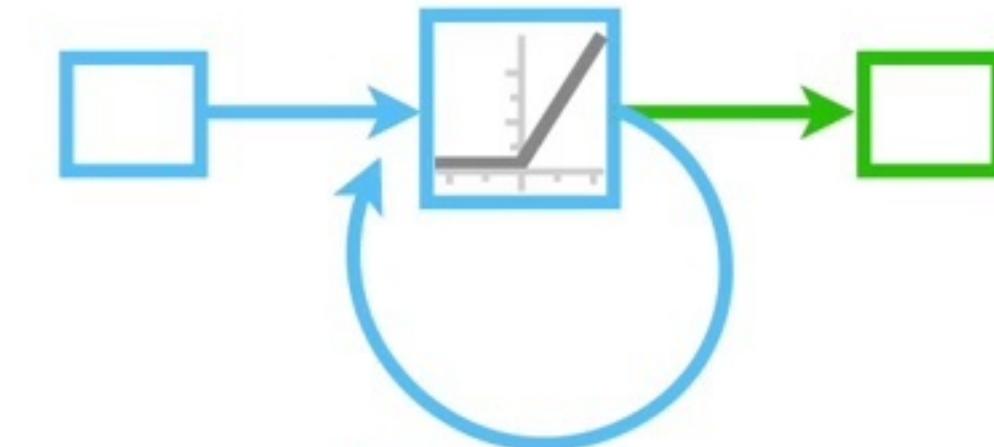


So, today, we're going to talk about **Long Short-Term Memory (LSTM)**, which is a type of **Recurrent Neural Network** that is designed to avoid the **exploding/vanishing gradient** problem.



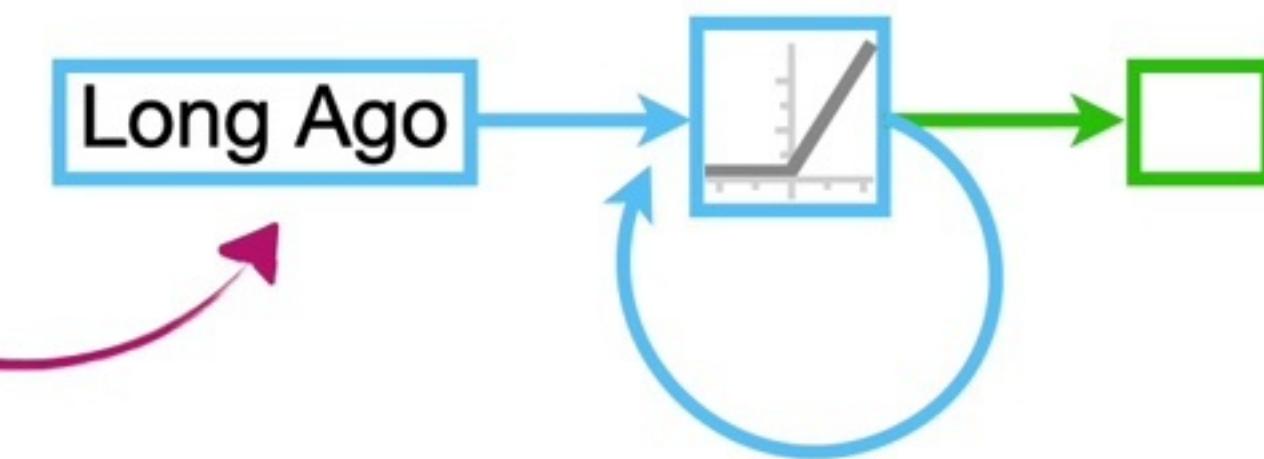


The main idea behind how
Long Short-Term Memory
(LSTM) works...

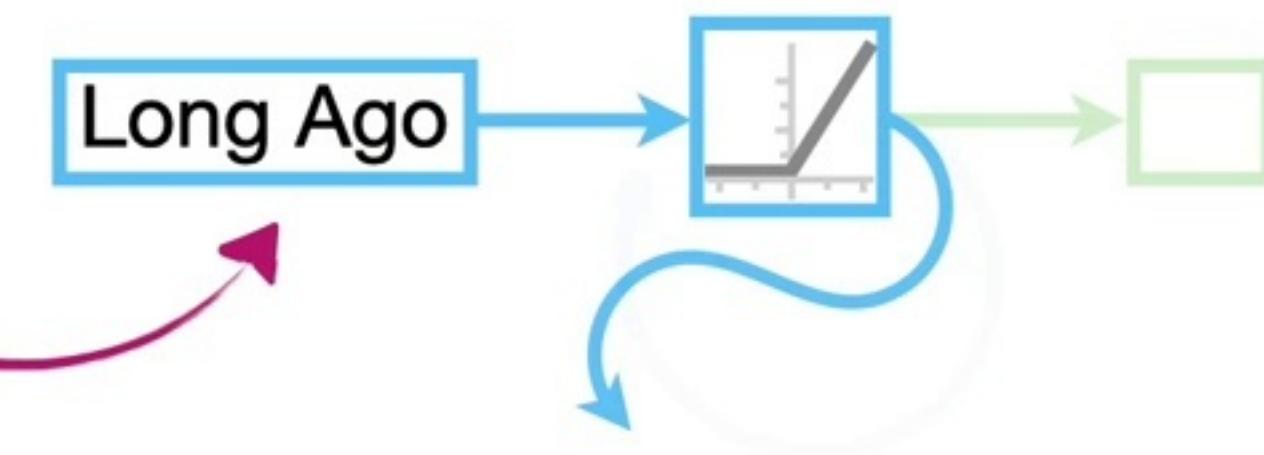


...is that instead of using the same feedback loop connection...





...for events that
happened long ago...



...for events that
happened long ago...

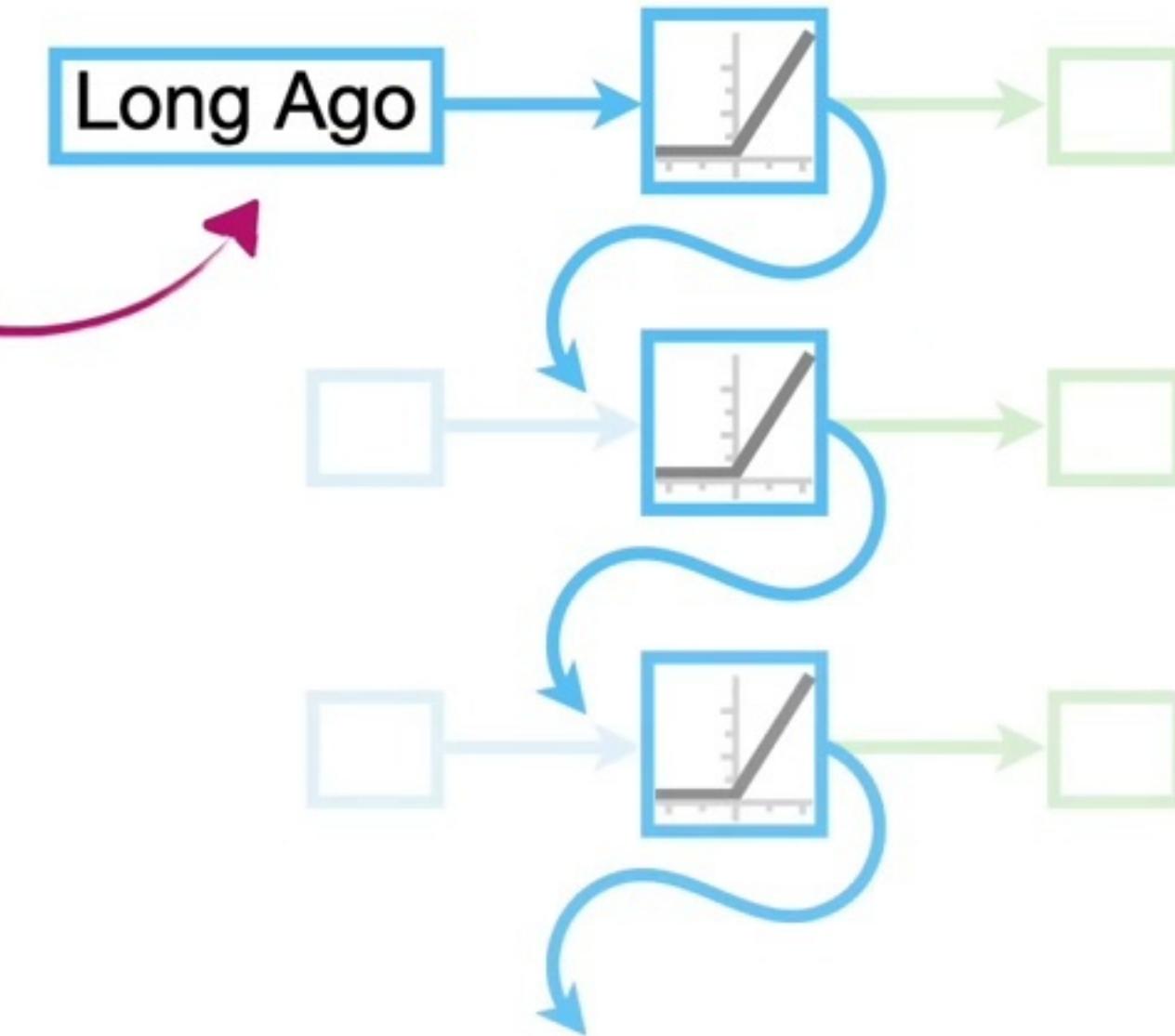


...for events that
happened long ago...





...for events that
happened long ago...

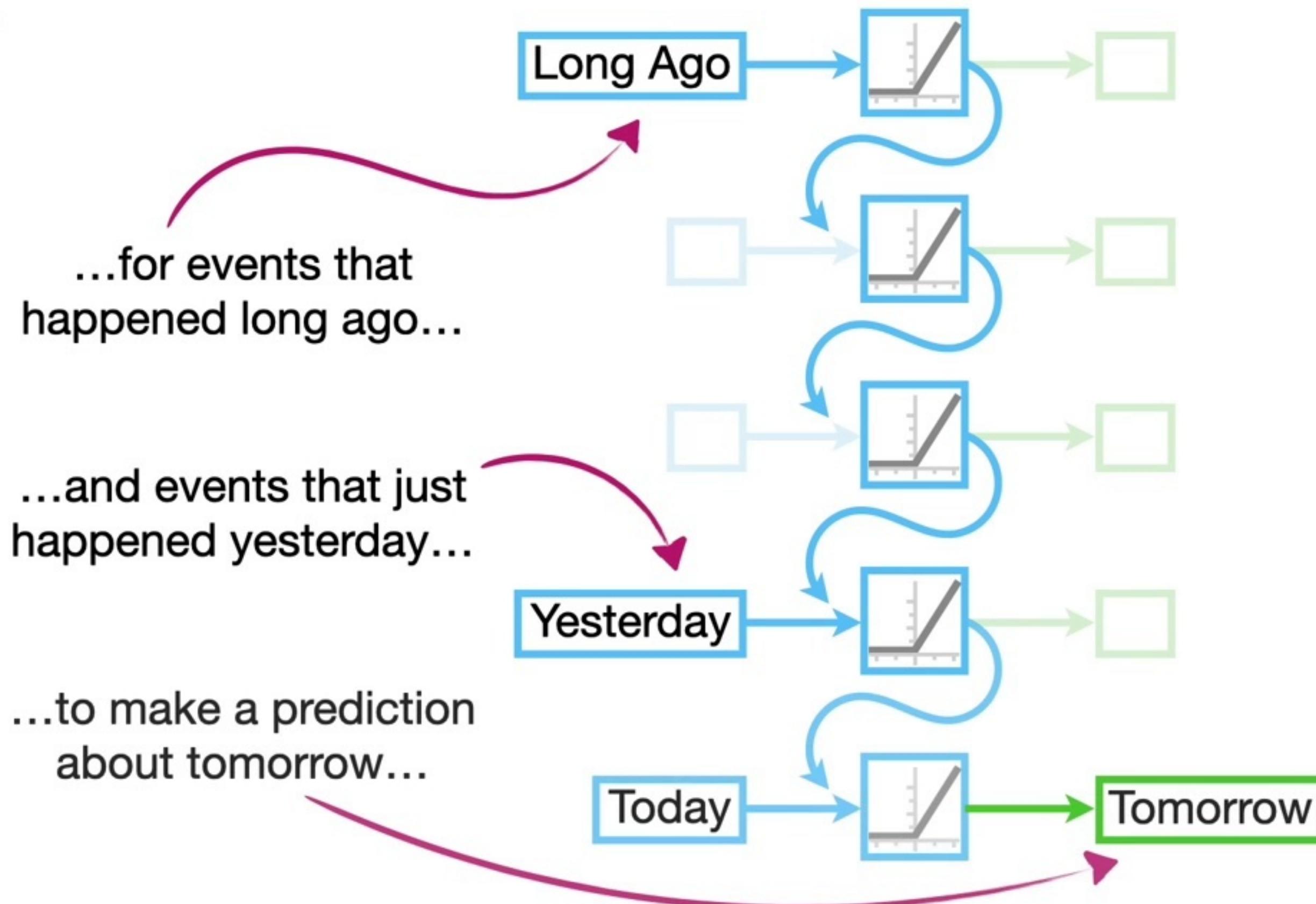




...for events that
happened long ago...

...and events that just
happened yesterday...







**...Long Short-Term
Memory uses two
separate paths to make
predictions about
tomorrow.**





One path is for **Long-Term Memories...**





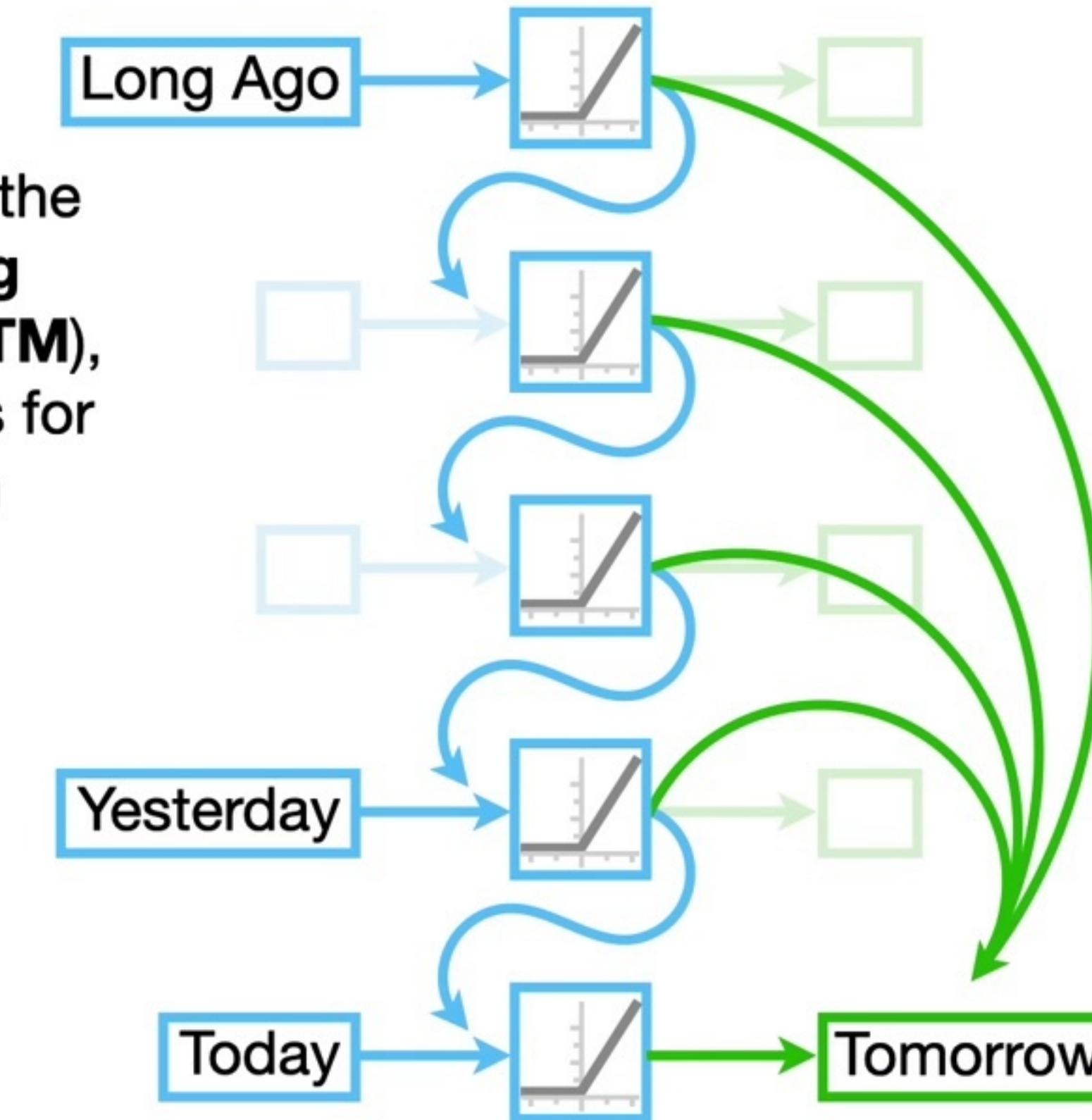
One path is for **Long-Term Memories**...

...and one is for **Short-Term Memories**.





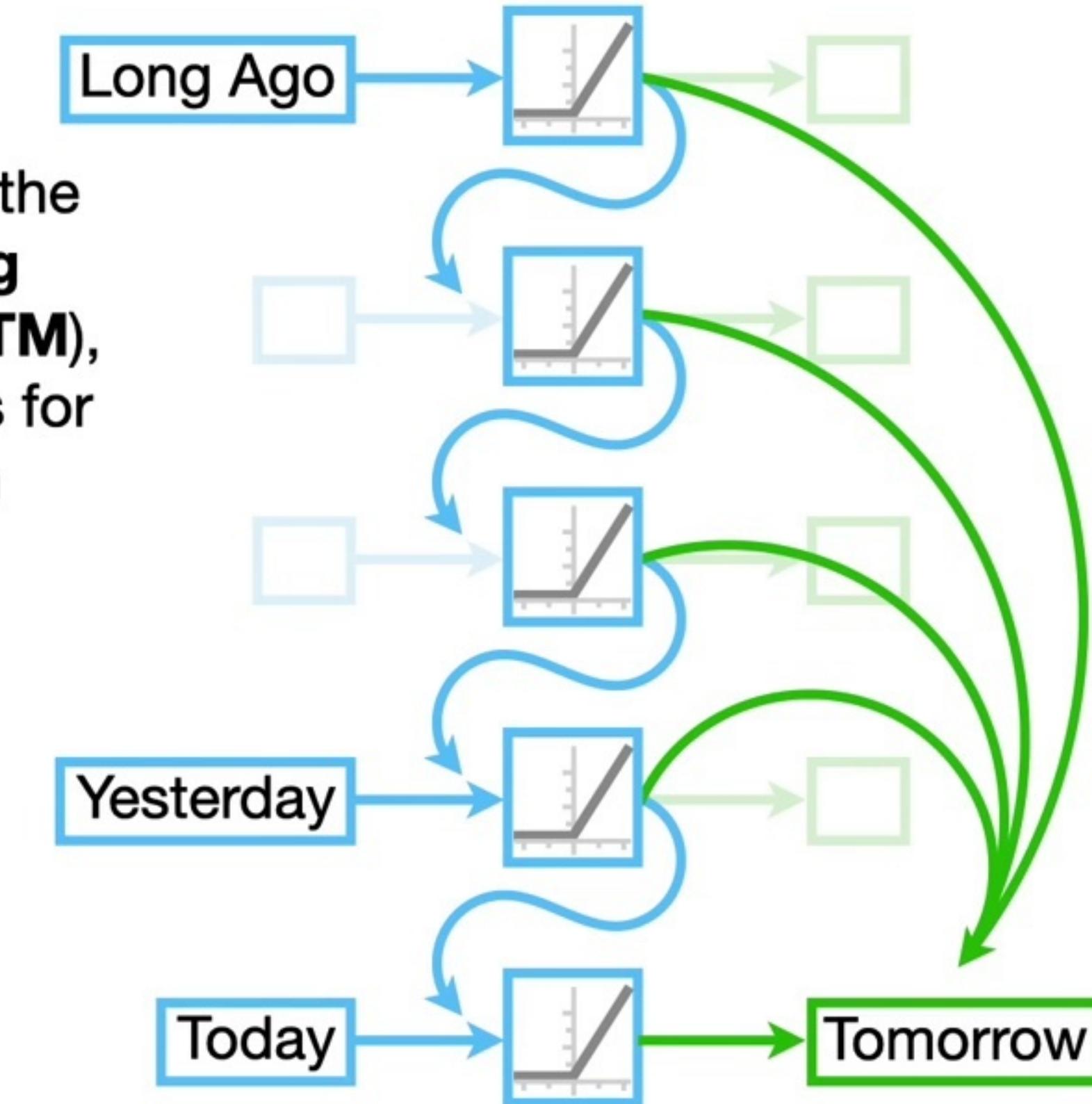
Now that we understand the **main idea** behind **Long Short-Term Memory (LSTM)**, that it uses different paths for **Long and Short-Term Memories...**





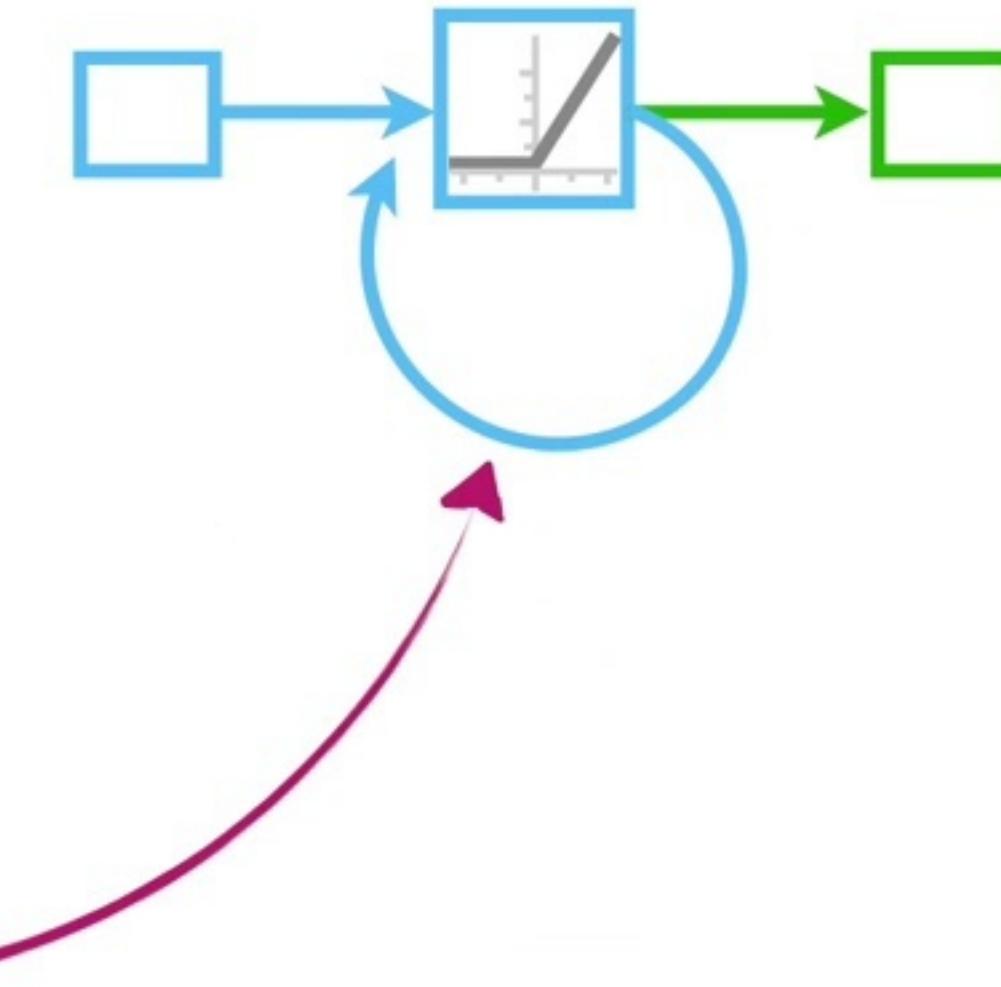
Now that we understand the **main idea** behind **Long Short-Term Memory (LSTM)**, that it uses different paths for **Long and Short-Term Memories...**

...let's talk about the **details**.



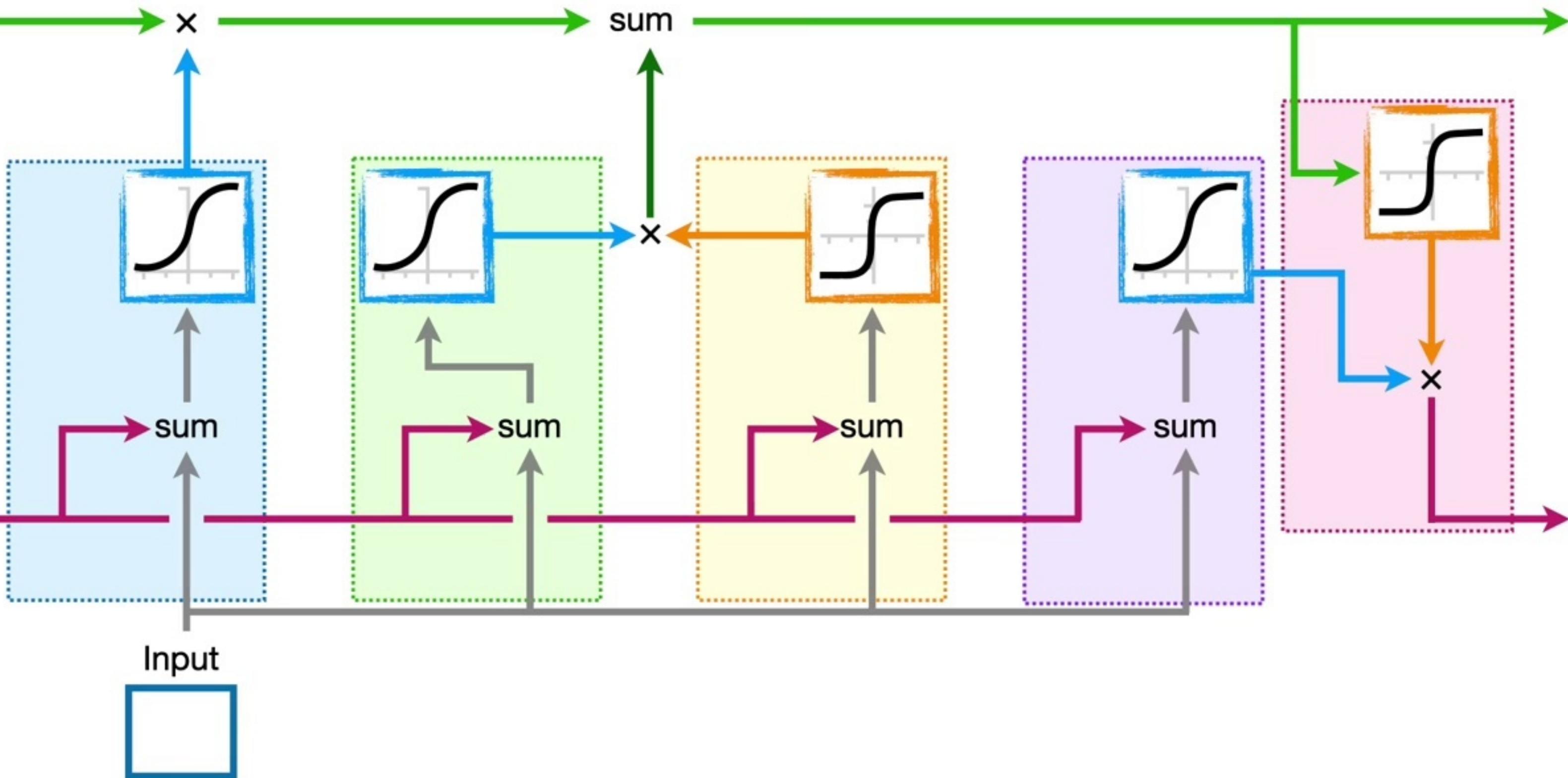


The bad news is that compared to a basic, vanilla **Recurrent Neural Network**, which unrolls from a relatively simple unit...



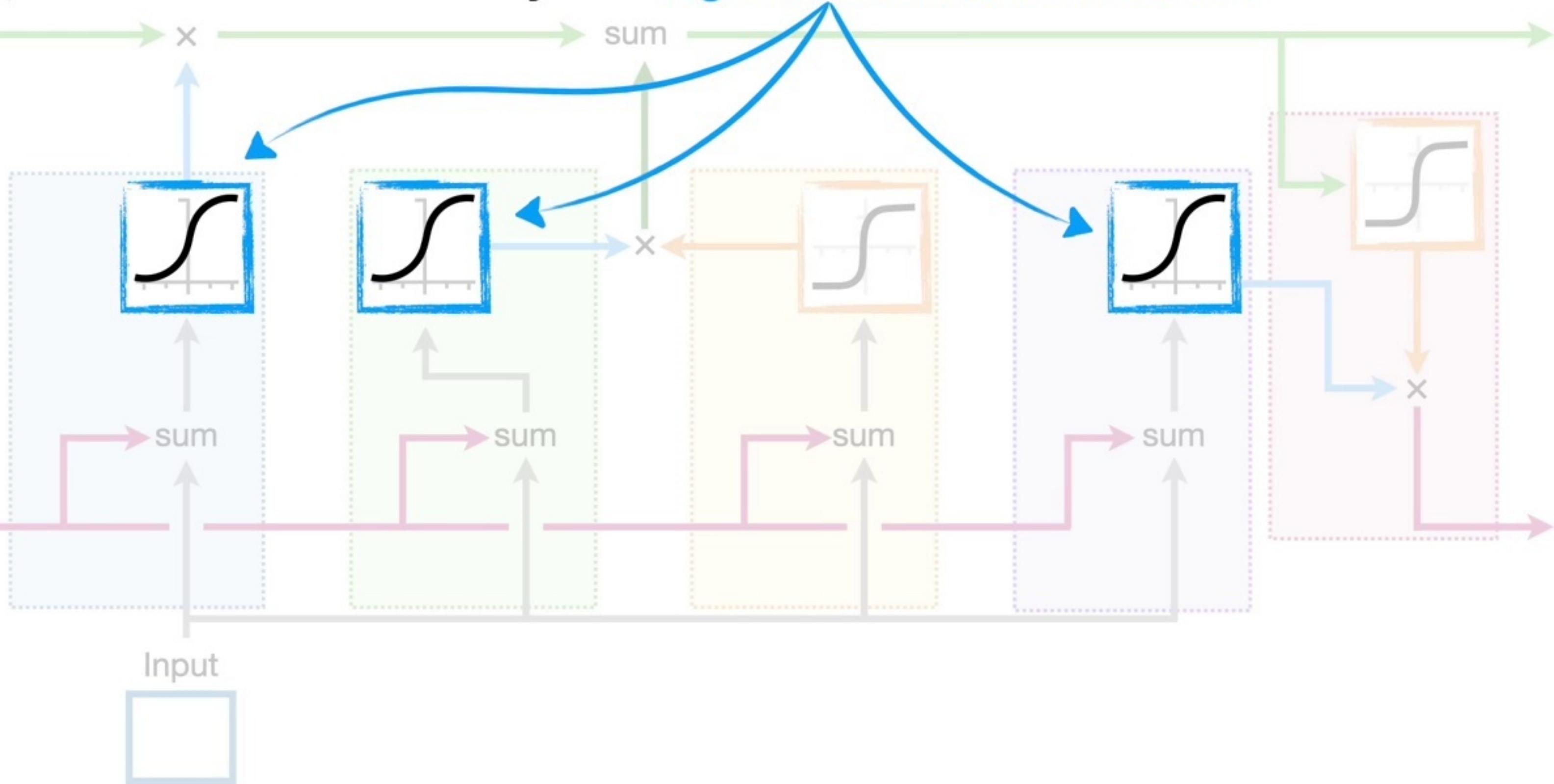


...Long Short-Term Memory is based
on a much more complicated unit.



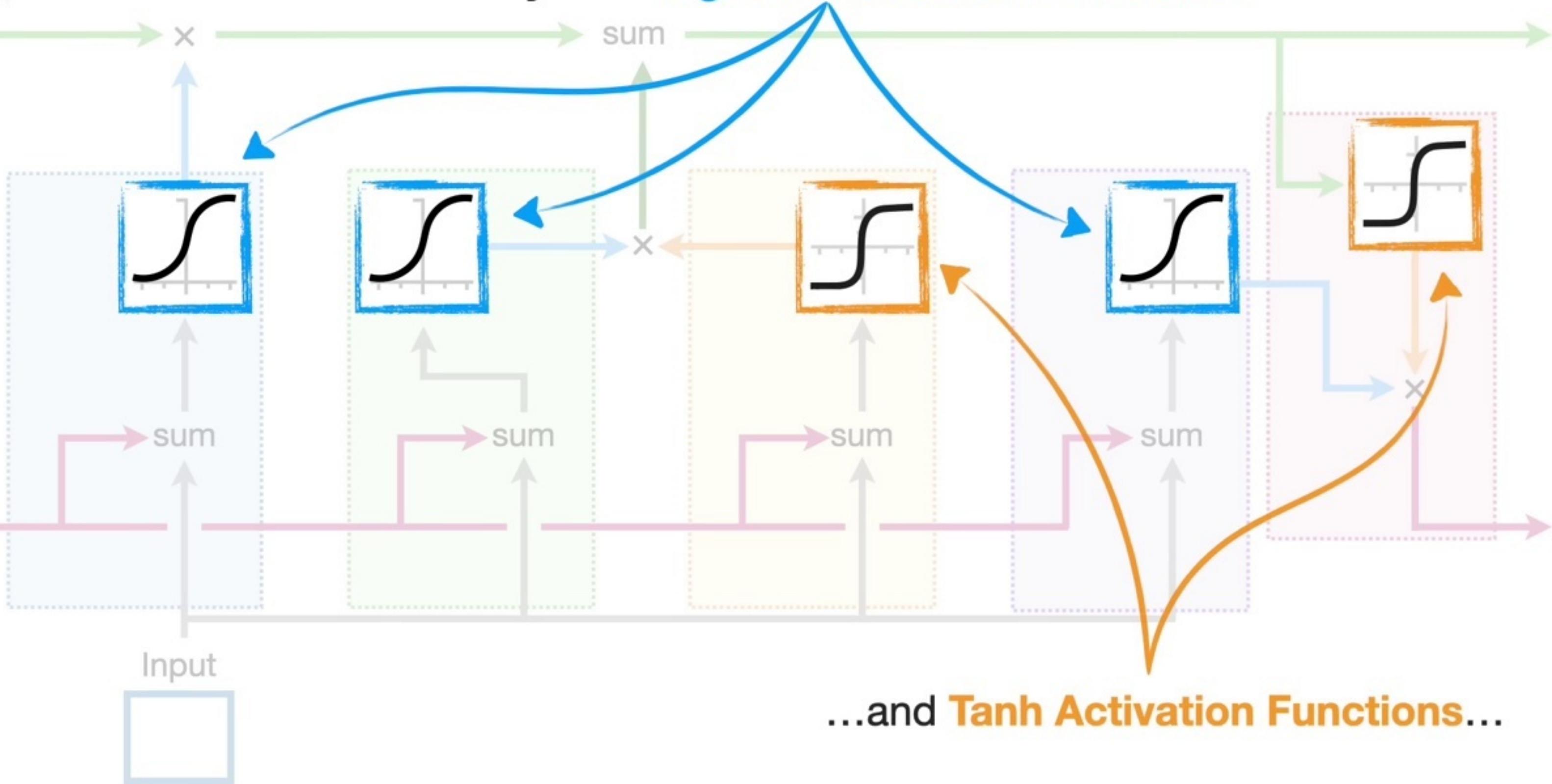


NOTE: Unlike the networks we've used before in this series, **Long Short-Term Memory** uses **Sigmoid Activation Functions**...





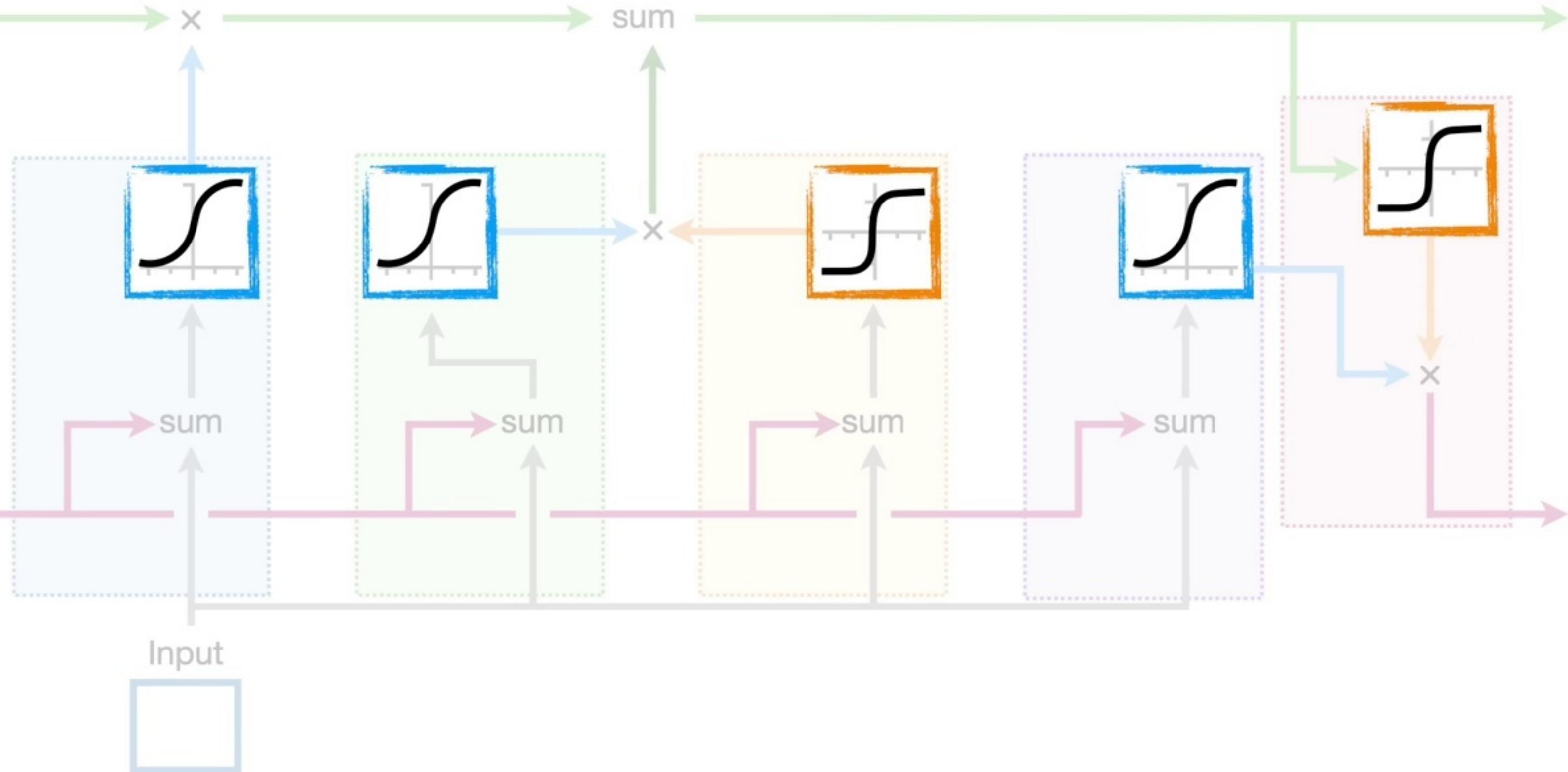
NOTE: Unlike the networks we've used before in this series, **Long Short-Term Memory** uses **Sigmoid Activation Functions...**



...and Tanh Activation Functions...

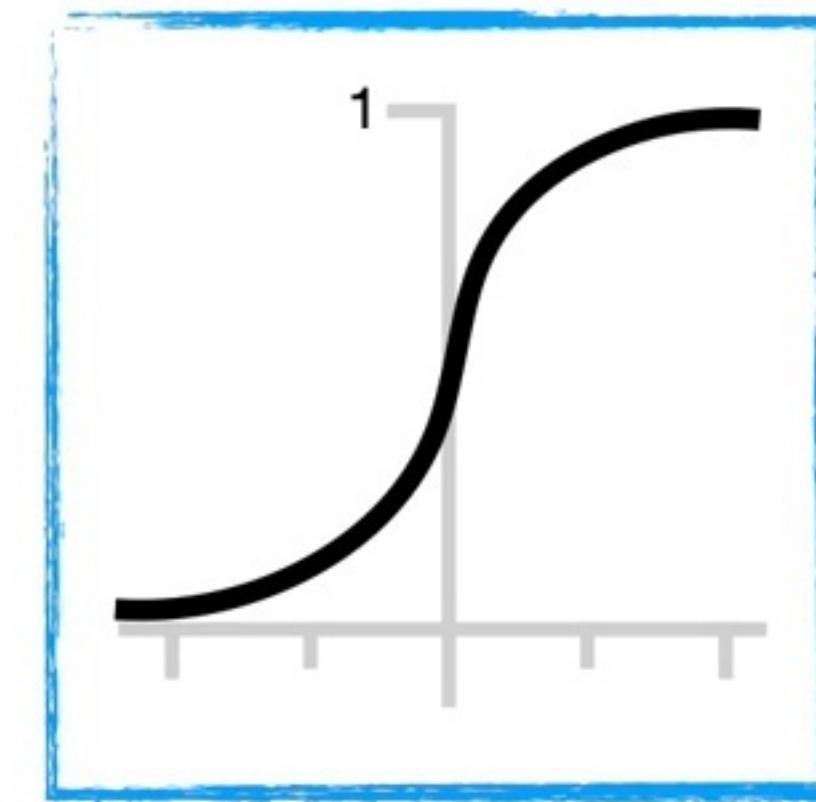


...so let's quickly talk about **Sigmoid** and **Tanh** Activation Functions.



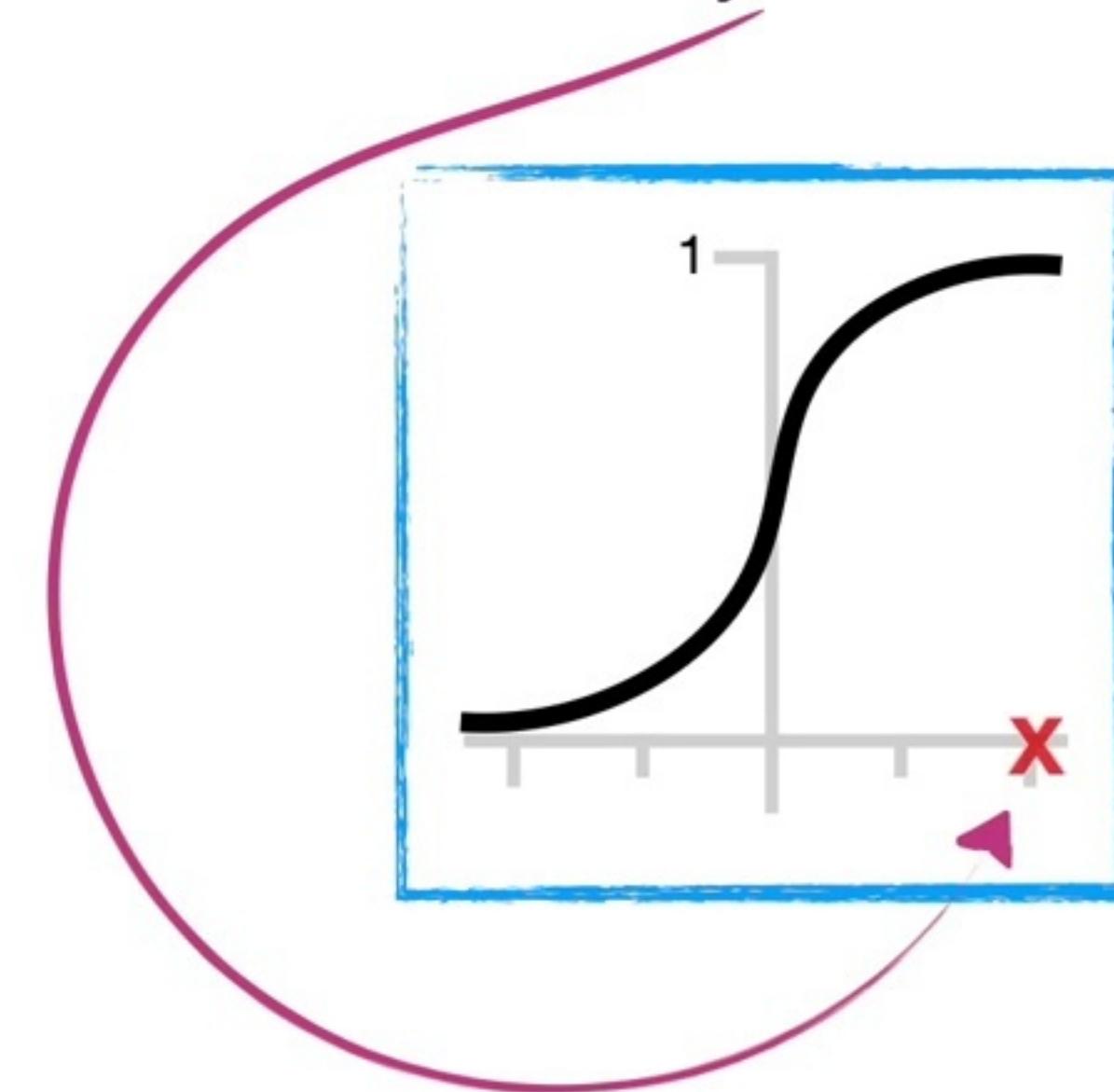


In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...



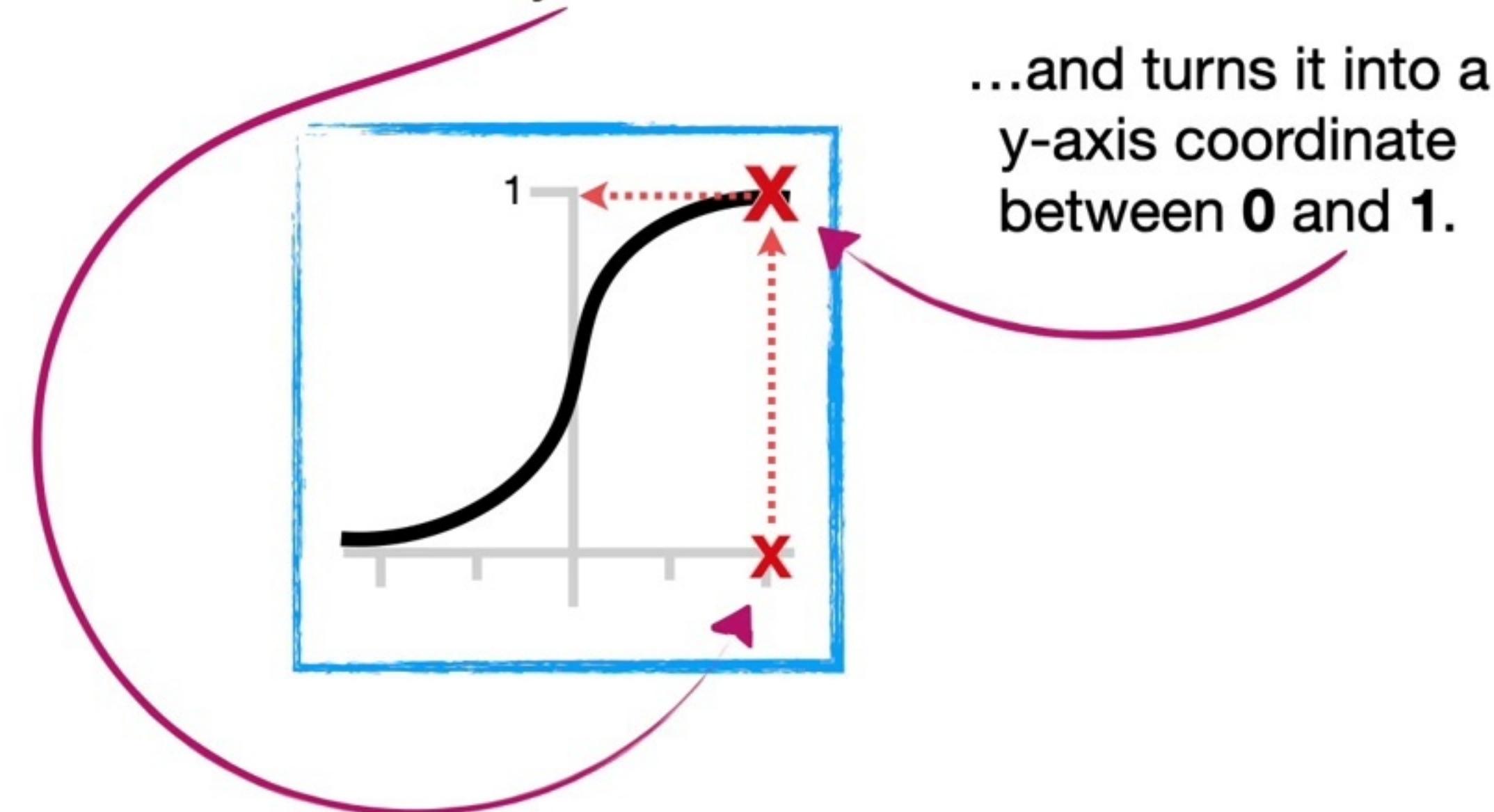


In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...





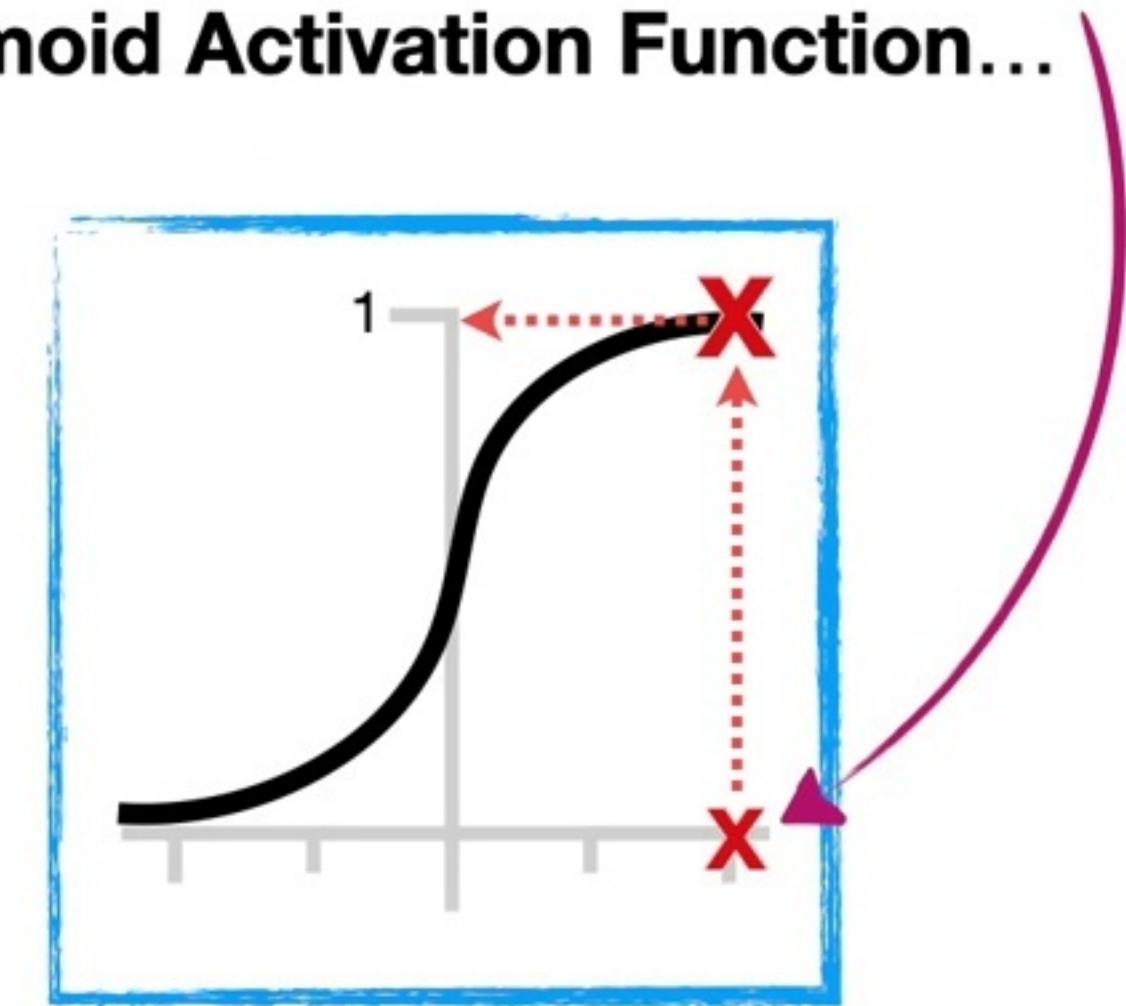
In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...



...and turns it into a y-axis coordinate between **0** and **1**.

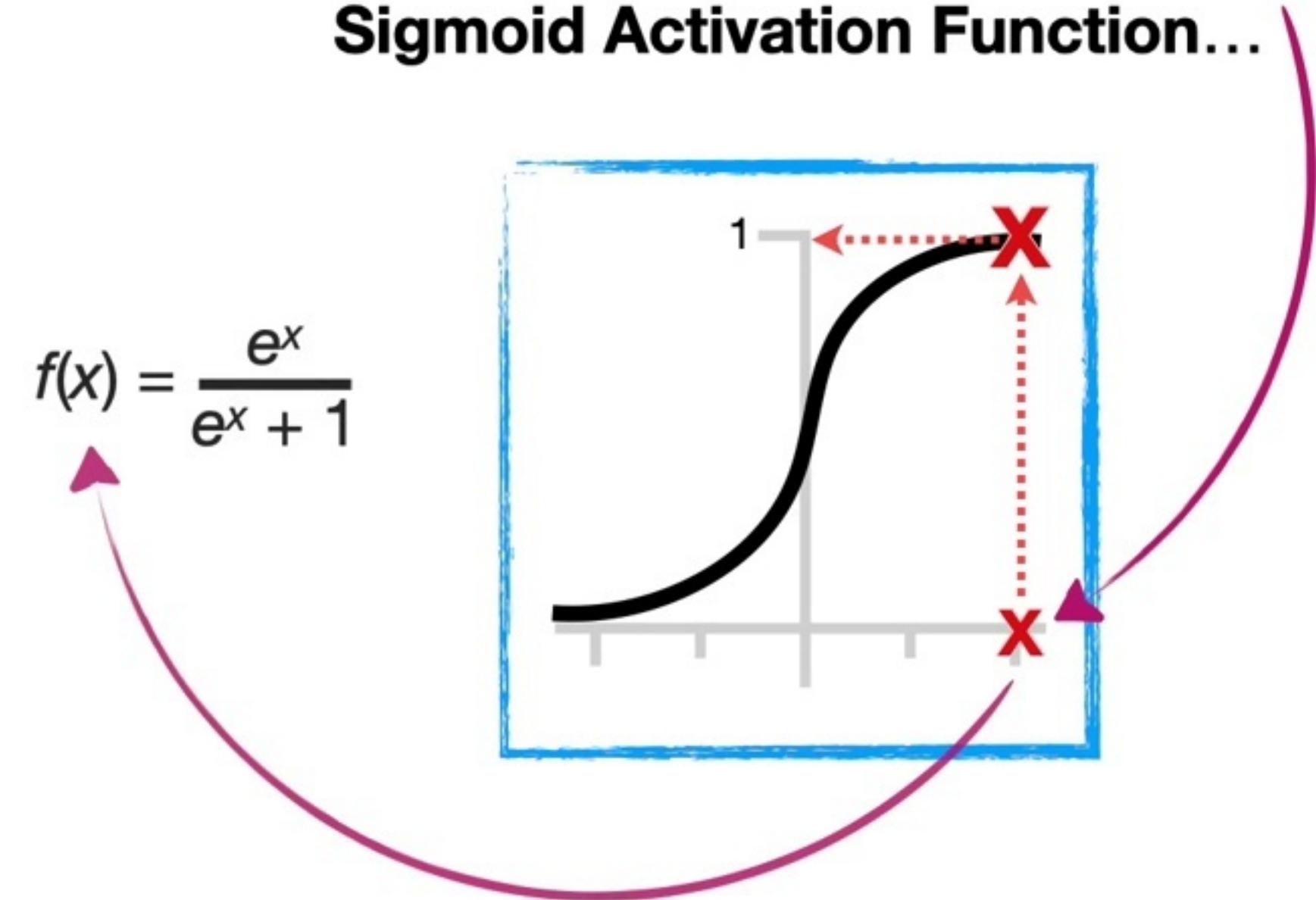


For example, when we plug in this x-axis coordinate, 10, into the equation for the **Sigmoid Activation Function**...





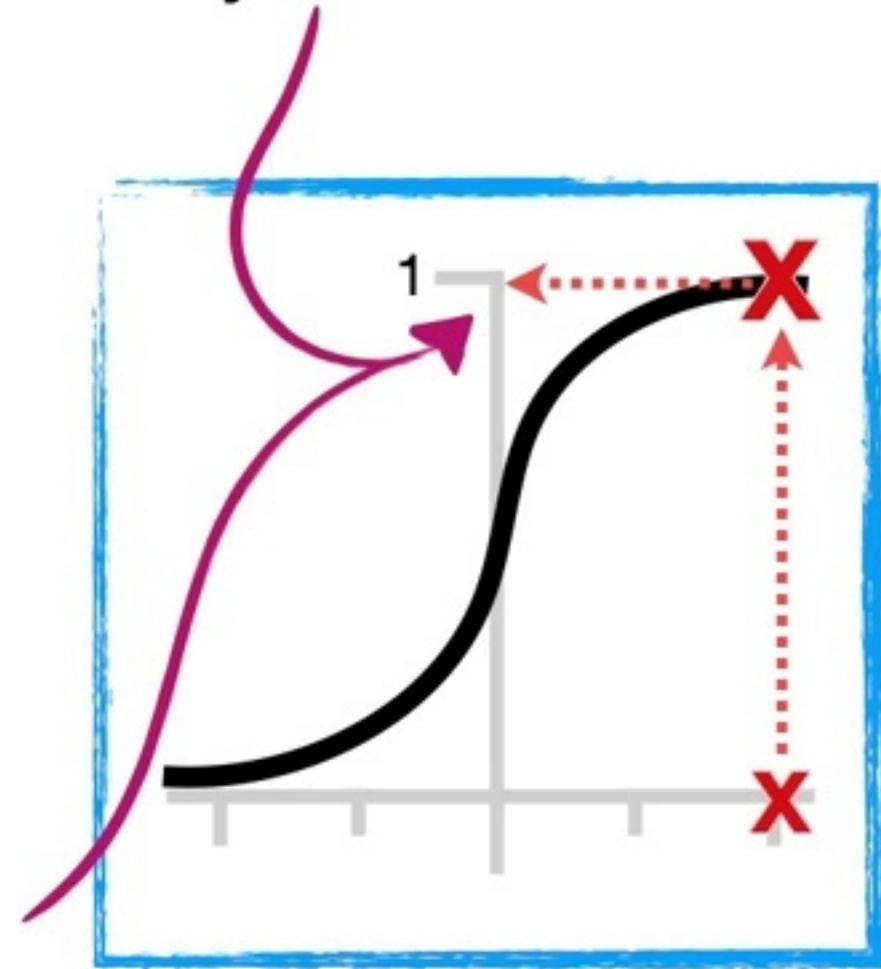
For example, when we plug in this x-axis coordinate, 10, into the equation for the **Sigmoid Activation Function**...





...we get **0.99995** as
the y-axis coordinate.

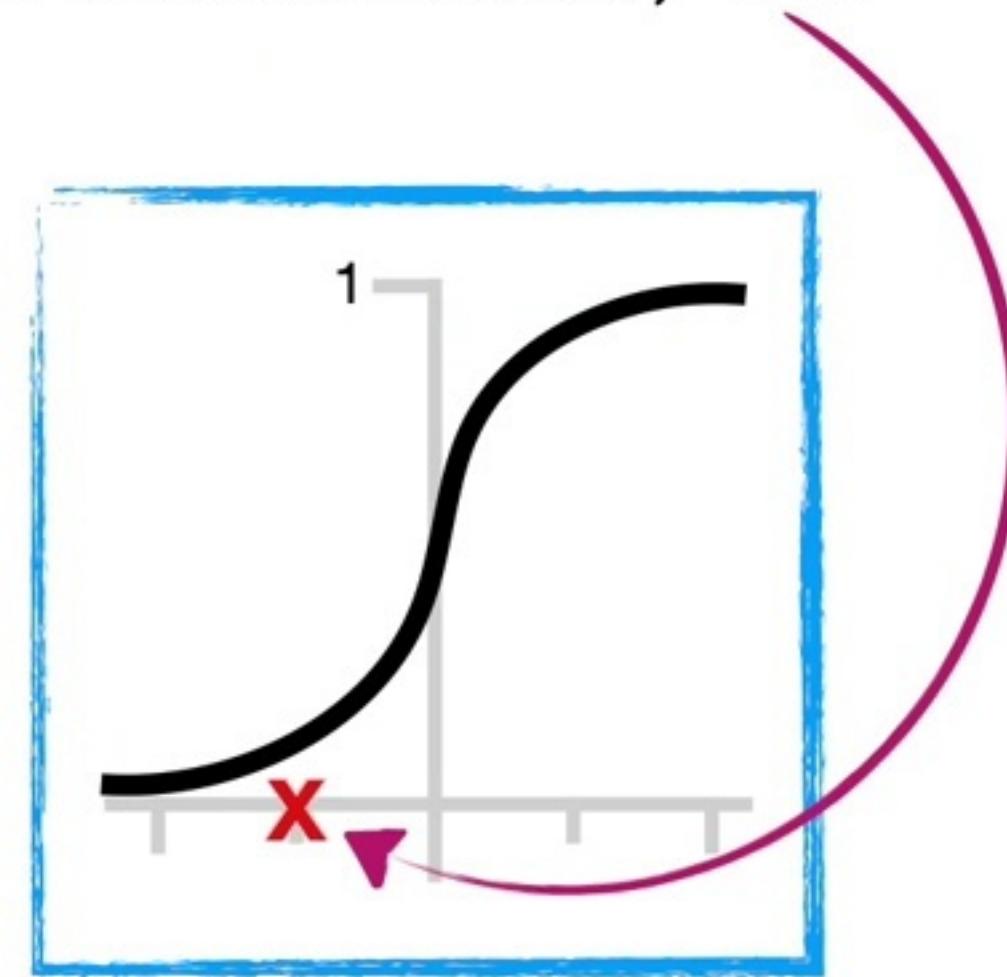
$$f(x) = \frac{e^x}{e^x + 1}$$
$$f(10) = \frac{e^{10}}{e^{10} + 1}$$
$$= 0.99995$$





And if we plug in this
x-axis coordinate, -5...

$$f(x) = \frac{e^x}{e^x + 1}$$



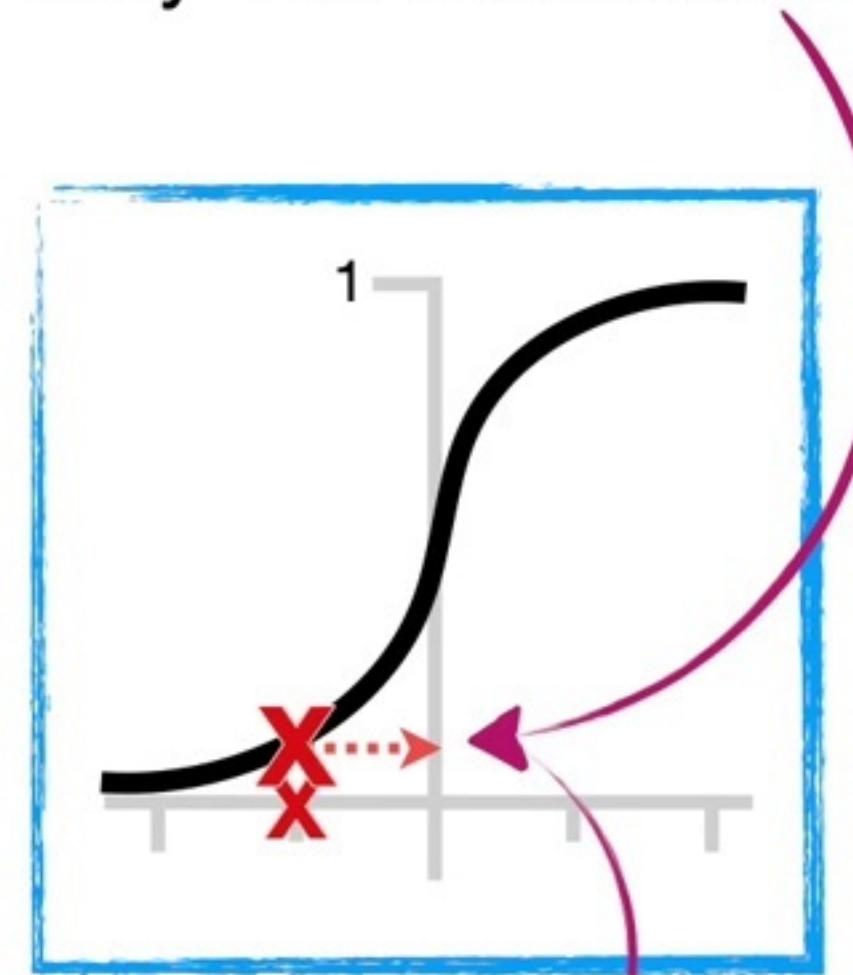


...then we get **0.01** as
the y-axis coordinate.

$$f(x) = \frac{e^x}{e^x + 1}$$

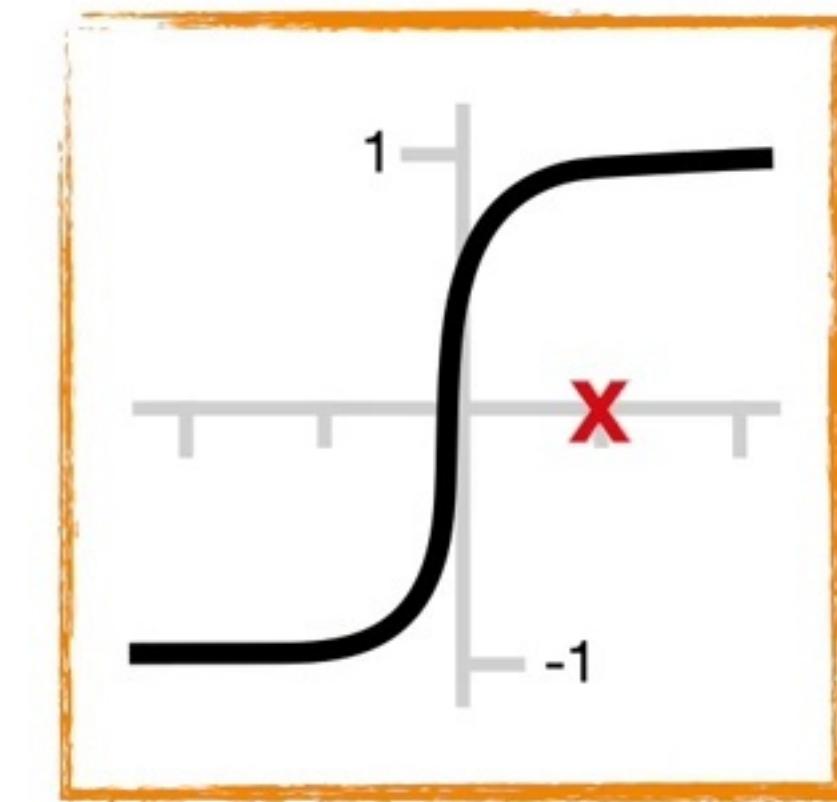
$$f(-5) = \frac{e^{-5}}{e^{-5} + 1}$$

$$= 0.01$$



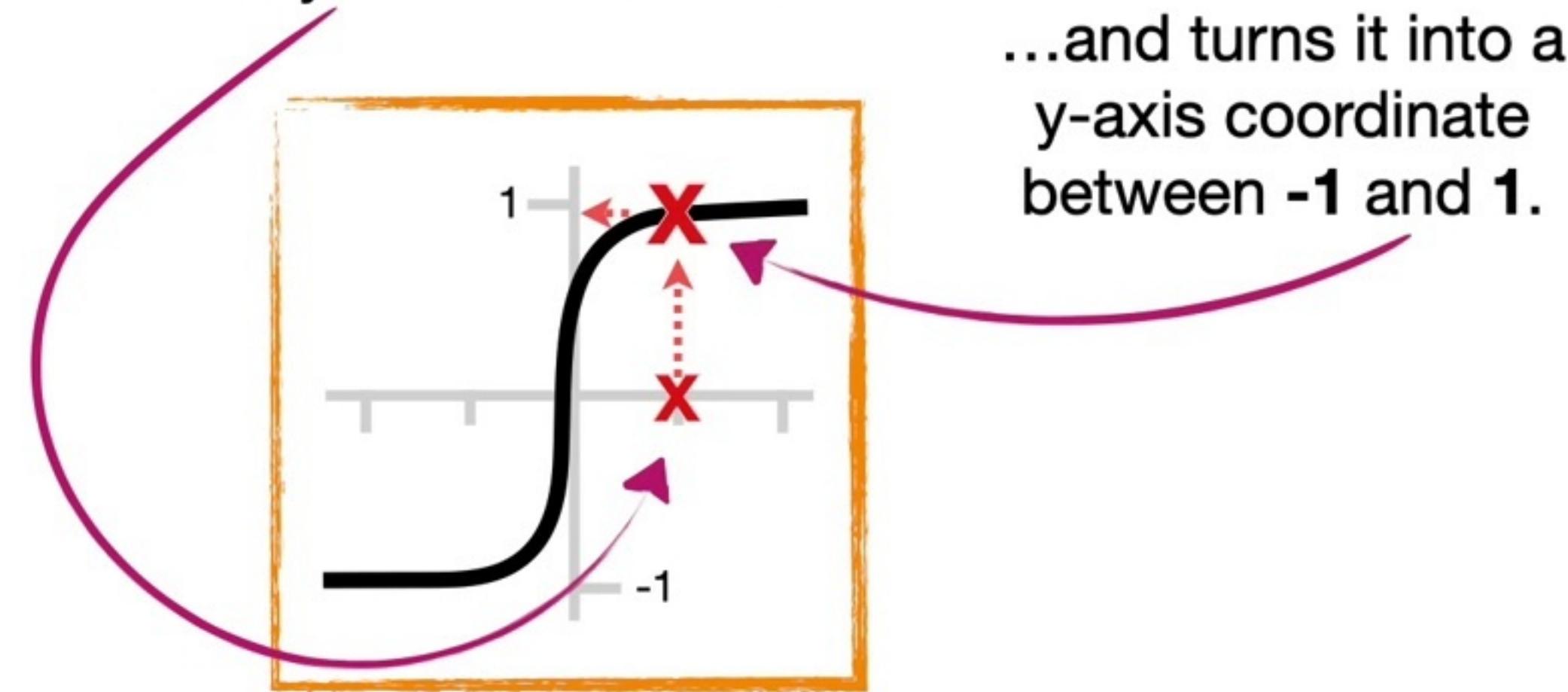


In contrast, the **tanh**, or **Hyperbolic Tangent, Activation Function** takes any x-axis coordinate...





In contrast, the **tanh**, or **Hyperbolic Tangent, Activation Function** takes any x-axis coordinate...



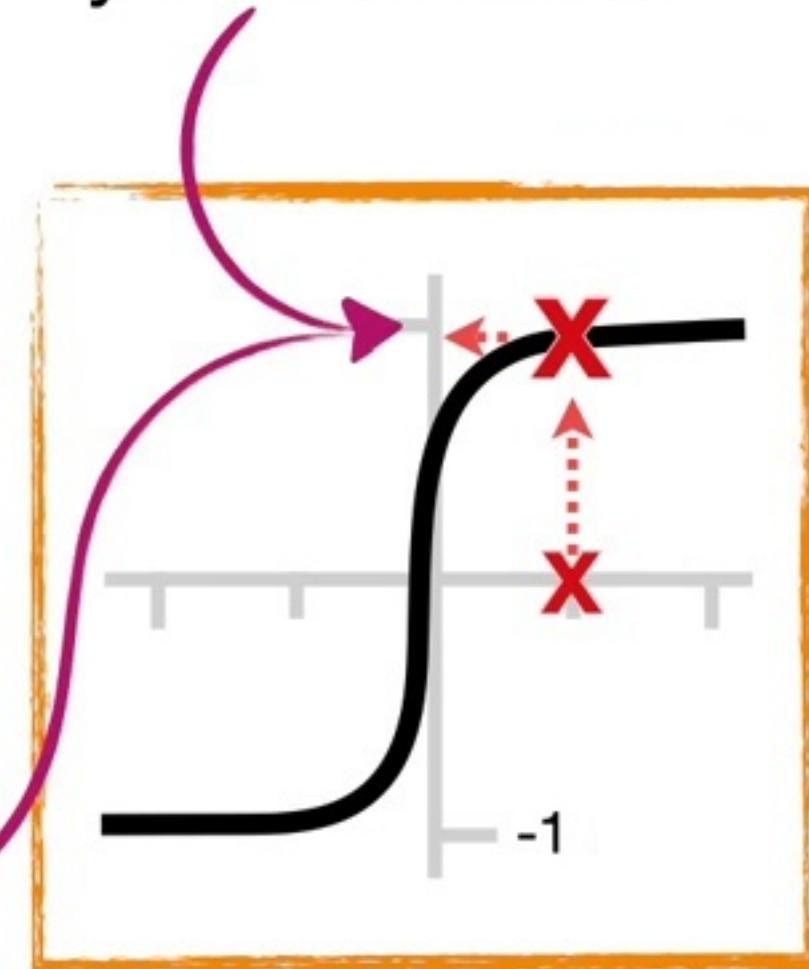


...we get **0.96** as the
y-axis coordinate.

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(2) = \frac{e^2 - e^{-2}}{e^2 + e^{-2}}$$

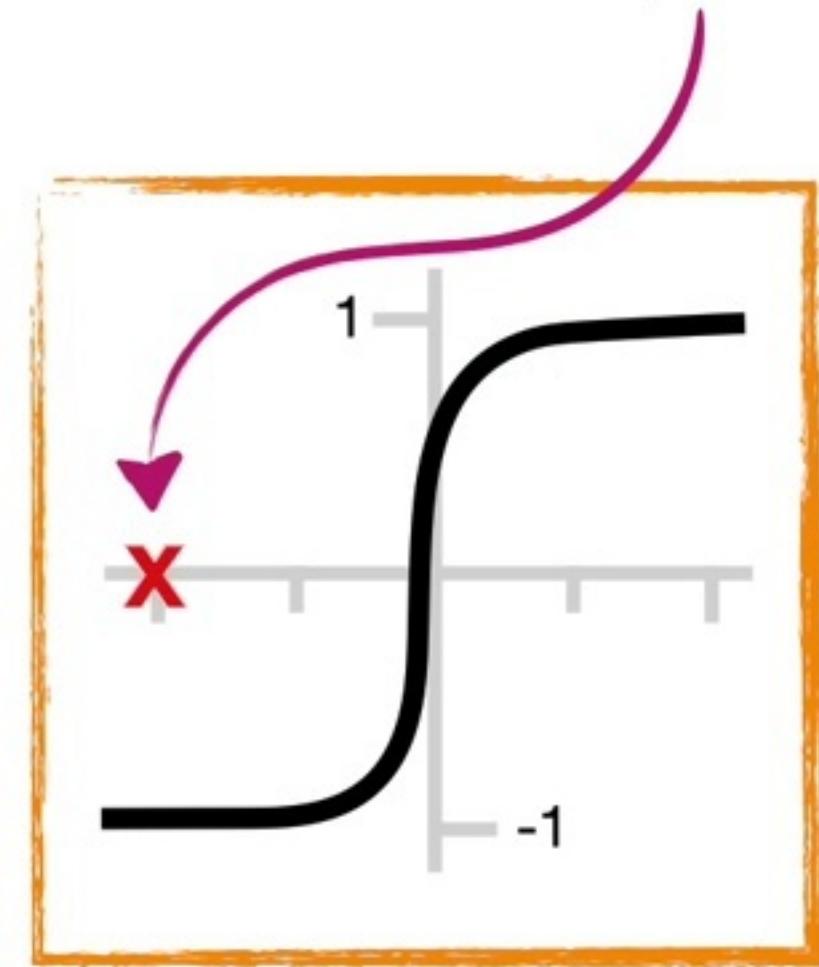
$$= 0.96$$





And if we plug in this
x-axis coordinate, -5...

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



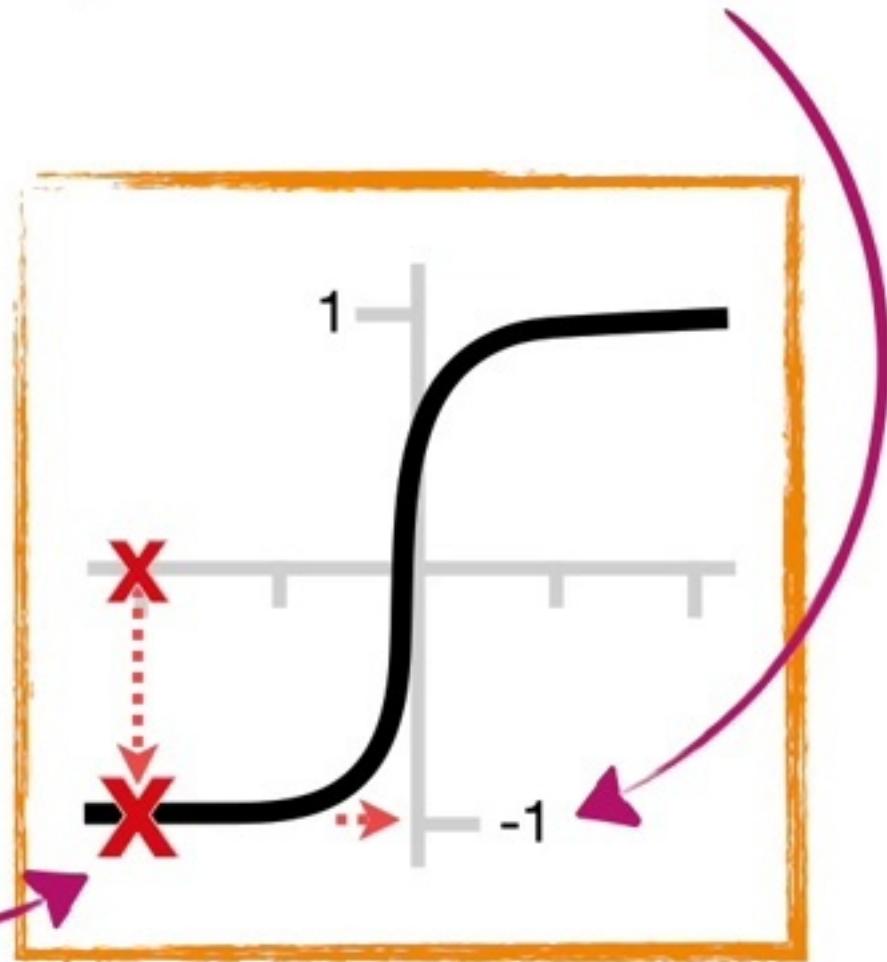


...we get -1 as the
y-axis coordinate.

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

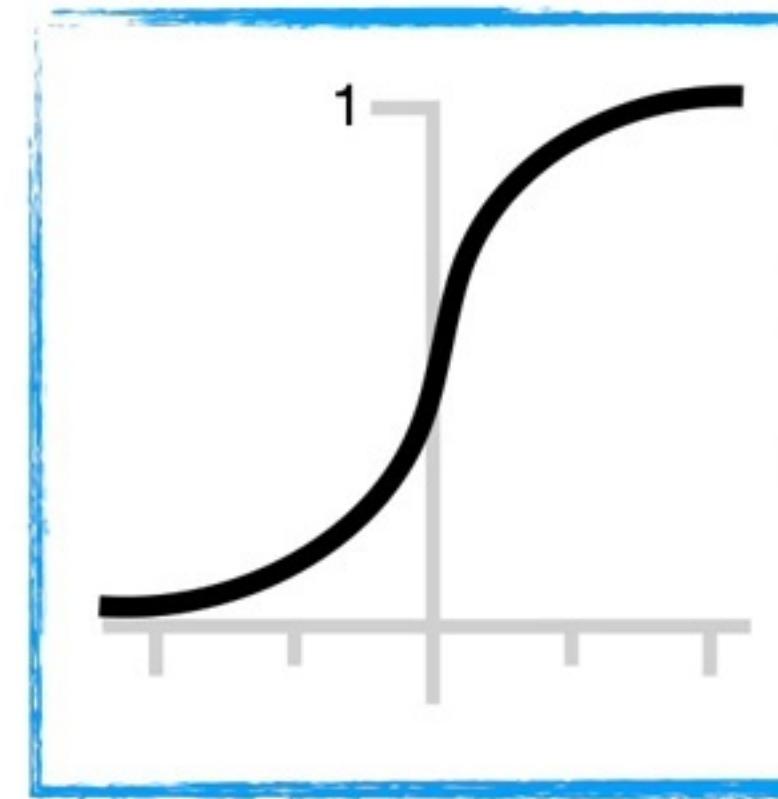
$$f(-5) = \frac{e^{-5} - e^5}{e^{-5} + e^5}$$

$$= -1$$

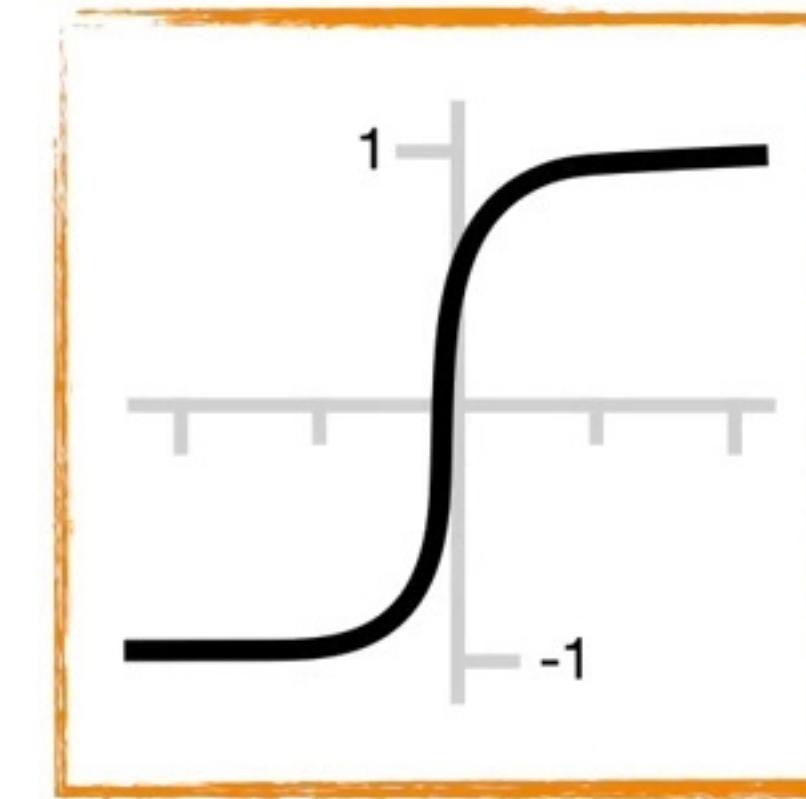




So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...

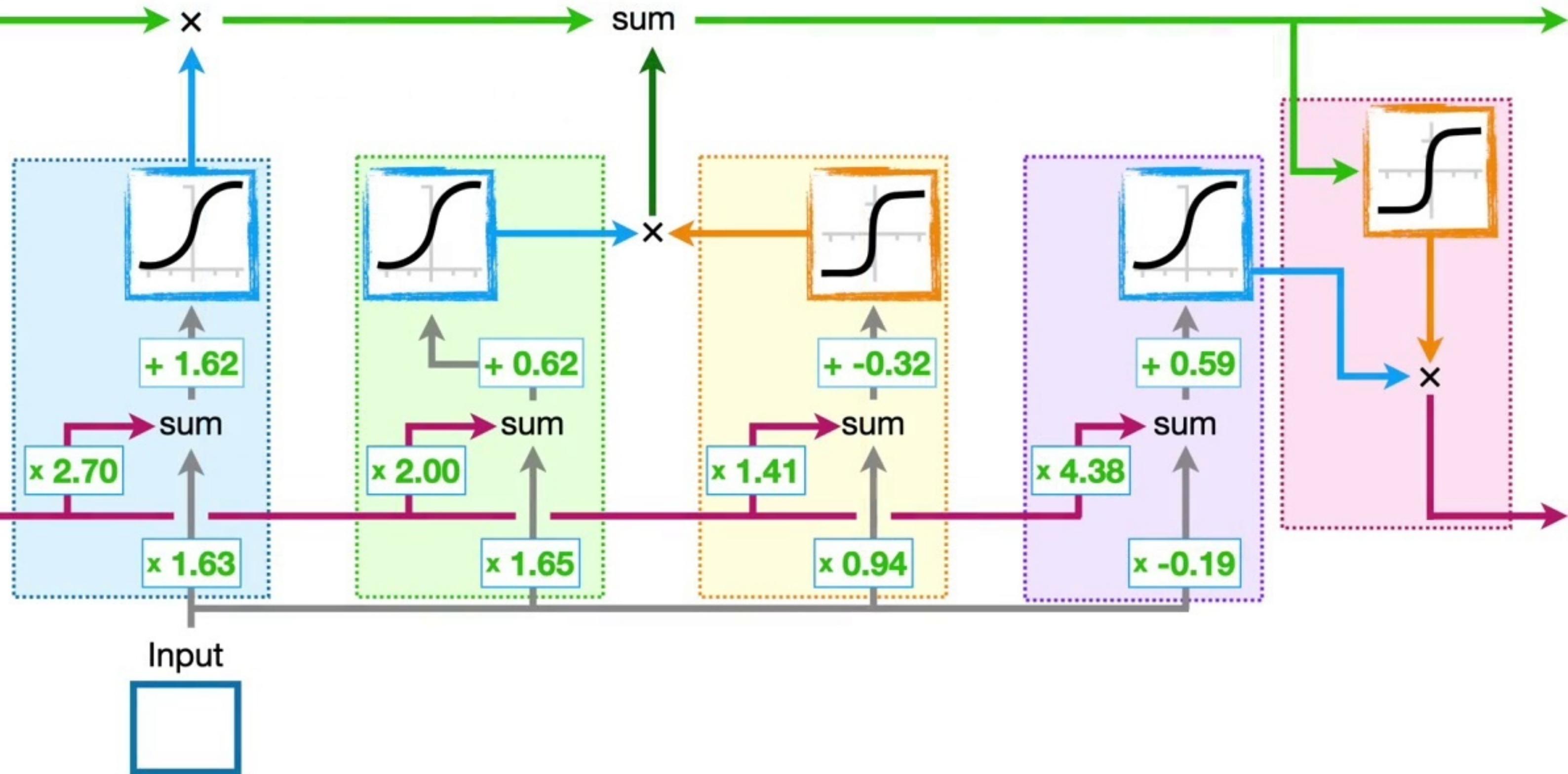


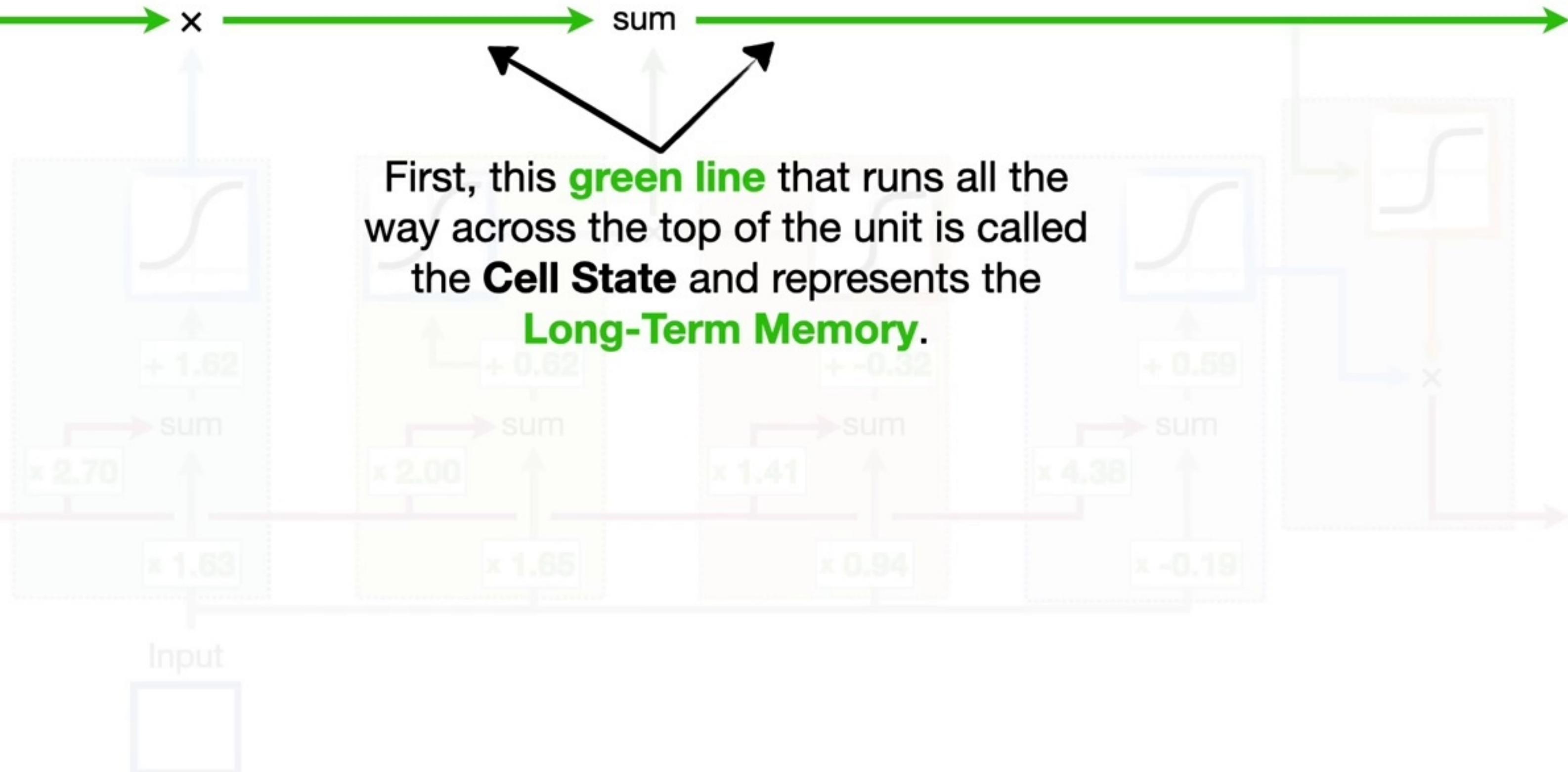
...and the **Tanh Activation Function** turns any input into a number between **-1** and **1**...

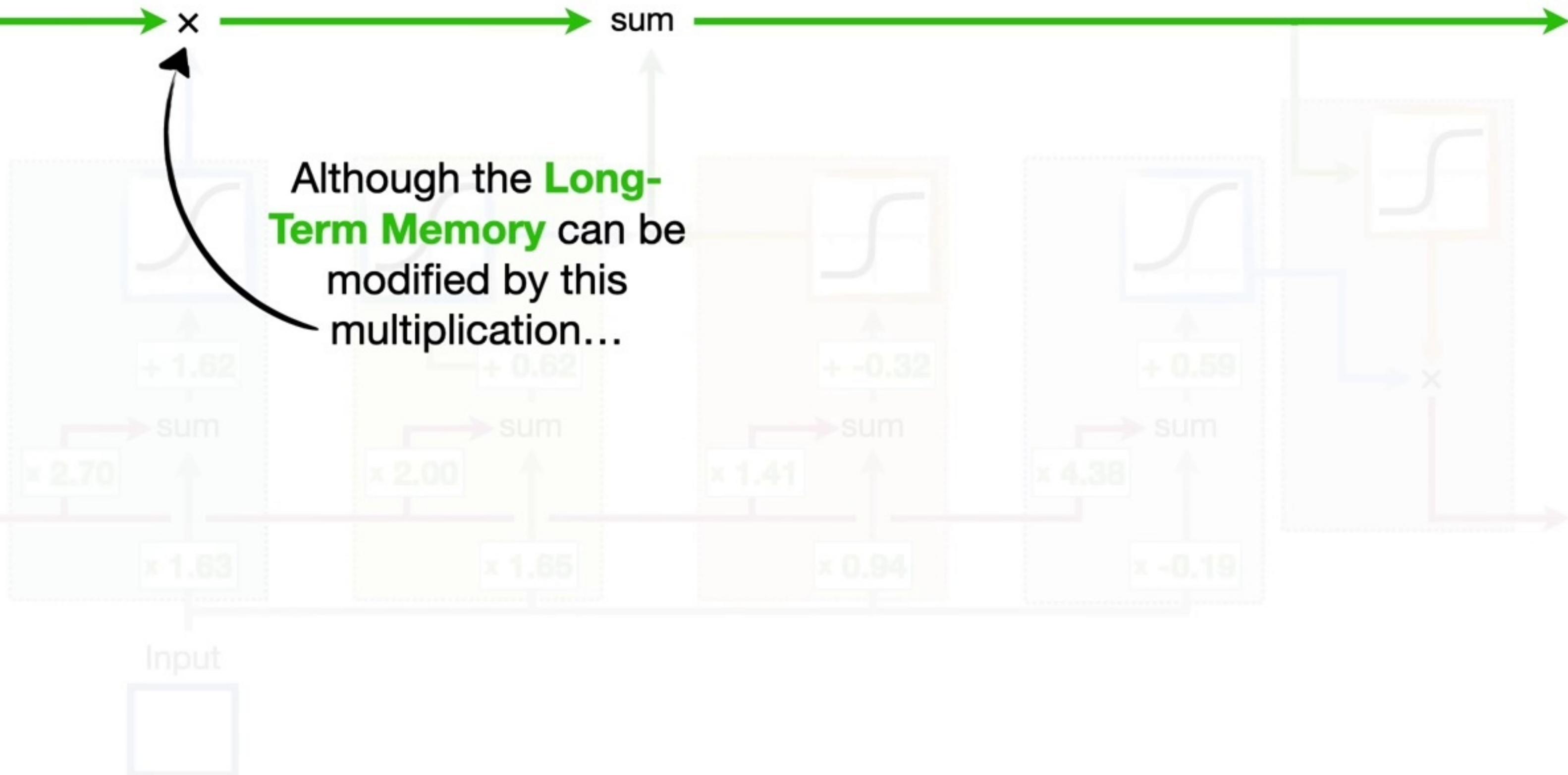


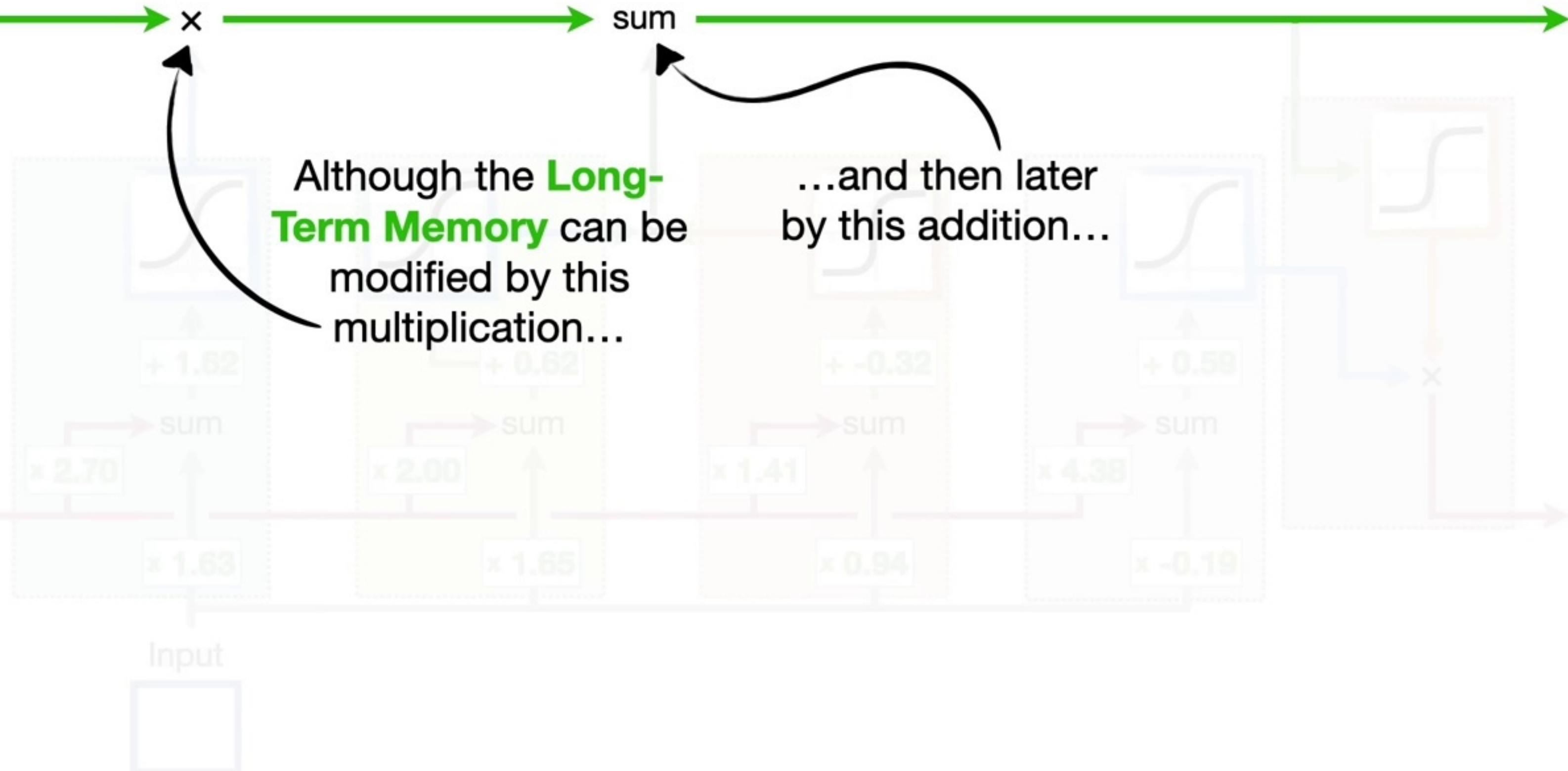


...let's talk about how the **Long Short-Term Memory** unit works.











x

sum

...you'll notice that there are
no **Weights** and **Biases** that
can modify it directly.



$$+ 1.62$$

sum

$$\times 2.70$$

$$\times 1.63$$

Input



$$+ 0.62$$

sum

$$\times 2.00$$

$$\times 1.65$$

sum

$$\times 1.41$$

sum

$$\times 0.94$$

$$+ -0.32$$

sum

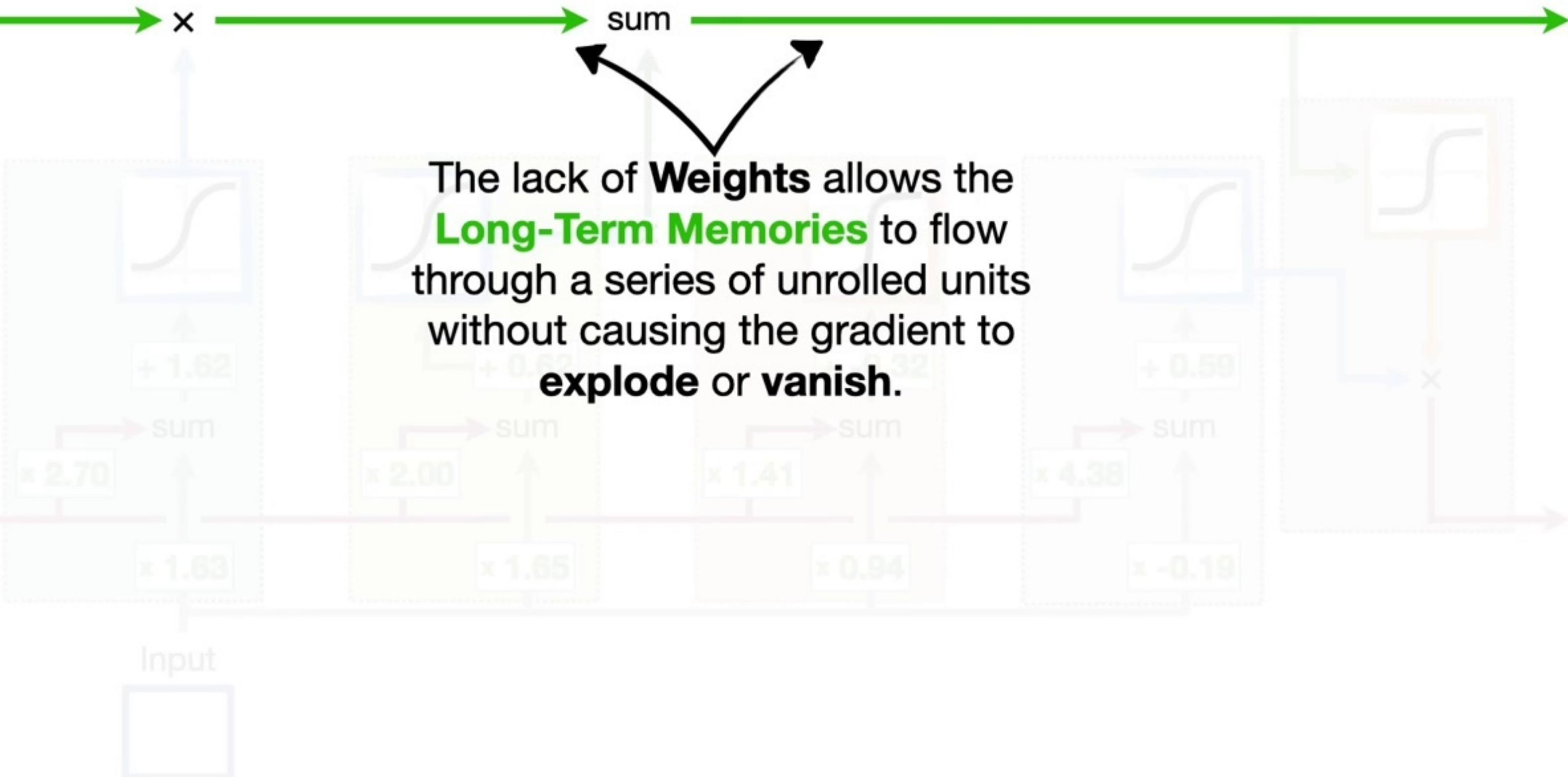
$$\times 4.38$$

$$+ 0.59$$

sum

$$\times -0.19$$

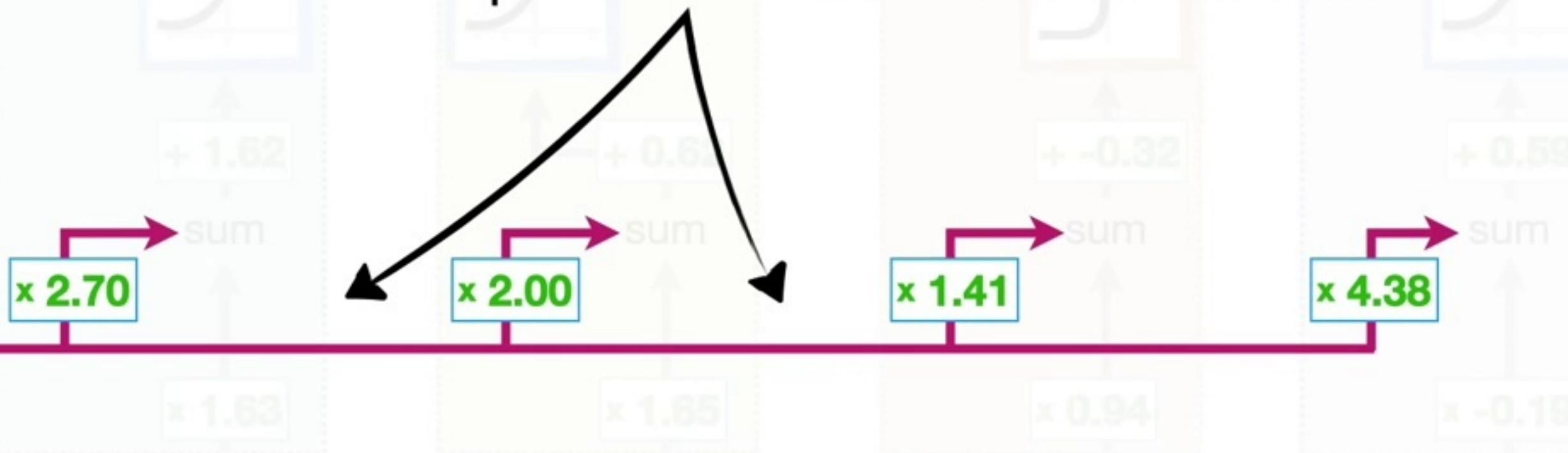




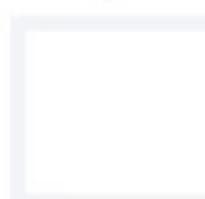


$$x \longrightarrow \text{sum}$$

Now, this **pink line**, called the **Hidden State**,
represents the **Short-Term Memories**.



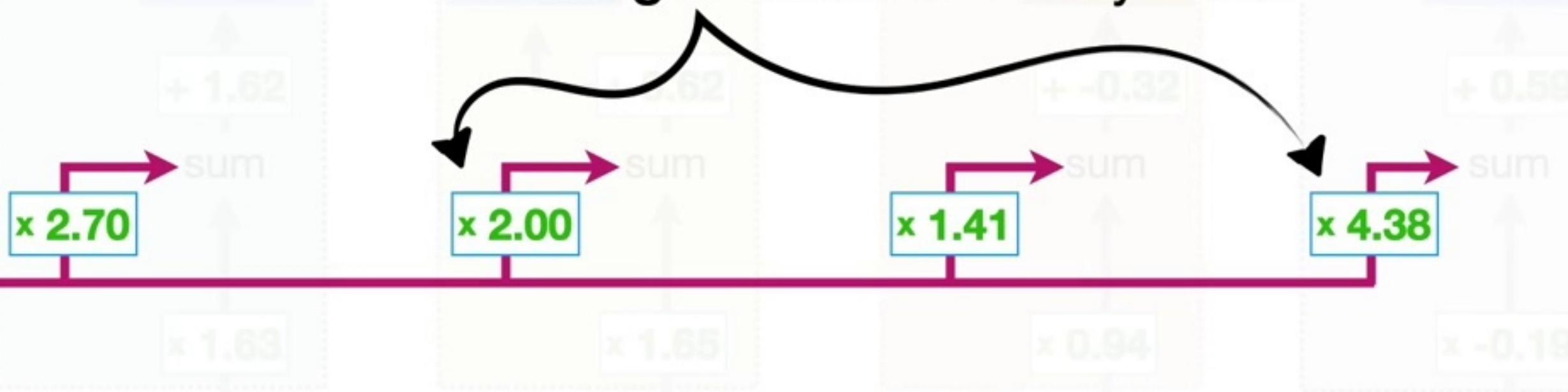
Input





x ————— sum

And, as we can see, the **Short-Term
Memories** are directly connected to
Weights that can modify them.





x → sum

To understand how the **Long** and **Short-Term Memories** interact and result in predictions, let's run some numbers through this unit.

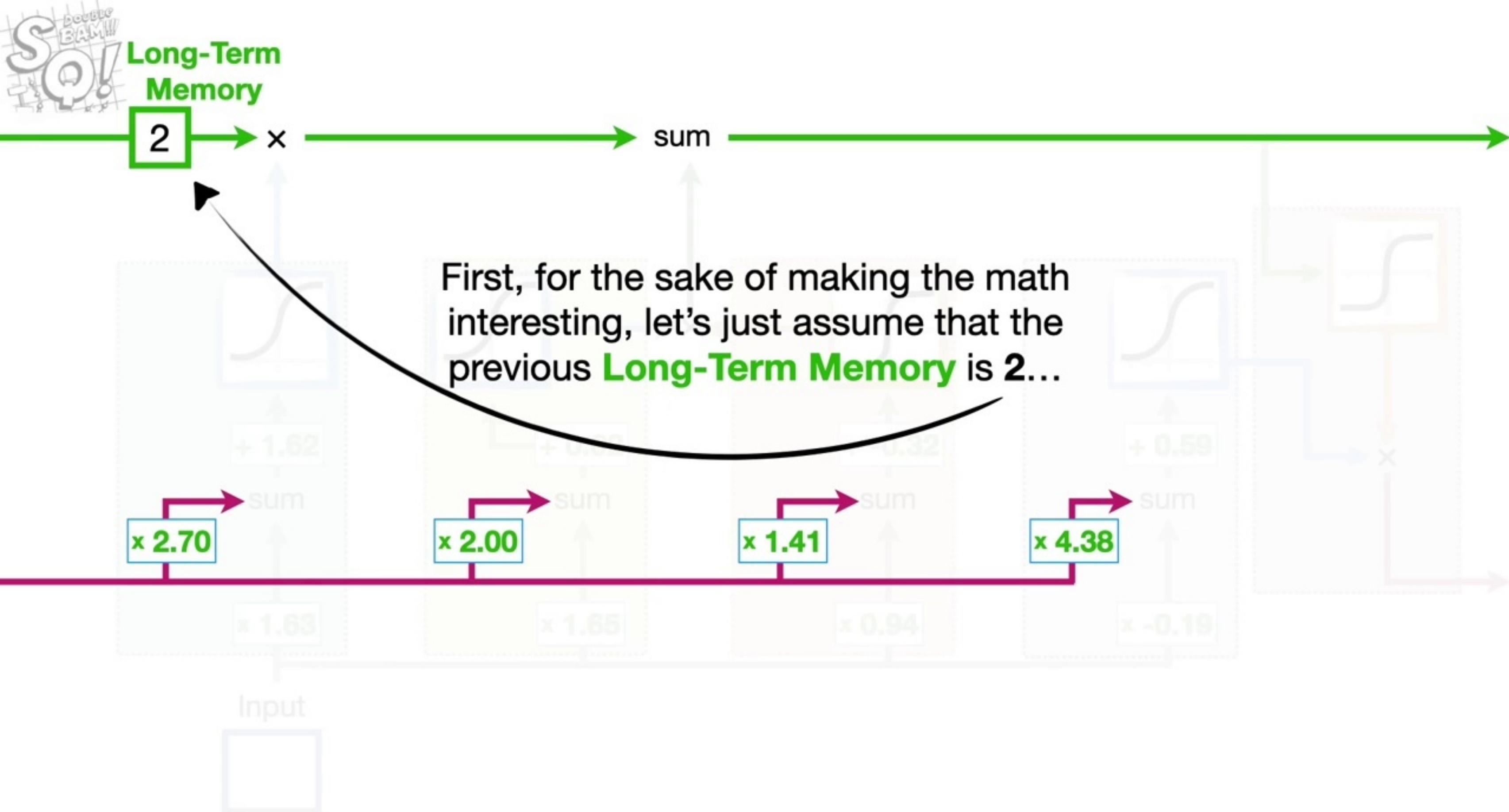
x 2.70

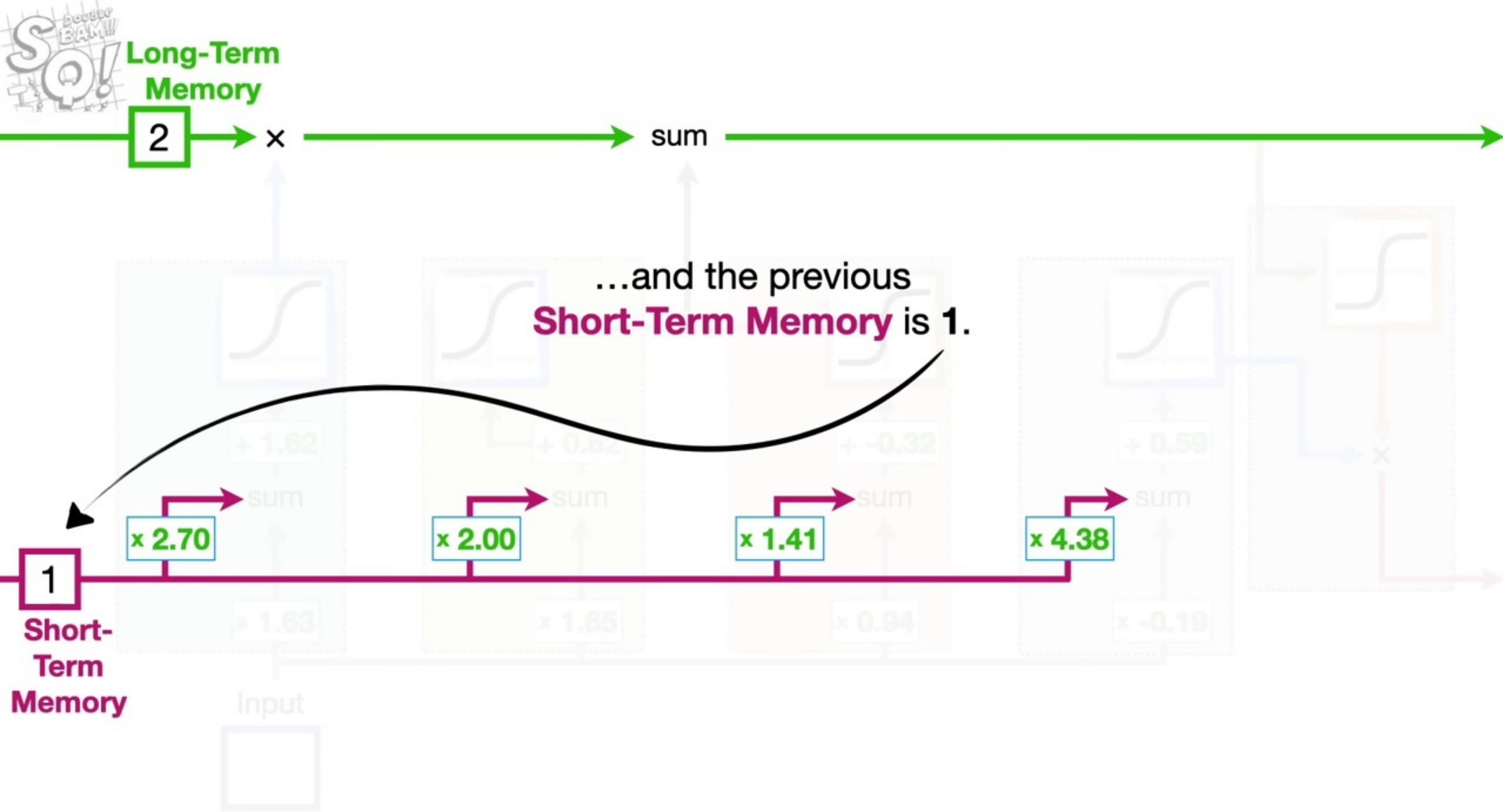
x 2.00

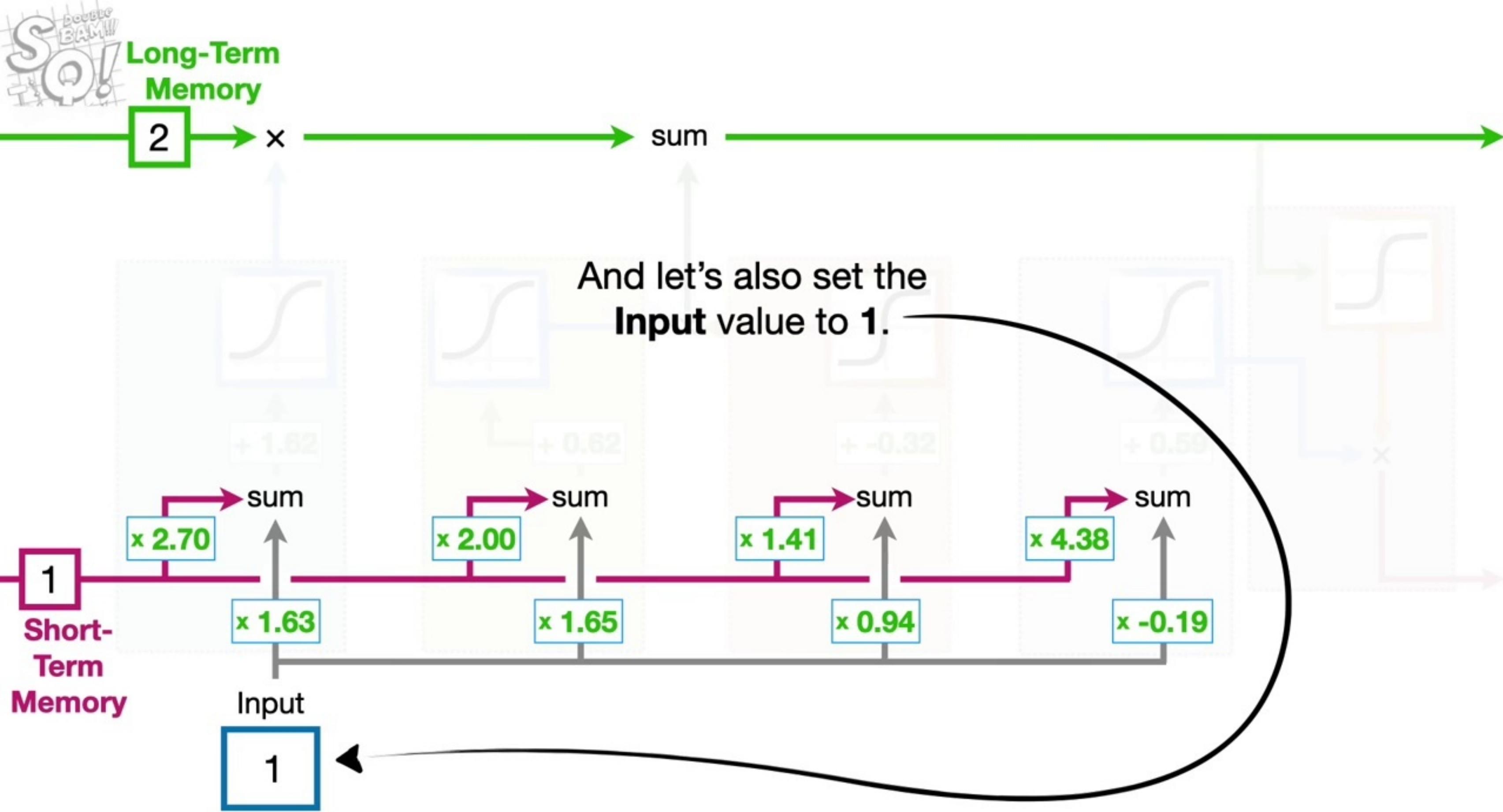
x 1.41

x 4.38

Input

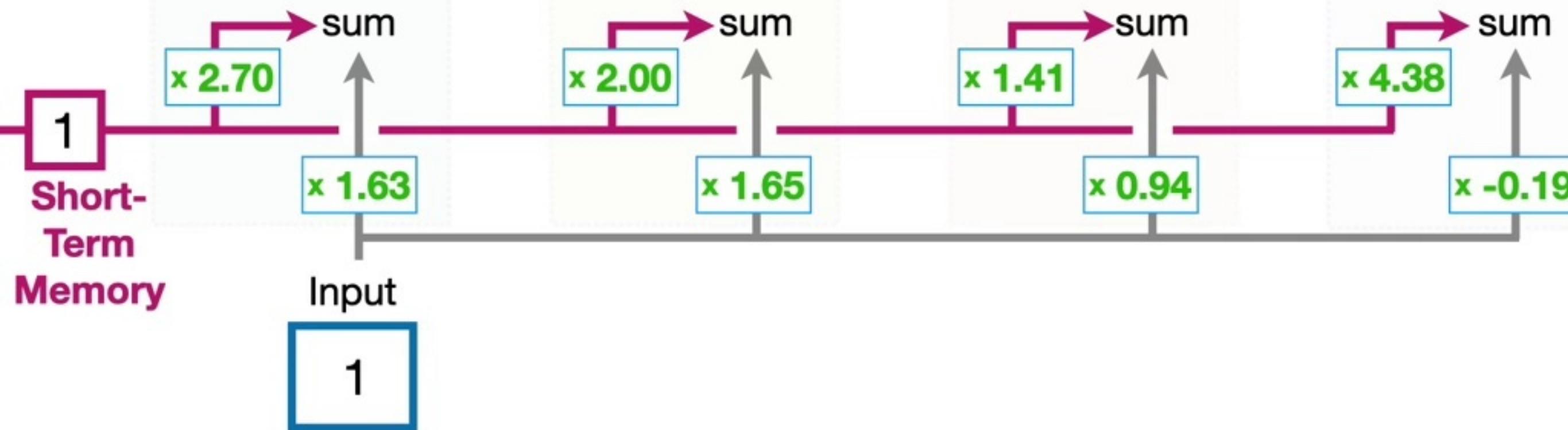


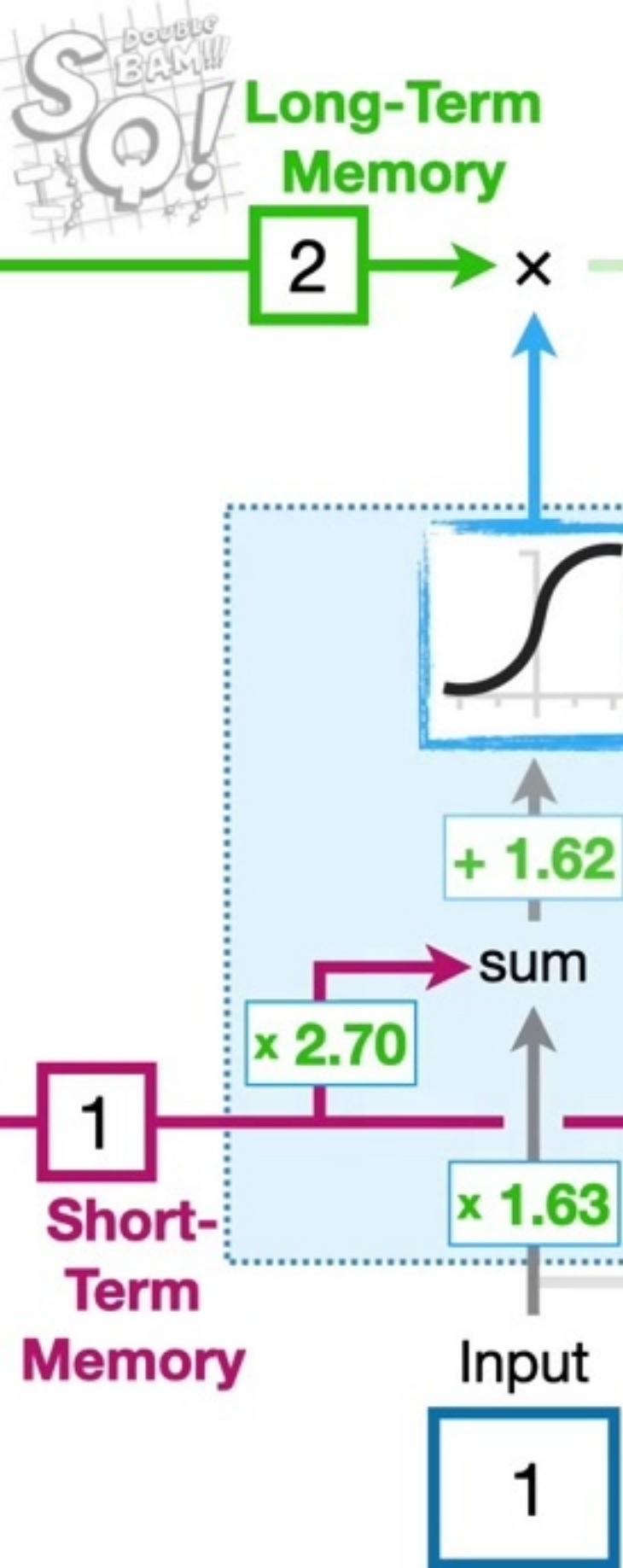




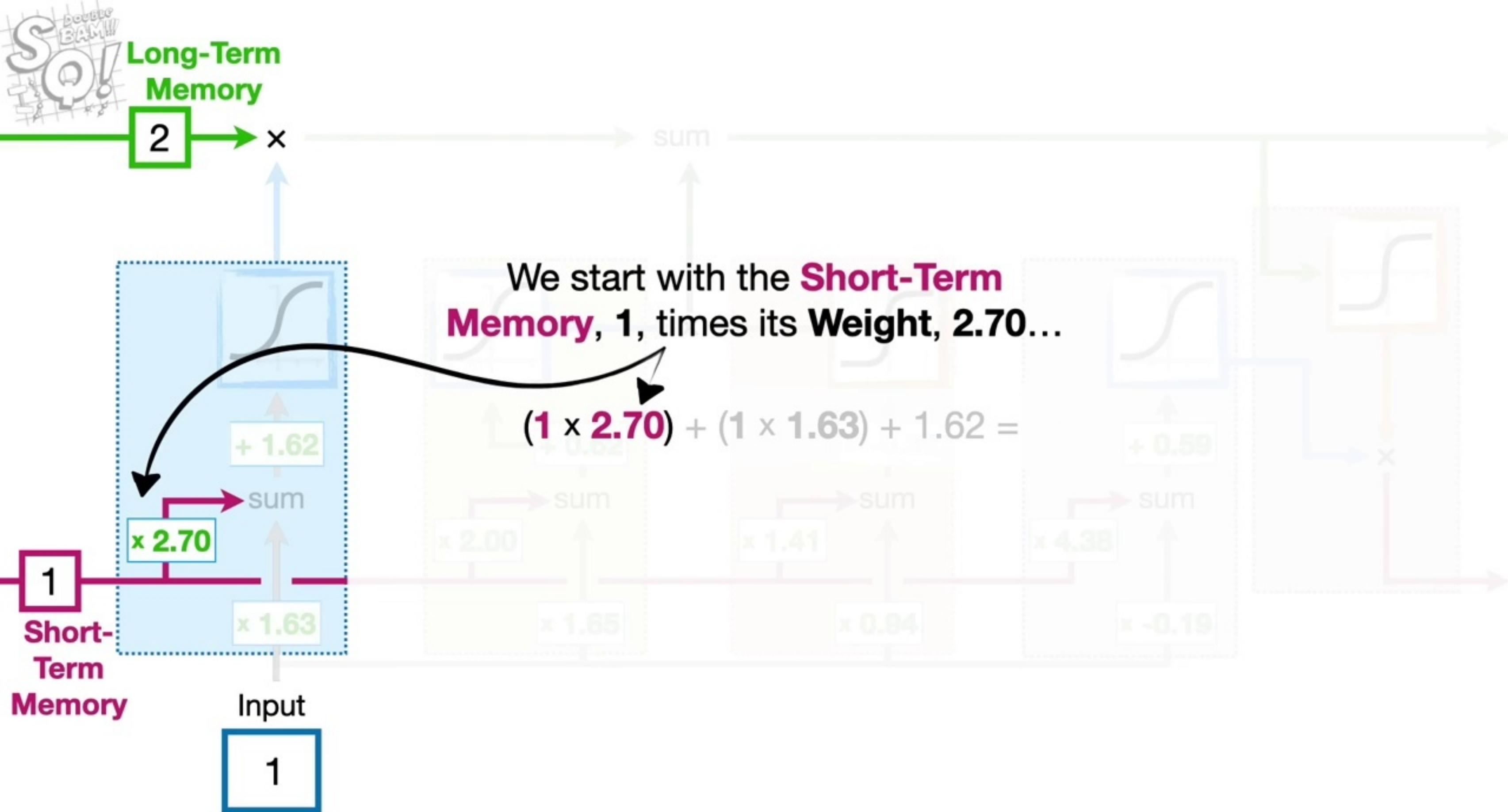


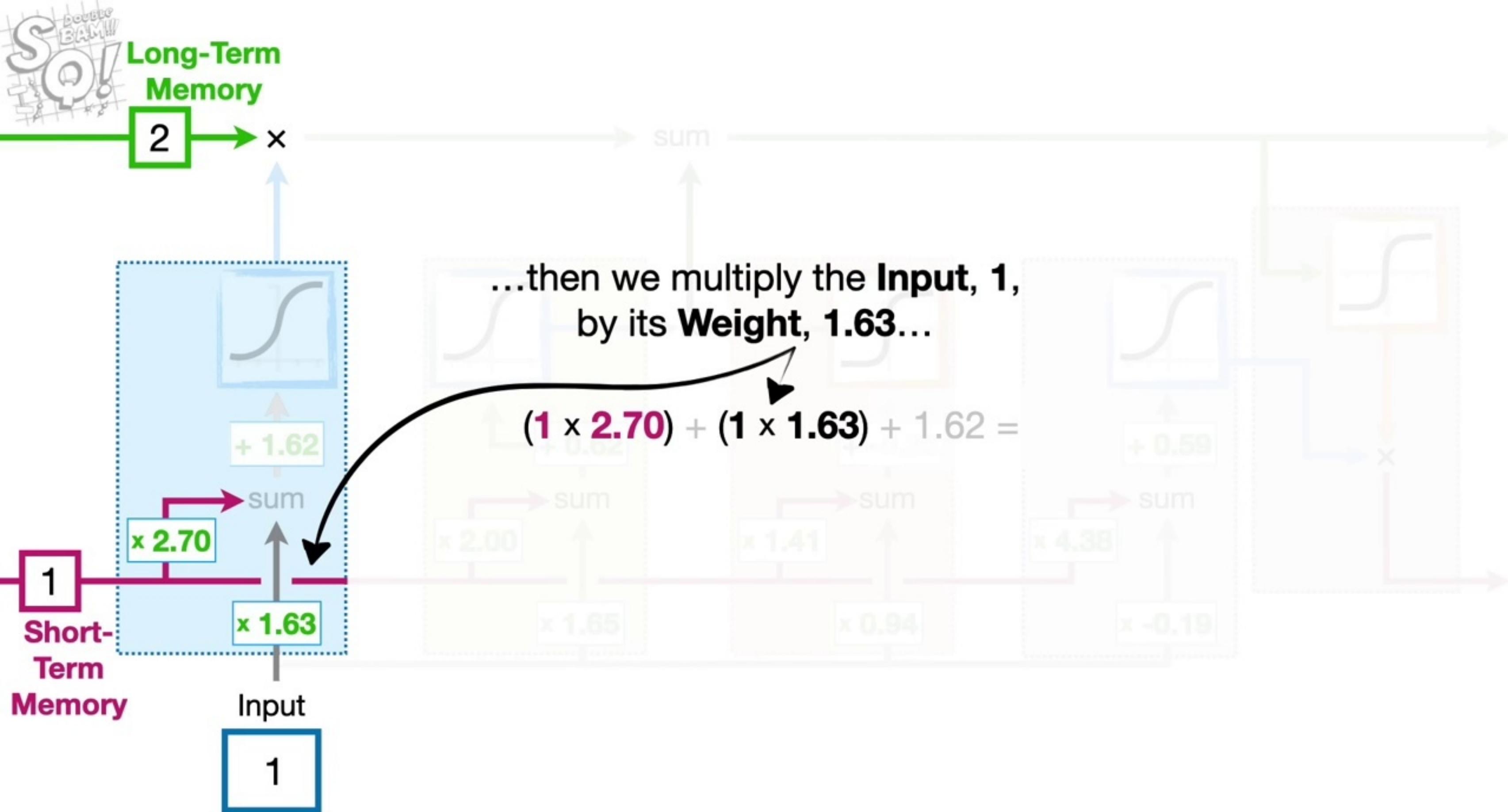
Now that we have plugged in some numbers, let's do the math to see what happens in the first stage of the **Long Short-Term Memory** unit.

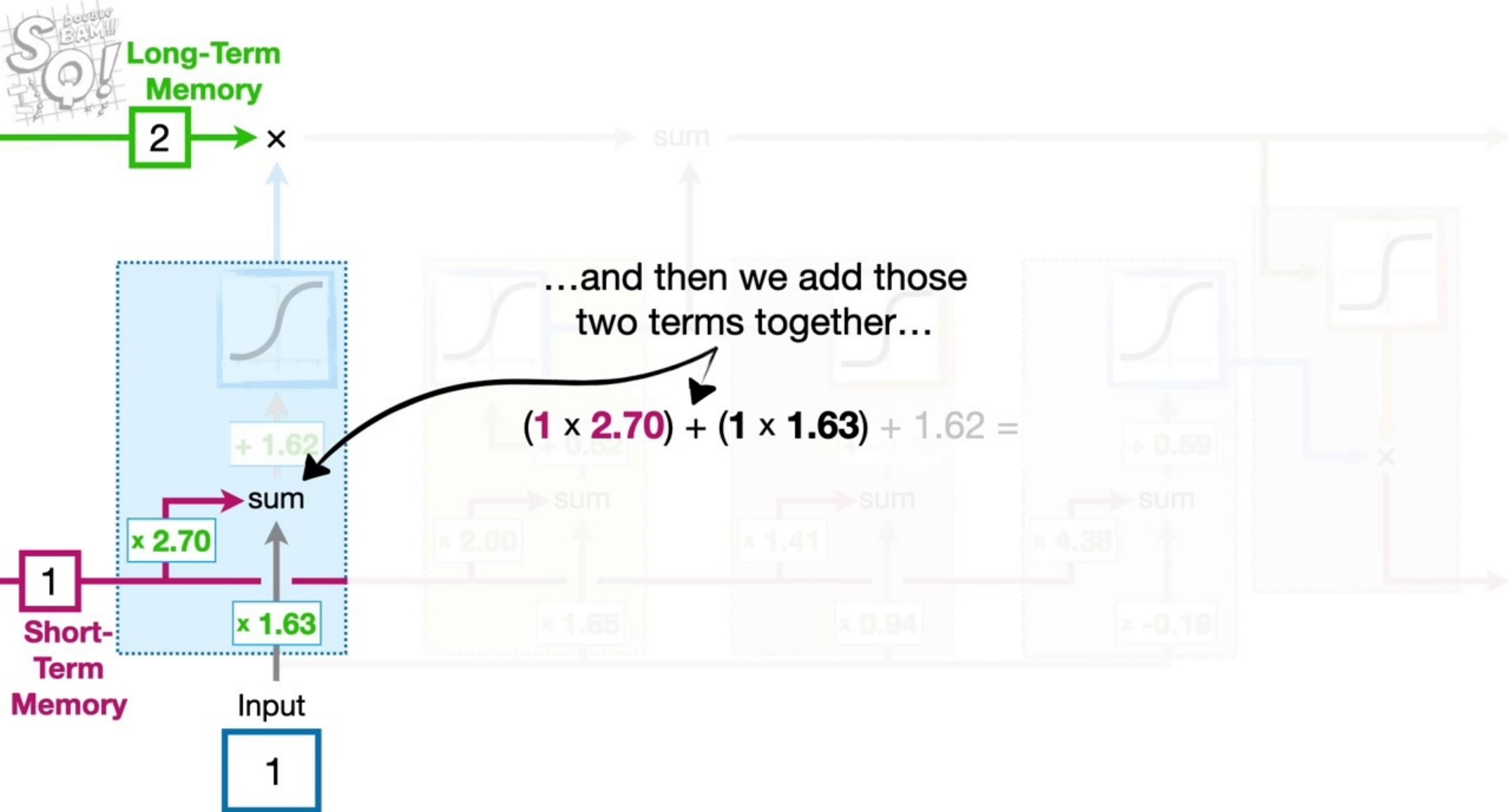


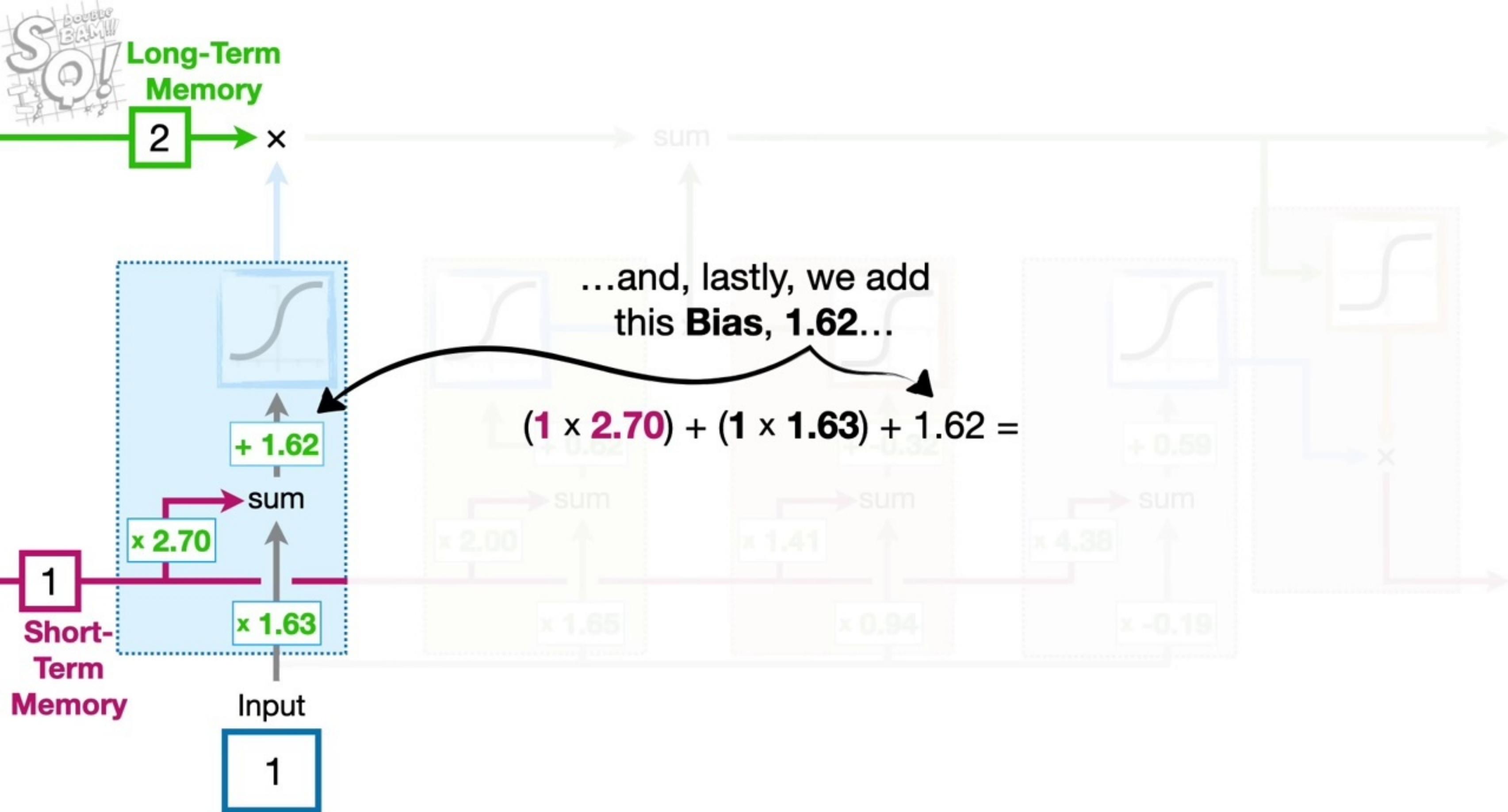


Now that we have plugged in some numbers, let's do the math to see what happens in the first stage of the **Long Short-Term Memory** unit.



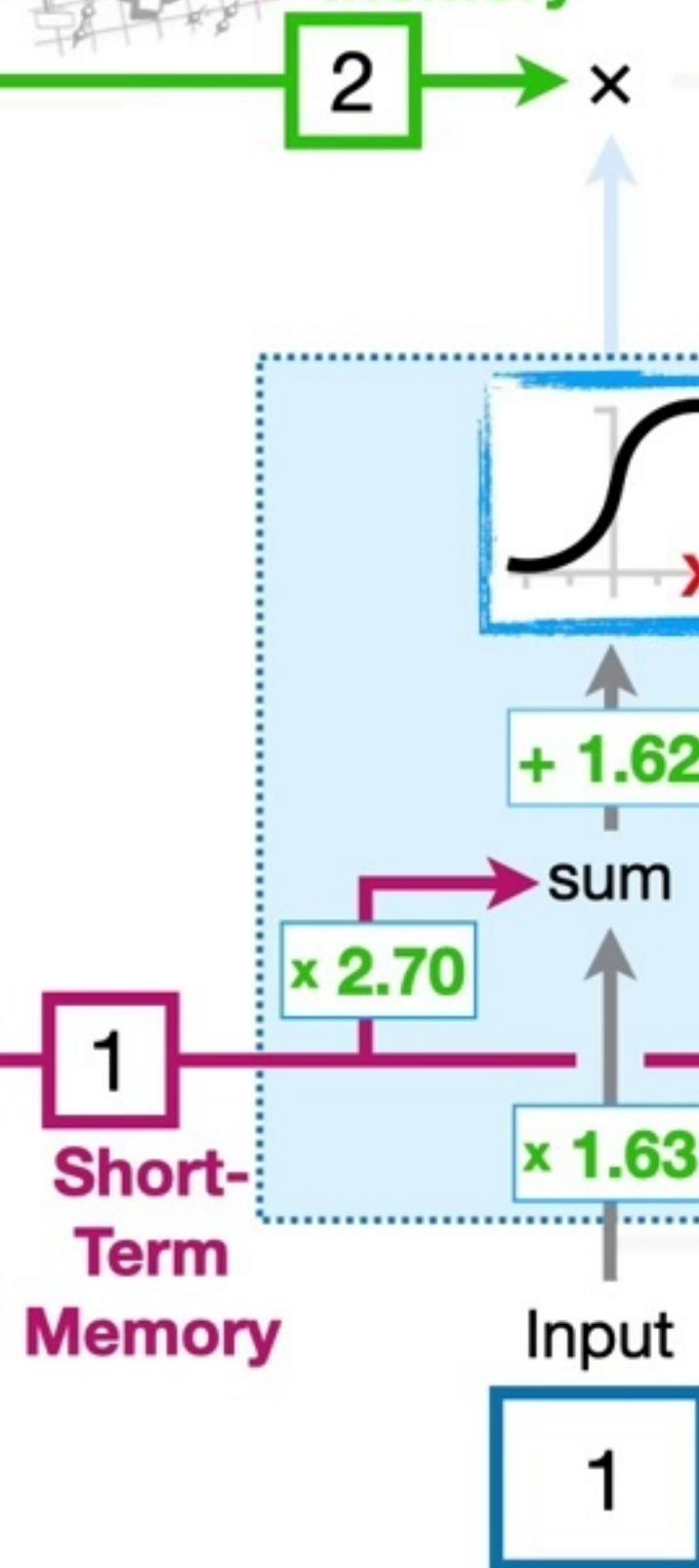






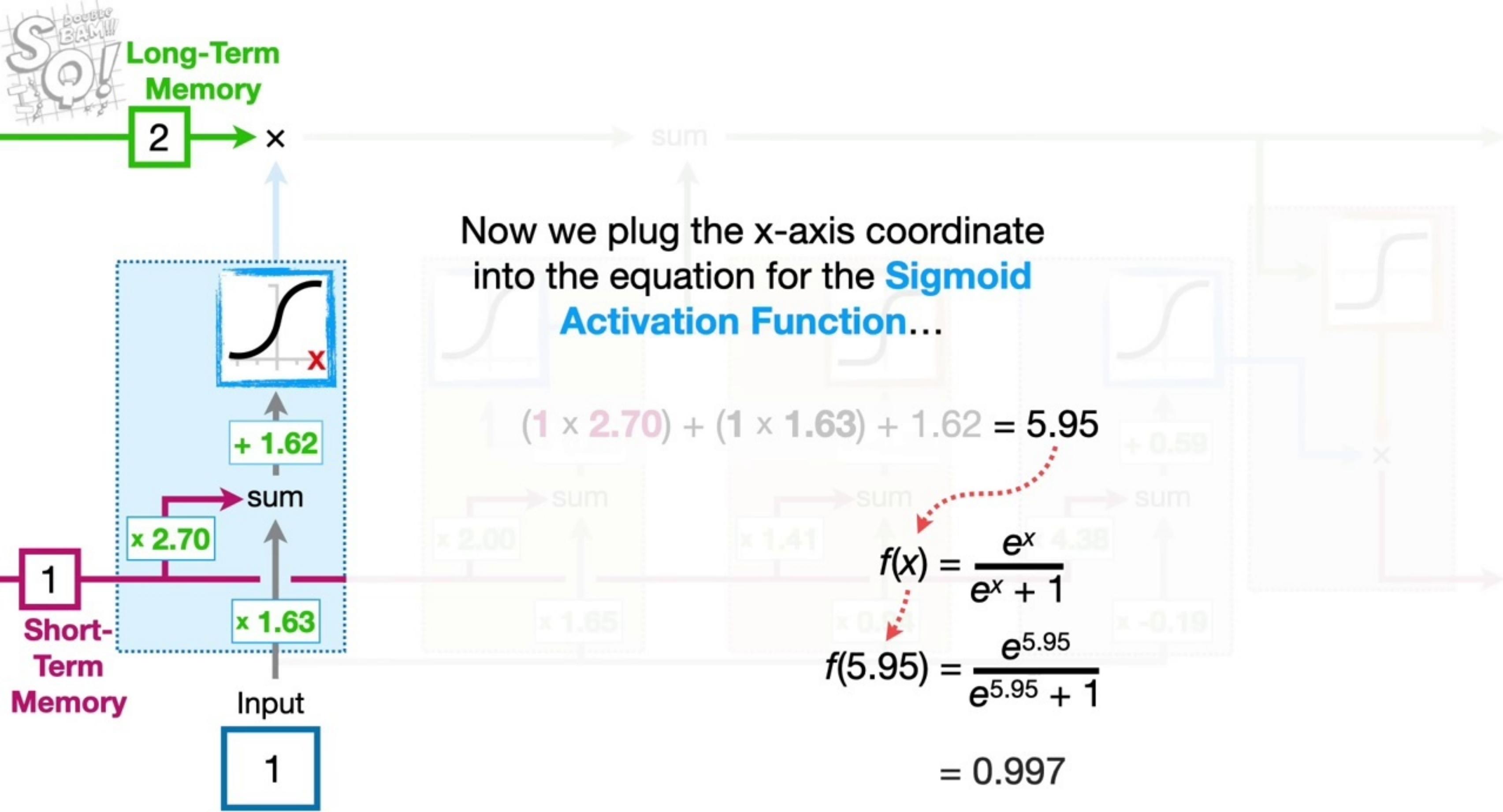
SQ!

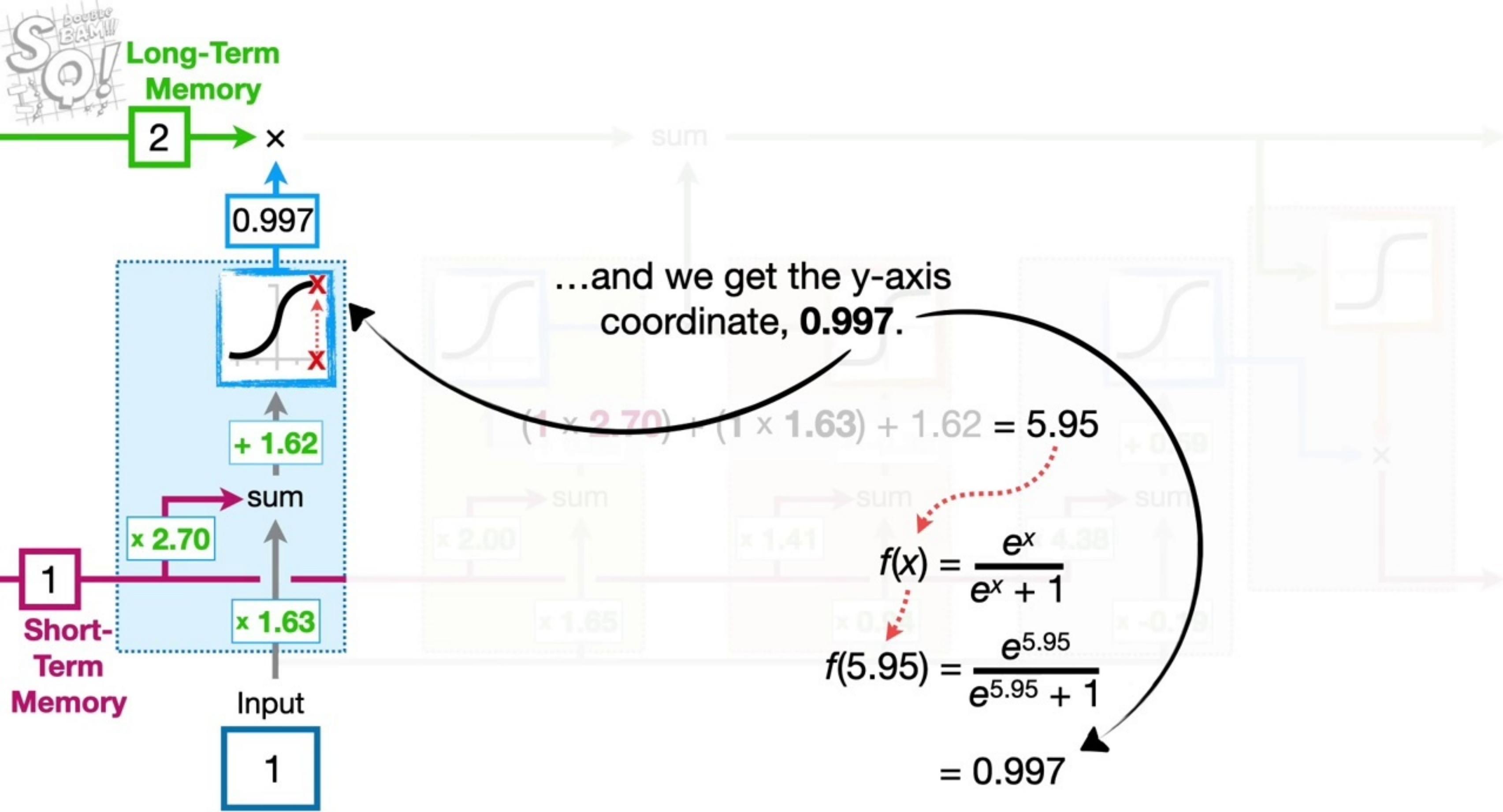
Long-Term Memory

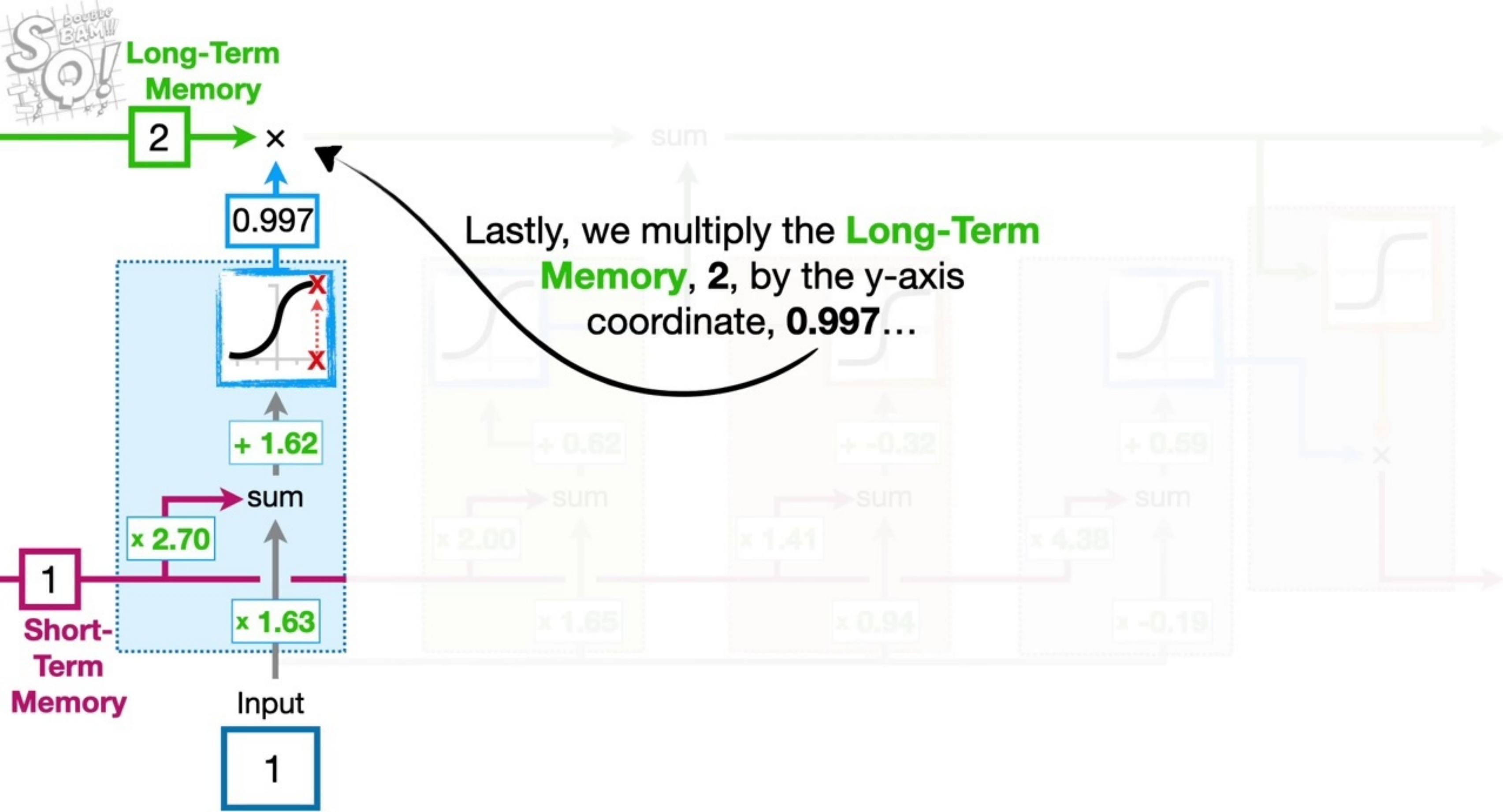


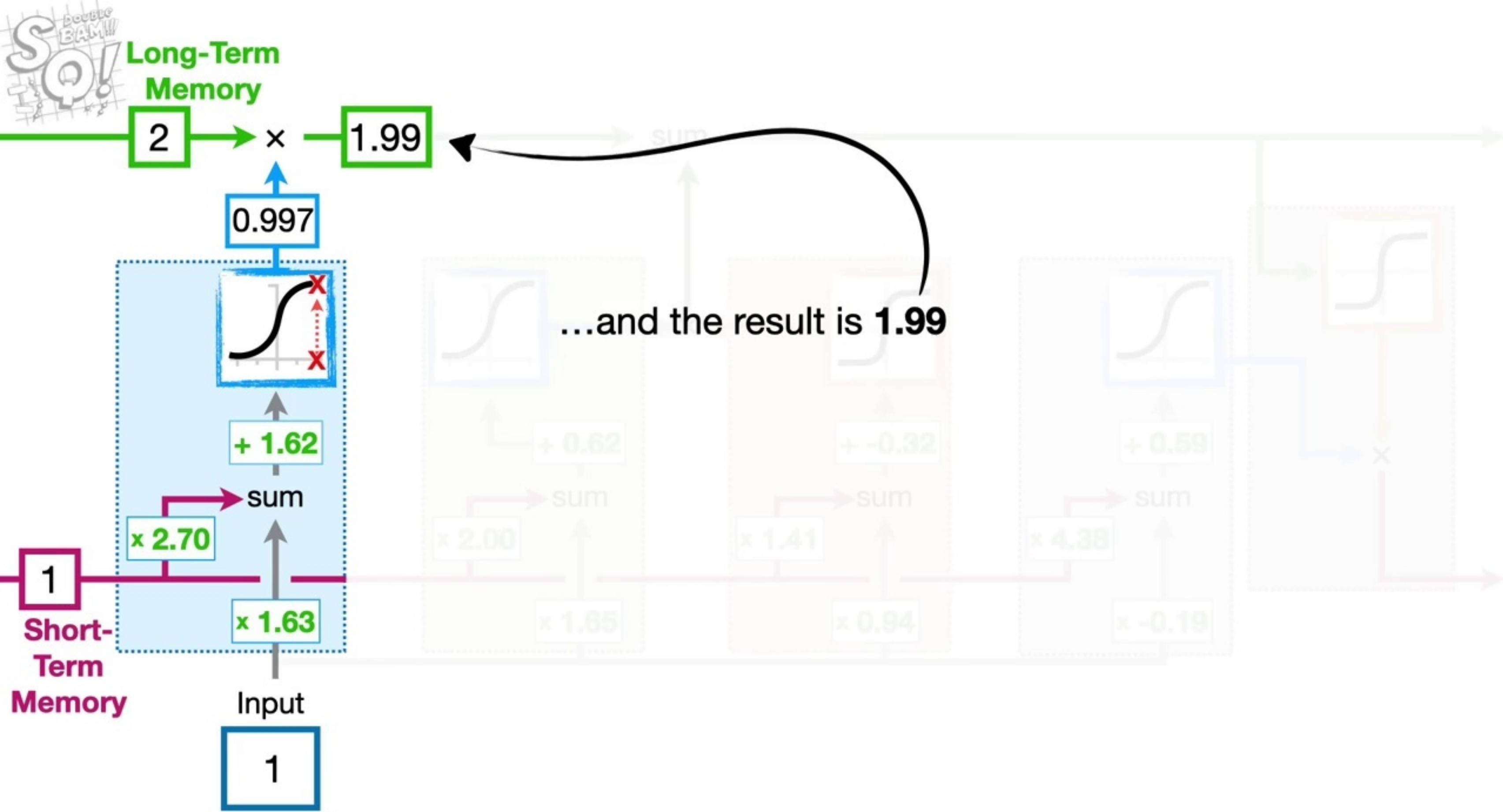
...to get **5.95**, an x-axis coordinate for the **Sigmoid Activation Function**.

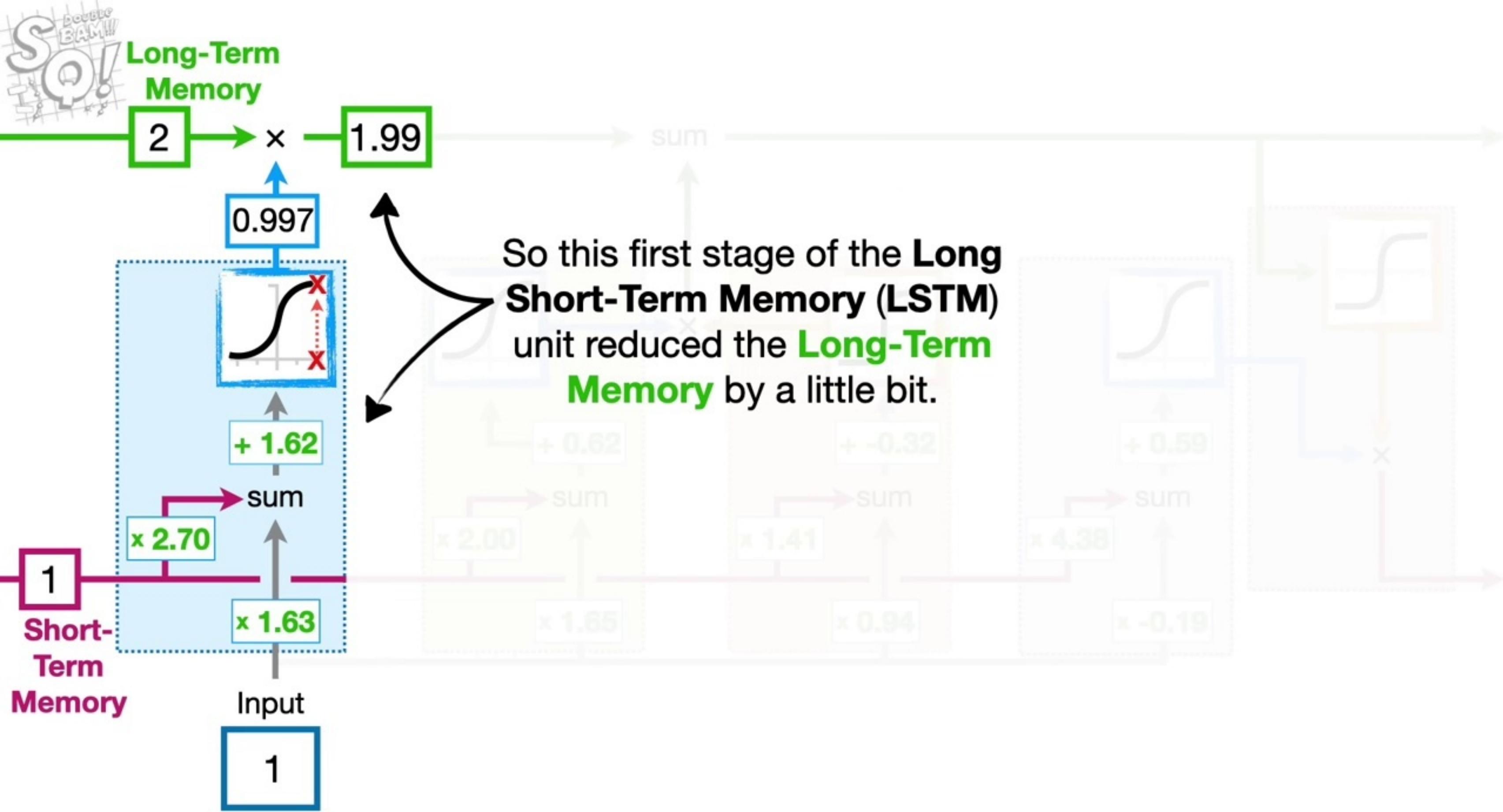
$$(1 \times 2.70) + (1 \times 1.63) + 1.62 = 5.95$$

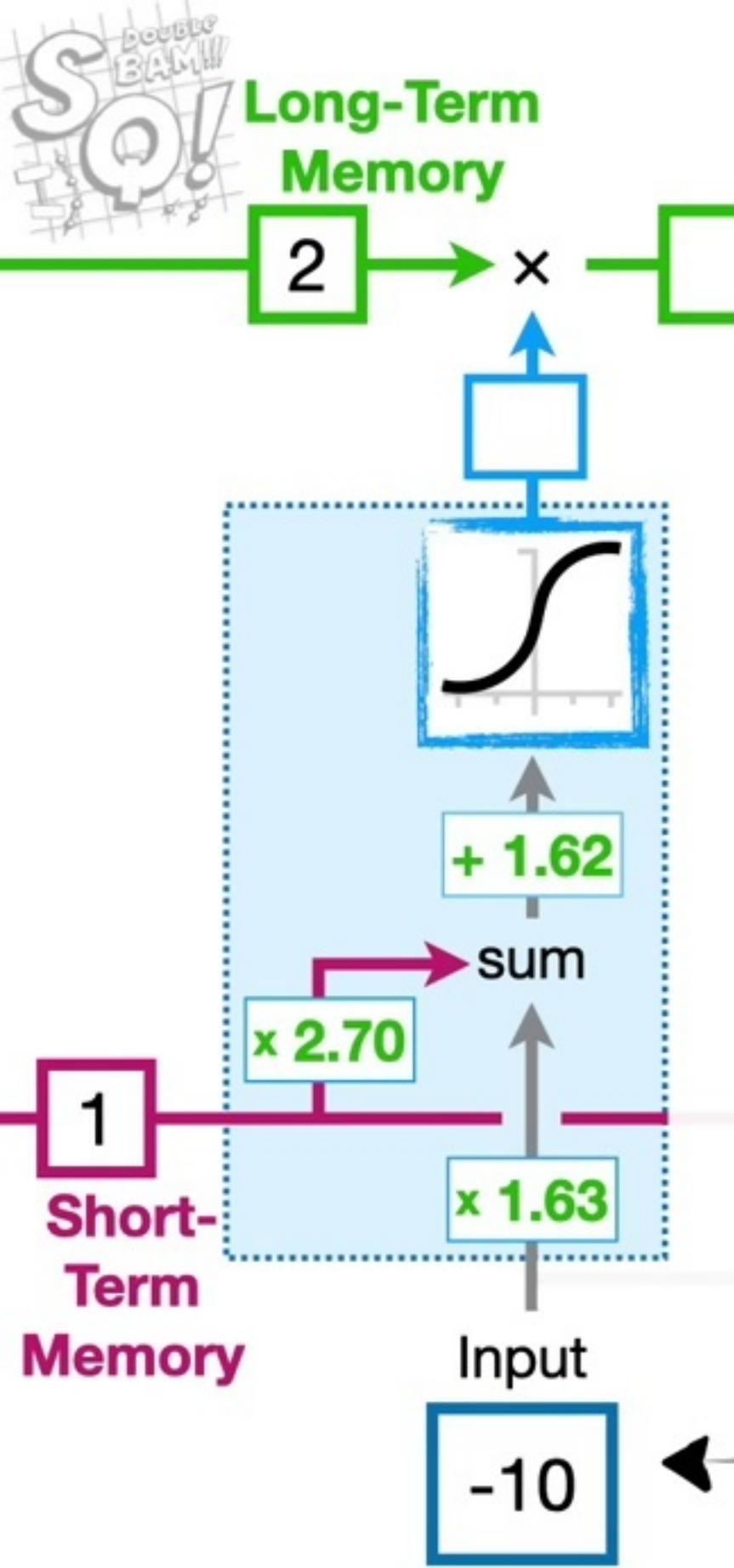




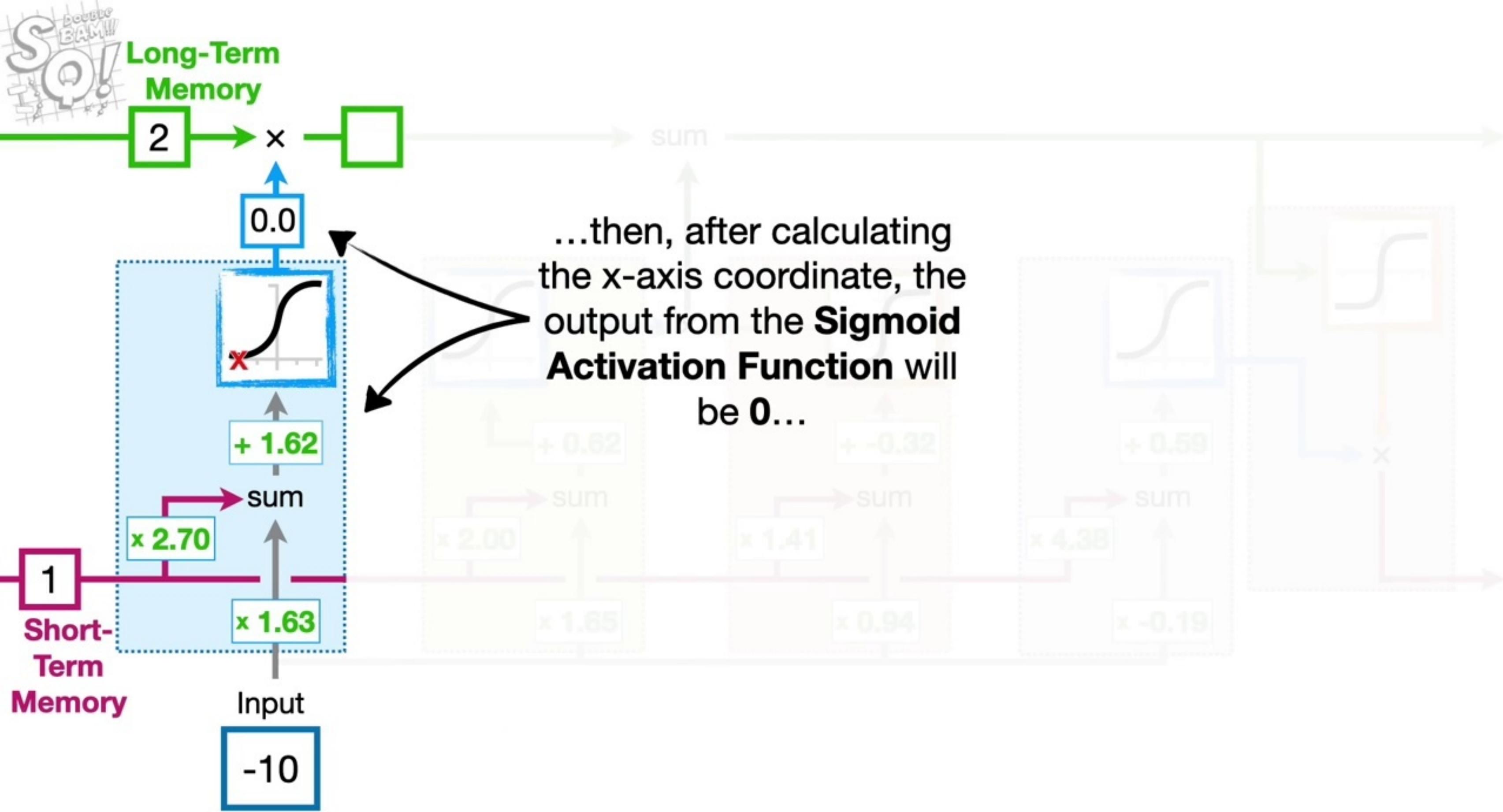


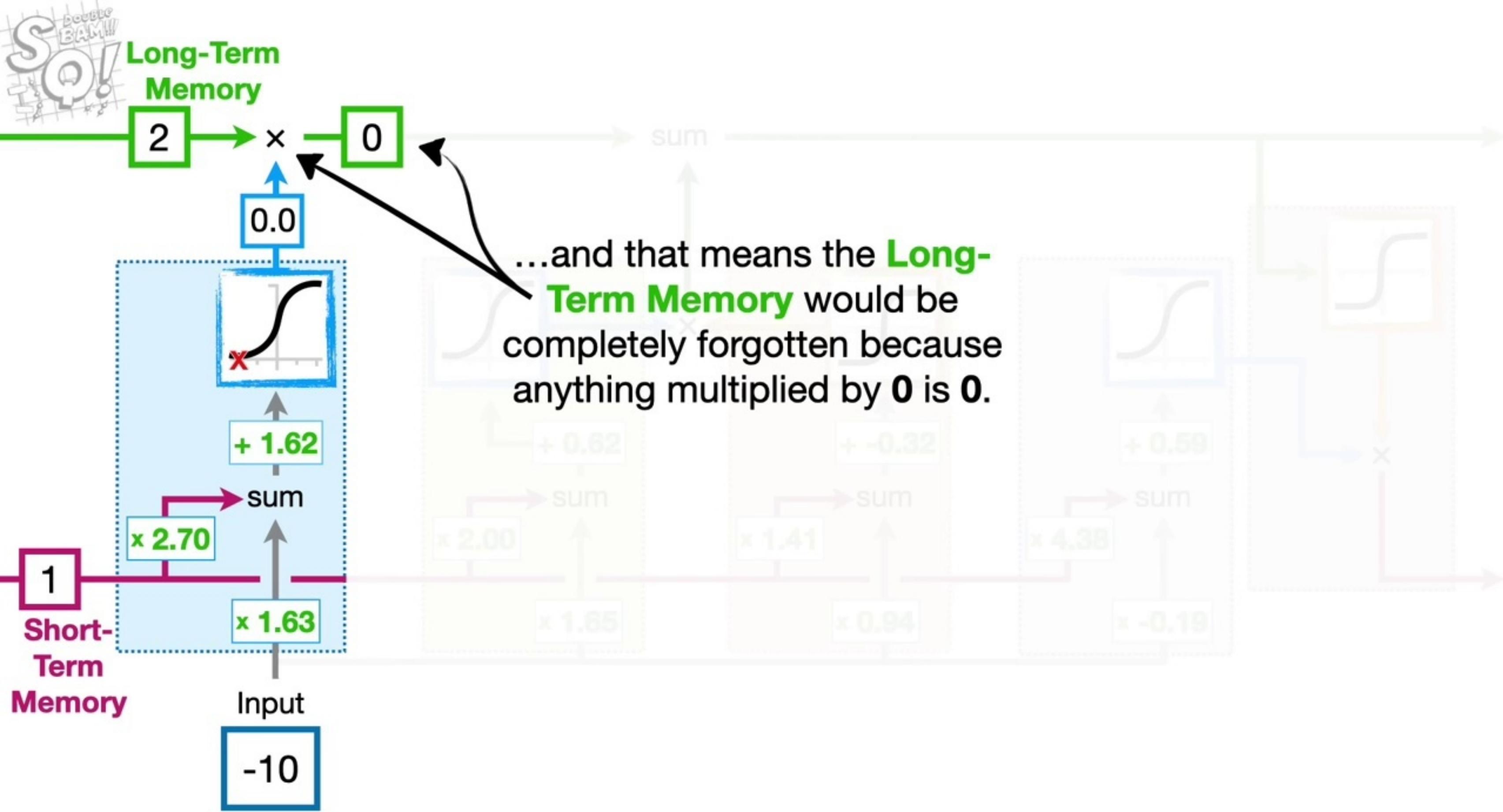


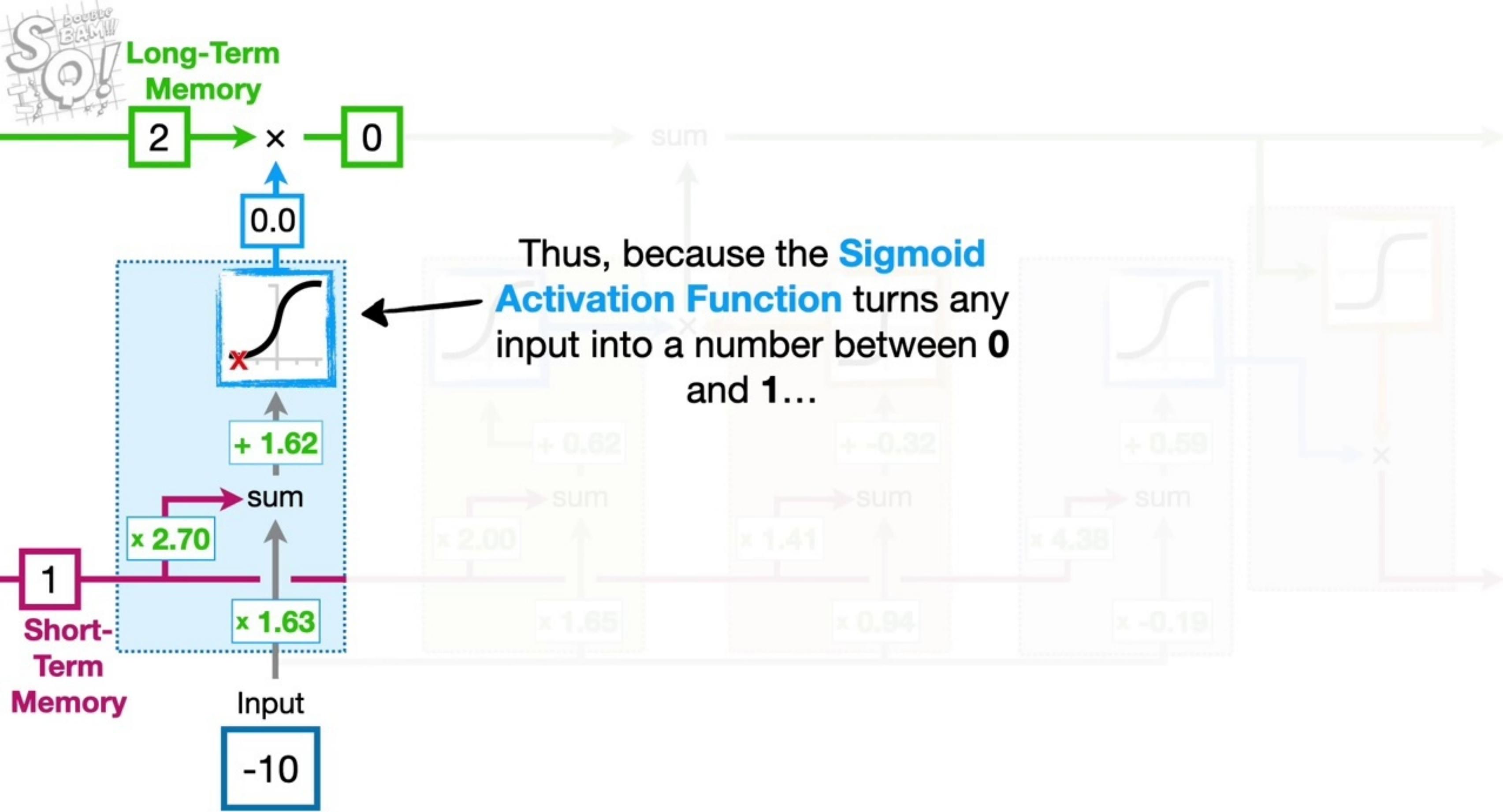


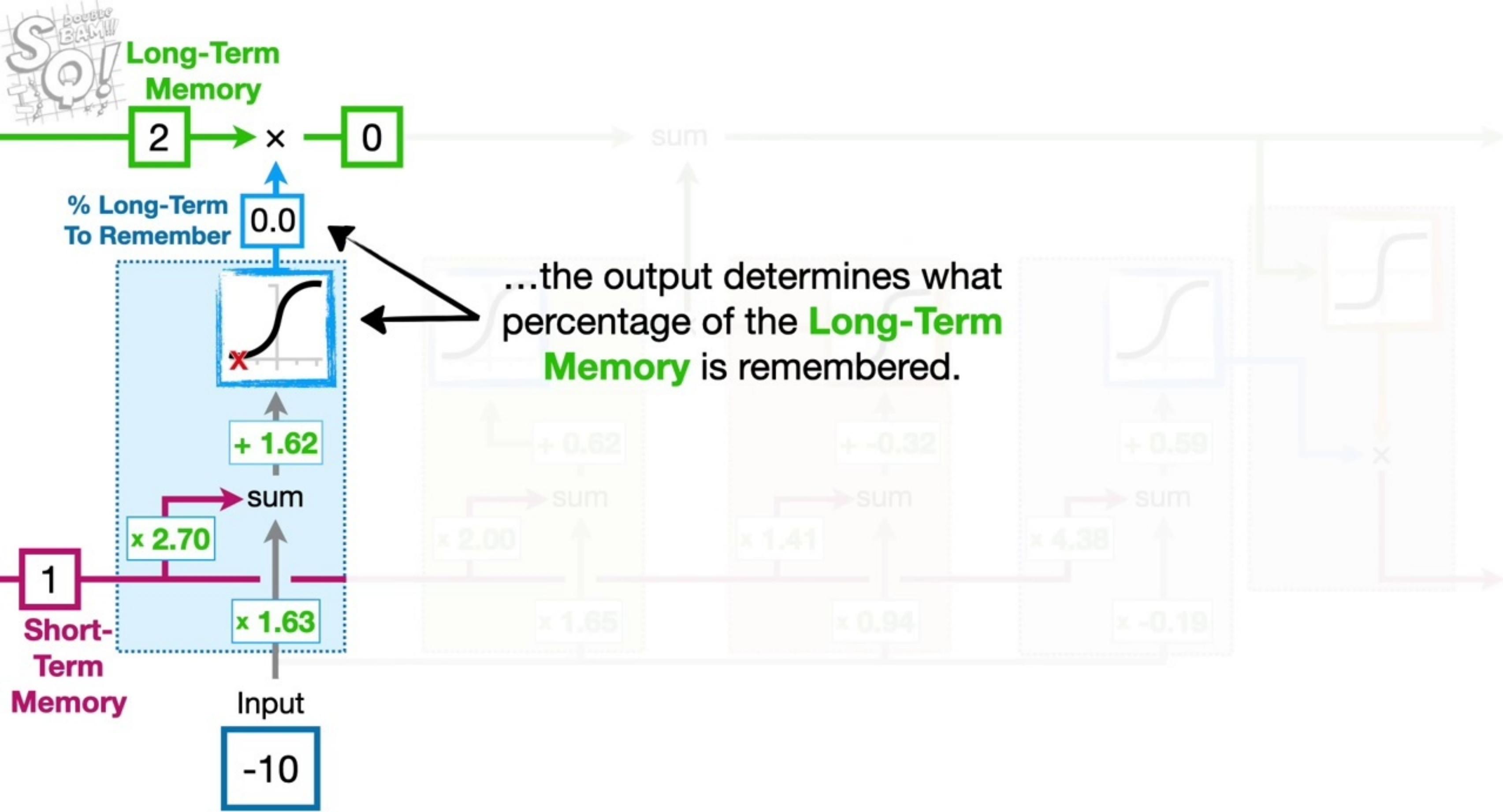


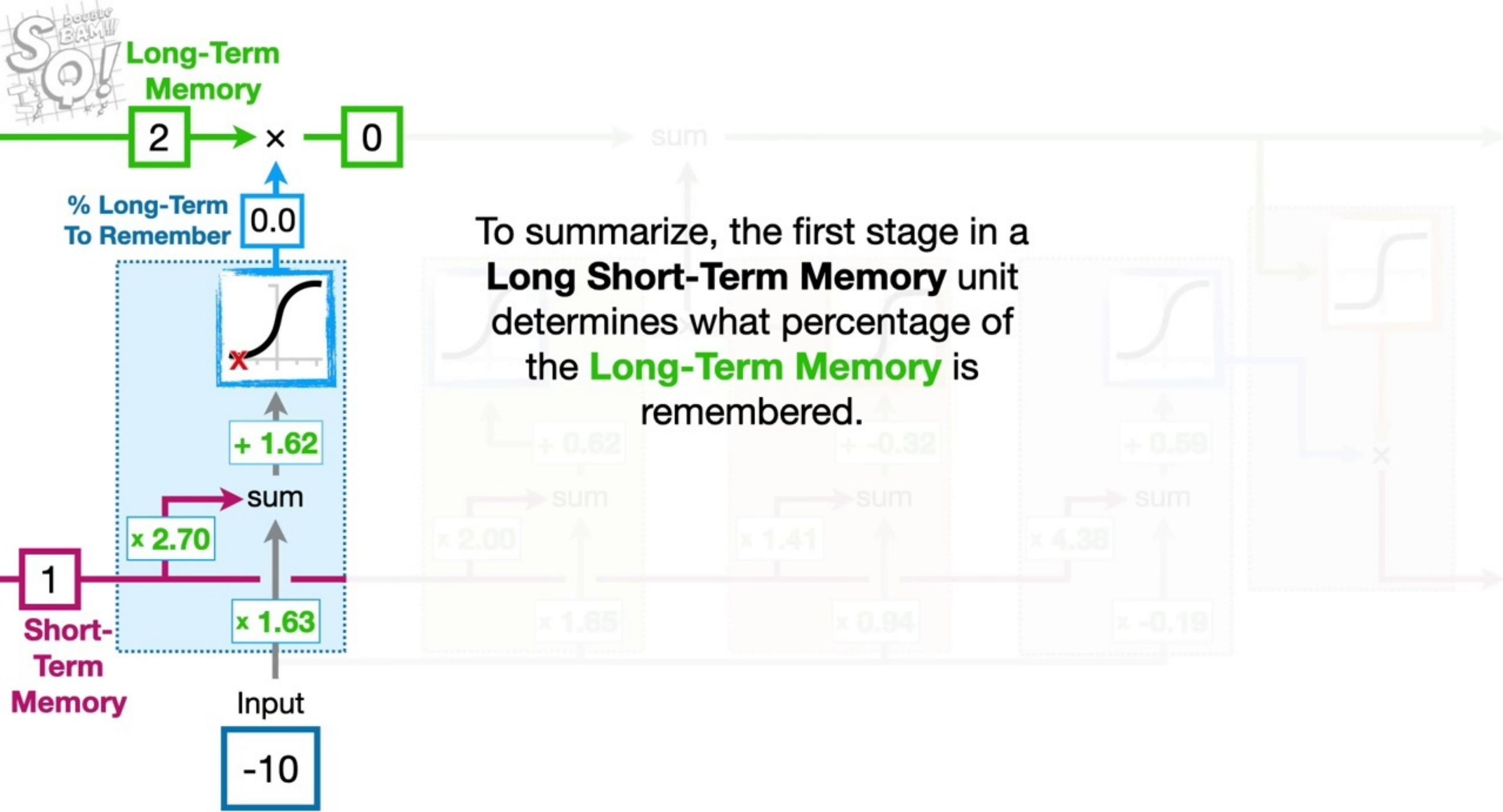
In contrast, if the **Input** to the **LSTM** was a relatively large negative number, like -10...

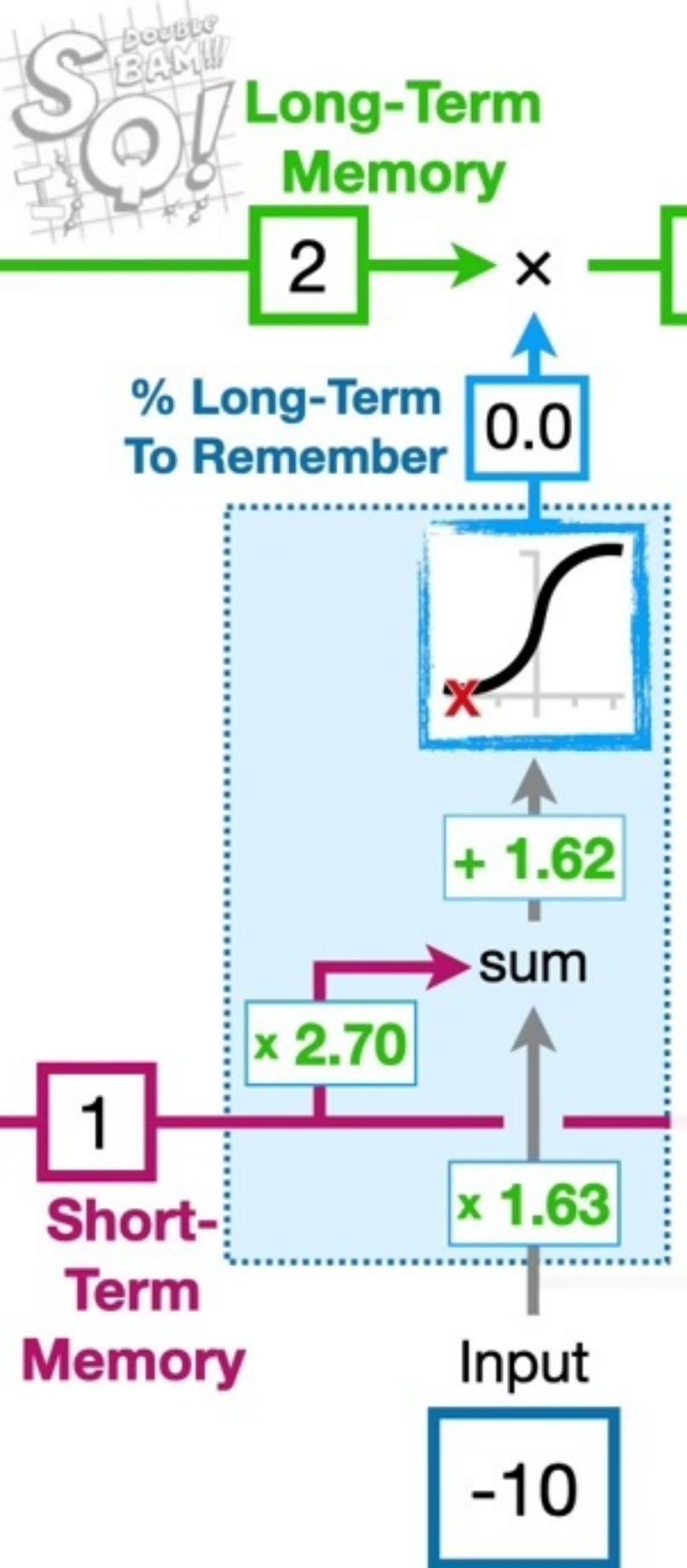








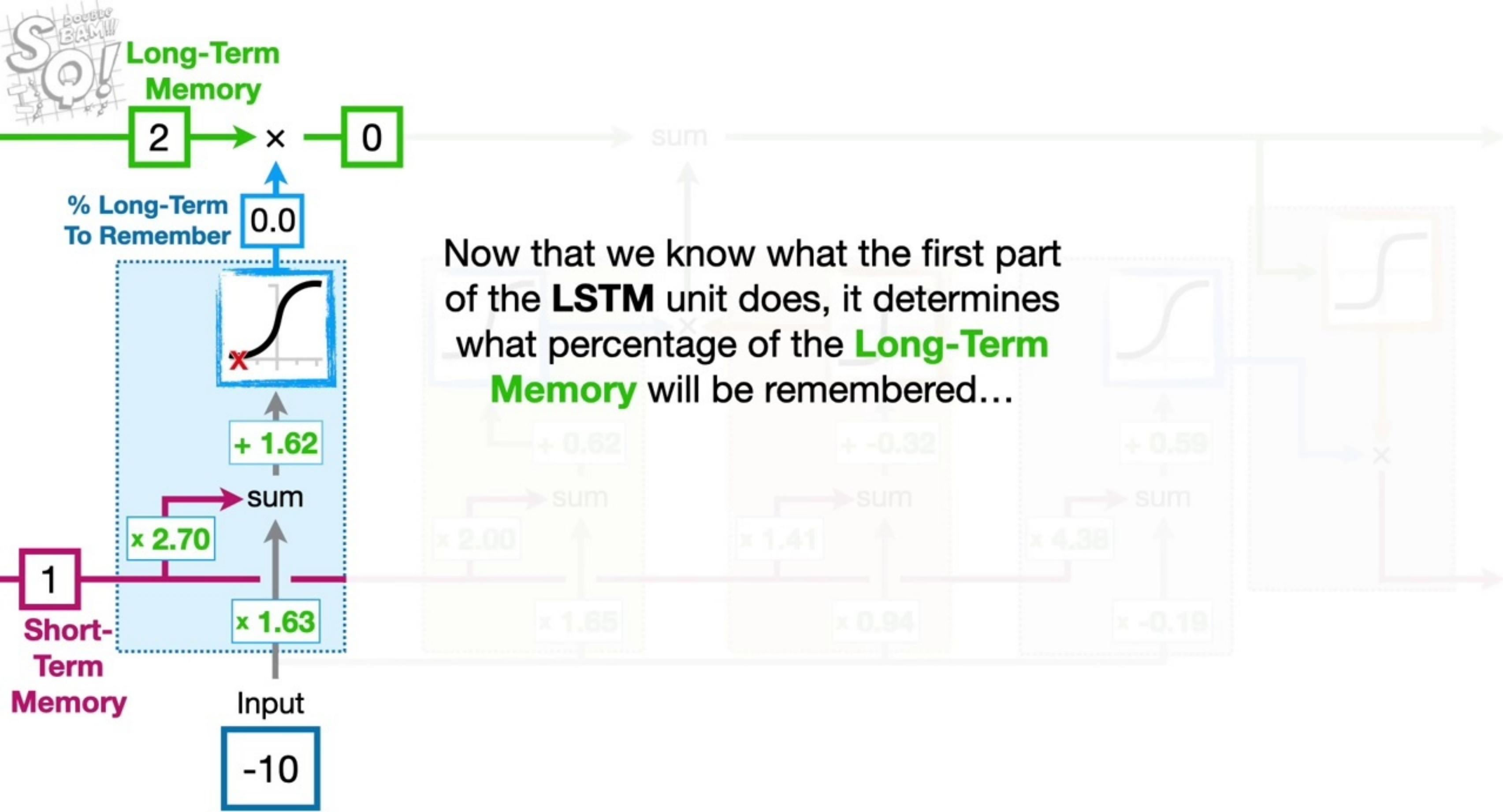


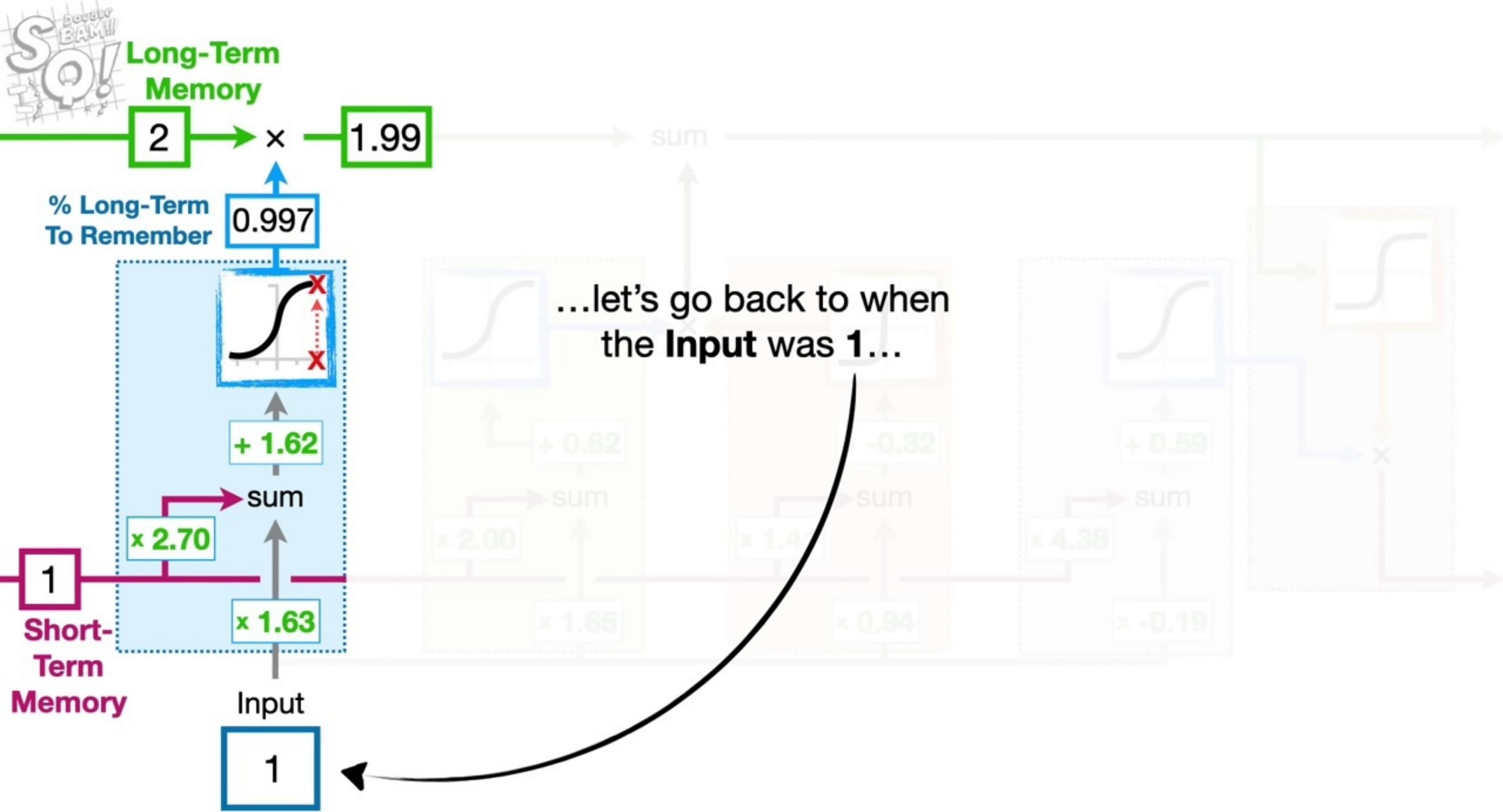


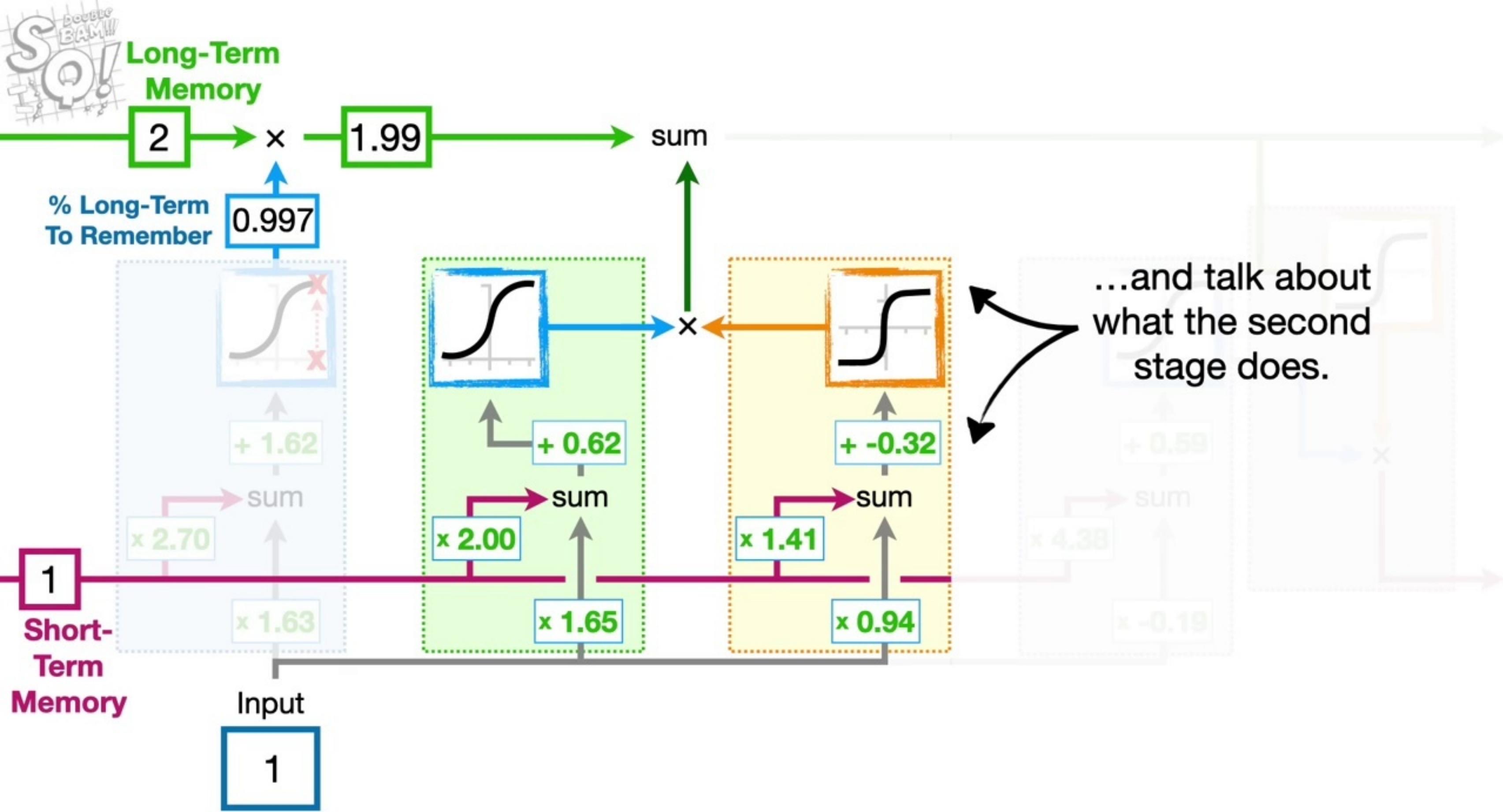
TERMINOLOGY ALERT!!!

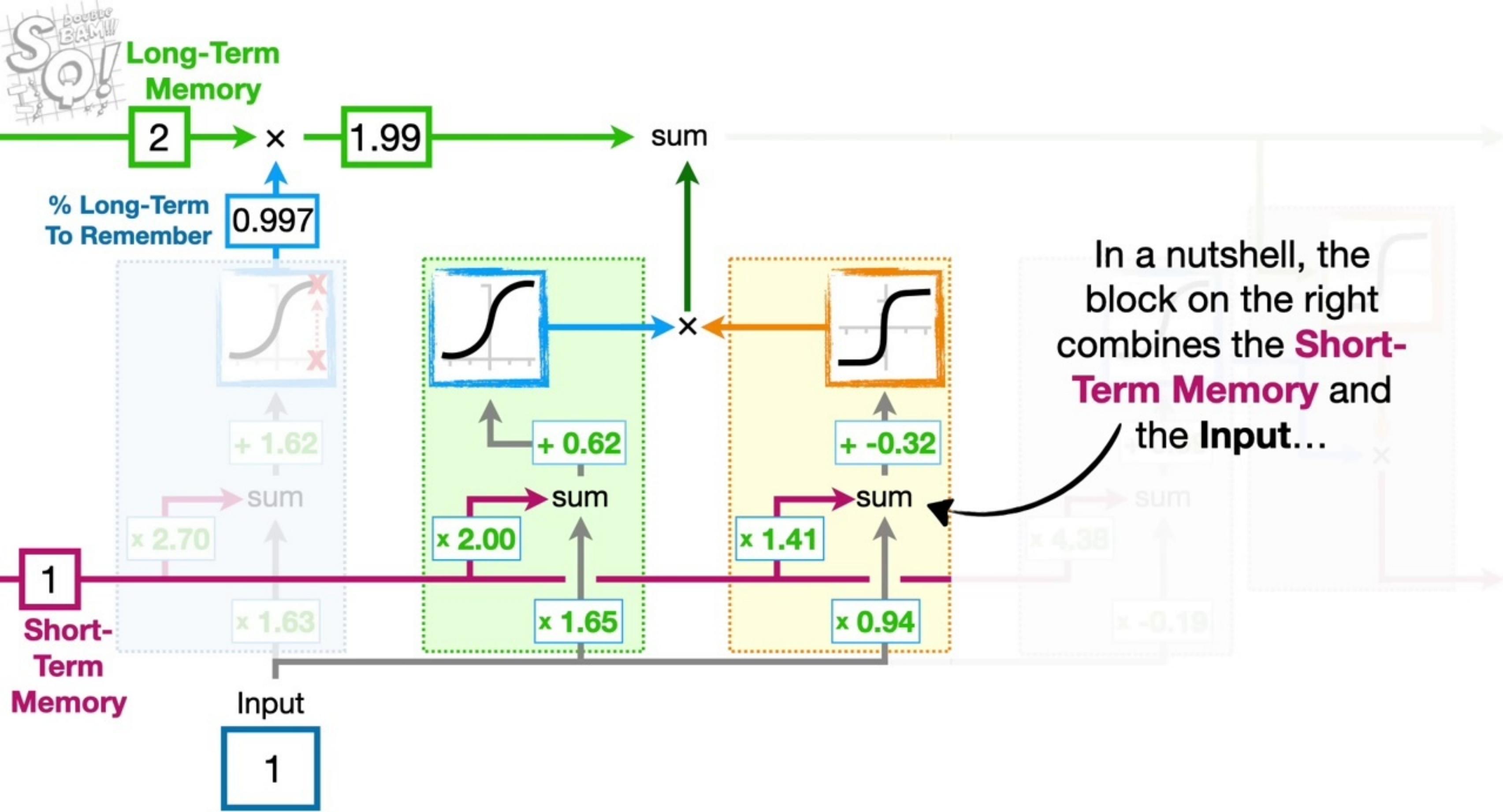
Even though this part of the **Long Short-Term Memory** unit determines what percentage of the **Long-Term Memory** will be remembered...

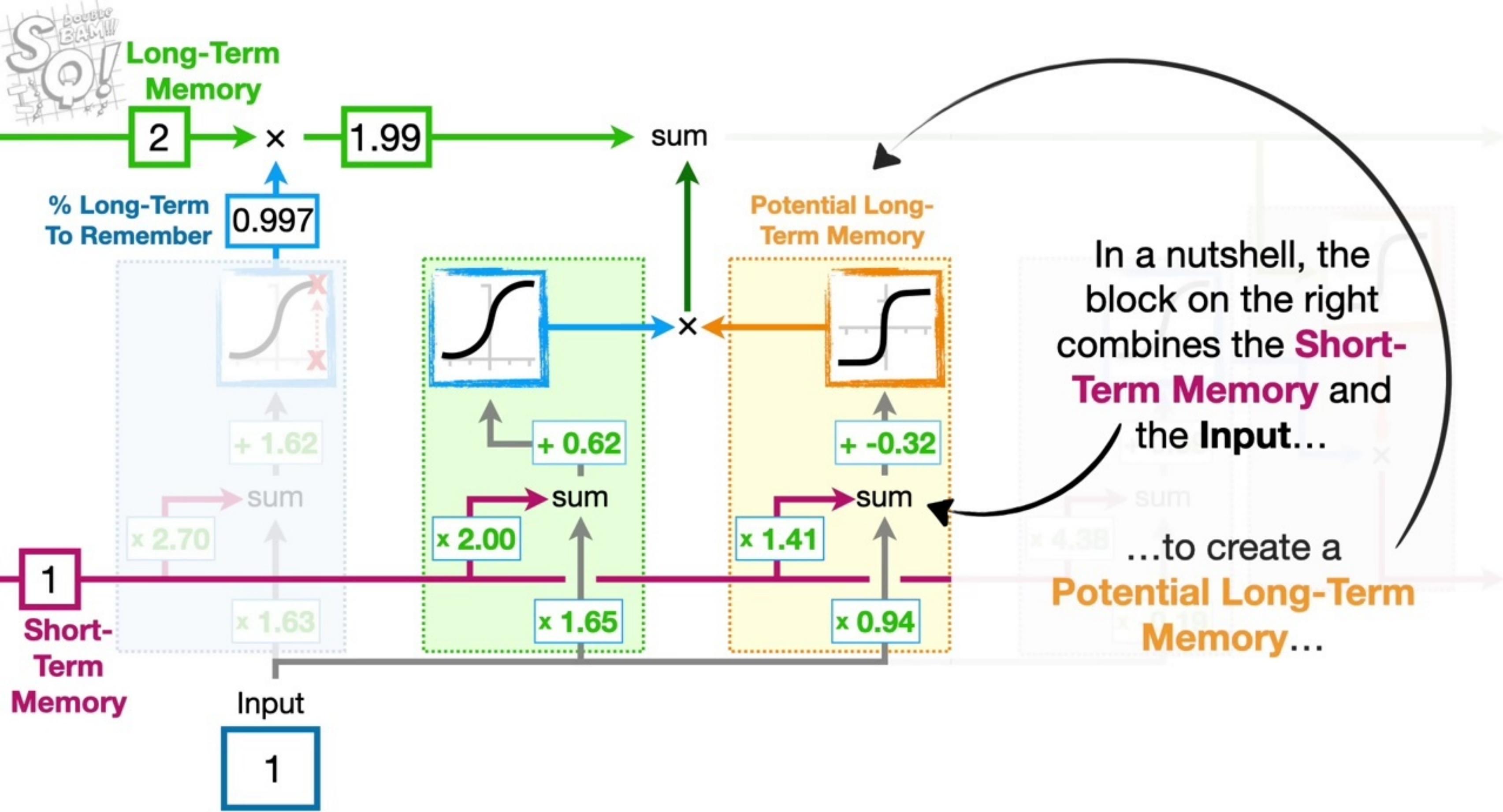
...it is usually called the **Forget Gate**.

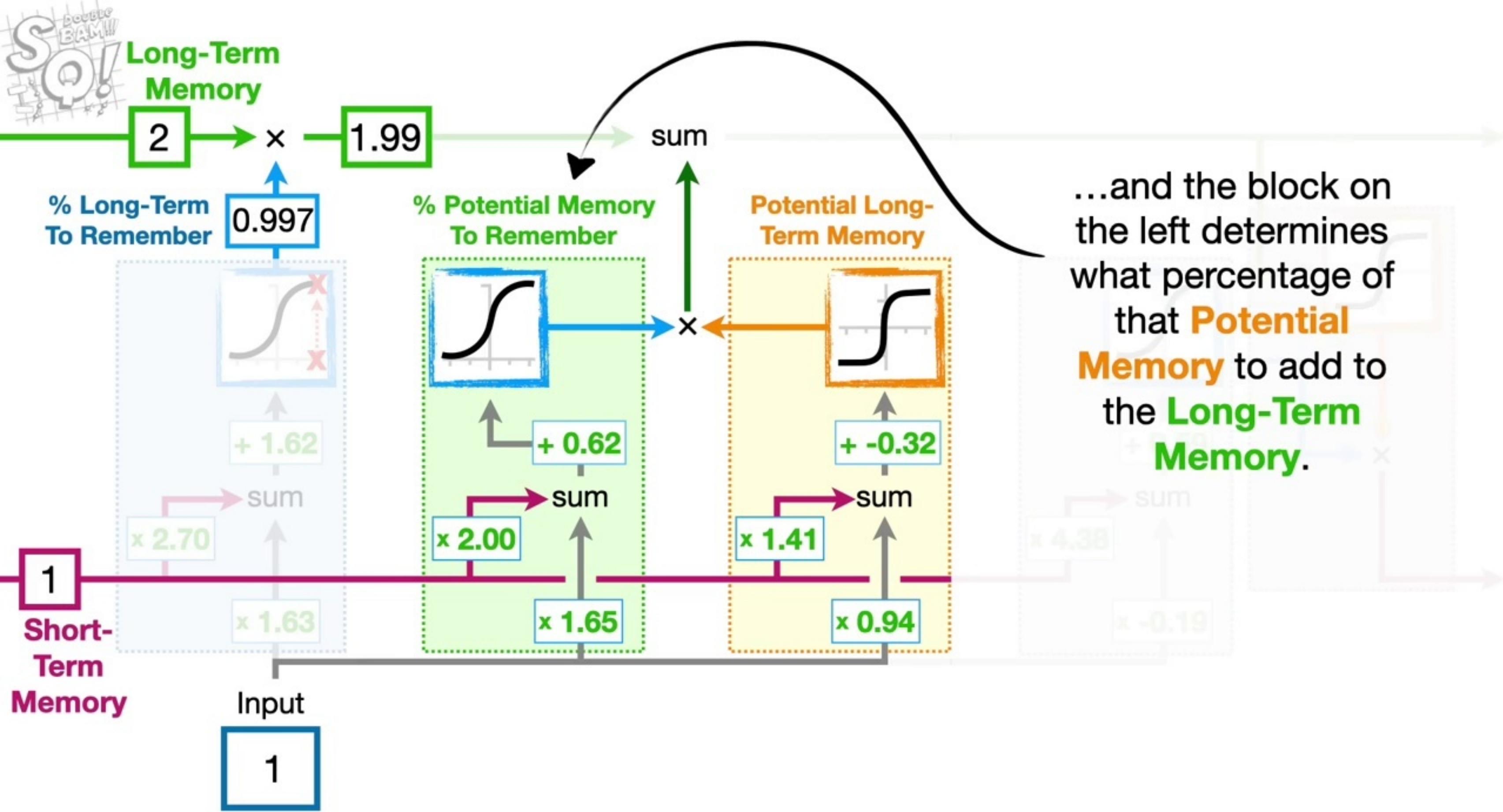


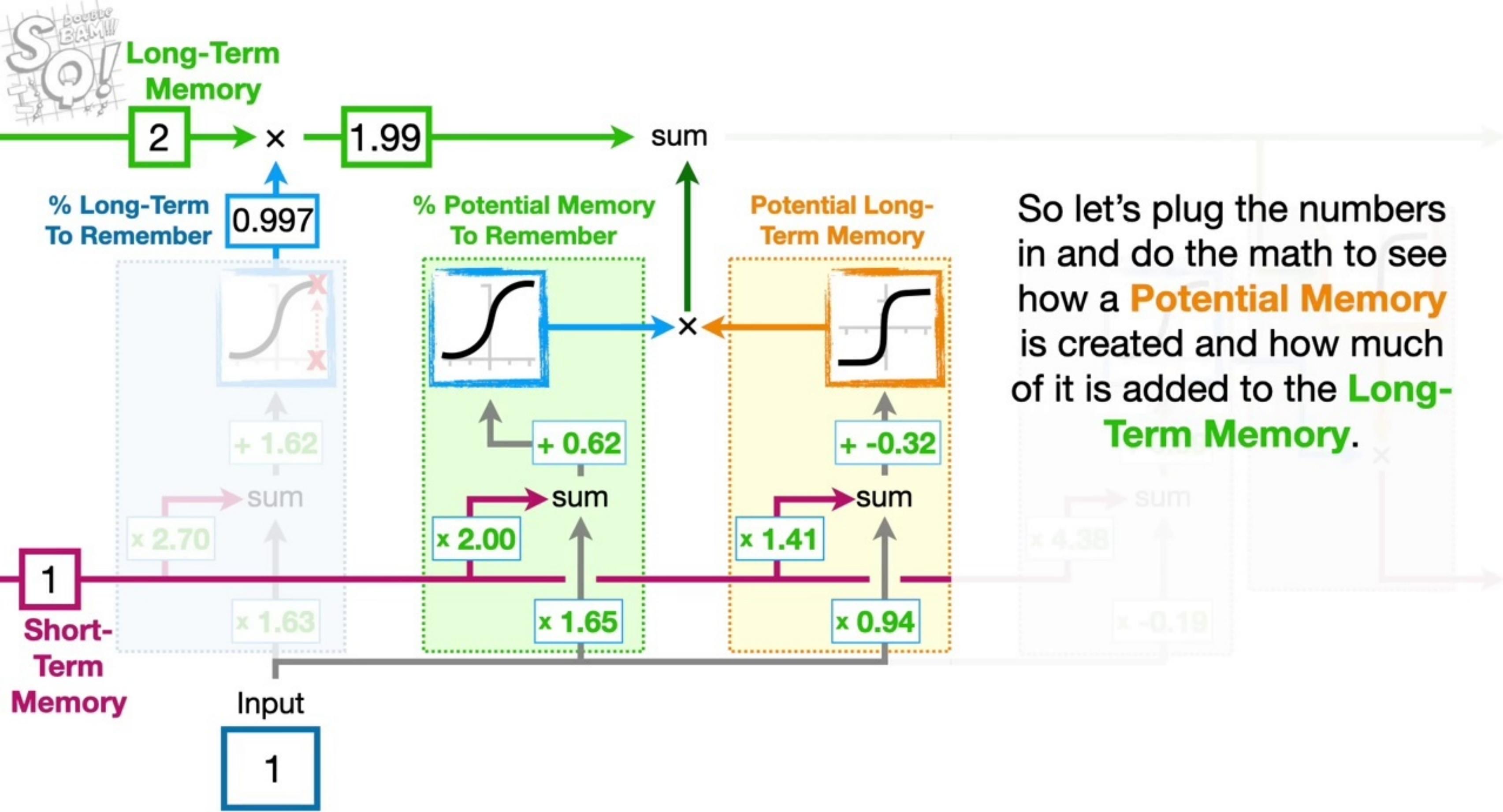




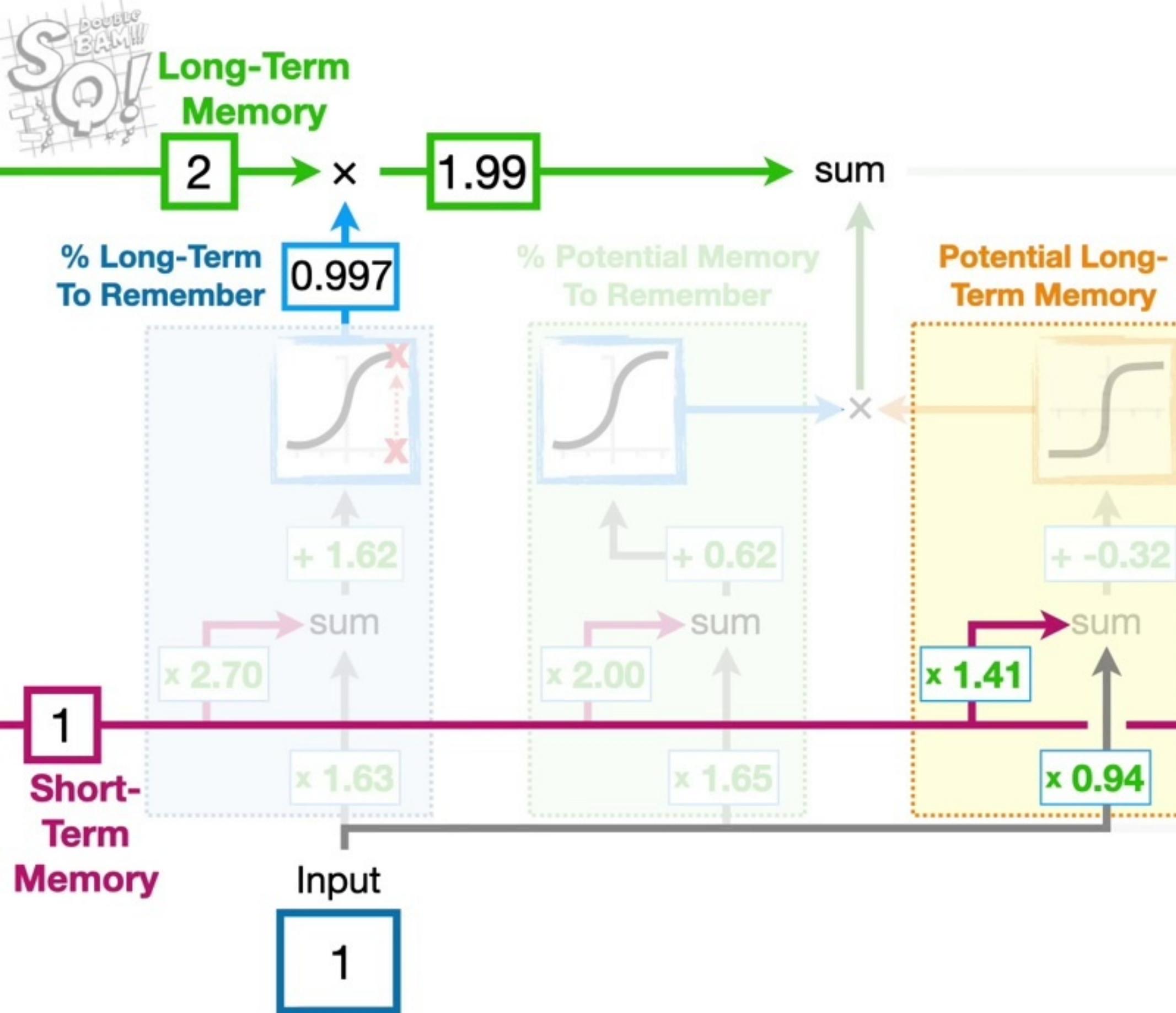




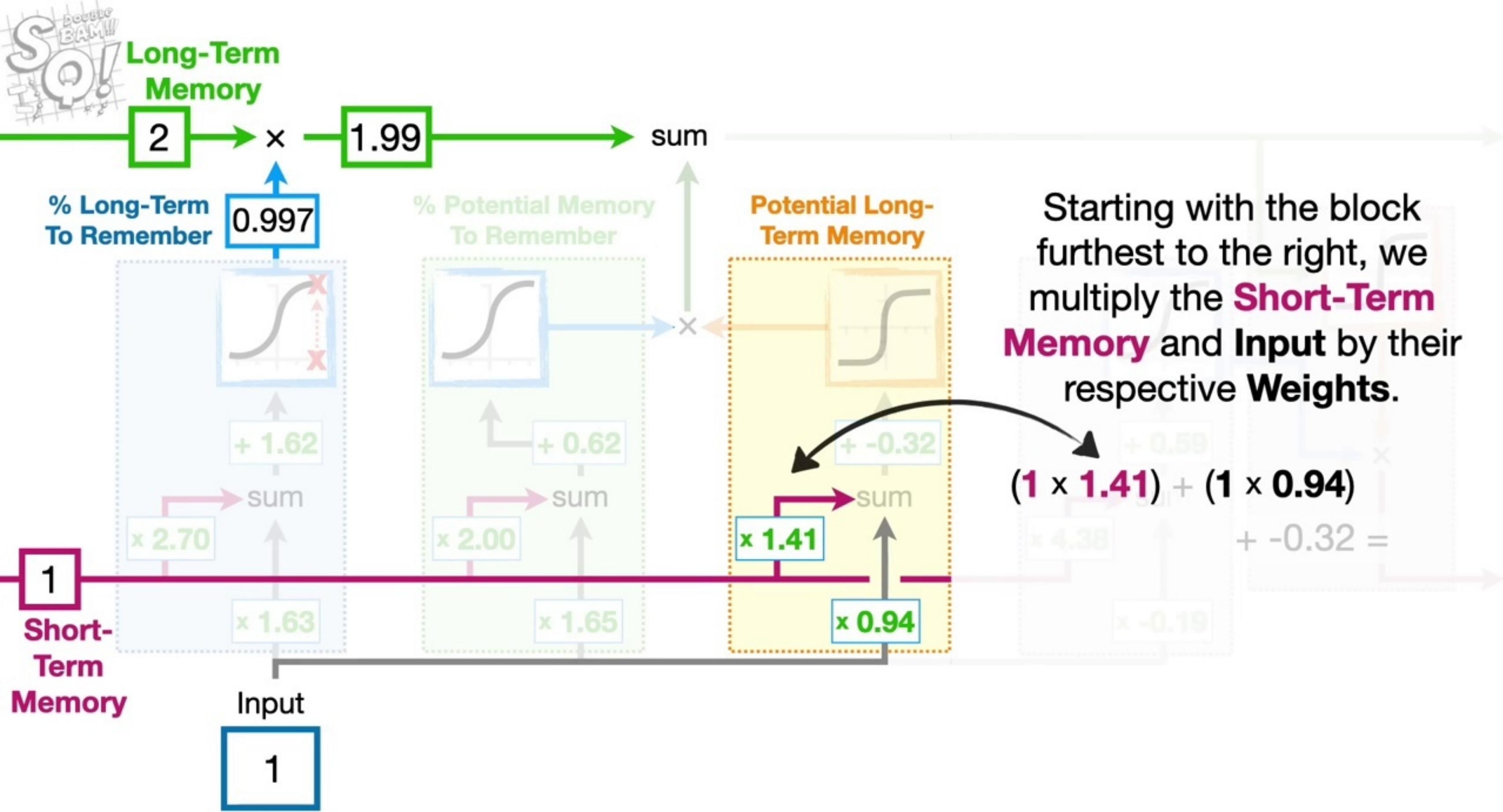


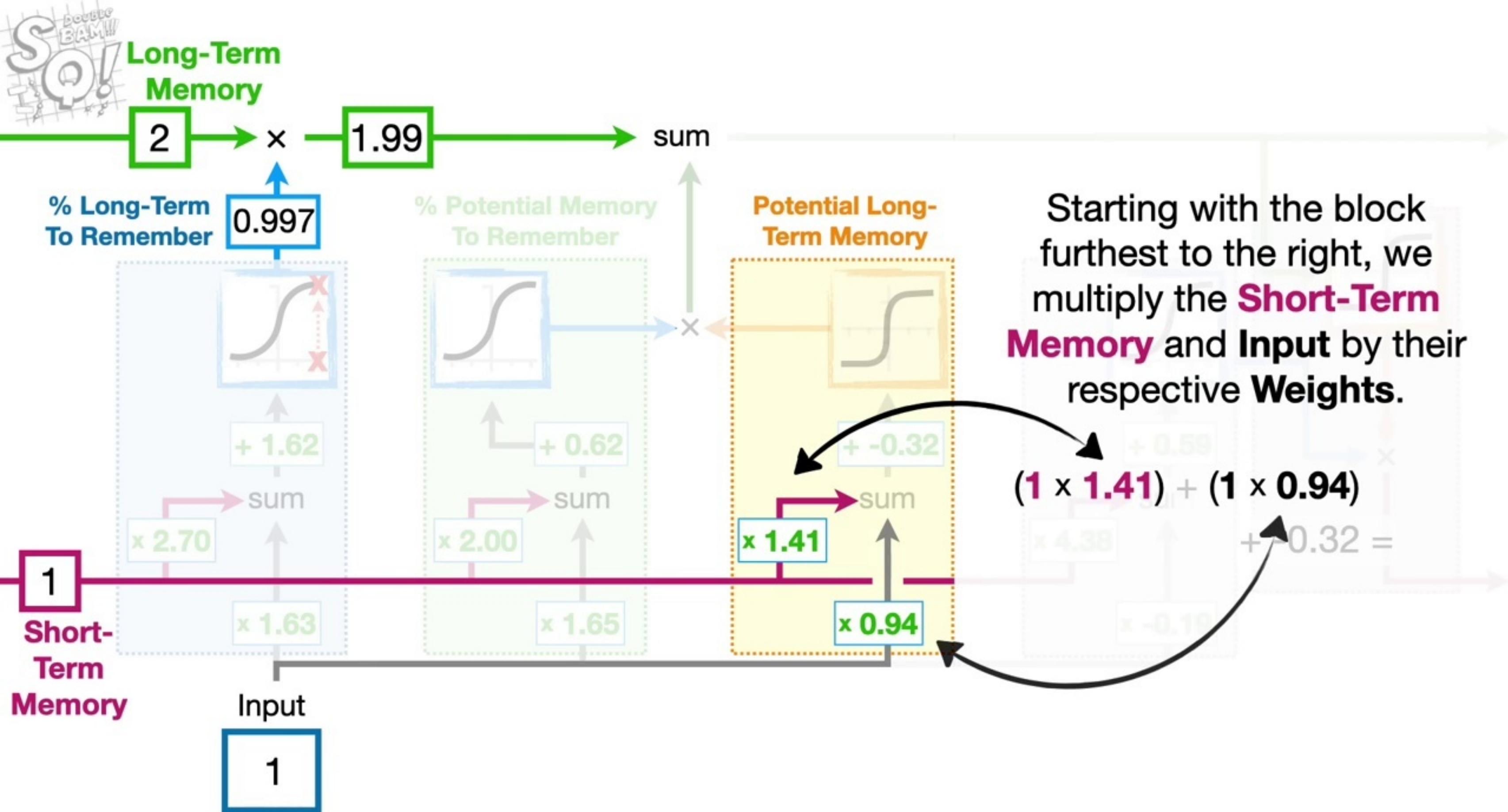


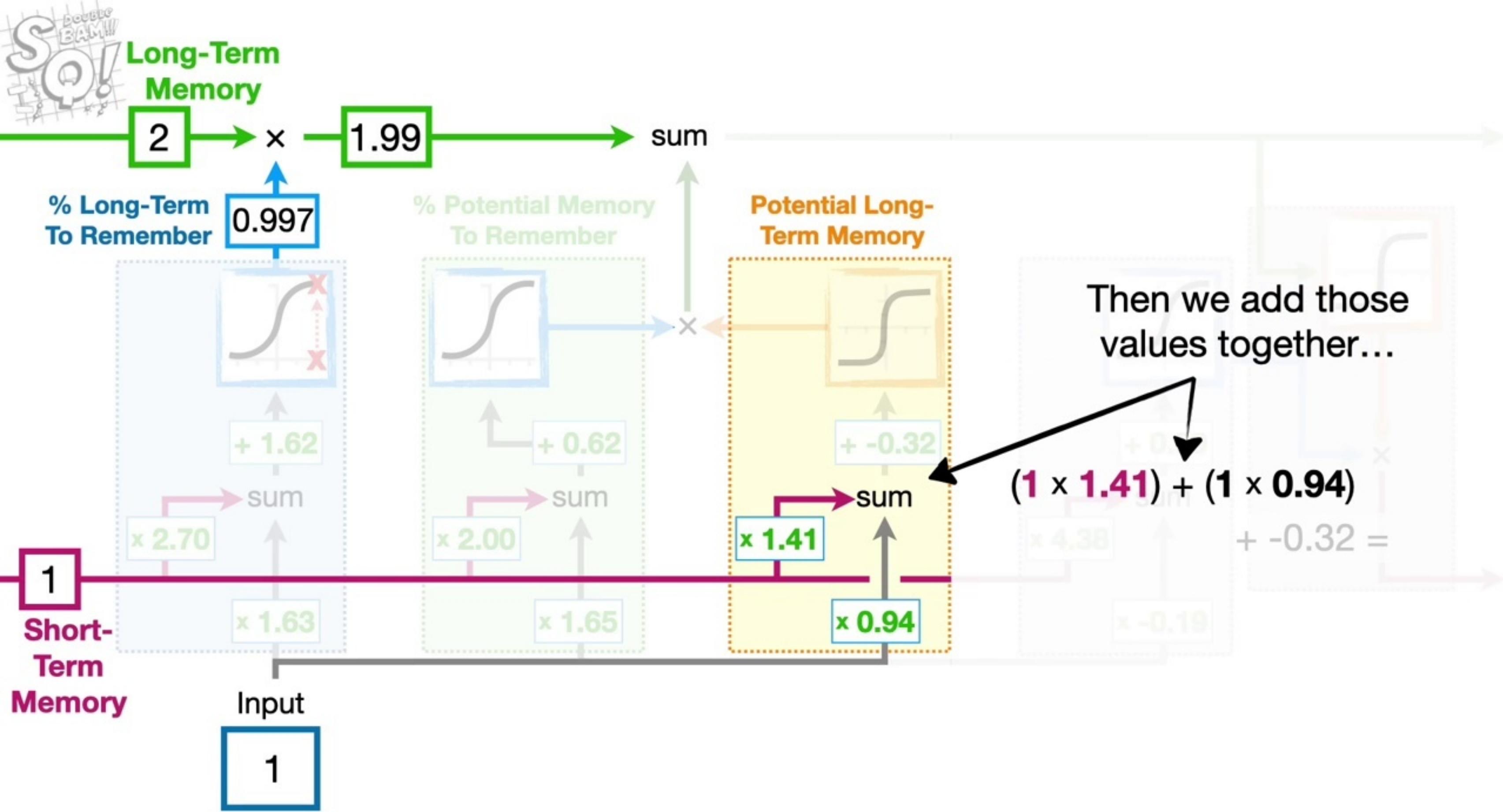
So let's plug the numbers in and do the math to see how a **Potential Memory** is created and how much of it is added to the **Long-Term Memory**.

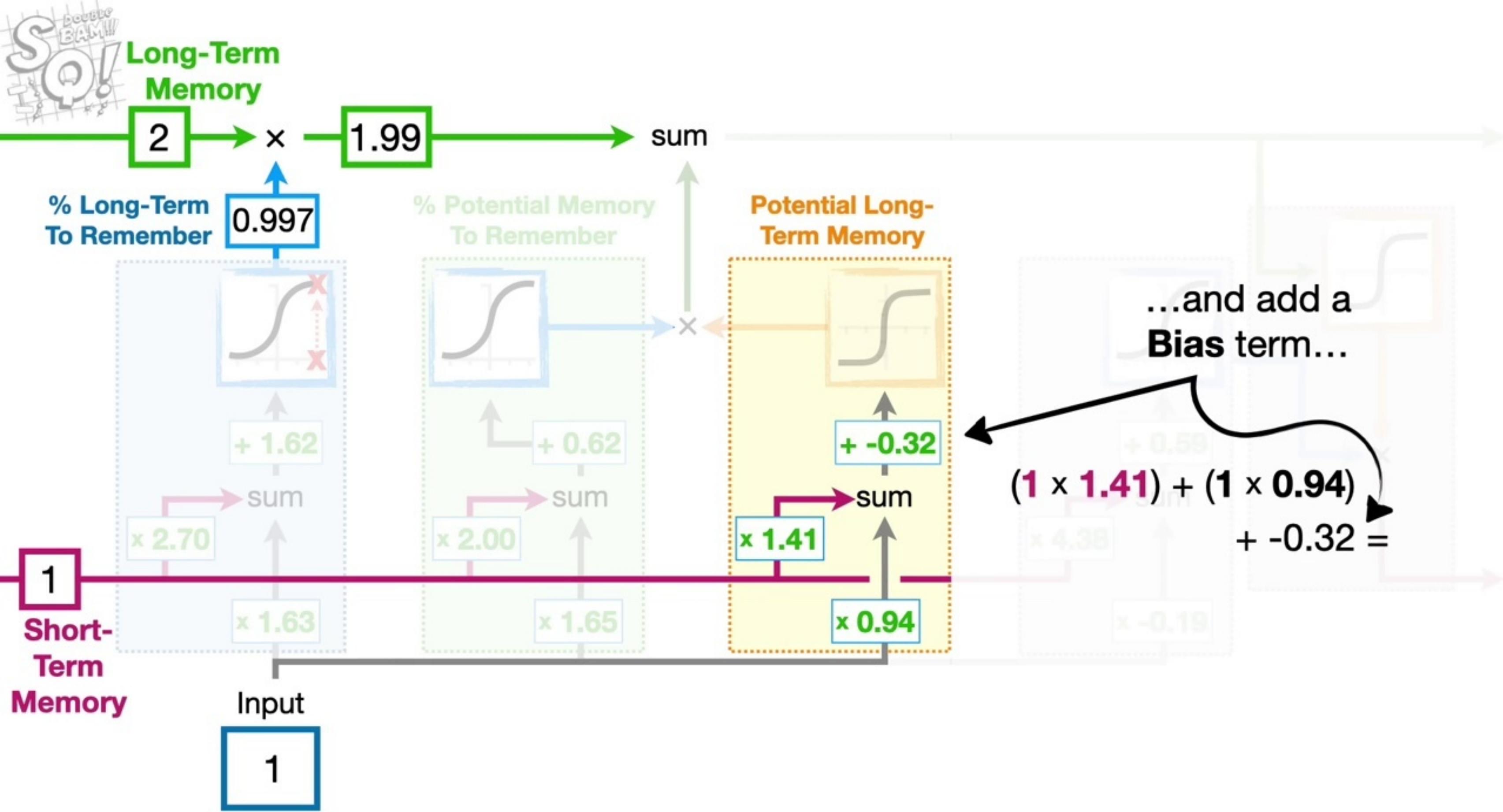


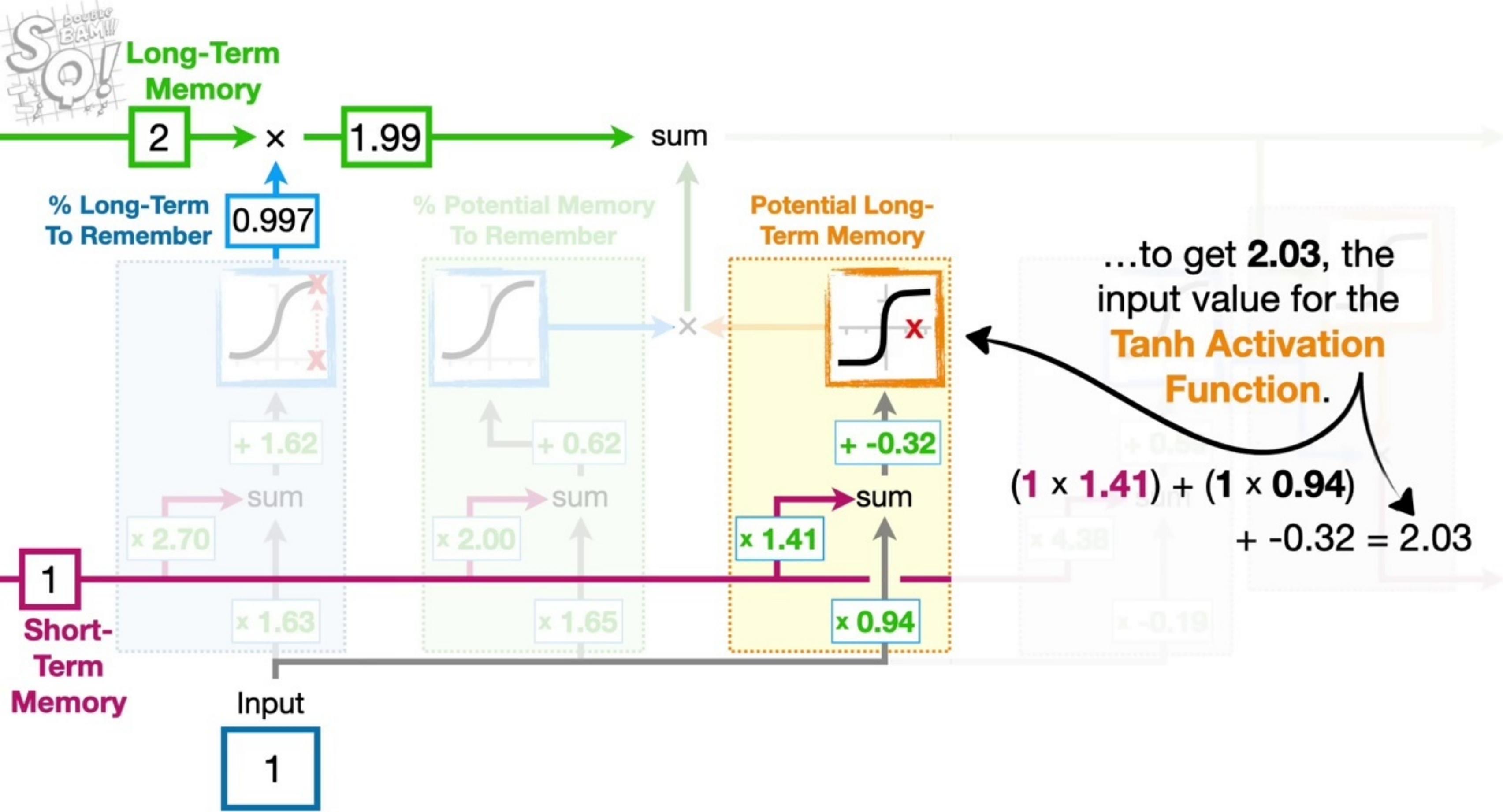
Starting with the block furthest to the right, we multiply the **Short-Term Memory** and **Input** by their respective **Weights**.

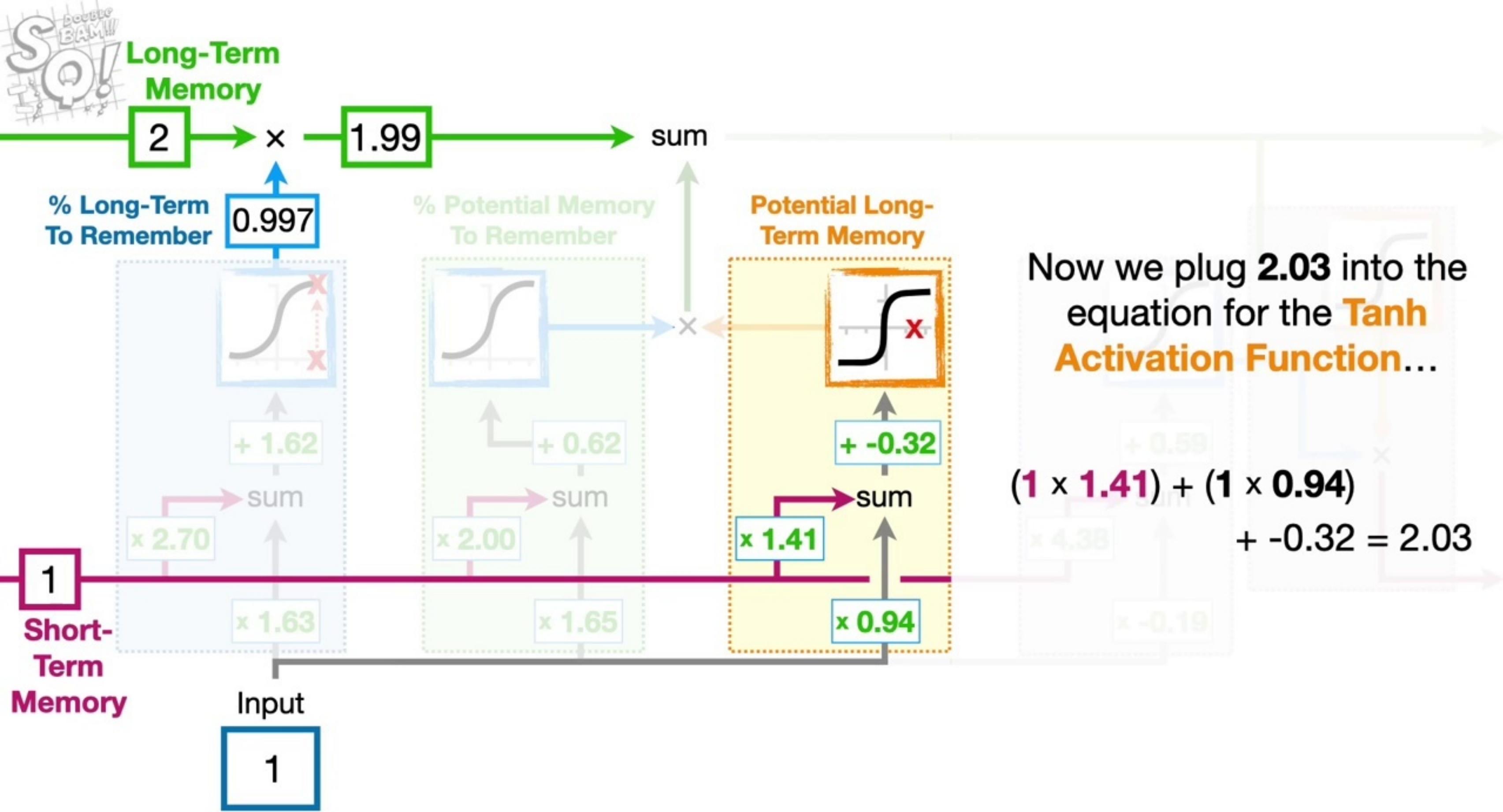


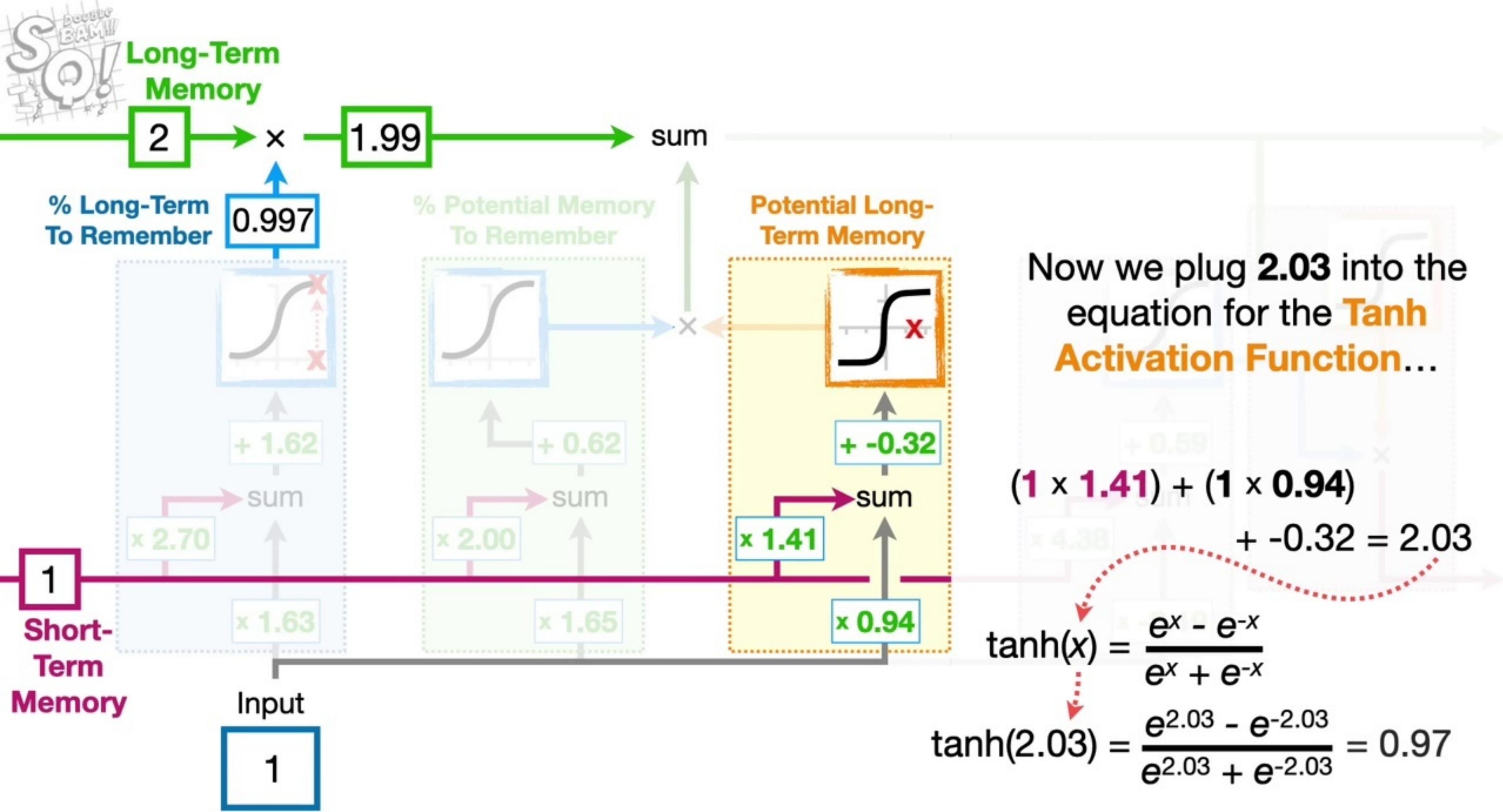


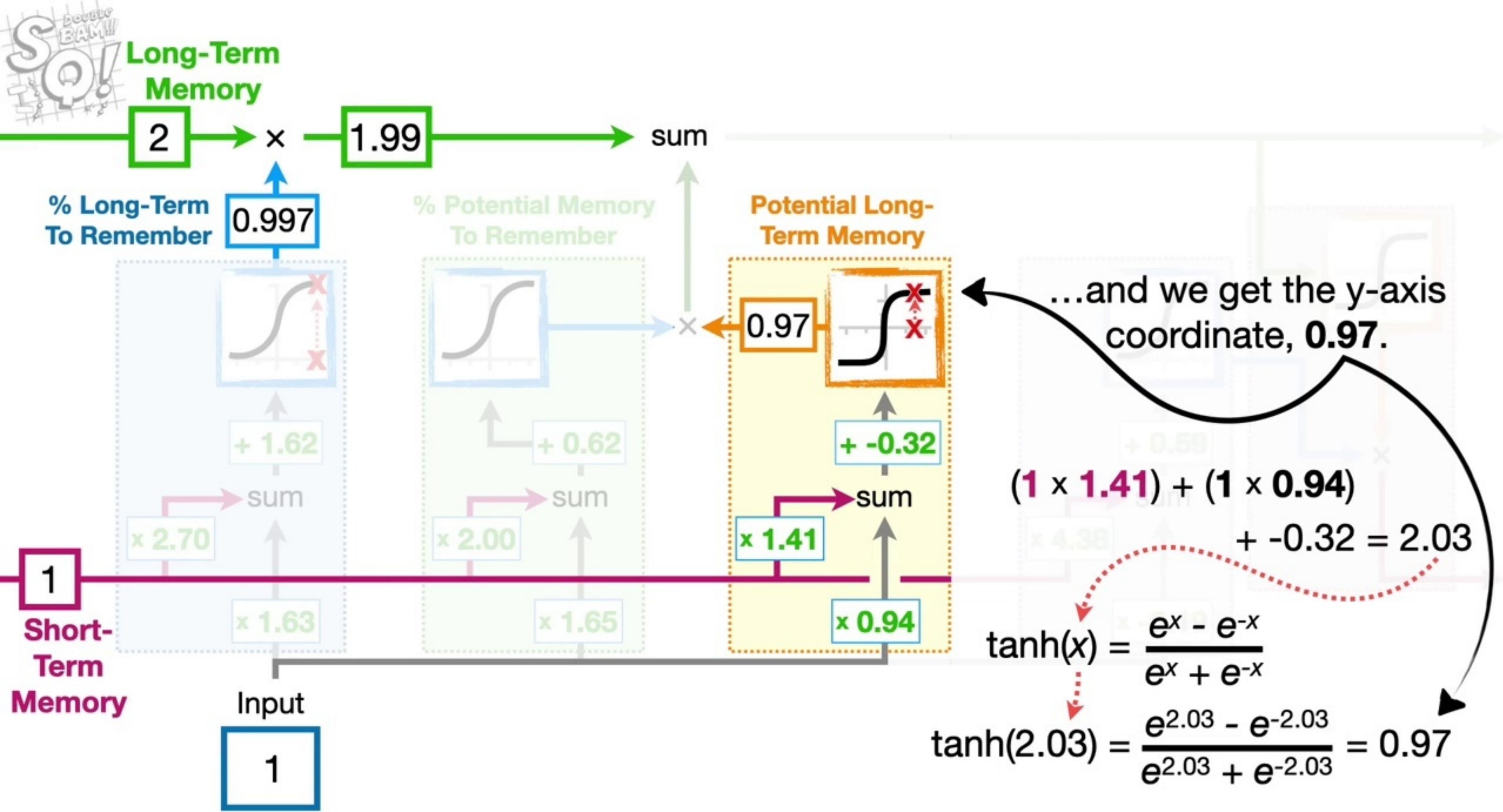


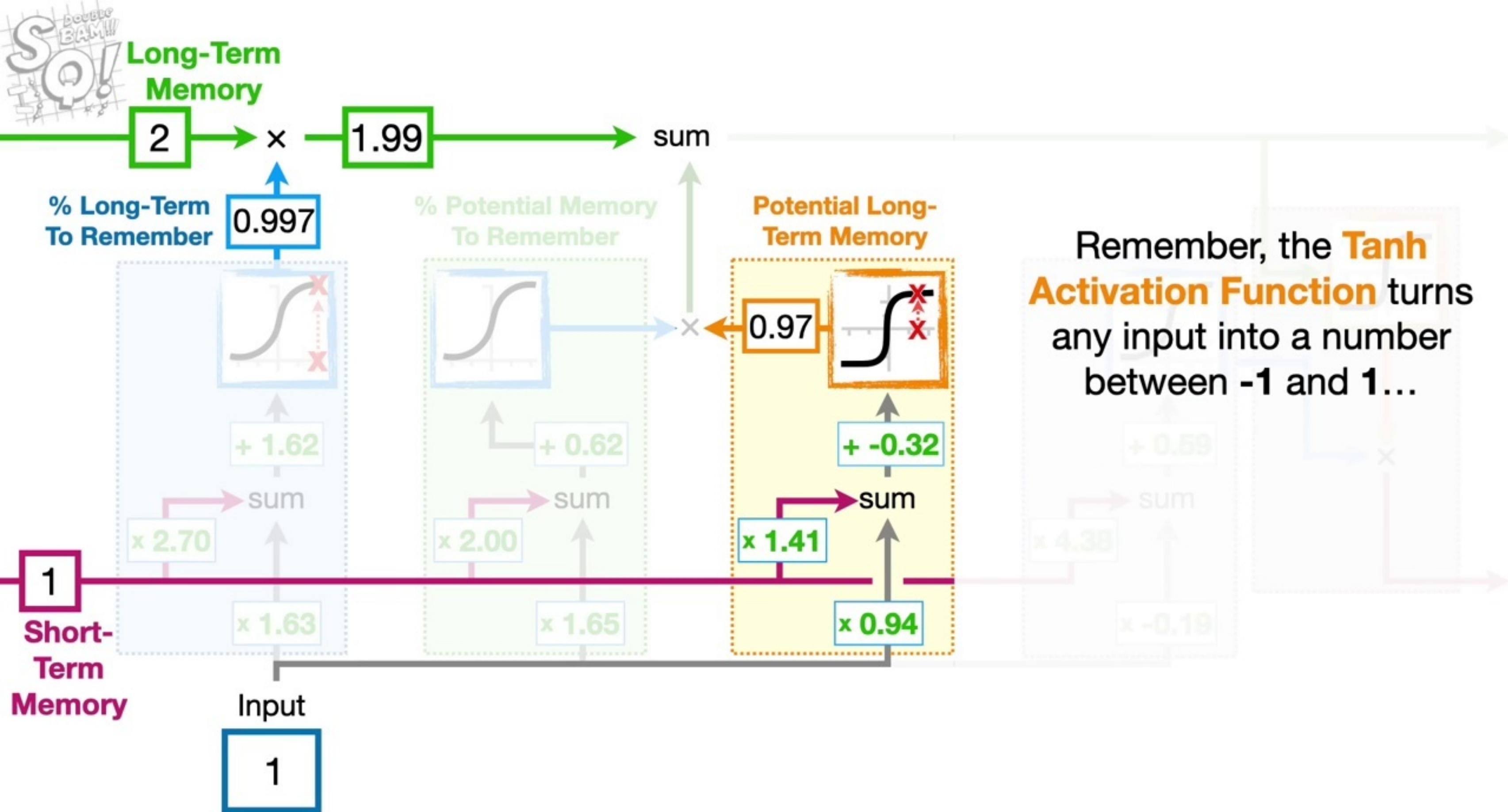


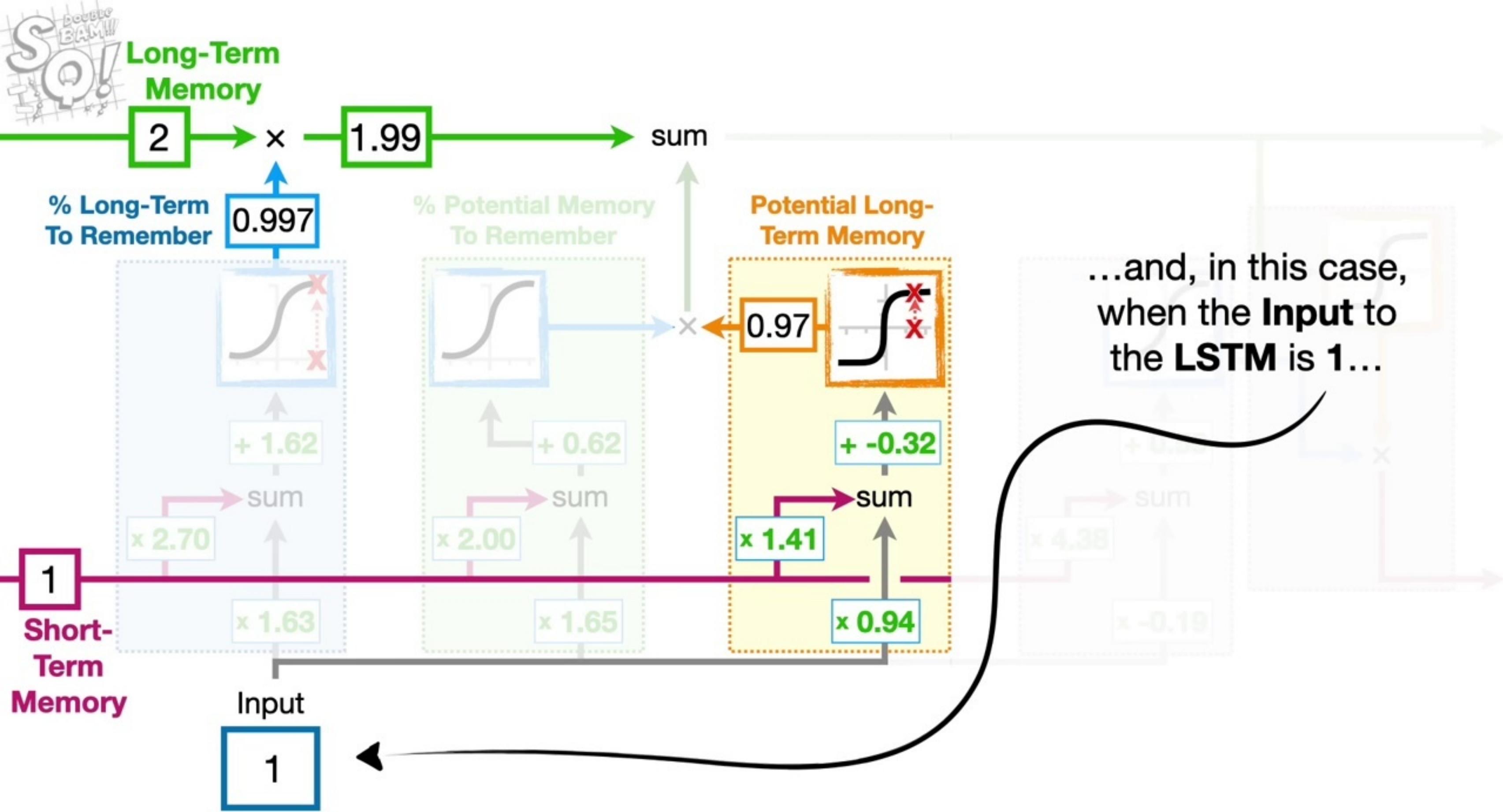


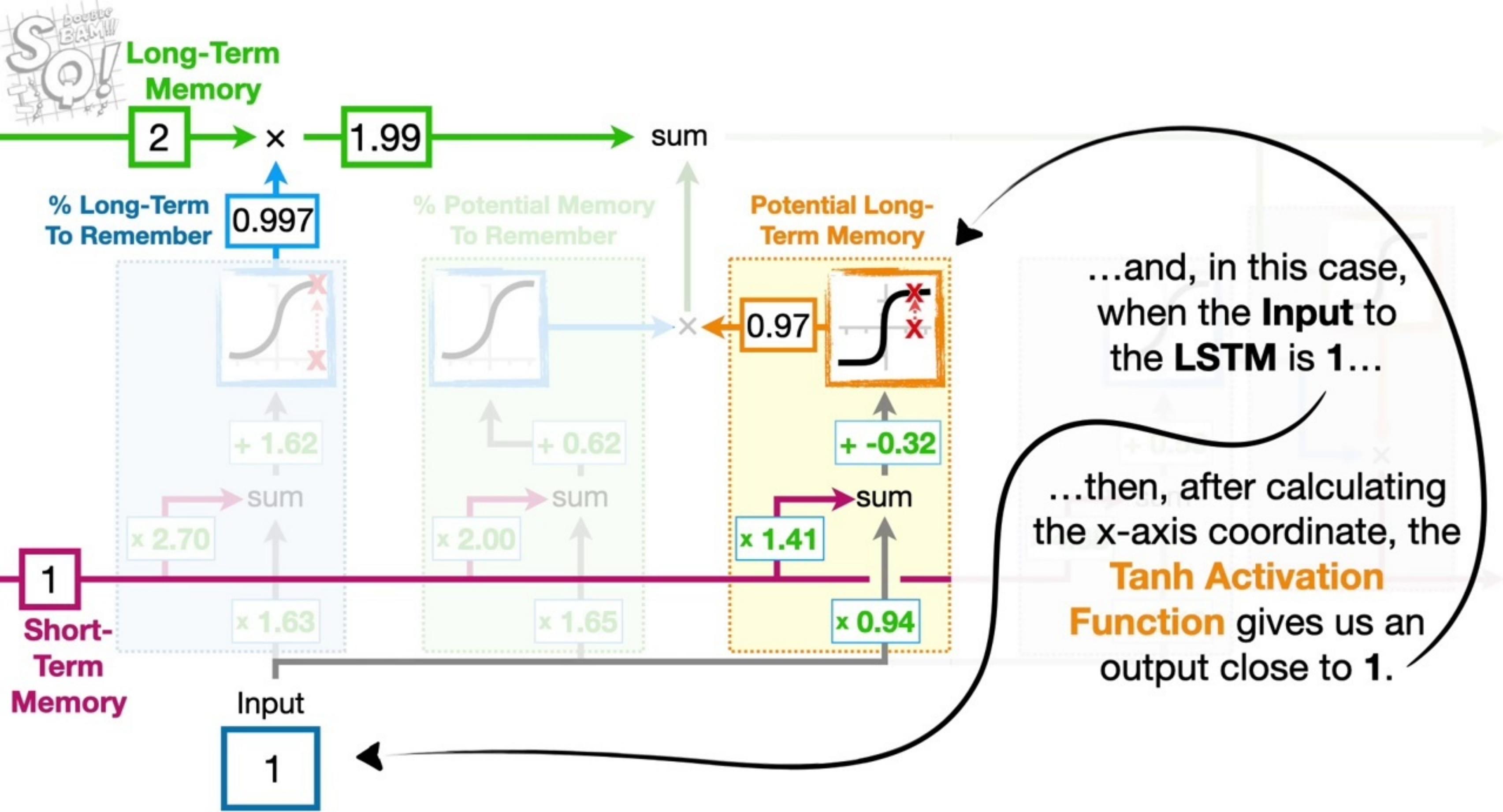


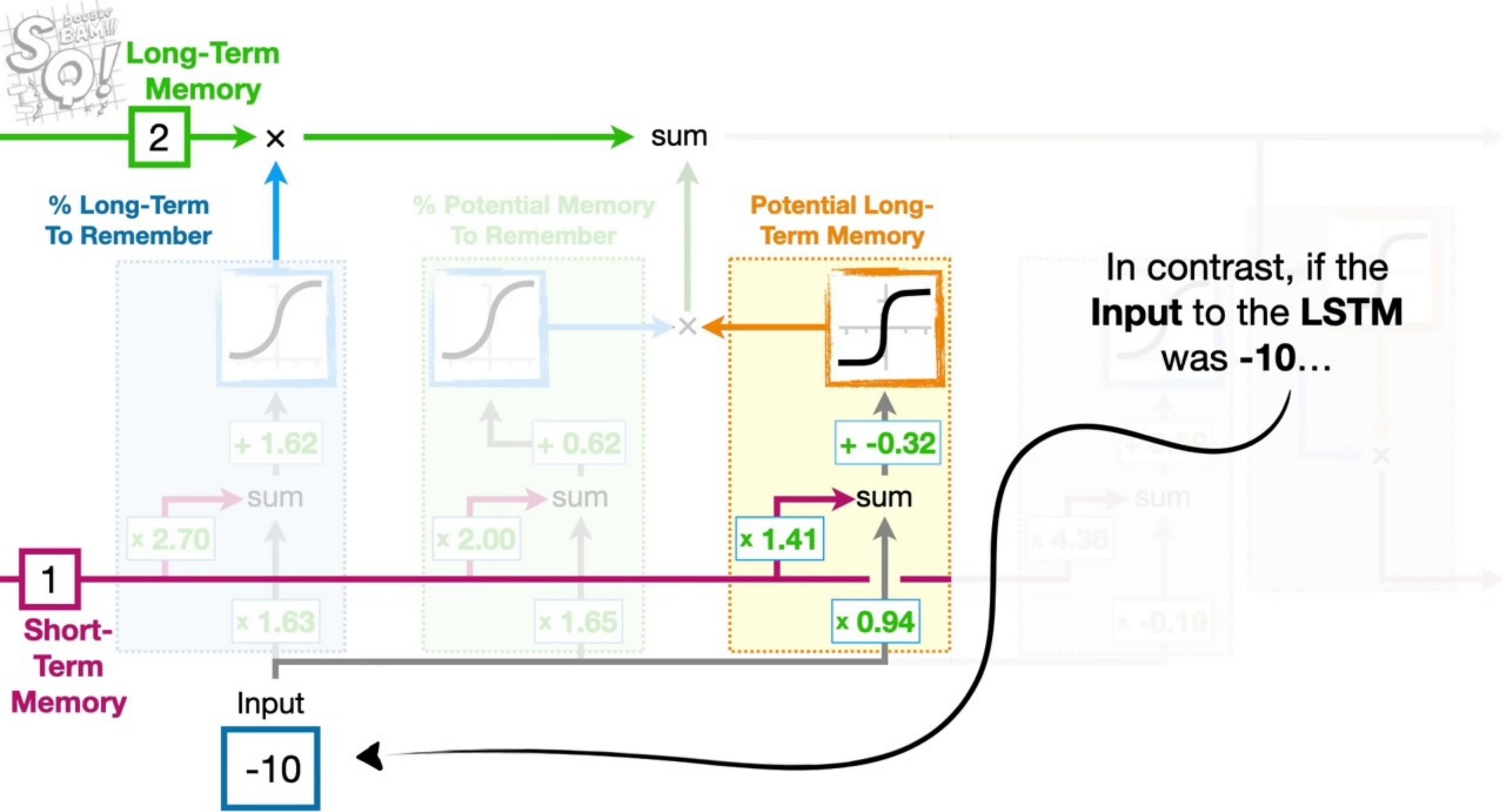


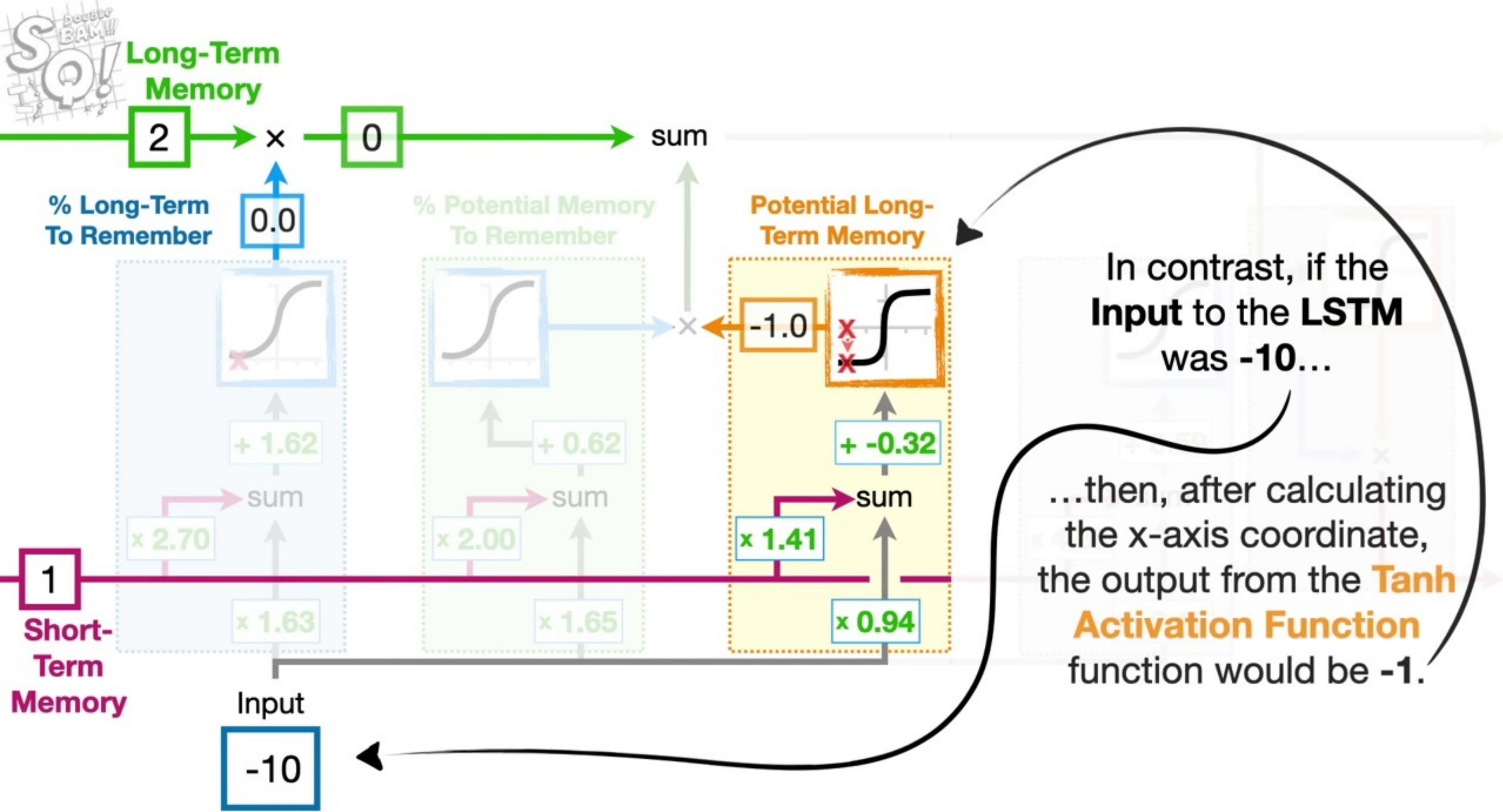


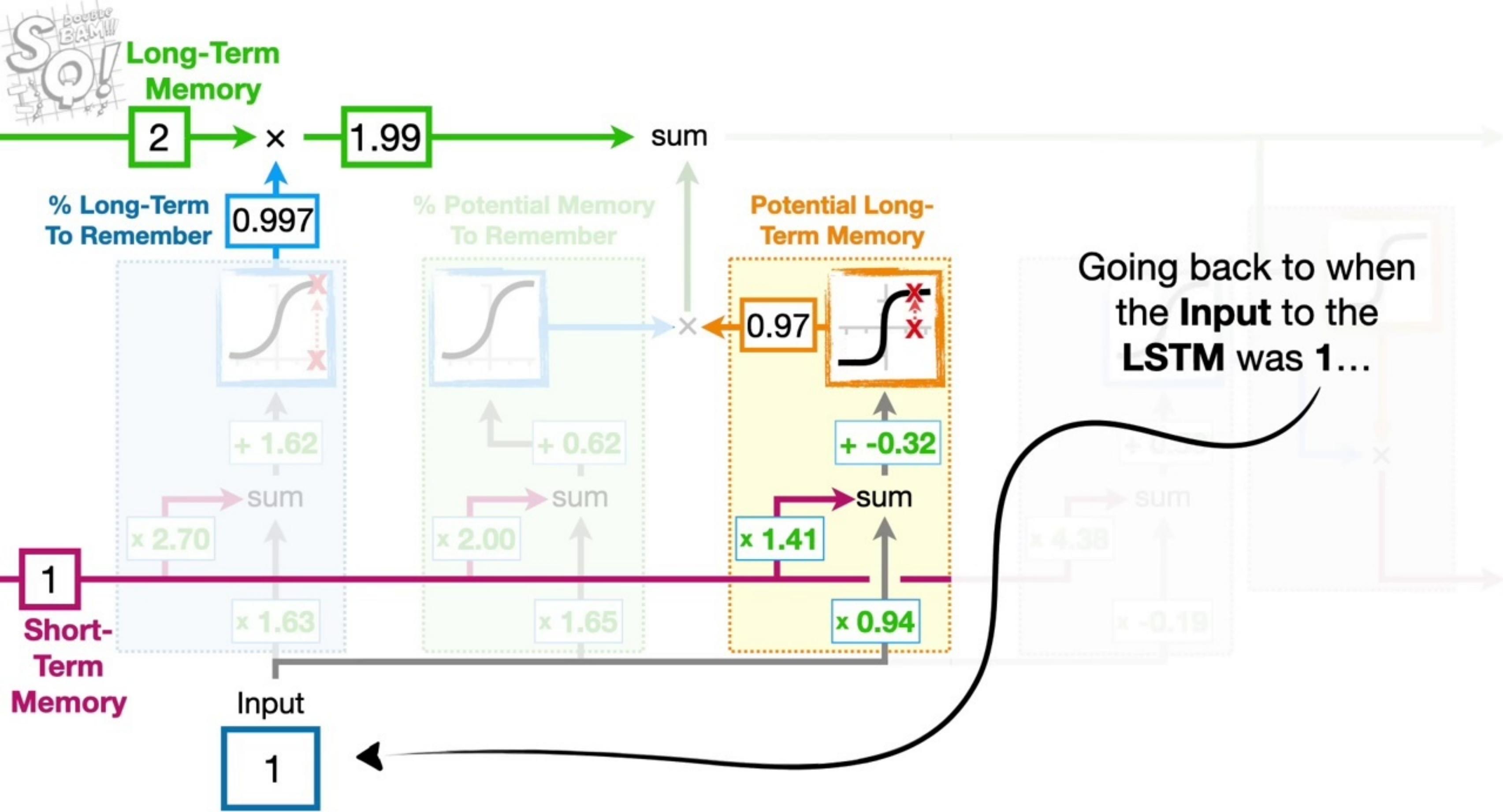


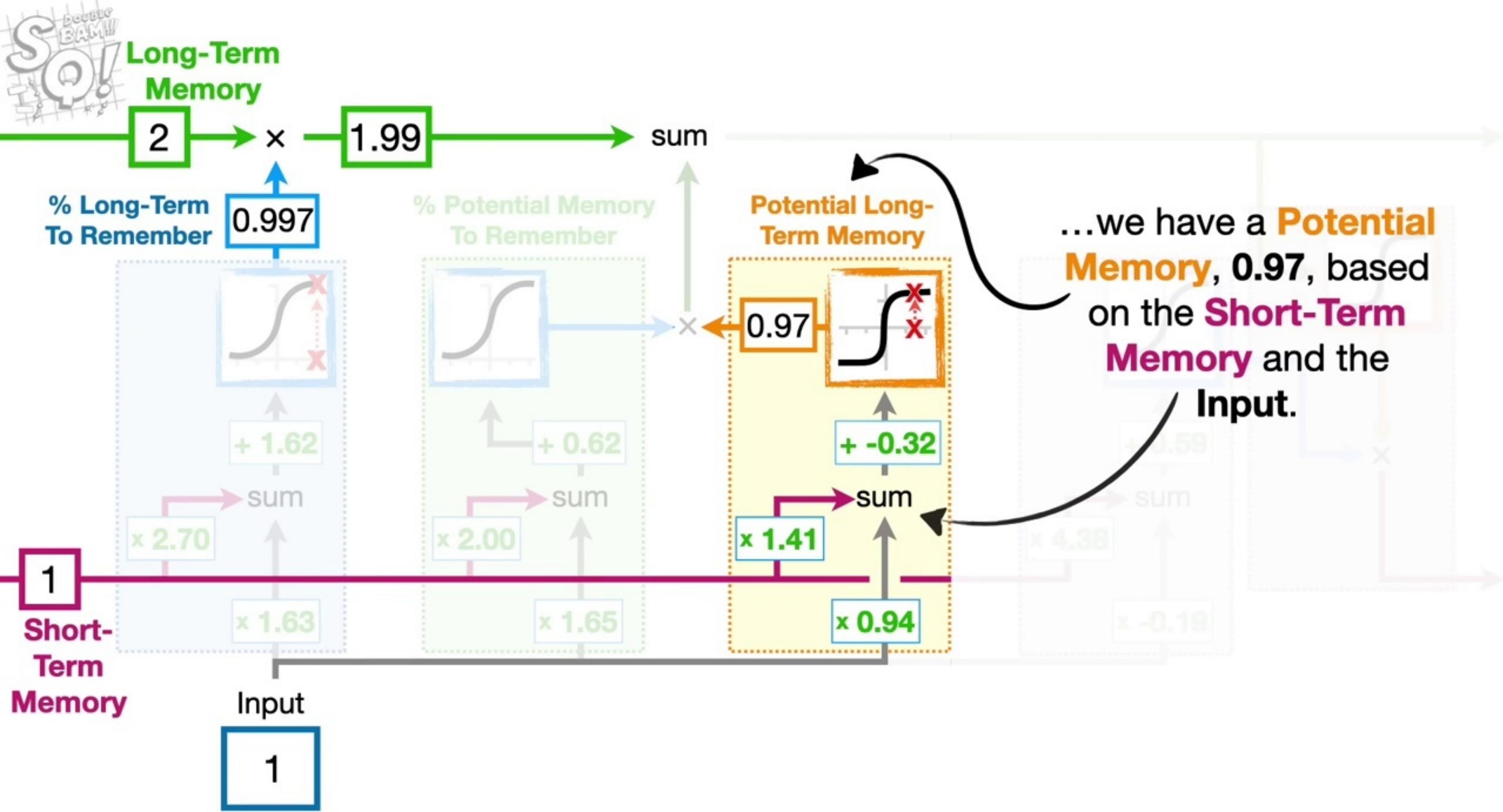


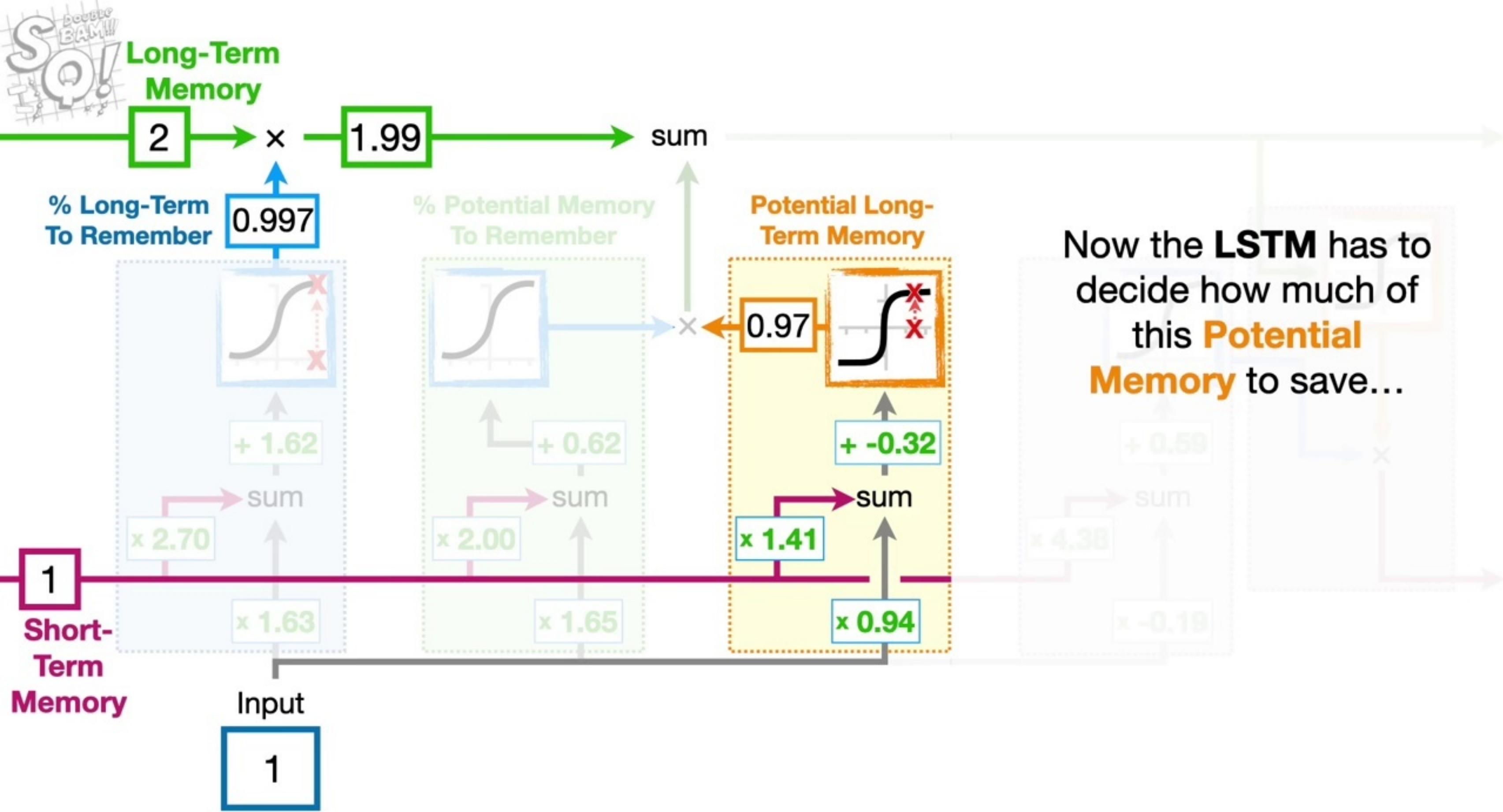


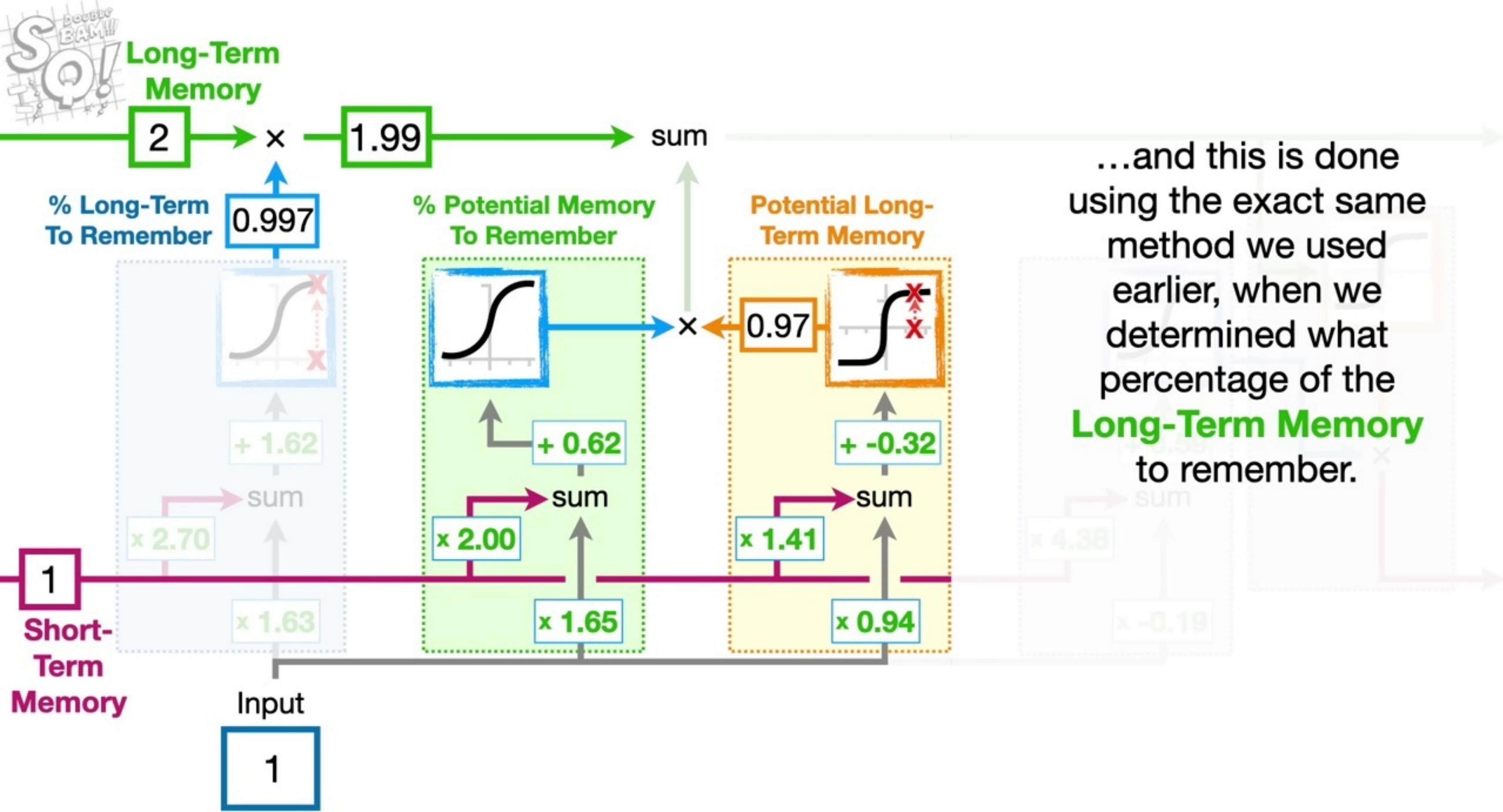




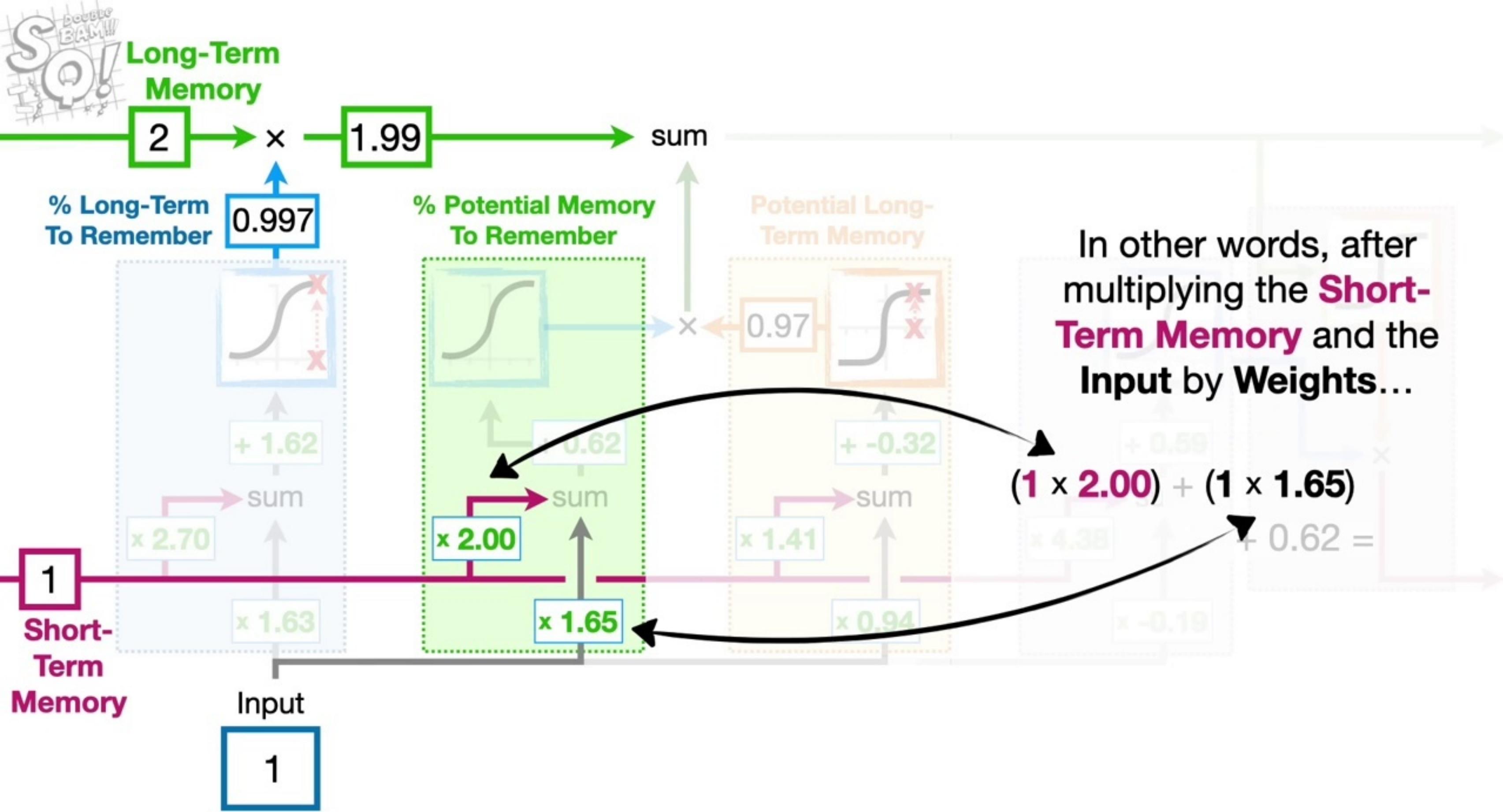


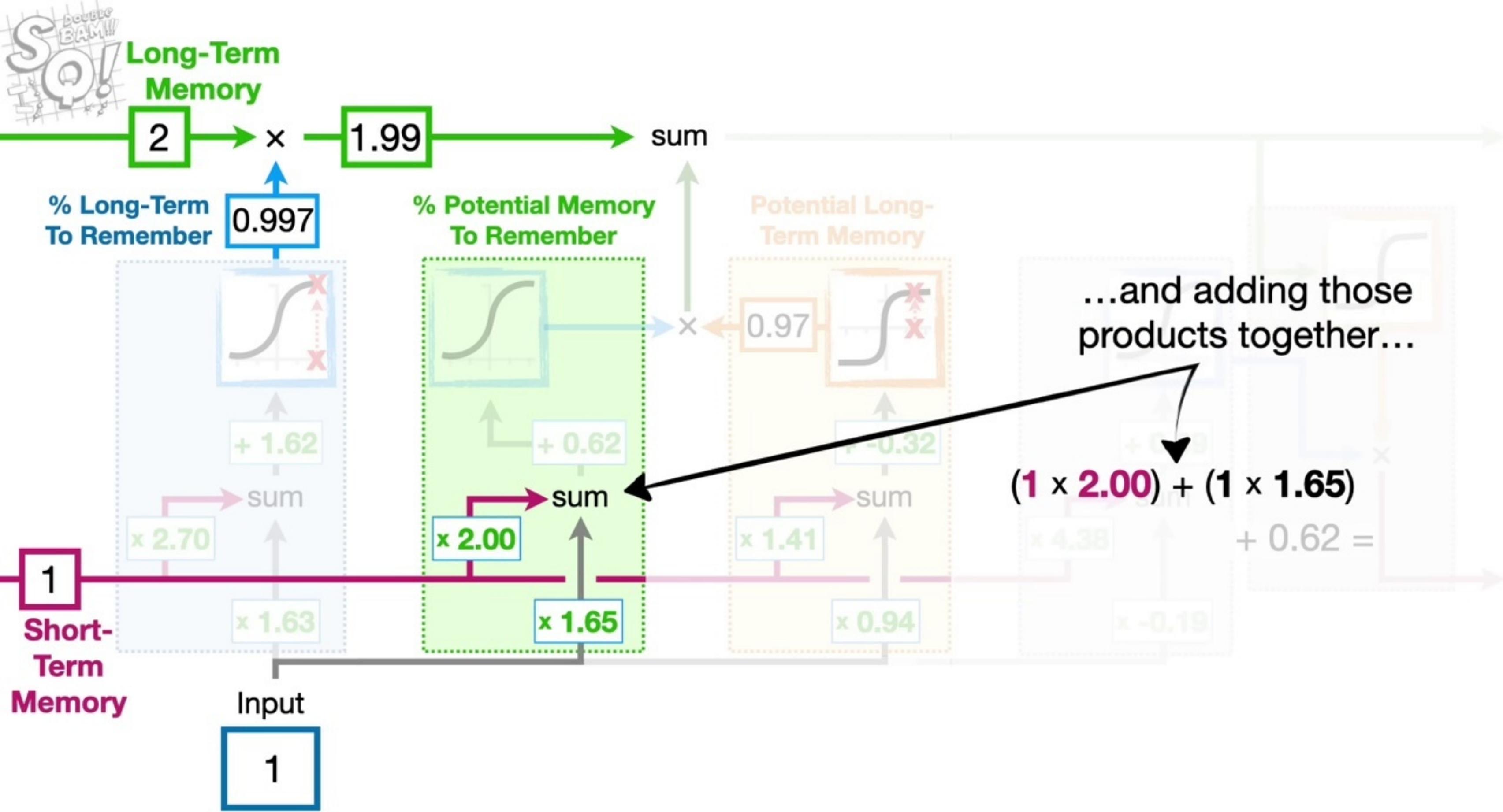


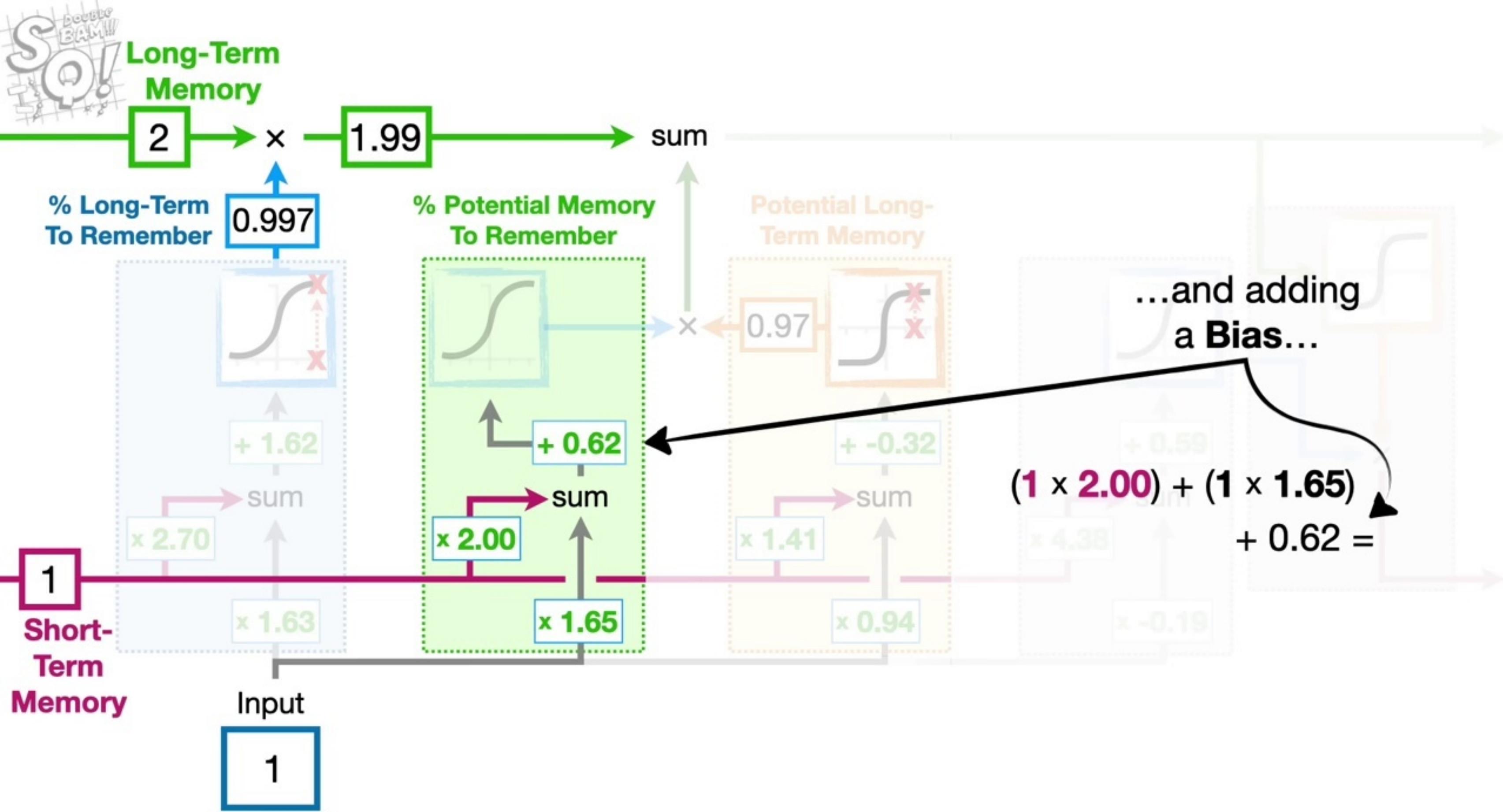


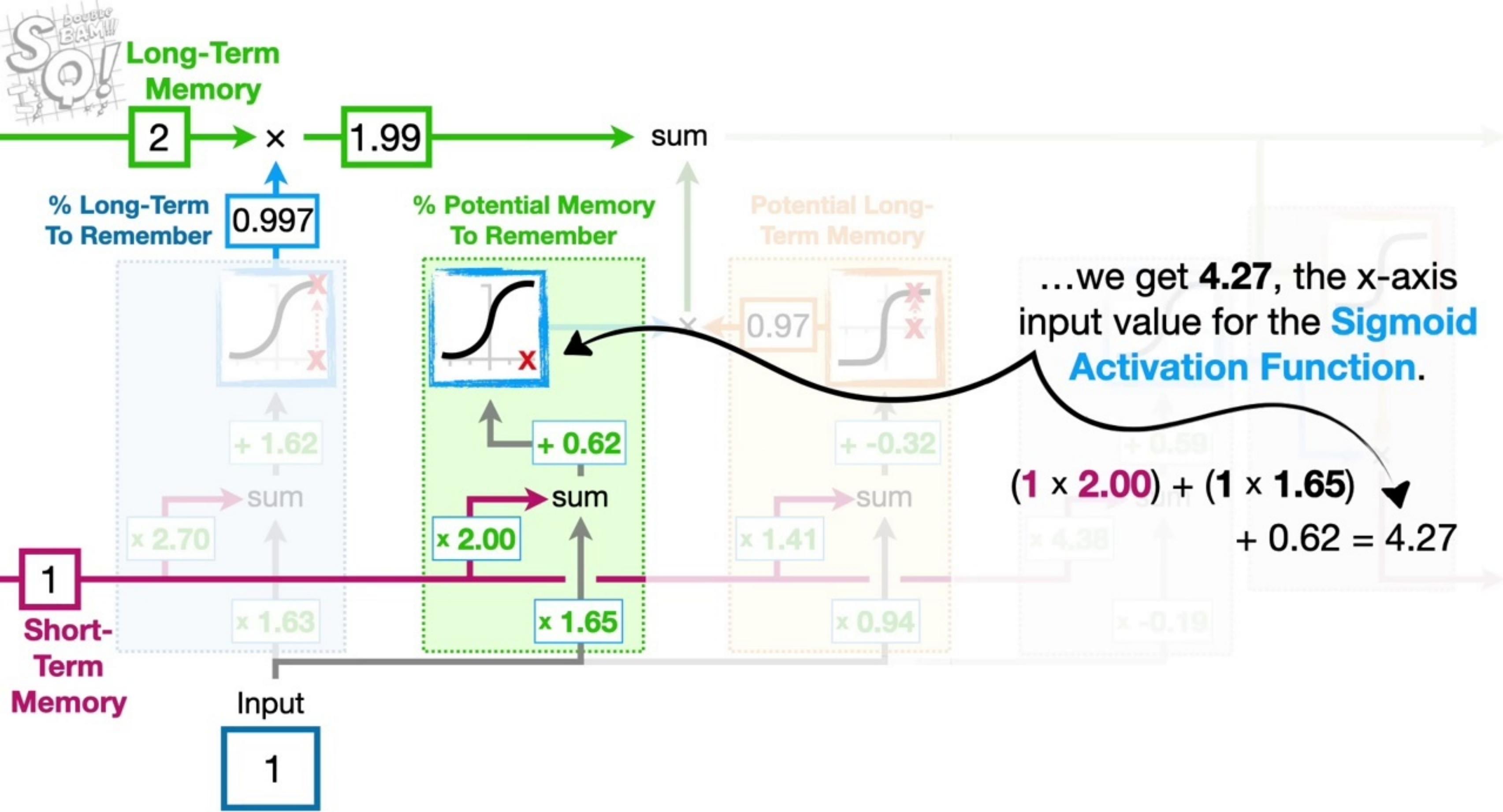


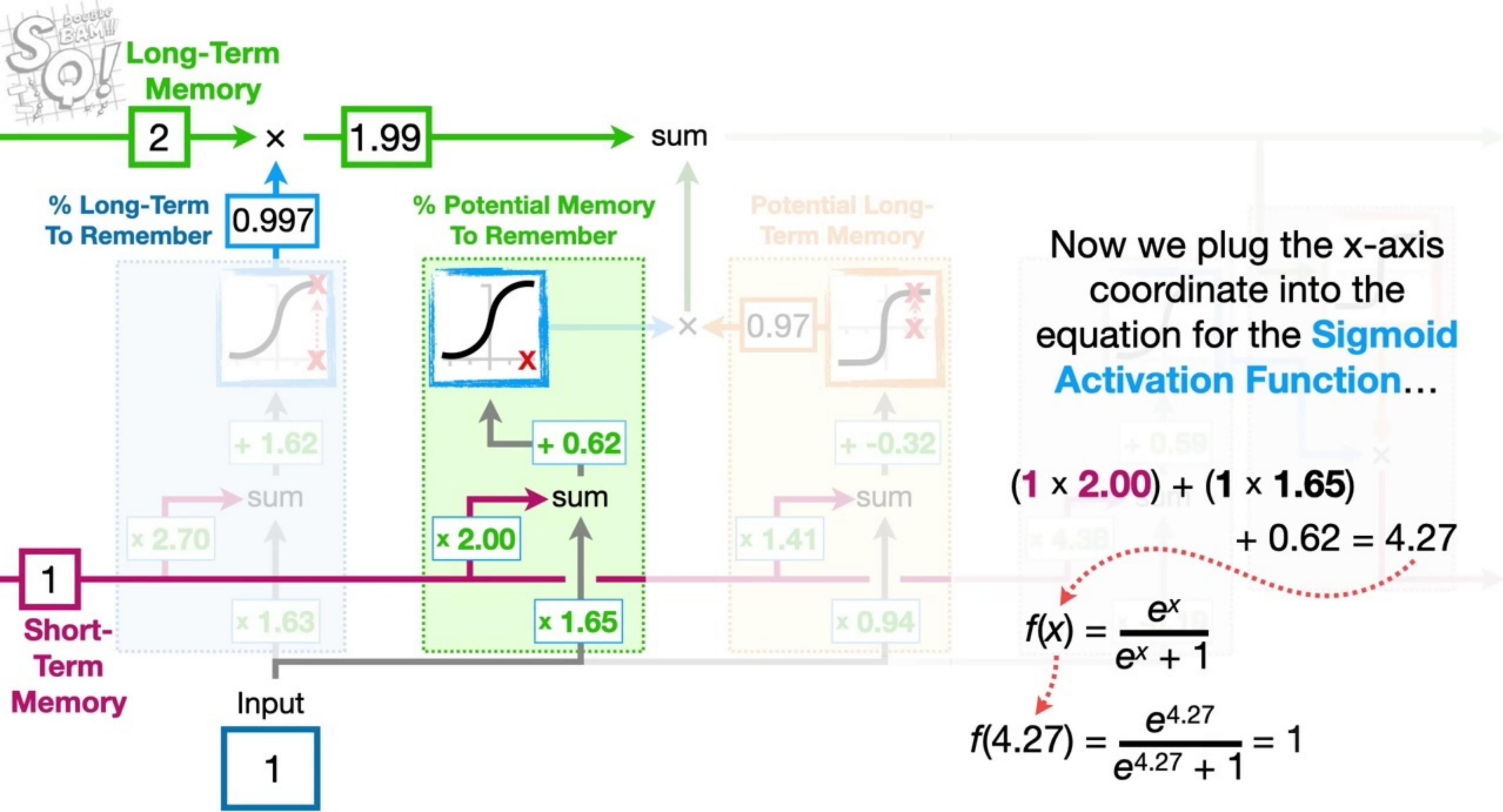
...and this is done using the exact same method we used earlier, when we determined what percentage of the **Long-Term Memory** to remember.

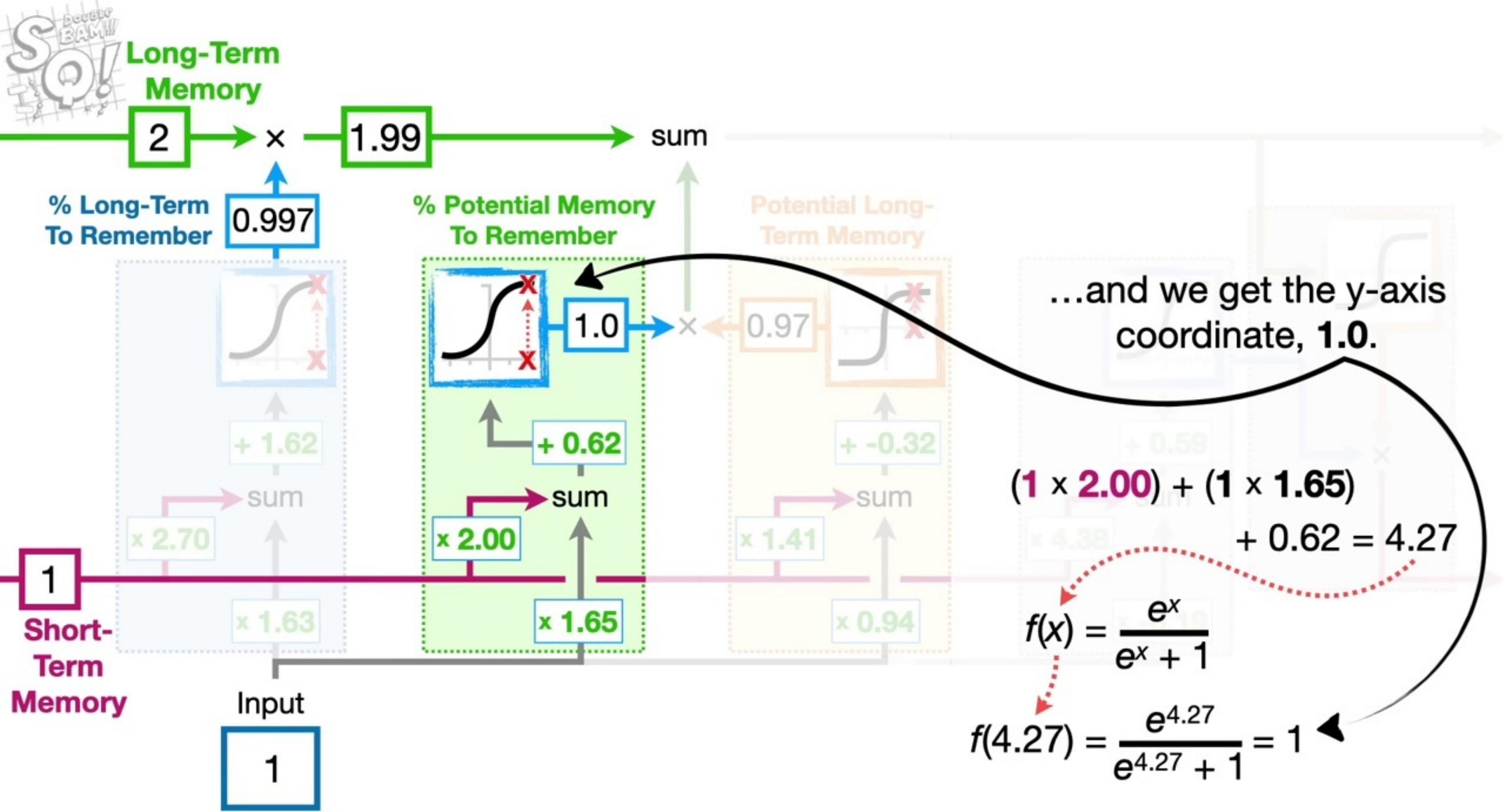


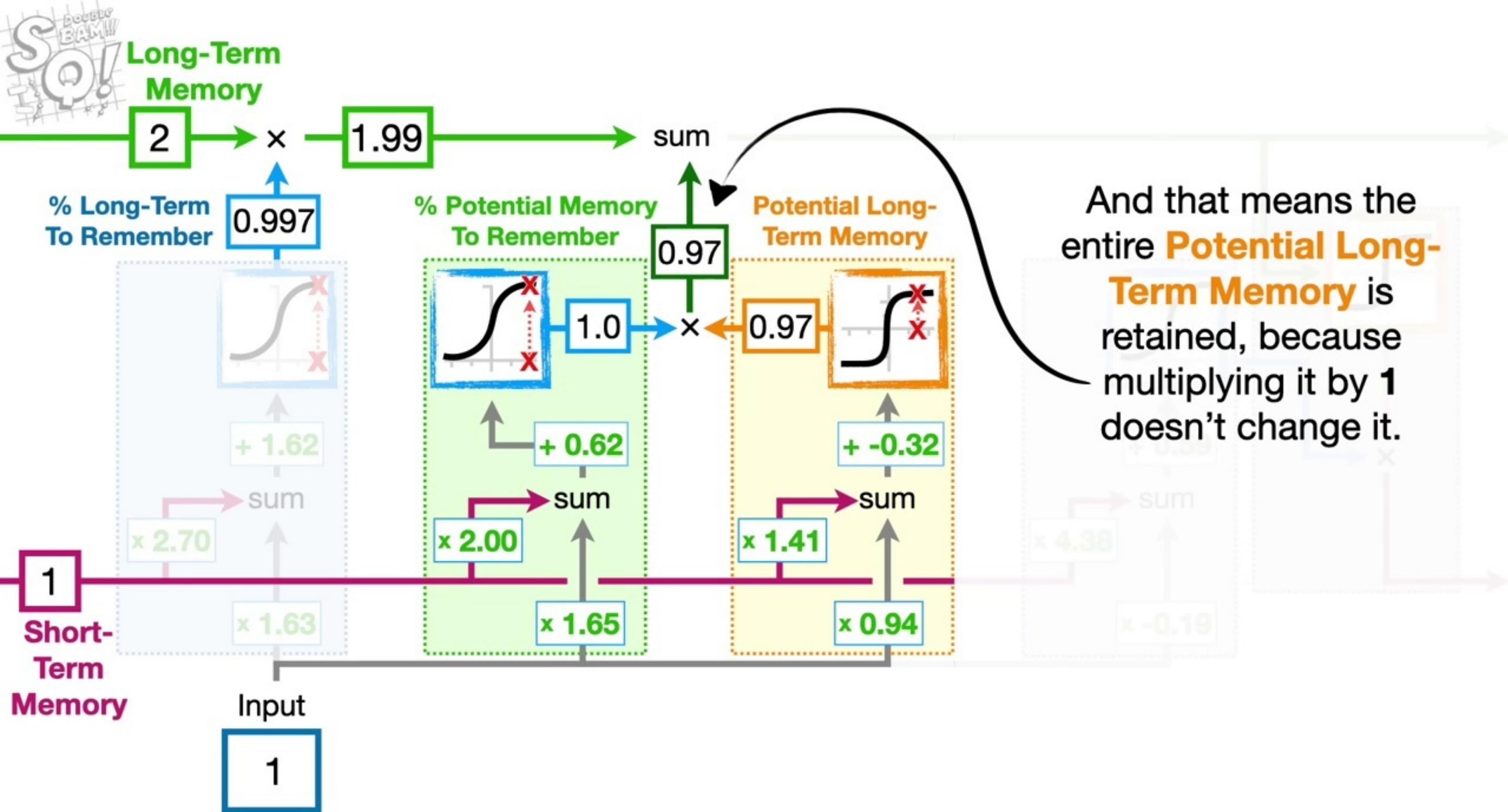


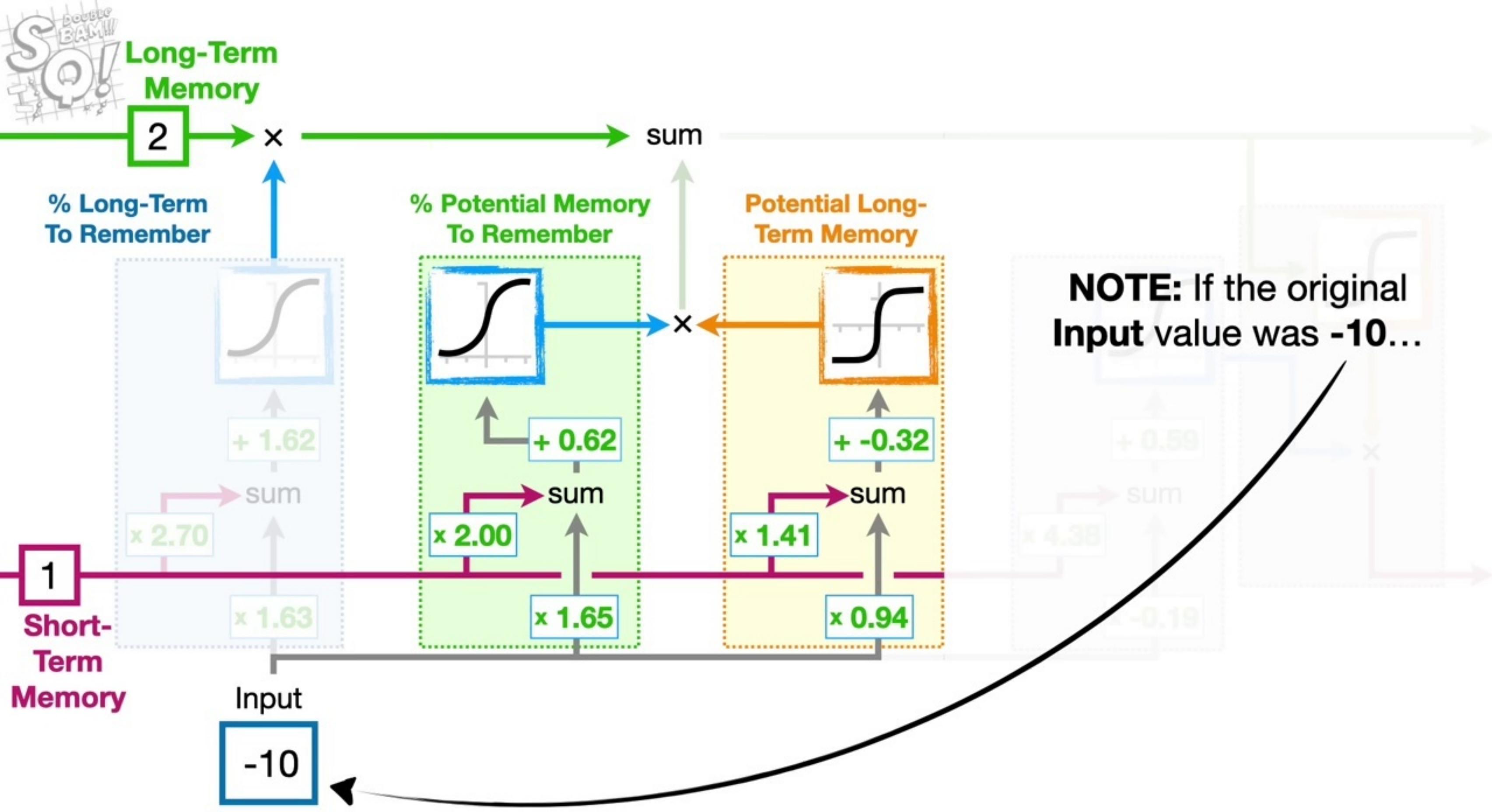


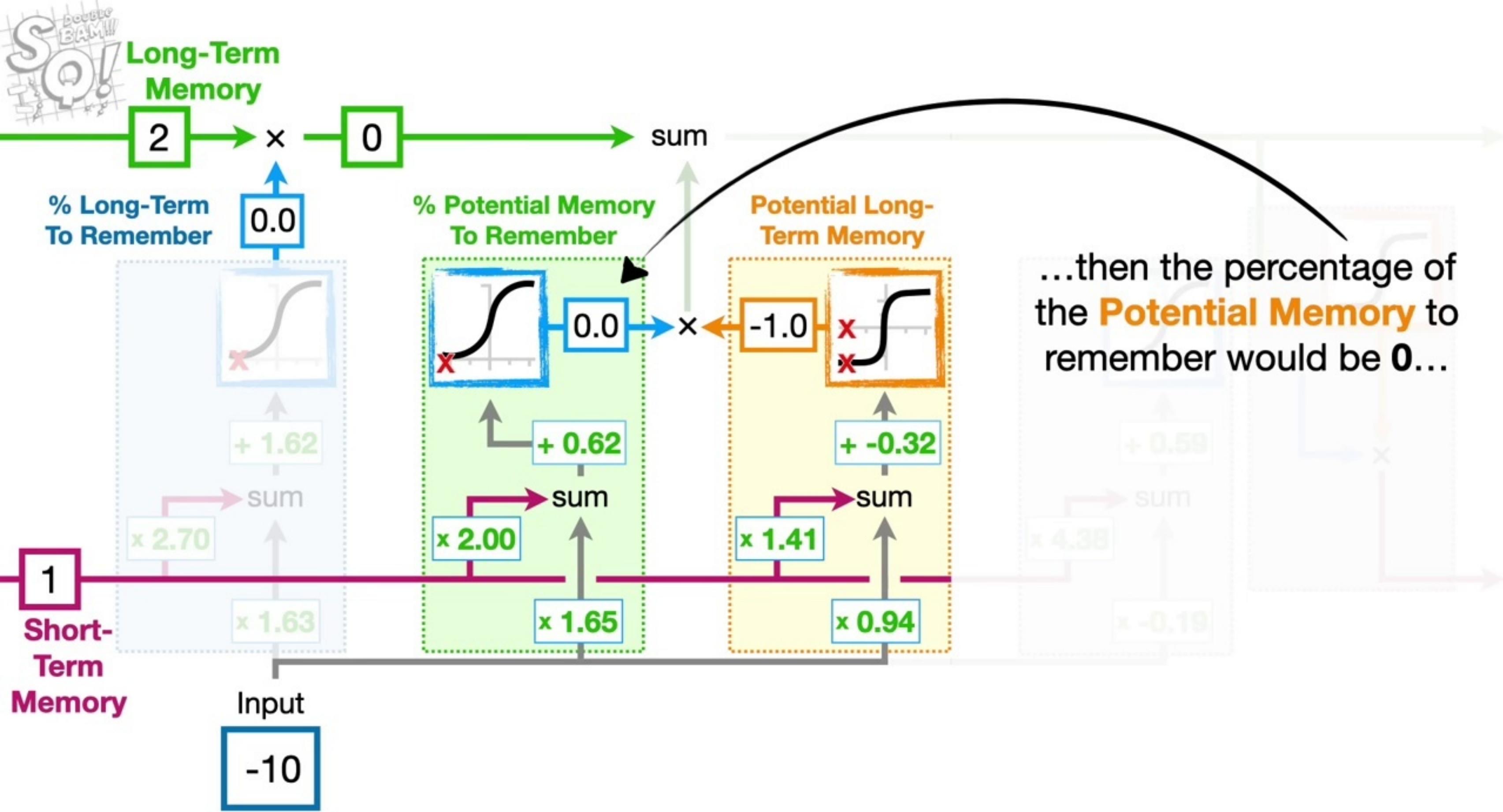


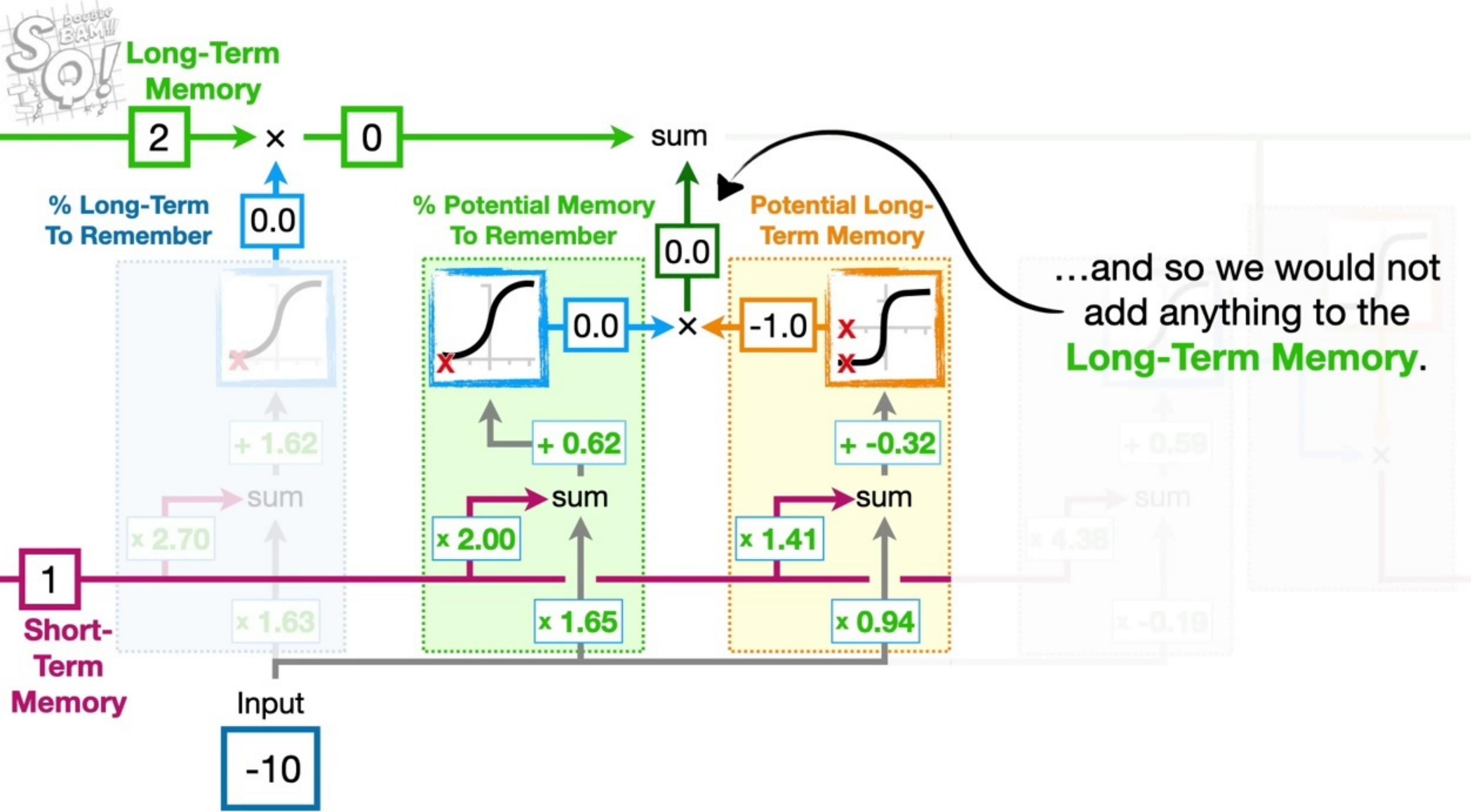


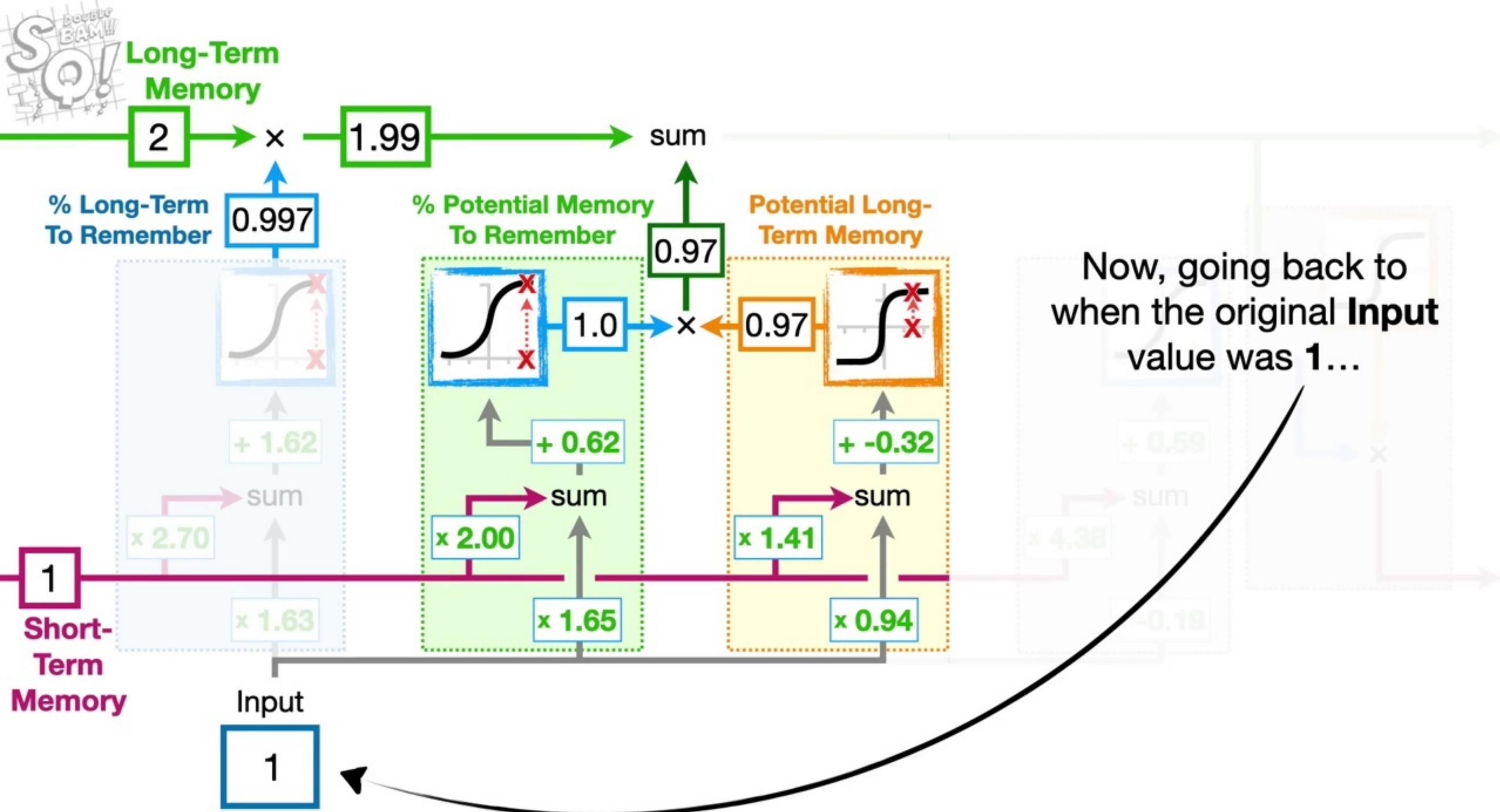


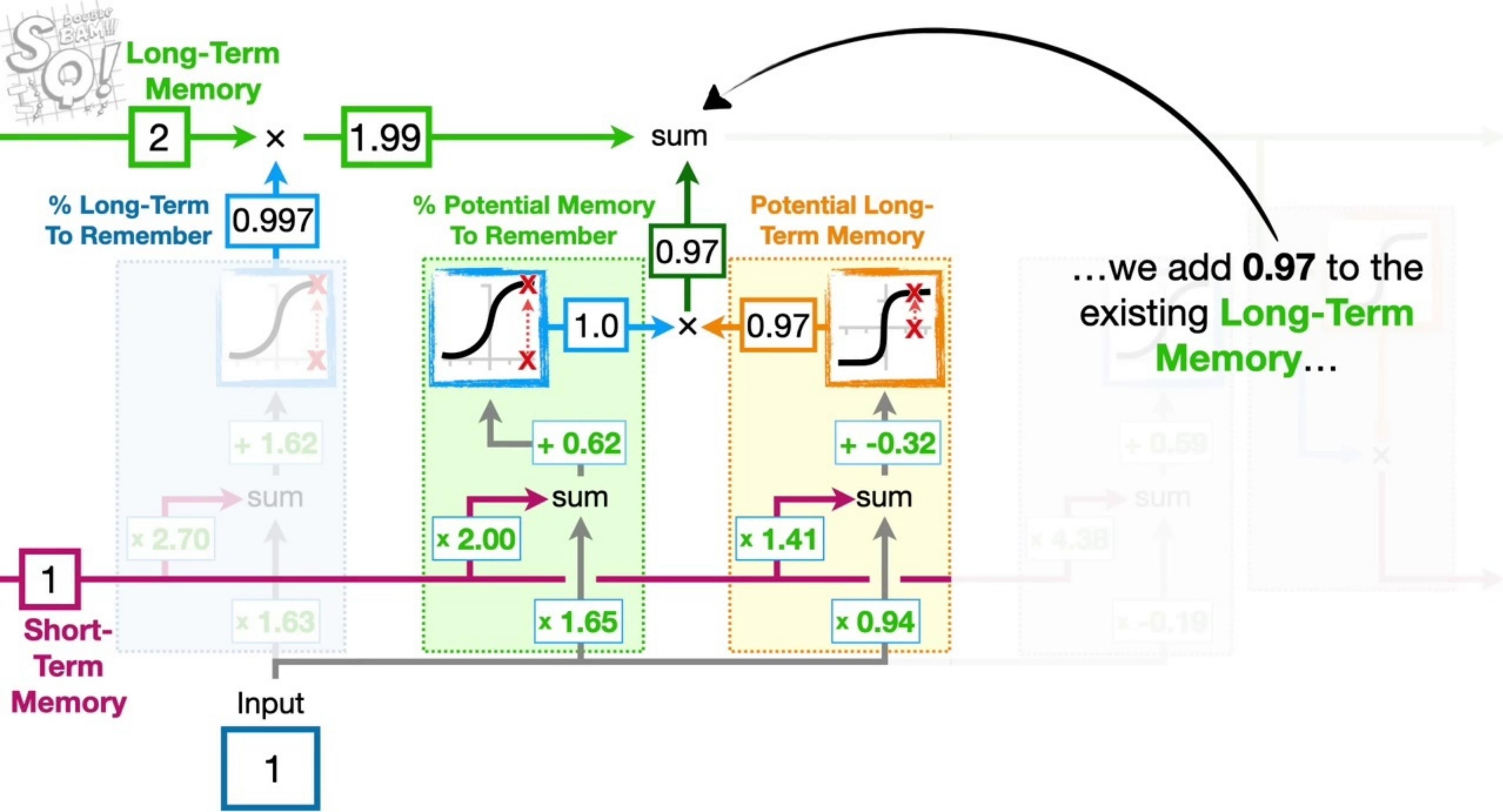


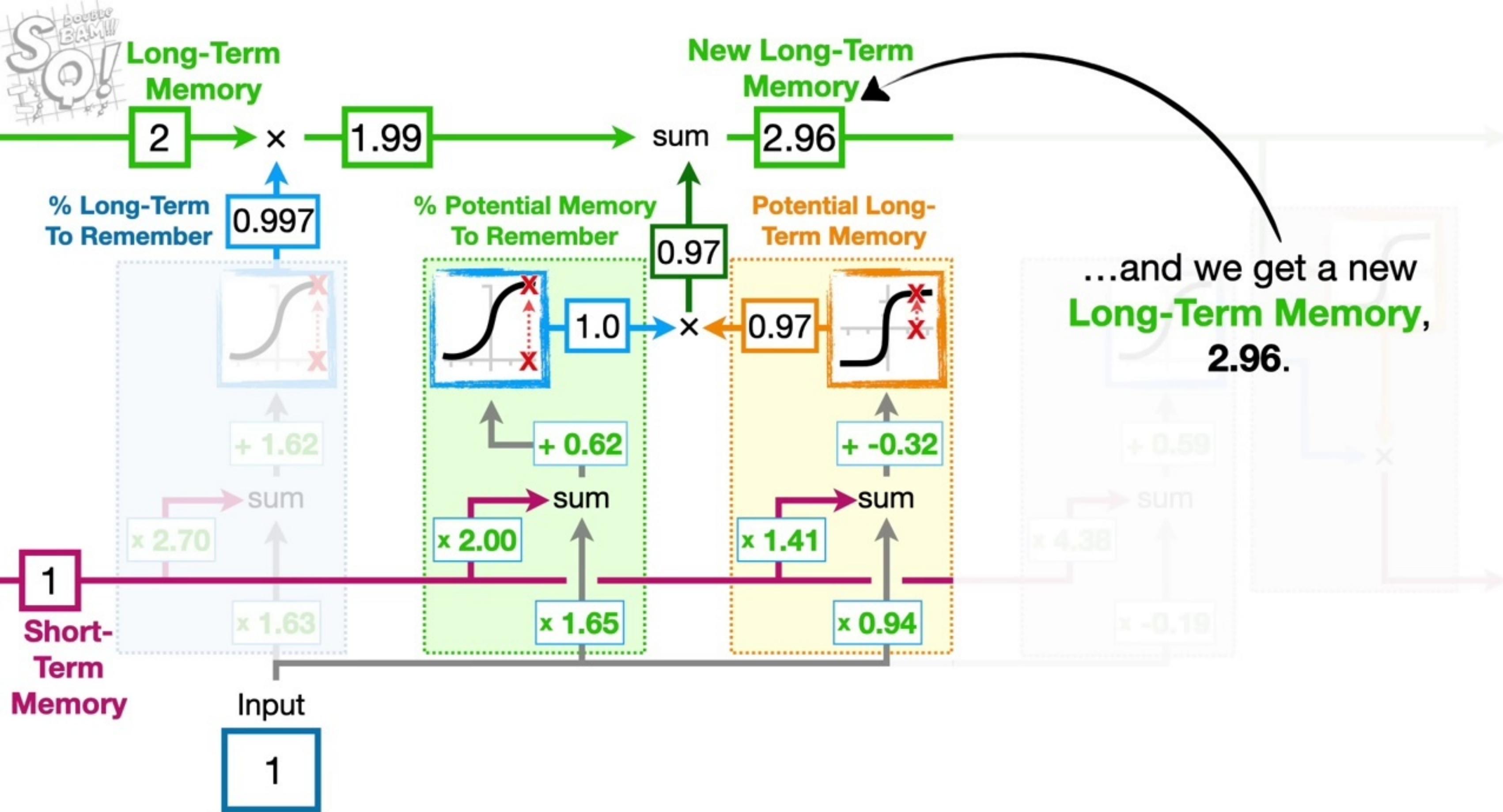


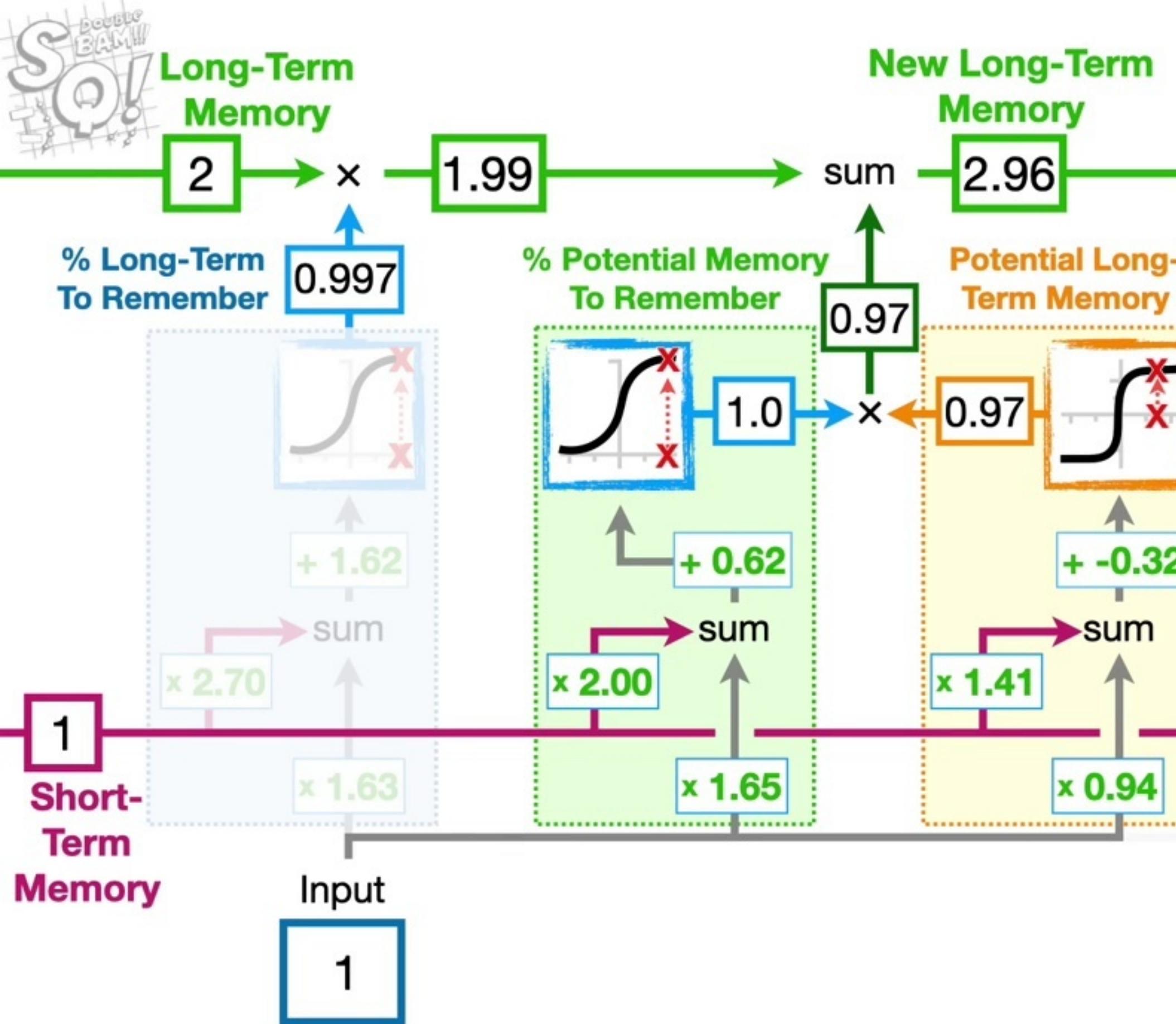








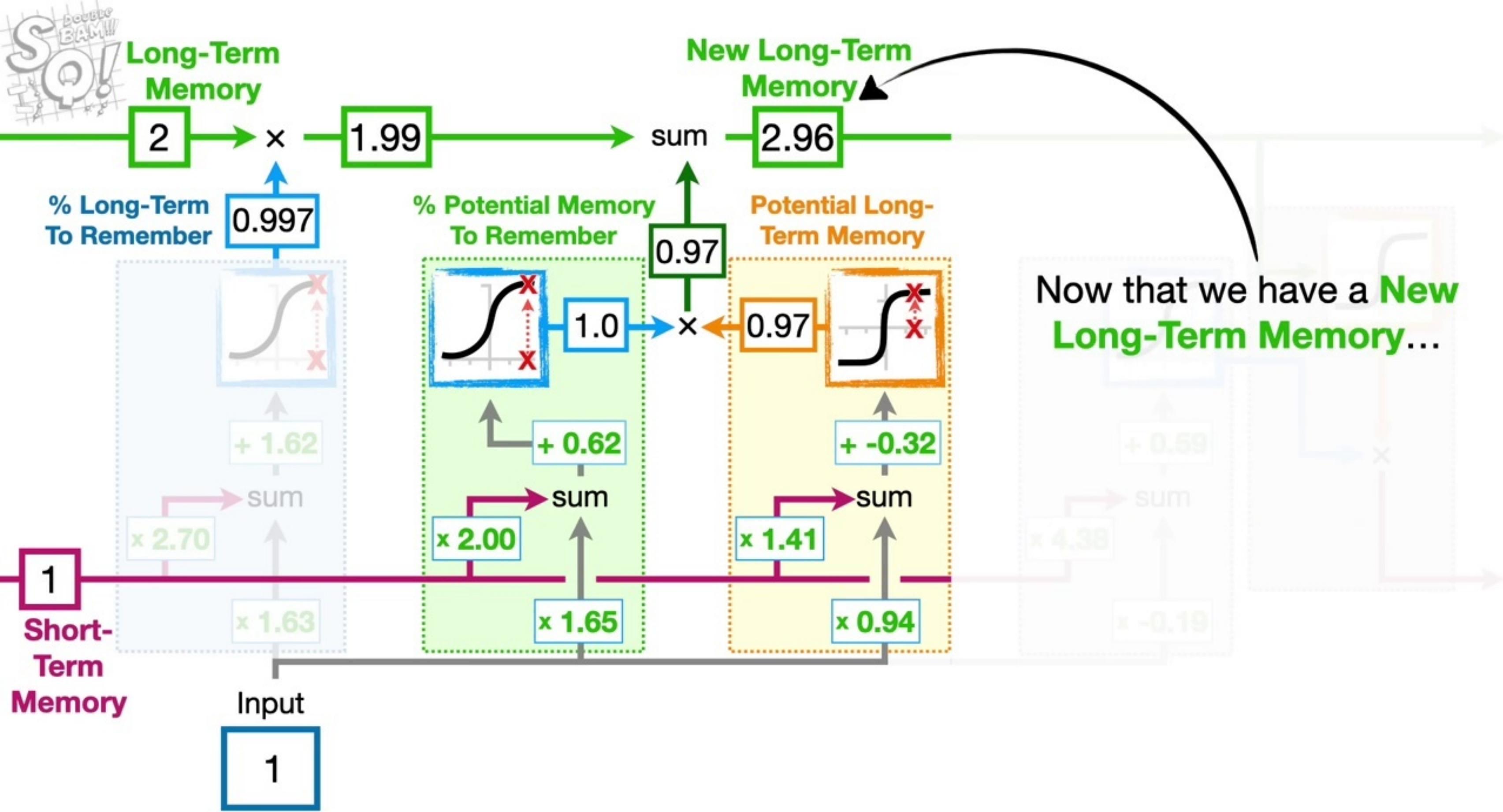


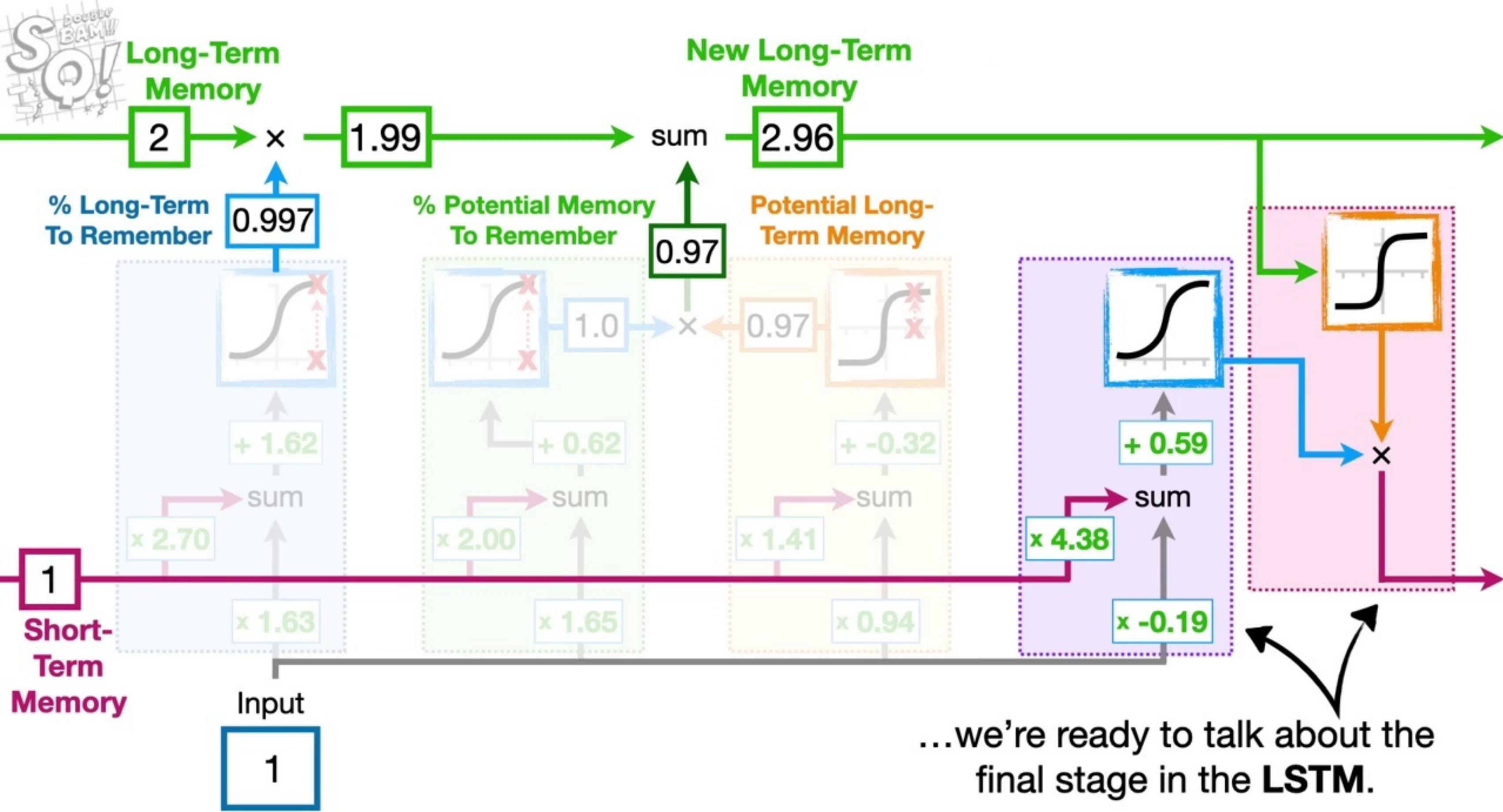


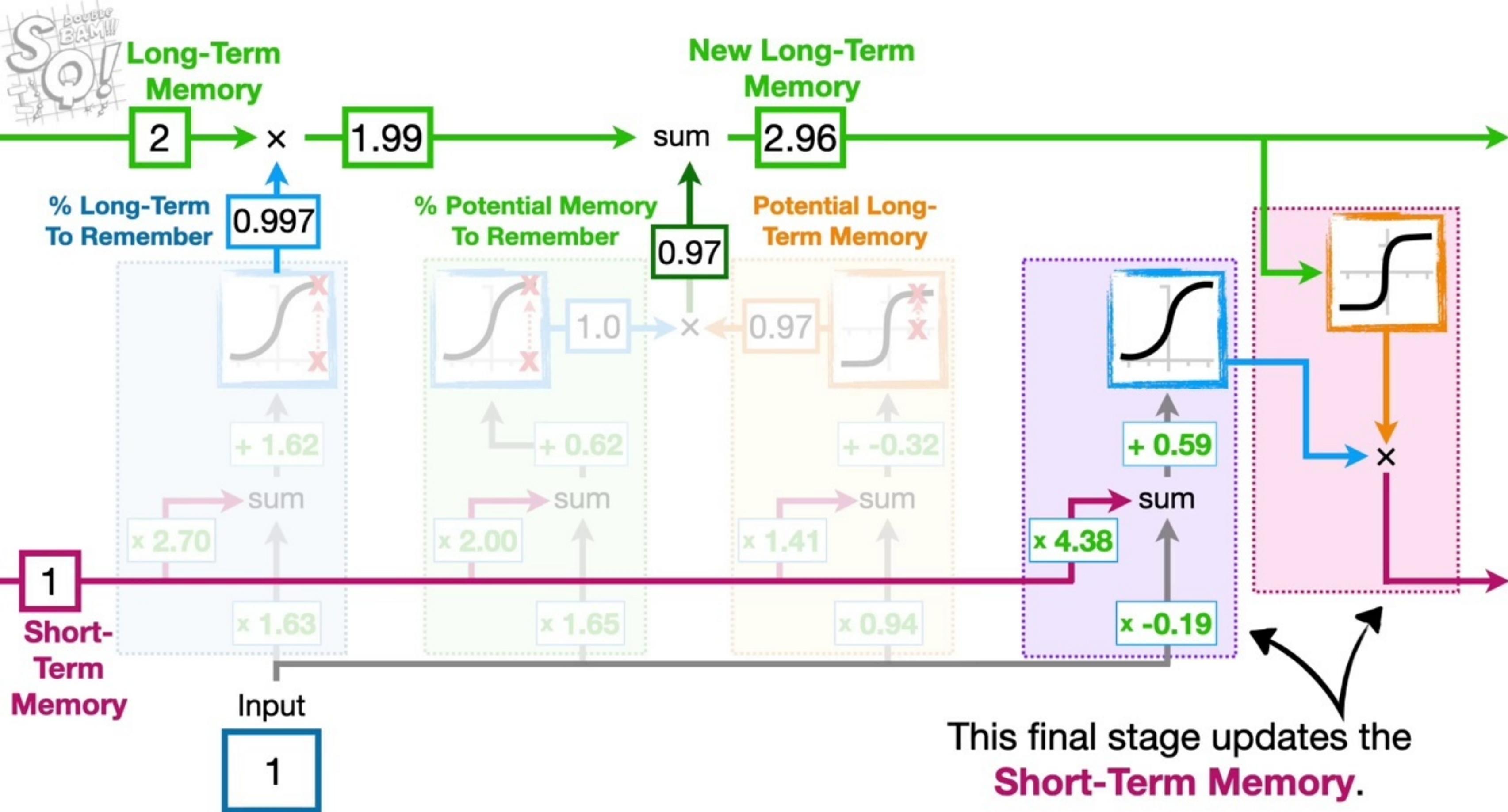
TERMINOLOGY ALERT!!!

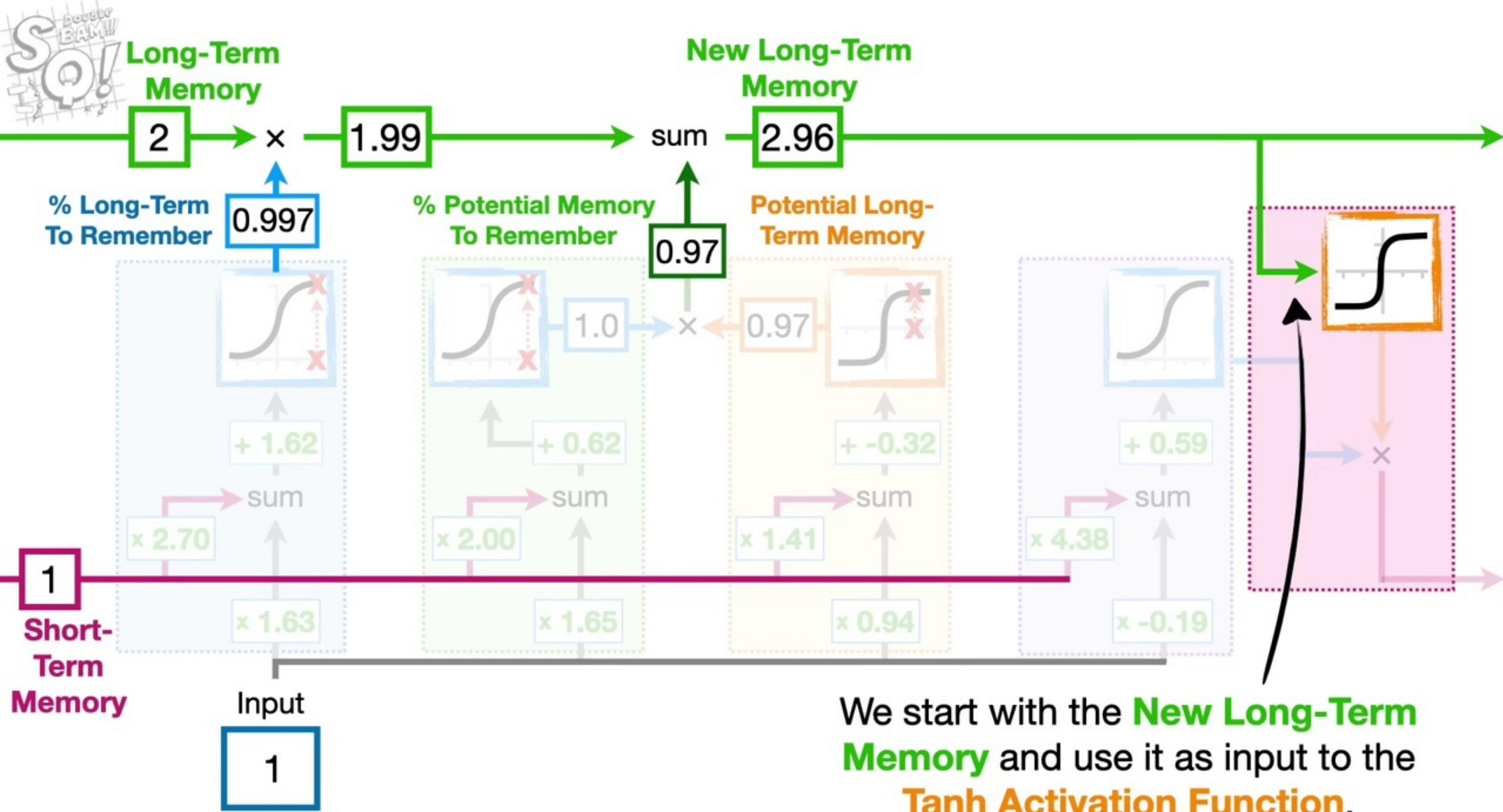
Even though this part of the **Long Short-Term Memory** unit determines how we should update the **Long-Term Memory**...

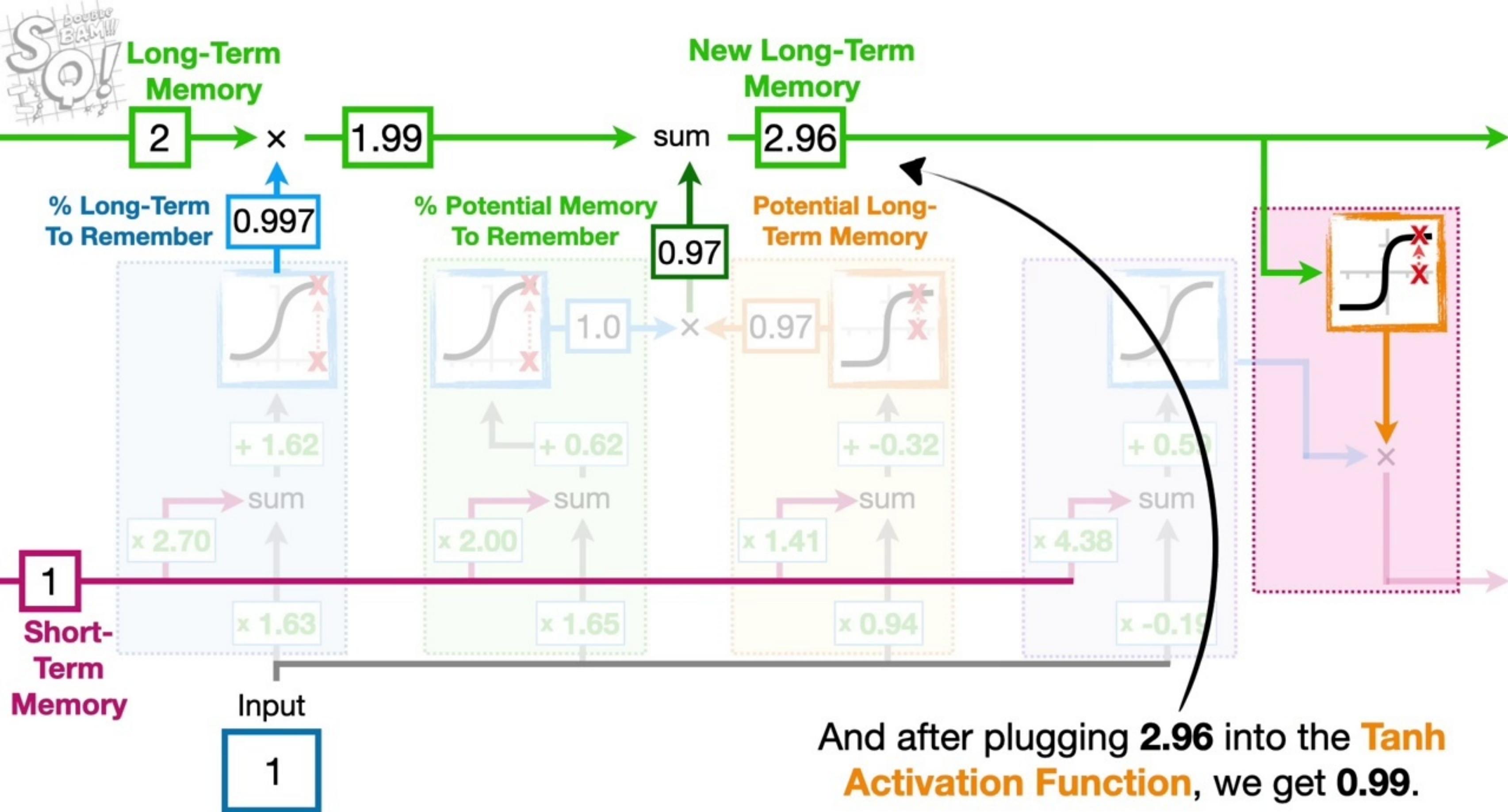
...it is usually called the **Input Gate**.

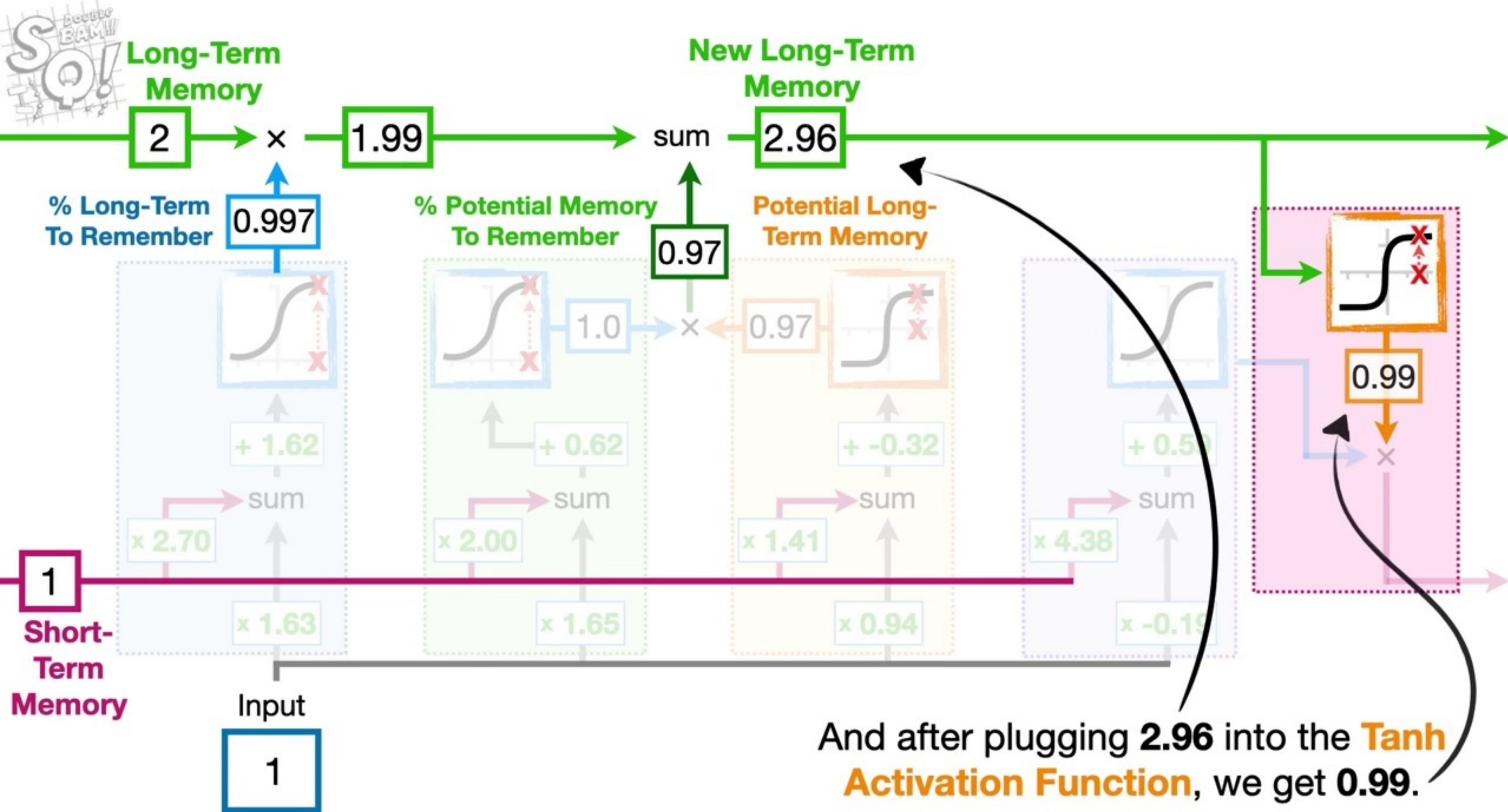




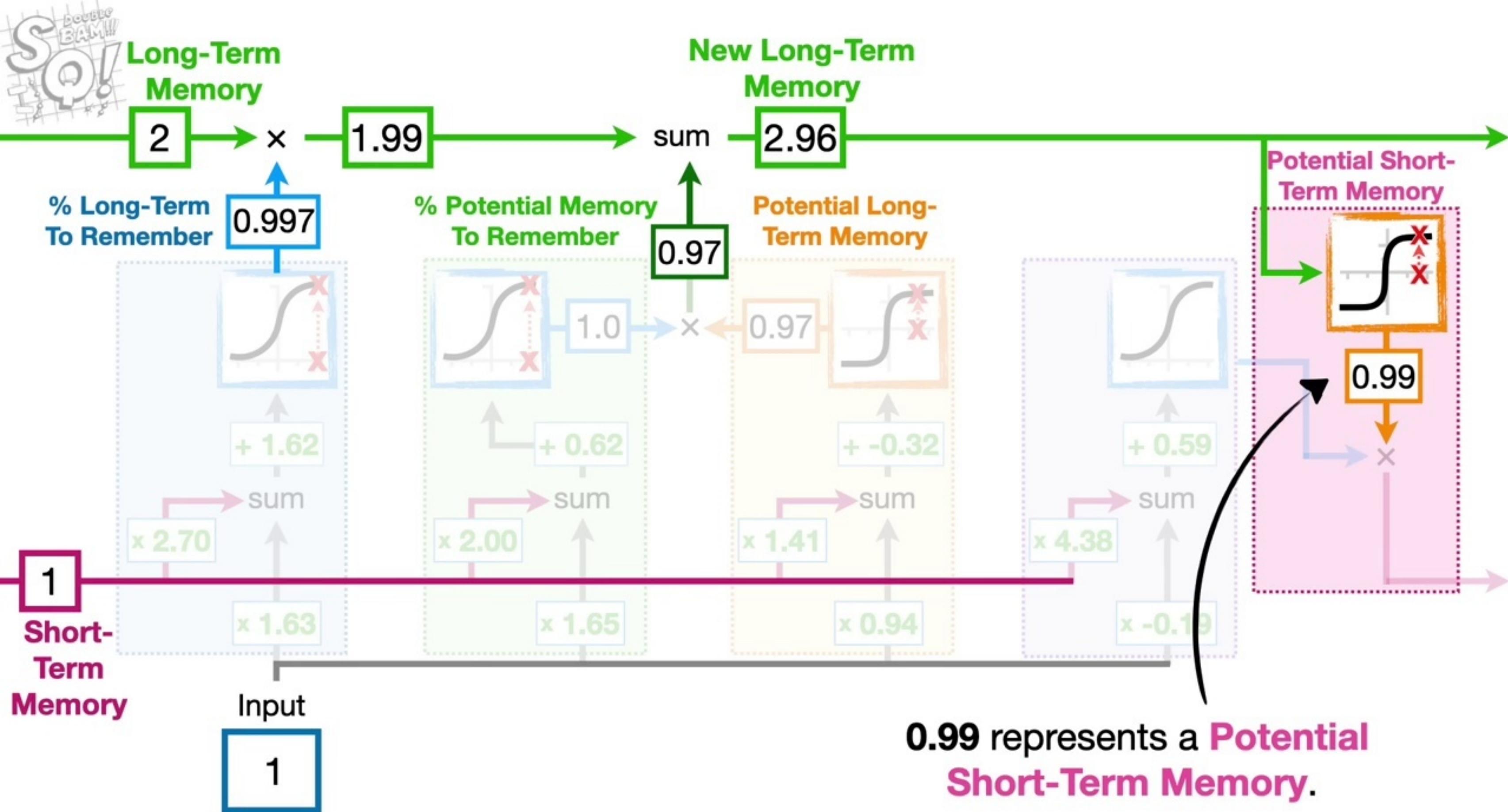


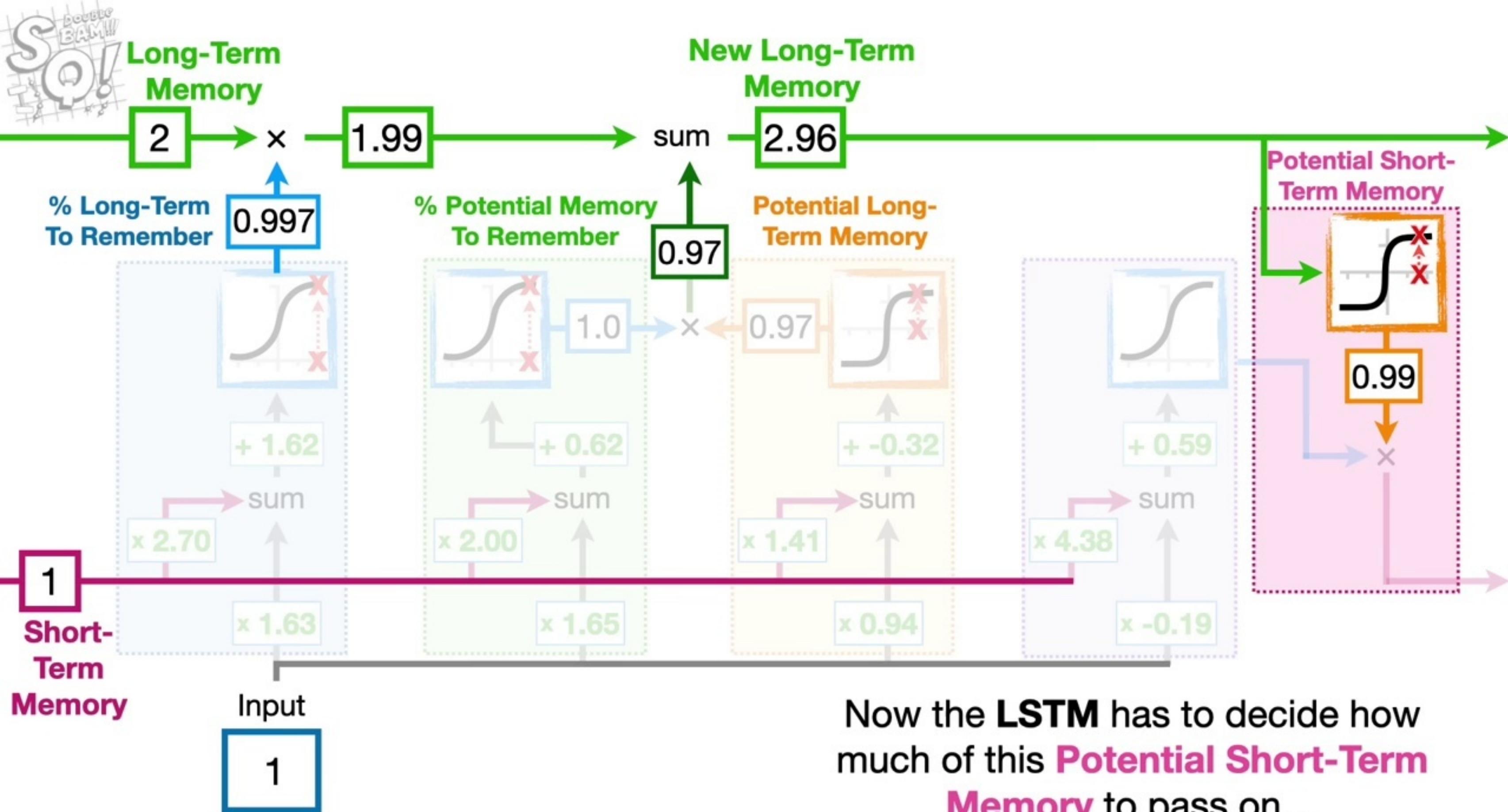




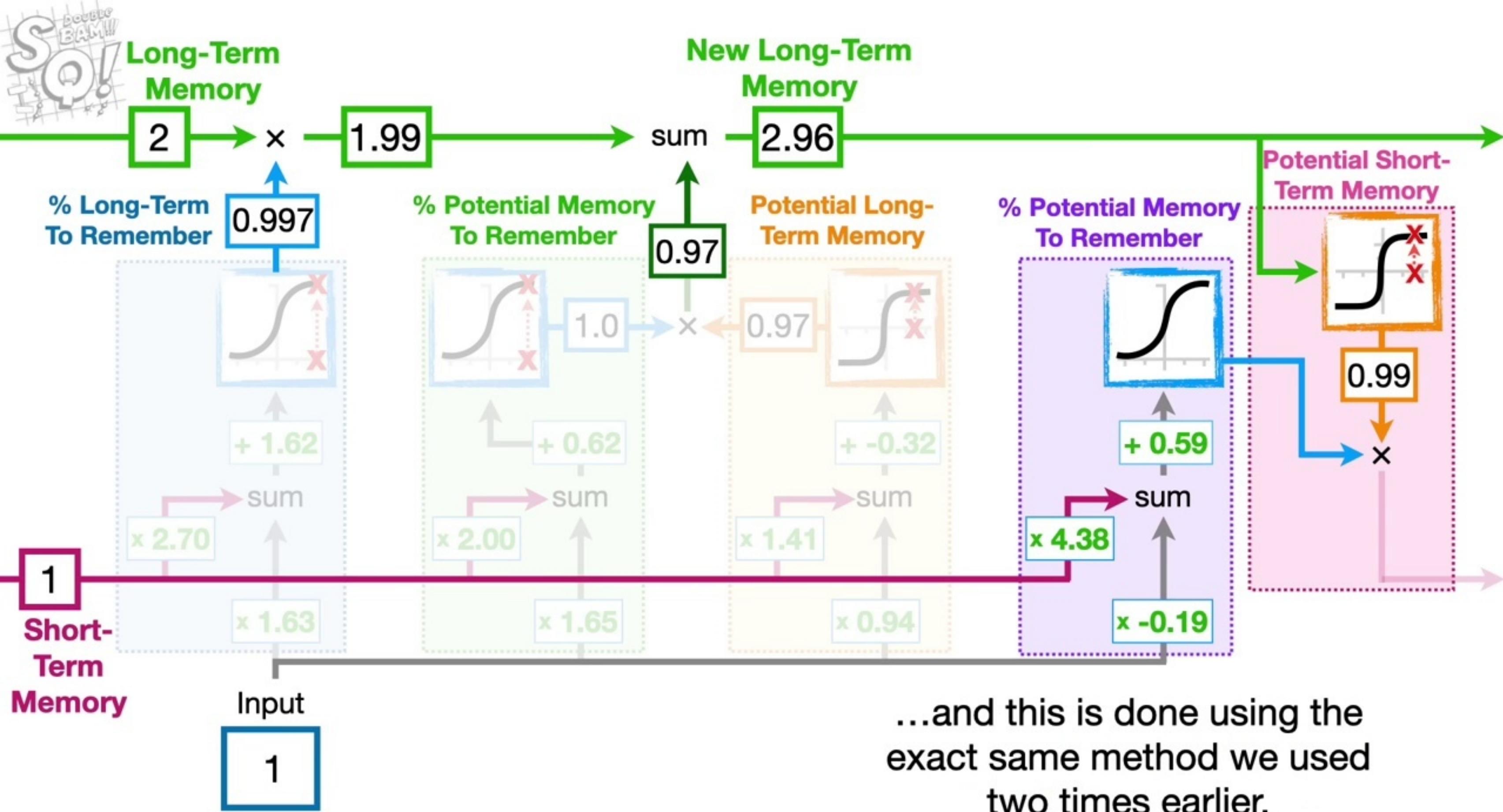


And after plugging **2.96** into the **Tanh Activation Function**, we get **0.99**.

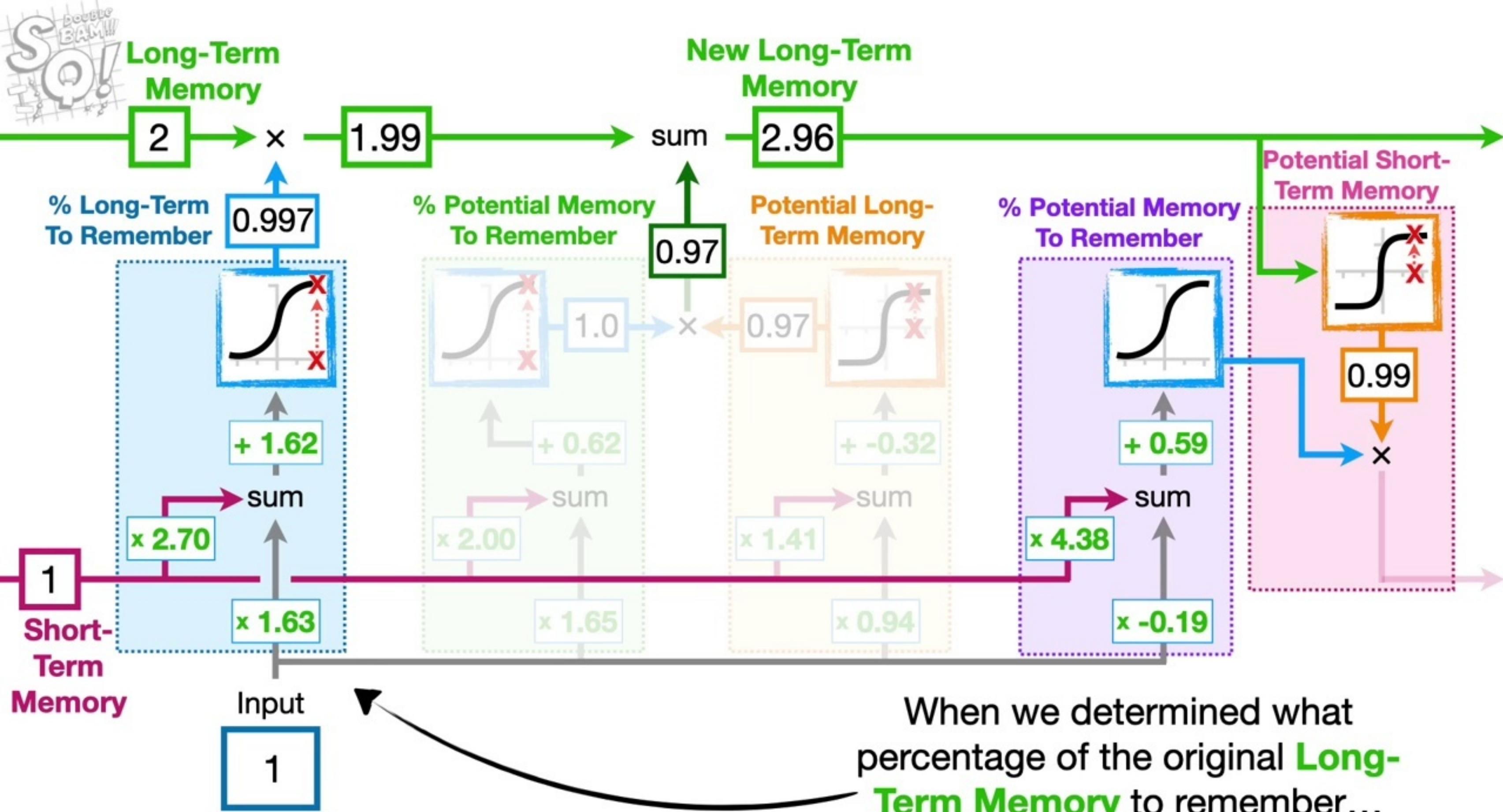


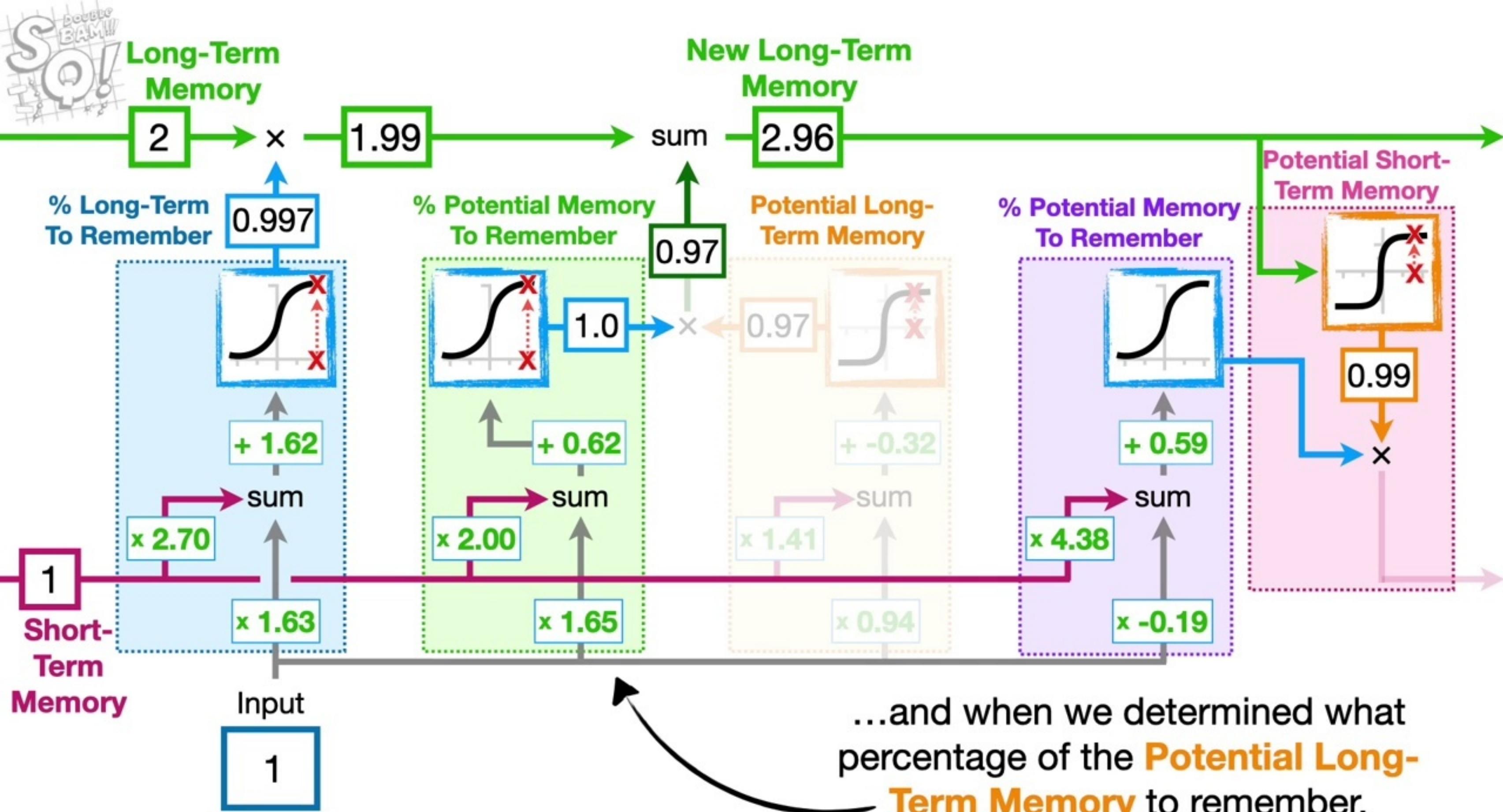


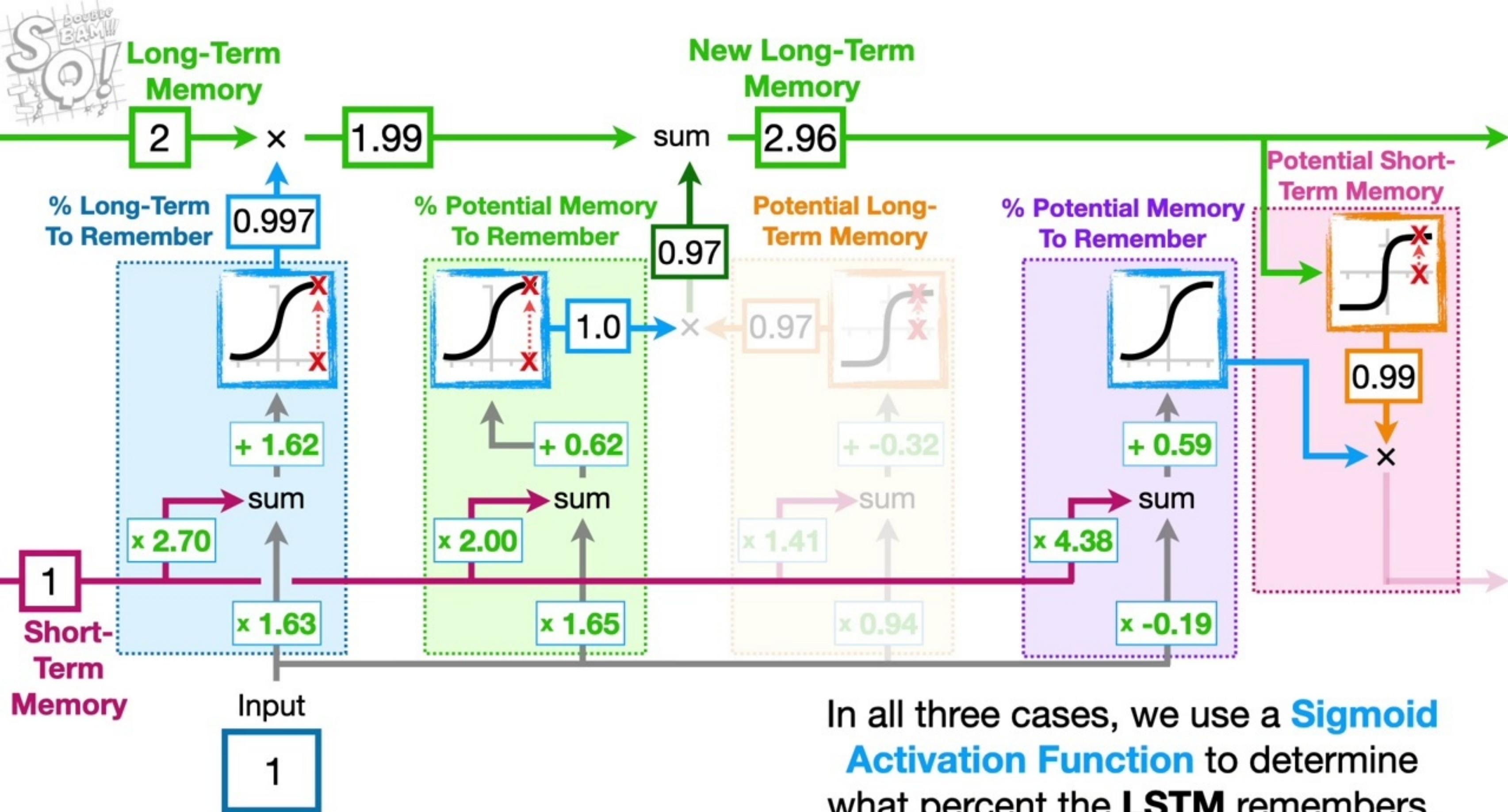
Now the **LSTM** has to decide how much of this **Potential Short-Term Memory** to pass on...



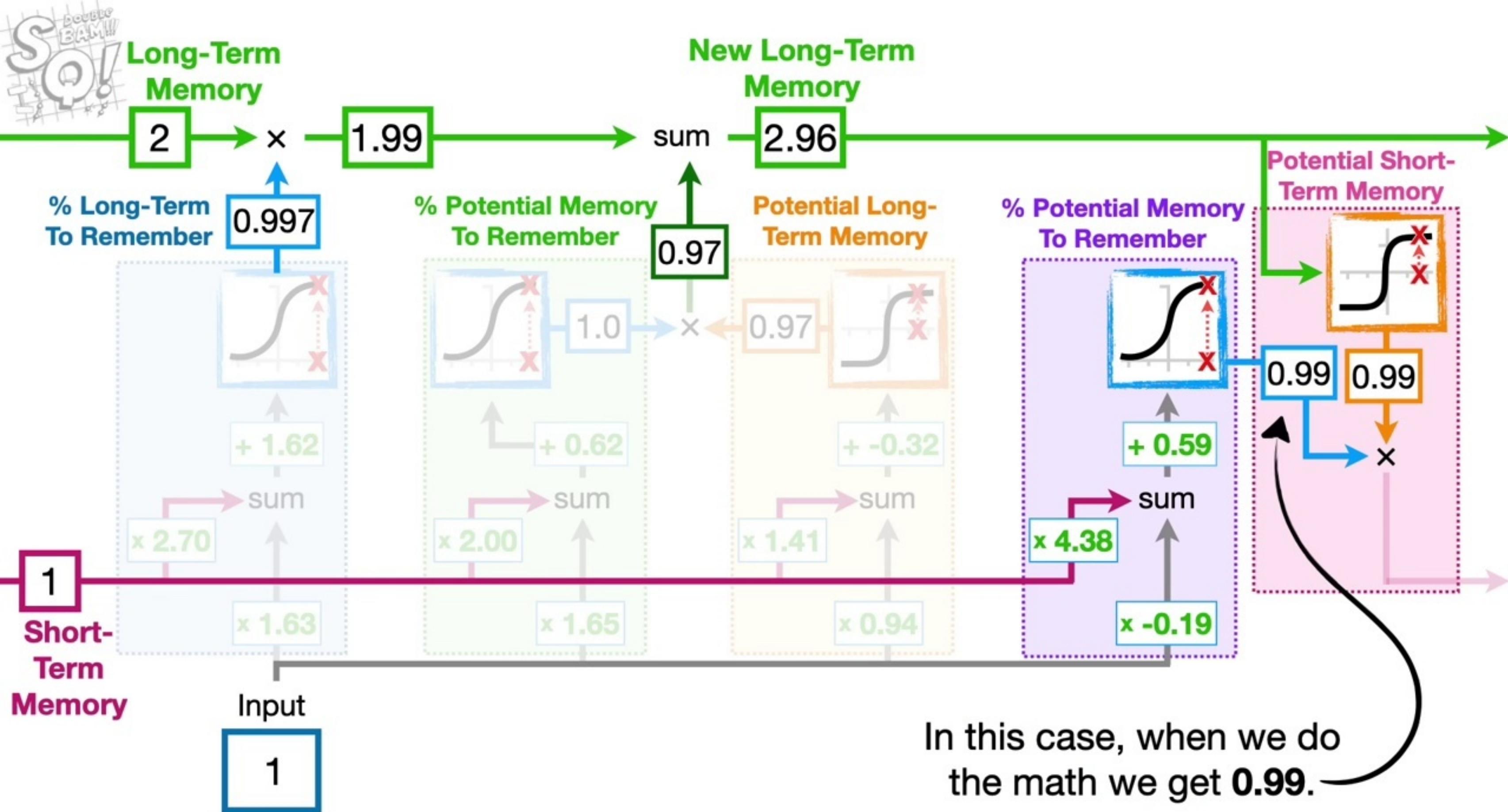
...and this is done using the exact same method we used two times earlier.

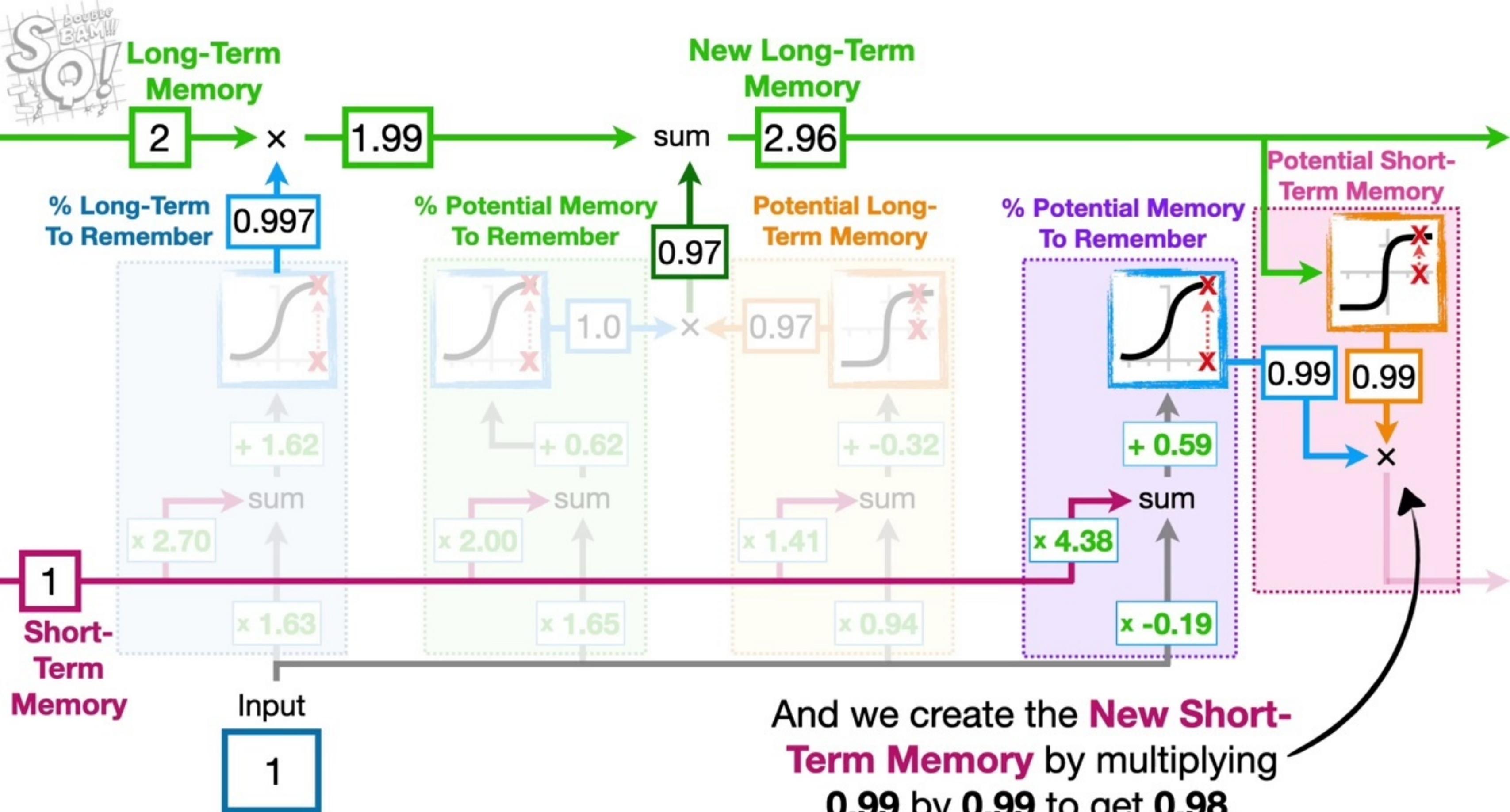


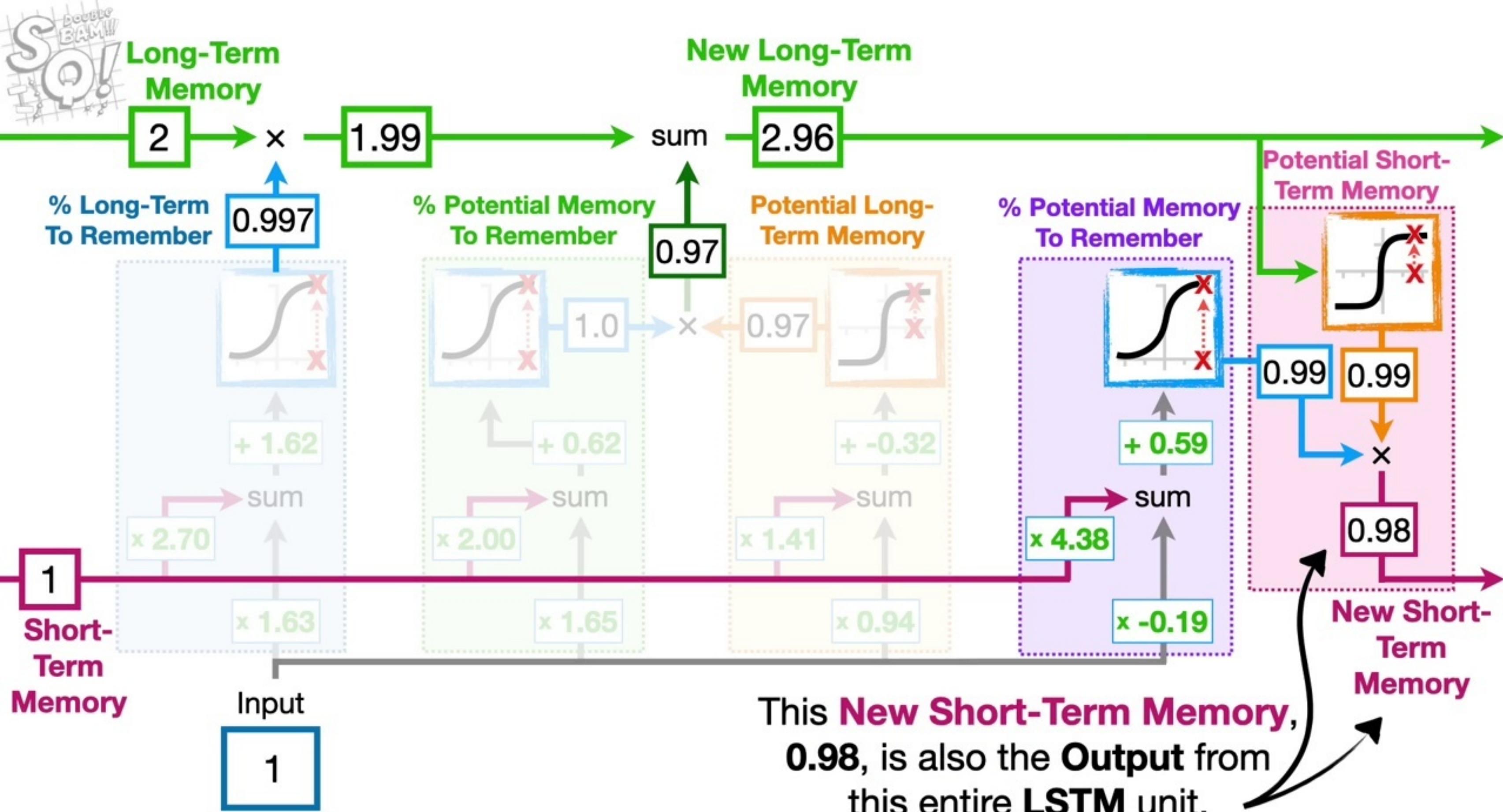


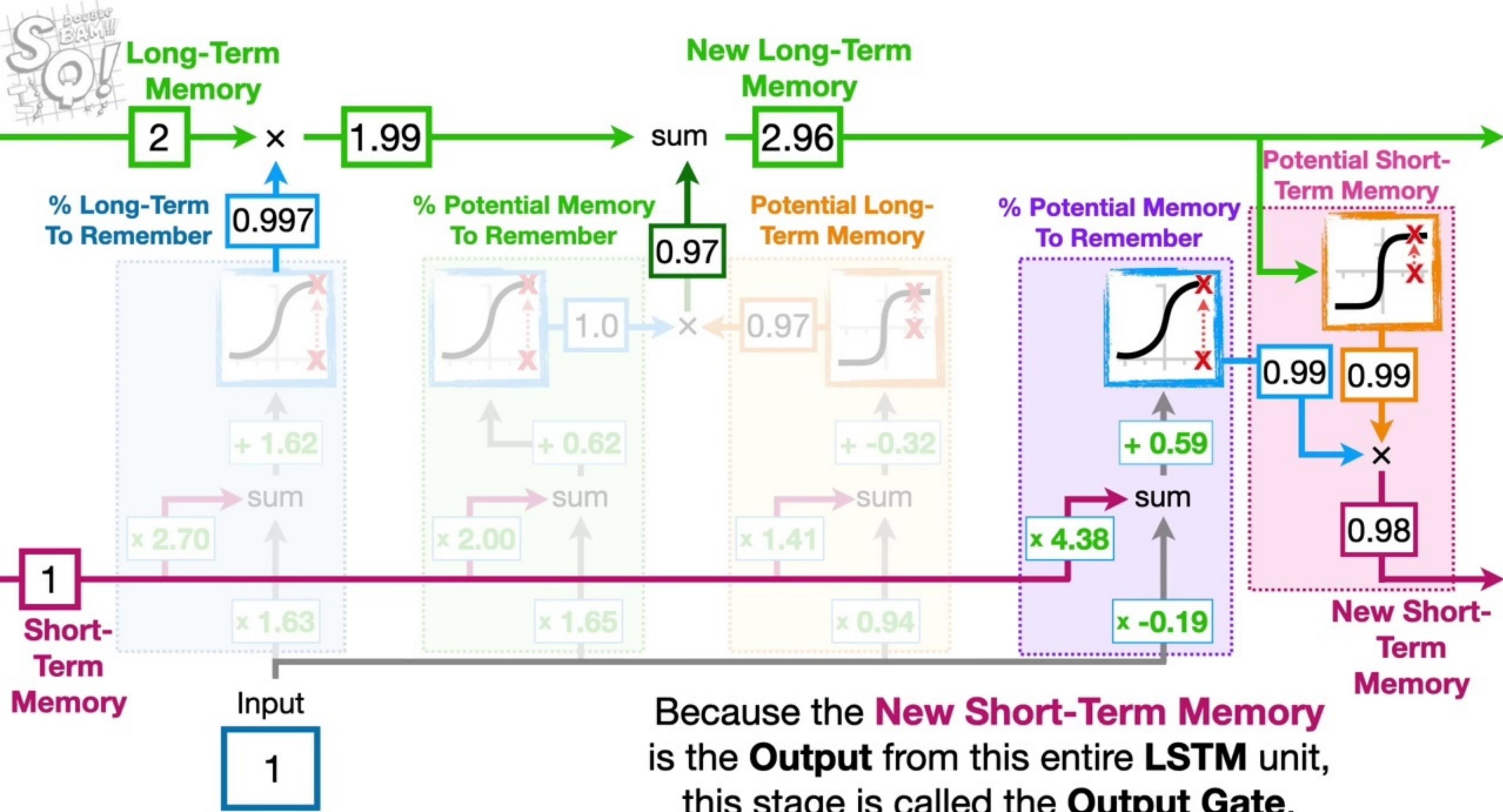


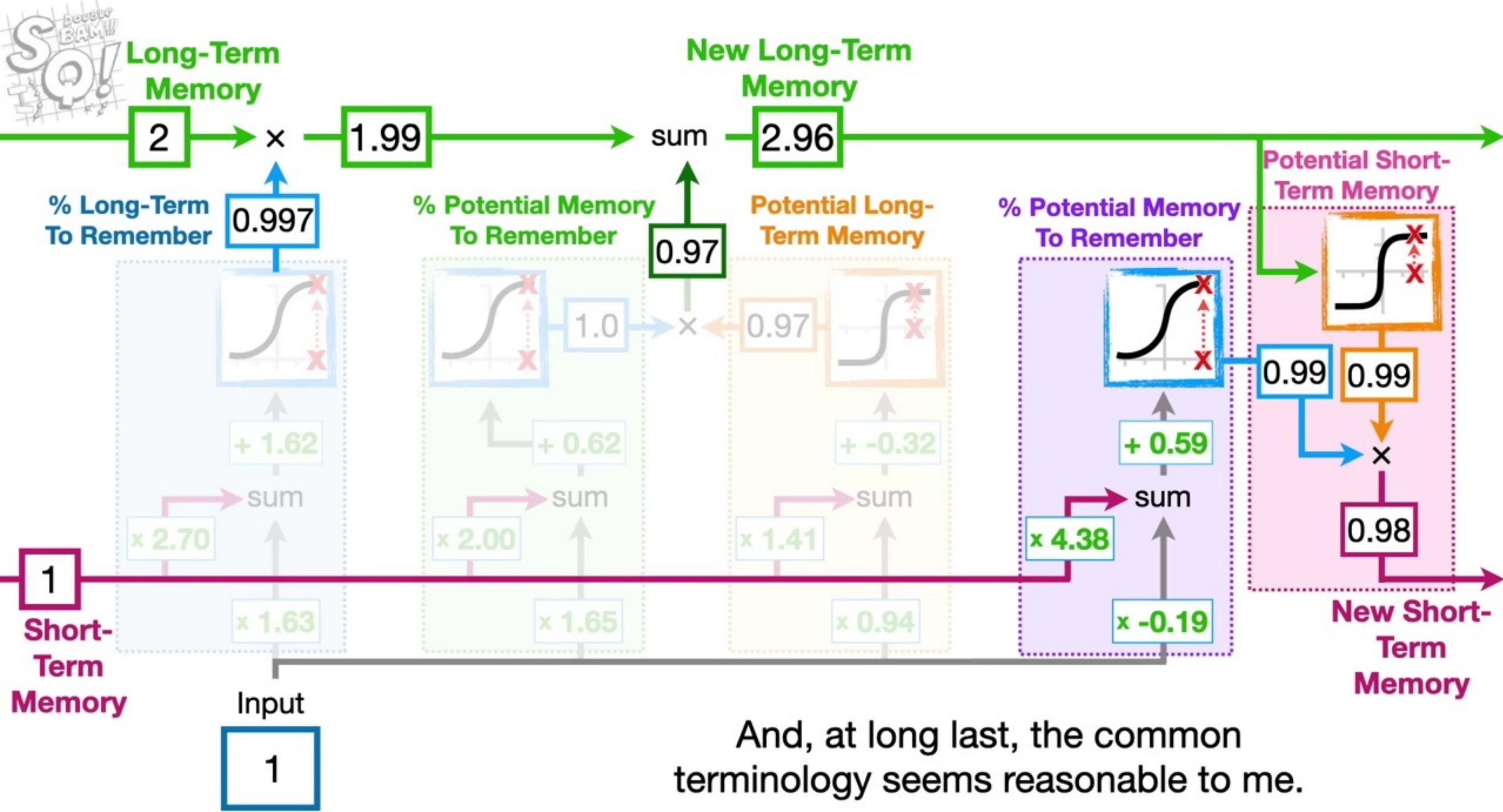
In all three cases, we use a **Sigmoid Activation Function** to determine what percent the **LSTM** remembers.

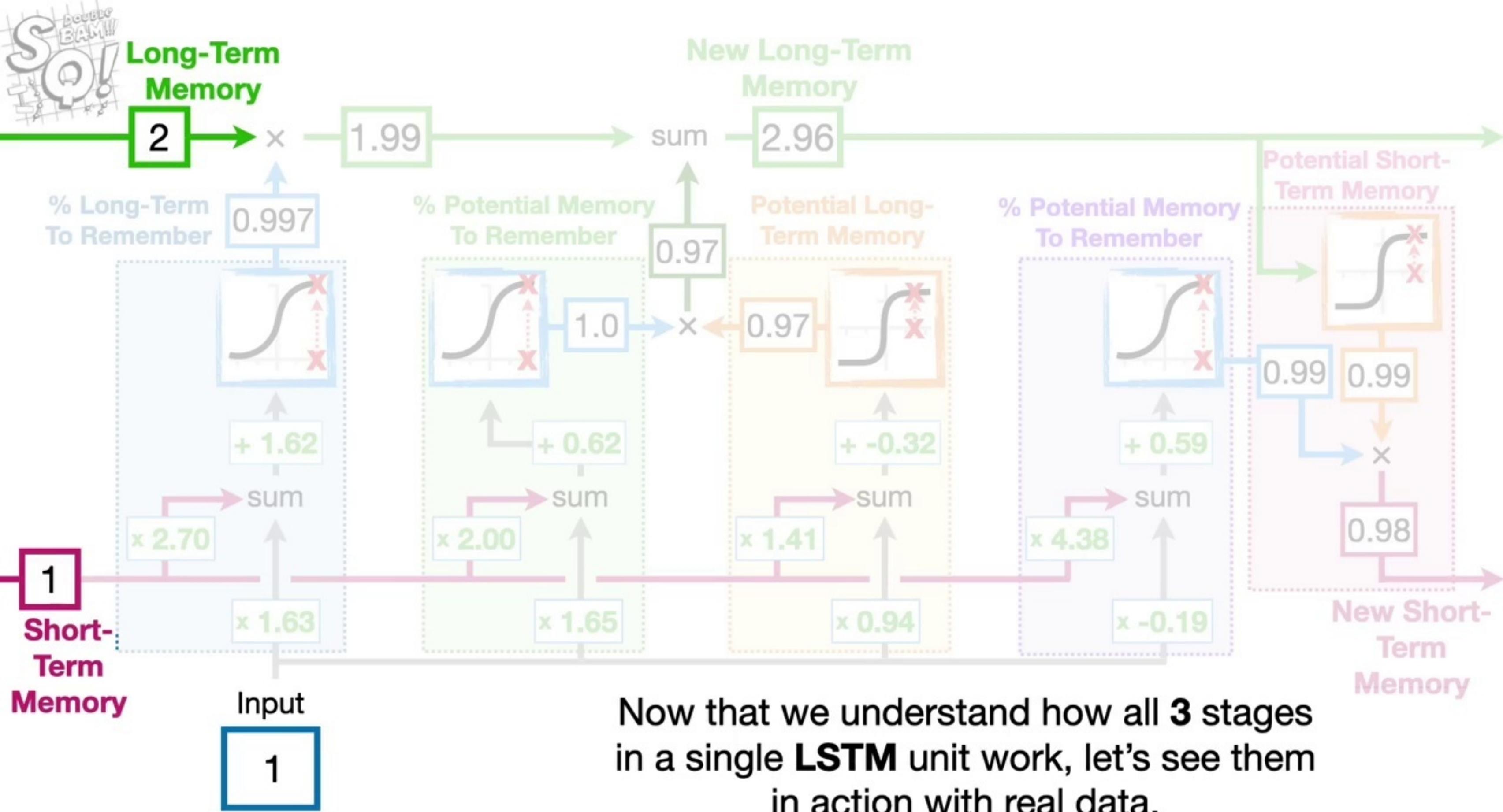






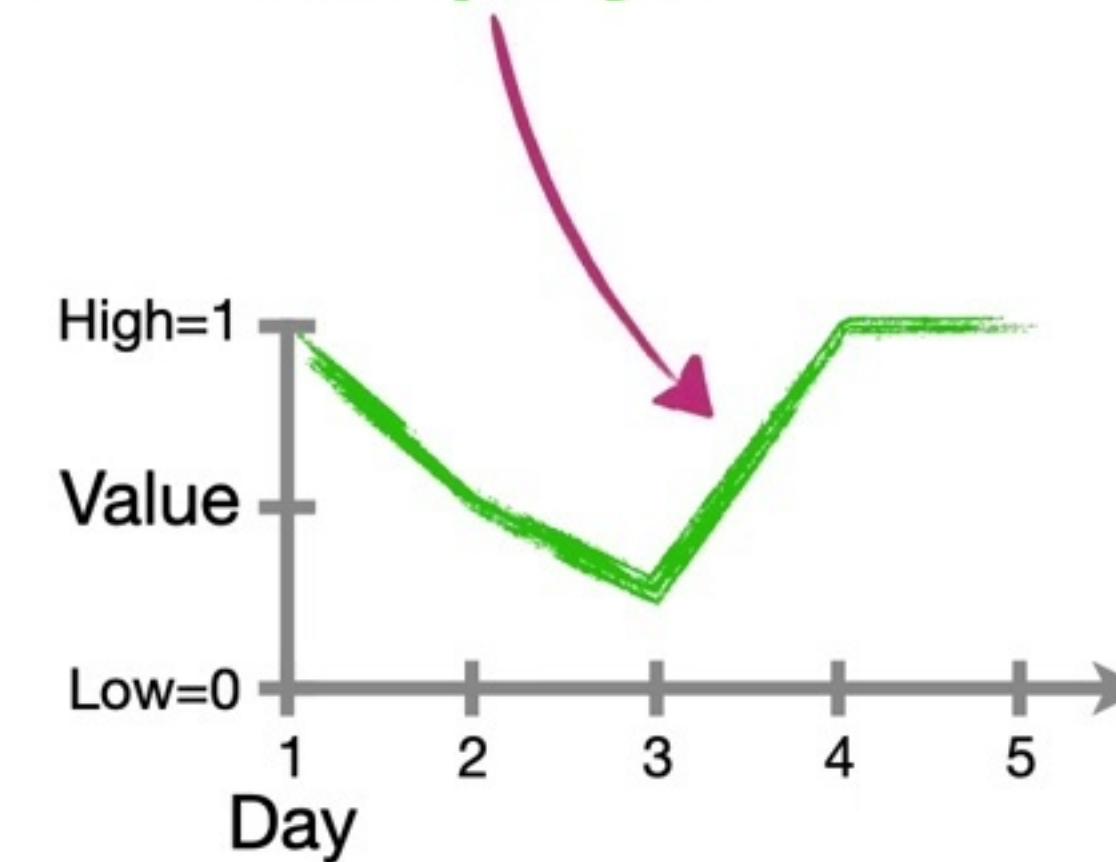
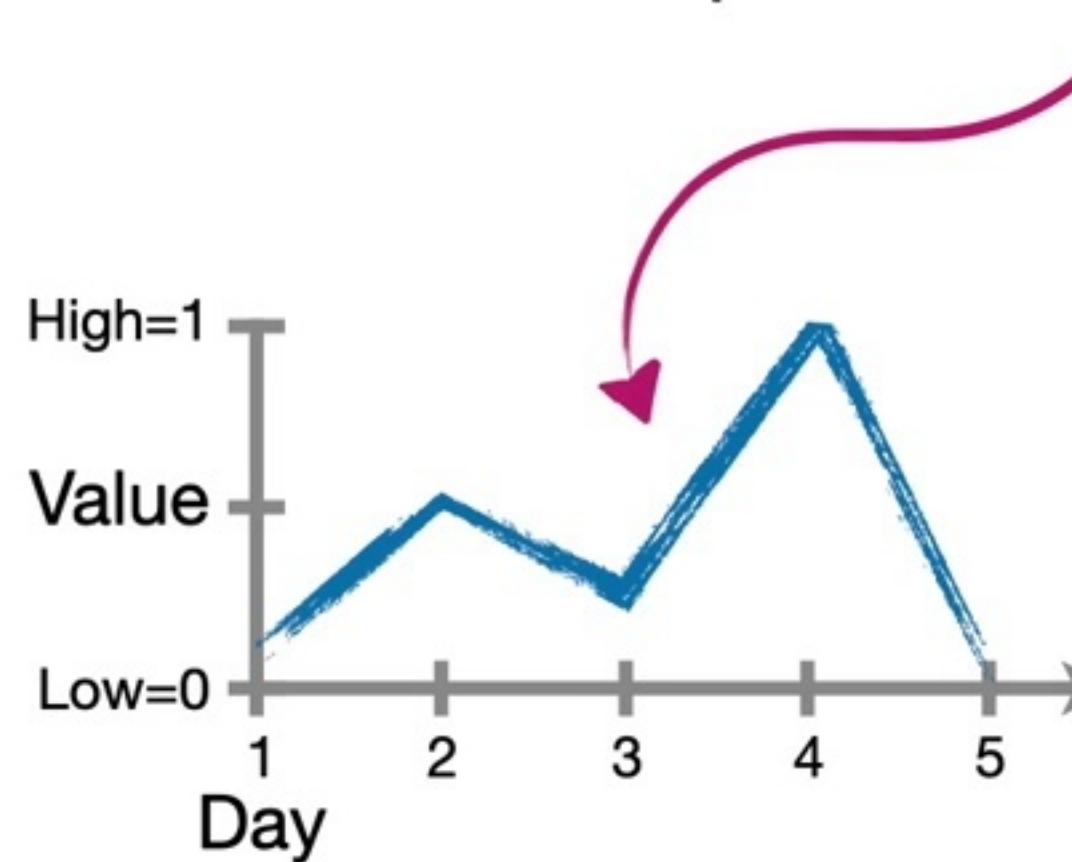








Here we have stock prices for two companies, **Company A** and **Company B**.

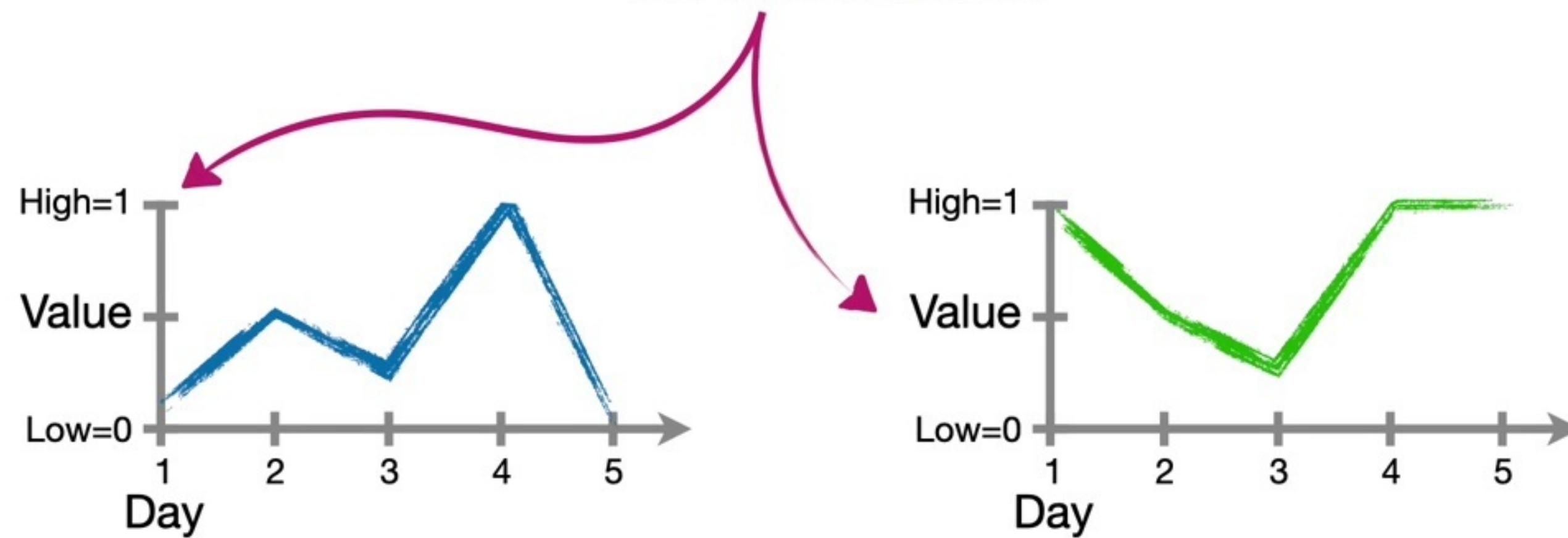


Company A =

Company B =



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

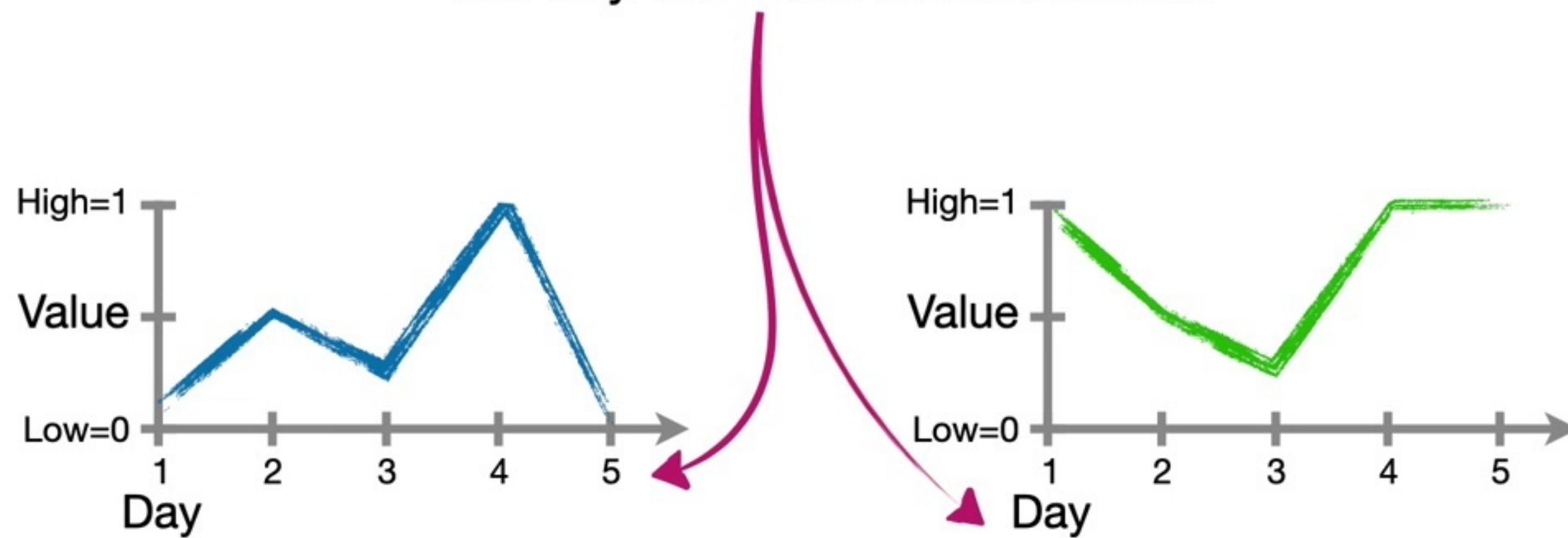


Company A =

Company B =



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

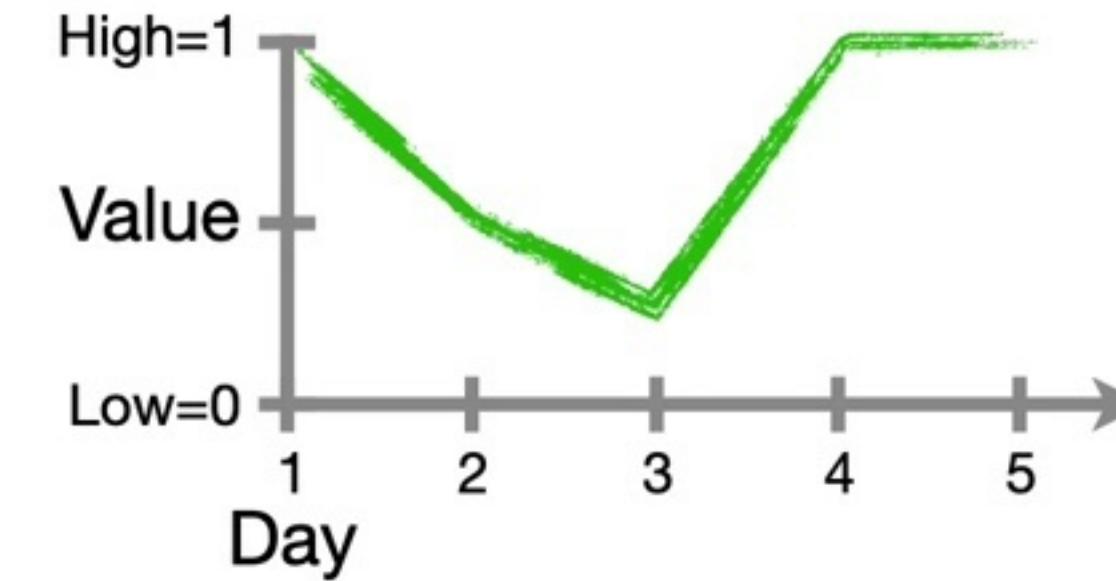
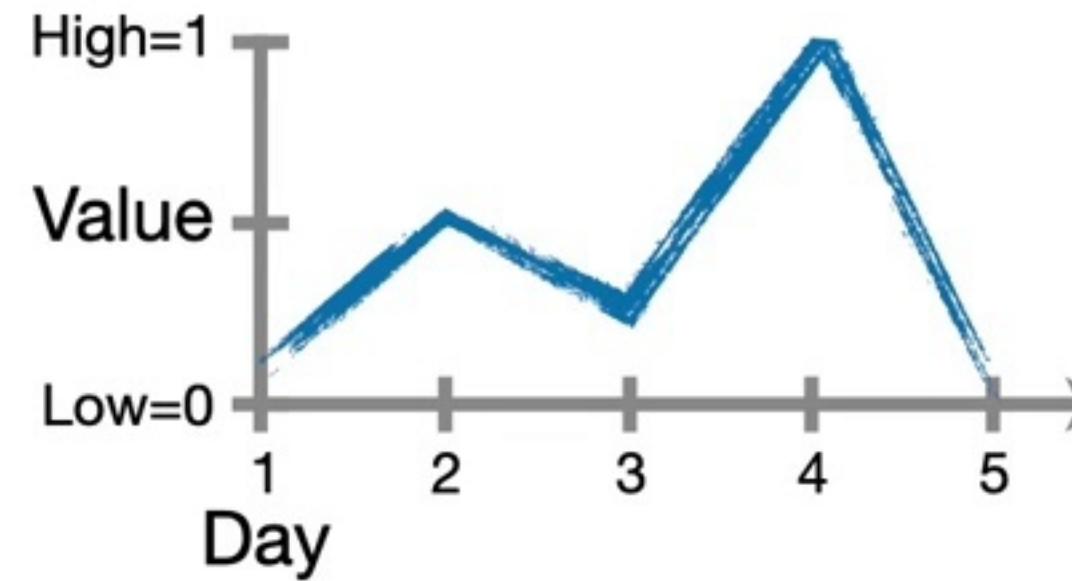


Company A =

Company B =



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

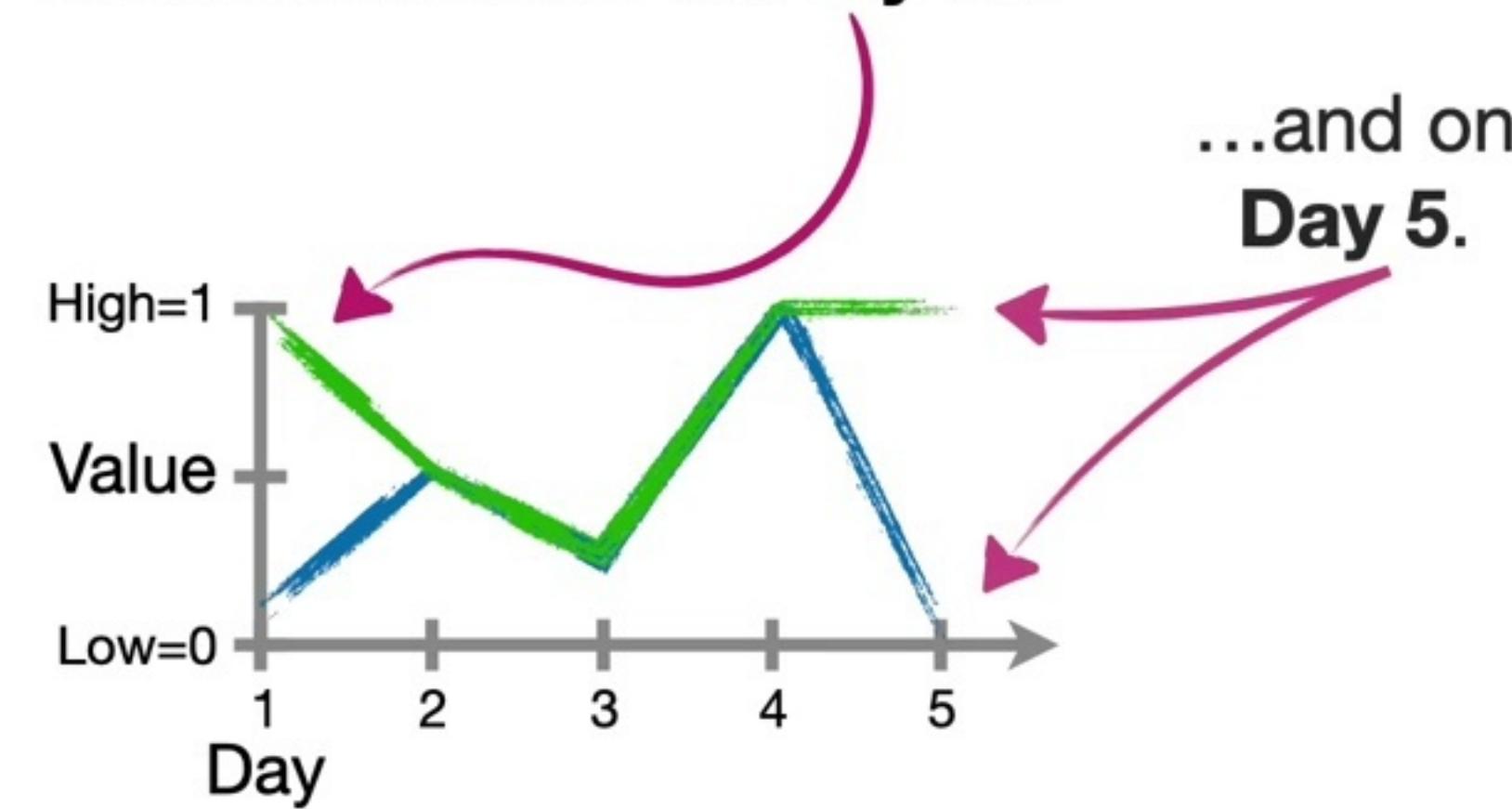


Company A =

Company B =



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

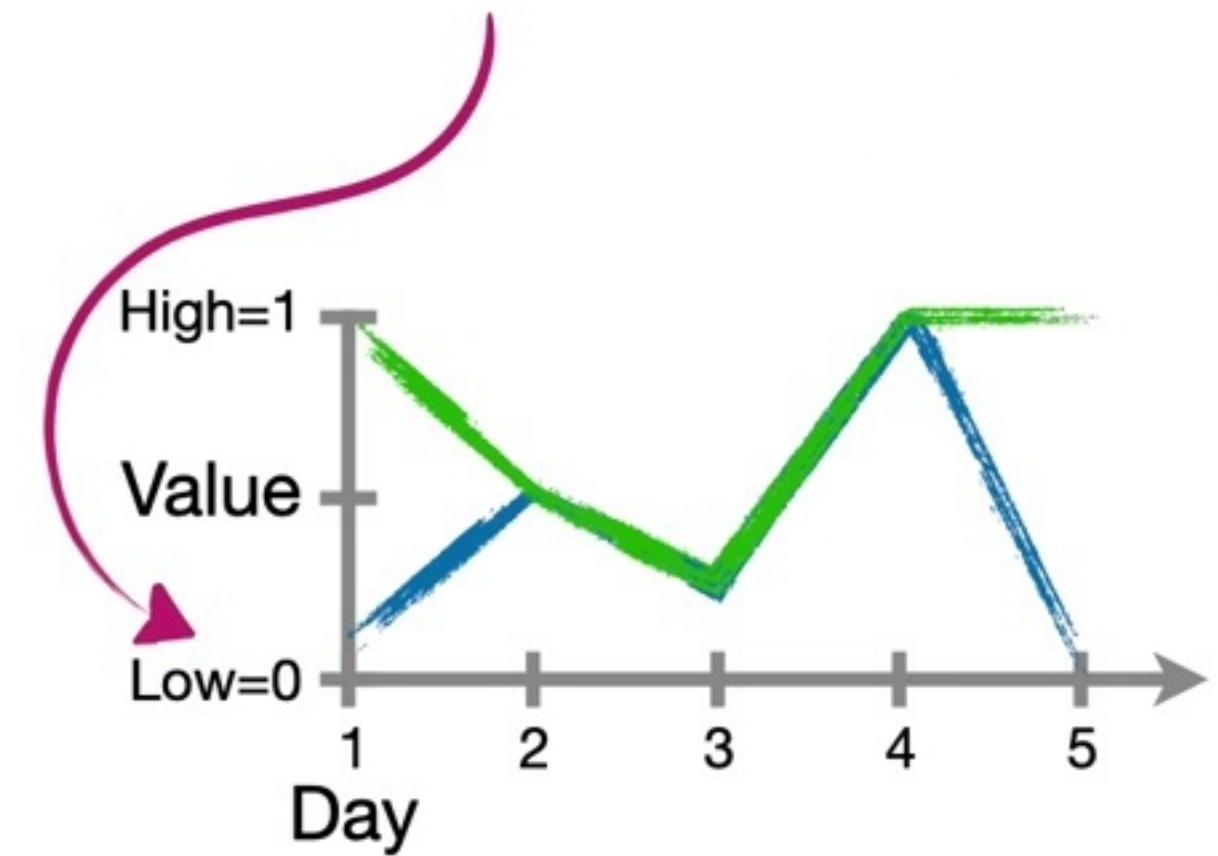


Company A =

Company B =



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

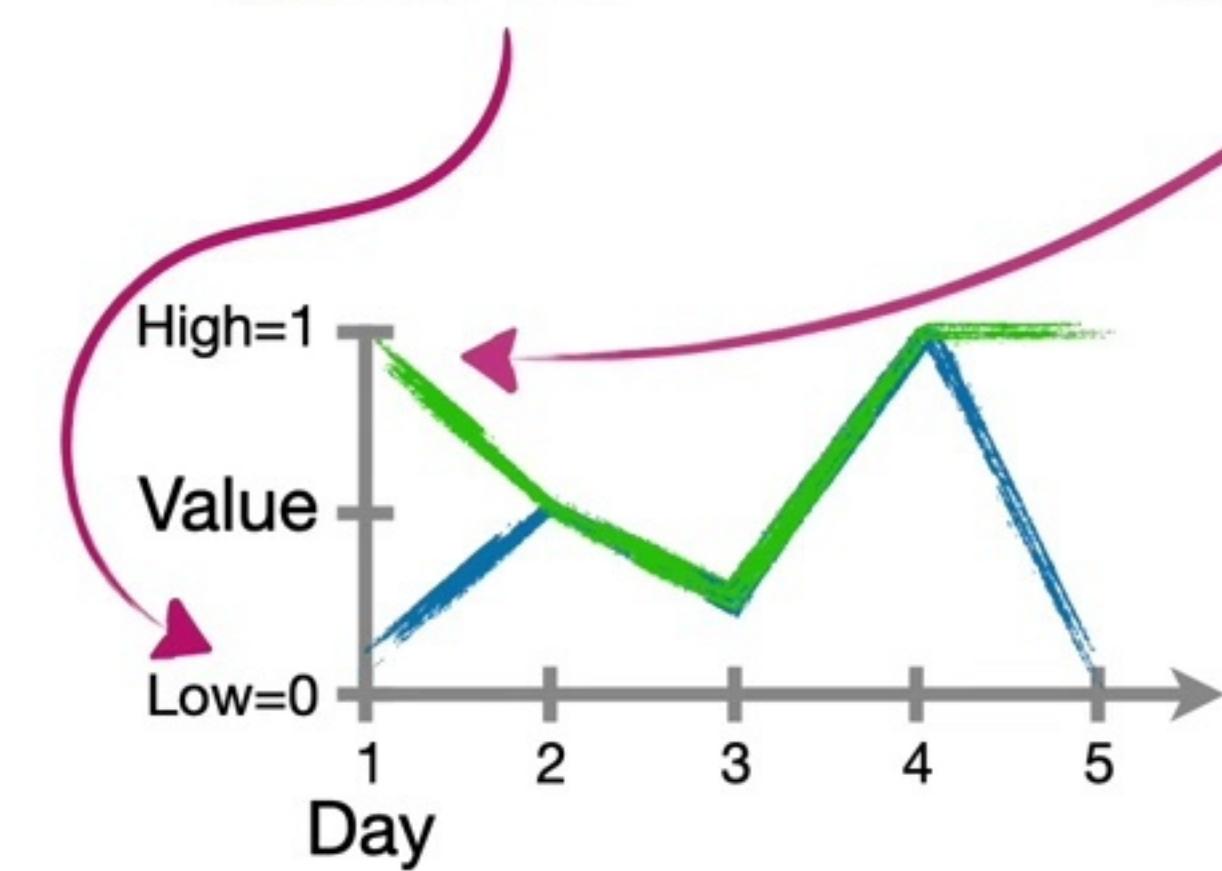


Company A =

Company B =



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

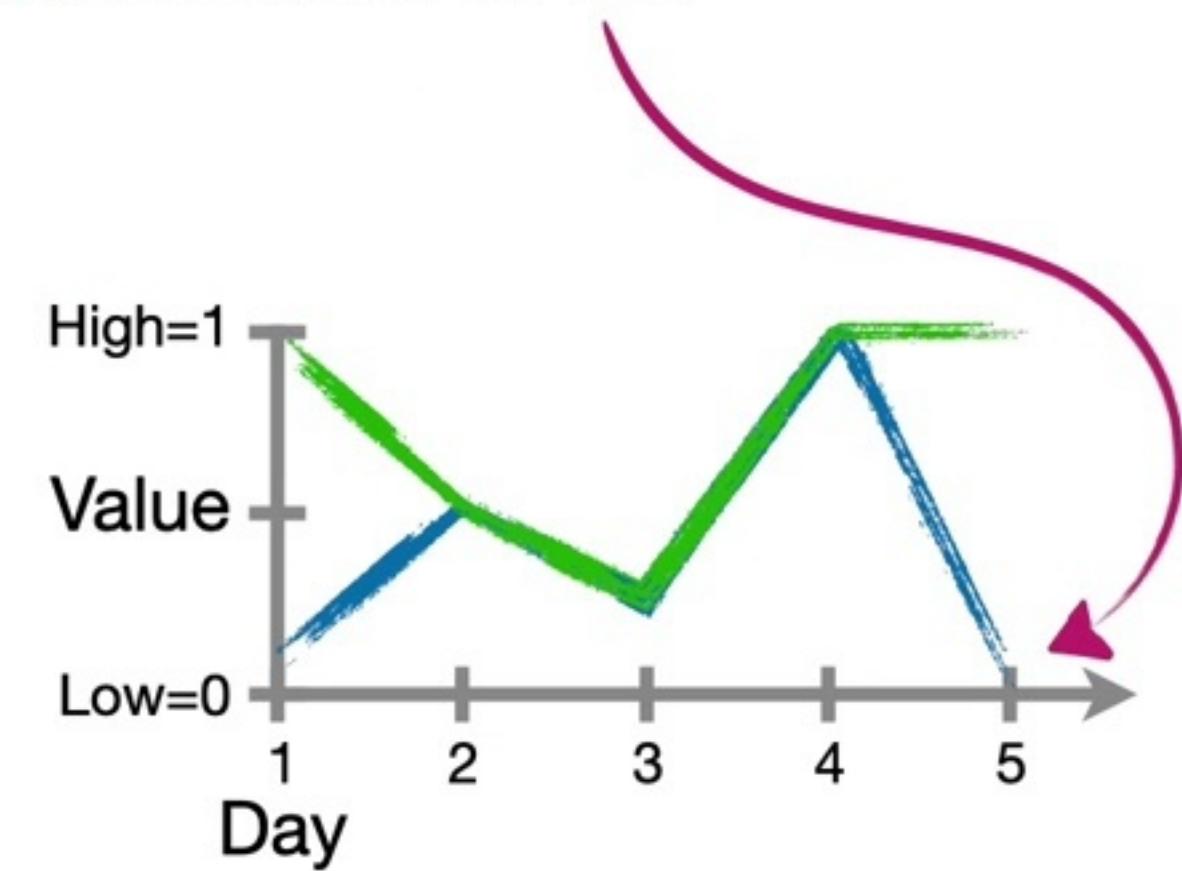


Company A =

Company B =



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

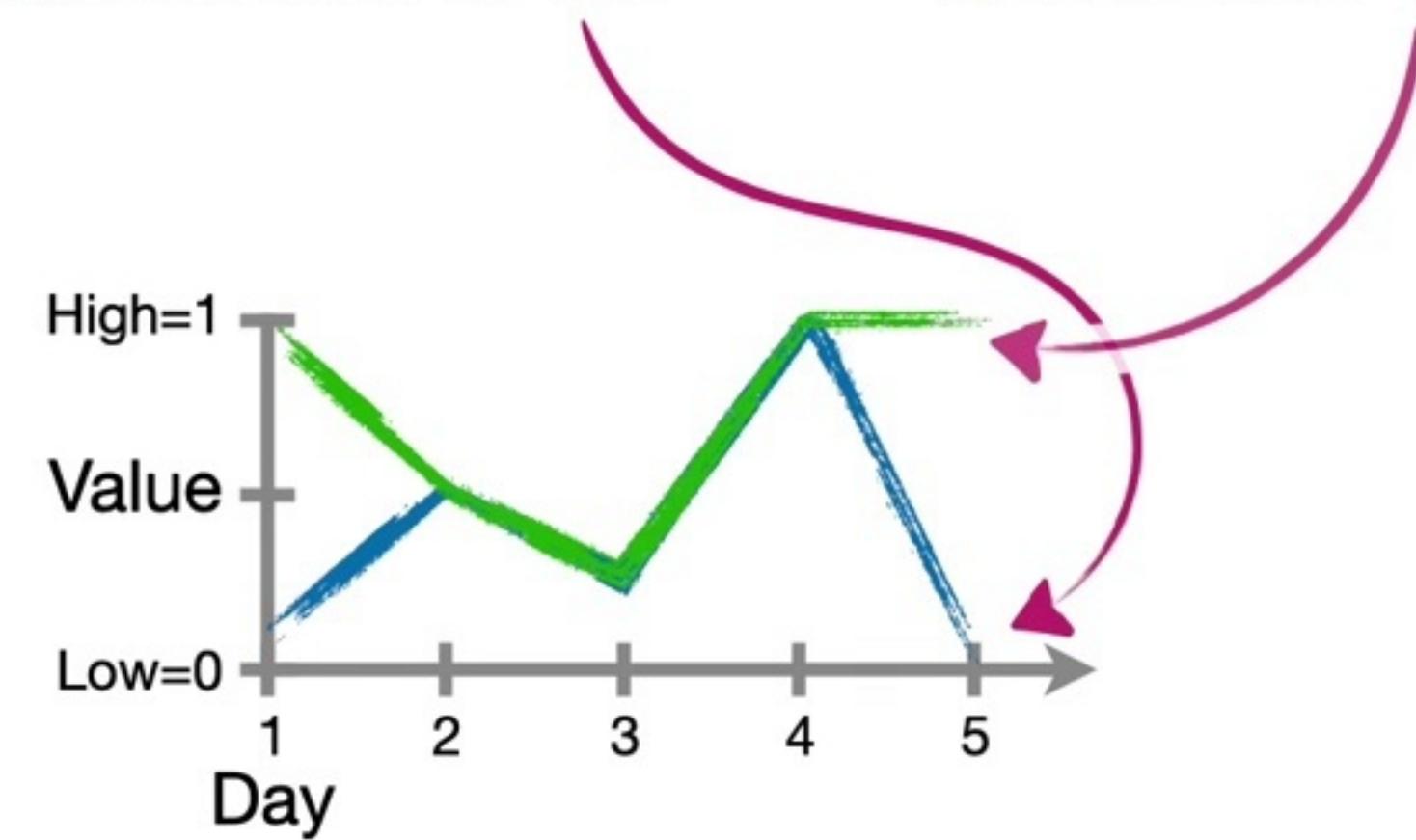


Company A =

Company B =



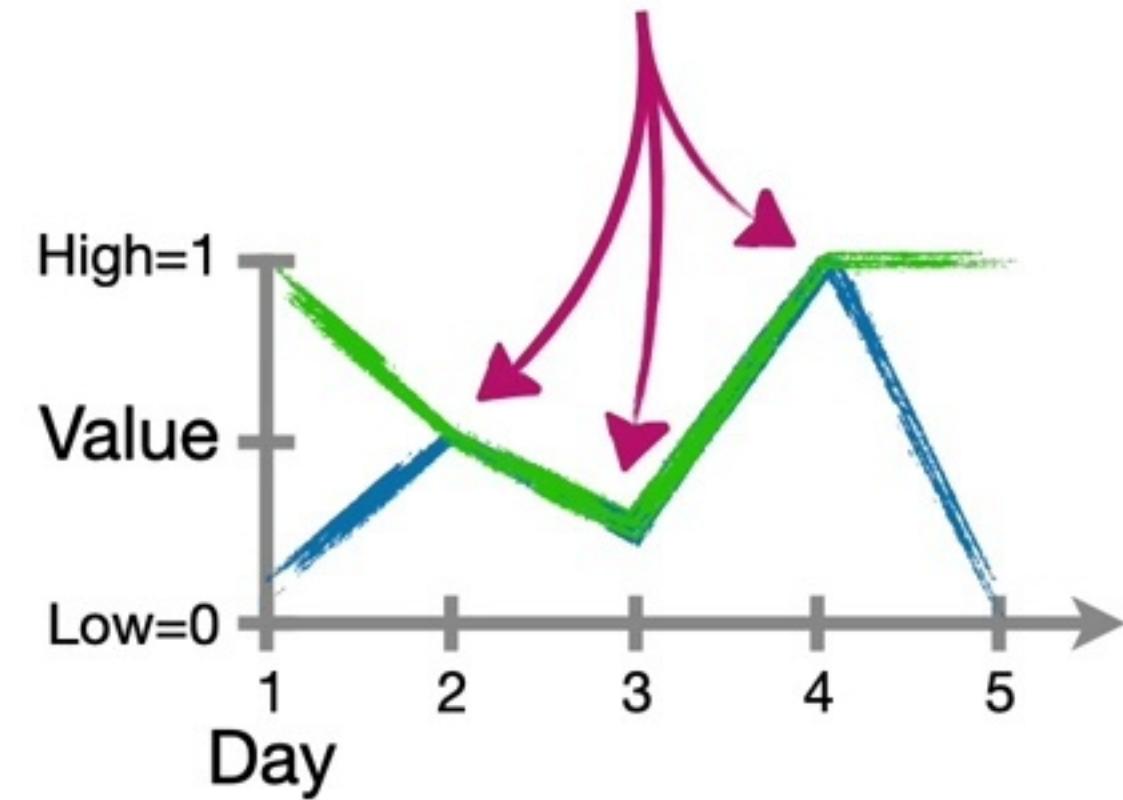
And on **Day 5**, **Company A** is returns to 0...
...and **Company B** returns to 1.



Company A =
Company B =



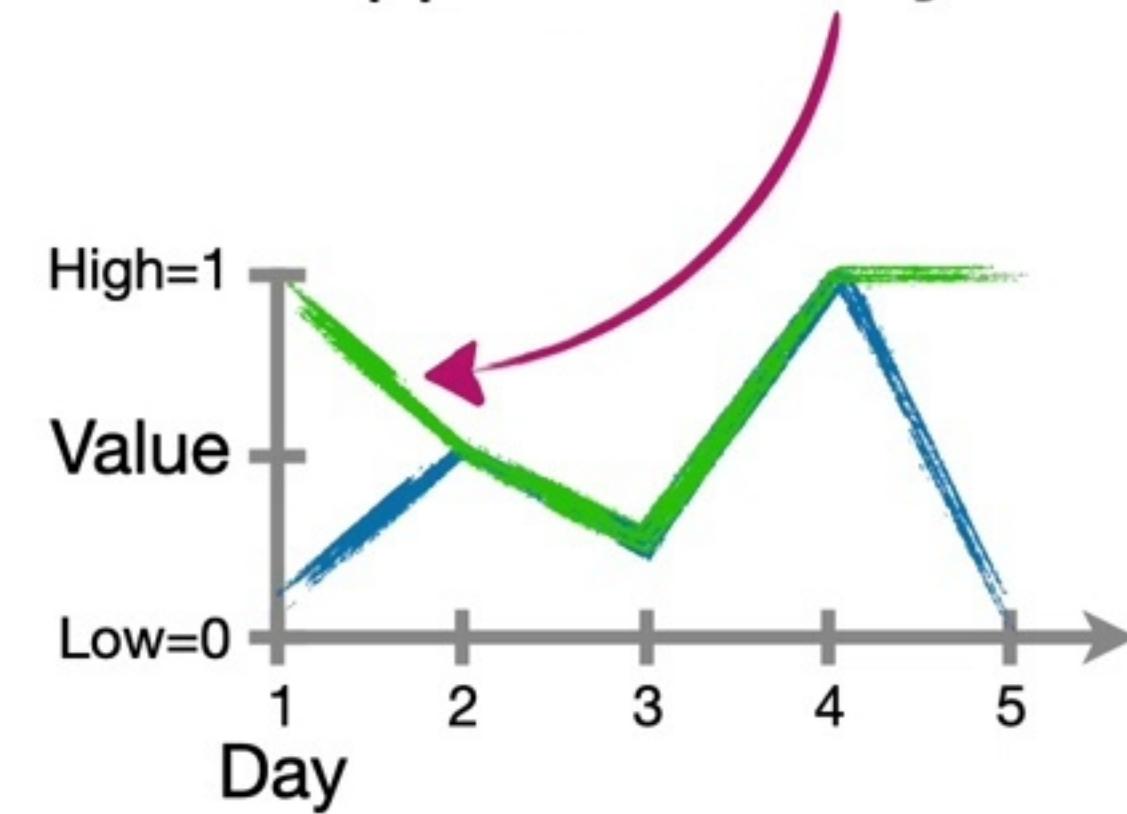
On all of the other days, **Days 2, 3 and 4**, both companies have the exact same values.



Company A =
Company B =



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

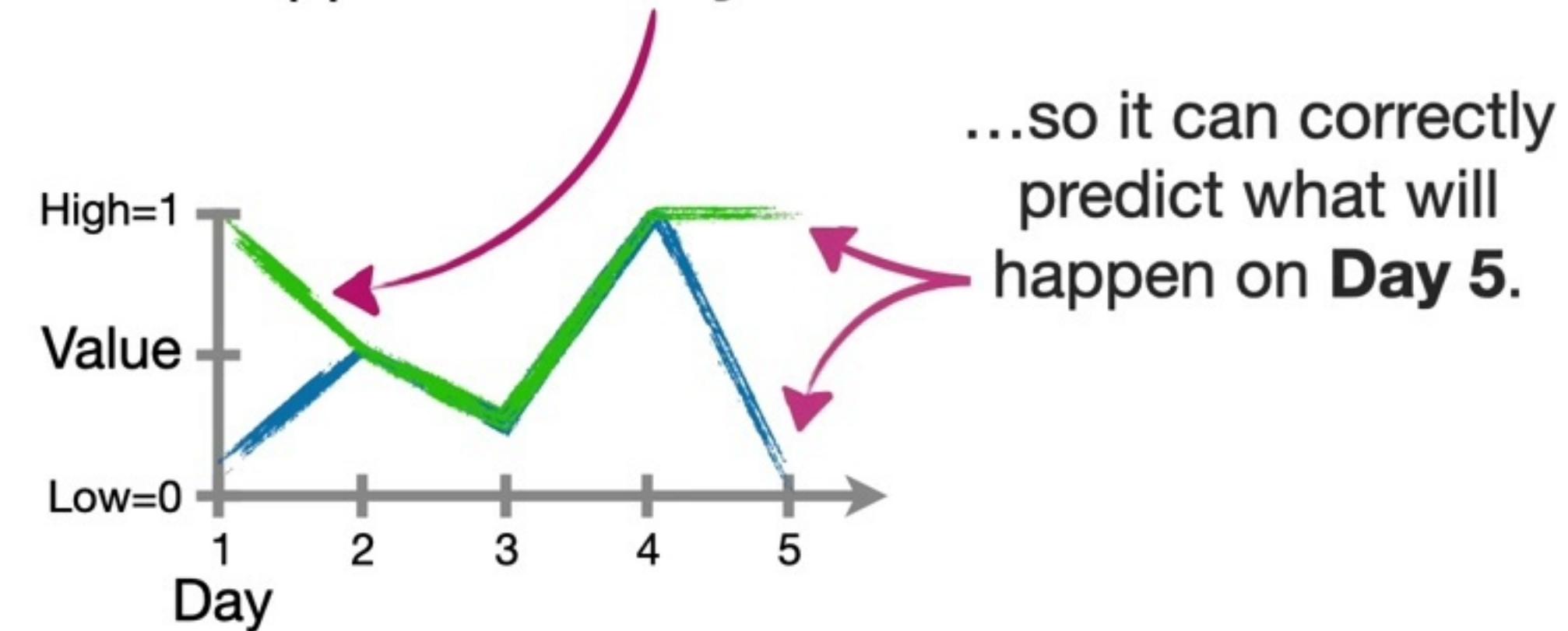


Company A = /checkmark

Company B = /checkmark



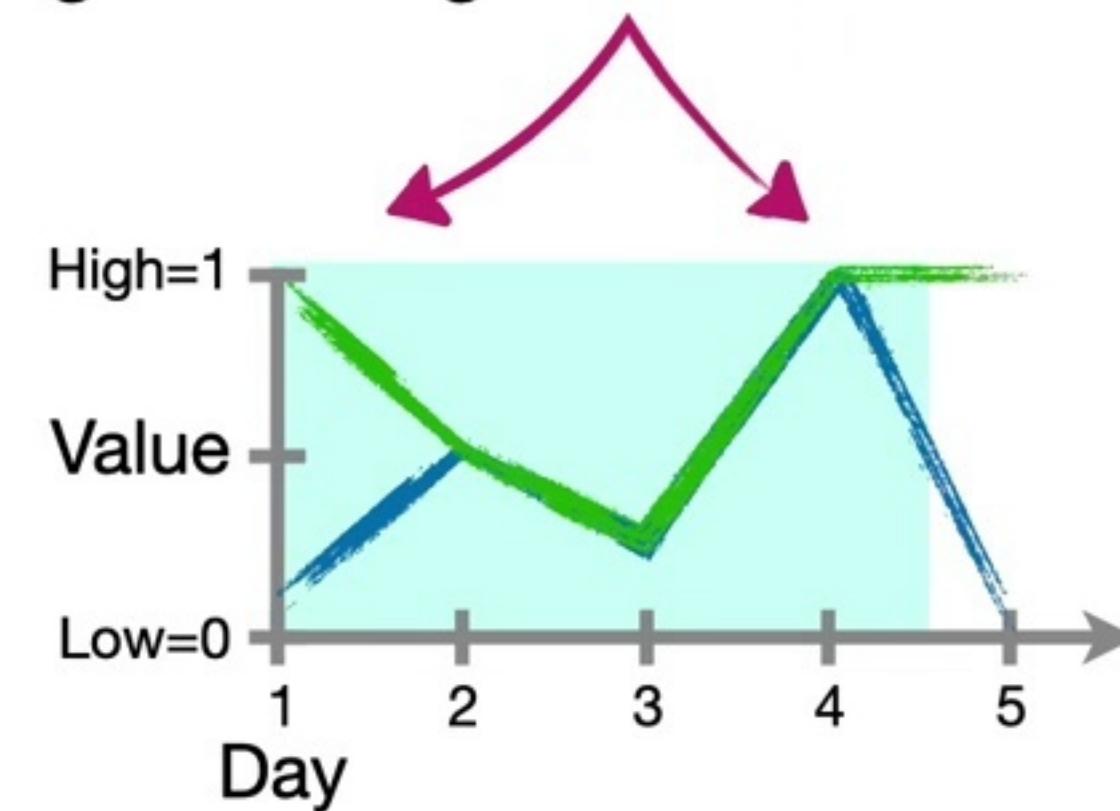
Given this sequential data, we want the **LSTM** to remember what happened on **Day 1**...



Company A =
Company B =



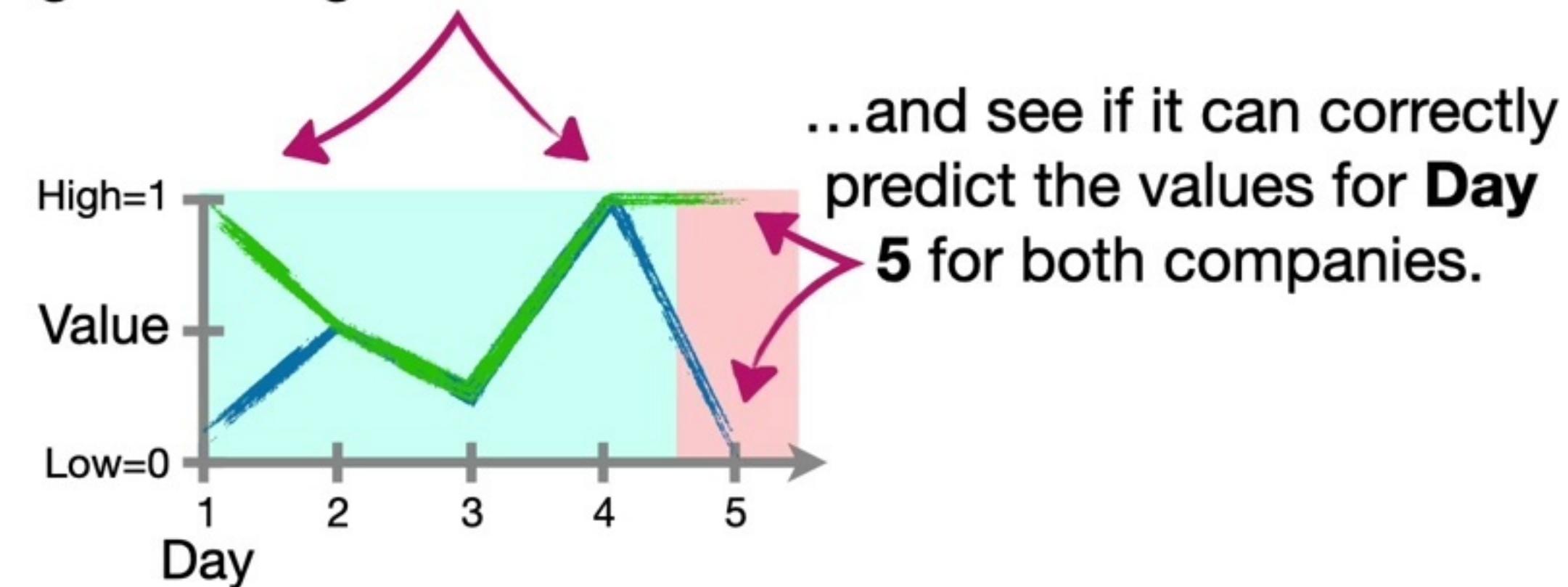
In other words, we're going to sequentially run the data from **Days 1** through **4** through an unrolled **LSTM**...



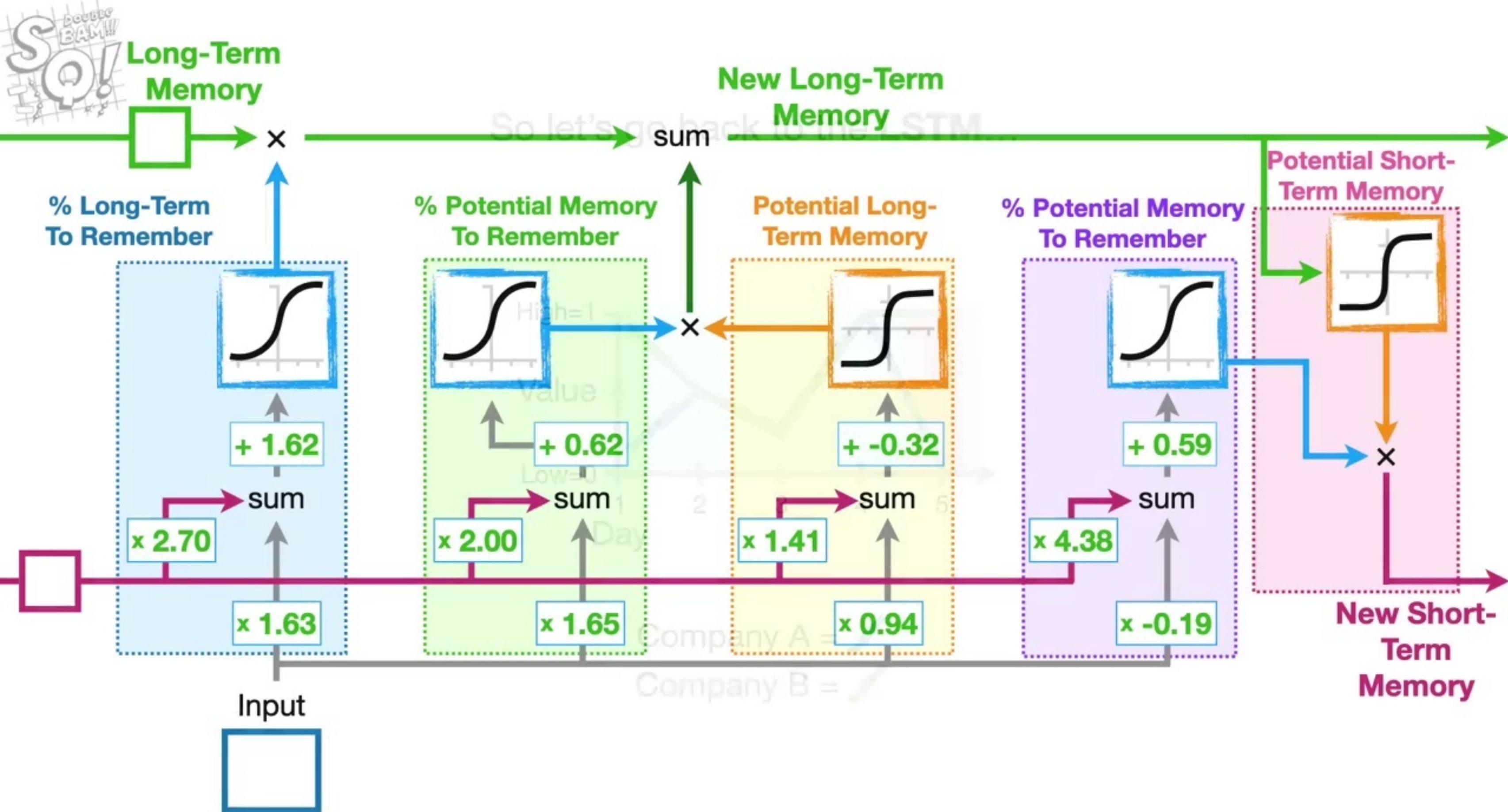
Company A = 
Company B = 

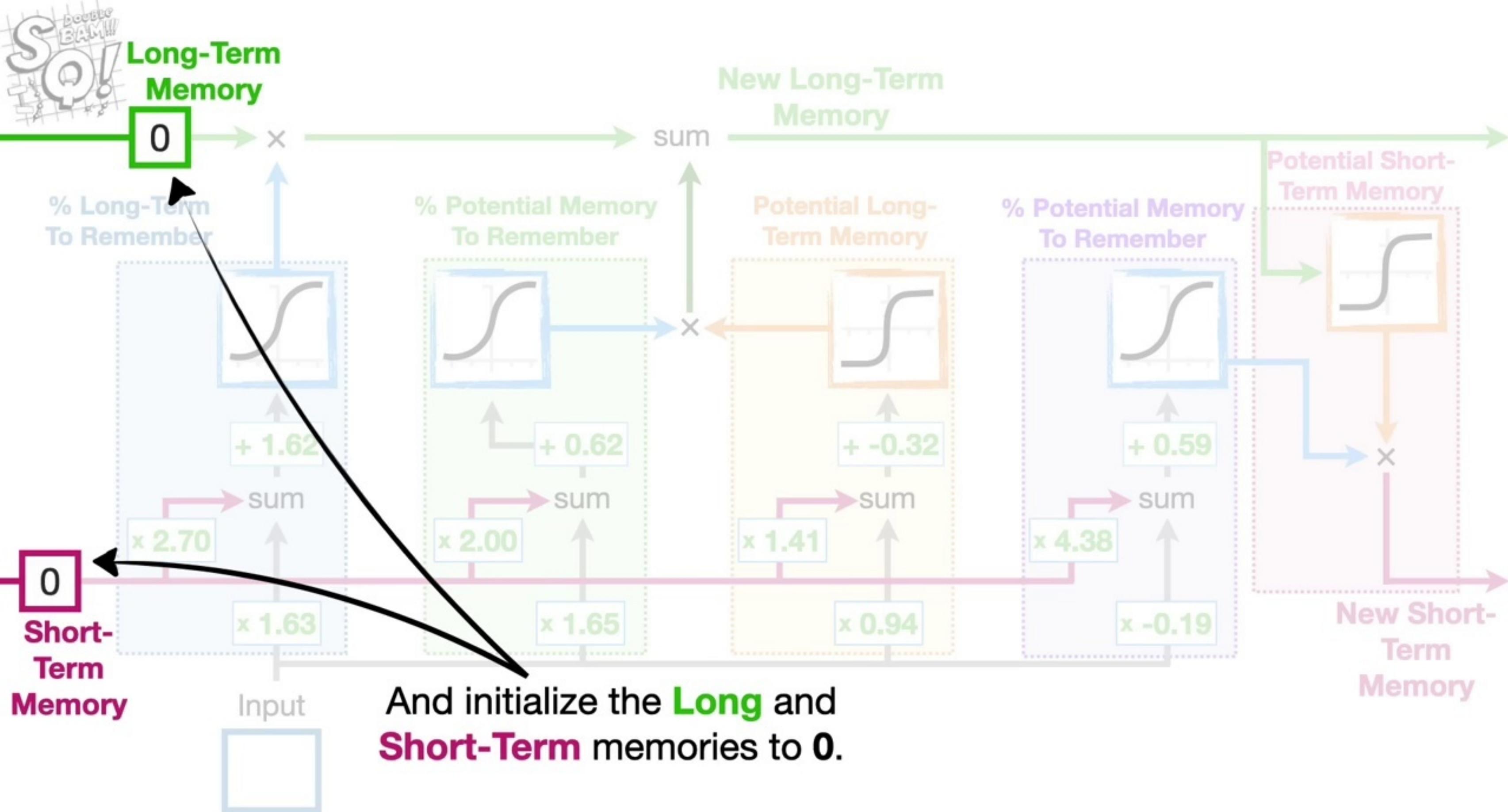


In other words, we're going to sequentially run the data from **Days 1** through **4** through an unrolled **LSTM**...

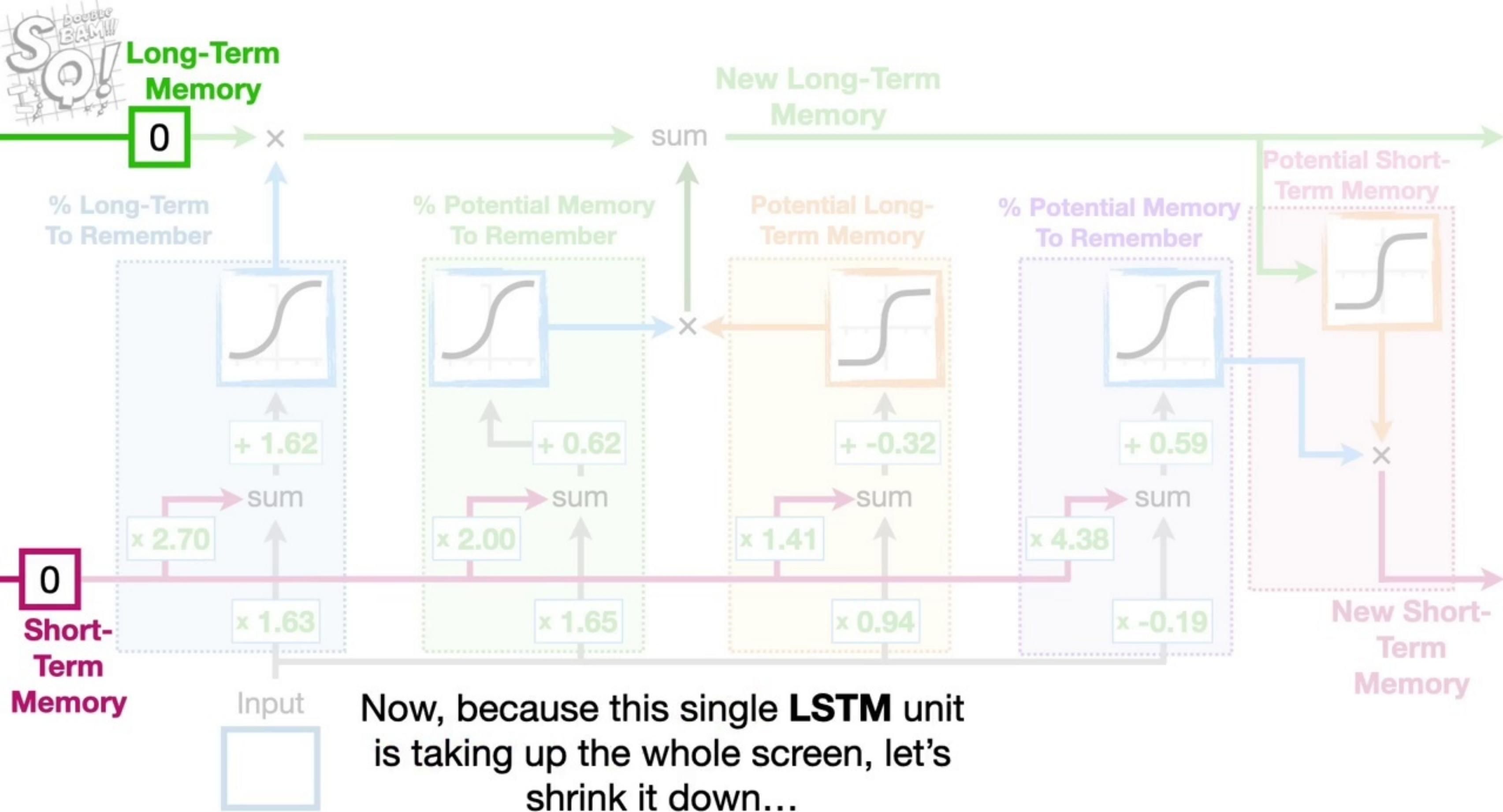


Company A = 
Company B = 





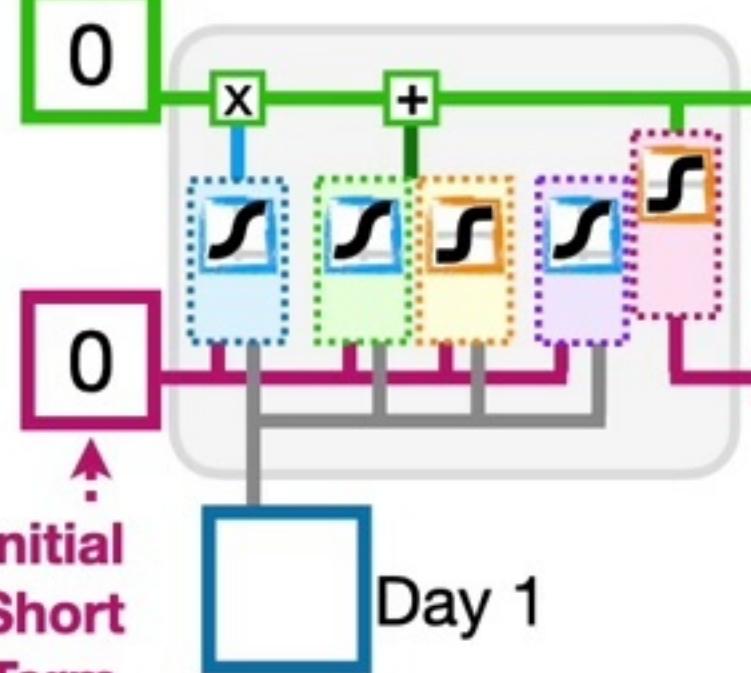
And initialize the **Long** and **Short-Term** memories to 0.



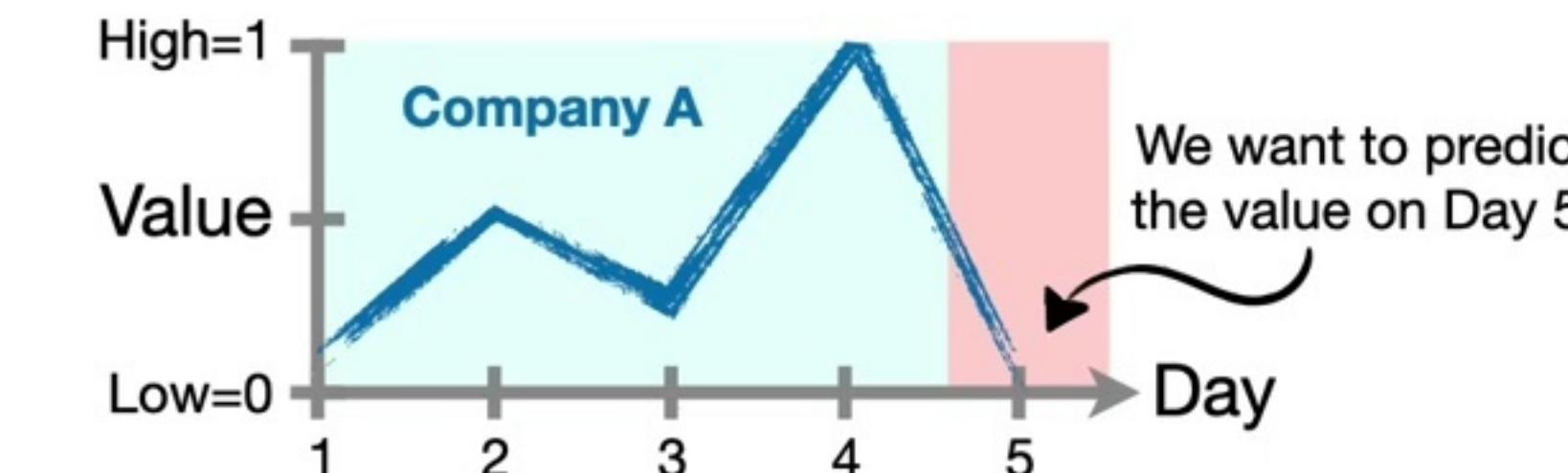
Now, because this single **LSTM** unit
is taking up the whole screen, let's
shrink it down...



Initial
Long Term
Memory



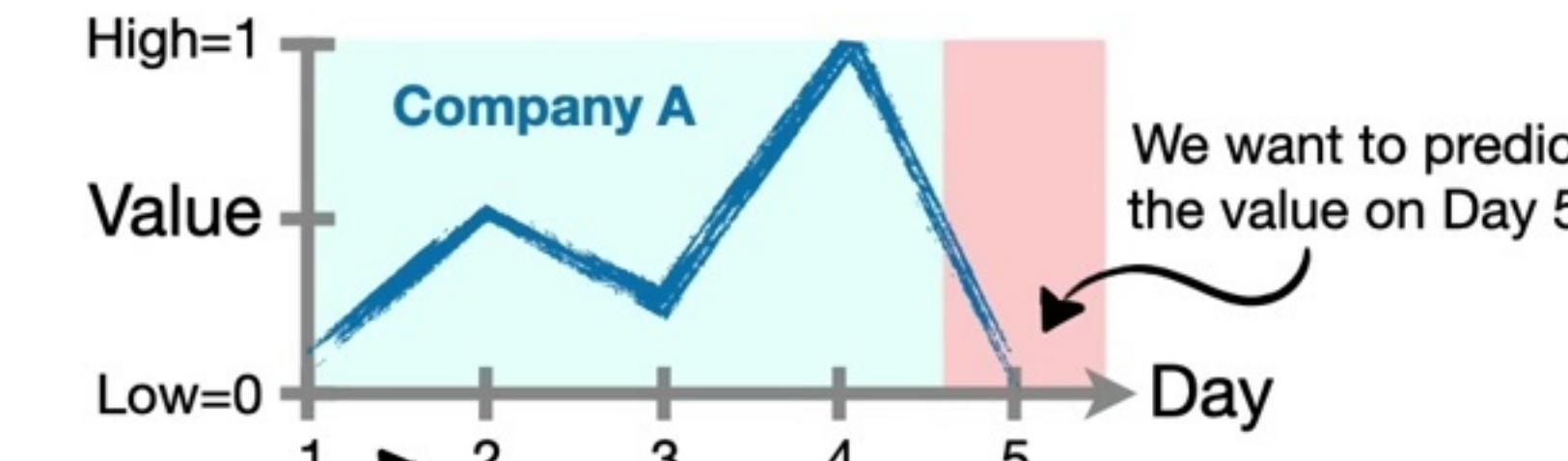
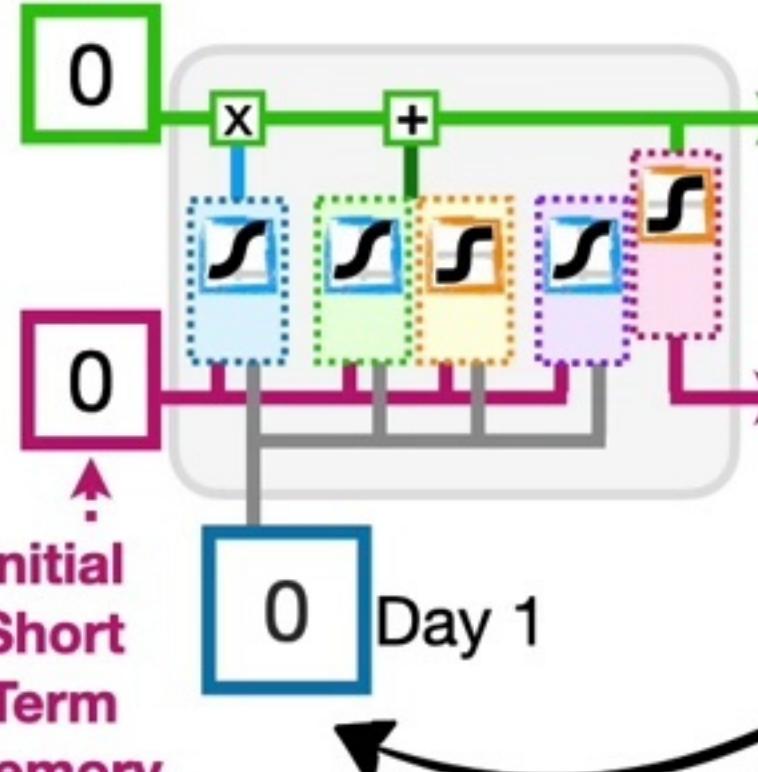
Initial
Short
Term
Memory



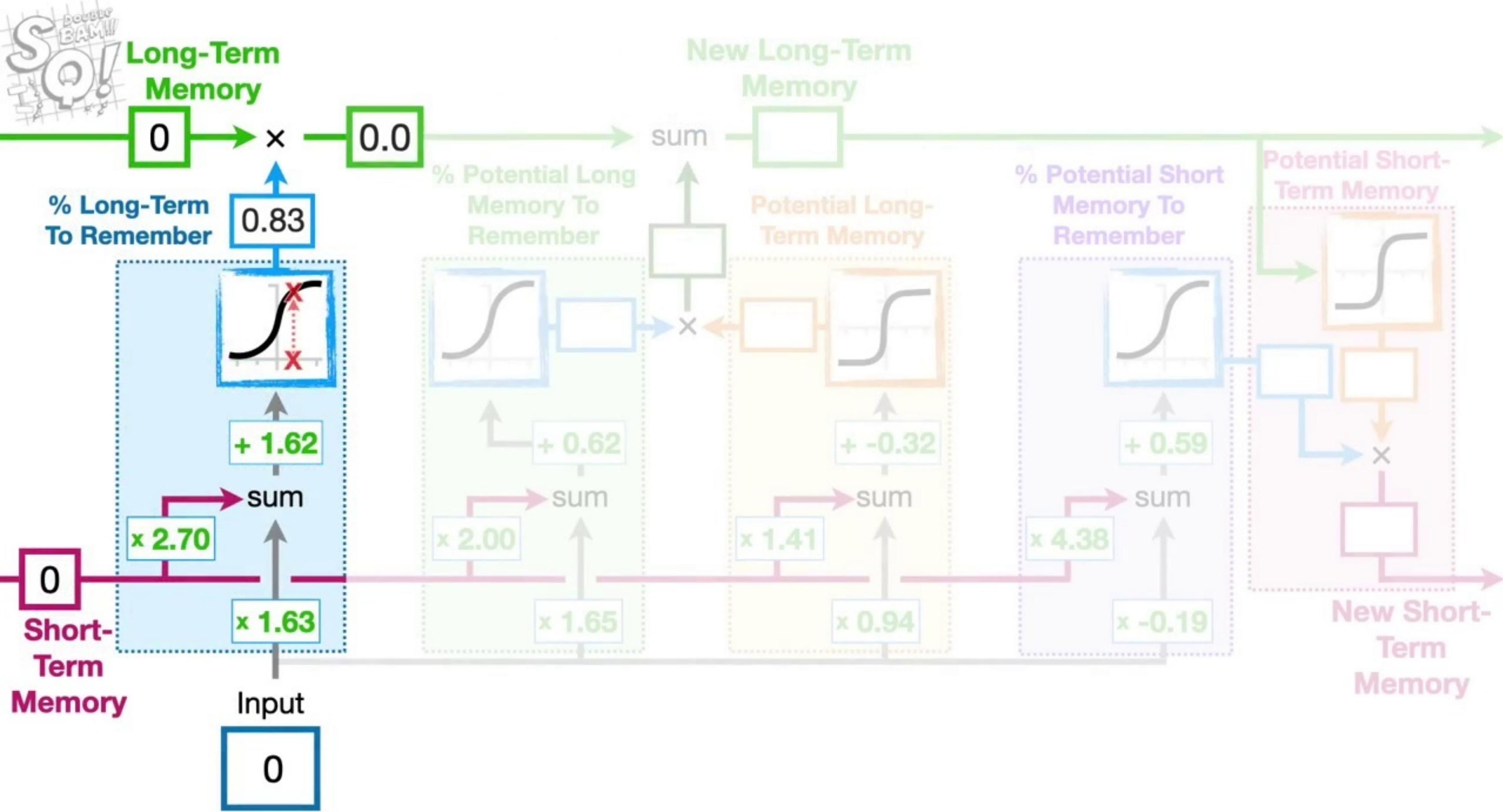
Now, if we want to sequentially run
Company A's values from Days 1
through 4 through this **LSTM**...

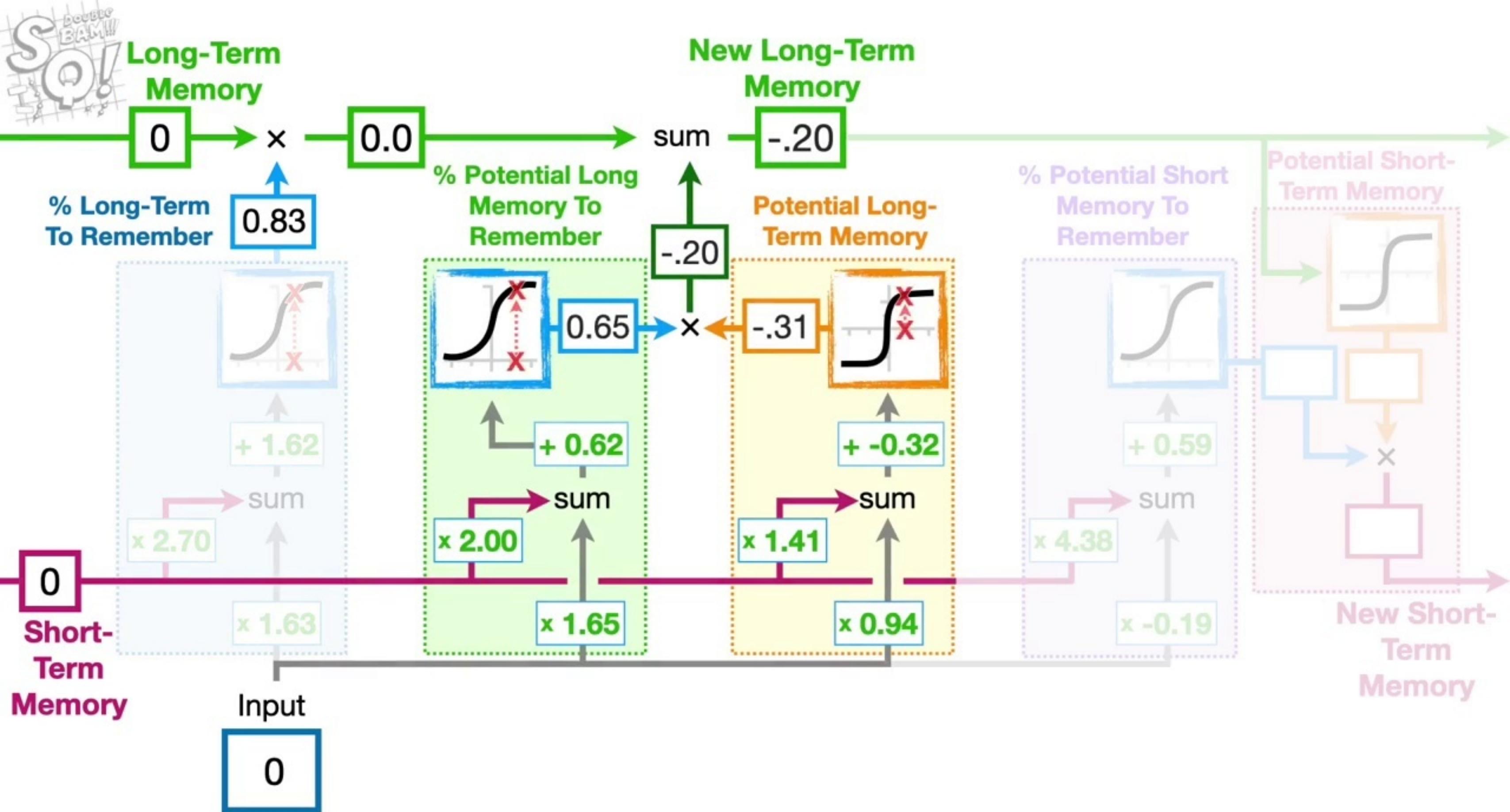


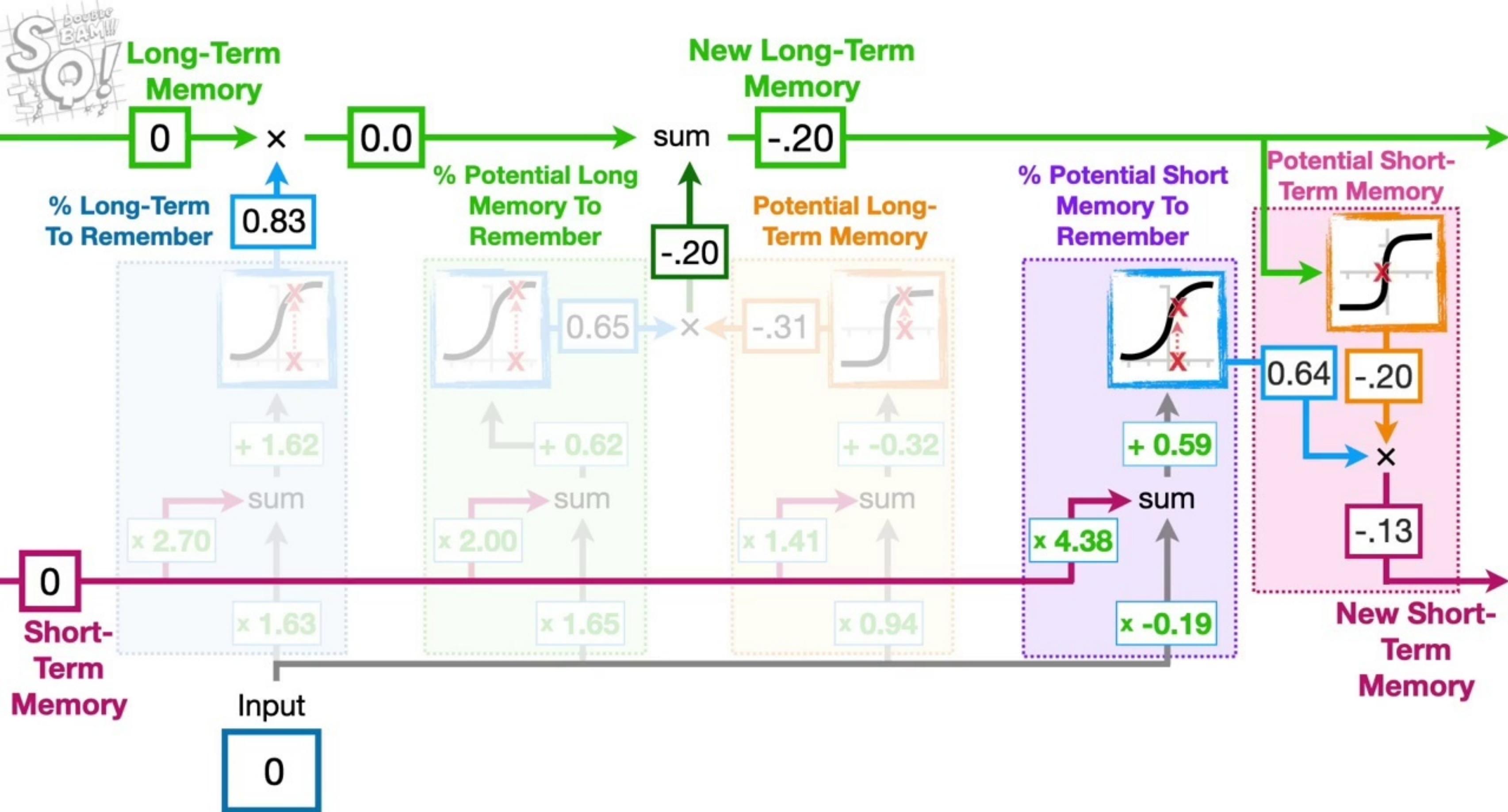
Initial
Long Term
Memory

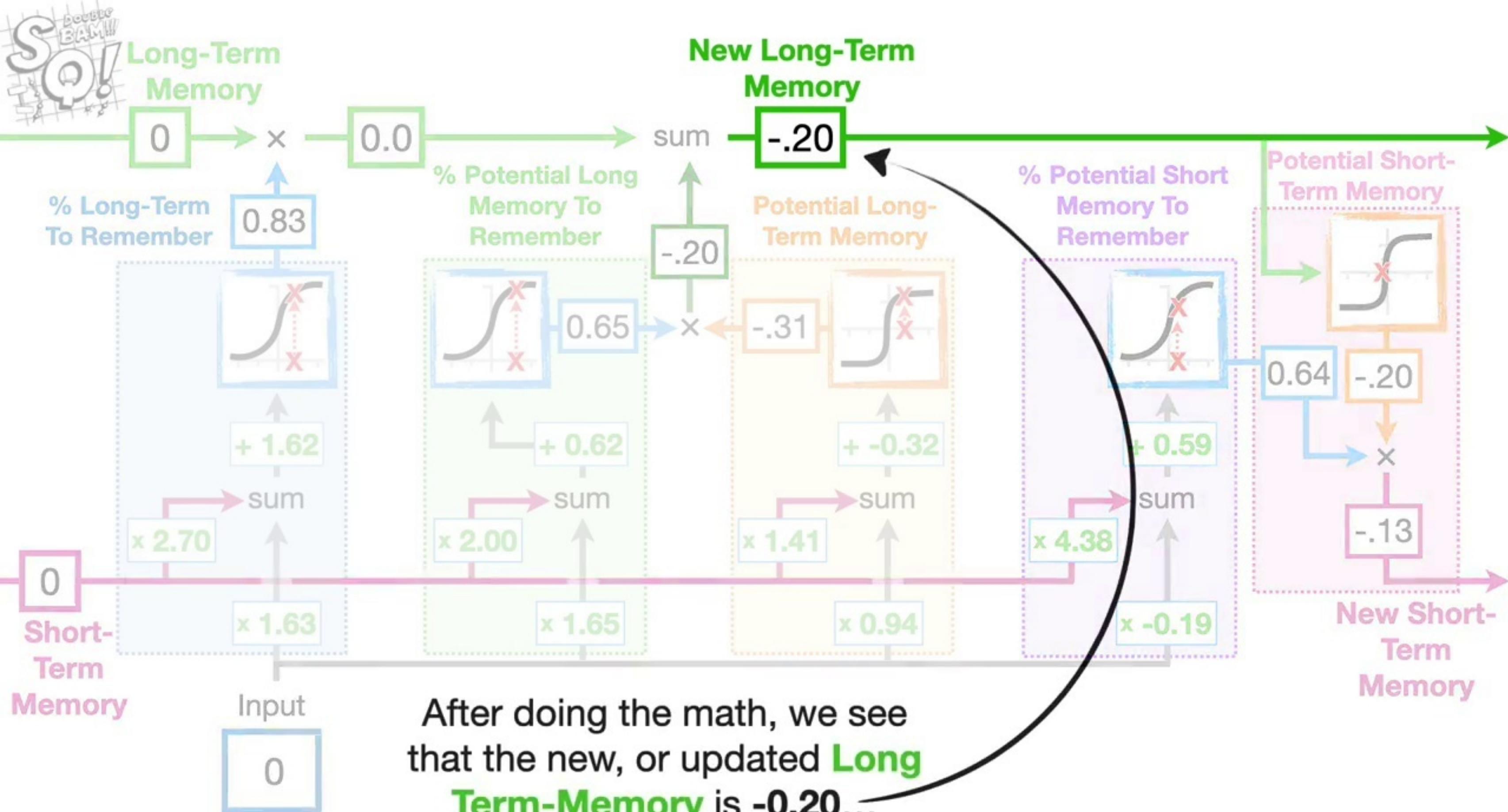


...then we'll start by plugging the value for **Day 1**, which is **0**, into the **Input**.

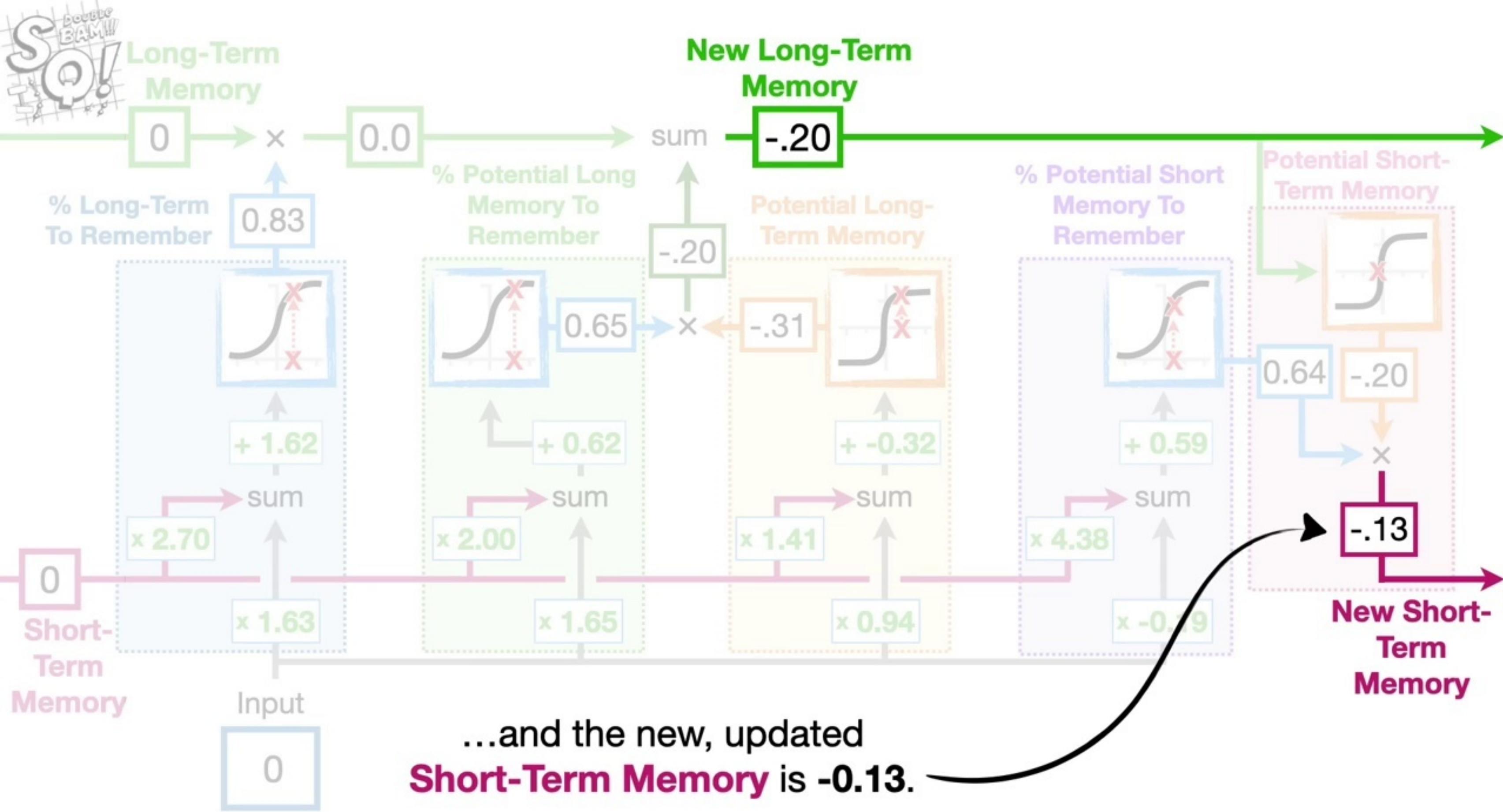








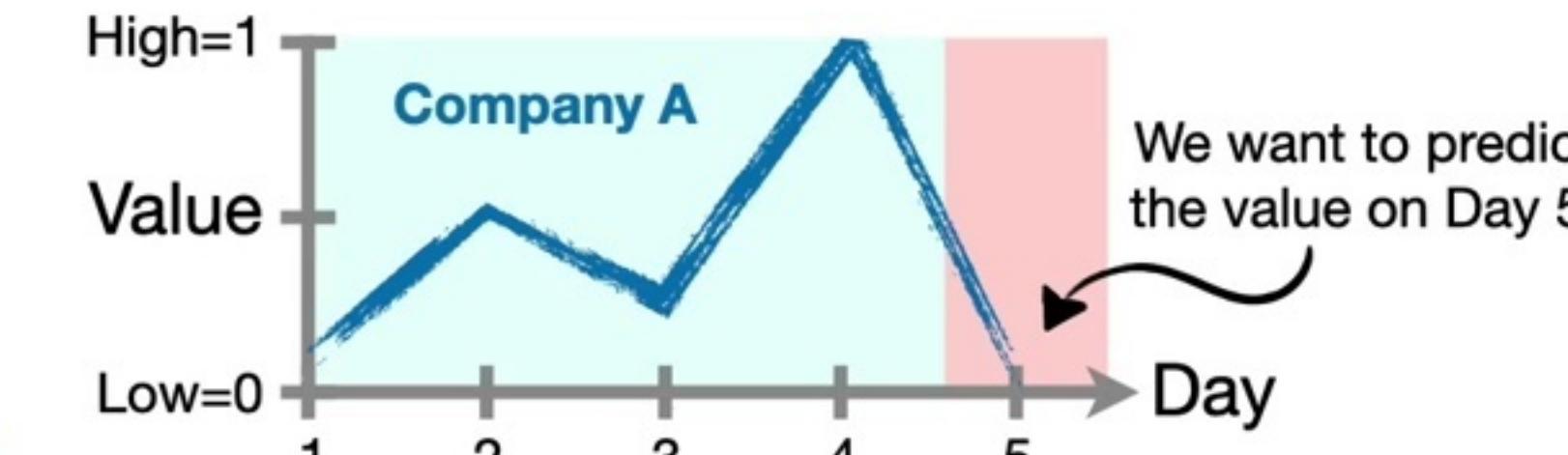
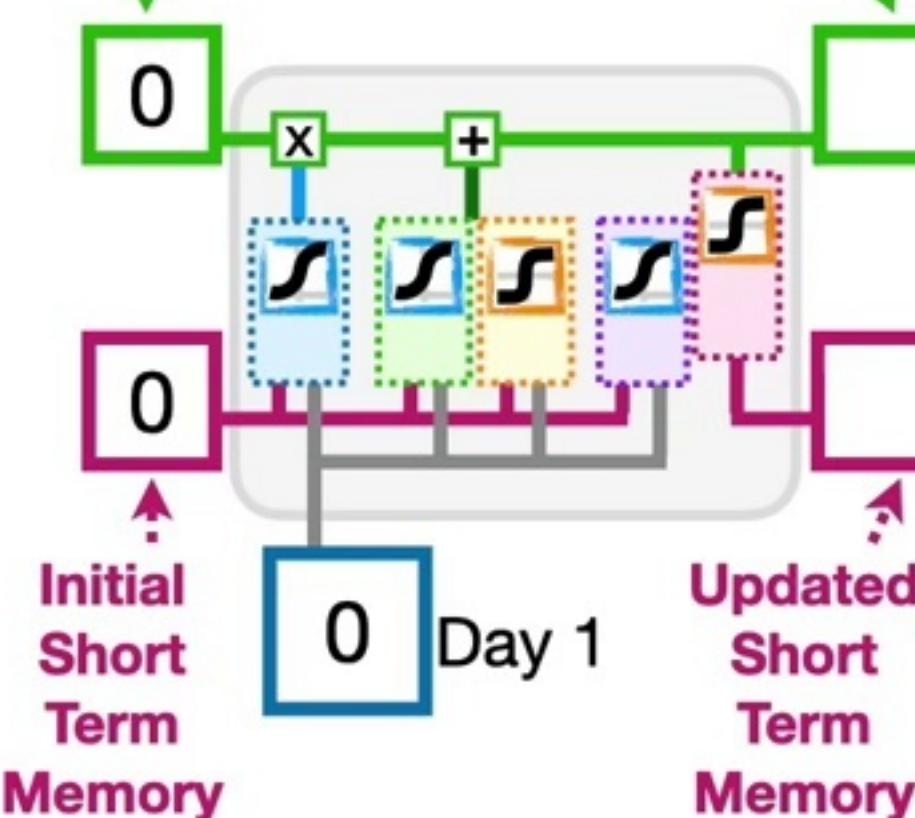
After doing the math, we see that the new, or updated **Long Term-Memory** is -0.20...



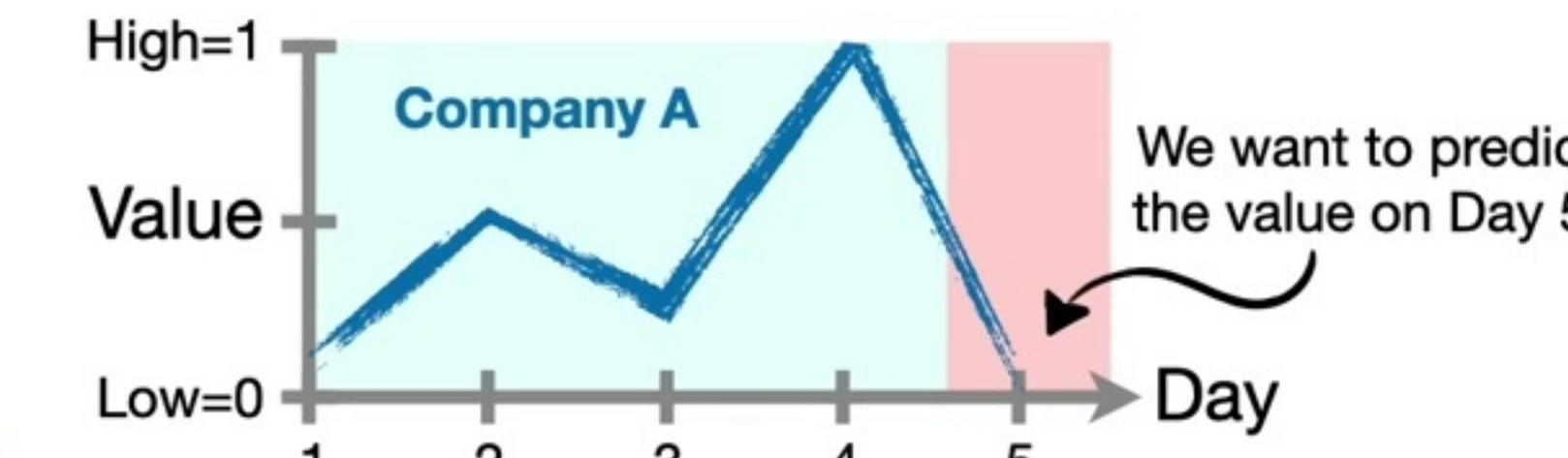
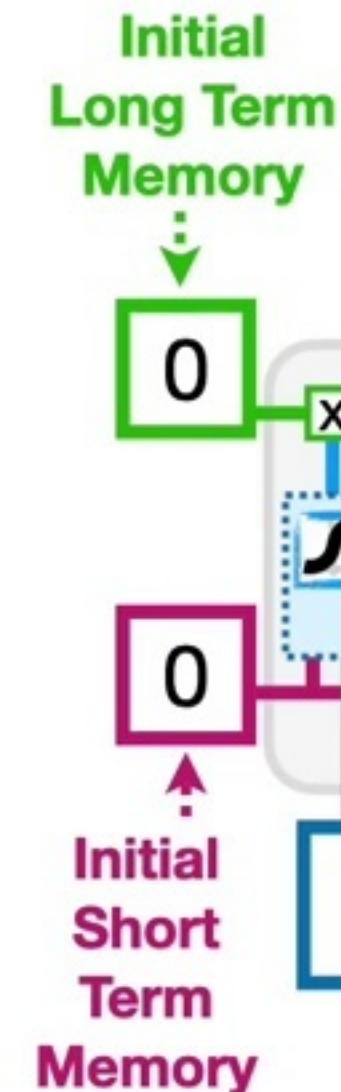


Initial
Long Term
Memory

Updated Long
Term Memory



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

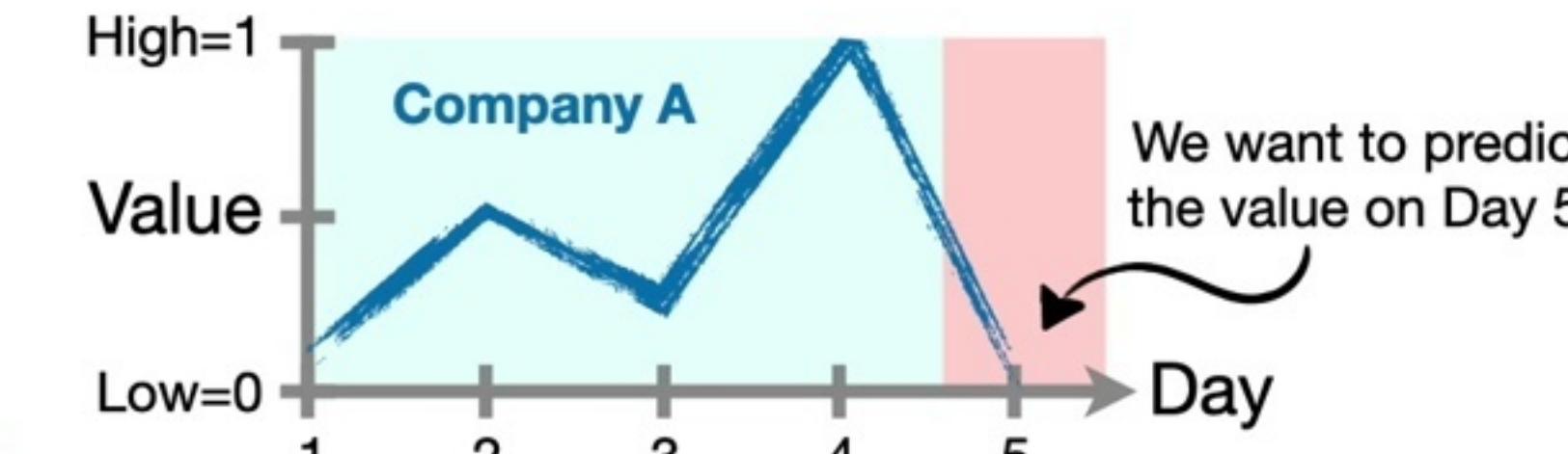
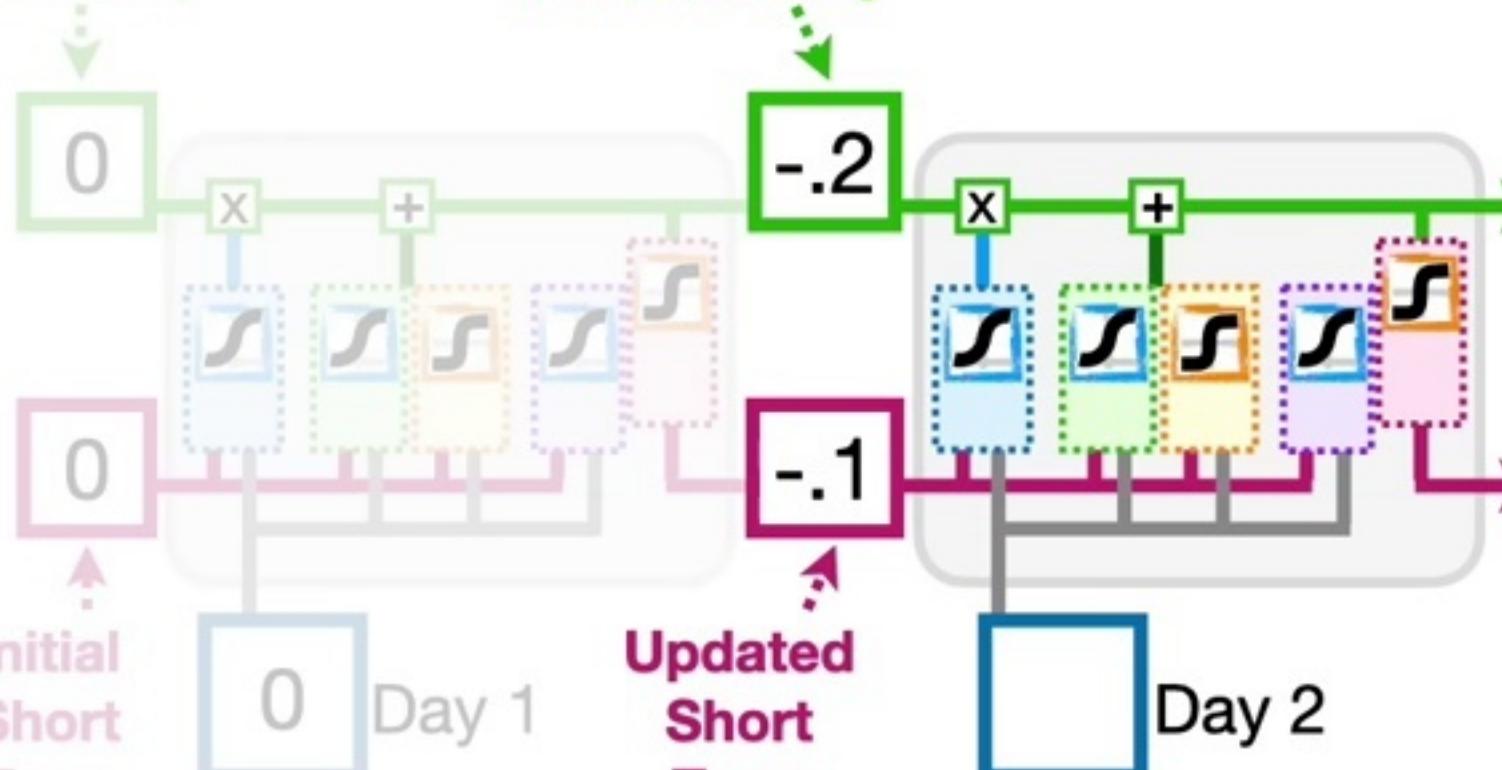


...and -0.1 (rounded) for the updated **Short-Term Memory**.



Initial

Long Term
Memory

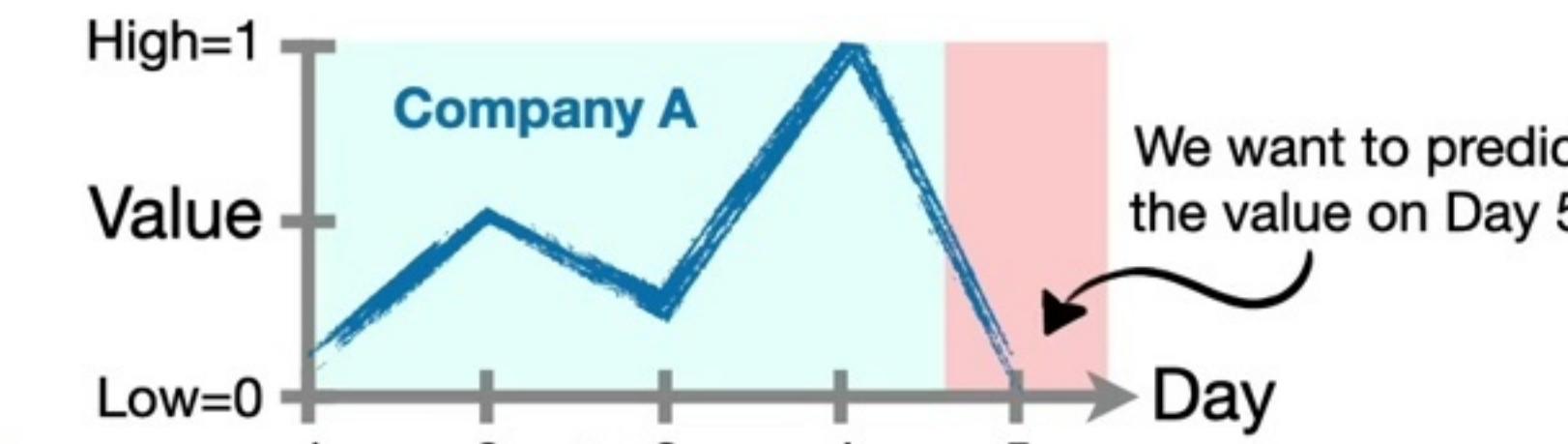
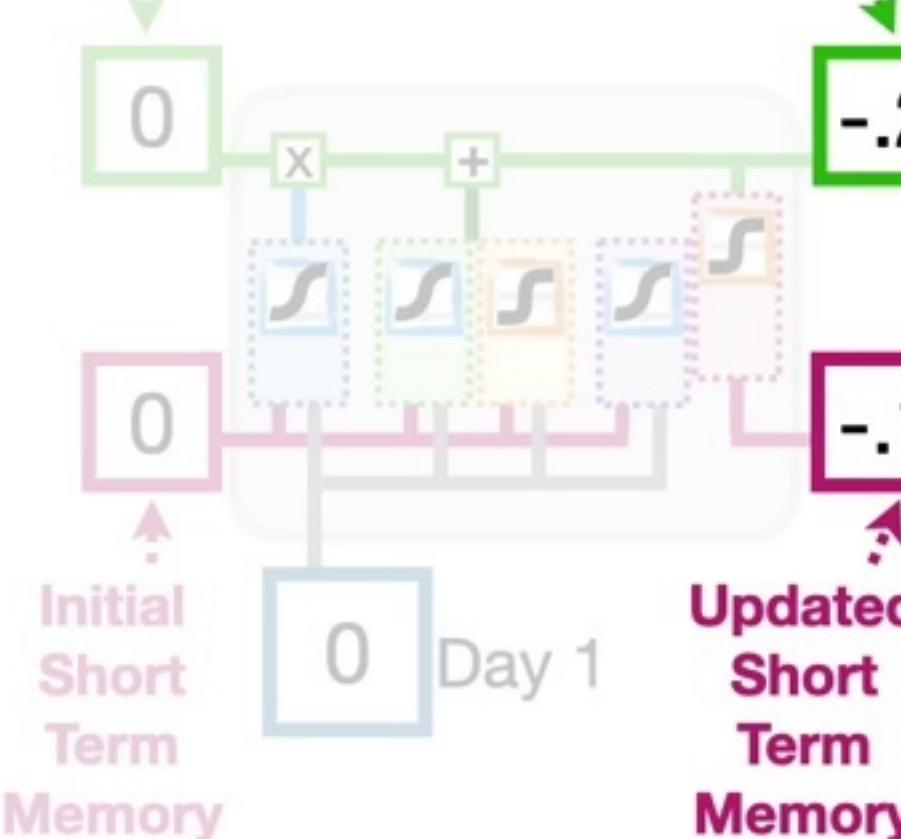


Now we unroll the **LSTM**, using
the updated memories...



Initial
Long Term
Memory

Updated Long
Term Memory



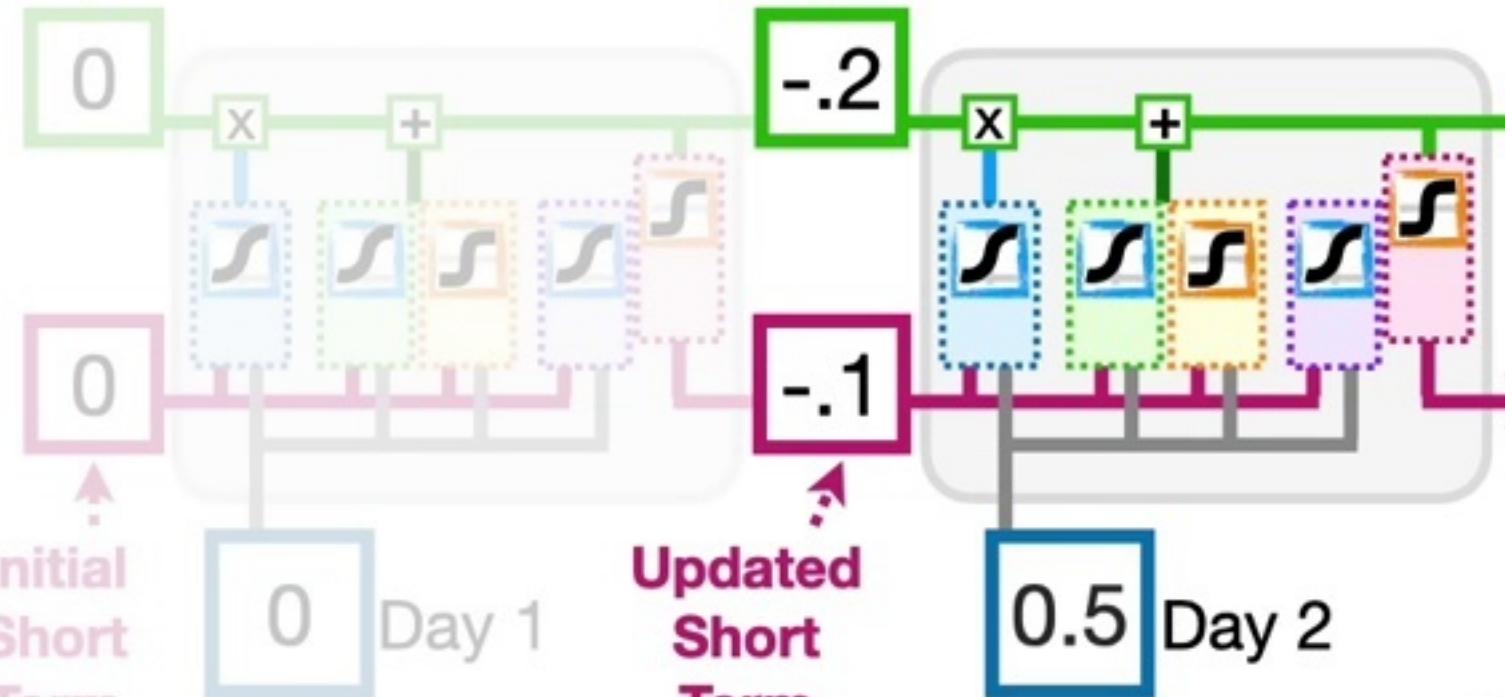
...and plug the value from
Day 2, 0.5, into the Input.



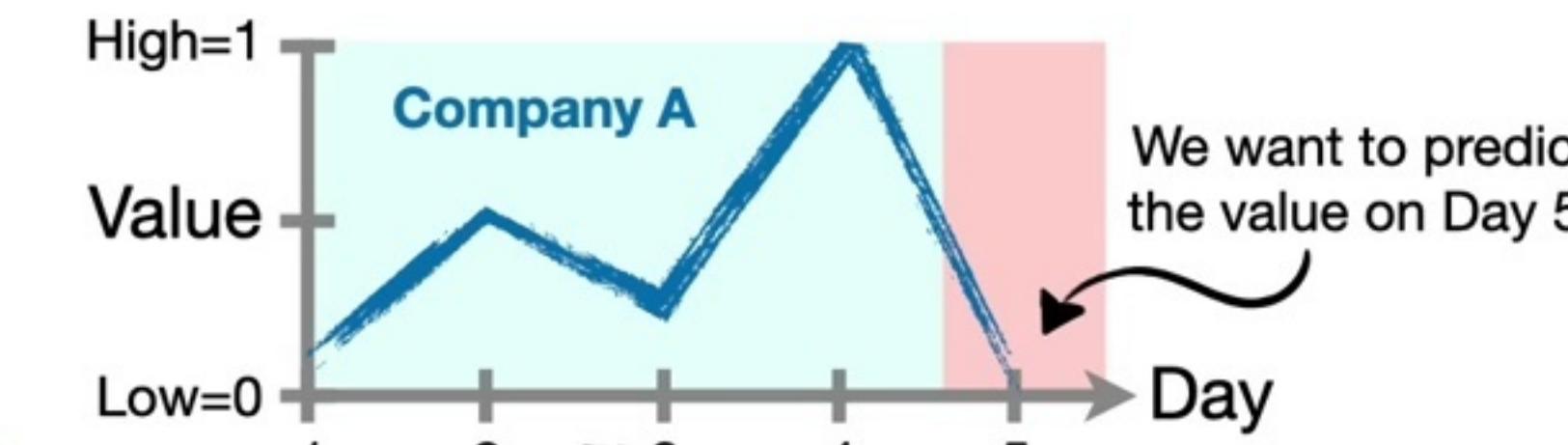
Initial
Long Term
Memory

Updated Long
Term Memory

Initial
Short
Term
Memory

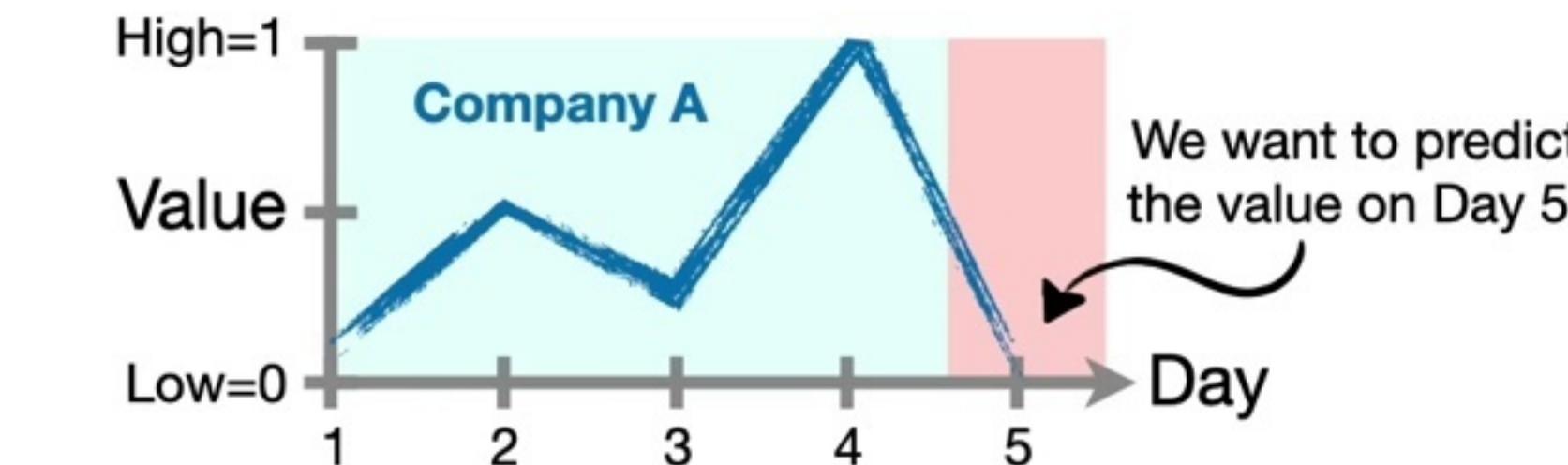
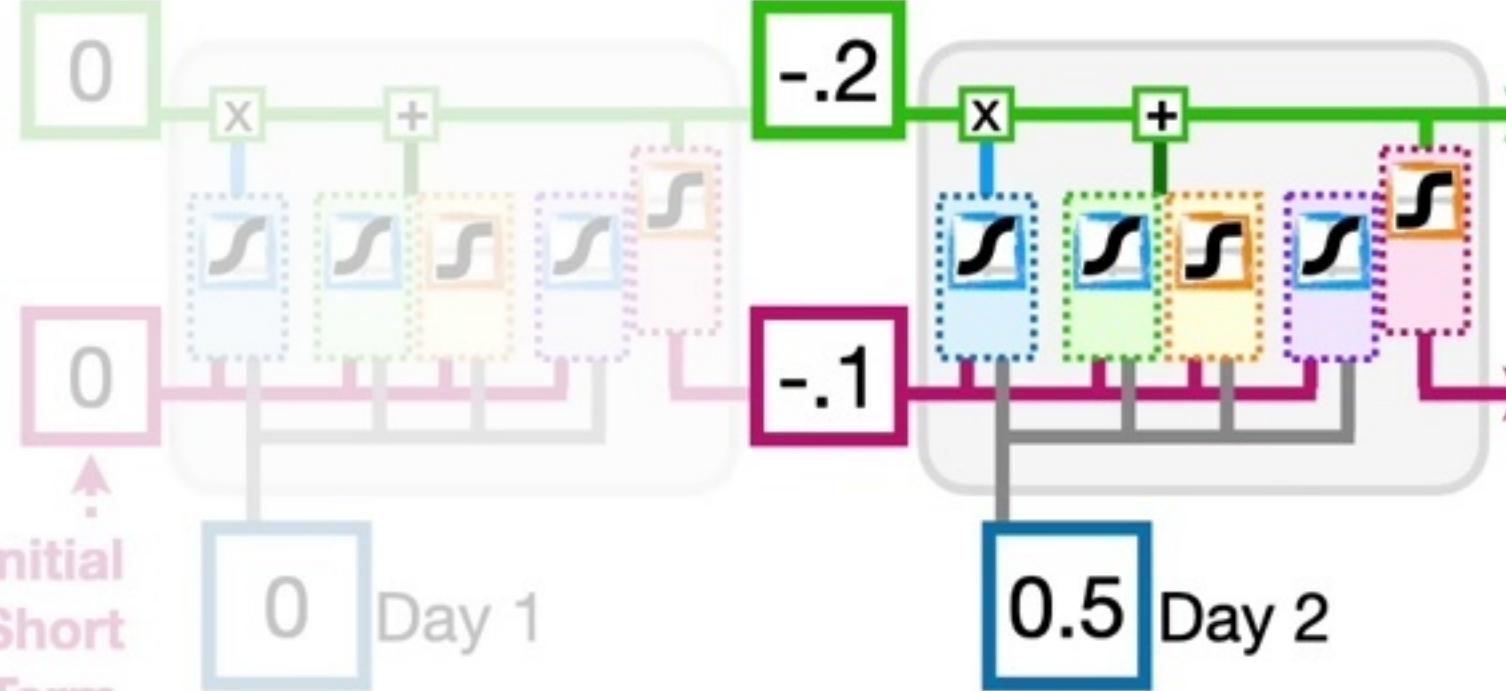


...and plug the value from
Day 2, 0.5, into the Input.





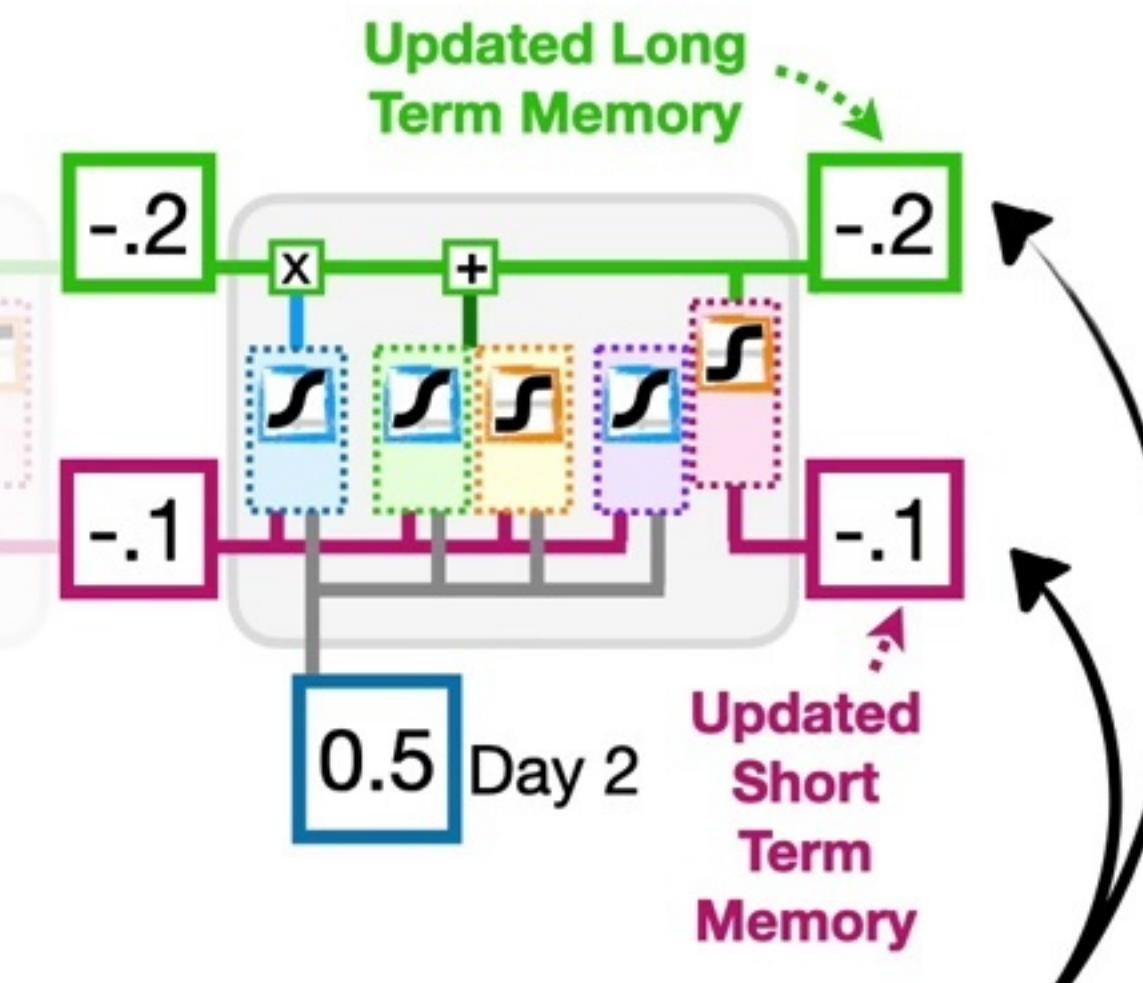
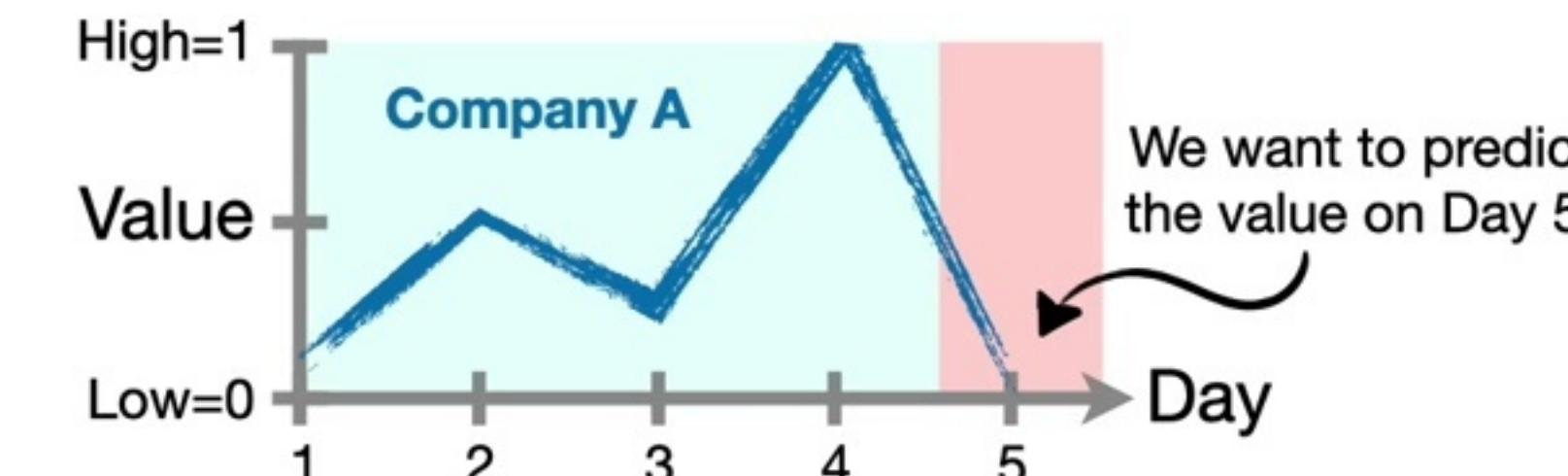
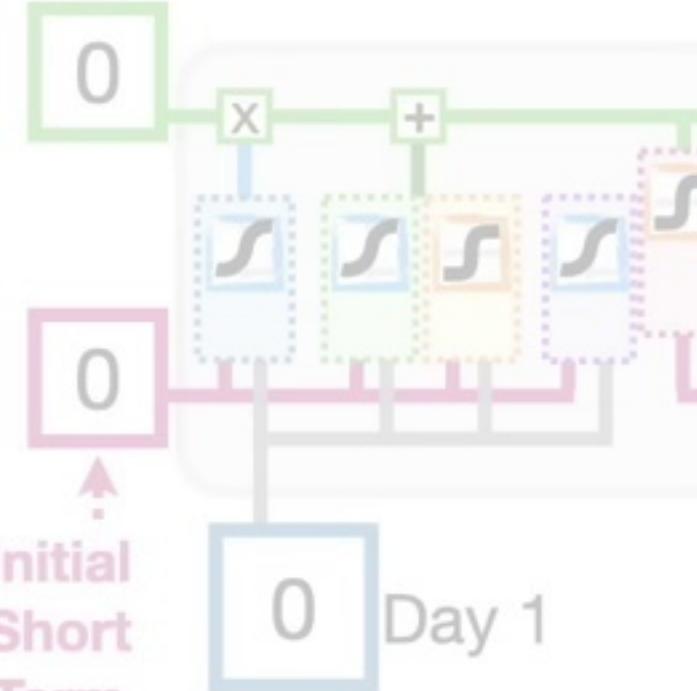
Initial
Long Term
Memory



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



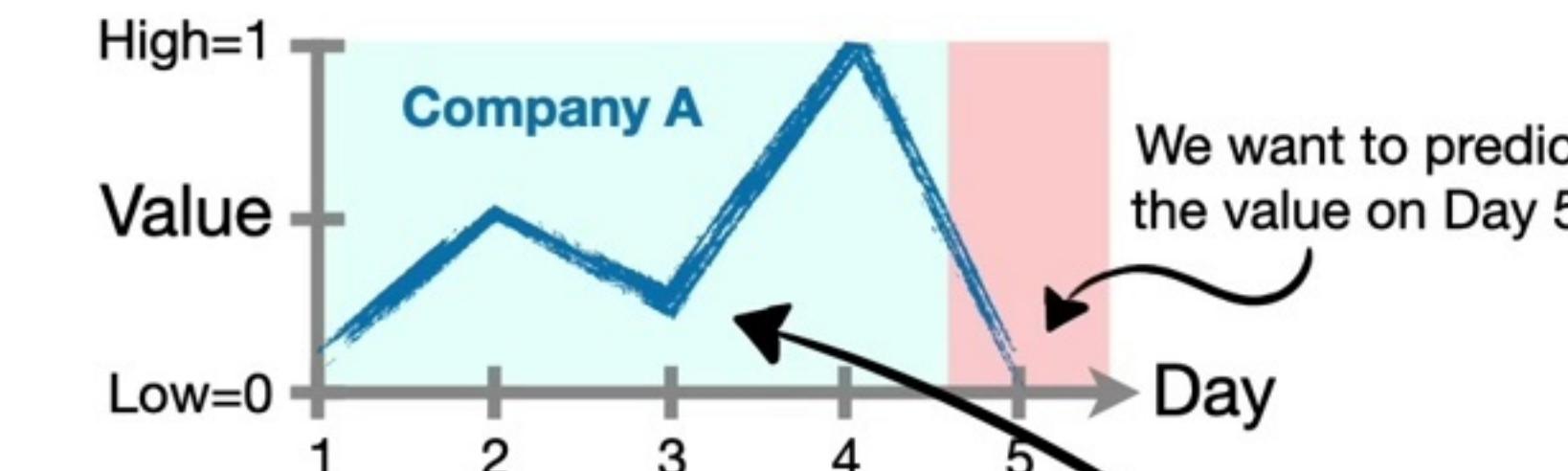
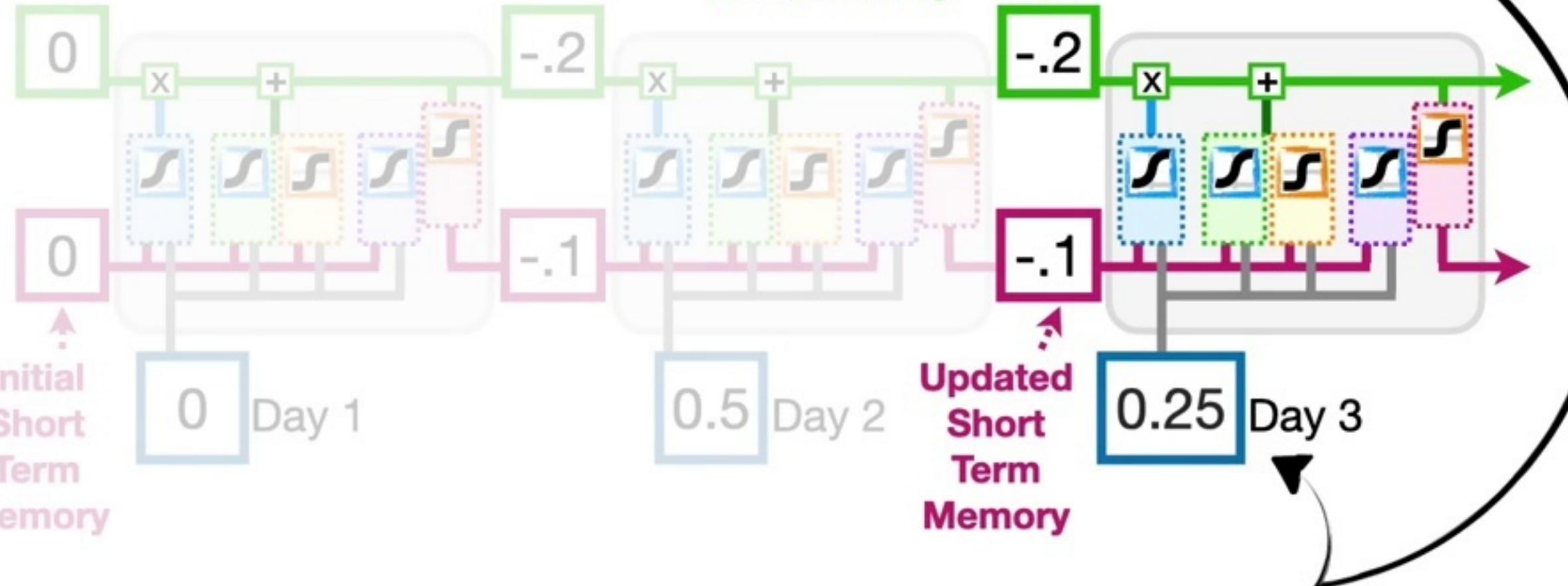
Initial
Long Term
Memory



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



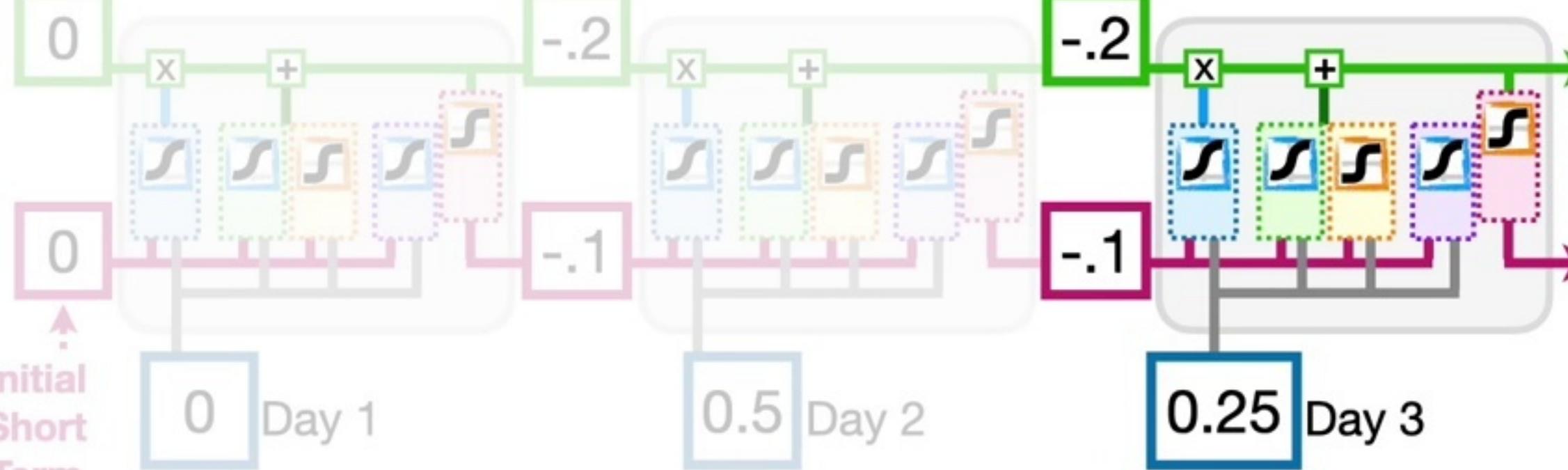
Initial
Long Term
Memory



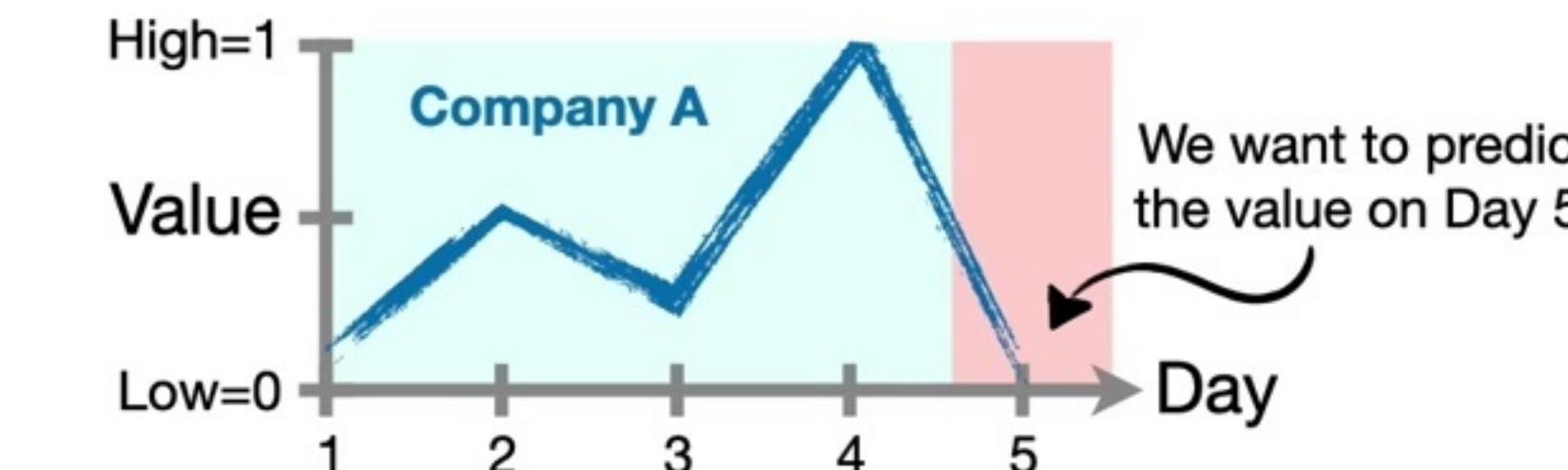
...and plug in the value for **Day 3**.

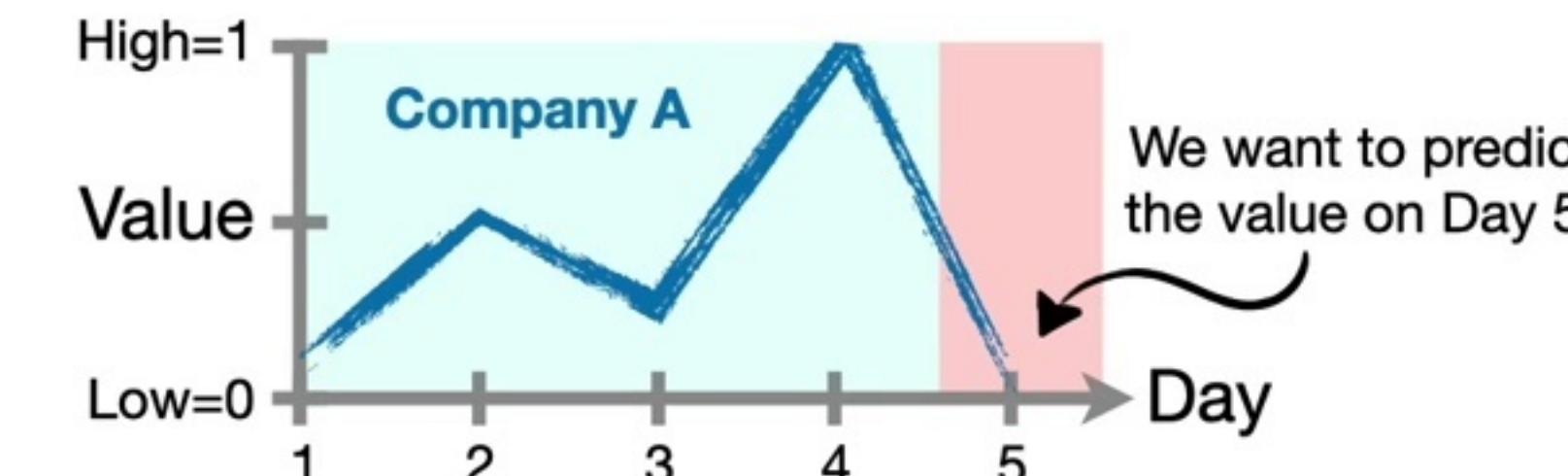


Initial
Long Term
Memory

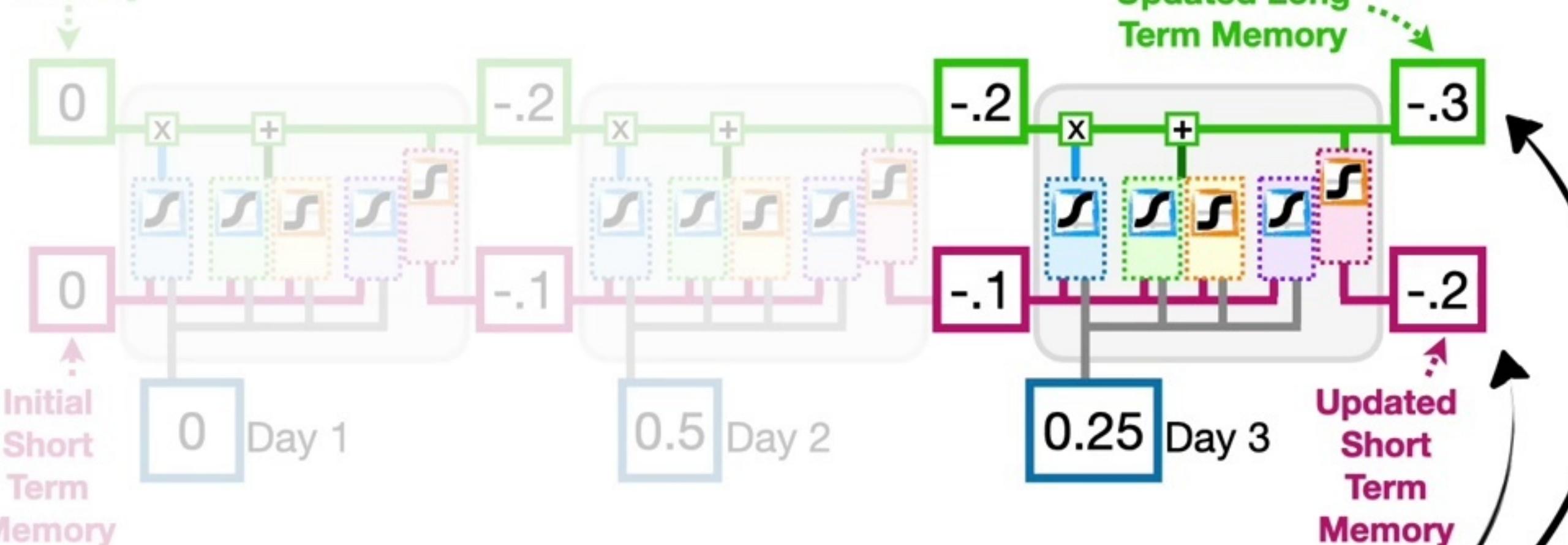


And then the **LSTM** does the
math, again, using the exact
same **Weights** and **Biases**...

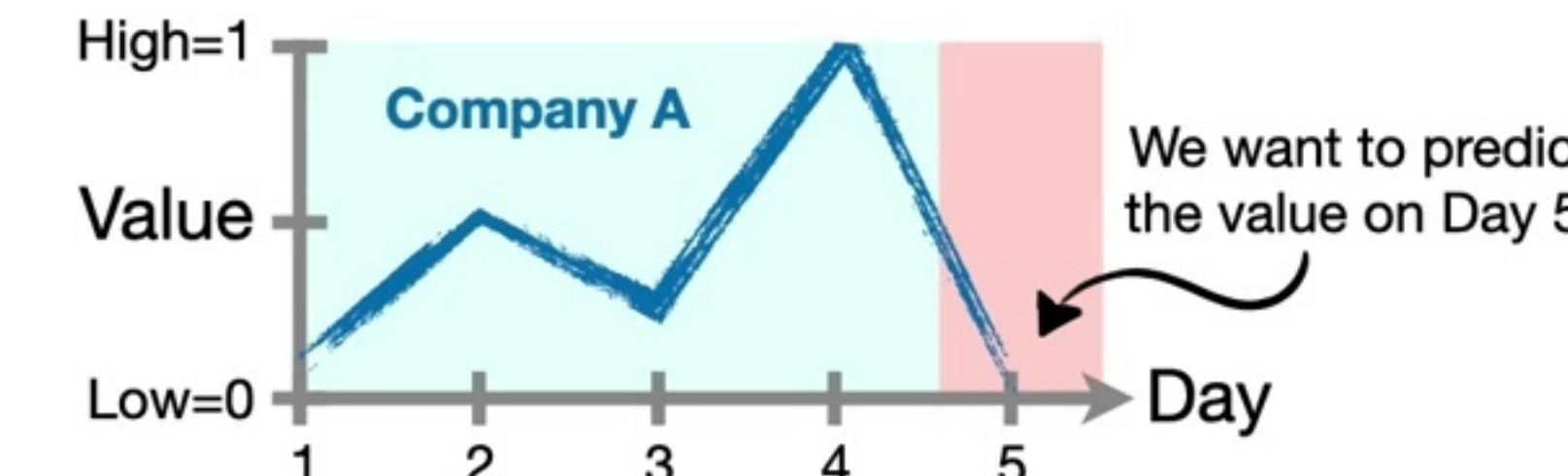




Initial
Long Term
Memory



...and gives us these
updated memories.



Initial
Long Term
Memory



0

-0.2

-0.2

-0.3



0

-0.1

-0.1

-0.2

Initial
Short
Term
Memory

0 Day 1

0.5 Day 2

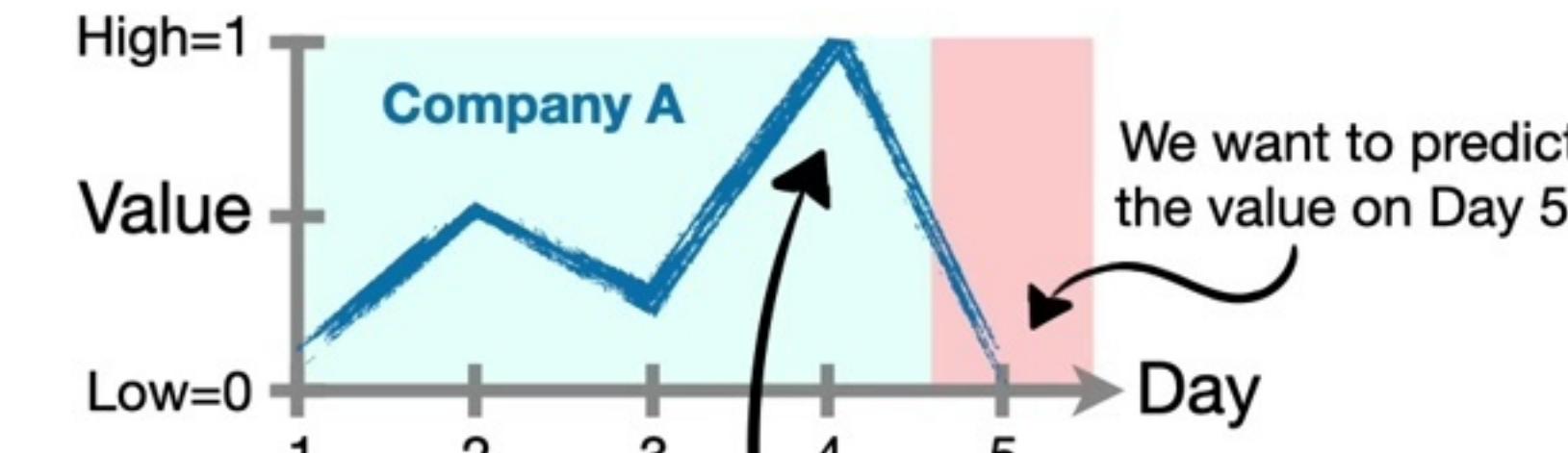
0.25 Day 3

Day 4

Updated Long
Term Memory

Updated
Short
Term
Memory

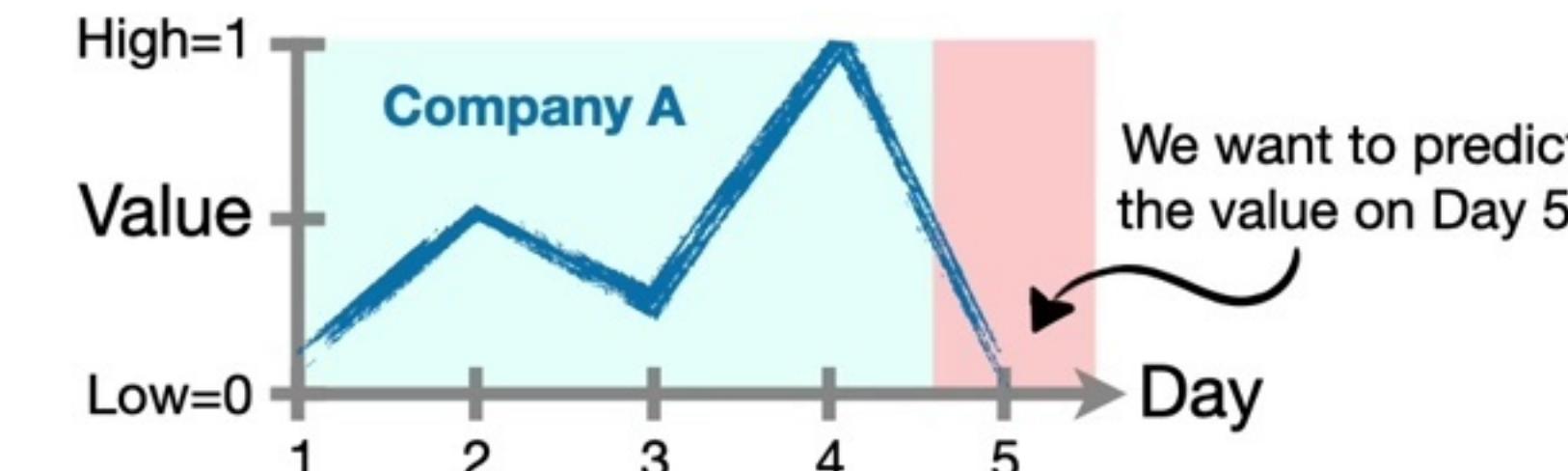
Then we unroll the **LSTM**
one last time...



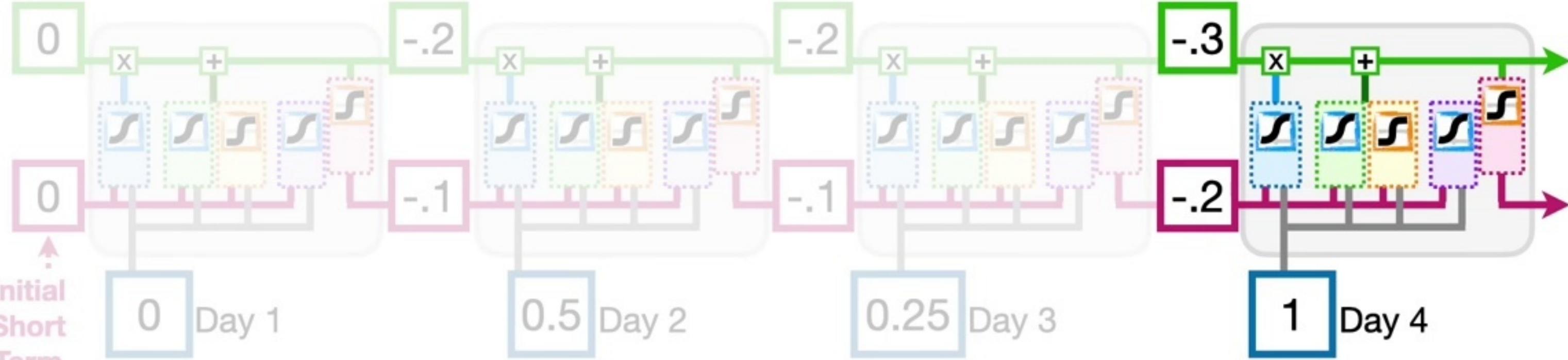
Initial
Long Term
Memory



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

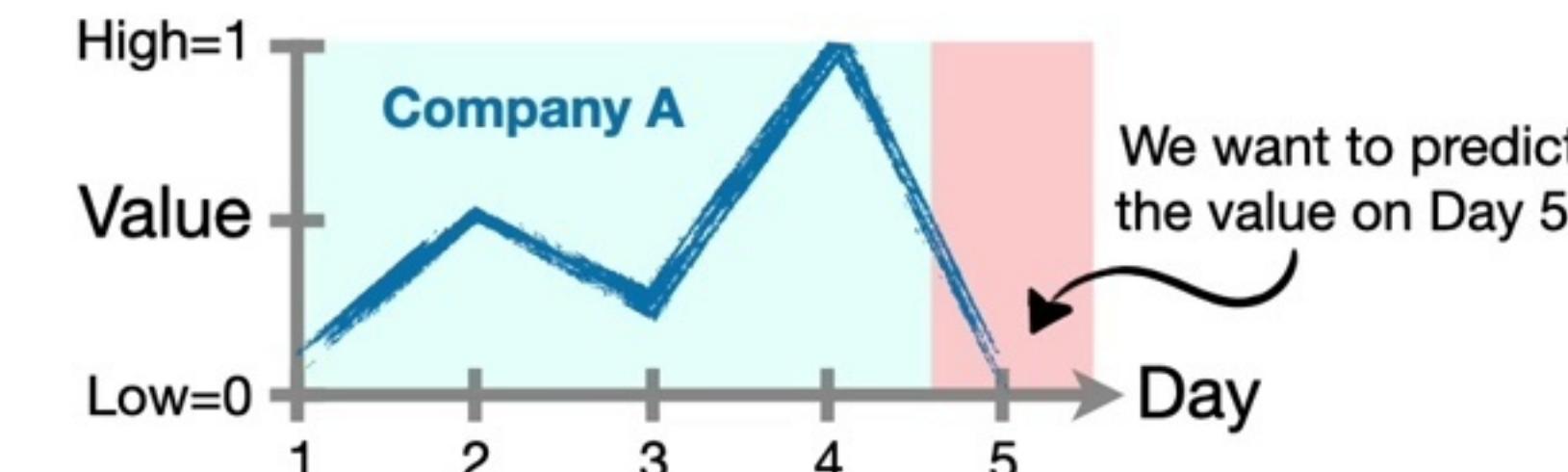


Initial
Long Term
Memory

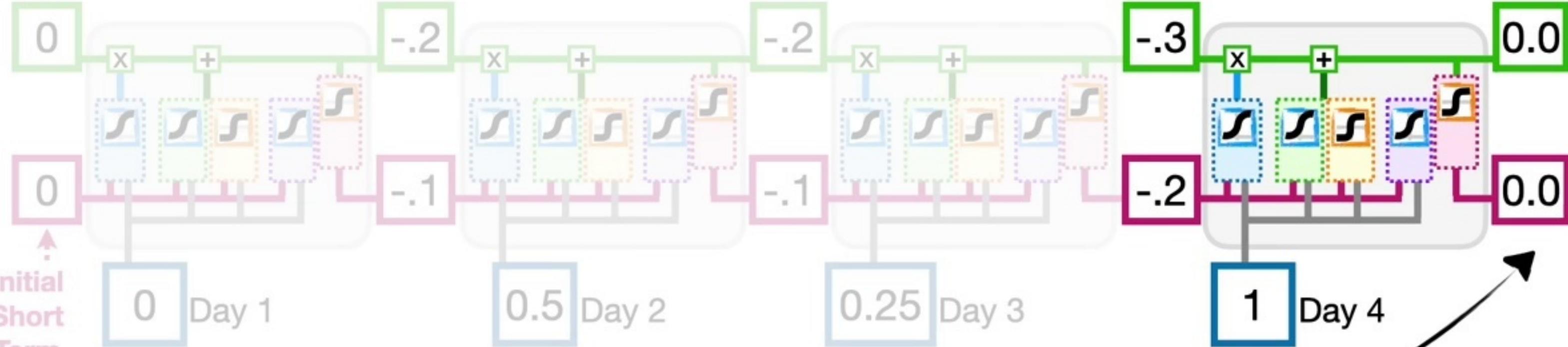


Initial
Short
Term
Memory

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



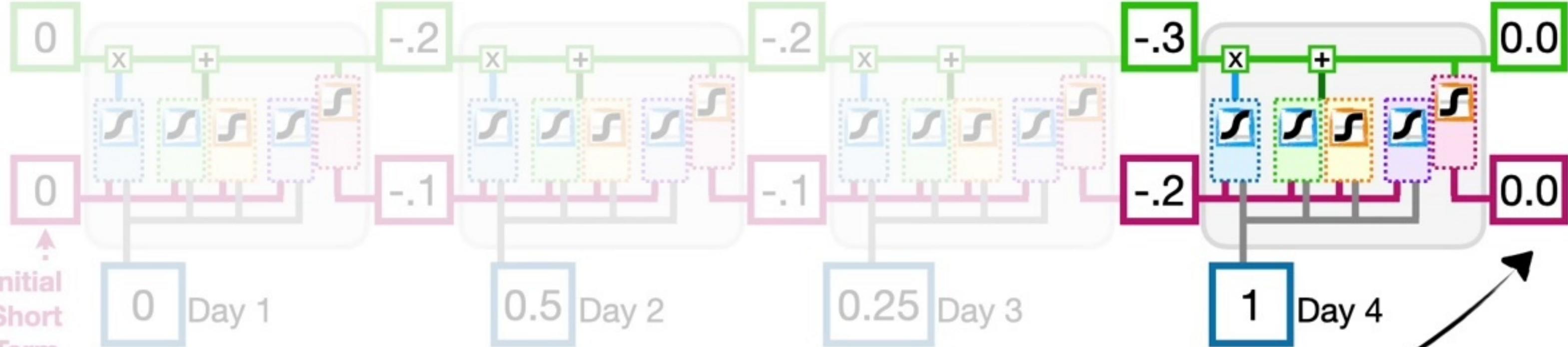
Initial
Long Term
Memory



...and gives us the
final memories.



Initial
Long Term
Memory

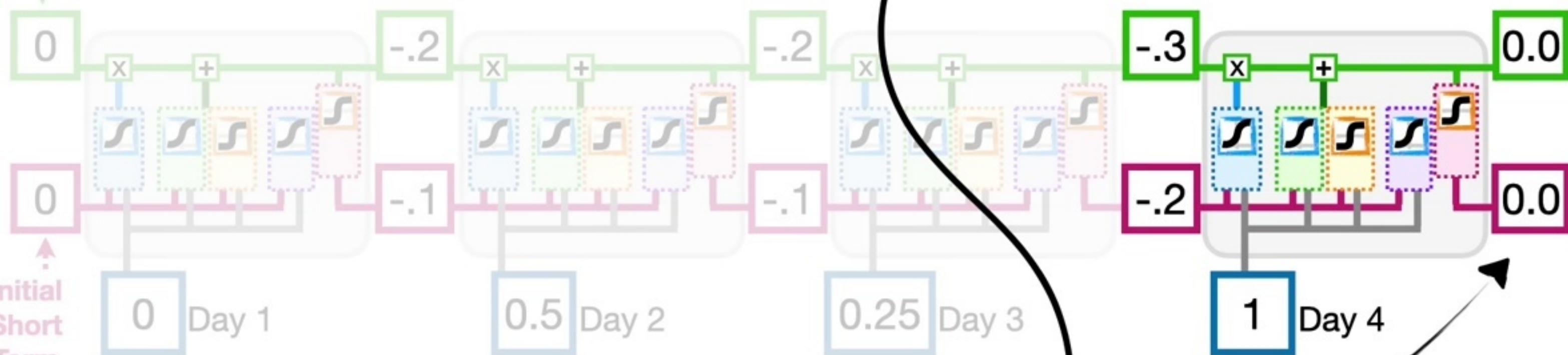


And the final **Short-Term
Memory**, 0.0, is the **Output**
from the unrolled **LSTM**.

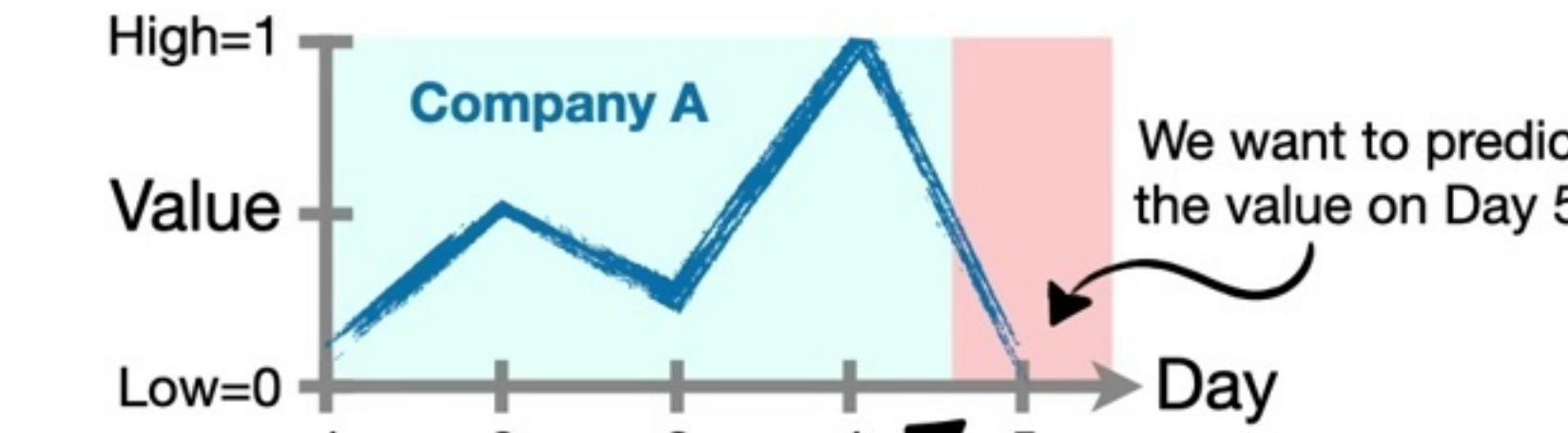
We want to predict
the value on Day 5.

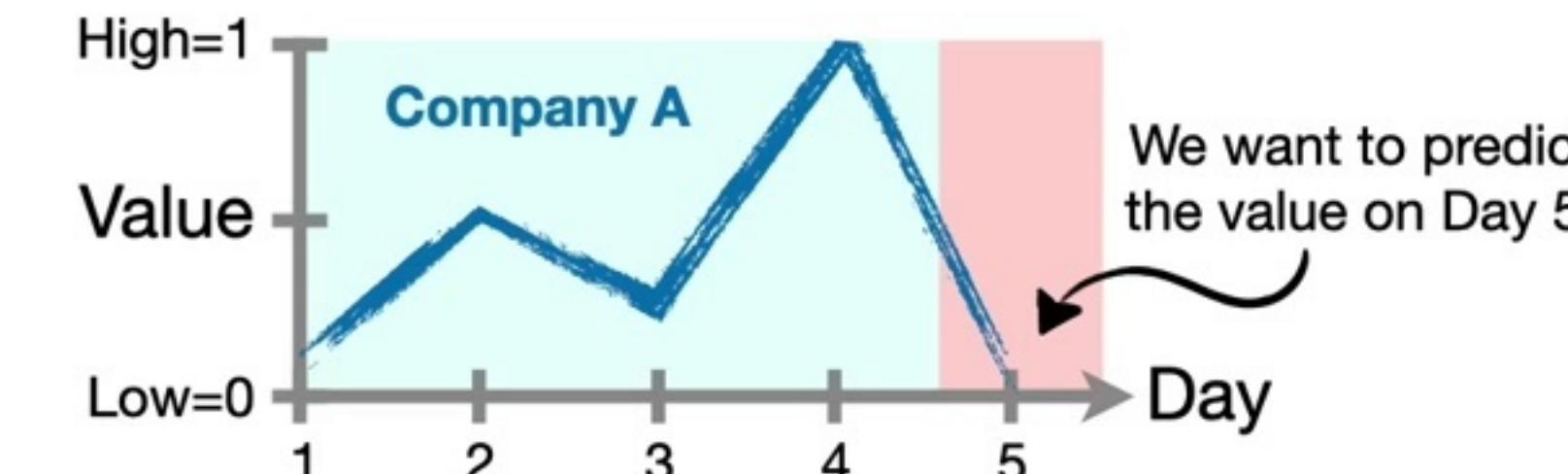


Initial
Long Term
Memory

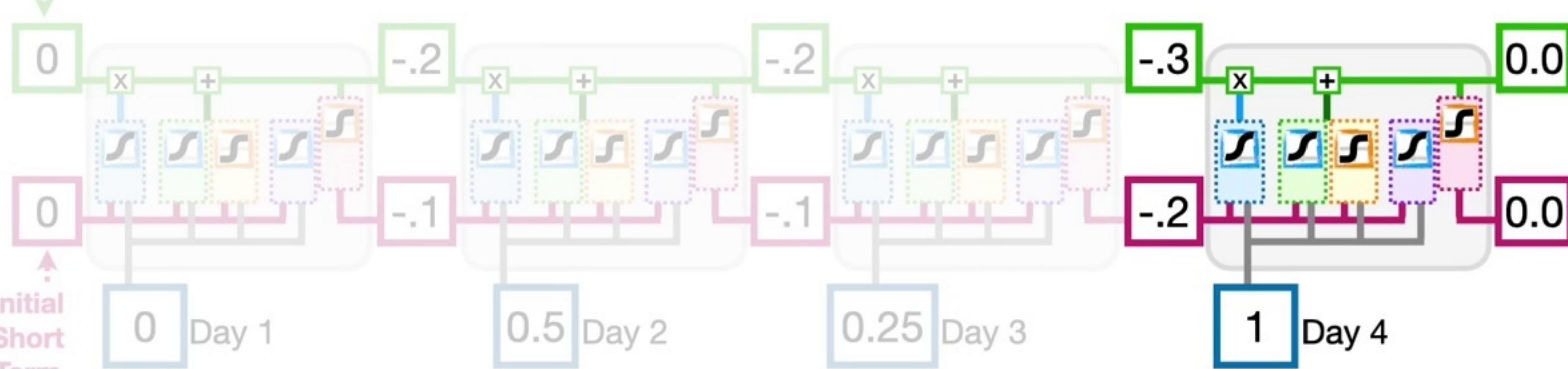


And that means that the **Output**
from the **LSTM** correctly predicts
Company A's value for Day 5.

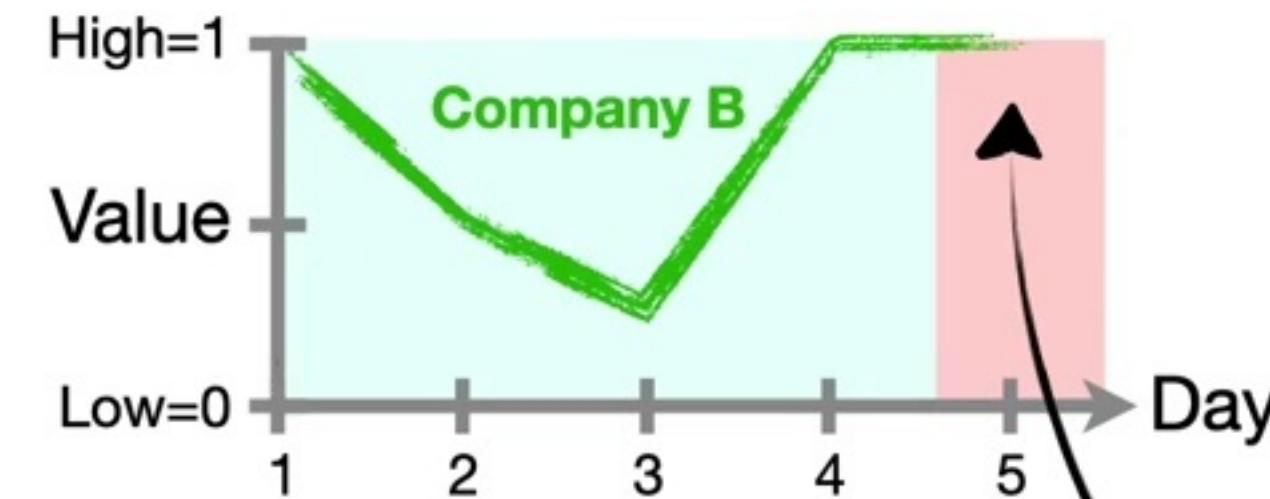




Initial
Long Term
Memory



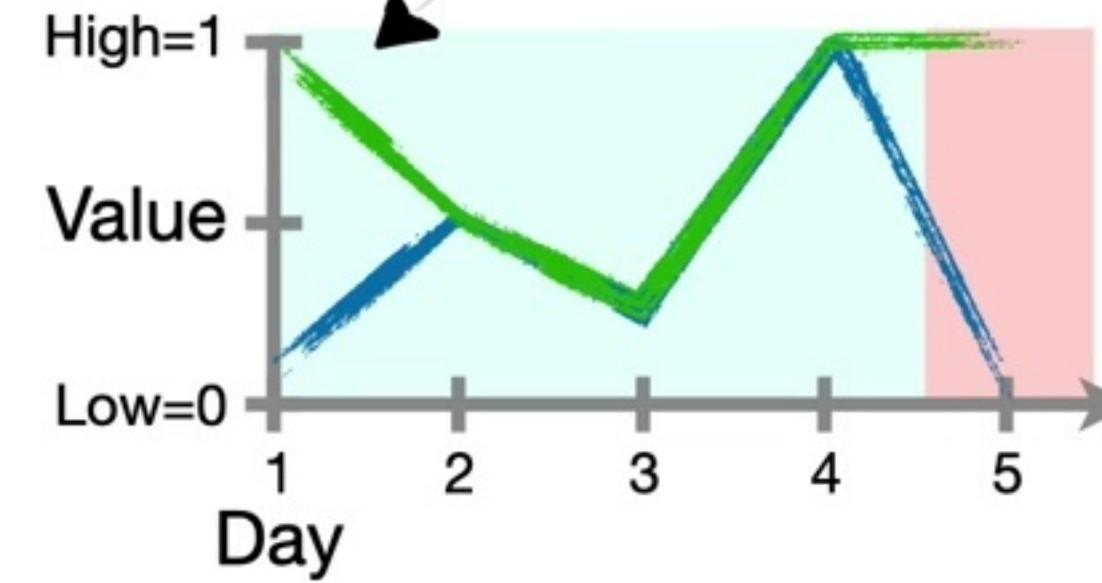
Now that we have shown that the
LSTM can correctly predict the value
on **Day 5** for **Company A**...



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



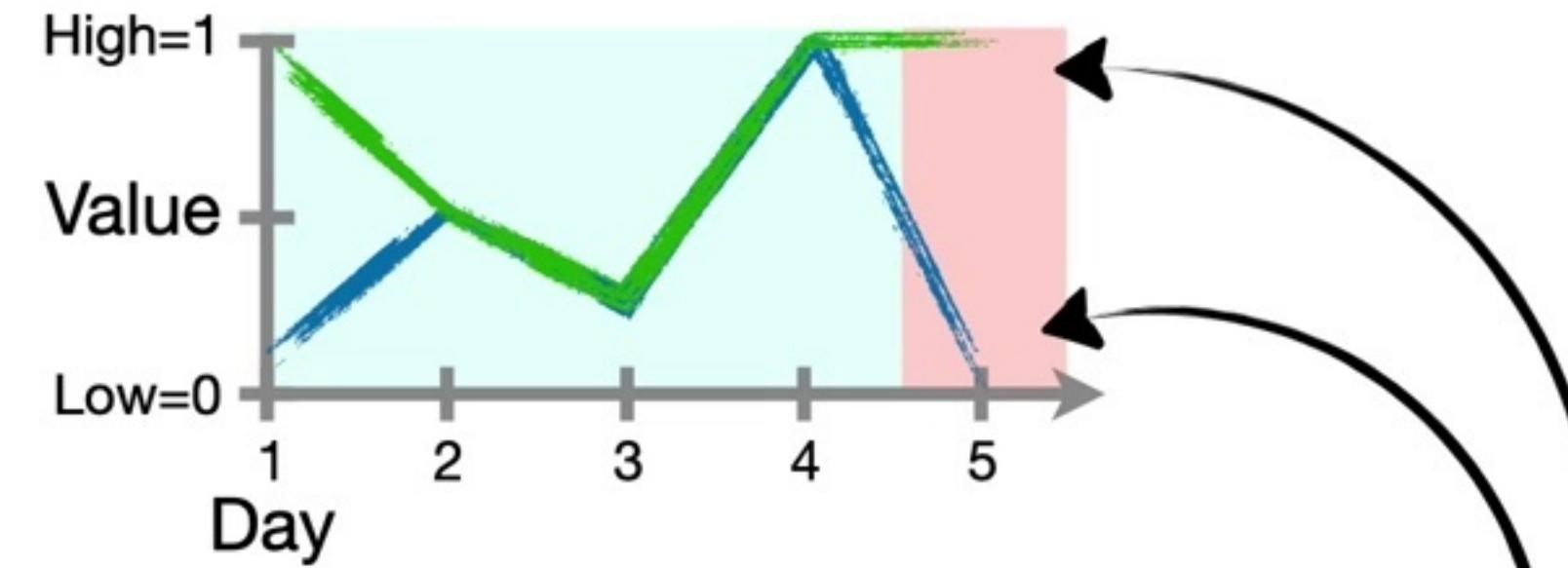
Company A =
Company B =



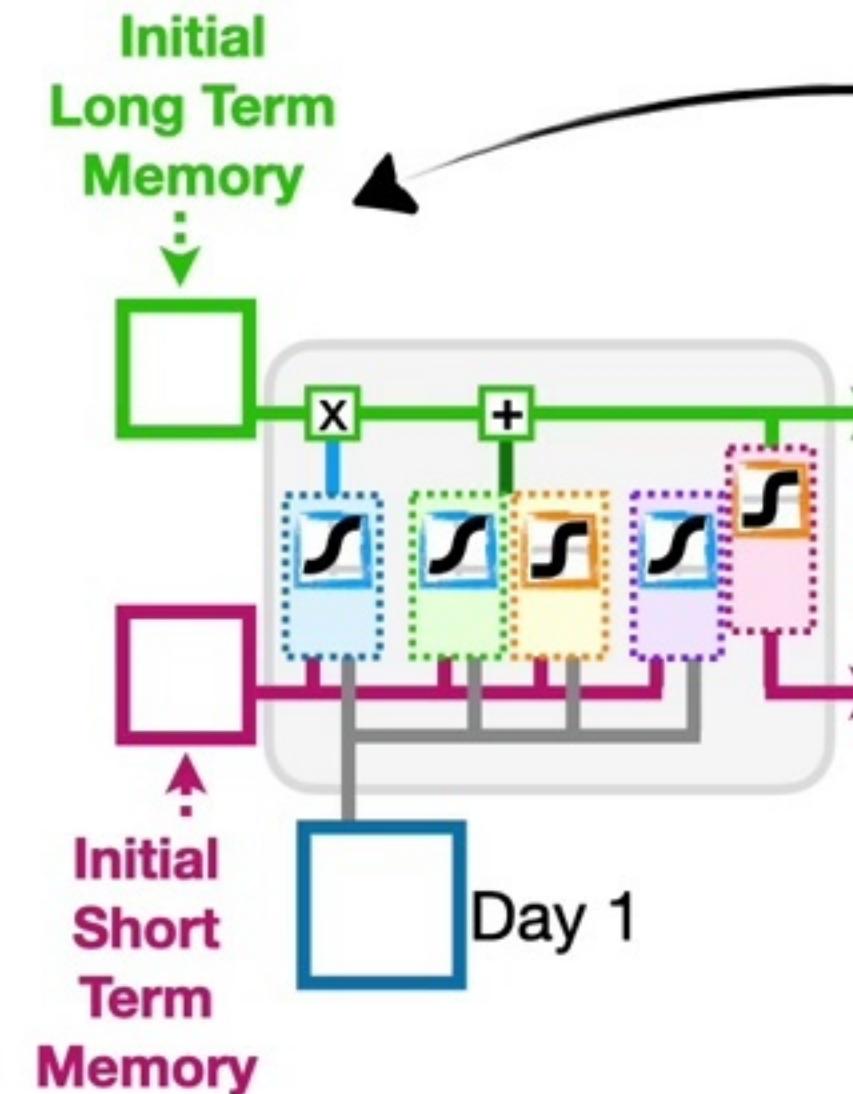
NOTE: Remember, on **Days 1 through 4**, the only difference between the companies occurs on **Day 1...**



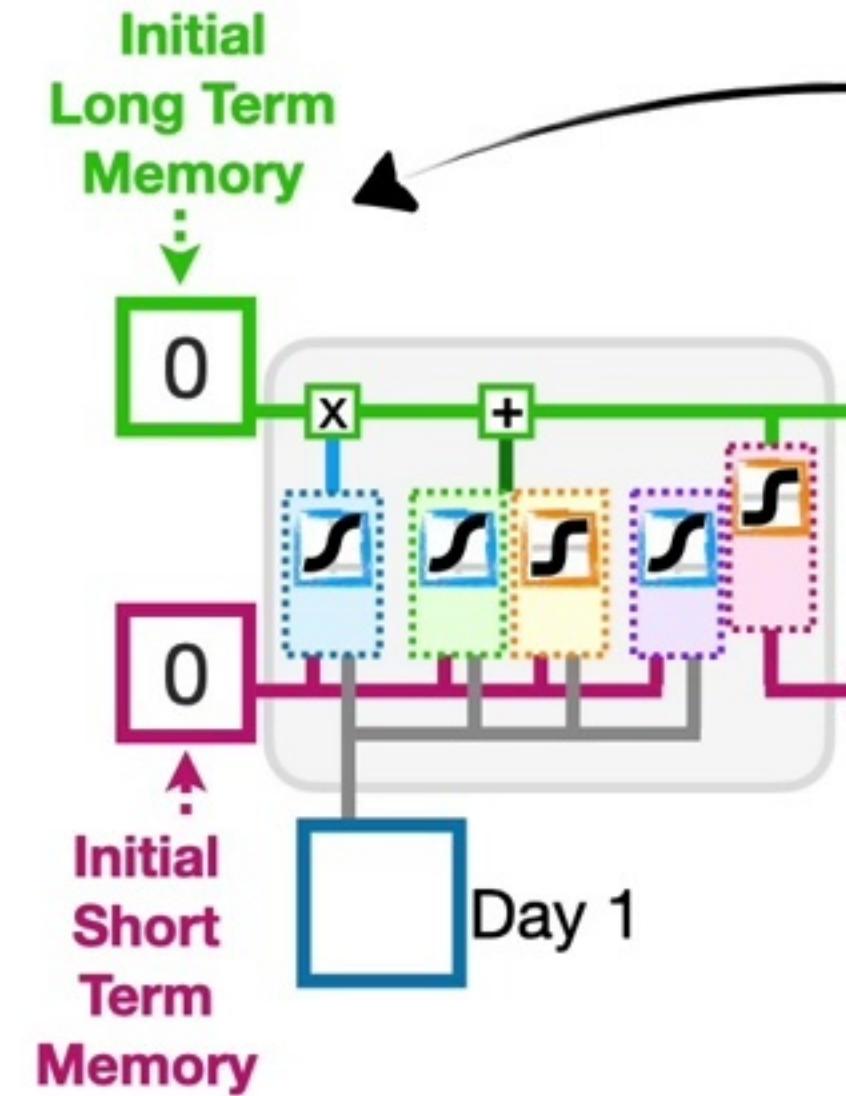
Company A =
Company B =



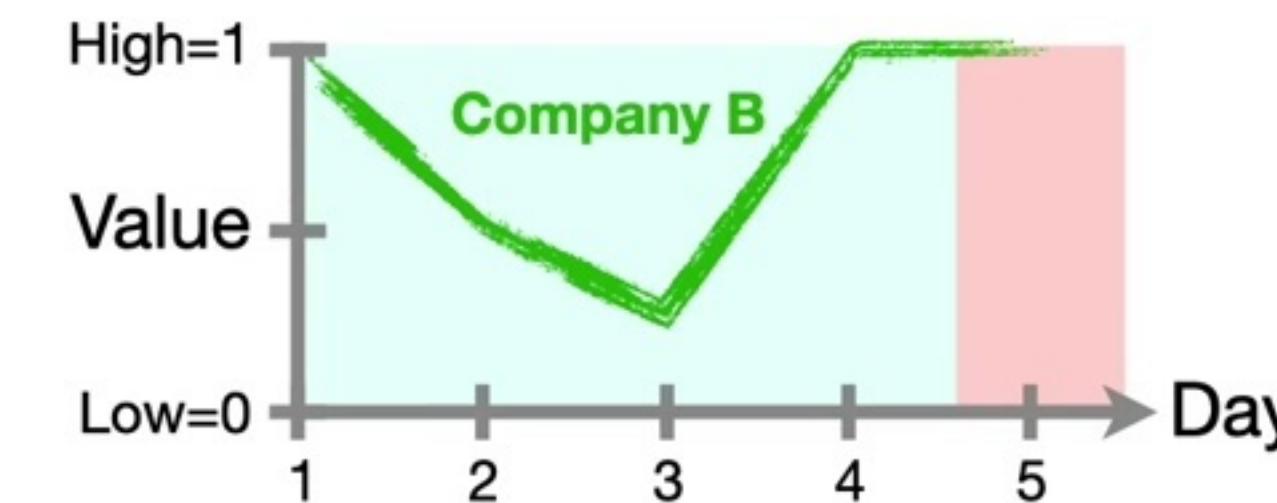
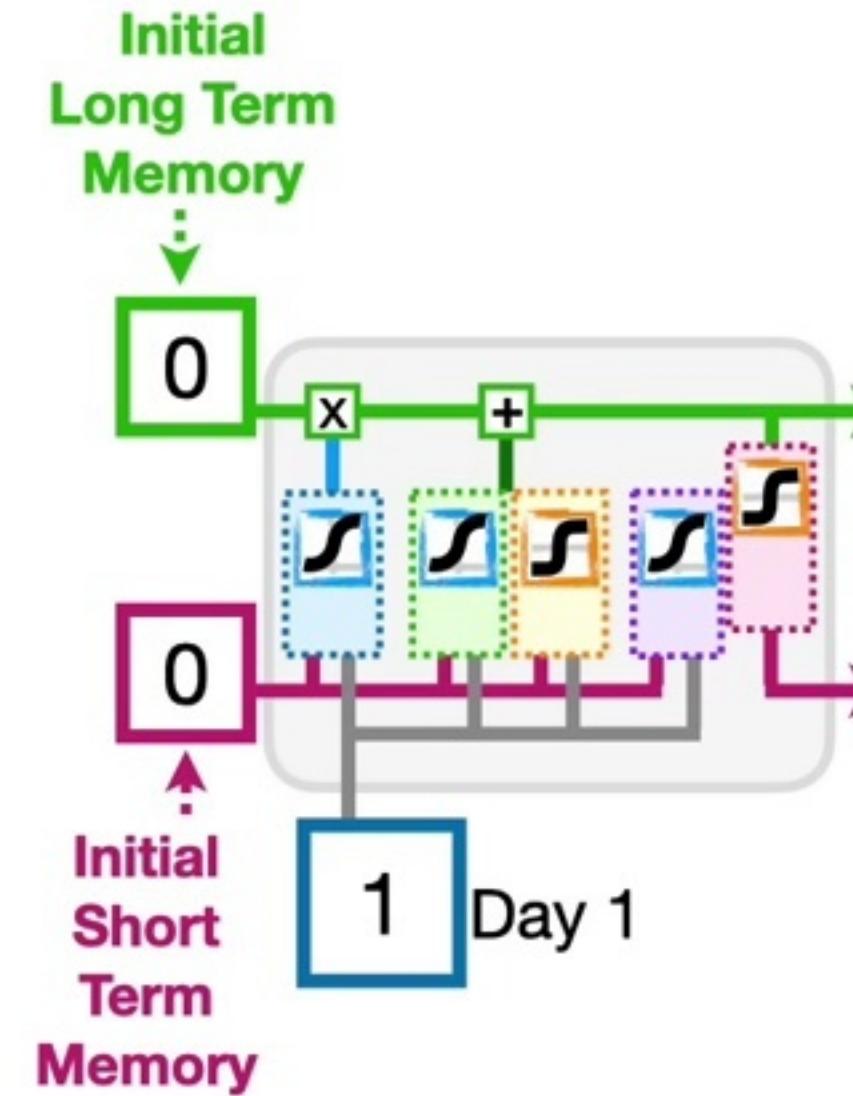
...and that means the **LSTM** has to remember
what happened on **Day 1** in order to correctly
predict the different output values on **Day 5**.



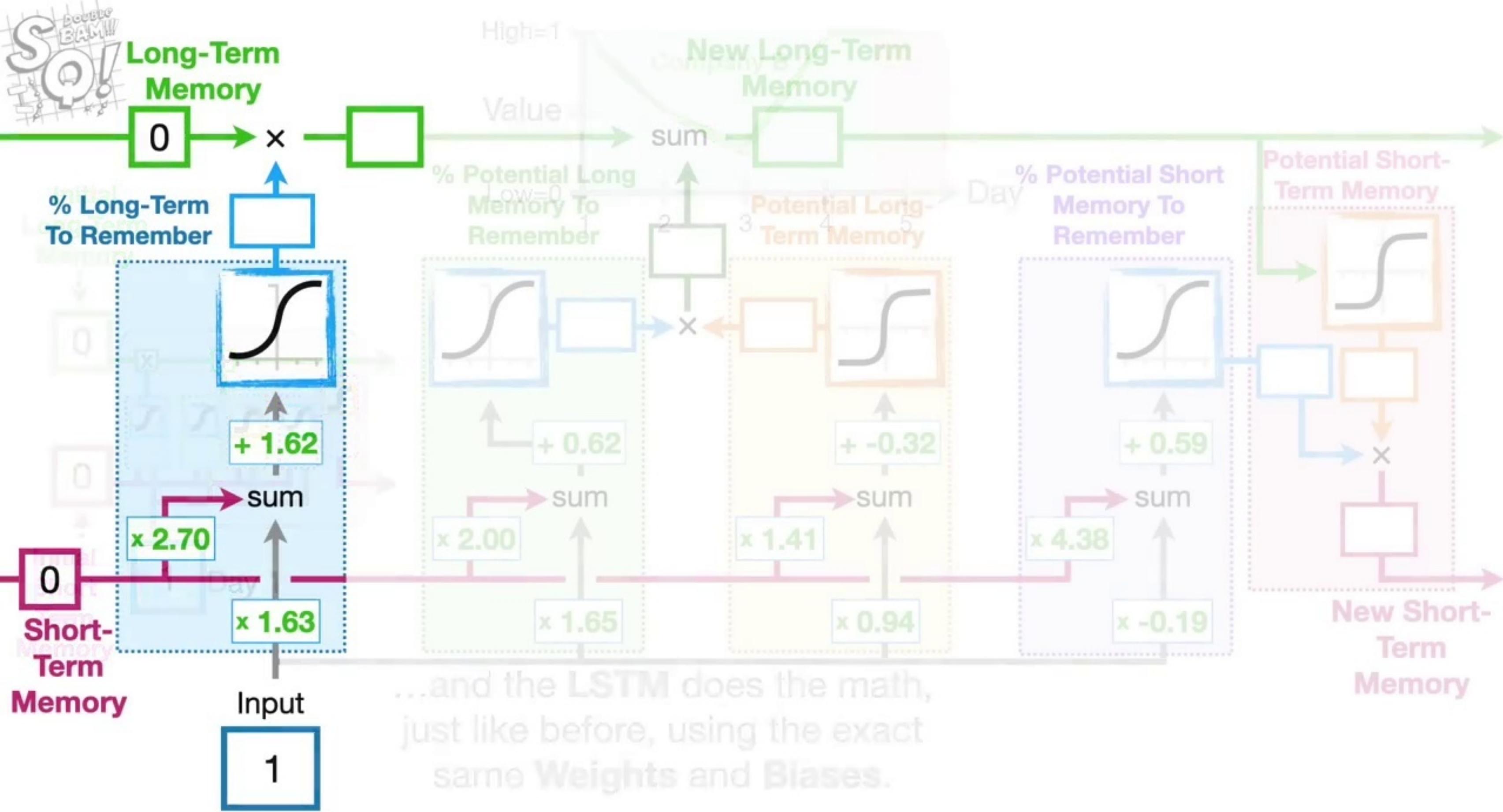
So let's start by initializing the **Long** and **Short-Term Memories** to 0.

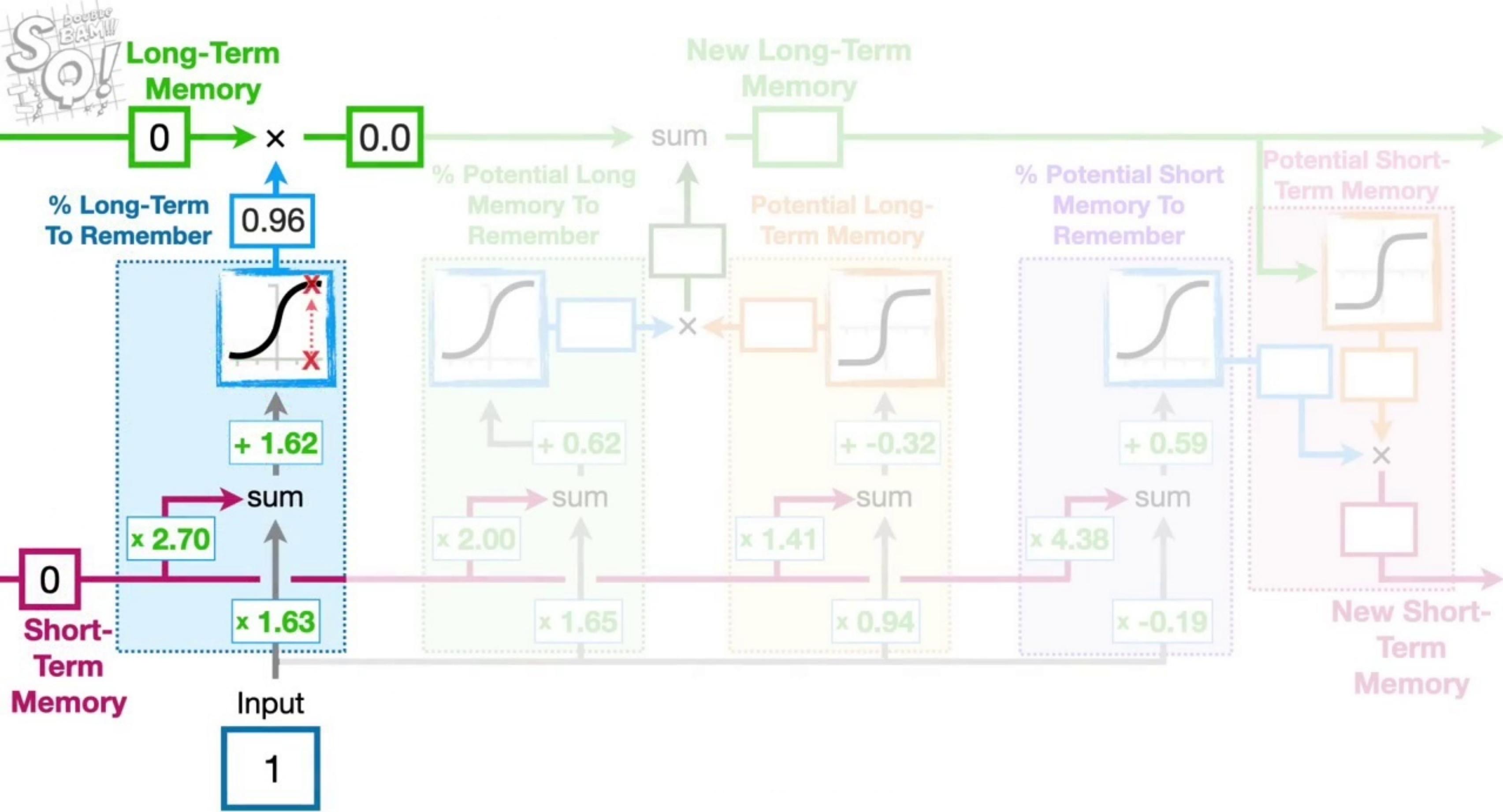


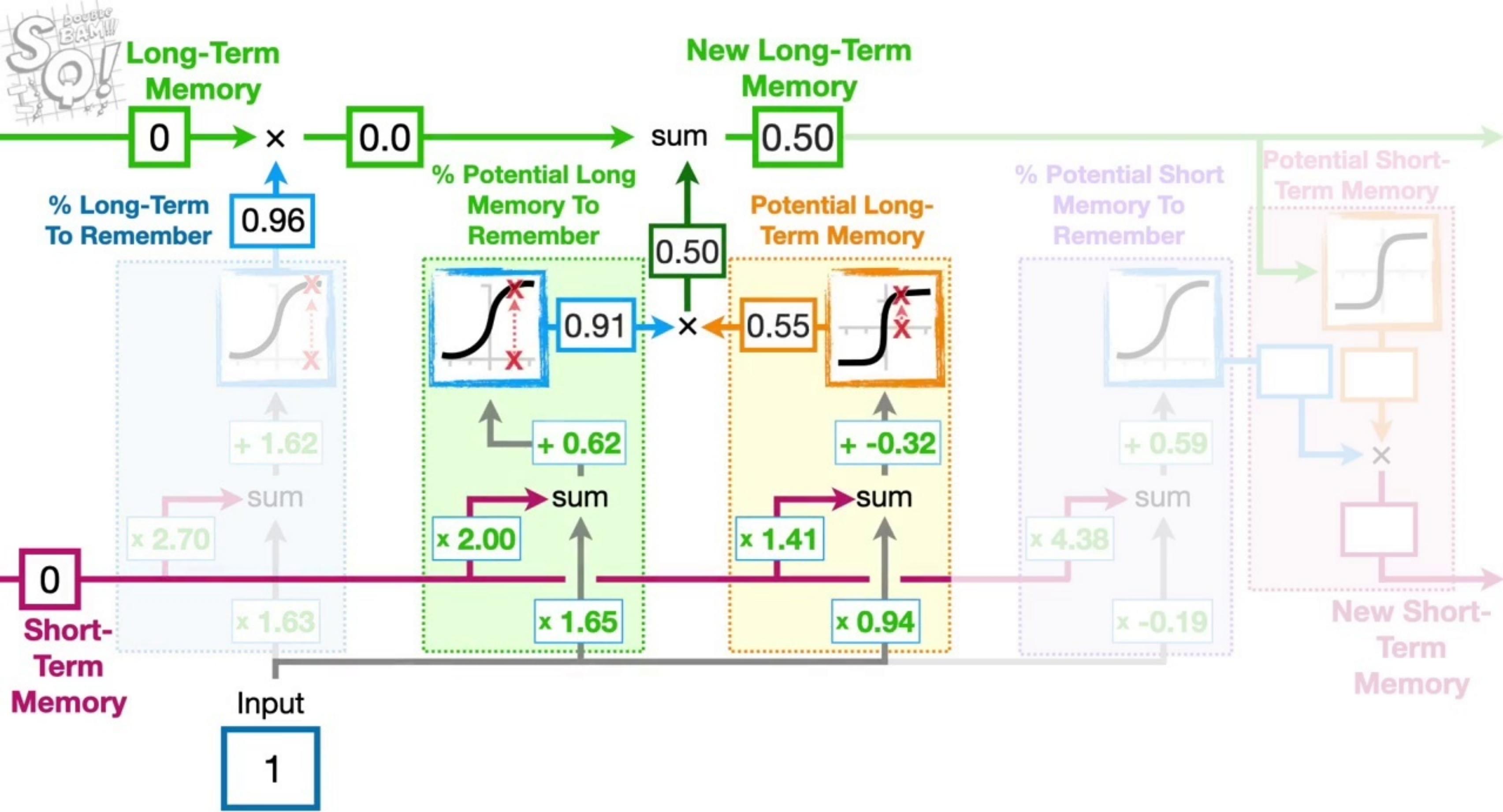
So let's start by initializing the **Long** and **Short-Term Memories** to 0.

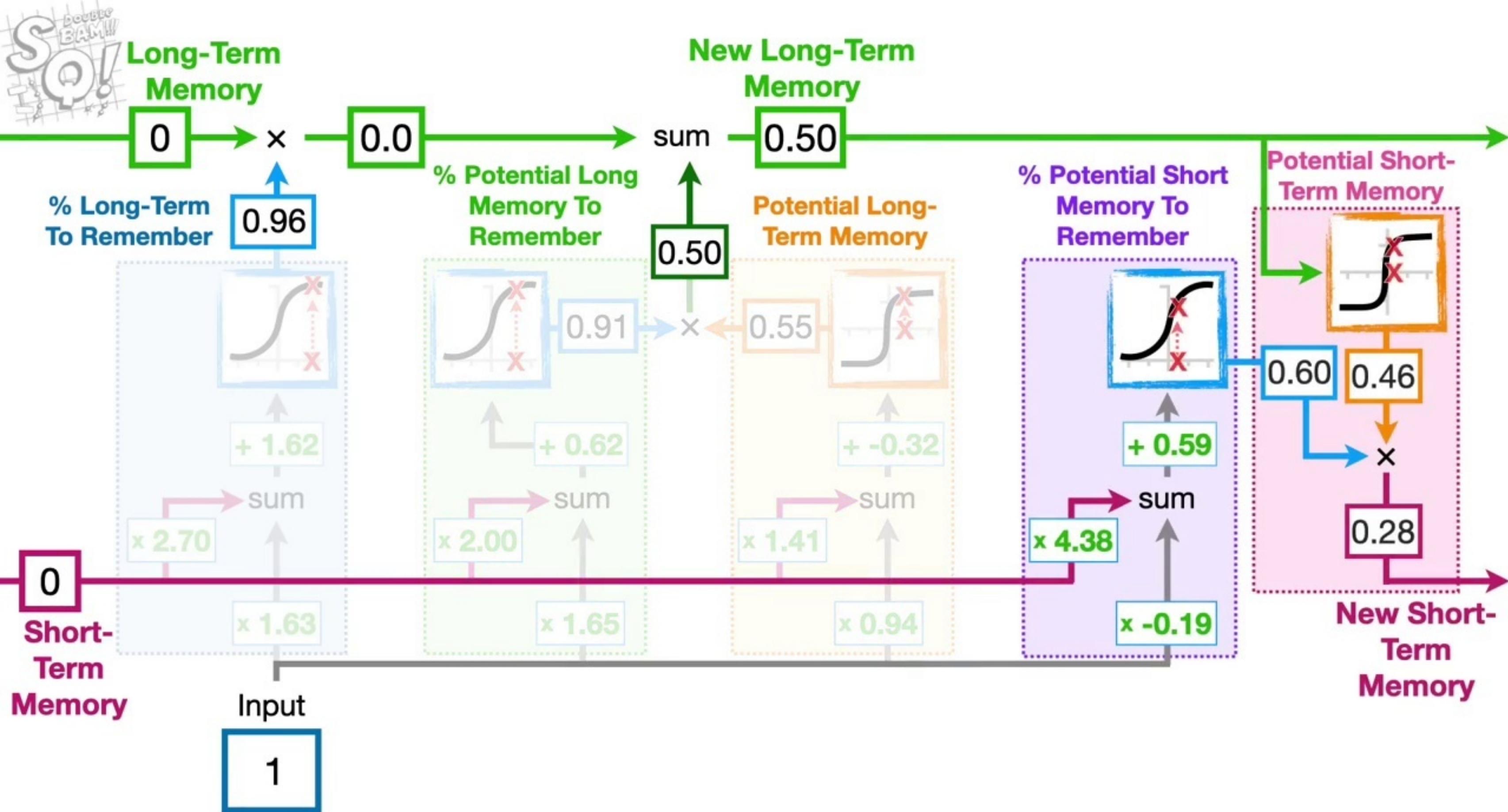


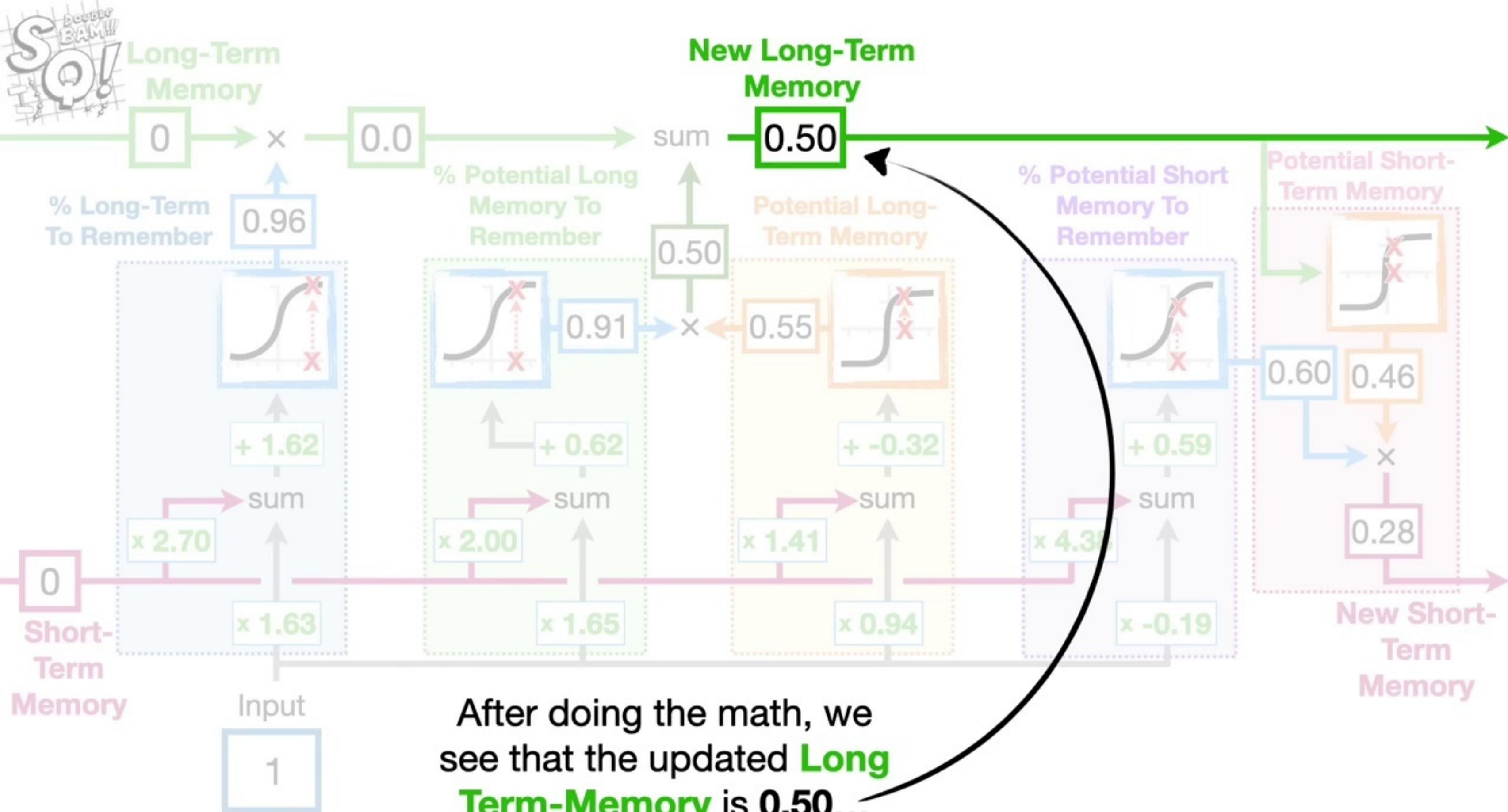
...and the **LSTM** does the math,
just like before, using the exact
same **Weights** and **Biases**.

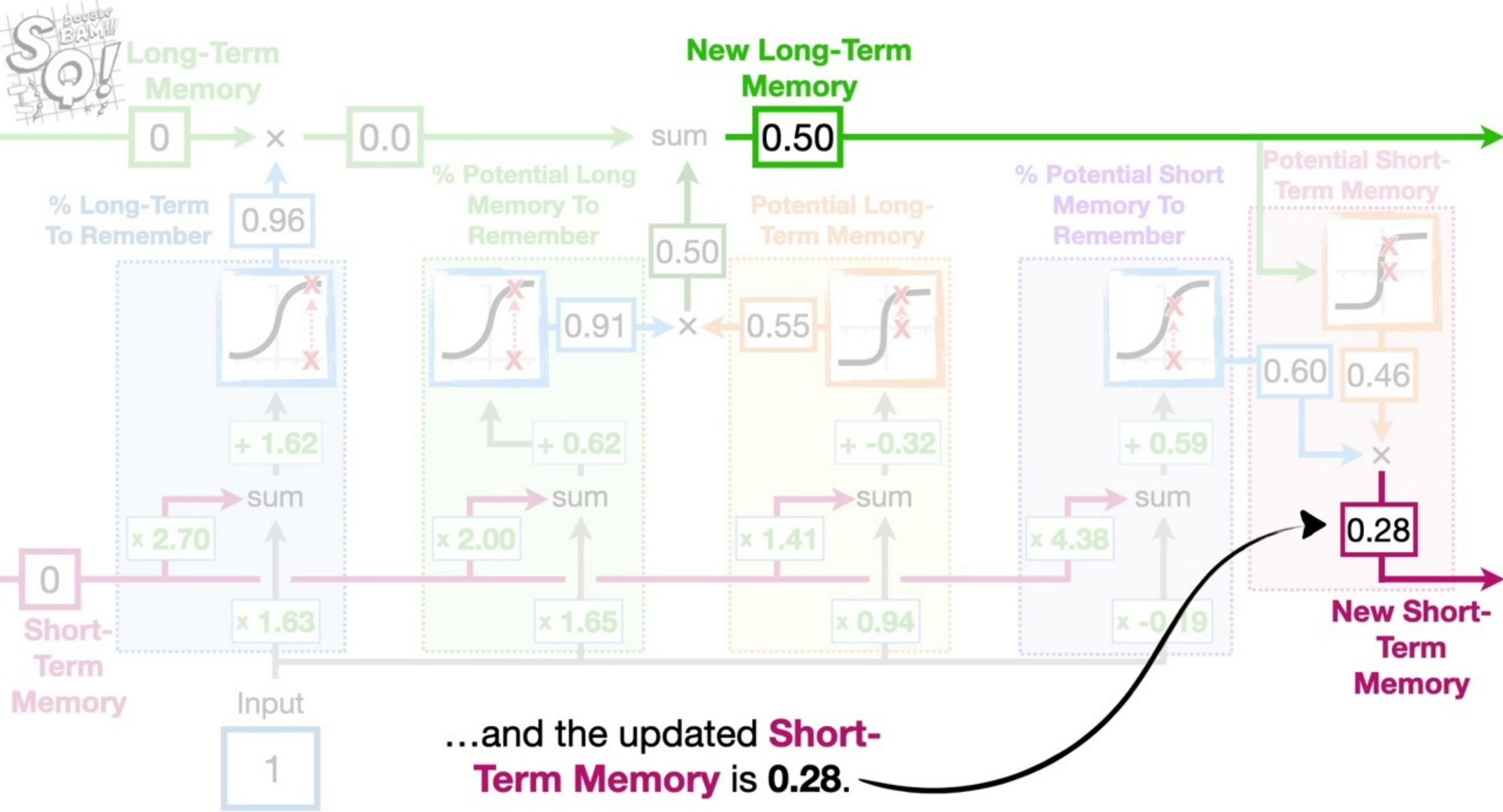


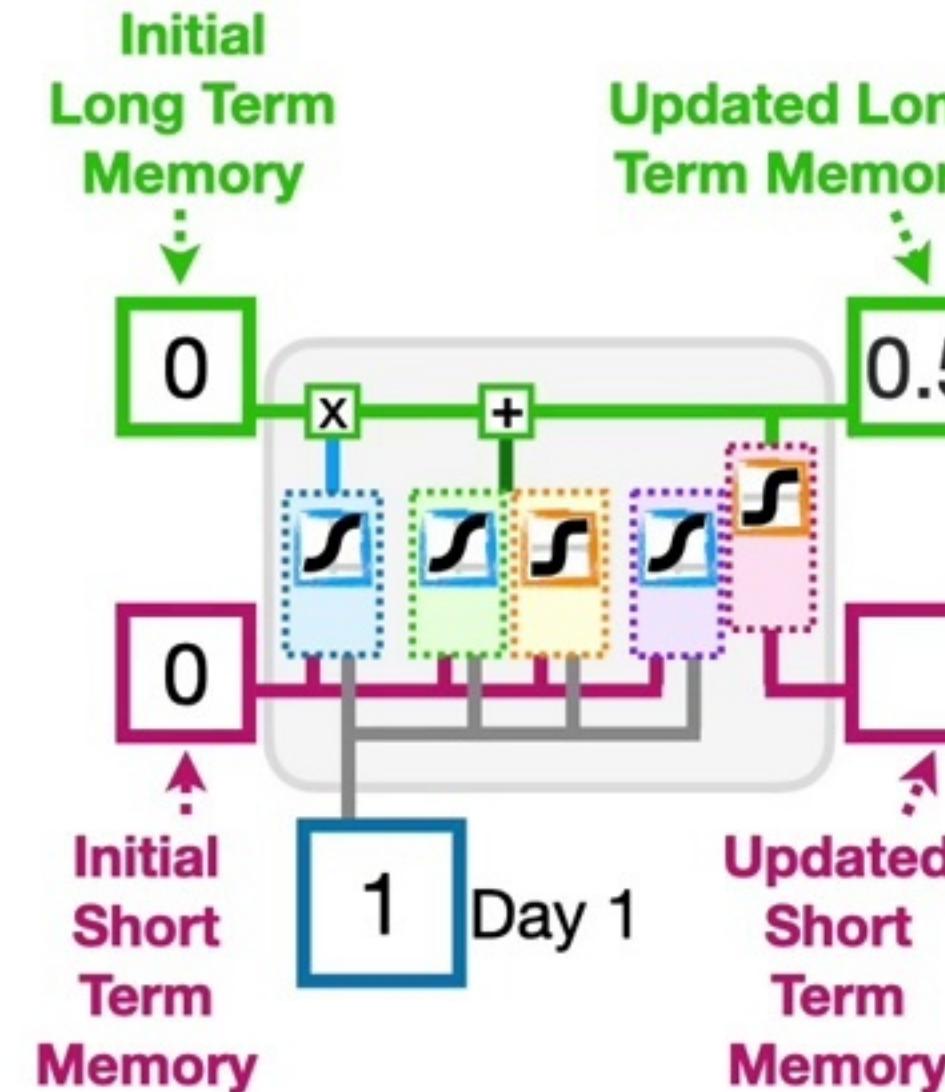
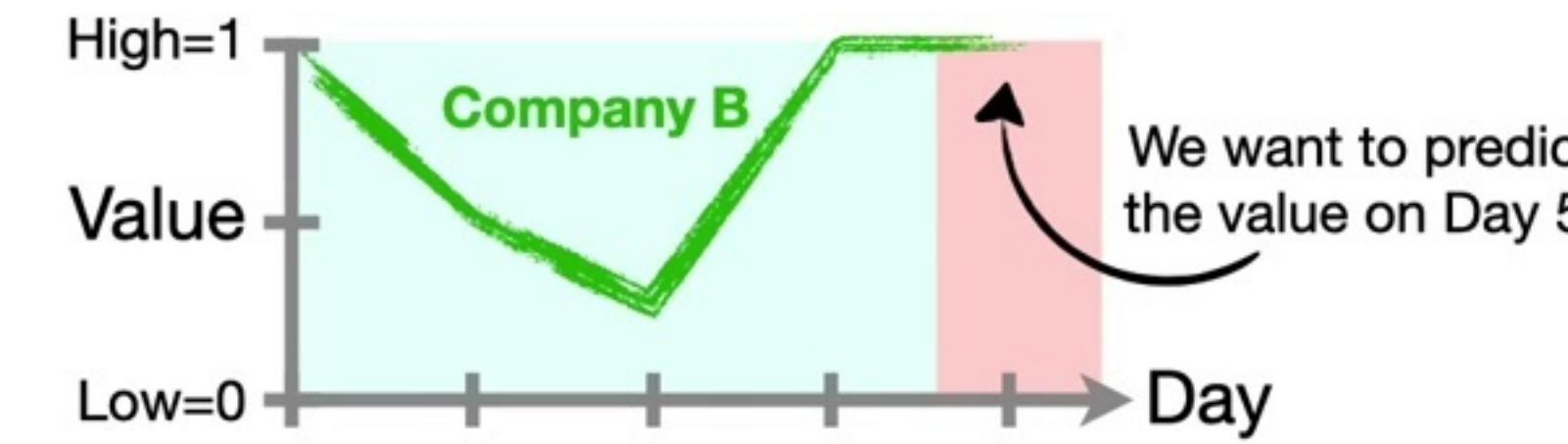




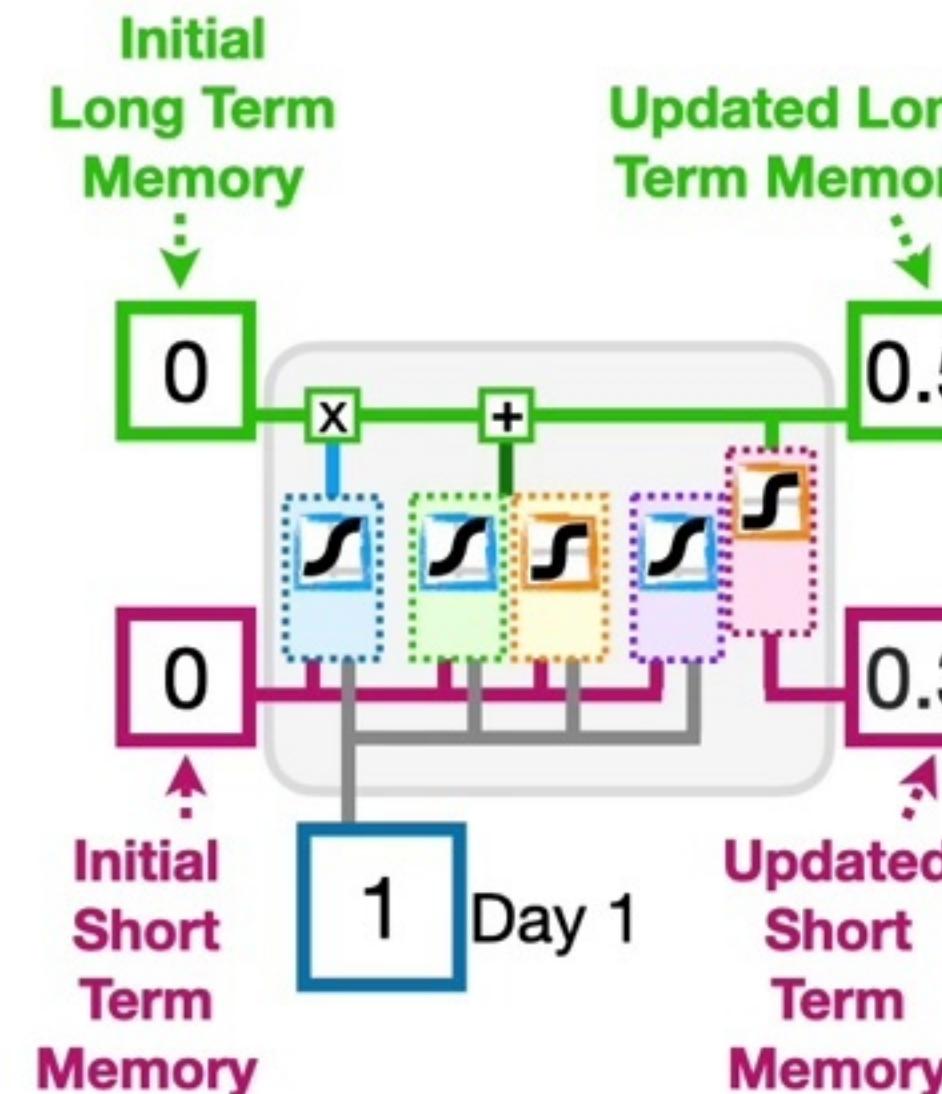




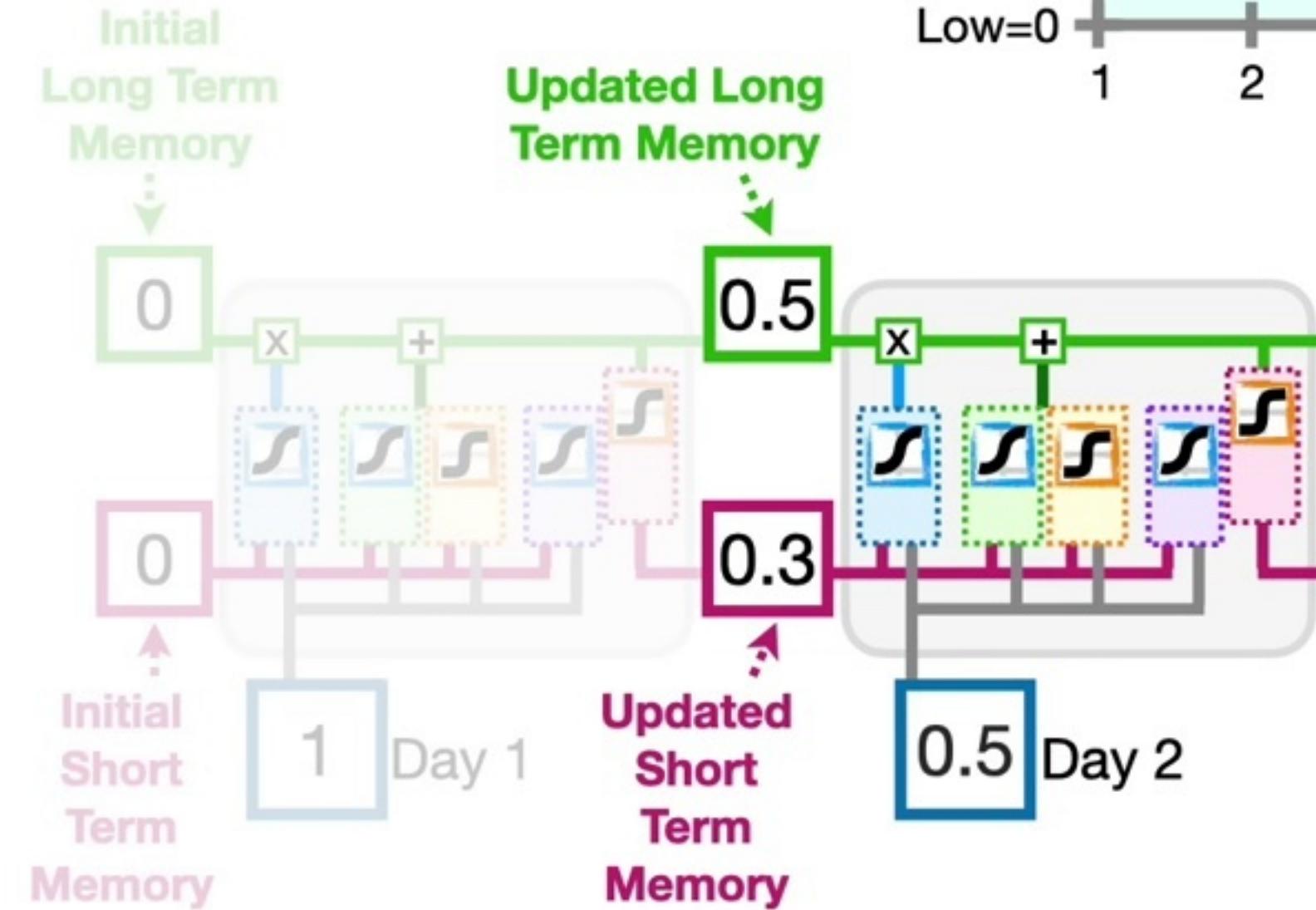
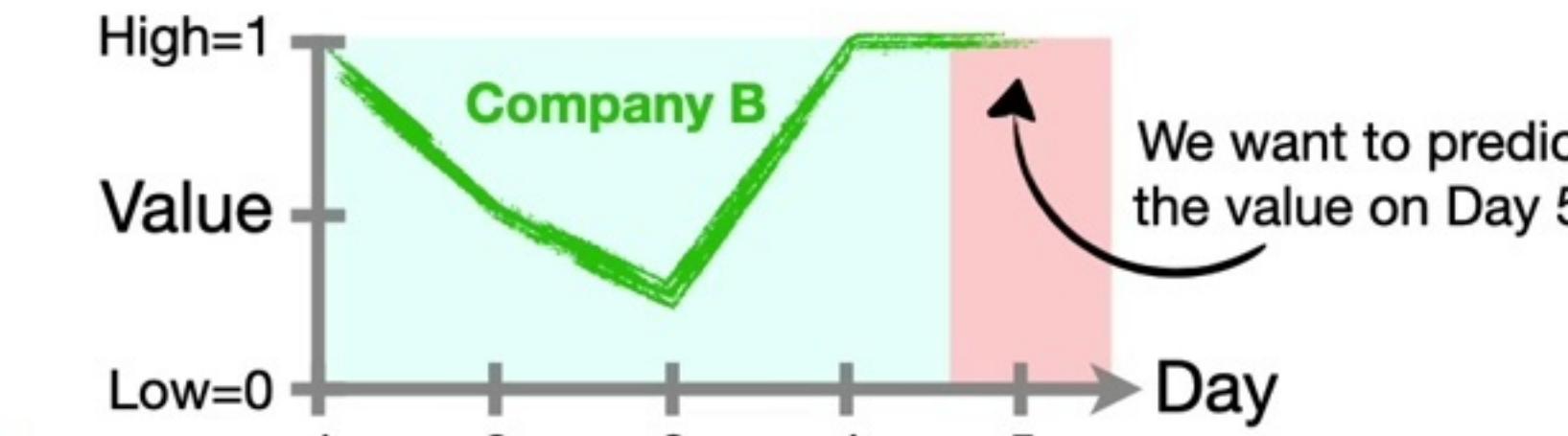




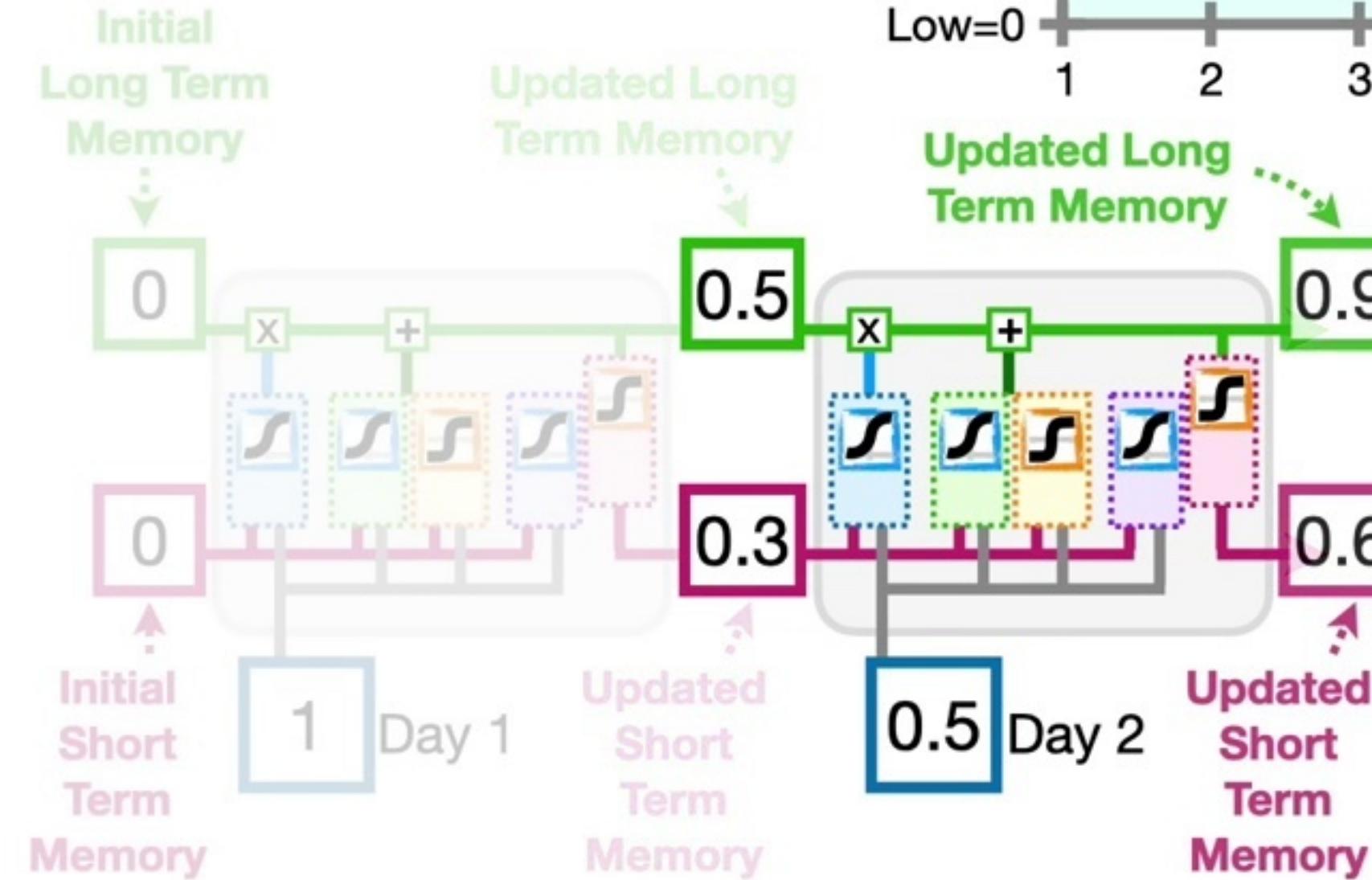
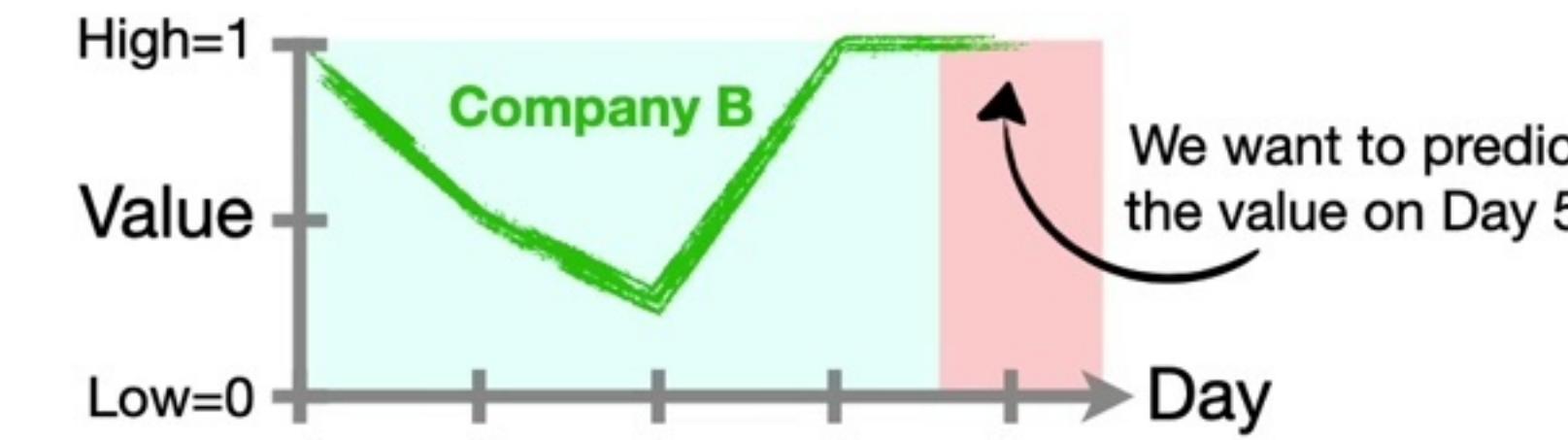
...so we plug in **0.5** for the
updated **Long-Term
Memory**...



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



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



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

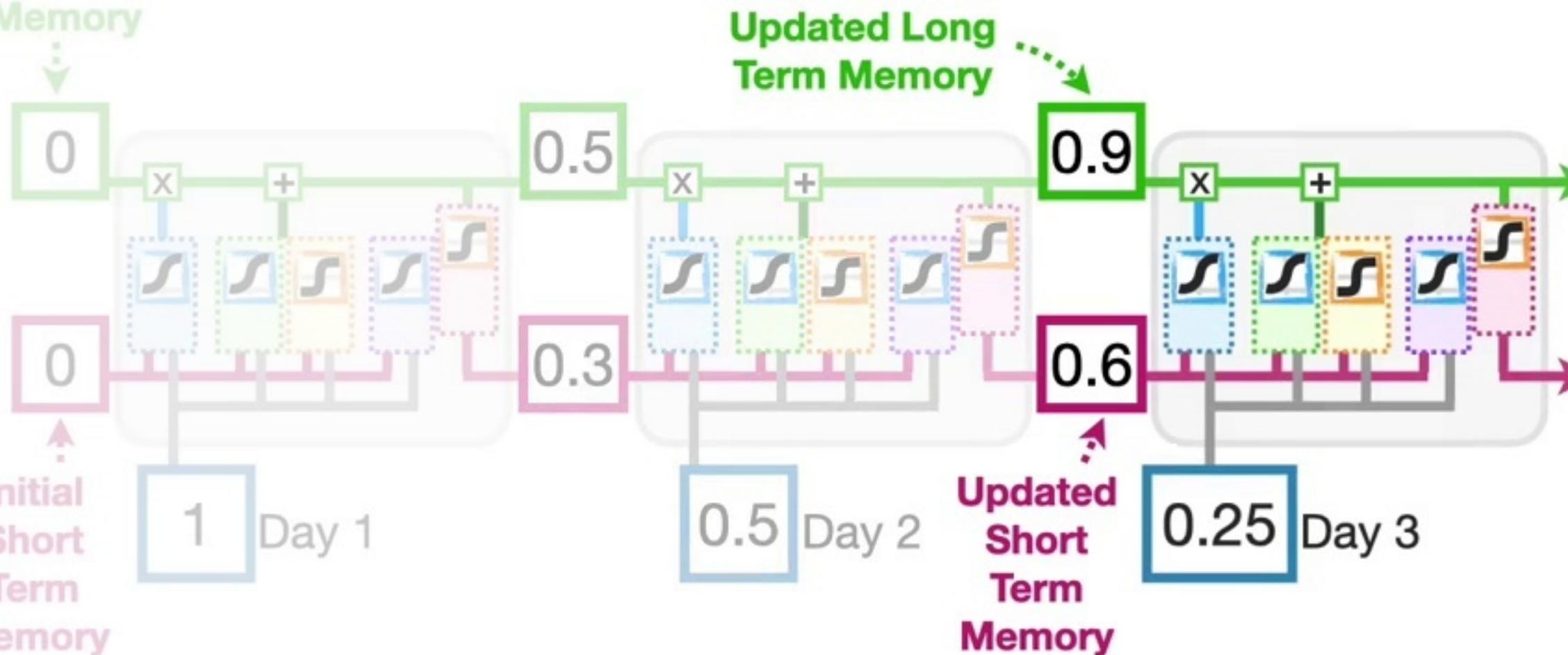


Initial Long Term Memory

↓

Initial Short Term Memory

↑



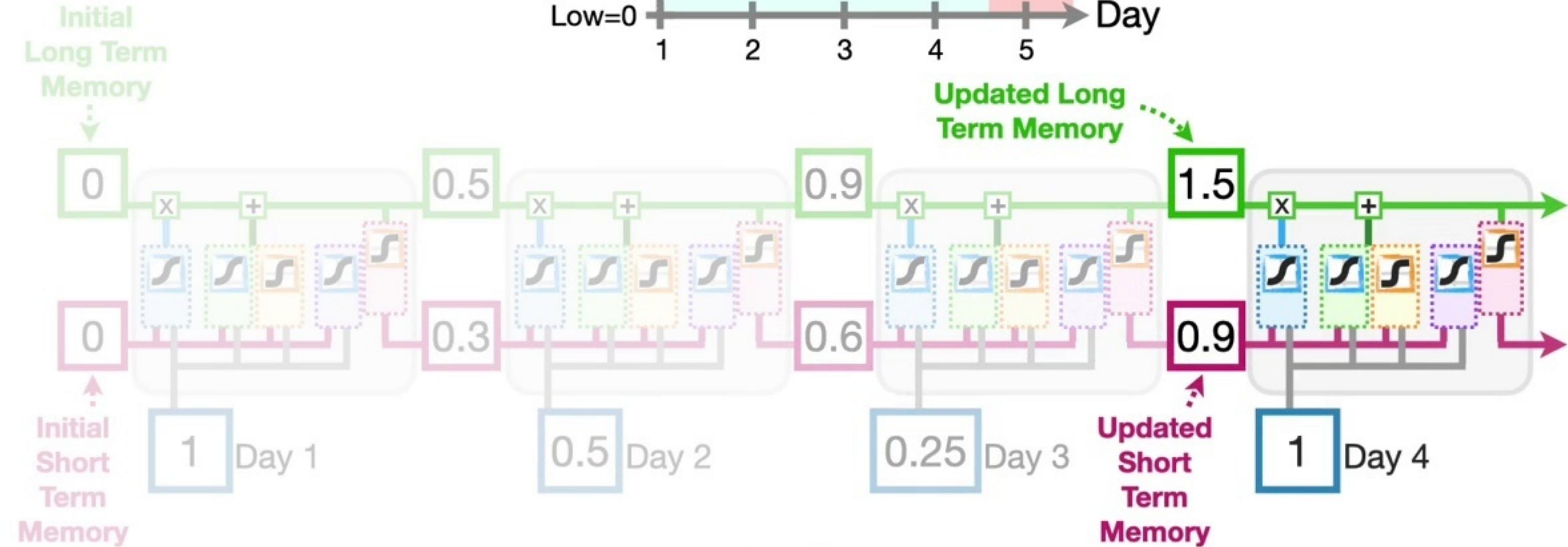
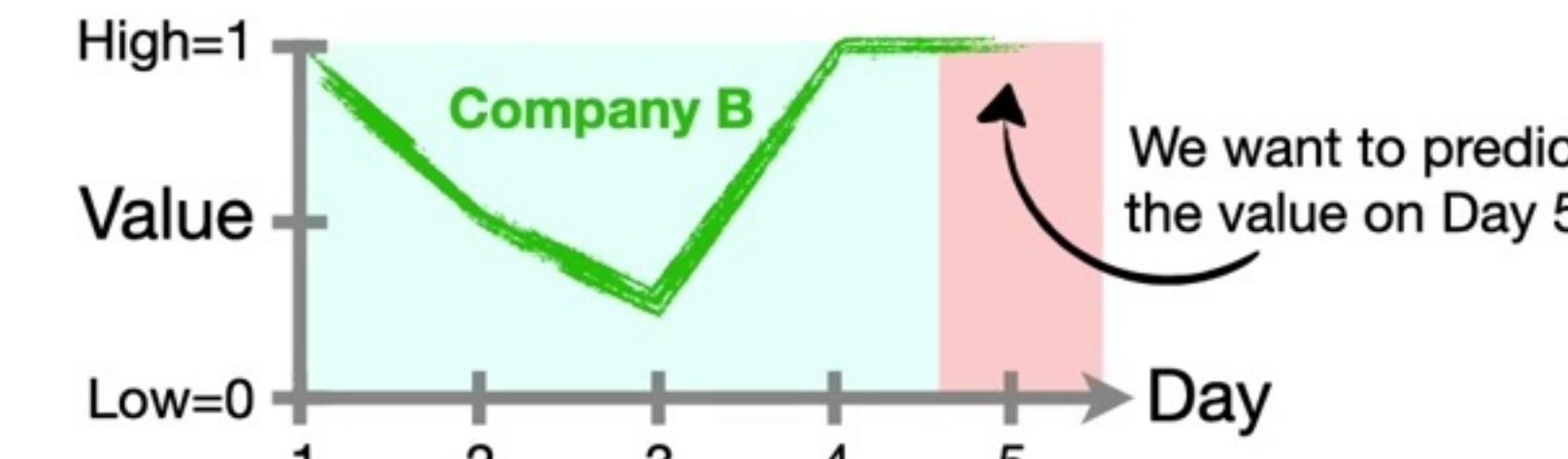
Now we unroll the **LSTM** and do the math with the remaining input values...



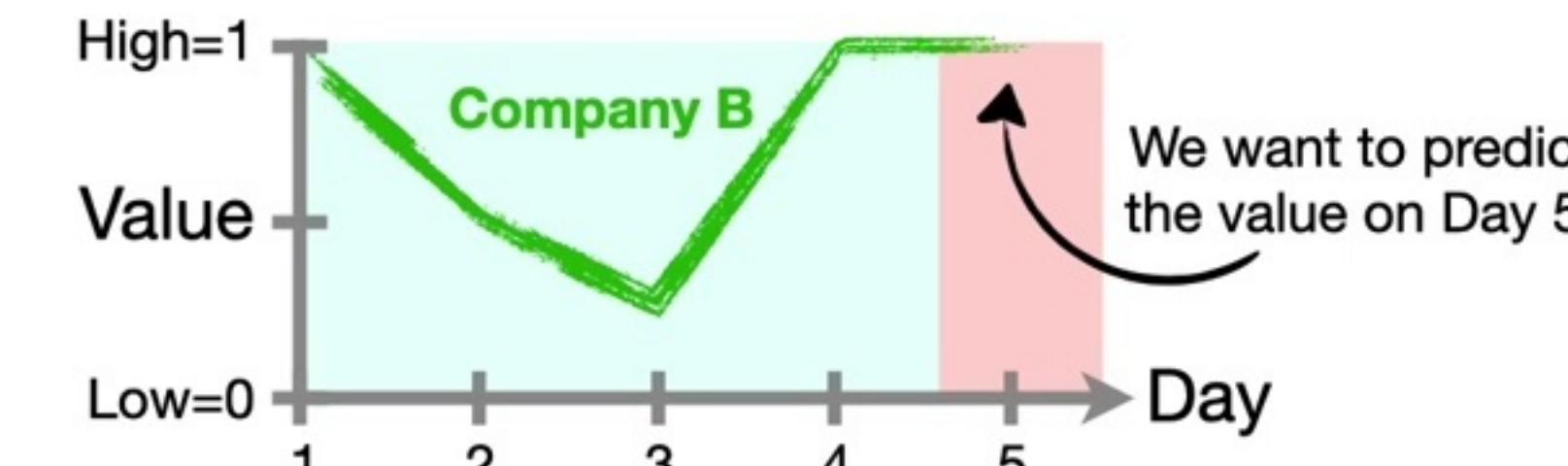
Initial
Long Term
Memory



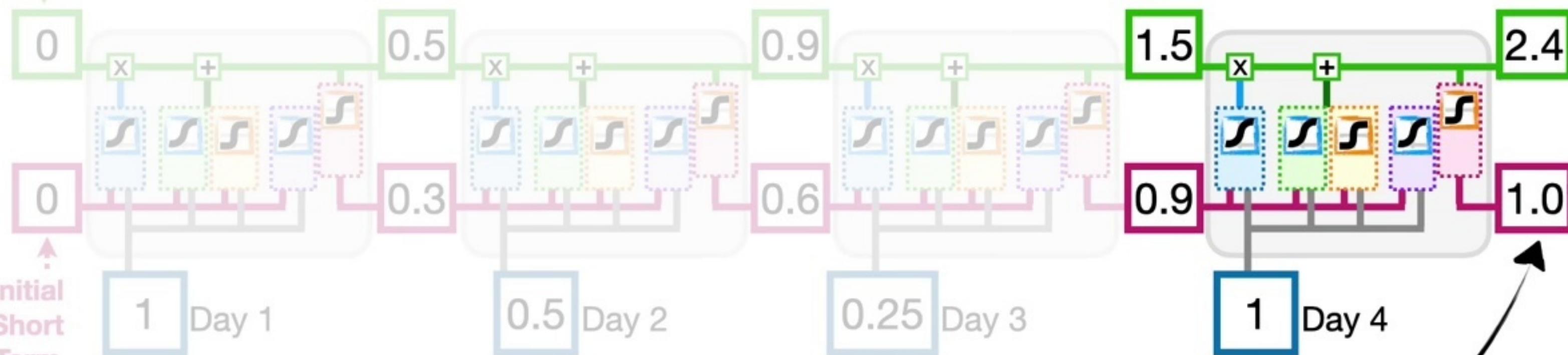
Now we unroll the **LSTM** and do the math with the remaining input values...



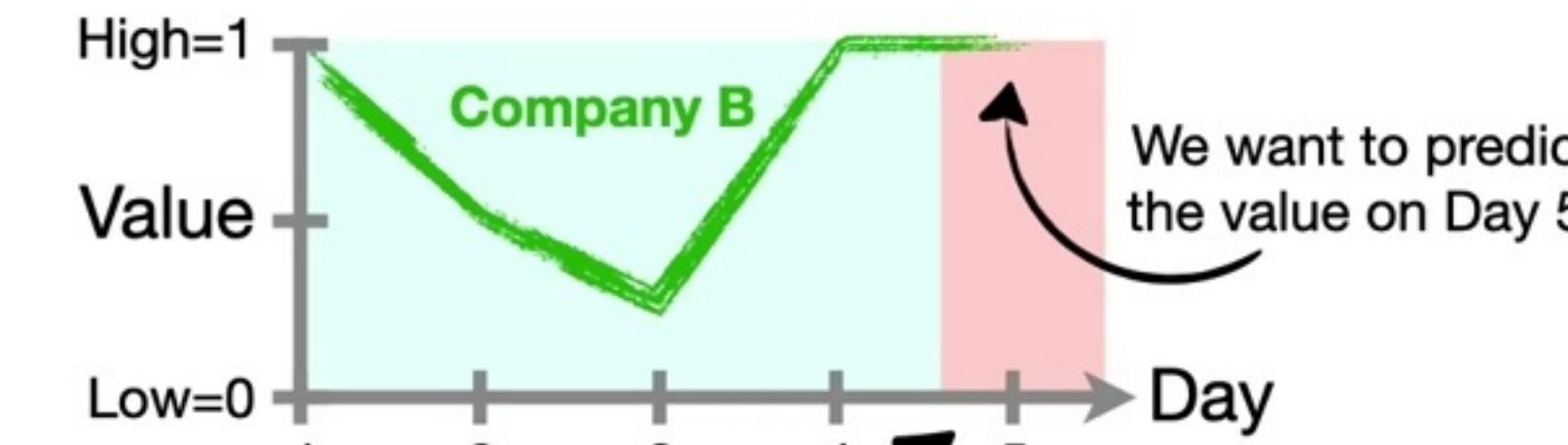
Now we unroll the **LSTM** and do the math with the remaining input values...



Initial
Long Term
Memory



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



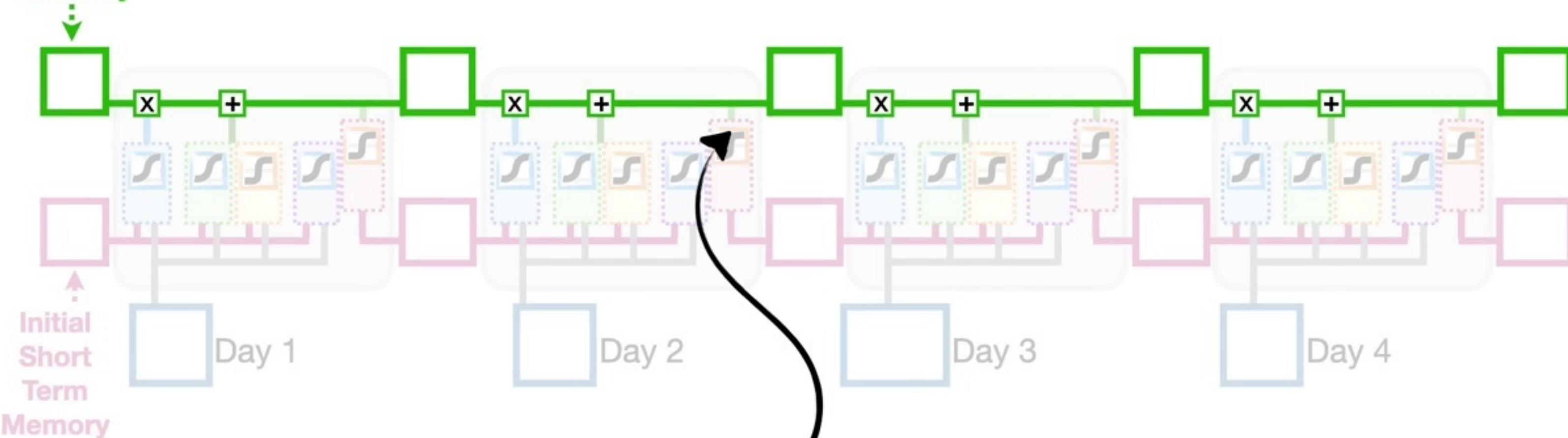
Initial
Long Term
Memory



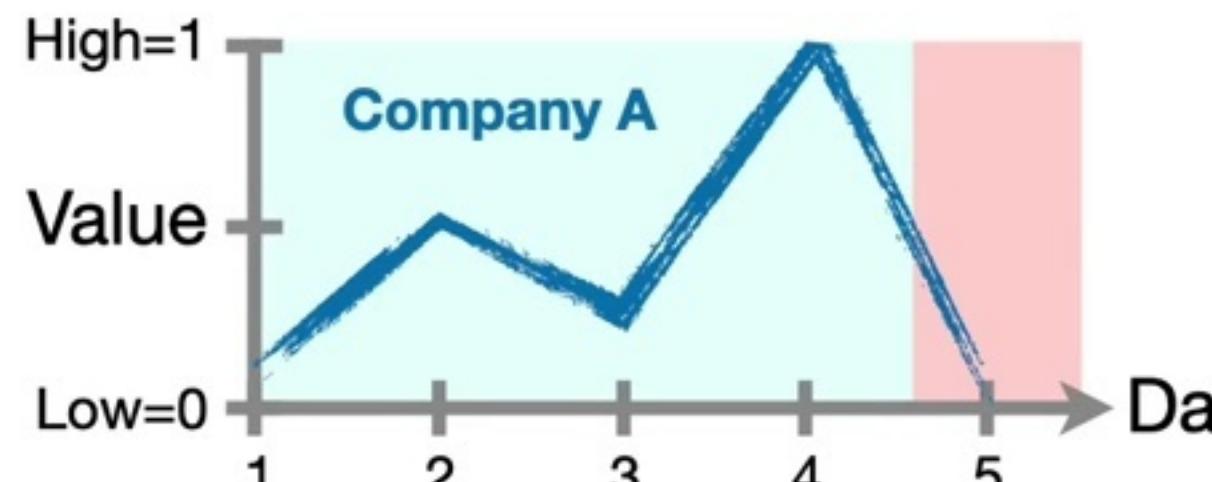
And that means that the **Output** from the **LSTM** correctly predicts **Company B's value for Day 5.**



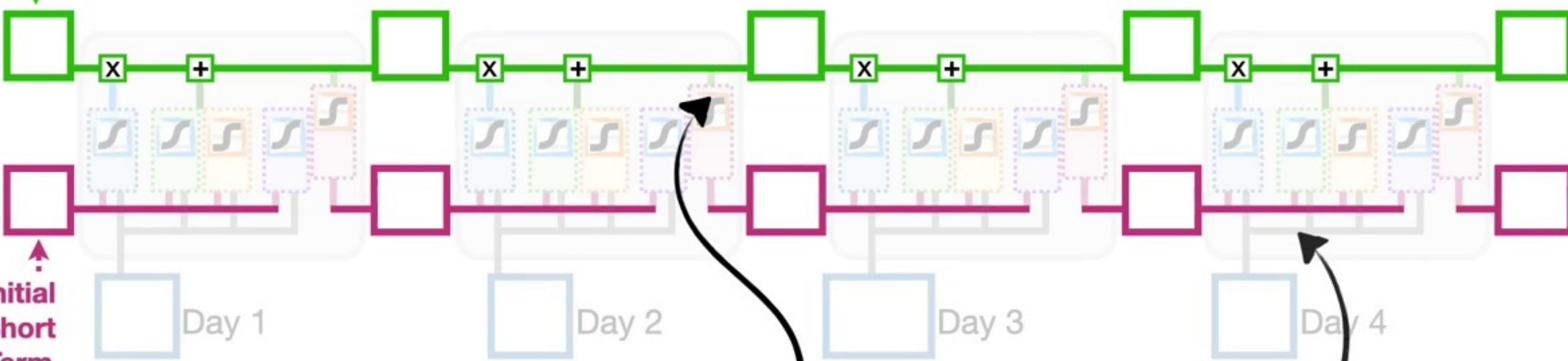
Initial
Long Term
Memory



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



Initial
Long Term
Memory



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

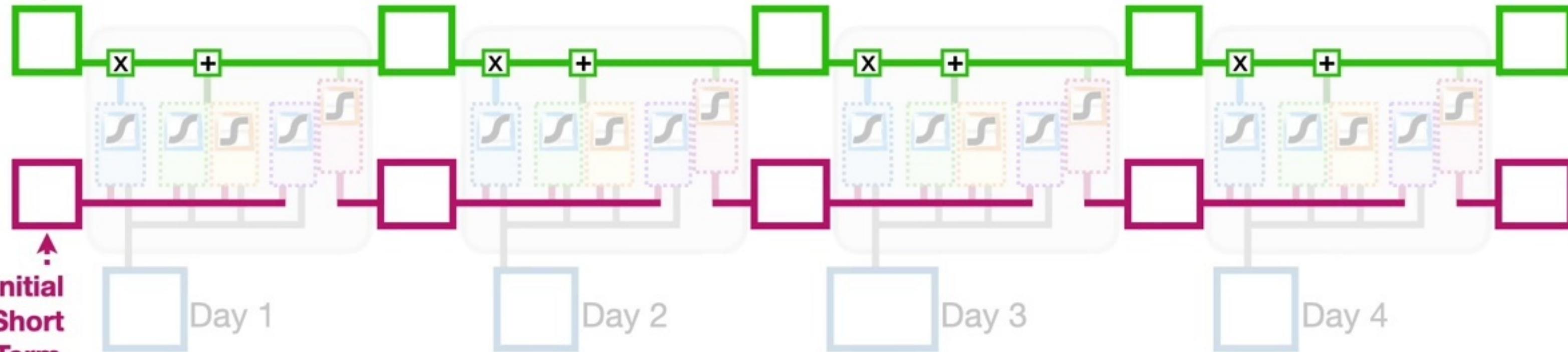
...and **Short-Term
Memories...**



Initial
Long Term
Memory



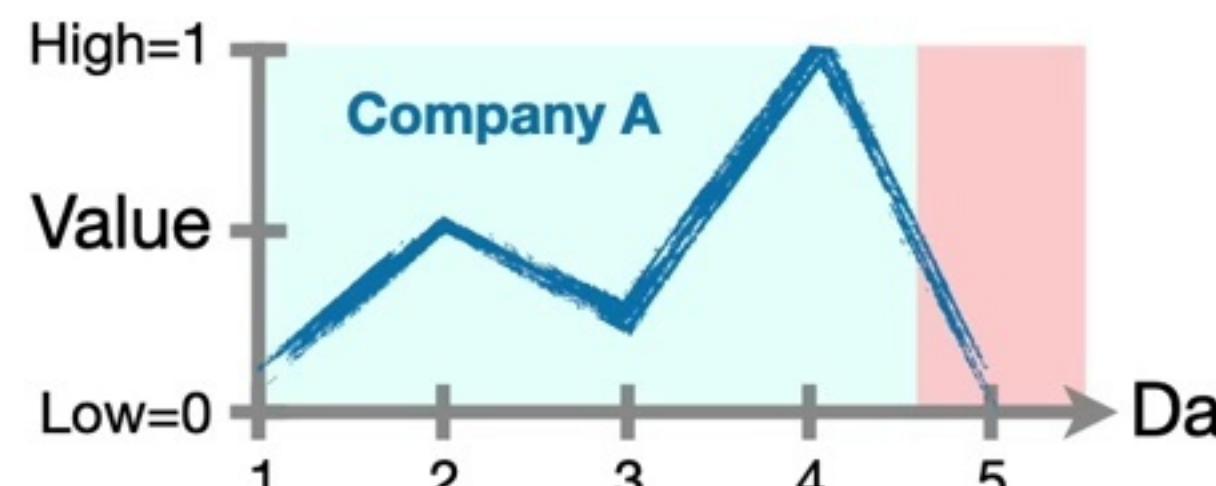
Initial
Short Term
Memory



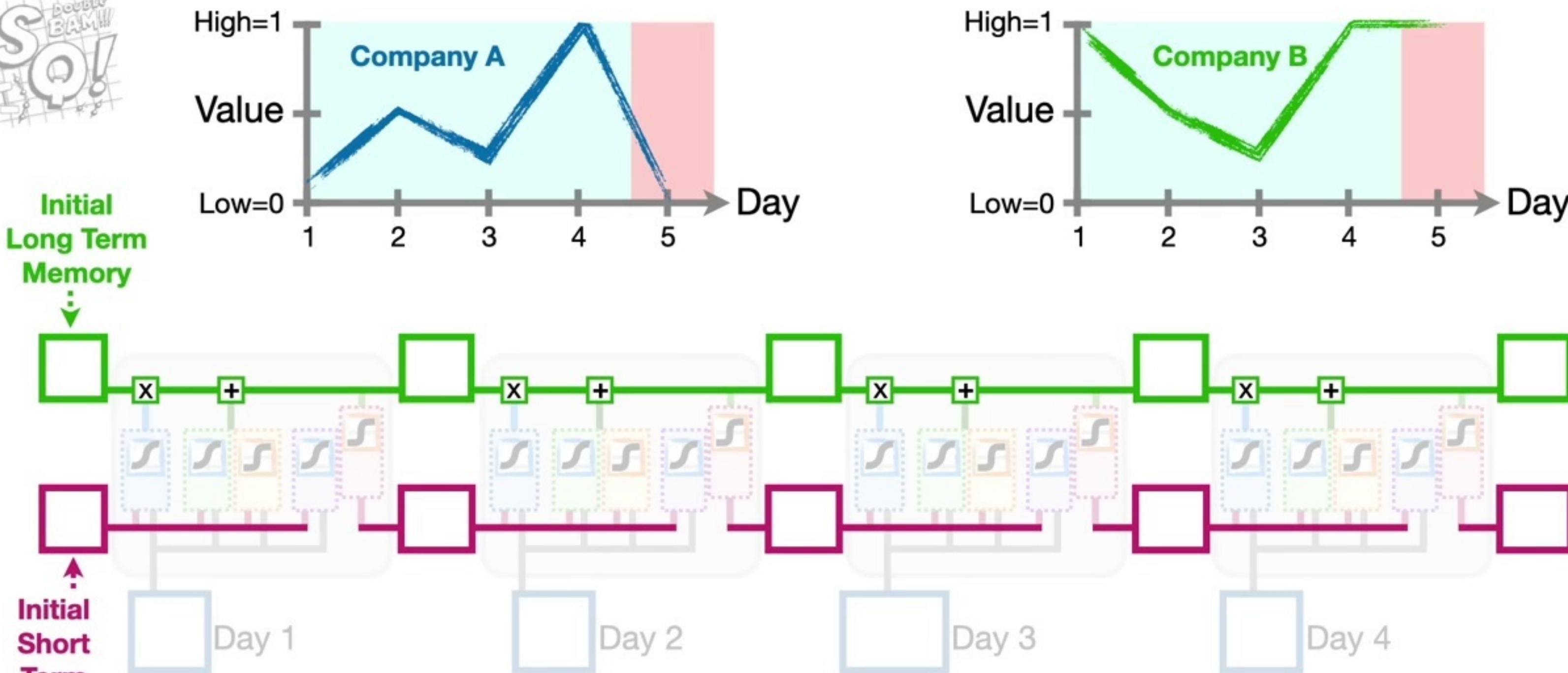
...Long Short-Term Memory
Networks avoid the exploding/
vanishing gradient problem...



Initial
Long Term
Memory



Initial
Short
Term
Memory



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