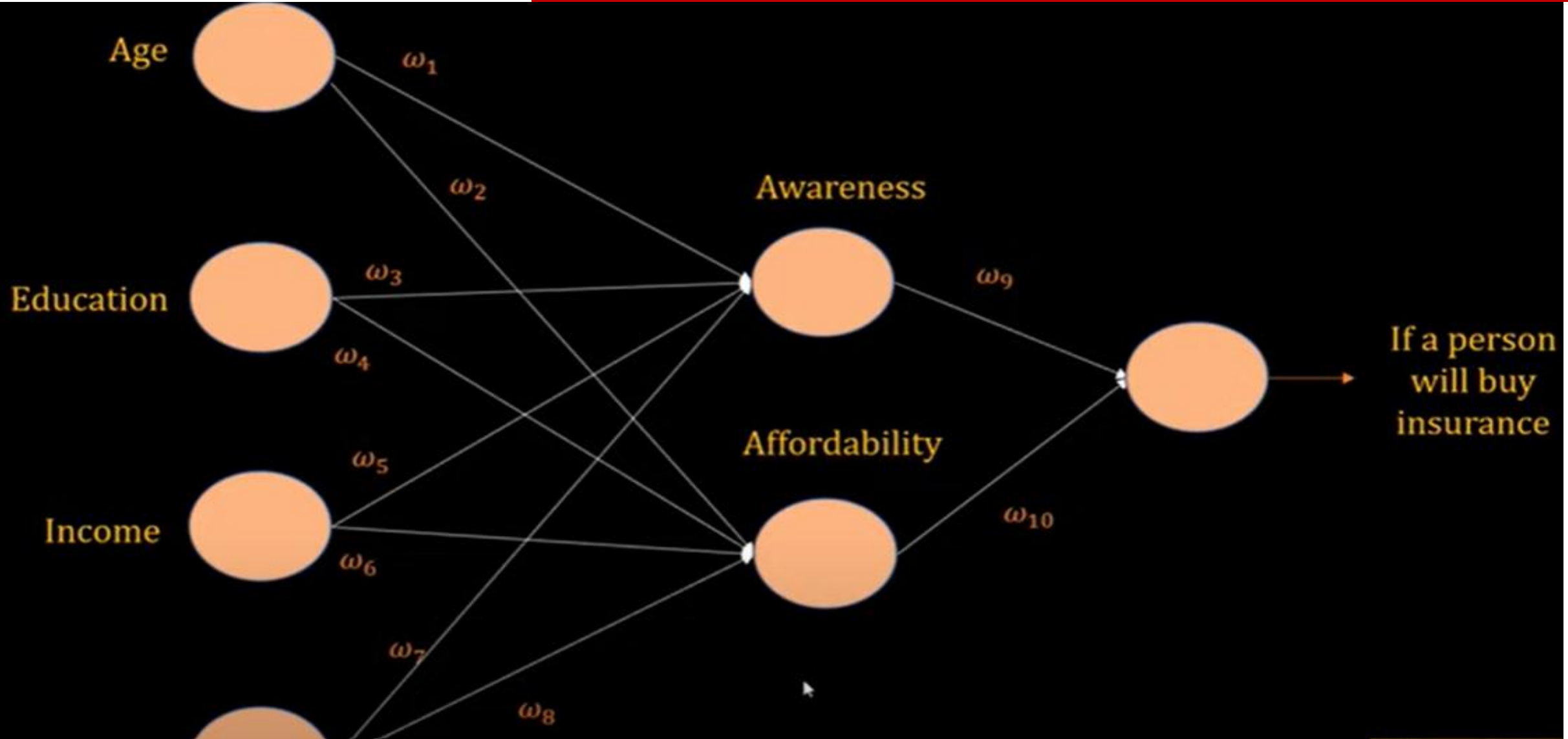
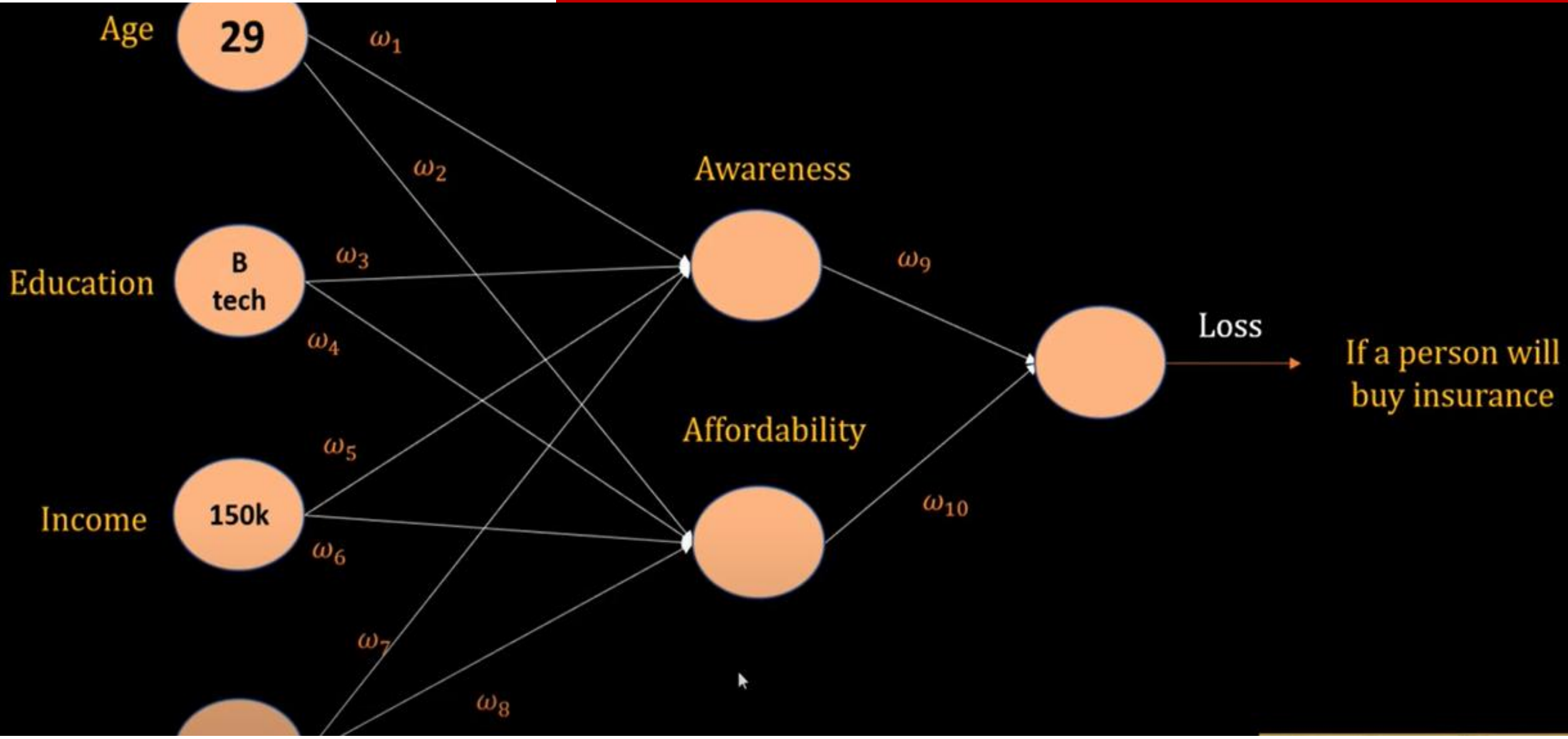
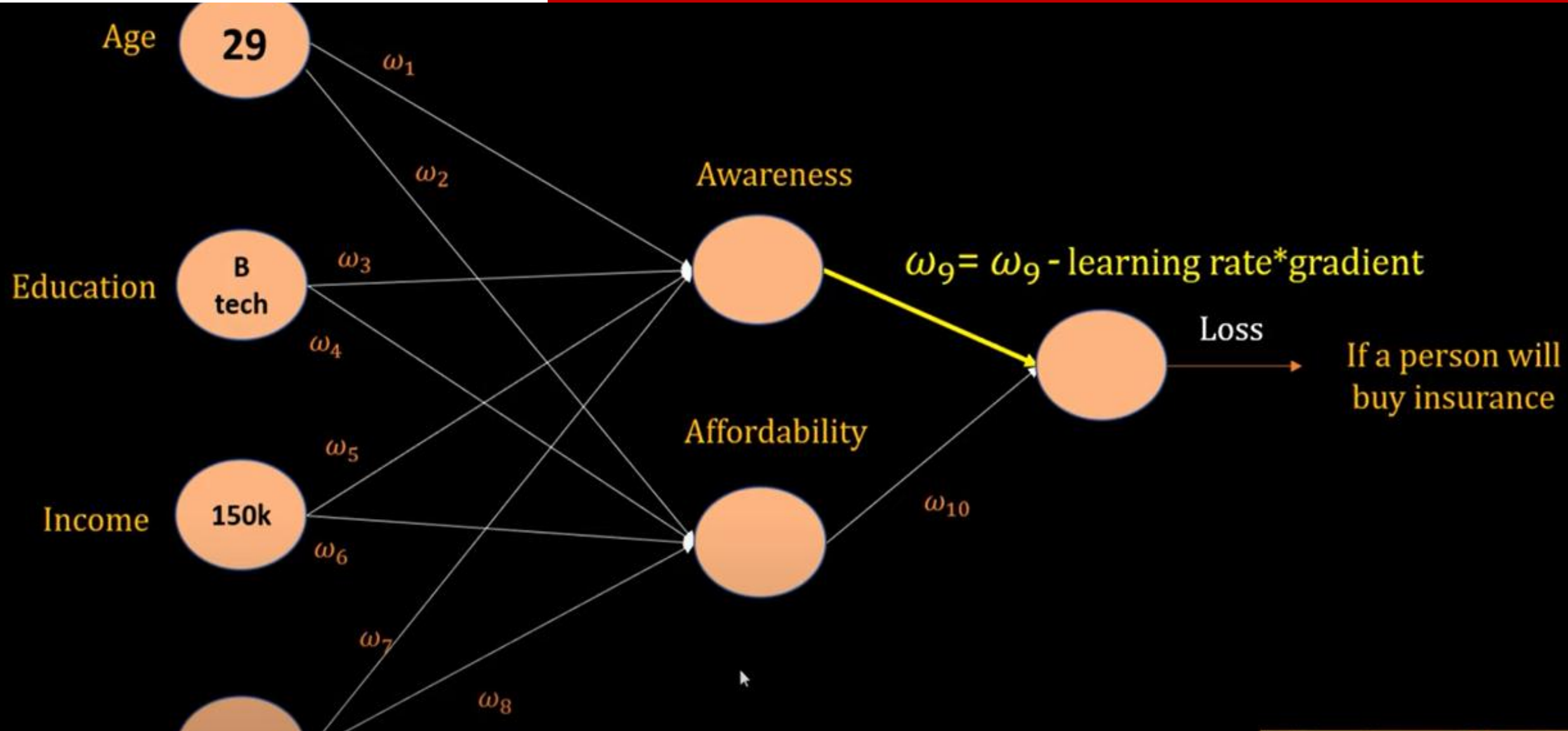


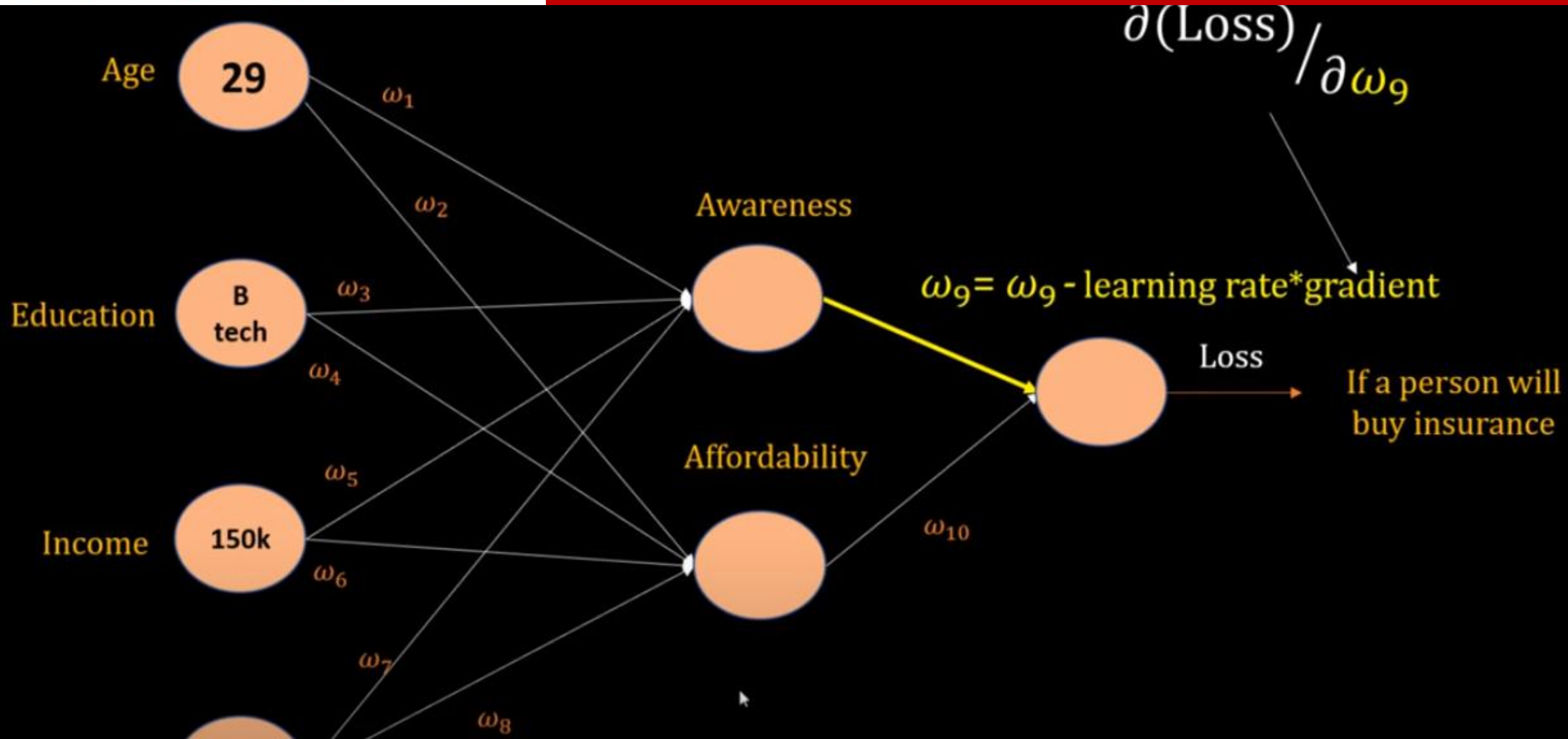
Long Short-Term Memory (LSTM)

Vanishing Gradient Issue

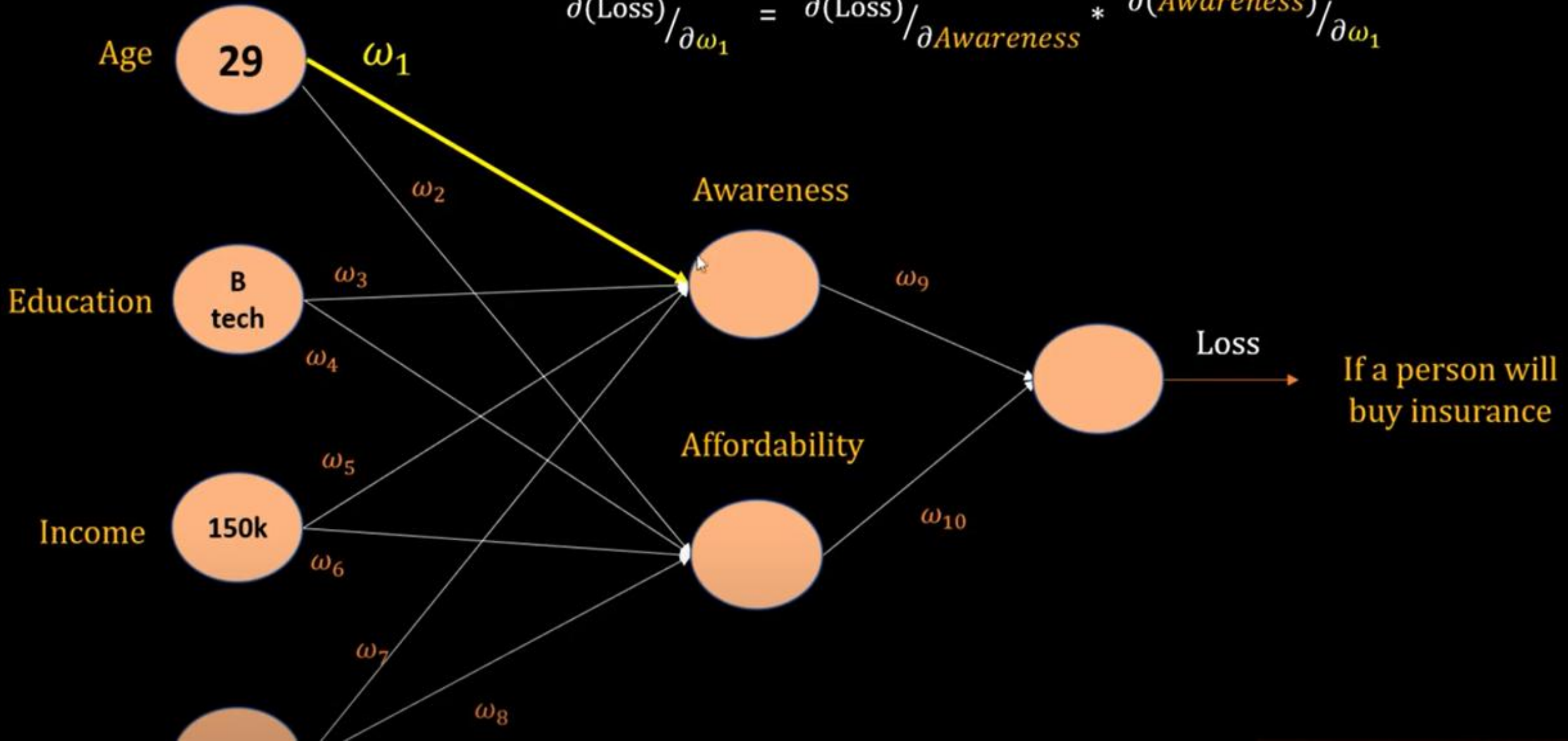








$$\frac{\partial(\text{Loss})}{\partial \omega_1} = \frac{\partial(\text{Loss})}{\partial \text{Awareness}} * \frac{\partial(\text{Awareness})}{\partial \omega_1}$$



$$\frac{\partial(\text{Loss})}{\partial \omega_1} = \frac{\partial(\text{Loss})}{\partial \text{Awareness}} * \frac{\partial(\text{Awareness})}{\partial \omega_1}$$

$$\frac{\partial(\text{Loss})}{\partial \omega_1} = \frac{\partial(\text{Loss})}{\partial \text{Awareness}} * \frac{\partial(\text{Awareness})}{\partial \omega_1}$$

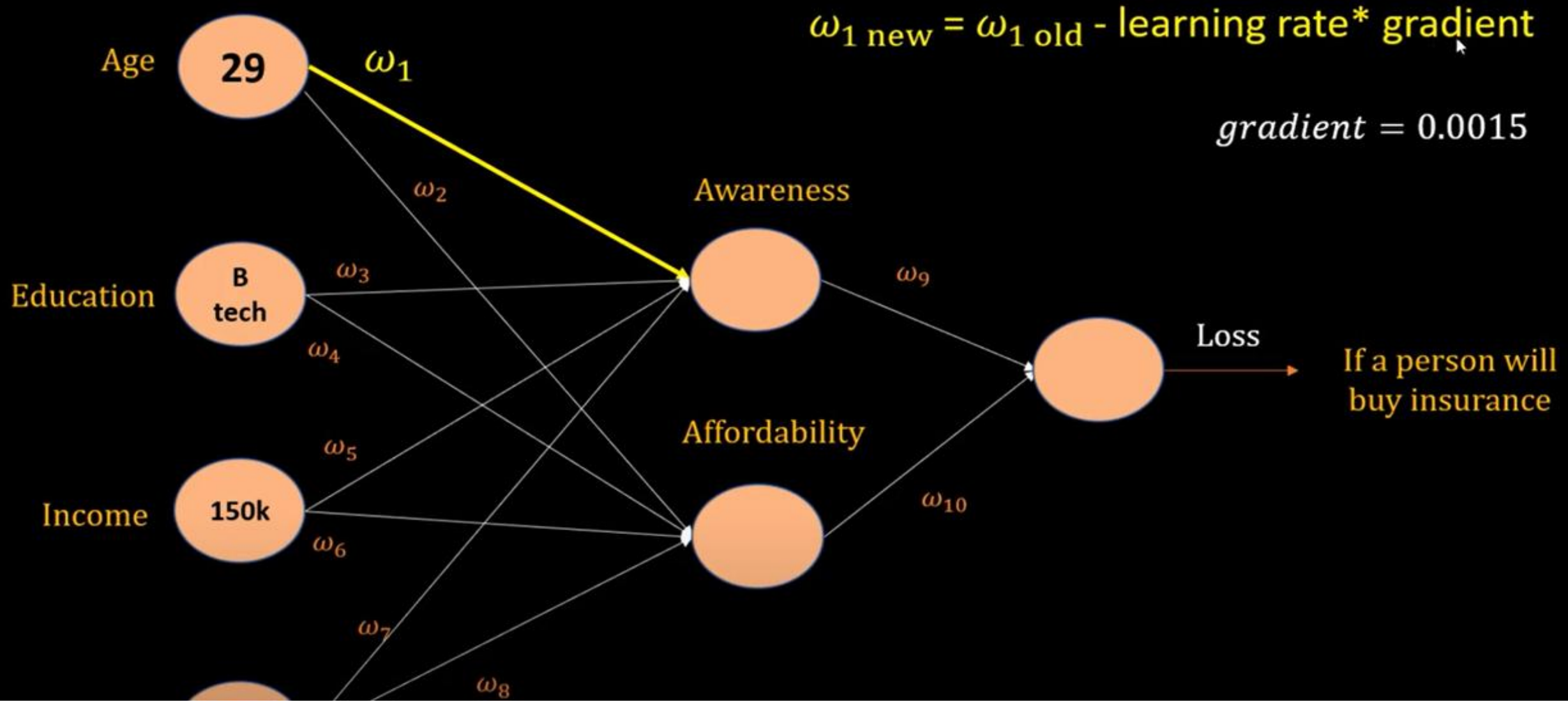
$$\text{gradient} = d1 * d2$$

$$\frac{\partial(\text{Loss})}{\partial \omega_1} = \frac{\partial(\text{Loss})}{\partial \text{Awareness}} * \frac{\partial(\text{Awareness})}{\partial \omega_1}$$

$$\text{gradient} = d1 * d2$$

$$\text{gradient} = 0.03 * 0.05$$

$$\text{gradient} = 0.0015$$



As number of hidden layers grow, gradient becomes very small and weights will hardly change . This will hamper the learning process.

Vanishing Gradients

When individual derivatives are large, the final derivative will also become huge and weights would change drastically.

Exploding Gradients

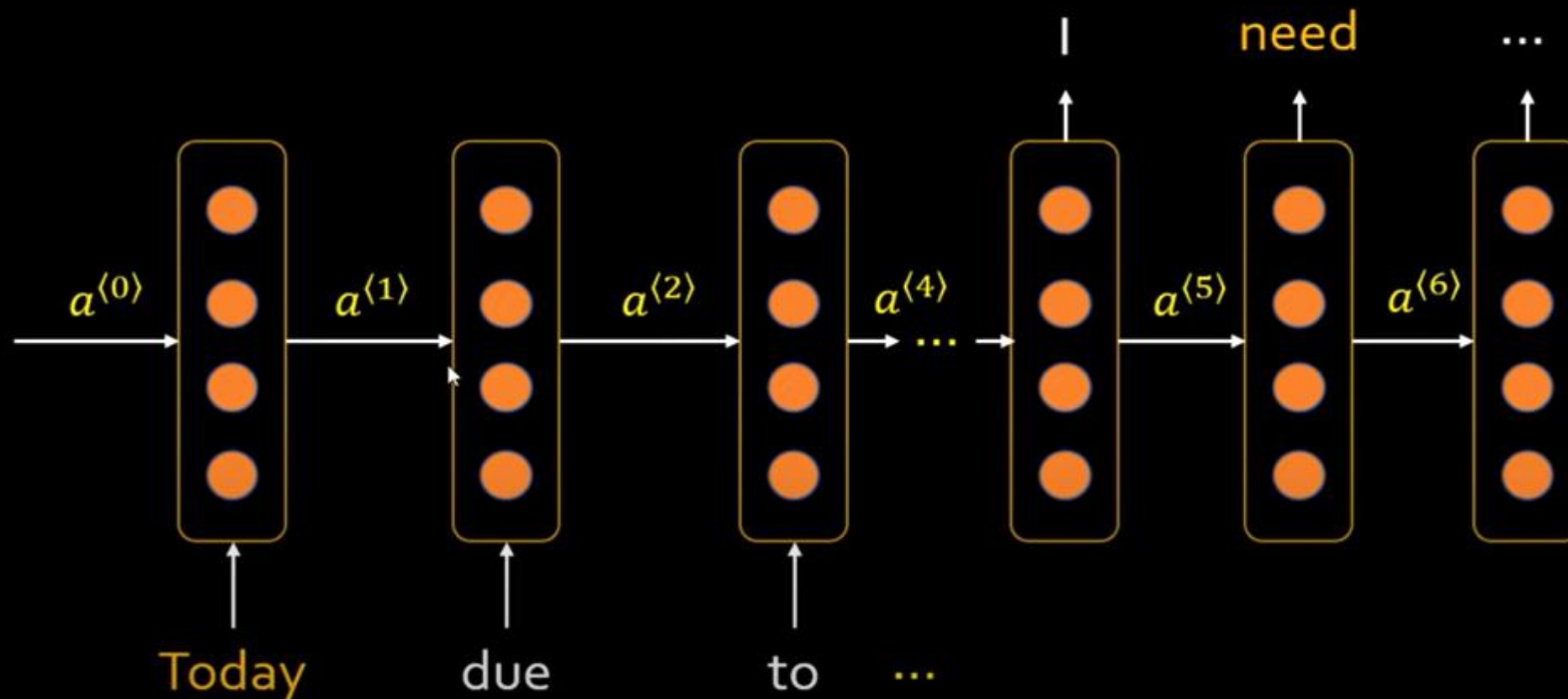
$$\text{gradient} = d1 * d2 * d3 * d4 * \dots * dn$$

Vanishing gradient problem is more prominent in very deep neural networks.

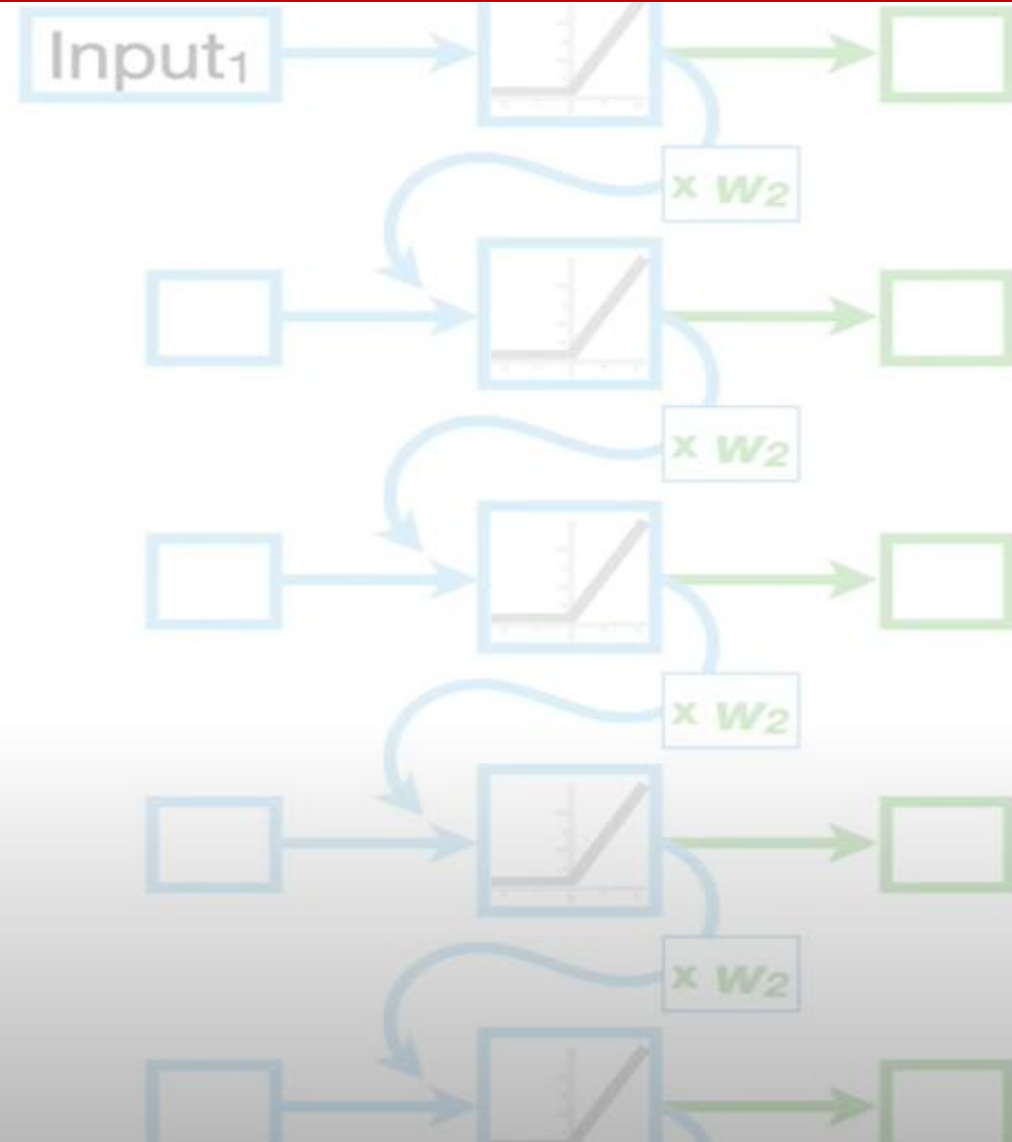
Today, due to my current job situation and family conditions, I need to take a loan.

Last year, due to my current job situation and family conditions, I had to take a loan.

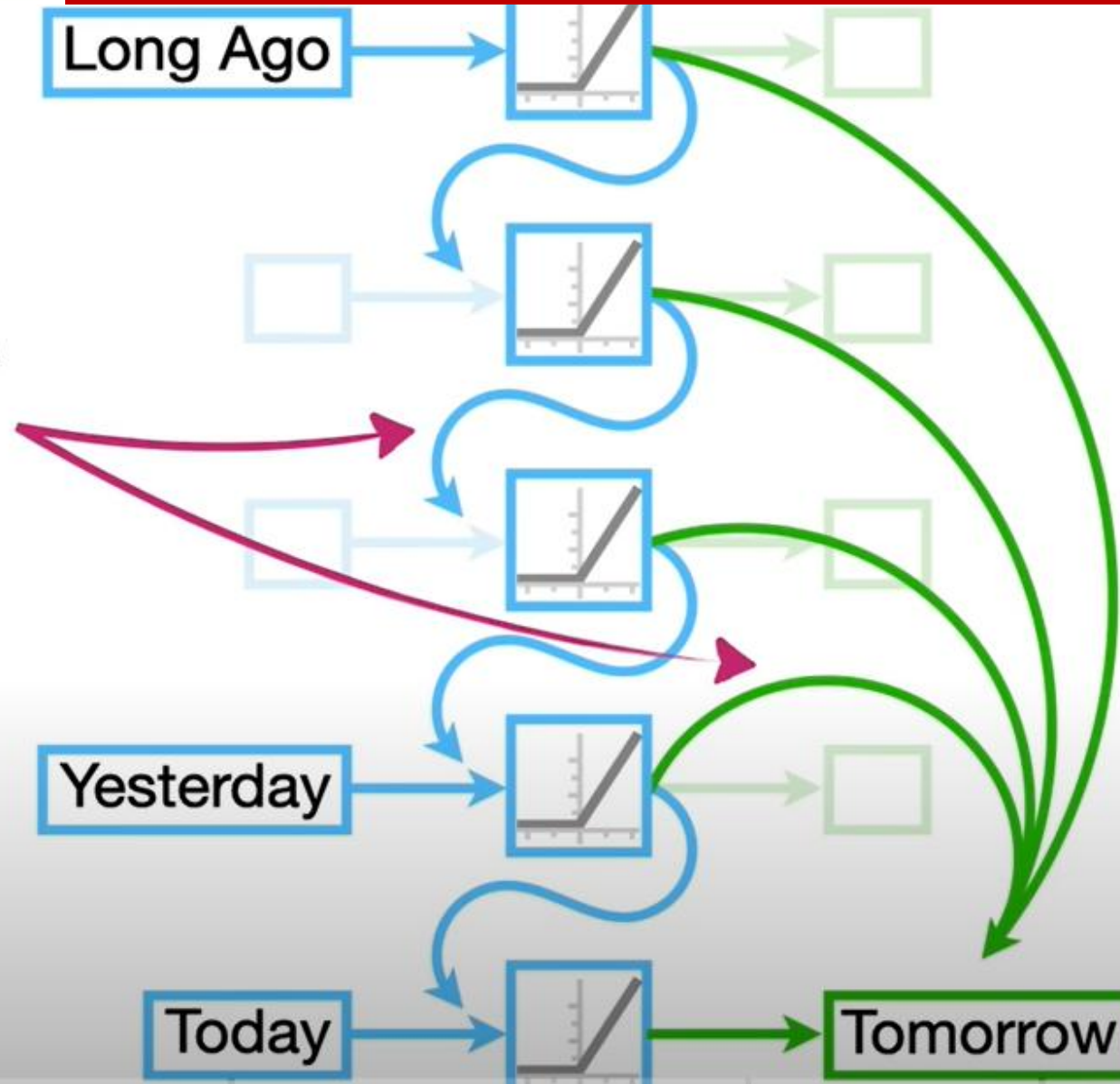
Today, due to my current job situation and family conditions, I need to take a loan.

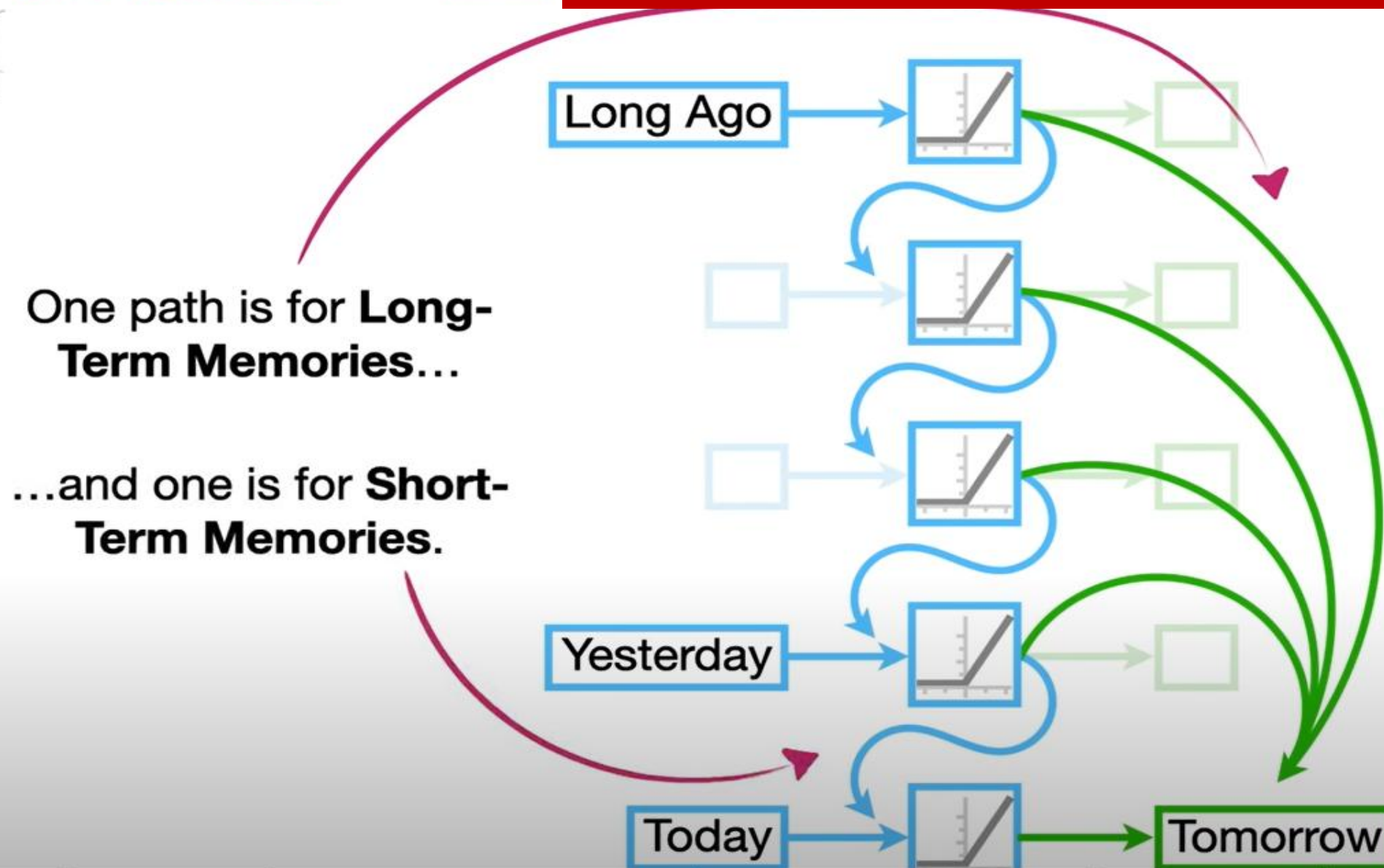


So, today, we're going to talk about **Long Short-Term Memory (LSTM)**, which is a type of **Recurrent Neural Network** that is designed to avoid the **exploding/vanishing** gradient problem.



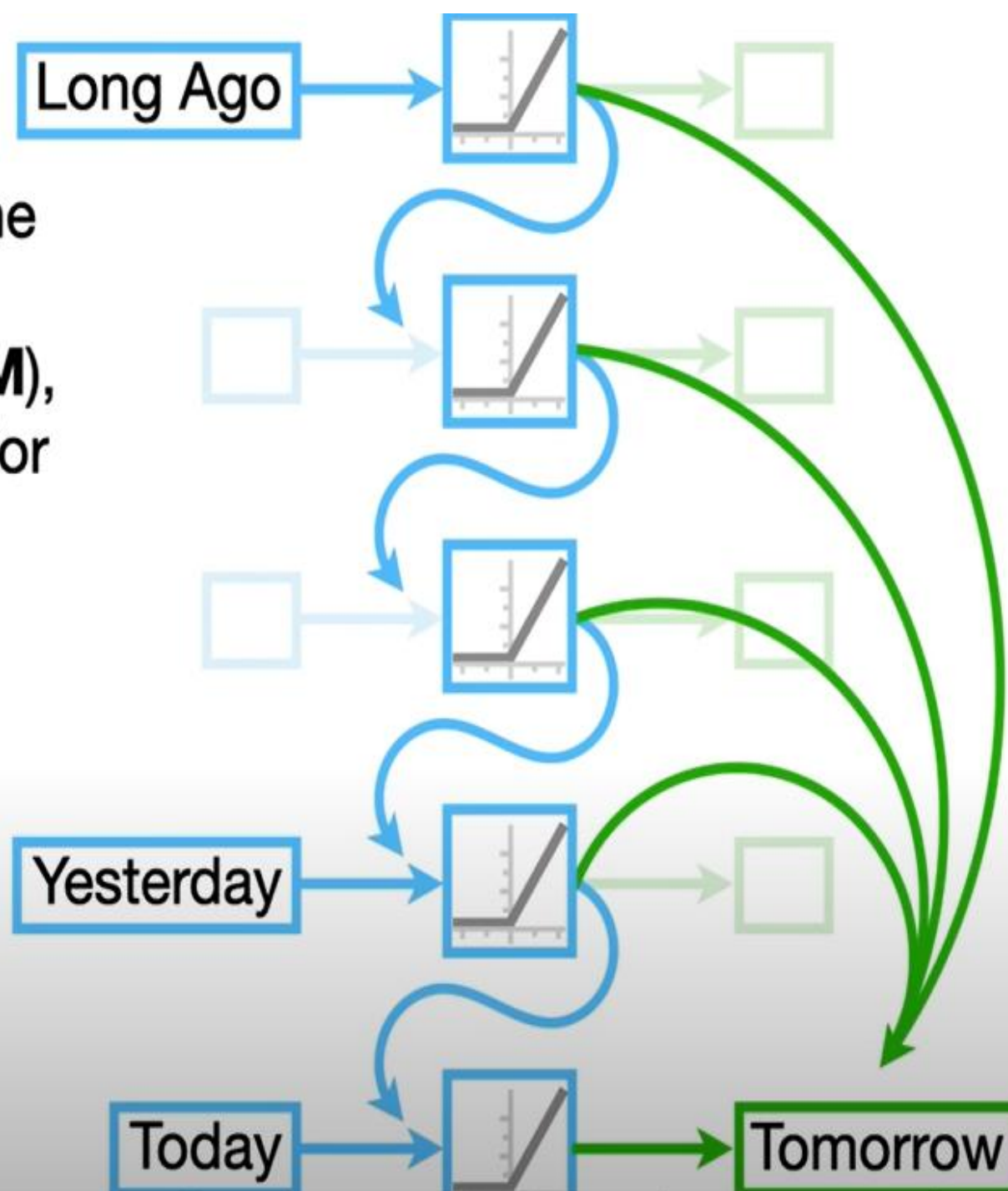
...**Long Short-Term Memory** uses two separate paths to make predictions about tomorrow.

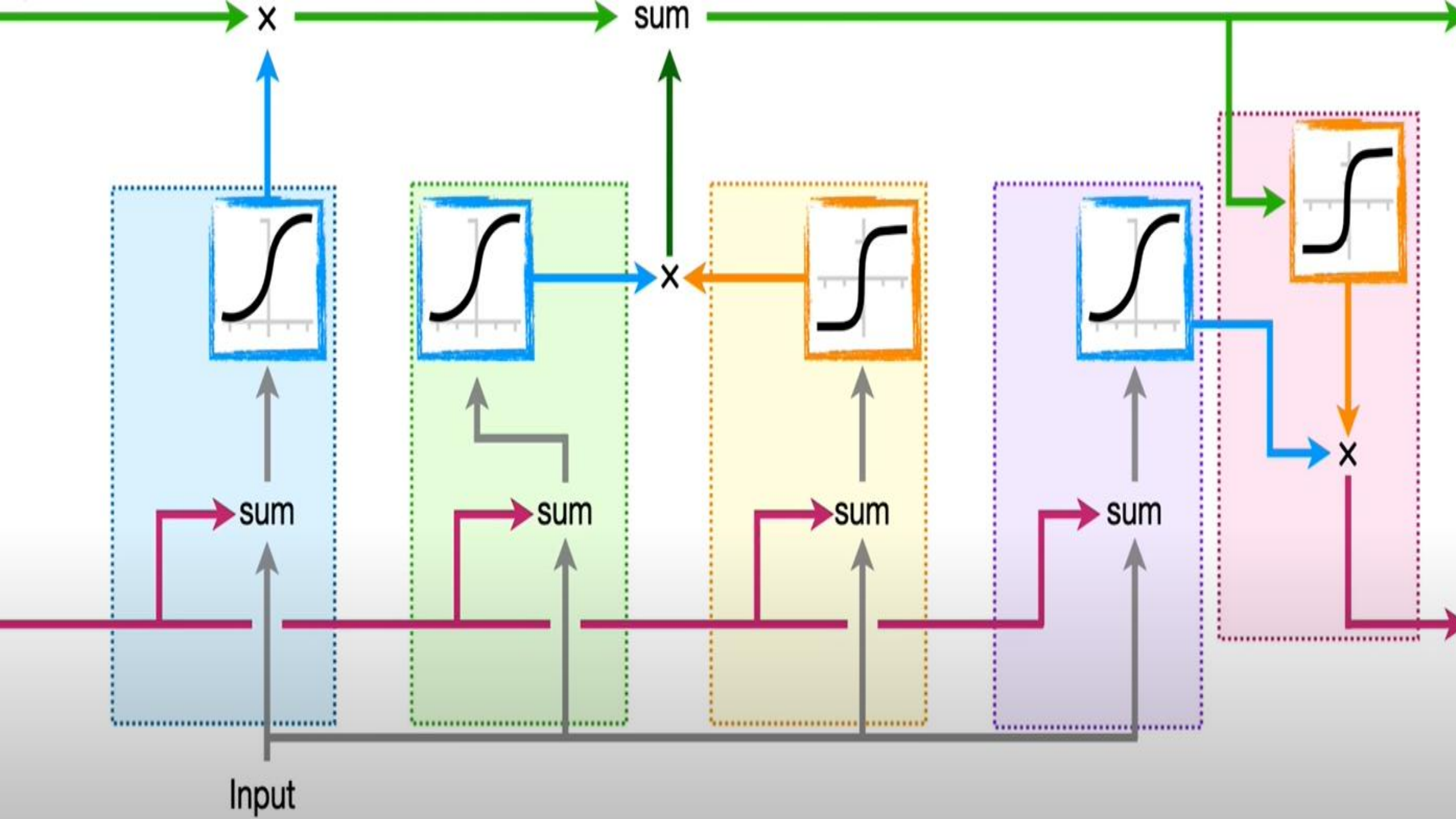




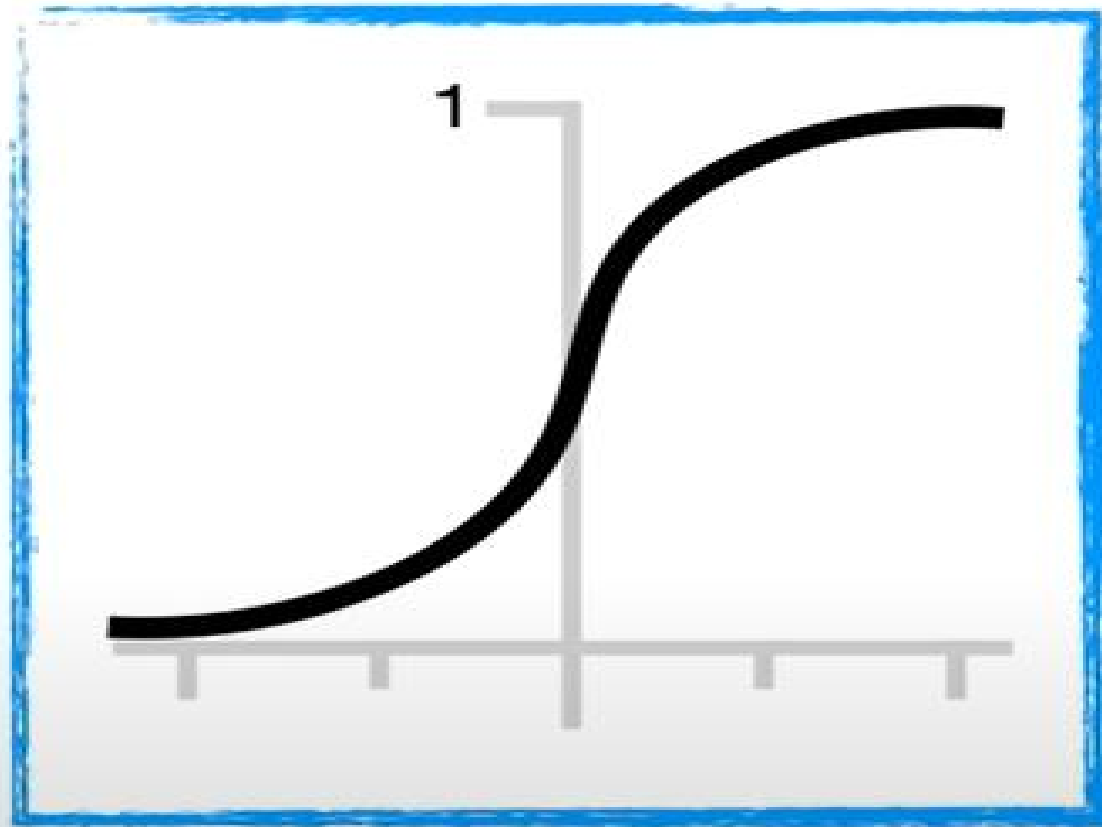


Now that we understand the
main idea behind **Long
Short-Term Memory (LSTM)**,
that it uses different paths for
Long and **Short-Term**
Memories...

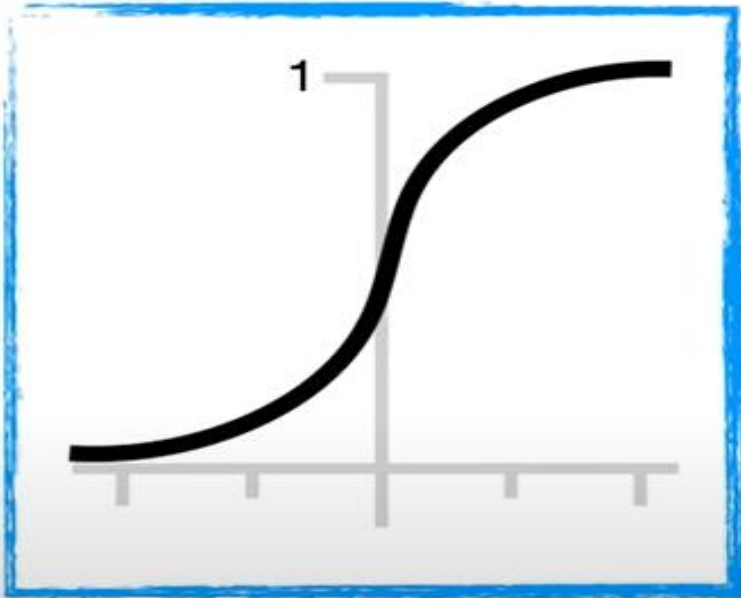




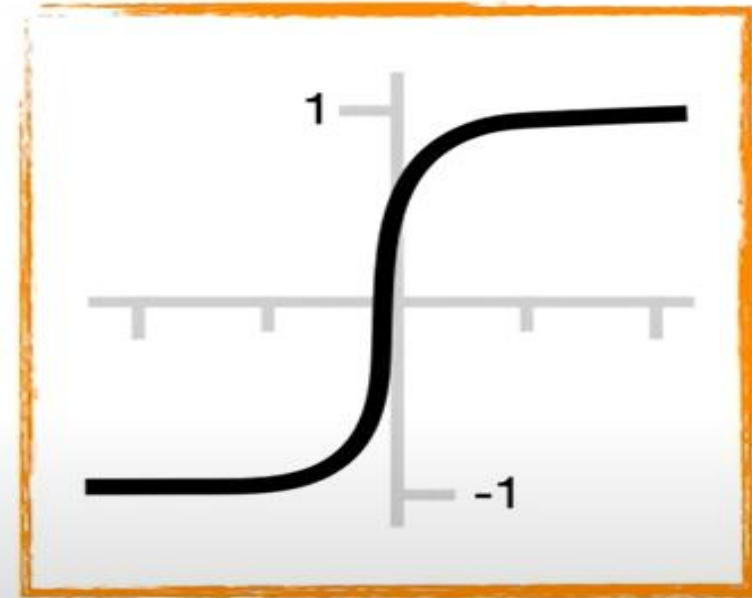
In a nutshell, the **Sigmoid Activation Function** takes any x-axis coordinate...

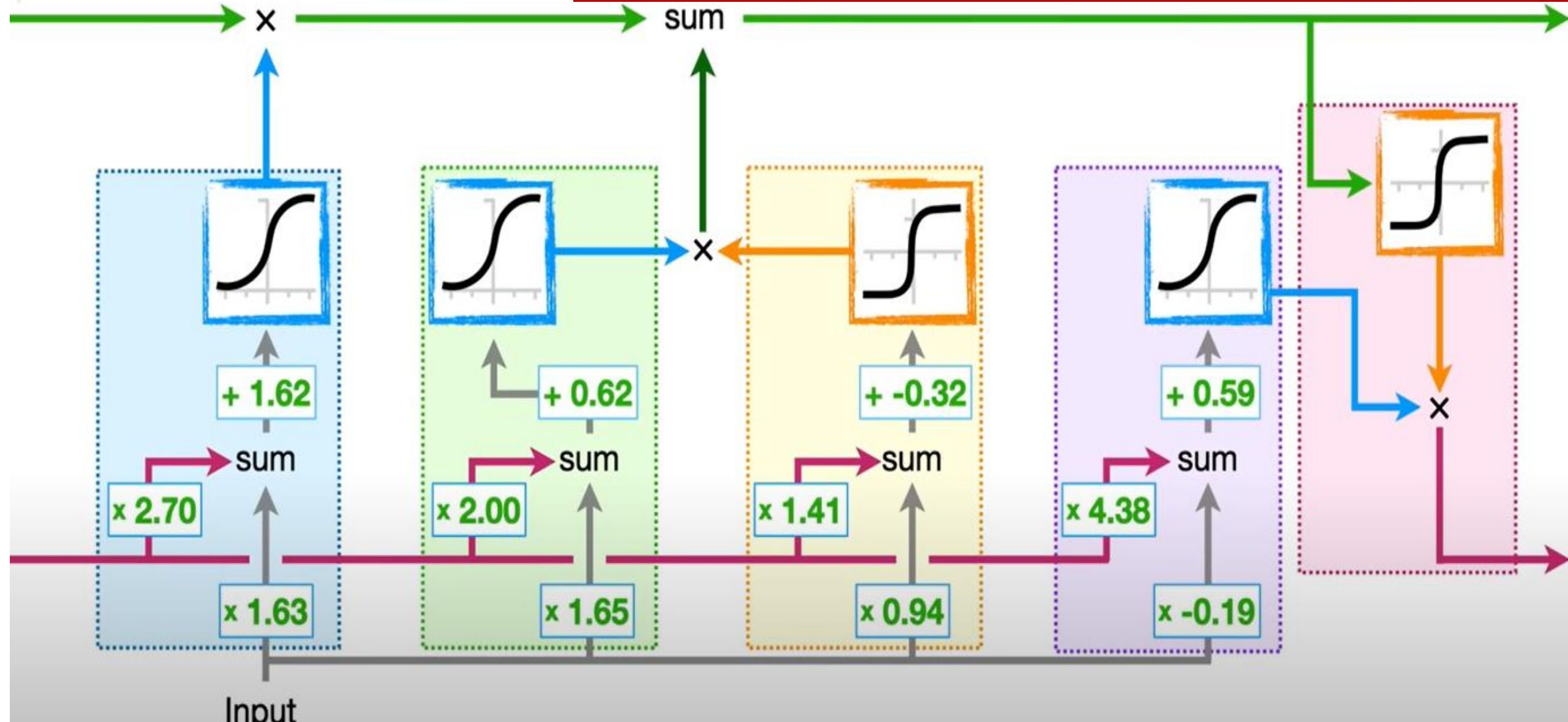


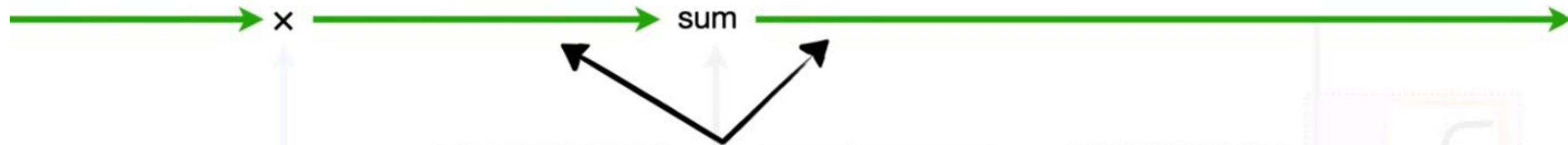
So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...



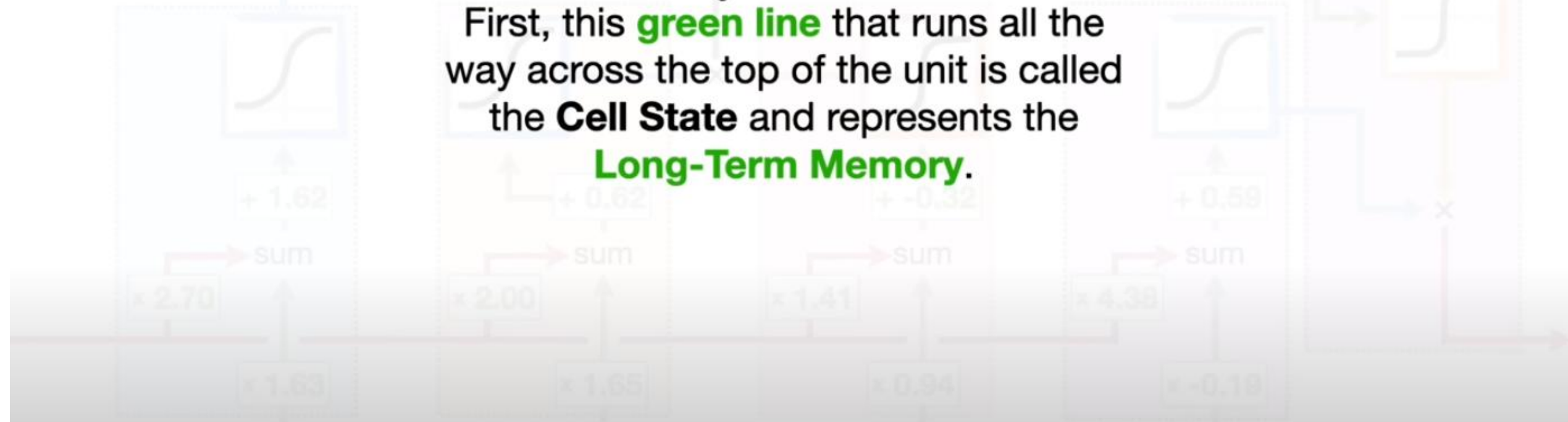
...and the **Tanh Activation Function** turns any input into a number between **-1** and **1**...



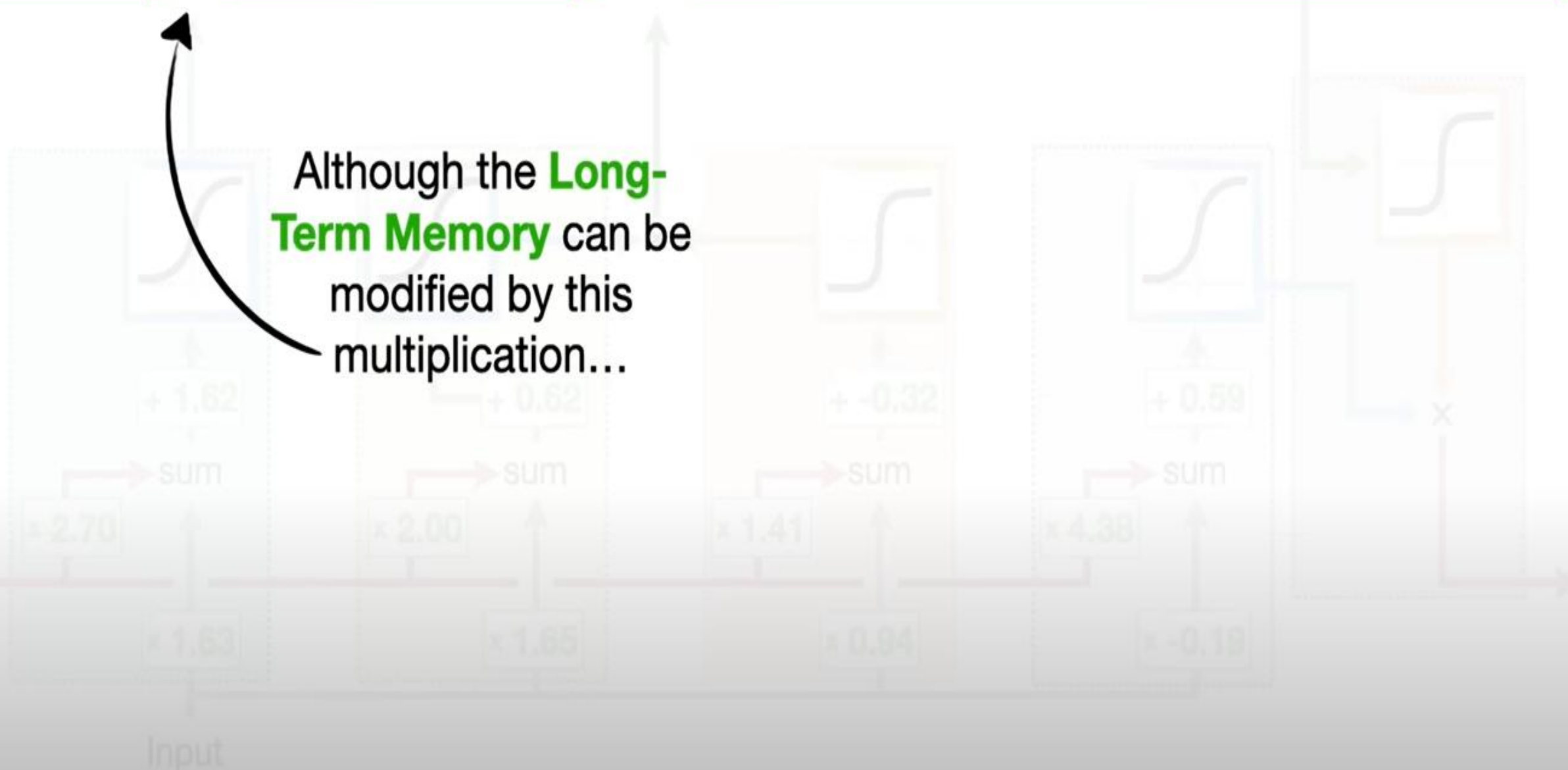




First, this **green line** that runs all the way across the top of the unit is called the **Cell State** and represents the **Long-Term Memory**.



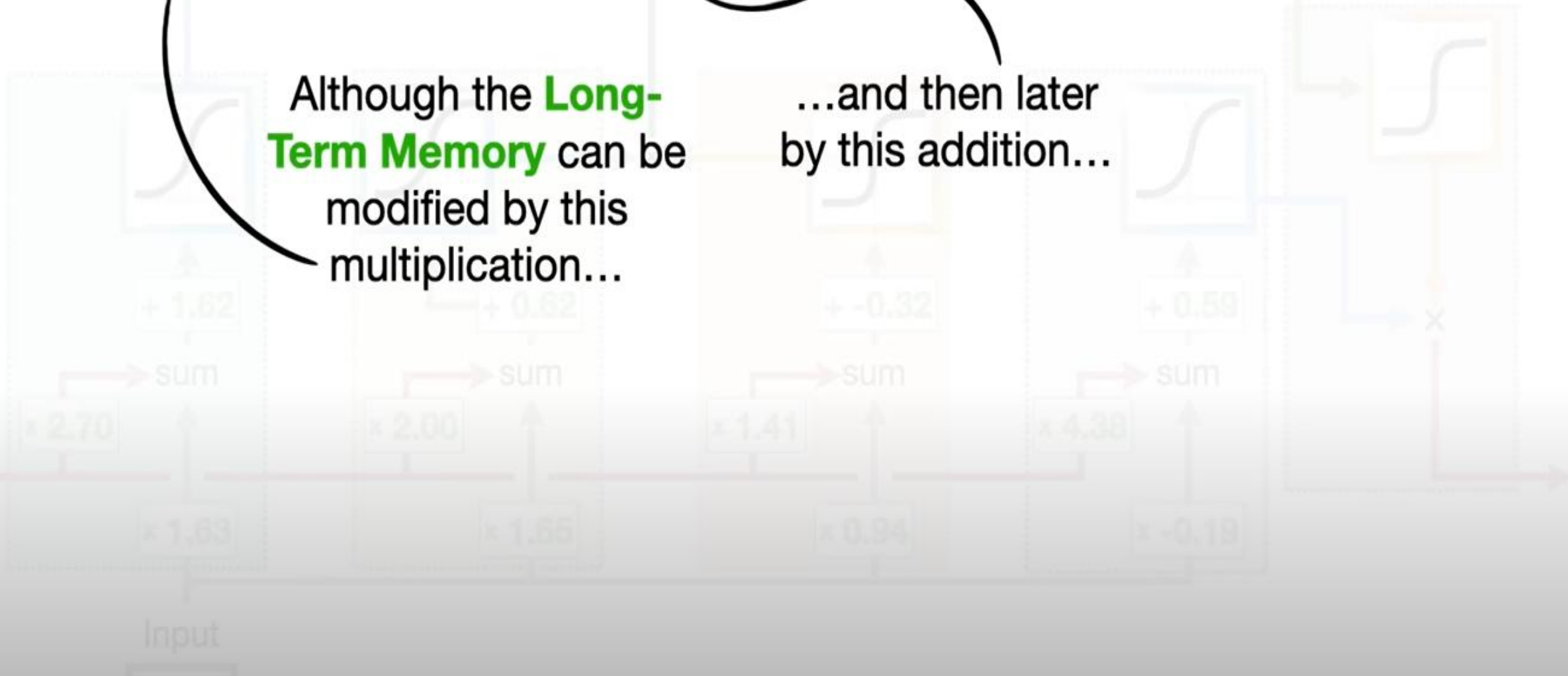
Although the **Long-Term Memory** can be modified by this multiplication...

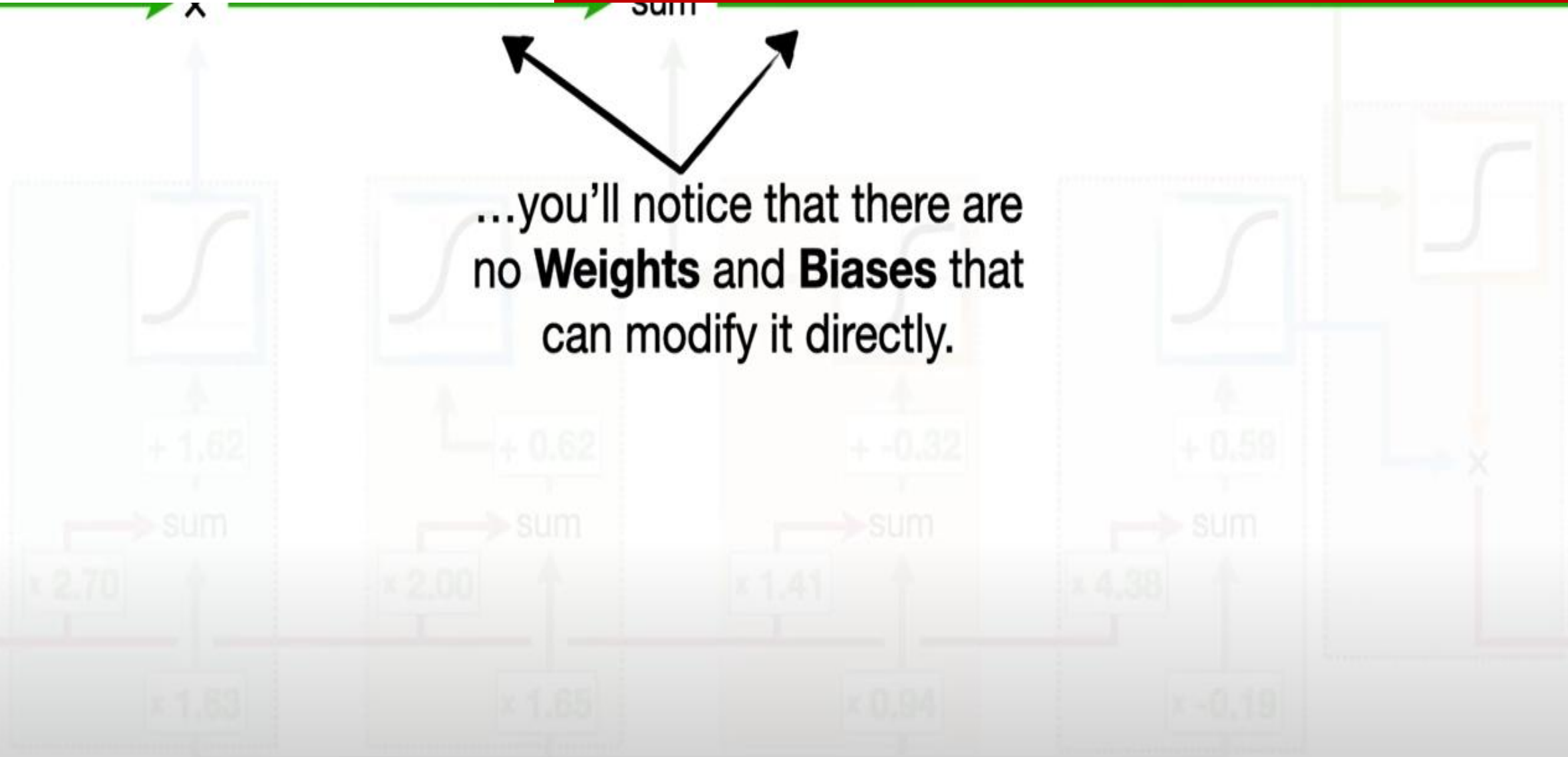




Although the **Long-Term Memory** can be modified by this multiplication...

...and then later by this addition...

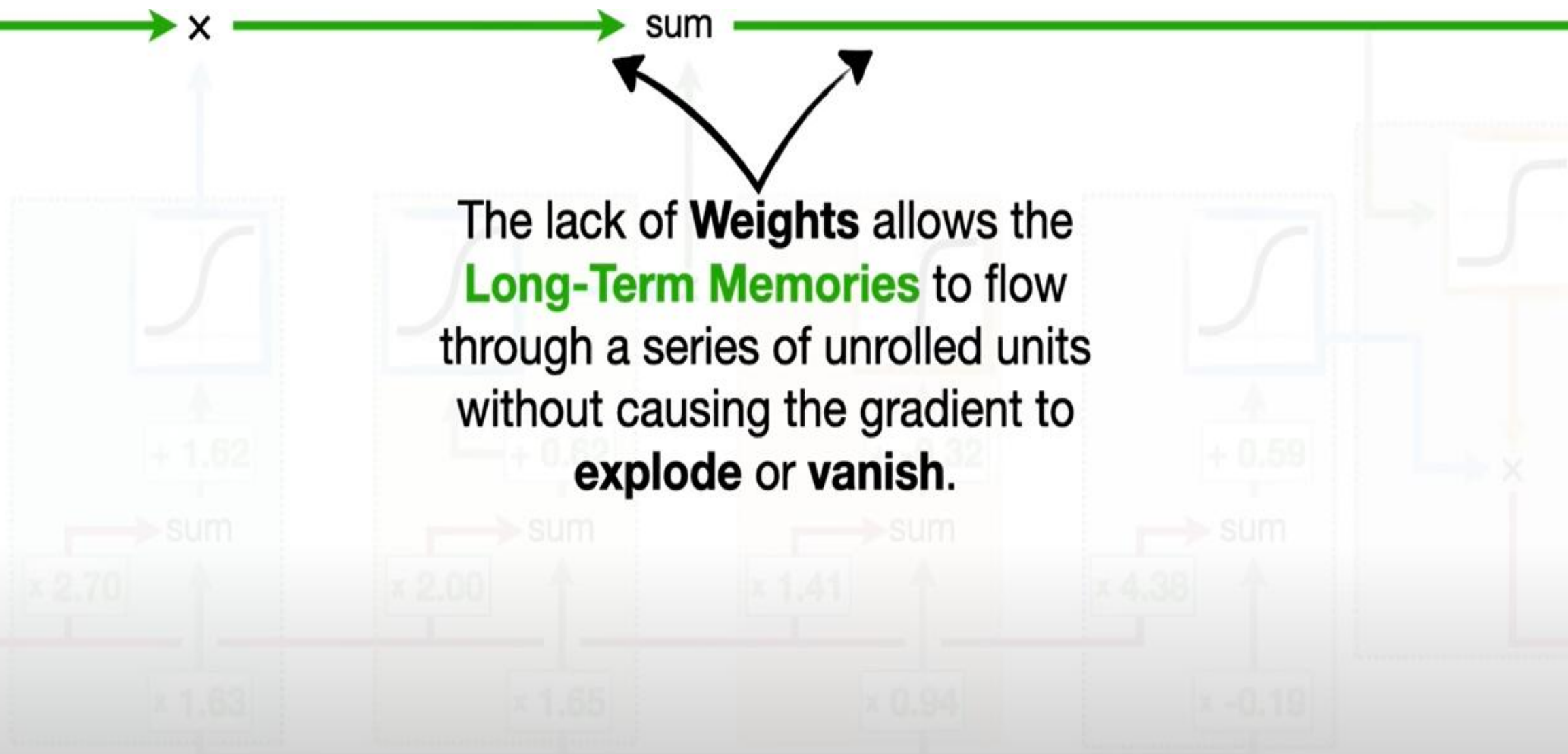




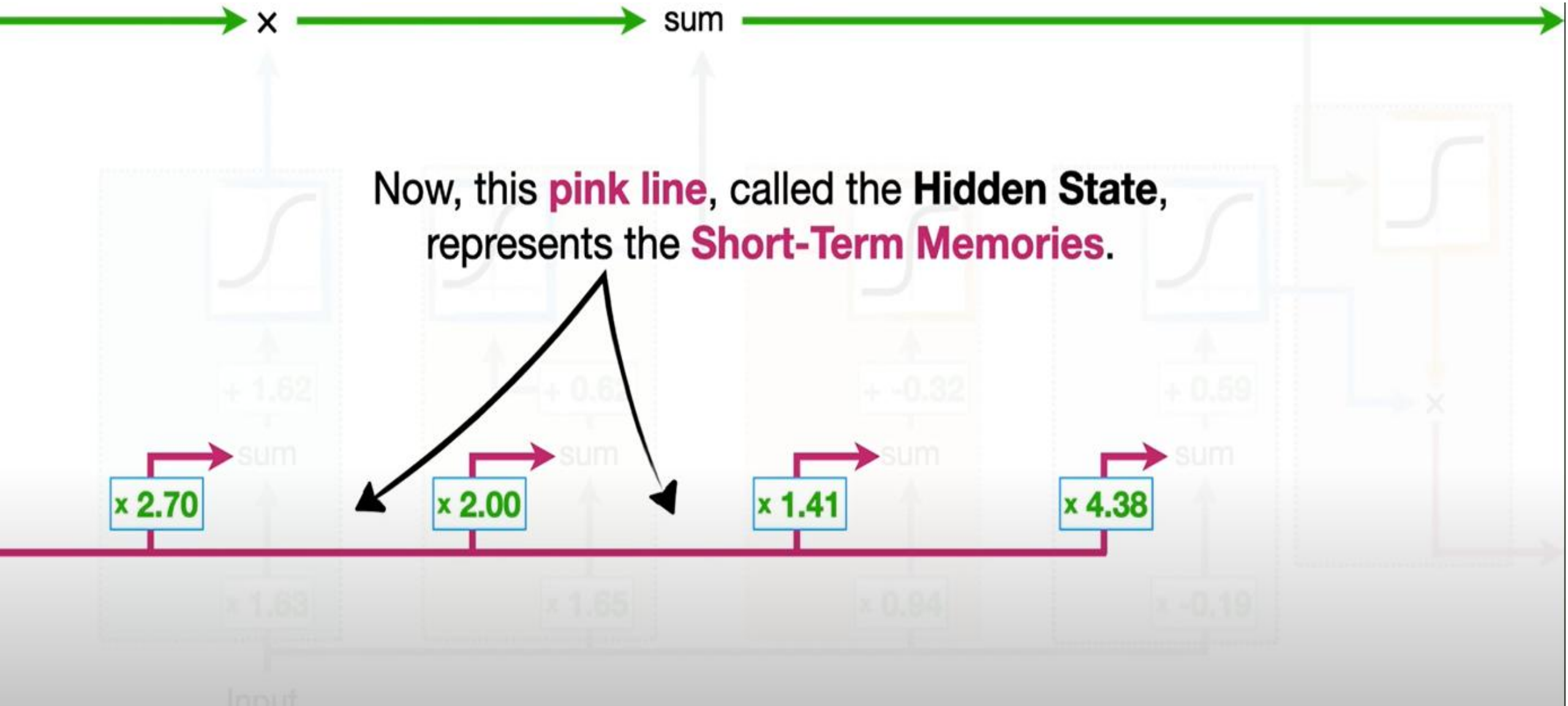
...you'll notice that there are
no **Weights** and **Biases** that
can modify it directly.

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**.





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



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

sum

x 2.00

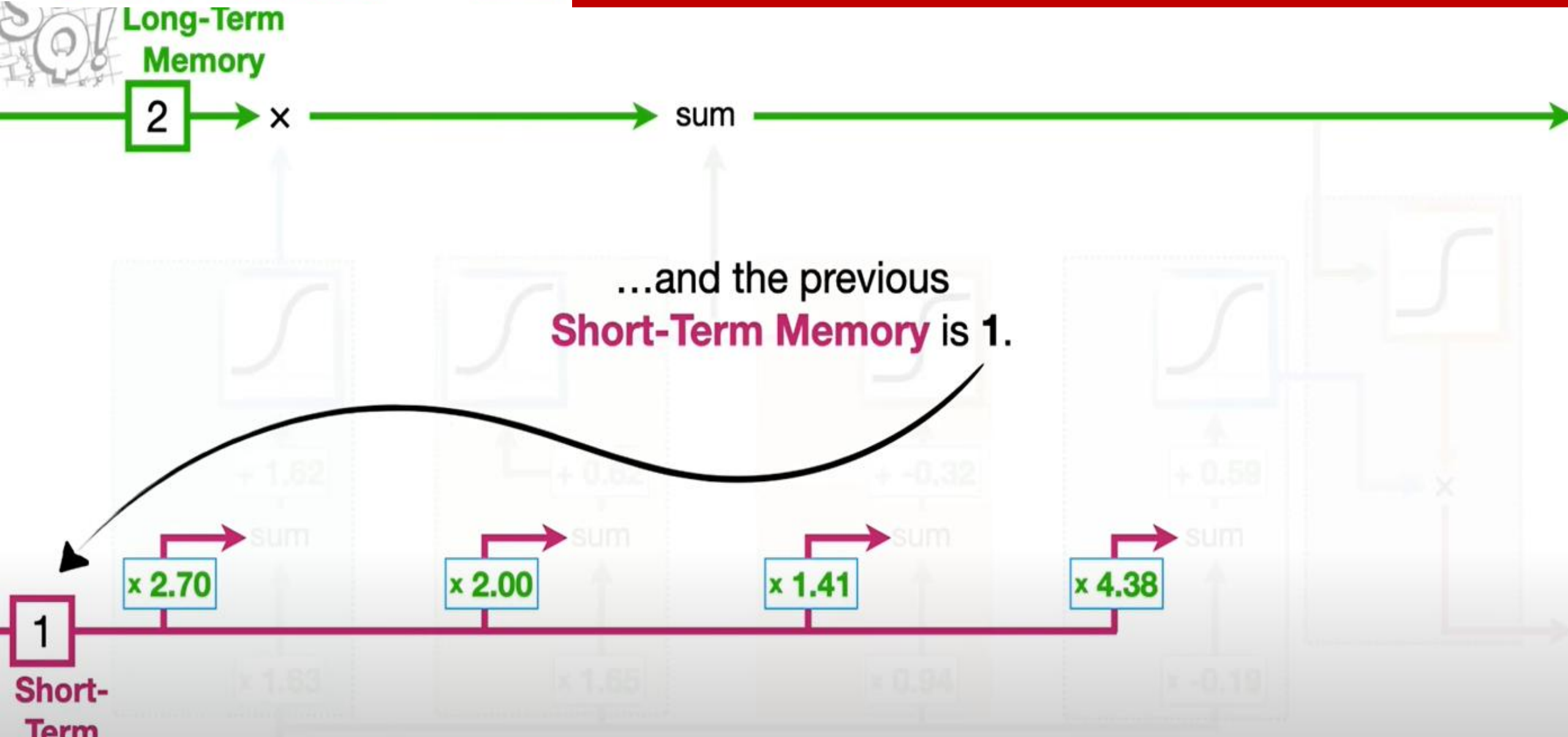
sum

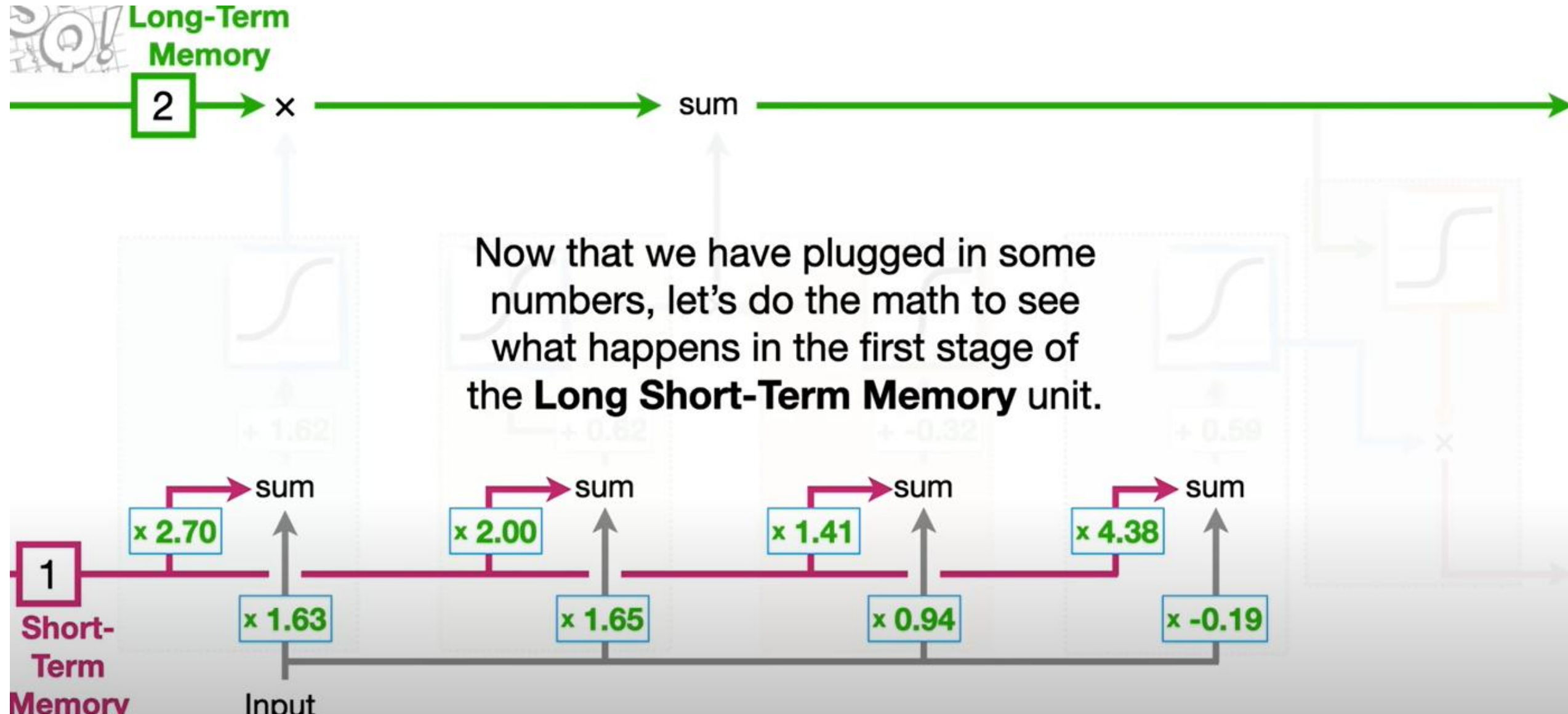
x 1.41

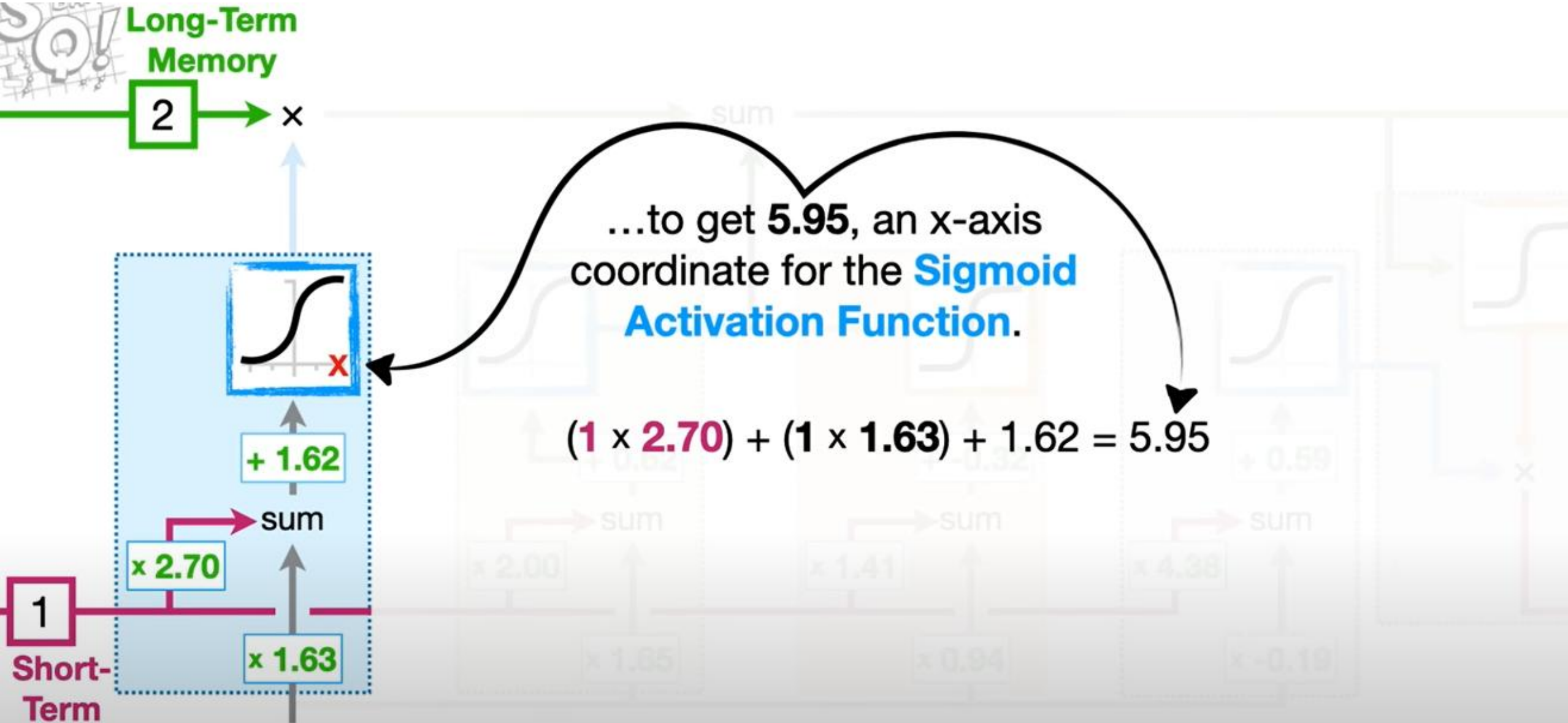
sum

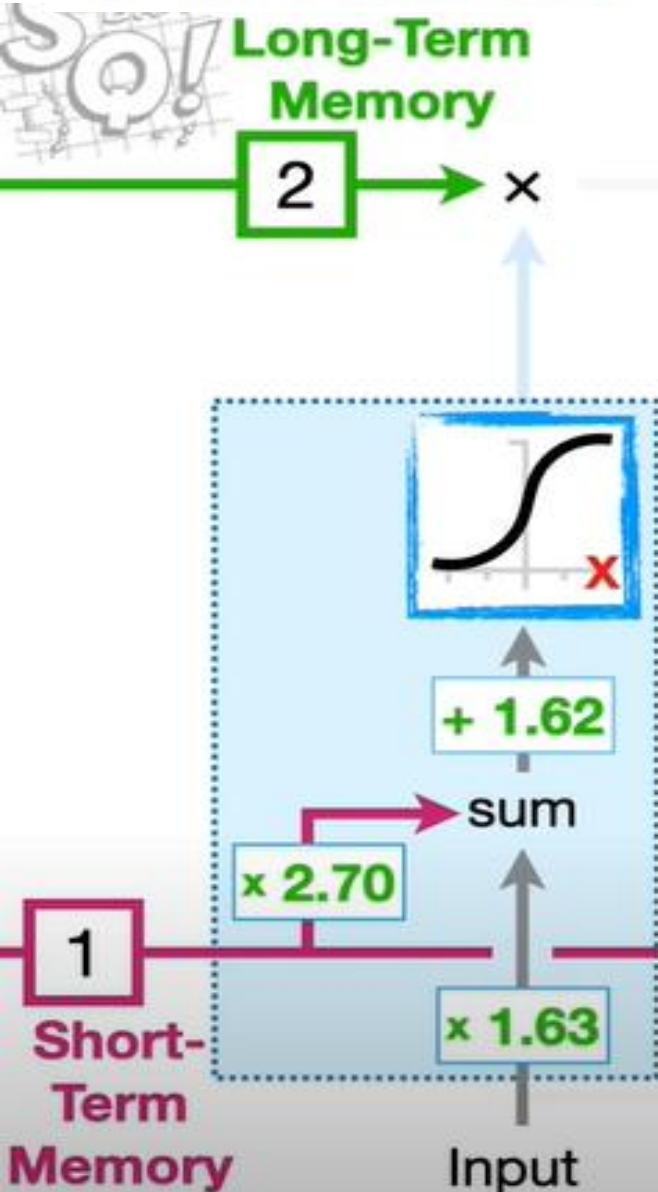
x 4.38

sum







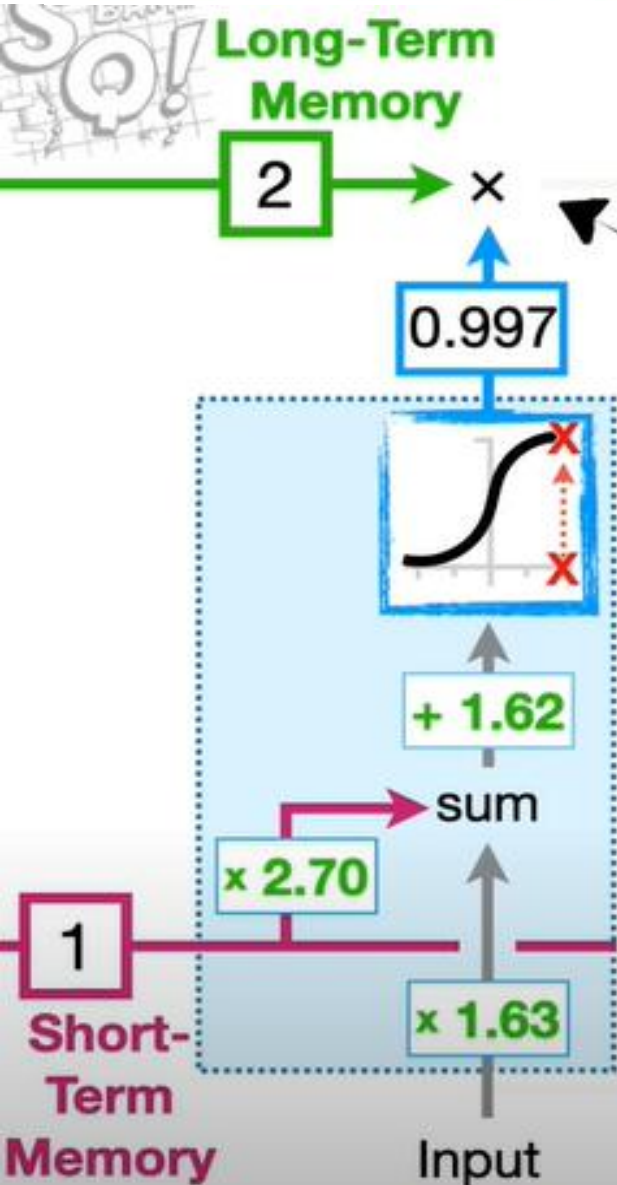


Now we plug the x-axis coordinate into the equation for the **Sigmoid Activation Function...**

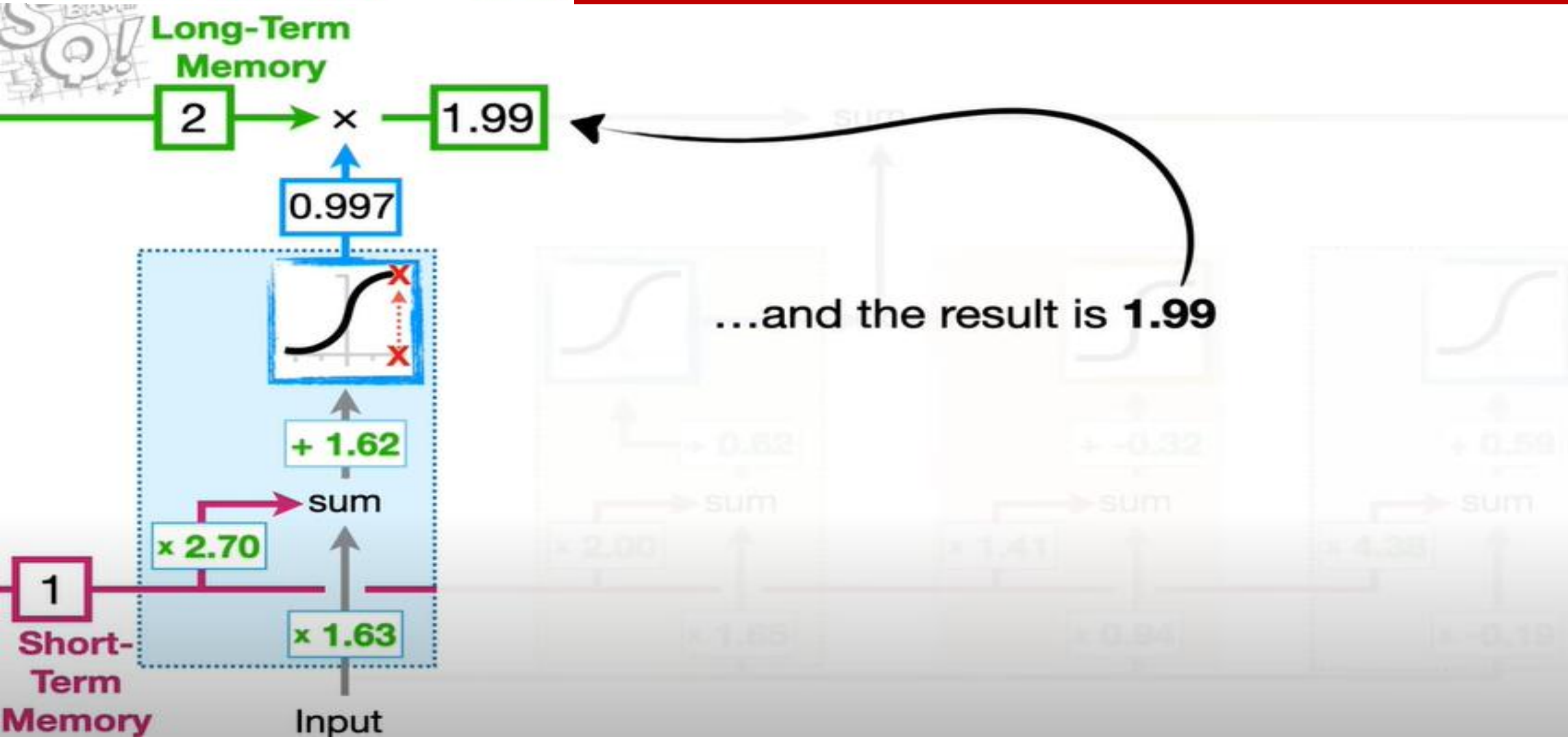
$$(1 \times 2.70) + (1 \times 1.63) + 1.62 = 5.95$$

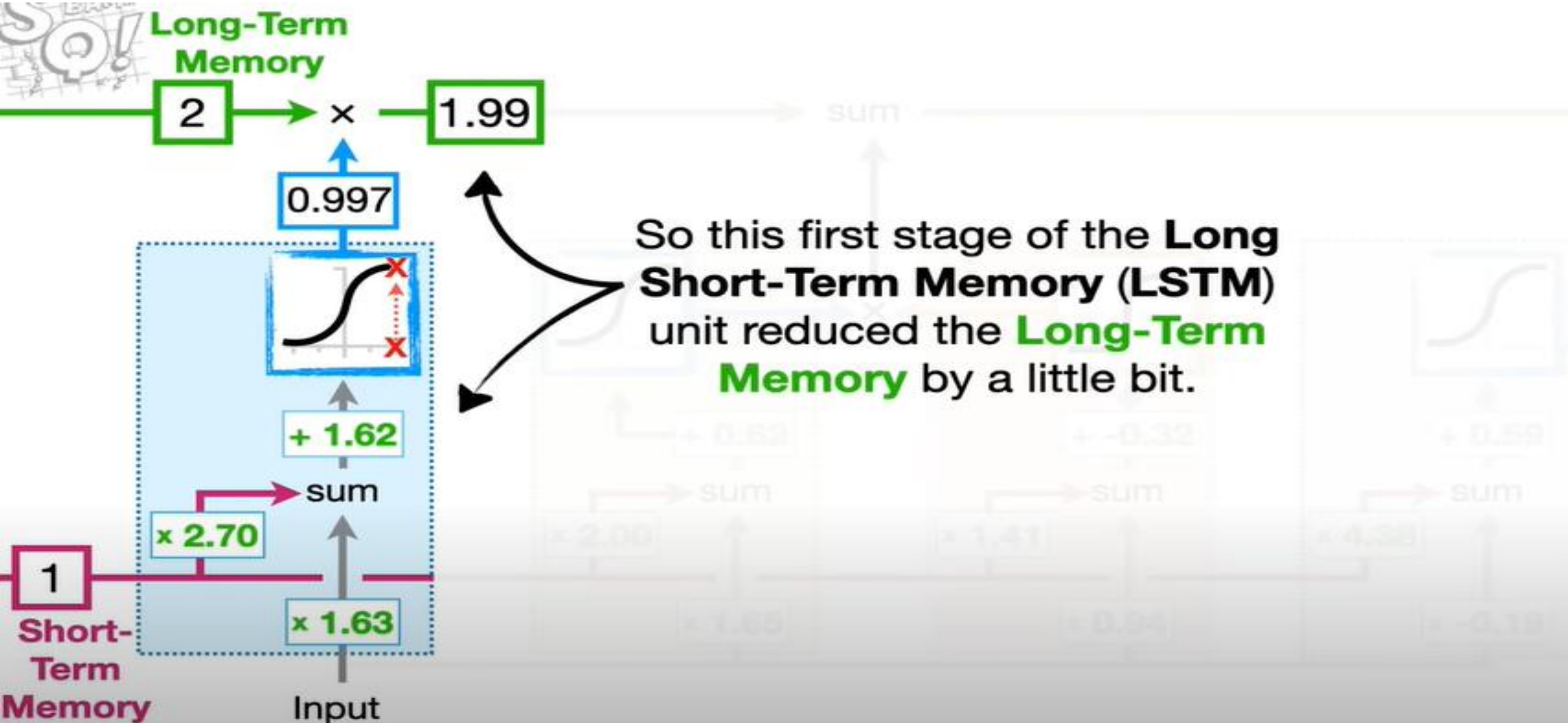
$$f(x) = \frac{e^x}{e^x + 1}$$

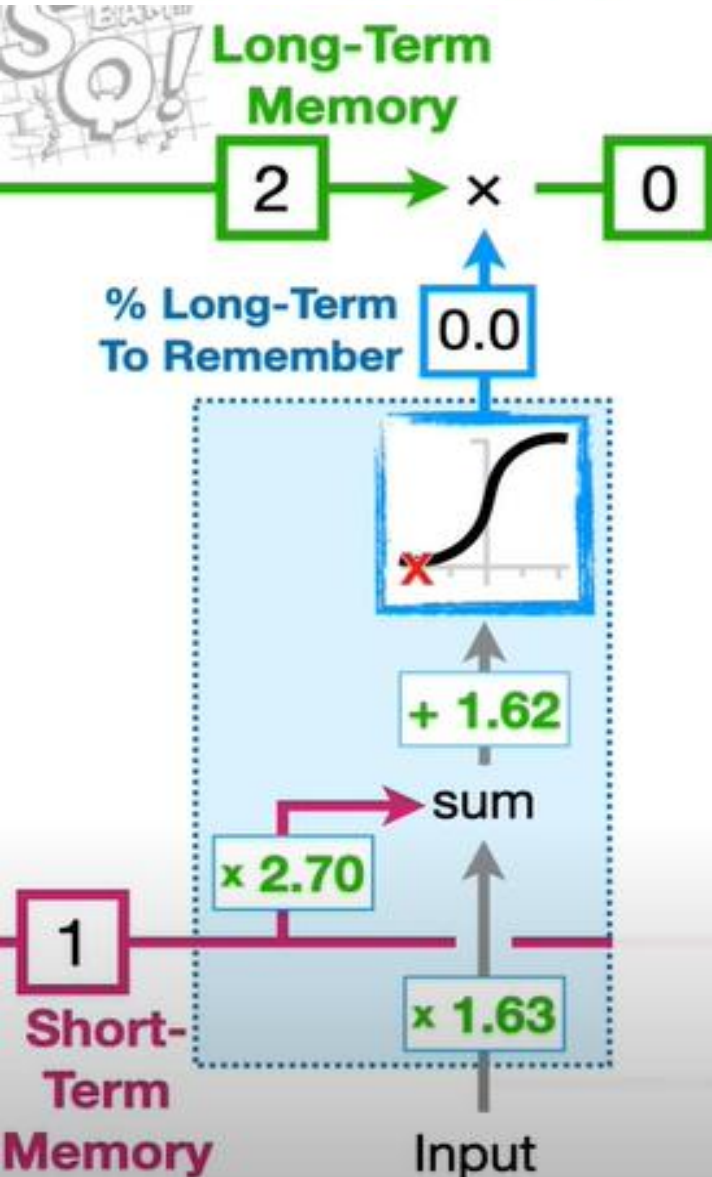
$$f(5.95) = \frac{e^{5.95}}{e^{5.95} + 1}$$



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



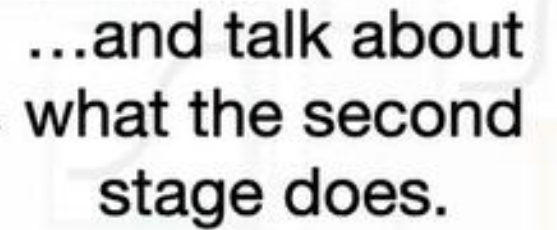


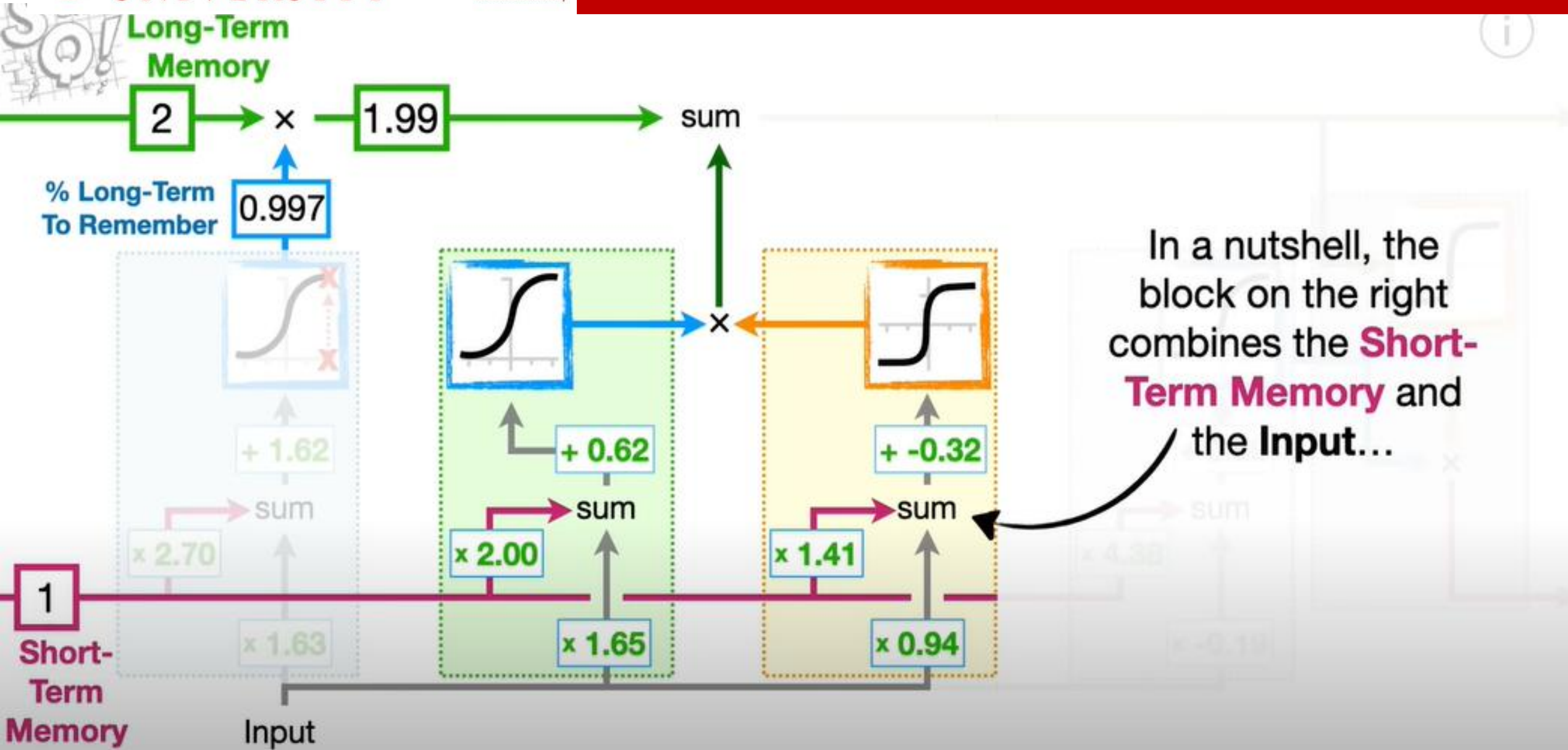


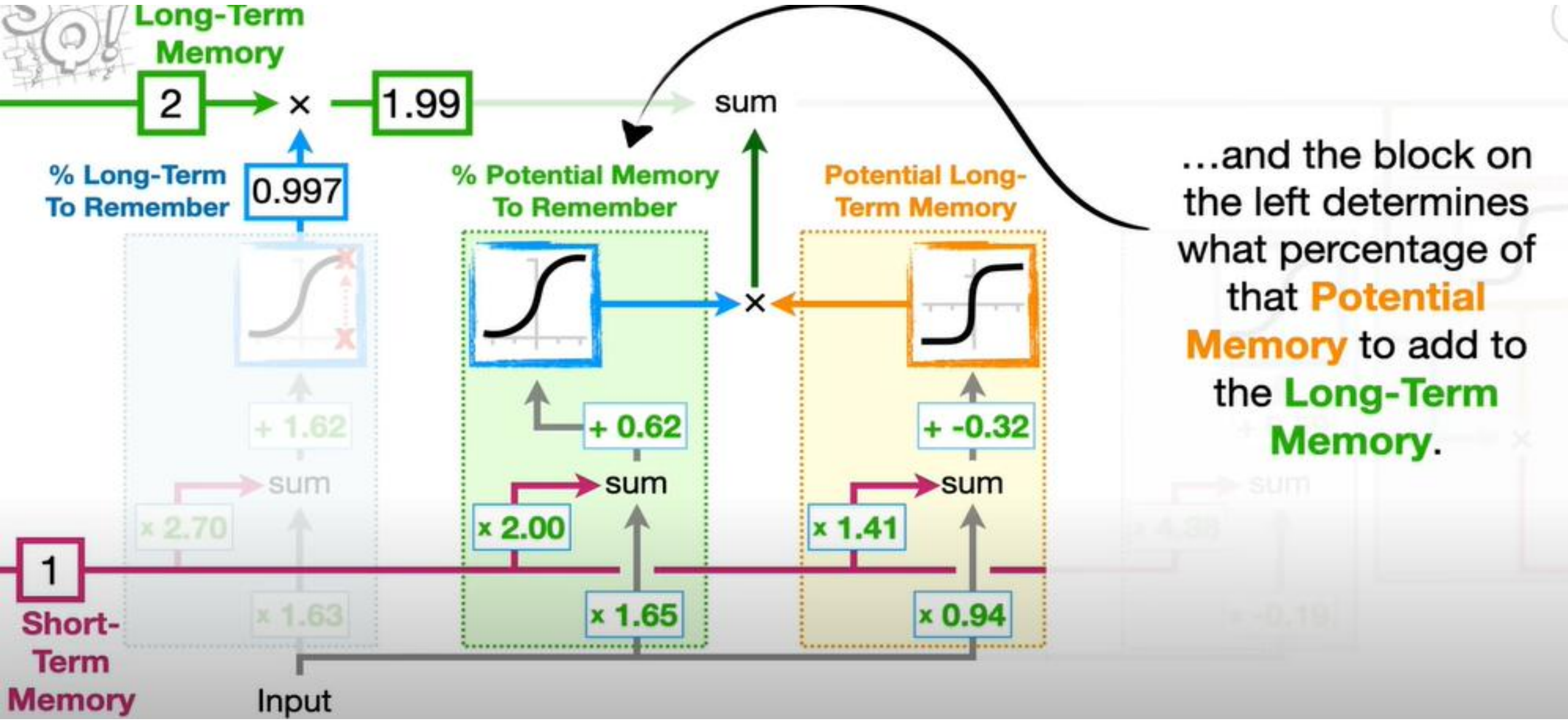
TERMINOLOGY ALERT!!!

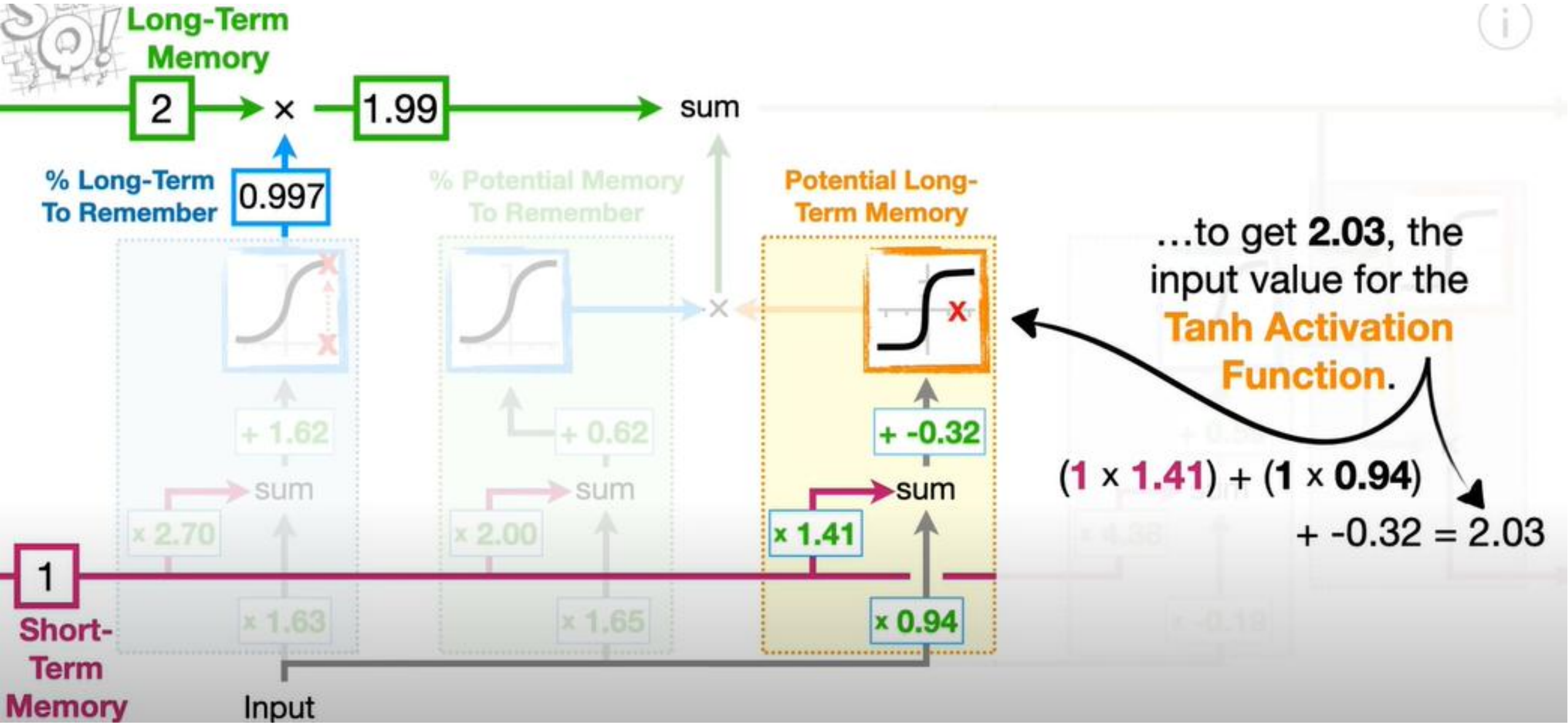
Even though this part of the **Long Short-Term Memory** unit determines what percentage of the **Long-Term Memory** will be remembered...

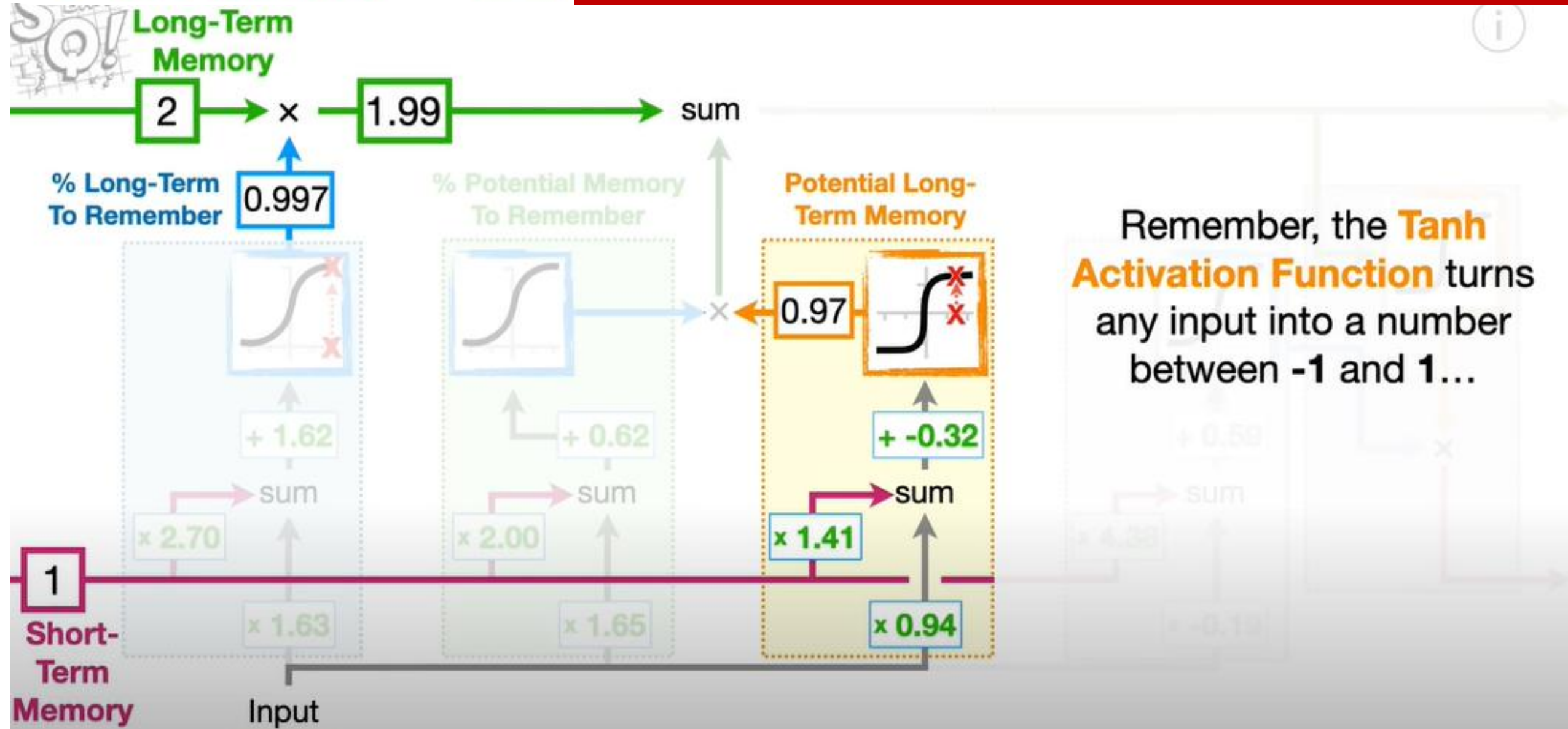
...it is usually called the **Forget Gate**.

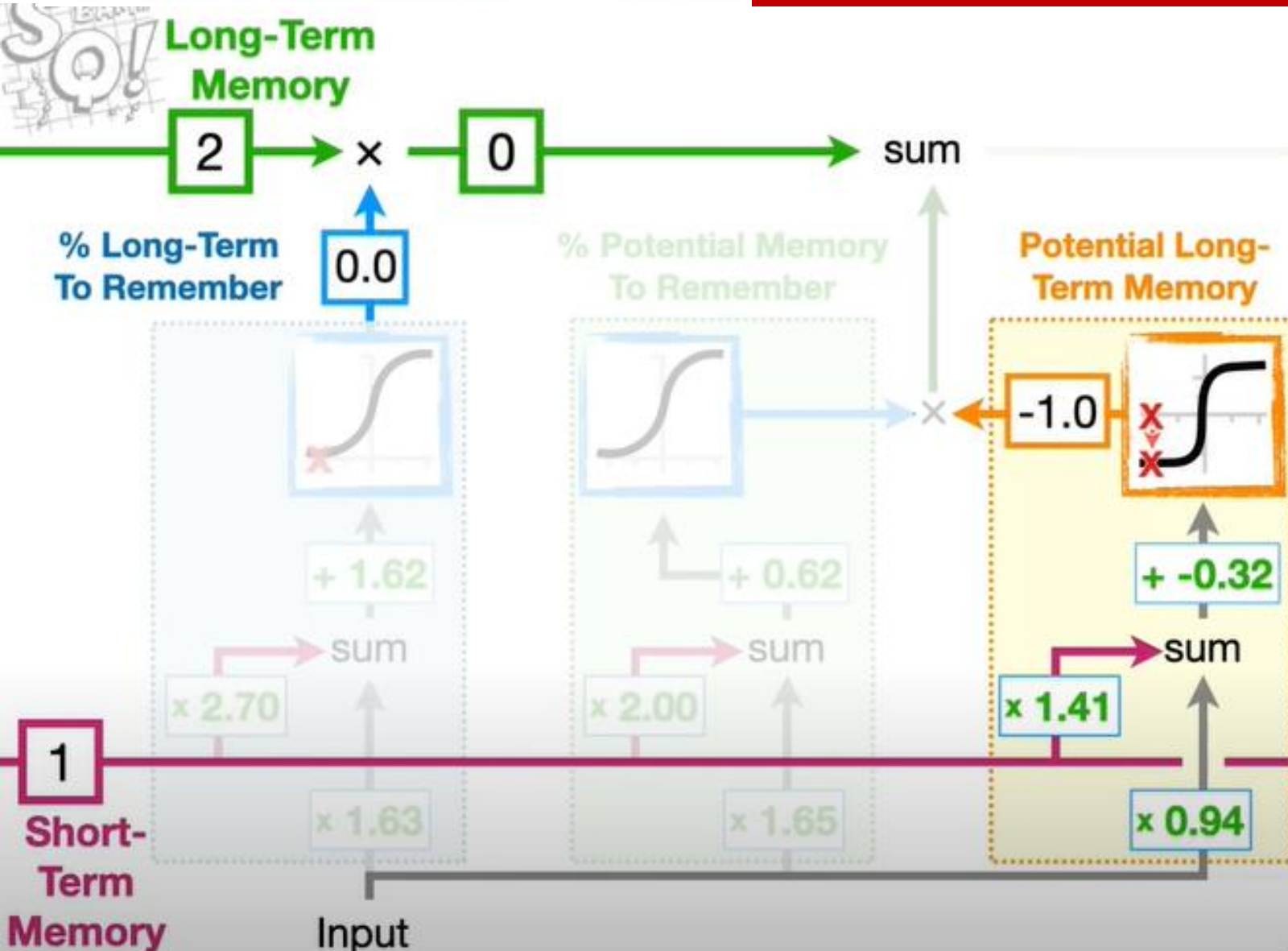






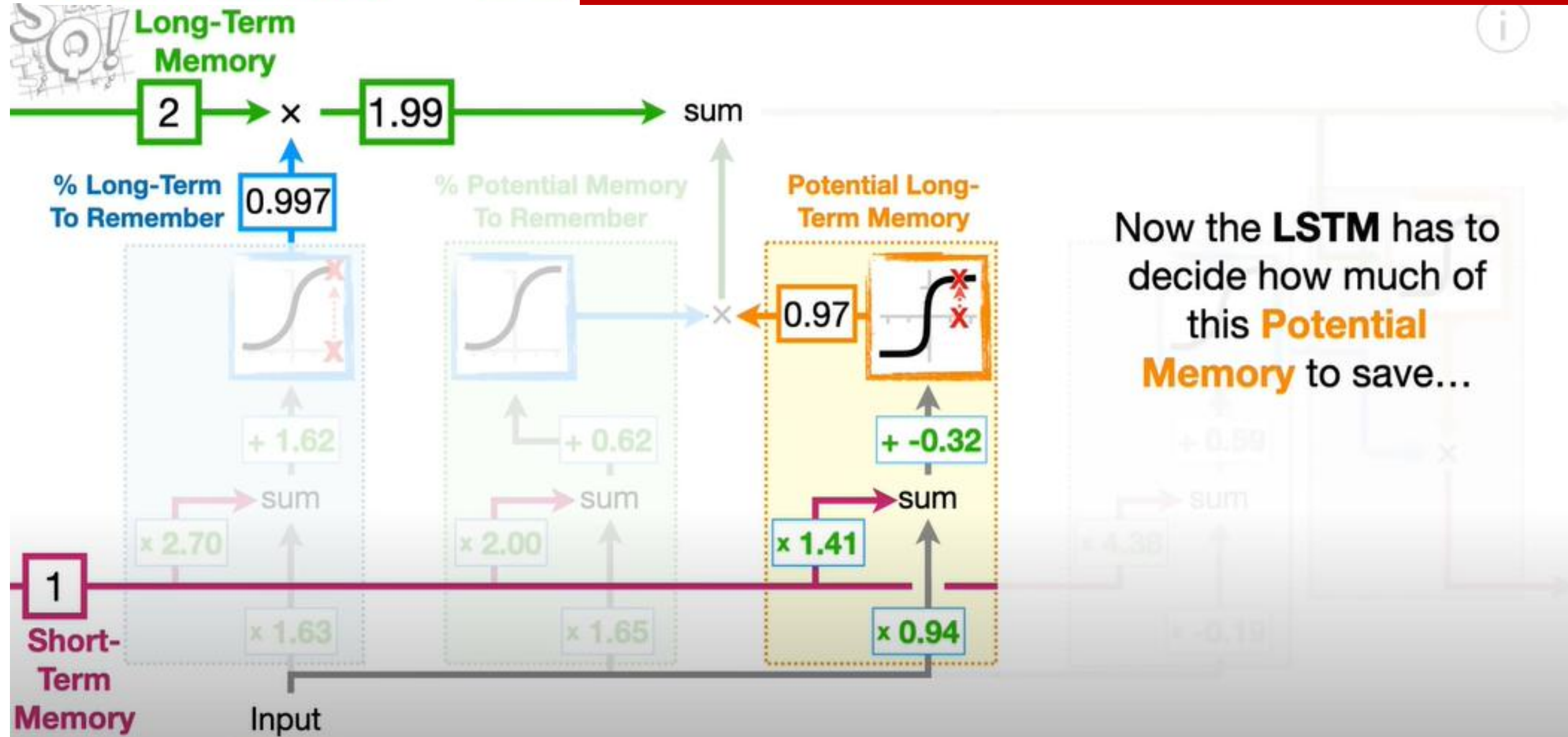


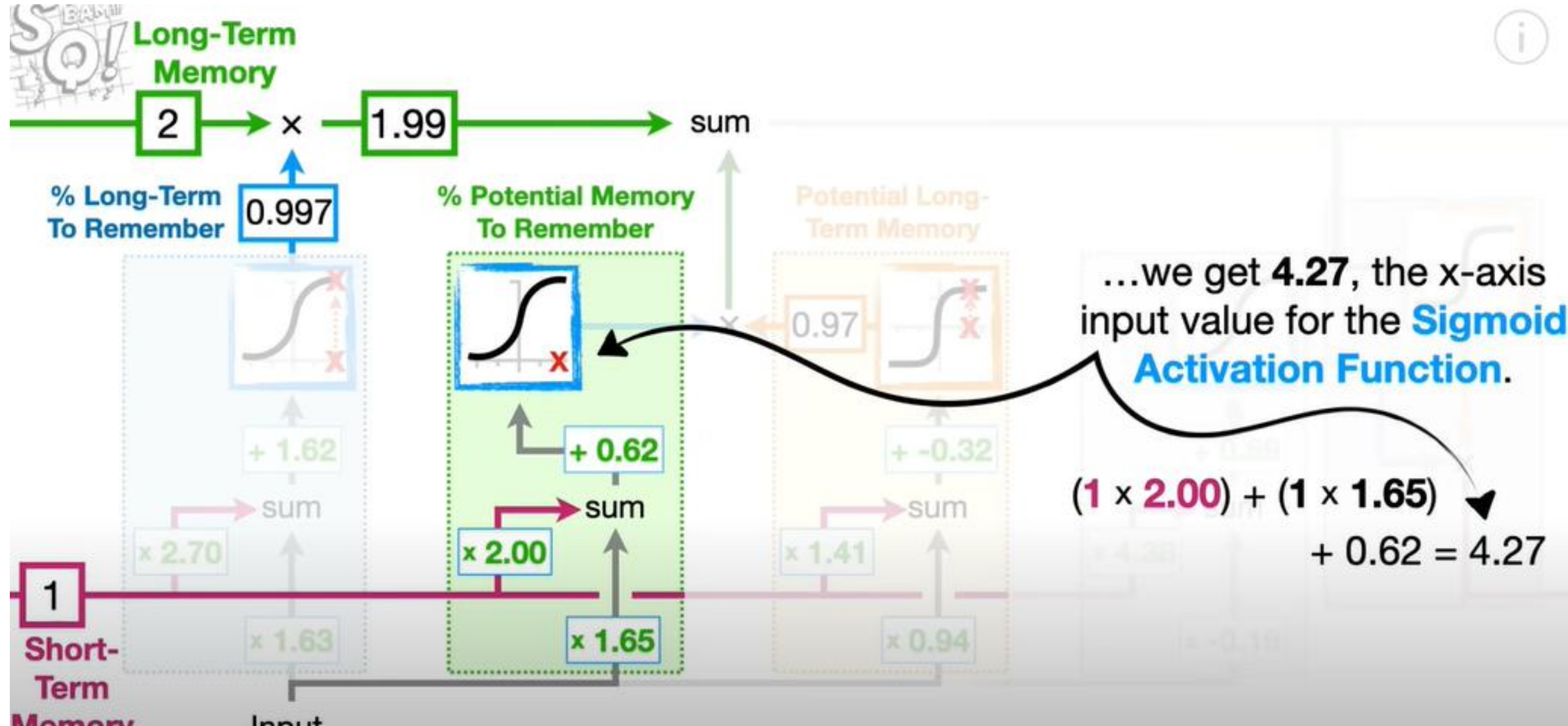


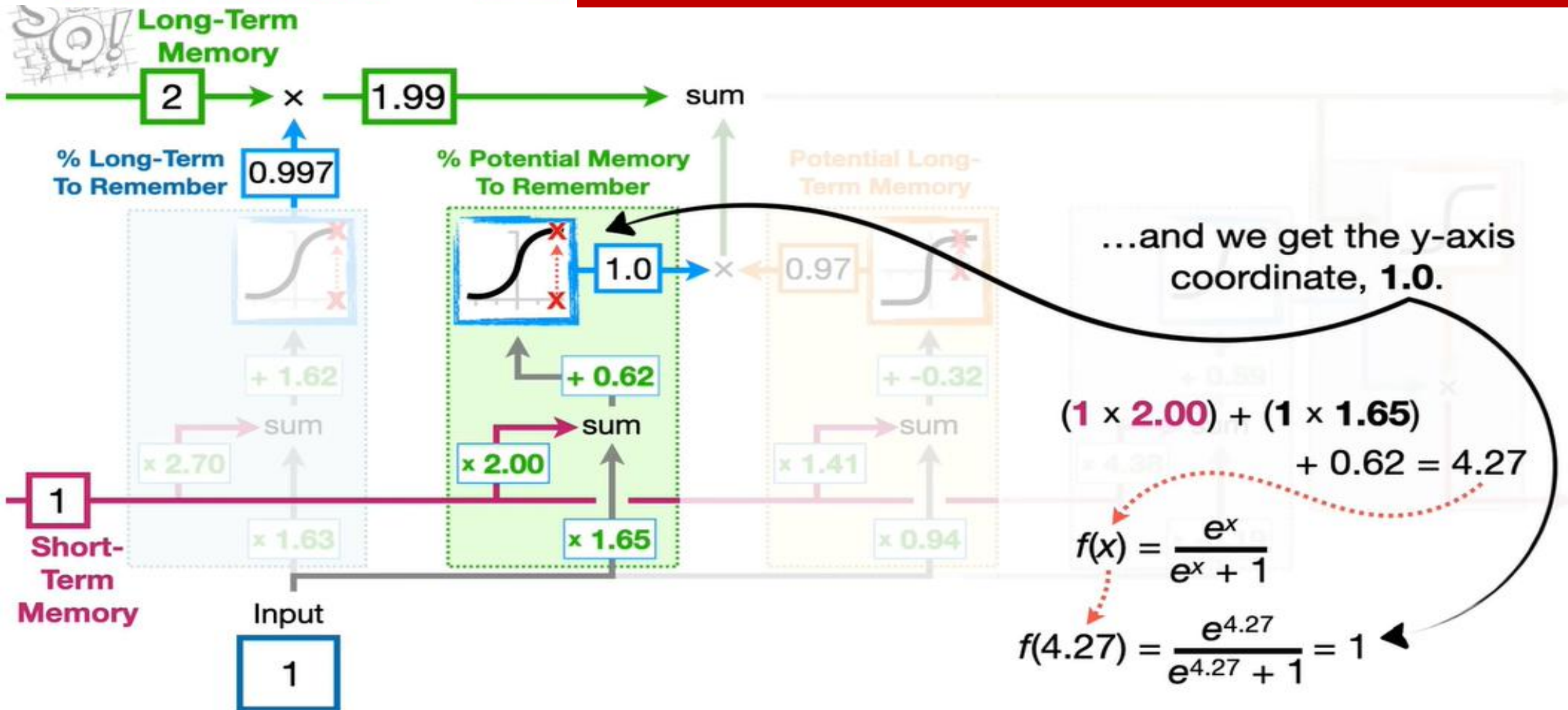


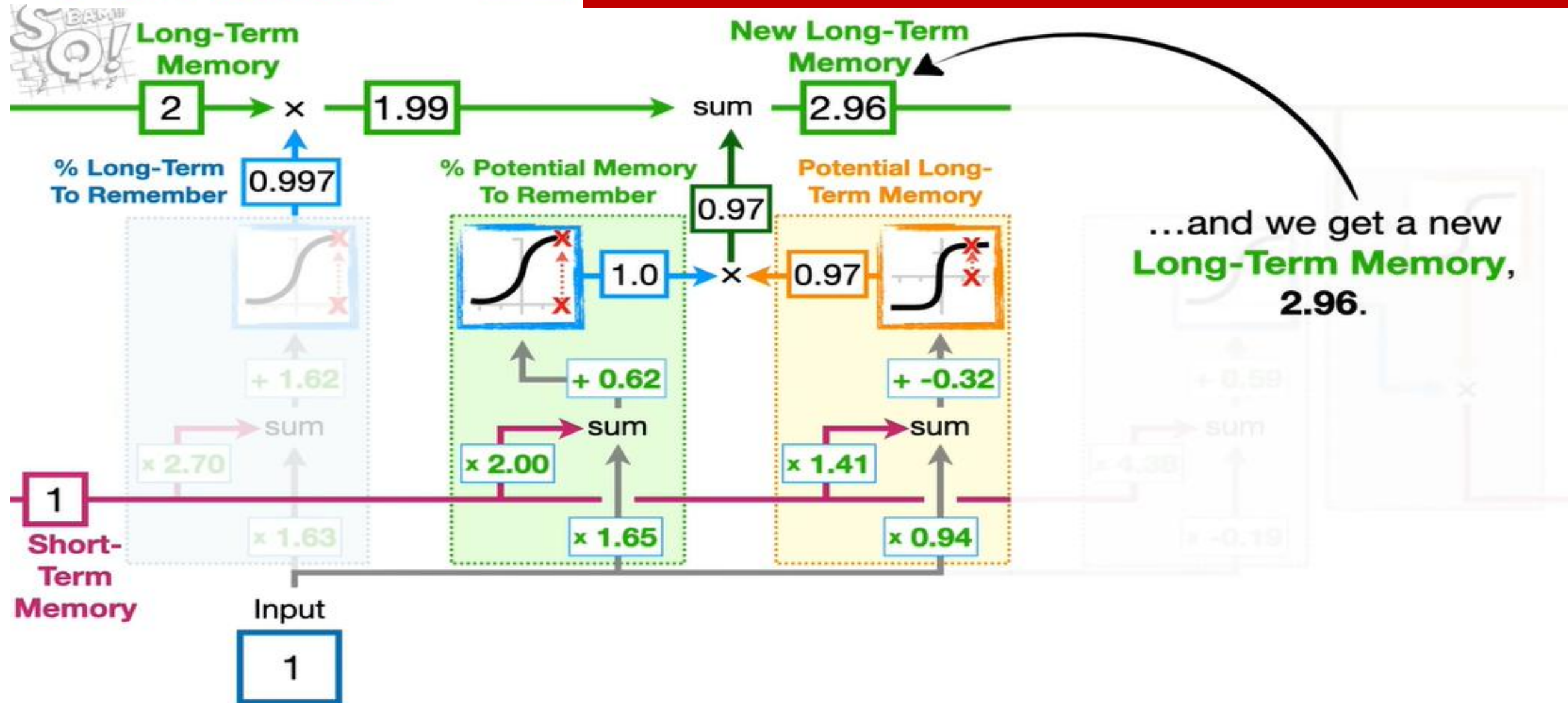
In contrast, if the **Input** to the **LSTM** was **-10...**

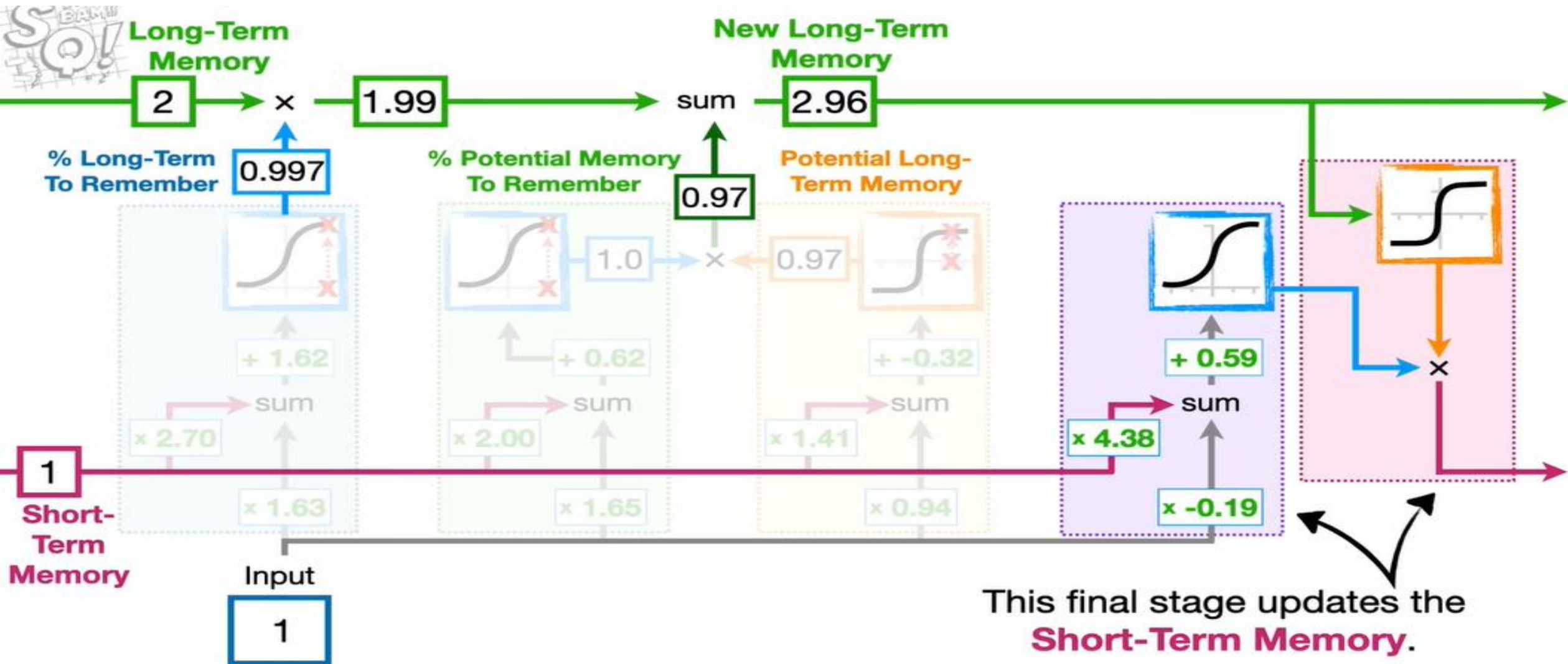
...then, after calculating the x-axis coordinate, the output from the **Tanh Activation Function** function would be **-1**.

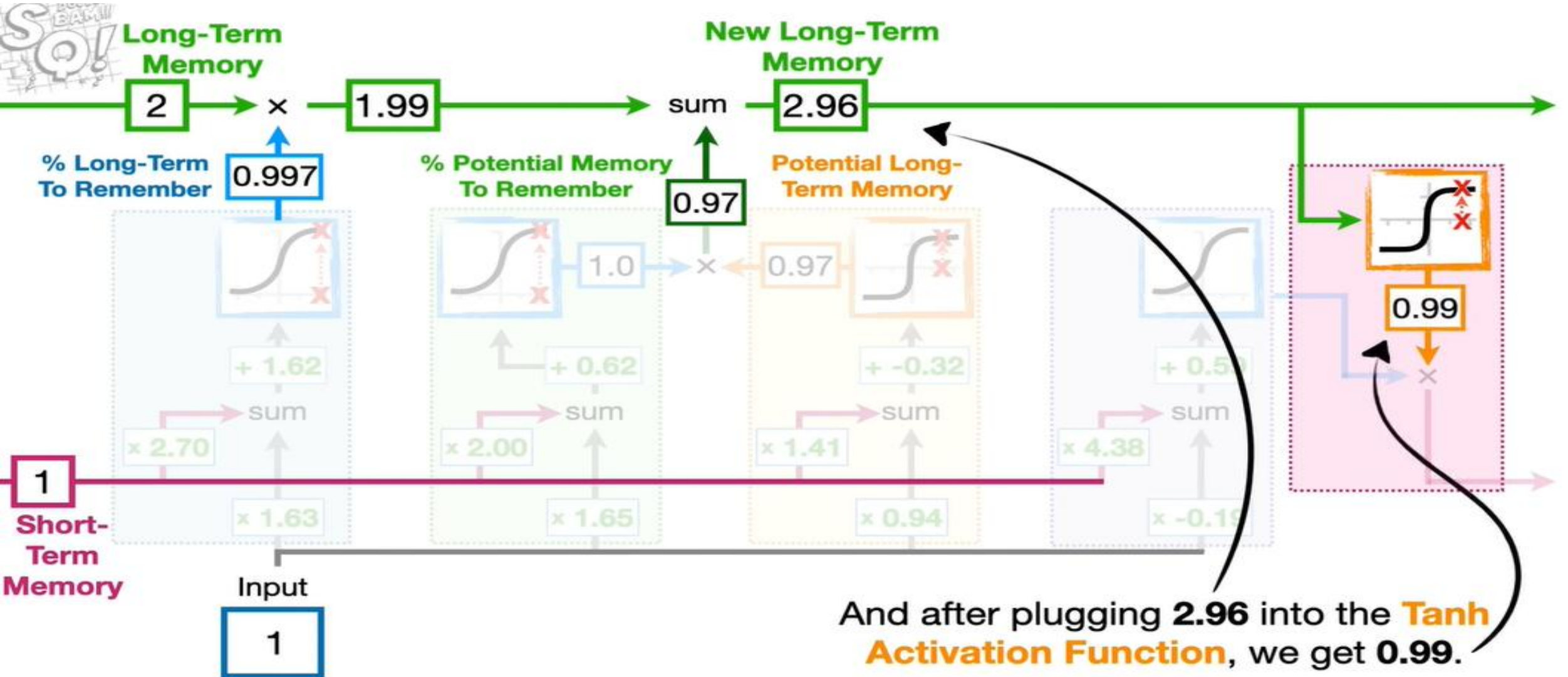


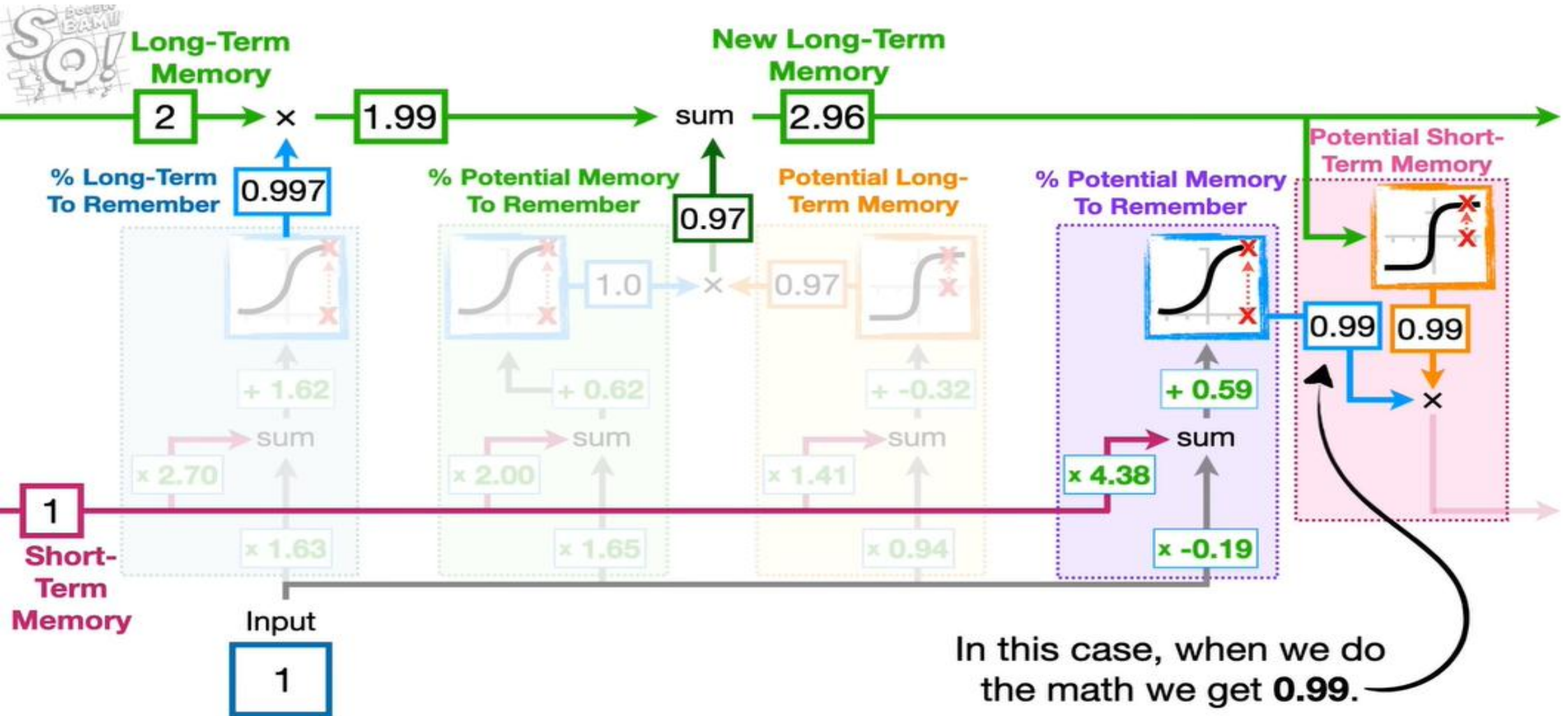


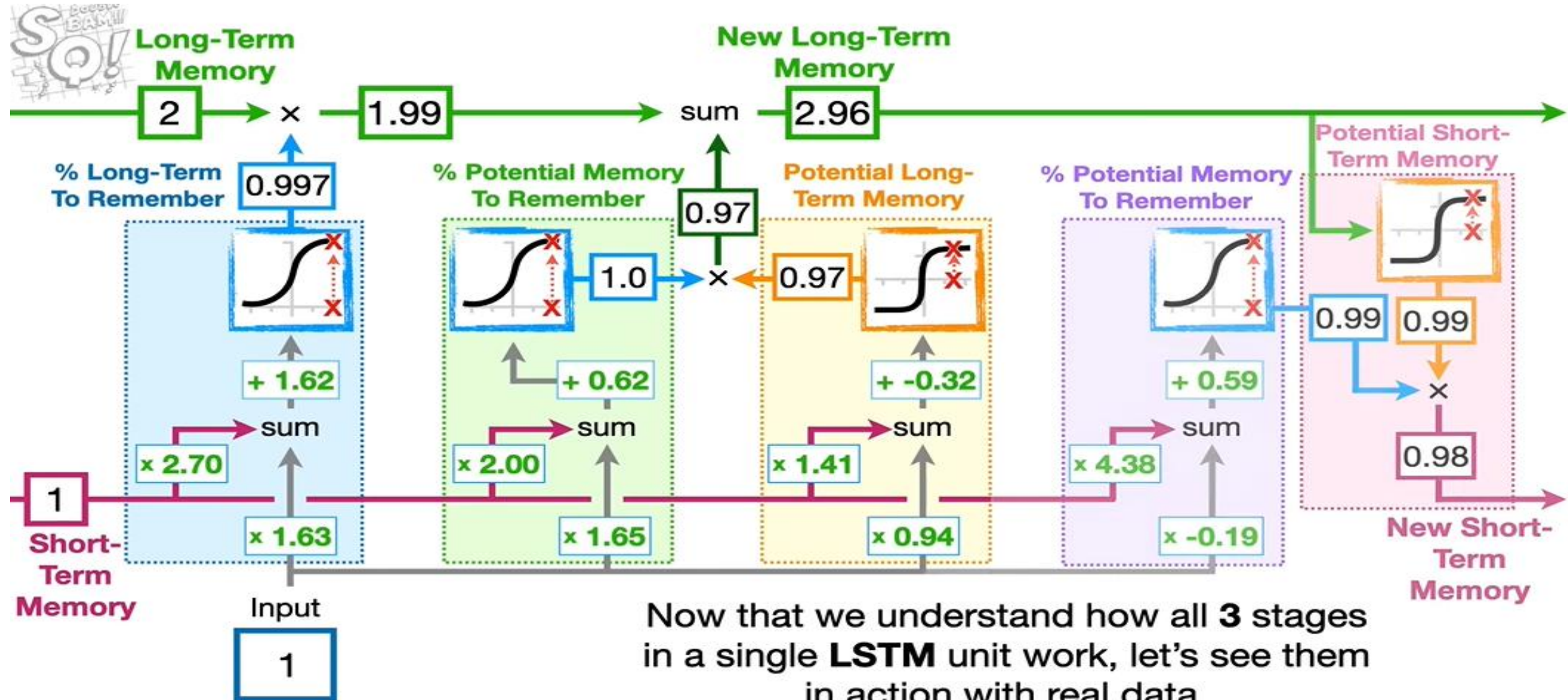












Thank You