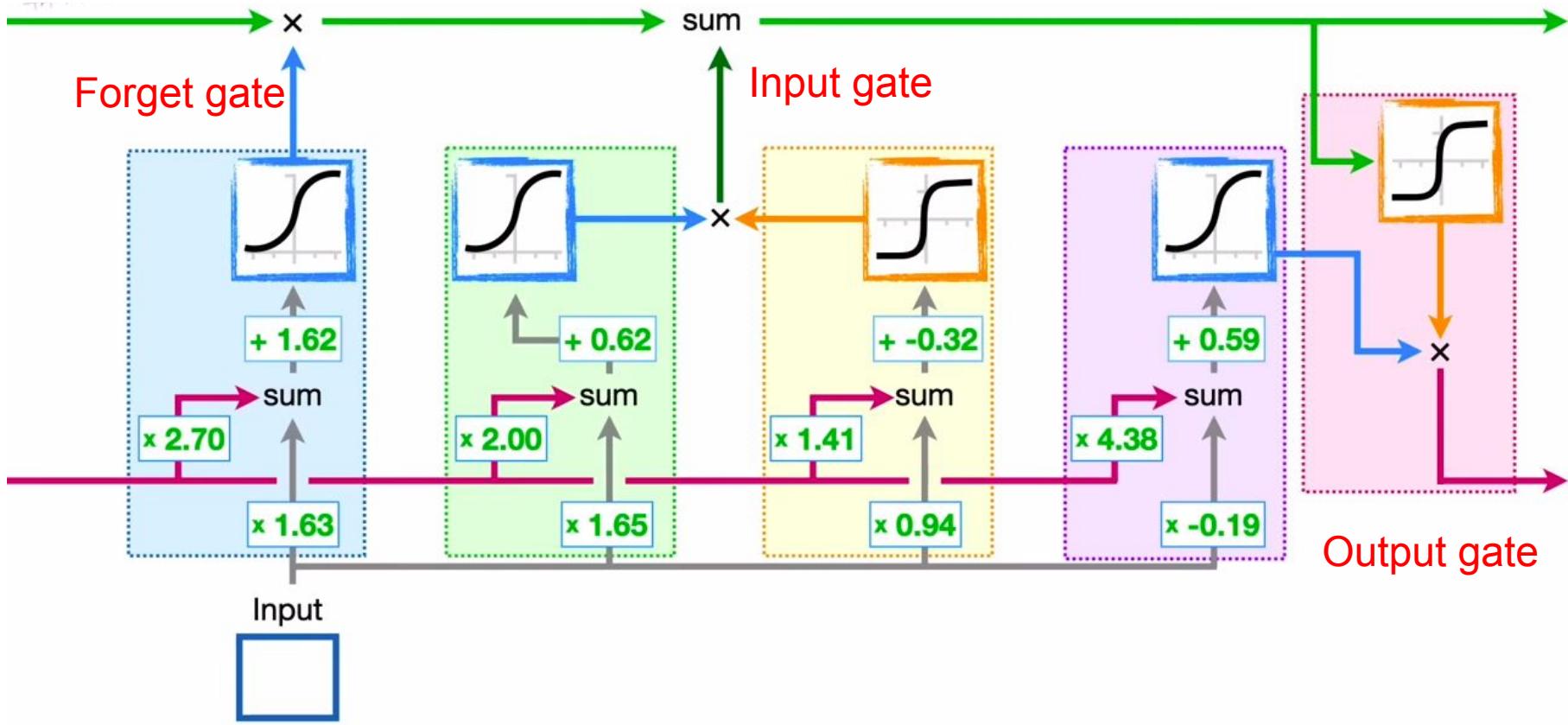


LSTM

Usman Nazir

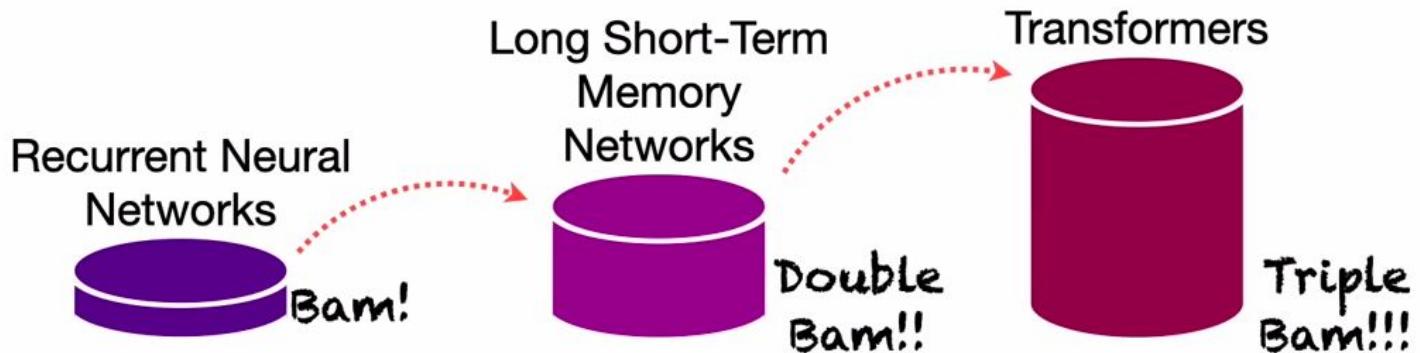


A Always

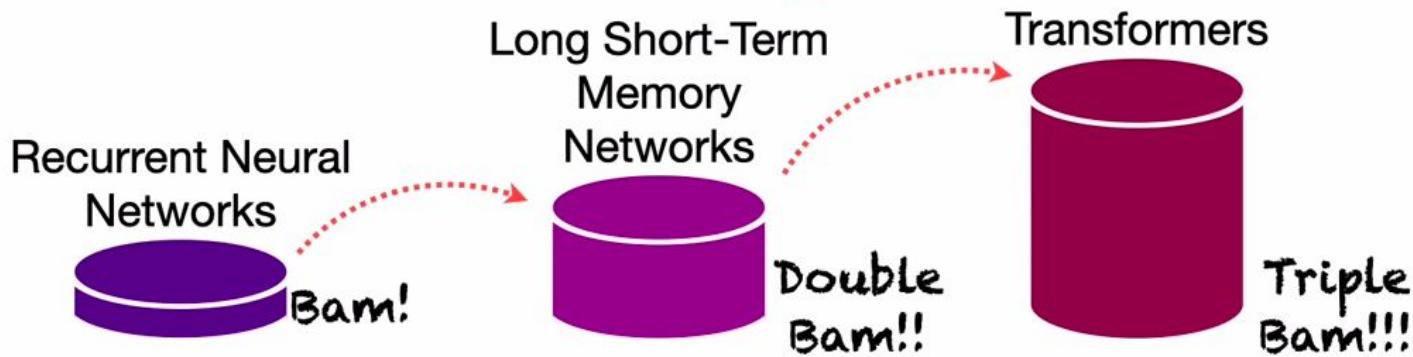
B Be

C Curious

ALSO NOTE: Although **Long Short-Term Memory (LSTM)** is totally awesome, it is also a stepping stone to learning about **Transformers**, which we will talk about in future **StatQuests**.

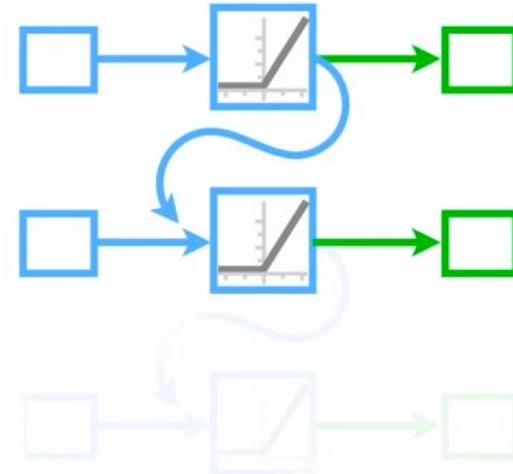


In other words, today we're taking the second step in our '**Quest**'.



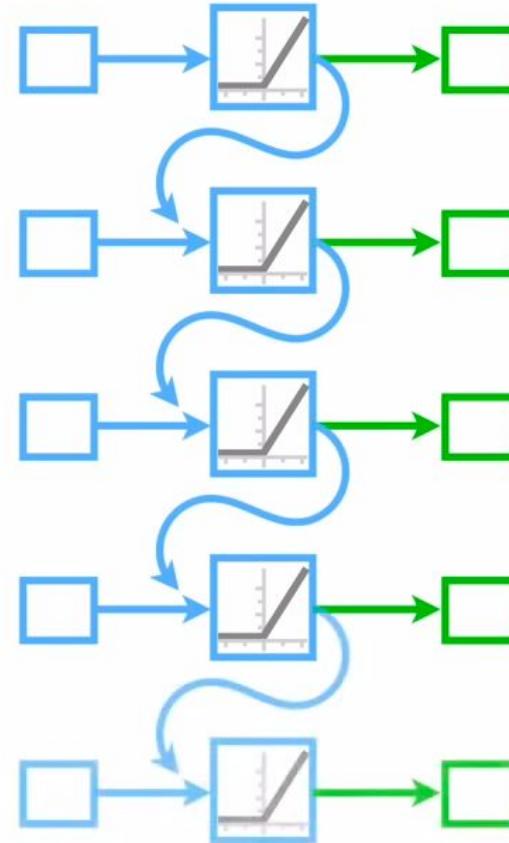


...to **unroll** a network that
works well with *different*
amounts of sequential data.



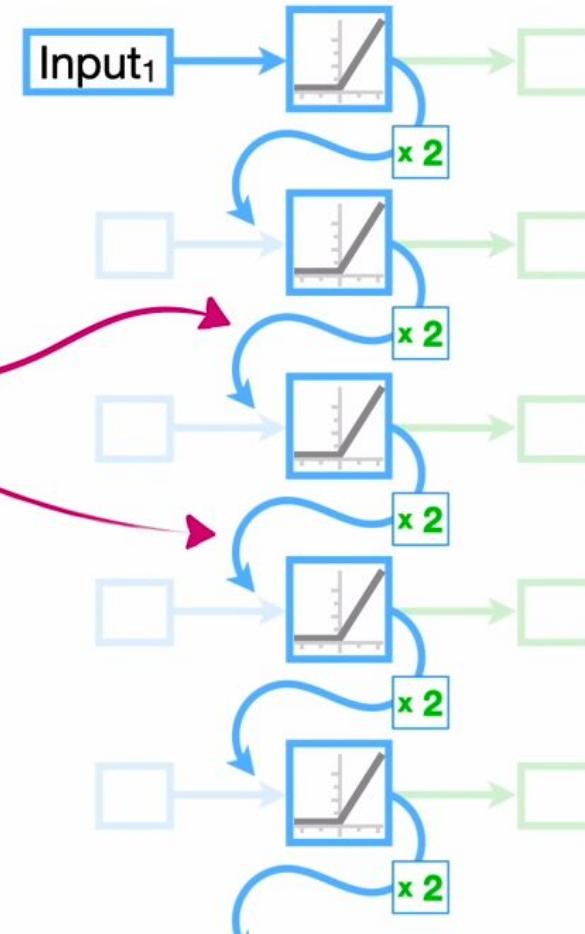


...to **unroll** a network that
works well with *different*
amounts of **sequential data**.





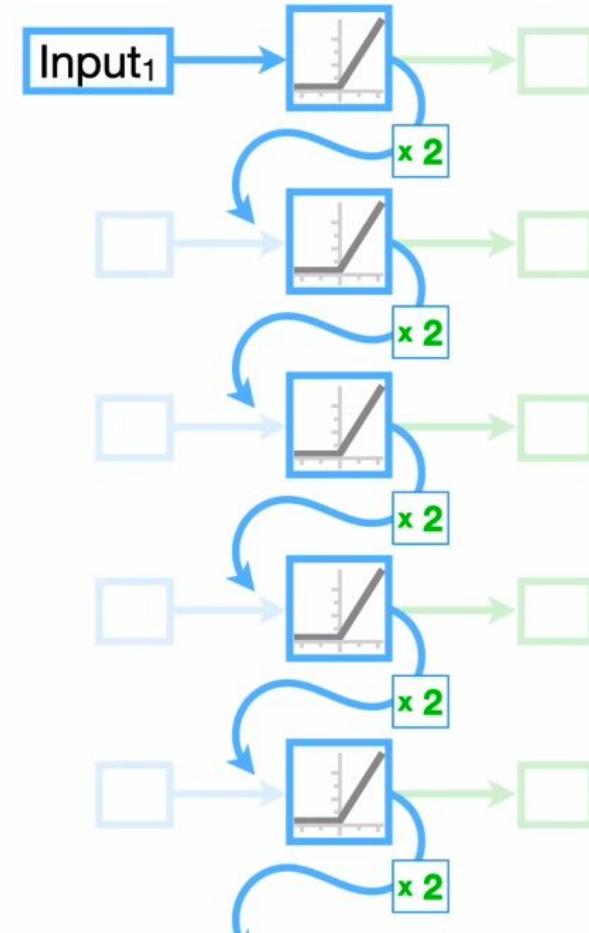
...when the **Weight** on the feedback loop is greater than 1, and in this example the **Weight** is 2...





...then, when we do the math,
we end up multiplying the
Input by the **Weight**, which in
this case is **2**, raised to the
number of times we unrolled
the network.

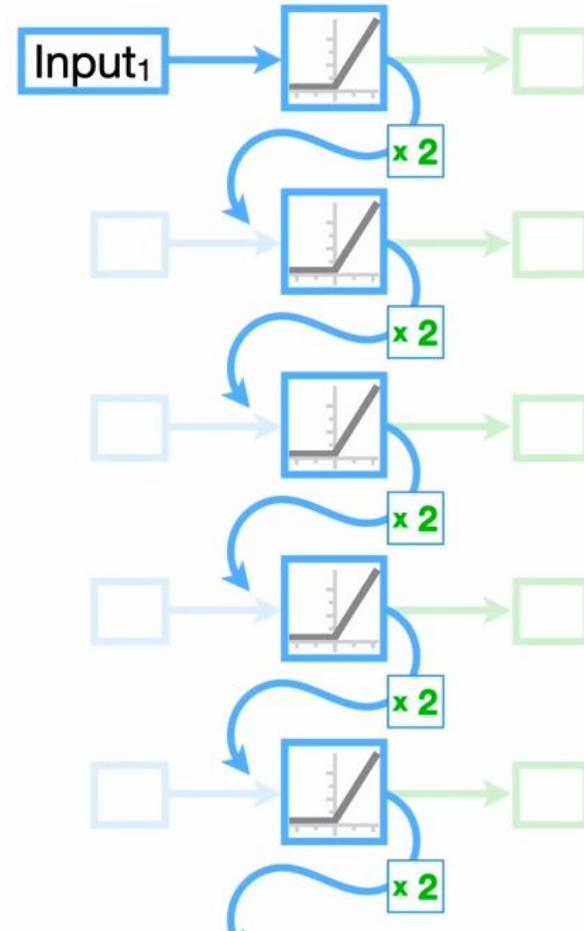
$$= \text{Input}_1 \times 2^{\text{Num. Unroll}}$$





And thus, if we had **50** sequential data points, like **50** days of stock market data, which isn't really that much...

= **Input₁** × **2Num. Unroll**

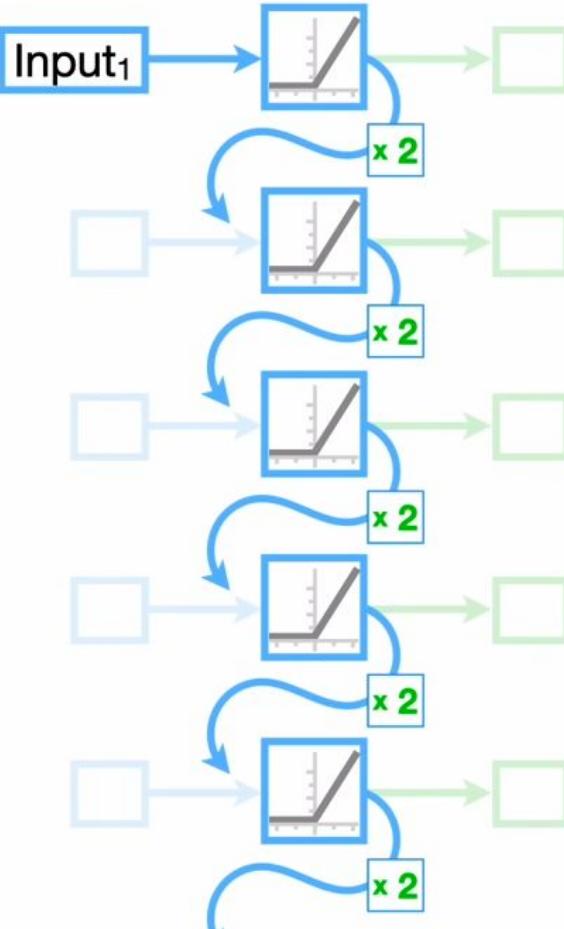




And thus, if we had **50** sequential data points, like **50** days of stock market data, which isn't really that much...

...then we would raise
2 by 50...

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

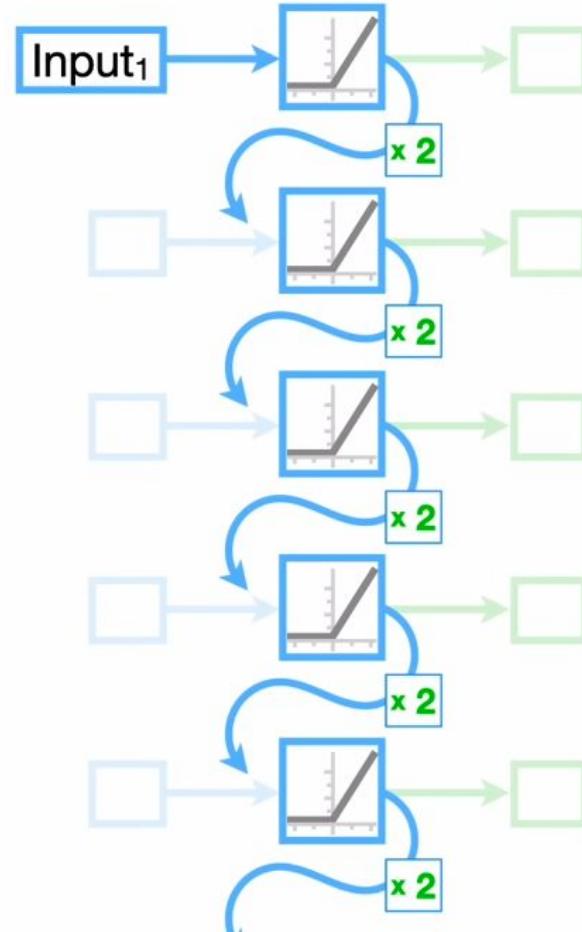




...and this **HUGE NUMBER**
would cause the gradient,
which we need for **Gradient
Descent**...

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

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

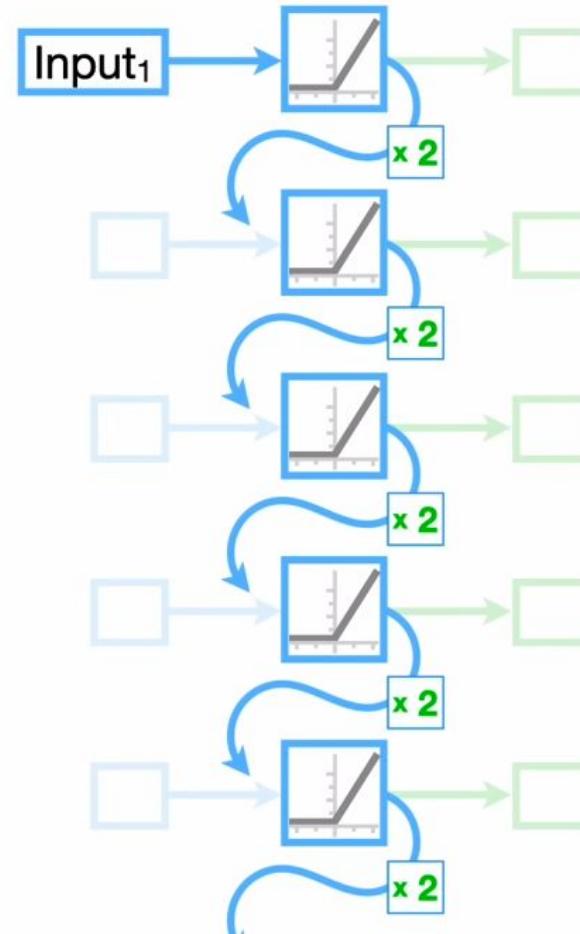




...to **EXPLODE!!!**

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

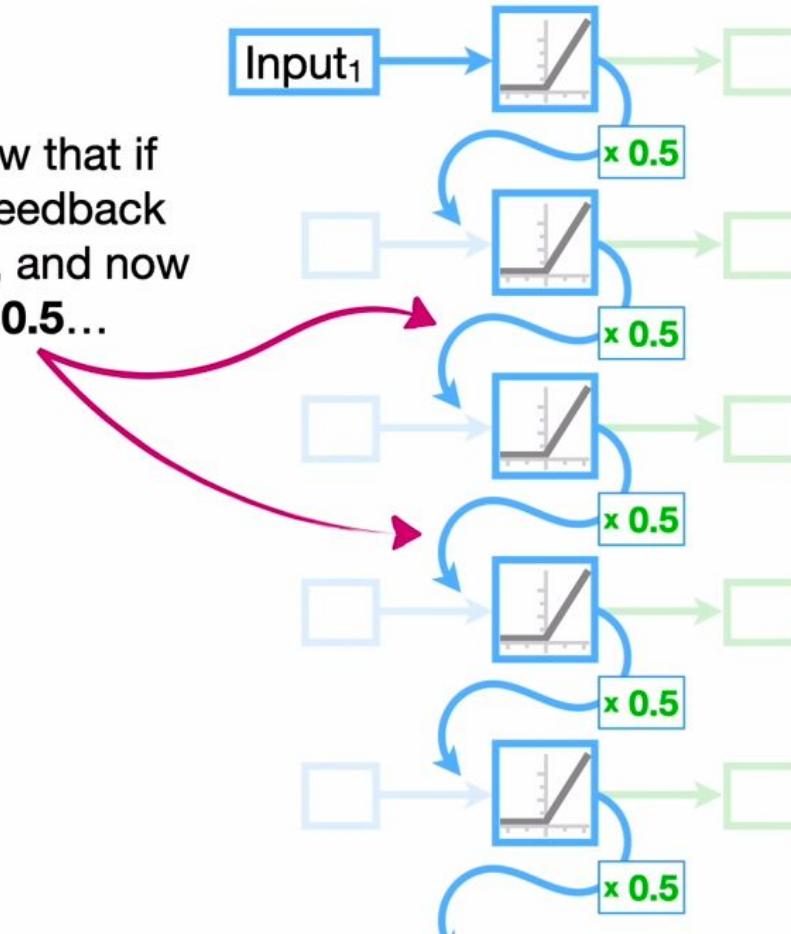
$$= \text{Input}_1 \times \text{A HUGE NUMBER}$$





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...

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

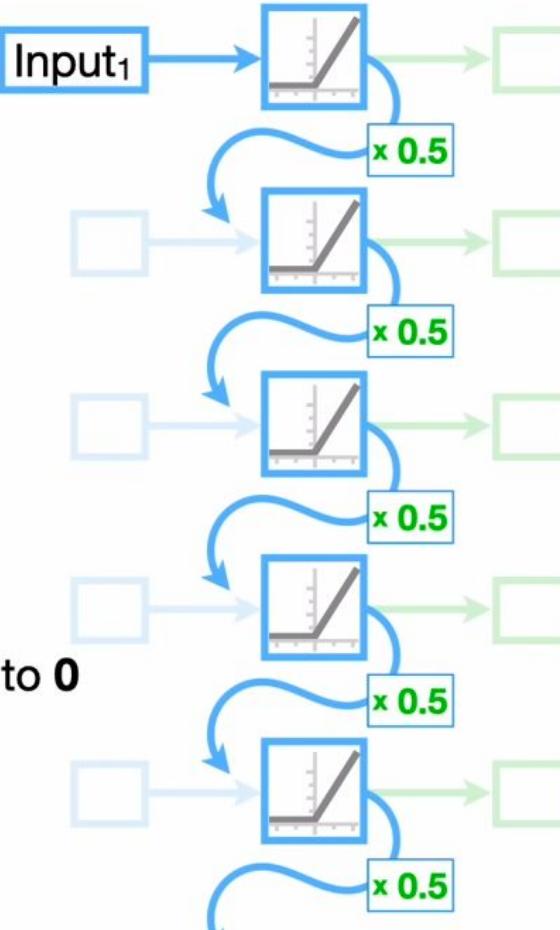




...and **0.5** raised to the **50th** power is a number super close to **0**...

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

= **Input₁** × a number super close to 0





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

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

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

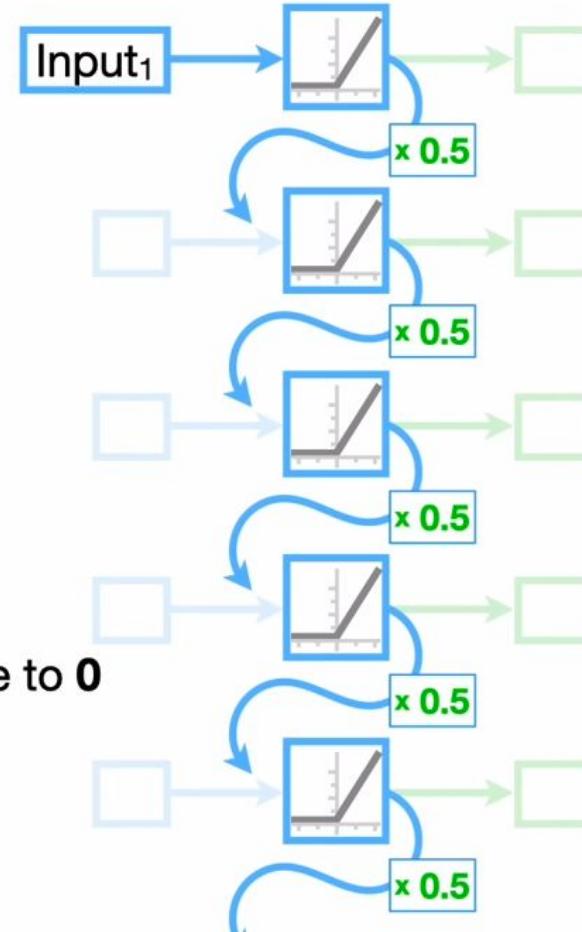




...to vanish.

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

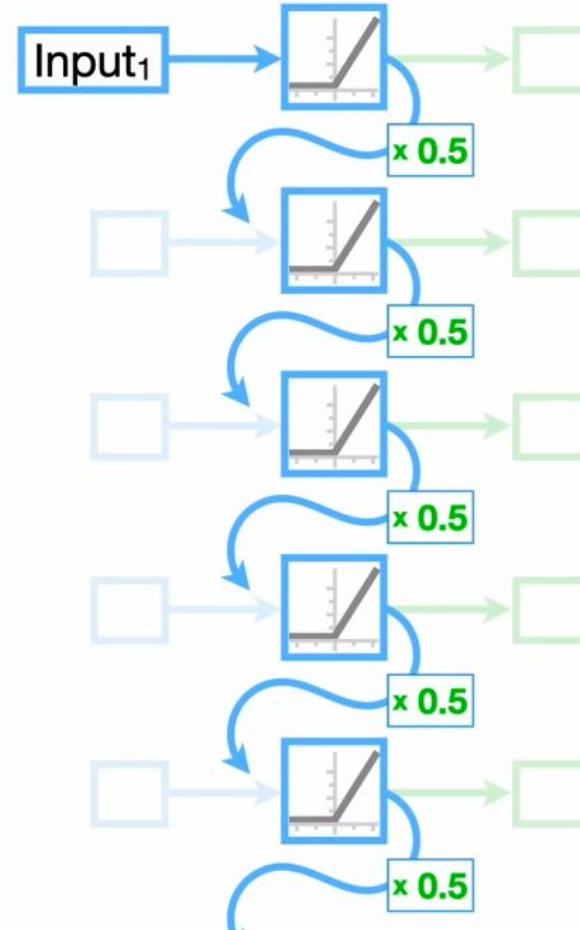
$$= \text{Input}_1 \times \text{a number super close to } 0$$





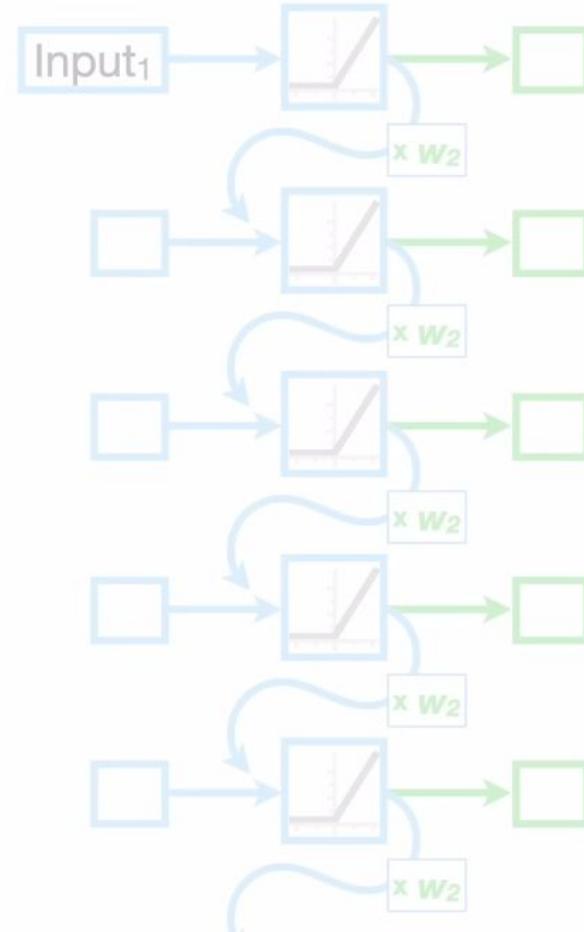
In summary, basic, vanilla **Recurrent Neural Networks** are hard to train because the gradients can **explode**, or **vanish**.

$$= \text{Input}_1 \times w_2^{\text{Num. Unroll}}$$



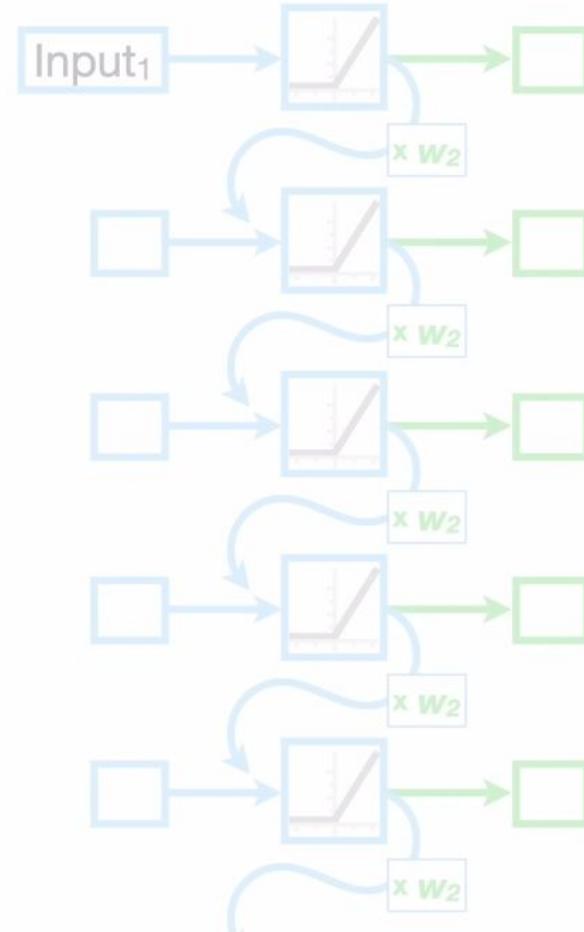


The good news is that it doesn't take much to extend the basic, vanilla **Recurrent Neural Network** so that we can avoid this problem.



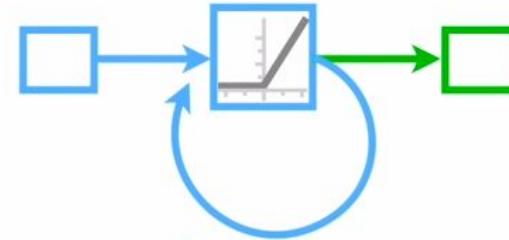


So, today, we're going to talk about **Long Short-Term Memory (LSTM)**, which is a type of **Recurrent Neural Network** that is designed to avoid the **exploding/vanishing gradient** problem.





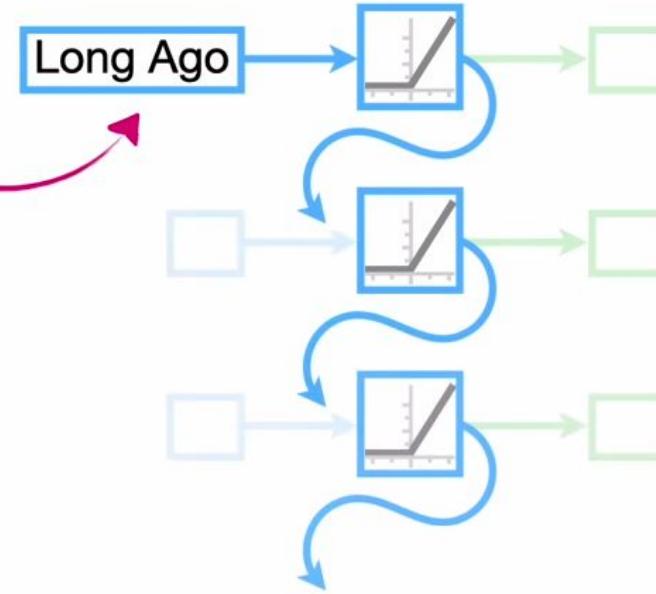
The main idea behind how
Long Short-Term Memory
(LSTM) works...



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



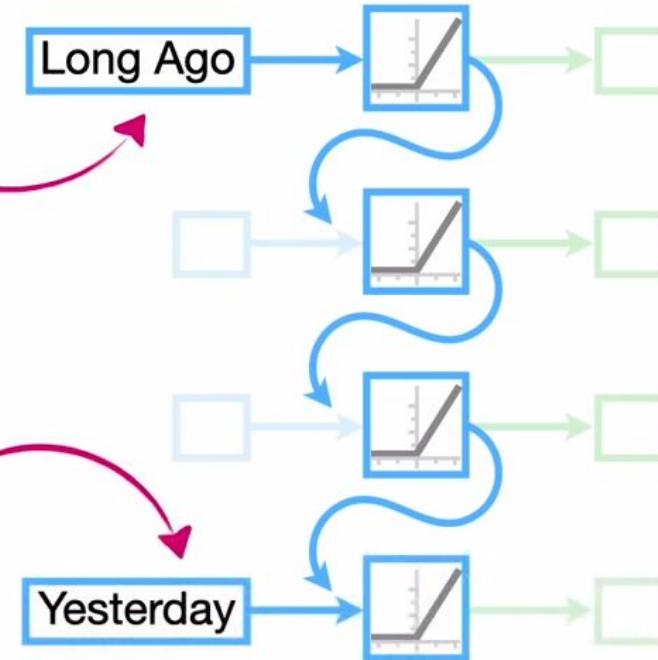
...for events that
happened long ago...





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

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

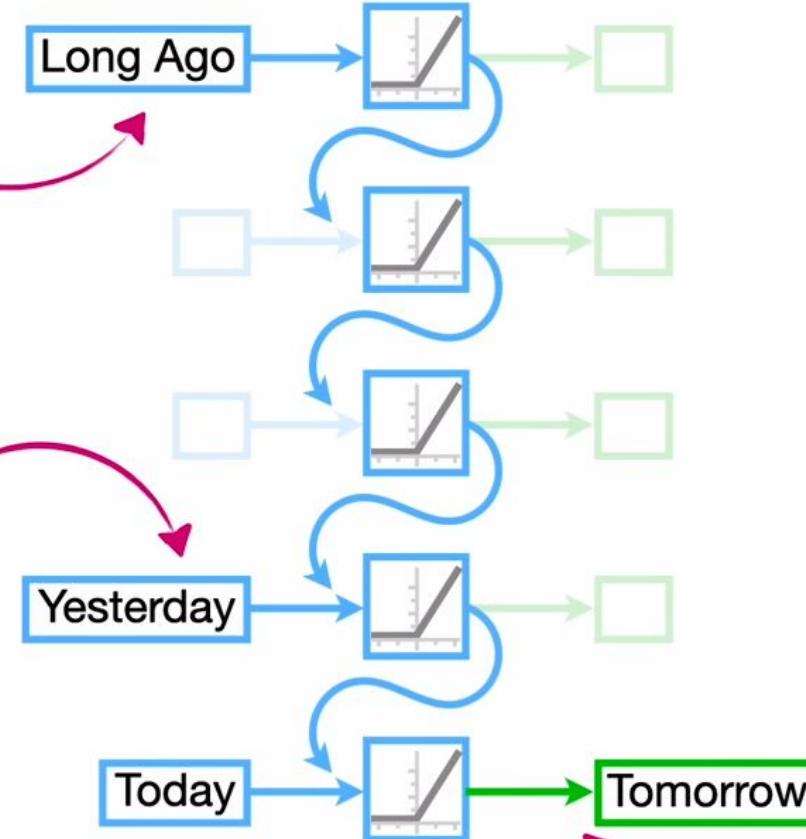




...for events that
happened long ago...

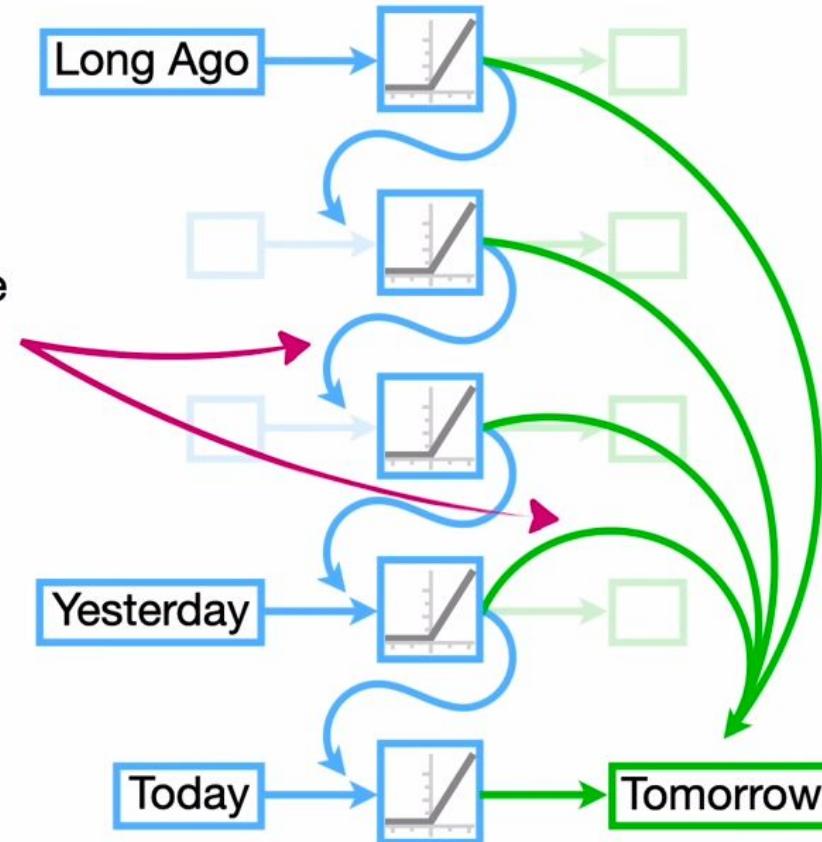
...and events that just
happened yesterday...

...to make a prediction
about tomorrow...



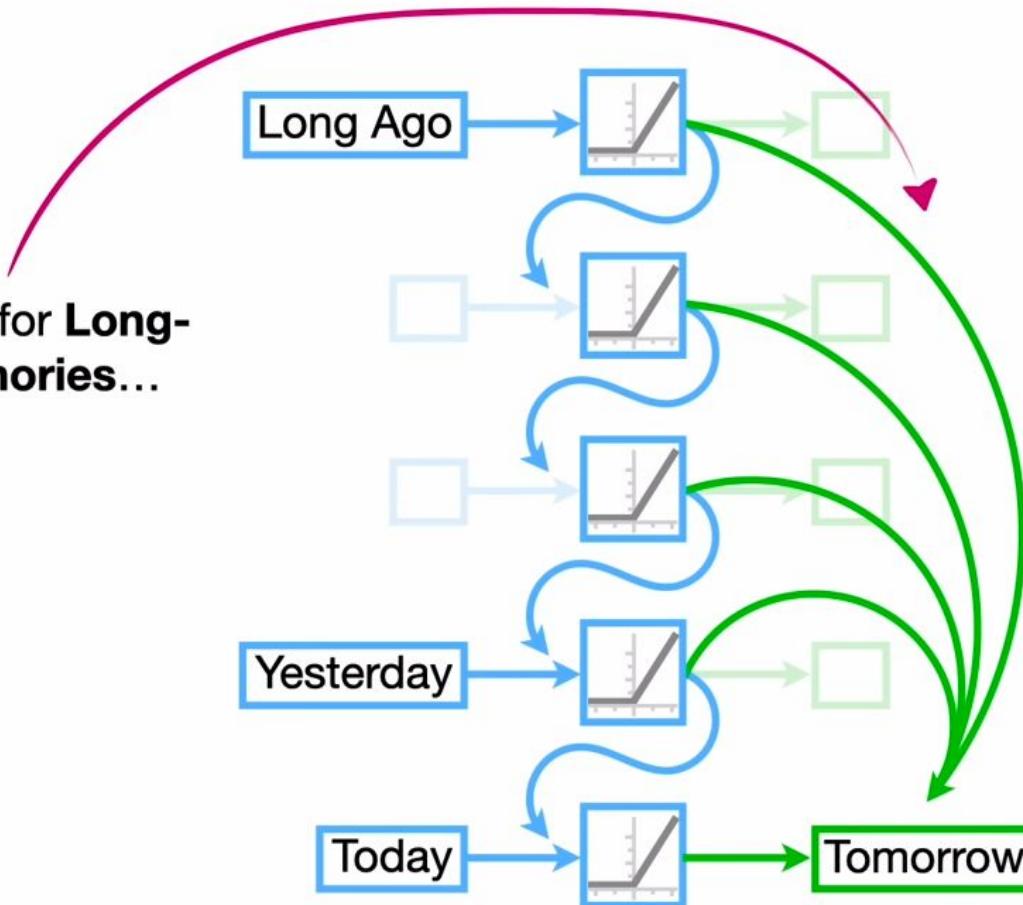


**...Long Short-Term
Memory uses two
separate paths to make
predictions about
tomorrow.**





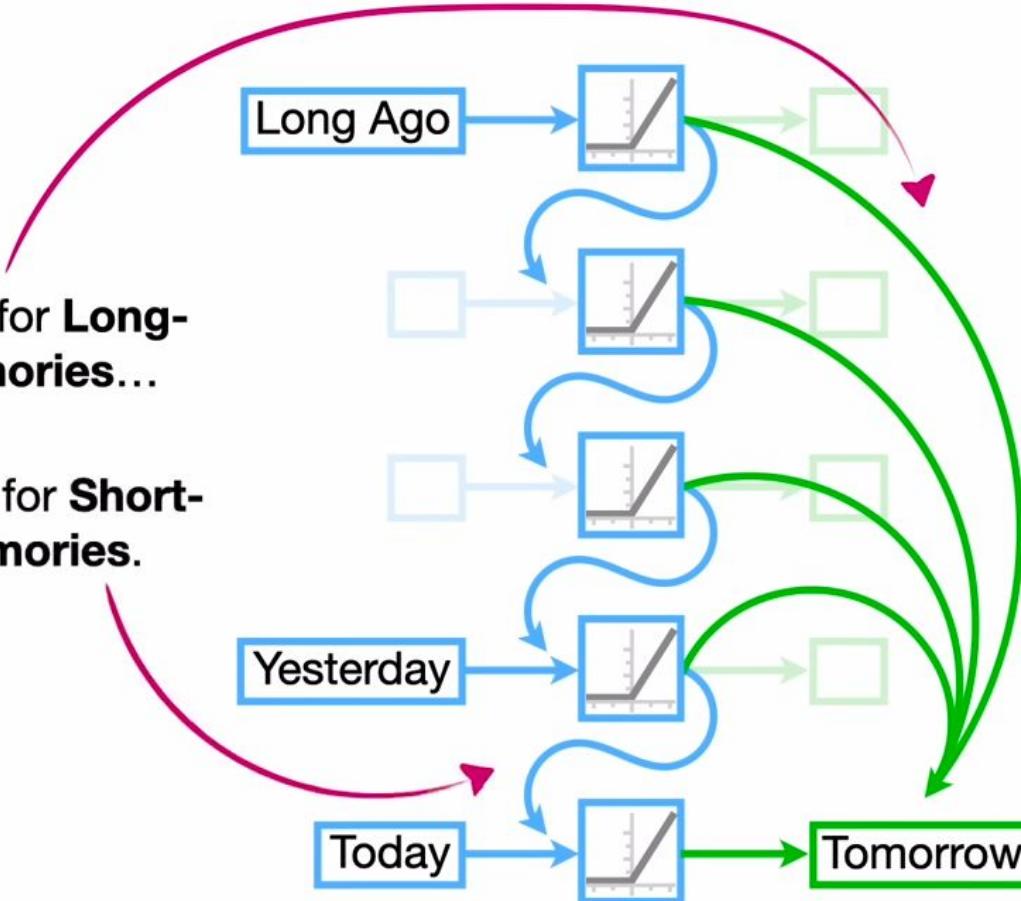
One path is for **Long-Term Memories...**





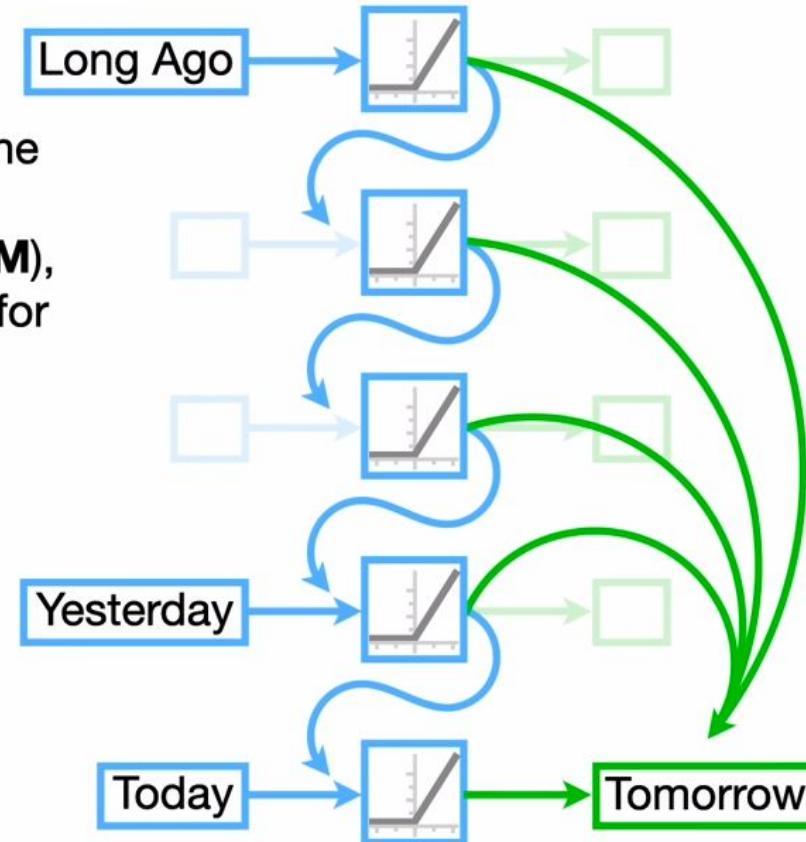
One path is for **Long-Term Memories**...

...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...**





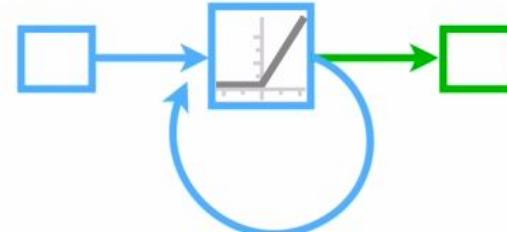
Now that we understand the **main idea** behind **Long Short-Term Memory (LSTM)**, that it uses different paths for **Long and Short-Term Memories...**

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



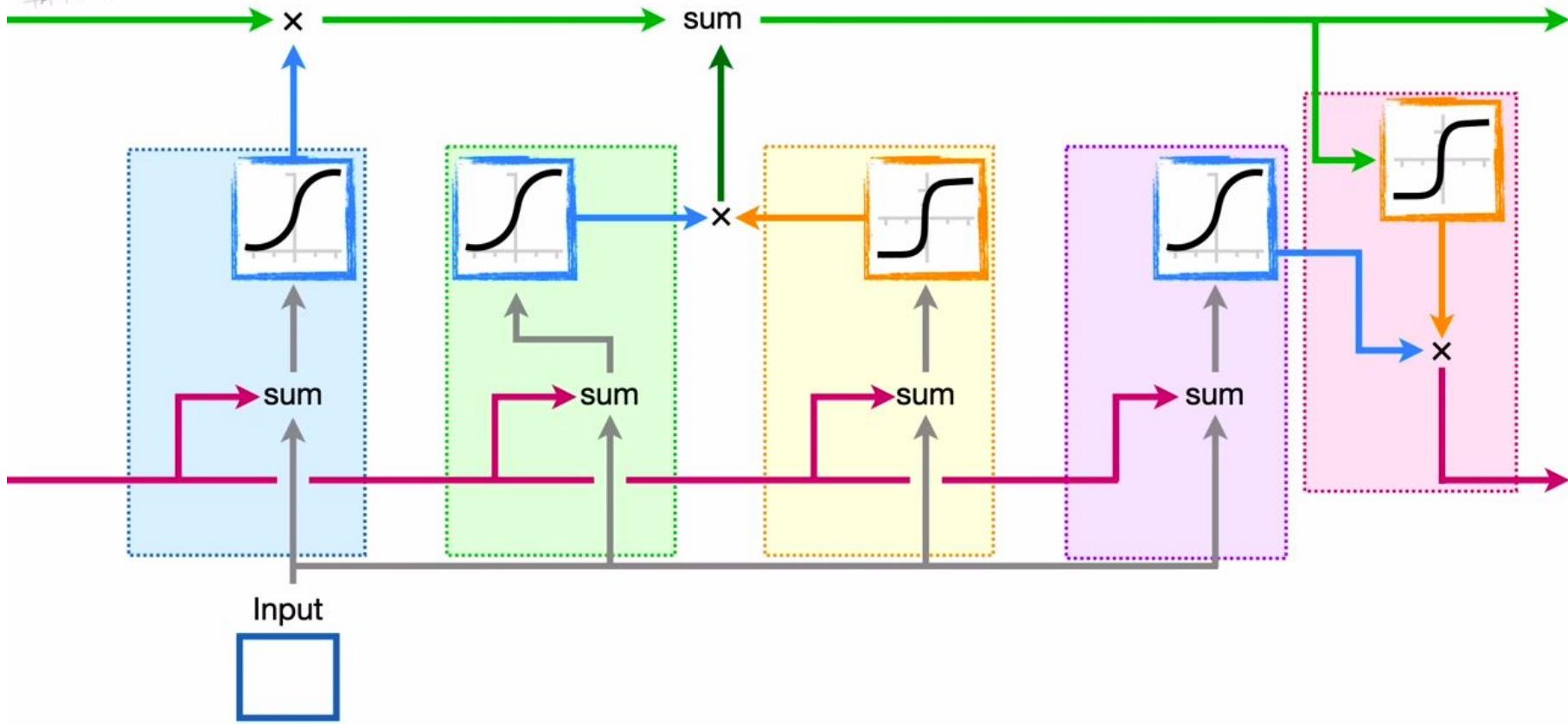


The bad news is that compared to a basic, vanilla **Recurrent Neural Network**, which unrolls from a relatively simple unit...



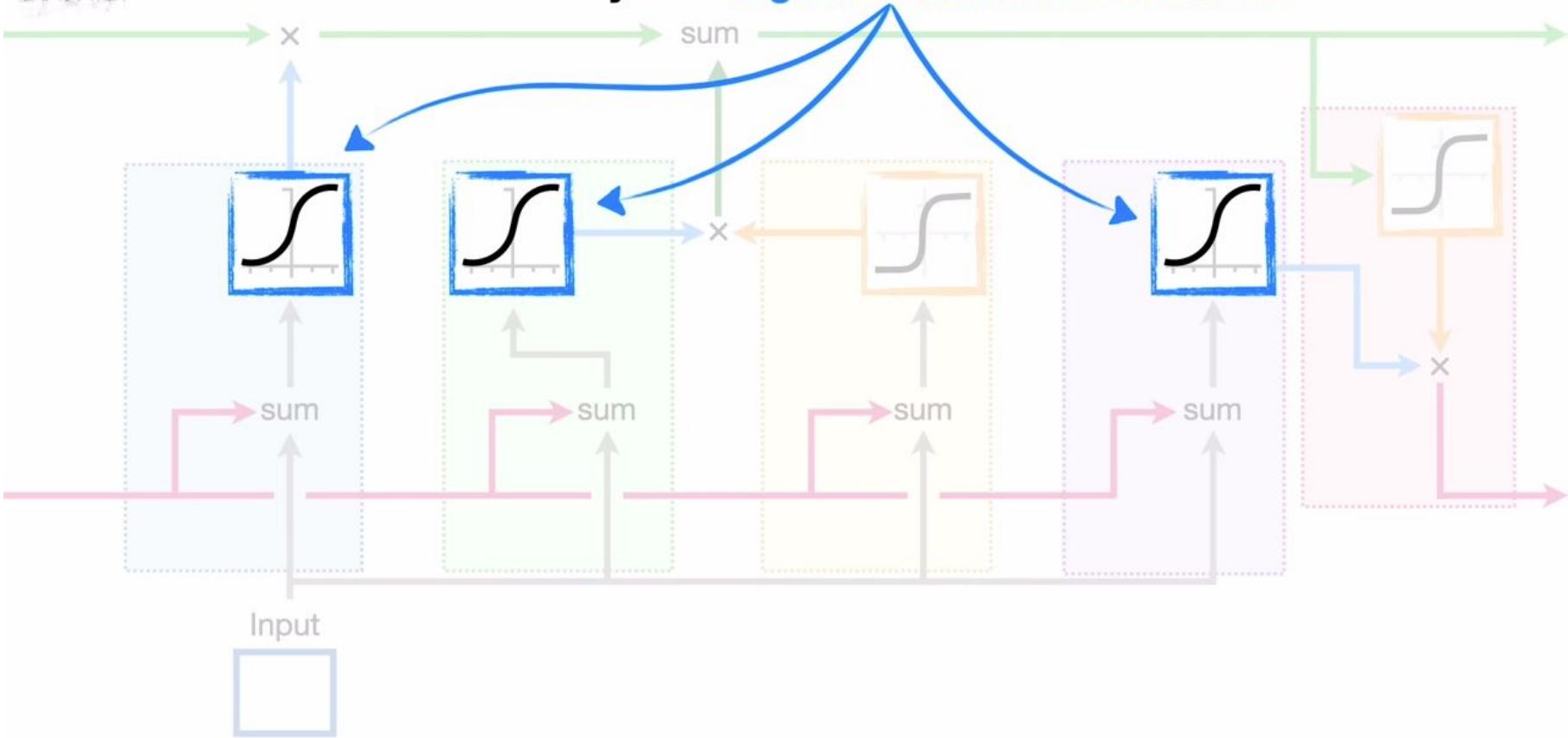


...Long Short-Term Memory is based
on a much more complicated unit.



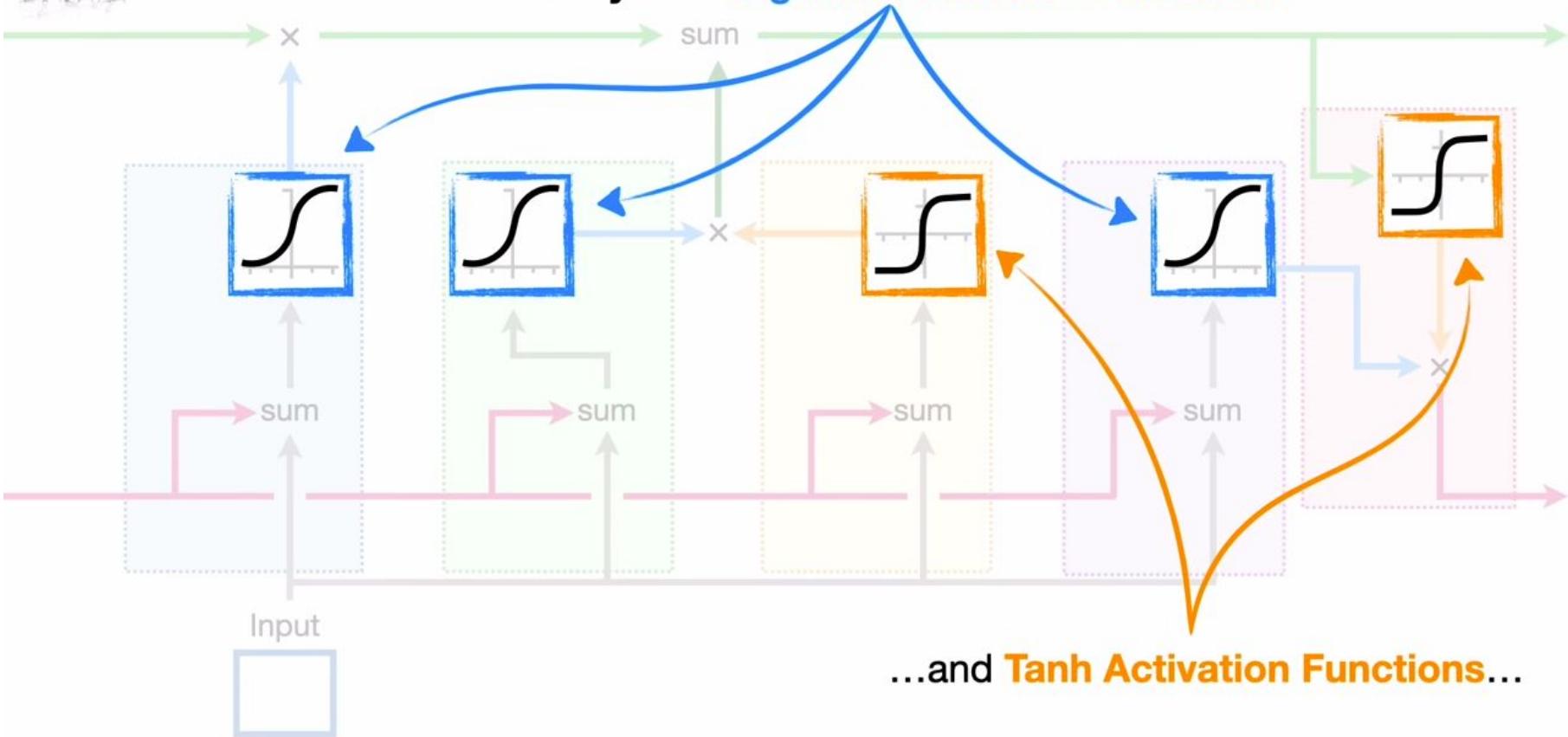


NOTE: Unlike the networks we've used before in this series, **Long Short-Term Memory** uses **Sigmoid Activation Functions**...





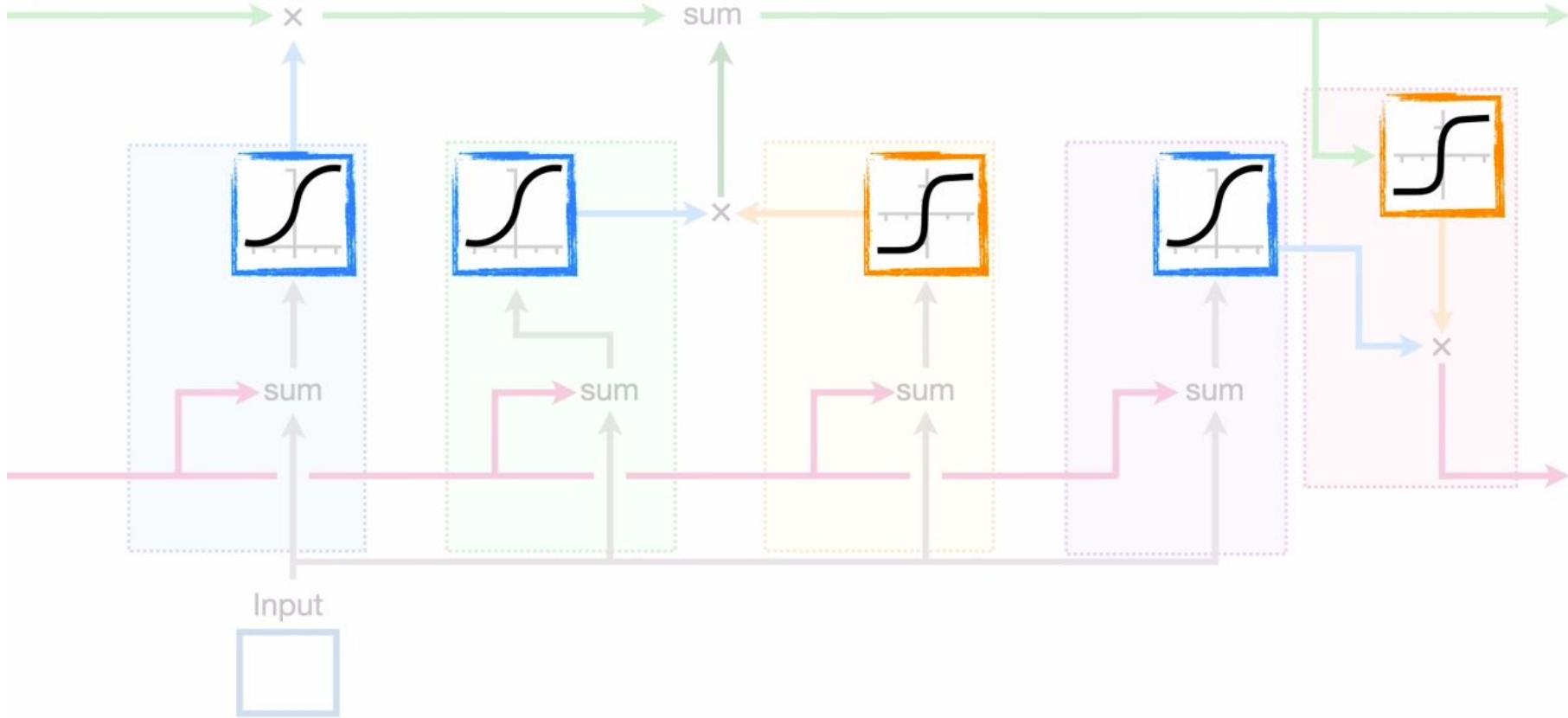
NOTE: Unlike the networks we've used before in this series, **Long Short-Term Memory** uses **Sigmoid Activation Functions...**



...and Tanh Activation Functions...

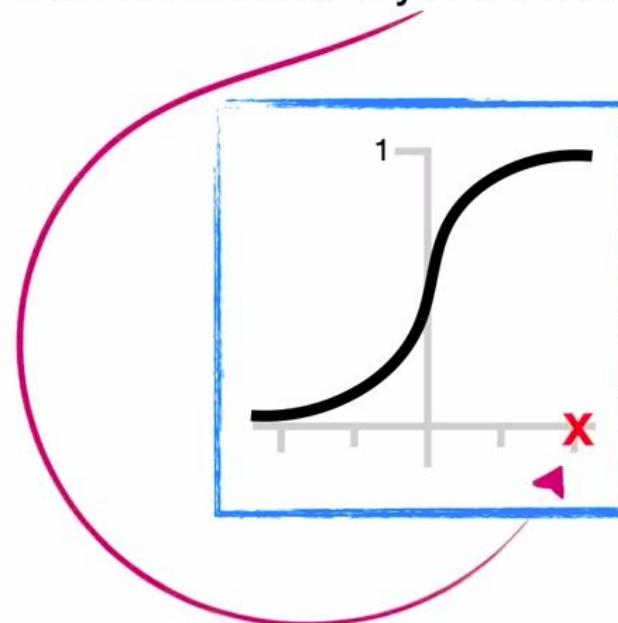


...so let's quickly talk about **Sigmoid** and **Tanh** Activation Functions.



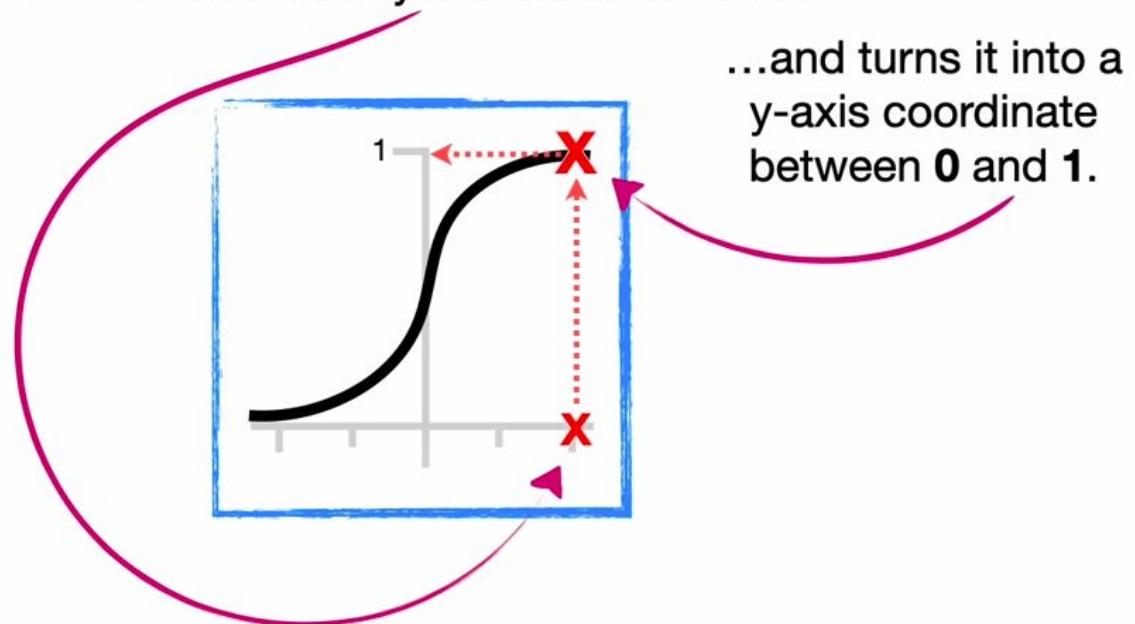


In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...



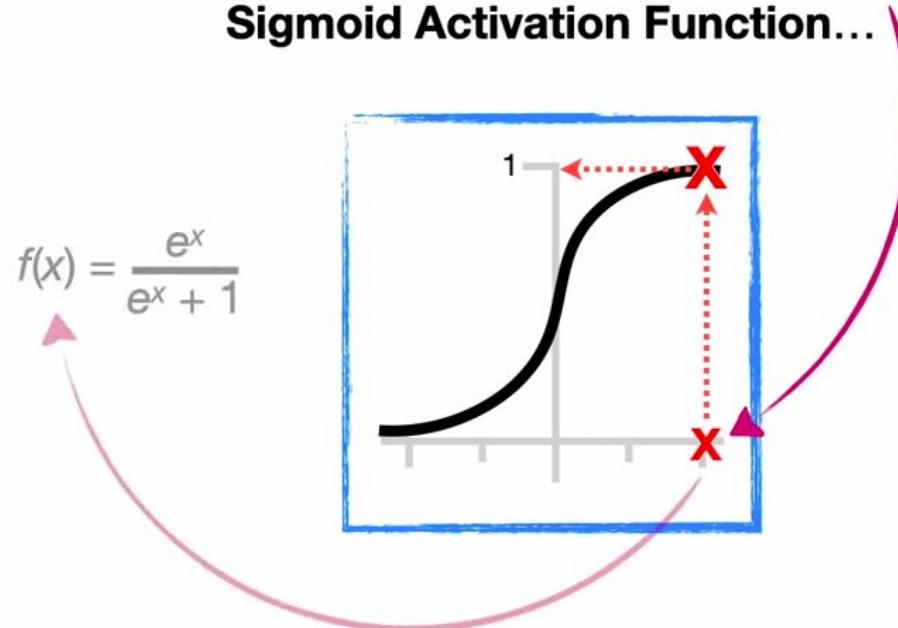


In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...





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

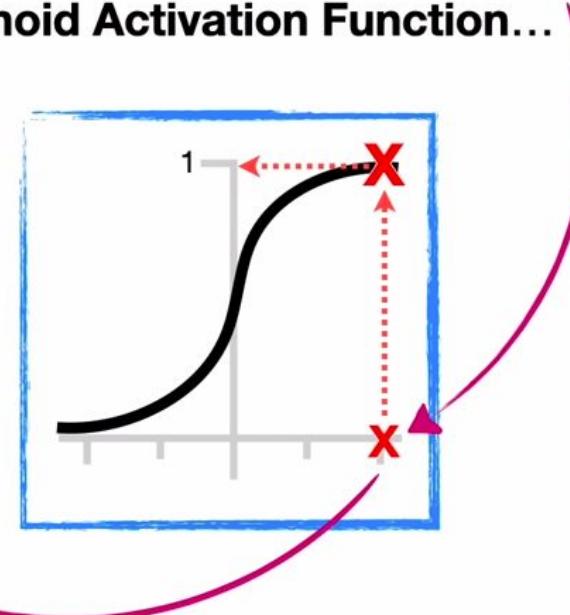




For example, when we plug in this x-axis coordinate, **10**, into the equation for the **Sigmoid Activation Function**...

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

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



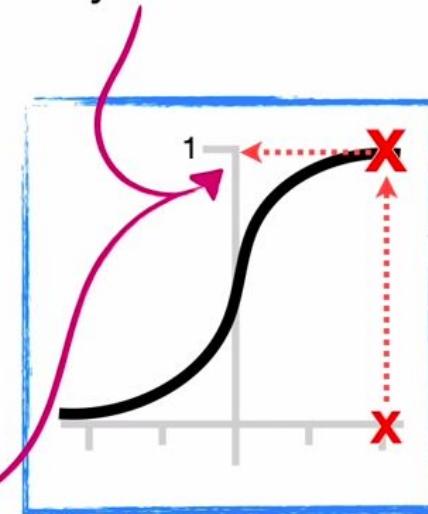


...we get **0.99995** as
the y-axis coordinate.

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

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

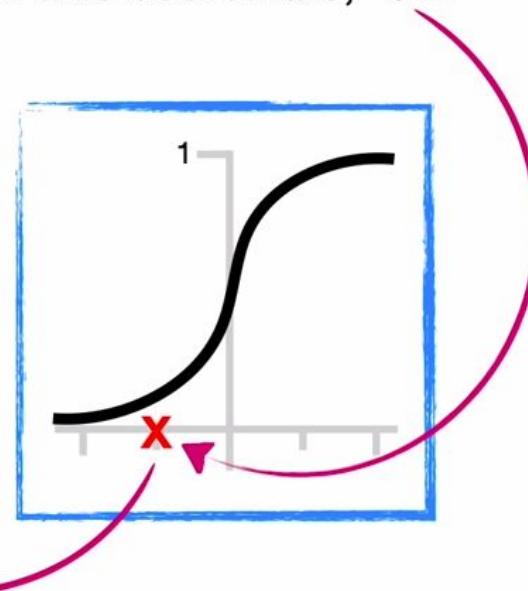
$$= 0.99995$$





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

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



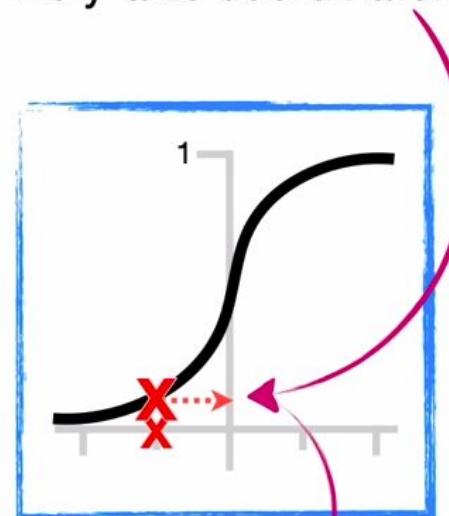


...then we get **0.01** as
the y-axis coordinate.

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

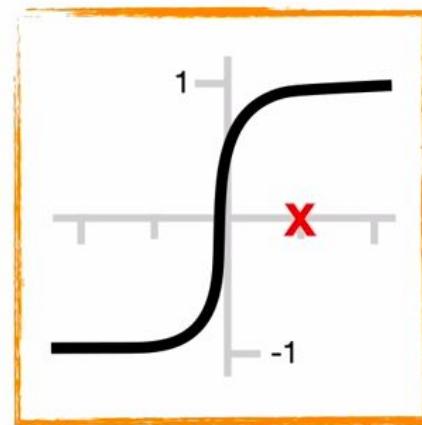
$$f(-5) = \frac{e^{-5}}{e^{-5} + 1}$$

$$= 0.01$$



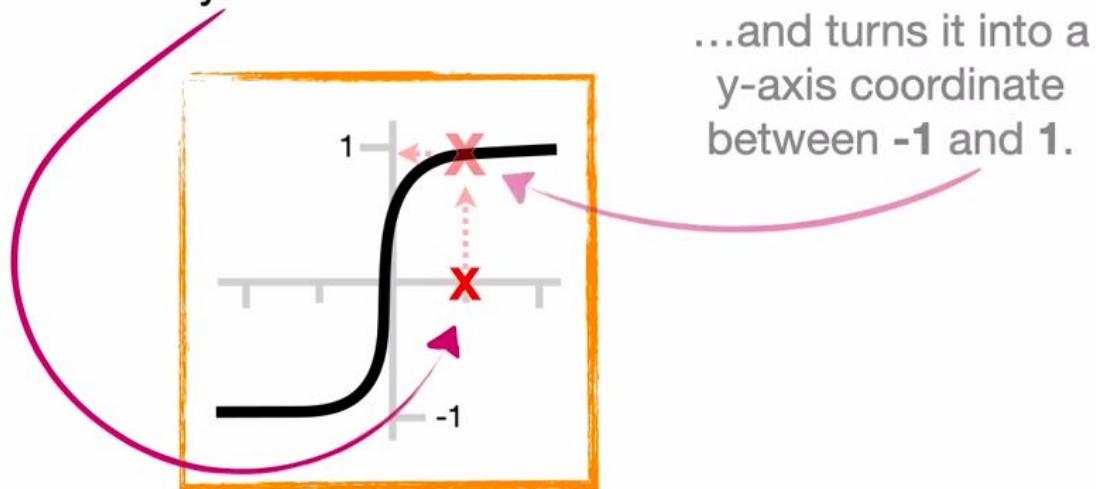


In contrast, the **tanh**, or **Hyperbolic Tangent, Activation Function** takes any x-axis coordinate...





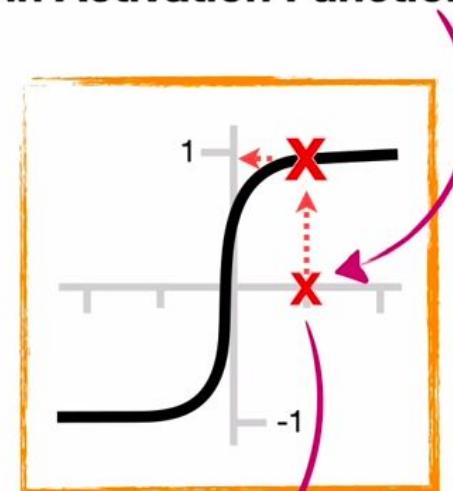
In contrast, the **tanh**, or **Hyperbolic Tangent, Activation Function** takes any x-axis coordinate...





For example, if we plug this x-axis coordinate, **2**, into the equation for the **Tanh Activation Function**...

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

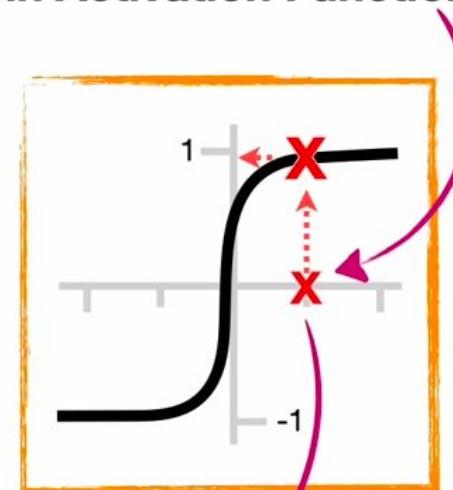




For example, if we plug this x-axis coordinate, **2**, into the equation for the **Tanh Activation Function**...

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

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



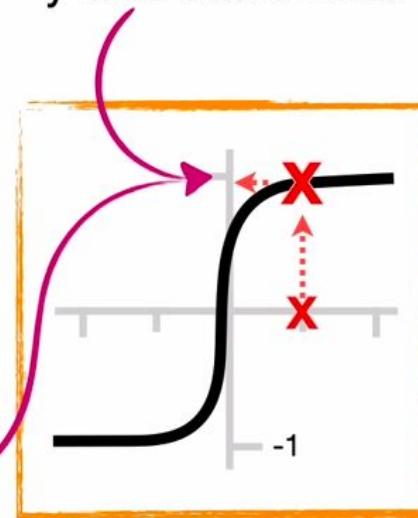


...we get **0.96** as the
y-axis coordinate.

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

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

$$= 0.96$$



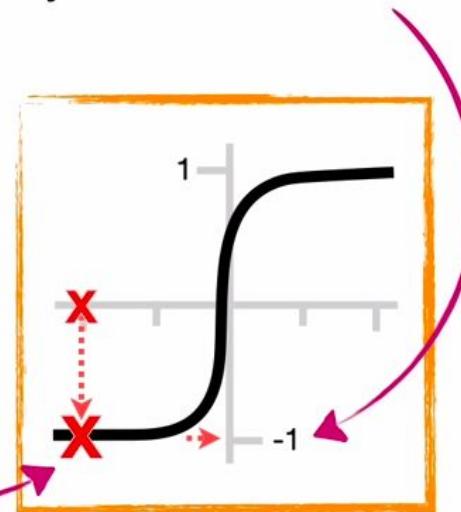


...we get **-1** as the
y-axis coordinate.

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

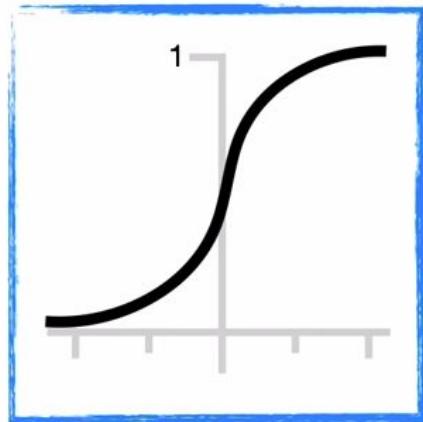
$$f(-5) = \frac{e^{-5} - e^5}{e^{-5} + e^5}$$

$$= -1$$



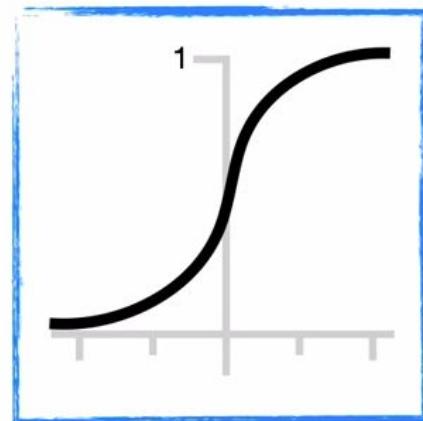


So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...

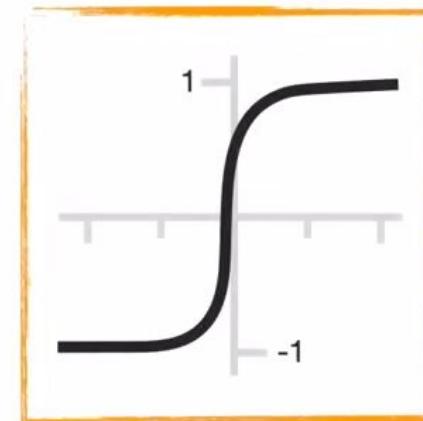




So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...

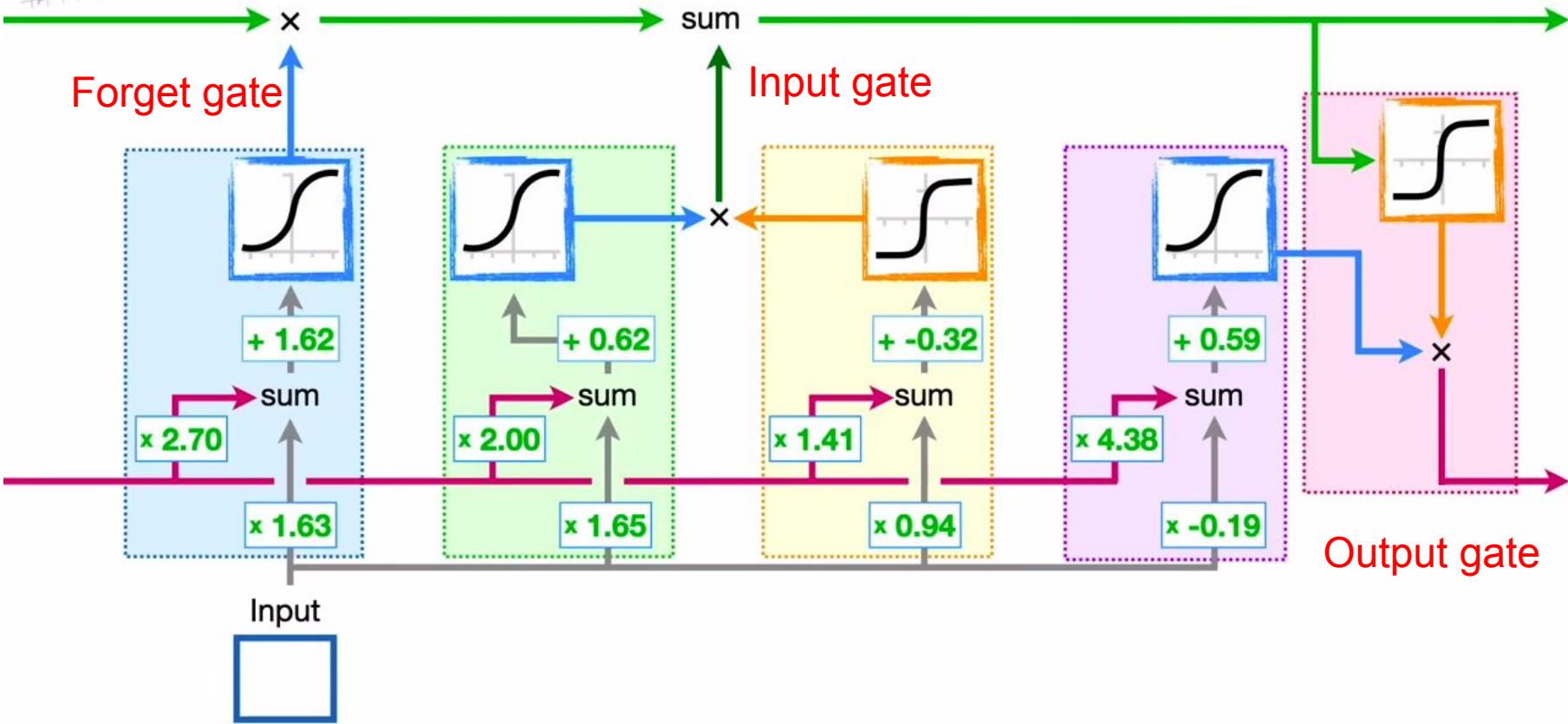


...and the **Tanh Activation Function** turns any input into a number between **-1** and **1**...



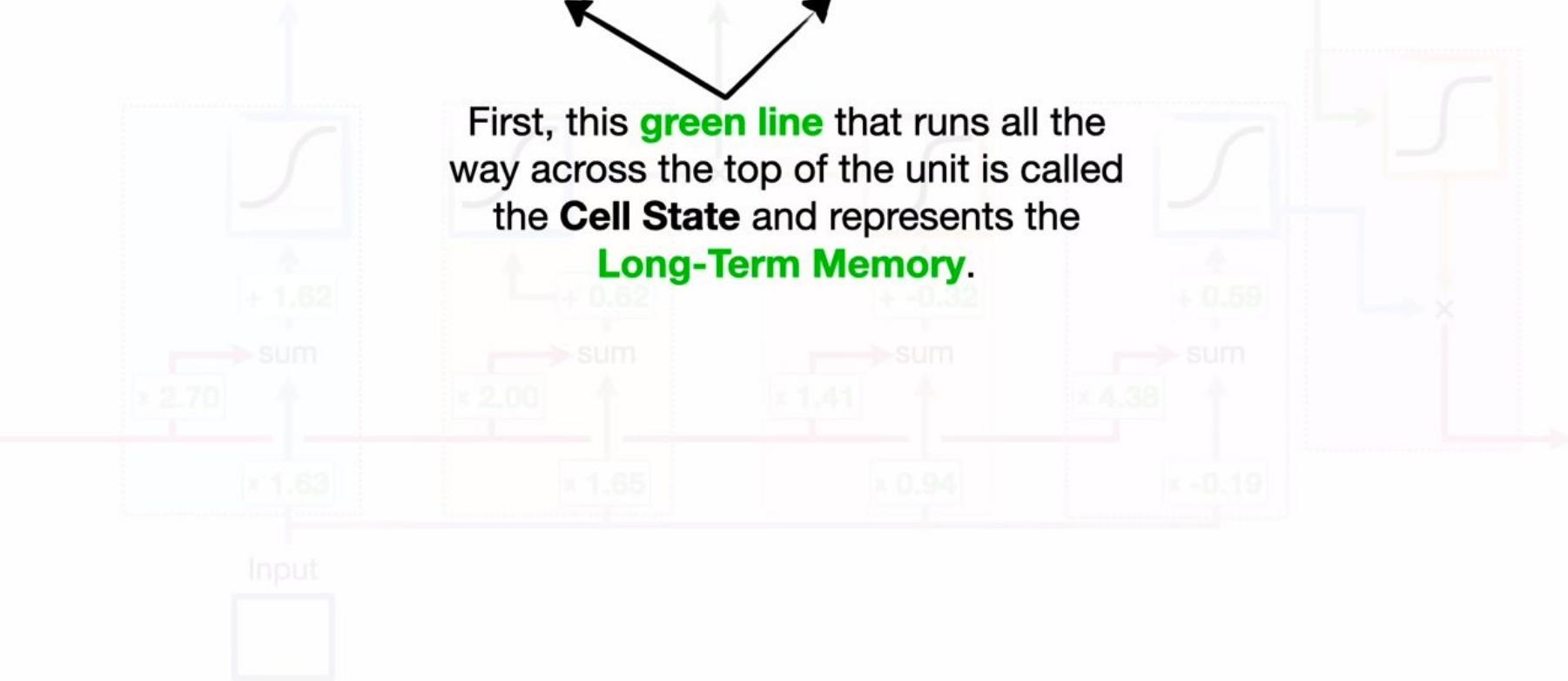


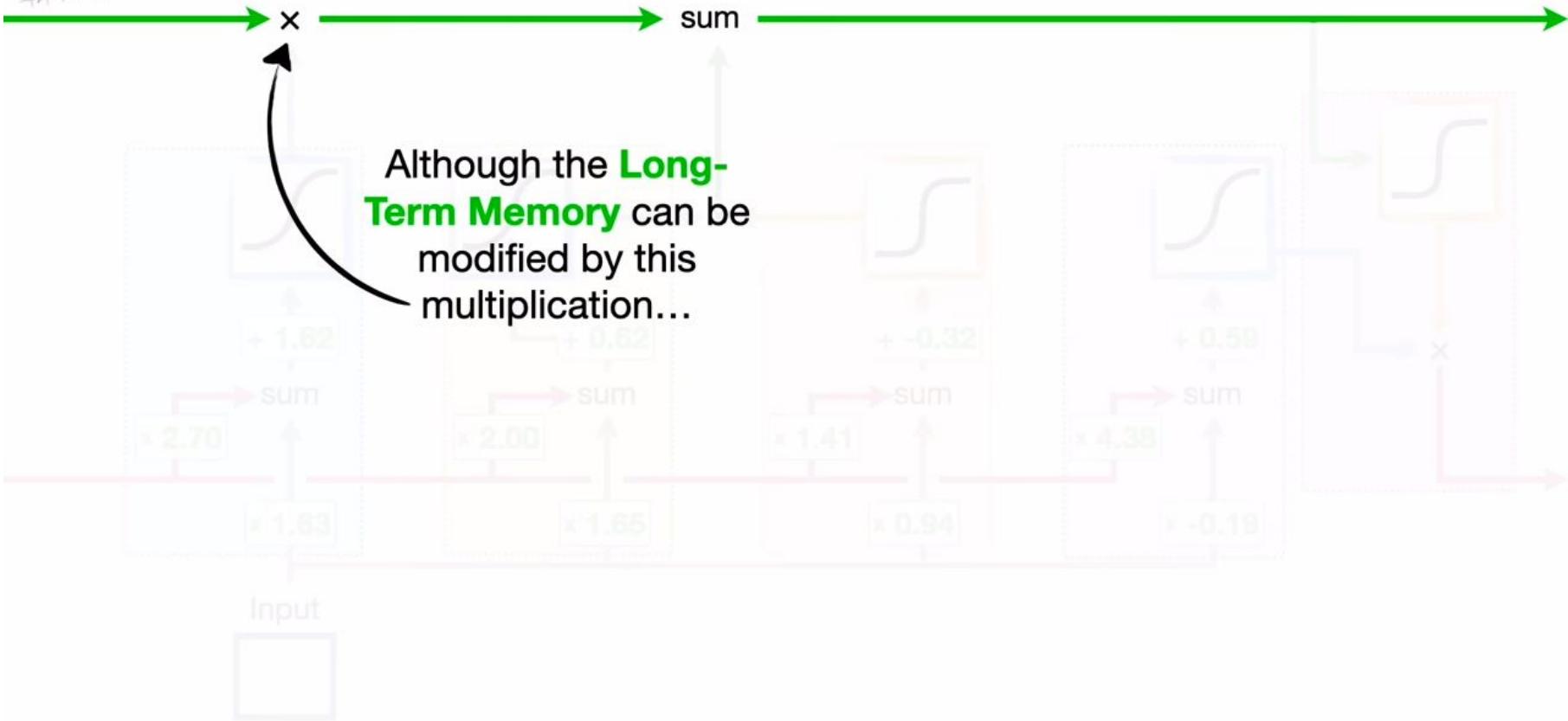
...let's talk about how the **Long Short-Term Memory** unit works.

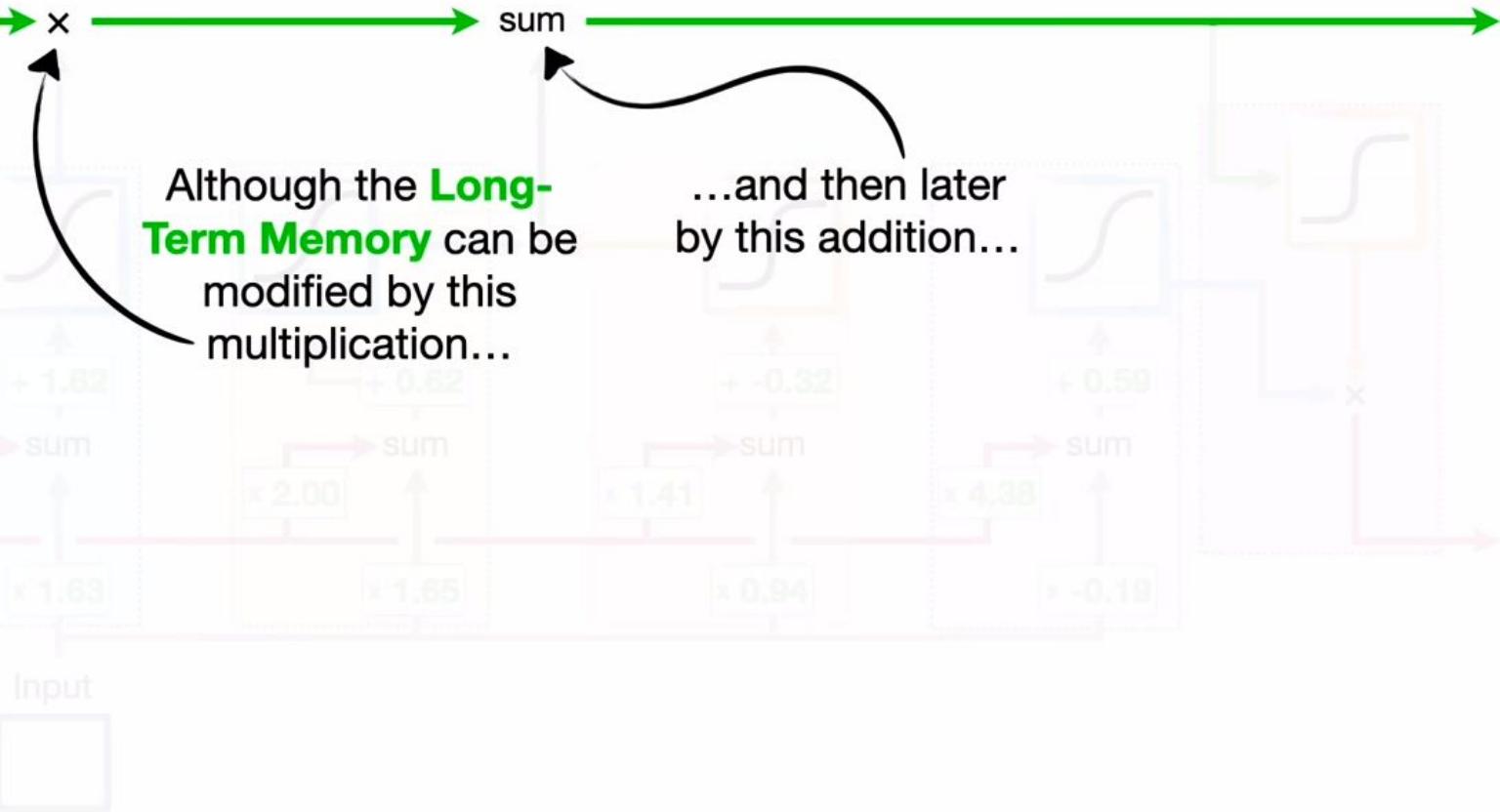




First, this **green line** that runs all the way across the top of the unit is called the **Cell State** and represents the **Long-Term Memory**.









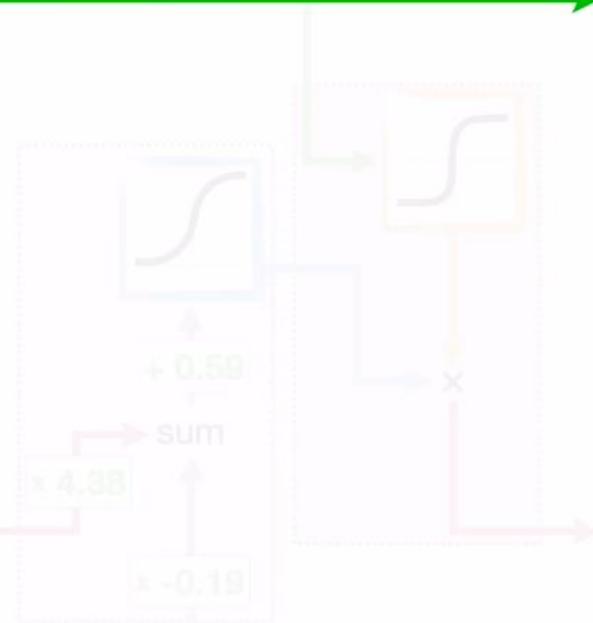
x

sum

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



Input

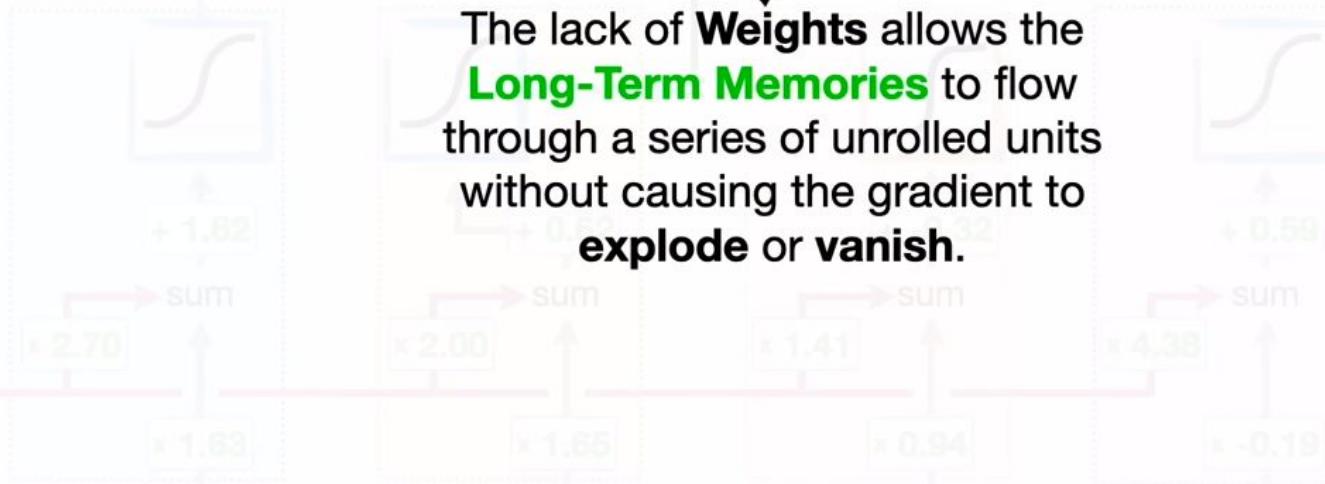


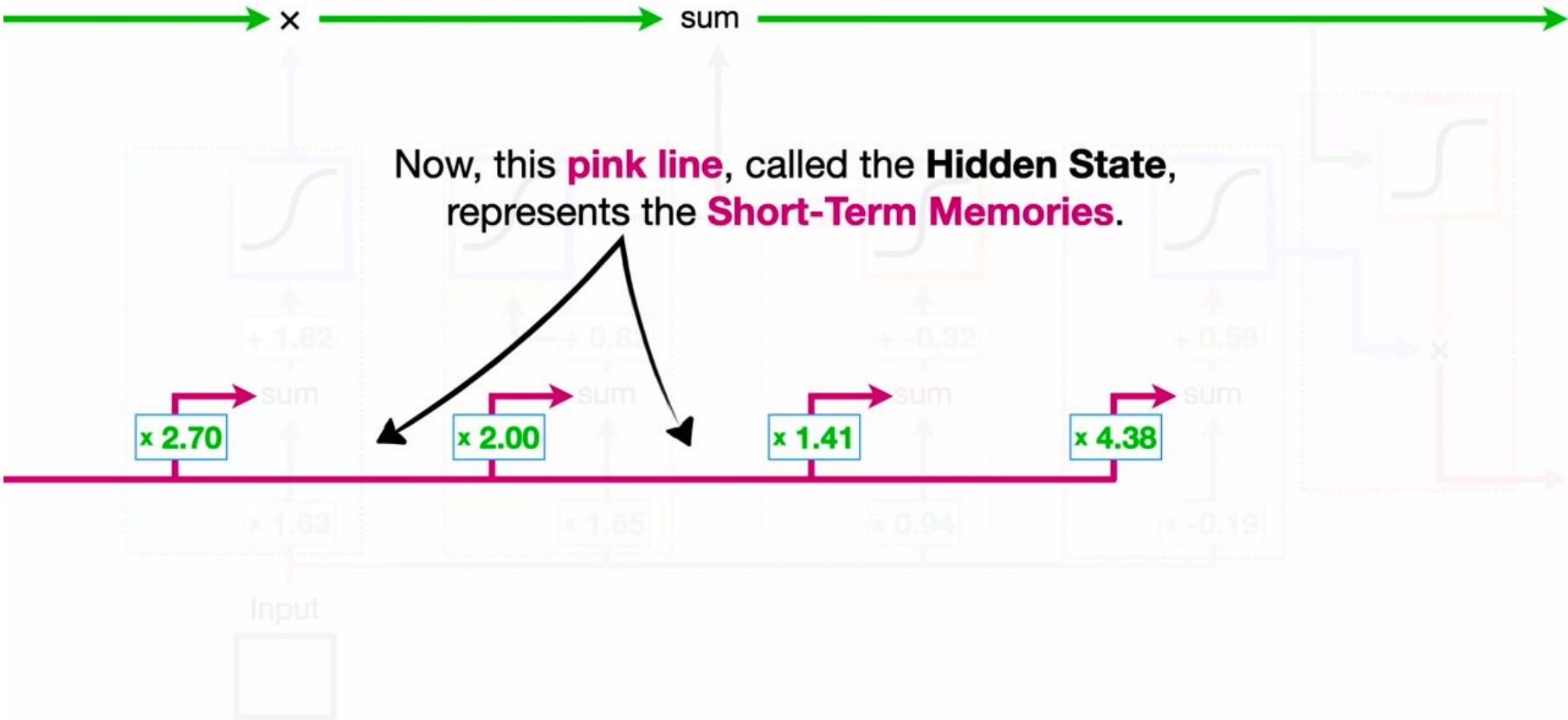


x

sum

The lack of **Weights** allows the
Long-Term Memories to flow
through a series of unrolled units
without causing the gradient to
explode or vanish.







→ x → sum →

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

x 2.70

x 2.00

x 1.41

x 4.38

Input



x 1.63

x 1.65

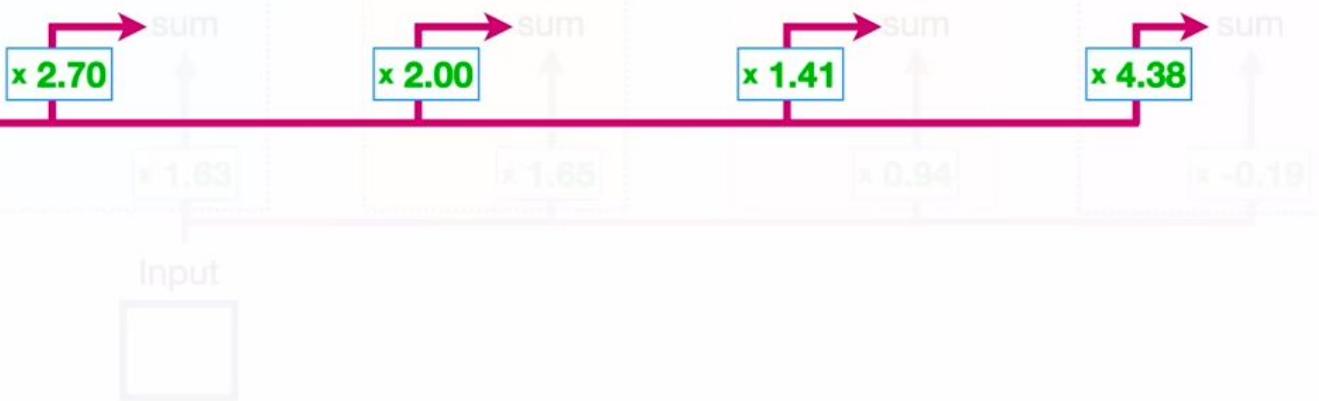
x 0.94

x -0.19

Input



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





Long-Term Memory

2

x

sum

First, for the sake of making the math interesting, let's just assume that the previous **Long-Term Memory** is 2...

x 2.70

x 2.00

x 1.41

x 4.38

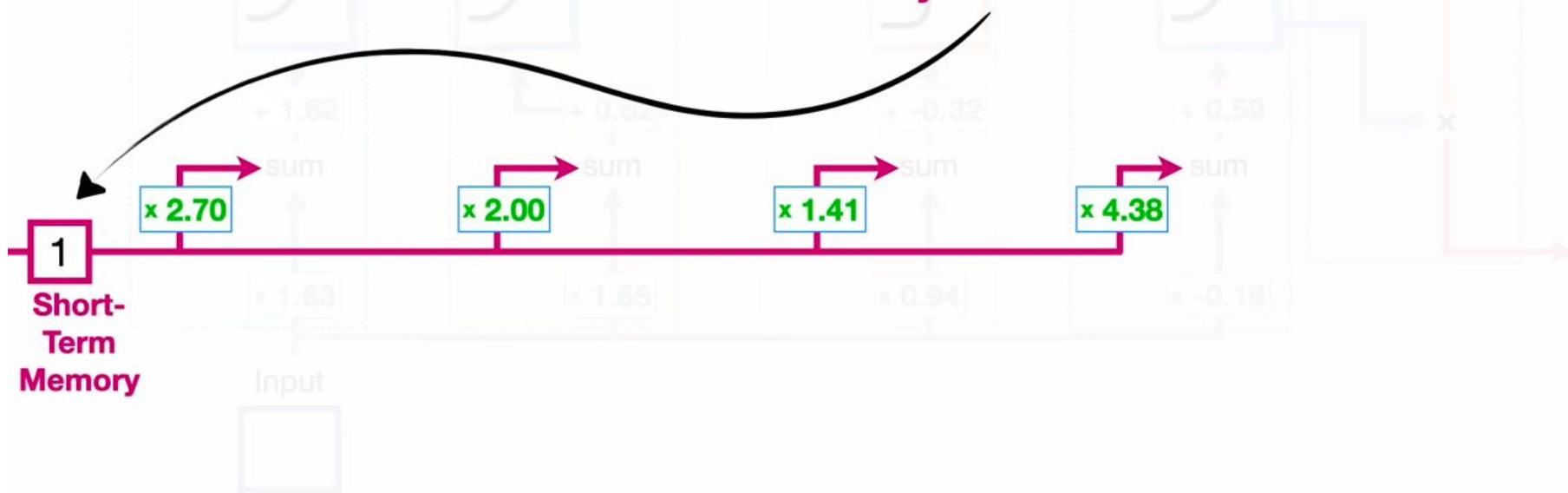
Input



Long-Term Memory



...and the previous
Short-Term Memory is 1.



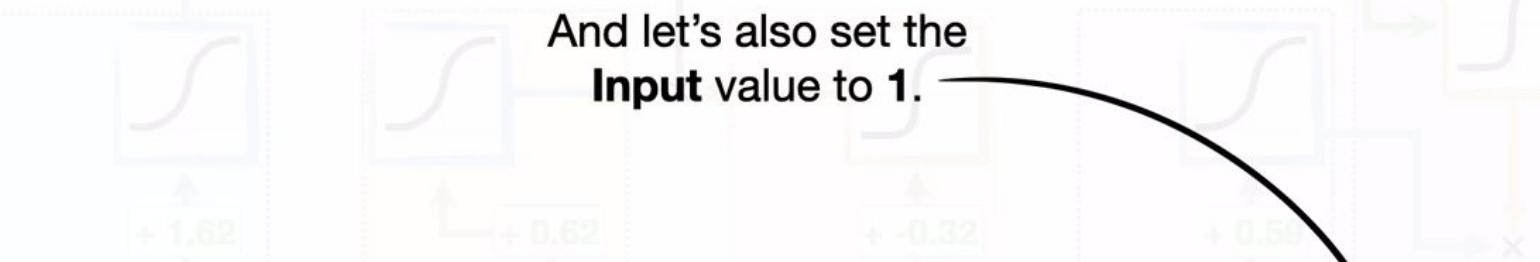


Long-Term Memory

2

x

sum



Short-Term Memory

1

Input

1

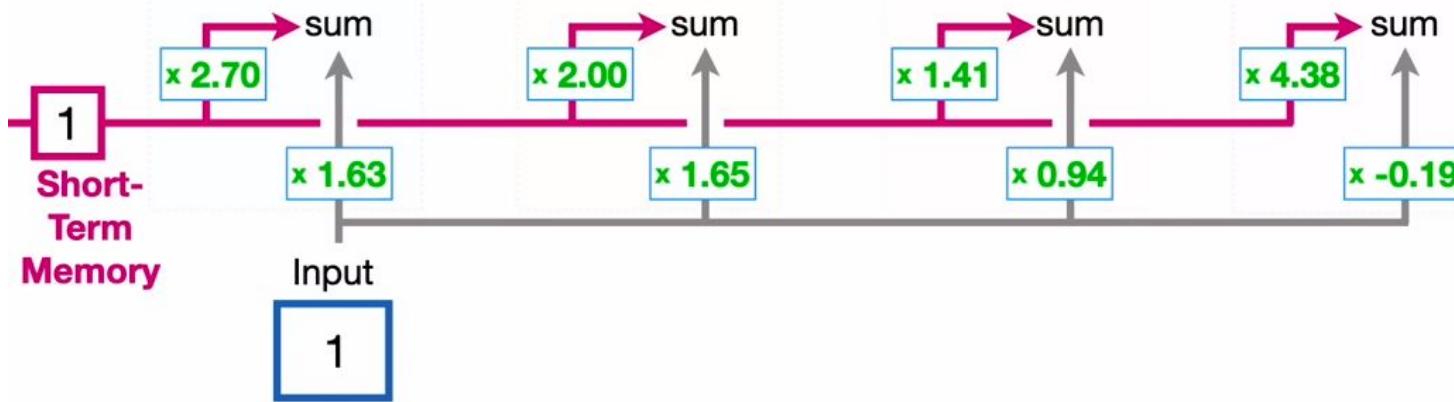


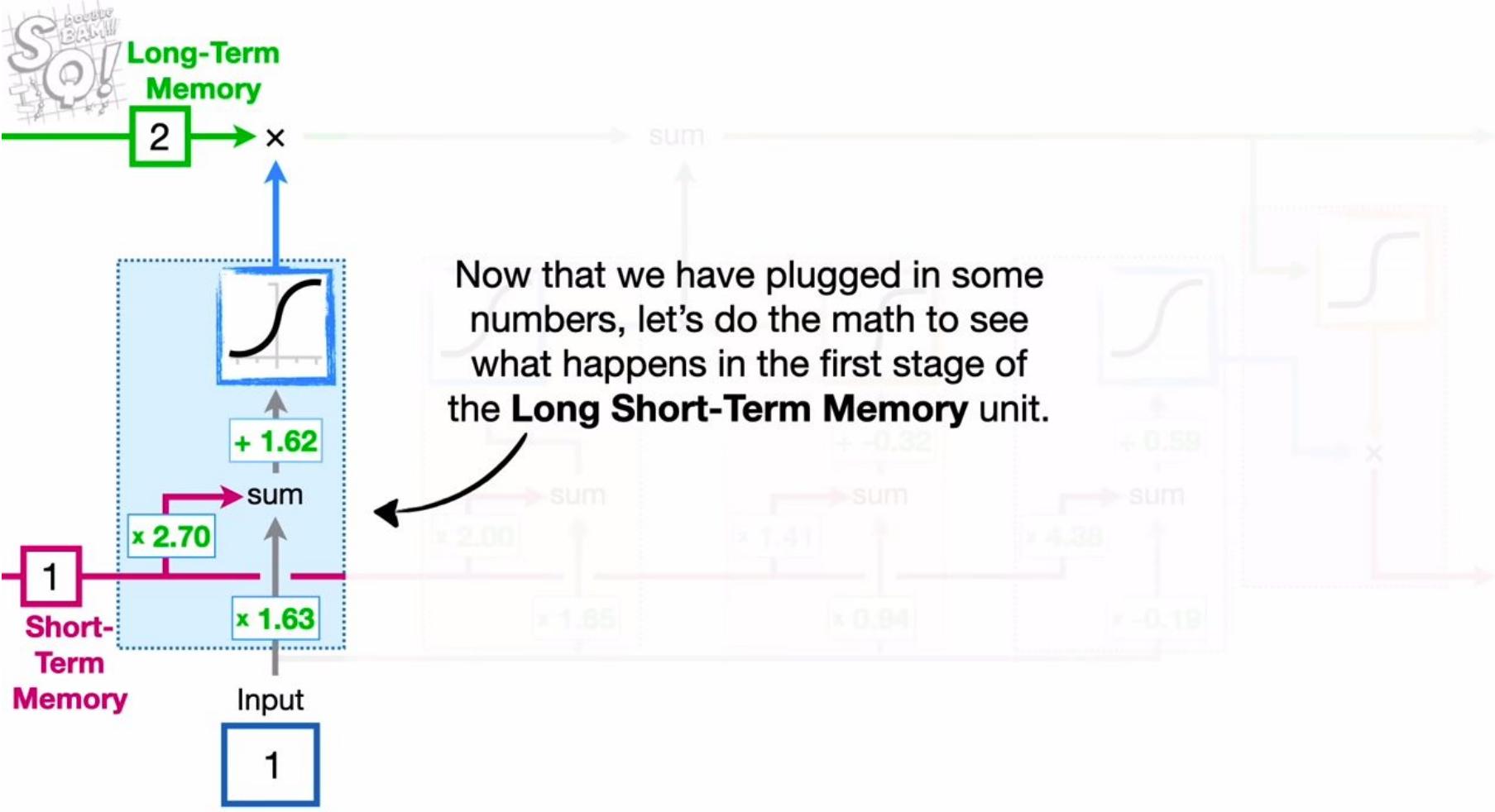
So!

Long-Term
Memory



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.

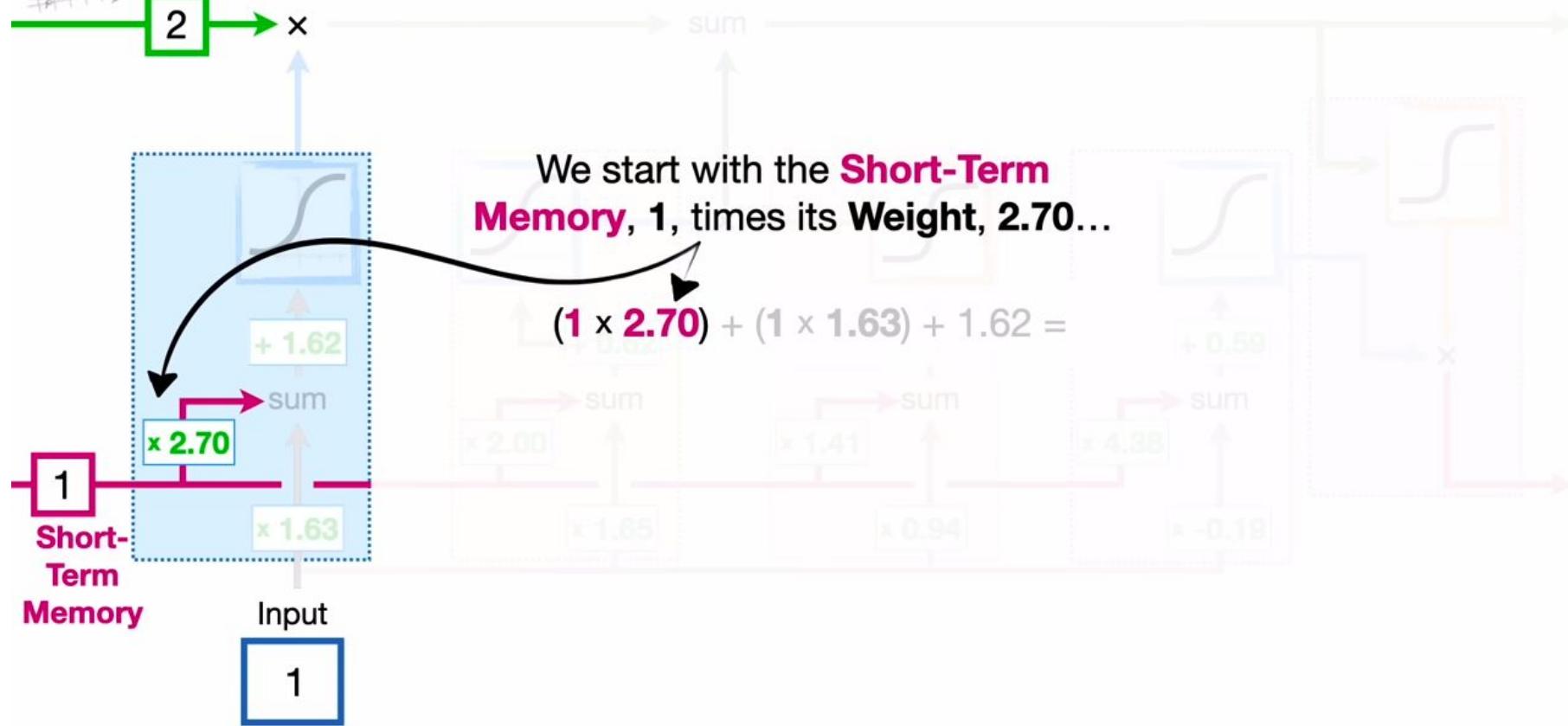






Long-Term
Memory

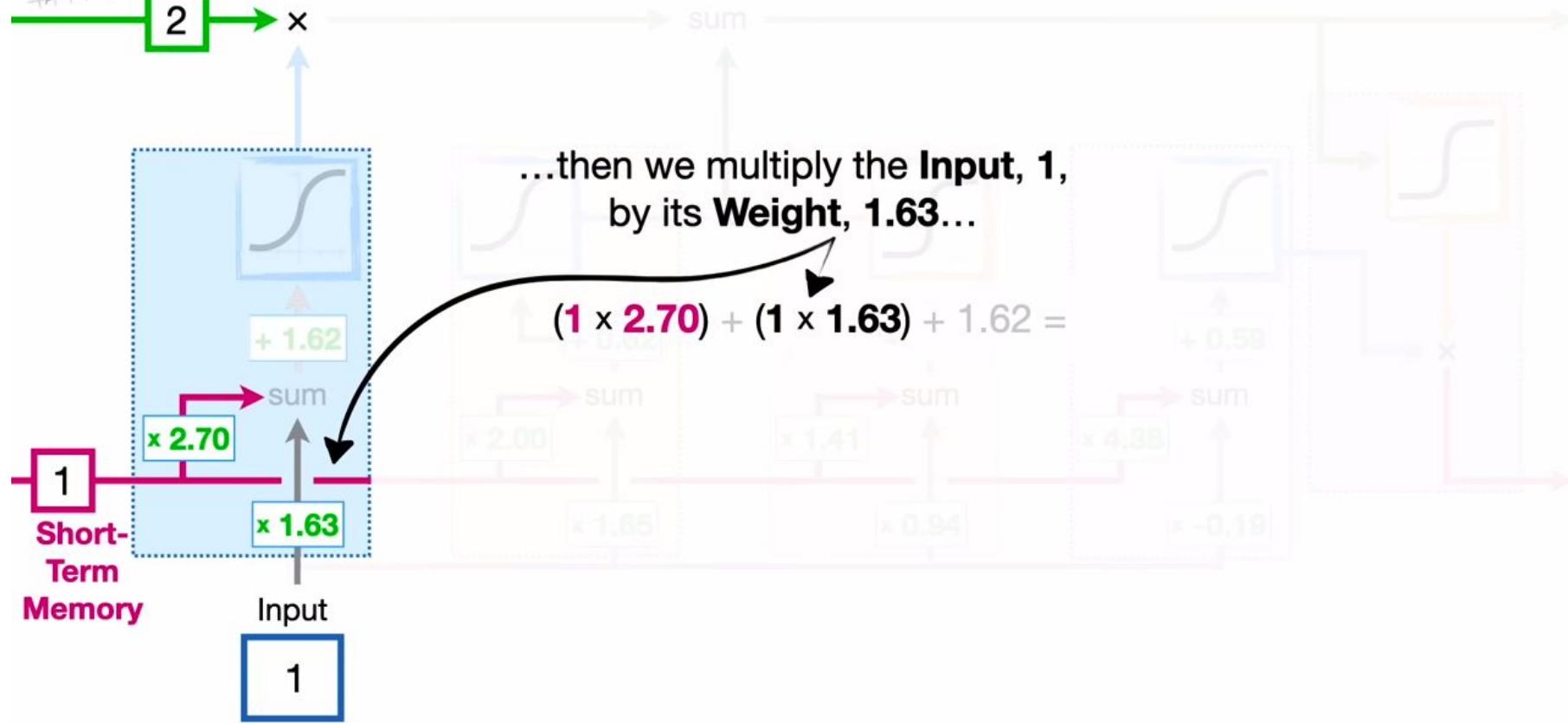
2

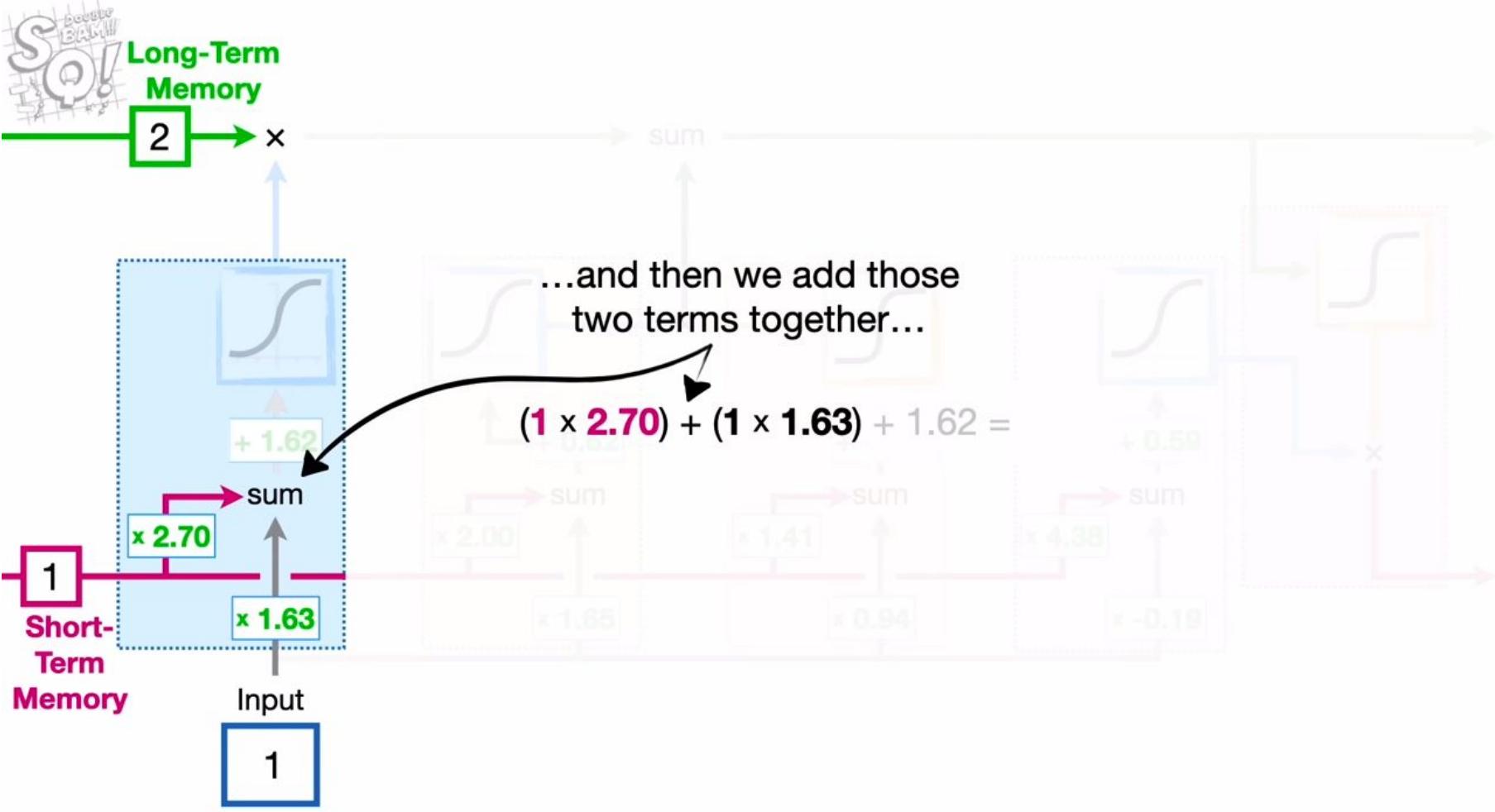


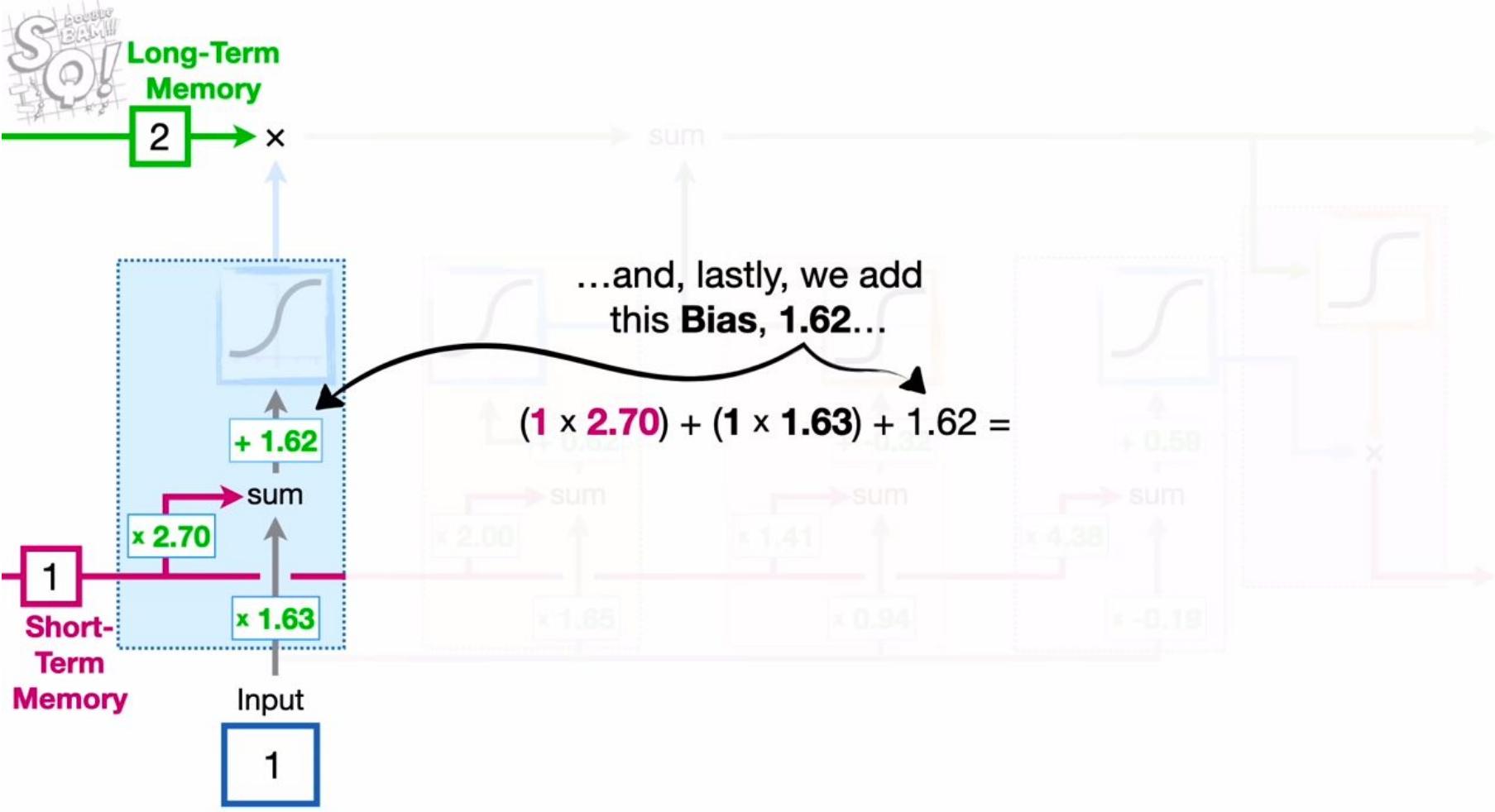


Long-Term Memory

2



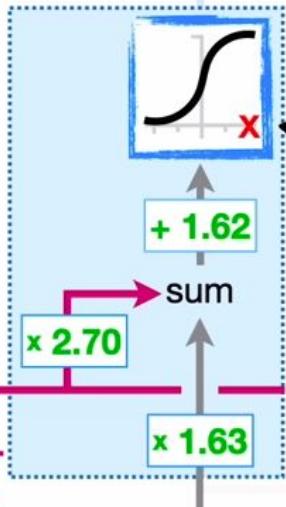






Long-Term
Memory

2 → x



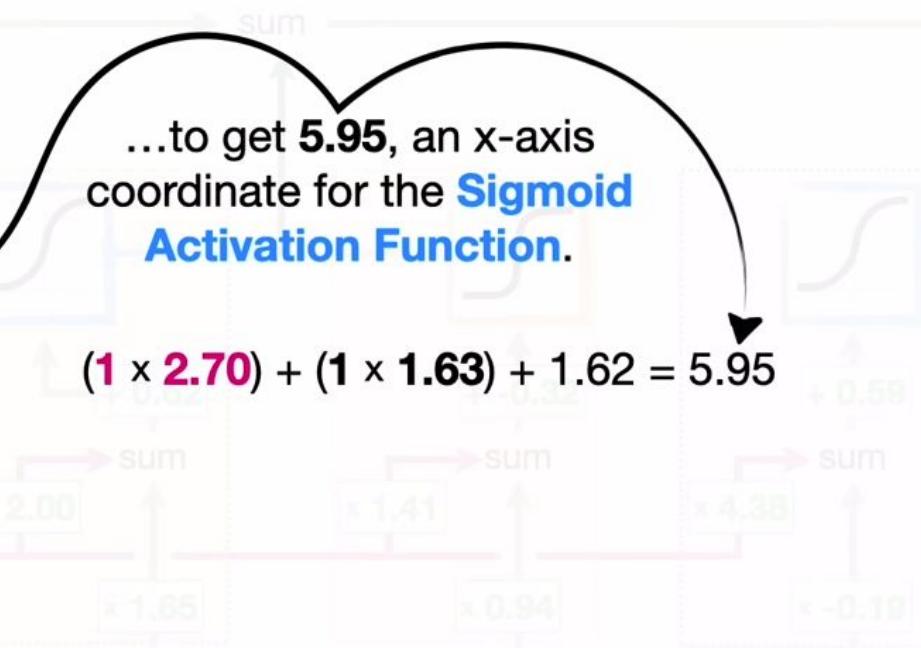
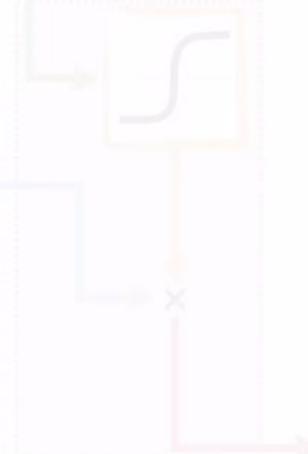
1
Short-
Term
Memory

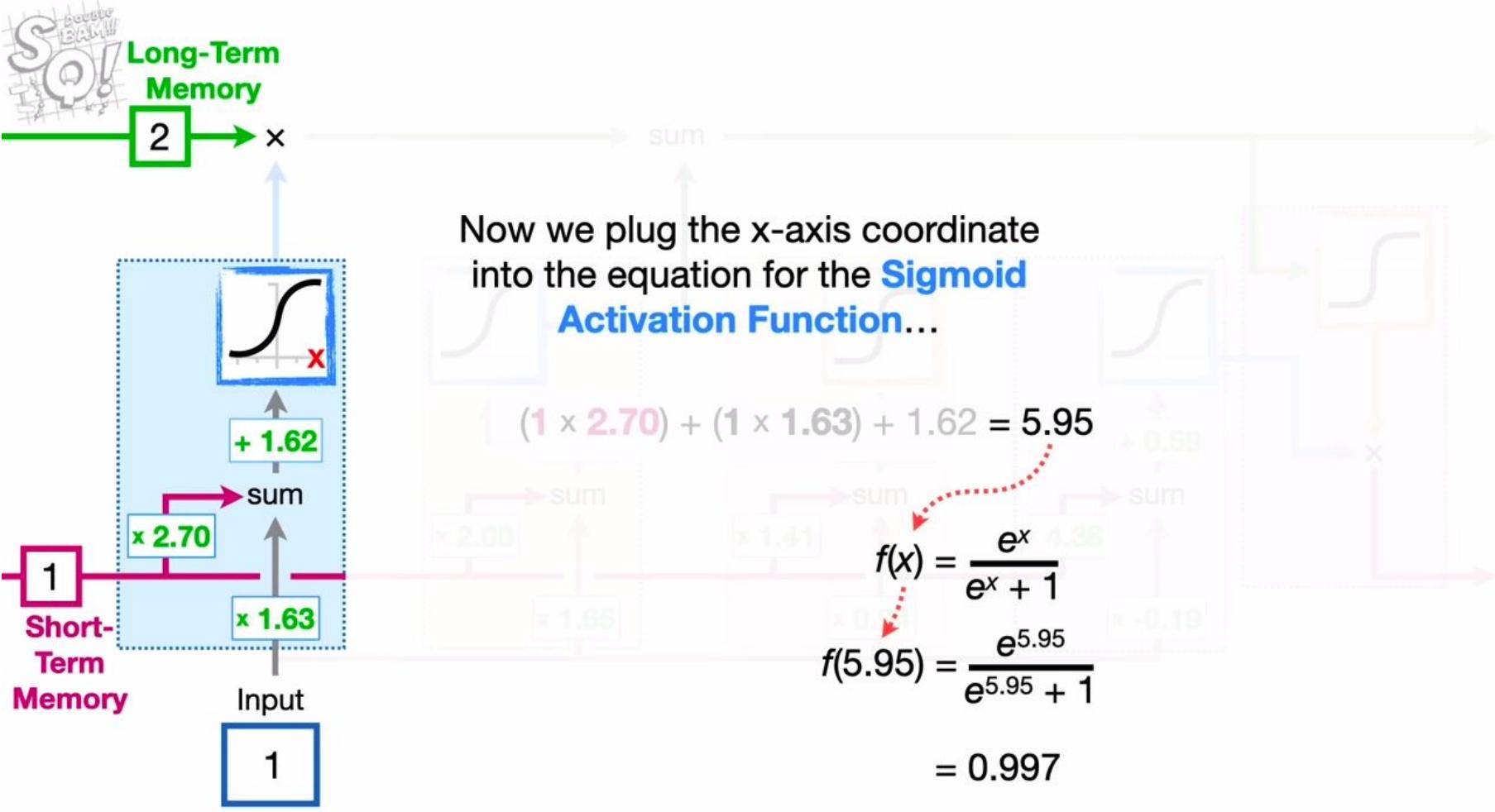
Input

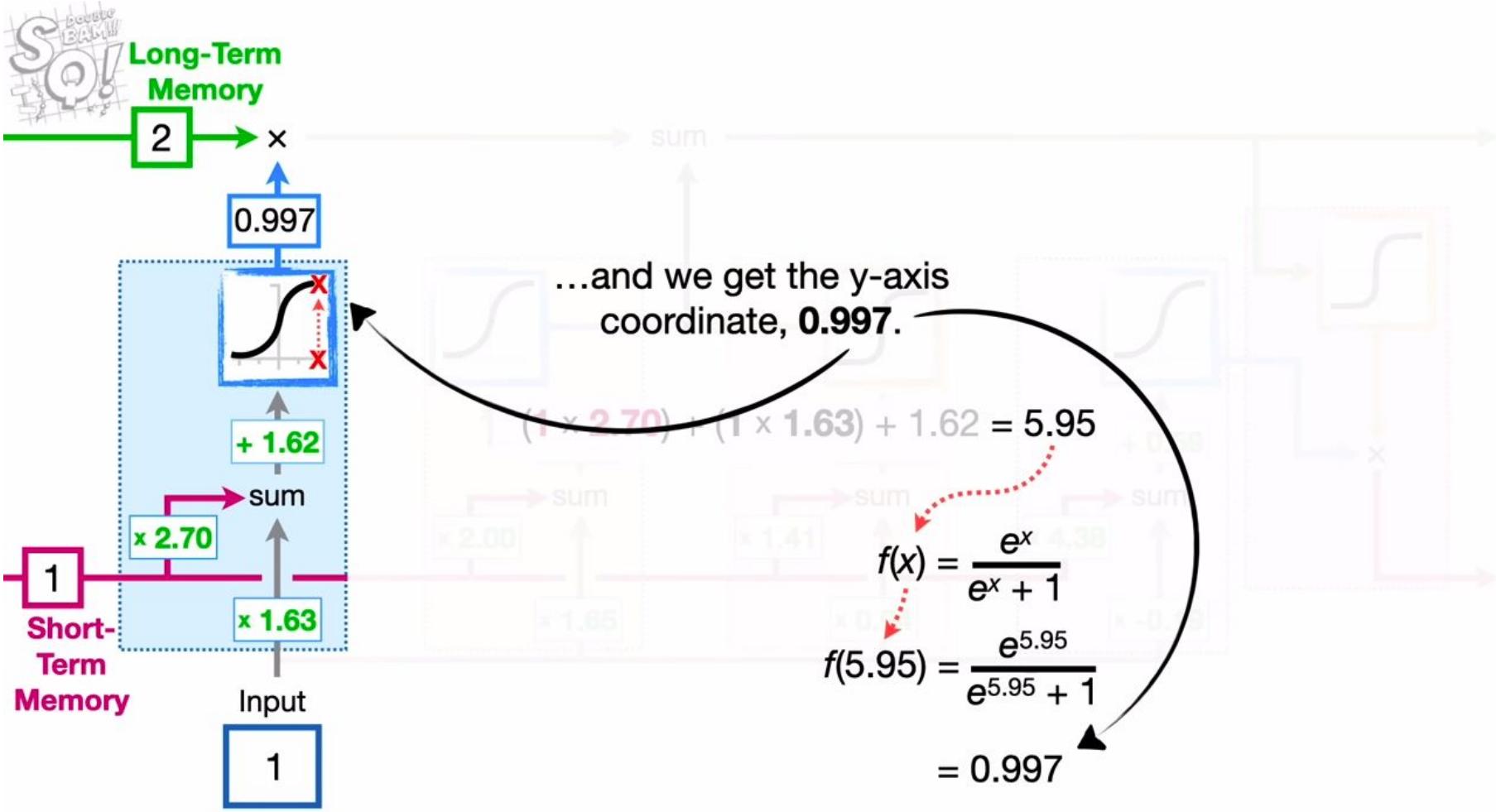
1

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

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





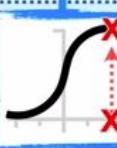




Long-Term Memory

2

0.997



+ 1.62

x 2.70

x 1.63

Input

1

Short-Term Memory

1

Lastly, we multiply the **Long-Term Memory**, 2, by the y-axis coordinate, 0.997...

sum

sum

sum

sum

sum

sum

sum

sum

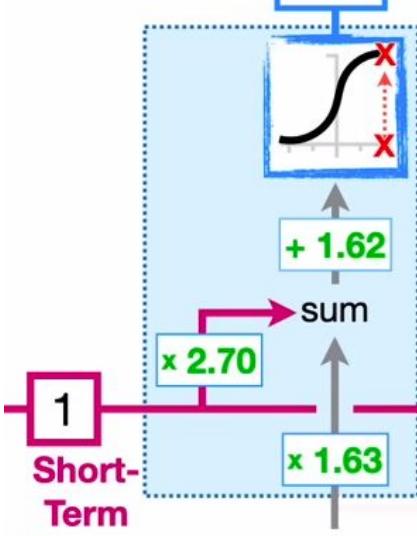


Long-Term Memory

$$2 \rightarrow \times \rightarrow 1.99$$

0.997

...and the result is 1.99



Short-Term Memory

1

1

1



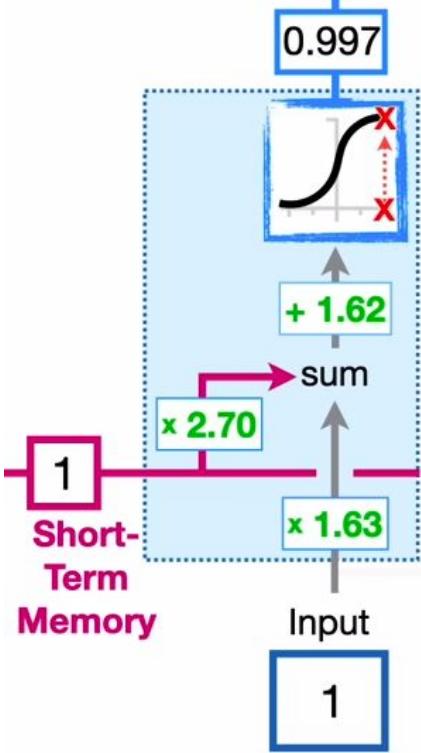


Long-Term
Memory

$$2 \rightarrow \times \rightarrow 1.99$$

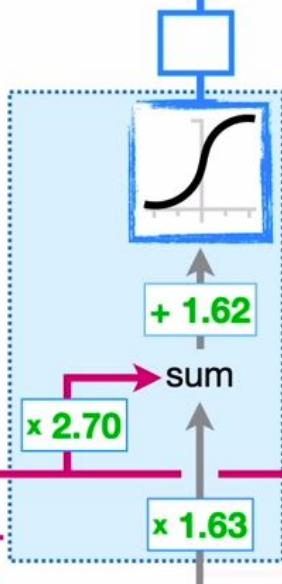
0.997

So this first stage of the **Long Short-Term Memory (LSTM)** unit reduced the **Long-Term Memory** by a little bit.

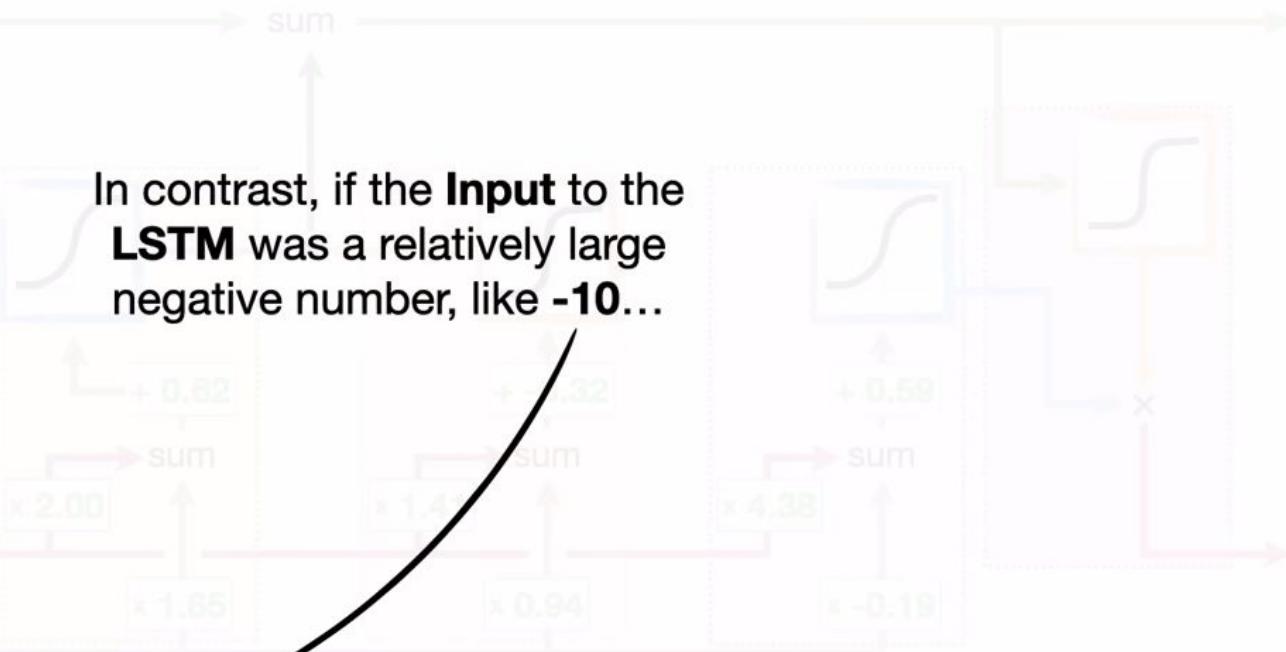
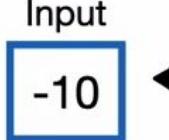




Long-Term Memory



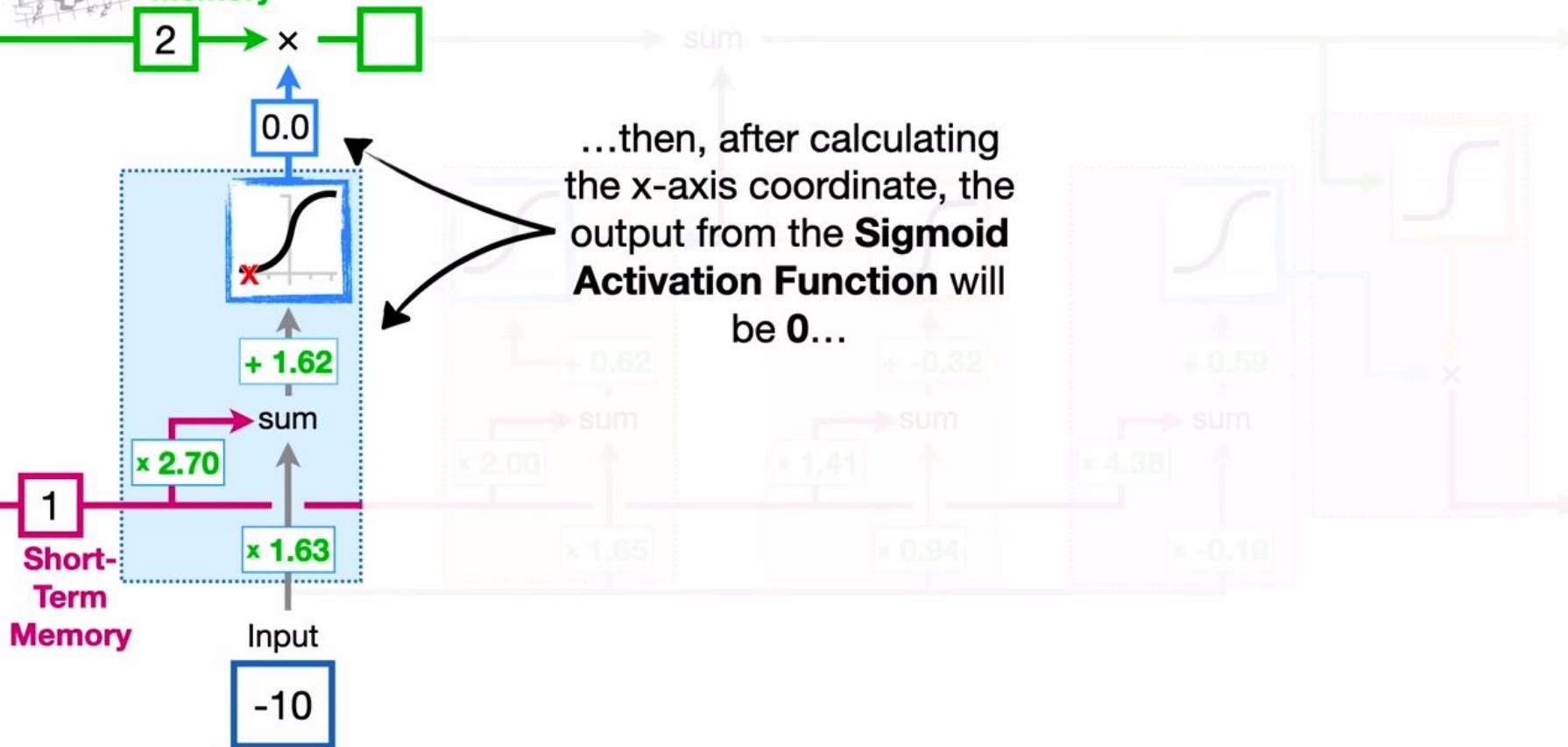
In contrast, if the **Input** to the **LSTM** was a relatively large negative number, like **-10**...



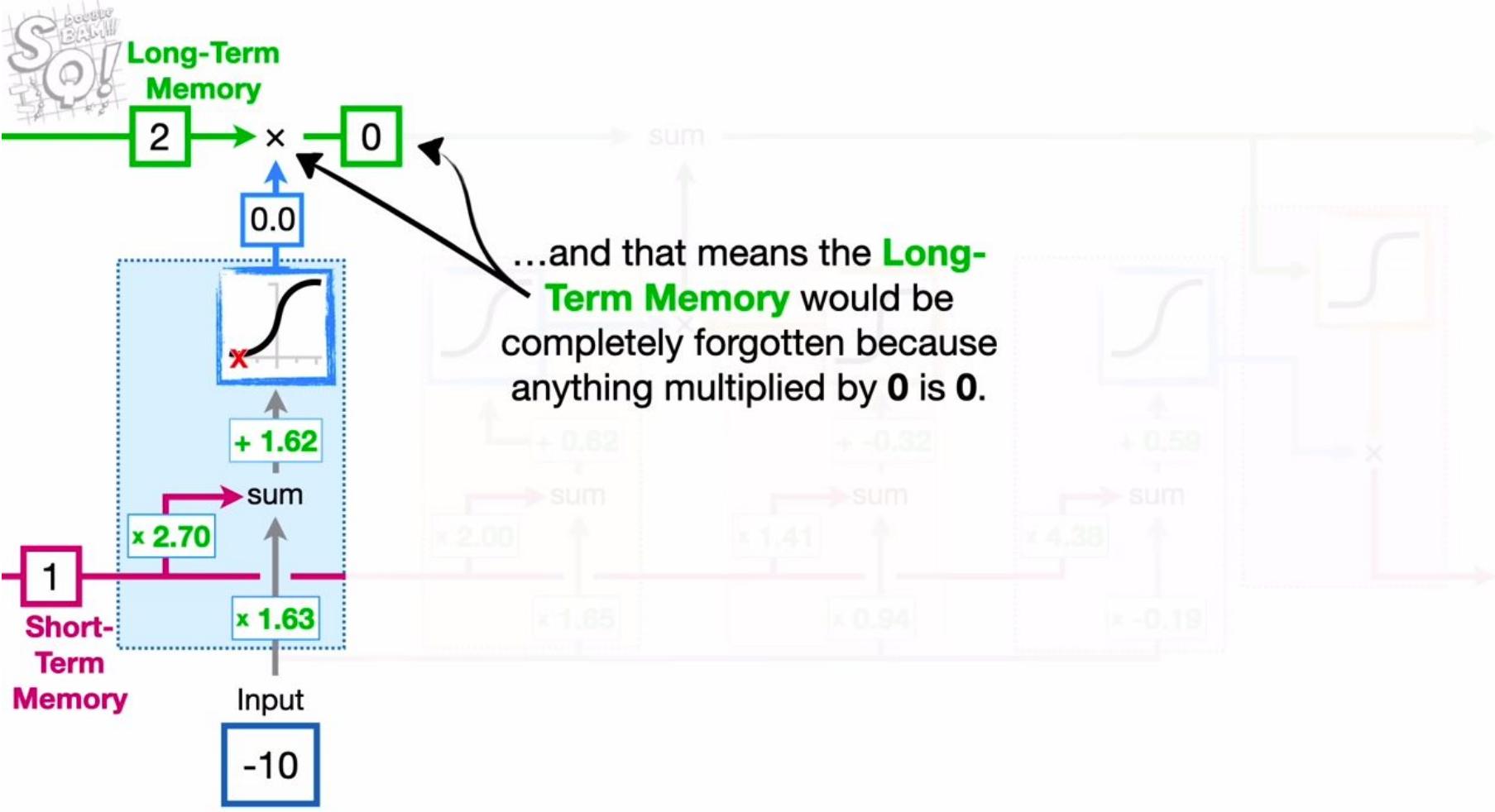


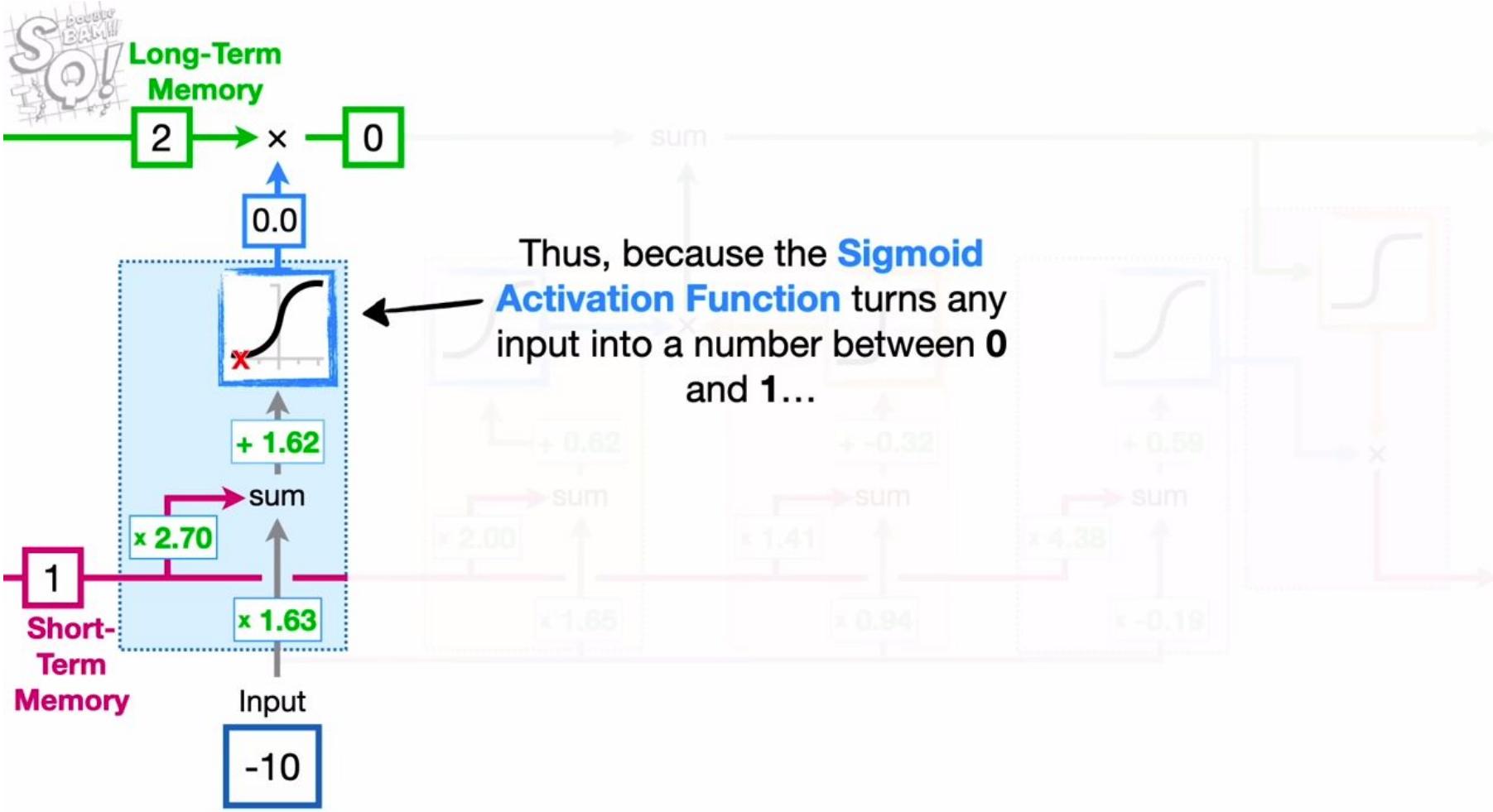
Long-Term Memory

2



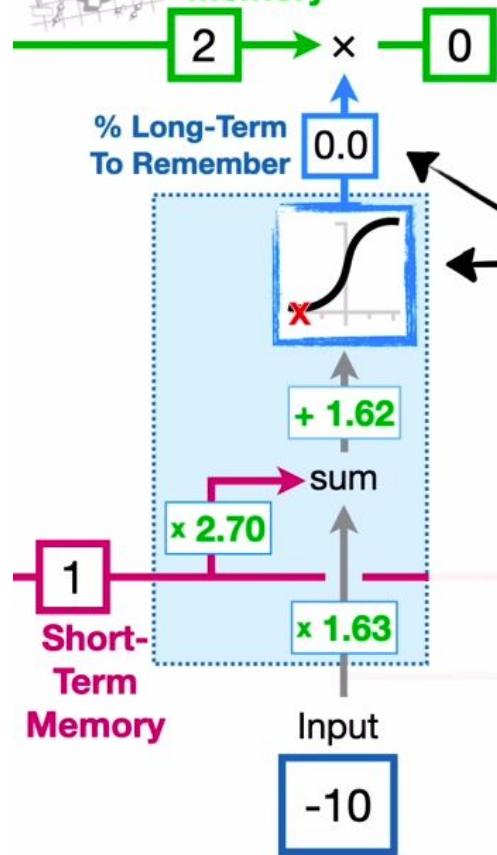
...then, after calculating
the x-axis coordinate, the
output from the **Sigmoid
Activation Function** will
be 0...



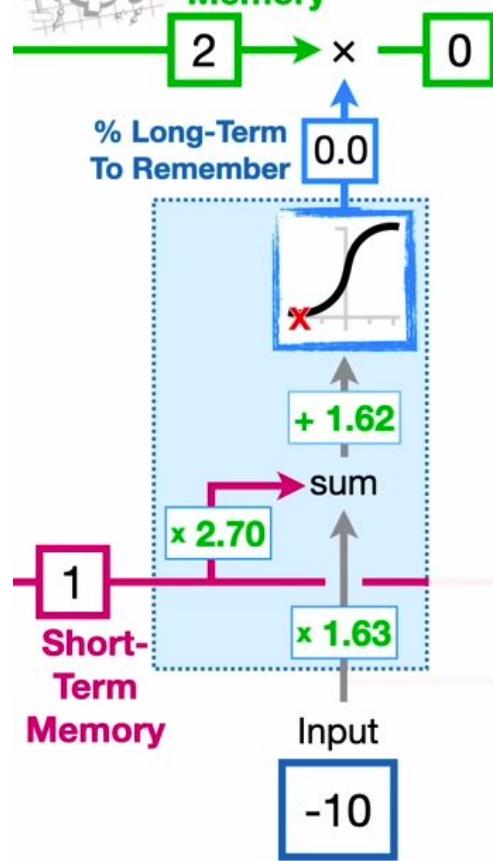




Long-Term Memory



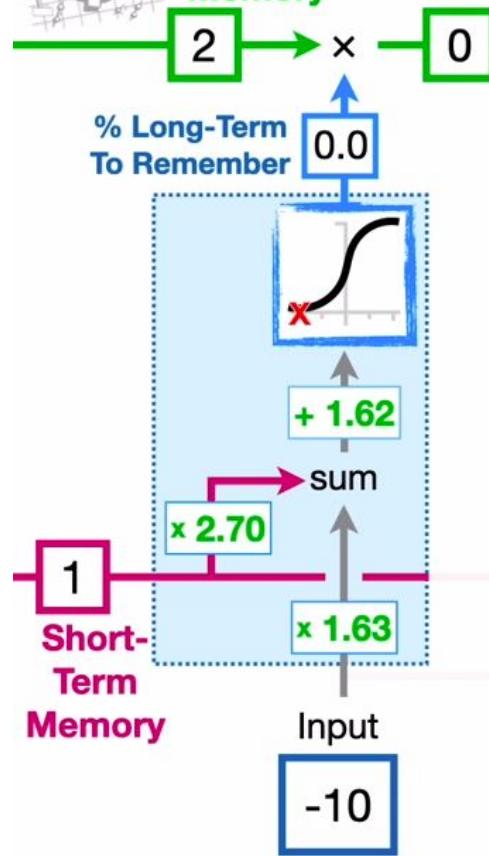
...the output determines what percentage of the **Long-Term Memory** is remembered.



To summarize, the first stage in a **Long Short-Term Memory** unit determines what percentage of the **Long-Term Memory** is remembered.

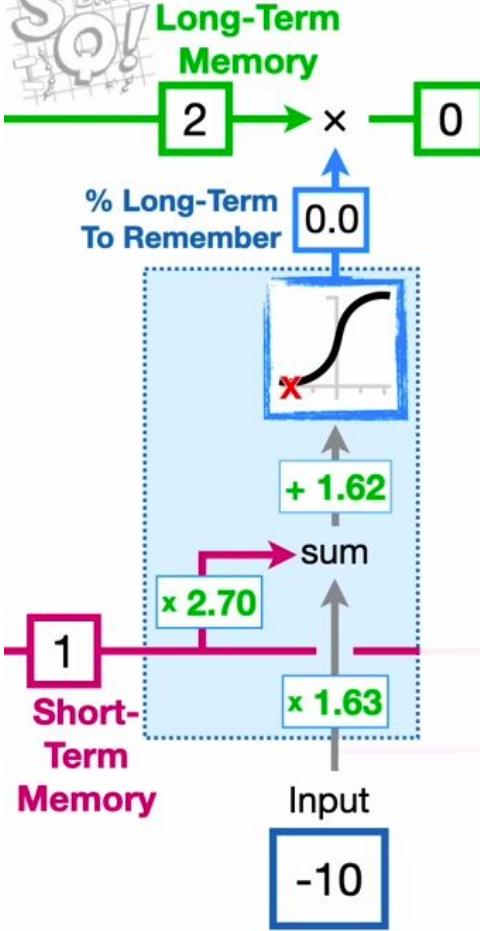


Long-Term Memory



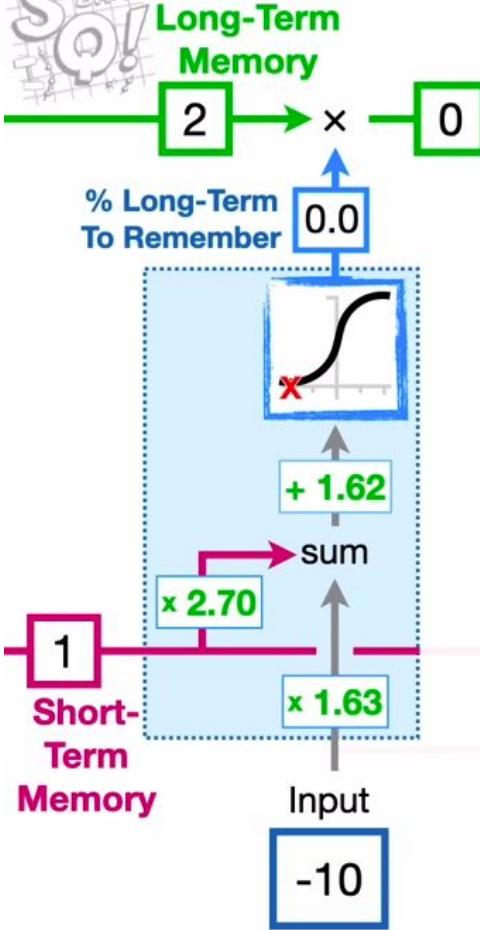
TERMINOLOGY ALERT!!!





TERMINOLOGY ALERT!!!

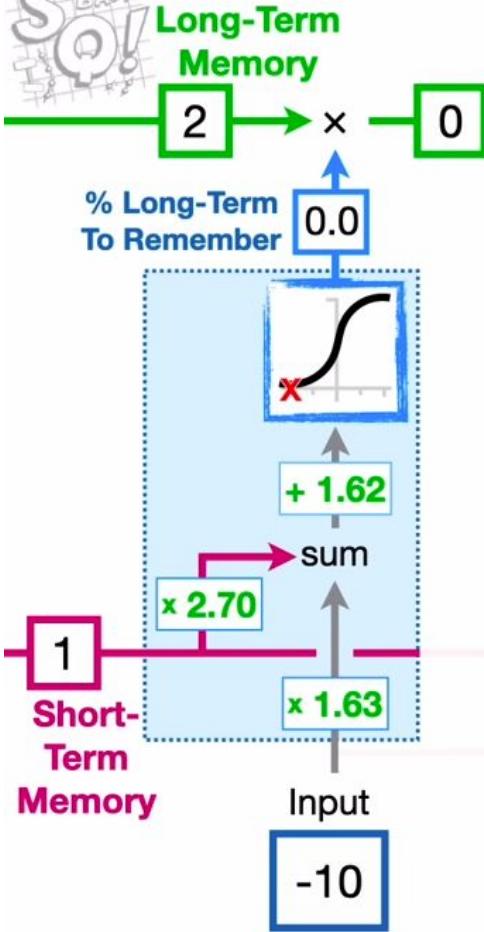
Even though this part of the **Long Short-Term Memory** unit determines what percentage of the **Long-Term Memory** will be remembered...



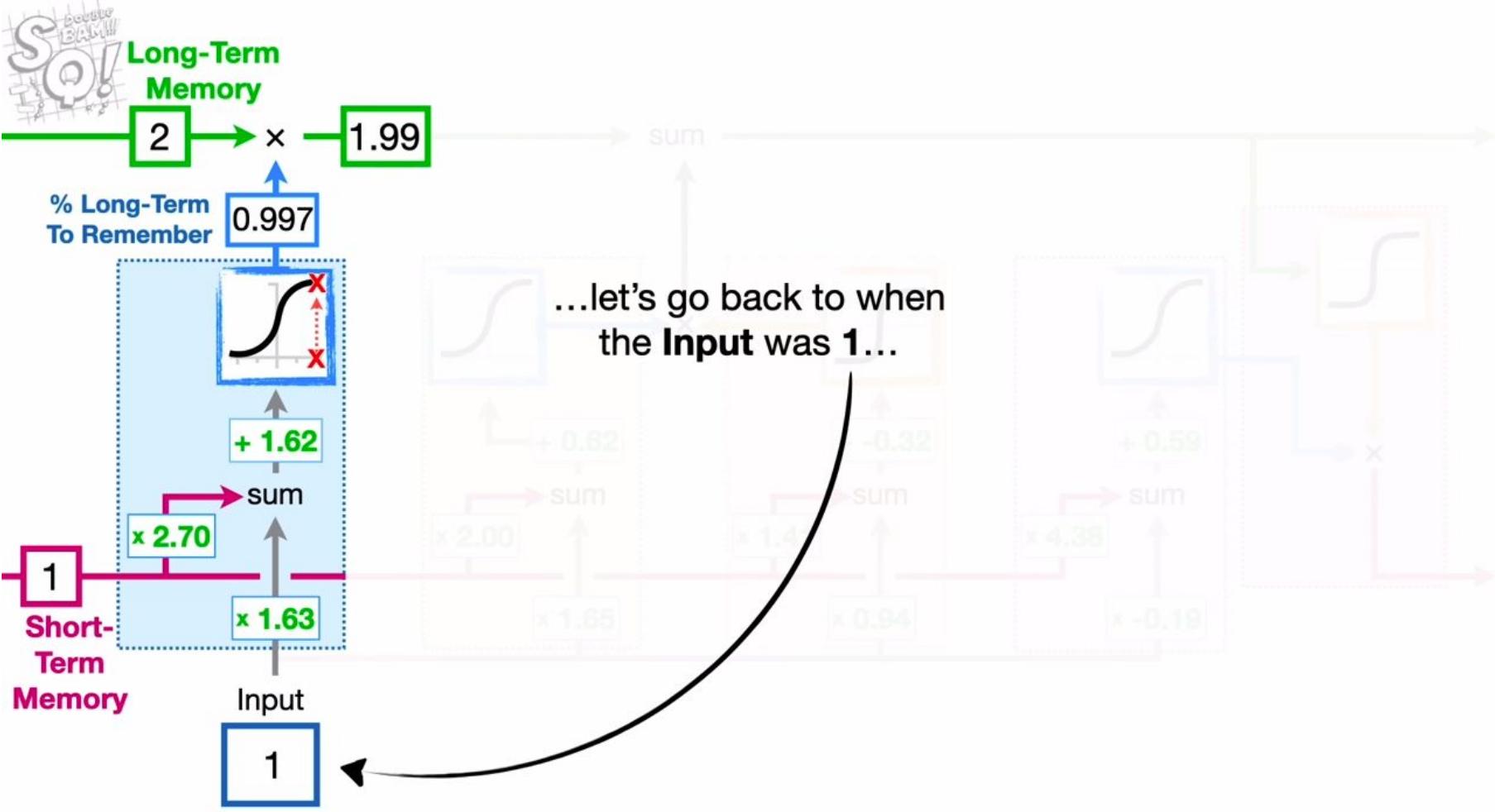
TERMINOLOGY ALERT!!!

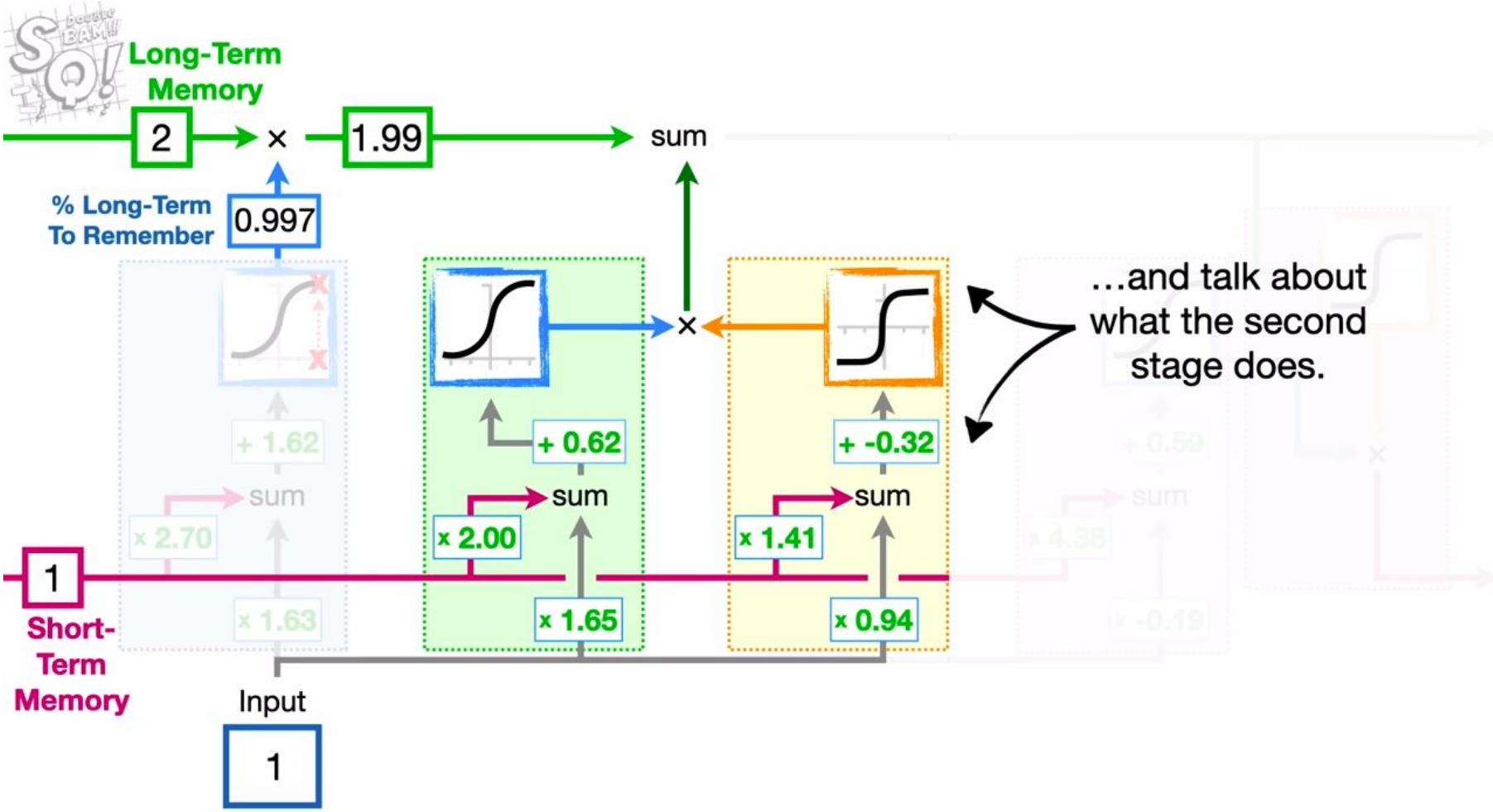
Even though this part of the **Long Short-Term Memory** unit determines what percentage of the **Long-Term Memory** will be remembered...

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



Now that we know what the first part of the **LSTM** unit does, it determines what percentage of the **Long-Term Memory** will be remembered...



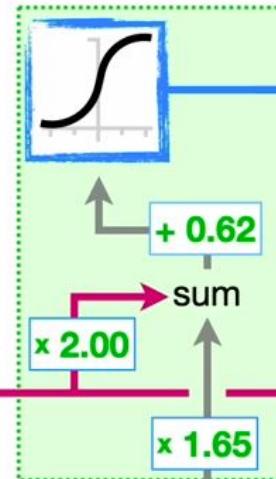




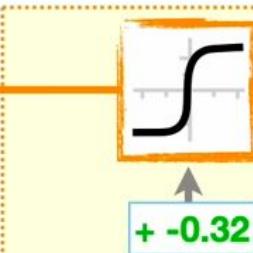
Long-Term Memory

$$2 \rightarrow x \quad 1.99 \rightarrow \text{sum}$$

% Long-Term To Remember
0.997



Potential Long-Term Memory



In a nutshell, the block on the right combines the **Short-Term Memory** and the **Input**...

...to create a **Potential Long-Term Memory**...



Long-Term Memory

$$2 \rightarrow \times \quad 1.99$$

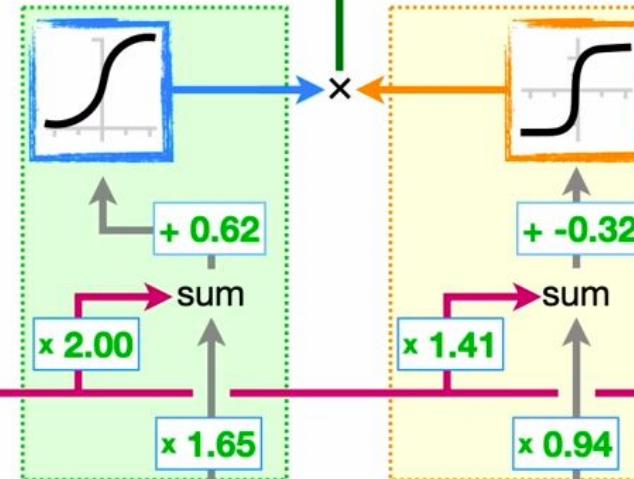
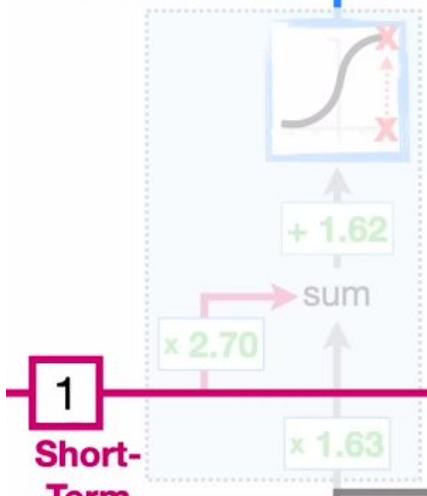
$$\% \text{ Long-Term} \rightarrow \text{To Remember} \quad 0.997$$

$$\% \text{ Potential Memory} \rightarrow \text{To Remember}$$

sum

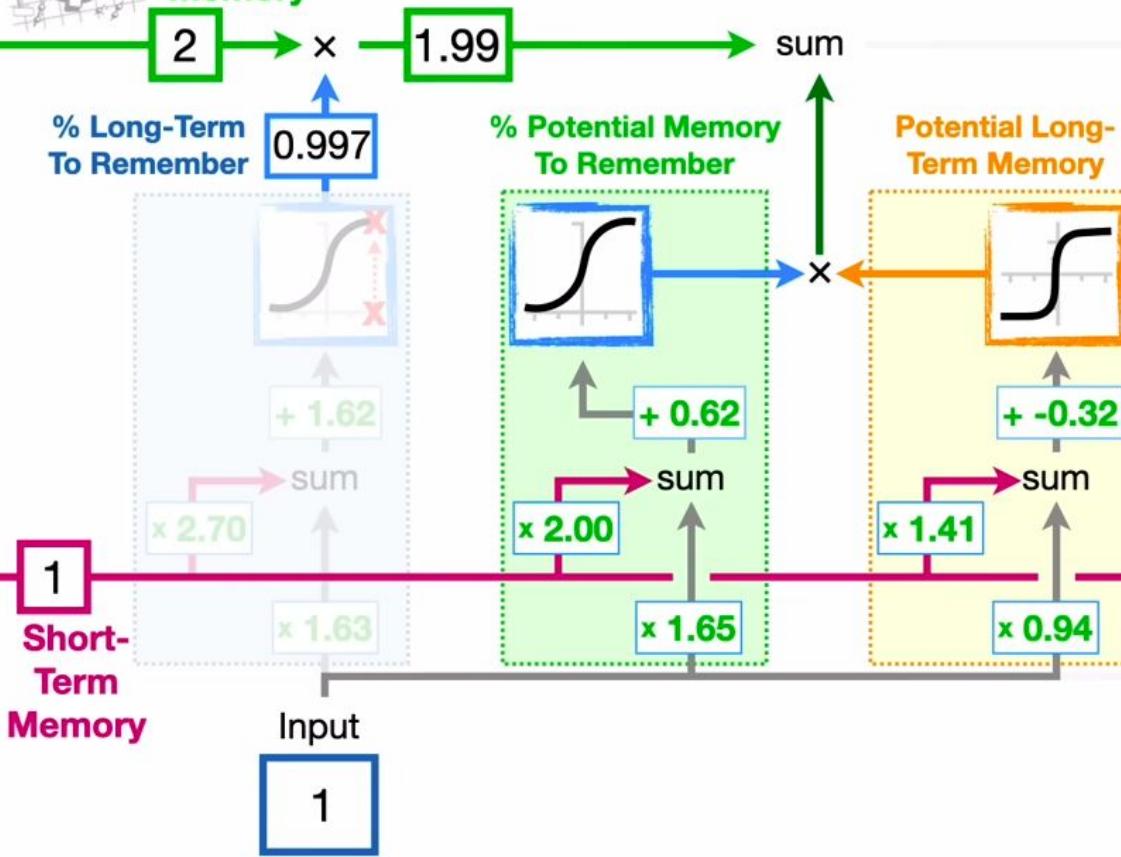
$$\text{Potential Long-Term Memory}$$

...and the block on the left determines what percentage of that **Potential Memory** to add to the **Long-Term Memory**.

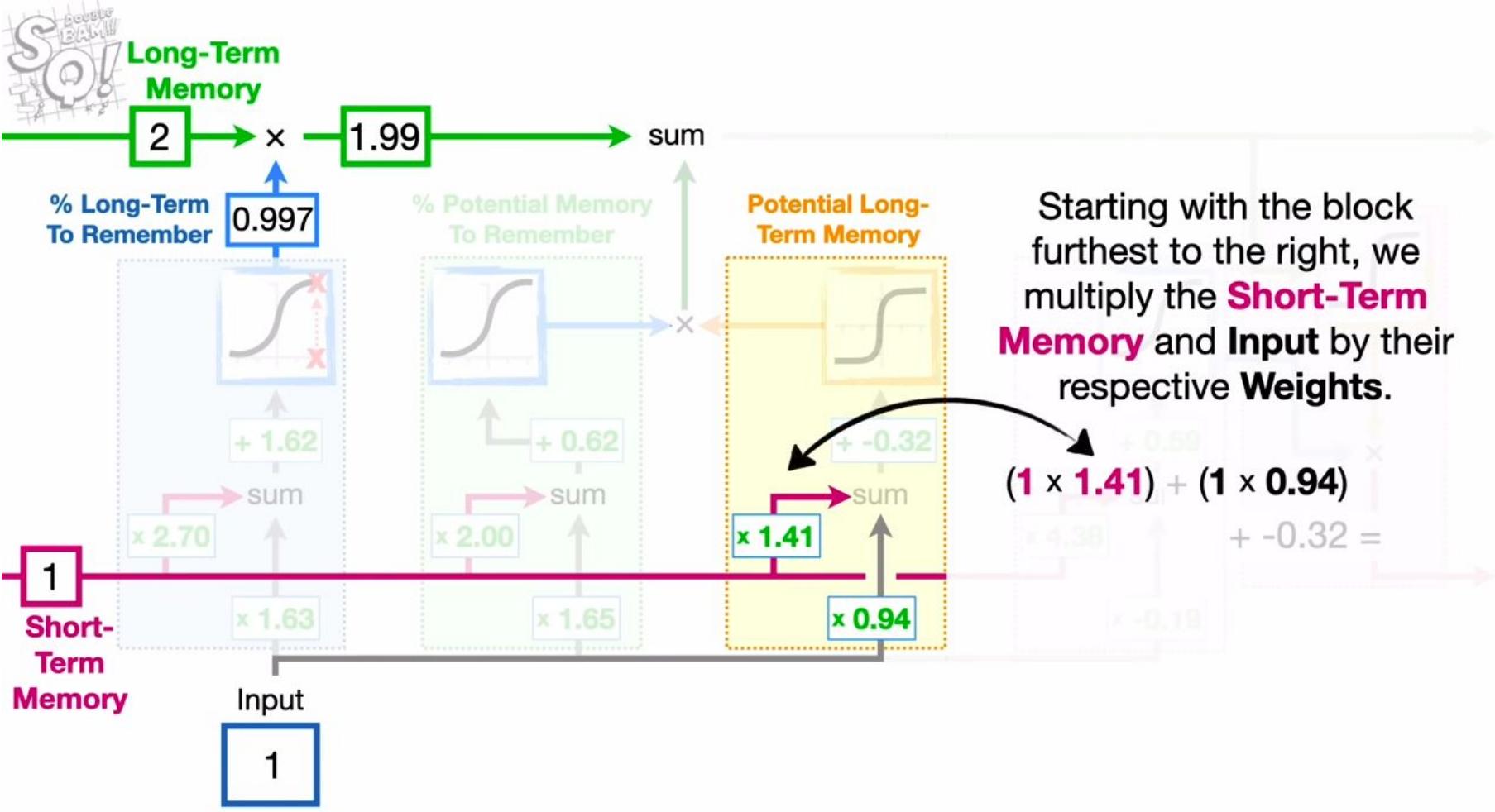


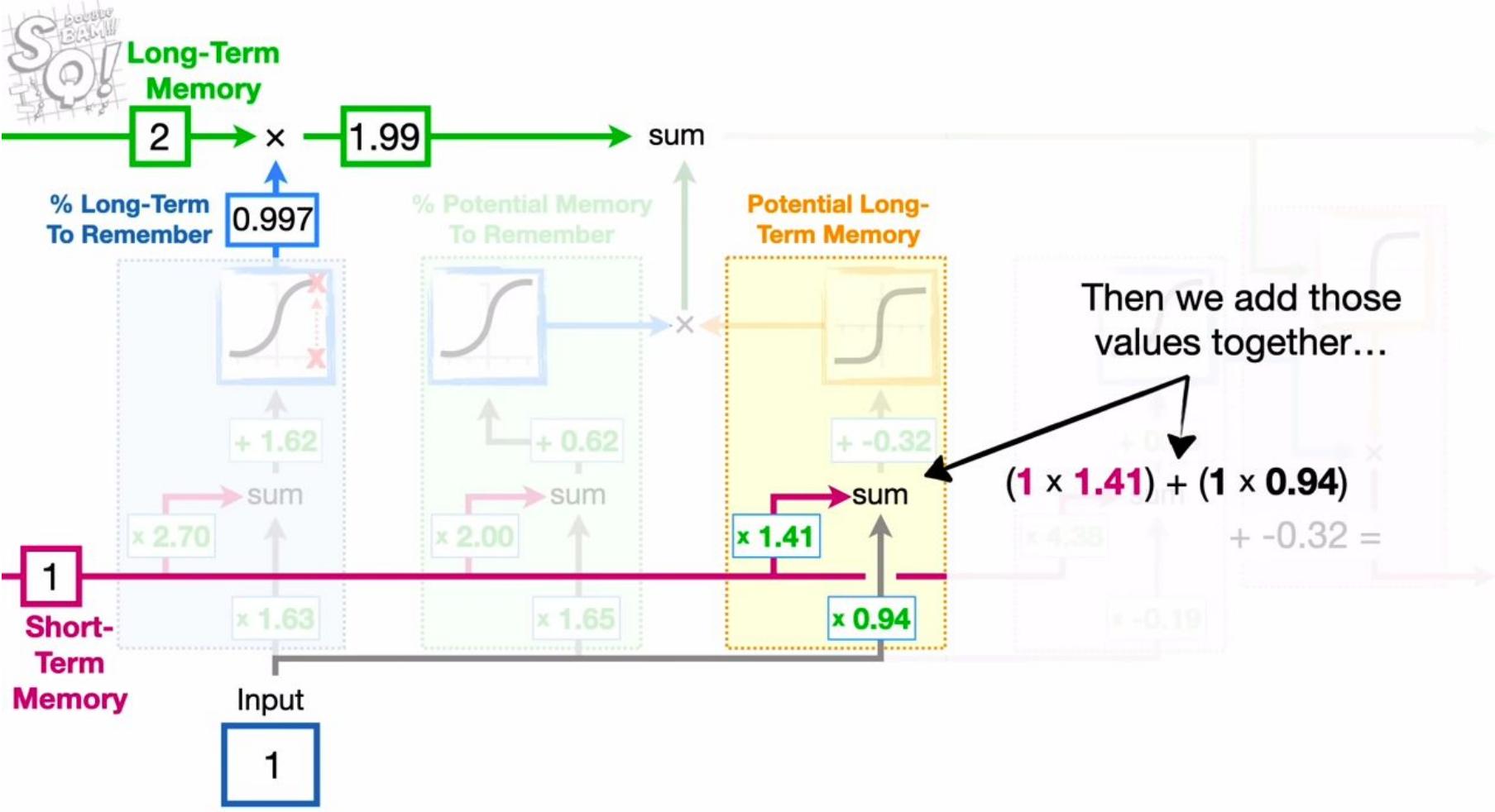


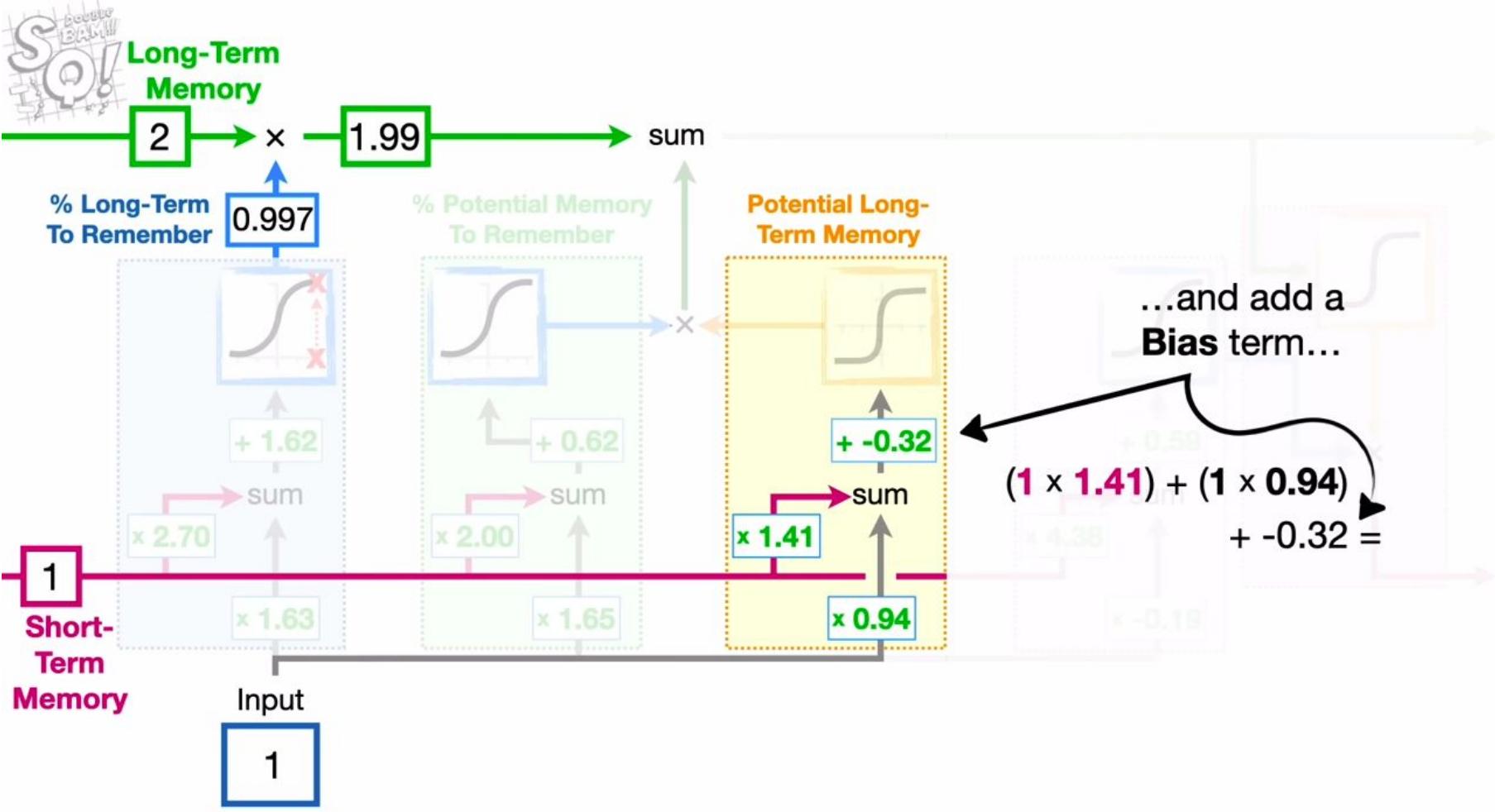
Long-Term Memory

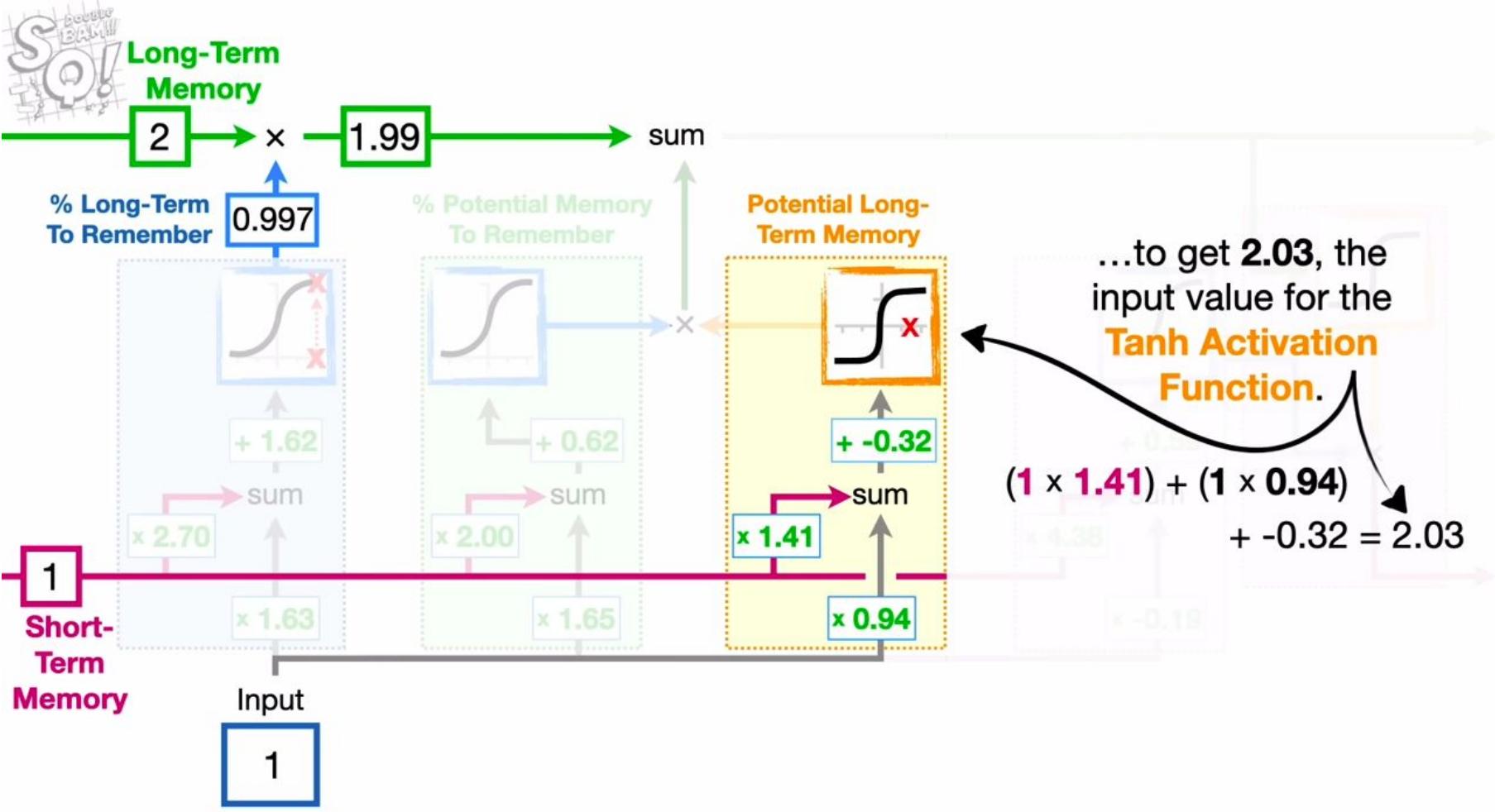


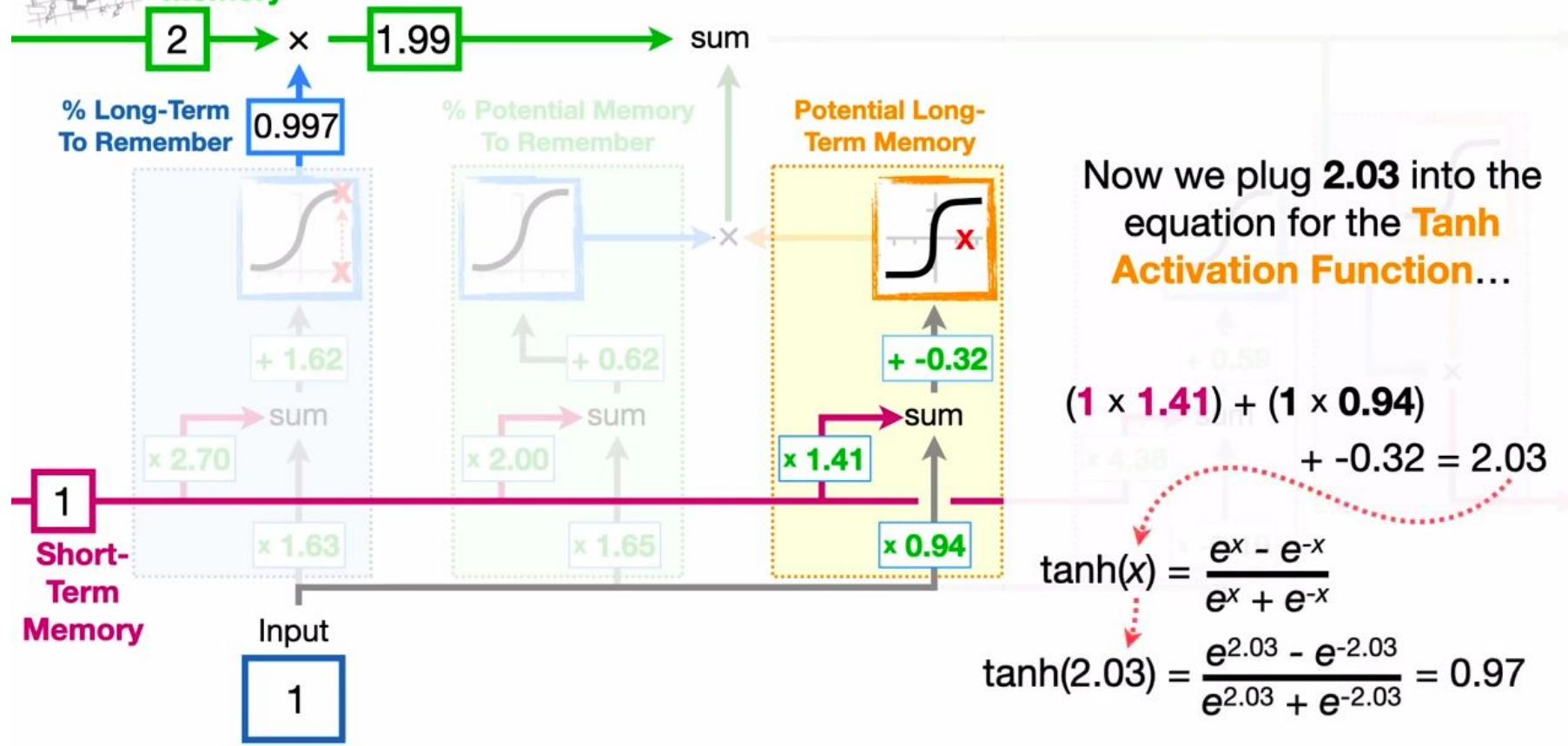
So let's plug the numbers in and do the math to see how a **Potential Memory** is created and how much of it is added to the **Long-Term Memory**.

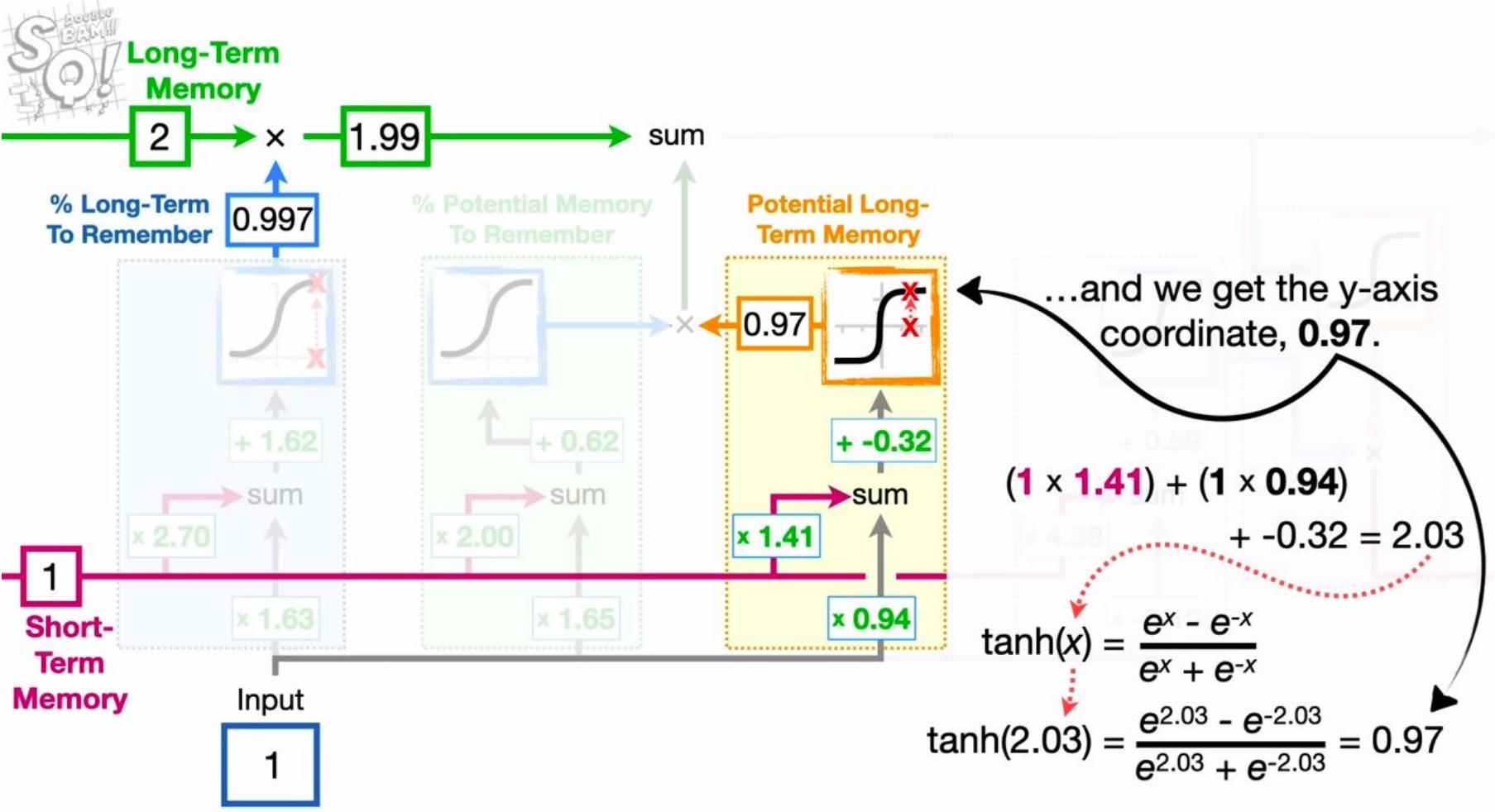


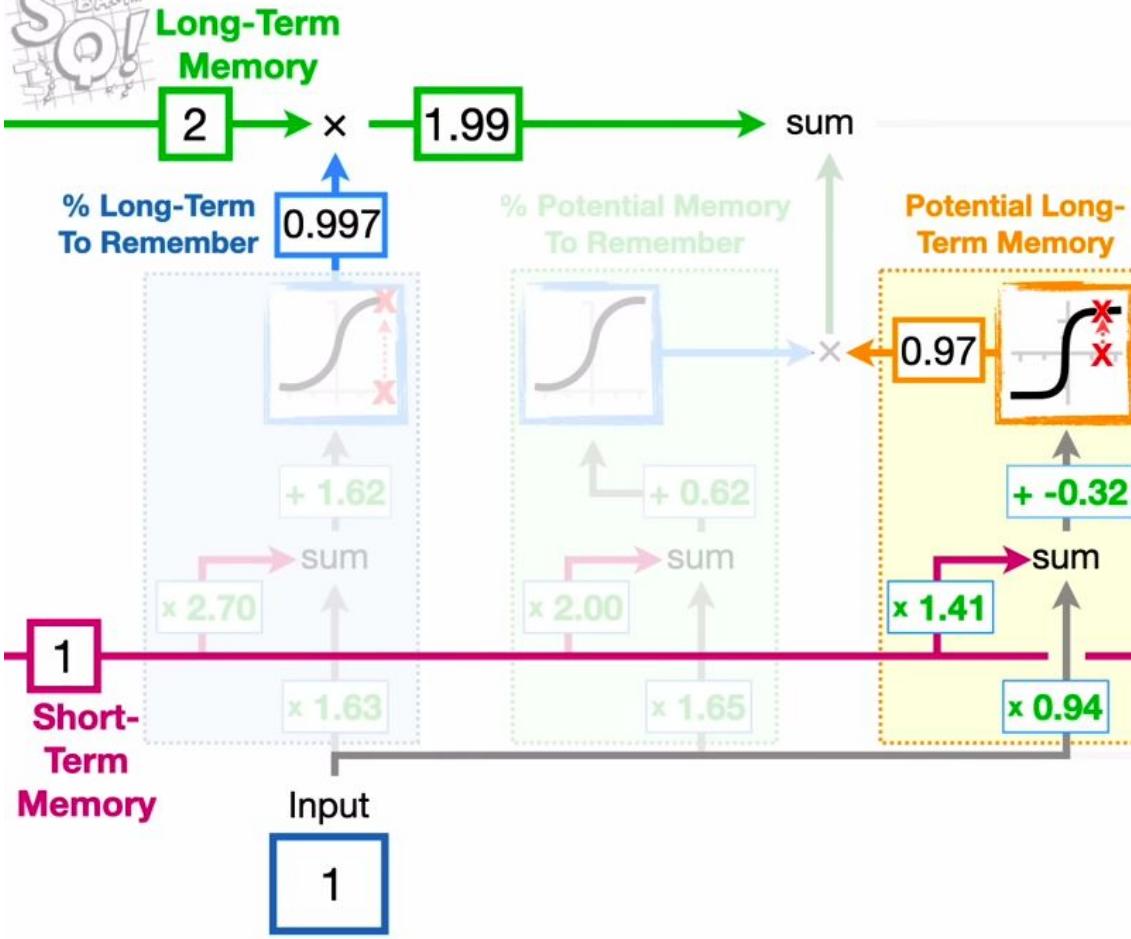




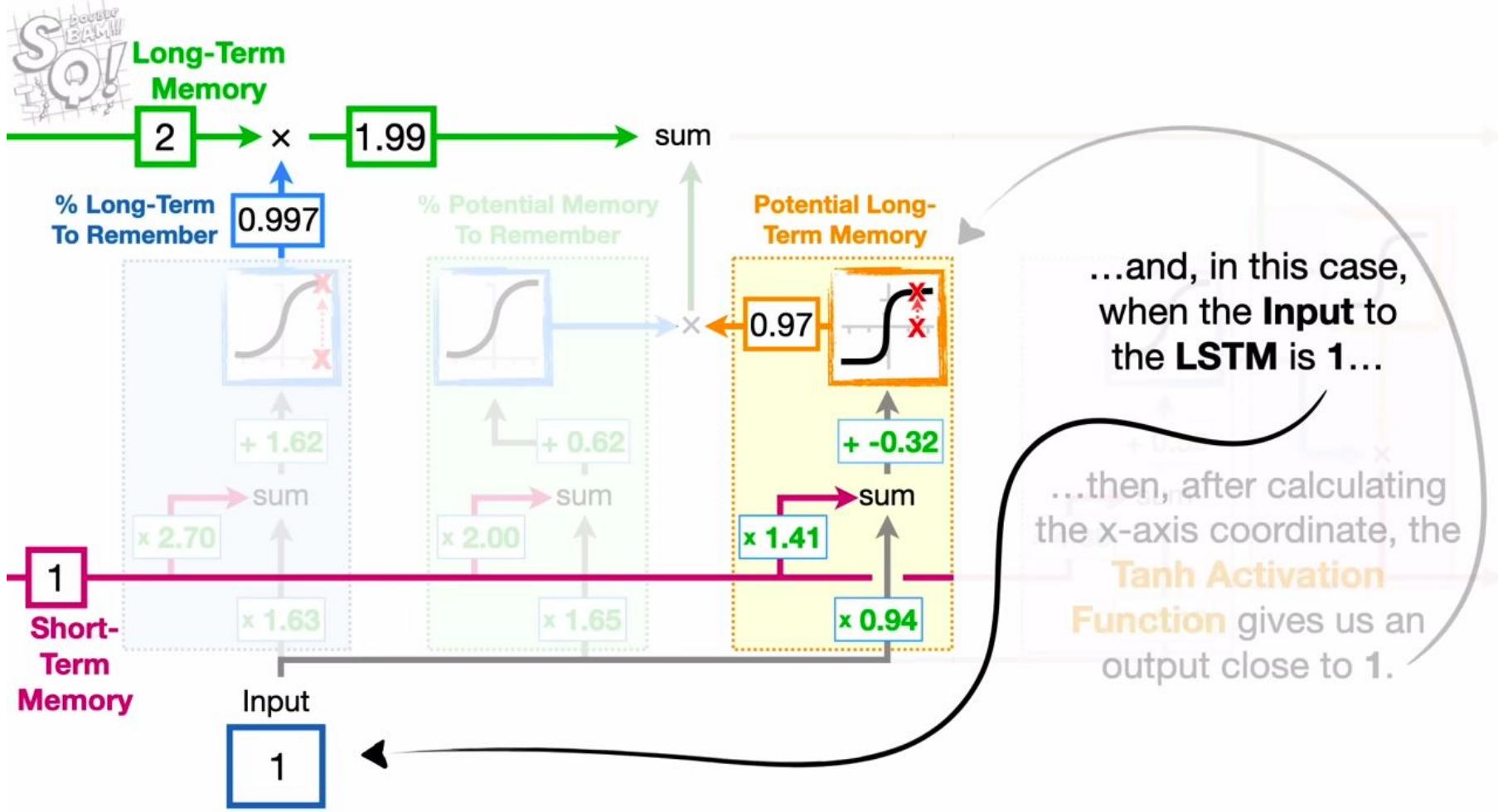


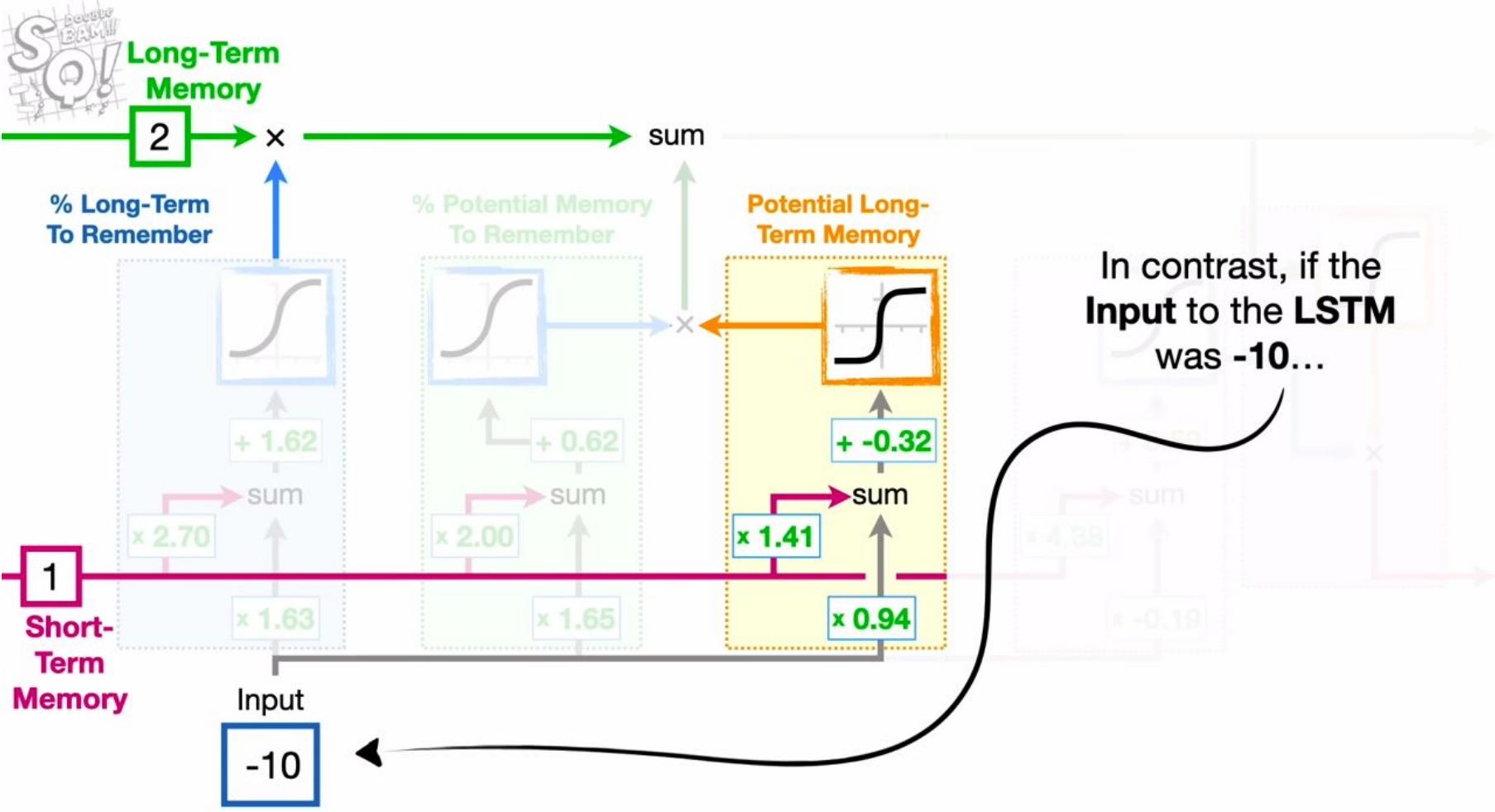






Remember, the **Tanh Activation Function** turns any input into a number between -1 and 1...







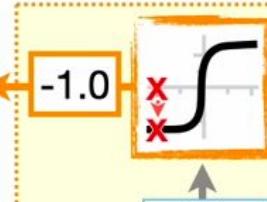
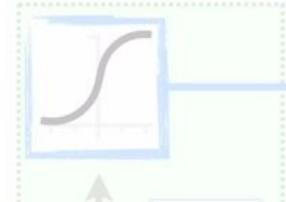
Long-Term Memory



% Long-Term
To Remember
0.0

% Potential Memory
To Remember

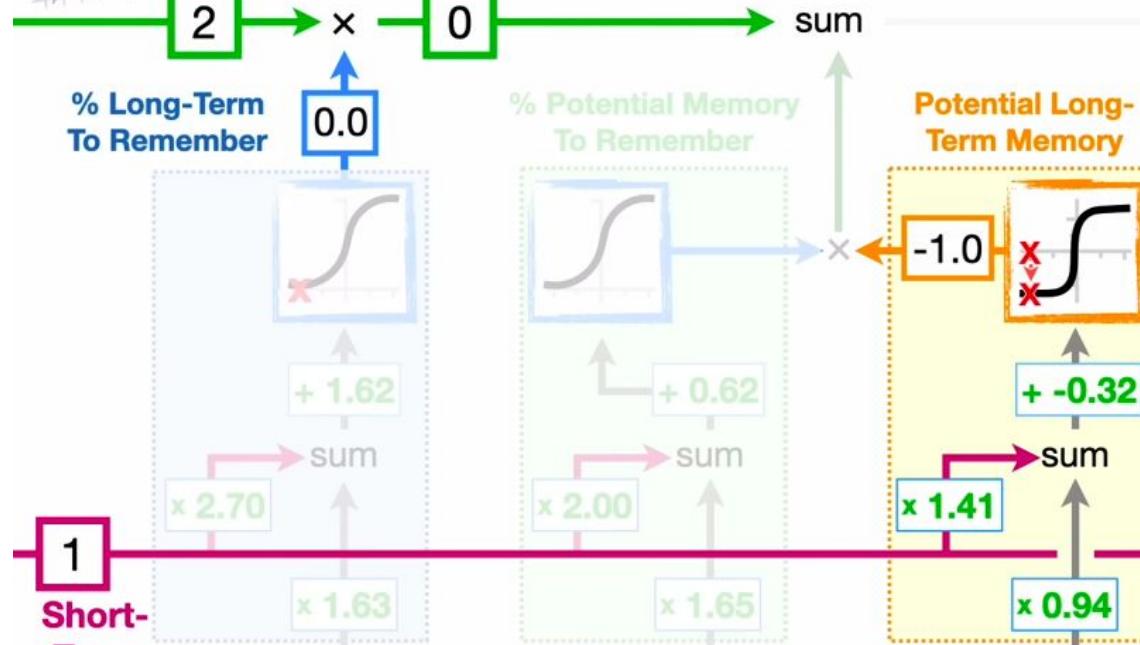
Potential Long-
Term Memory



1
Short-
Term
Memory

Input

-10



In contrast, if the
Input to the LSTM
was -10...

...then, after calculating
the x-axis coordinate,
the output from the Tanh
Activation Function
function would be -1.



Long-Term Memory

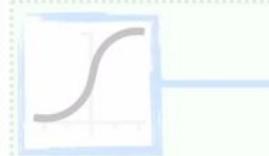
2

x 1.99

sum

% Long-Term
To Remember

0.997

% Potential Memory
To RememberPotential Long-
Term Memory

Short-Term Memory

1

x 2.70

sum

x 1.63

Input

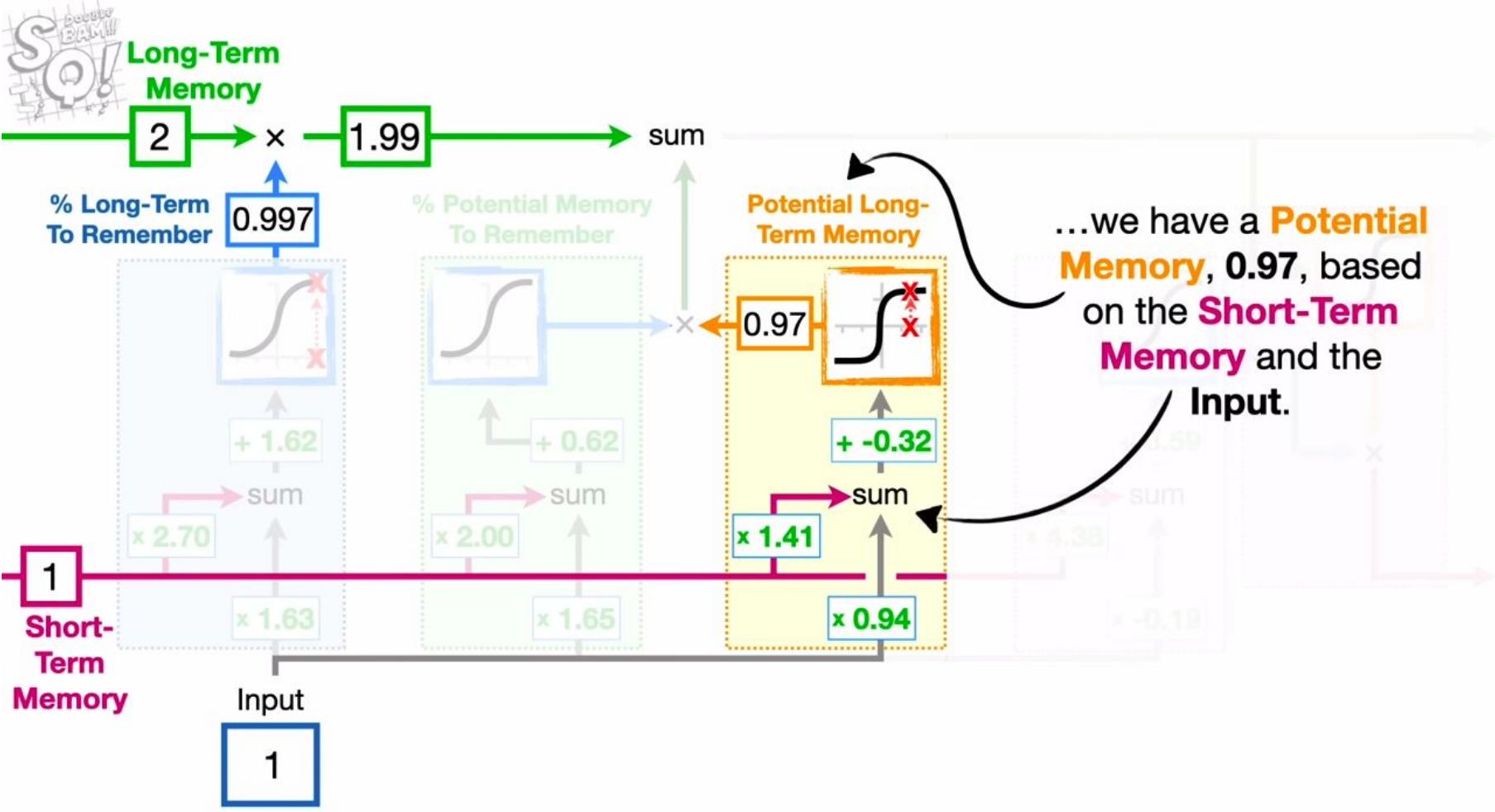
1

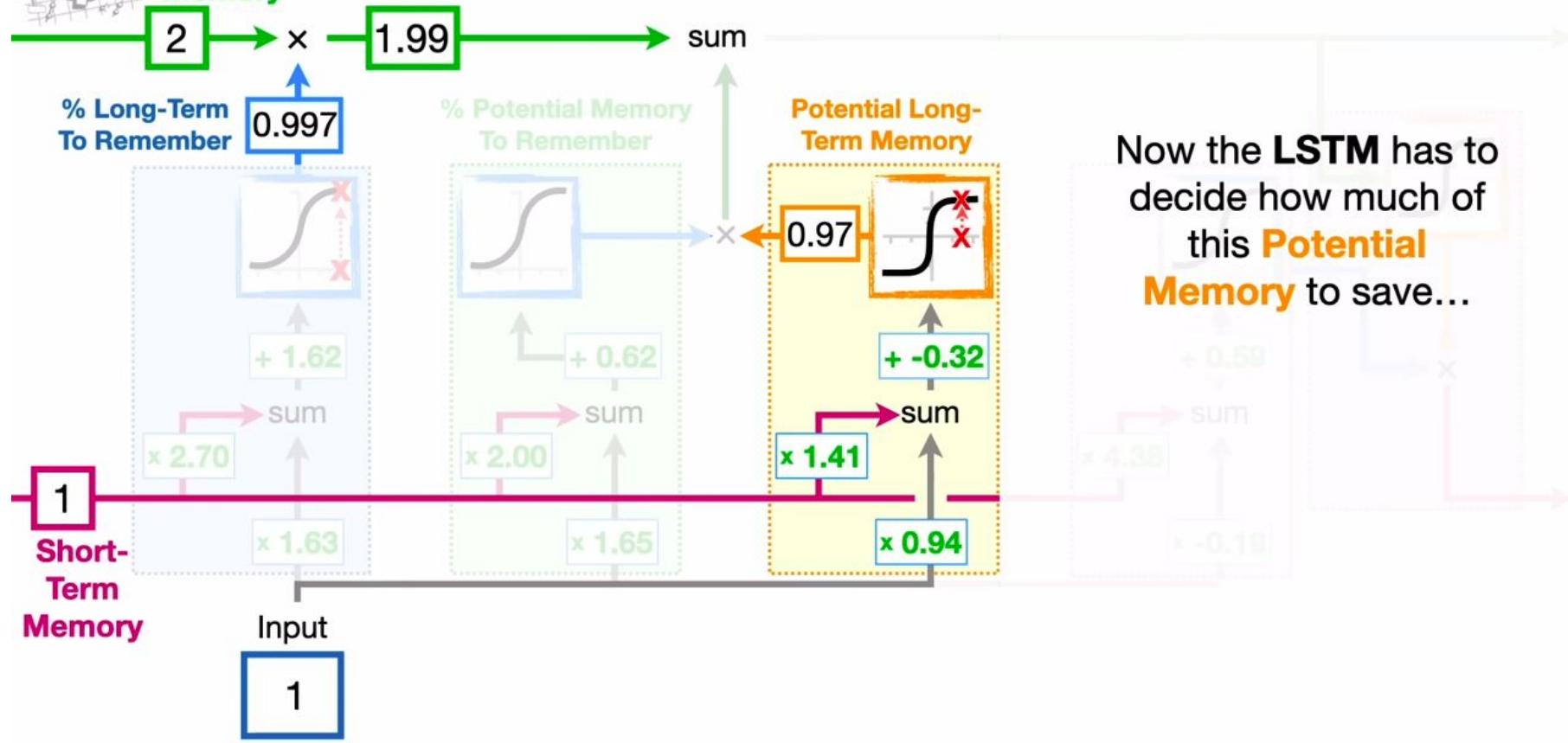
x 1.41

sum

x 0.94

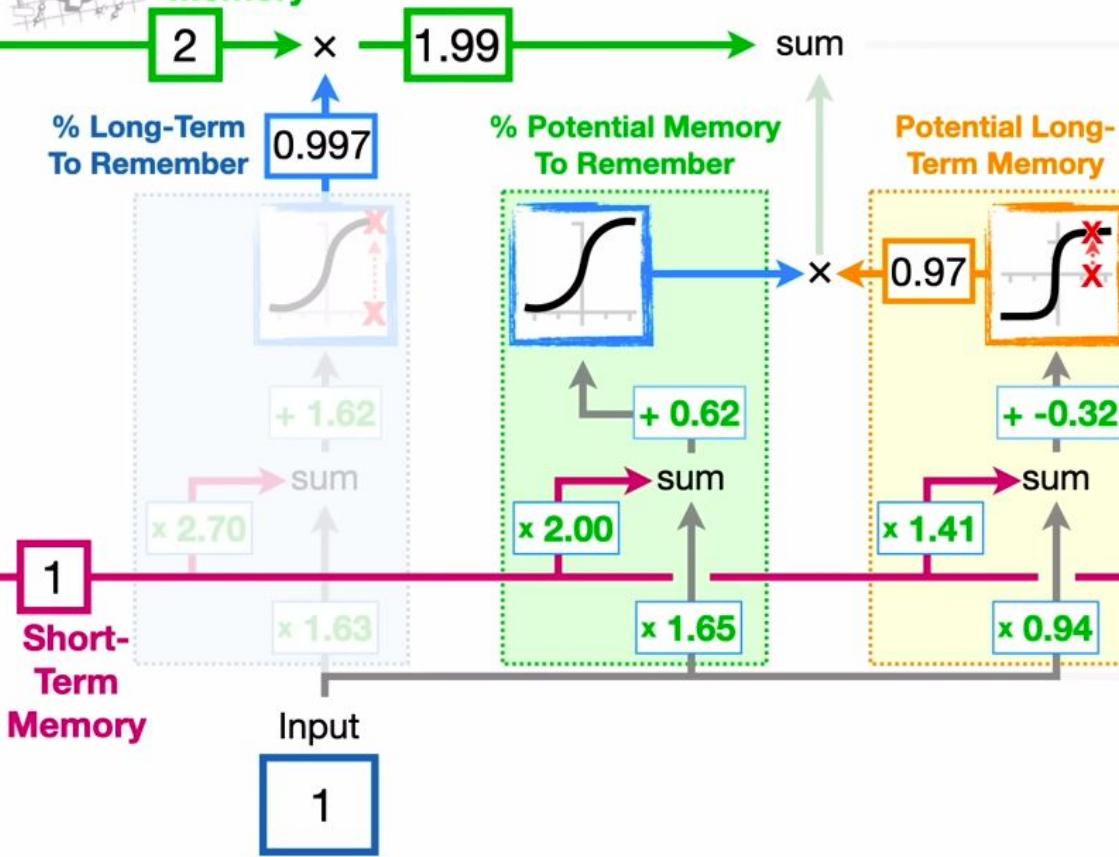
Going back to when
the **Input** to the
LSTM was 1...







Long-Term Memory



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



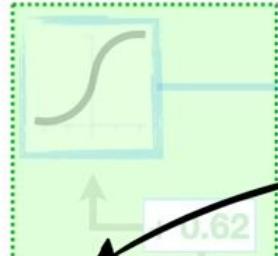
Long-Term Memory

$$2 \rightarrow x \quad 1.99 \rightarrow \text{sum}$$

$$\% \text{ Long-Term} \rightarrow \text{To Remember} \quad 0.997$$

$$\% \text{ Potential Memory} \rightarrow \text{To Remember} \quad 0.62$$

$$\text{Potential Long-Term Memory} \quad 0.97$$

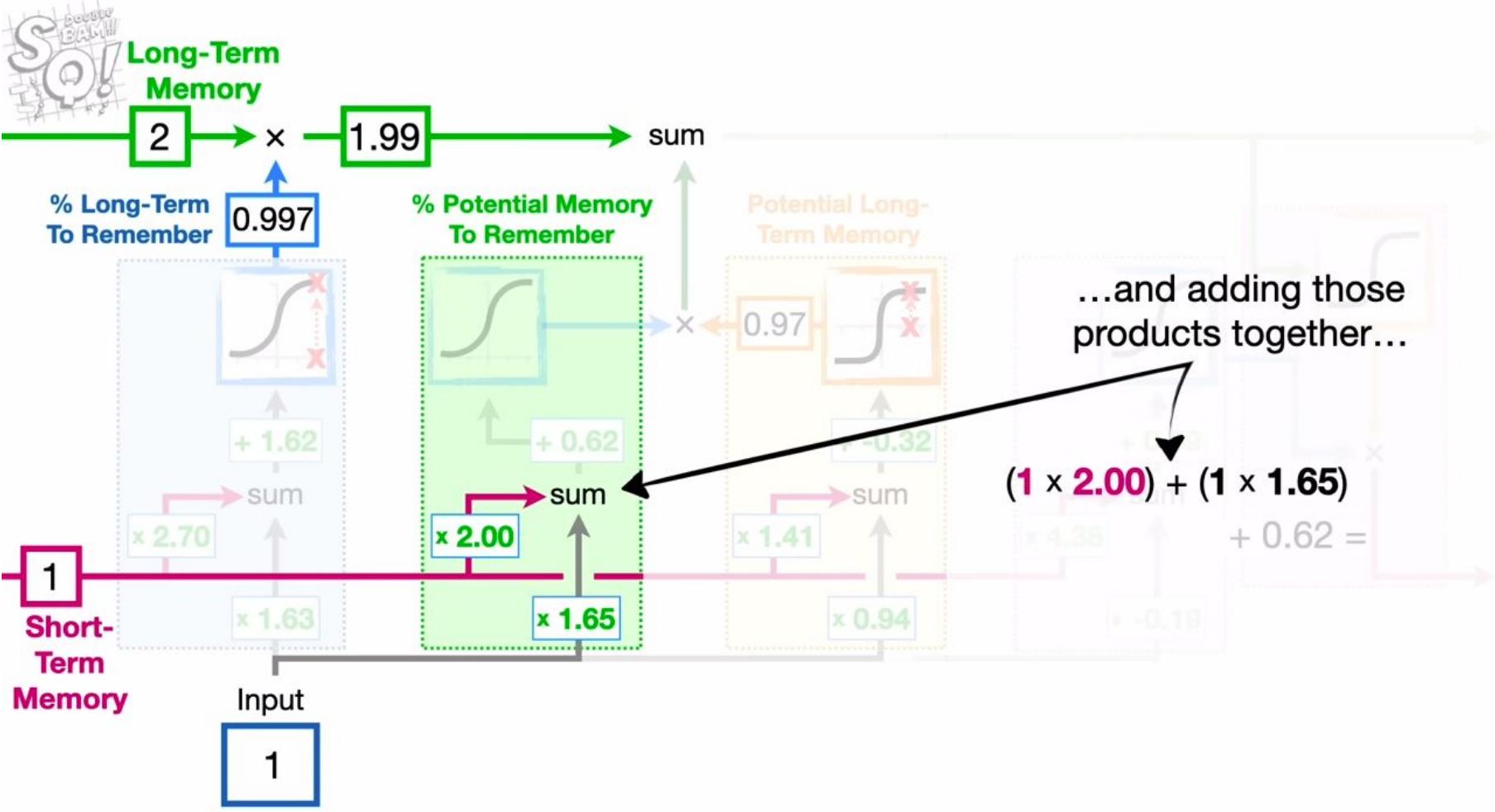


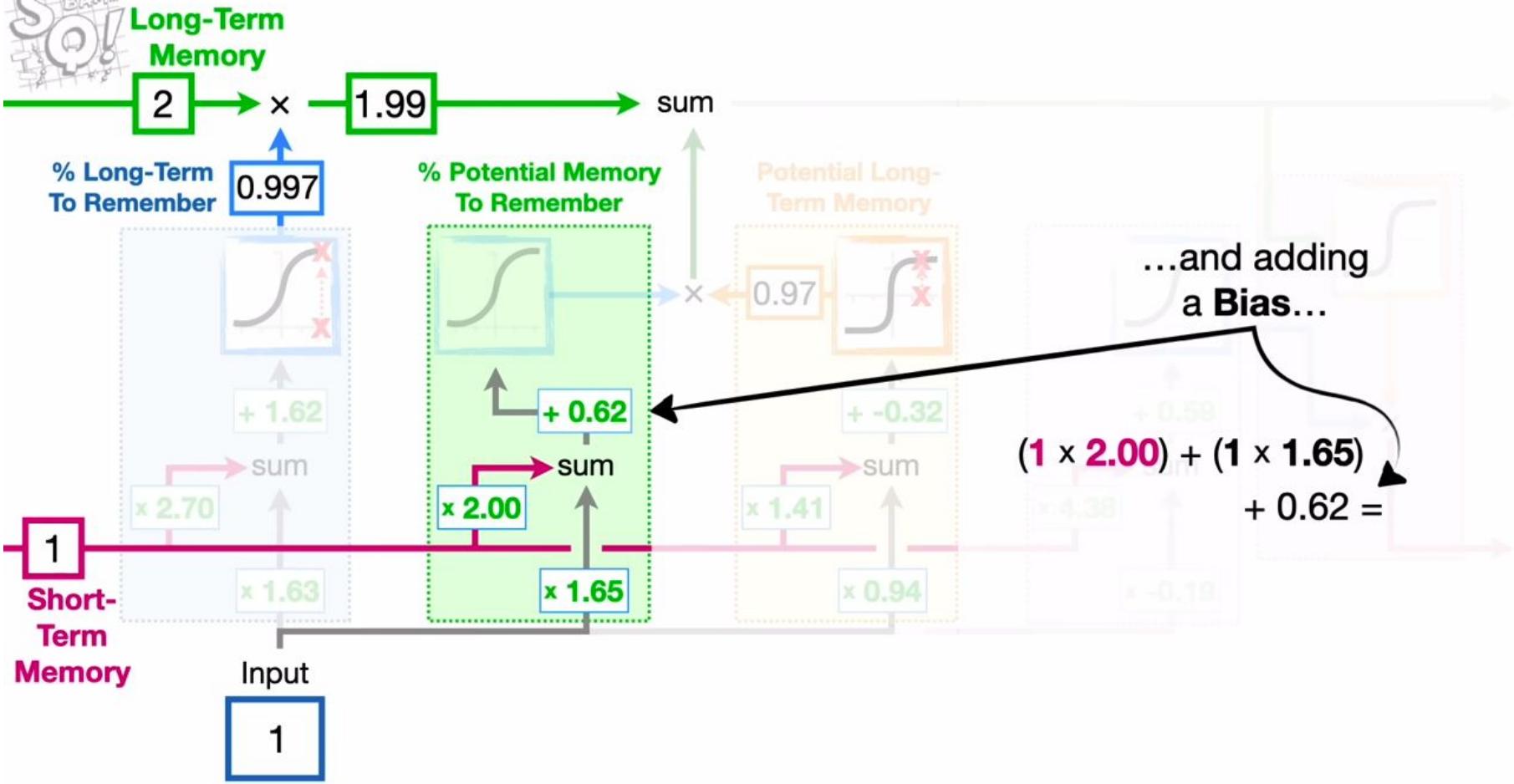
$$\begin{array}{l} 1 \\ \text{Short-Term} \\ \text{Memory} \\ \text{Input} \\ 1 \end{array}$$

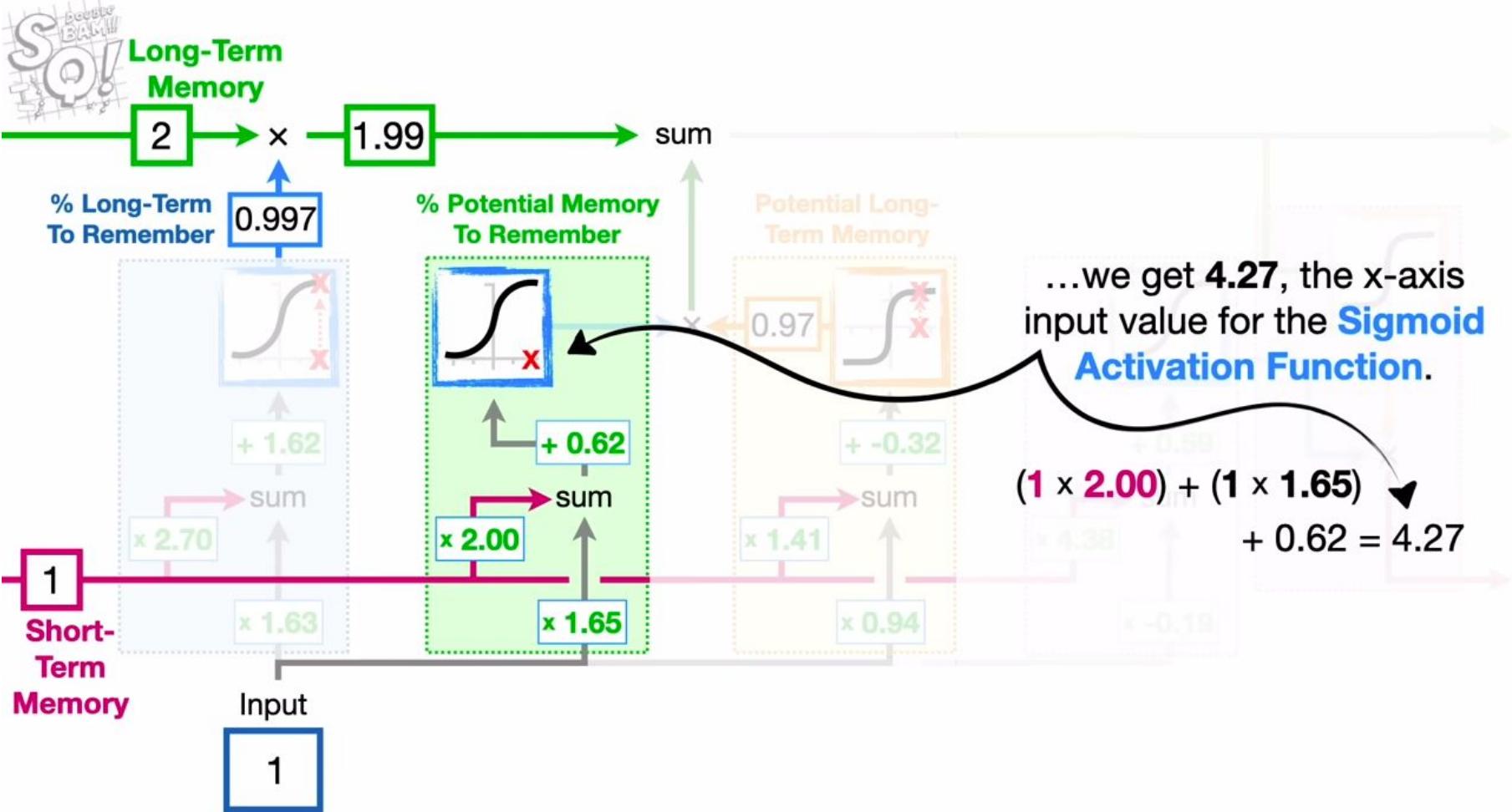
In other words, after multiplying the **Short-Term Memory** and the **Input** by **Weights**...

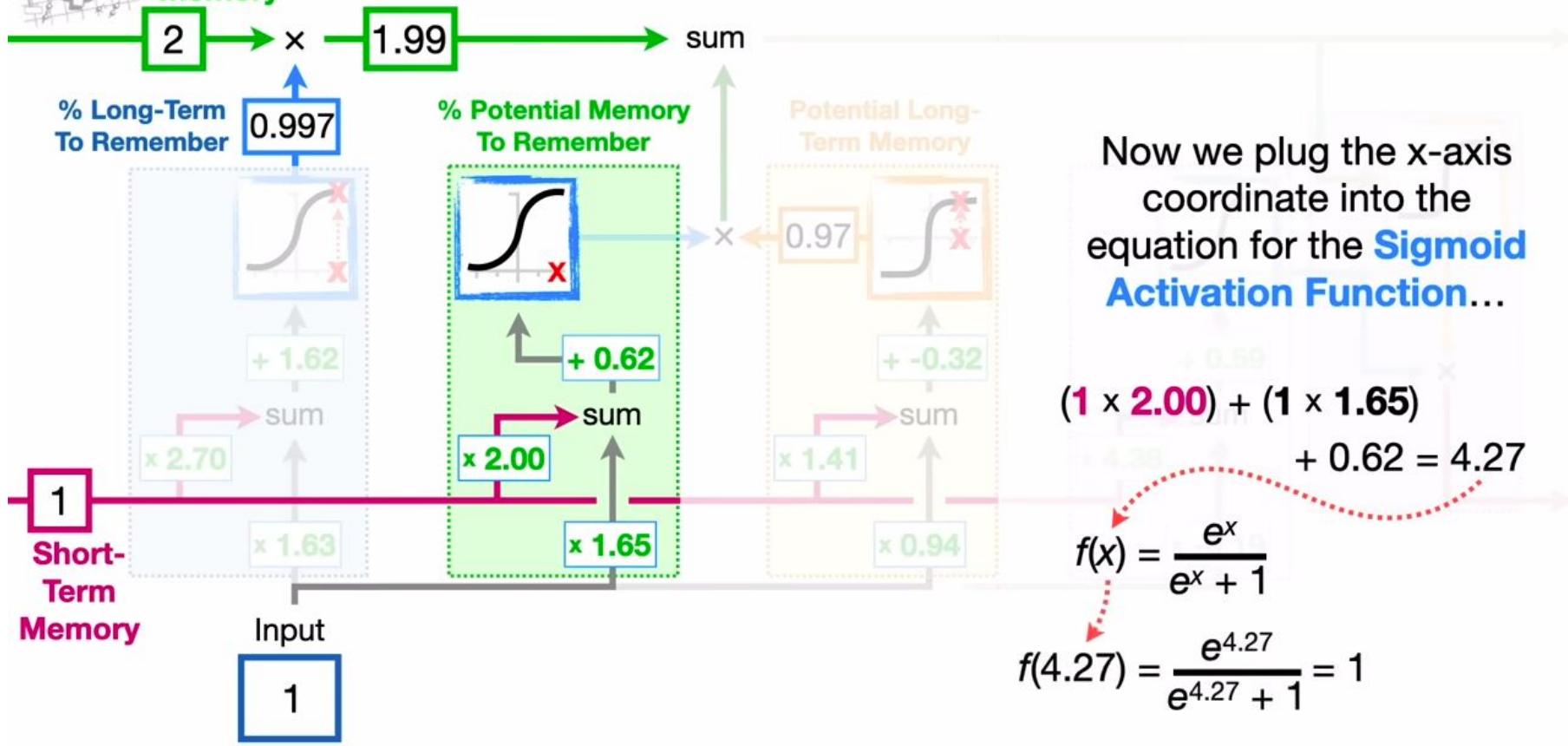
$$(1 \times 2.00) + (1 \times 1.65)$$

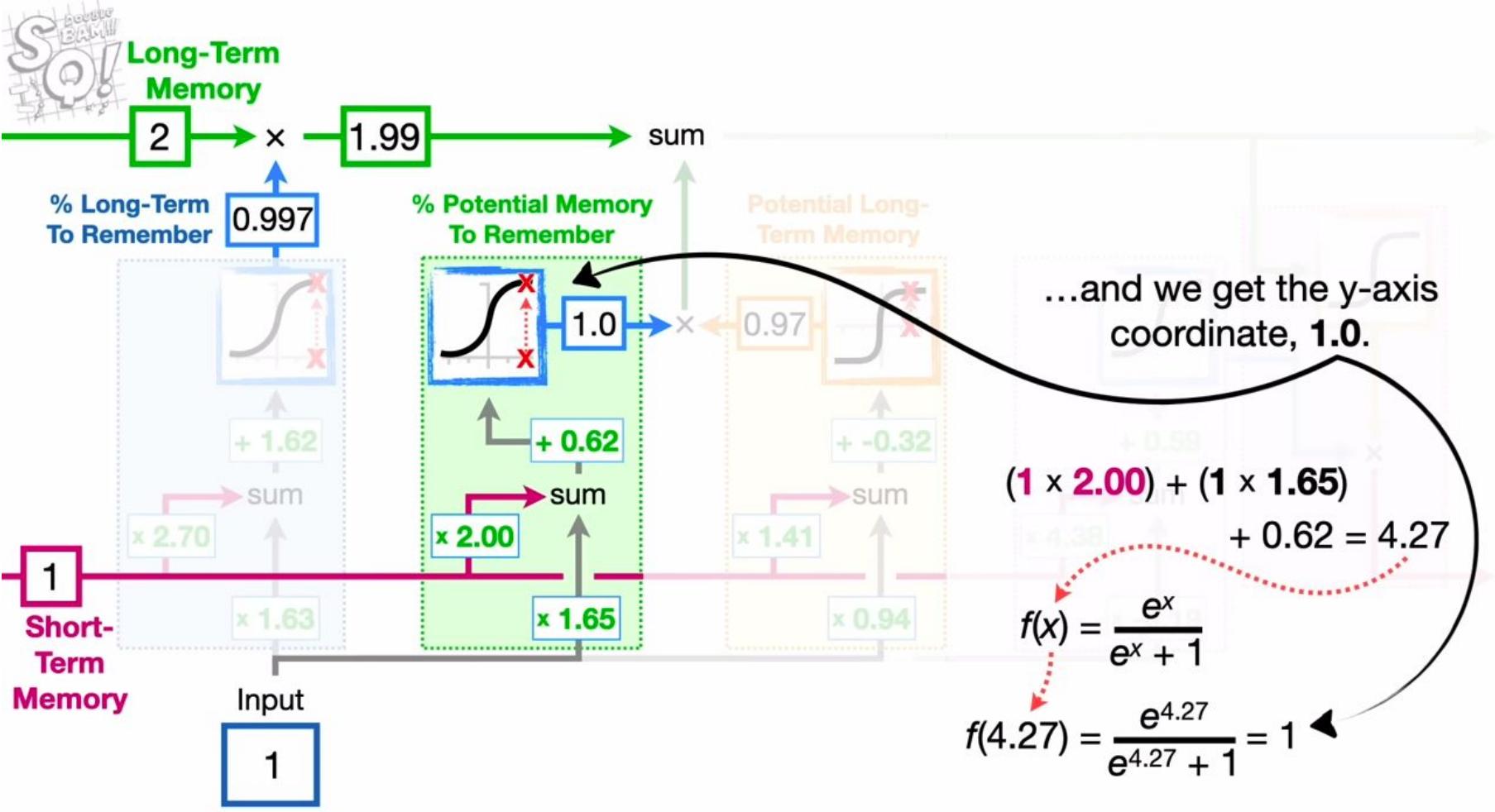
$$+ 0.62 =$$





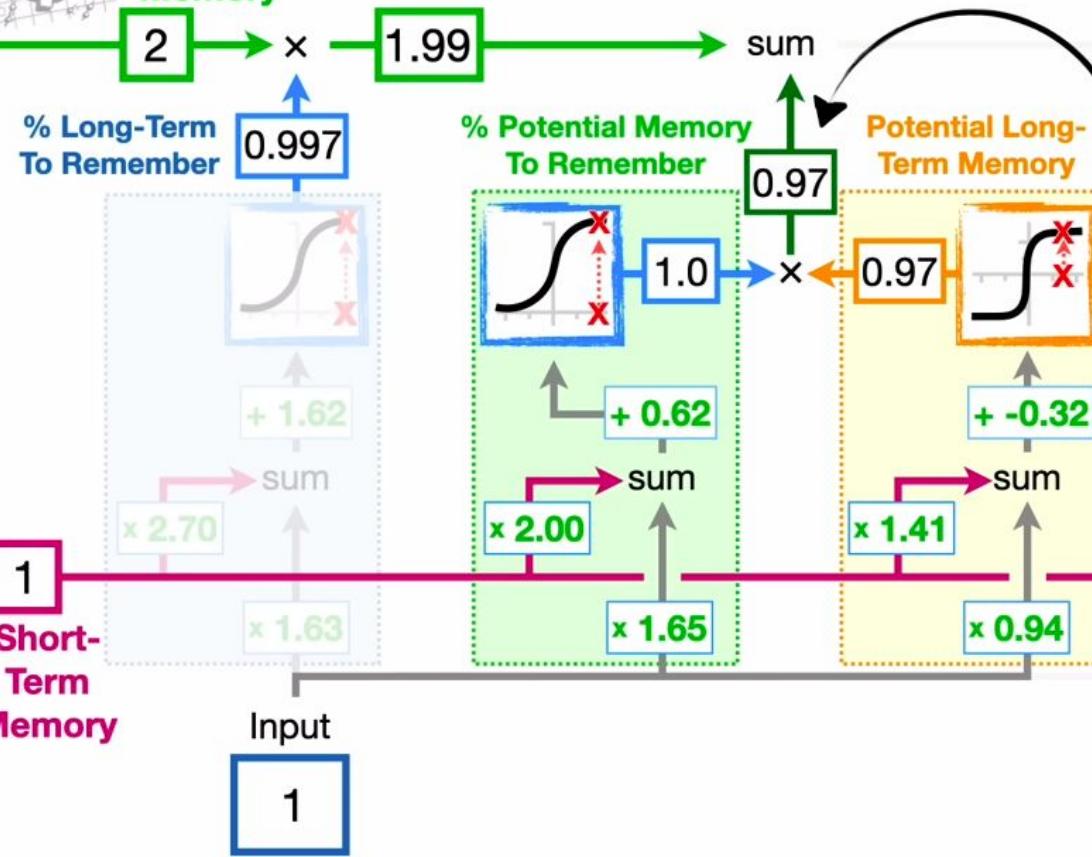








Long-Term Memory



And that means the entire **Potential Long-Term Memory** is retained, because multiplying it by **1** doesn't change it.



Long-Term Memory

2

x

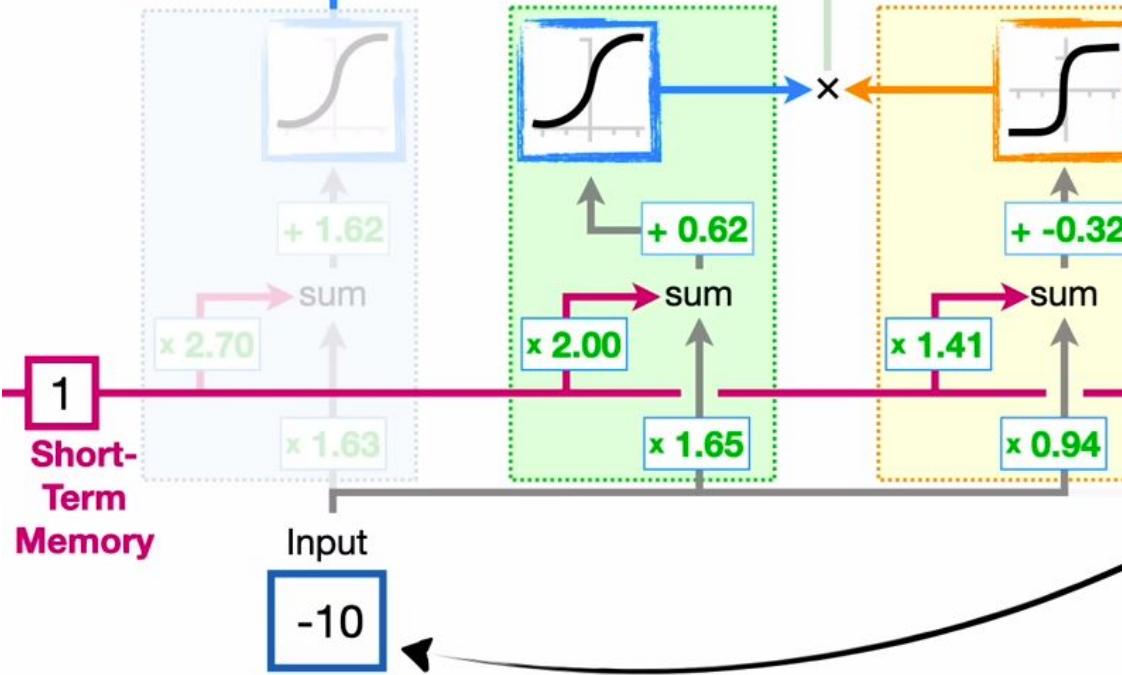
sum

% Long-Term
To Remember

% Potential Memory
To Remember

Potential Long-
Term Memory

NOTE: If the original
Input value was -10...





Long-Term Memory

2

0

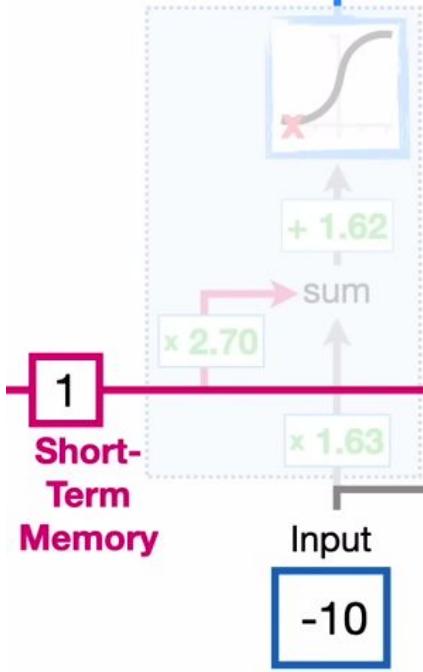
sum

% Long-Term
To Remember

0.0

% Potential Memory
To RememberPotential Long-
Term Memory

...then the percentage of
the **Potential Memory** to
remember would be 0...

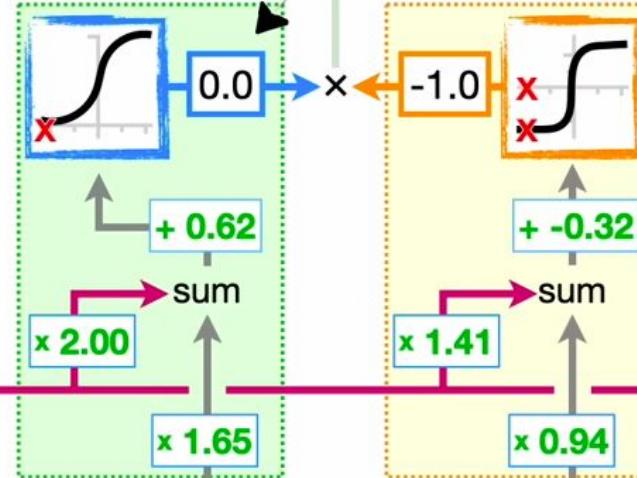


Short-Term Memory

1

Input

-10





Long-Term Memory

2

0

sum

% Long-Term
To Remember

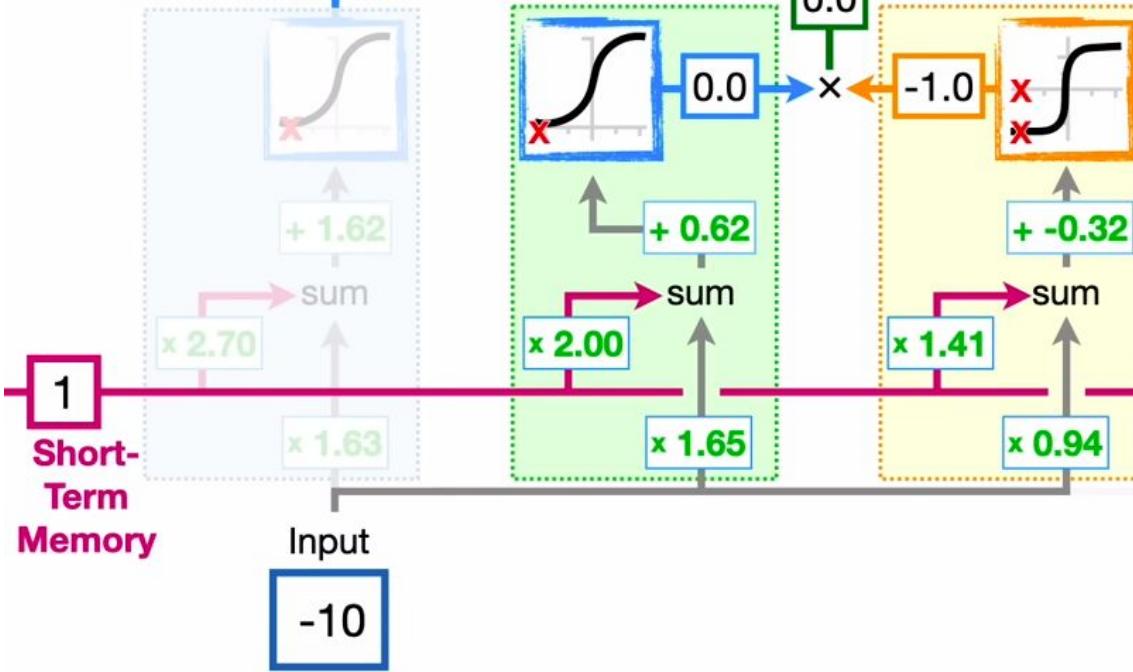
0.0

% Potential Memory
To Remember

0.0

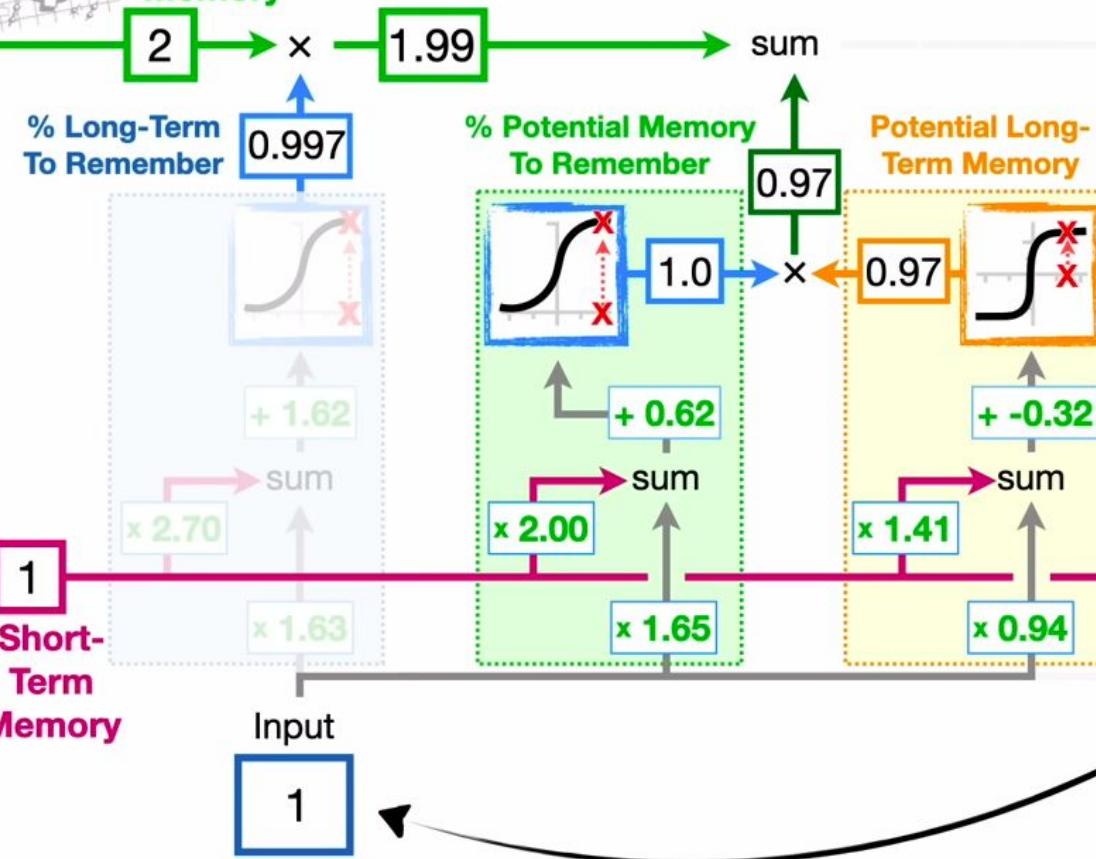
Potential Long-
Term Memory

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

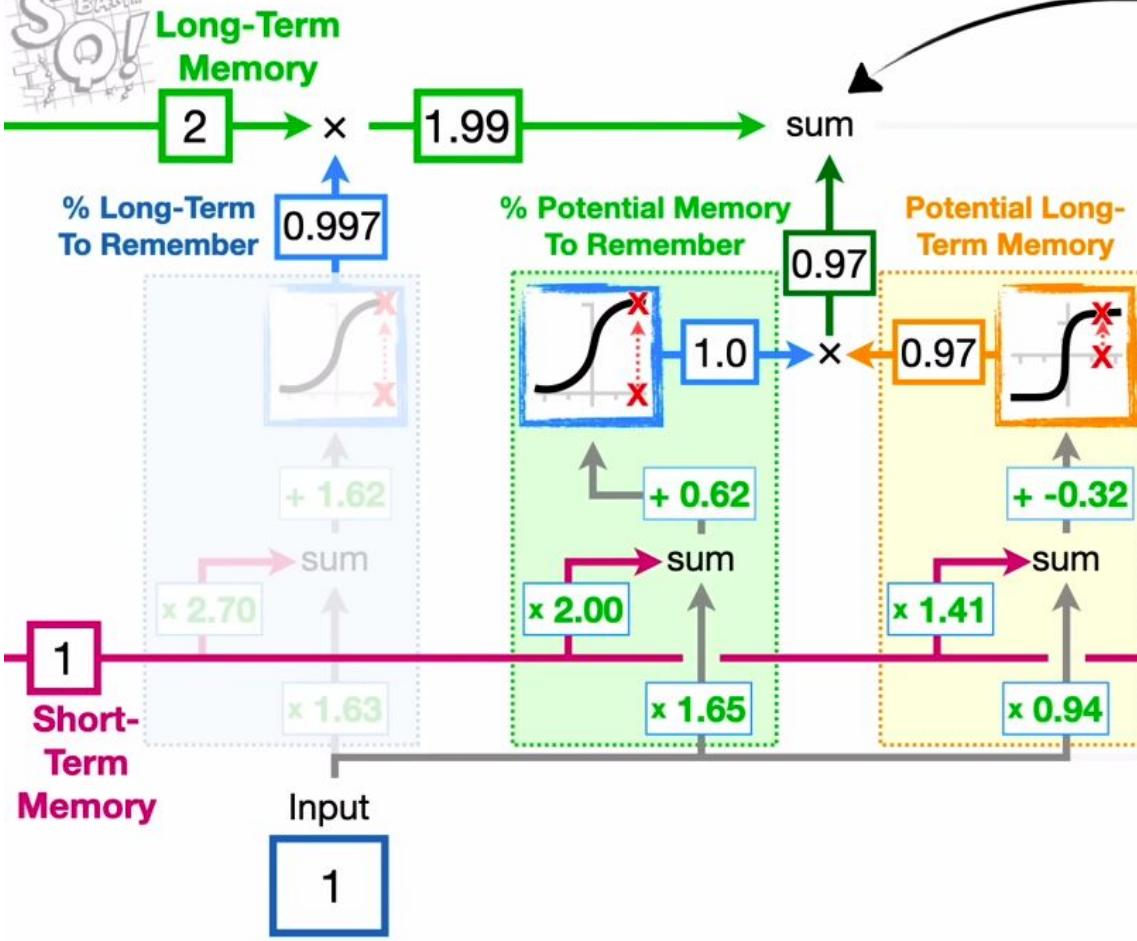




Long-Term Memory



Now, going back to
when the original **Input**
value was 1...



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



Long-Term
Memory

2

\times
1.99

New Long-Term
Memory

2.96

% Long-Term
To Remember

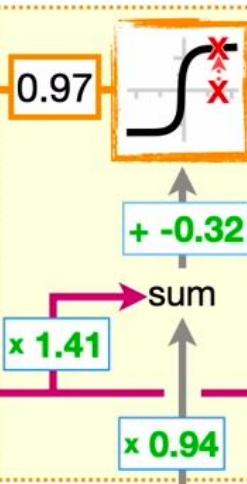
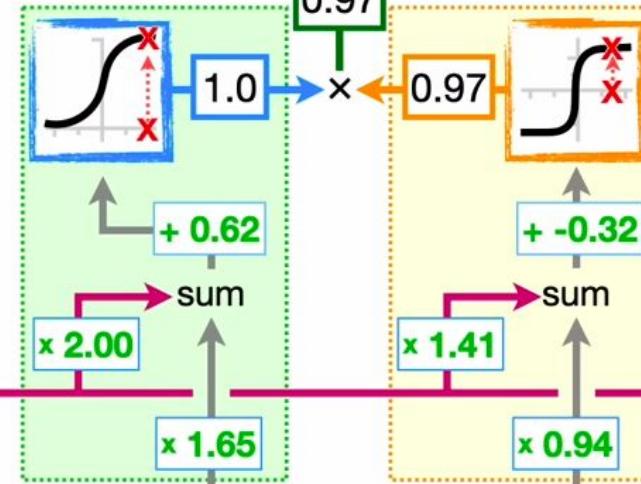
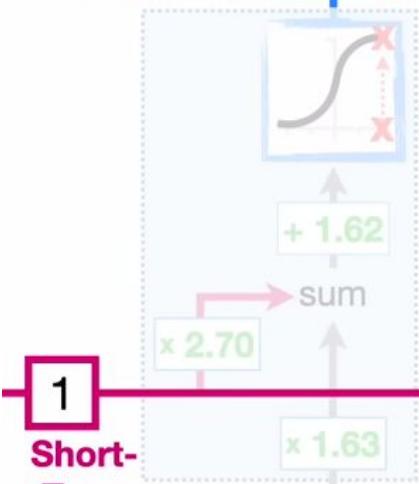
0.997

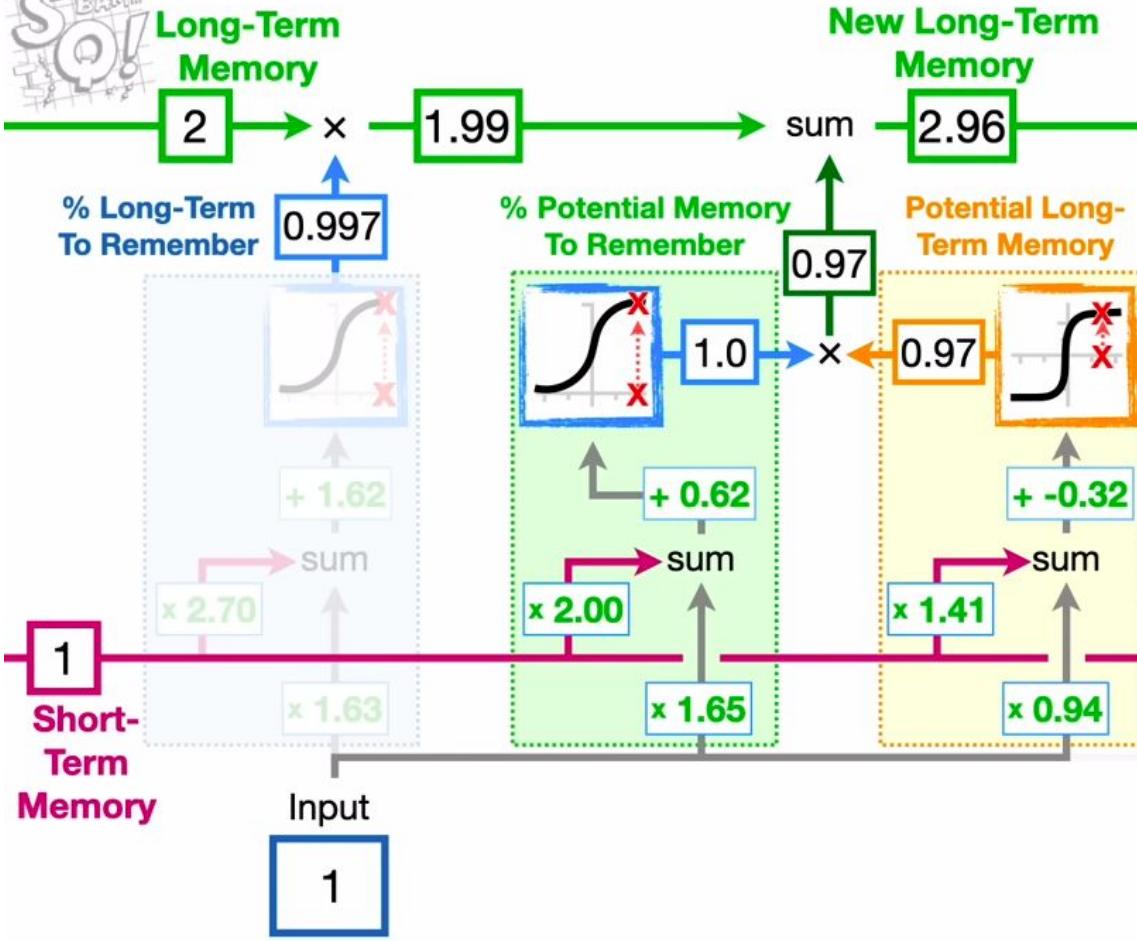
% Potential Memory
To Remember

Potential Long-
Term Memory

0.97

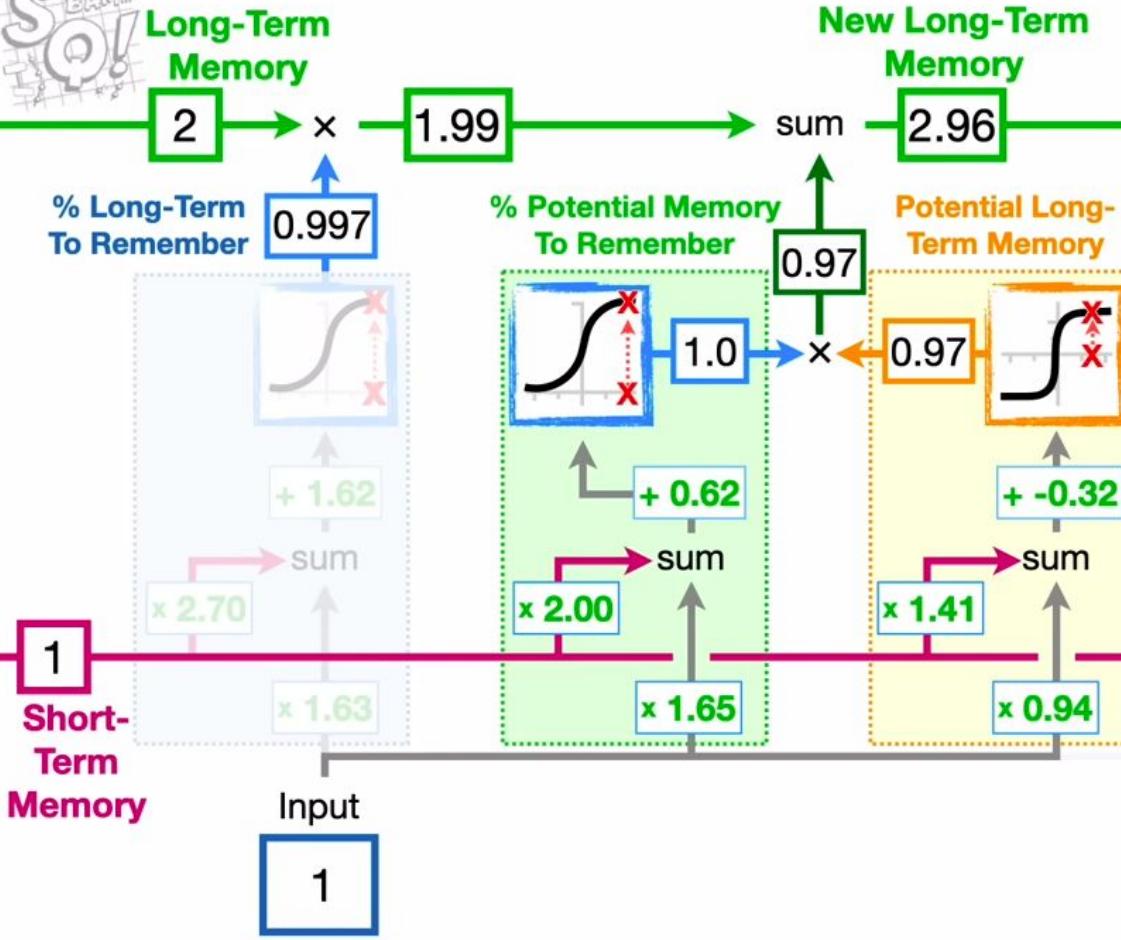
...and we get a new
Long-Term Memory,
2.96.





TERMINOLOGY ALERT!!!

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



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**.



Long-Term
Memory

2

\times
1.99

New Long-Term
Memory

2.96

% Long-Term
To Remember

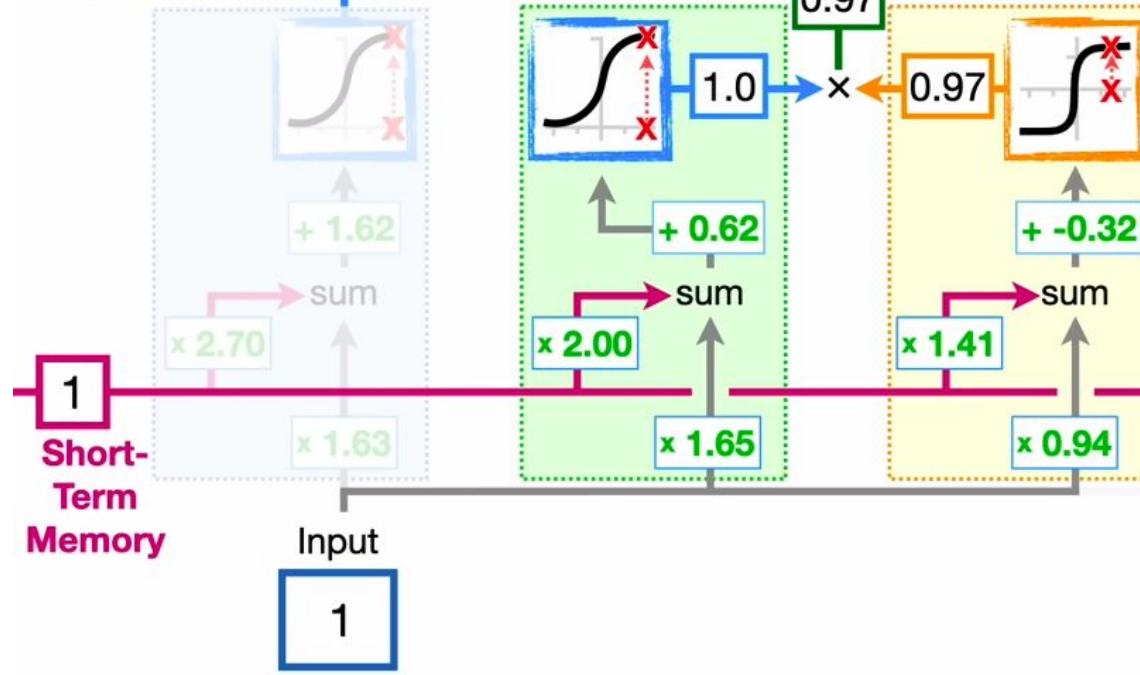
0.997

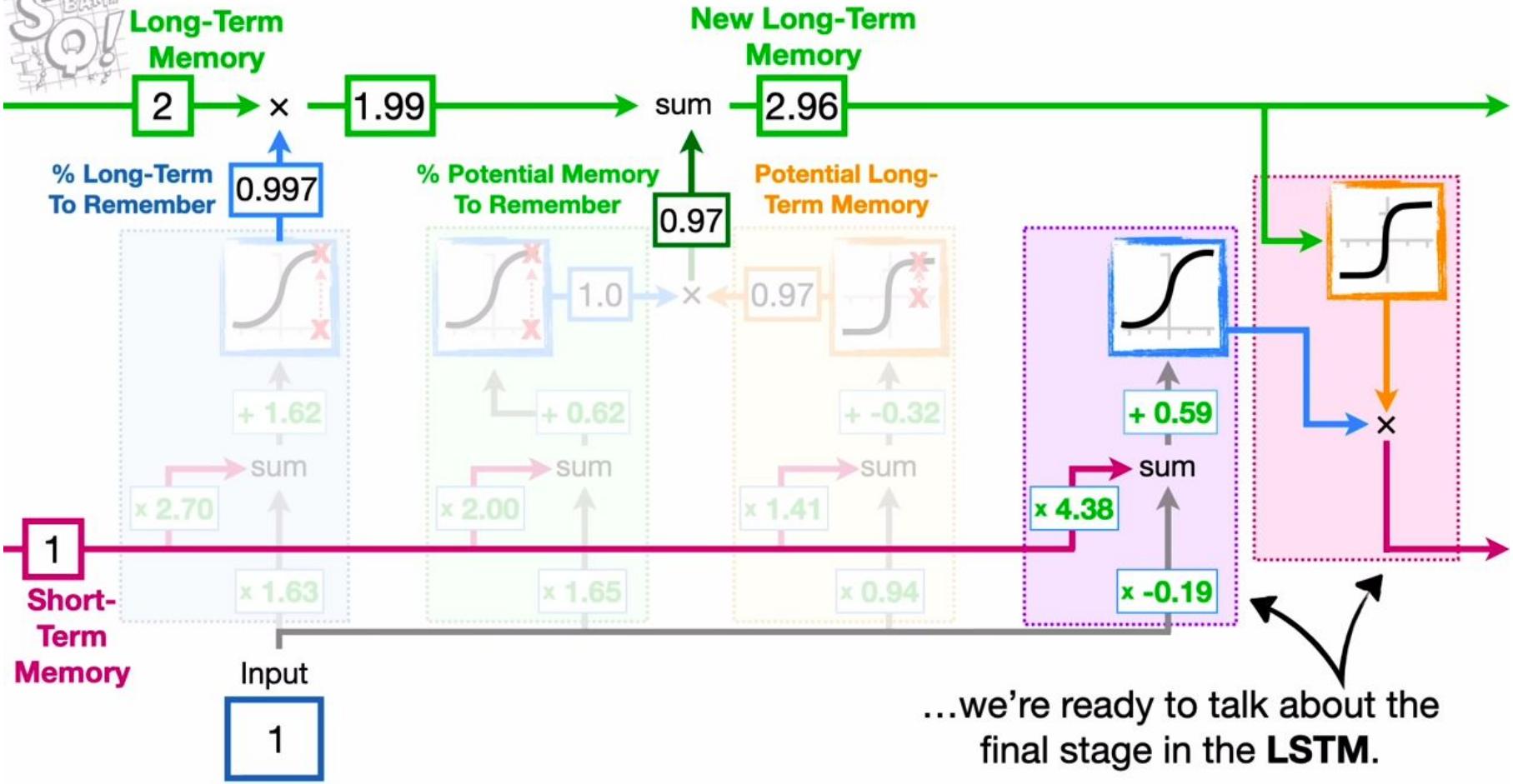
% Potential Memory
To Remember

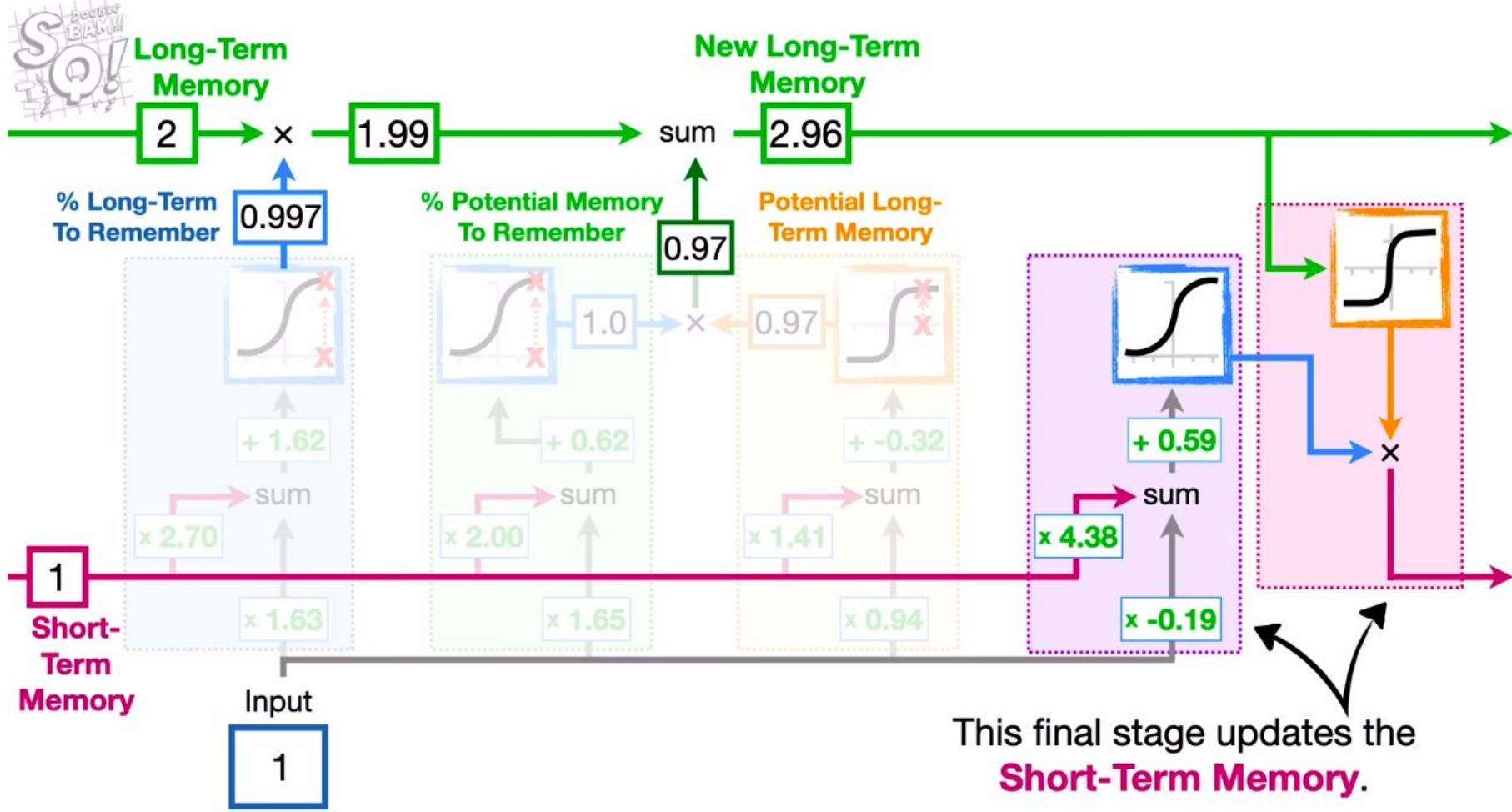
Potential Long-
Term Memory

0.97

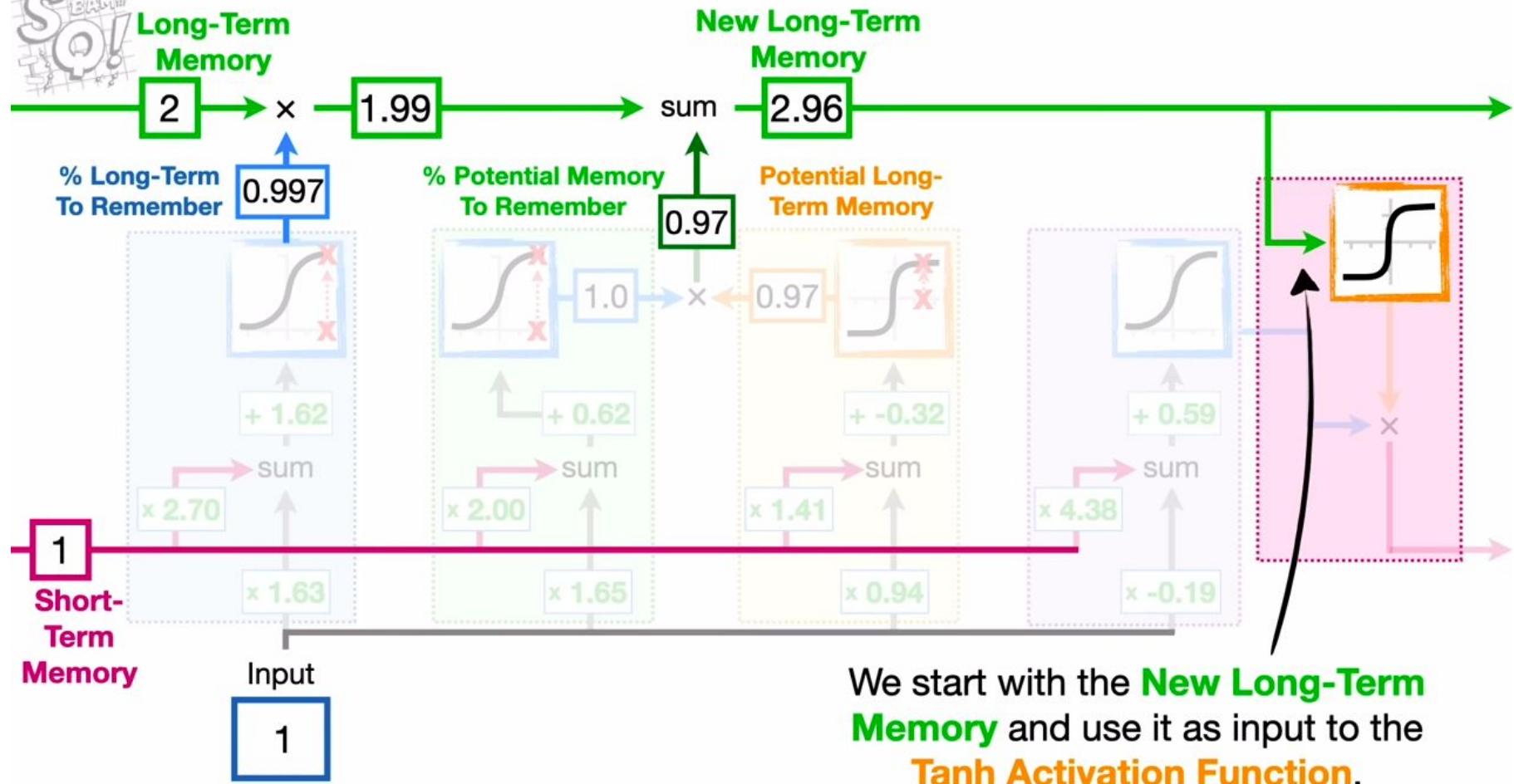
Now that we have a New
Long-Term Memory...



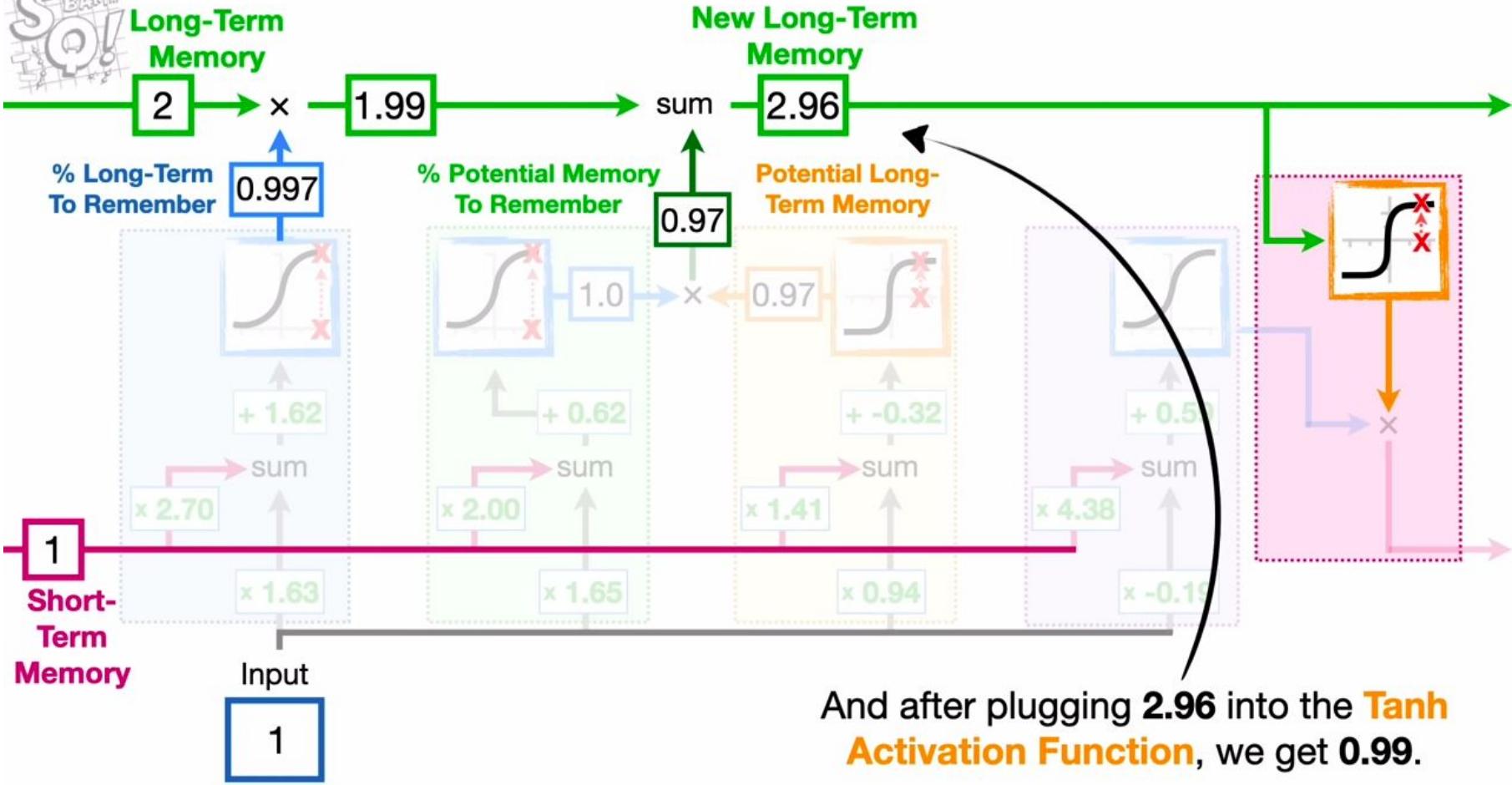




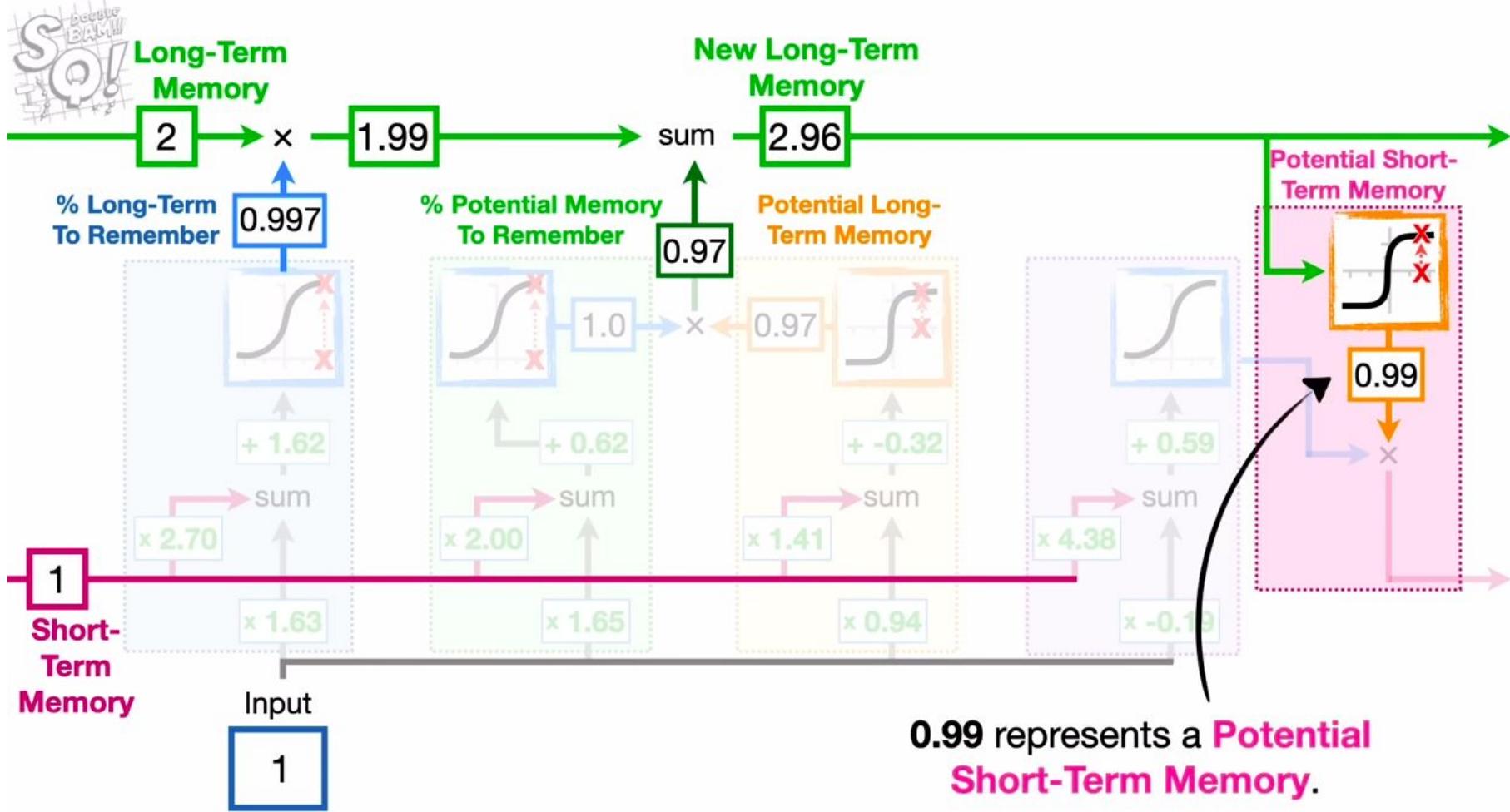
This final stage updates the **Short-Term Memory**.



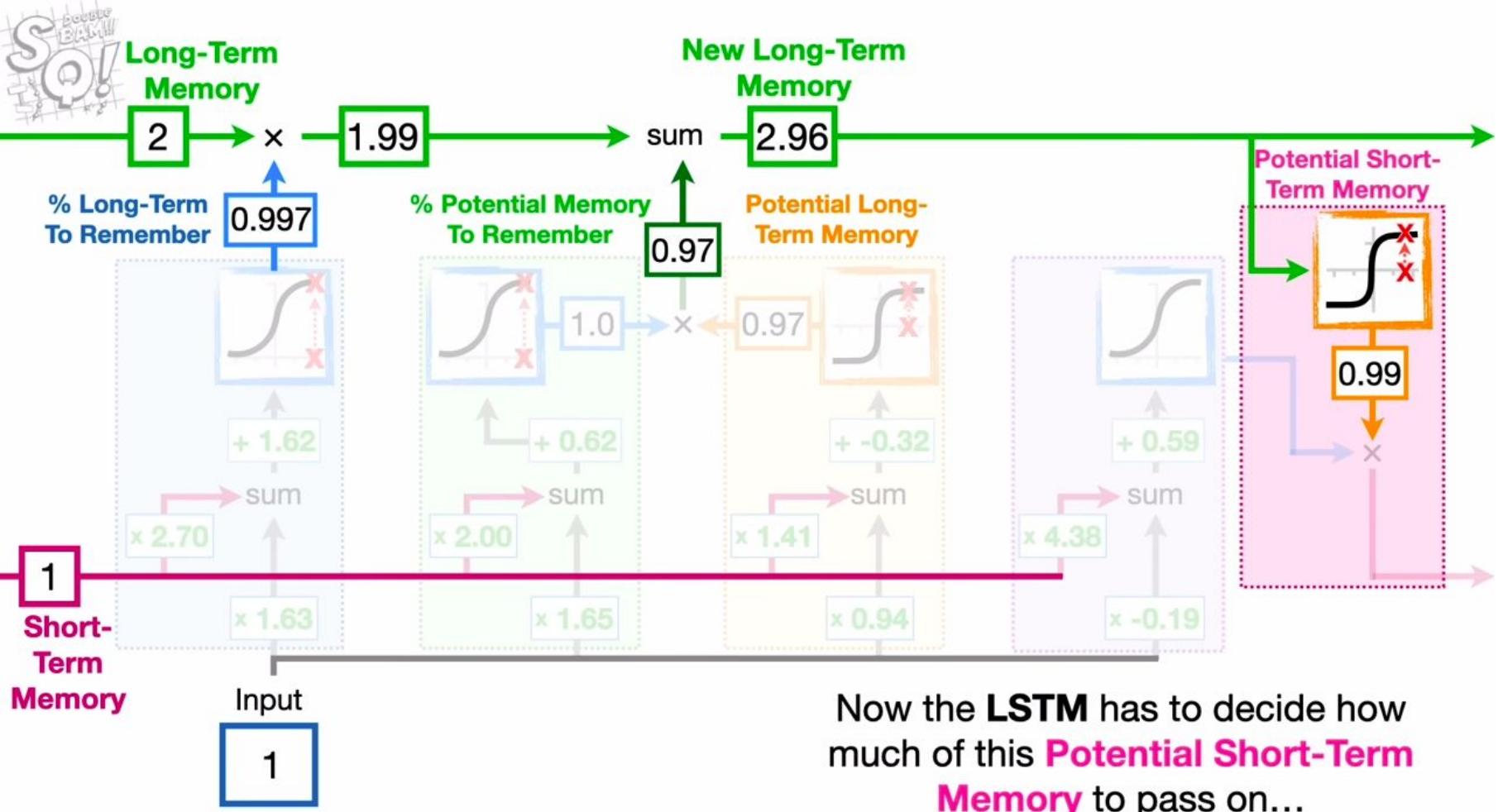
We start with the **New Long-Term Memory** and use it as input to the **Tanh Activation Function**.



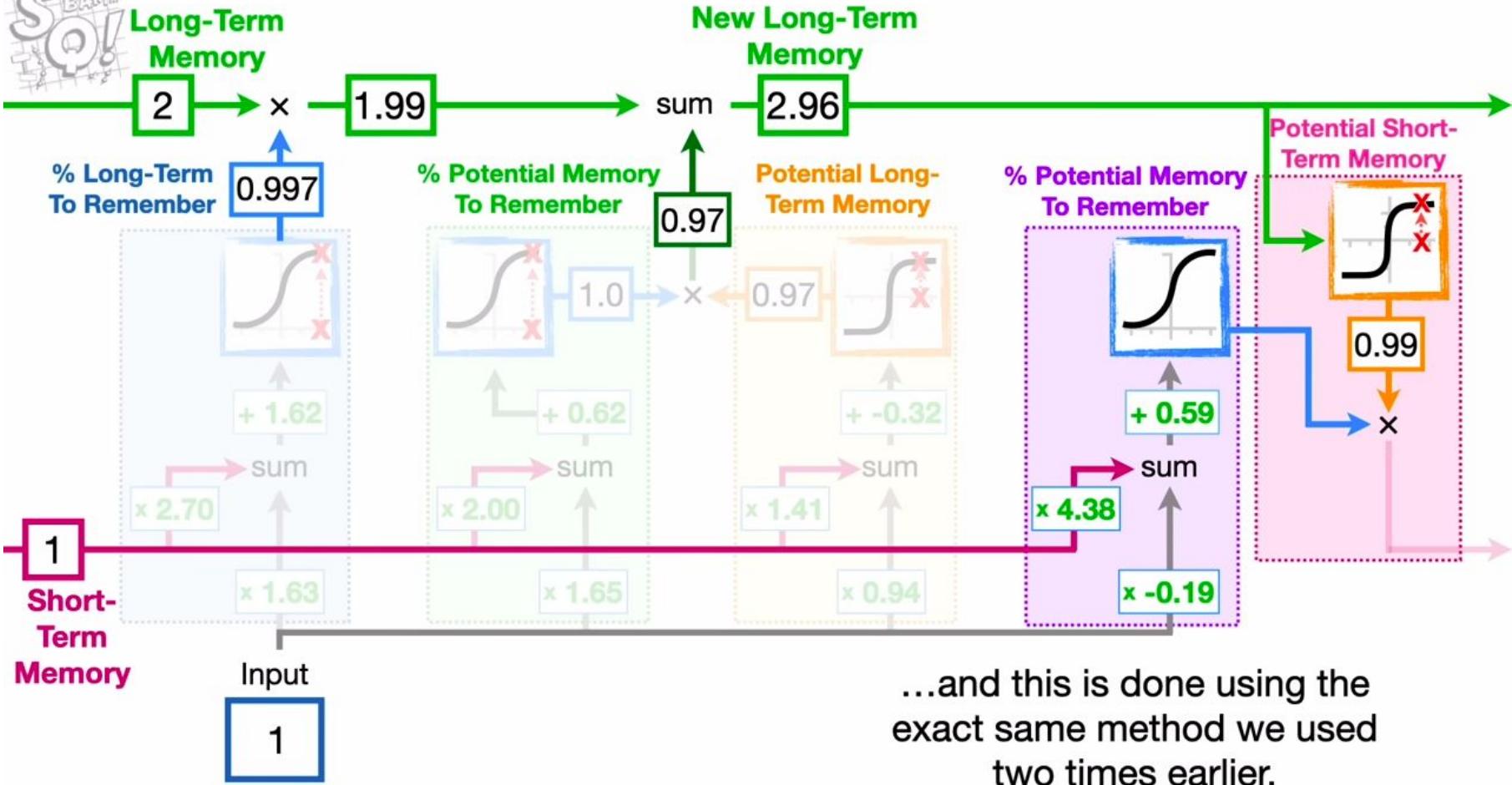
And after plugging 2.96 into the **Tanh Activation Function**, we get 0.99.

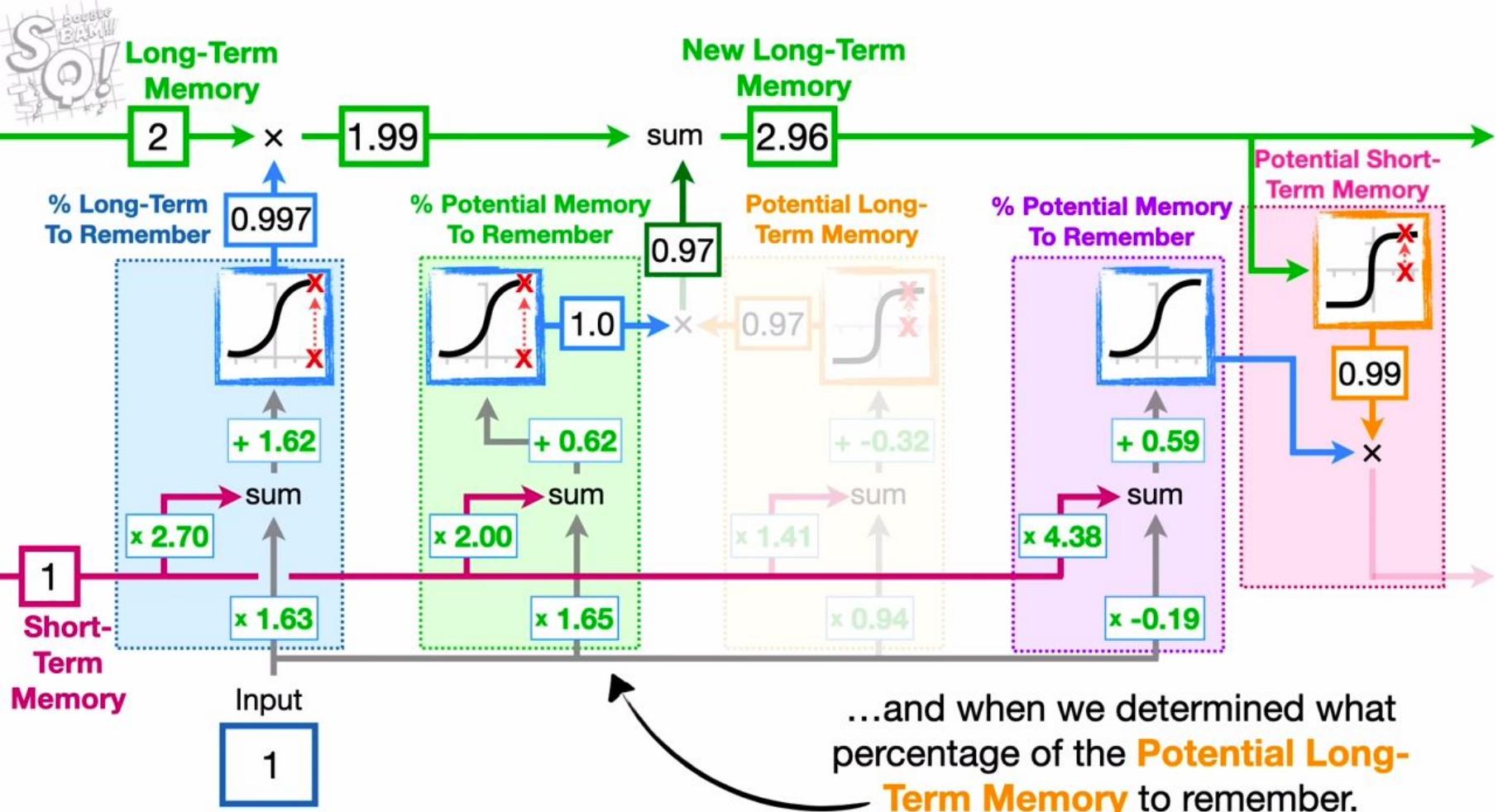


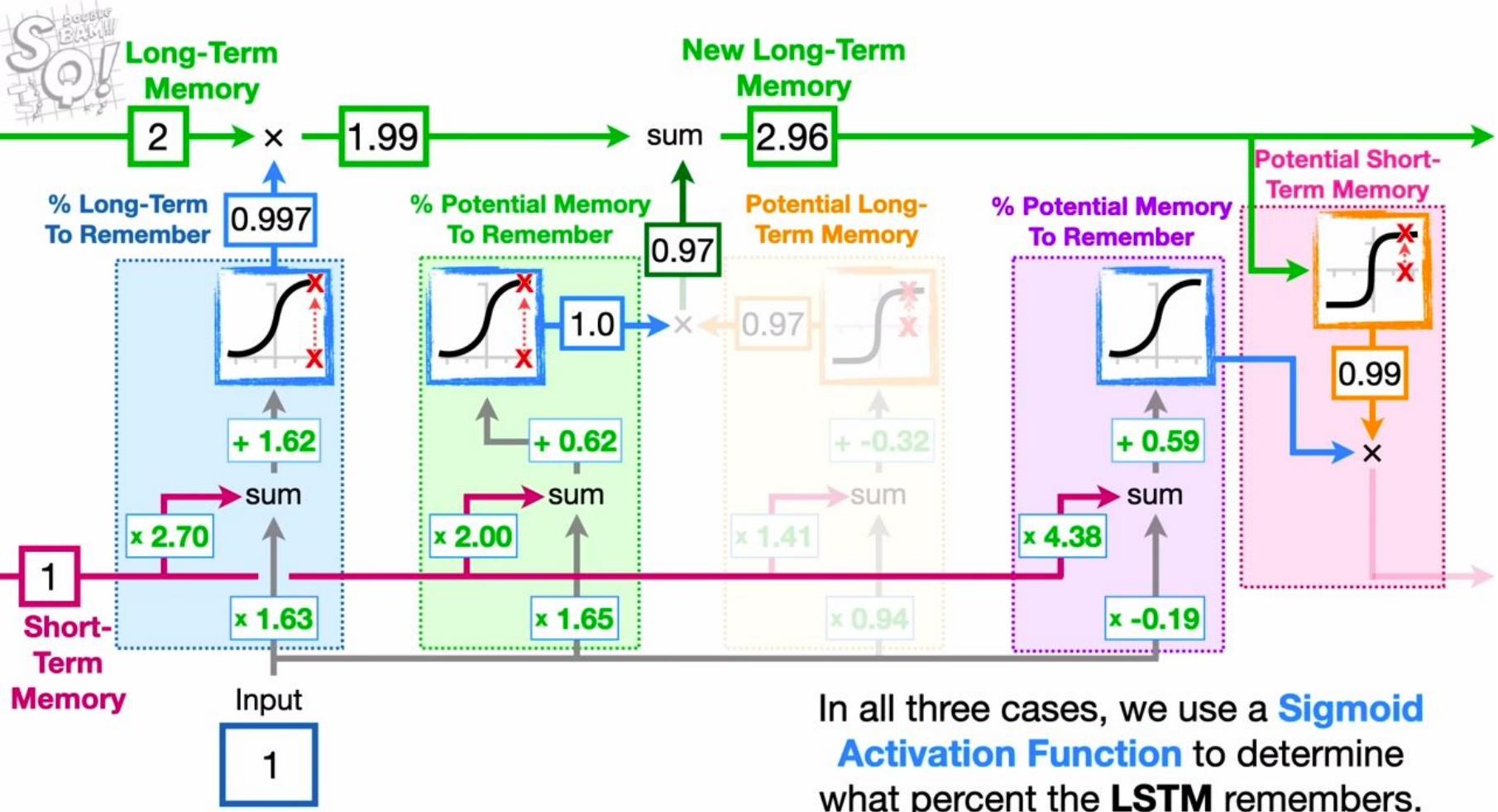
0.99 represents a Potential Short-Term Memory.



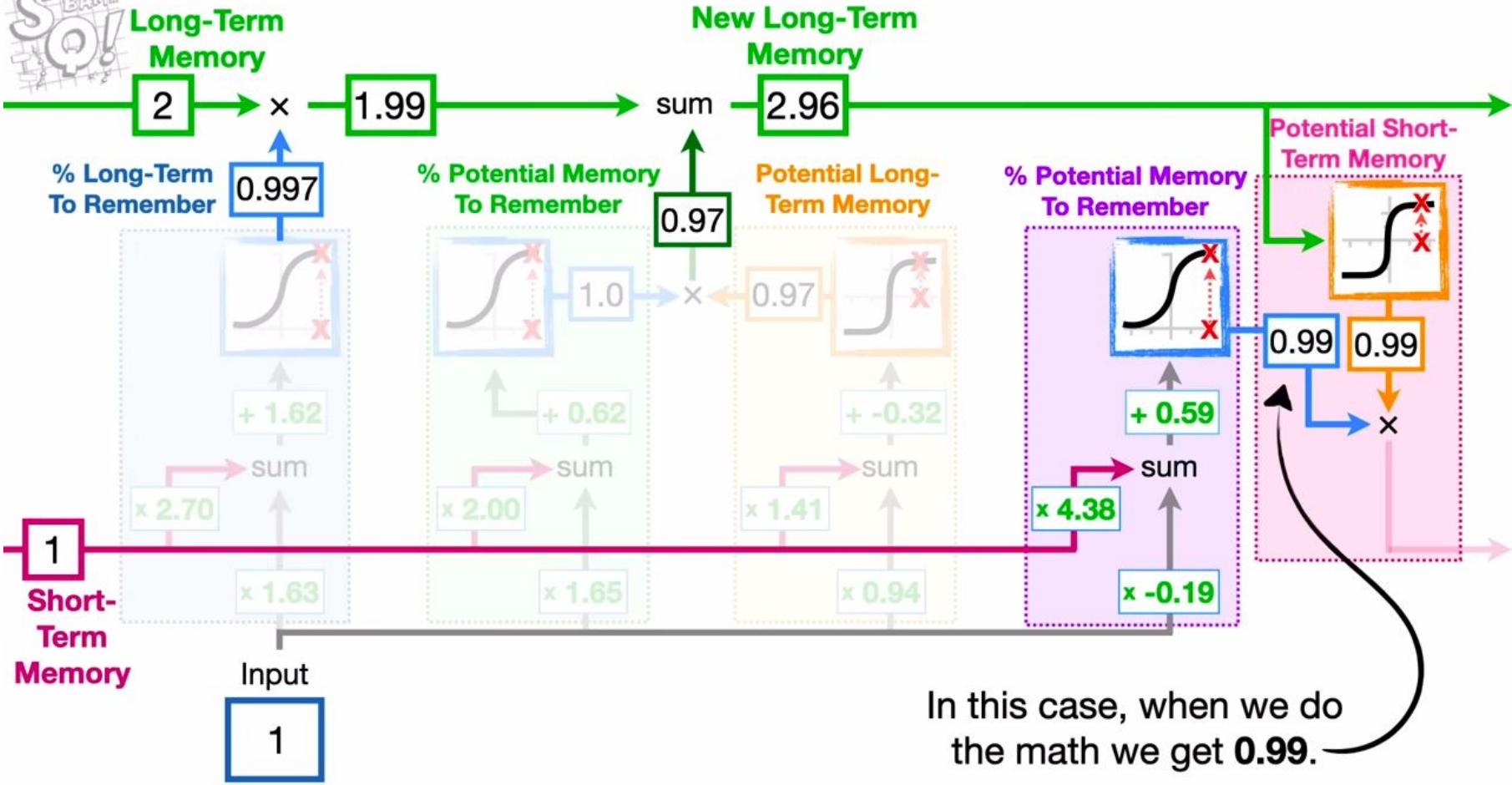
Now the **LSTM** has to decide how much of this **Potential Short-Term Memory** to pass on...

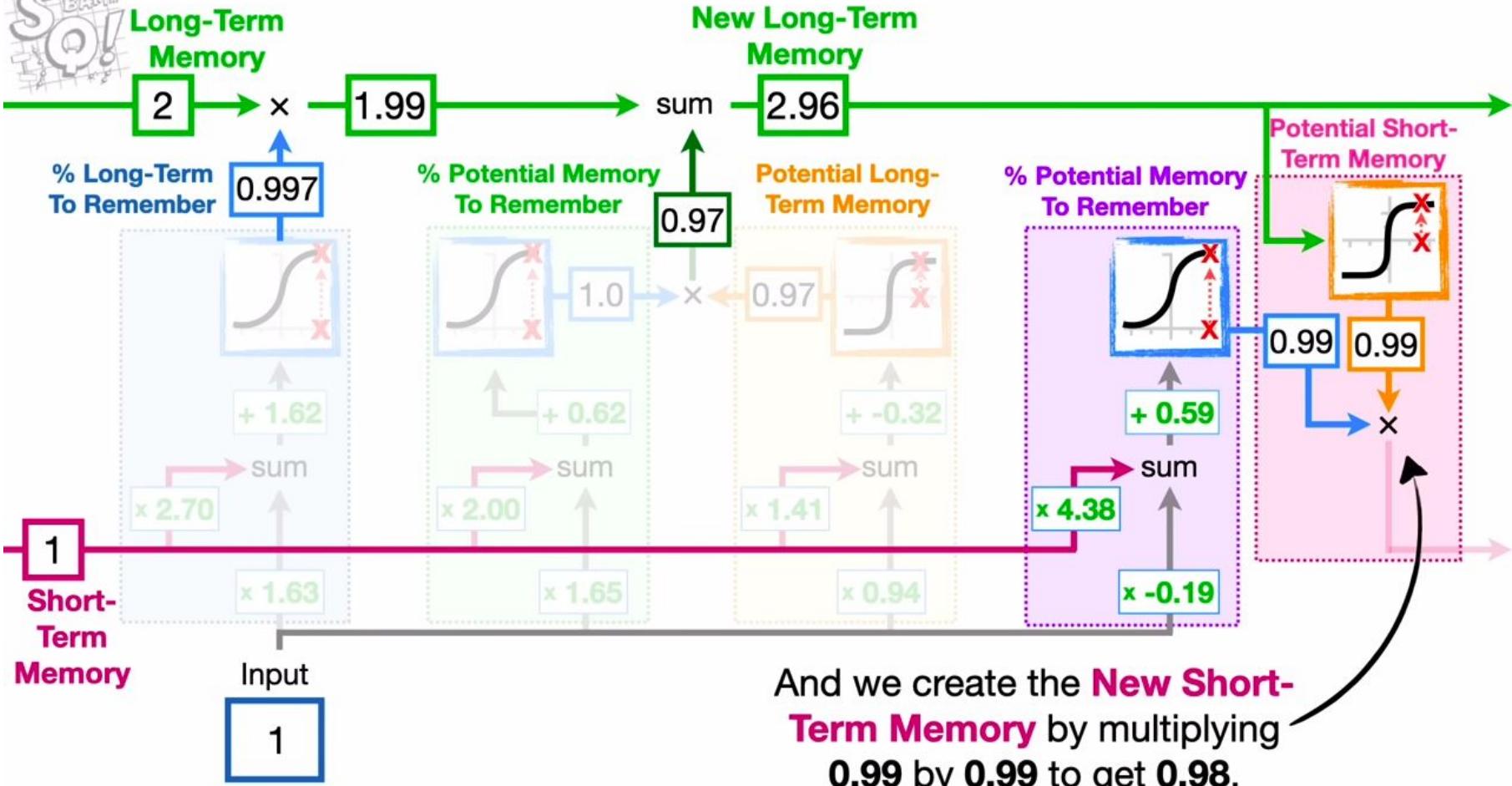




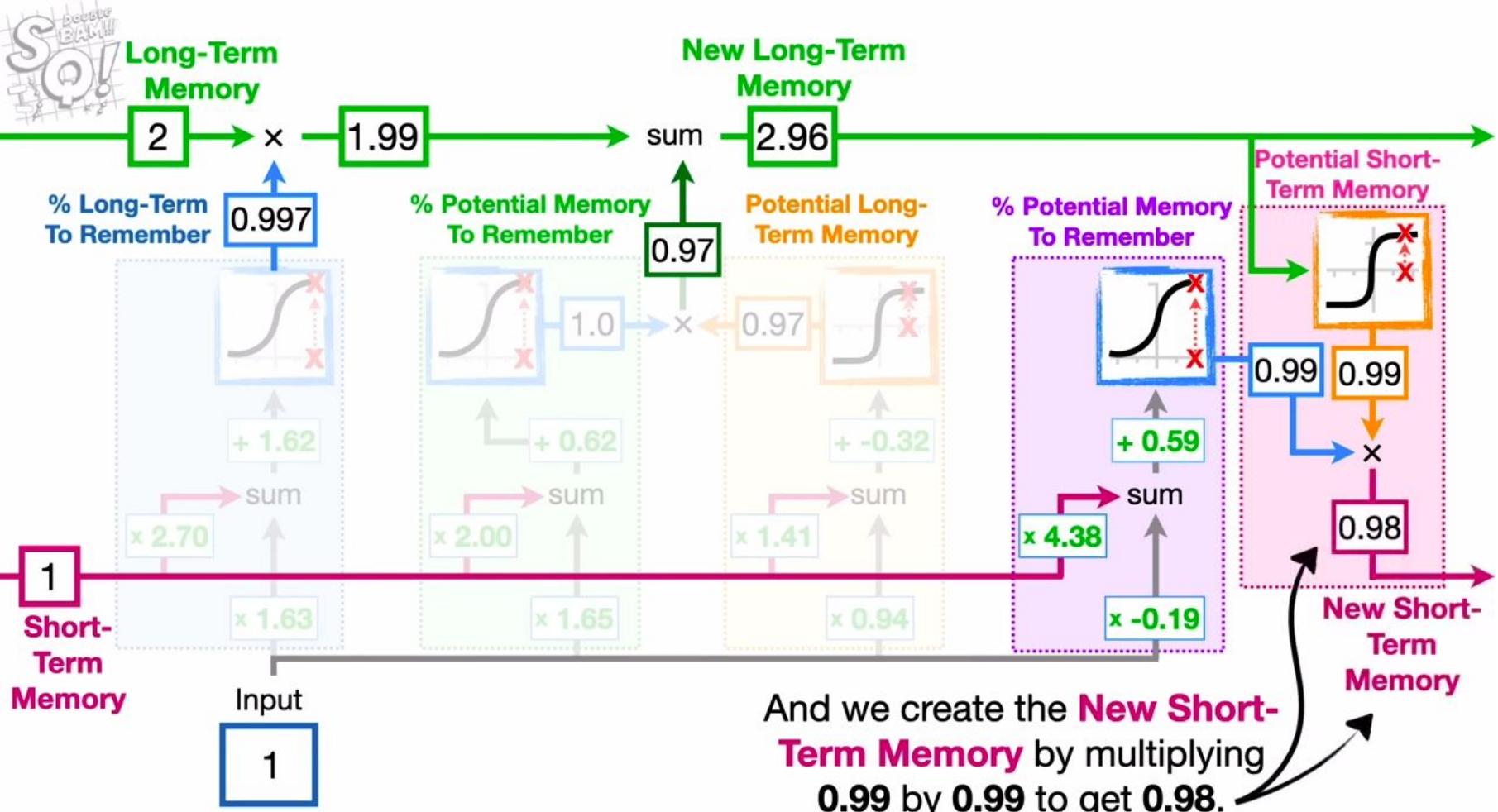


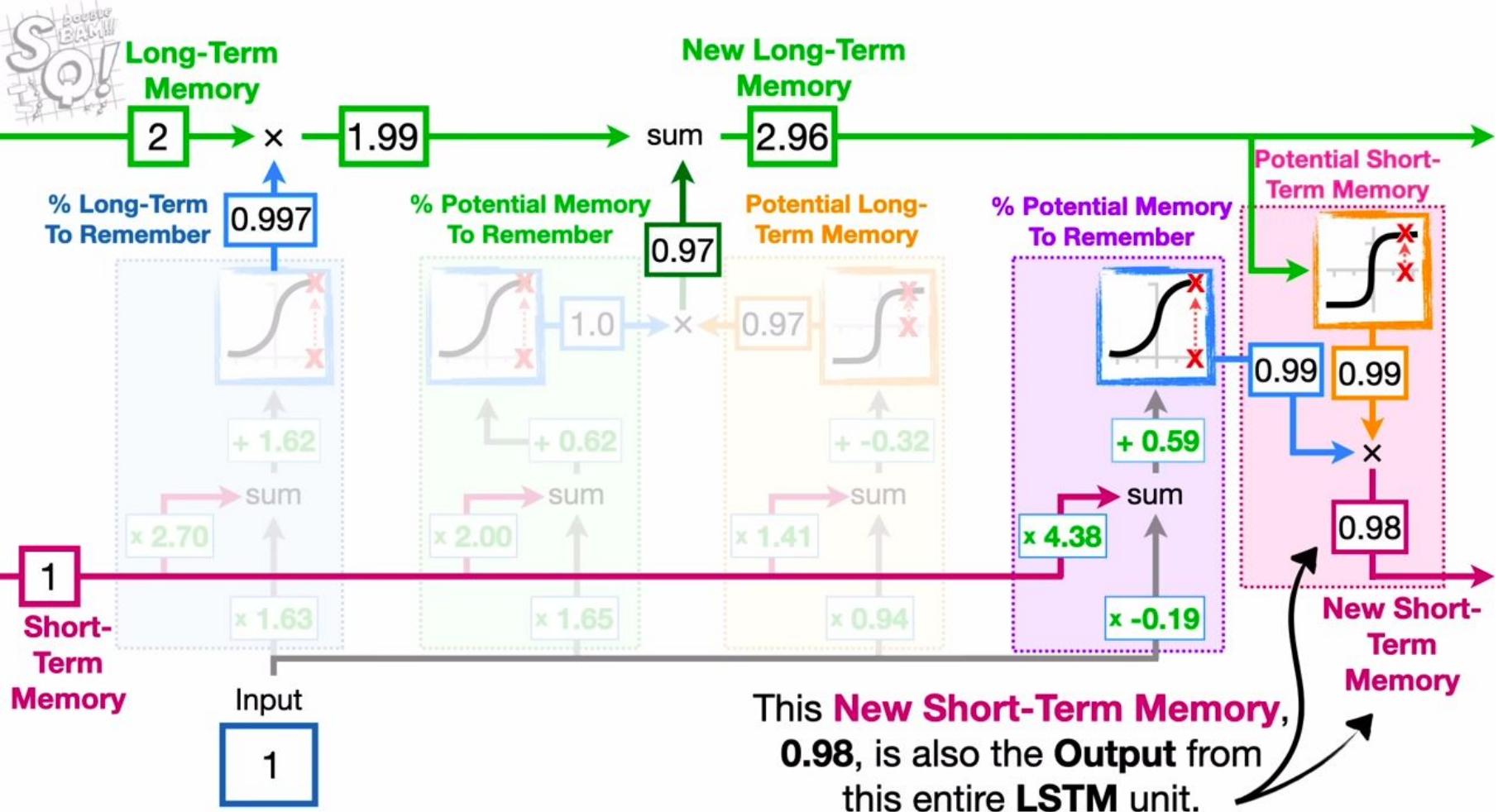
In all three cases, we use a **Sigmoid Activation Function** to determine what percent the **LSTM** remembers.

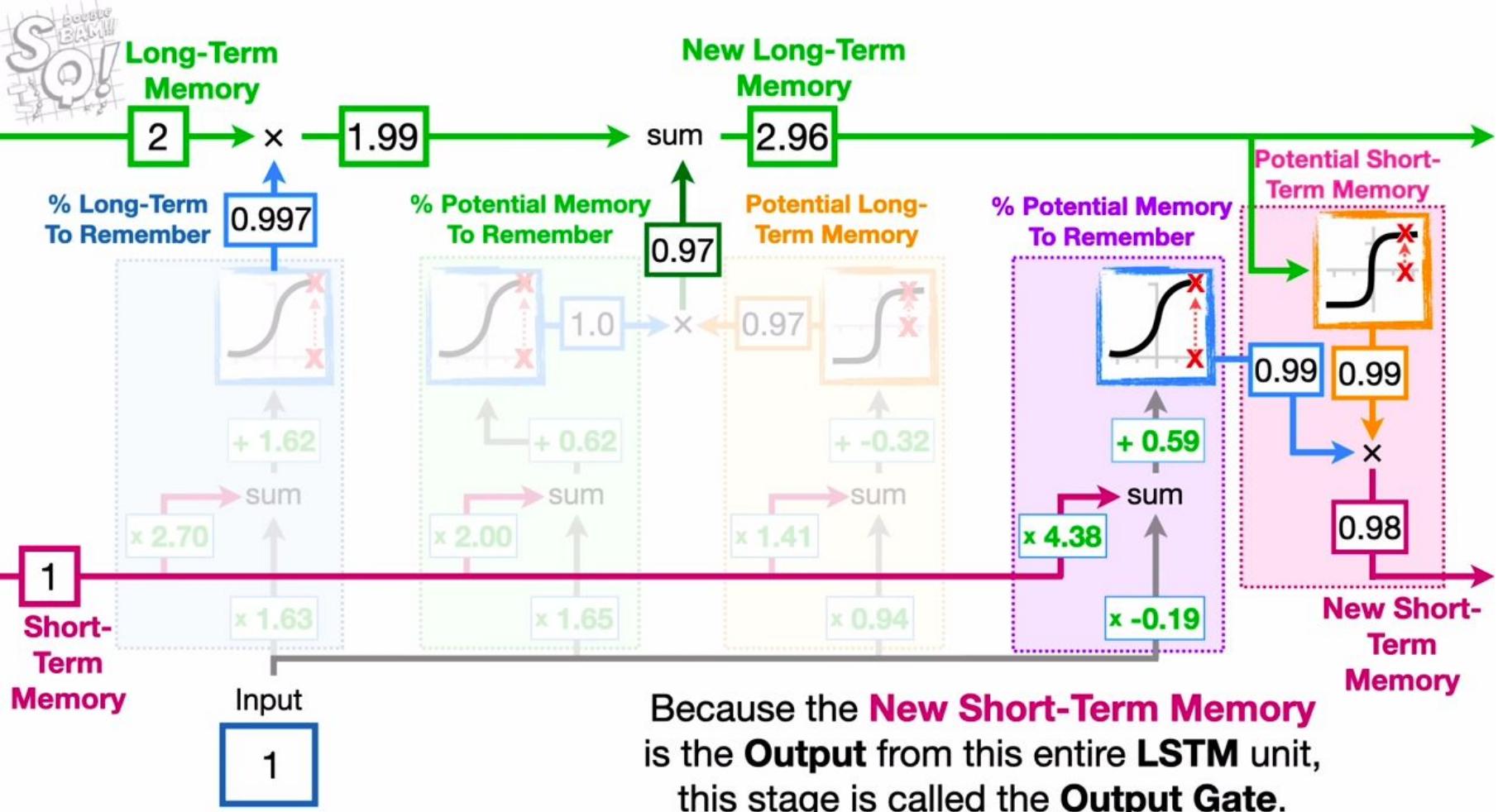


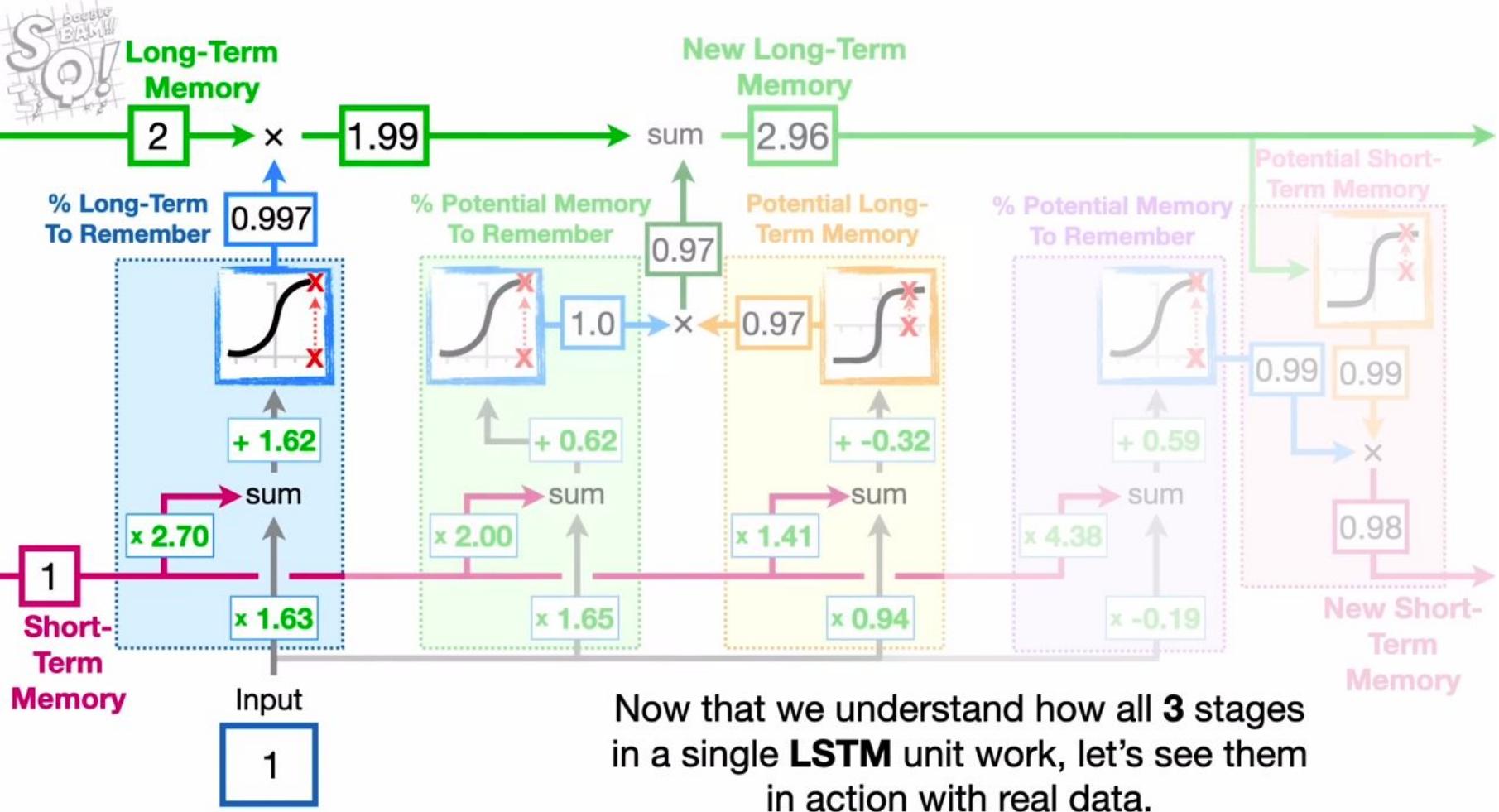


And we create the **New Short-Term Memory** by multiplying 0.99 by 0.99 to get 0.98.

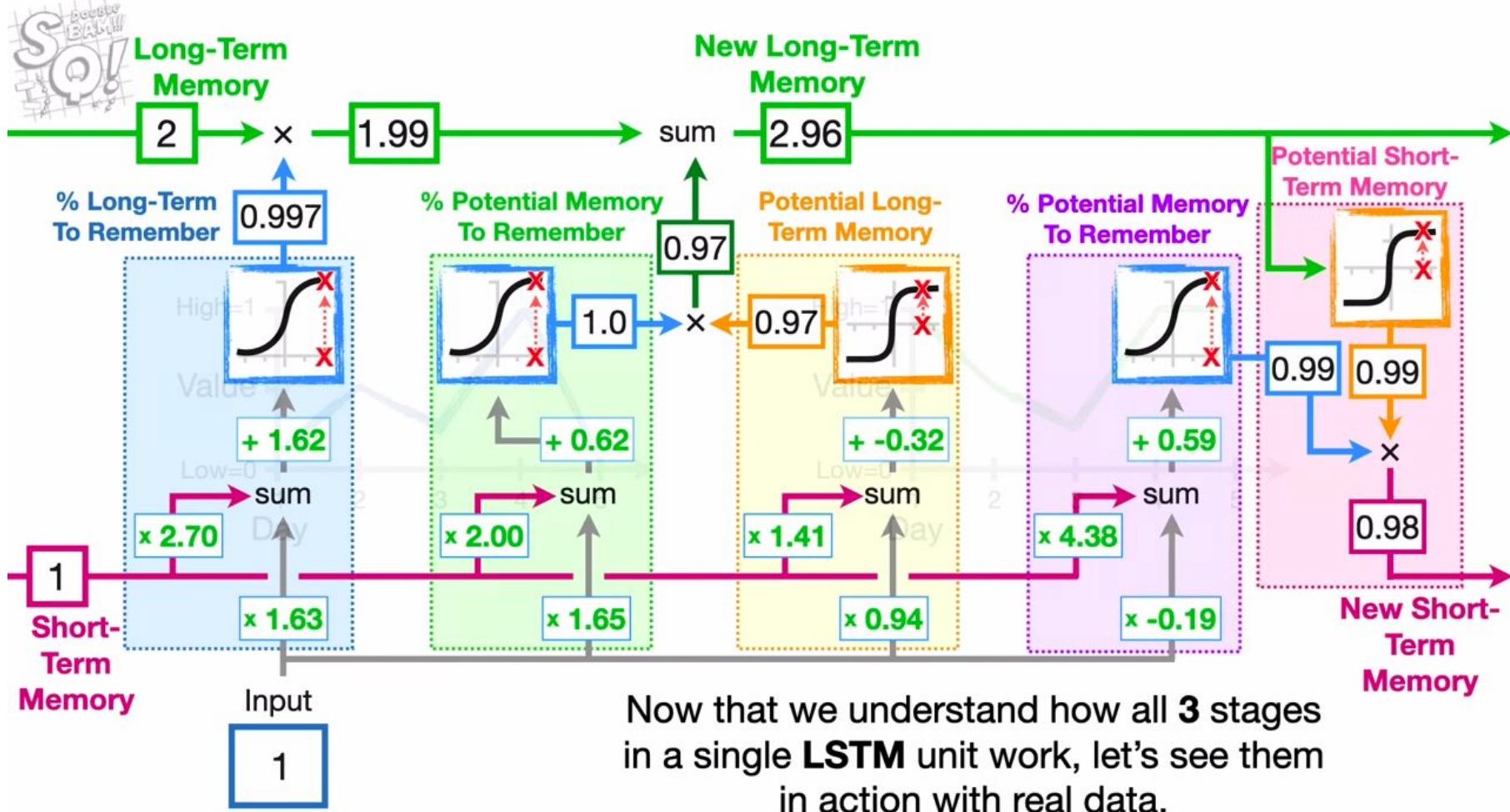








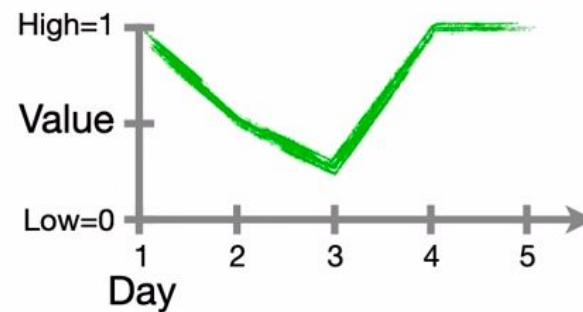
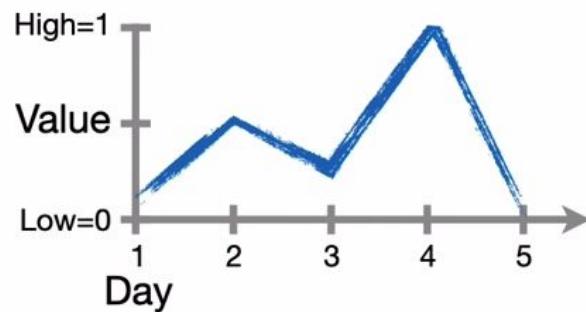
Now that we understand how all **3** stages in a single **LSTM** unit work, let's see them in action with real data.



Now that we understand how all **3** stages in a single **LSTM** unit work, let's see them in action with real data.



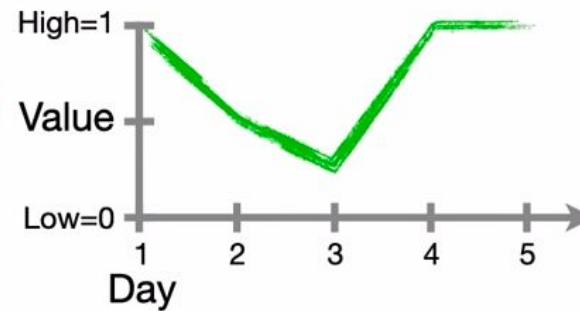
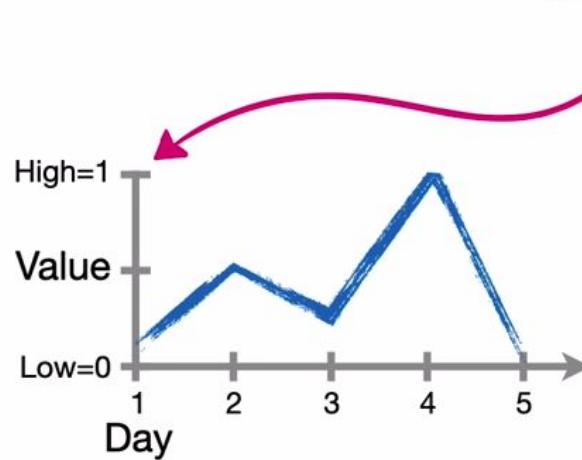
Here we have stock prices for two companies, **Company A** and **Company B**.



Company A =
Company B =



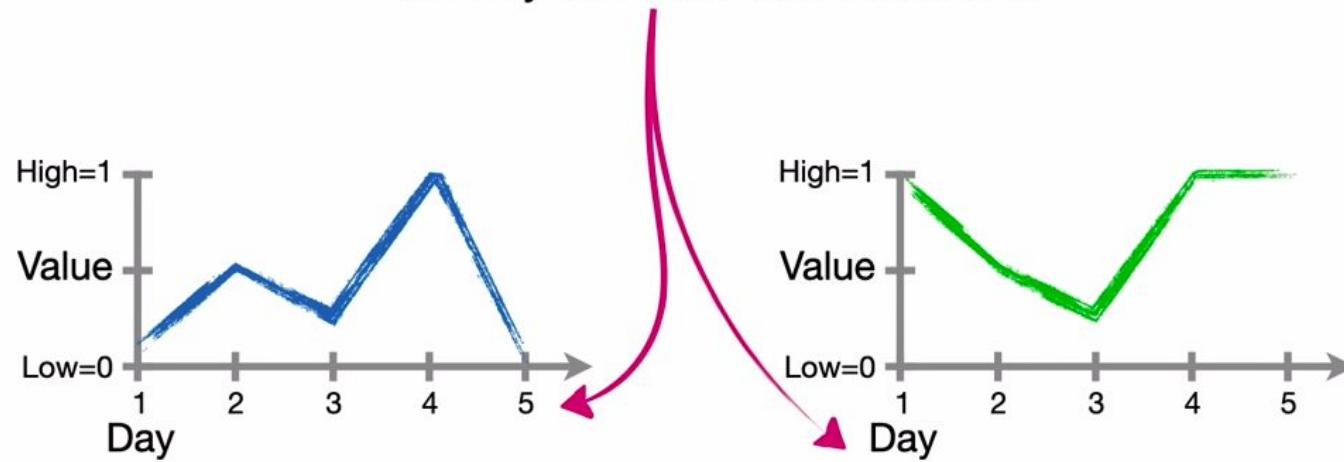
On the y-axis, we have
the stock value...



Company A =
Company B =



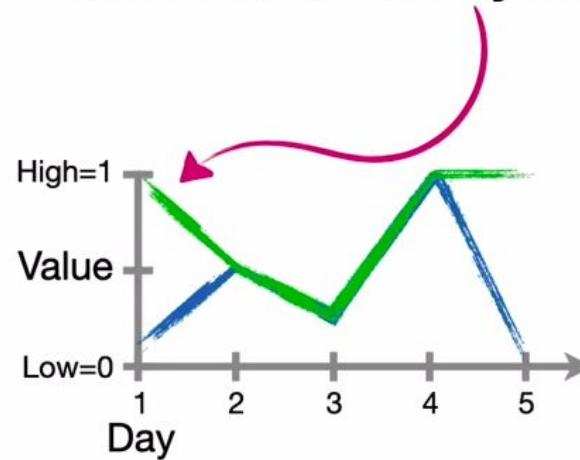
...and on the x-axis, we have
the day the value was recorded.



Company A =
Company B =



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

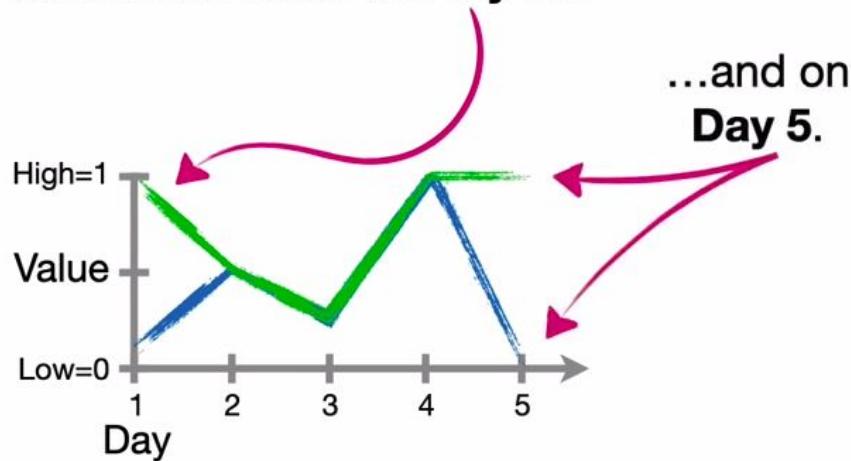


Company A =

Company B =



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

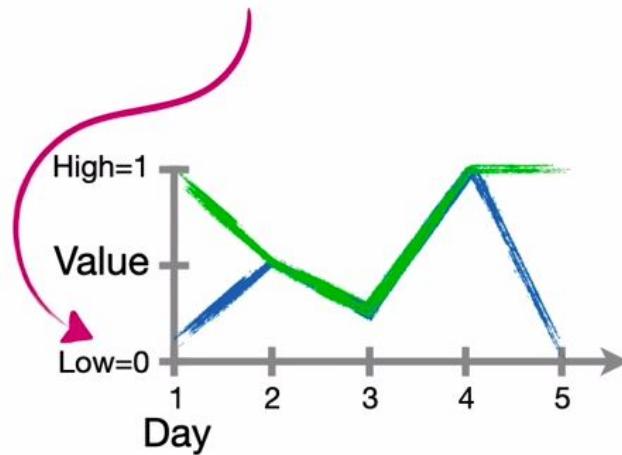


Company A =

Company B =



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

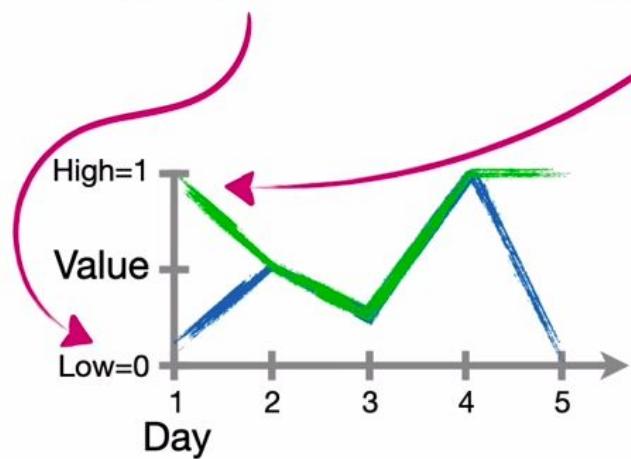


Company A = /

Company B = /



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

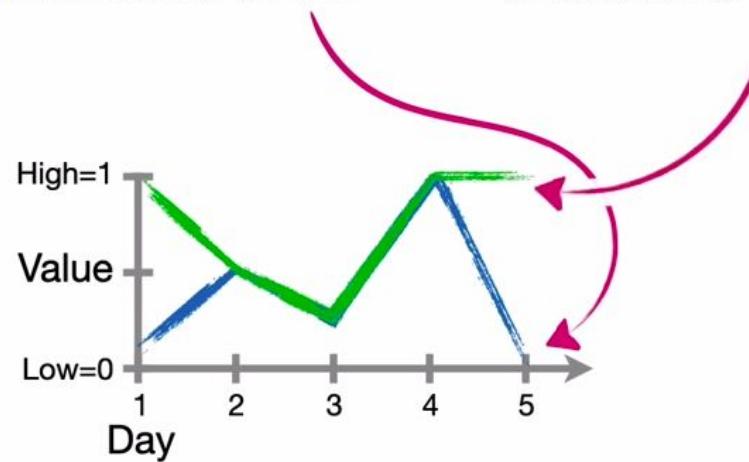


Company A =

Company B =



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

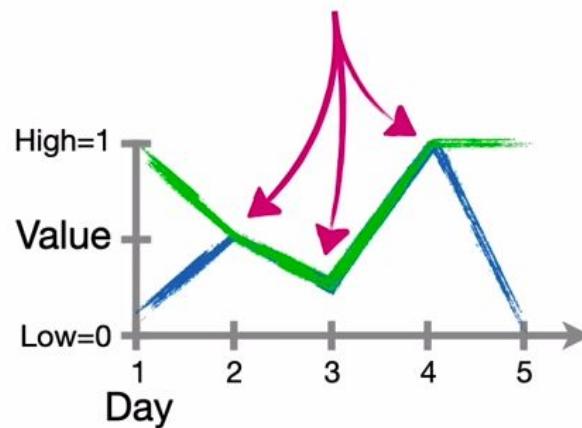


Company A =

Company B =



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

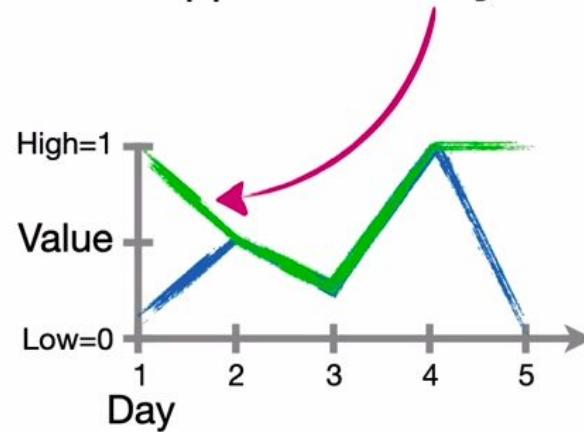


Company A =

Company B =



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

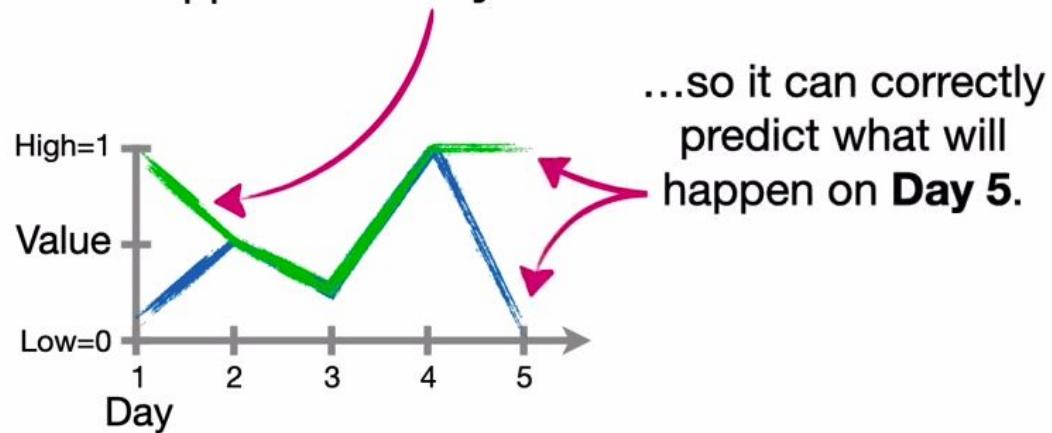


Company A =

Company B =



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

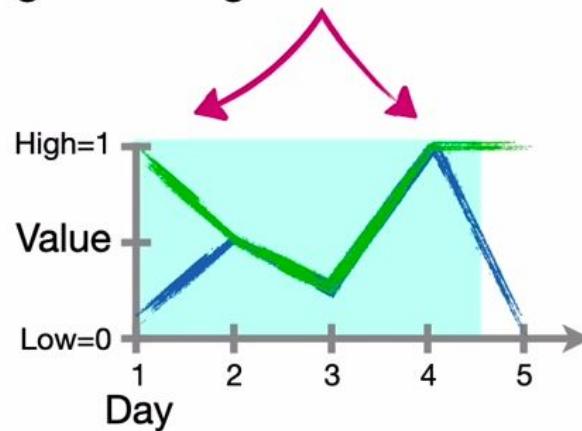


Company A =

Company B =



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

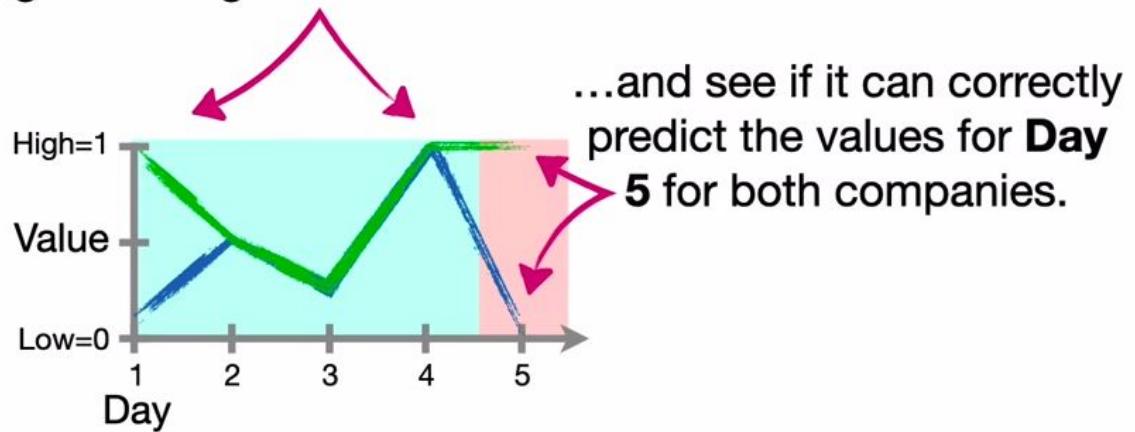


Company A =

Company B =

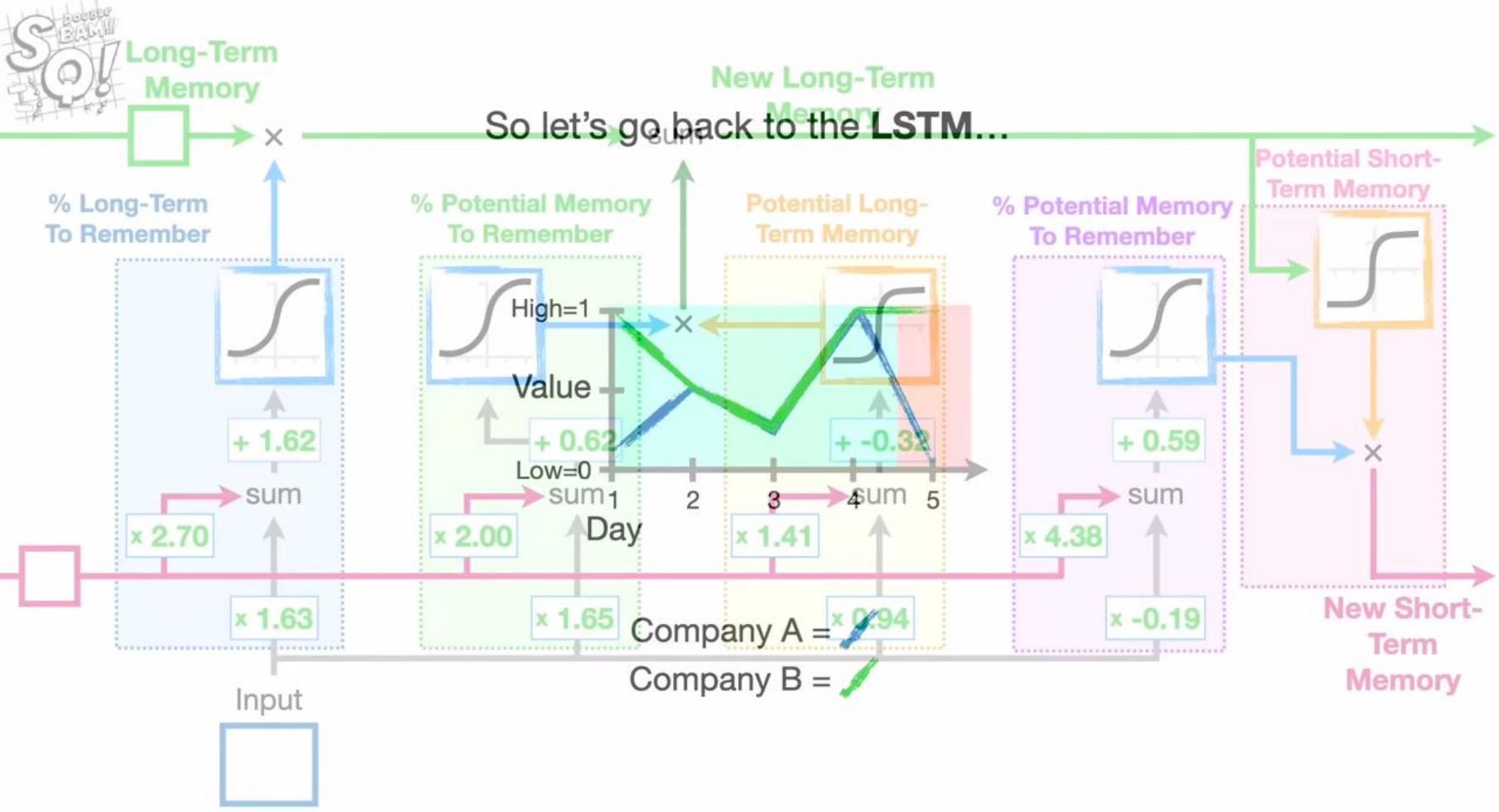


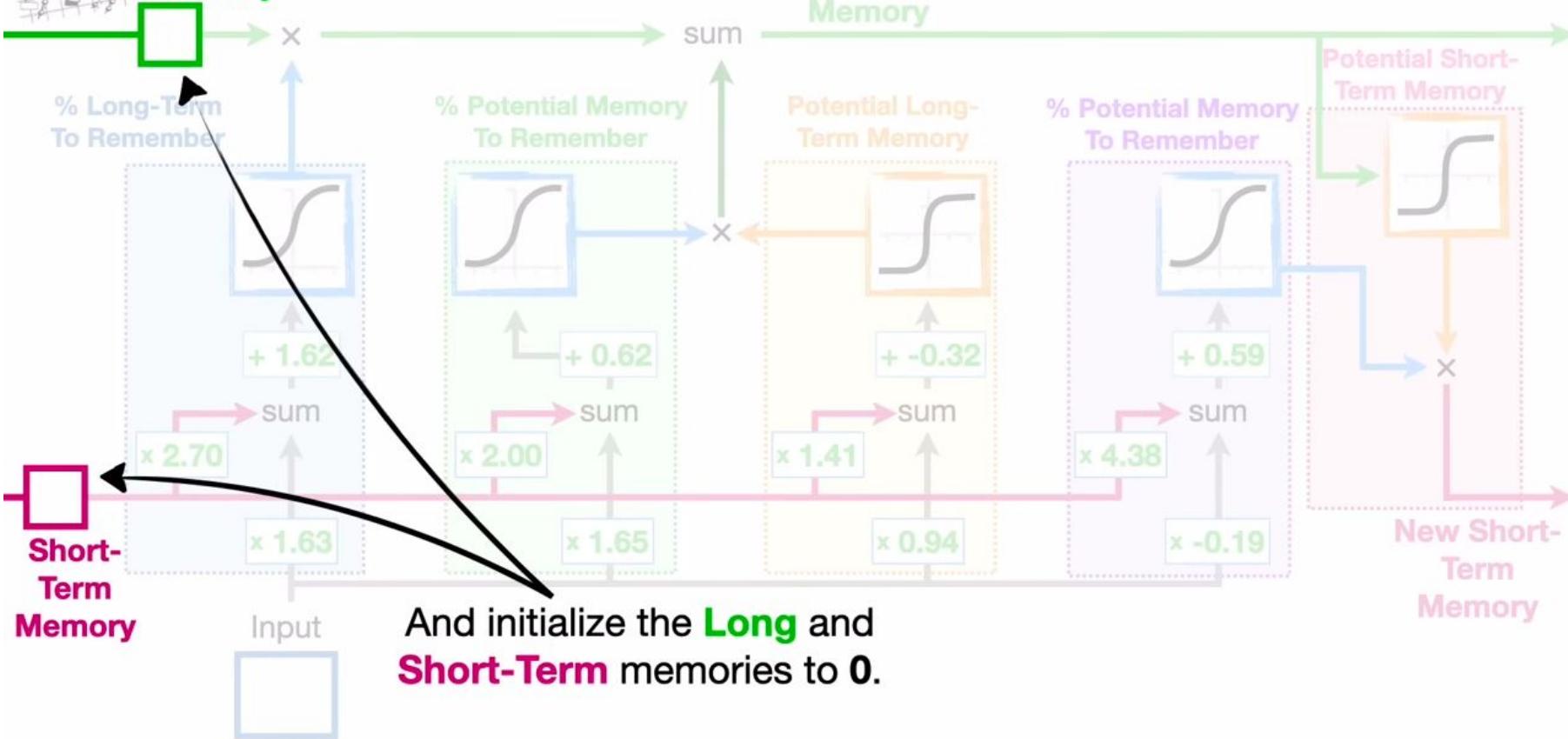
In other words, we're going to sequentially run the data from **Days 1** through **4** through an unrolled **LSTM**...



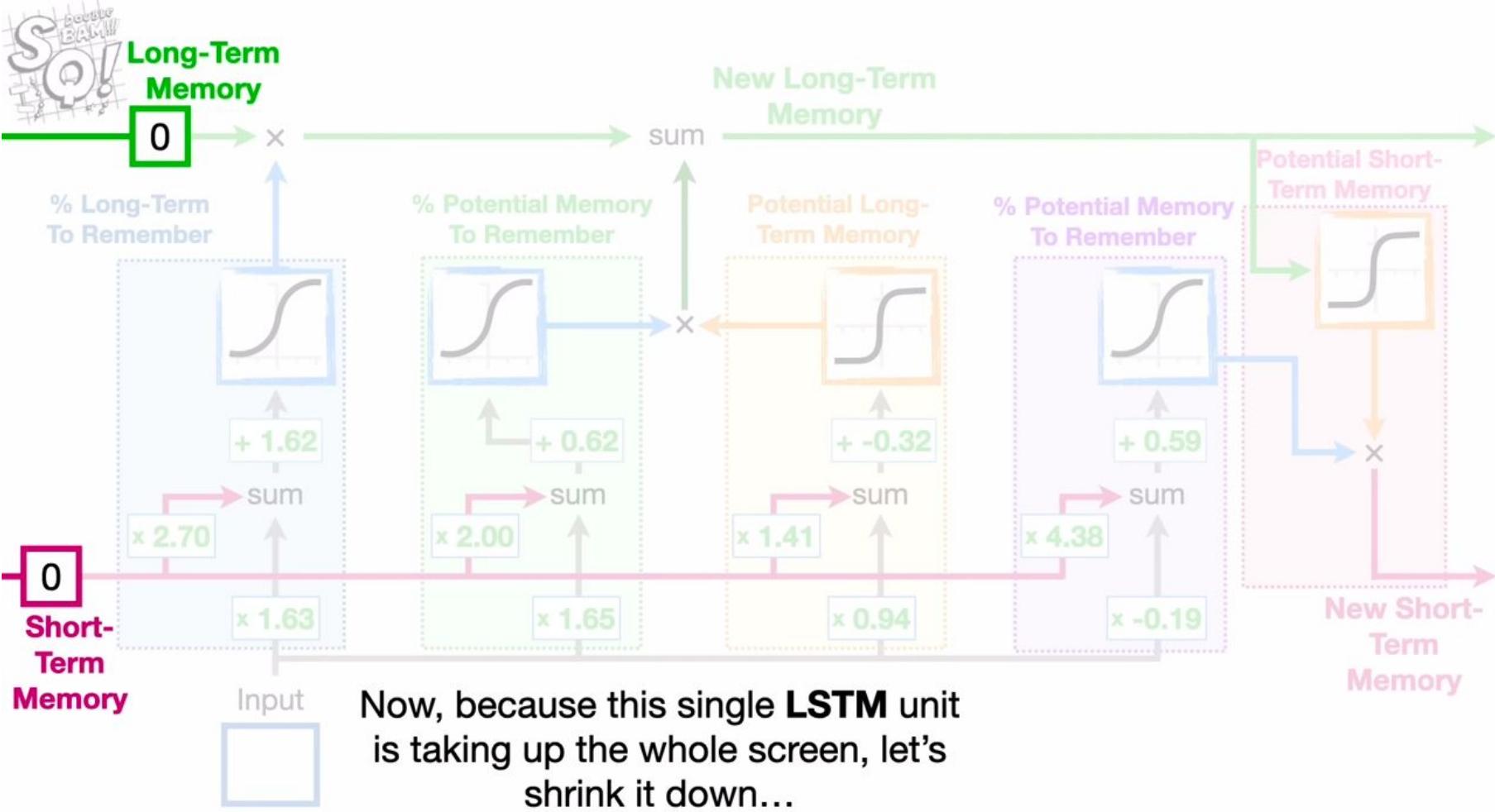
Company A =

Company B =





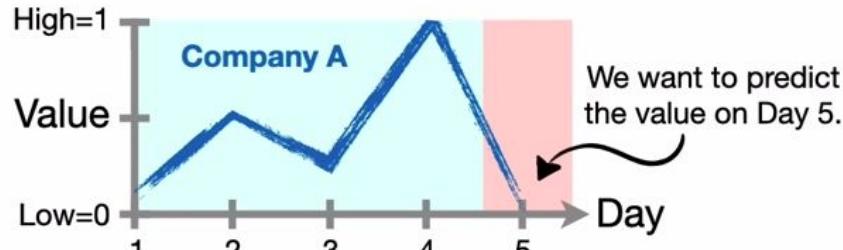
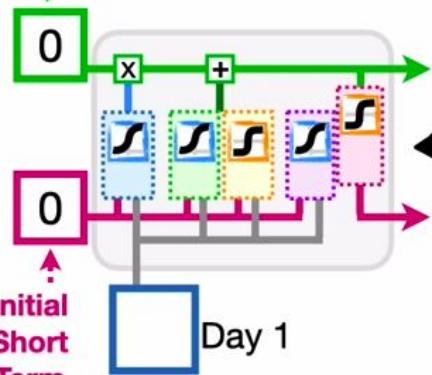
And initialize the **Long** and
Short-Term memories to 0.



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



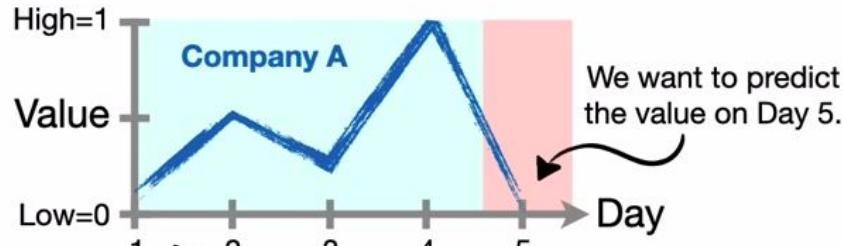
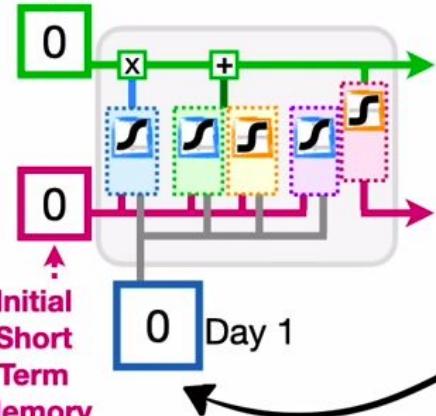
Initial
Long Term
Memory



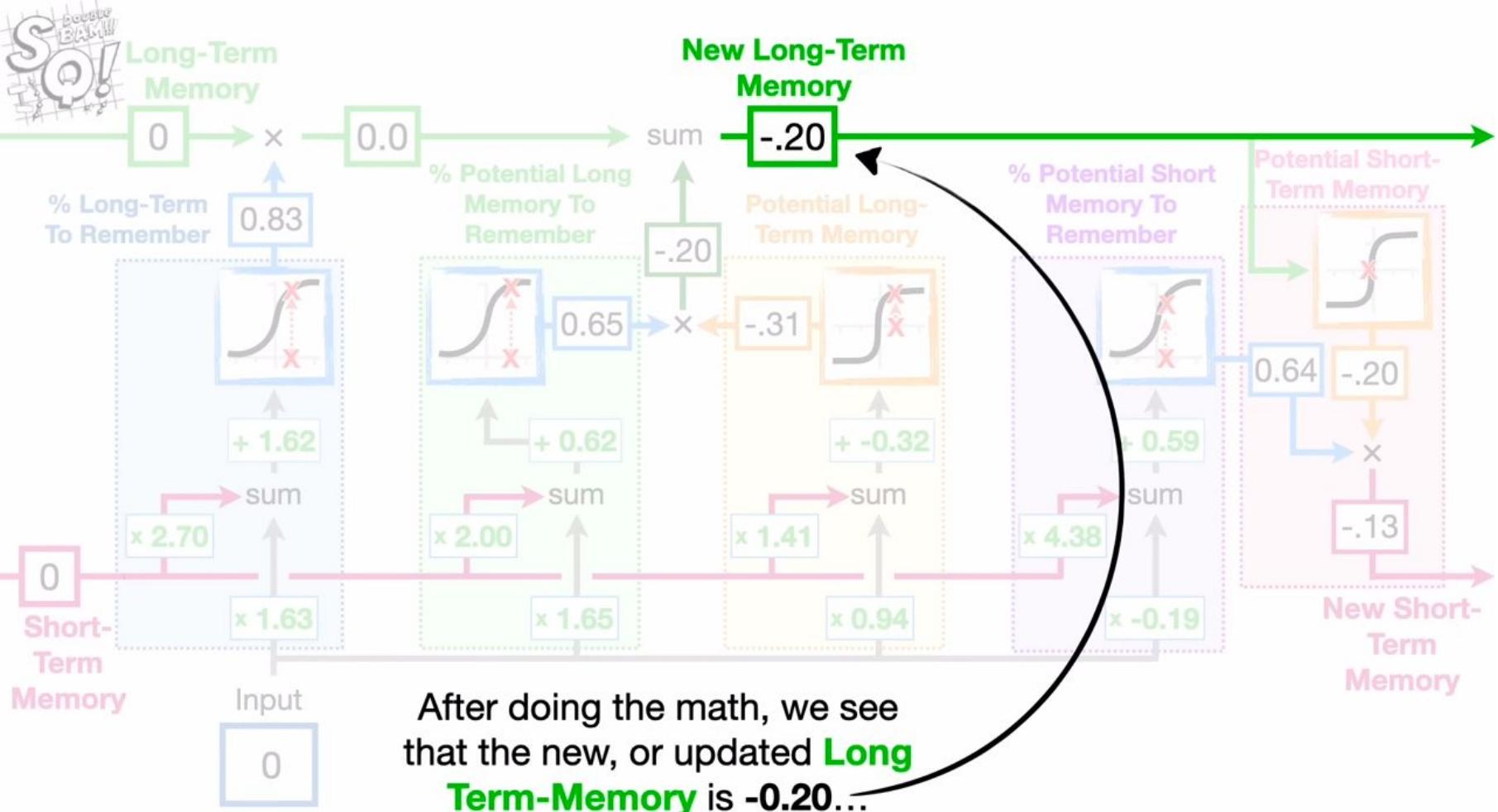
Now, if we want to sequentially run
Company A's values from **Days 1**
through **4** through this **LSTM**...

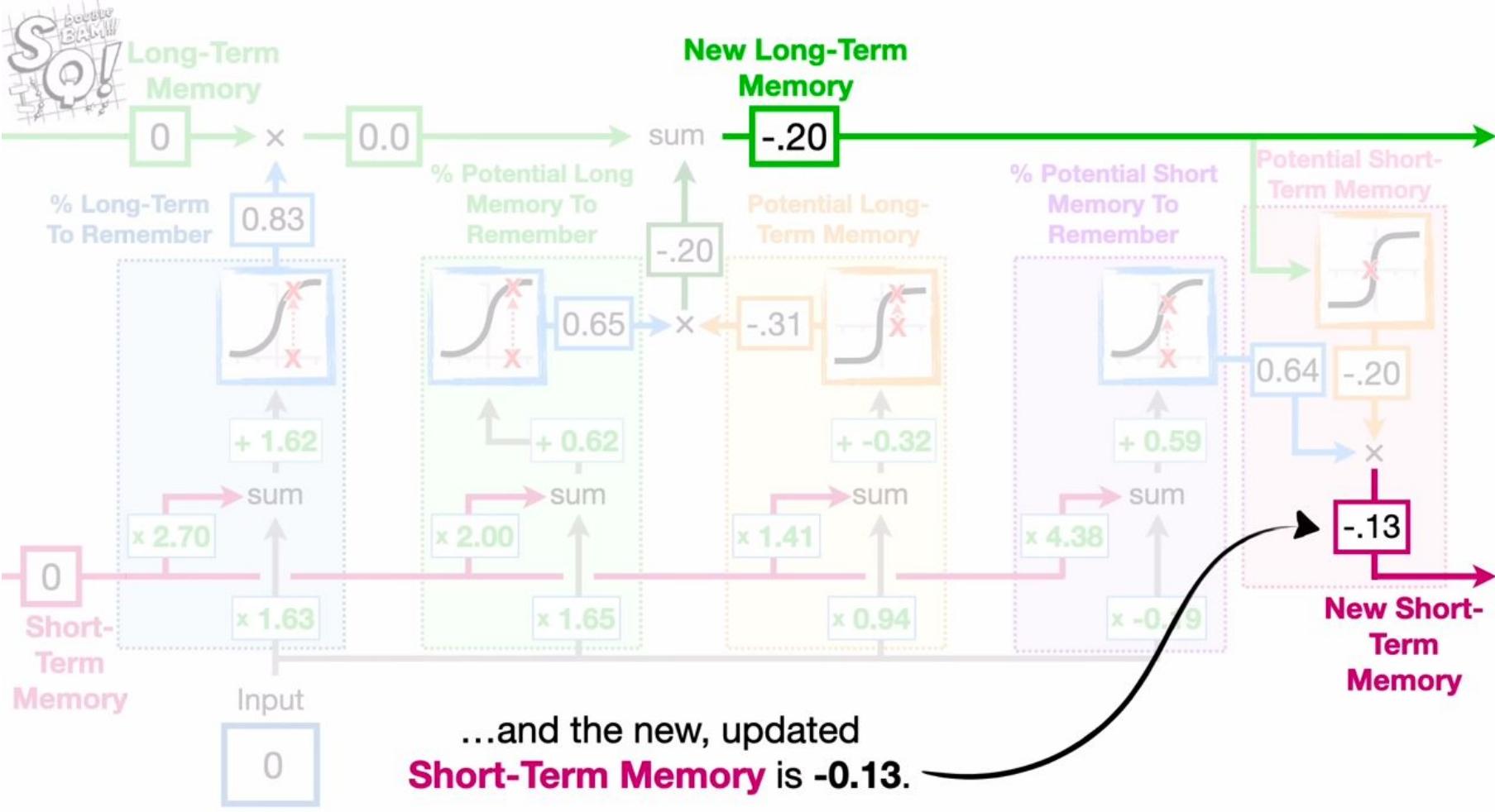


Initial
Long Term
Memory



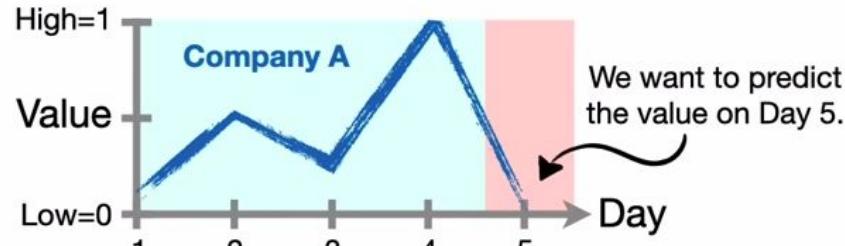
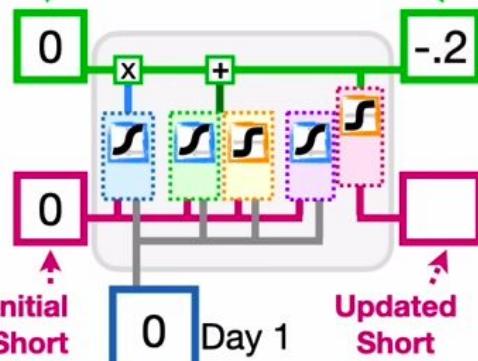
...then we'll start by plugging the value for **Day 1**, which is **0**, into the **Input**.



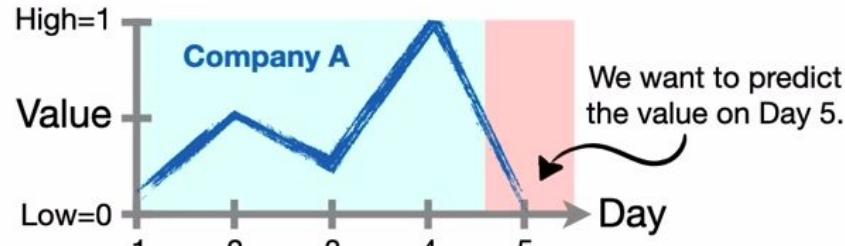
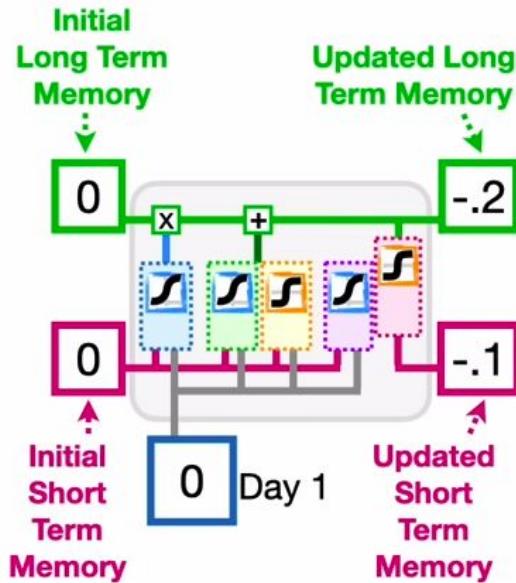




Initial
Long Term
Memory



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

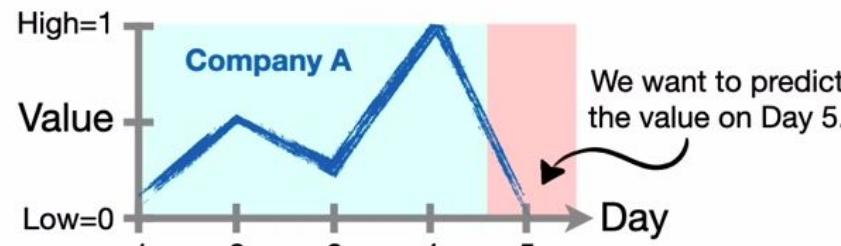
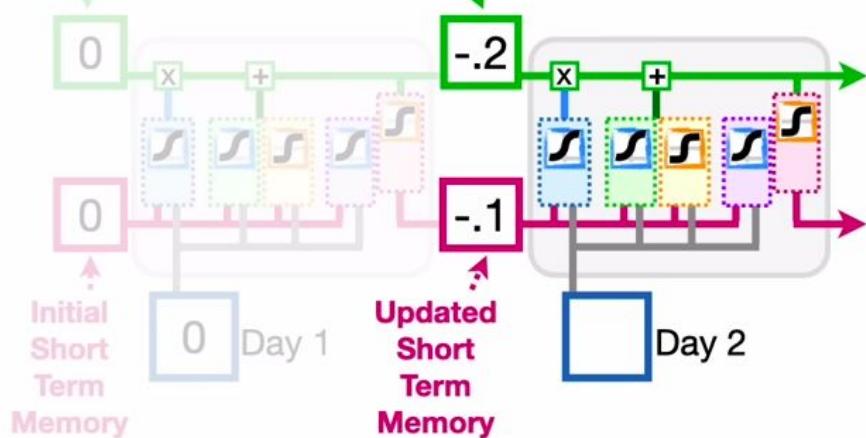


...and **-0.1** (rounded) for the
updated **Short-Term**
Memory.



Initial
Long Term
Memory

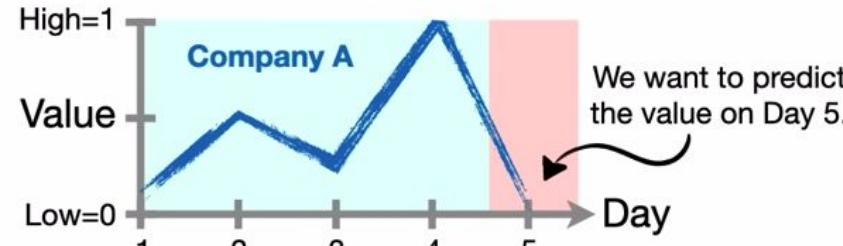
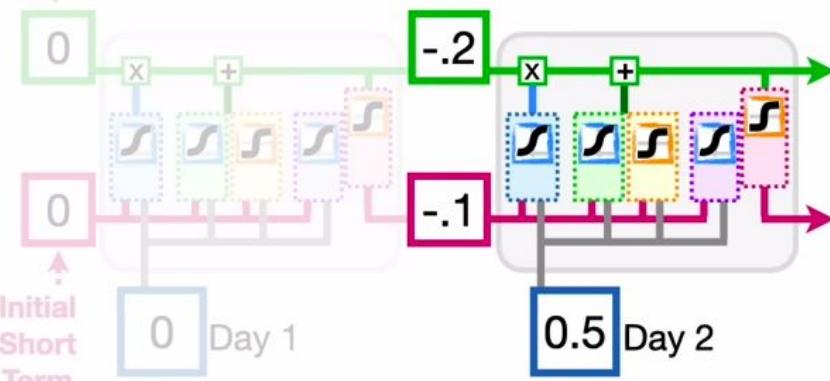
Updated Long
Term Memory



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



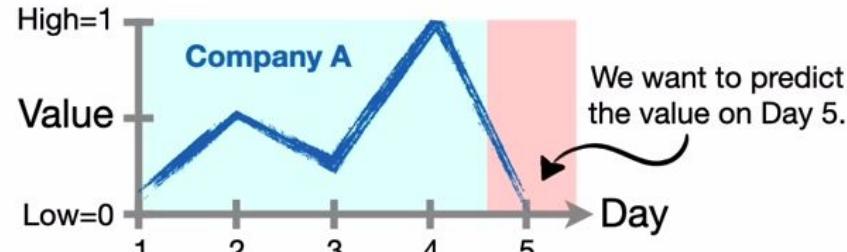
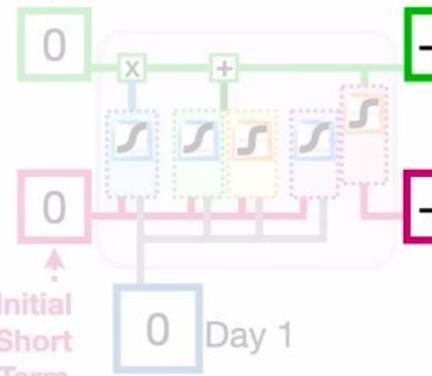
Initial
Long Term
Memory



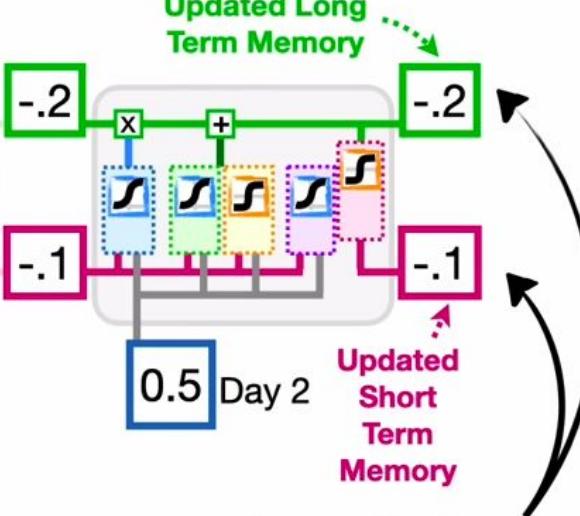
Then the **LSTM** does it's math,
using the exact same **Weights**
and **Biases** as before...



Initial
Long Term
Memory



Updated Long
Term Memory

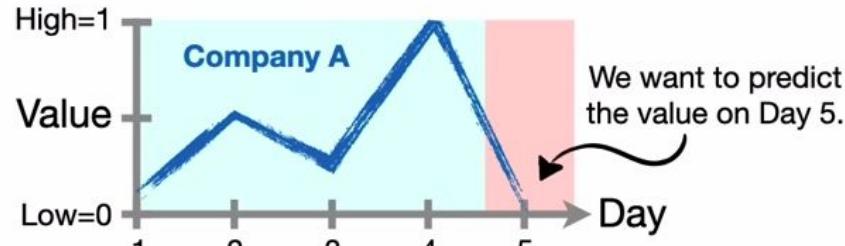
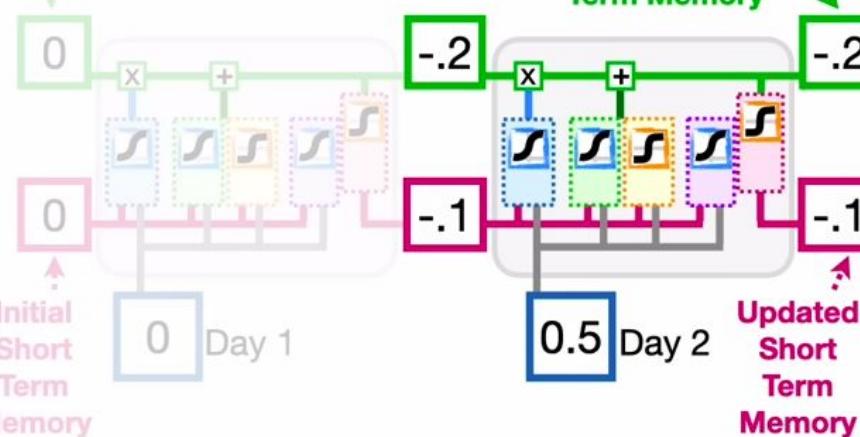


Updated
Short
Term
Memory

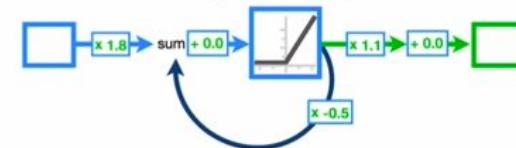
...and we end up with these updated
Long and Short-Term Memories.



Initial
Long Term
Memory



Recurrent Neural Networks (RNNs)...

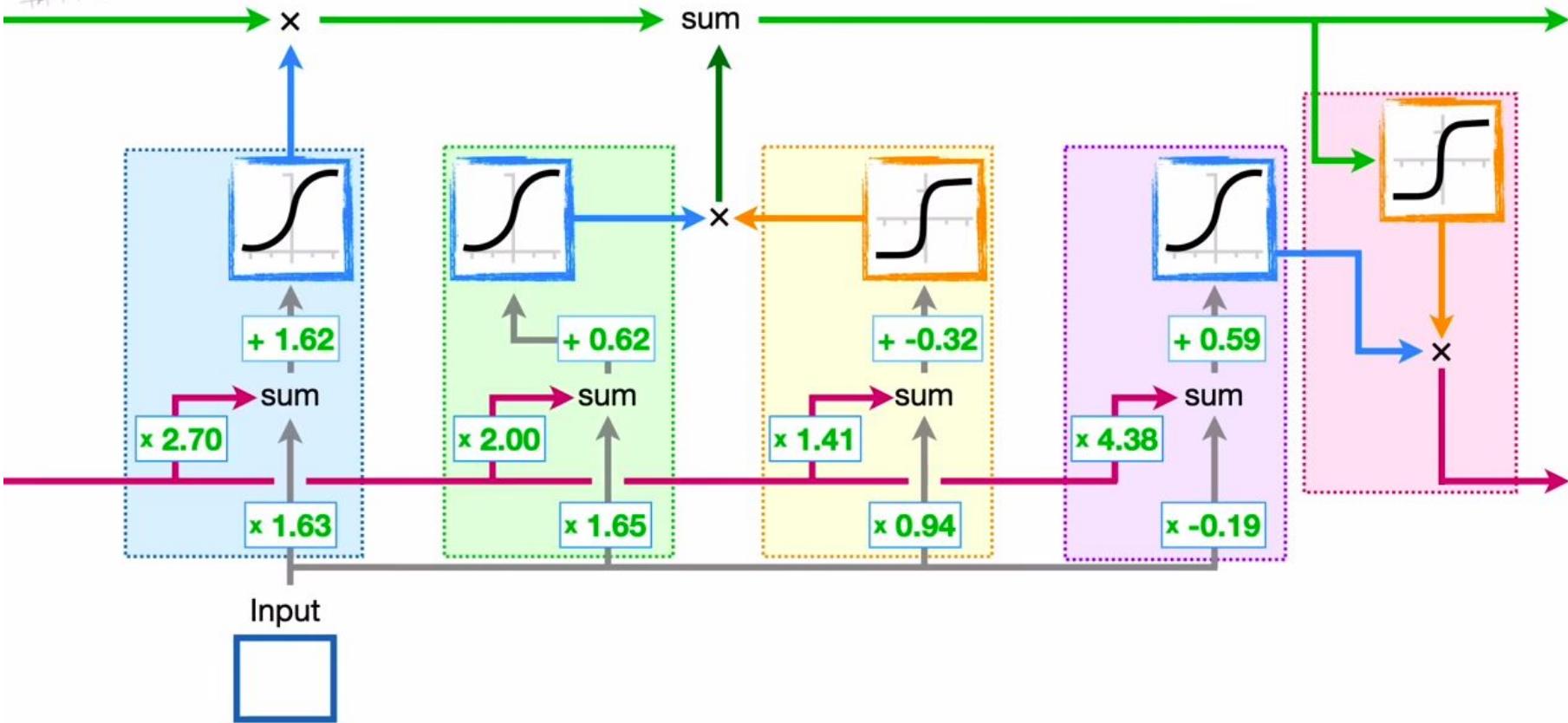


...Clearly Explained!!!

...the reason the **LSTM** reuses the exact same **Weights** and **Biases** is so it can handle data sequences of different lengths.

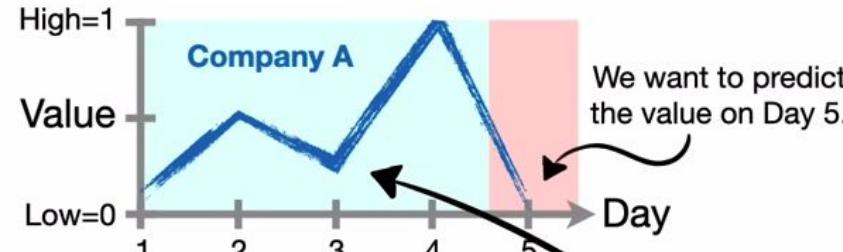
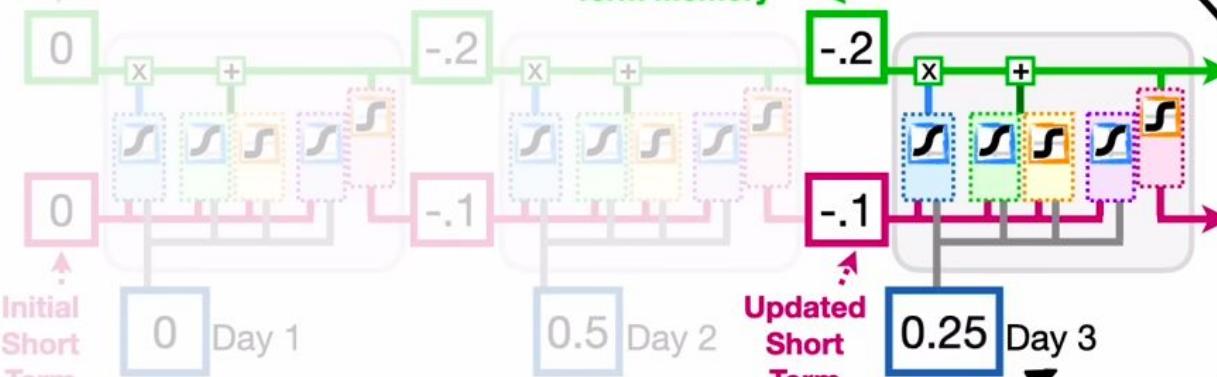


...let's talk about how the **Long Short-Term Memory** unit works.





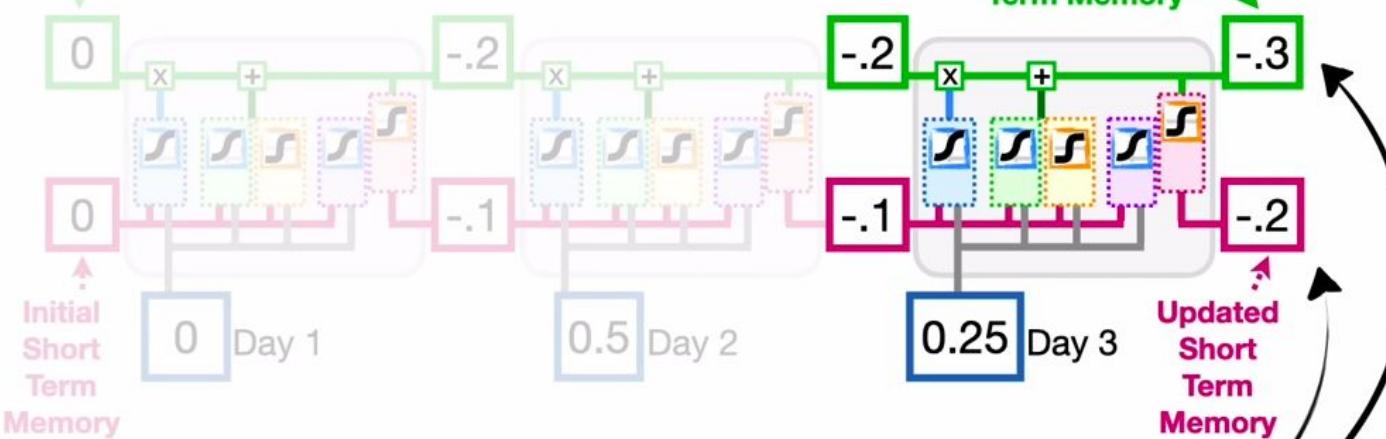
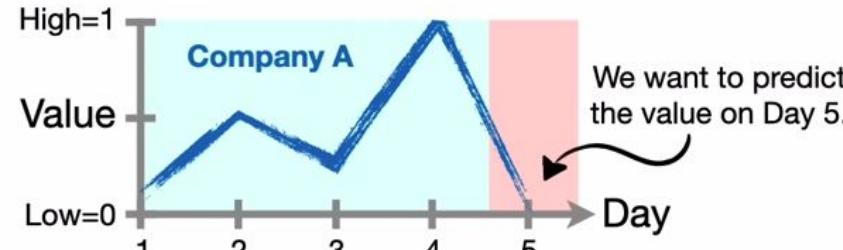
Initial
Long Term
Memory



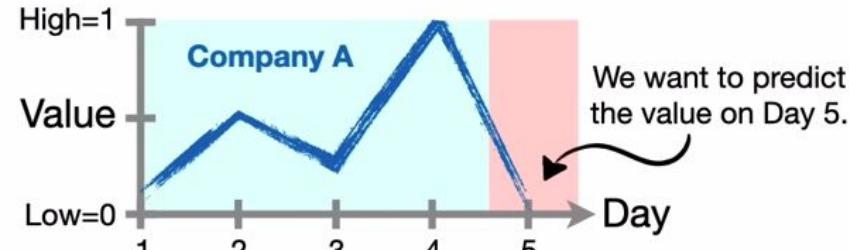
...and plug in the value for **Day 3**.



Initial
Long Term
Memory



...and gives us these
updated memories.



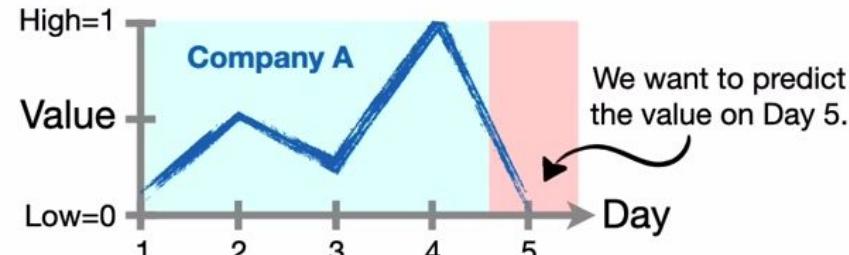
Initial
Long Term
Memory

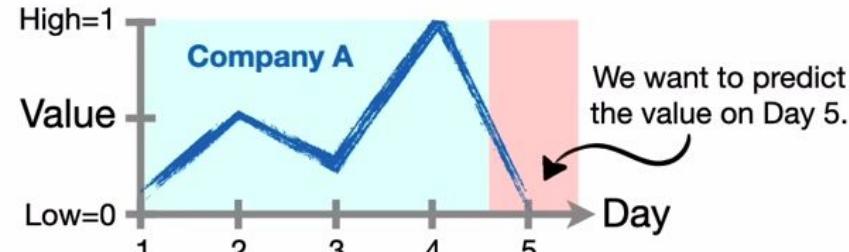


And the **LSTM** does the math,
again, using the exact same
Weights and Biases...

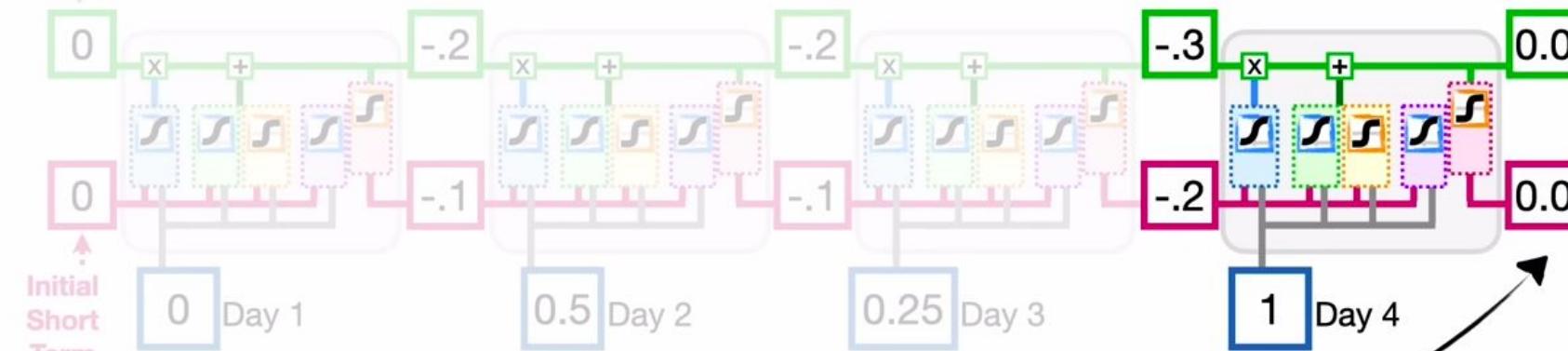


Initial
Long Term
Memory





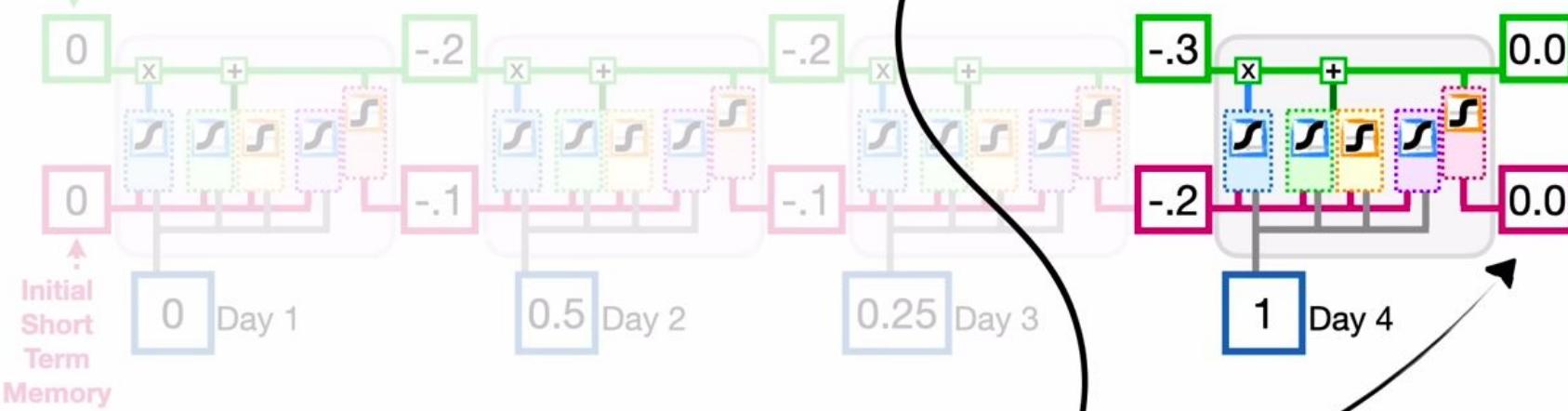
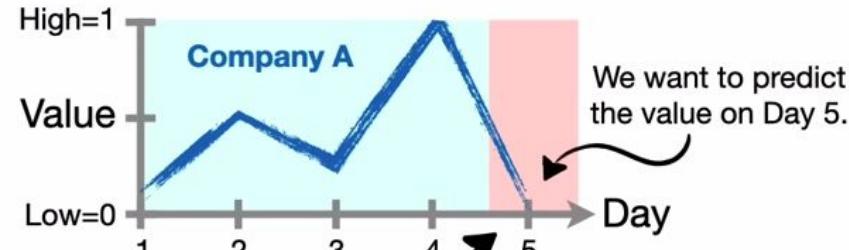
Initial
Long Term
Memory



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



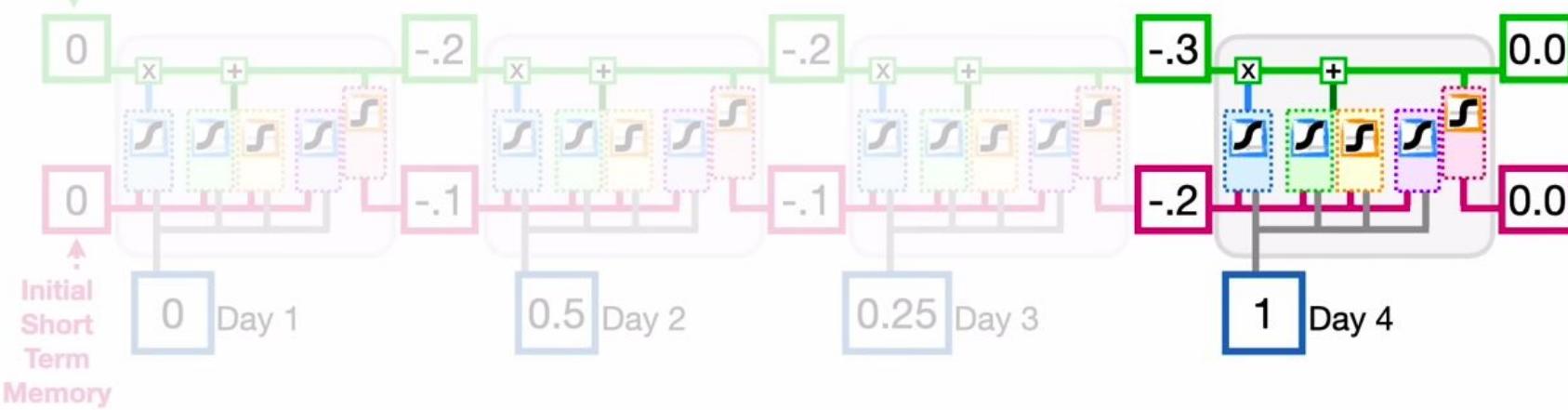
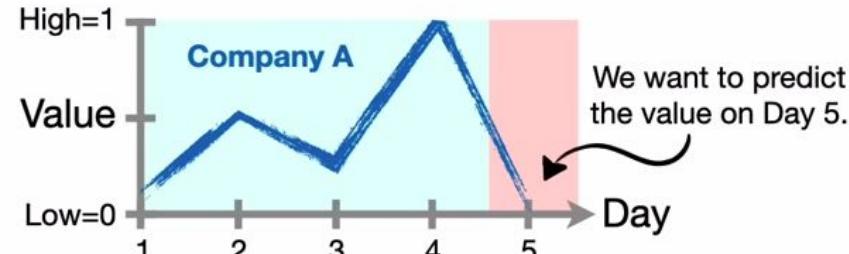
Initial
Long Term
Memory



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



Initial
Long Term
Memory

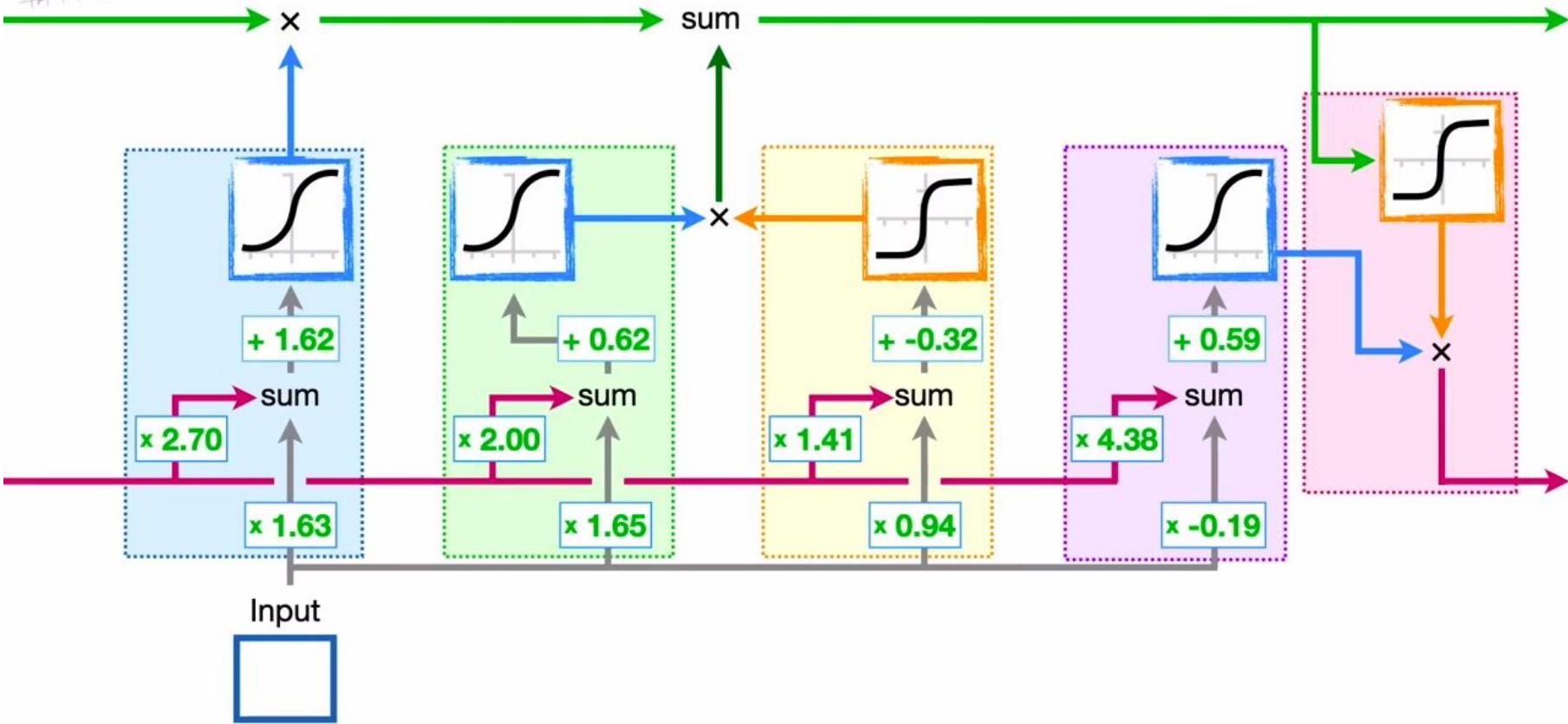


Now that we have shown that the
LSTM can correctly predict the value
on **Day 5** for **Company A**...

Quiz



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



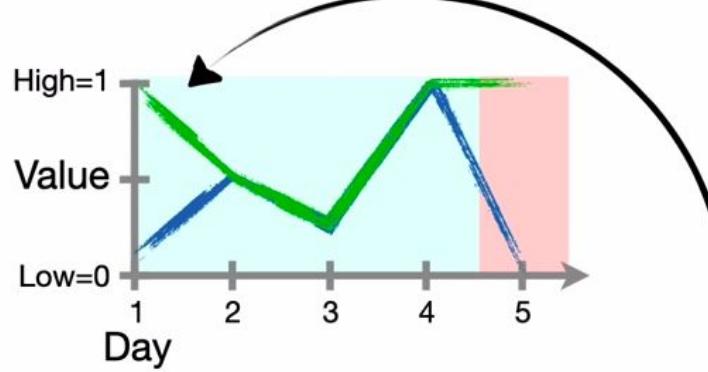


...let's show how the same **LSTM**, with the same **Weights** and **Biases**, can correctly predict the value on **Day 5** for **Company B**.

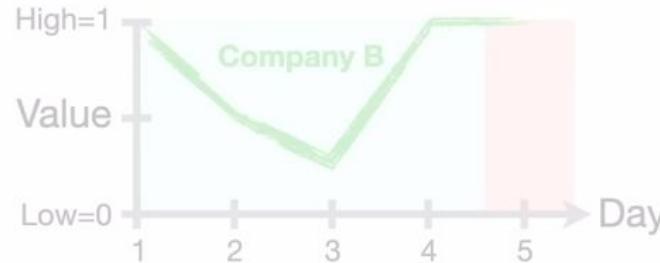


Company A =

Company B =

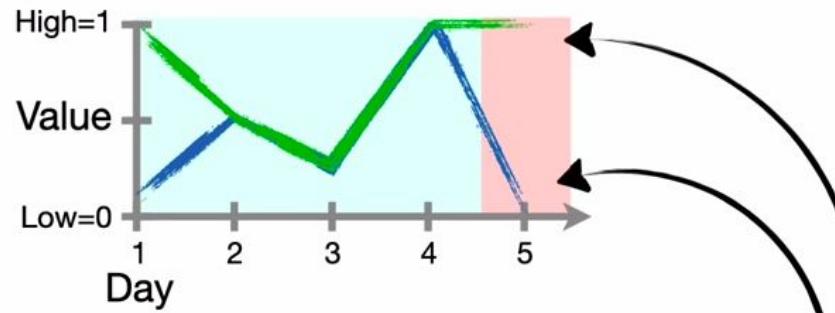


NOTE: Remember, on **Days 1 through 4**, the only difference between the companies occurs on **Day 1...**

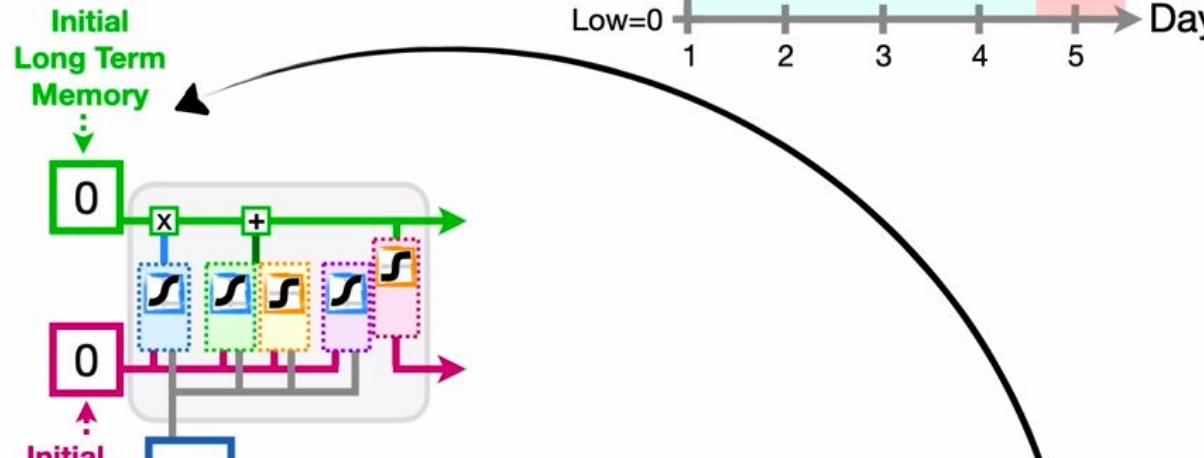


Company A =

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**.



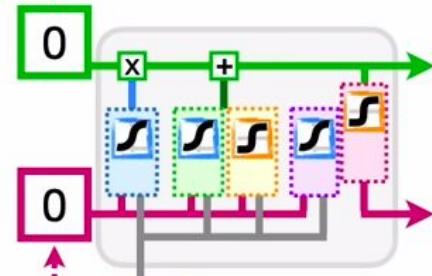
Initial
Short
Term
Memory

Day 1

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



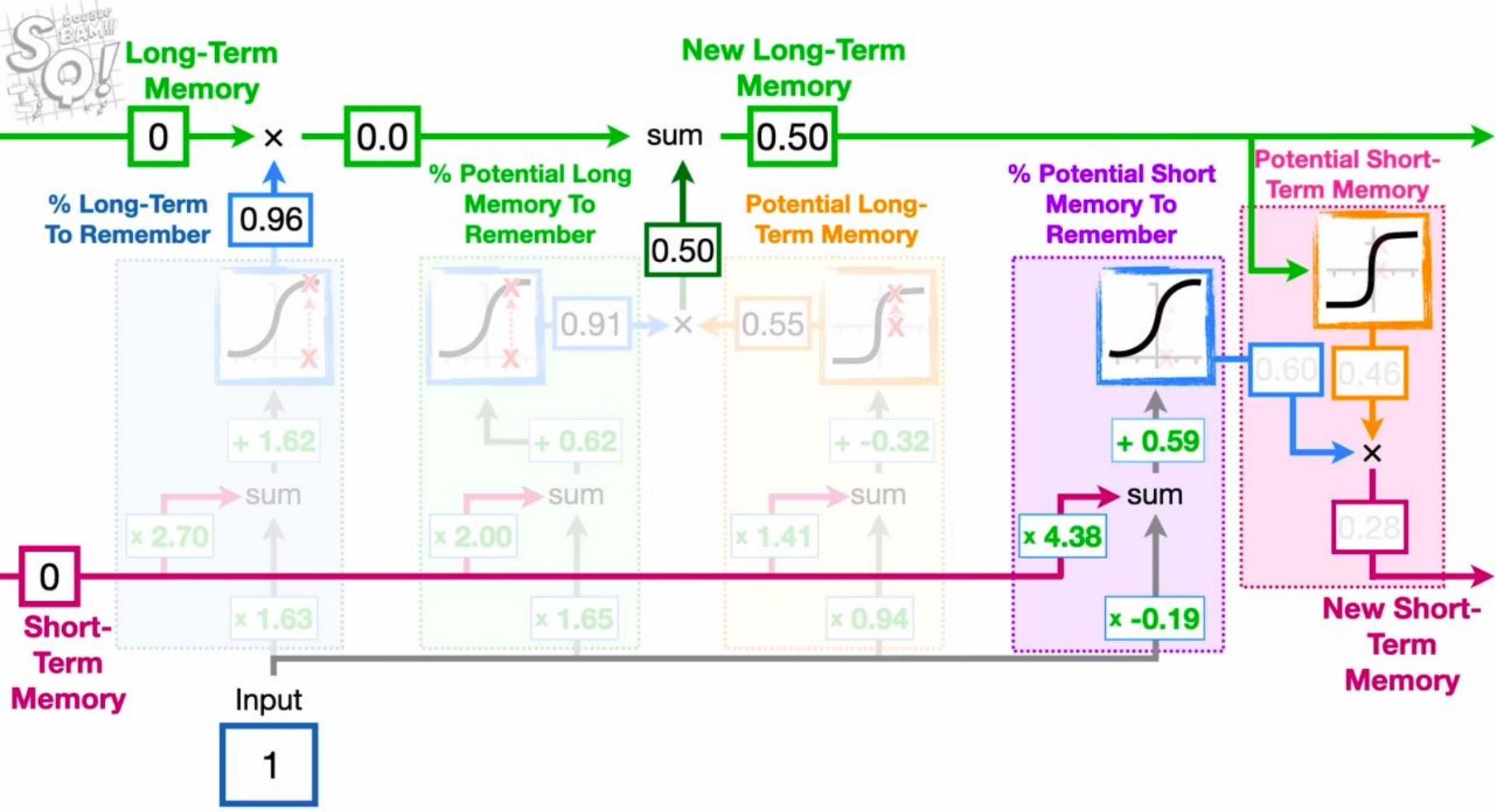
Initial
Long Term
Memory

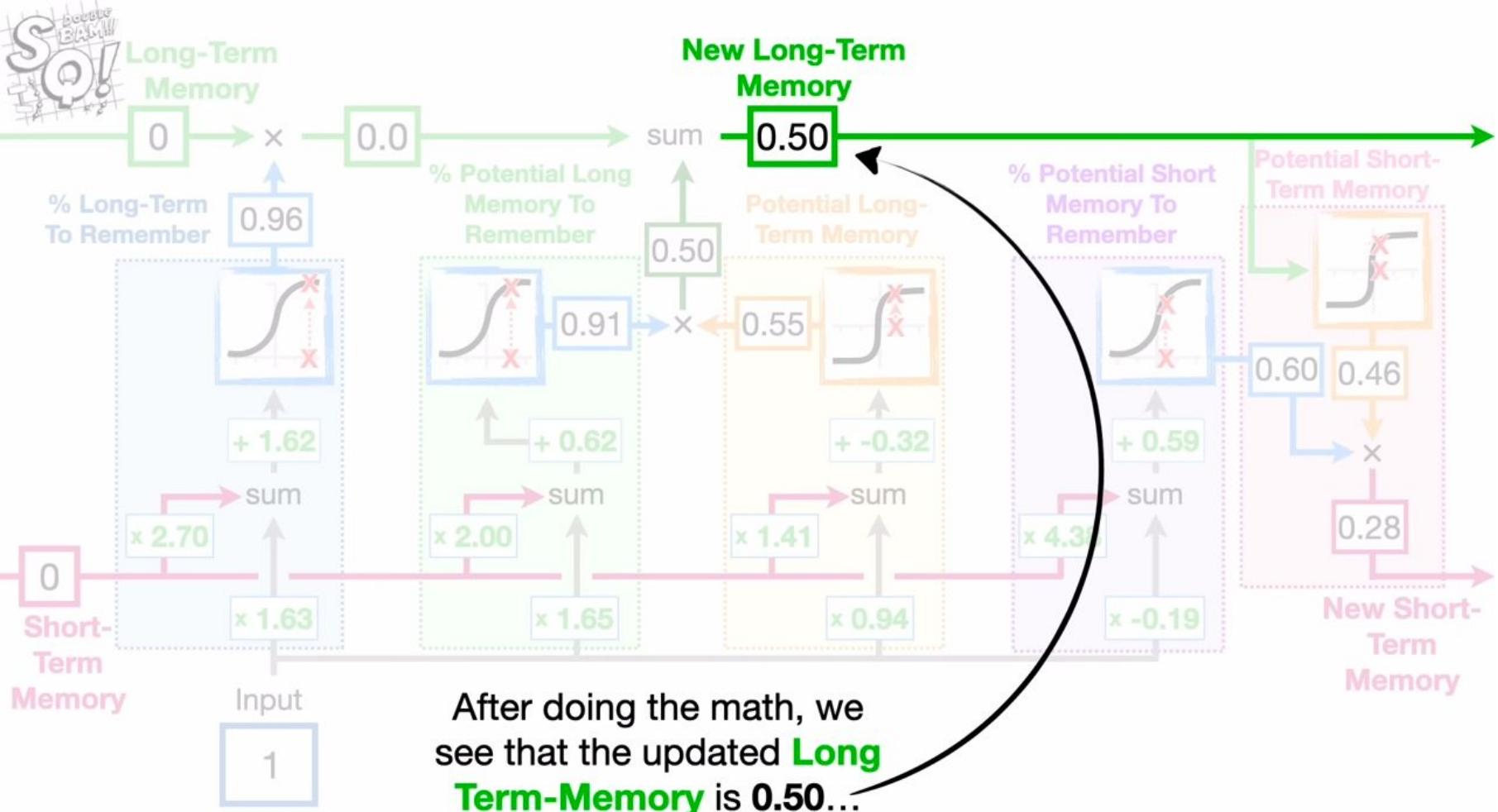


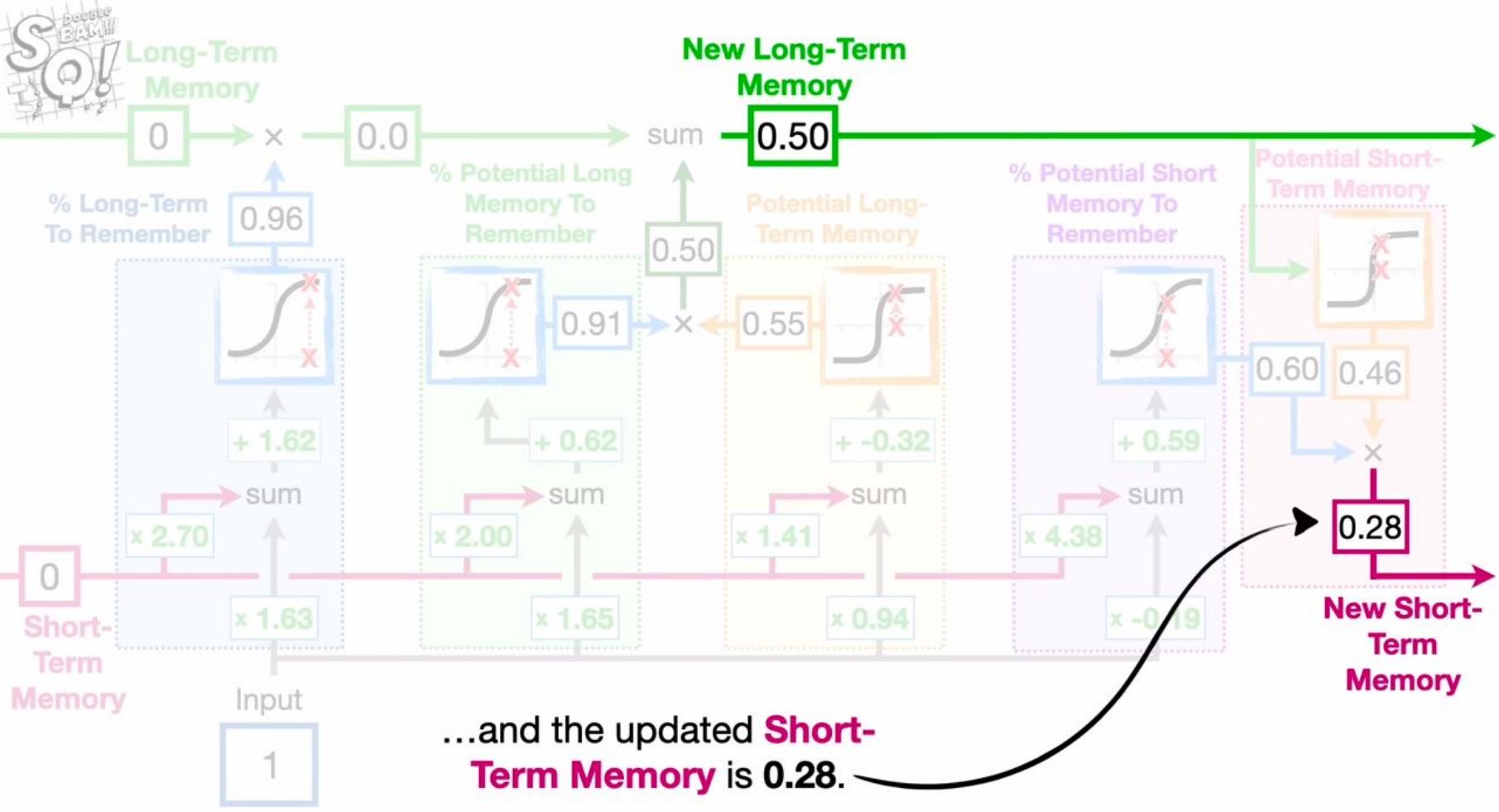
Initial
Short
Term
Memory

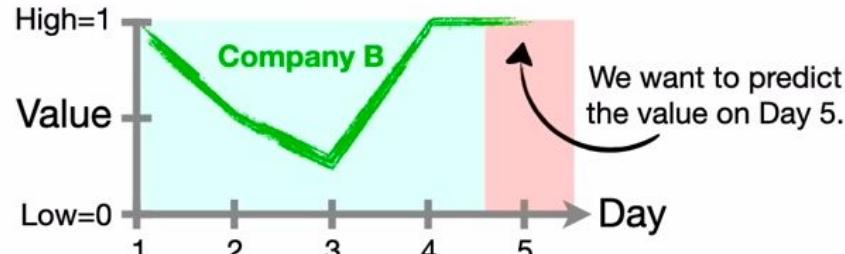
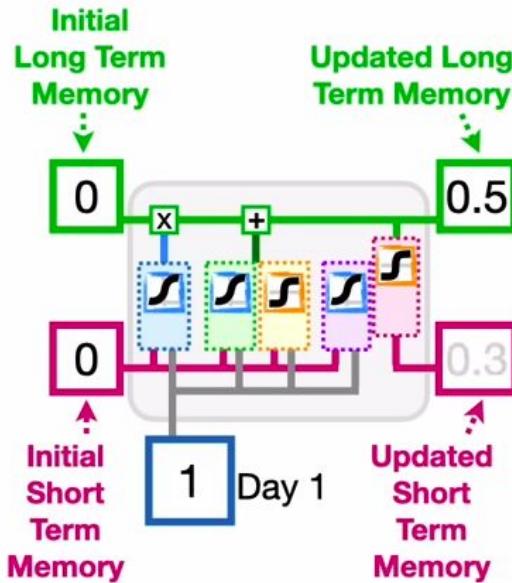


Now let's plug in the
value for **Day 1** from
Company B, 1...

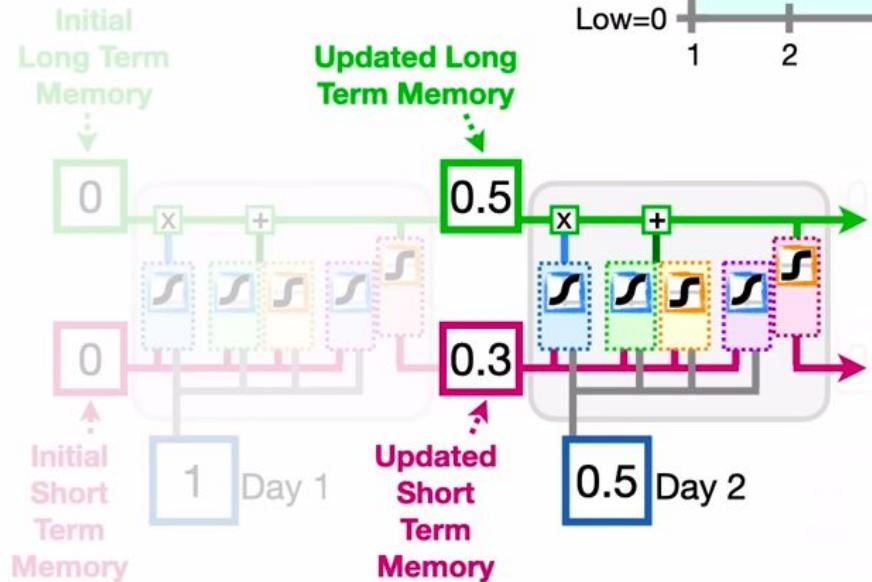




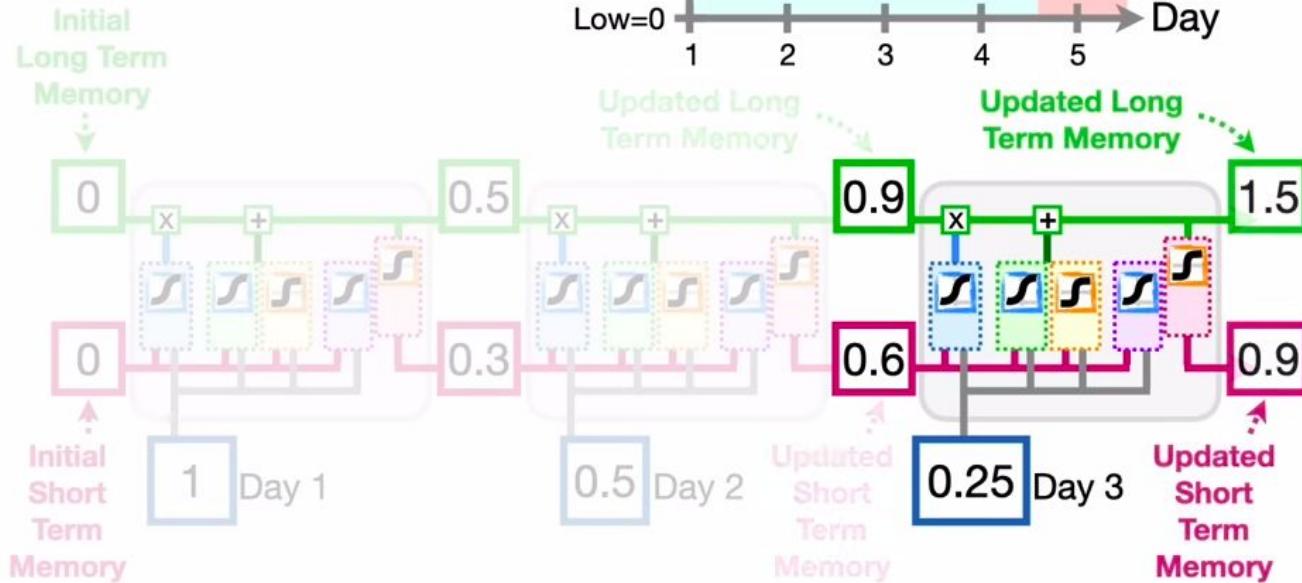
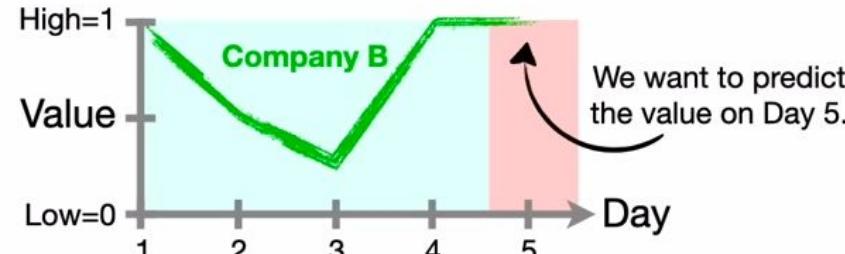




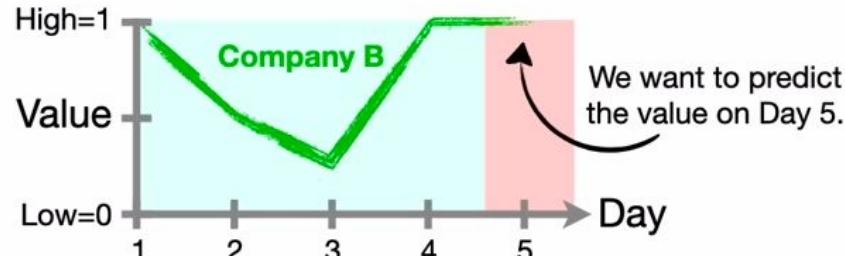
...and 0.3 (rounded) for the updated **Short-Term Memory**.



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...



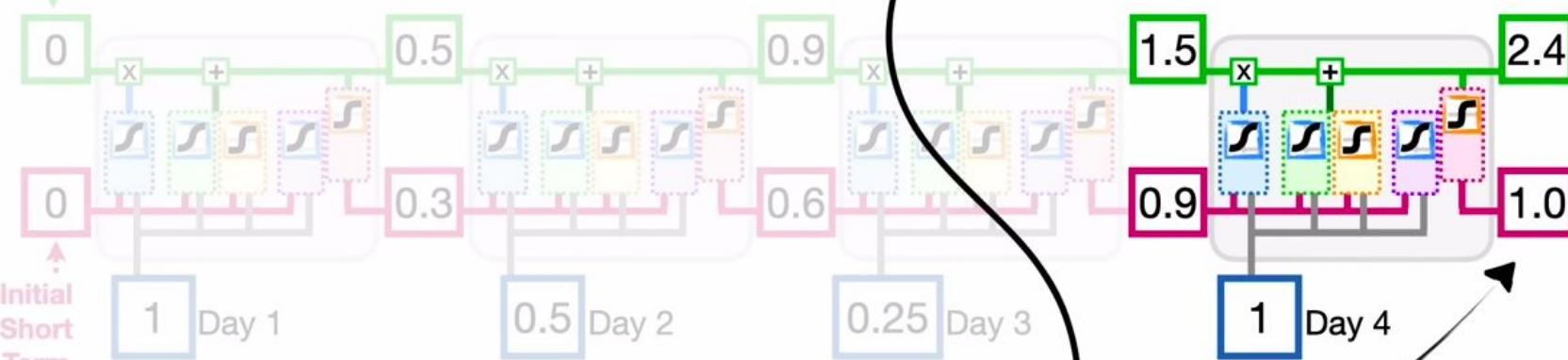
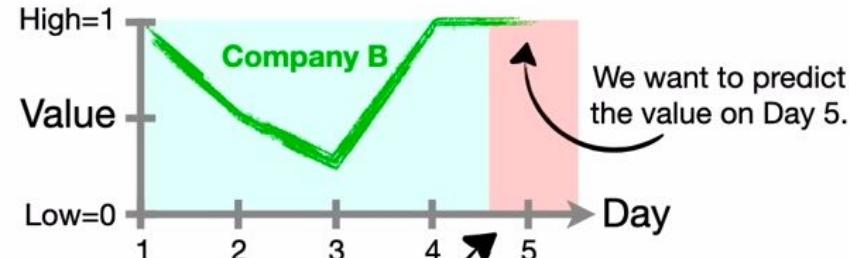
Initial
Long Term
Memory



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



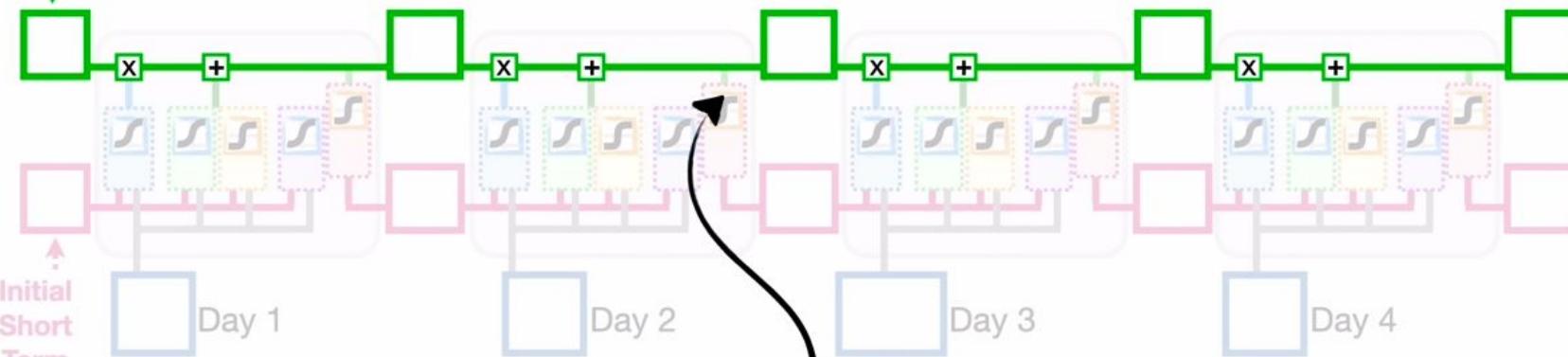
Initial
Long Term
Memory



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



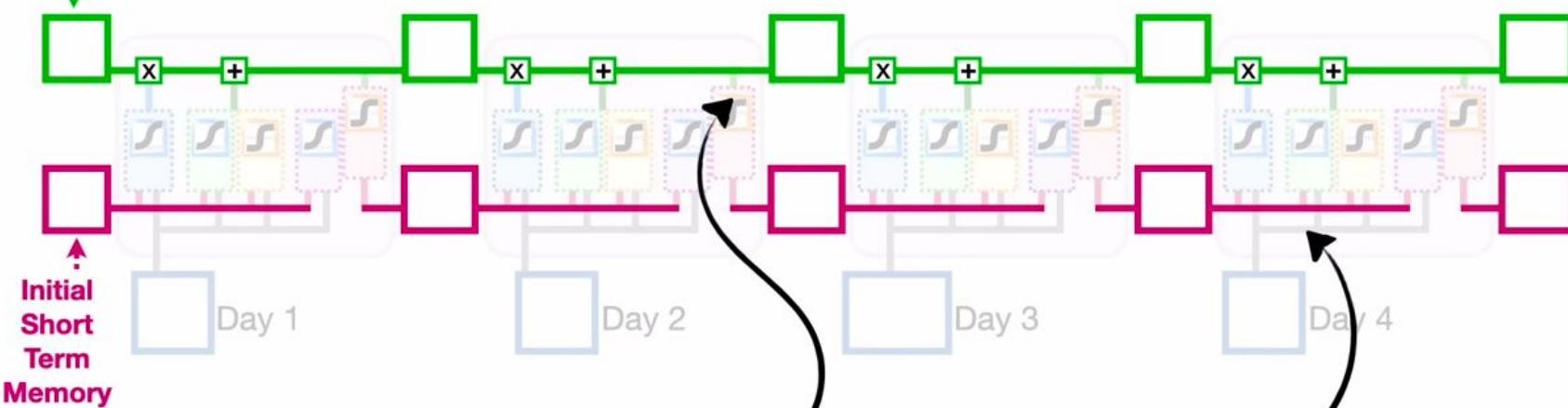
Initial
Long Term
Memory



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



Initial
Long Term
Memory

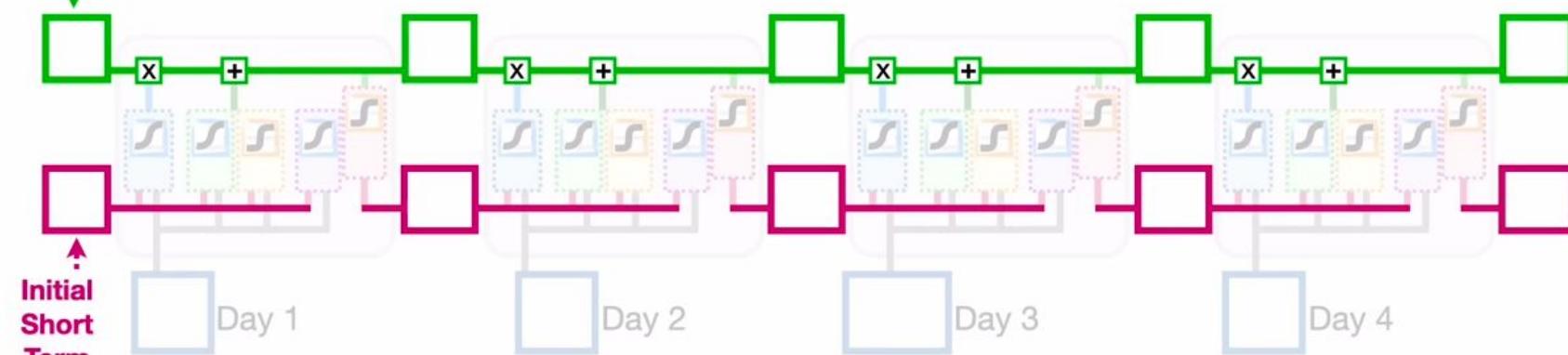


In summary, using
separate paths for **Long-
Term Memories...**

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



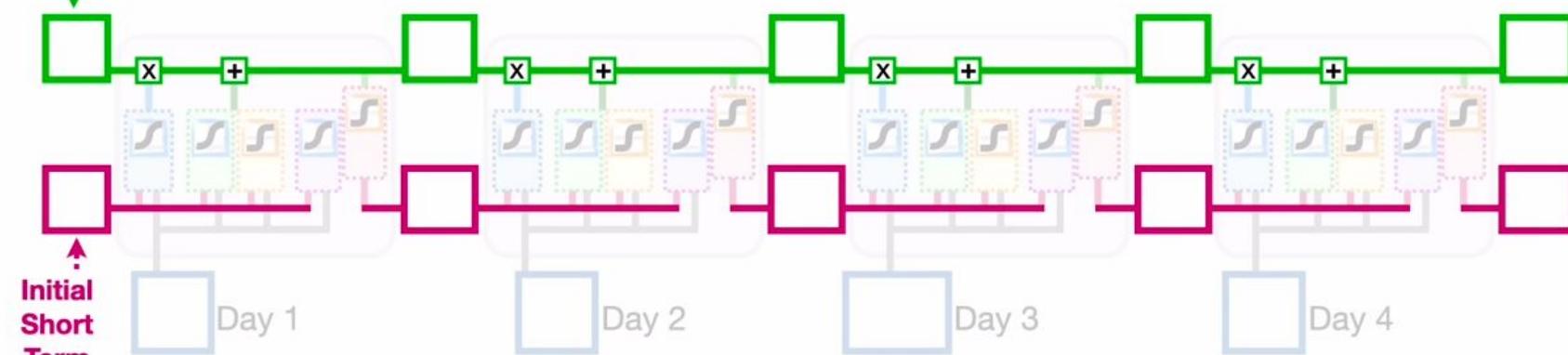
Initial
Long Term
Memory



...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**.



Initial
Long Term
Memory



At first I was scared of
how complicated the
LSTM was, but now I
understand.

