

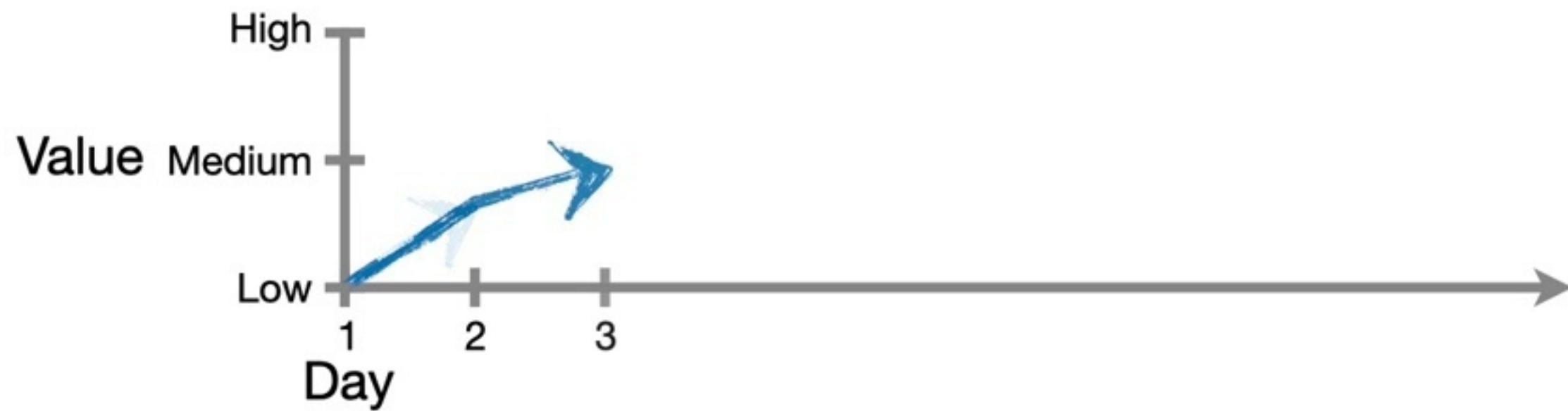


When we look at stock
prices, they tend to
change over time.



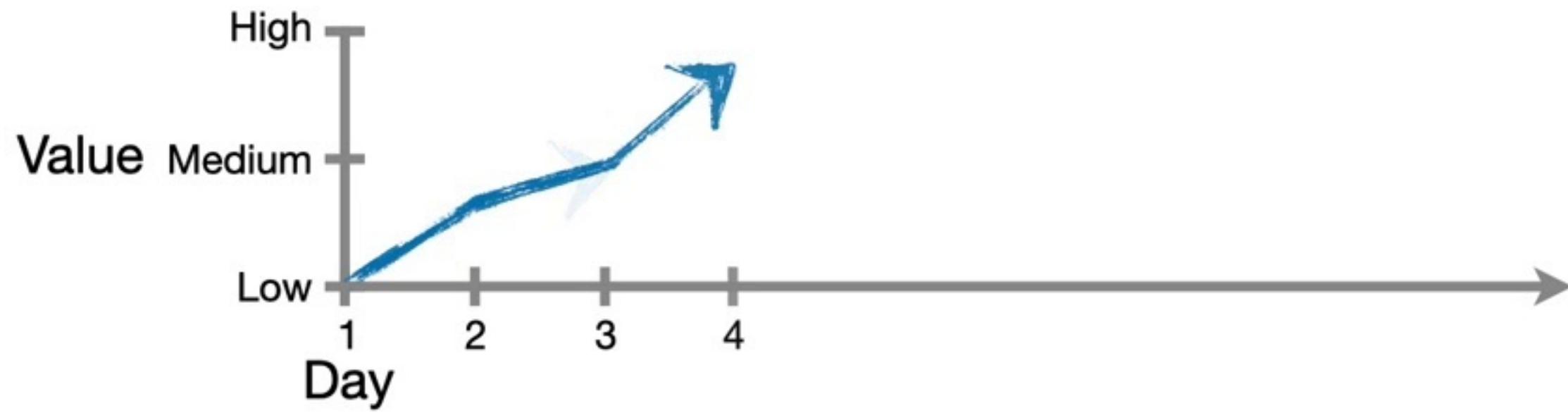


When we look at stock prices, they tend to change over time.



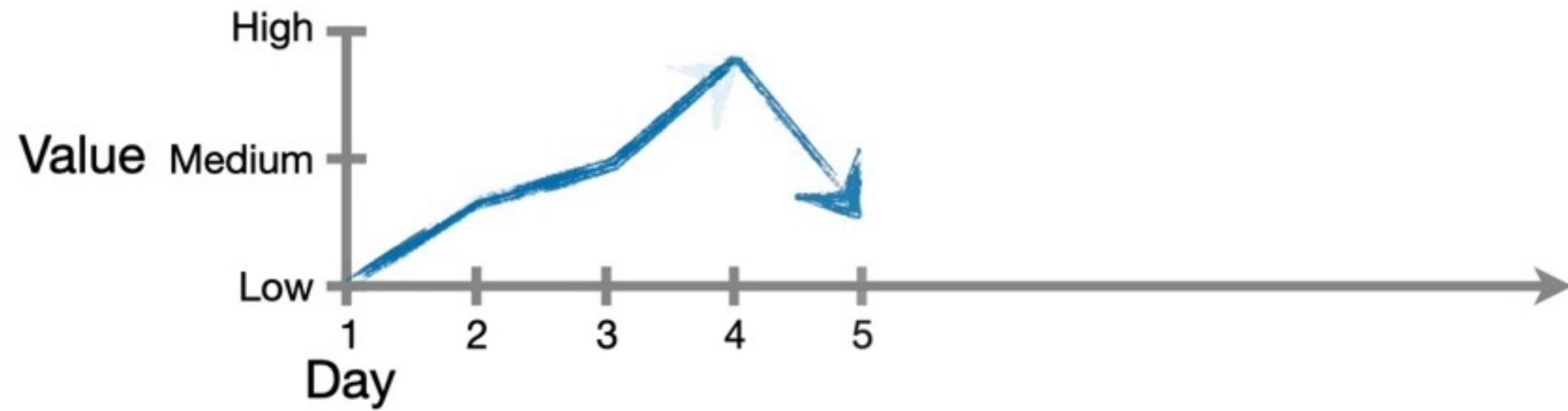


When we look at stock
prices, they tend to
change over time.



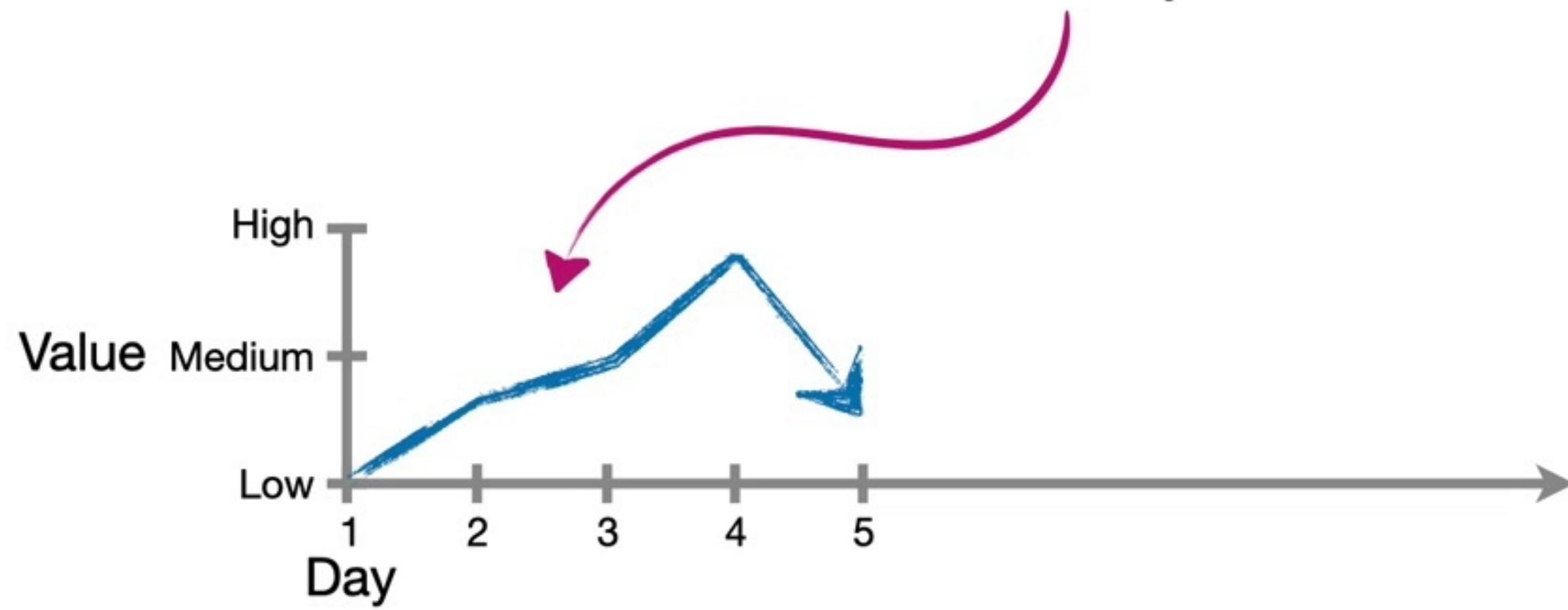


When we look at stock
prices, they tend to
change over time.



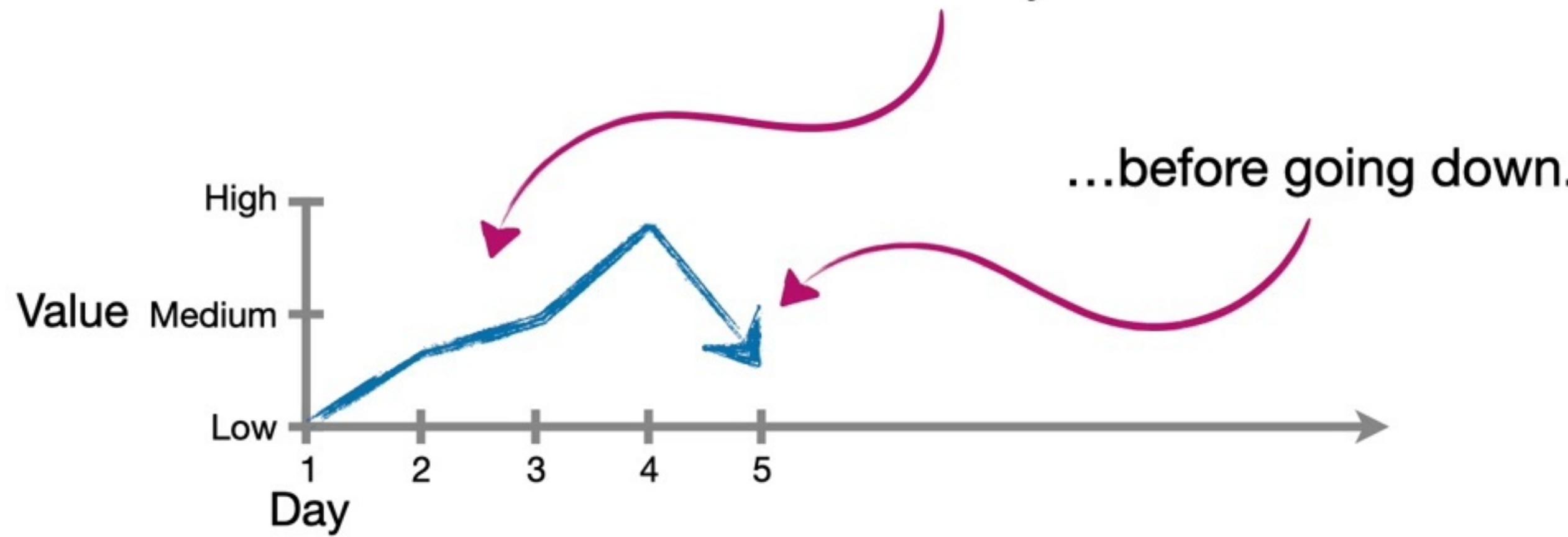


For example, the price
of this stock went up Hooray!
for 4 days...



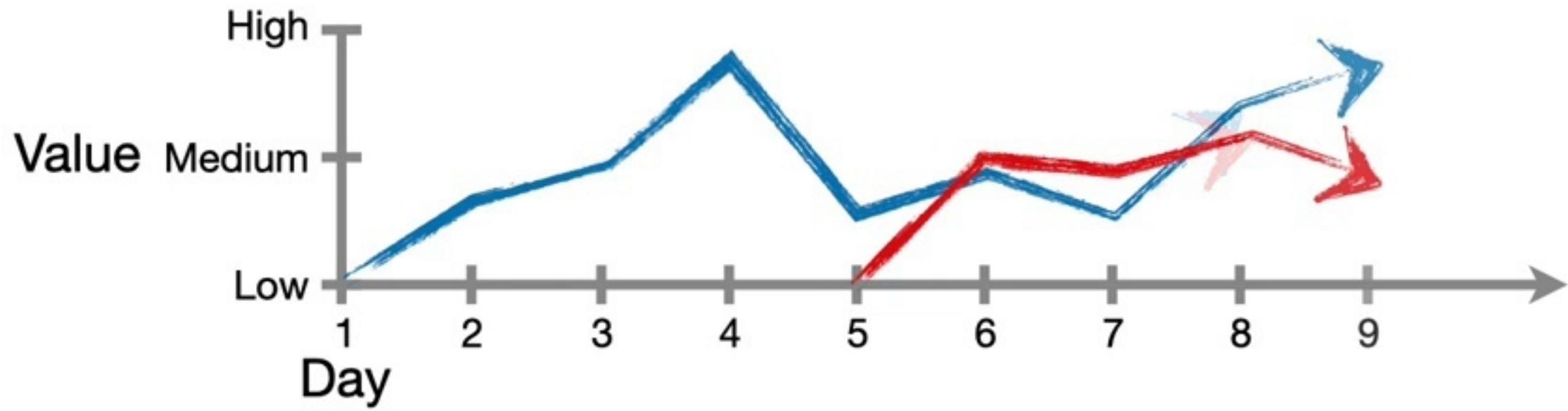


For example, the price
of this stock went up
for 4 days...



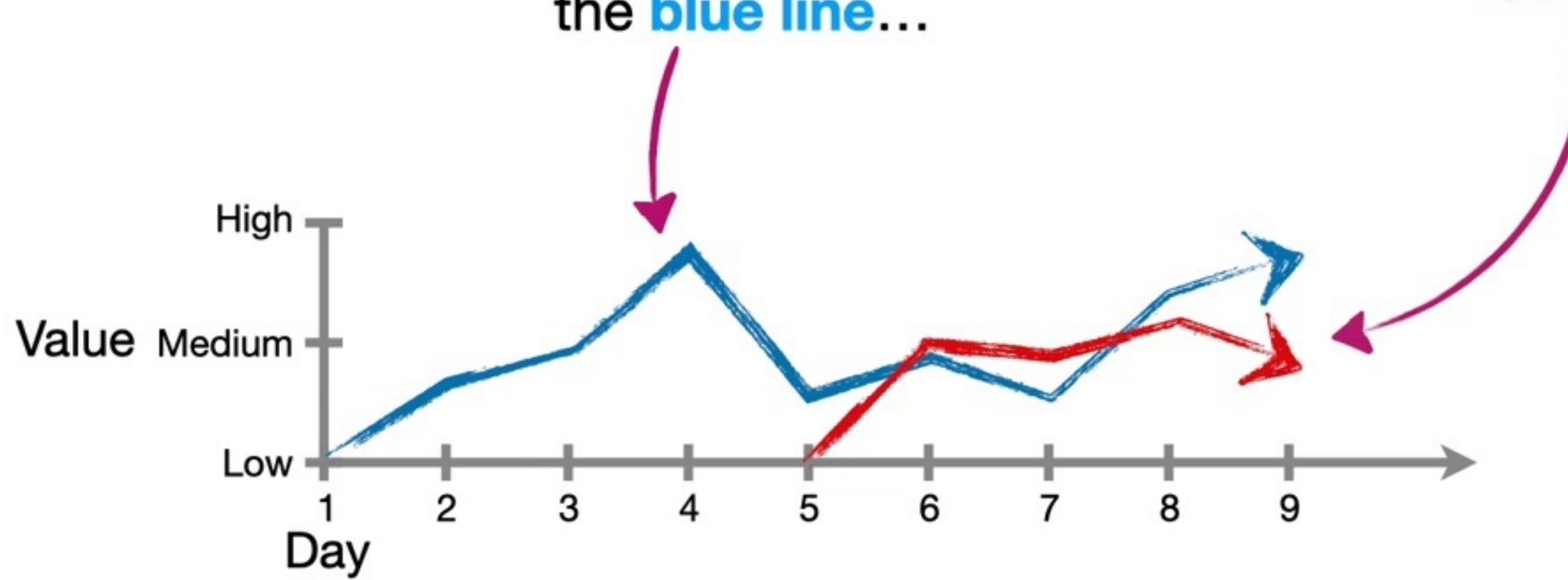


Also, the longer a company has been traded on the stock market, the more data we'll have for it.





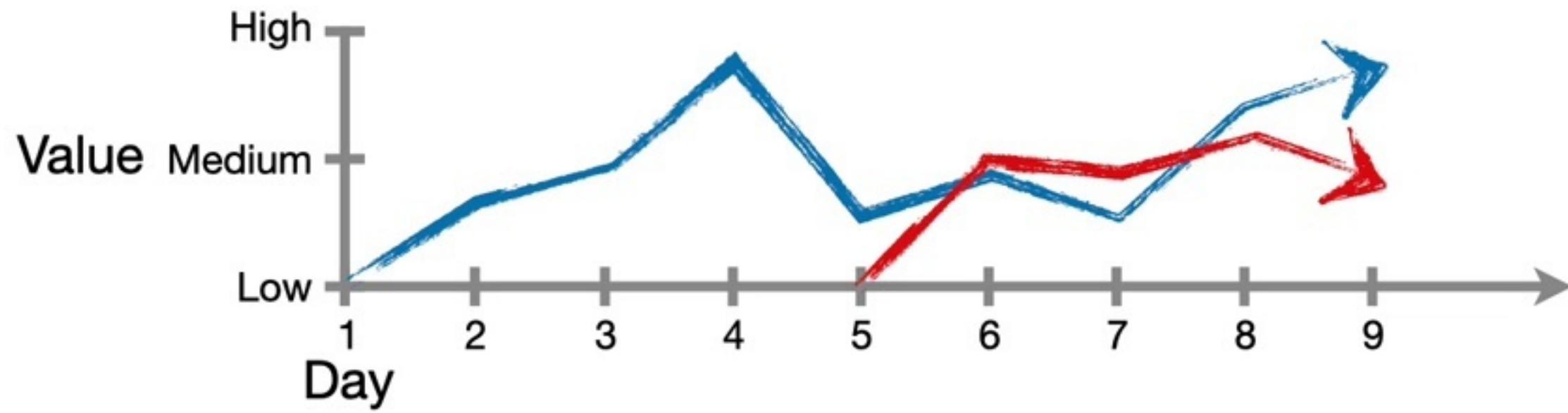
For example, we have more time points for the company represented by the **blue line**...



...then we have for the company represented by the **red line**.

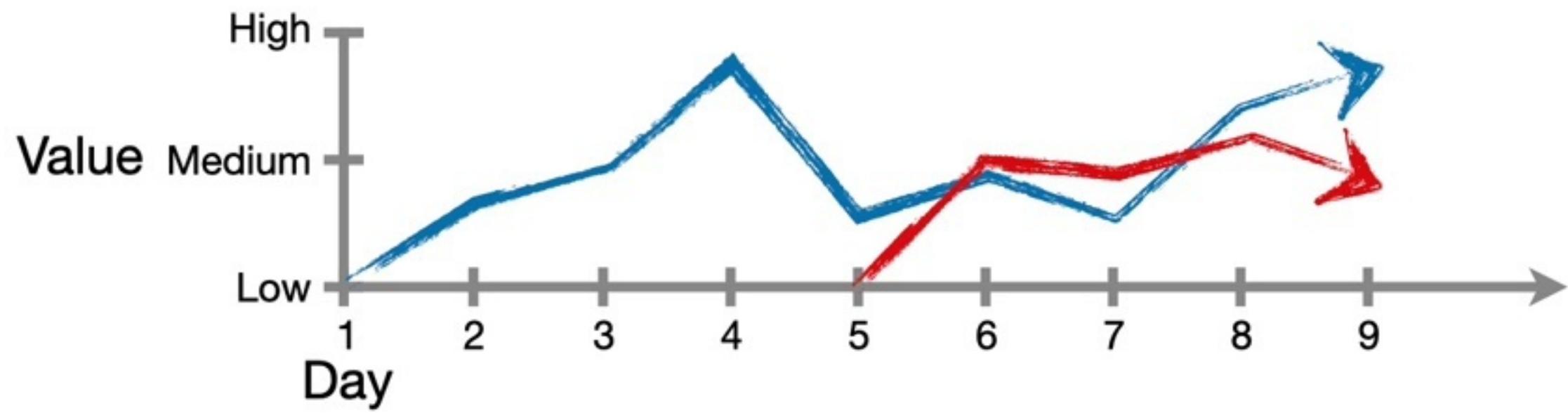


What this means is, if we want to use a neural network to predict stock prices...



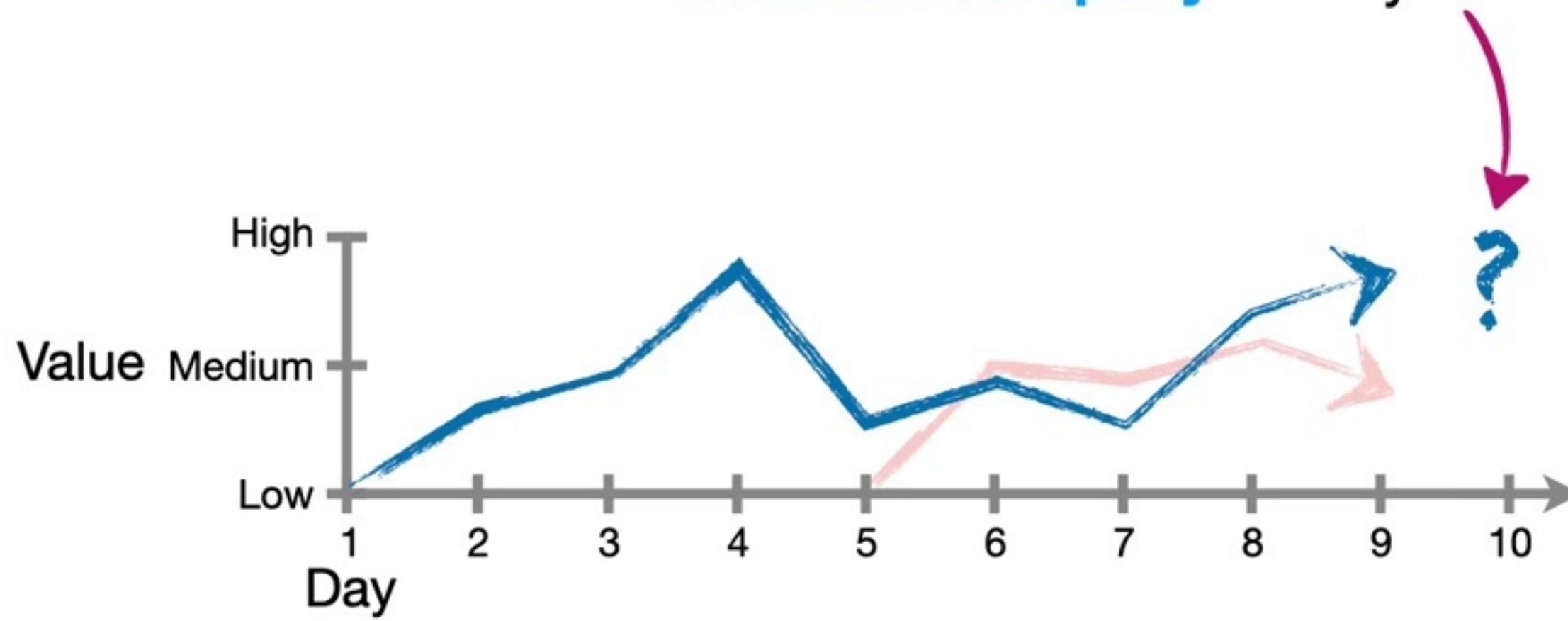


...then we need a neural network that works with different amounts of sequential data.



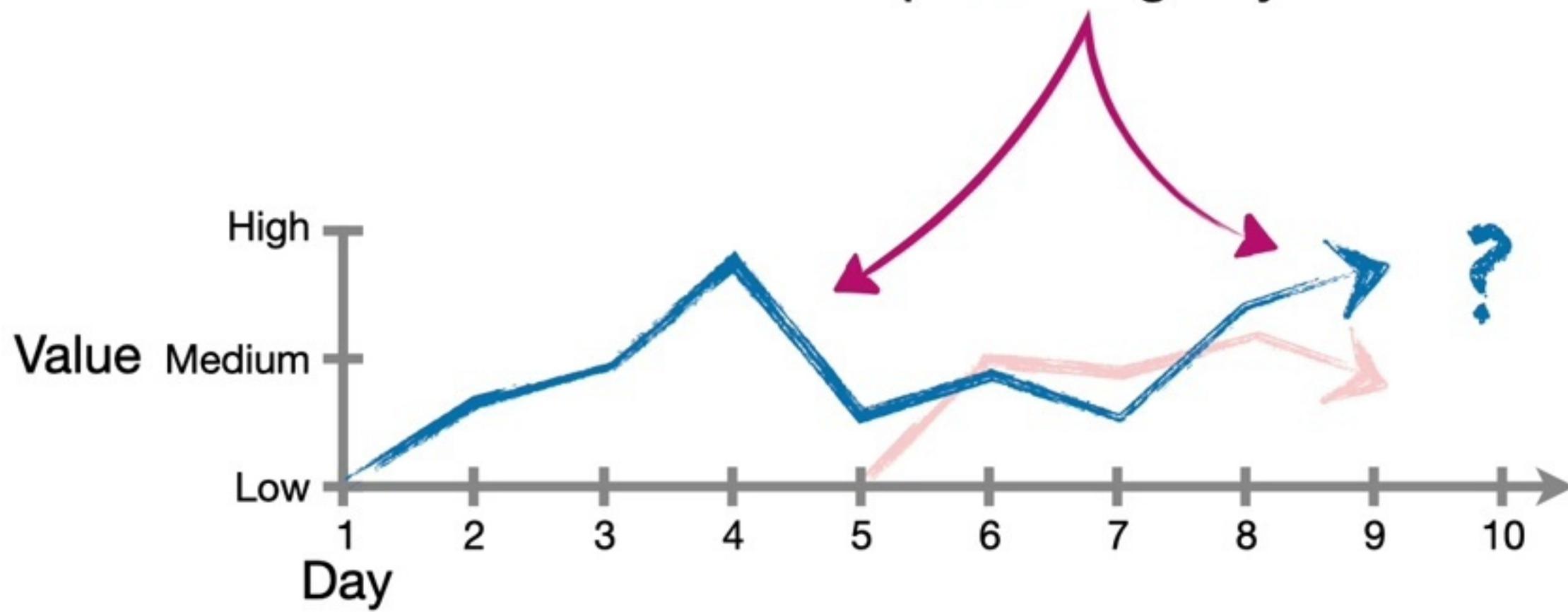


In other words, if we want to predict the stock price for the
blue line company on day 10...



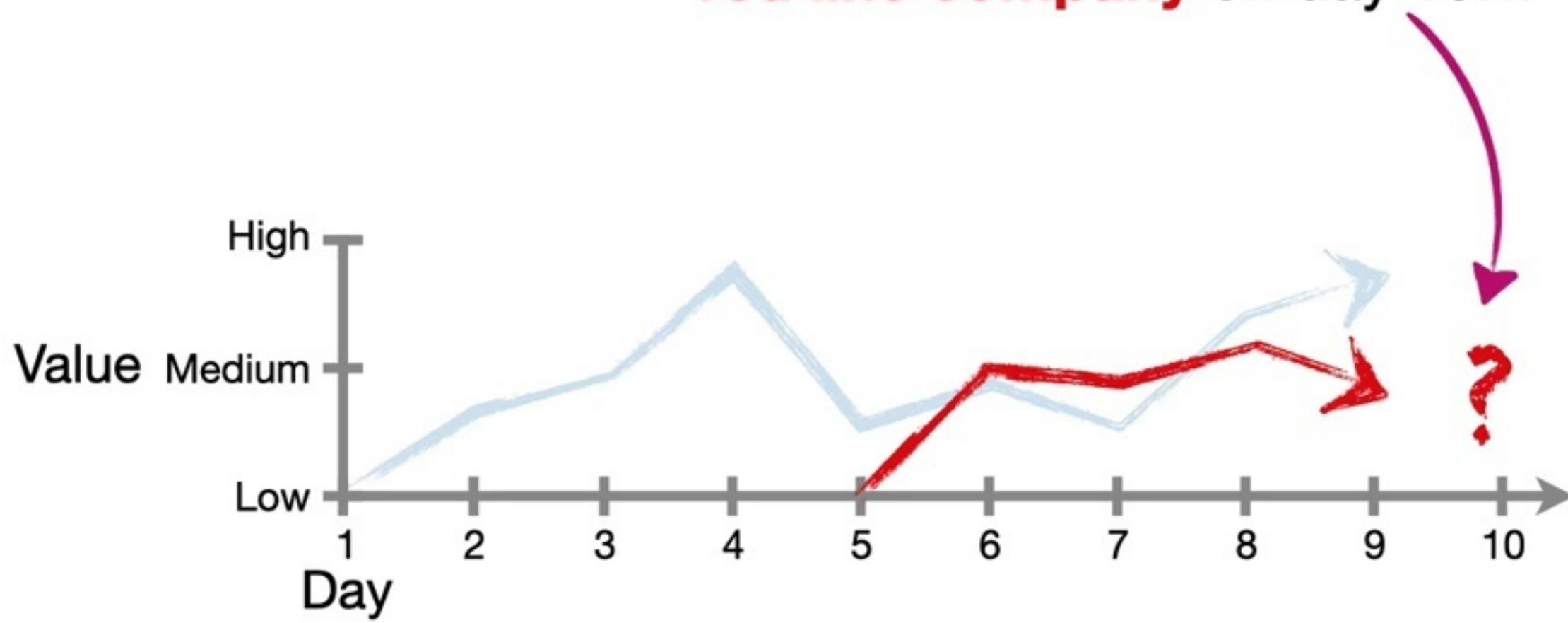


...then we might want to use
the data from all **9** of the
preceding days.



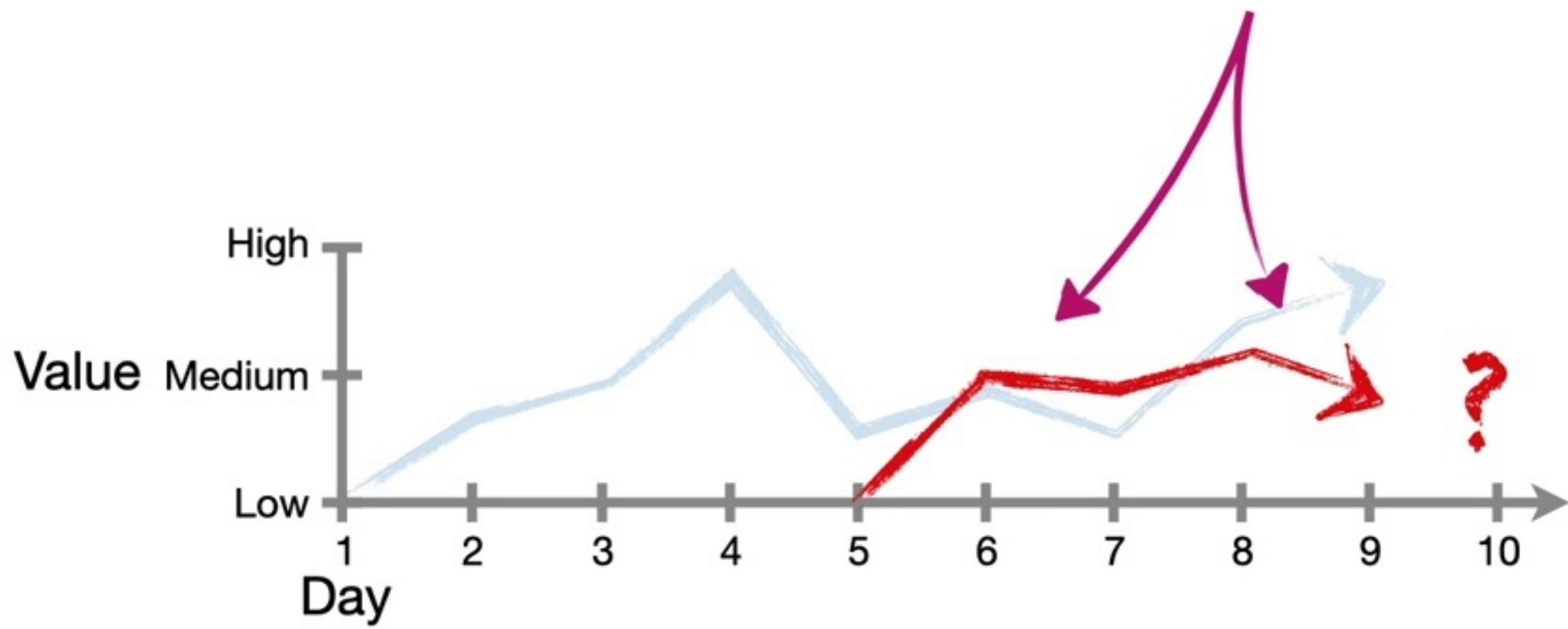


In contrast, if we wanted to predict the stock price for the
red line company on day 10...



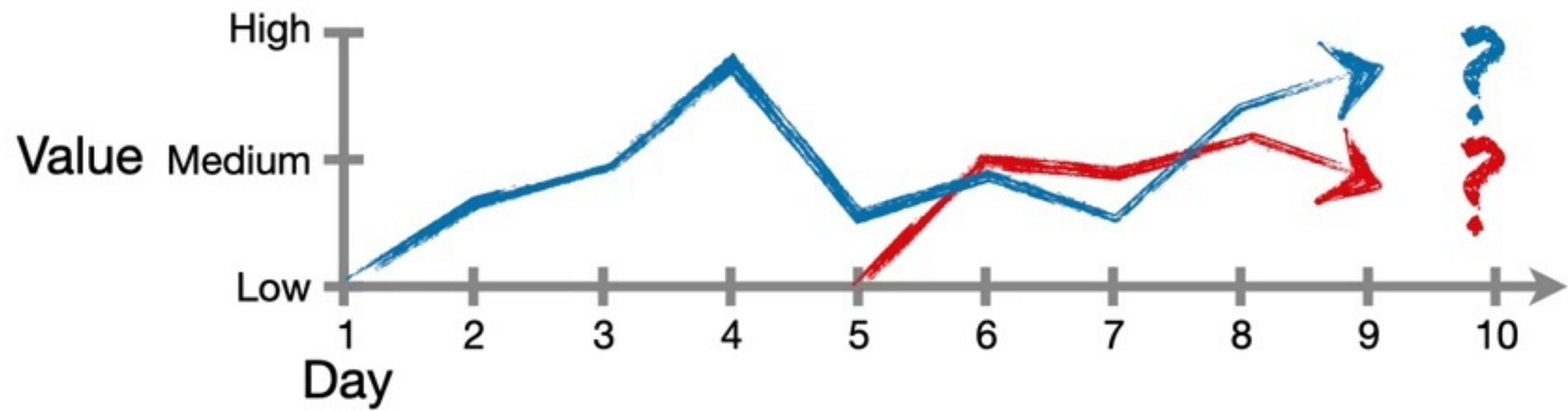


...then we would only have data
for the preceding **5** days.



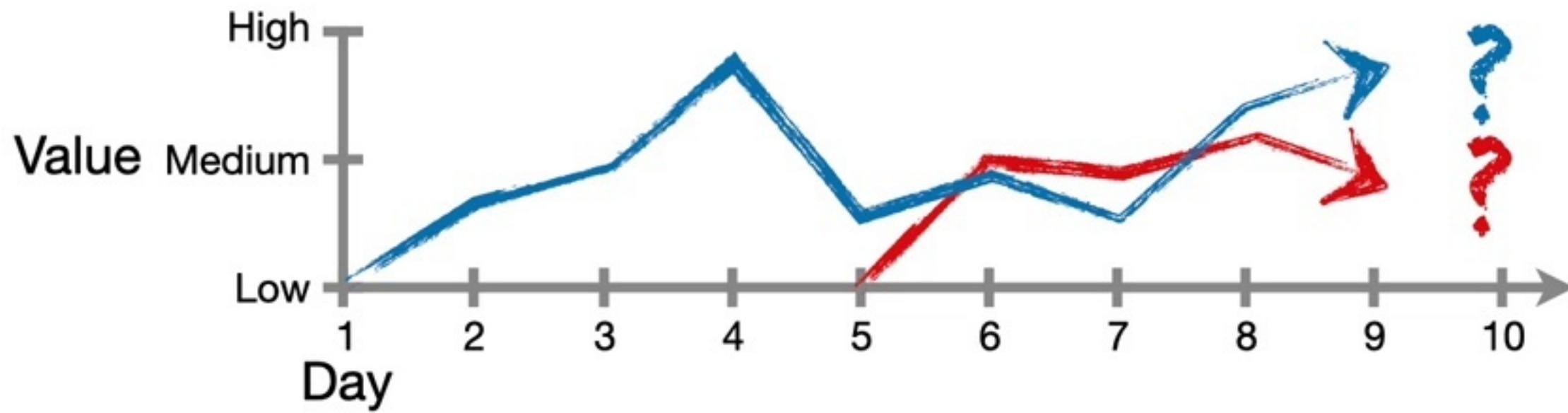


So we need the neural network
to be flexible in terms of how
much sequential data we use
to make a prediction.



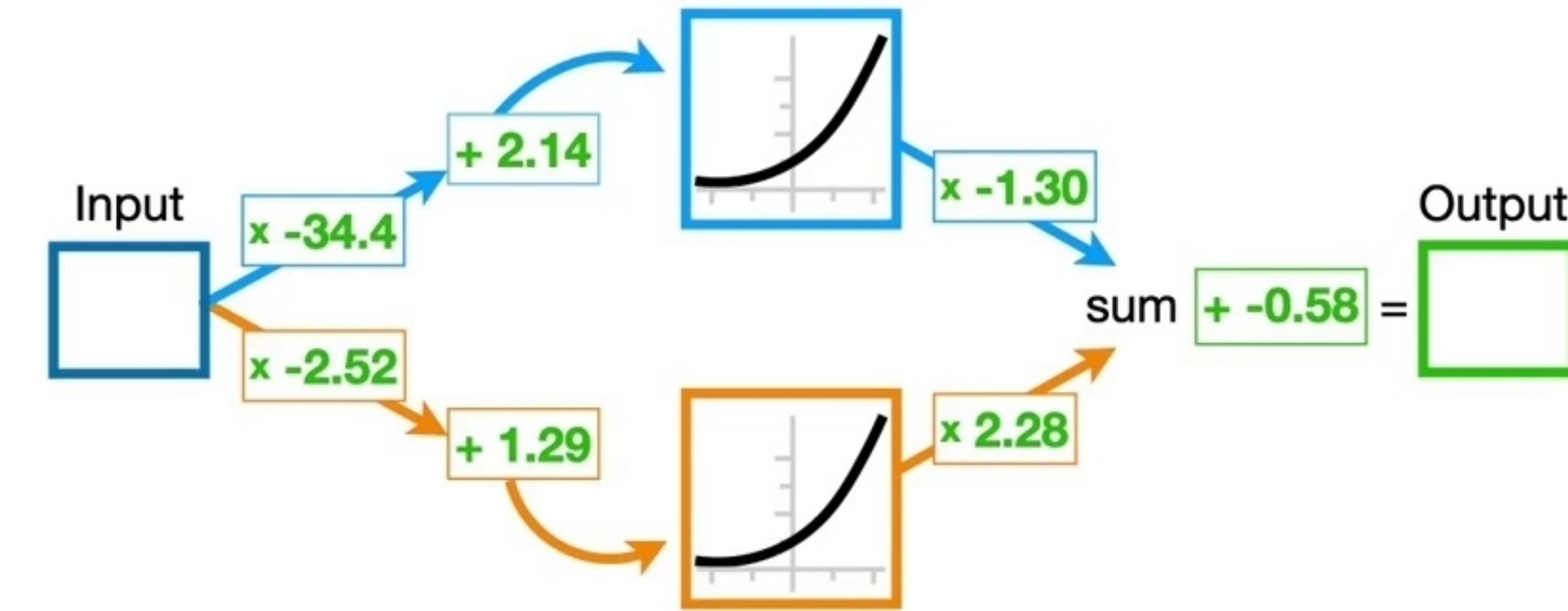


This is a big difference compared
to the other neural networks
we've looked at in this series.





For example, in **Neural Networks Clearly Explained**, we examined a neural network that made predictions using 1 input value, no more and no less.



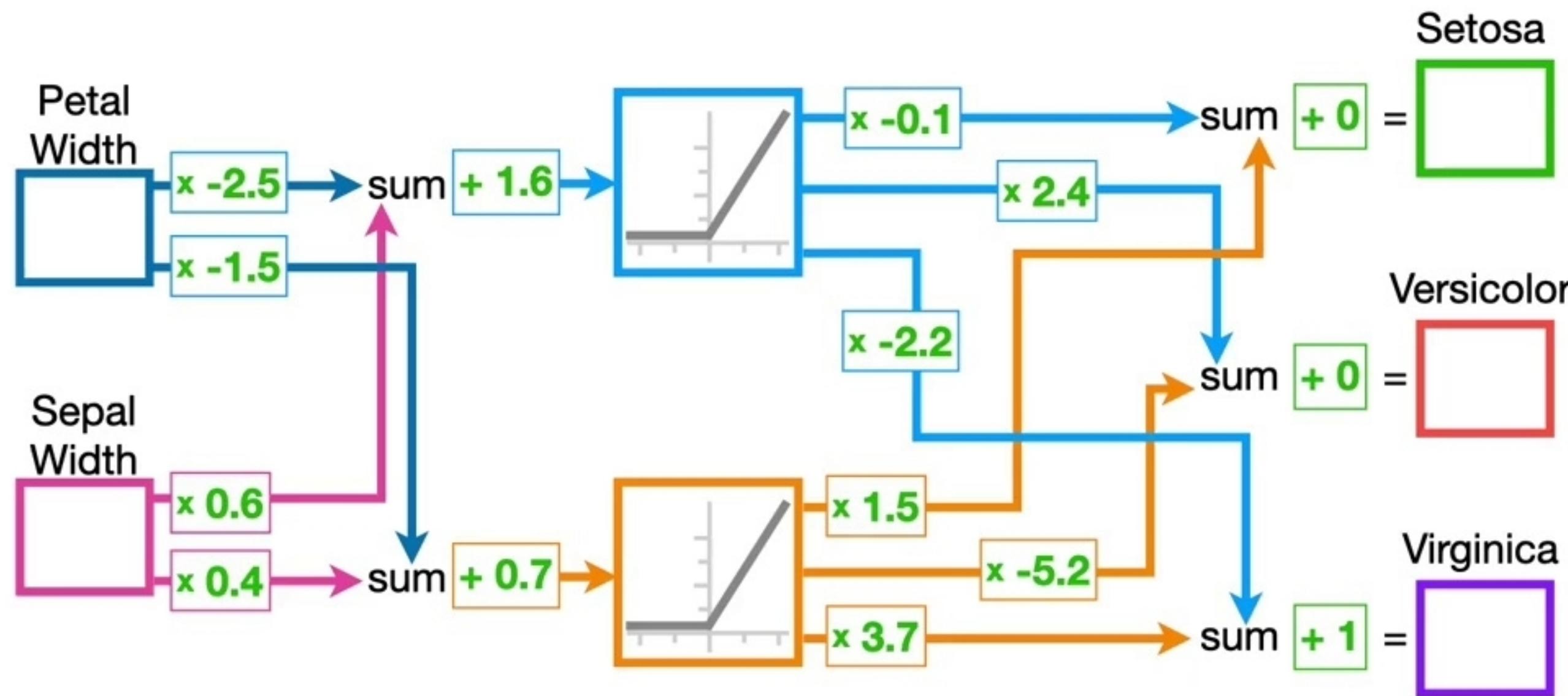


For example, in **Neural Networks Clearly Explained**, we examined a neural network that made predictions using **1** input value, no more and no less.



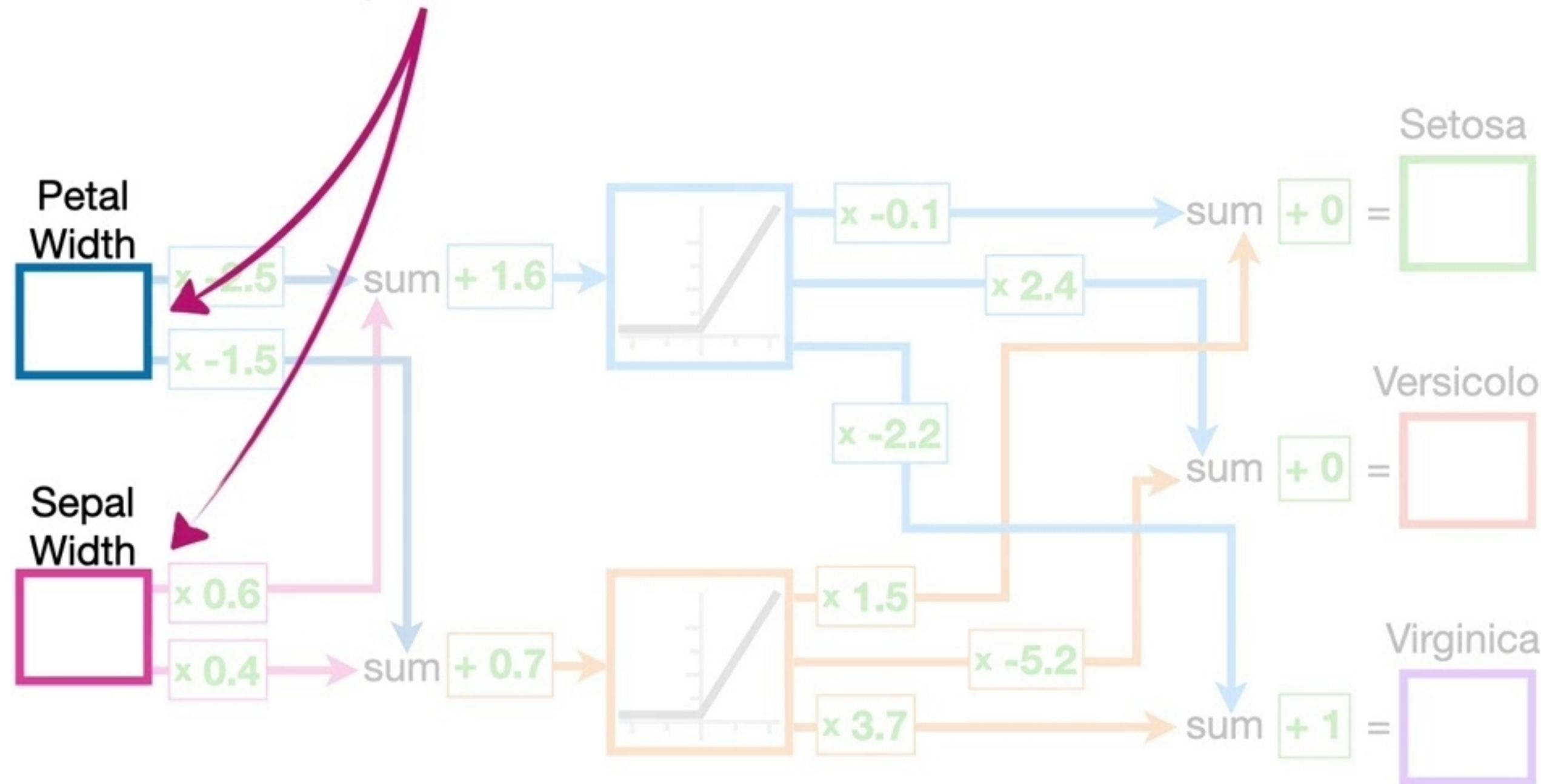


And if you saw the **StatQuest** on
Neural Networks with Multiple Inputs
and Outputs, you saw this neural
network that made predictions using **2**
input values, no more and no less.





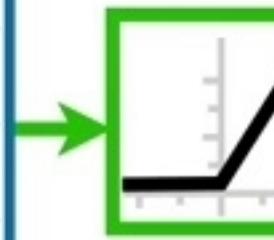
And if you saw the **StatQuest** on
Neural Networks with Multiple Inputs
and Outputs, you saw this neural
network that made predictions using **2**
input values, no more and no less.





Feature Map

1	-1	-2	-1
-1	-2	-1	-2
-2	-1	-2	-1
-1	-2	-1	1

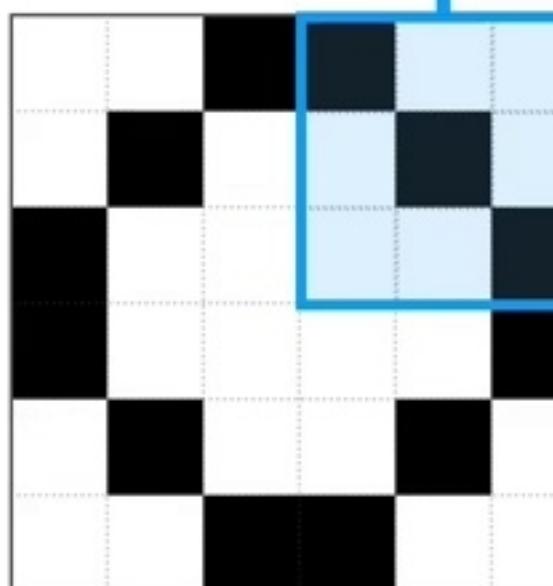


1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	1

$$+ -2$$

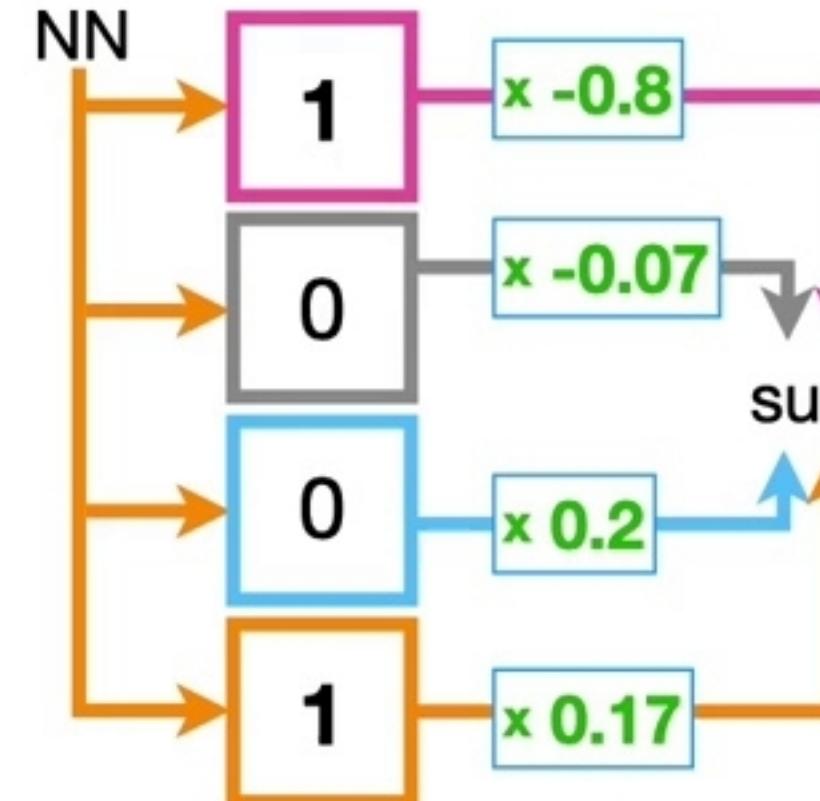


Filter
(Convolution)

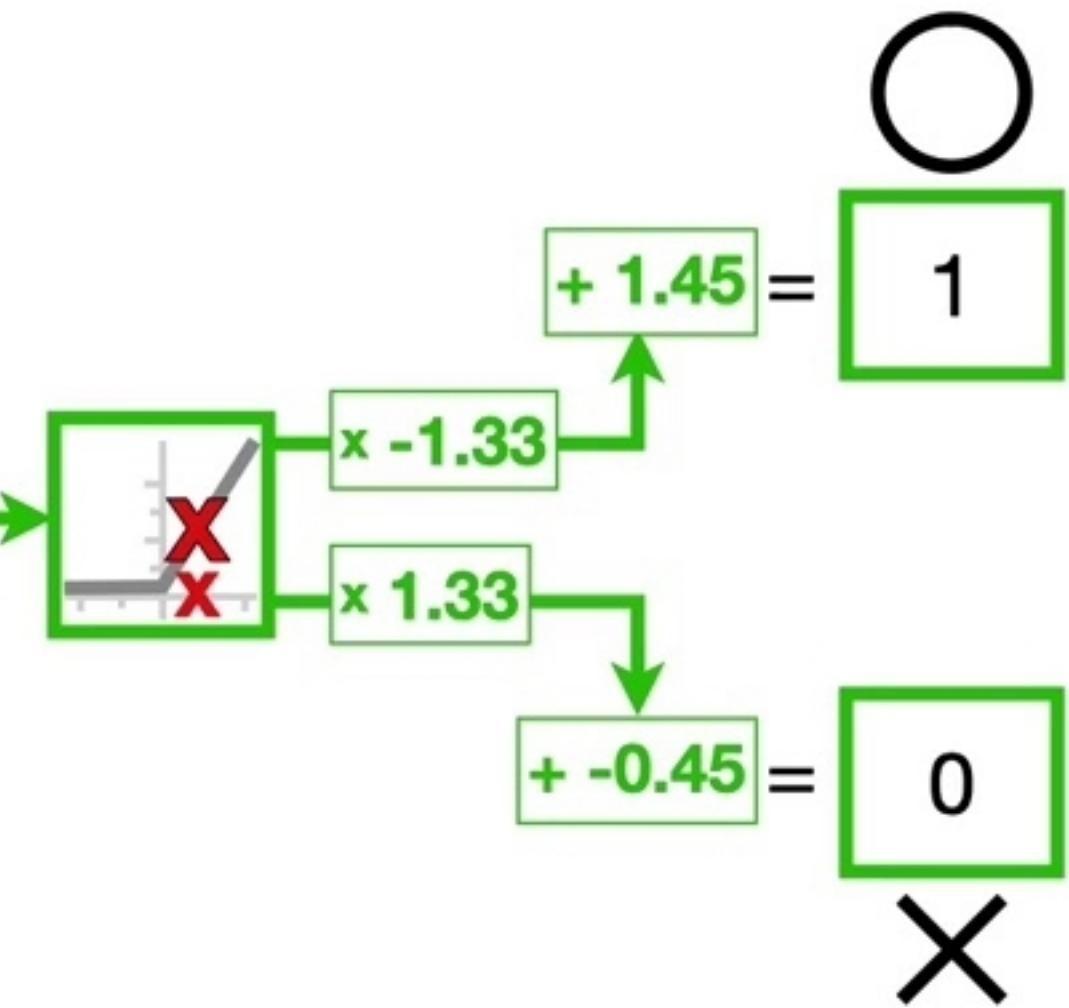


1	0
0	1

Input to NN

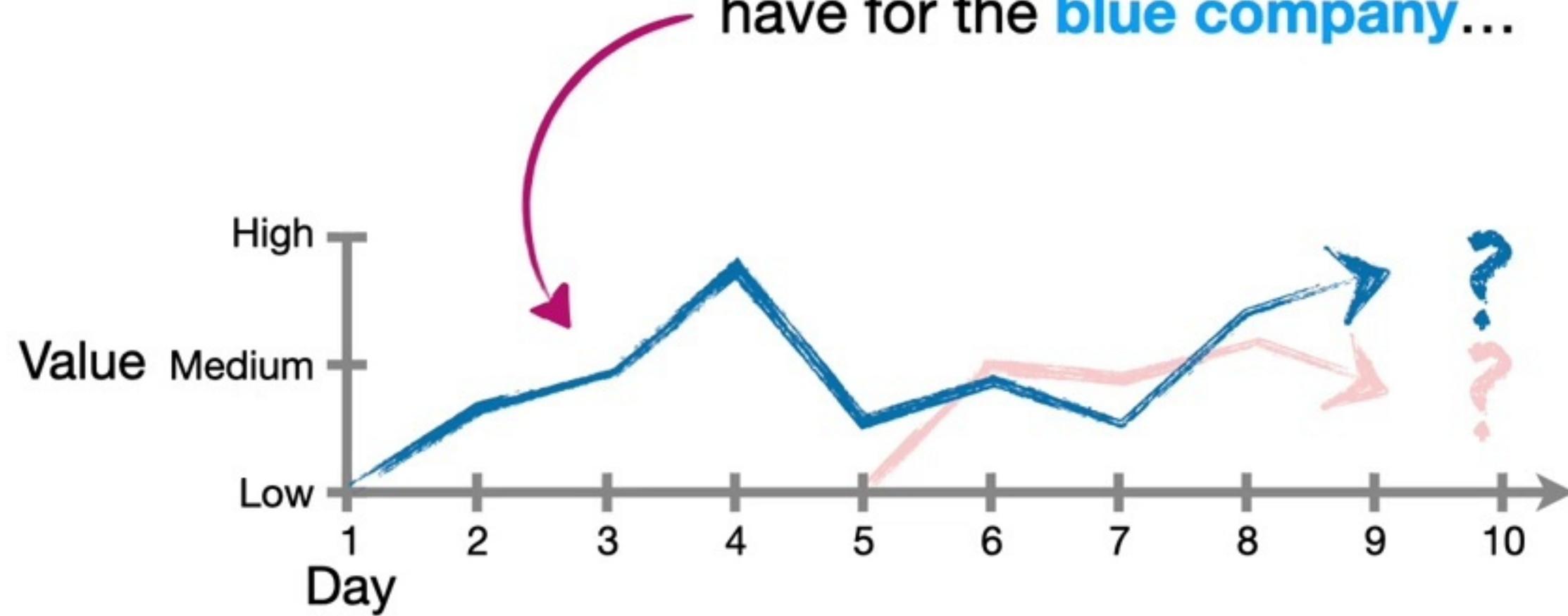


Max Pooling



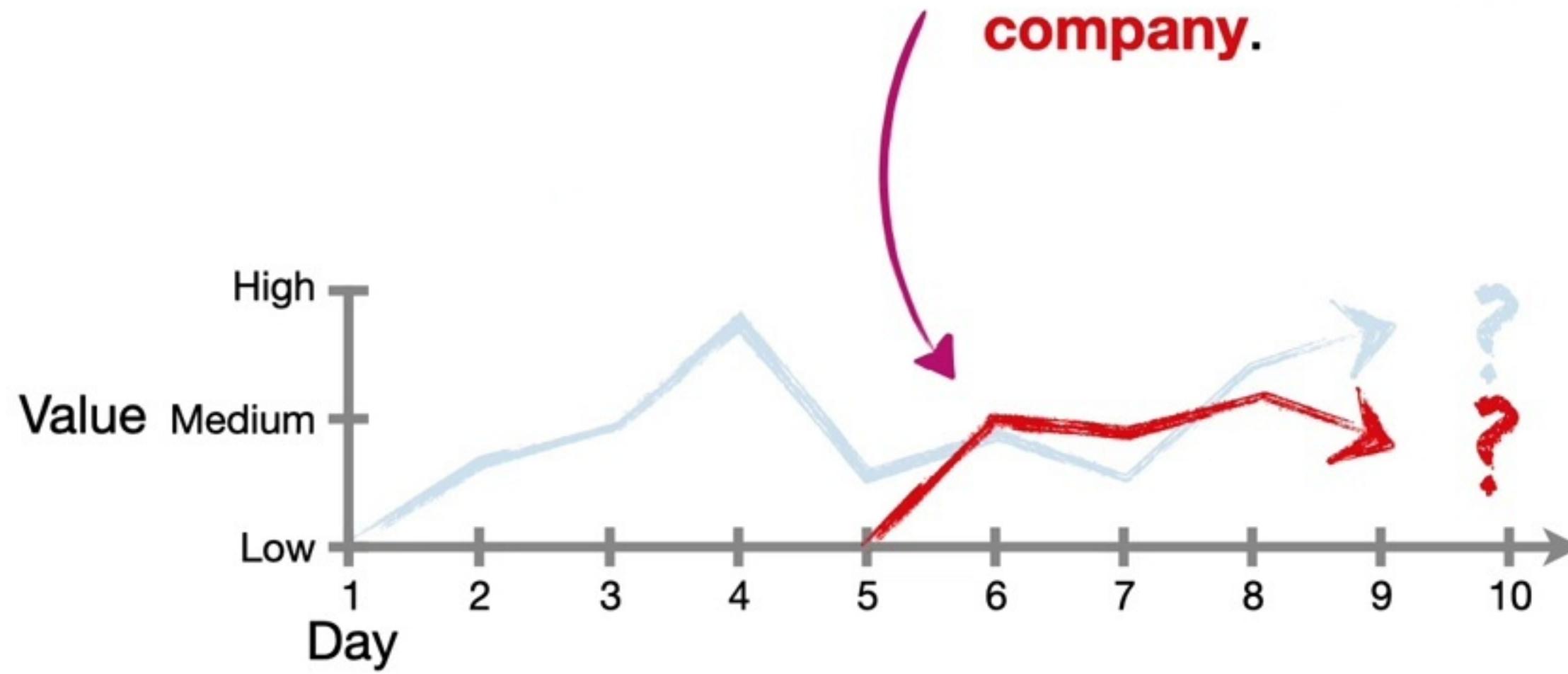


However, now we need a neural network that can make a prediction using the **9** values we have for the **blue company**...



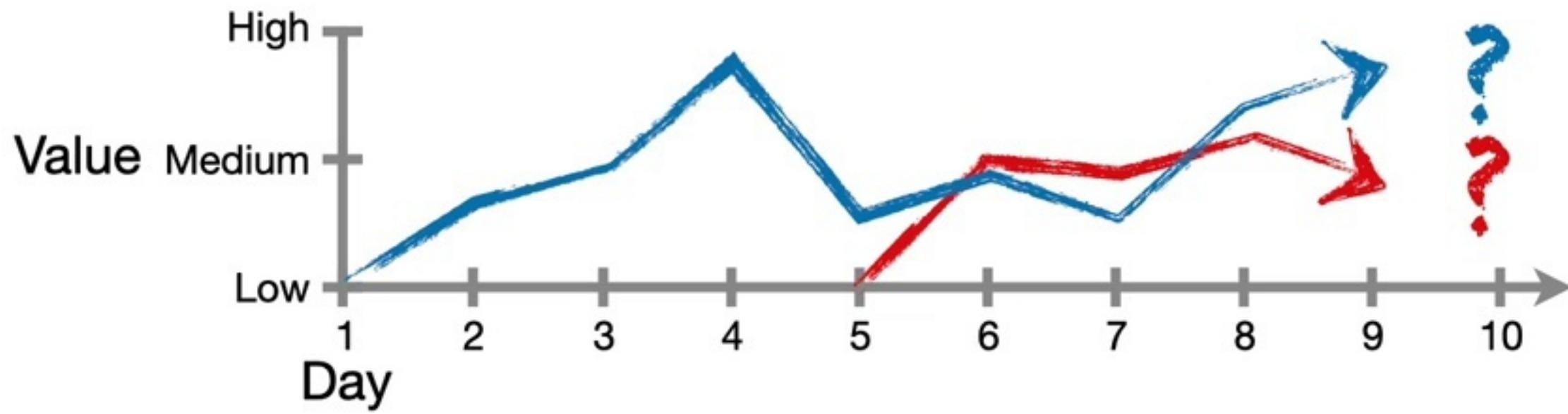


...and make a prediction using
the **5** values we have for the **red**
company.



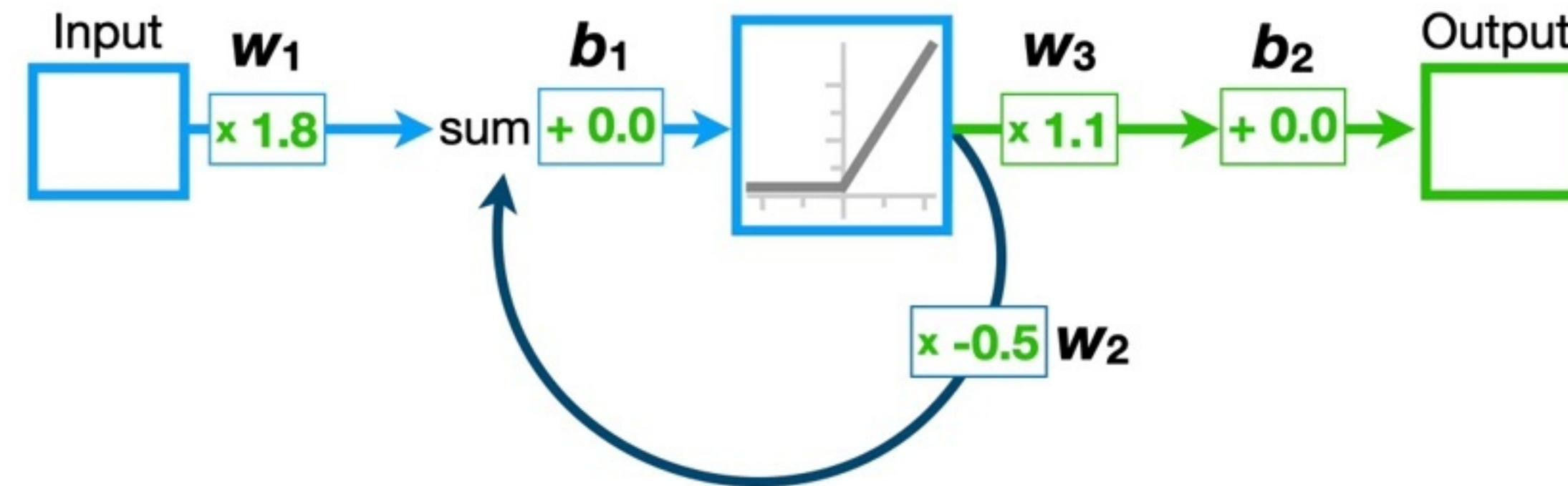


The good news is that one way to deal with the problem of having different amounts of input values is to use a **Recurrent Neural Network (RNN)**.



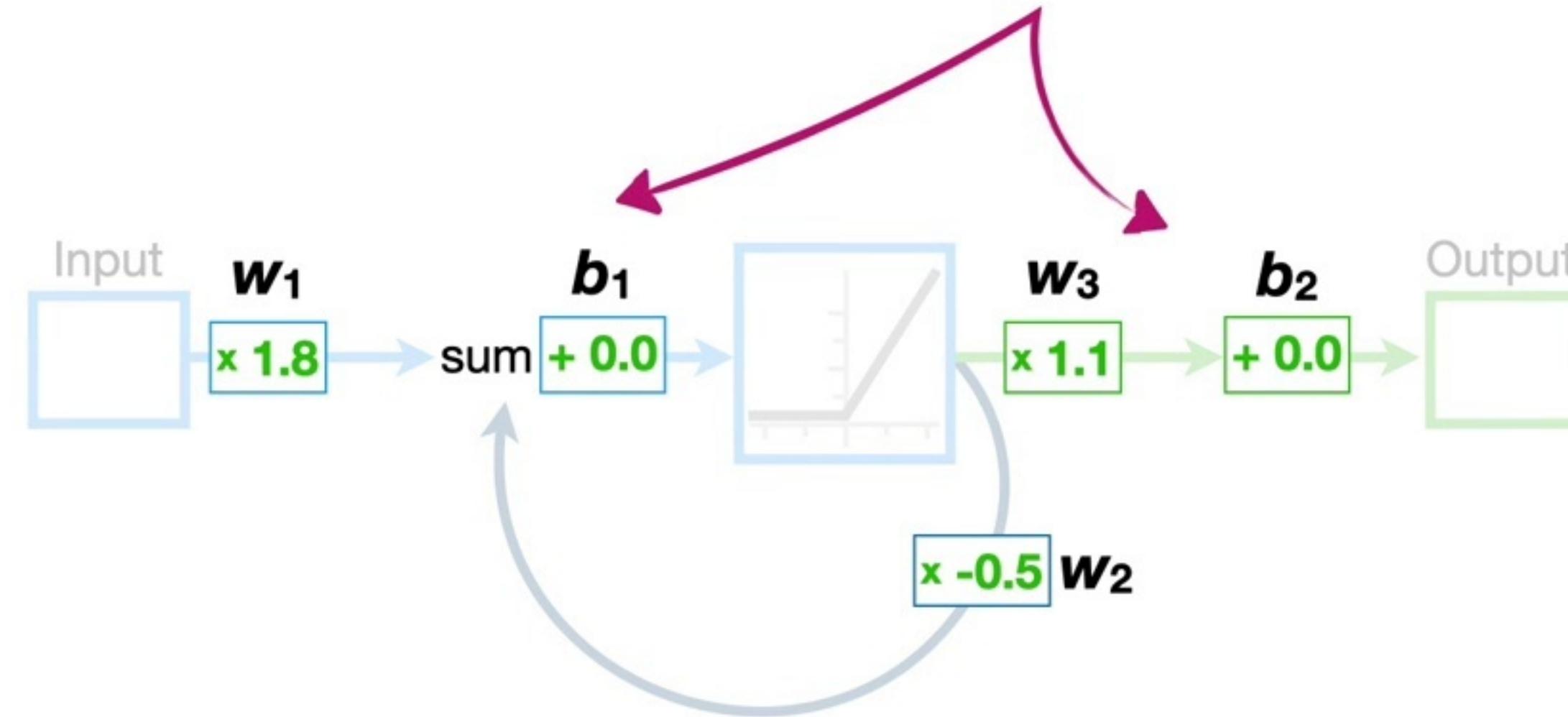


Just like the other neural networks that we've seen before, **Recurrent Neural Networks have weights, biases,**



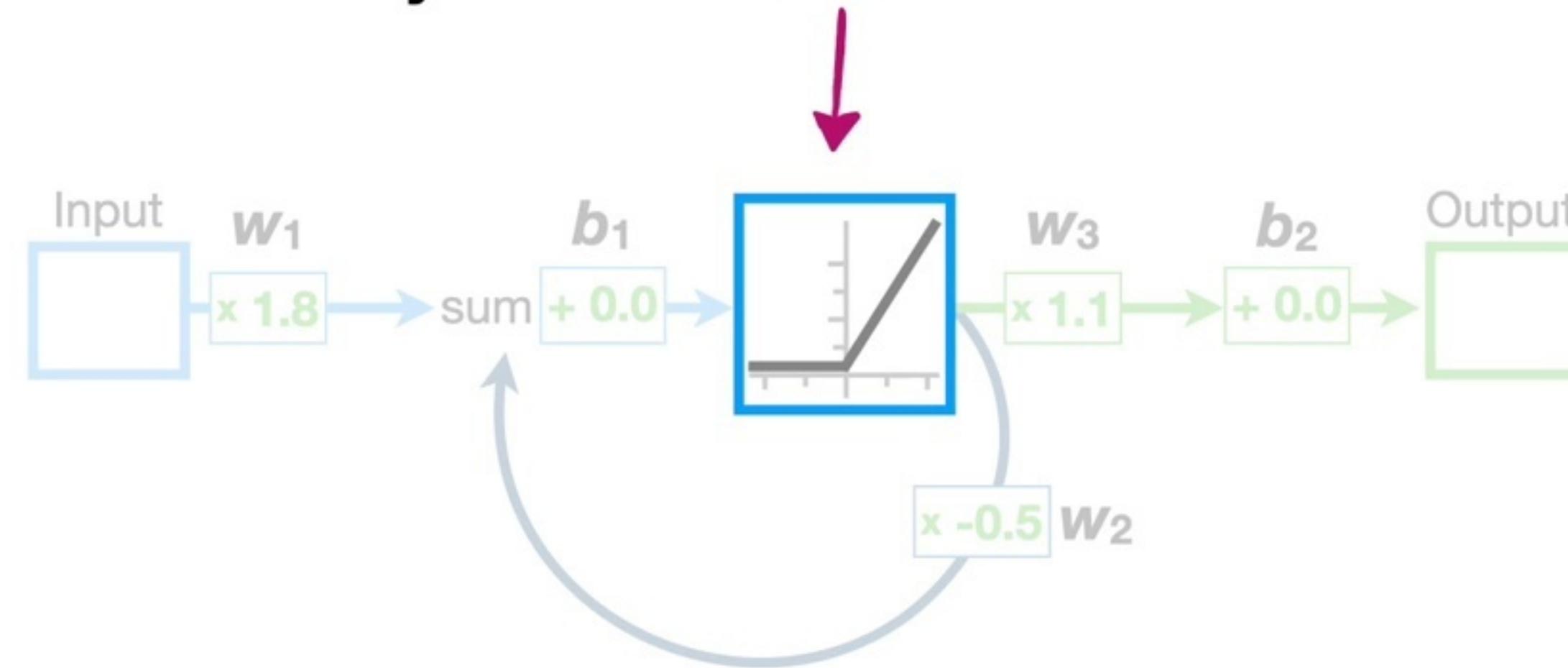


Just like the other neural networks that we've seen before, **Recurrent Neural Networks have weights, biases,**



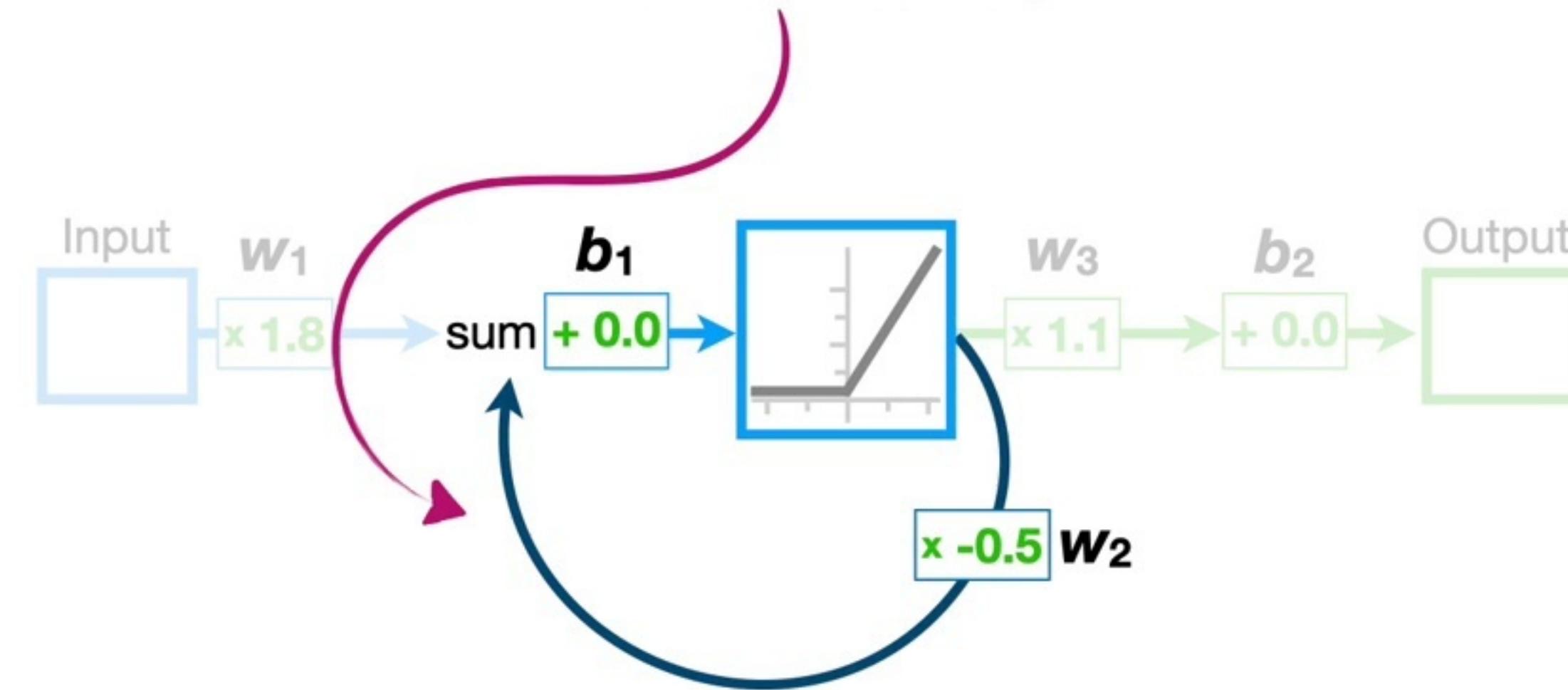


Just like the other neural networks that we've seen before, **Recurrent Neural Networks have weights, biases, layers and activation functions.**



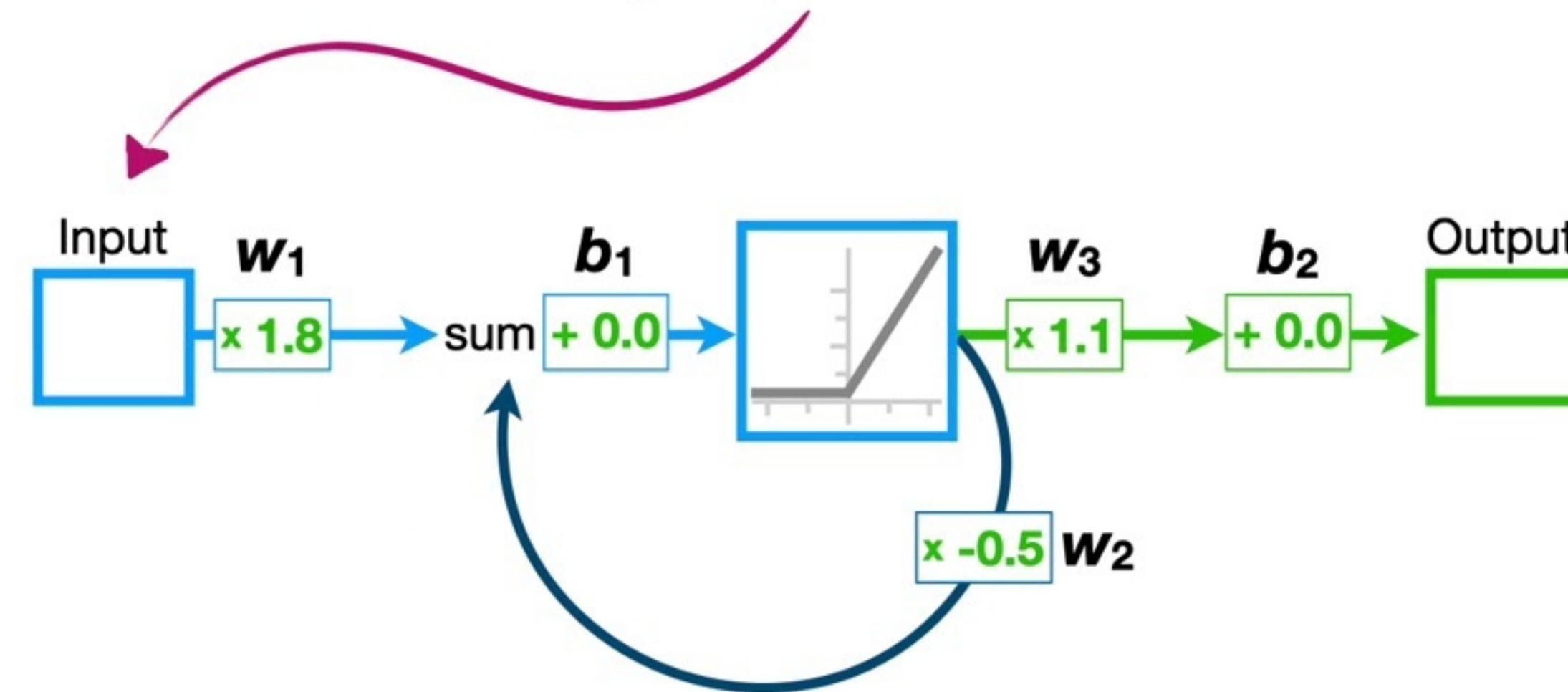


The big difference is that
Recurrent Neural Networks also
have **feedback loops**.



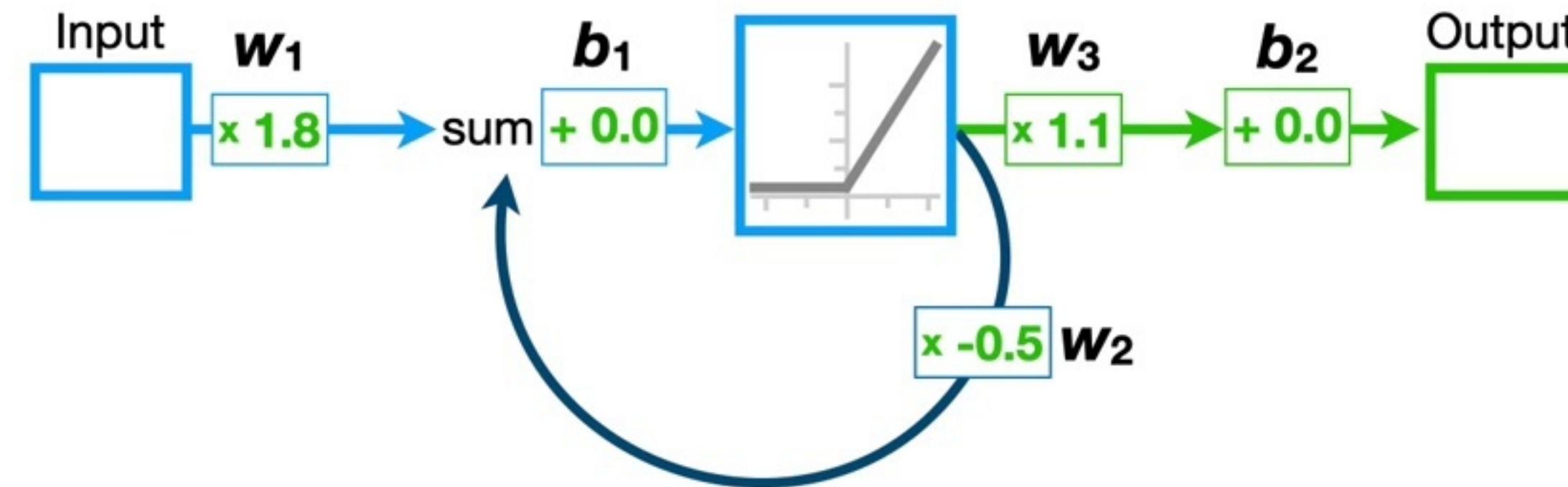


And, although this neural network
may look like it only takes a
single input value...



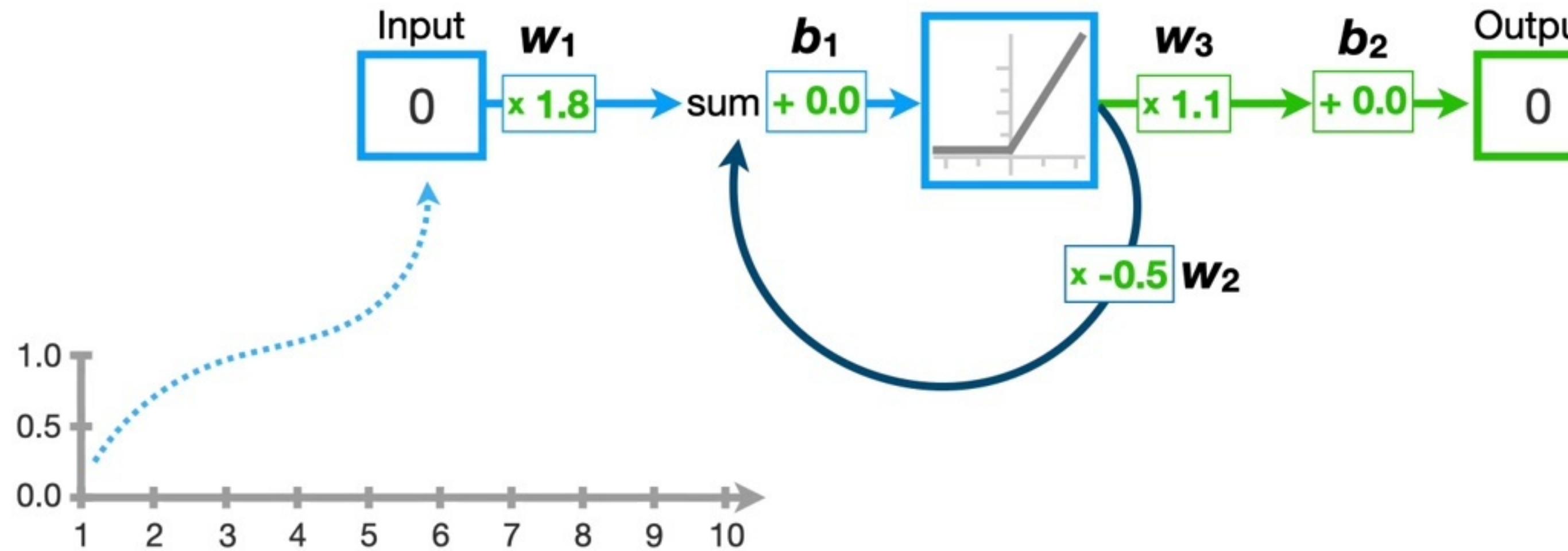


...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.



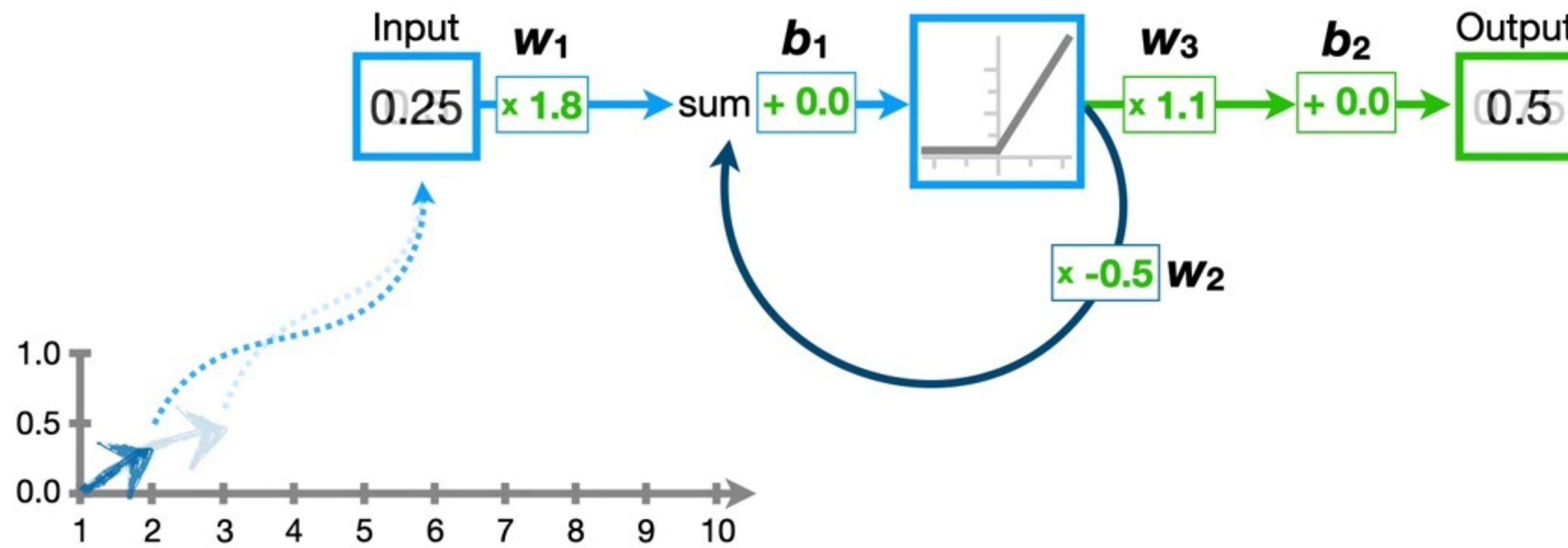


...the **feedback loop** makes it possible to use *sequential* input values, like stock market prices collected over time, to make predictions.



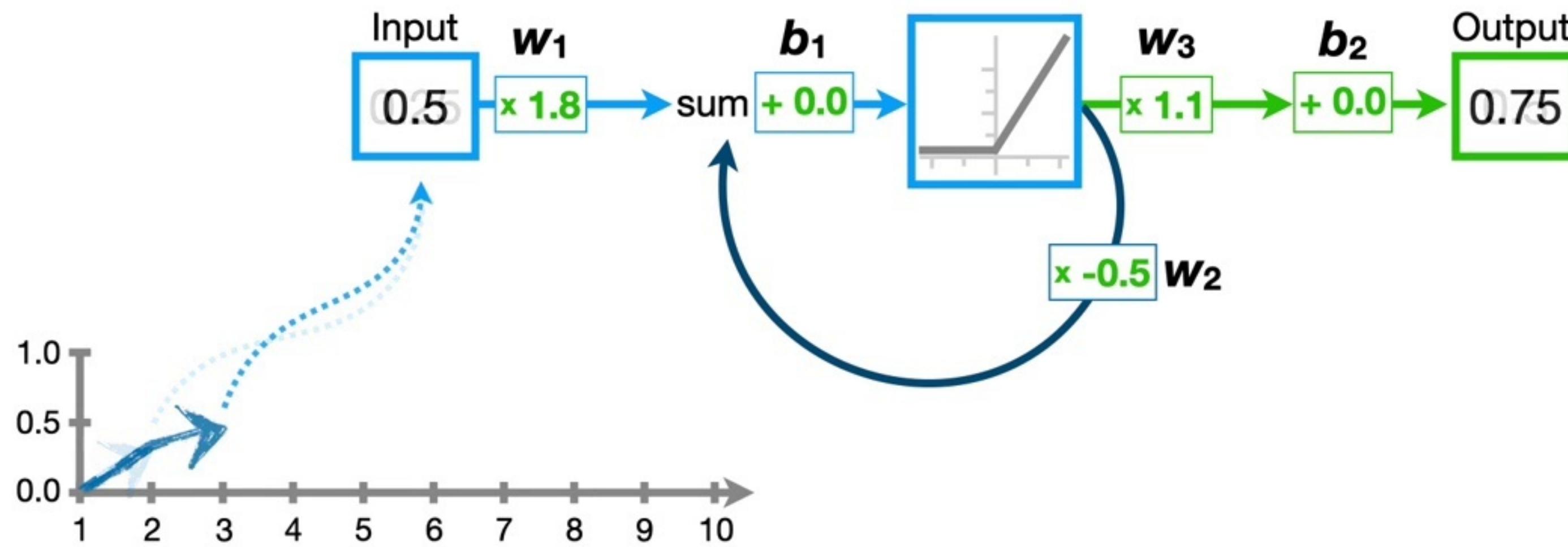


...the **feedback loop** makes it possible to use *sequential* input values, like stock market prices collected over time, to make predictions.



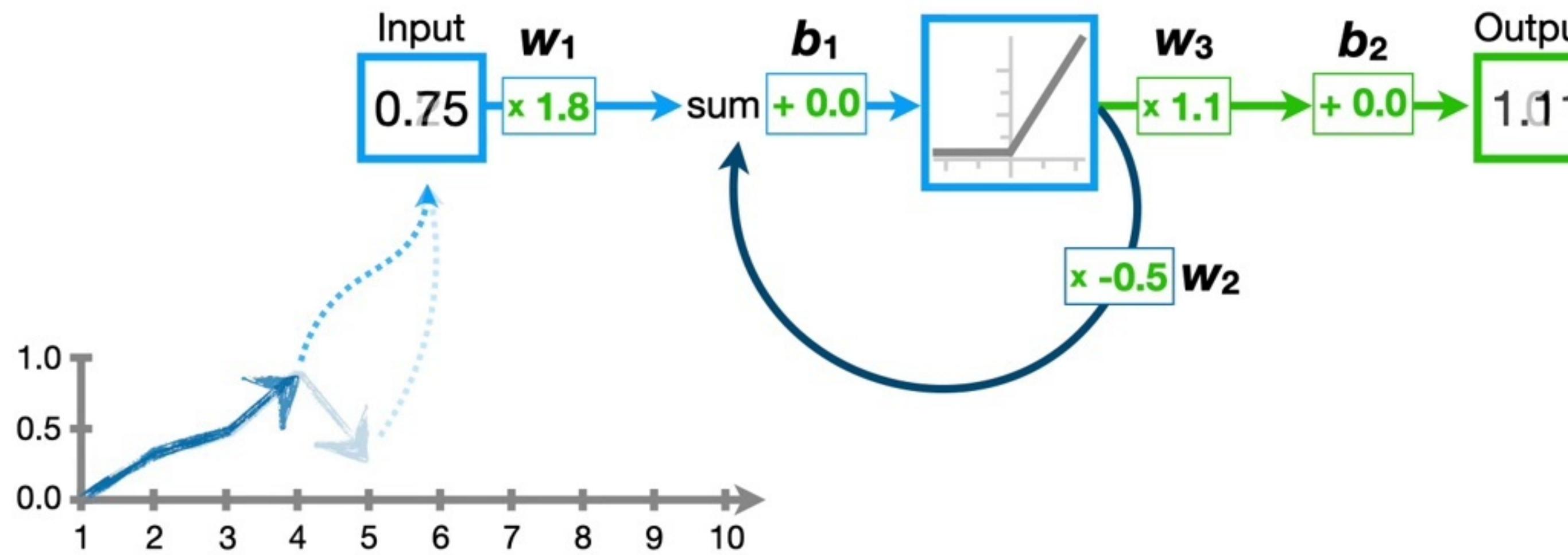


...the **feedback loop** makes it possible to use *sequential* input values, like stock market prices collected over time, to make predictions.



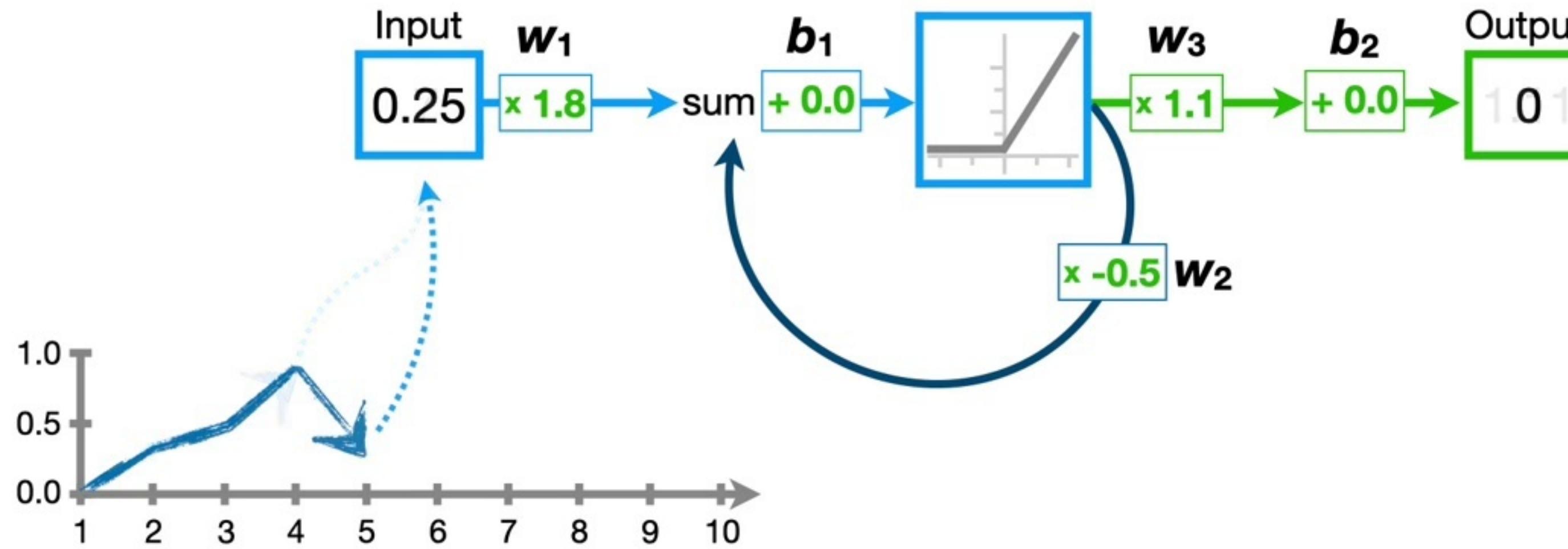


...the **feedback loop** makes it possible to use *sequential* input values, like stock market prices collected over time, to make predictions.



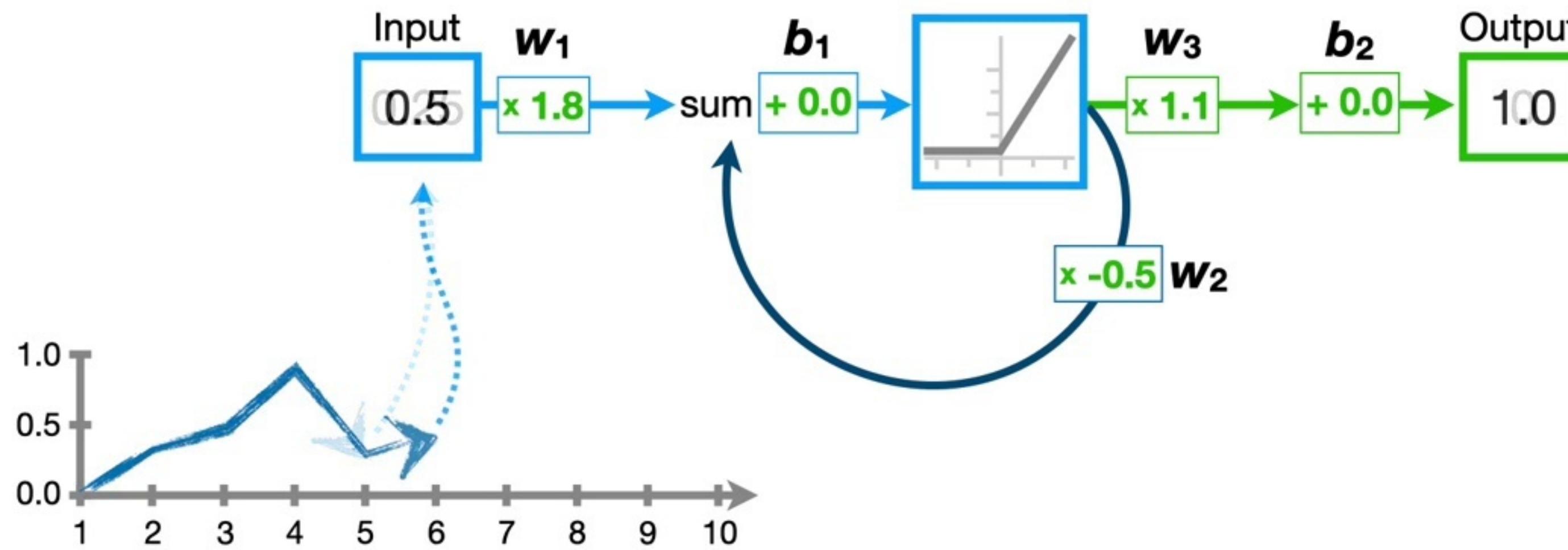


...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.



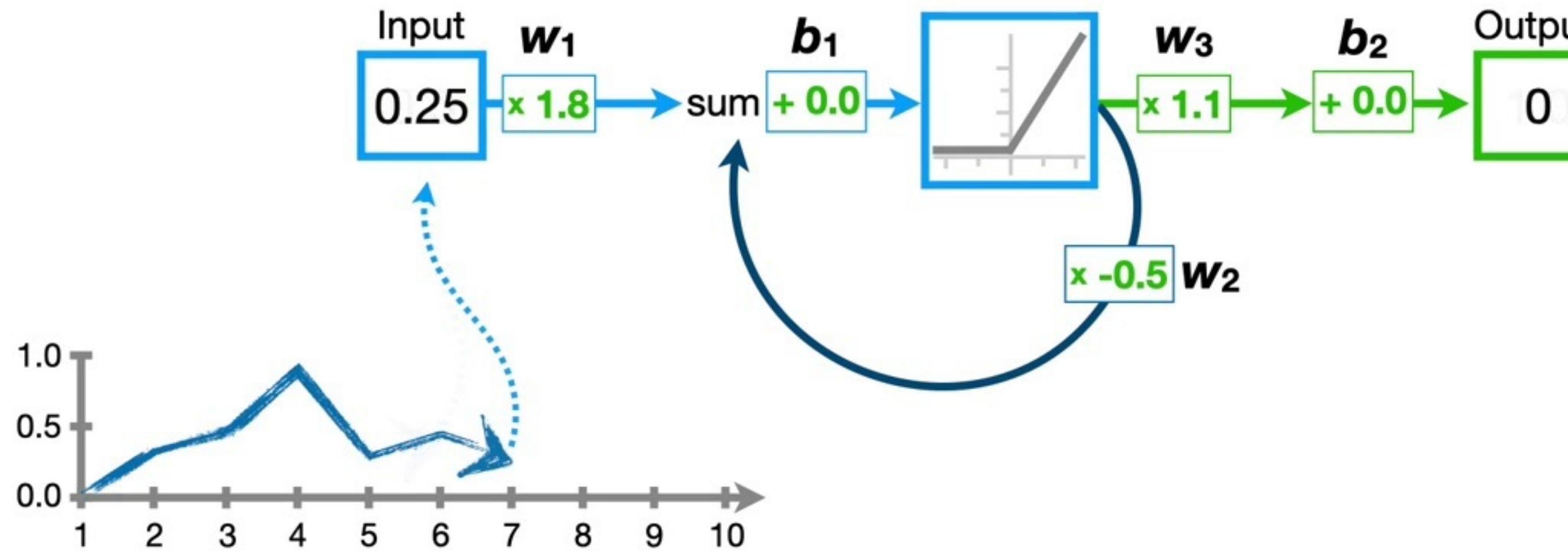


...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.



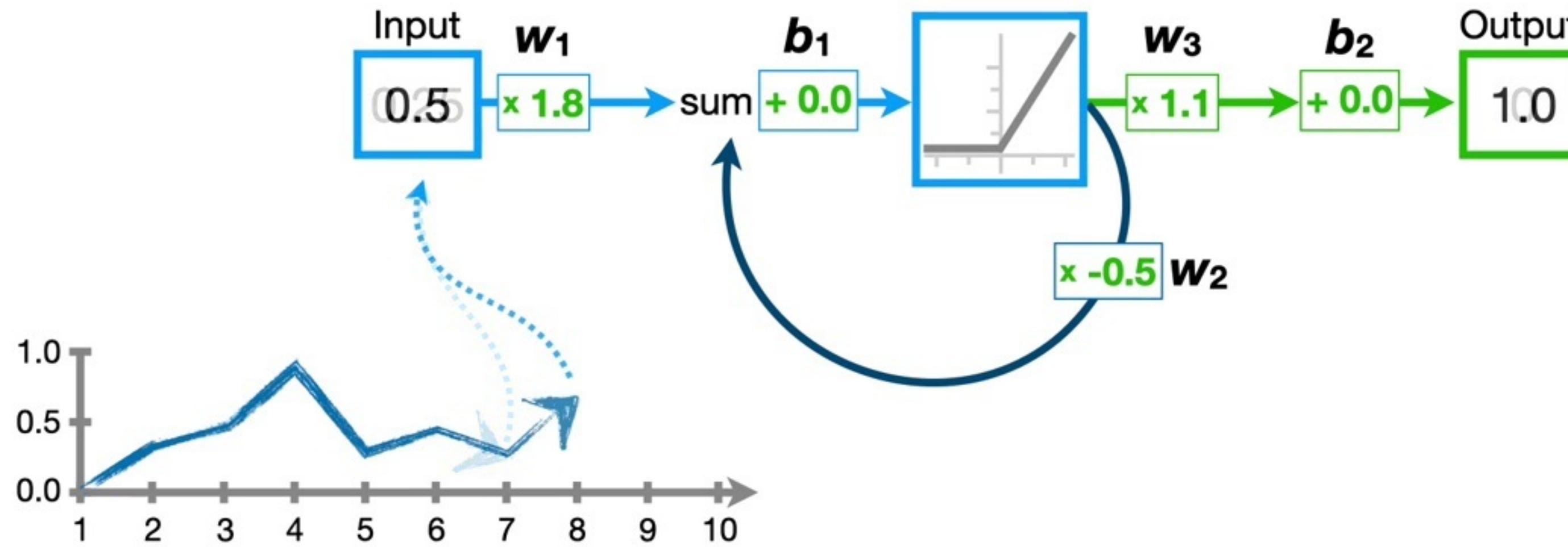


...the **feedback loop** makes it possible to use *sequential* input values, like stock market prices collected over time, to make predictions.



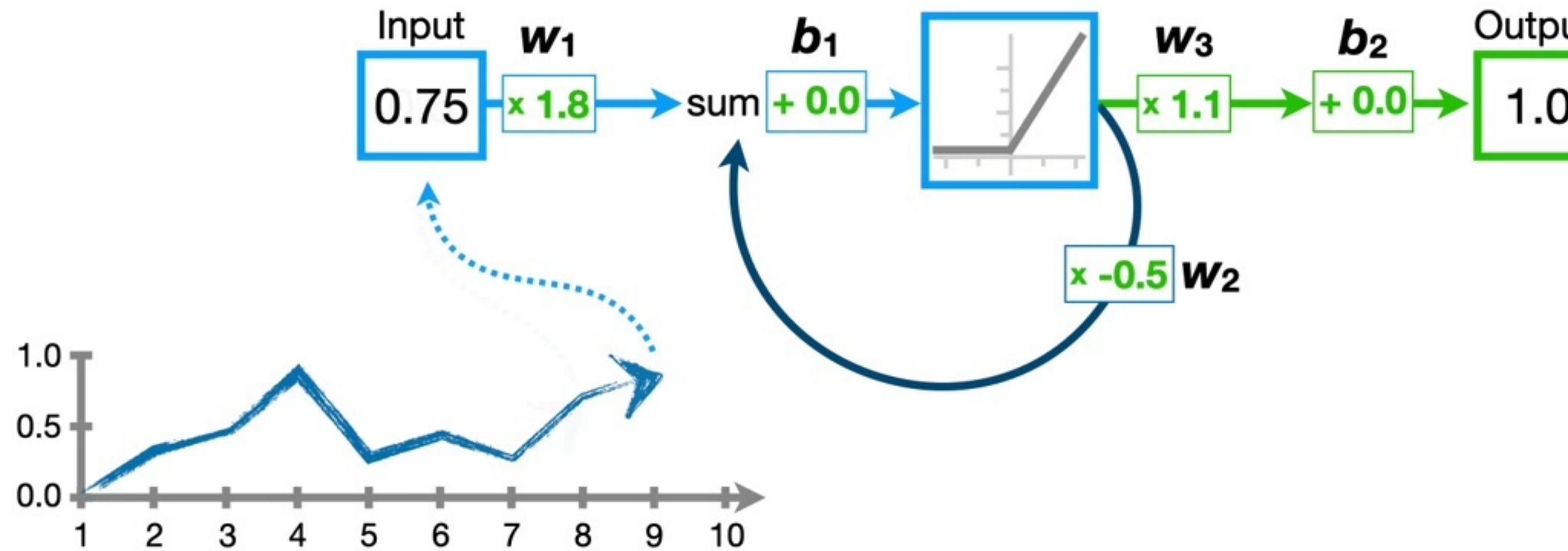


...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.



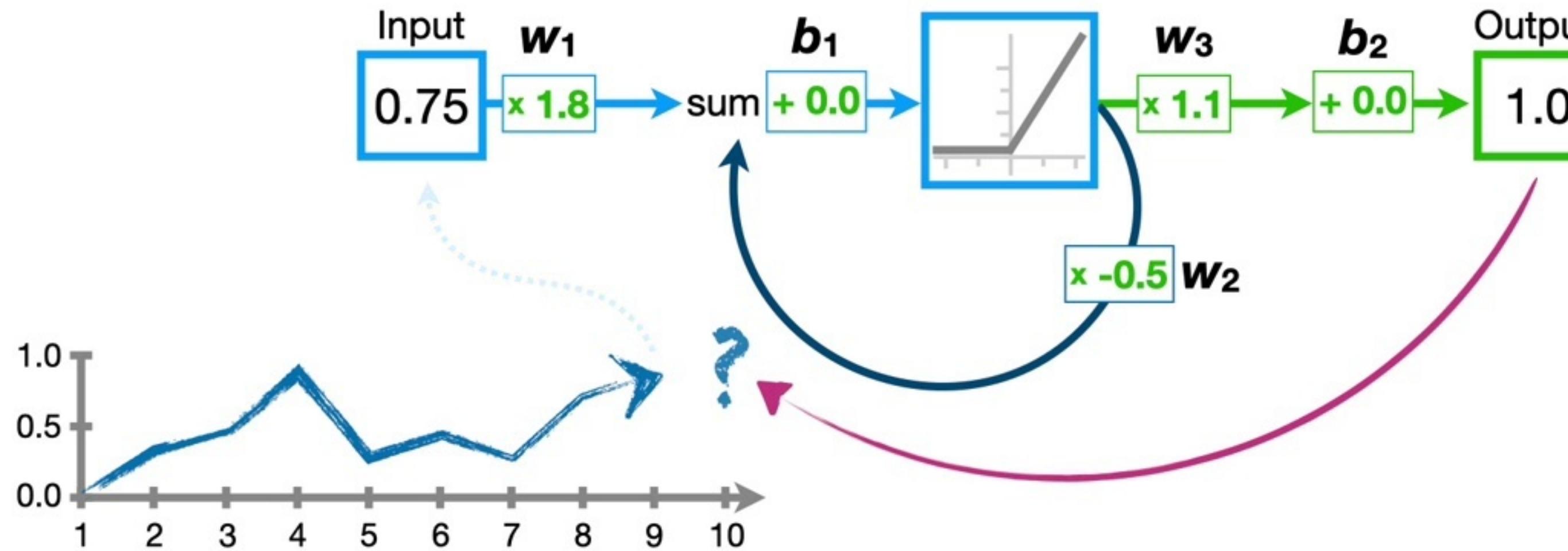


...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.



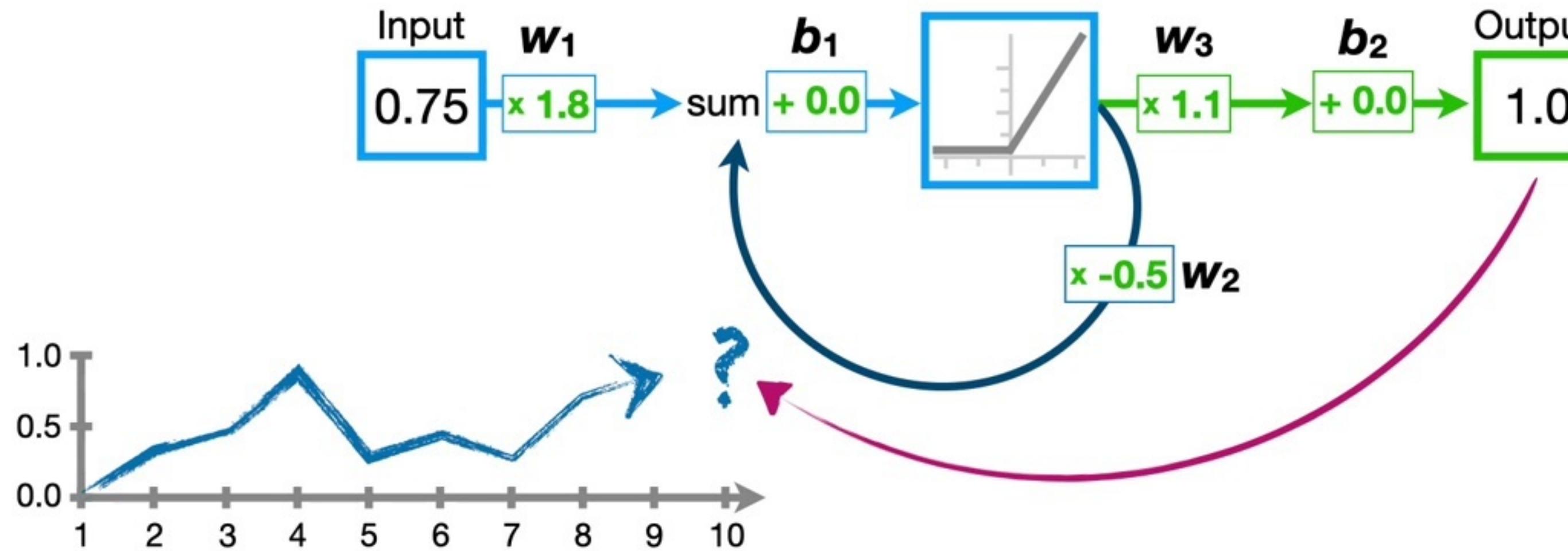


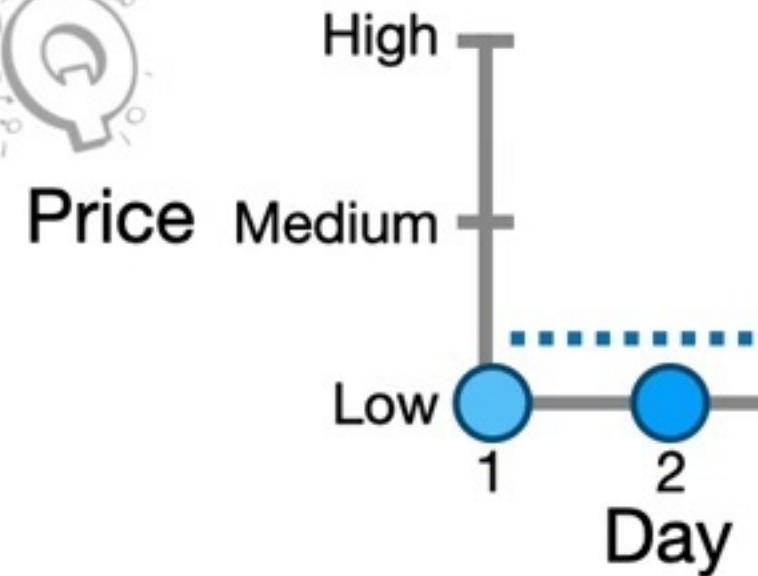
...the **feedback loop** makes it possible
to use *sequential* input values, like
stock market prices collected over
time, to make predictions.





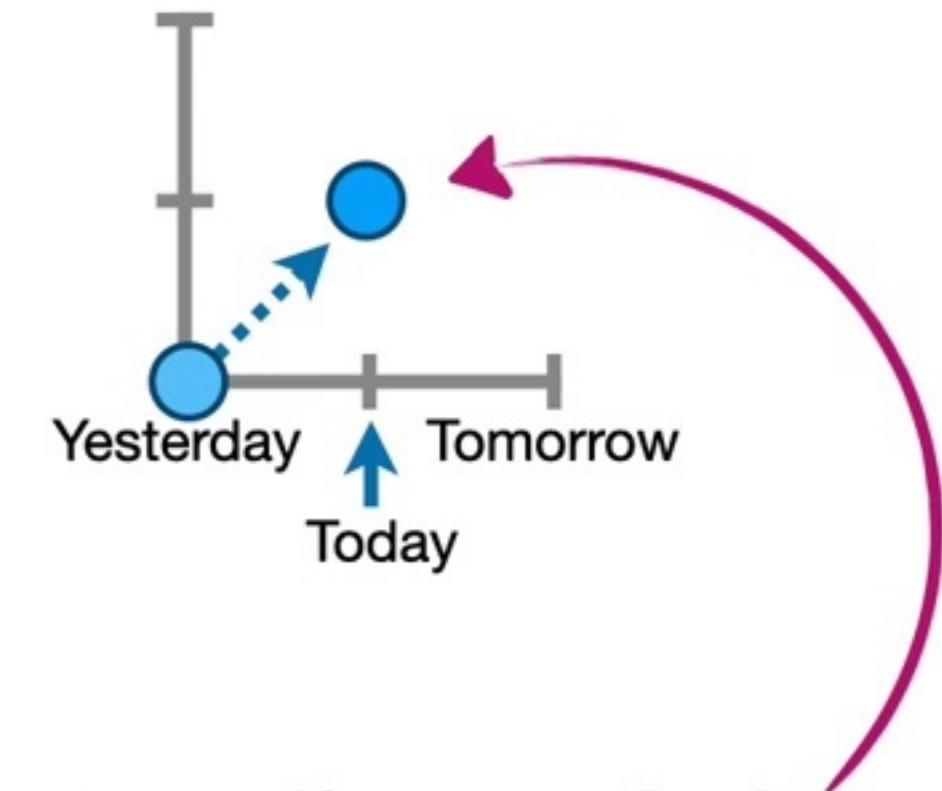
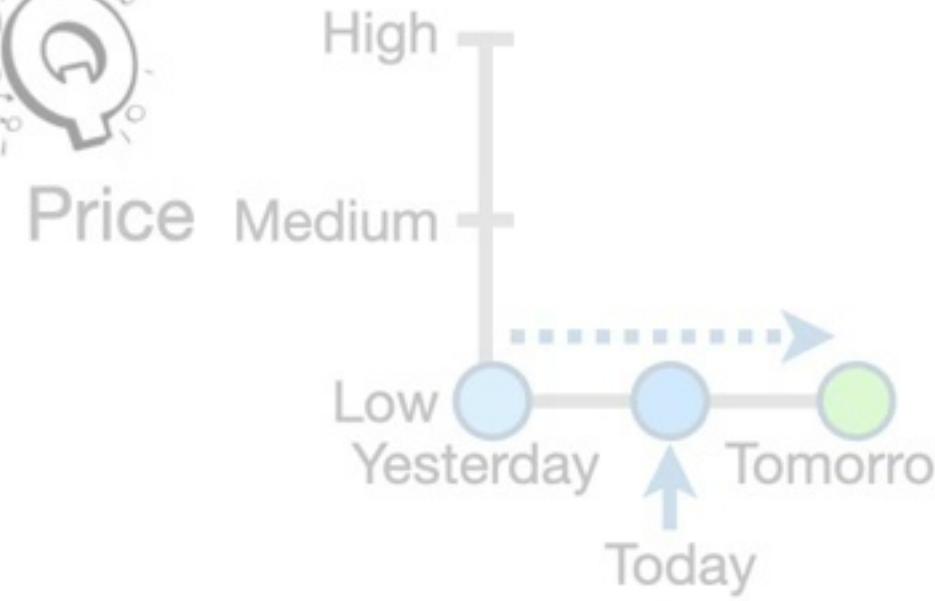
To understand how, exactly, this
Recurrent Neural Network can make
predictions with sequential input values...



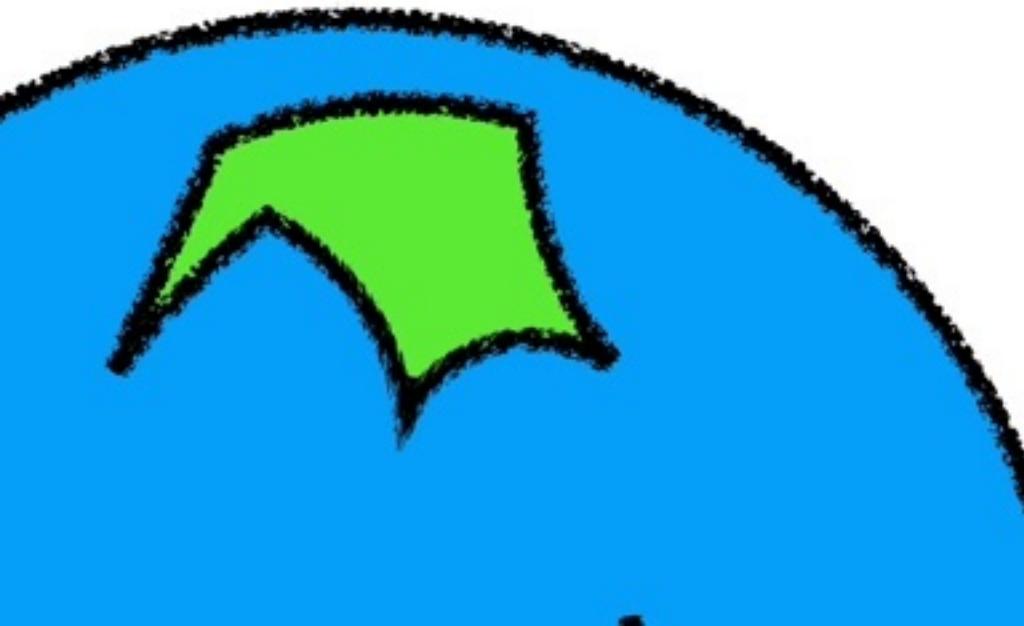


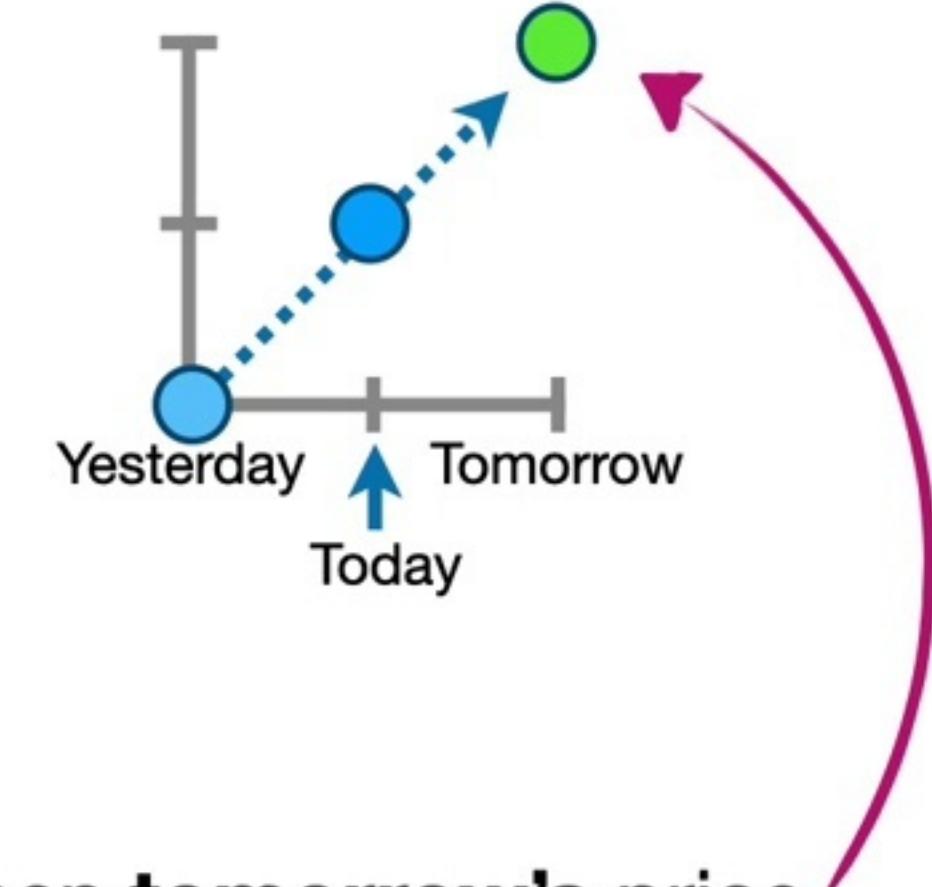
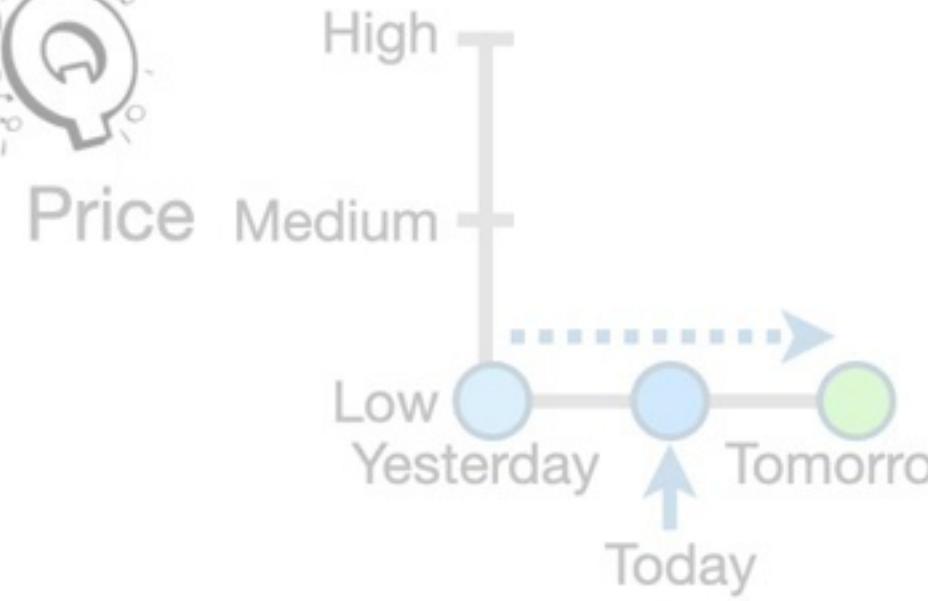
...then, more often than not, the price remains low on the next day.





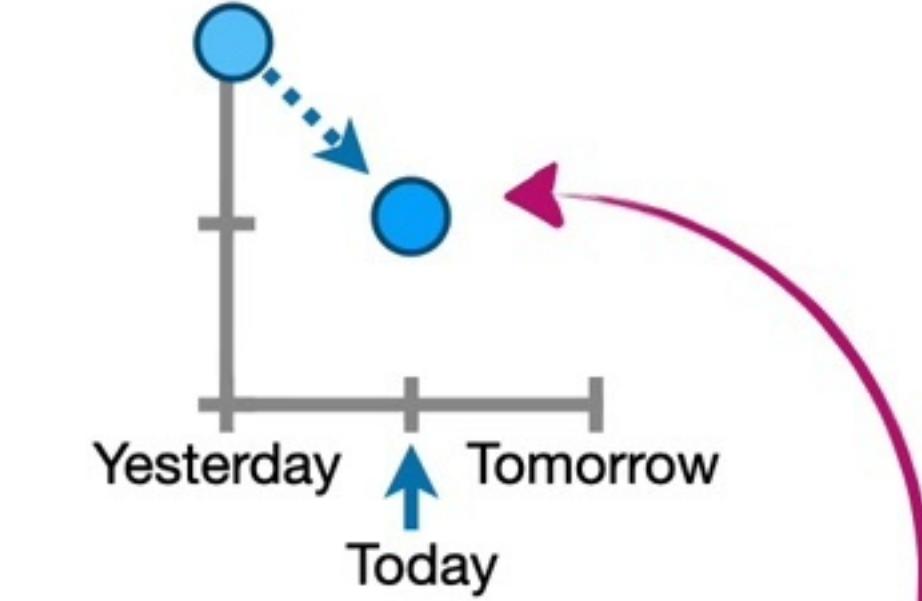
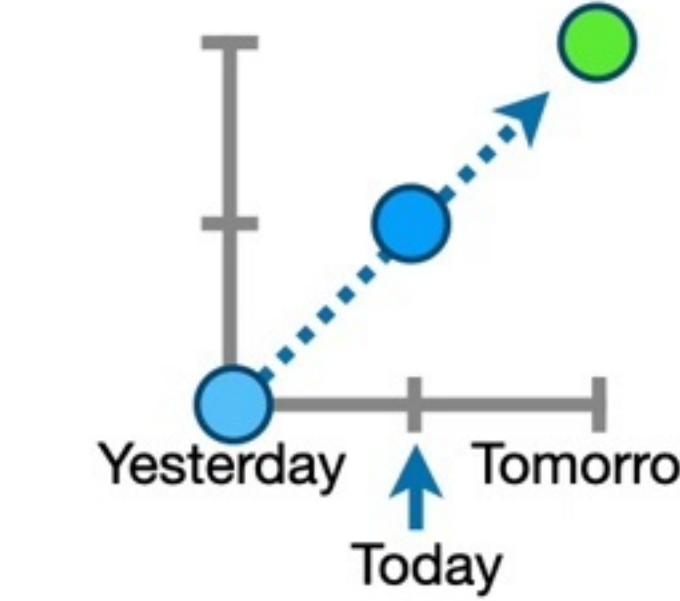
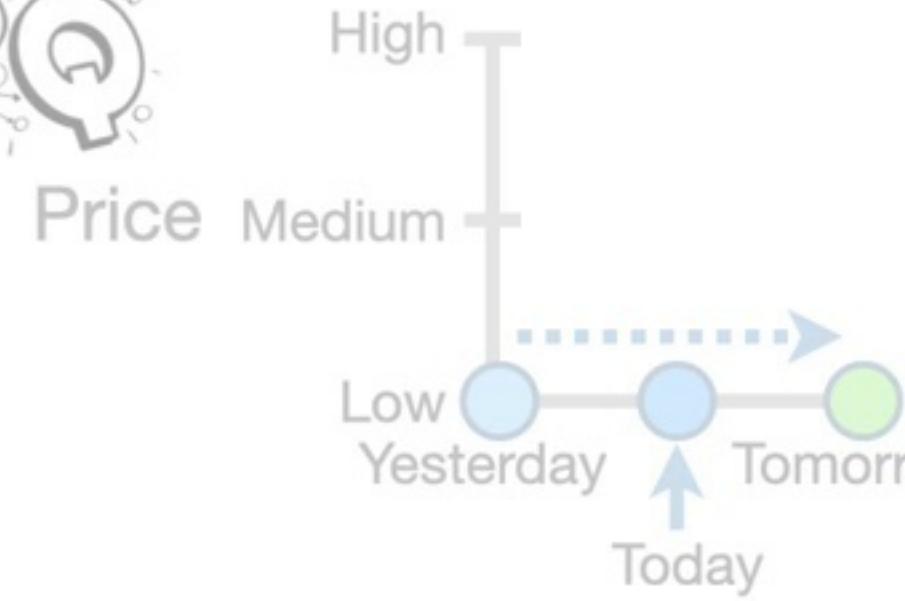
In contrast, if **yesterday's** price was low and **today's** price is medium...





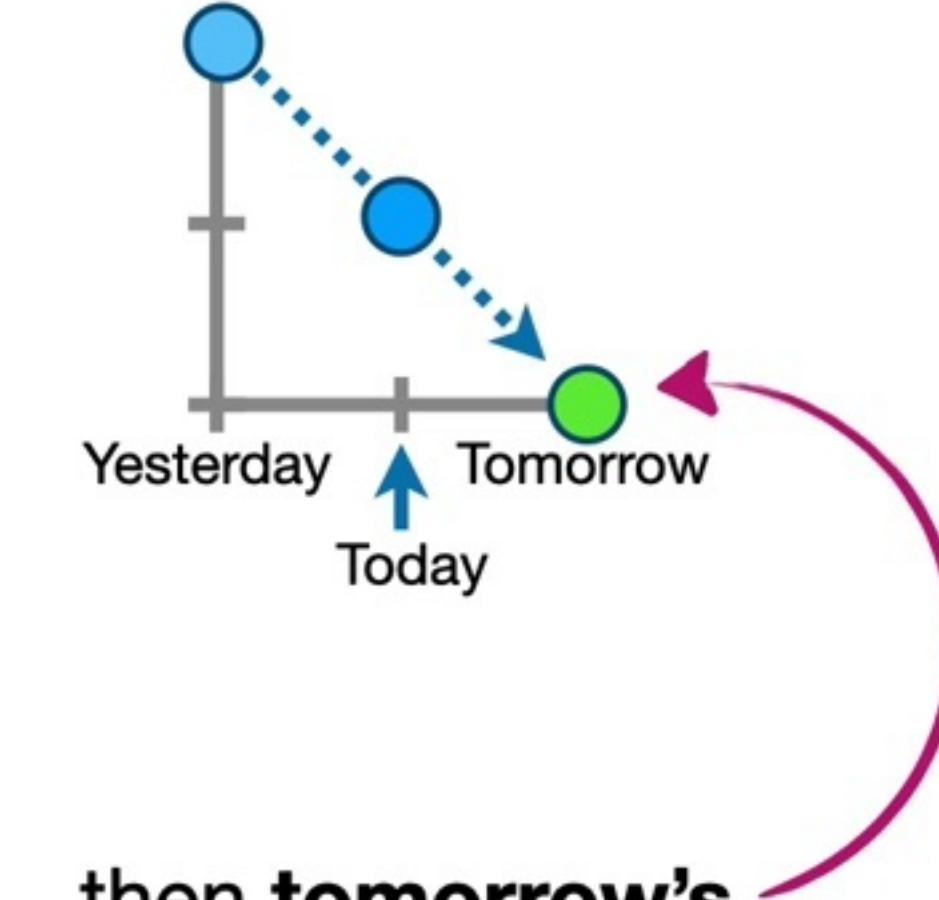
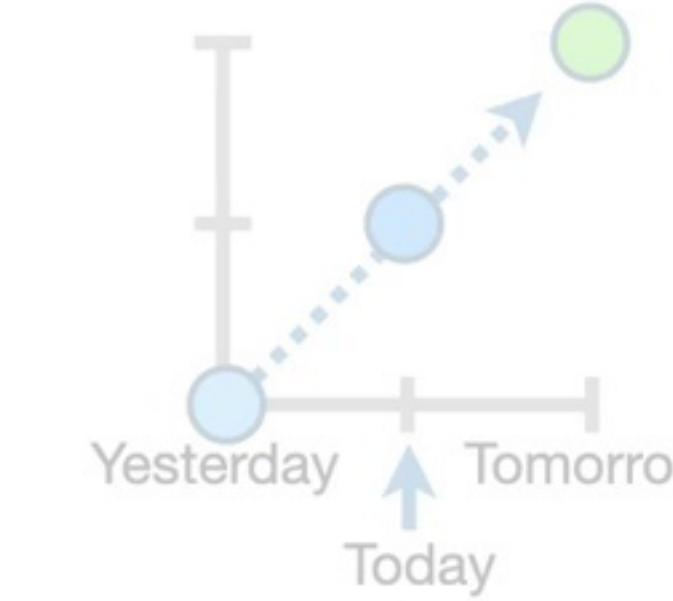
...then **tomorrow's** price
should be even higher.



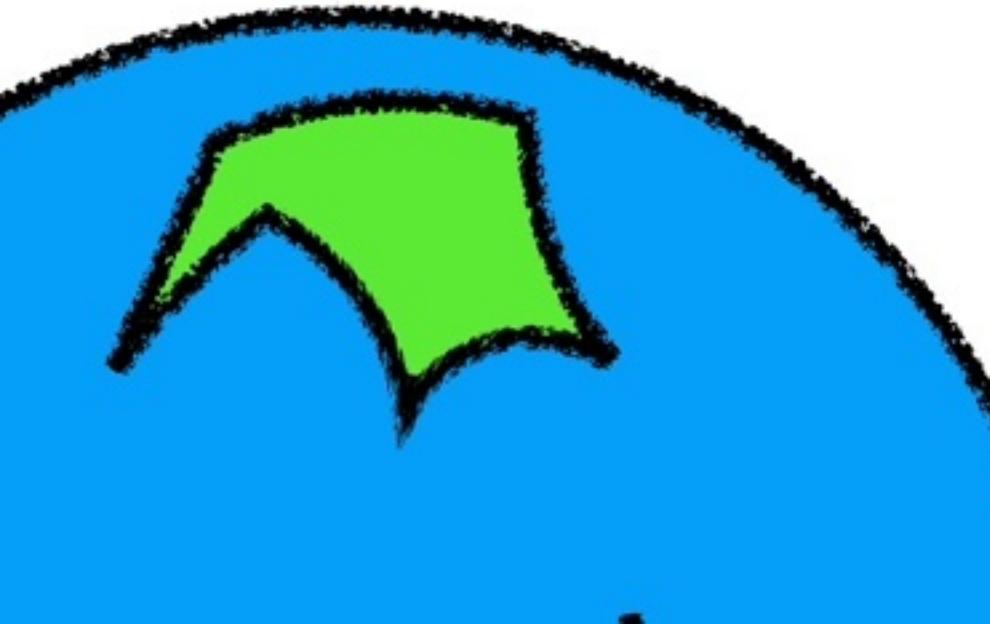


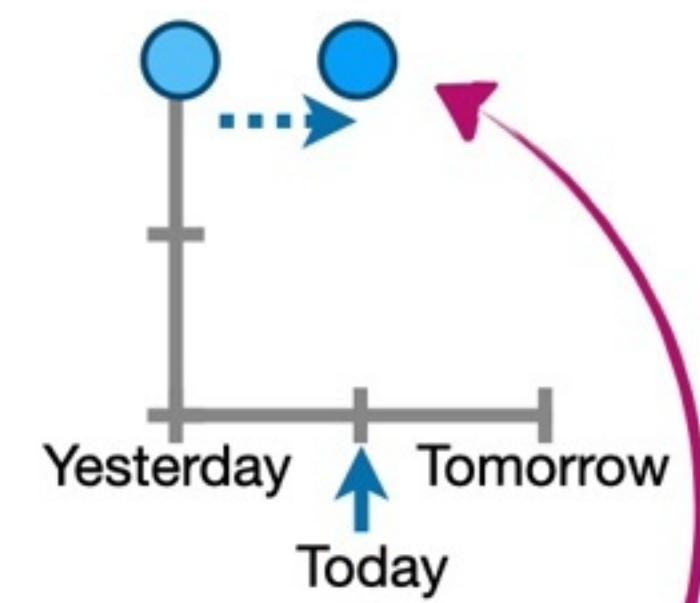
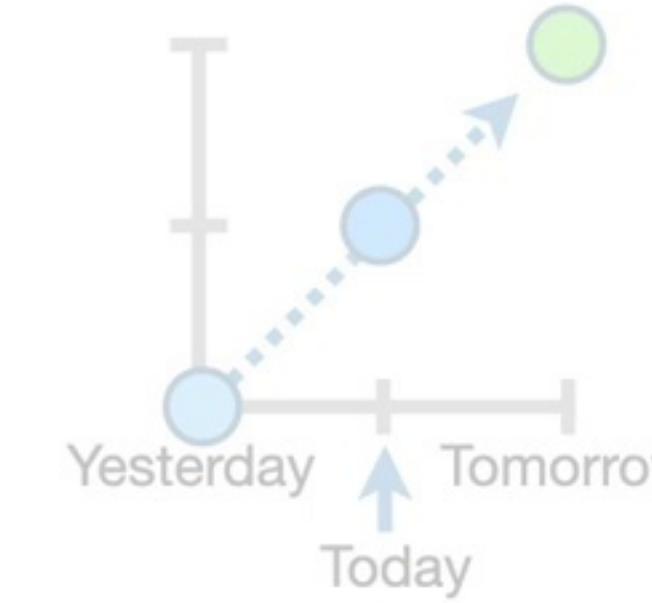
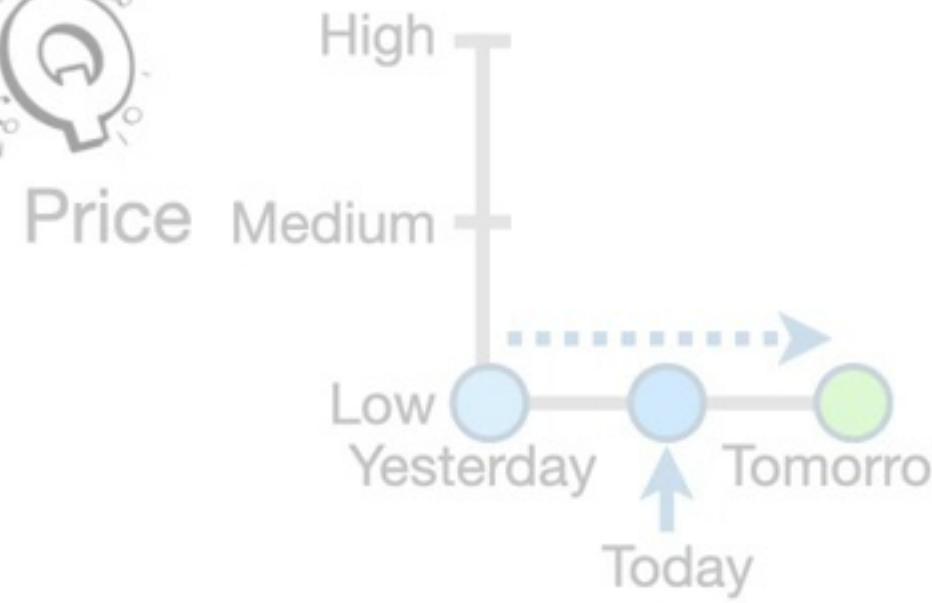
And when the price decreases from high to medium...





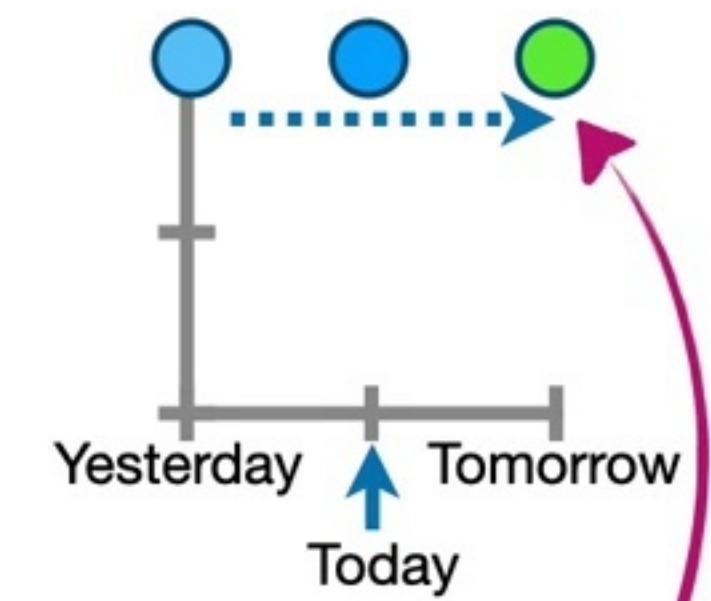
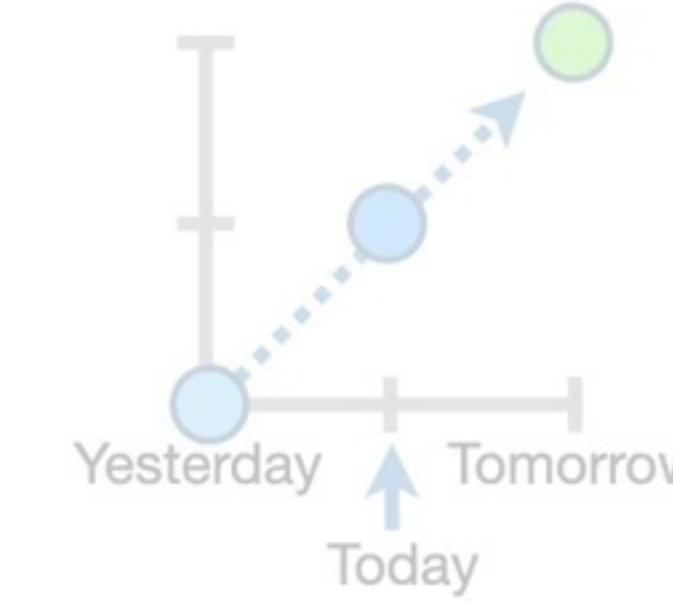
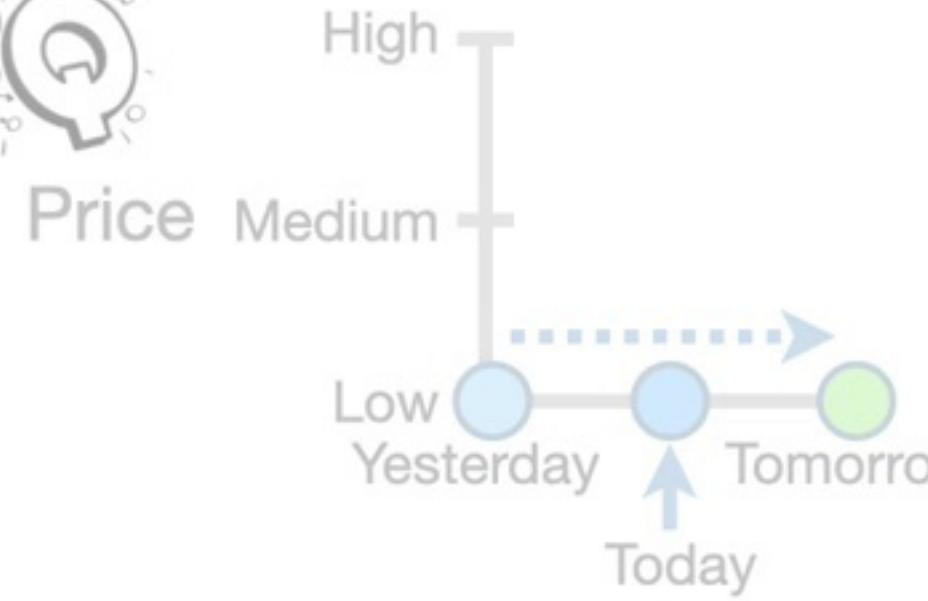
...then **tomorrow's**
price will be even lower.



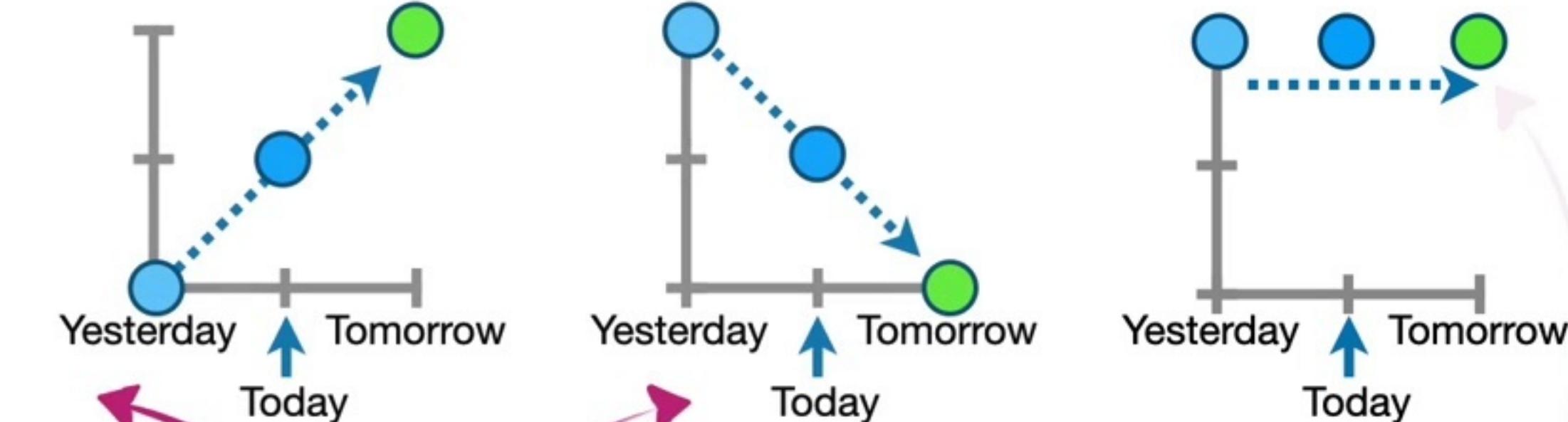
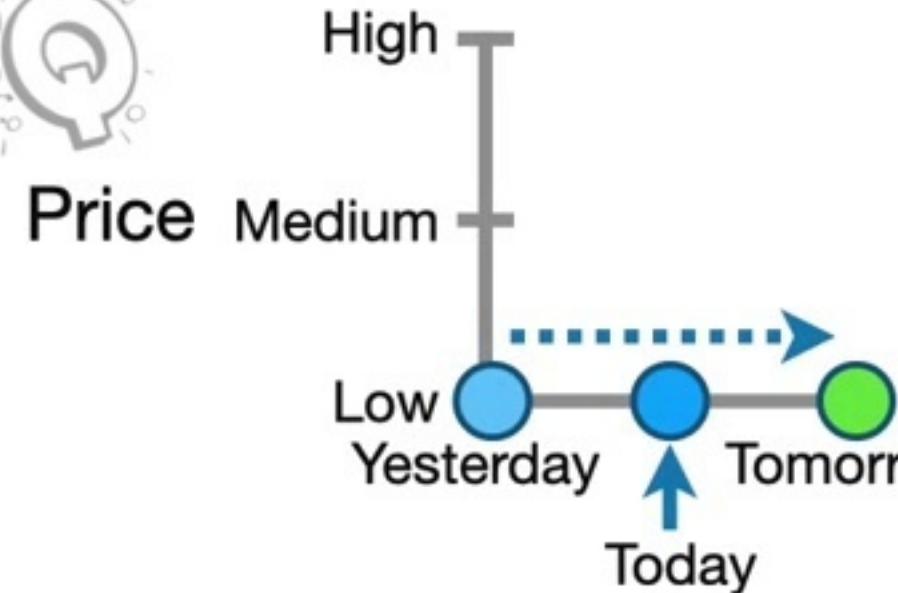


Lastly, if the price stays high for two days in a row...



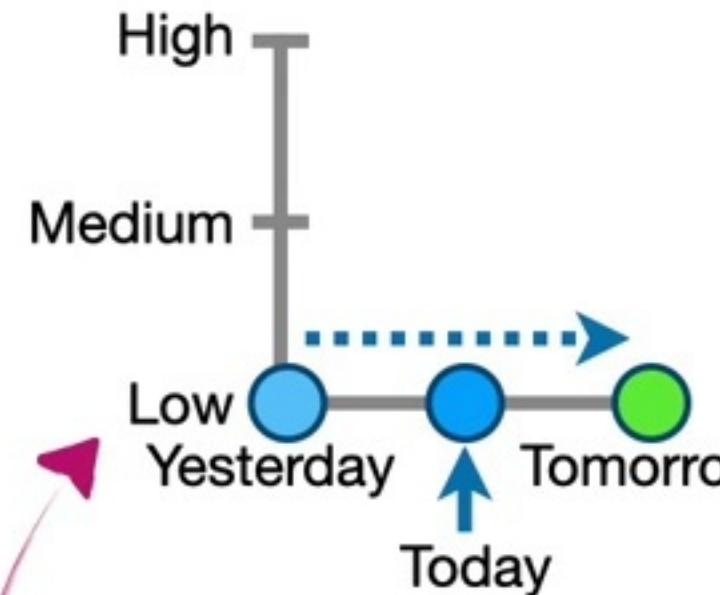


...then the price
will be high
tomorrow.

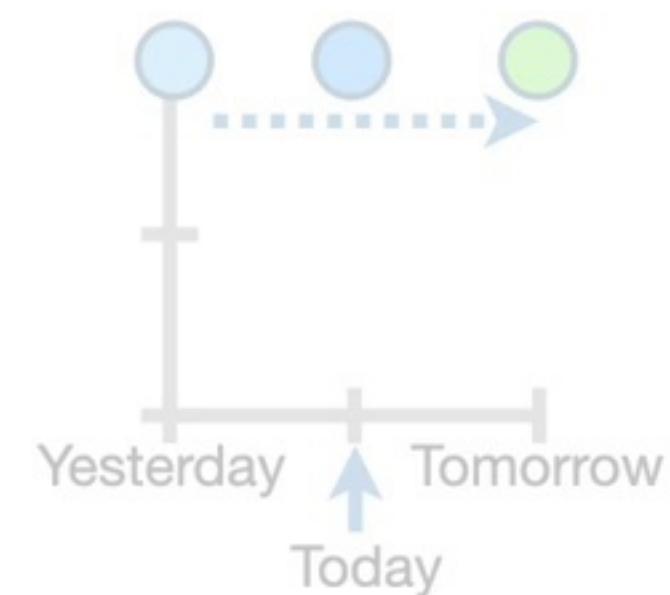
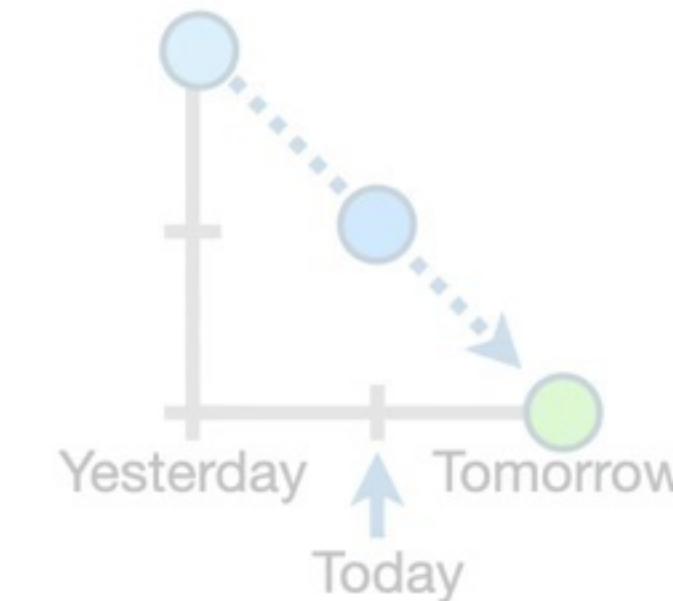
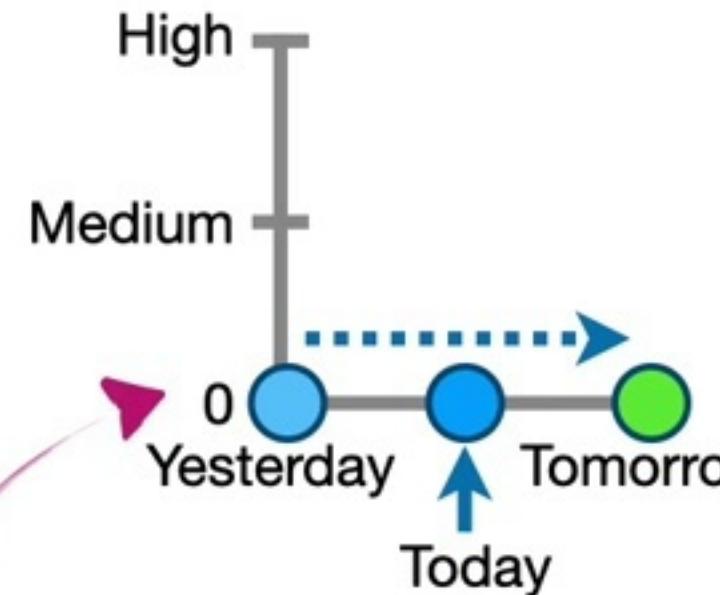


Now that we see the
general trends in stock
prices in **StatLand**...

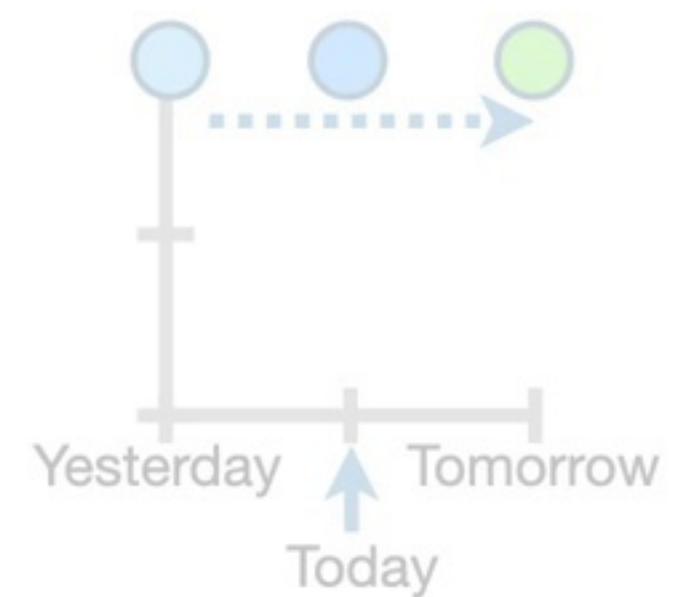
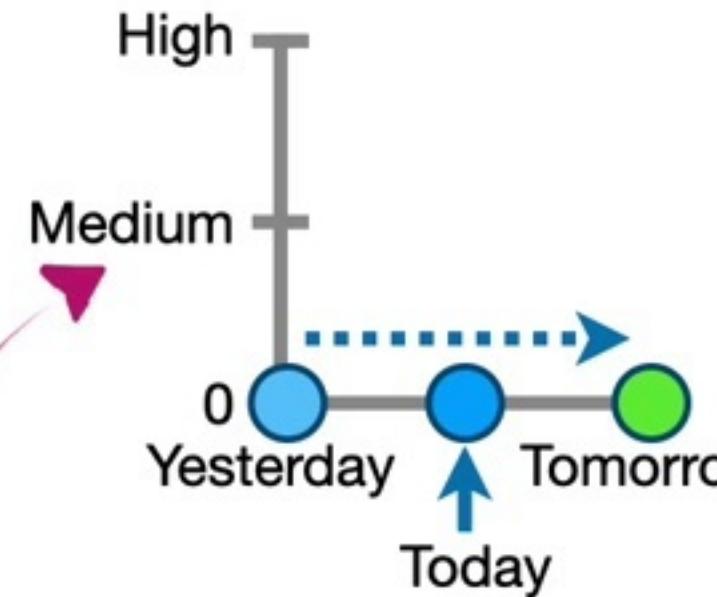
...then the price
will be high
tomorrow.



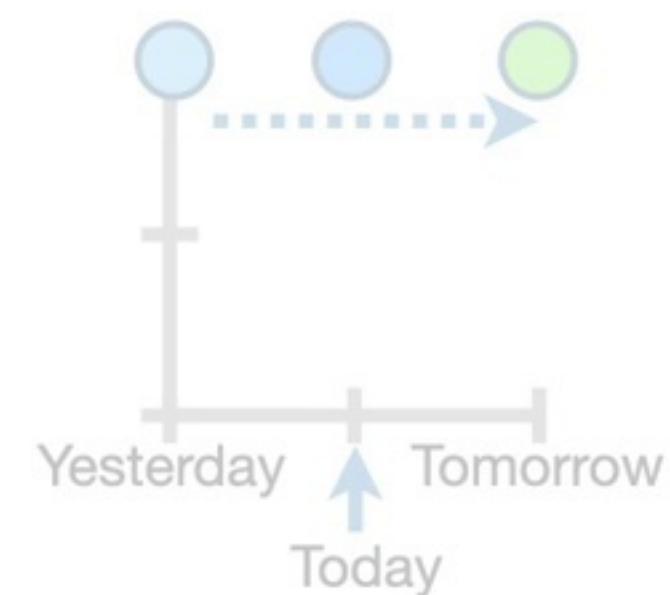
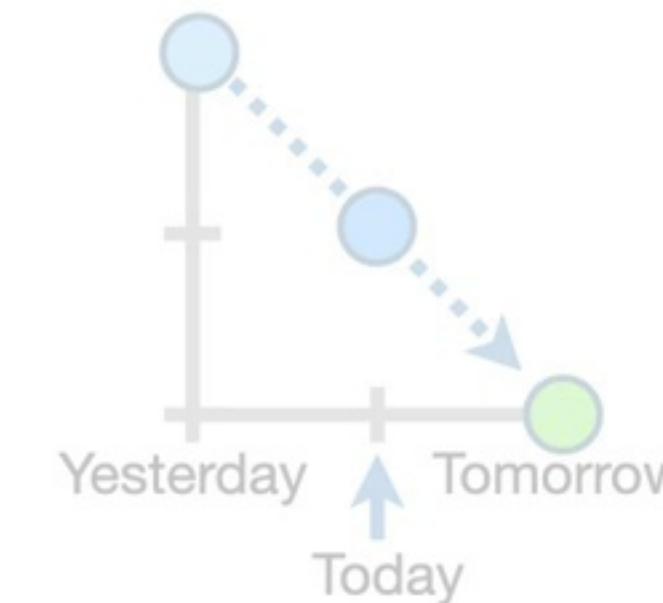
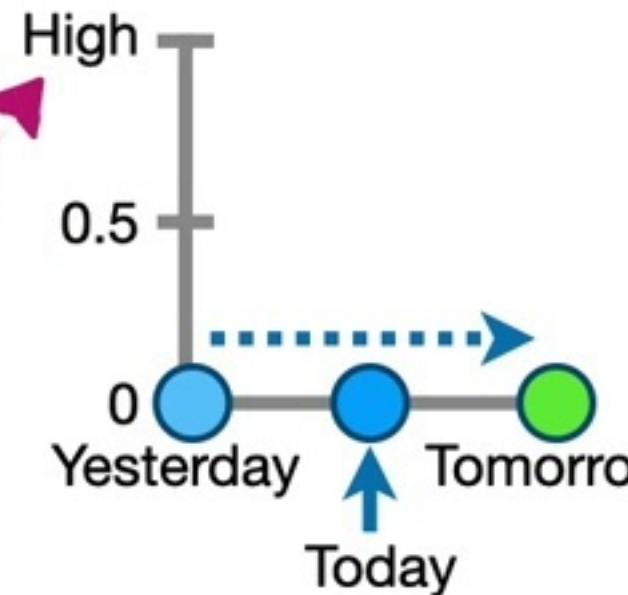
The first thing we'll do scale the
prices so that **Low = 0...**



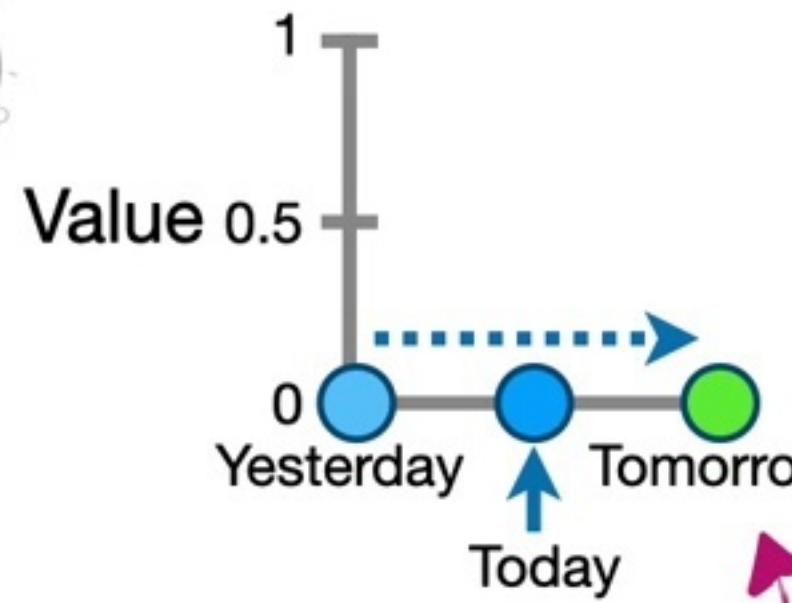
The first thing we'll do scale the
prices so that **Low = 0**...



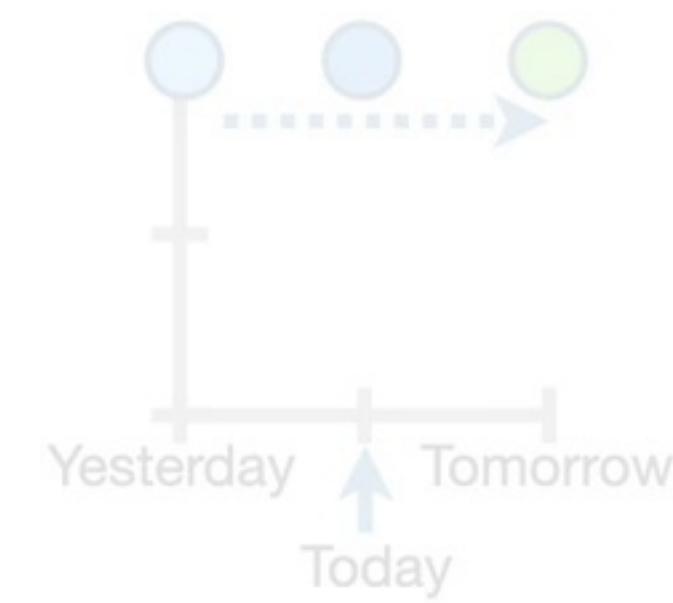
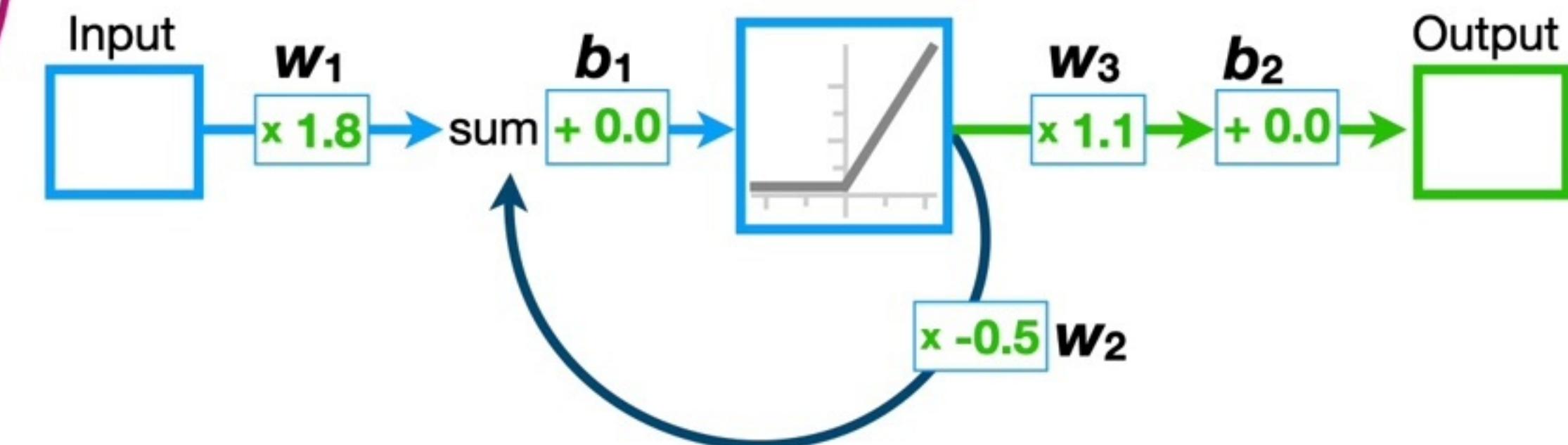
...Medium = 0.5...

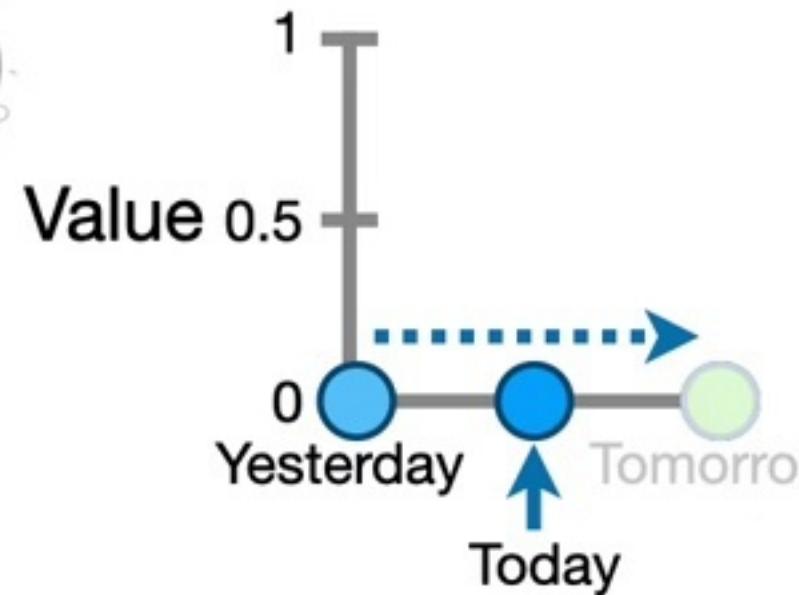


...and **High** = 1.

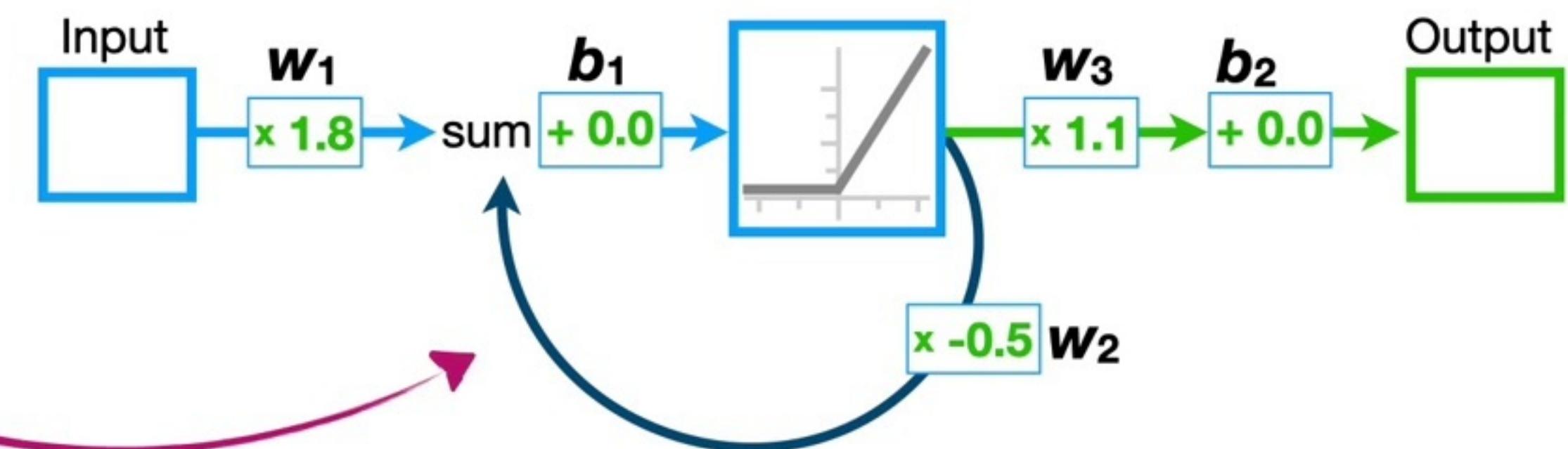


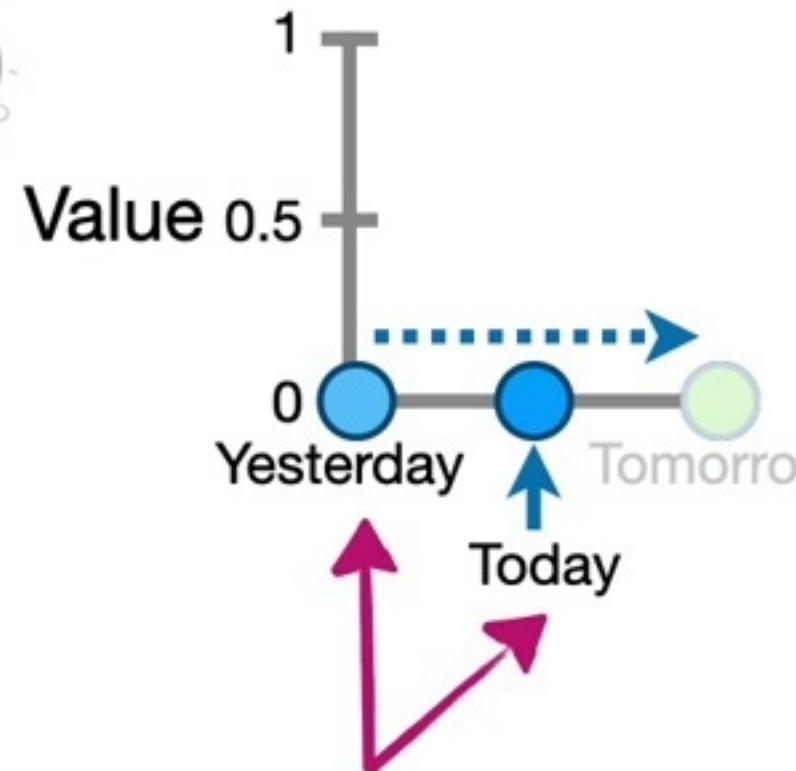
...and see if it can
correctly predict
tomorrow's value.



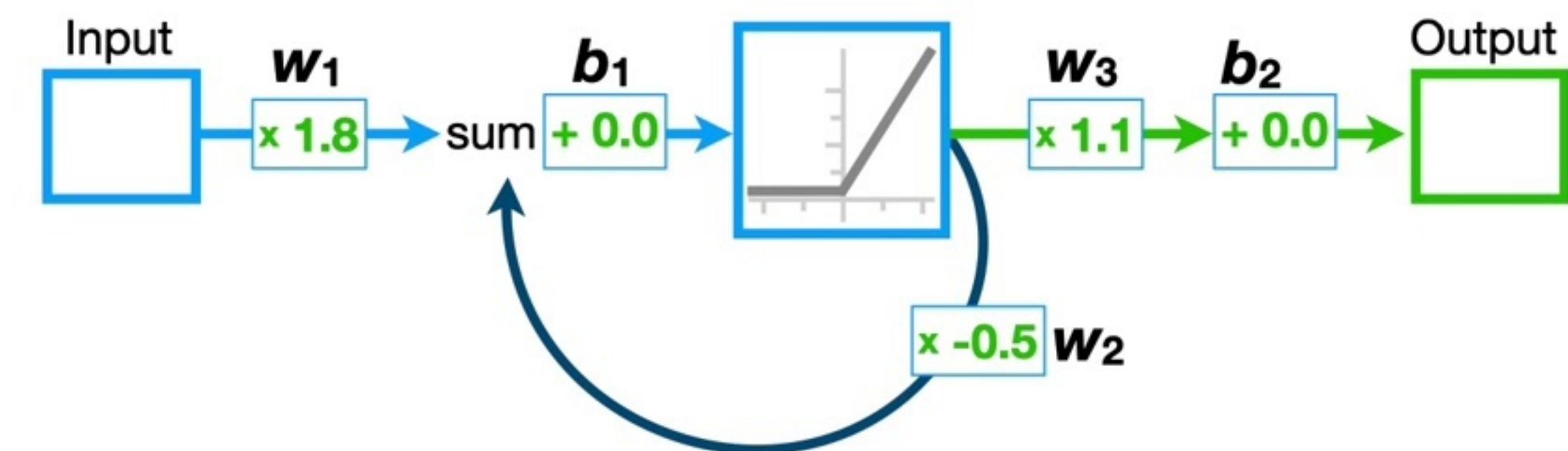


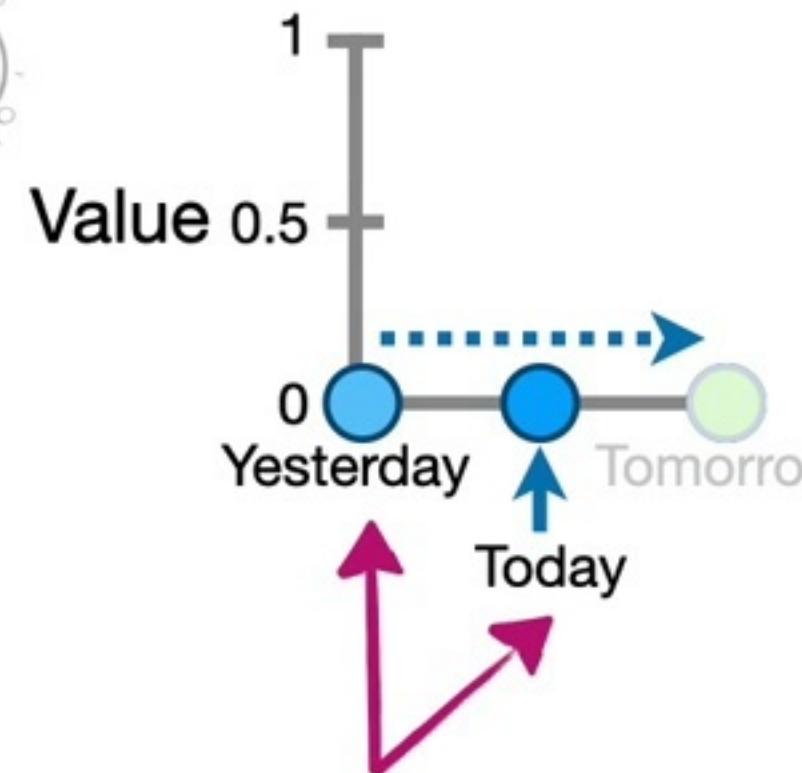
Now, because the recurrent neural network has a **feedback loop**...



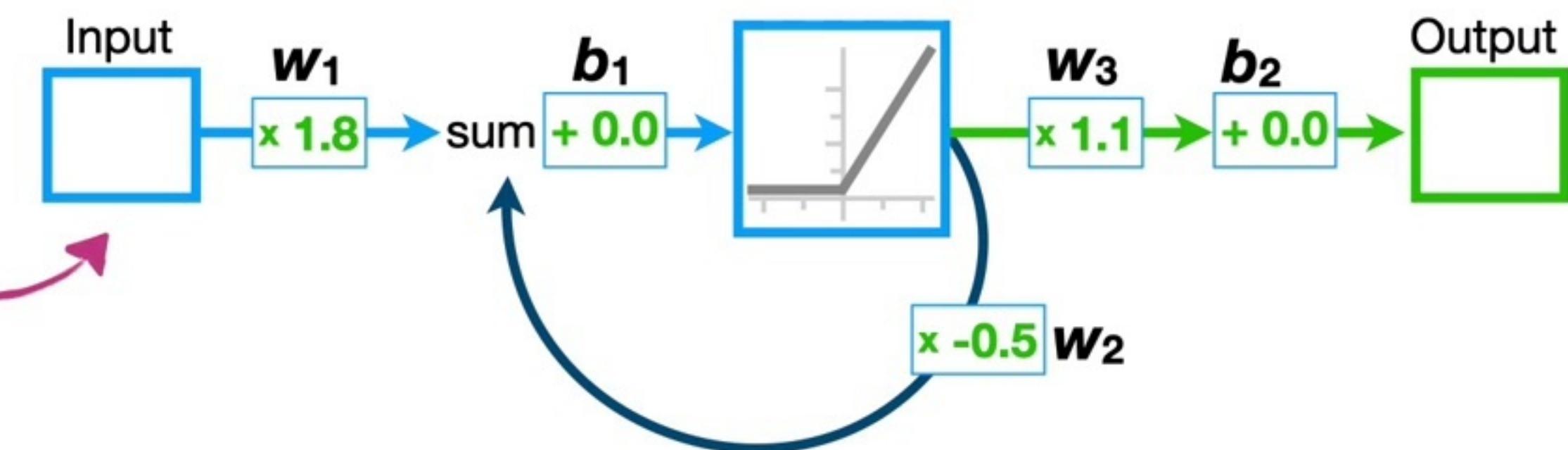


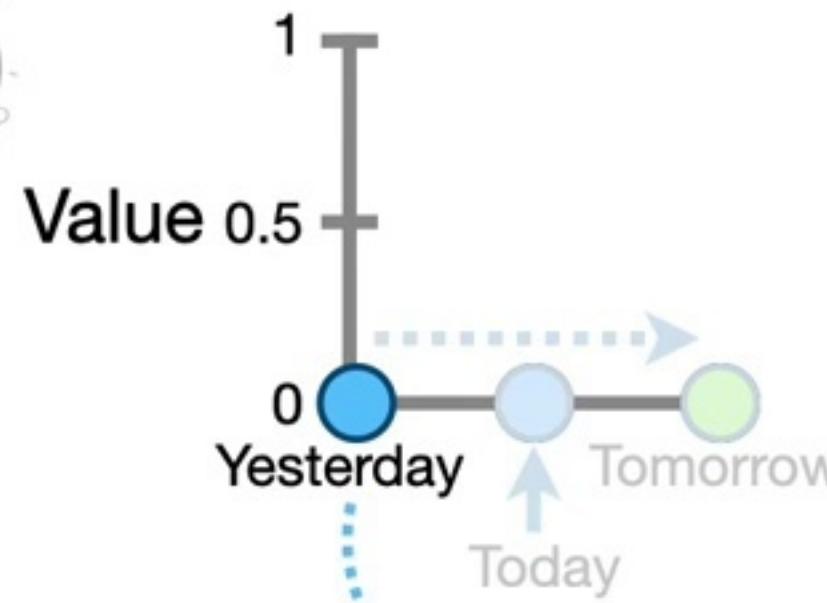
...we can enter
yesterday and **today's**
values into the input
sequentially.



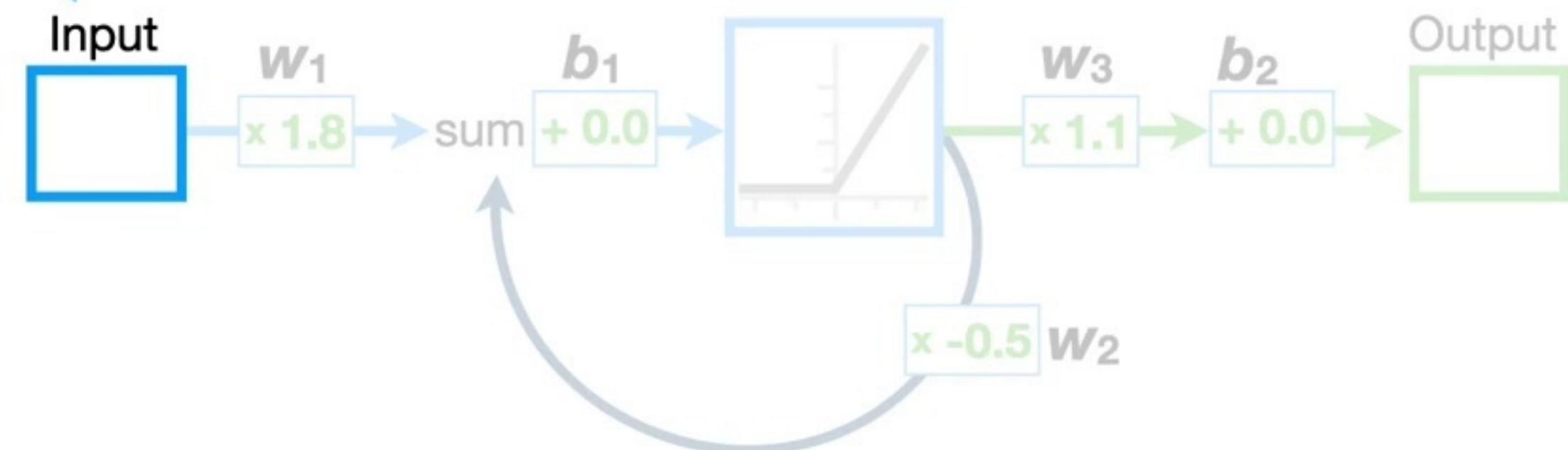


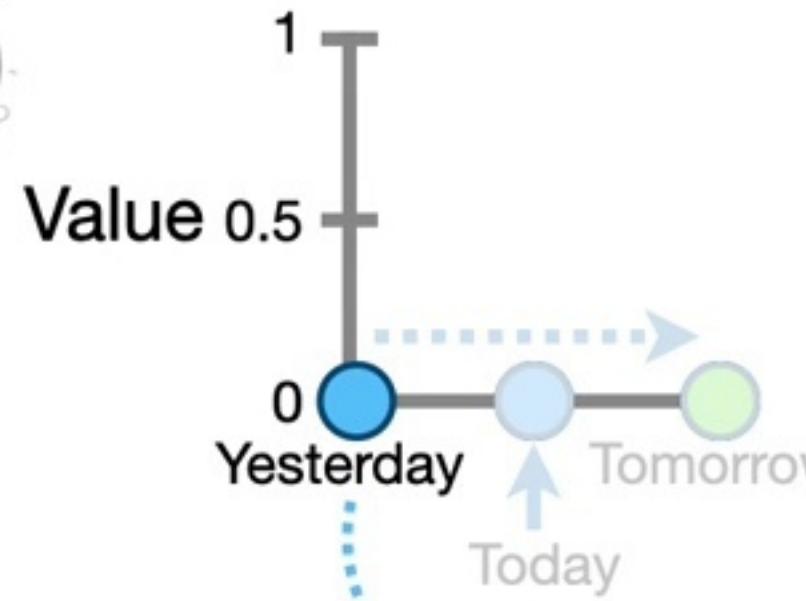
...we can enter
yesterday and **today's**
values into the input
sequentially.



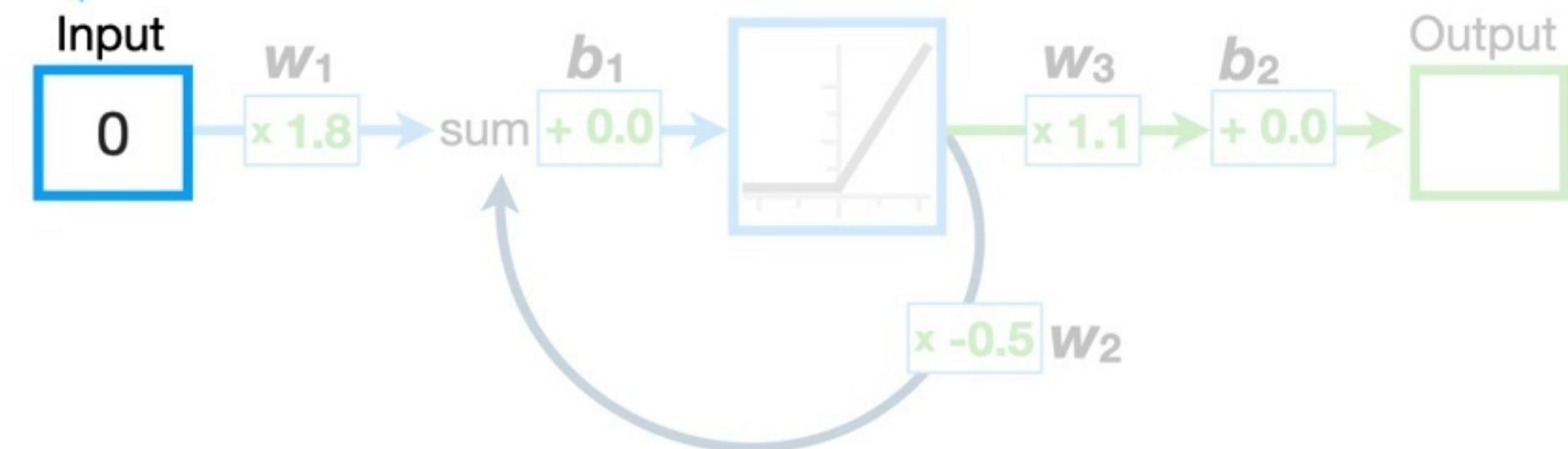


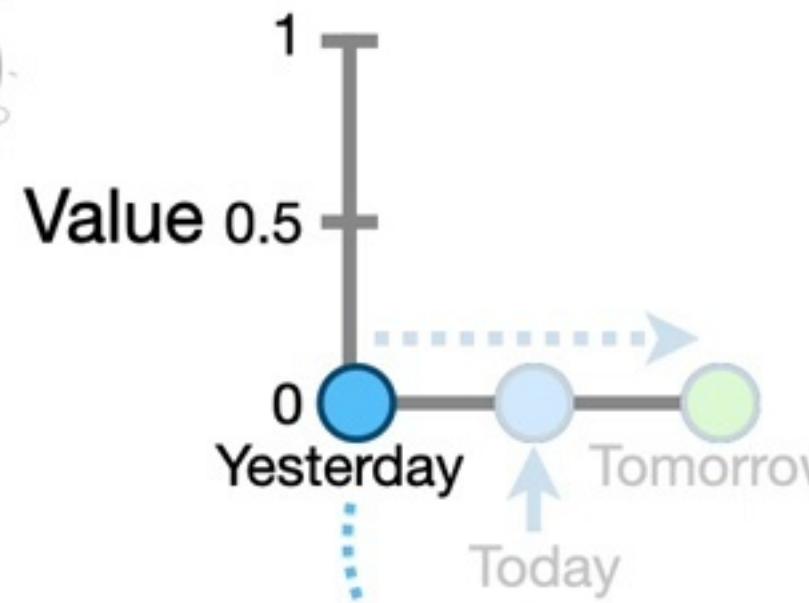
We'll start by plugging **yesterday's** value into the input.



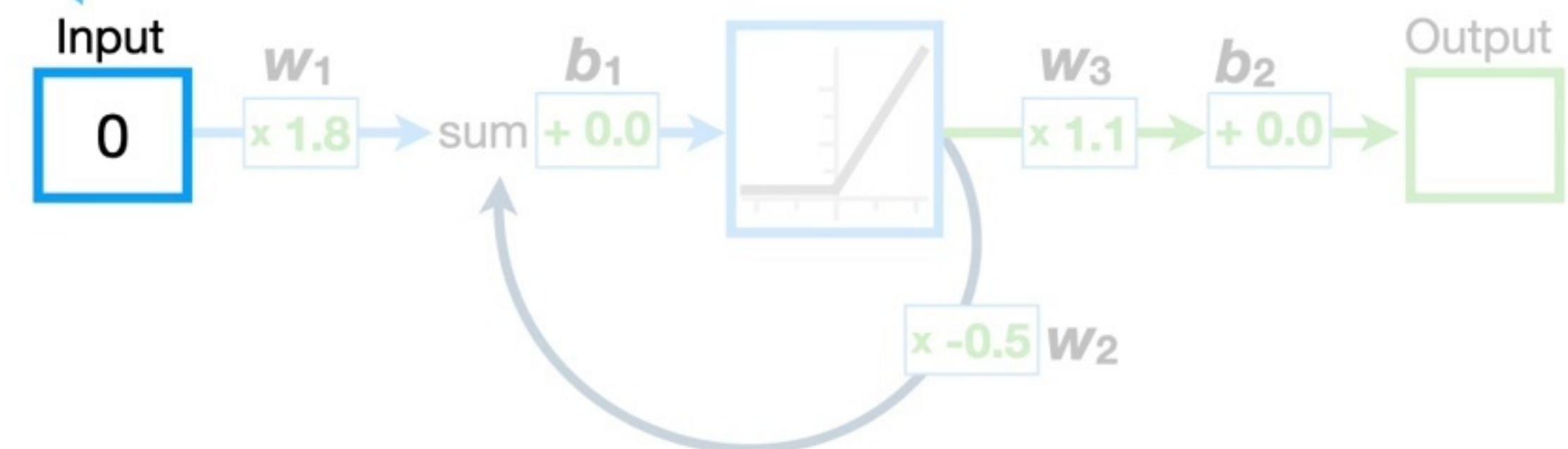


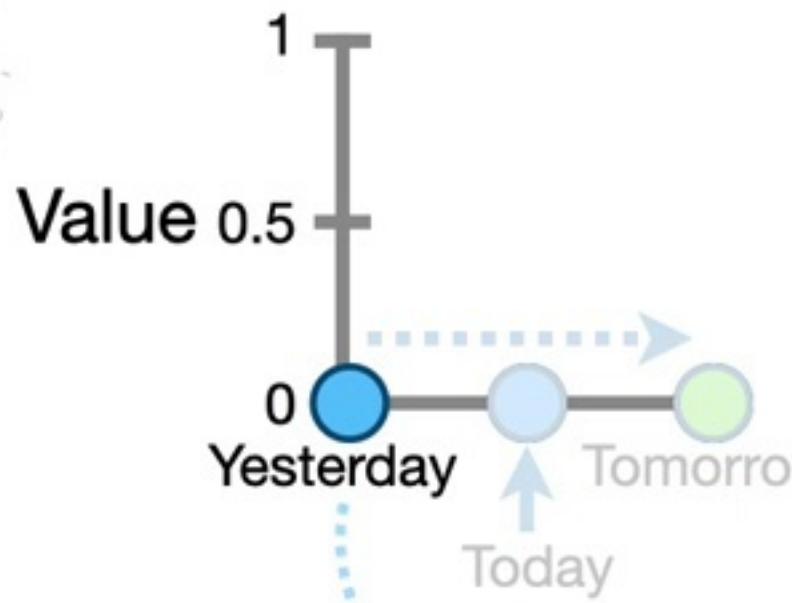
We'll start by plugging **yesterday's** value into the input.



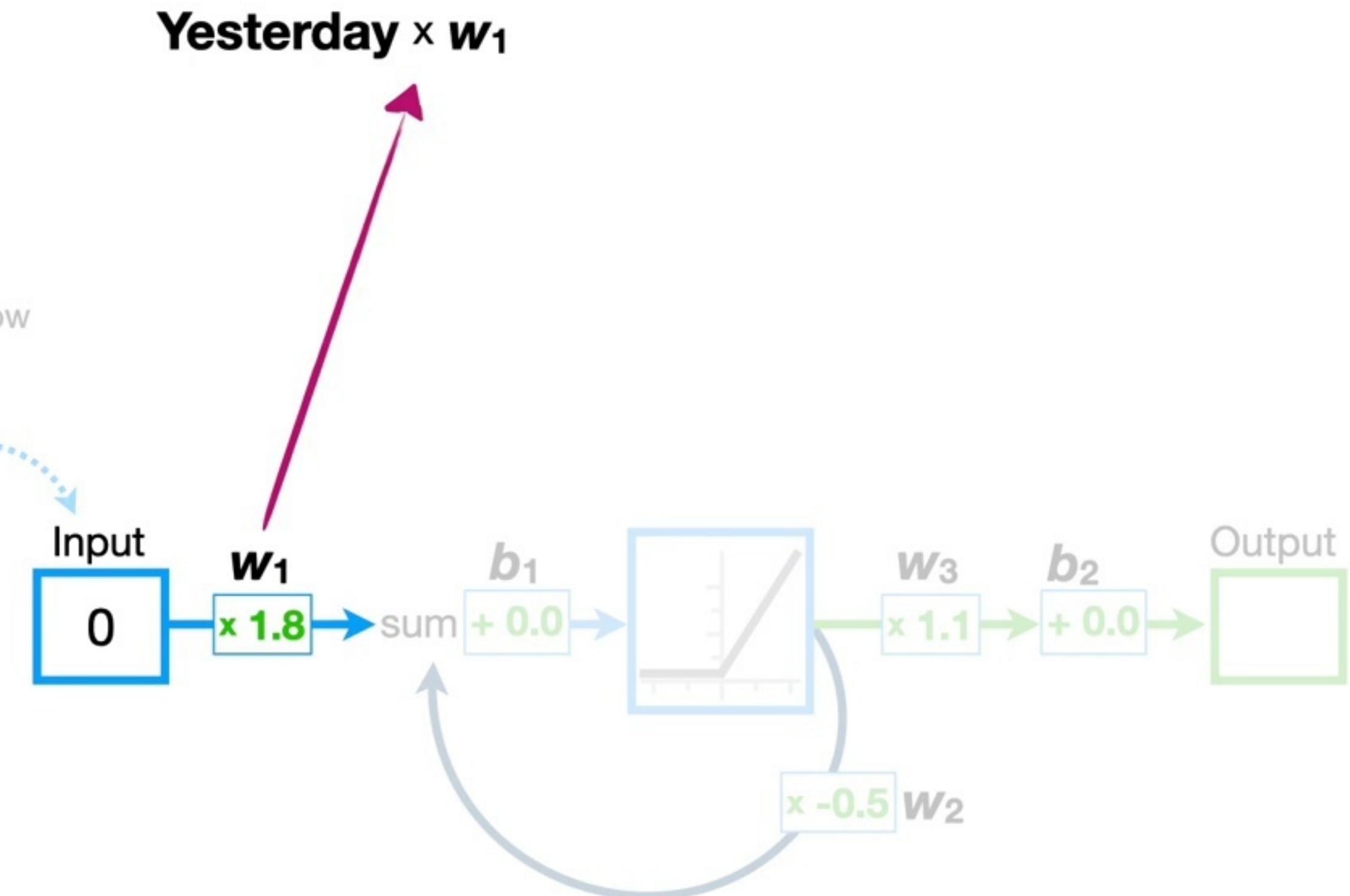


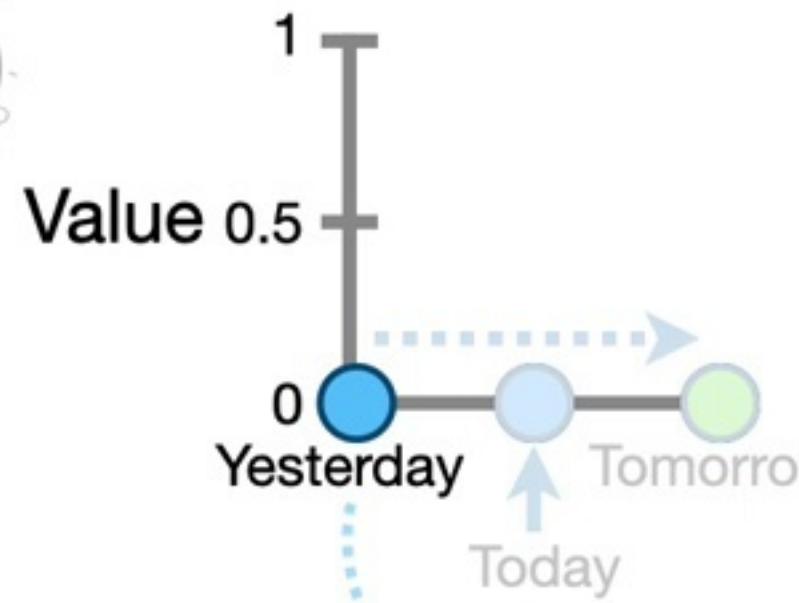
Now we can do the
math just like we
would for any other
neural network.



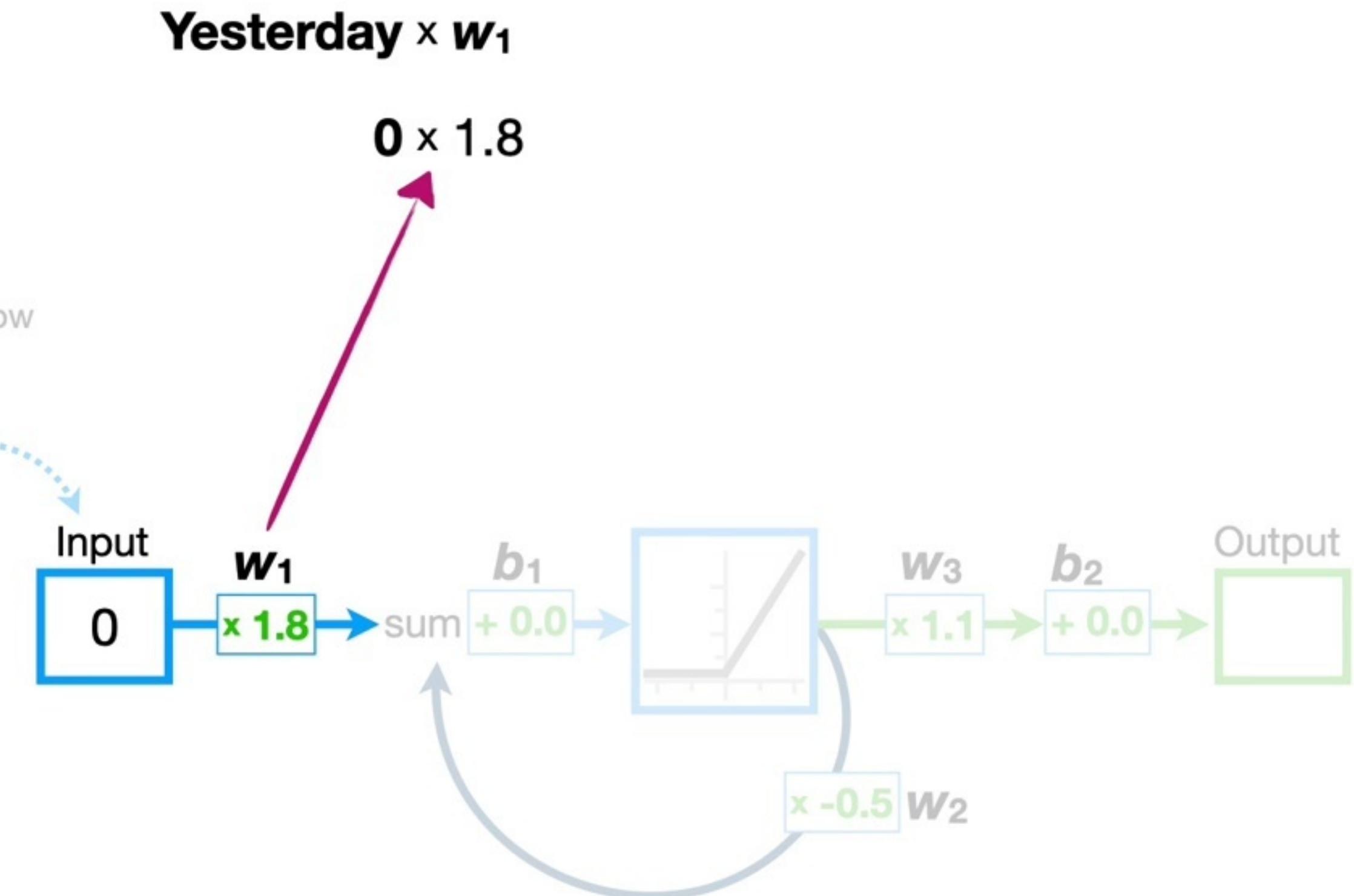


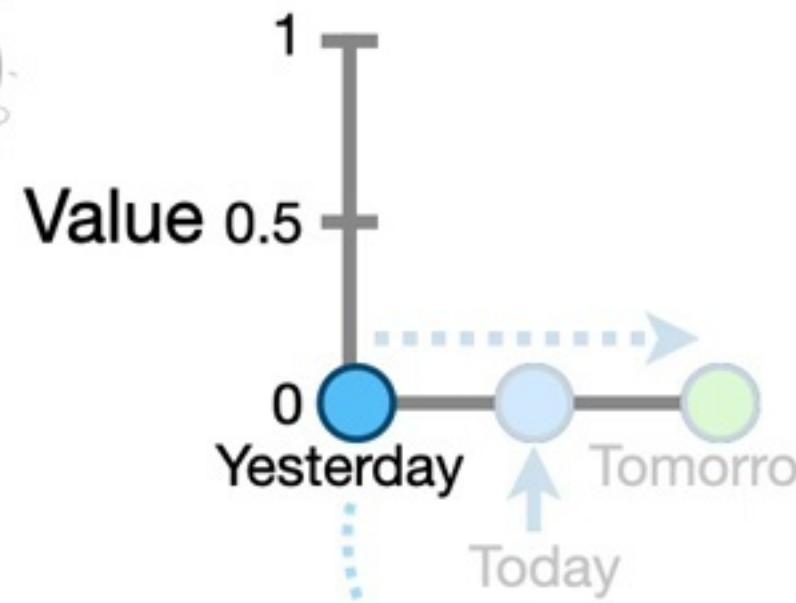
Now we can do the math just like we would for any other neural network.



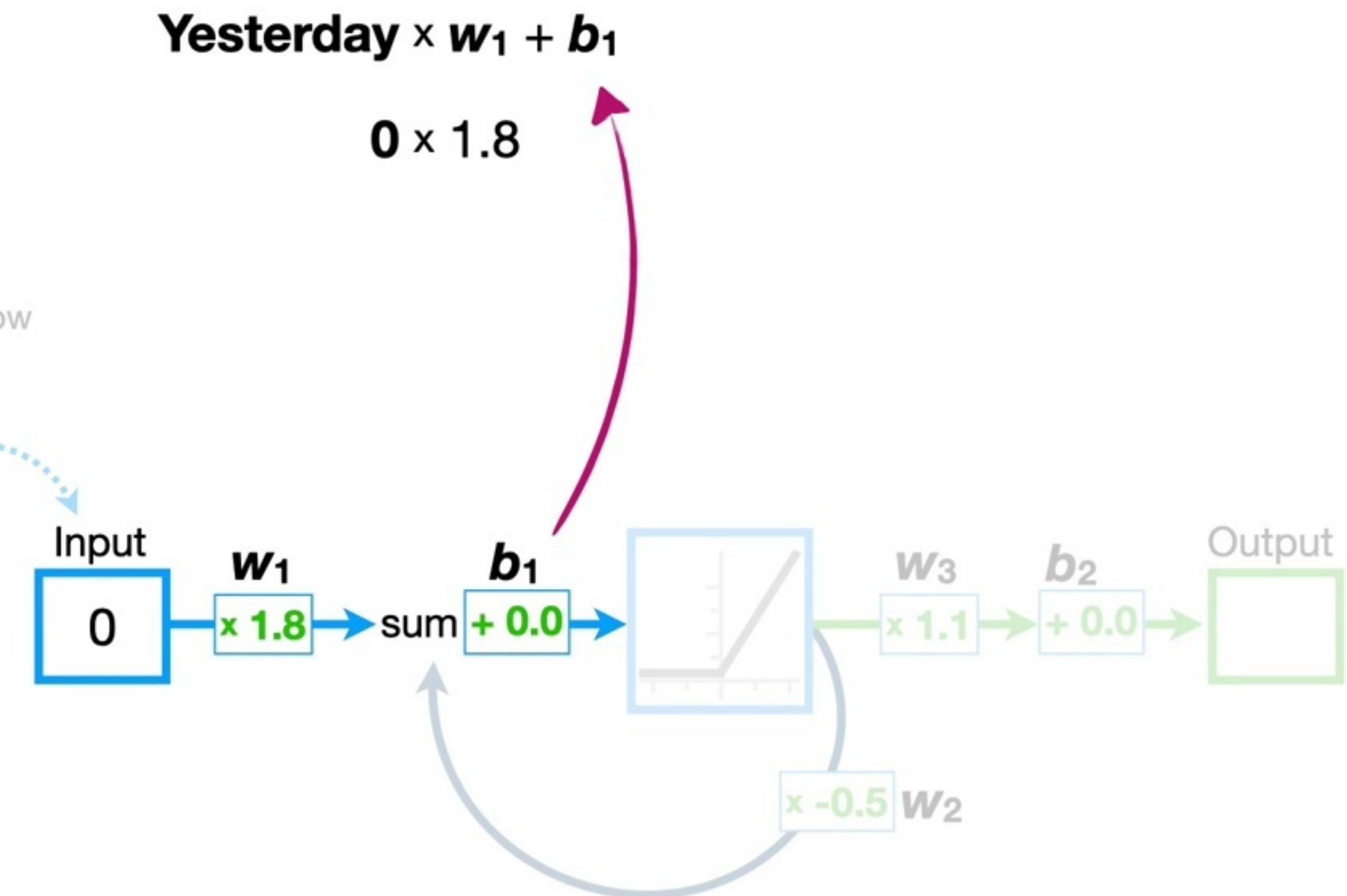


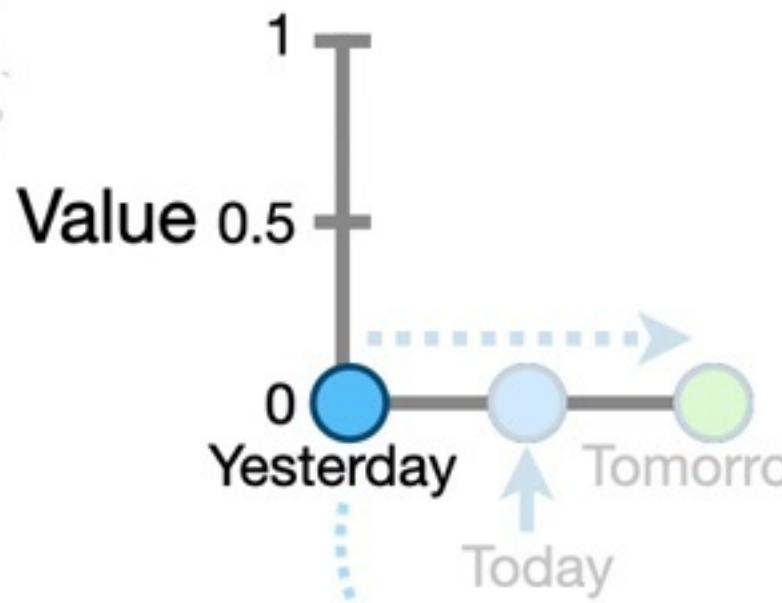
Now we can do the
math just like we
would for any other
neural network.





Now we can do the math just like we would for any other neural network.

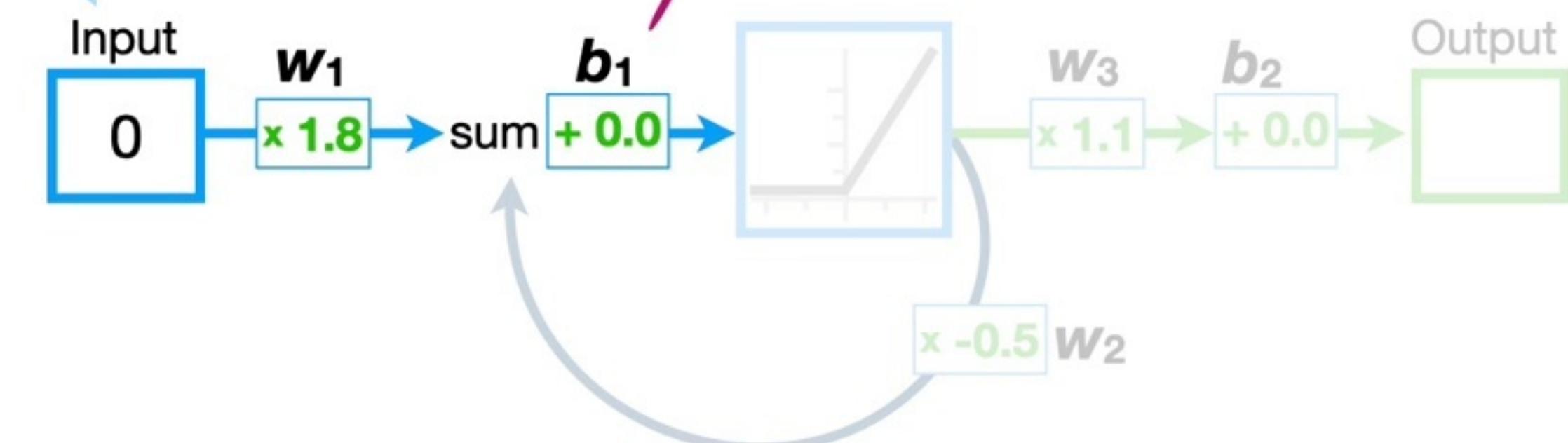


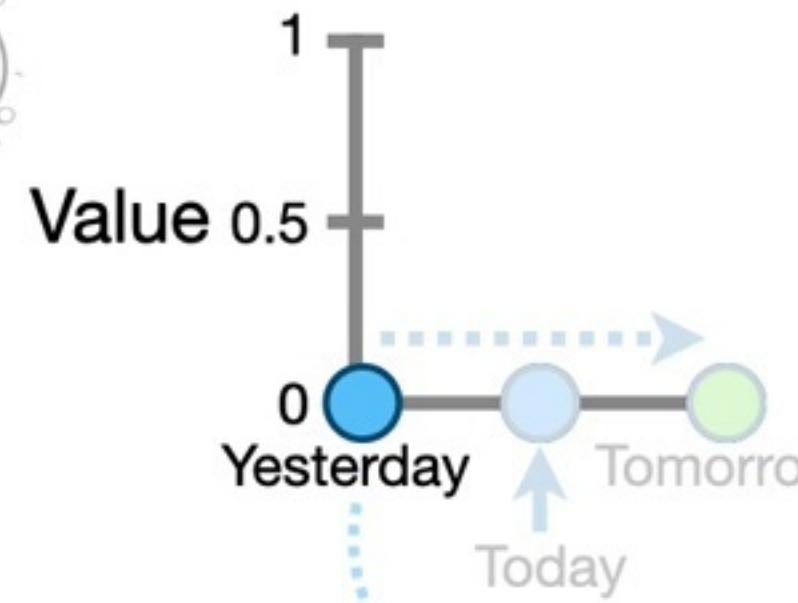


$$\text{Yesterday} \times w_1 + b_1$$

$$0 \times 1.8 + 0.0$$

Now we can do the math just like we would for any other neural network.

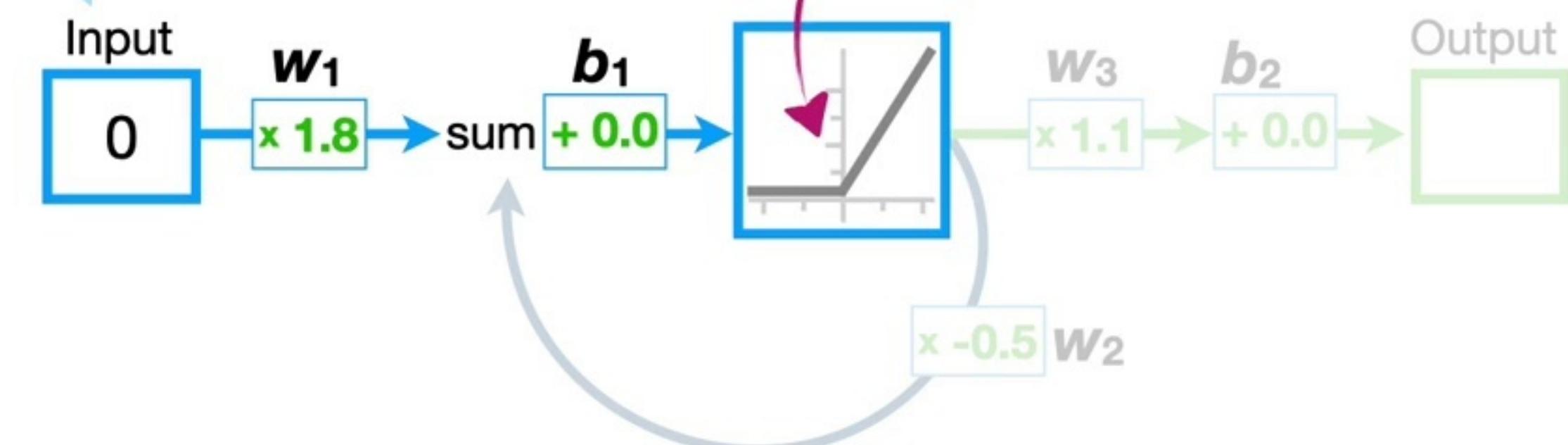


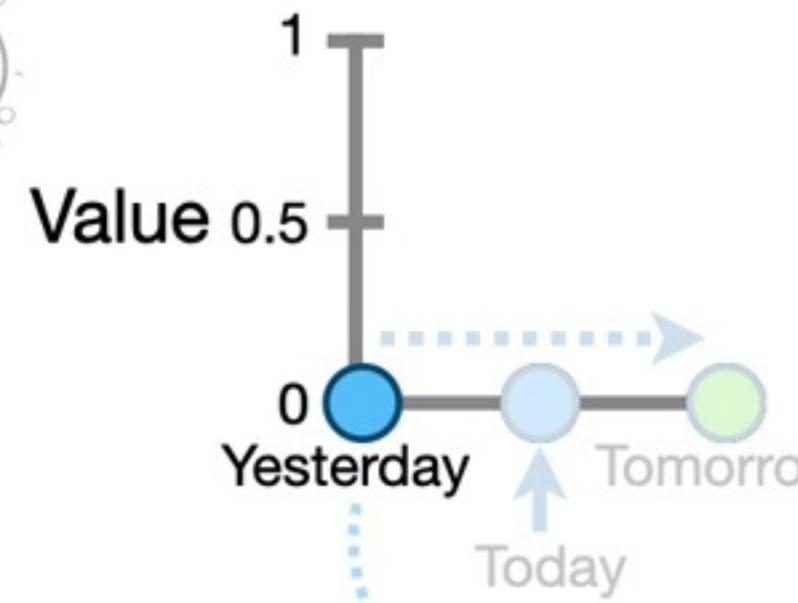


Yesterday $\times w_1 + b_1$ = x-axis coordinate

$$0 \times 1.8 + 0.0$$

Now we can do the math just like we would for any other neural network.

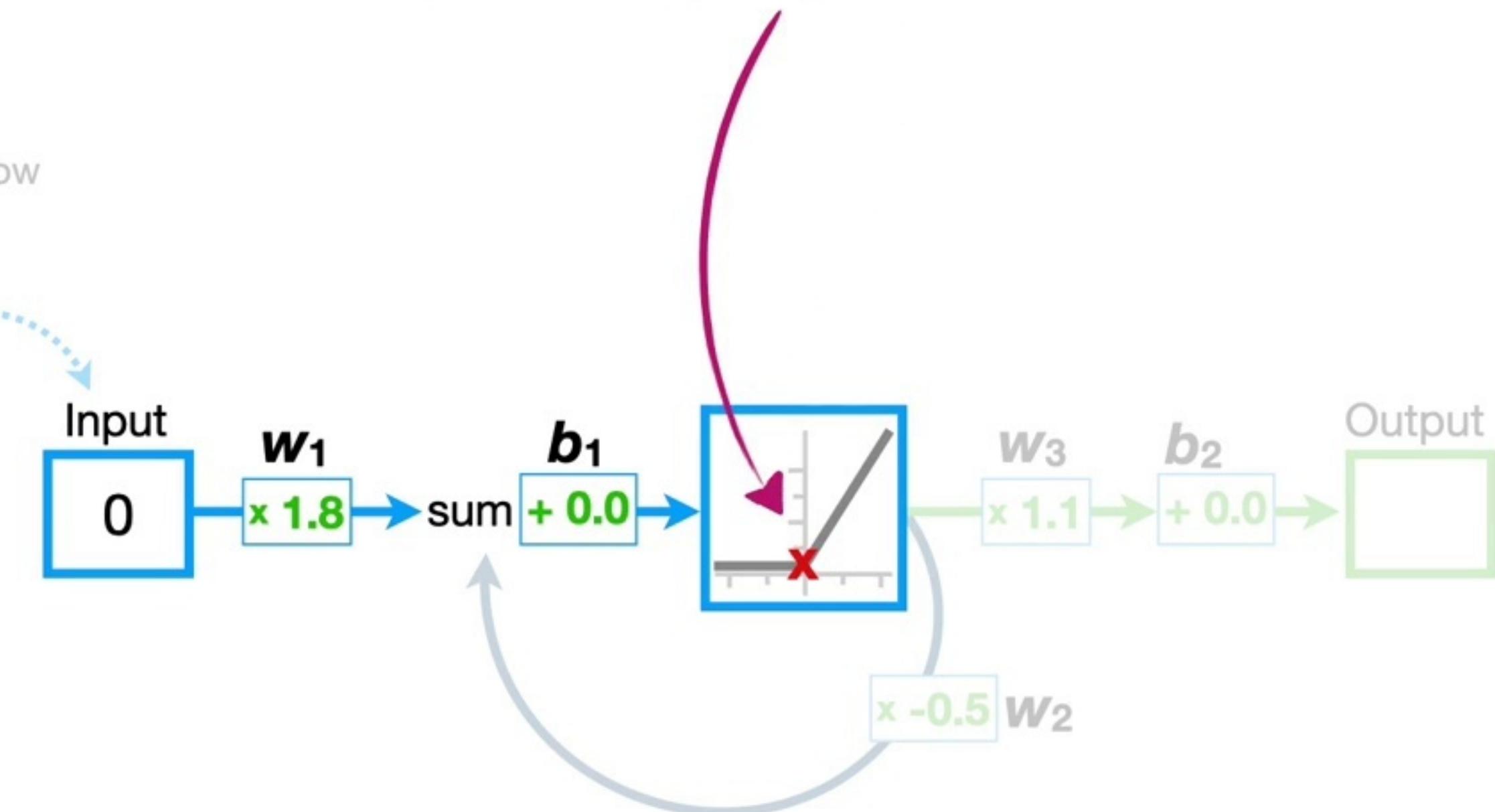


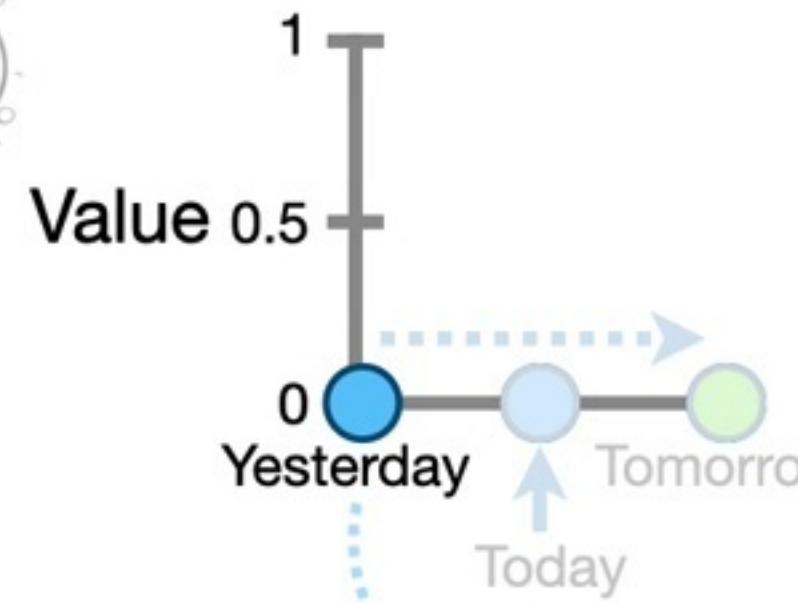


Yesterday $\times w_1 + b_1$ = x-axis coordinate

$$0 \times 1.8 + 0.0 = 0$$

Now we can do the math just like we would for any other neural network.



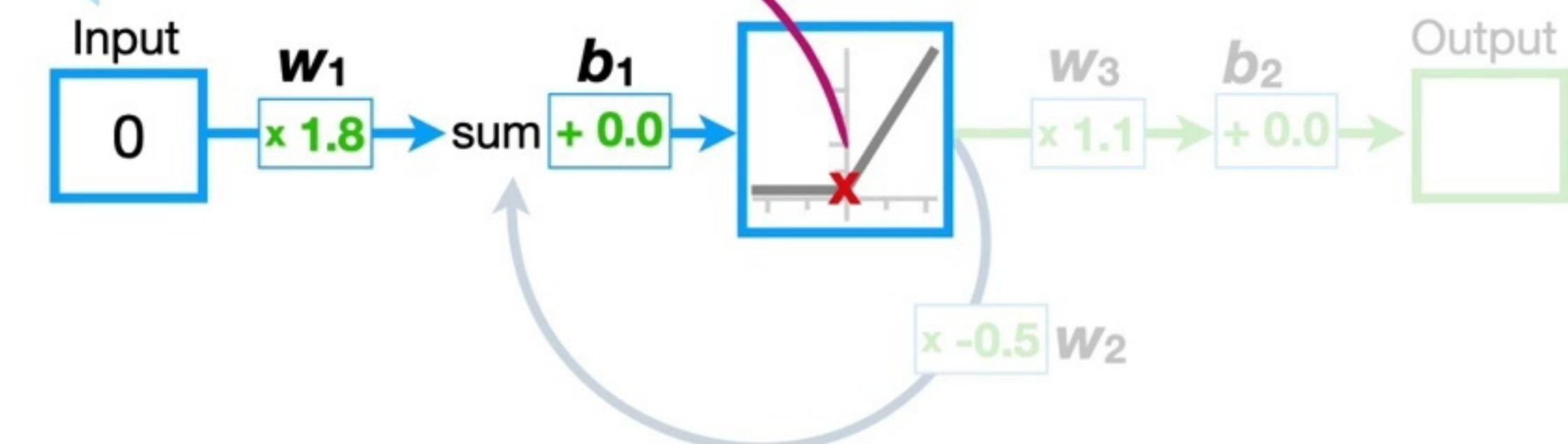


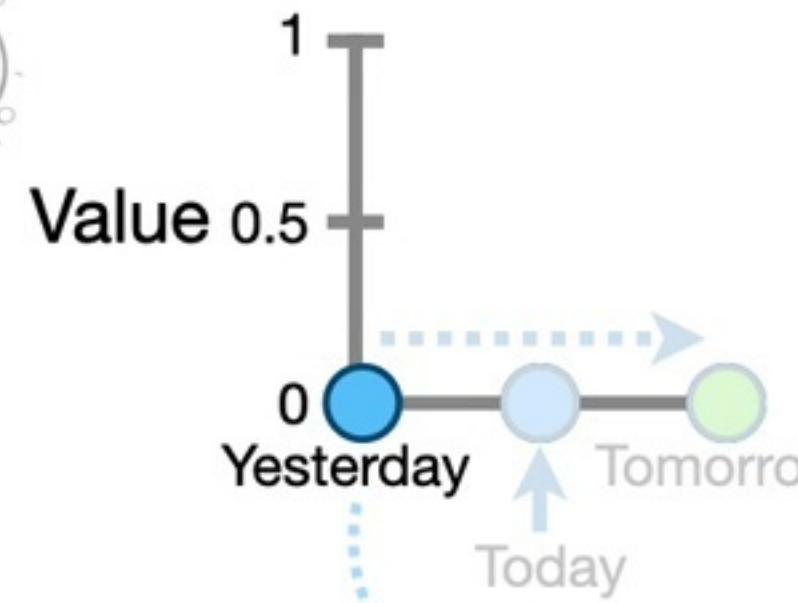
$\text{Yesterday} \times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$$f(x) = \max(0, x) = \text{y-axis coordinate}$$

Now we can do the math just like we would for any other neural network.



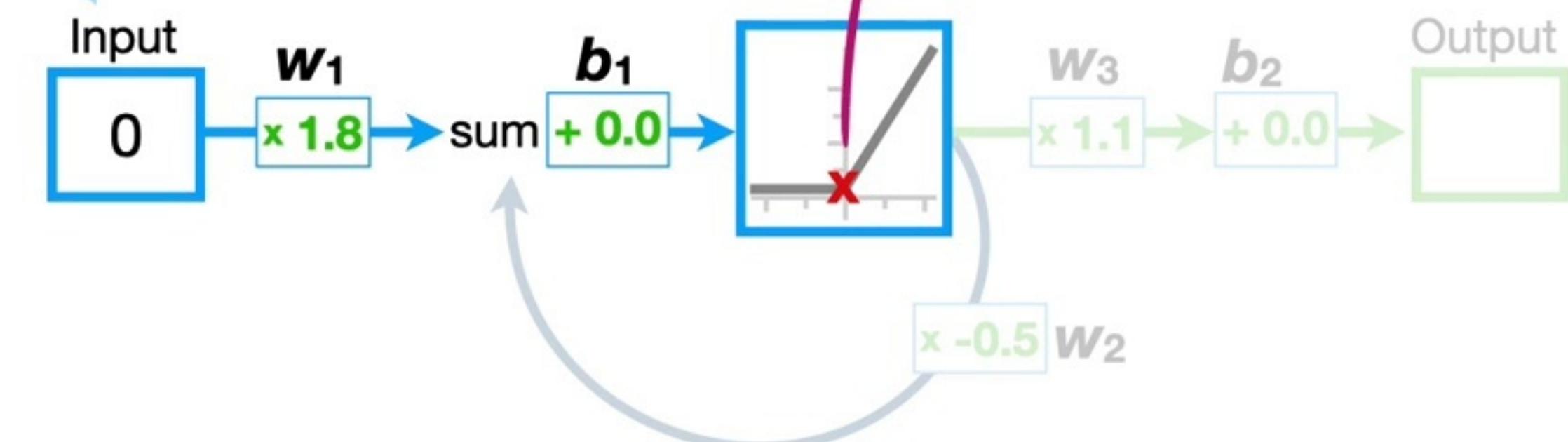


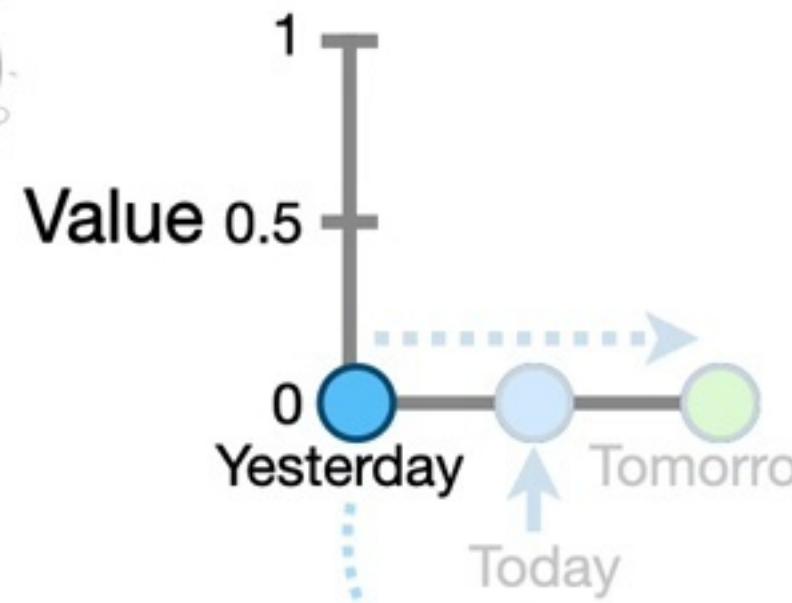
Yesterday $\times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(x) = \max(0, x) = \text{y-axis coordinate}$

Now we can do the
math just like we
would for any other
neural network.



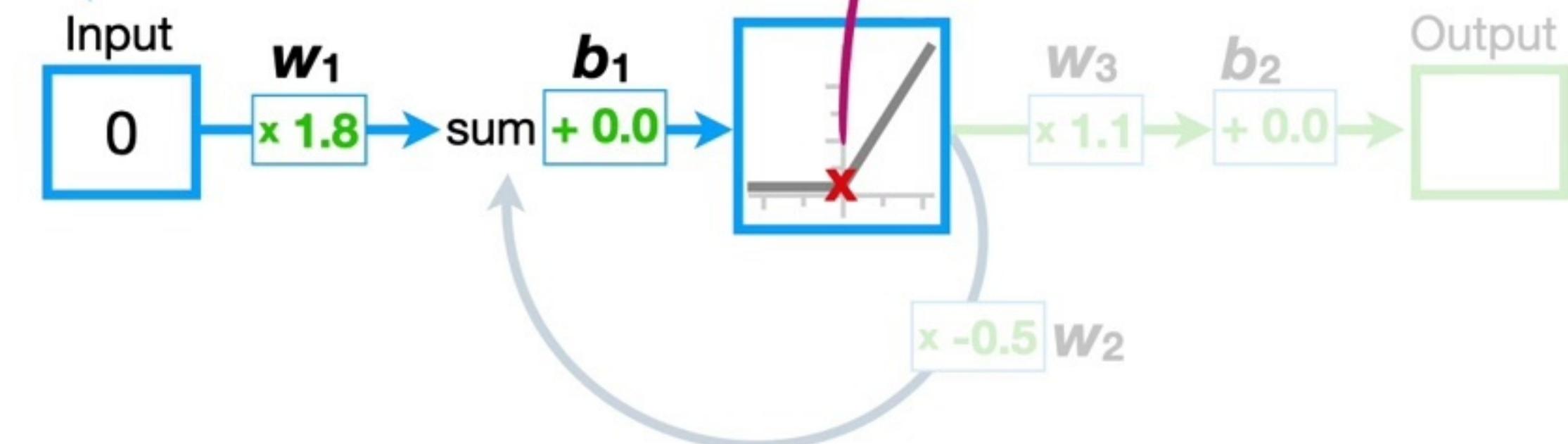


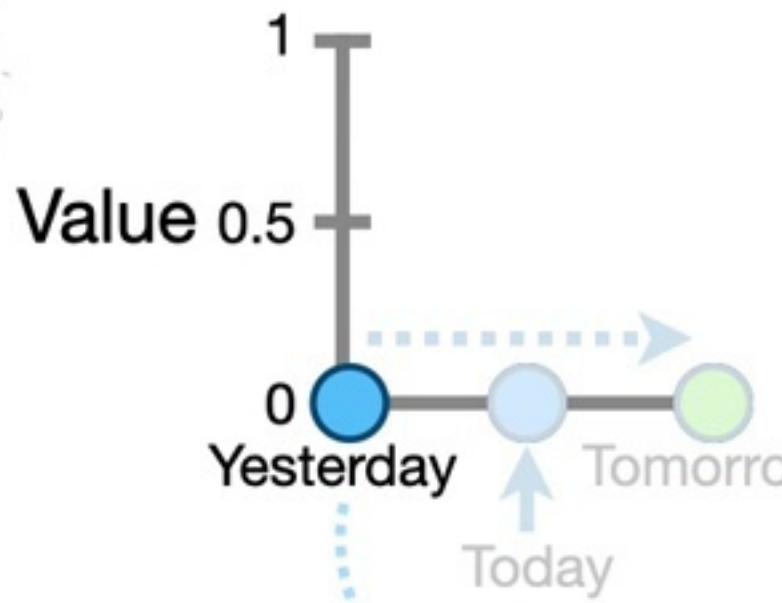
Yesterday $\times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(0) = \max(0, 0) = \text{y-axis coordinate}$

Now we can do the
math just like we
would for any other
neural network.



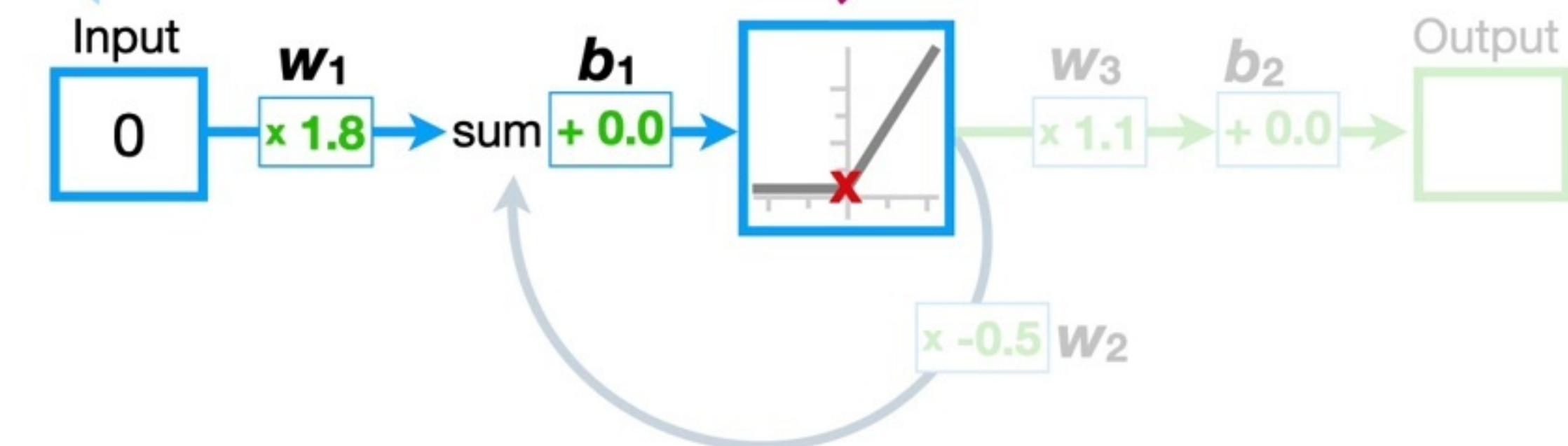


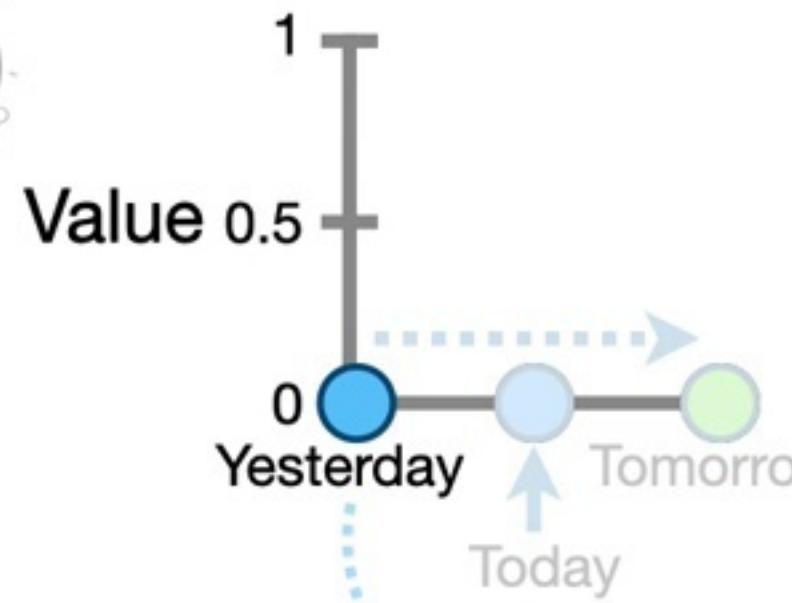
Yesterday $\times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(0) = \max(0, 0) = \text{y-axis coordinate}$

Now we can do the
math just like we
would for any other
neural network.



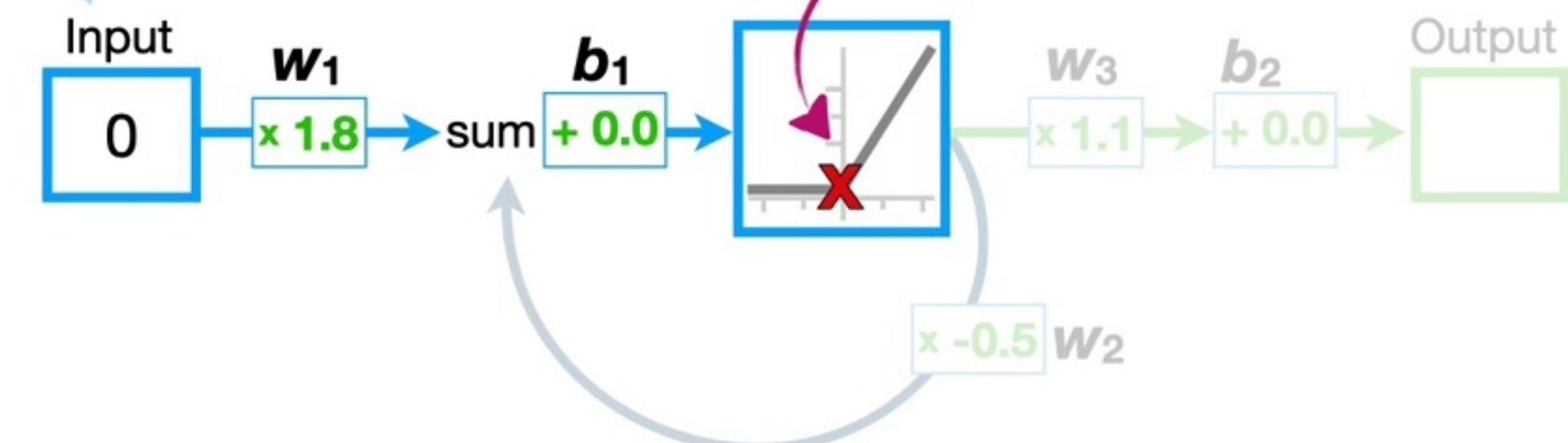


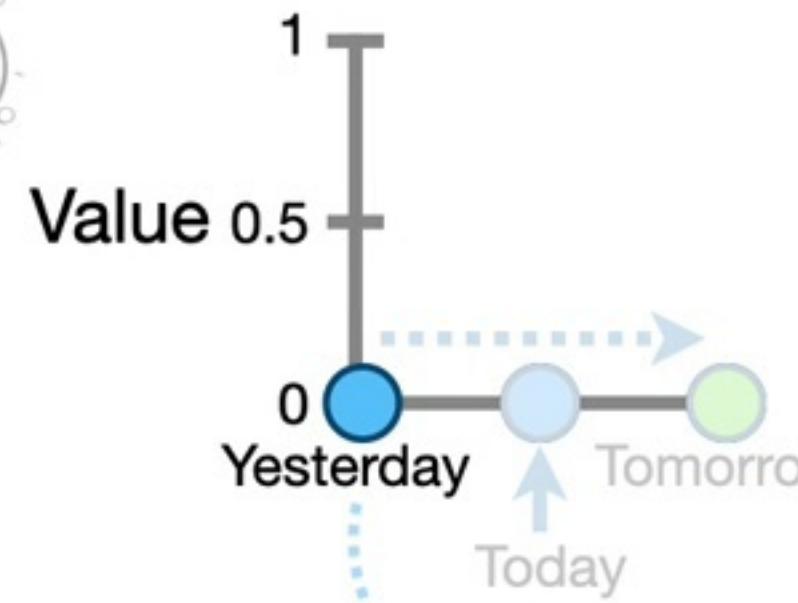
Yesterday $\times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$$f(0) = \max(0, 0) = 0$$

Now we can do the
math just like we
would for any other
neural network.

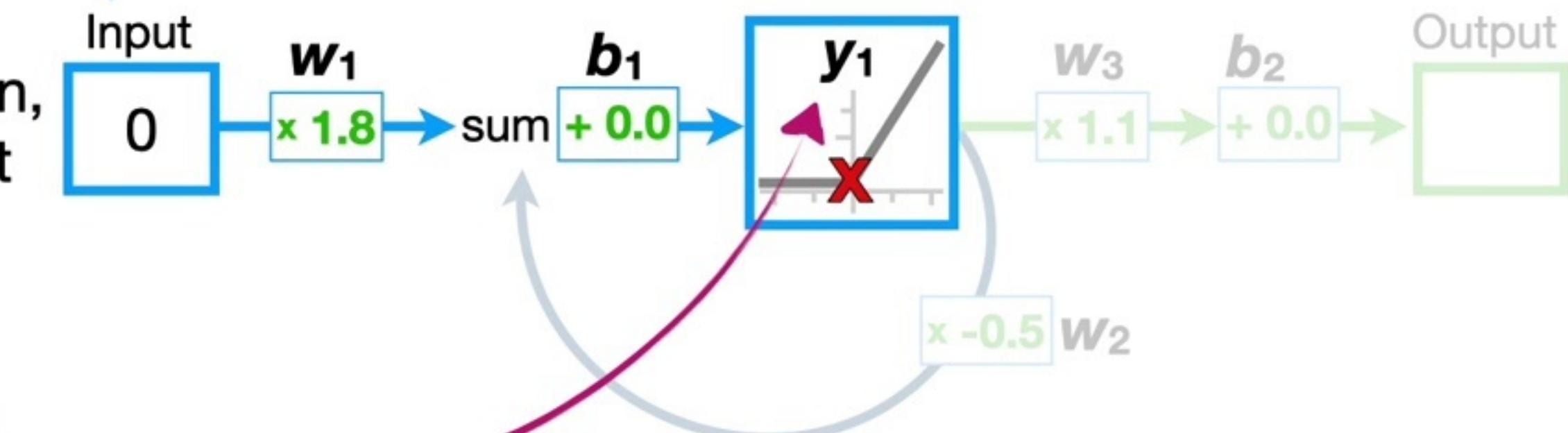


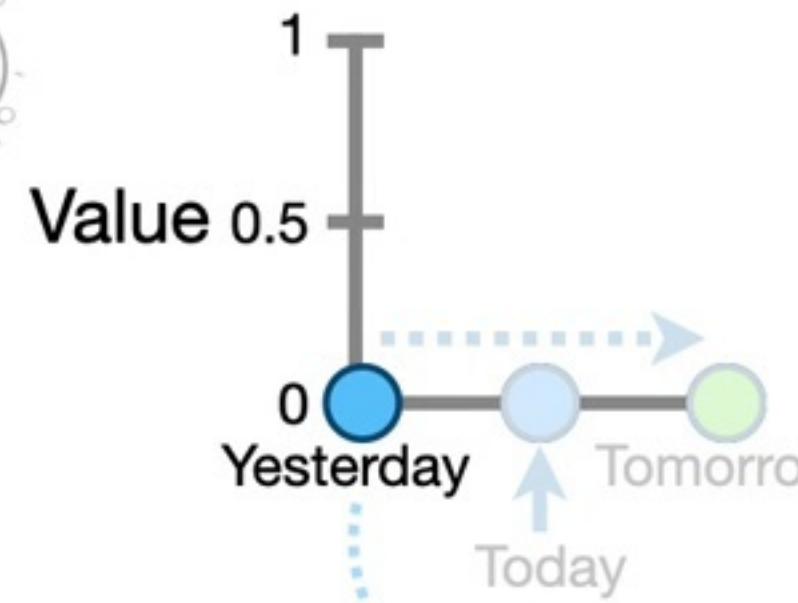


$$f(0) = \max(0, 0) = 0 = y_1$$

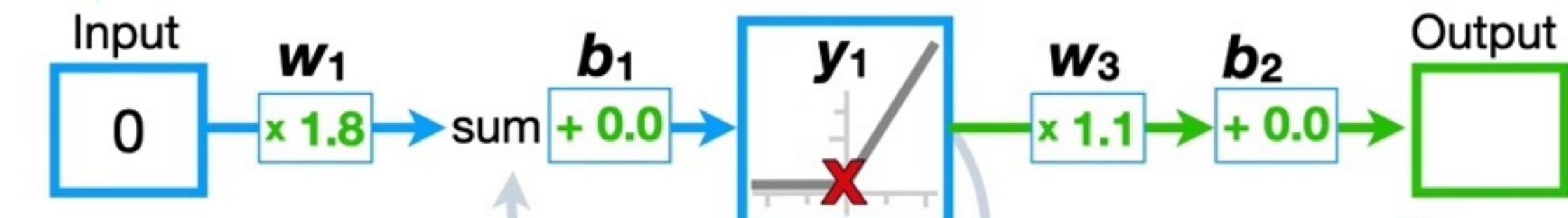
At this point, the output from the activation function, the y-axis coordinate that we'll call y_1 ...

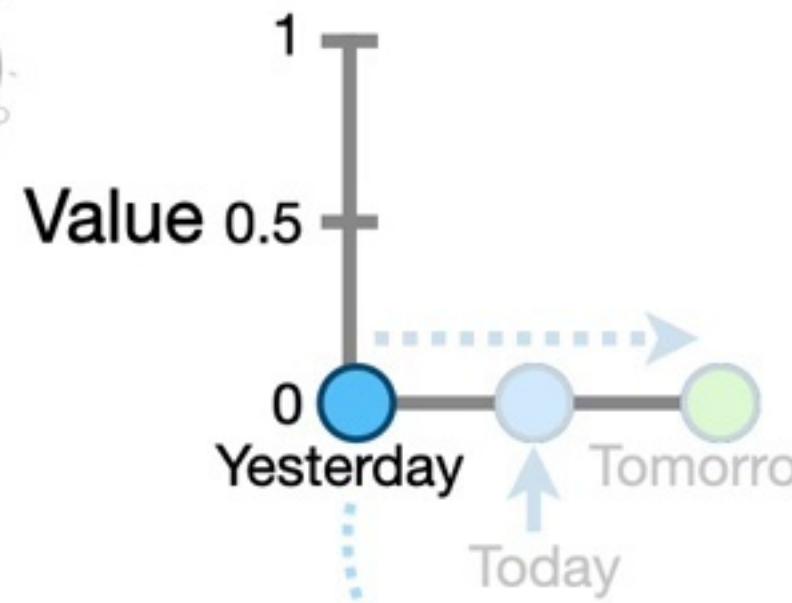
...can go two places.



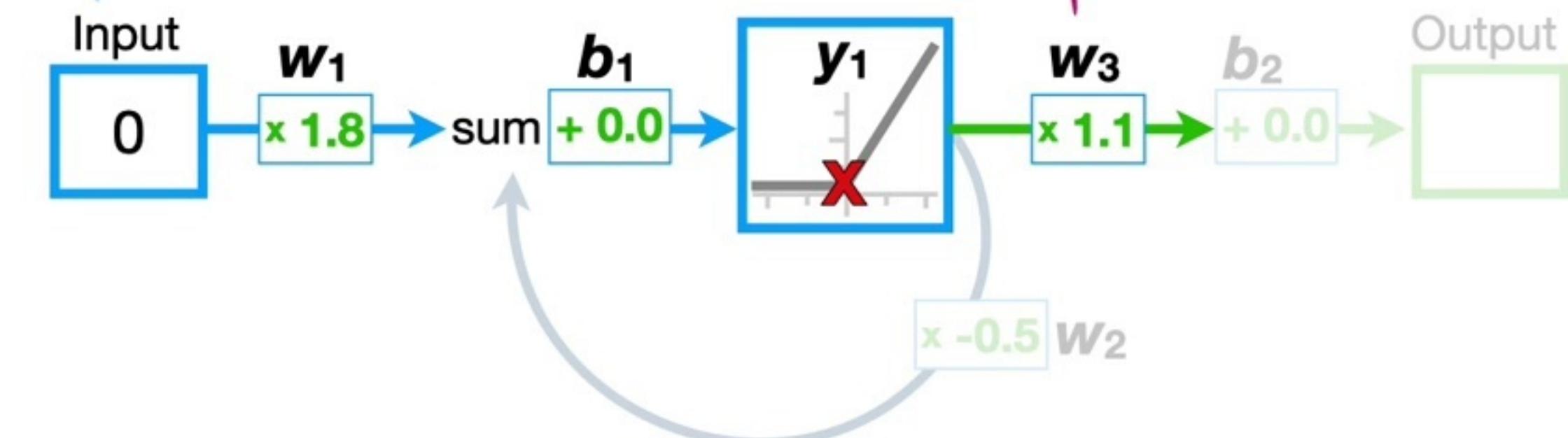


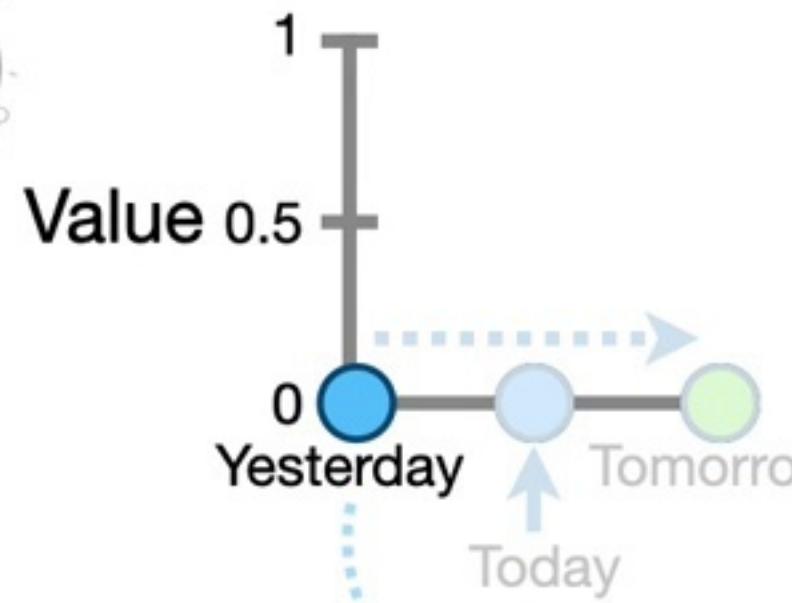
First, y_1 can go towards the output.



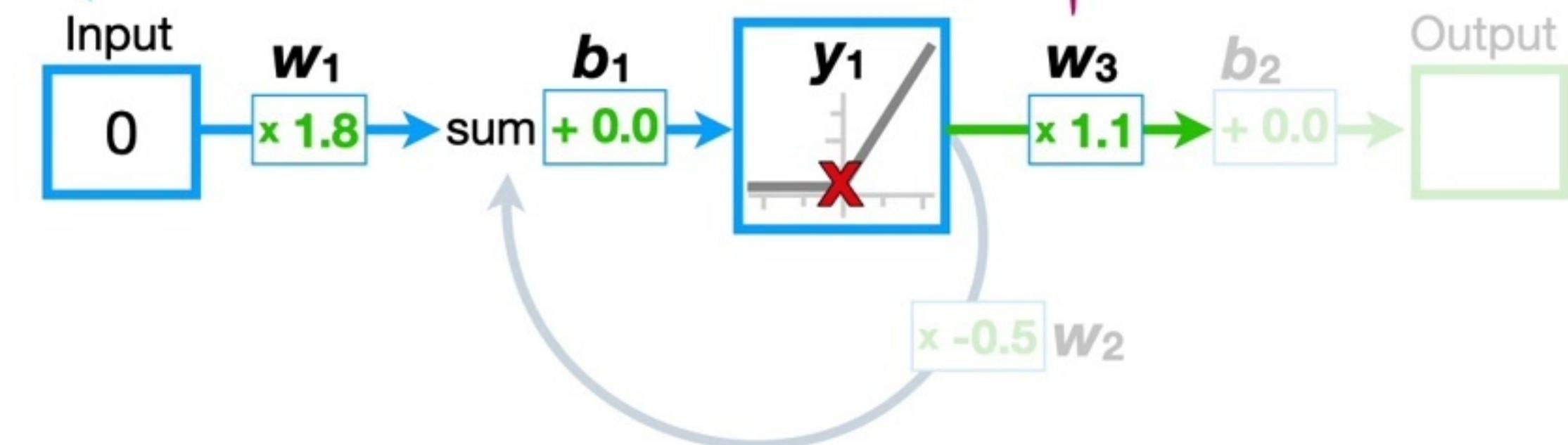


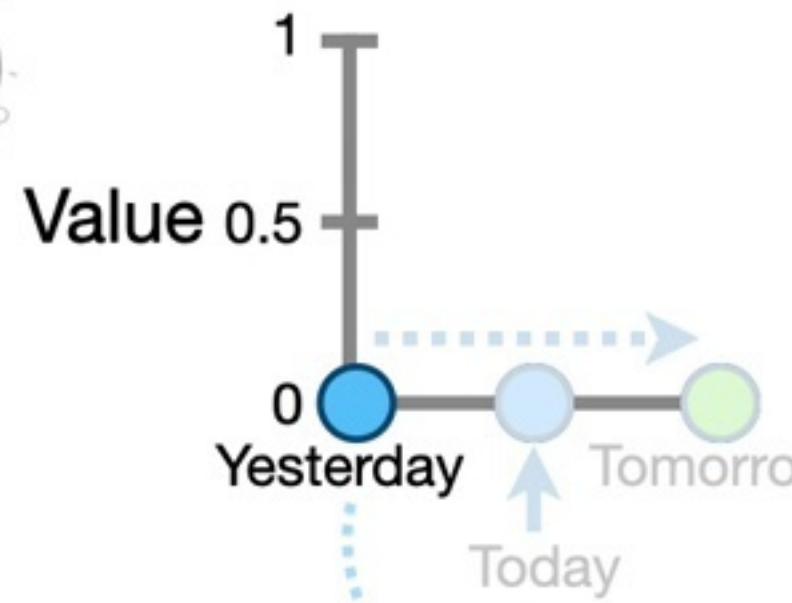
And if we go that way
and do the math...



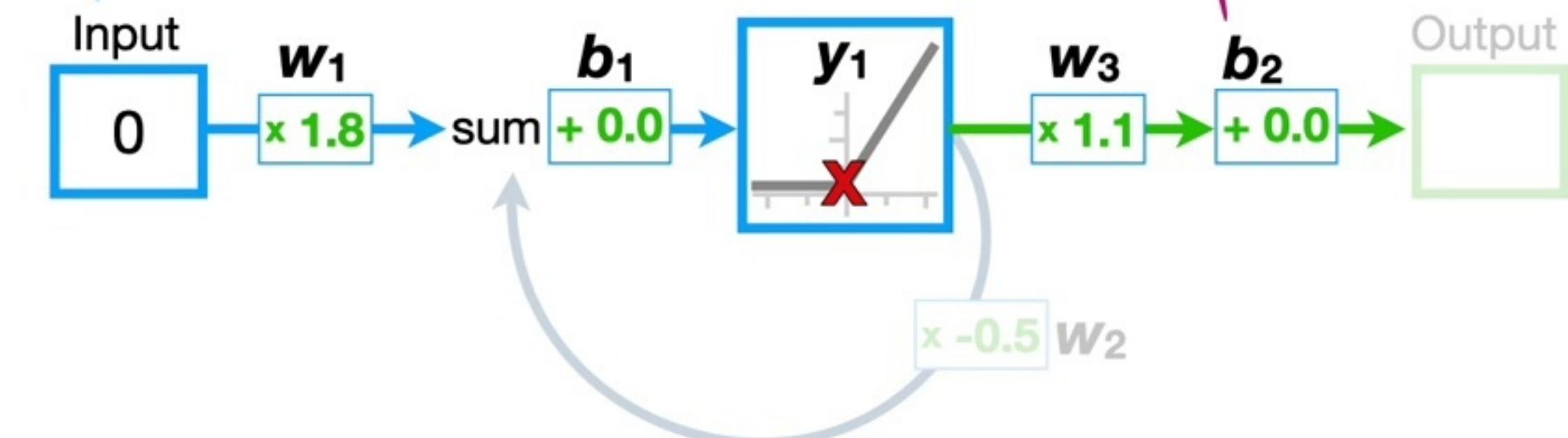


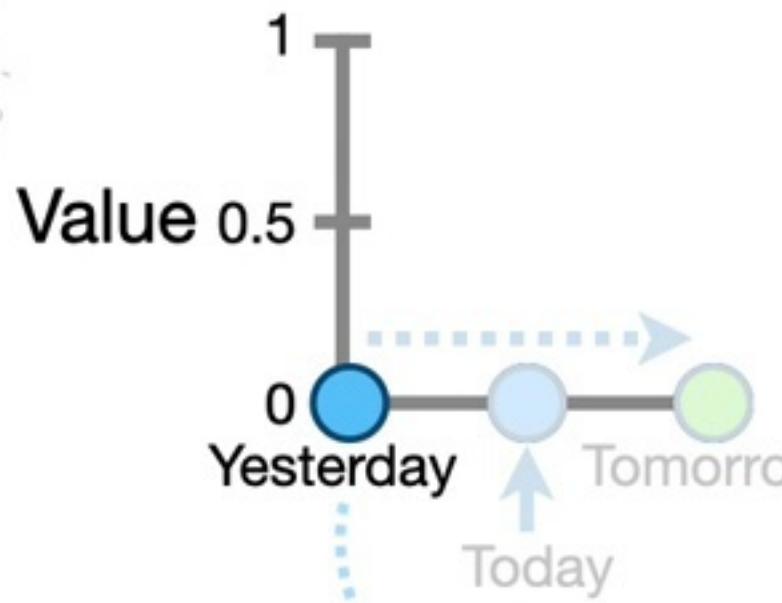
And if we go that way
and do the math...





And if we go that way
and do the math...

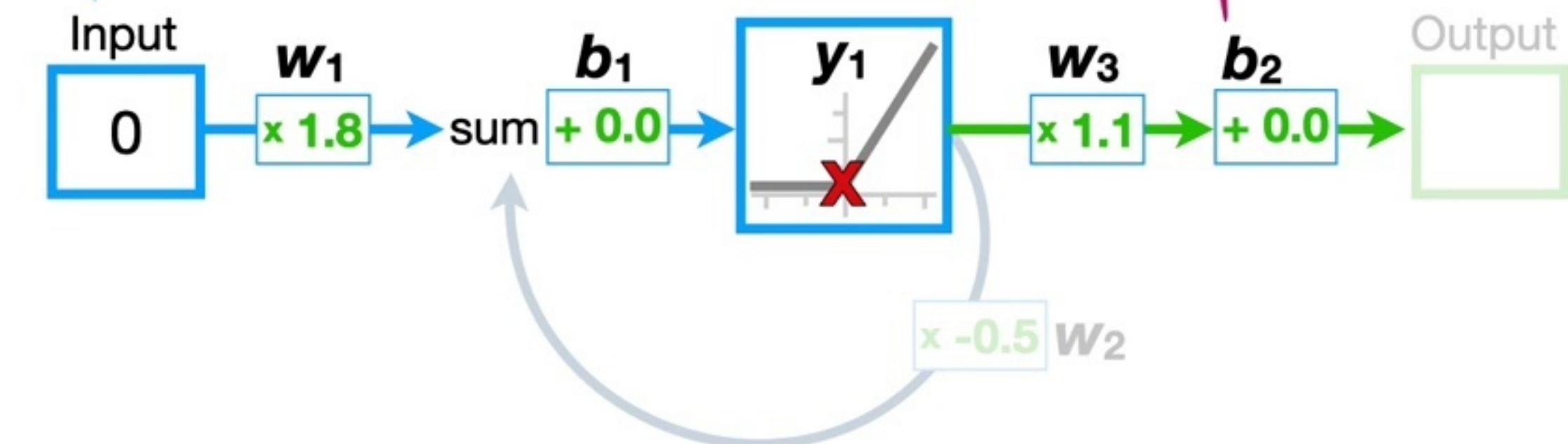


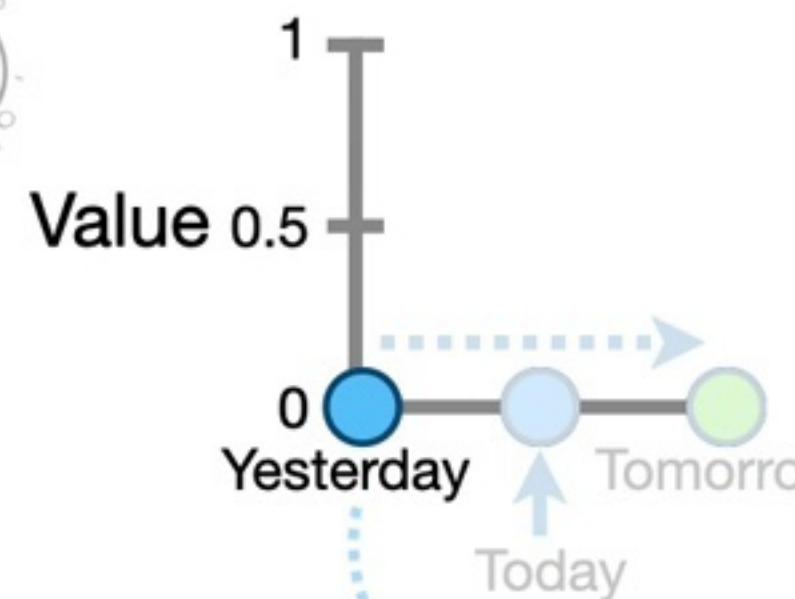


And if we go that way
and do the math...

$$y_1 \times w_3 + b_2$$

$$0 \times 1.1 + 0.0$$

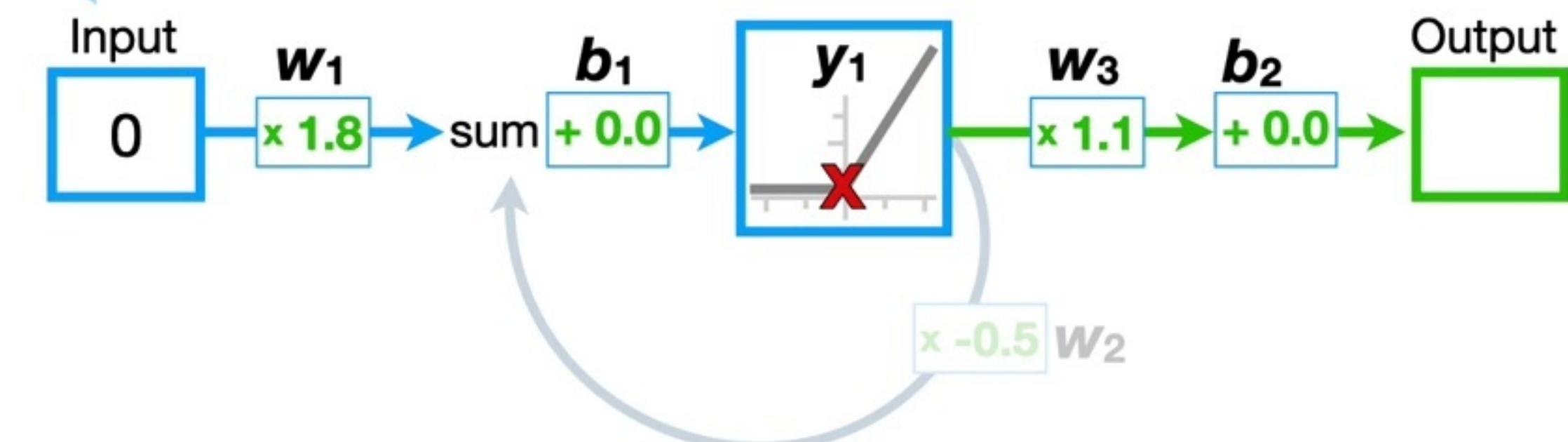


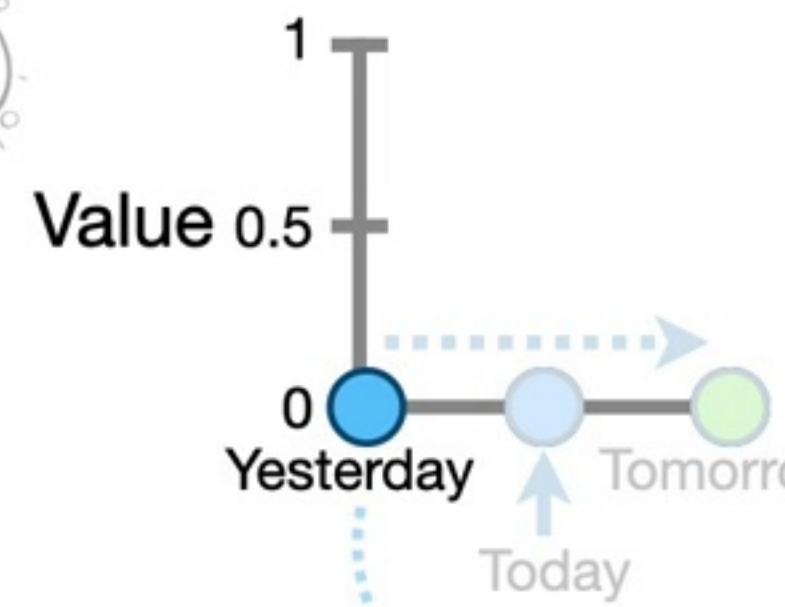


...then the output is the predicted value for **today**.

$y_1 \times w_3 + b_2 = \text{The Predicted Value for Today}$

$$0 \times 1.1 + 0.0$$

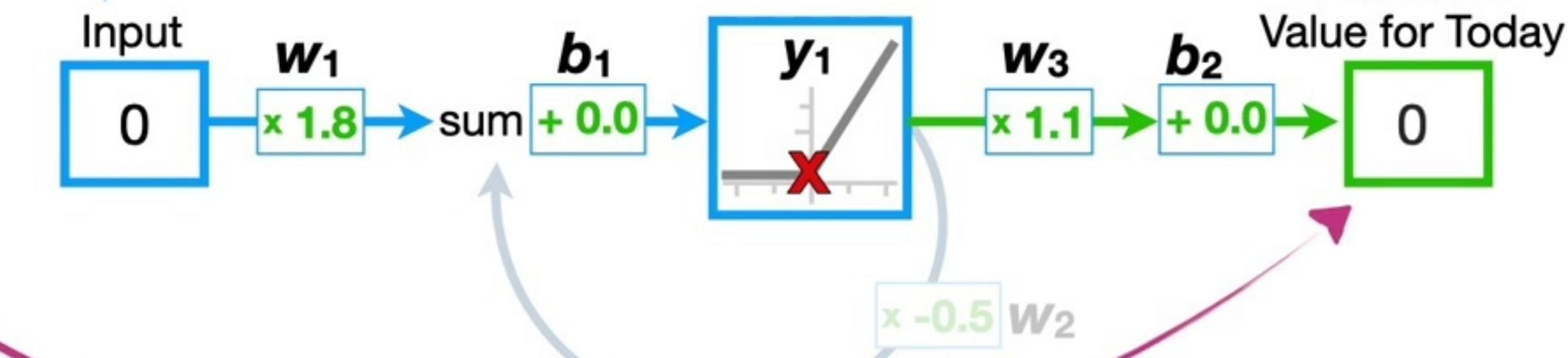


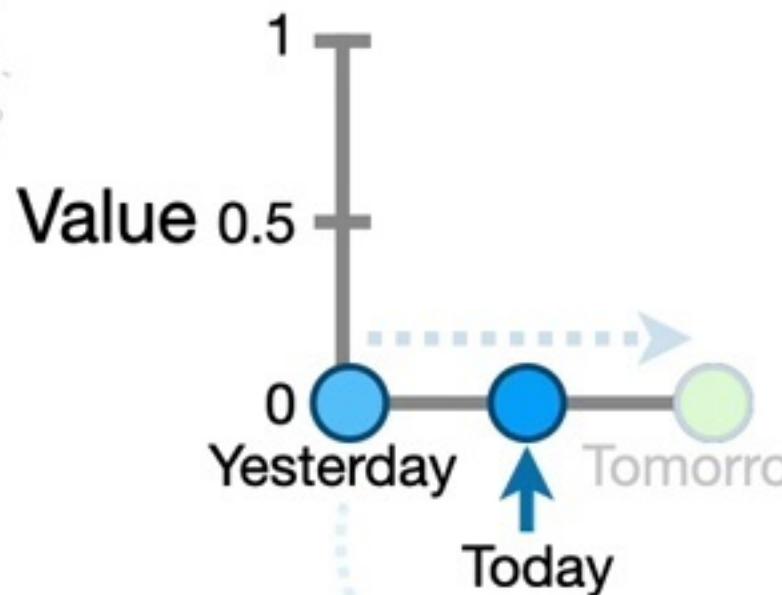


...then the output is the predicted value for today.

$y_1 \times w_3 + b_2 = \text{The Predicted Value for Today}$

$$0 \times 1.1 + 0.0 = 0$$

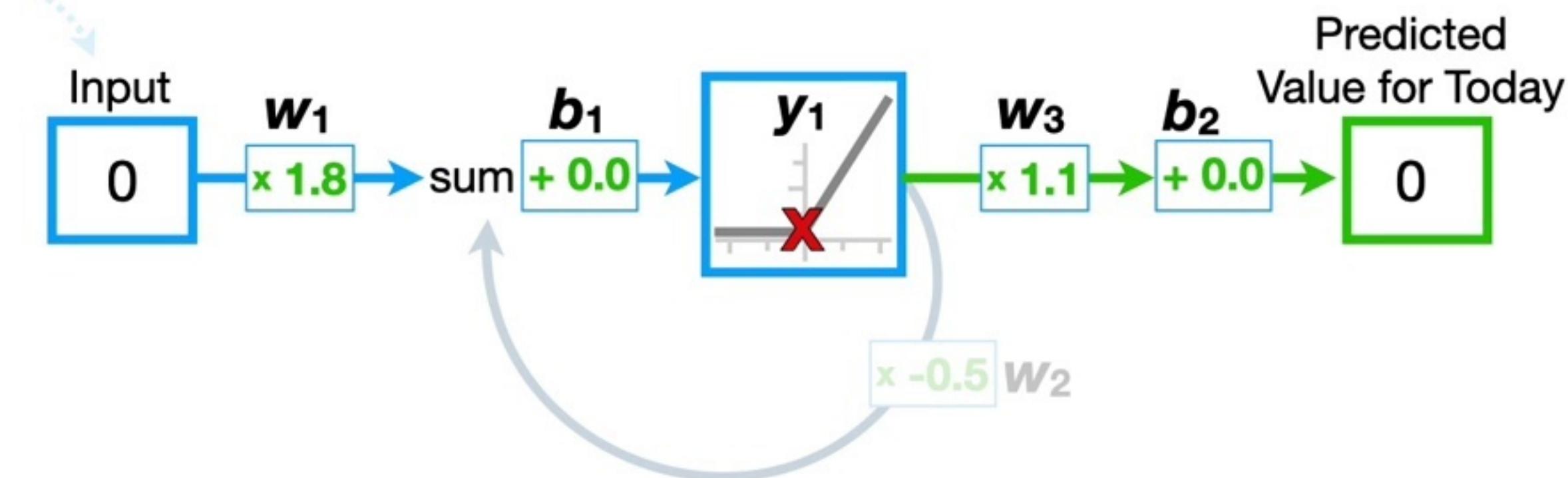


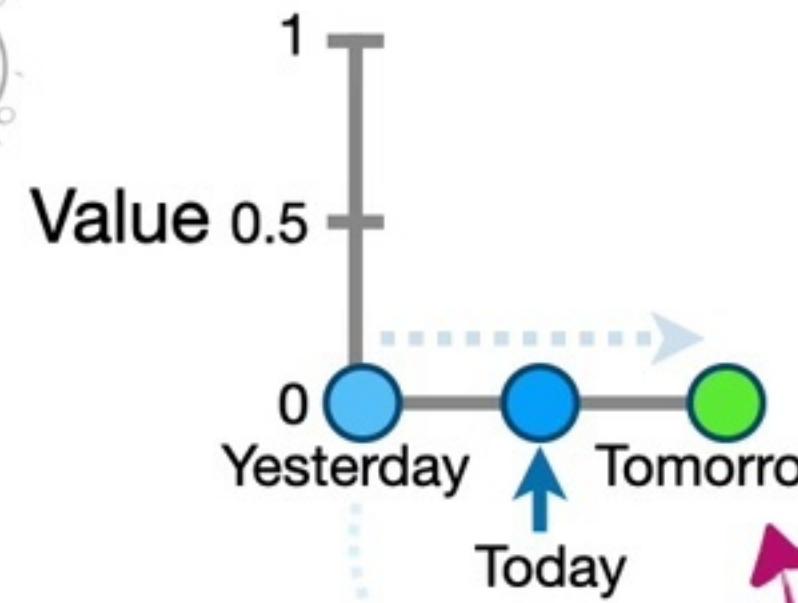


However, we're not interested in the *predicted* value for **today** because we already have the *actual* value for **today**.

$y_1 \times w_3 + b_2 = \text{The Predicted Value for Today}$

$$0 \times 1.1 + 0.0 = 0$$

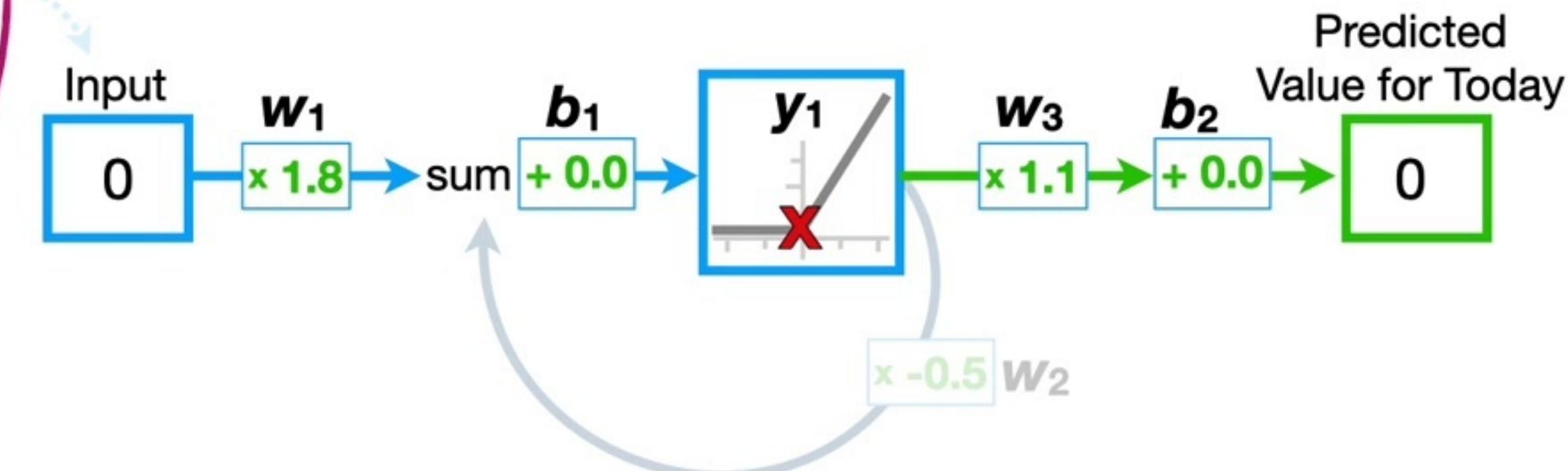


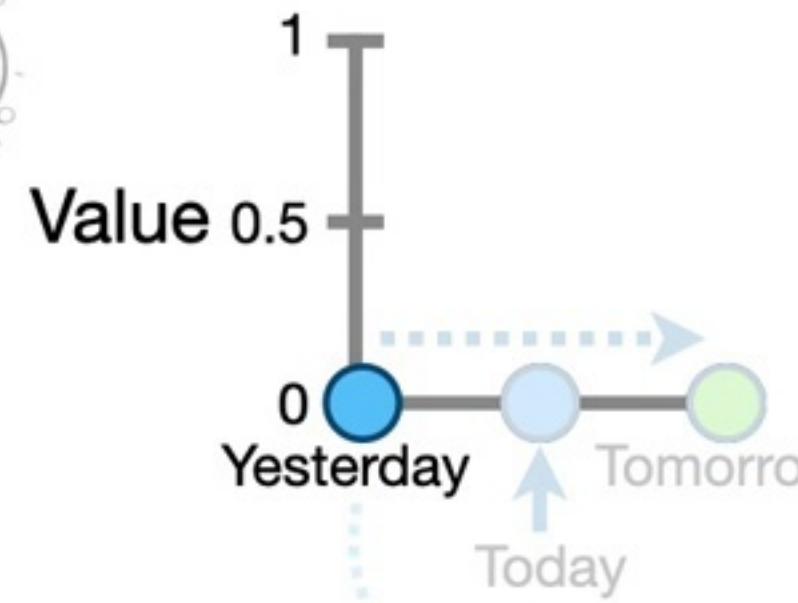


Instead, we want to use both **yesterday** and **today's** value to predict **tomorrow's** value.

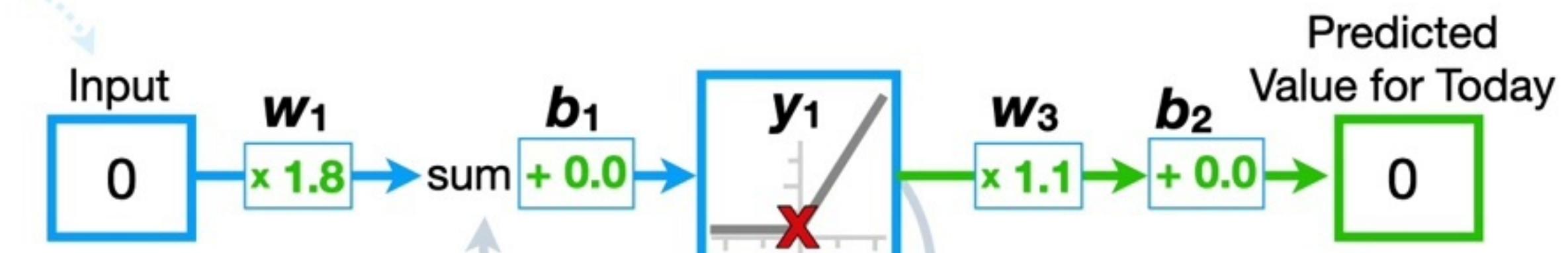
$y_1 \times w_3 + b_2 = \text{The Predicted Value for Today}$

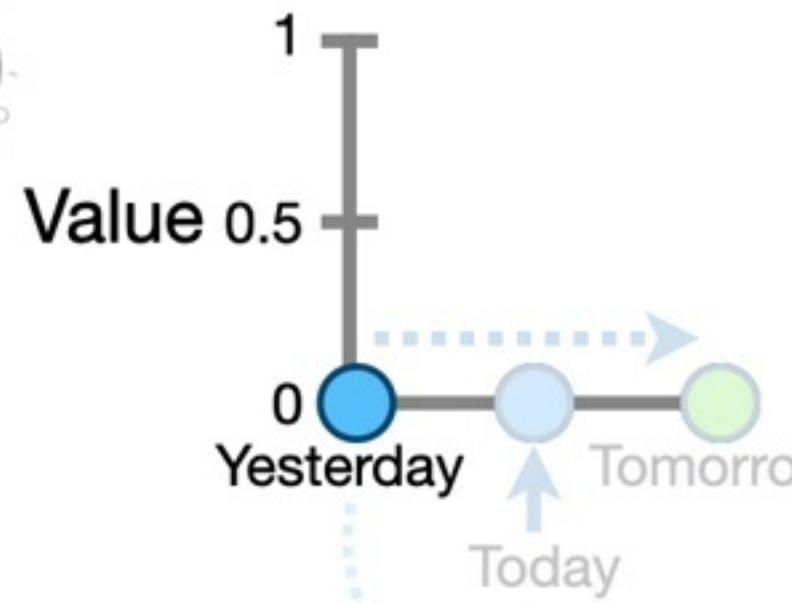
$$0 \times 1.1 + 0.0 = 0$$



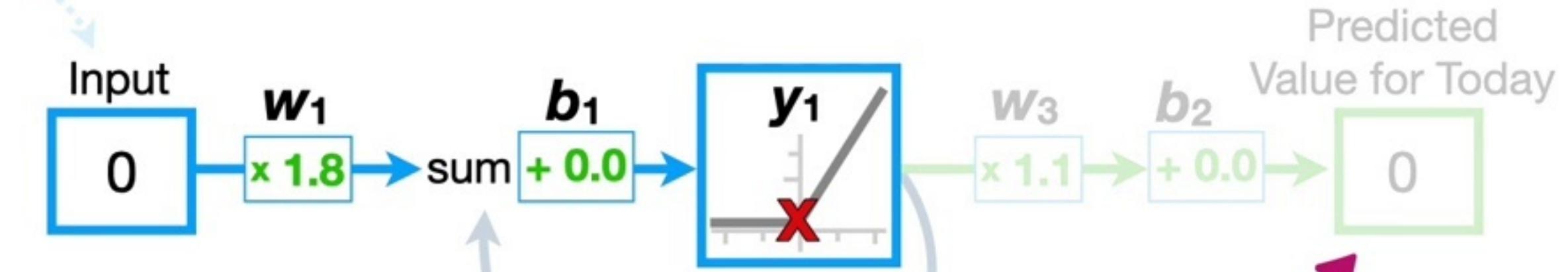


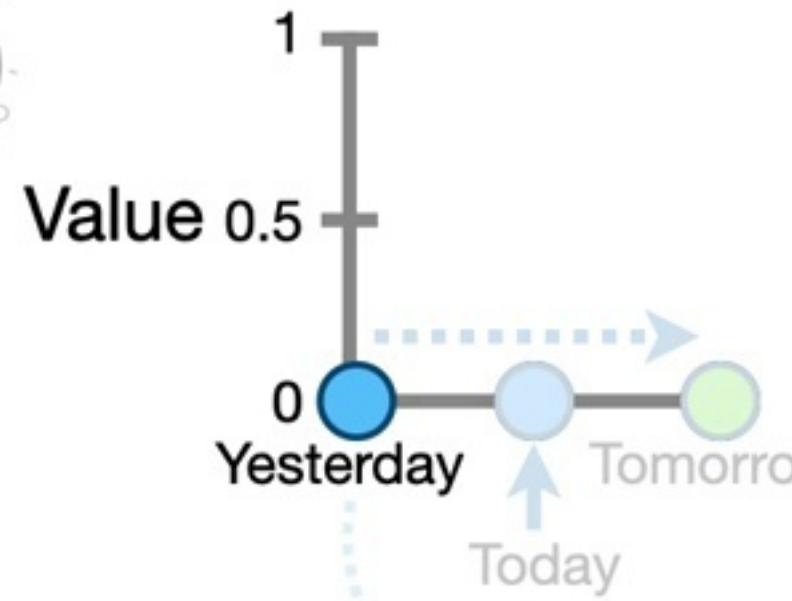
So, for now, we'll ignore this output...



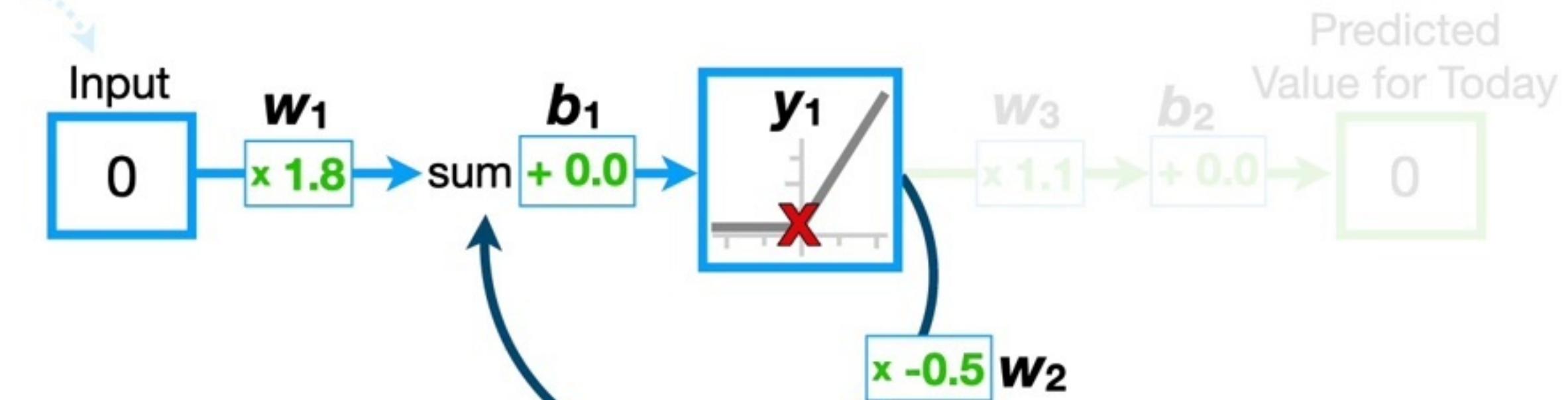


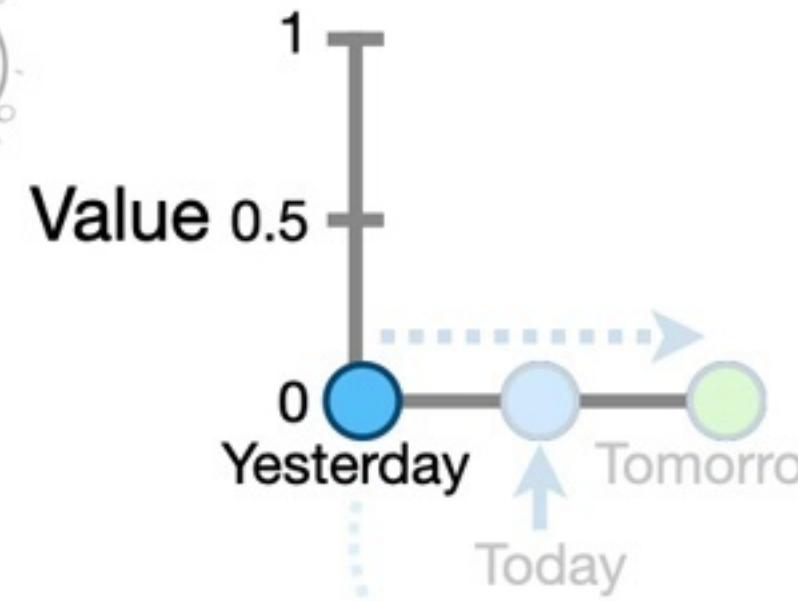
So, for now, we'll ignore this output...



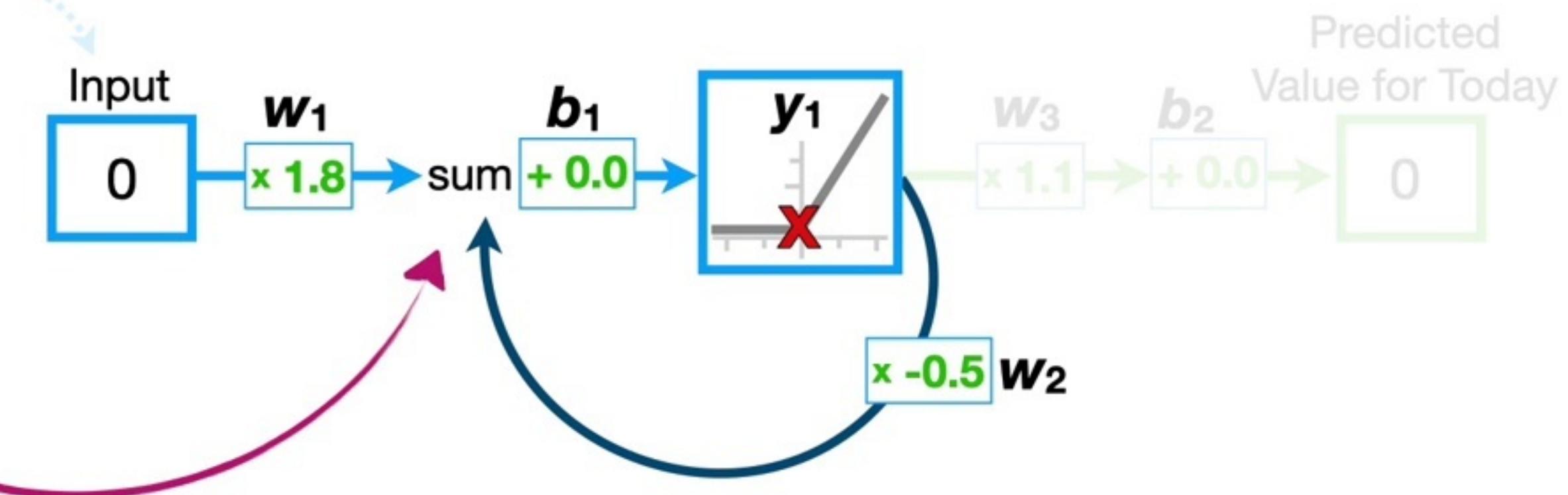


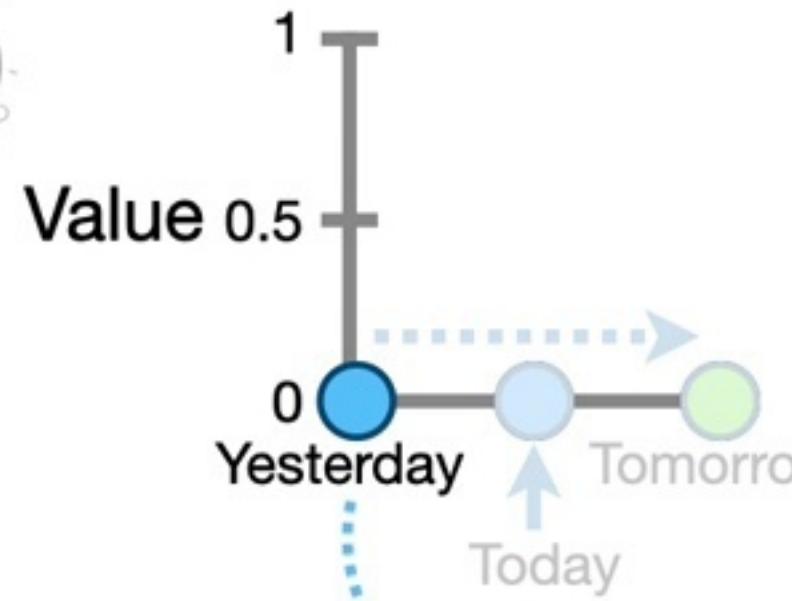
...and, instead, focus
on what happens with
this **feedback loop**.



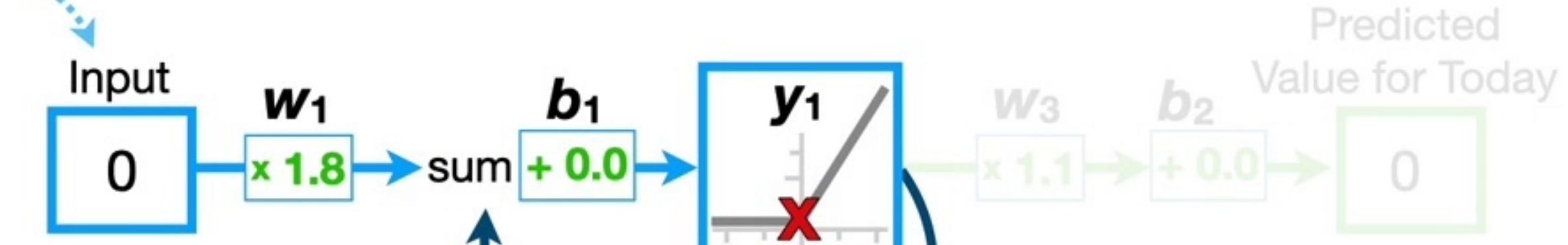


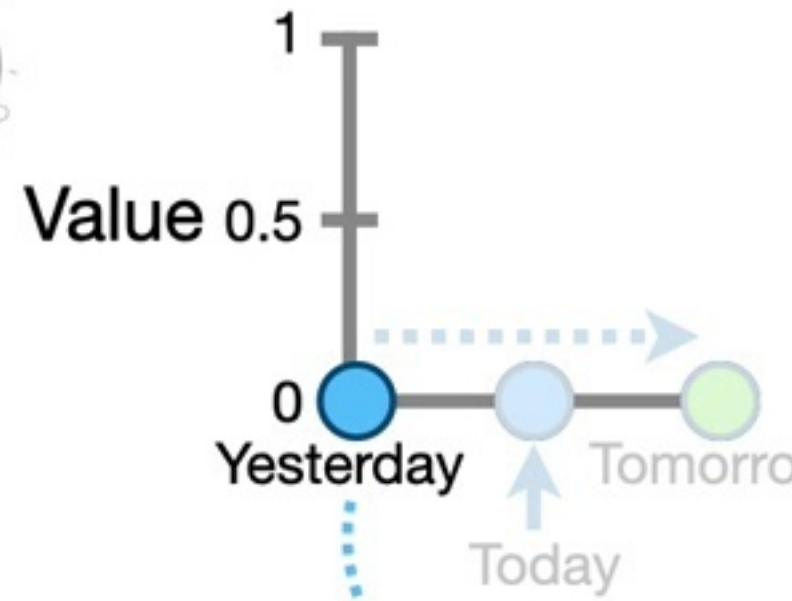
The key to understanding how the feedback loop works is this summation.



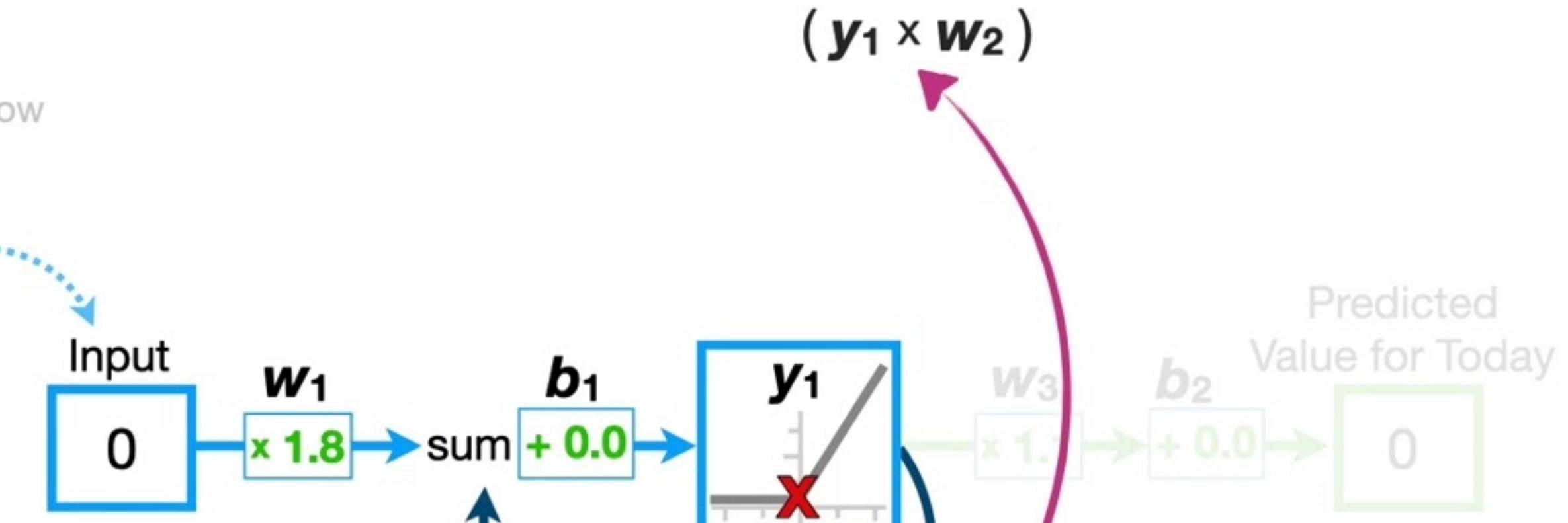


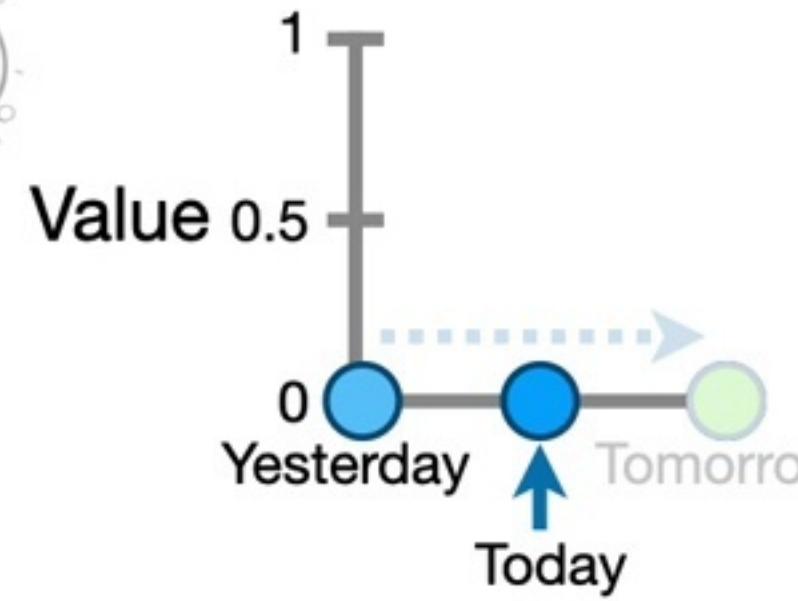
The summation allows us to add $y_1 \times w_2$, which is based on **yesterday's** value...



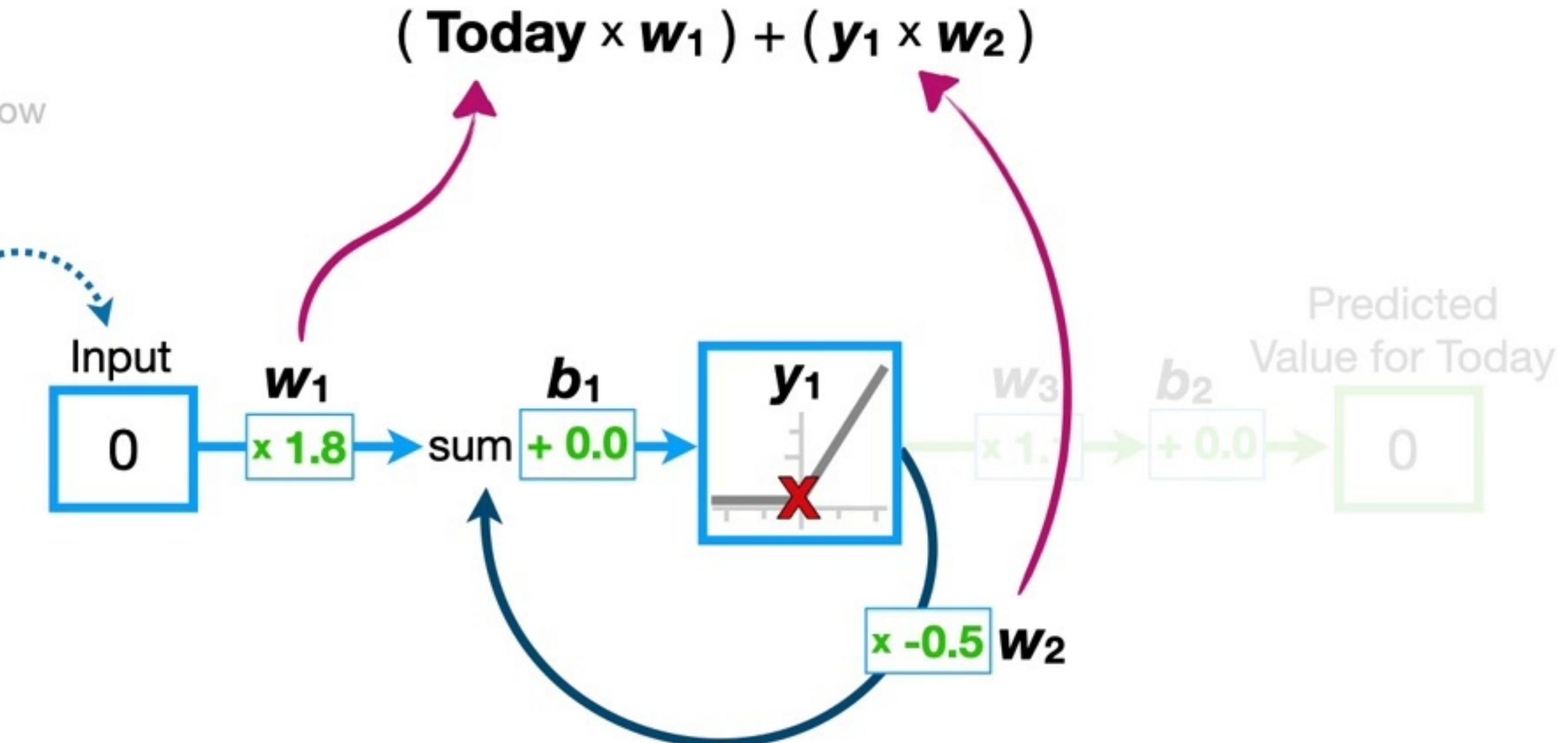


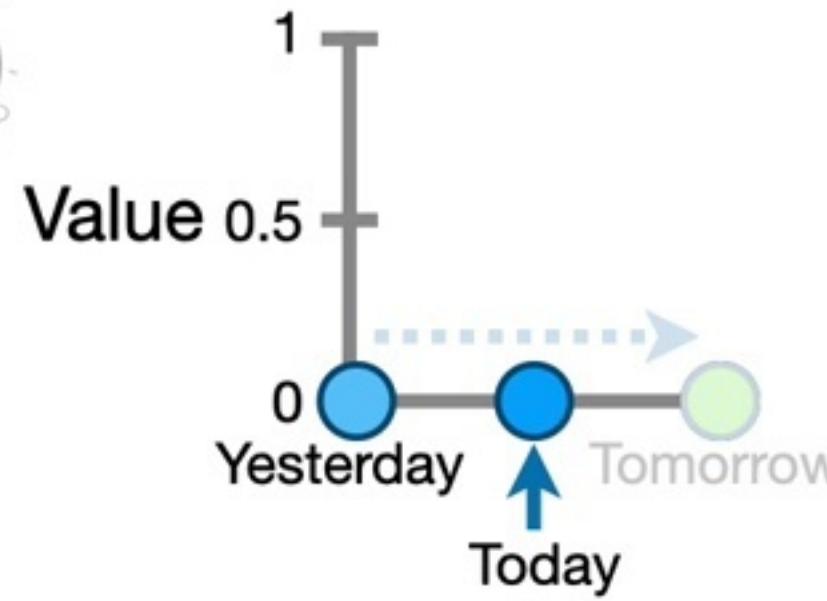
The summation allows us to add $y_1 \times w_2$, which is based on **yesterday's** value...



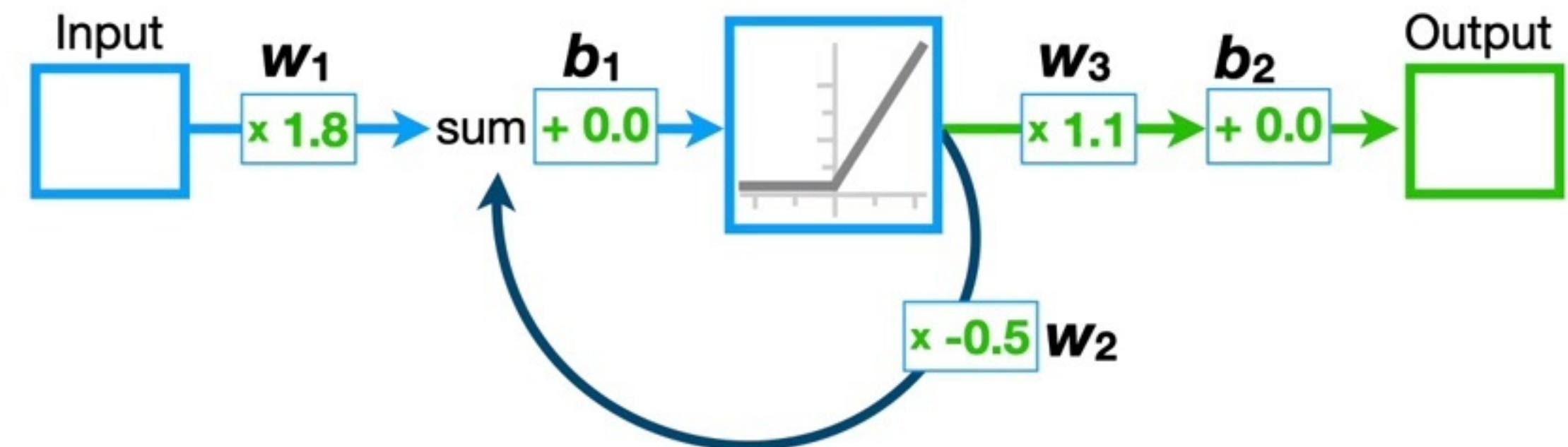


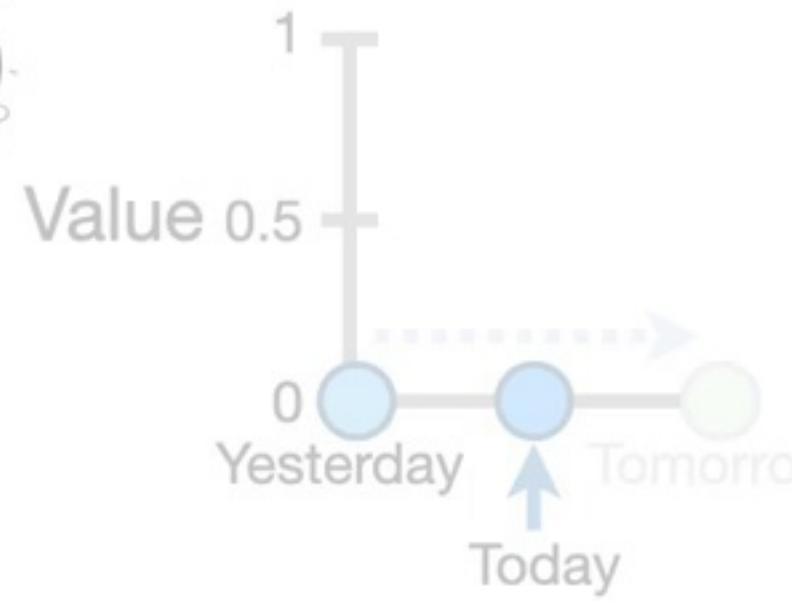
In other words, the **feedback loop** allows both **yesterday** and **today's** values to influence the prediction.



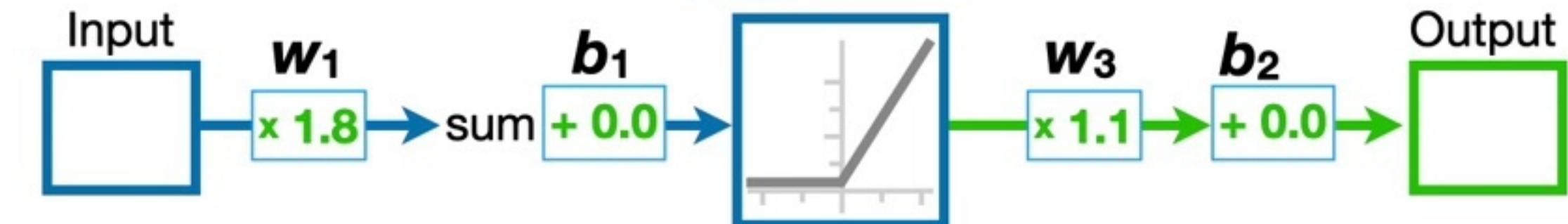
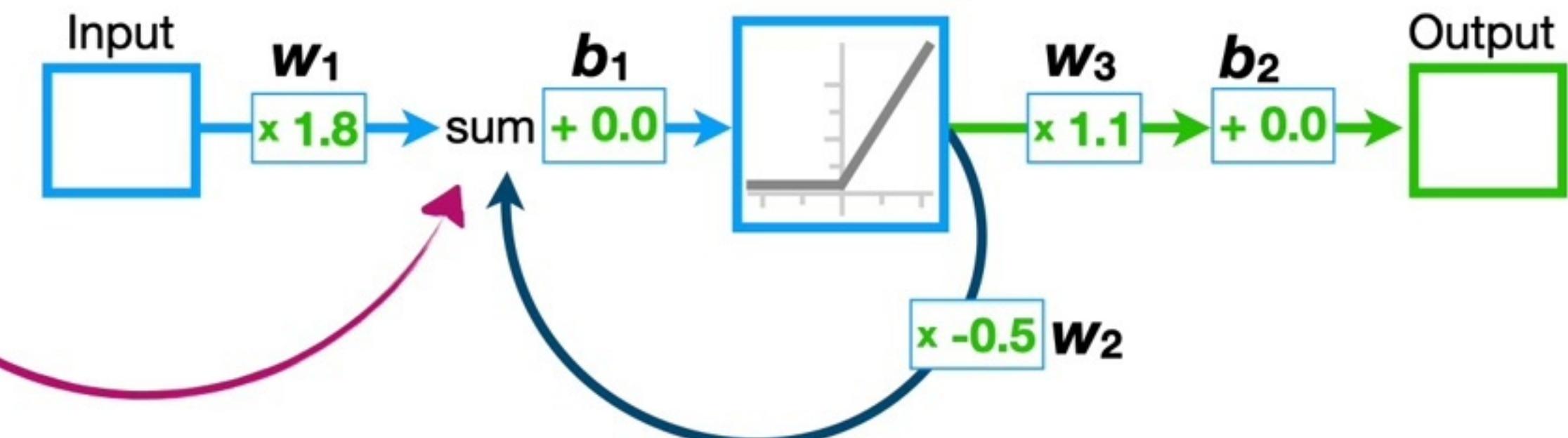


...we can **unroll** the feedback loop by making a copy of the neural network for each input value.



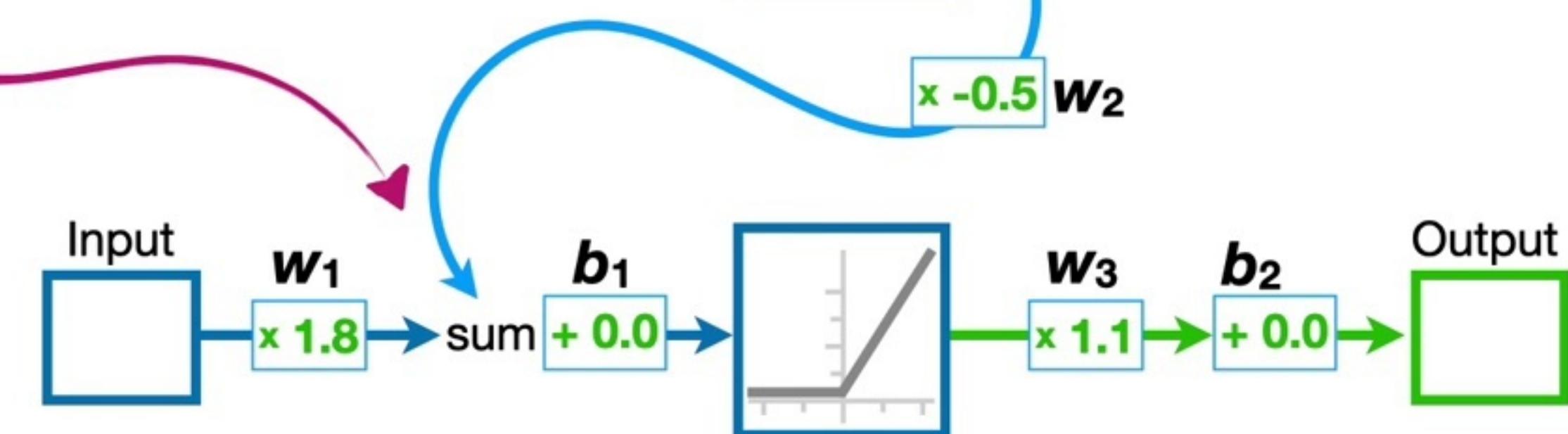
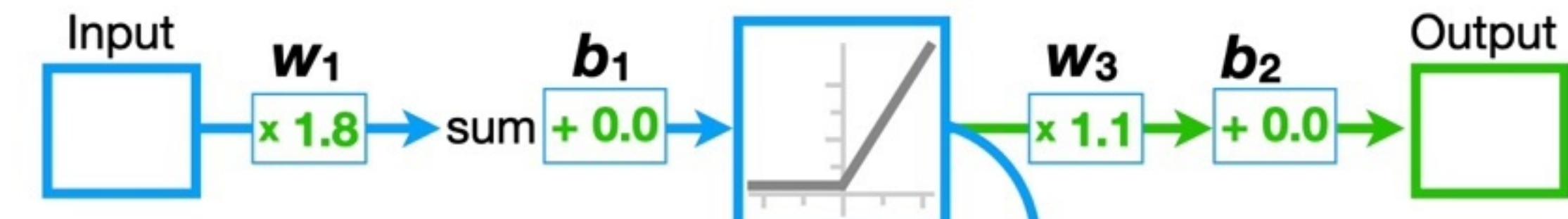


Now, instead of pointing the **feedback loop** to the sum in the first copy...



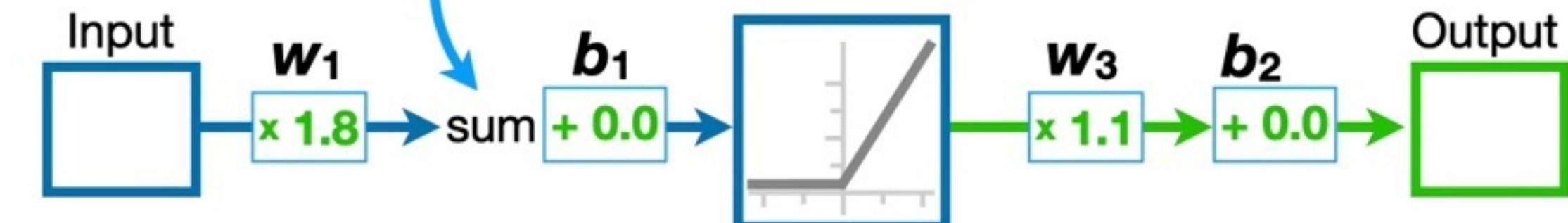


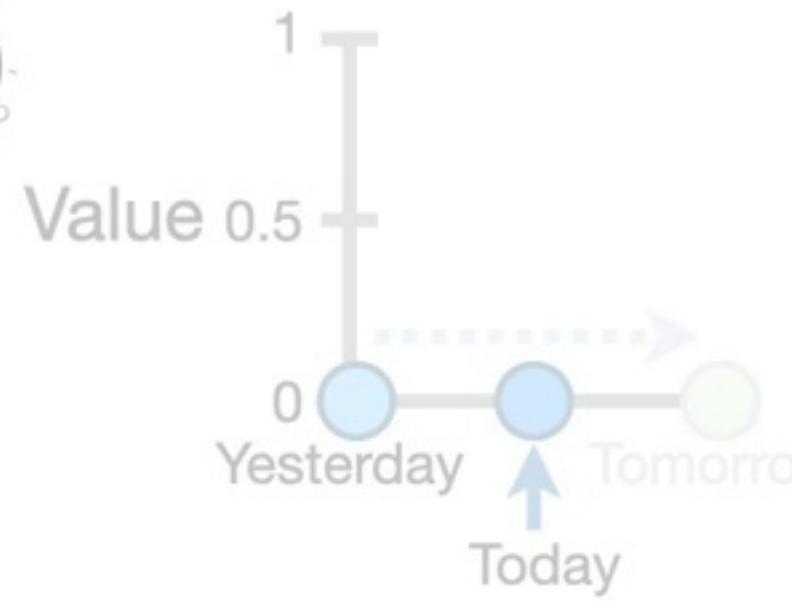
...we can point it to
the sum in the second
copy.



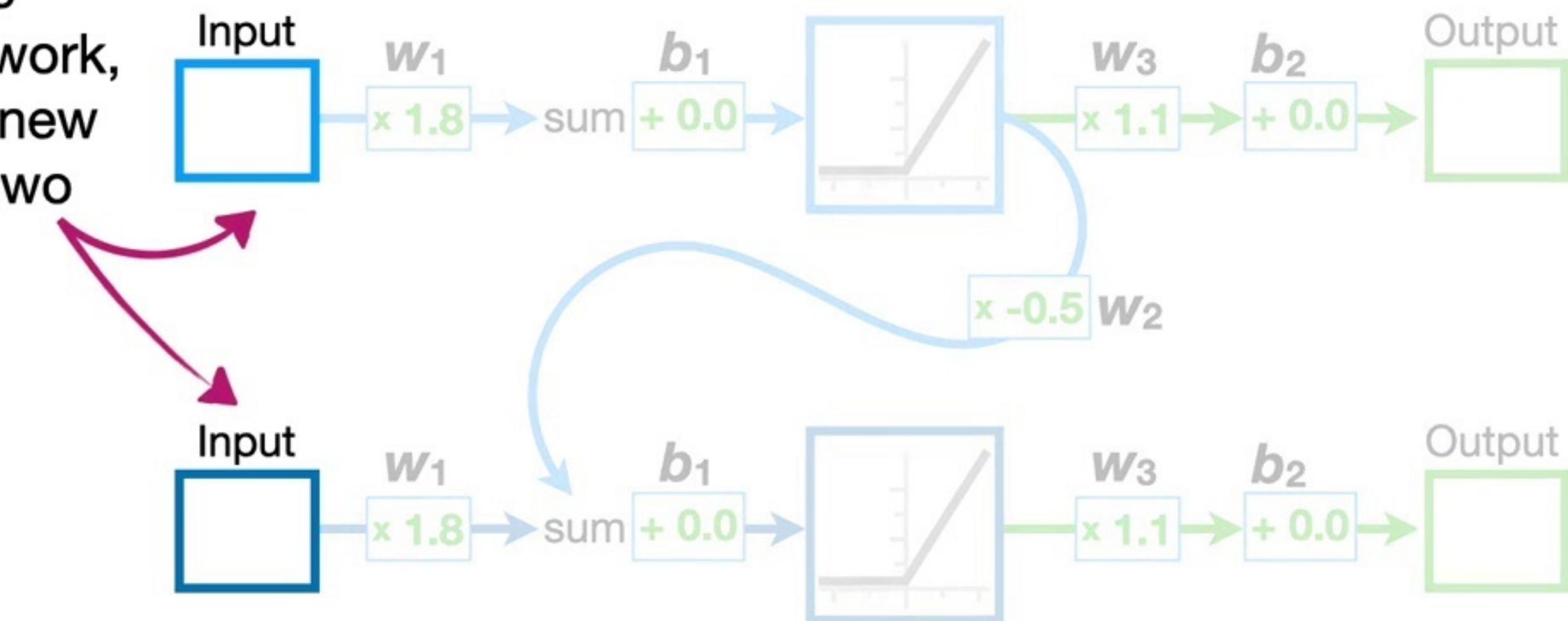


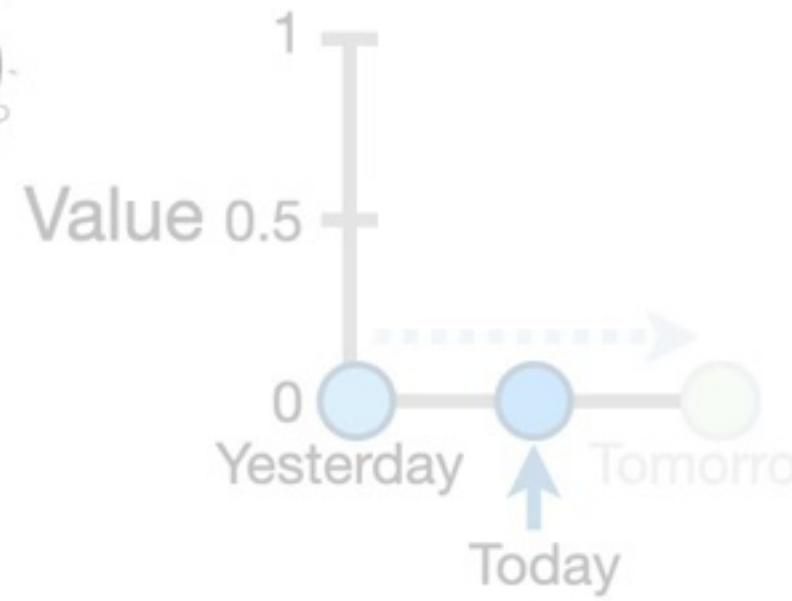
By **unrolling** the recurrent neural network, we end up with a new network that has two inputs...





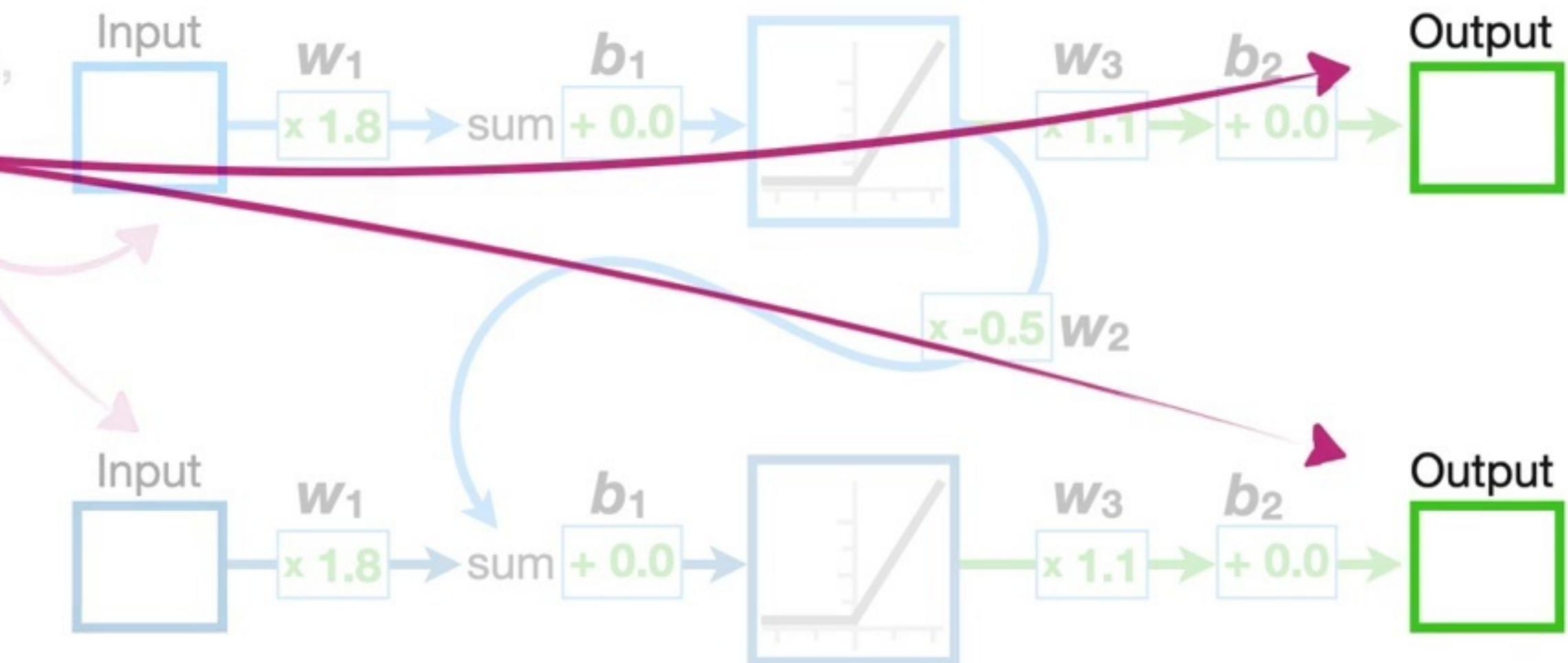
By unrolling the recurrent neural network, we end up with a new network that has two inputs...

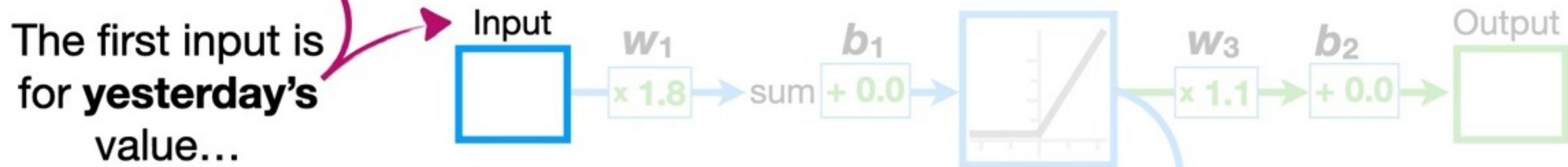
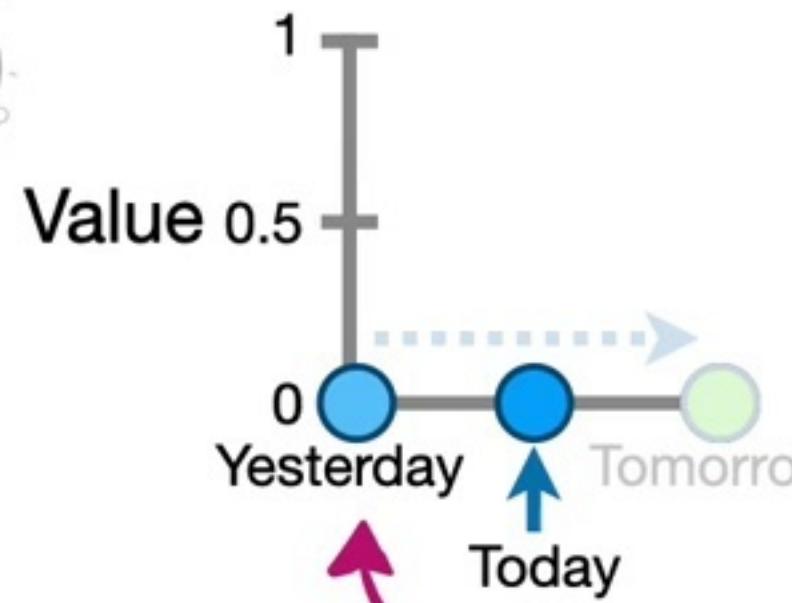


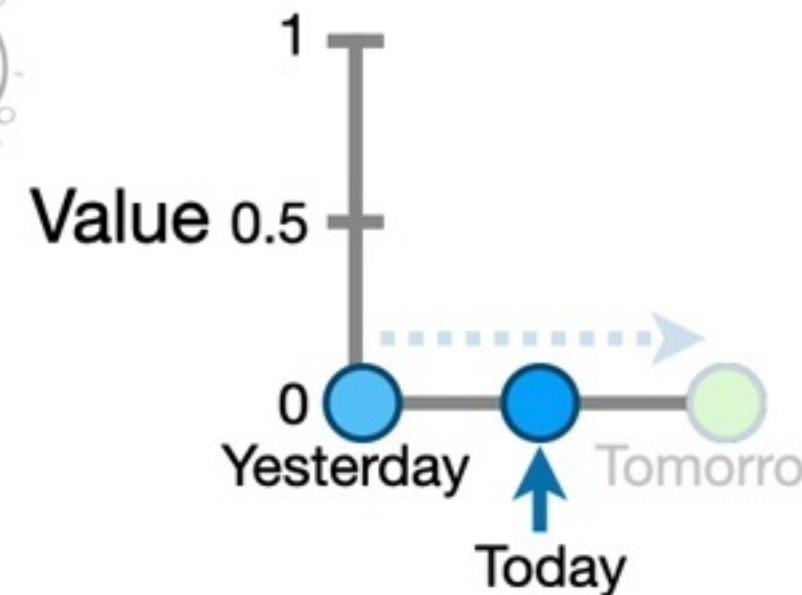


By unrolling the recurrent neural network,
we end up with a new network that has two inputs...

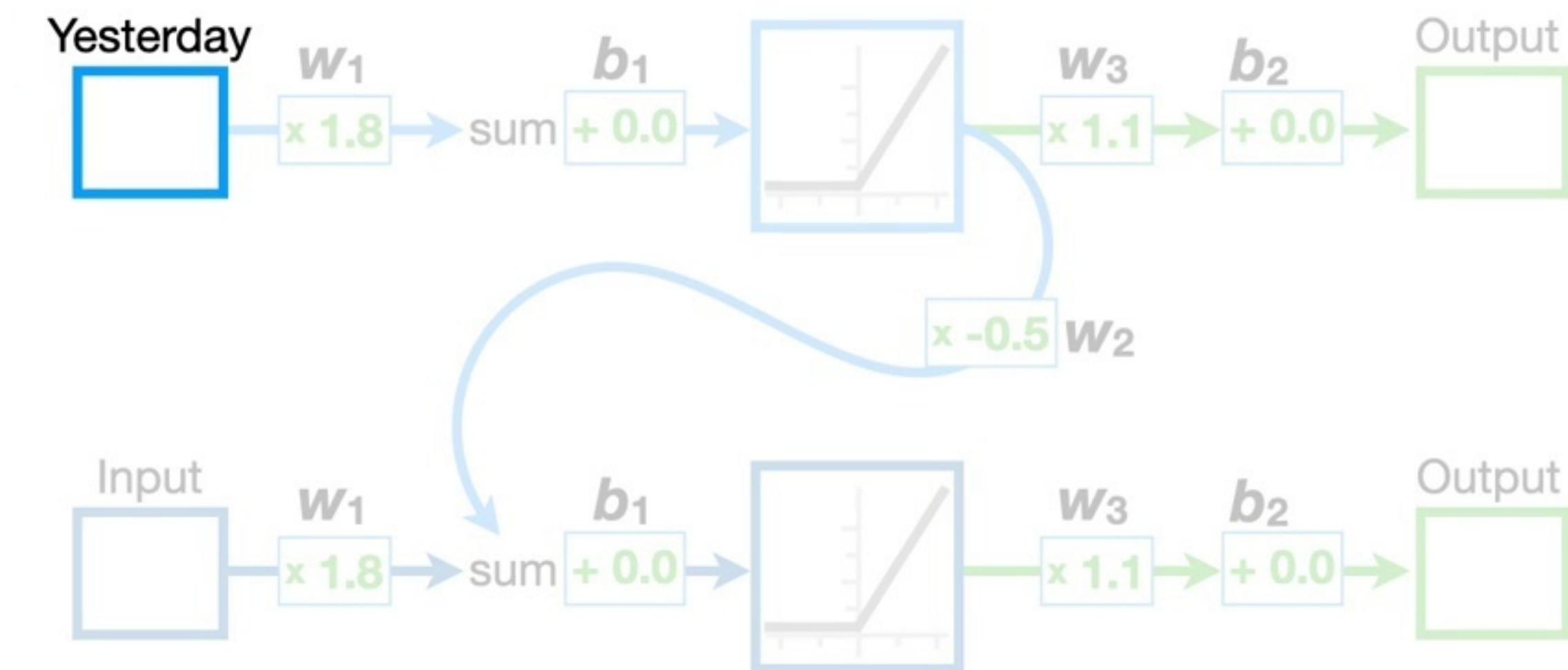
...and two outputs.

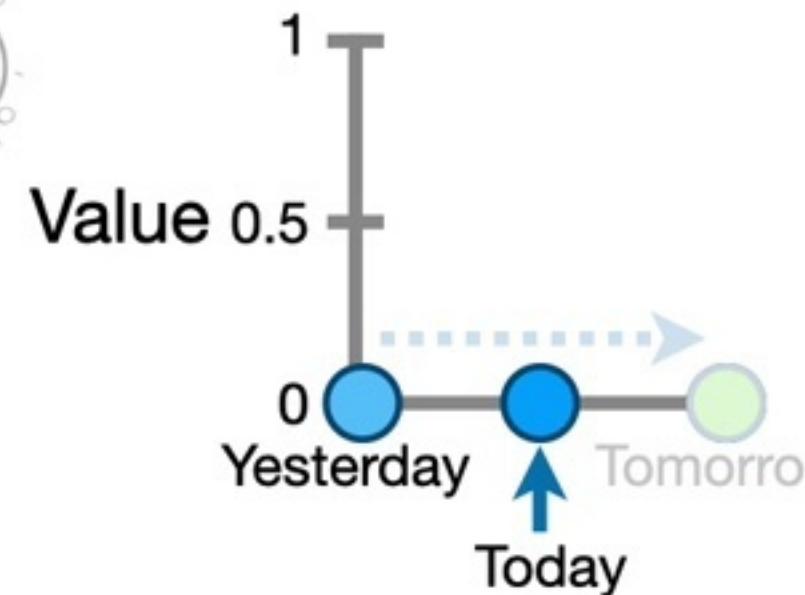




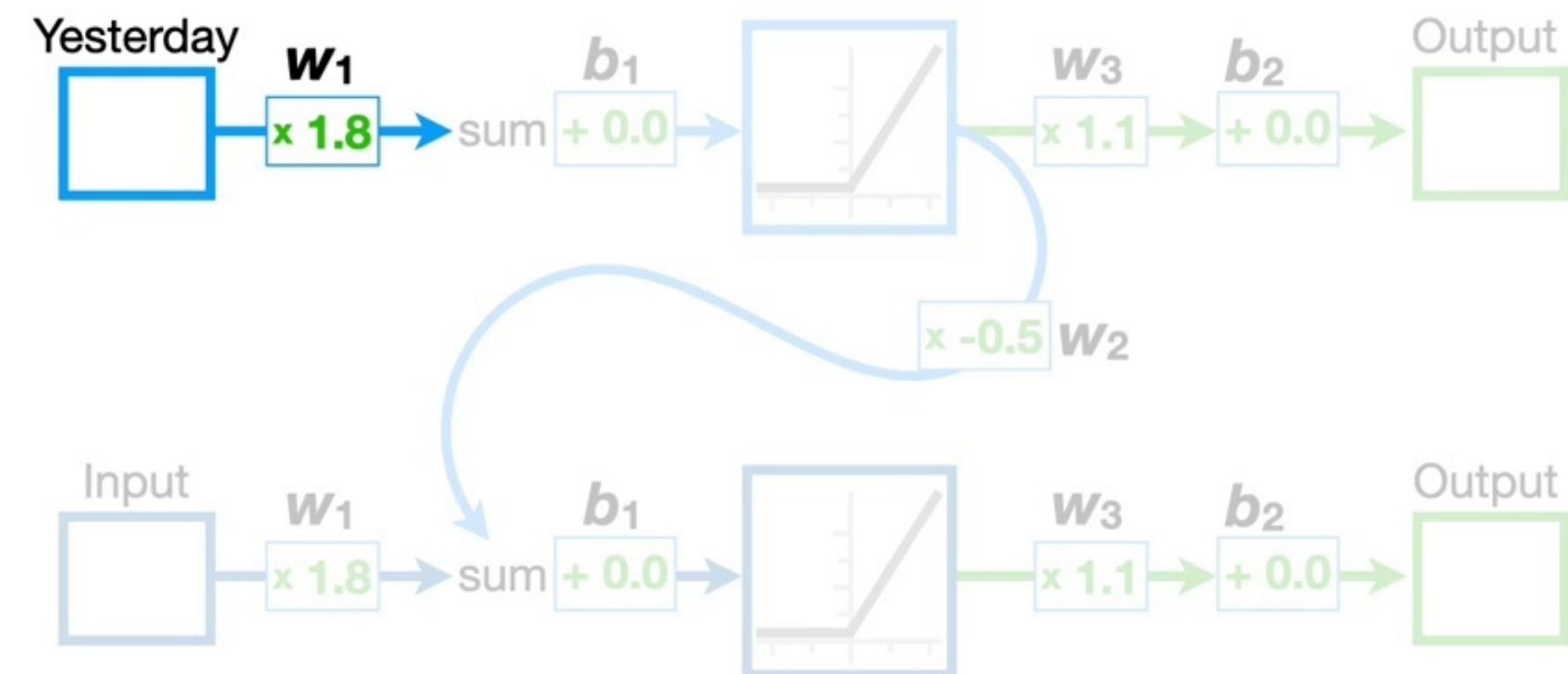


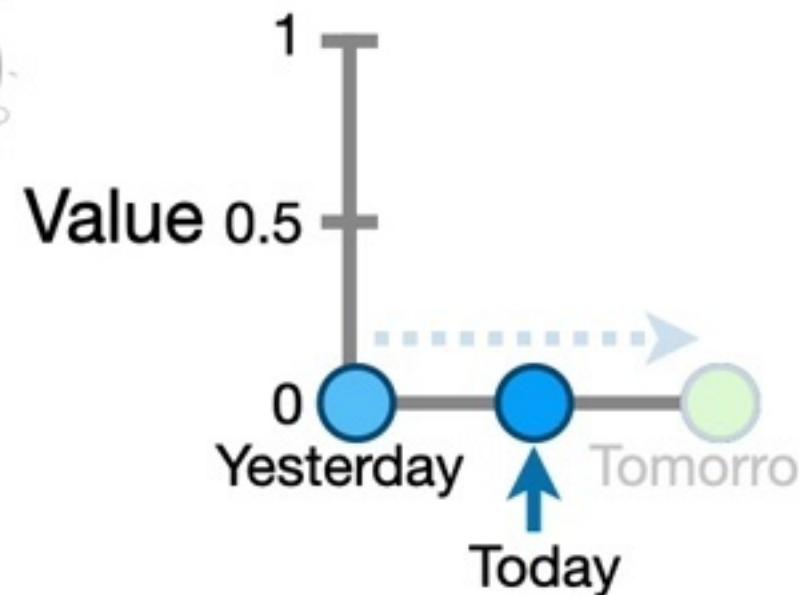
...and if we do
the math straight
through to the
first output like
we did earlier...



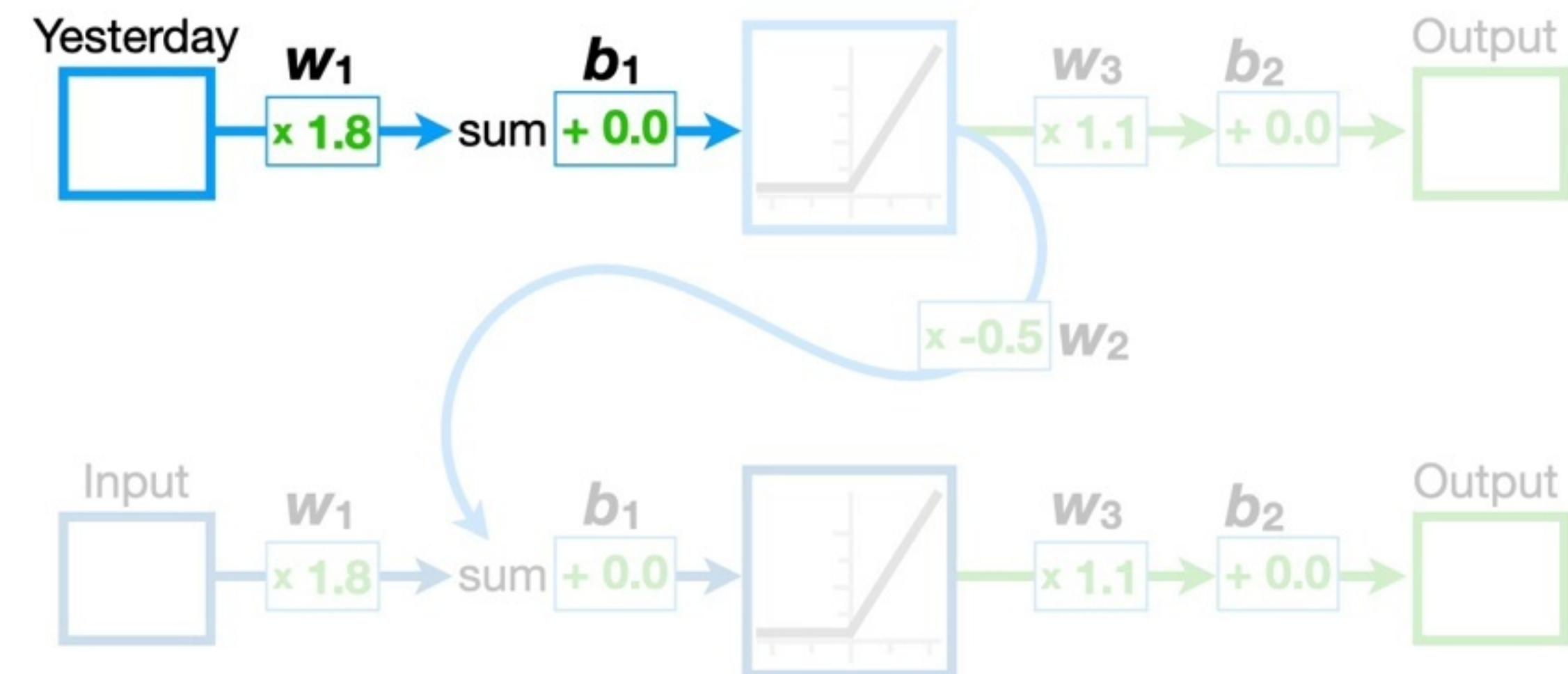


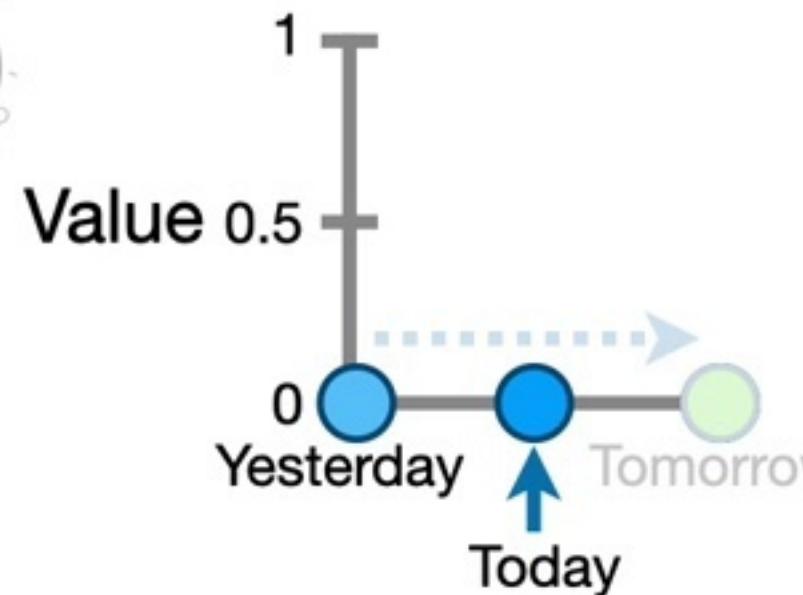
...and if we do
the math straight
through to the
first output like
we did earlier...



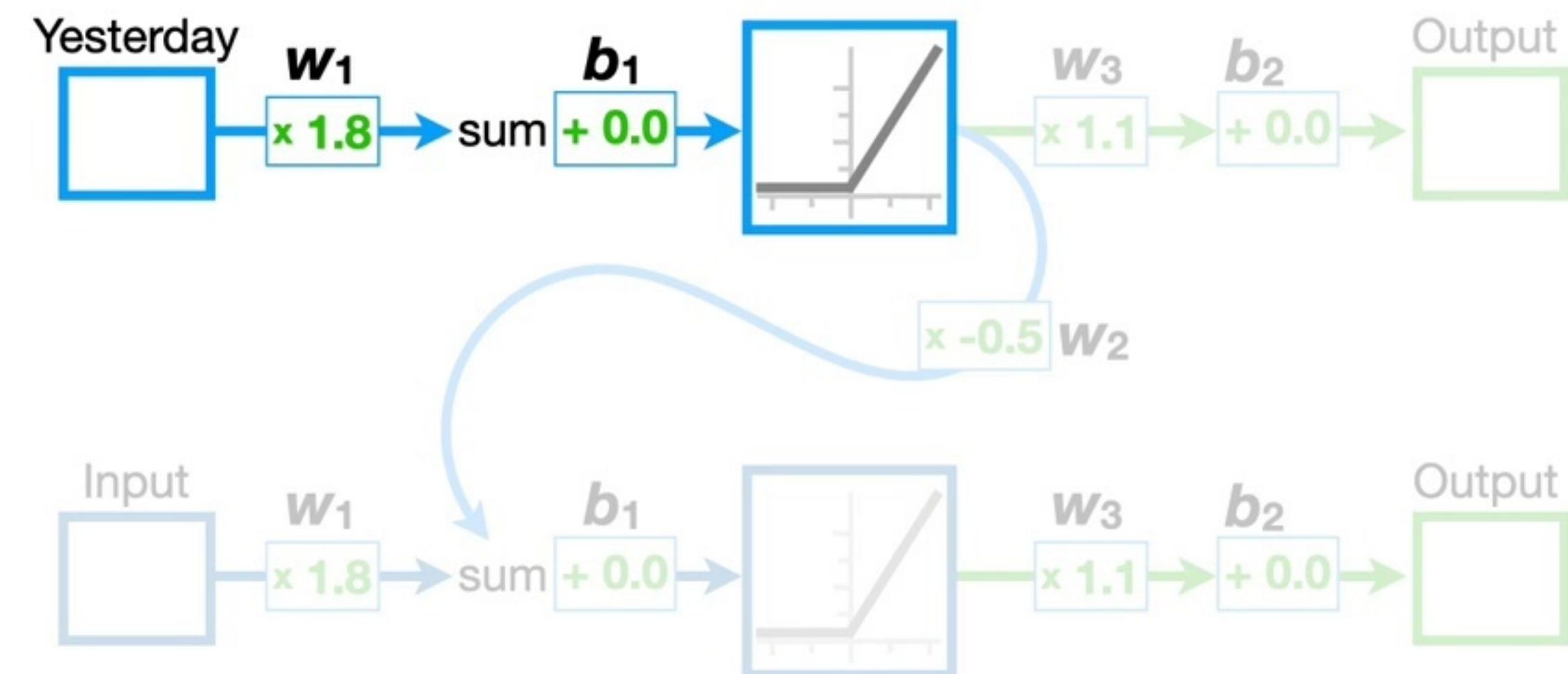


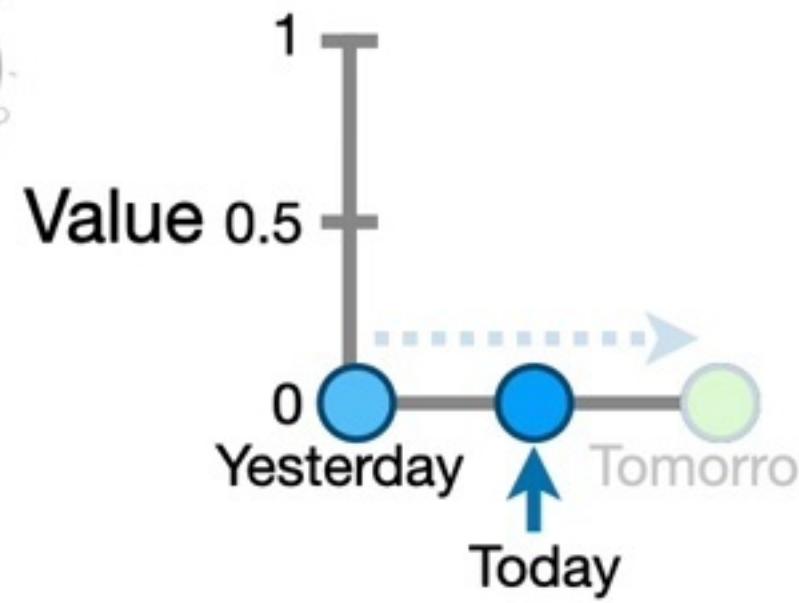
...and if we do
the math straight
through to the
first output like
we did earlier...



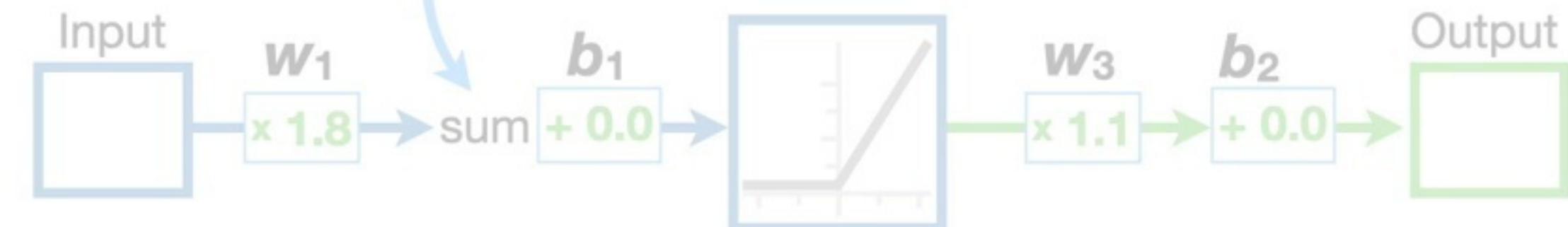


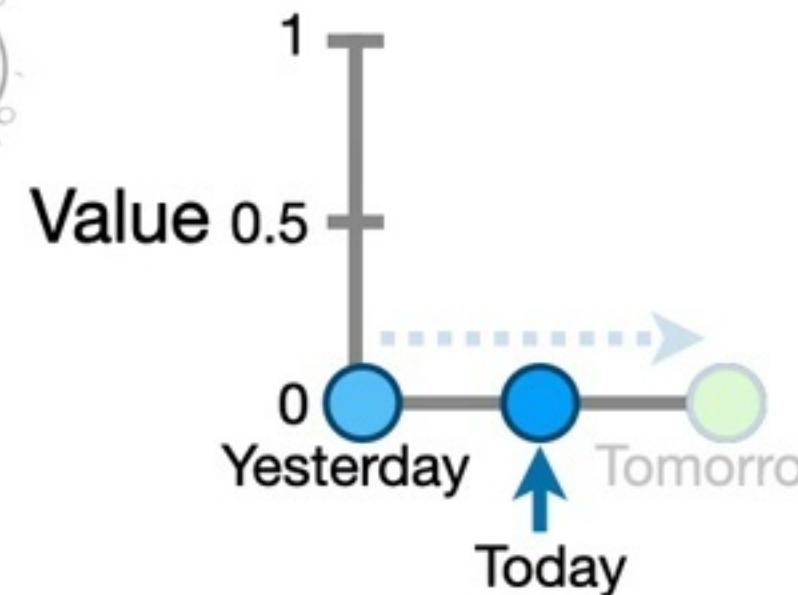
...and if we do
the math straight
through to the
first output like
we did earlier...



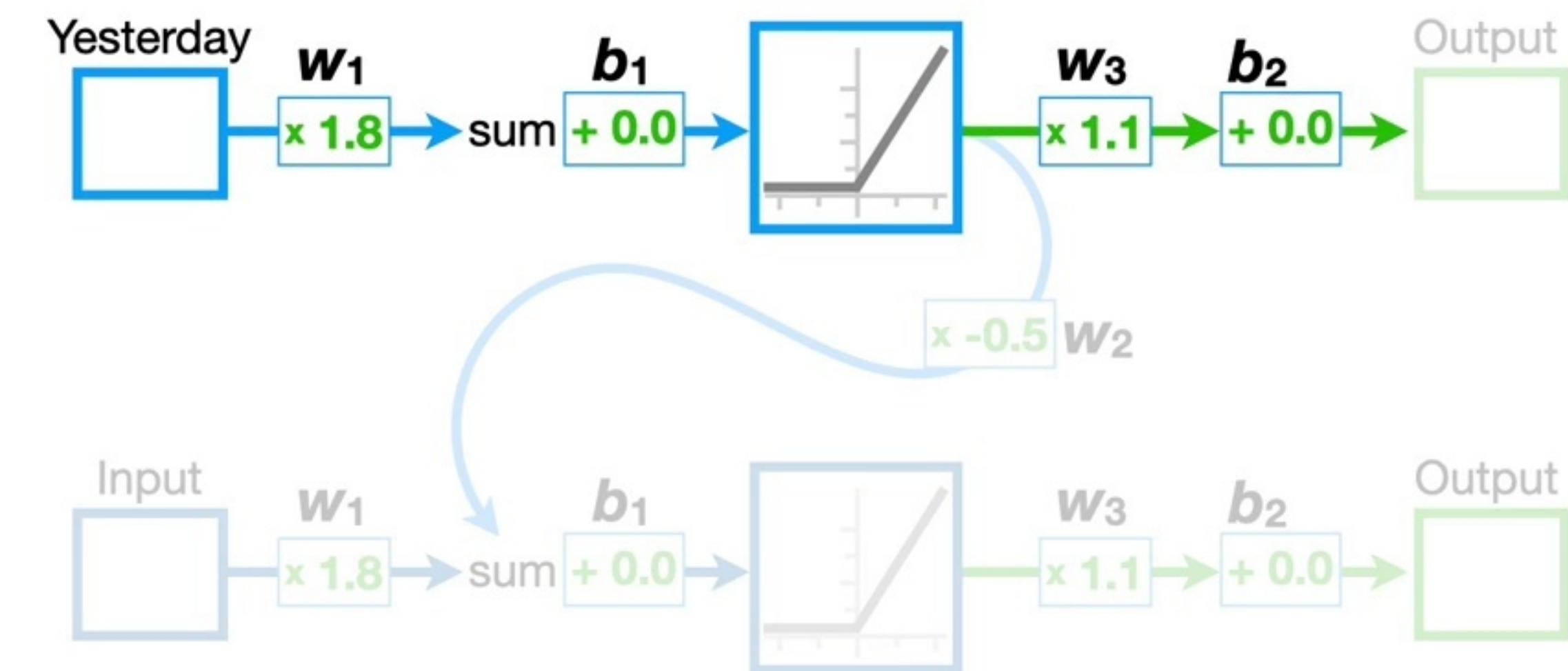


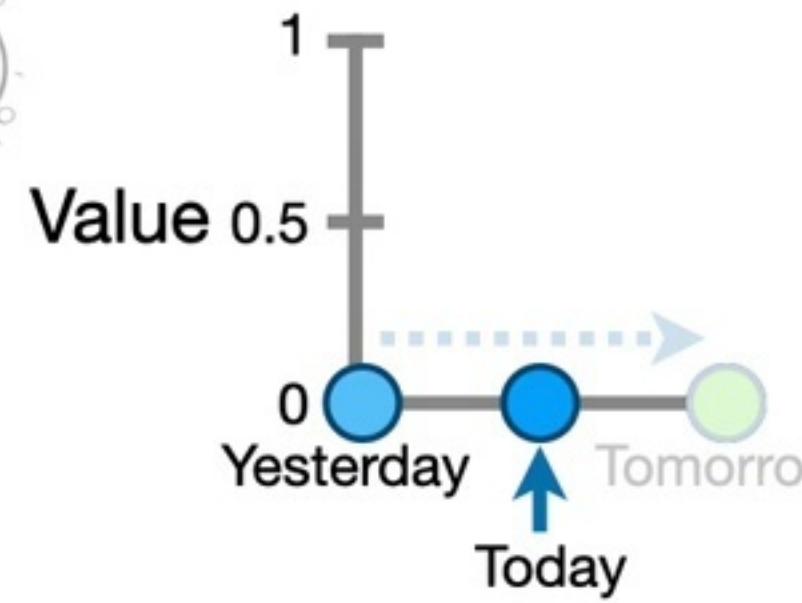
...and if we do
the math straight
through to the
first output like
we did earlier...



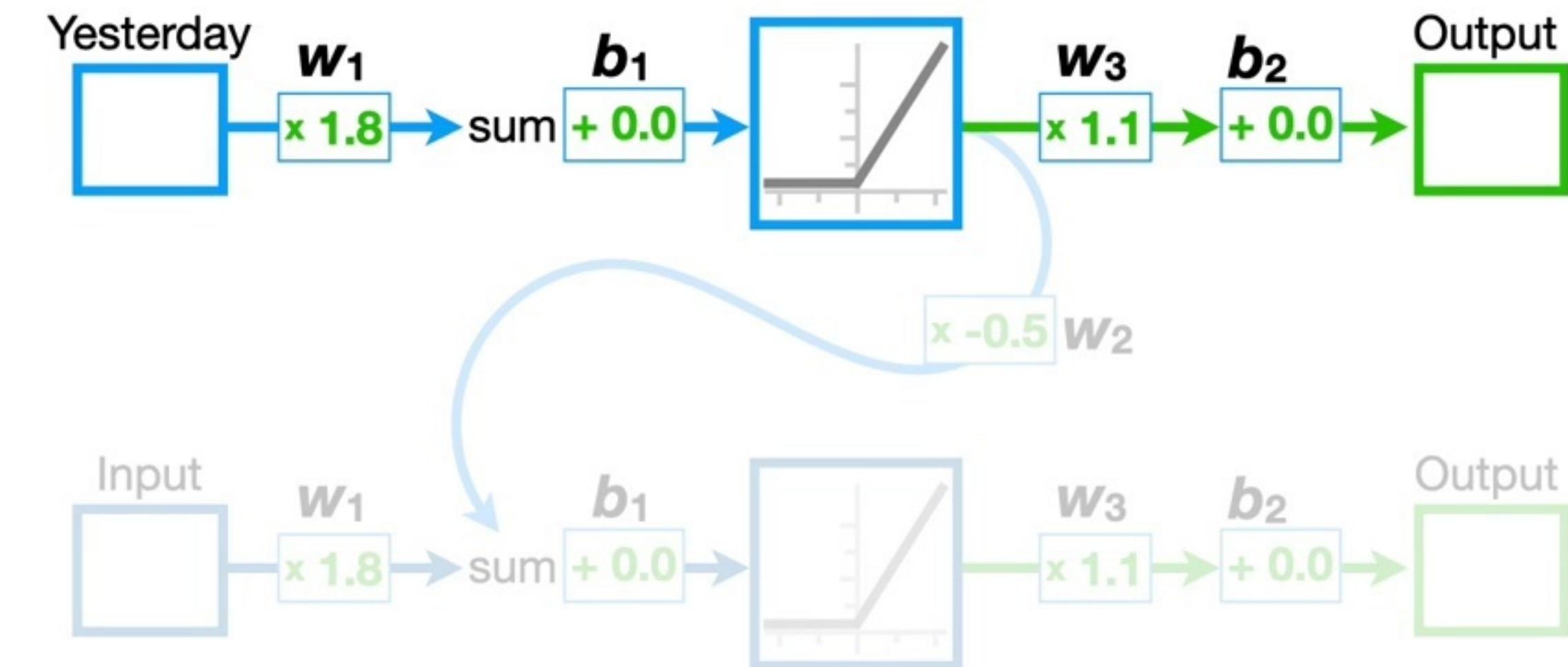


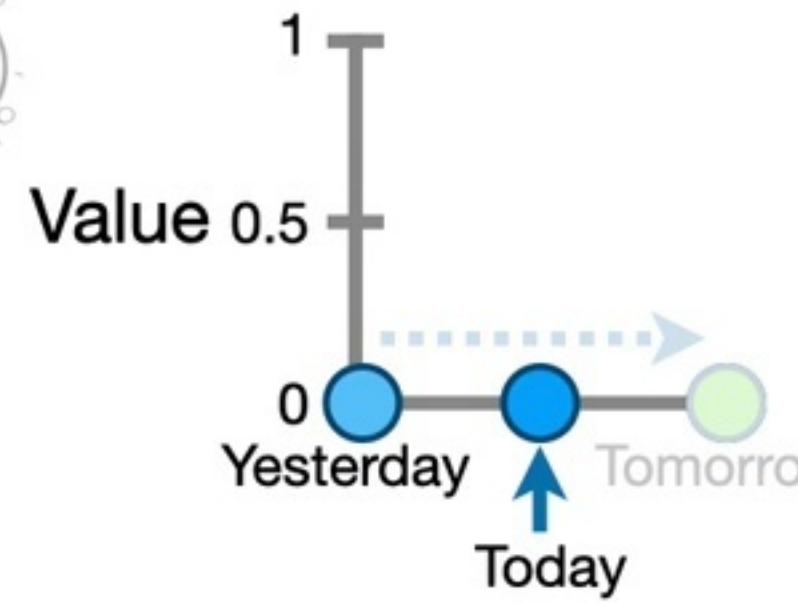
...and if we do
the math straight
through to the
first output like
we did earlier...



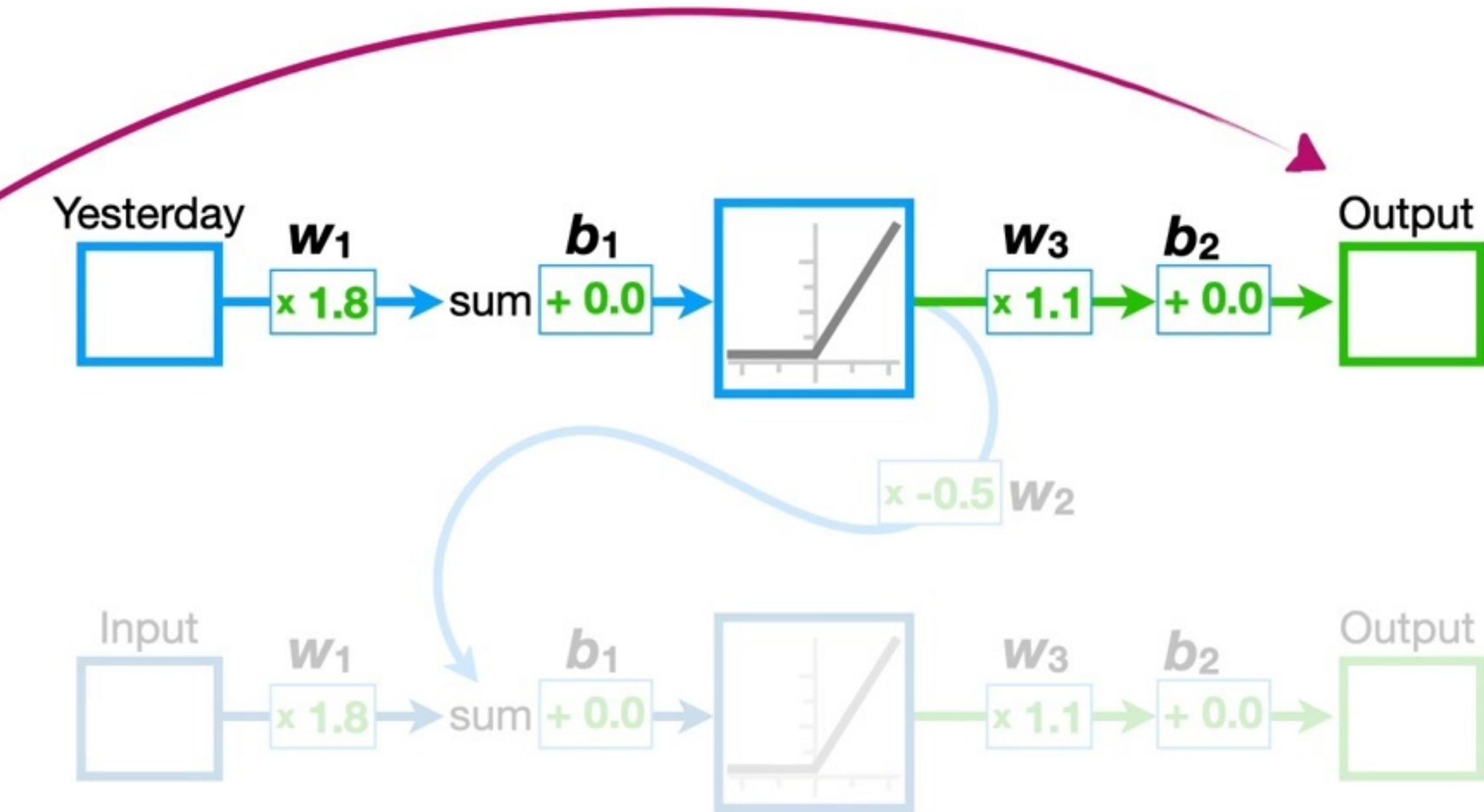


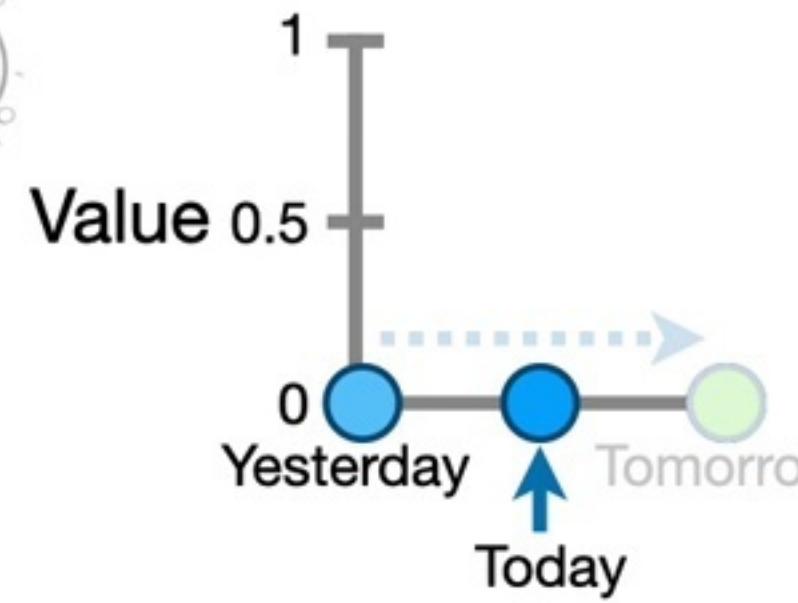
...and if we do
the math straight
through to the
first output like
we did earlier...



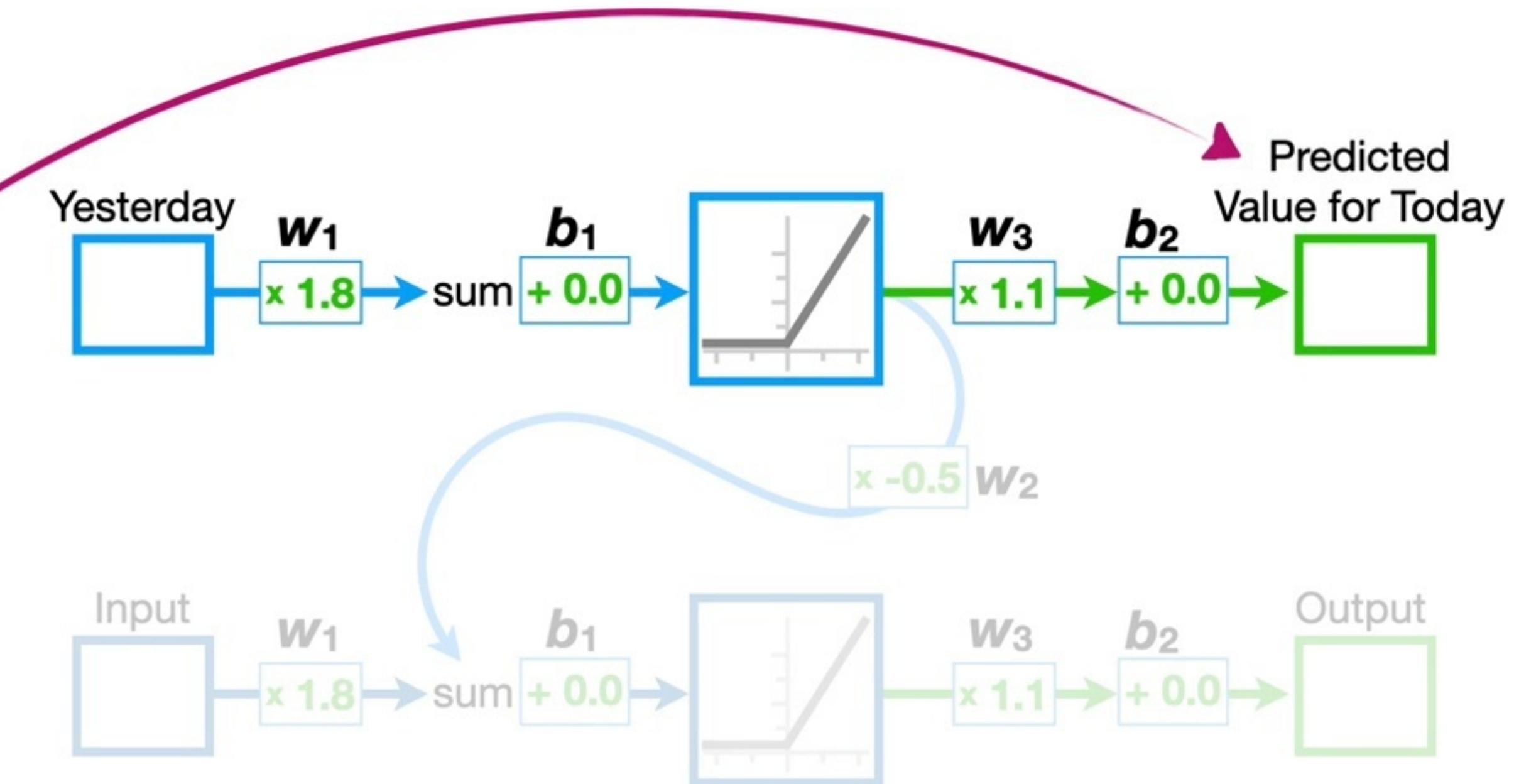


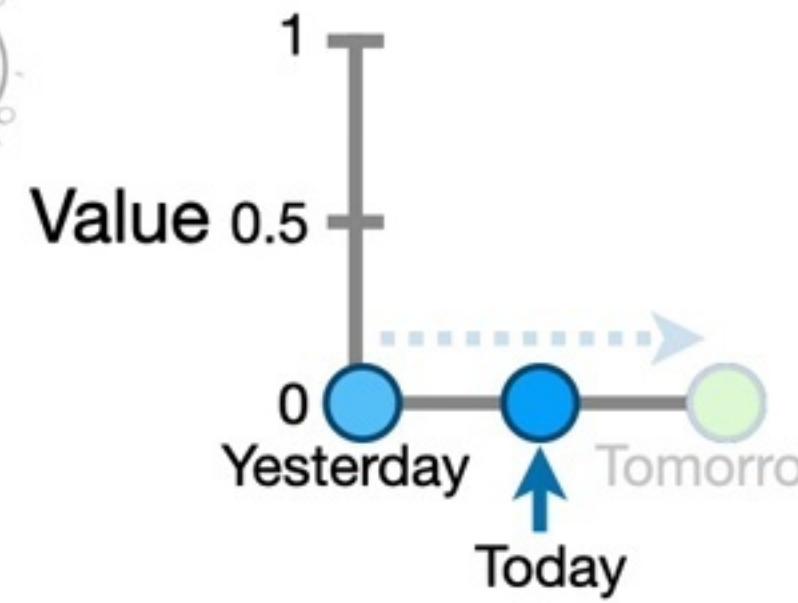
...we get the predicted value for **today**.



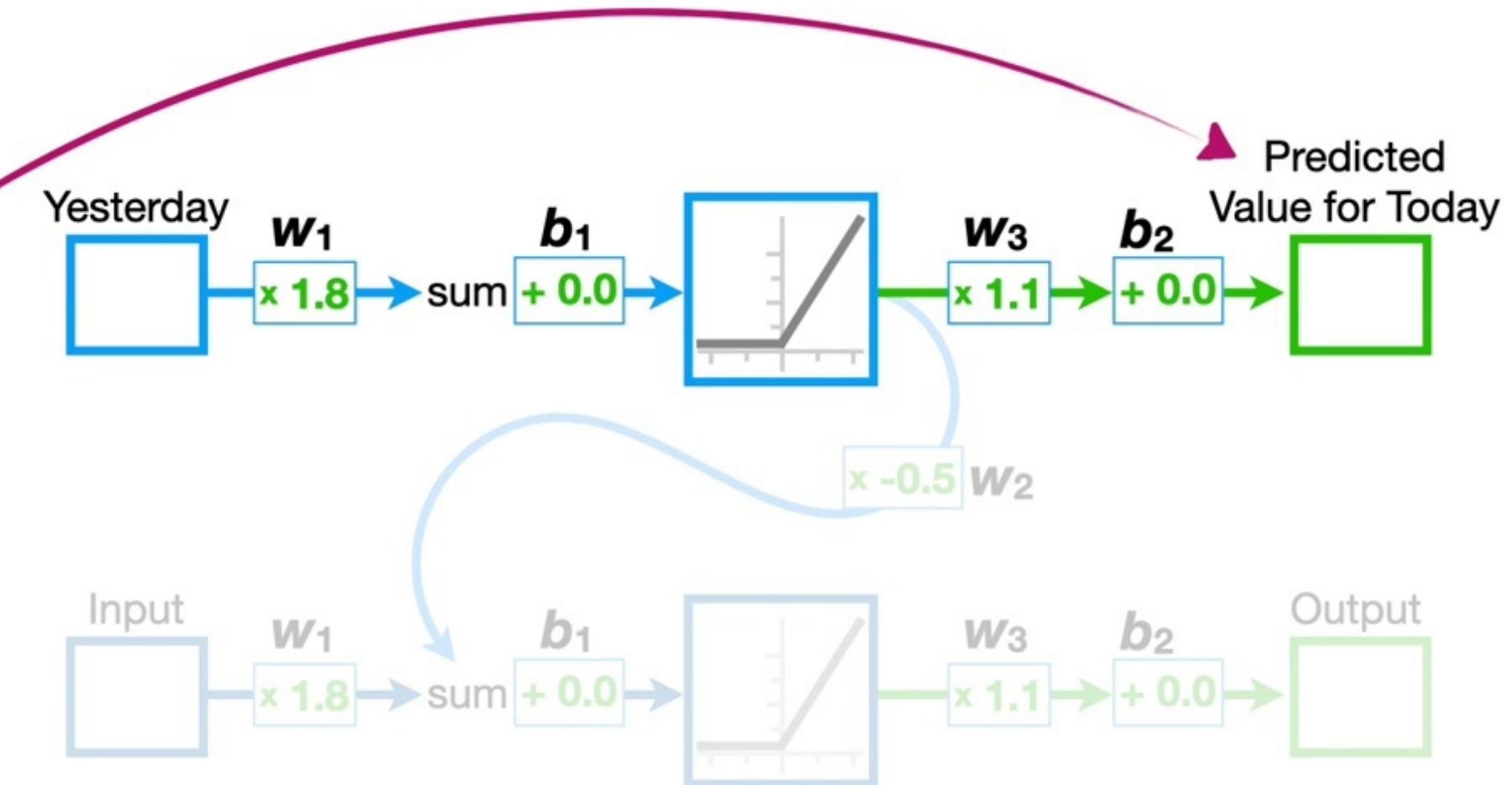


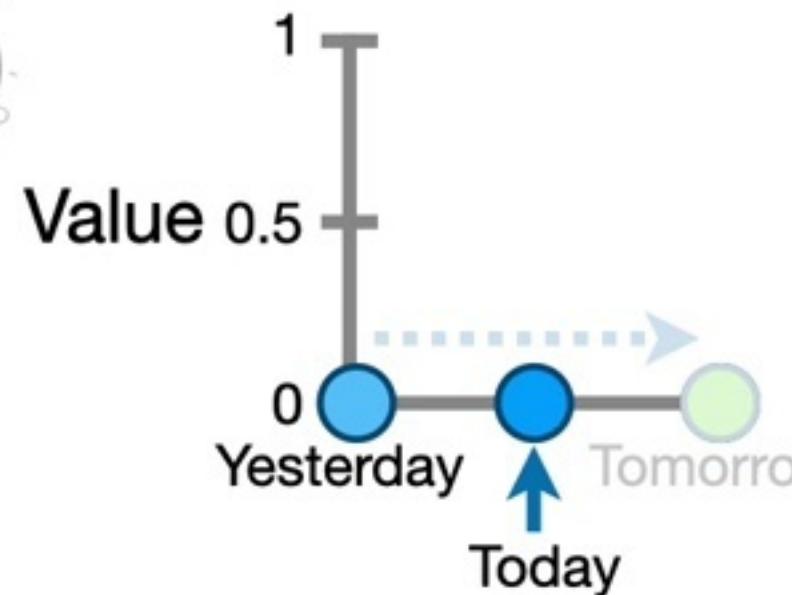
...we get the predicted value for **today**.



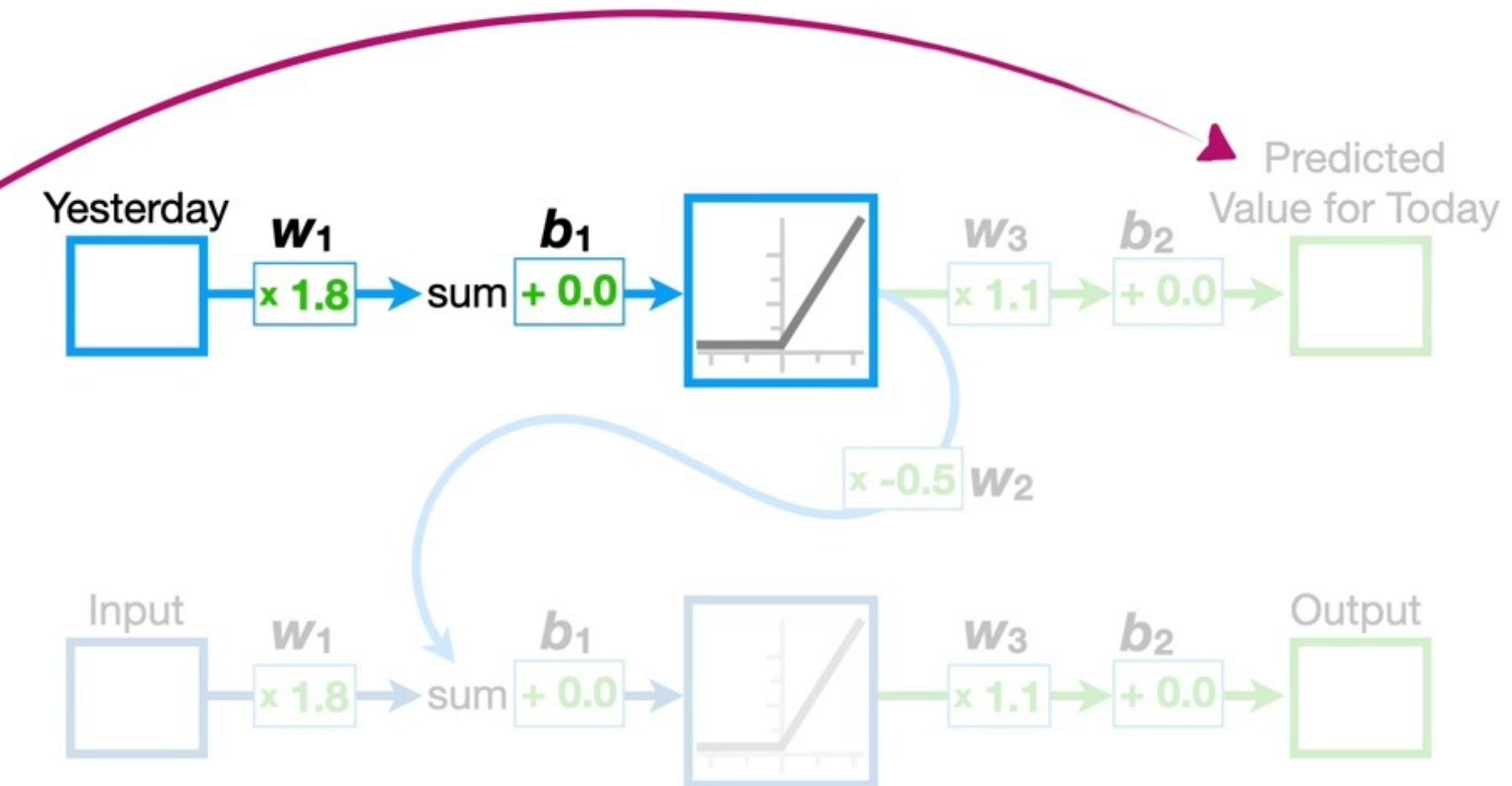


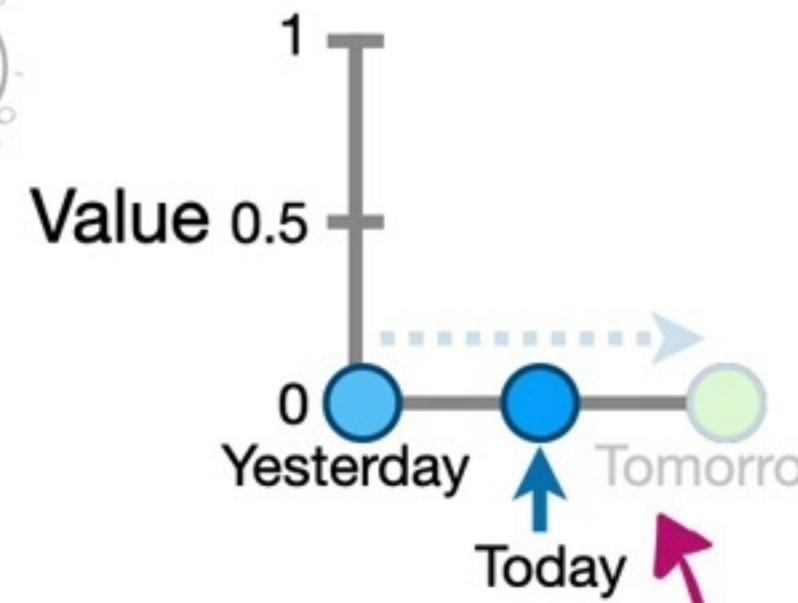
However, as we saw earlier, we can ignore this output.



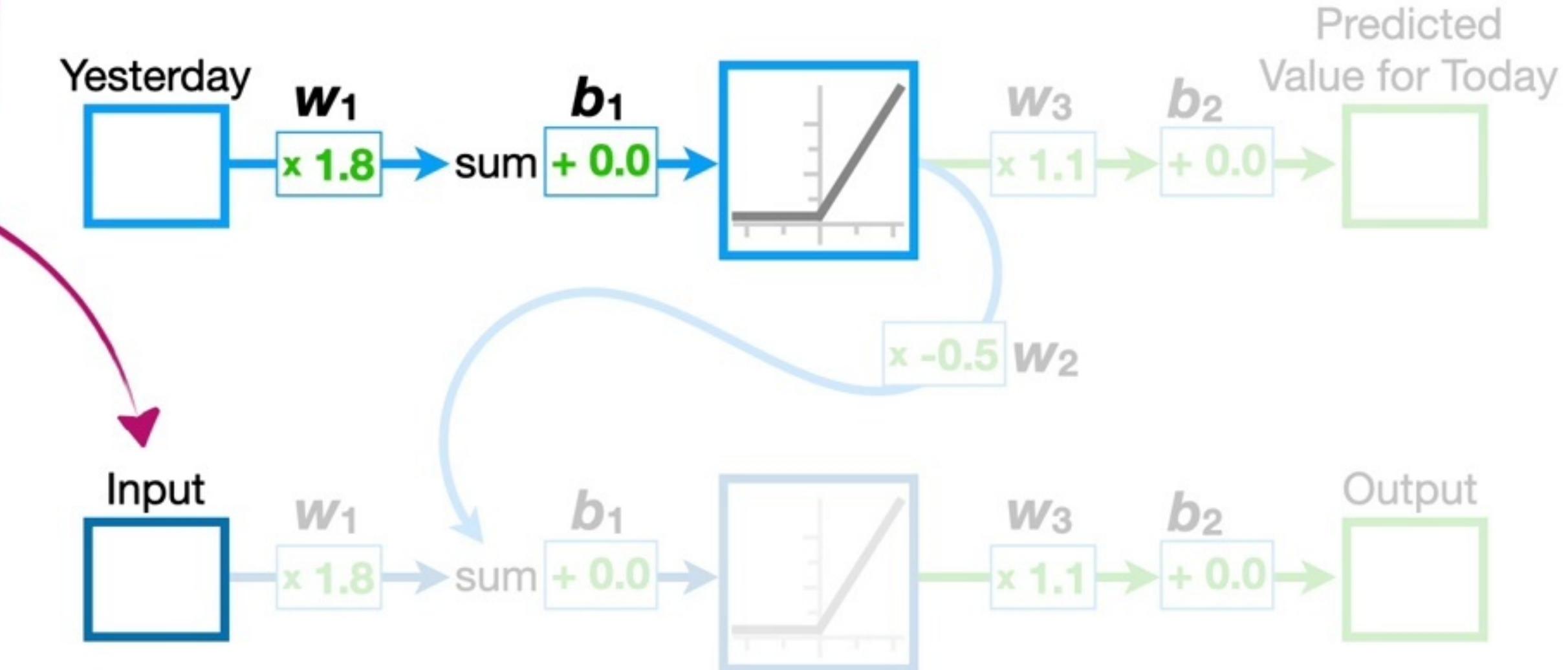


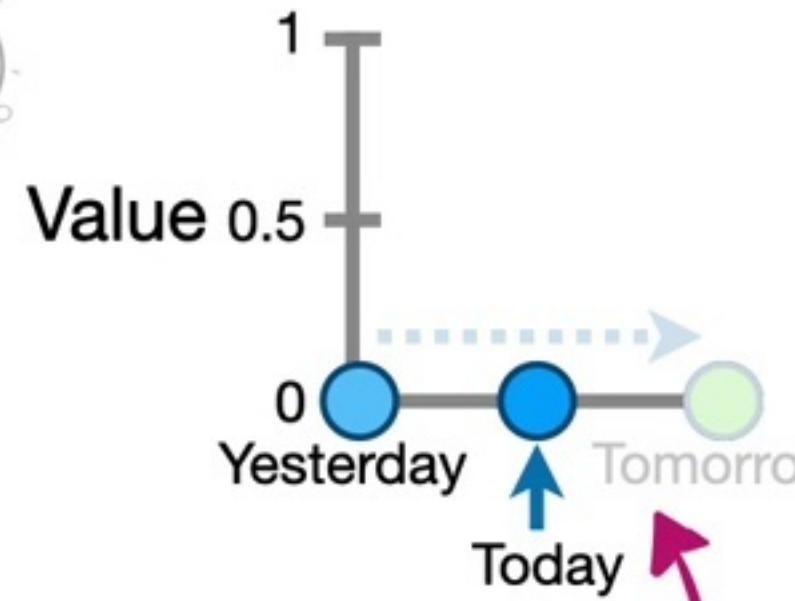
However, as we saw earlier, we can ignore this output.



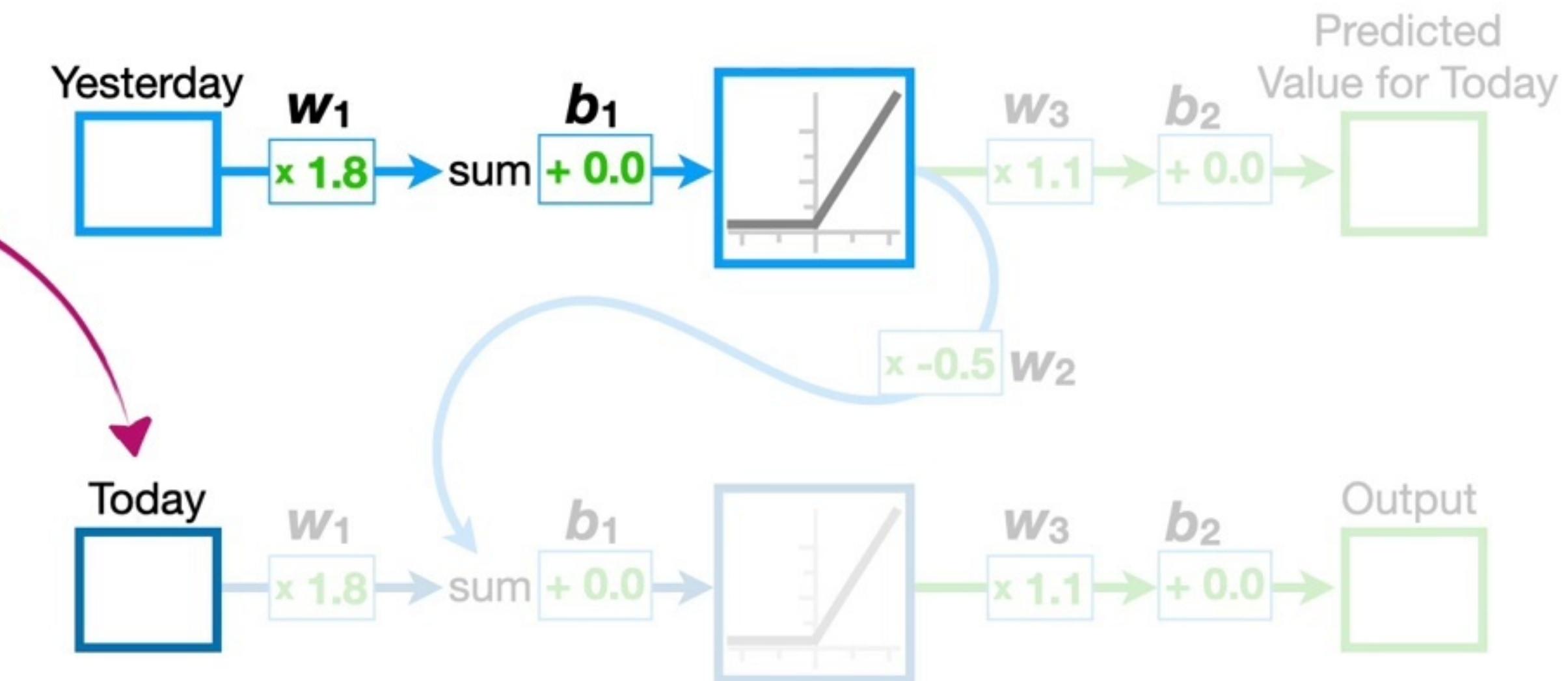


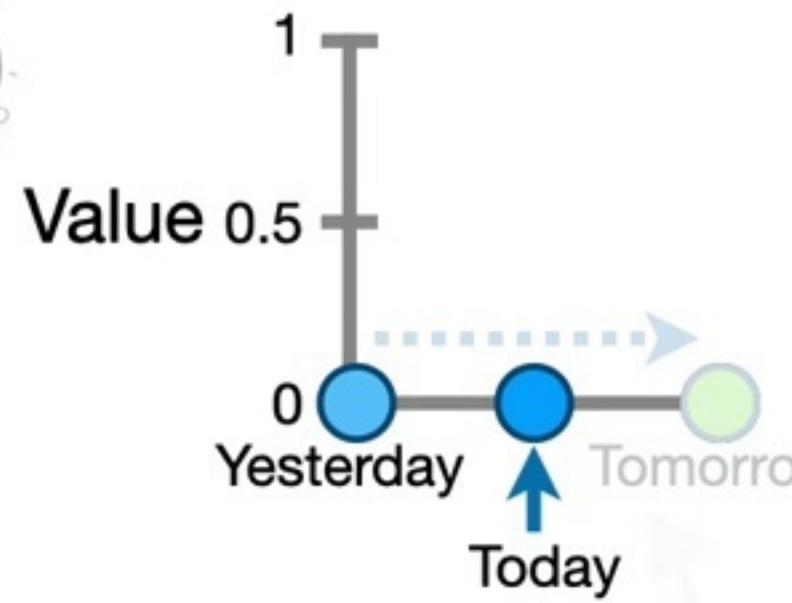
The second
input is for
today's value...



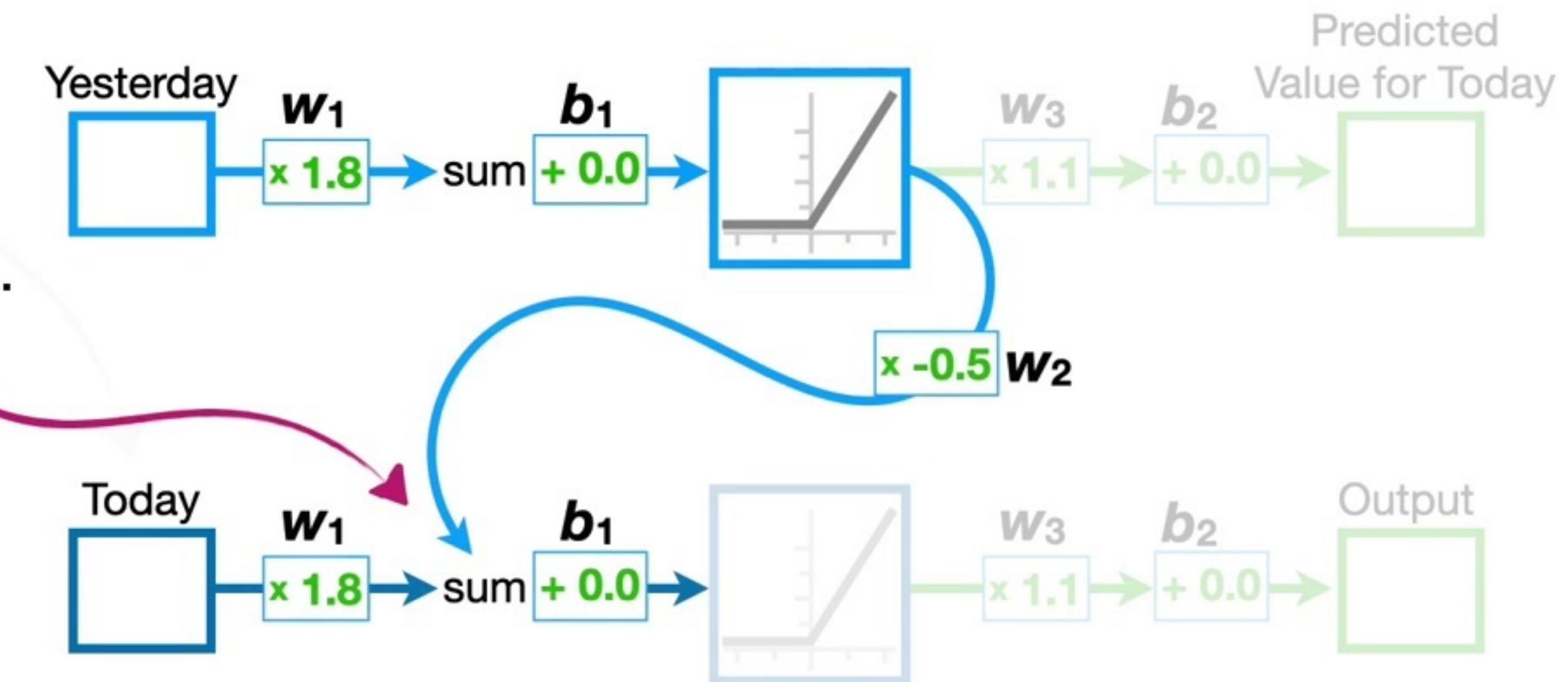


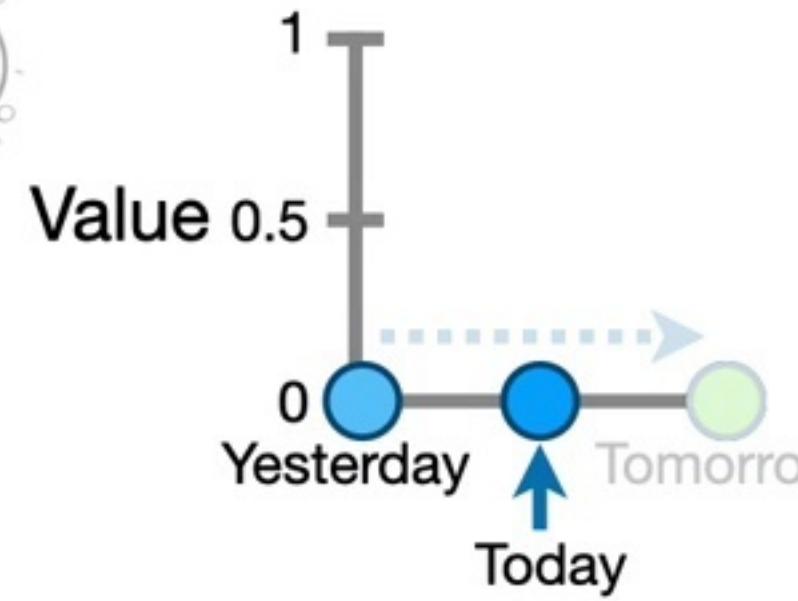
The second
input is for
today's value..



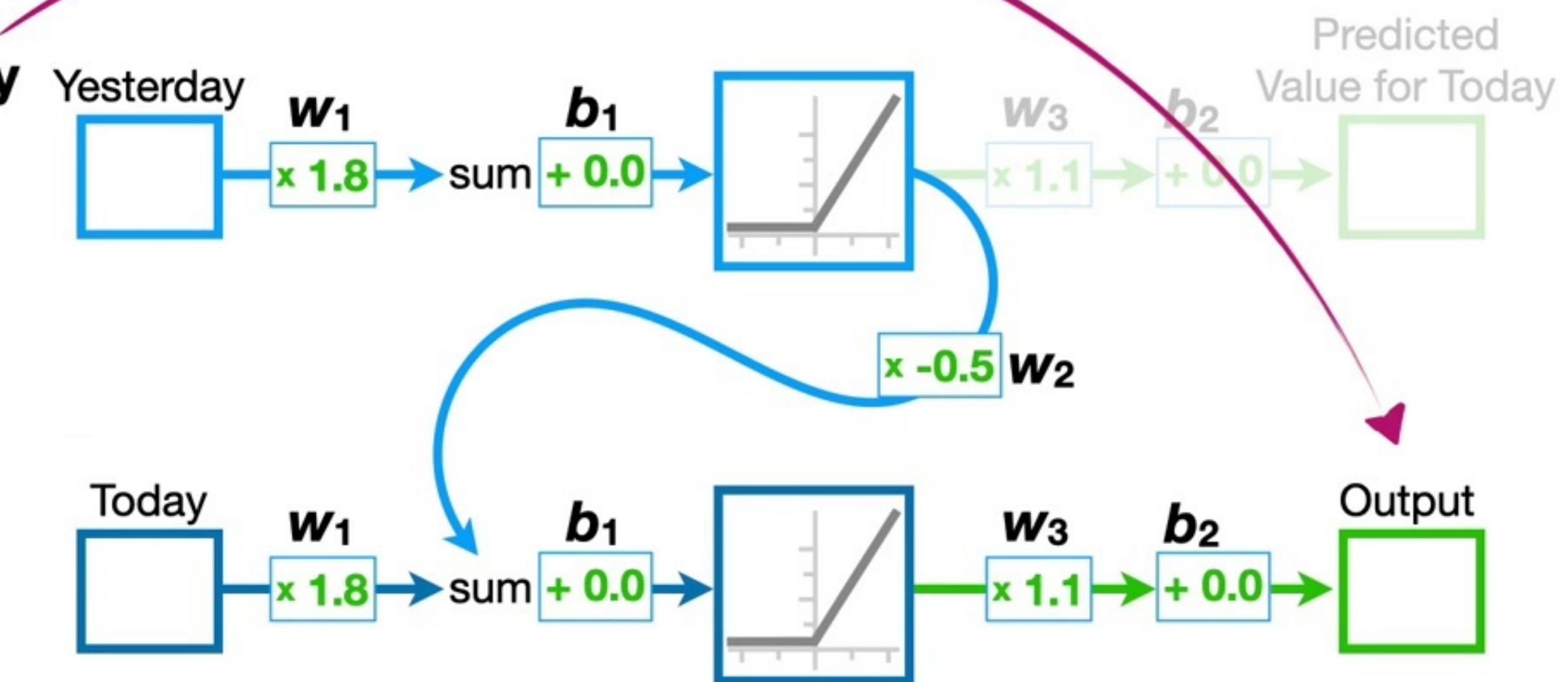


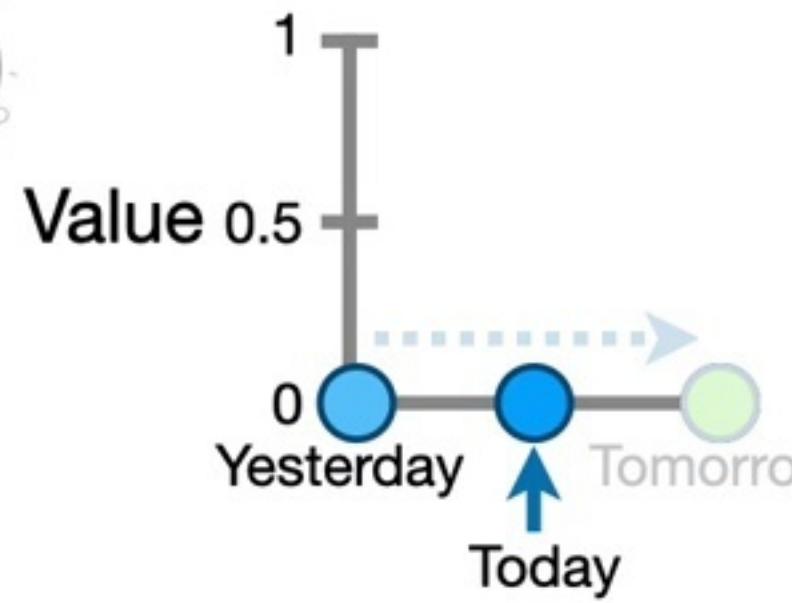
...and the connection
between the first
activation function and
the second summation...



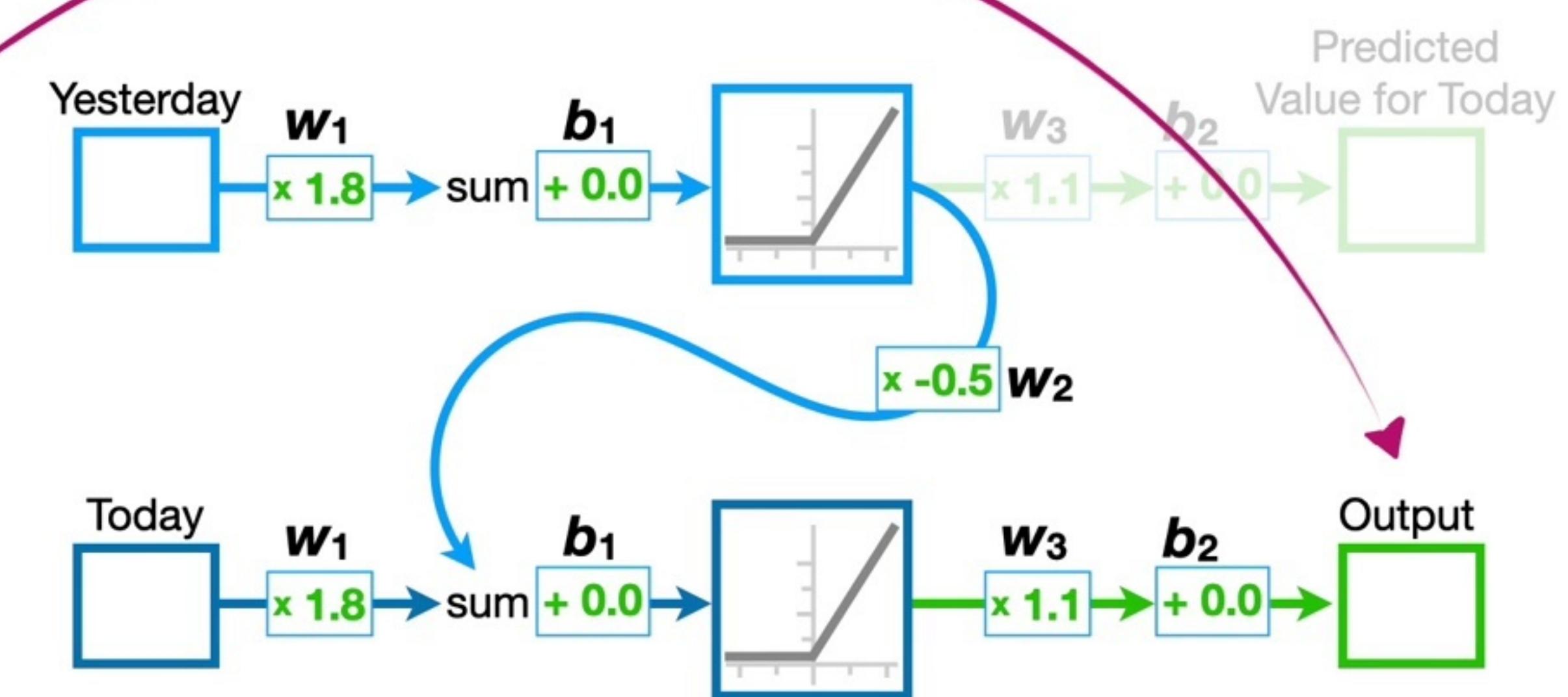


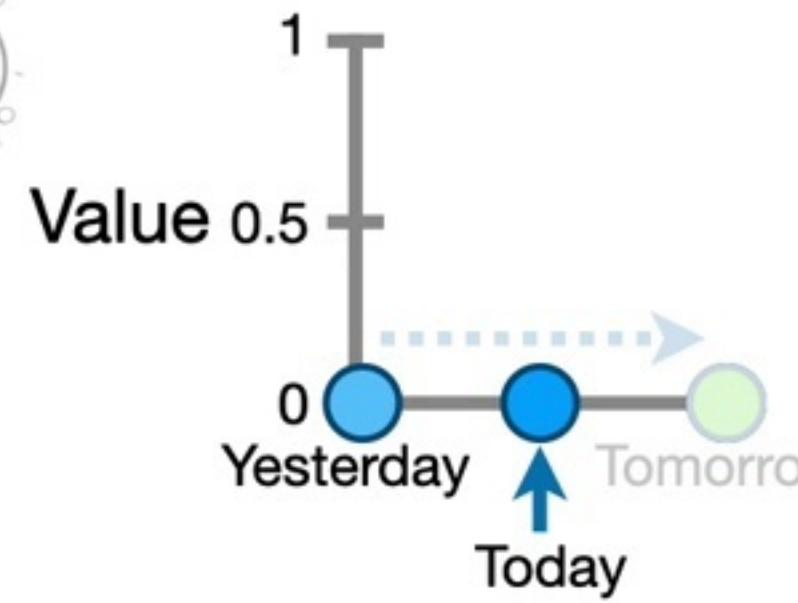
...allows both **yesterday** and **today's** values to influence the final output...



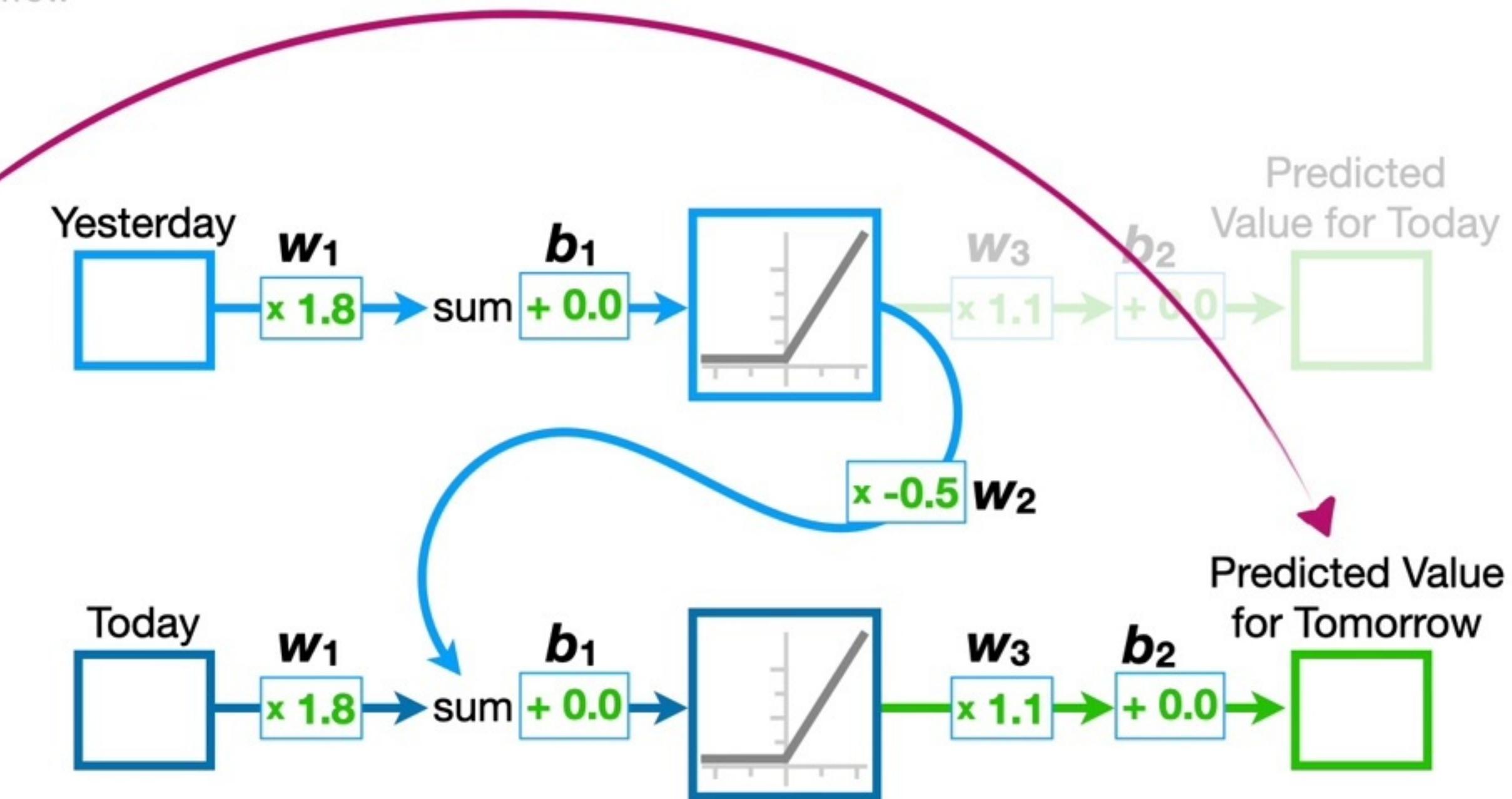


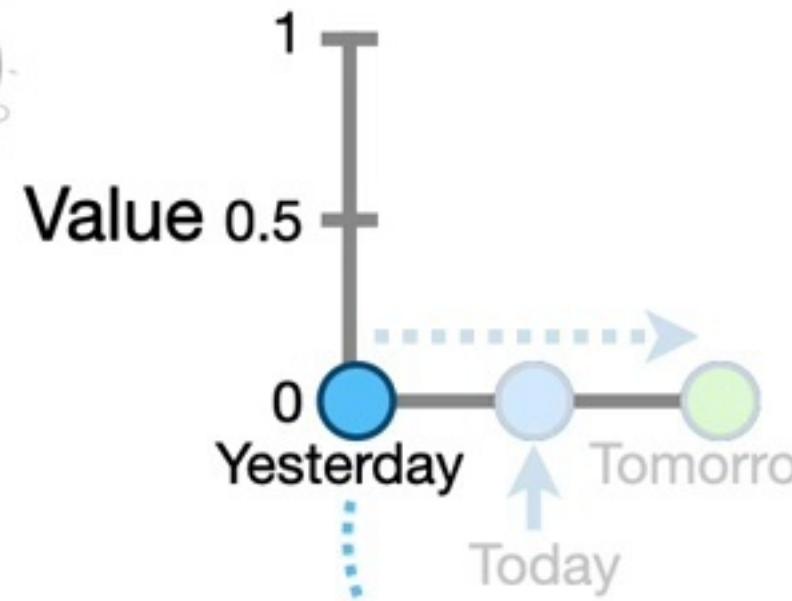
...which gives us the predicted value for tomorrow.



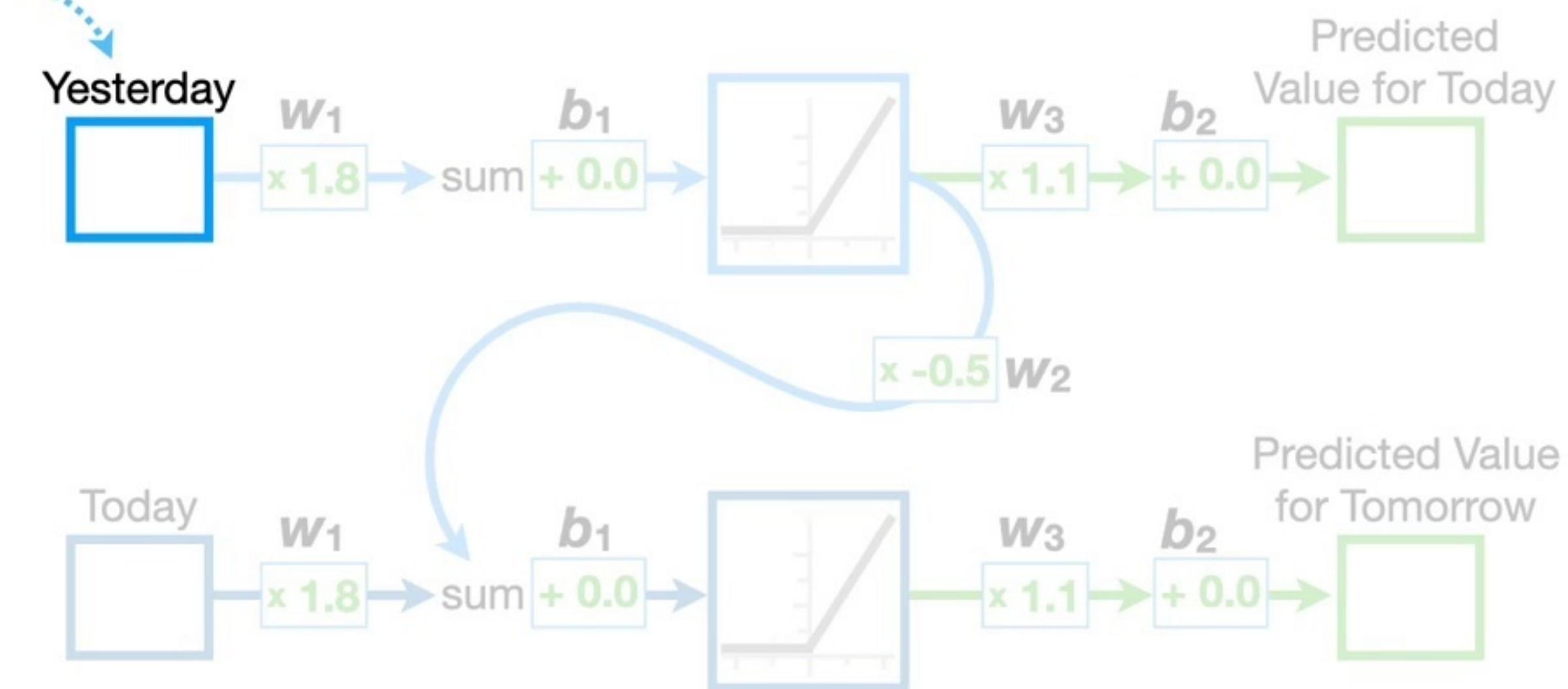


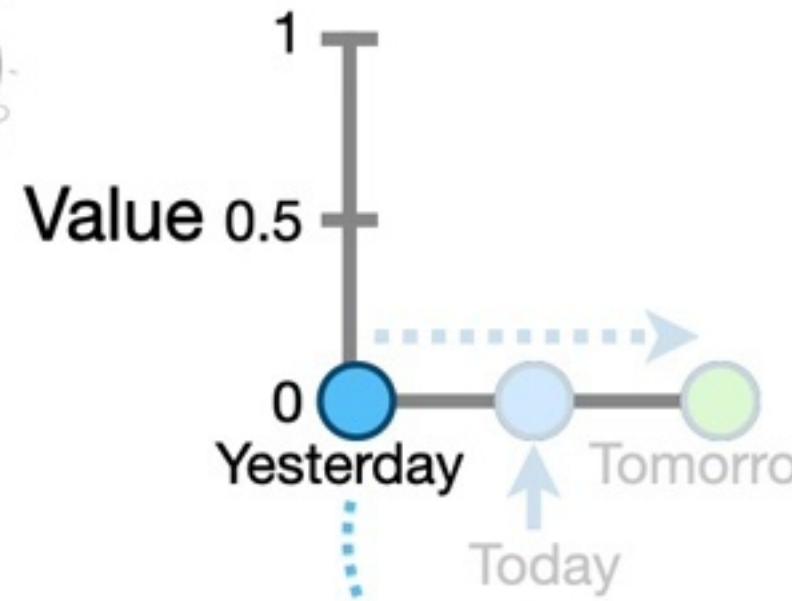
...which gives us the predicted value for tomorrow.





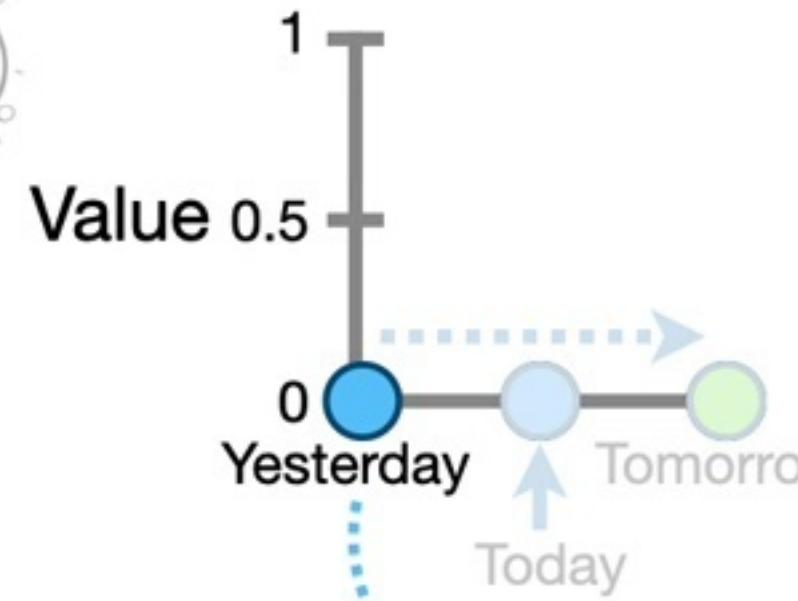
Now, when we put
yesterday's value into
the first input...





Now, when we put
yesterday's value into
the first input...



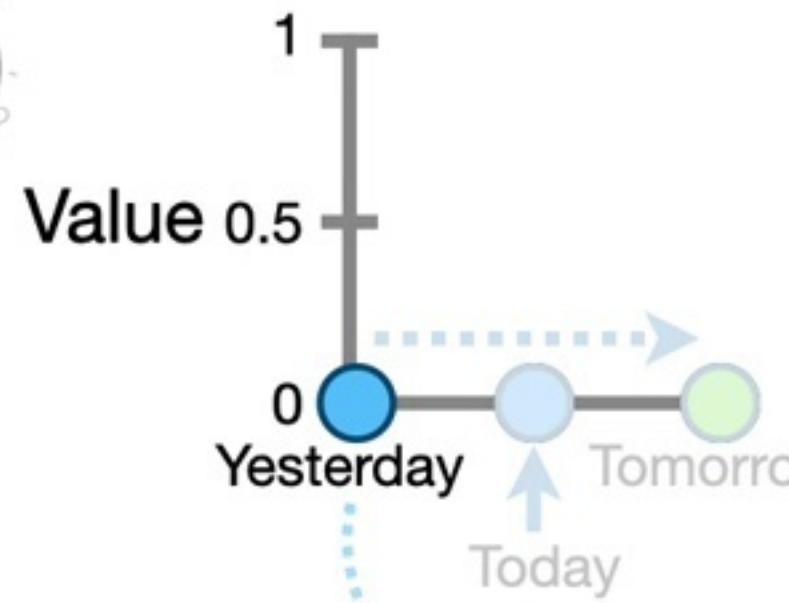


...and we do the math
just like before...



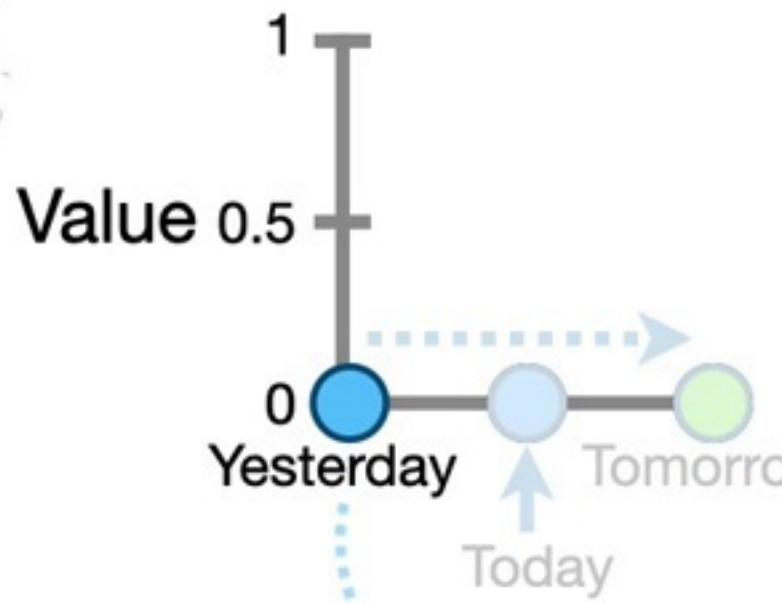


Yesterday $\times w_1$



...and we do the math just like before...





...and we do the math just like before...

Yesterday $\times w_1$

$$0 \times 1.8$$

Yesterday

w₁
0

$$\times 1.8$$

sum

$$+ 0.0$$



w₃

$$\times 1.1$$

+ 0.0

Predicted Value for Today

b₂

$$+ 0.0$$

Predicted Value for Tomorrow

b₂

b₂

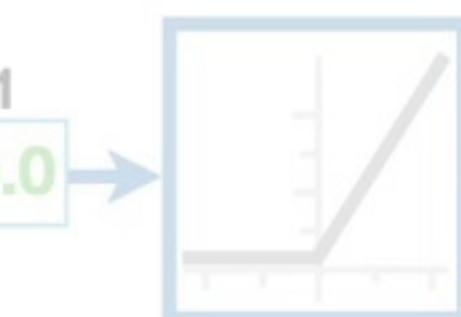
Today

w₁
x 1.8

$$\times 1.8$$

sum

$$+ 0.0$$

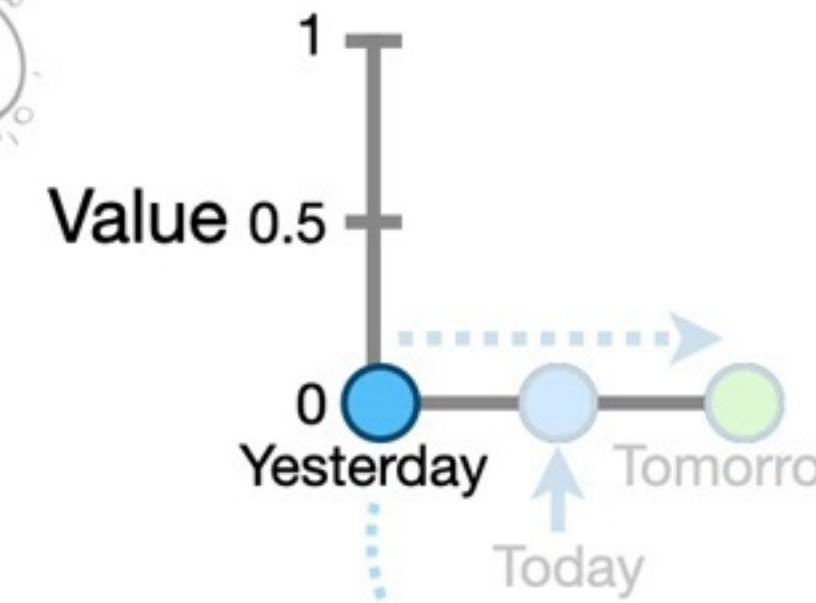


w₃

$$\times 1.1$$

+ 0.0

Predicted Value for Tomorrow



...and we do the math
just like before...

Yesterday $\times w_1 + b_1$

$$0 \times 1.8$$

Yesterday

$$0$$

$$w_1$$

$$\times 1.8$$

sum

$$+ 0.0$$

$$b_1$$

$$+ 0.0$$

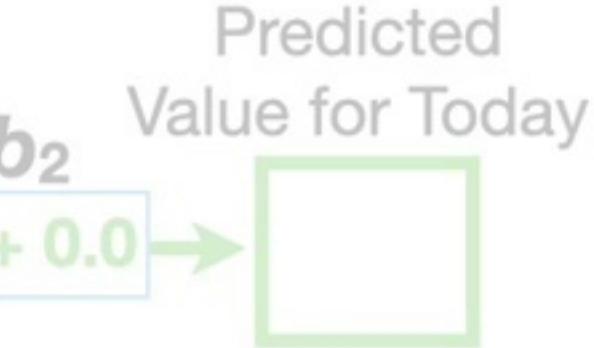


$$w_3$$

$$\times 1.1$$

$$b_2$$

$$+ 0.0$$



Predicted Value
for Tomorrow

Today

$$\square$$

$$w_1$$

$$\times 1.8$$

sum

$$+ 0.0$$

$$b_1$$

$$+ 0.0$$



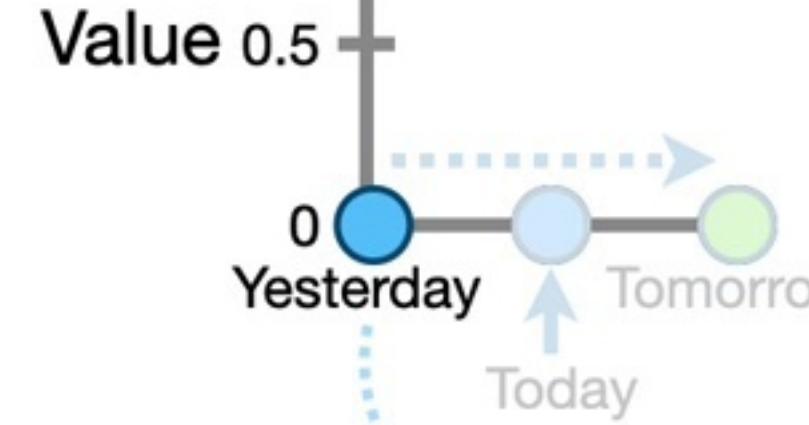
$$w_3$$

$$\times 1.1$$

$$b_2$$

$$+ 0.0$$





Yesterday $\times w_1 + b_1$

$$0 \times 1.8 + 0.0$$

...and we do the math
just like before...

Yesterday

$$0 \rightarrow w_1$$

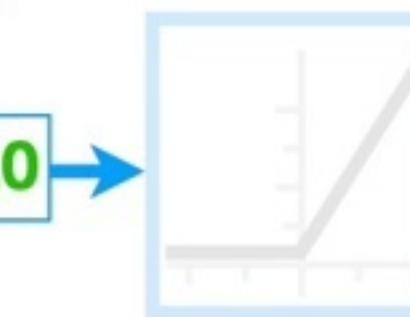
$\times 1.8$

sum

$$+ 0.0$$

b_1

$+ 0.0$

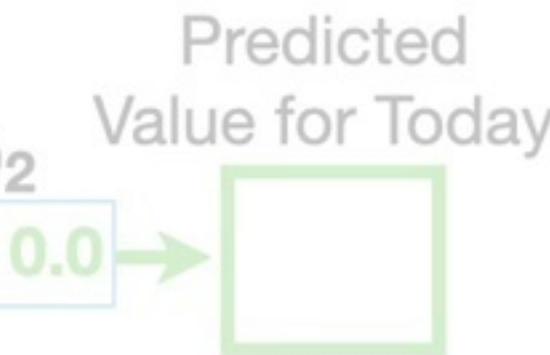


w_3

$$\times 1.1$$

b_2

$$+ 0.0$$



Today

$$\rightarrow w_1$$

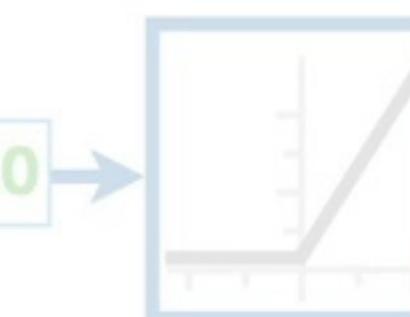
$$\times 1.8$$

sum

$$+ 0.0$$

b_1

$$+ 0.0$$



w_3

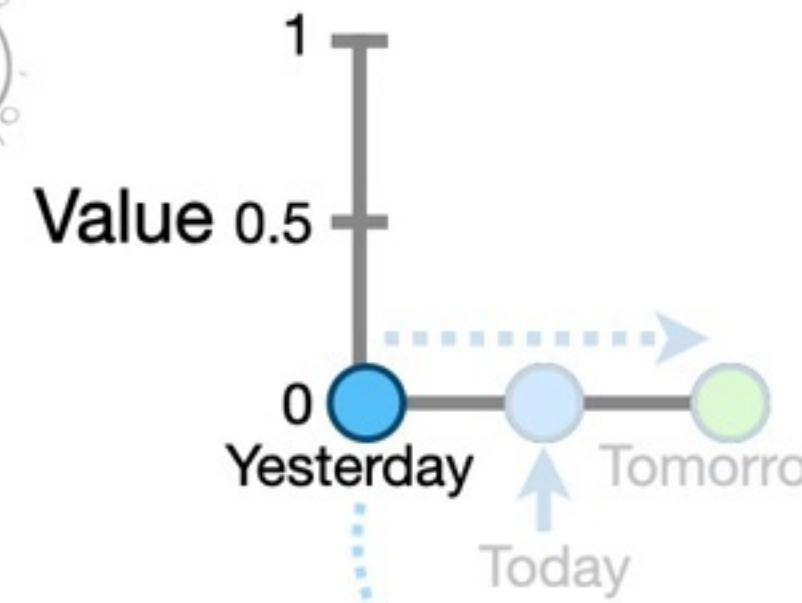
$$\times 1.1$$

b_2

$$+ 0.0$$



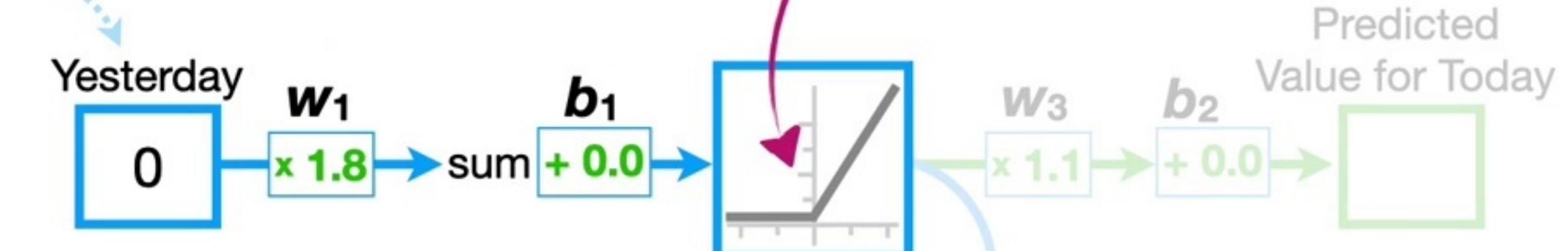
$\times -0.5$ w_2

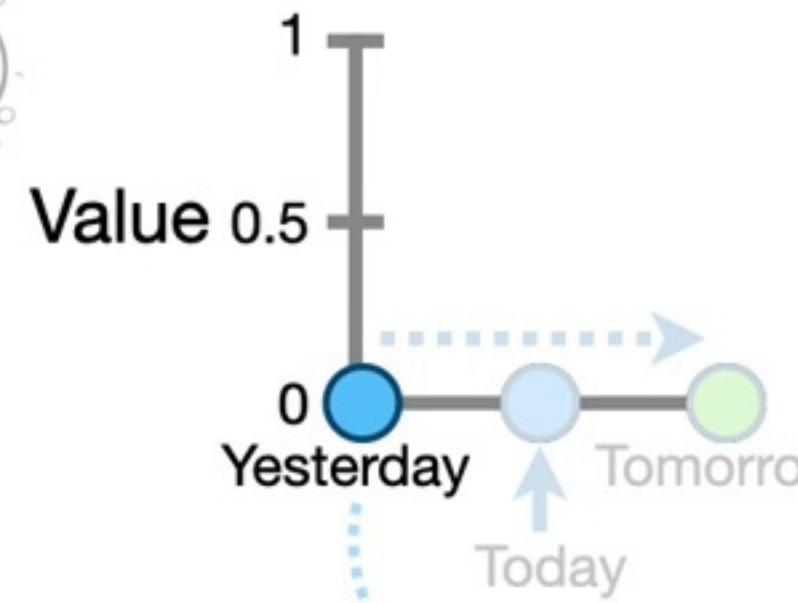


...and we do the math
just like before...

Yesterday $\times w_1 + b_1$ = x-axis coordinate

$$0 \times 1.8 + 0.0$$

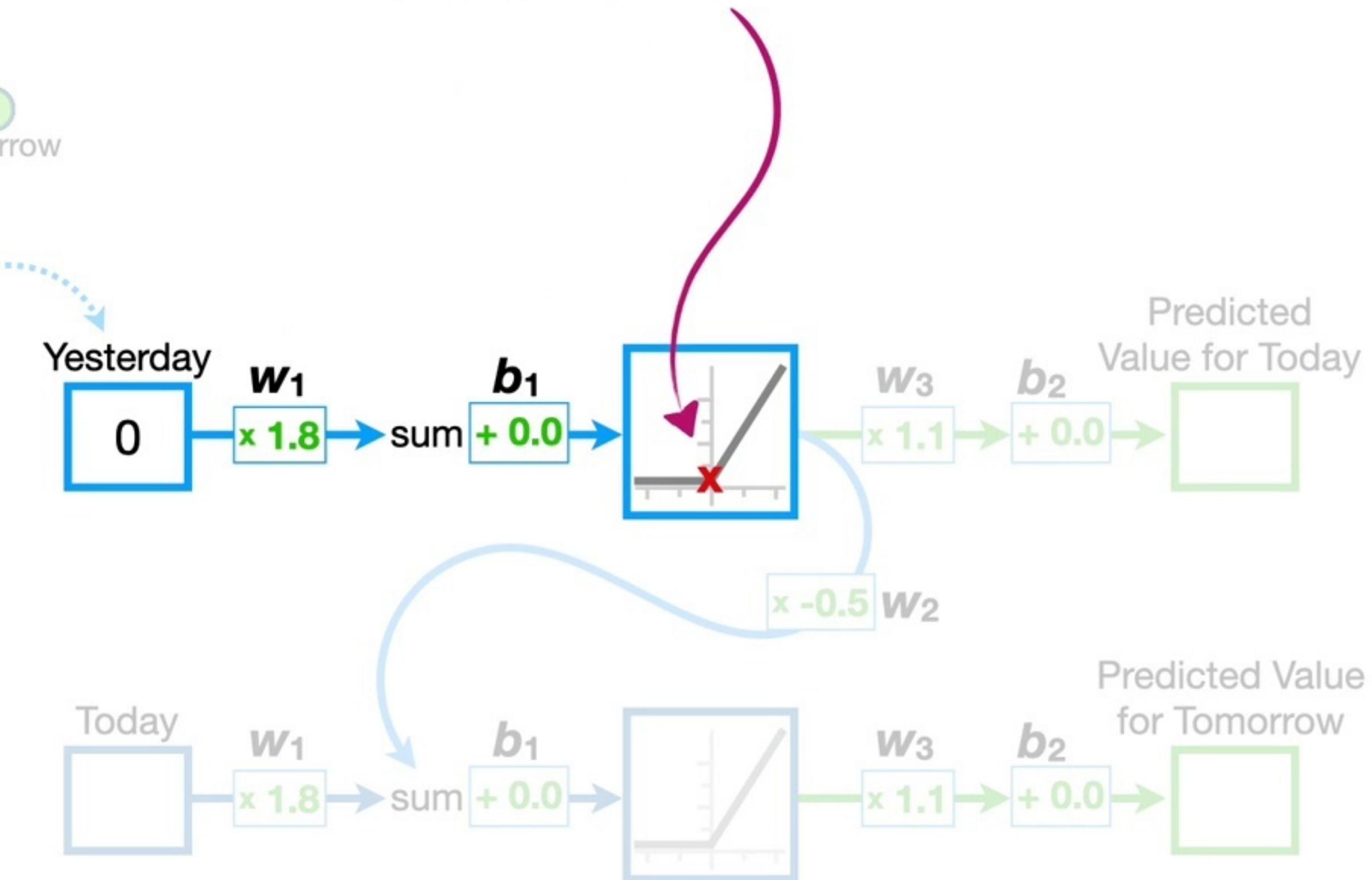


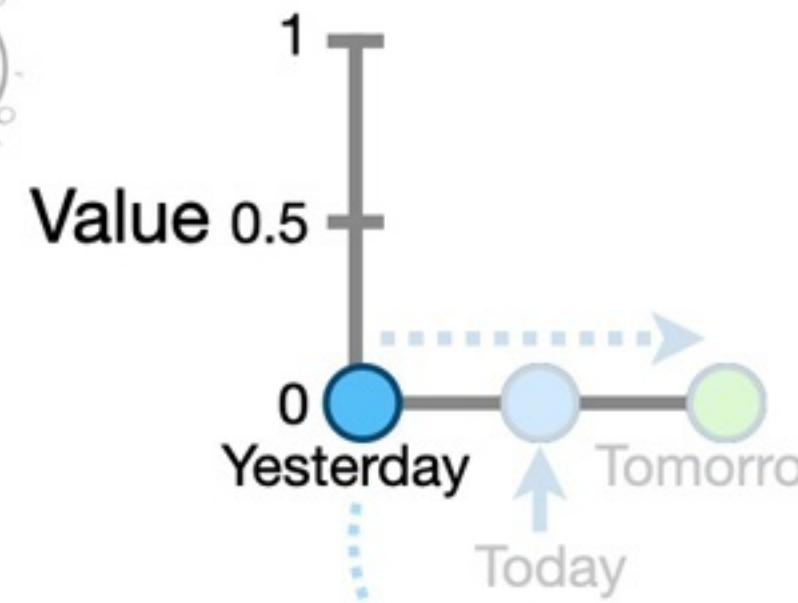


...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$



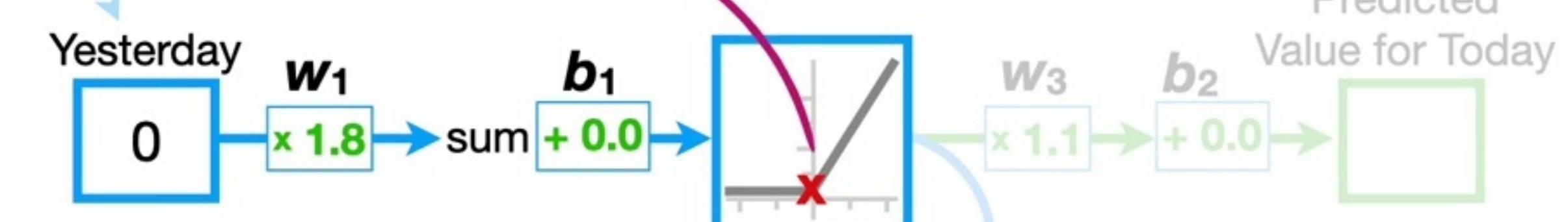


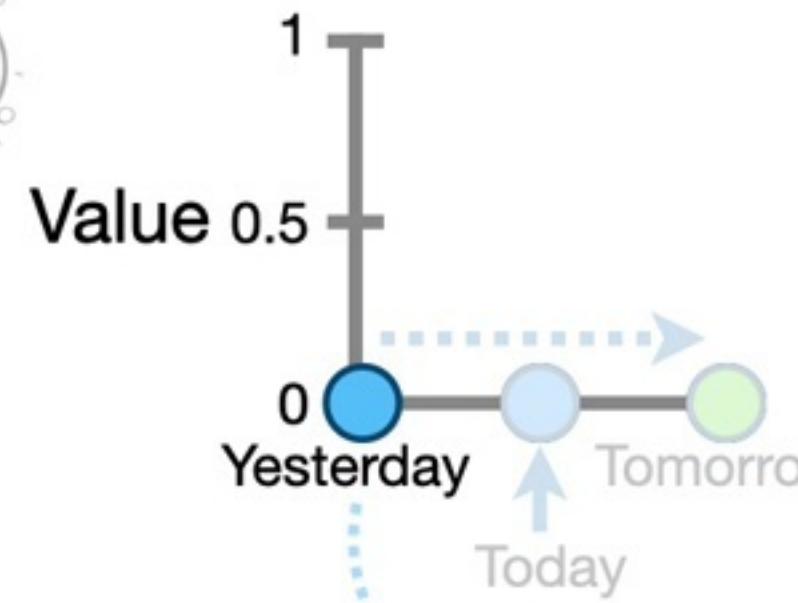
...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = x\text{-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$$f(x) = \max(0, x) = y\text{-axis coordinate}$$



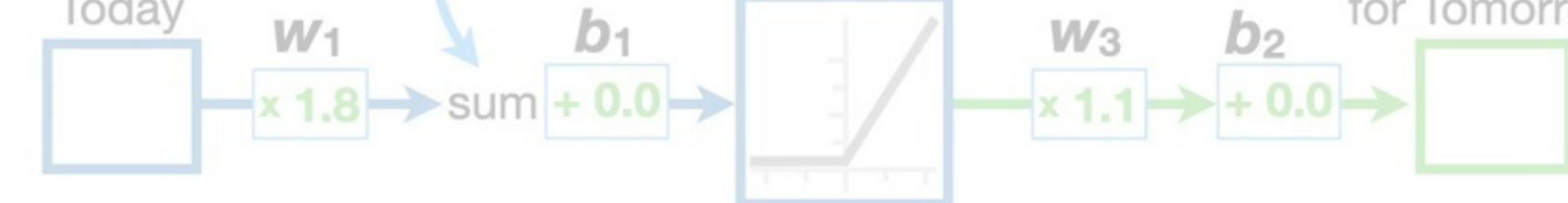


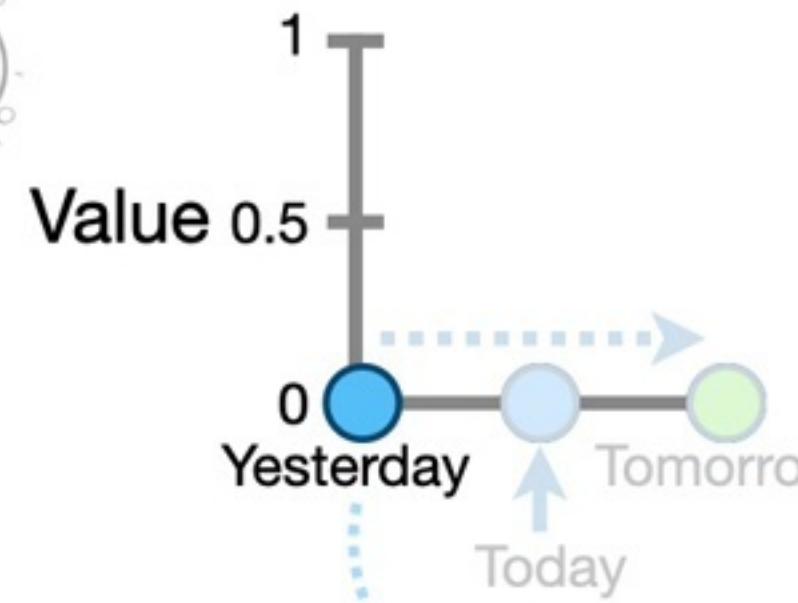
...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = x\text{-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(x) = \max(0, x) = y\text{-axis coordinate}$



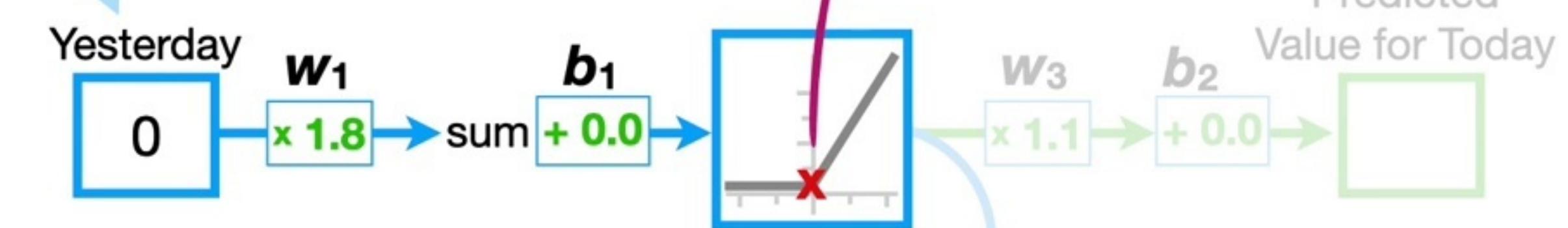


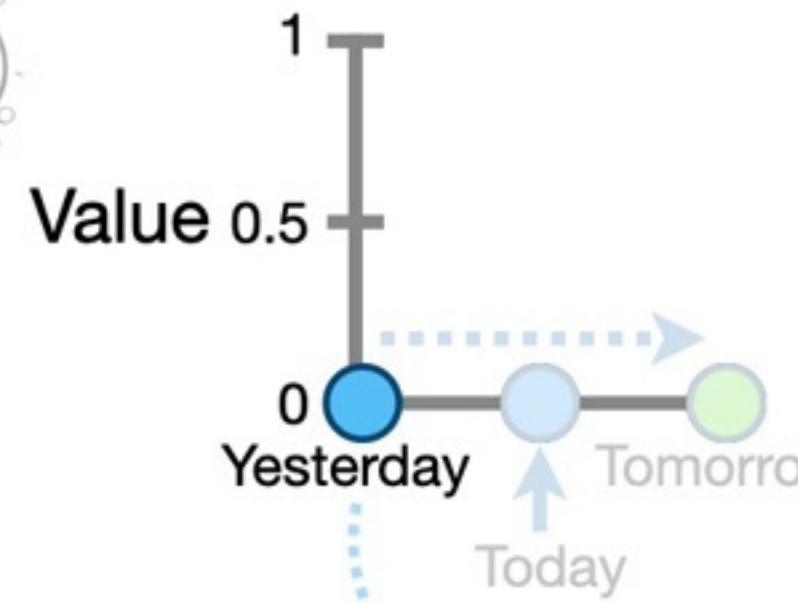
...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = x\text{-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(0) = \max(0, 0) = y\text{-axis coordinate}$



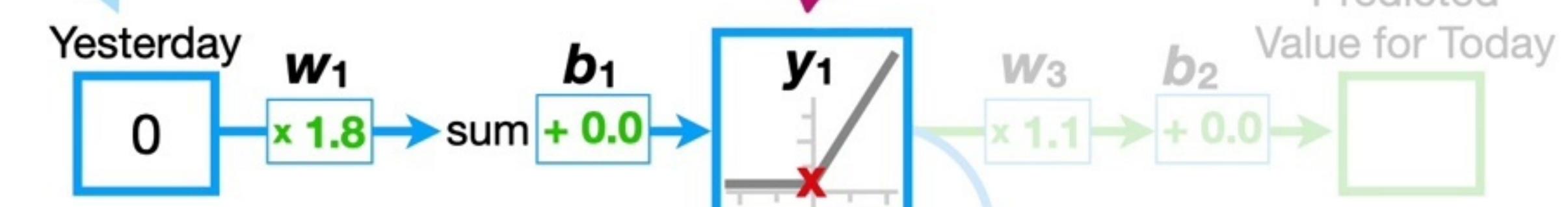


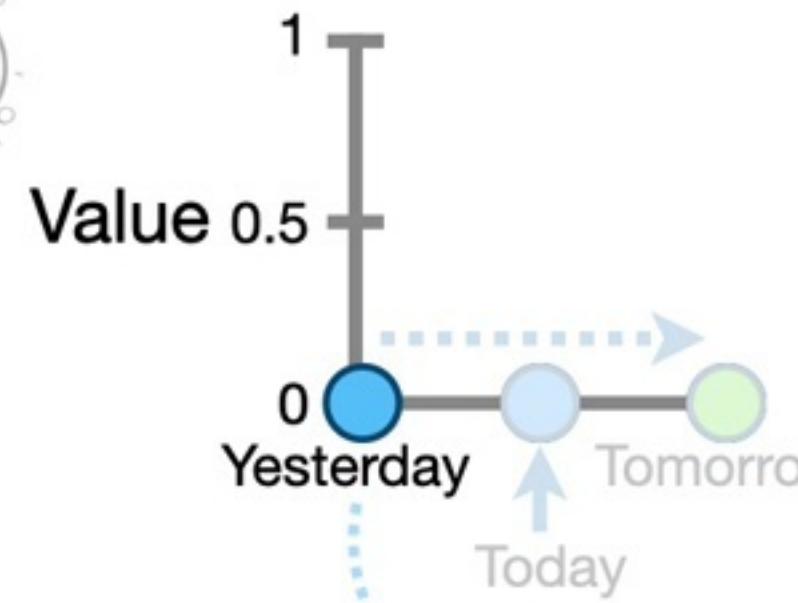
...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = x\text{-axis coordinate}$

$$0 \times 1.8 + 0.0 = 0$$

$f(0) = \max(0, 0) = y\text{-axis coordinate}$



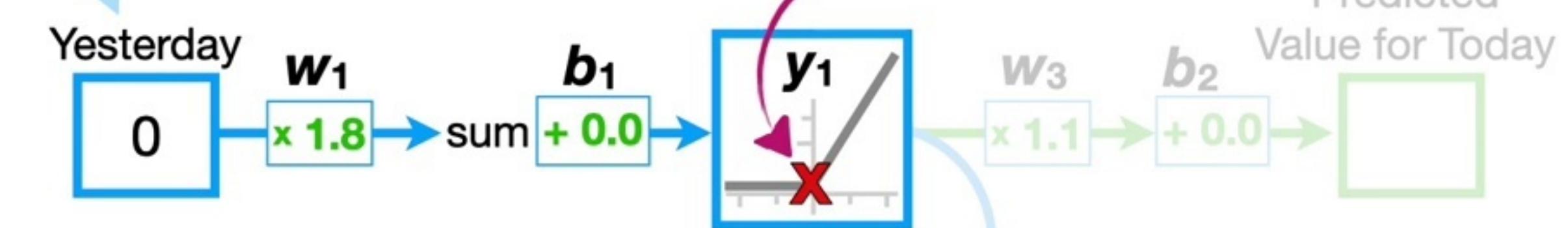


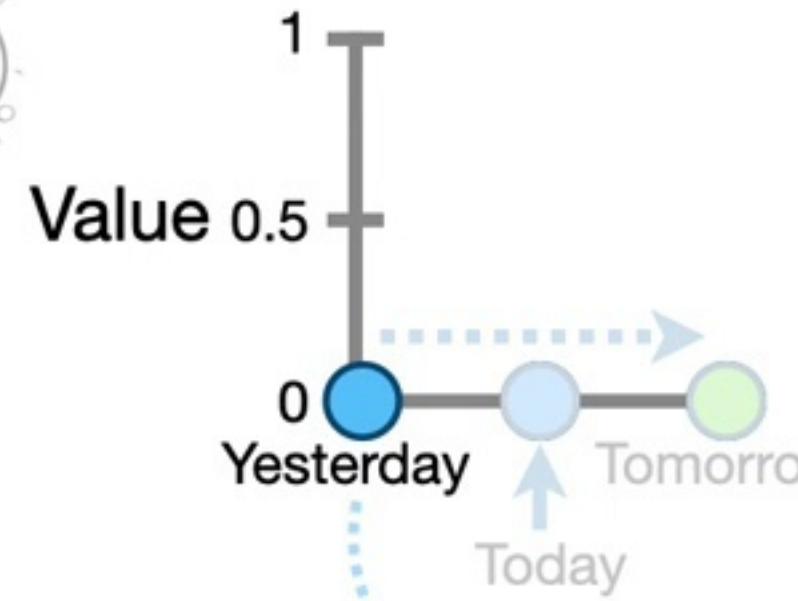
...and we do the math
just like before...

Yesterday $\times w_1 + b_1 = x\text{-axis coordinate}$

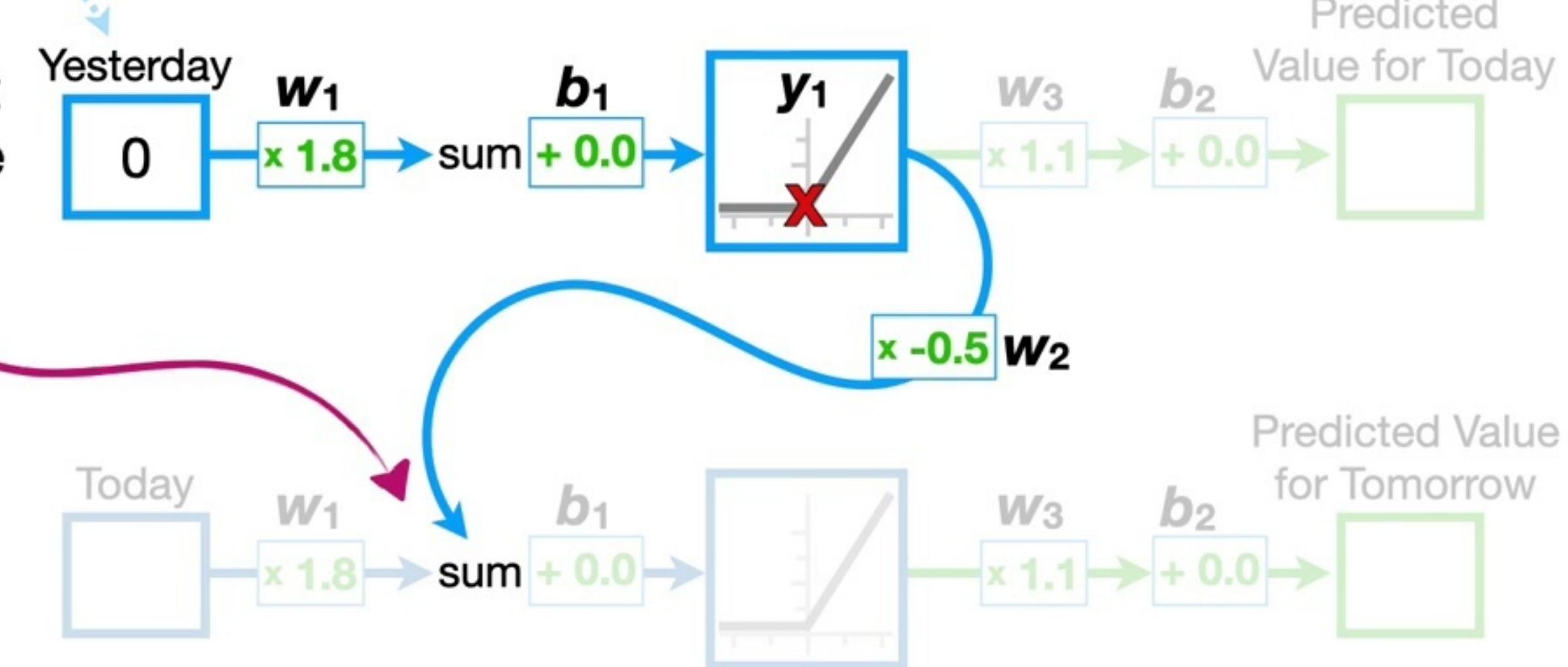
$$0 \times 1.8 + 0.0 = 0$$

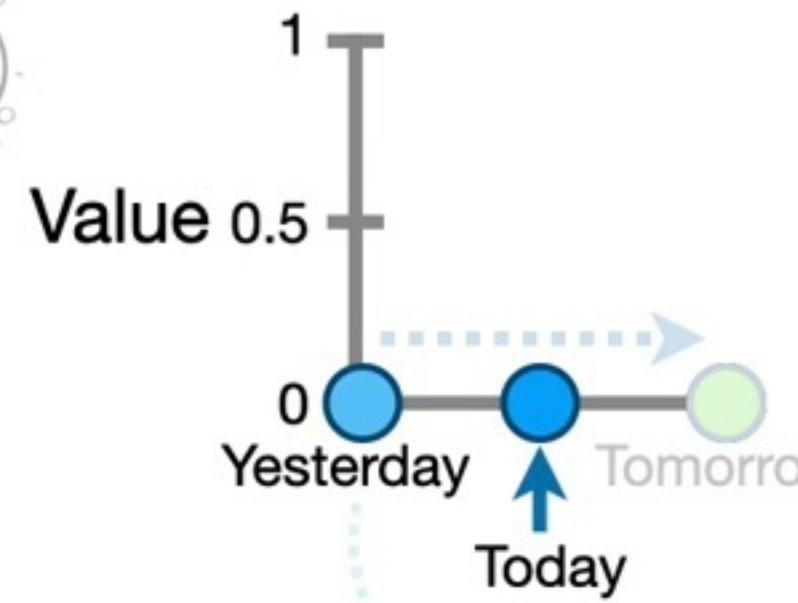
$$f(0) = \max(0, 0) = 0$$



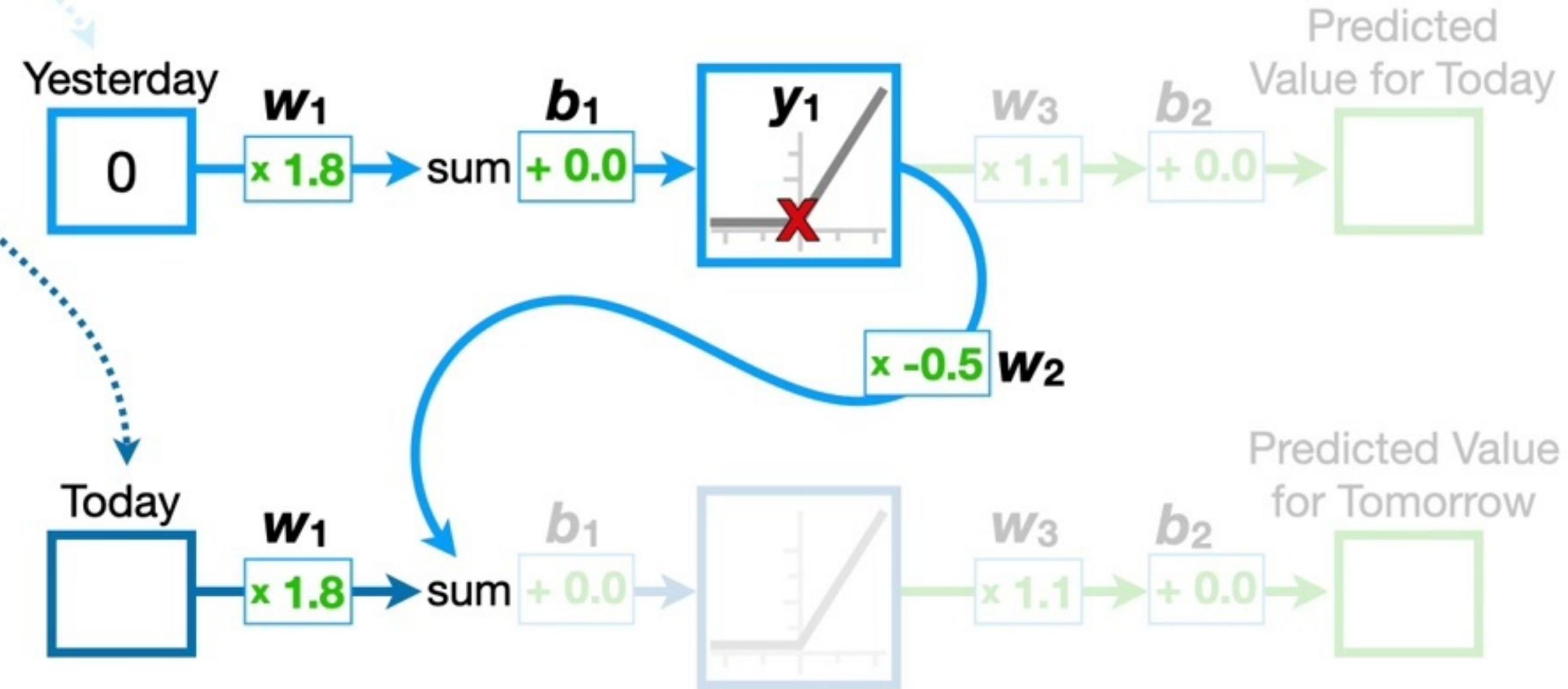


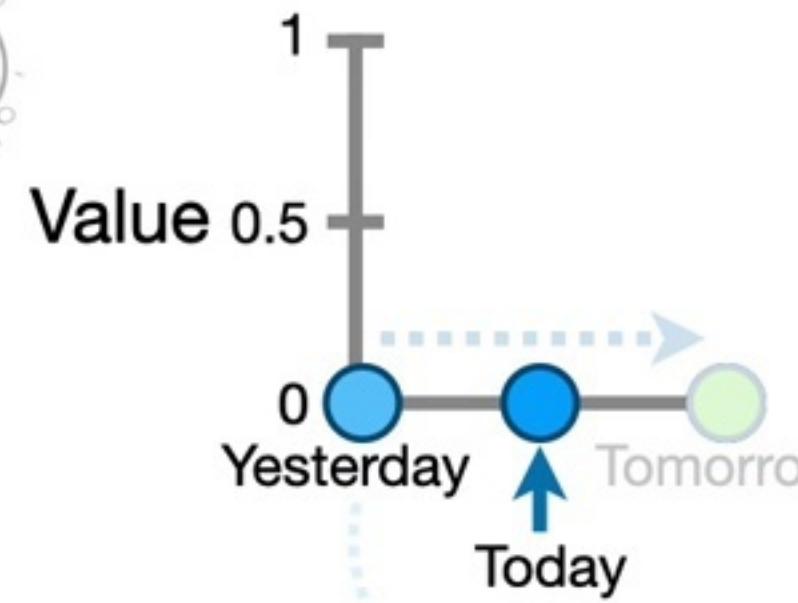
...then we follow the connection from the first activation function to the summation in the second copy of the neural network.



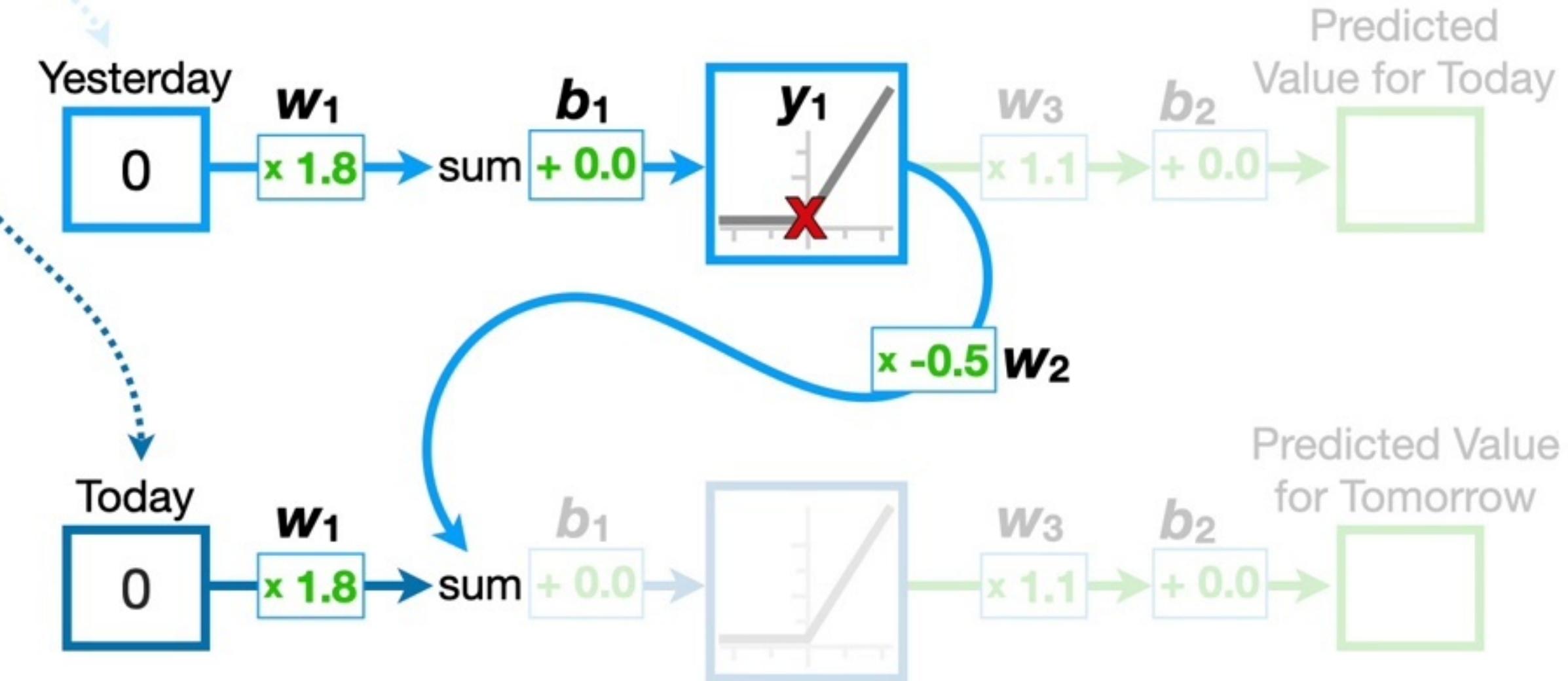


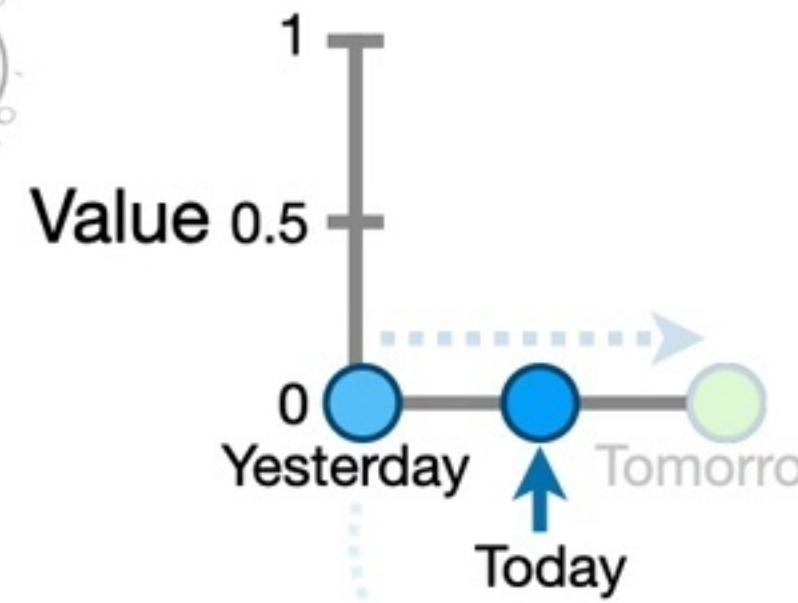
Now we put
today's value
into the second
input...



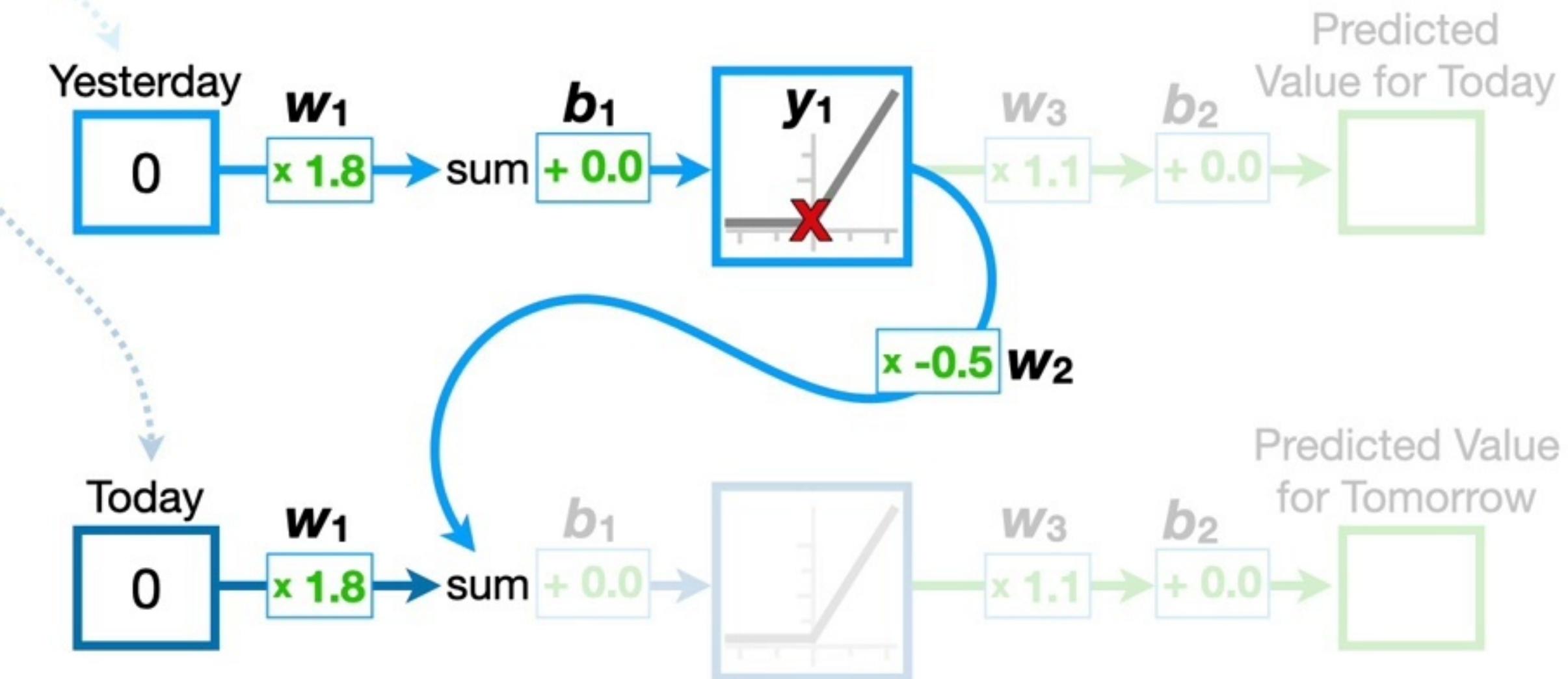


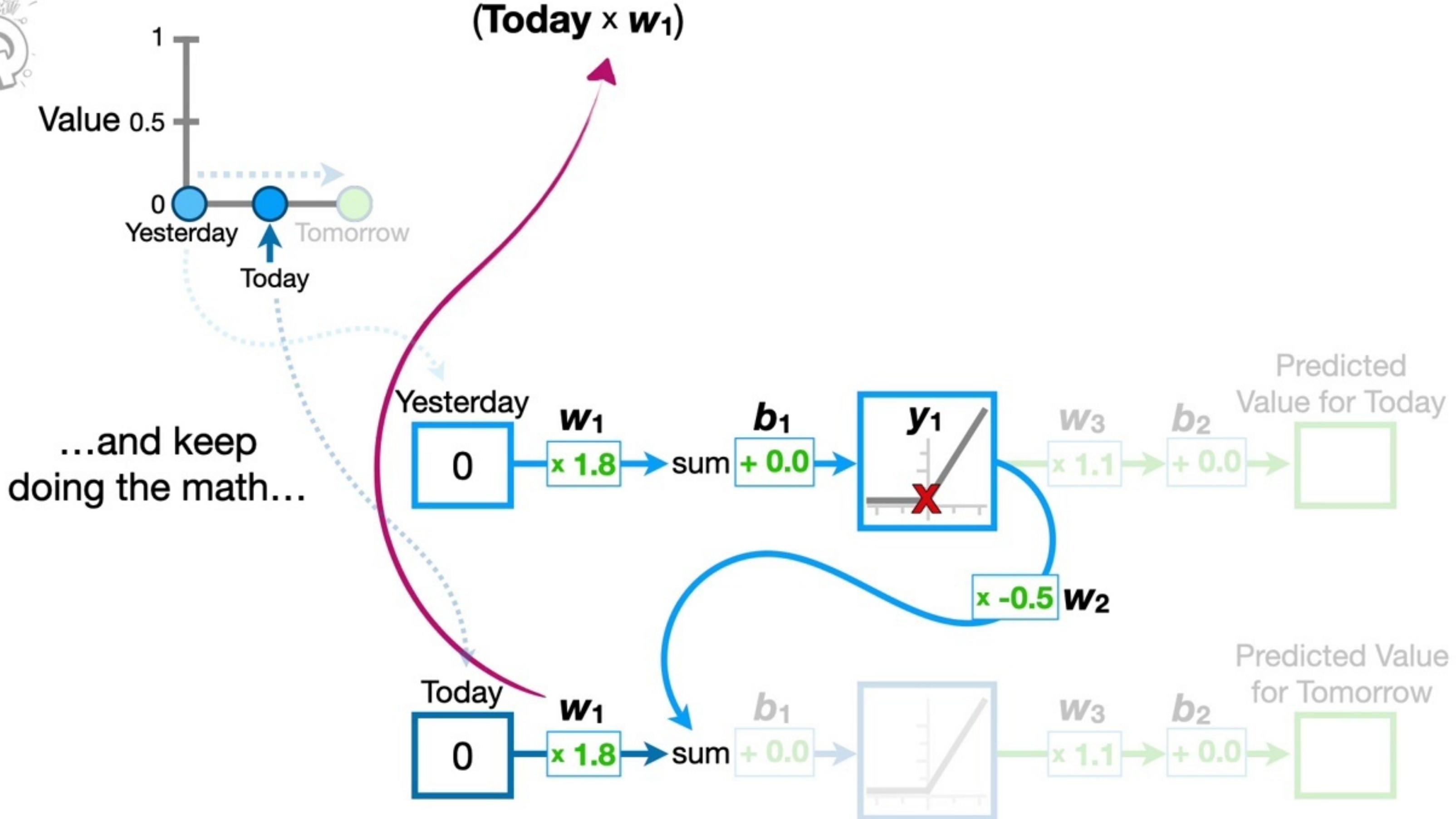
Now we put
today's value
into the second
input...

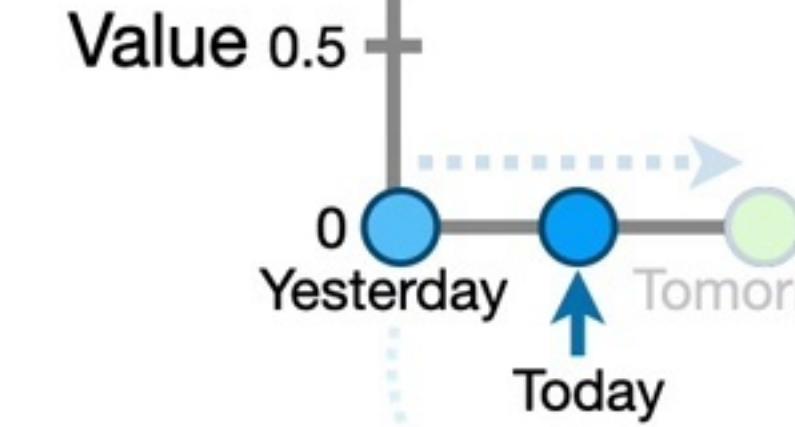




...and keep
doing the math...

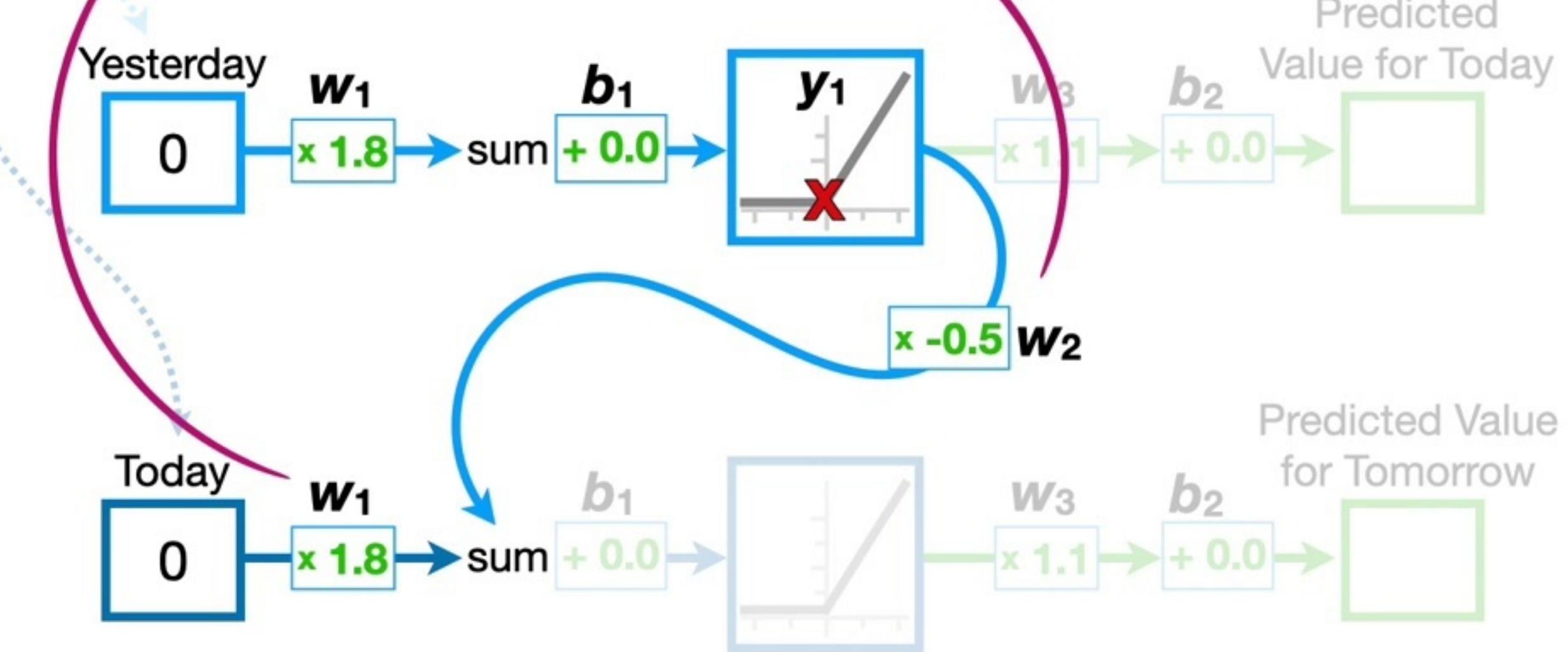






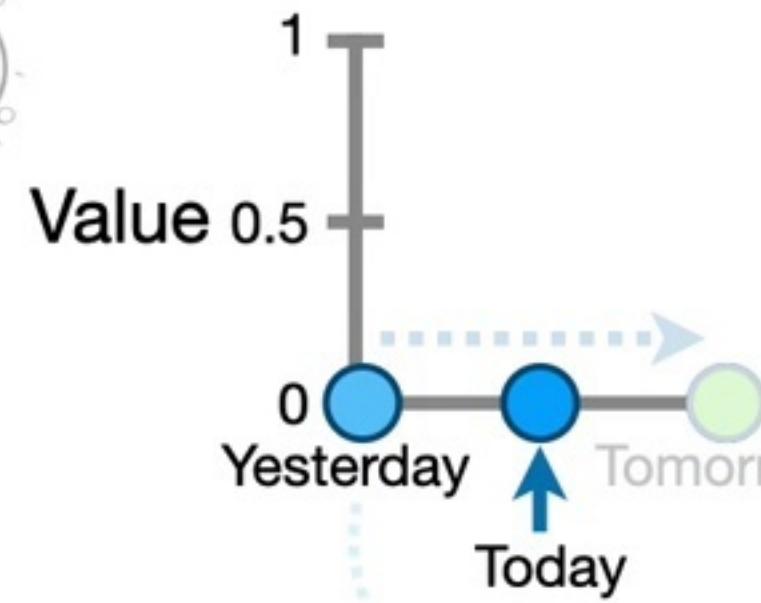
$$(\text{Today} \times w_1) + (y_1 \times w_2)$$

...and keep
doing the math...





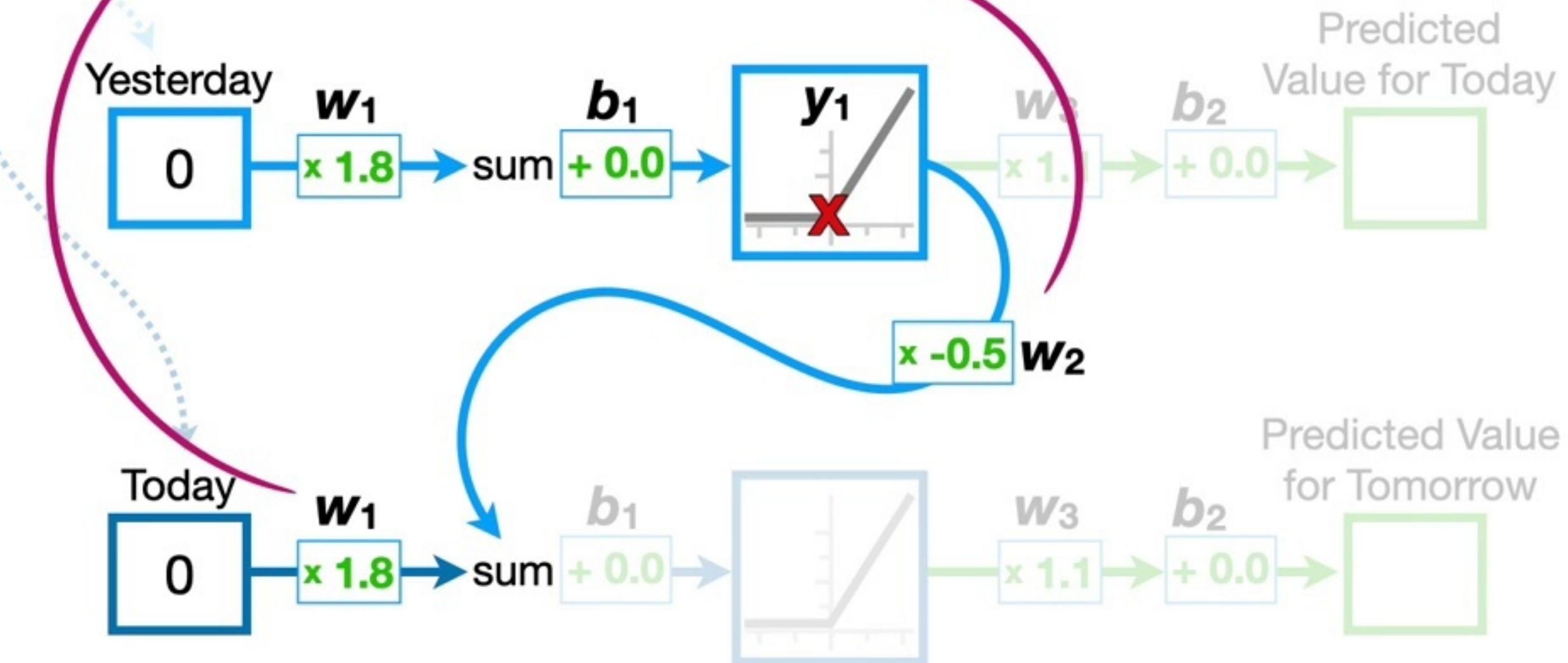
$$(\text{Today} \times w_1) + (y_1 \times w_2)$$

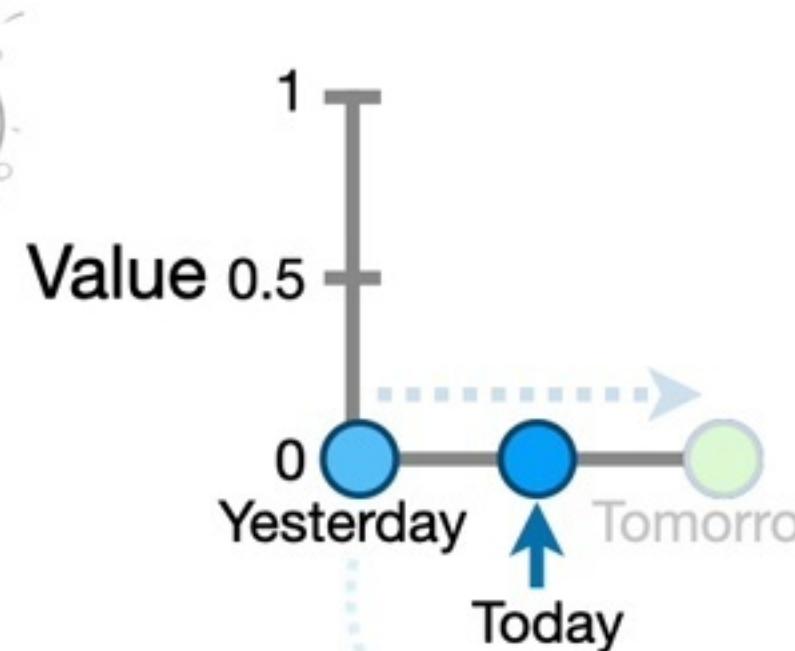


...and keep
doing the math.

$$0 \times 1.8 + (0 \times -0.5)$$

$$0 \times 1.8 + (0 \times -0.5)$$

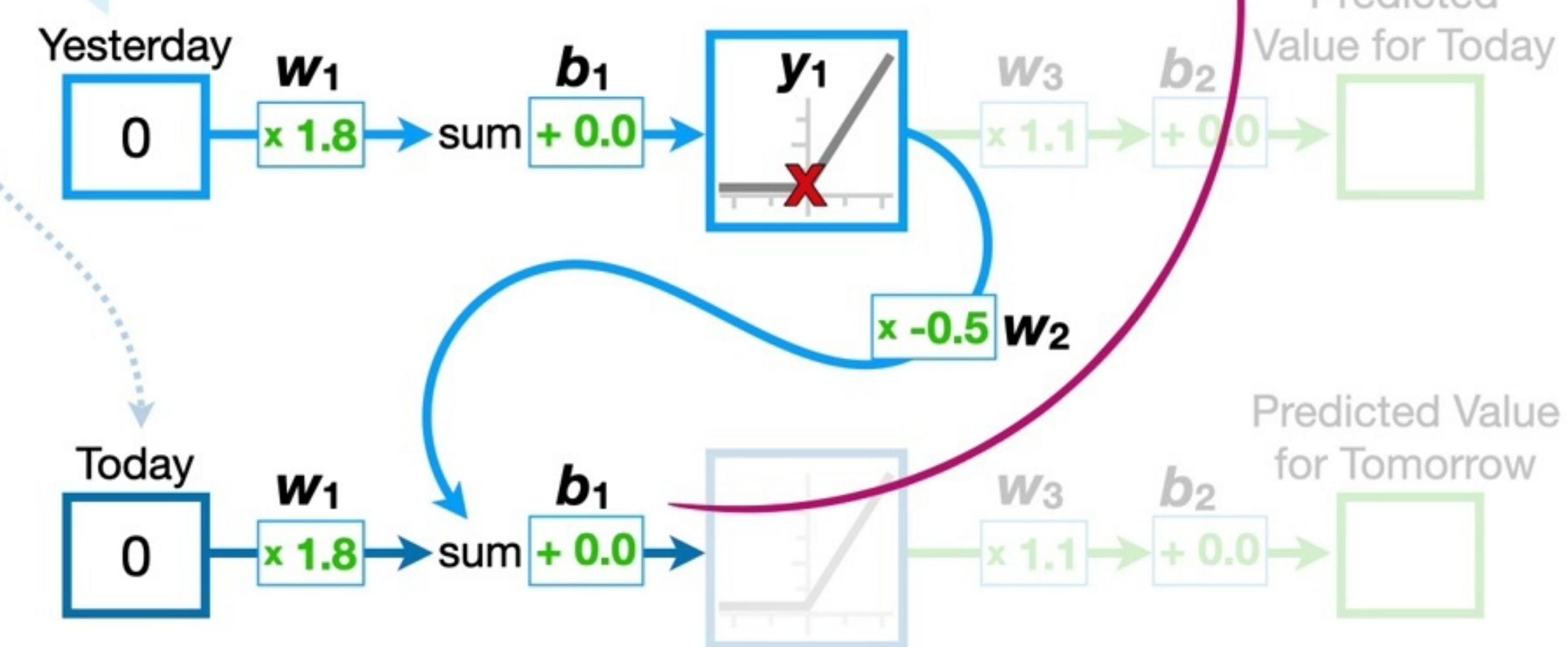


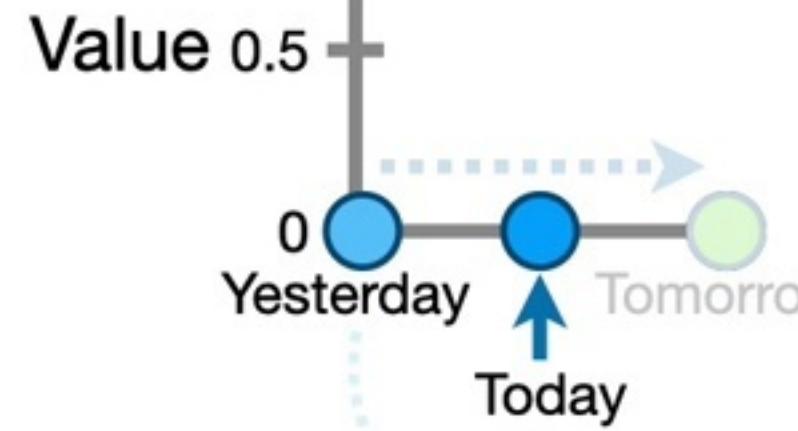


...and keep
doing the math..

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b$$

$$0 \times 1.8 + (0 \times -0.5)$$

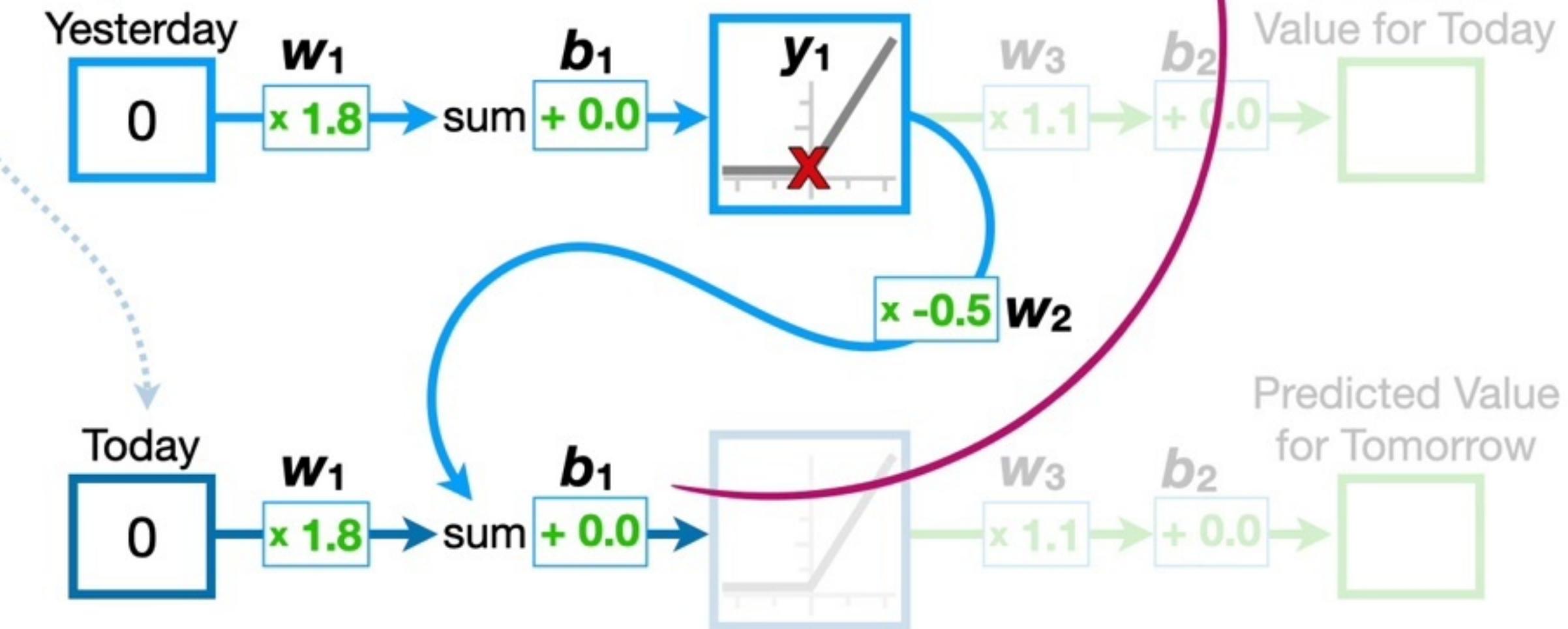


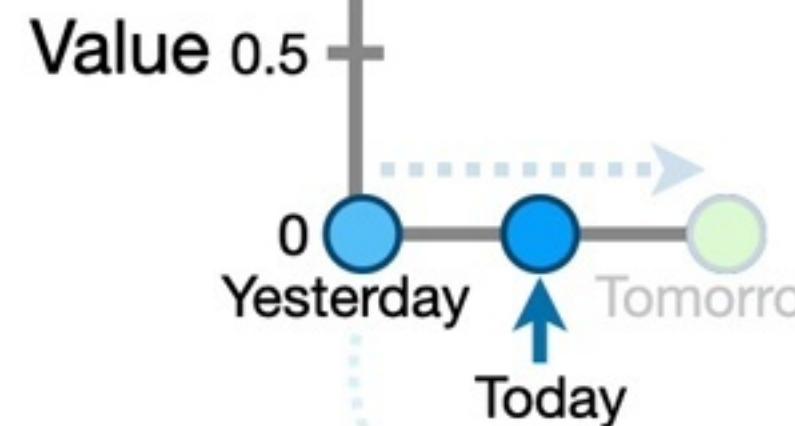


...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0$$

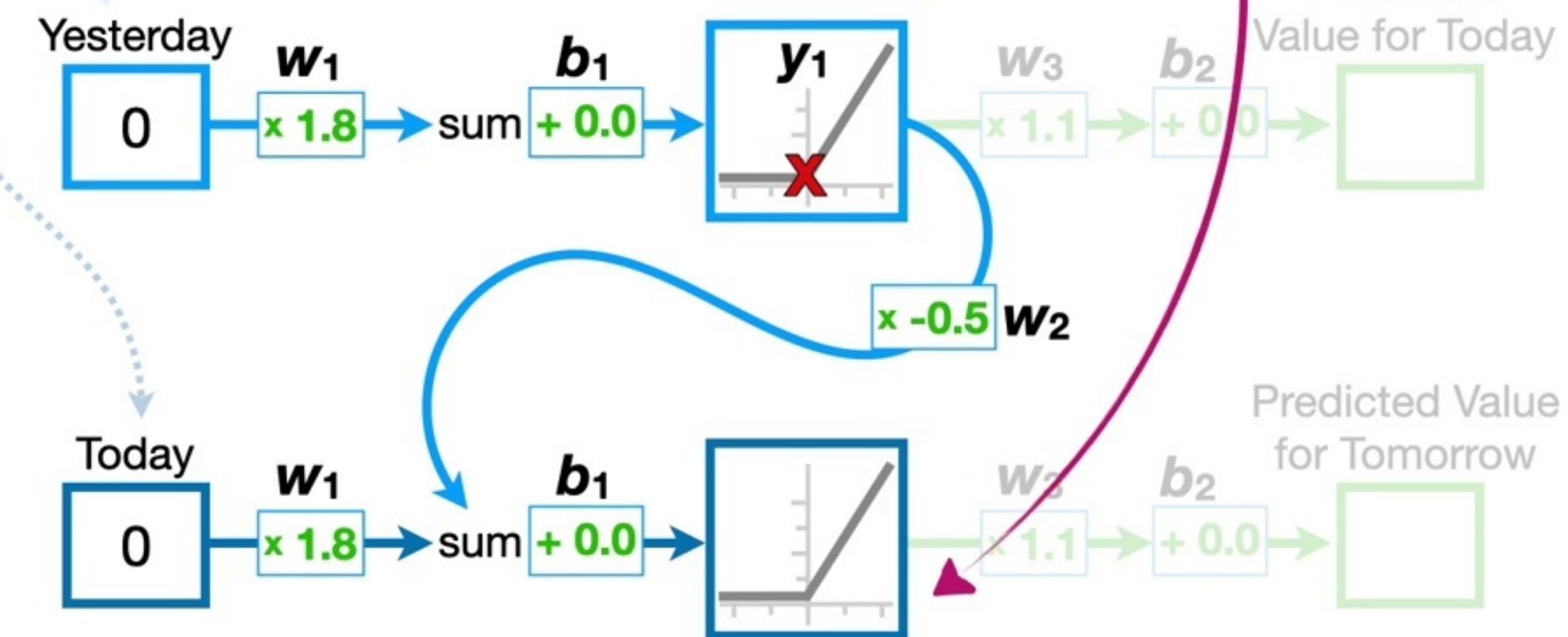


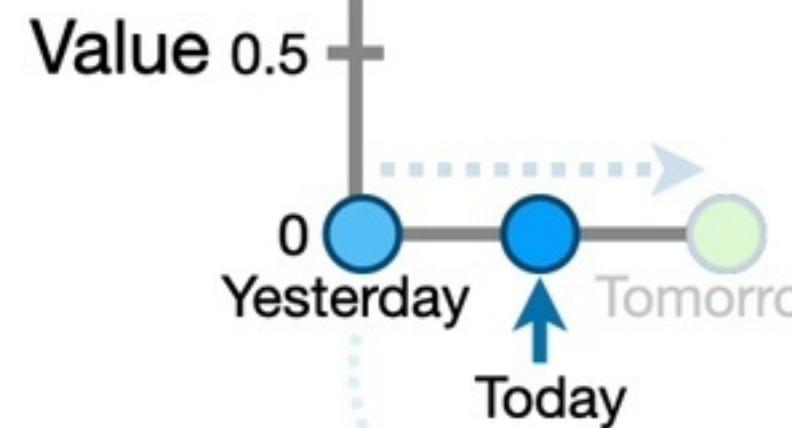


...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = x\text{-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0$$

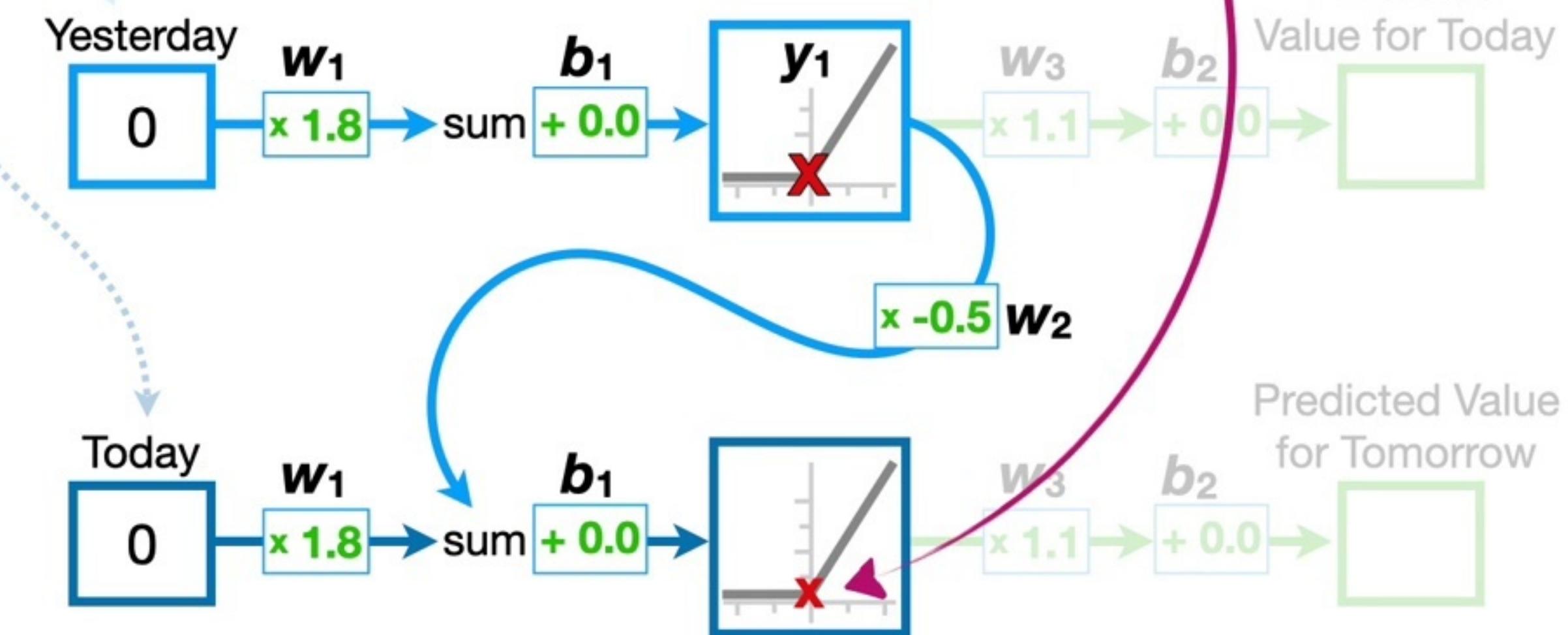


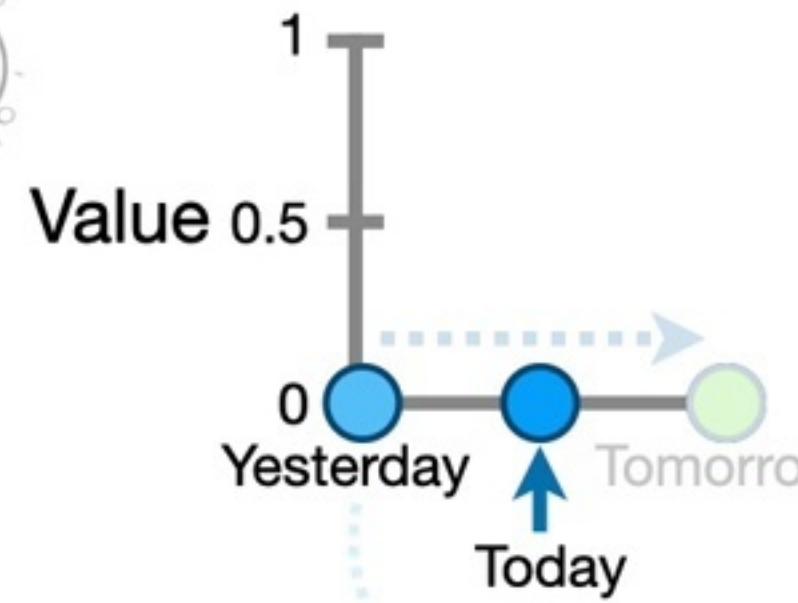


...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = x\text{-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$



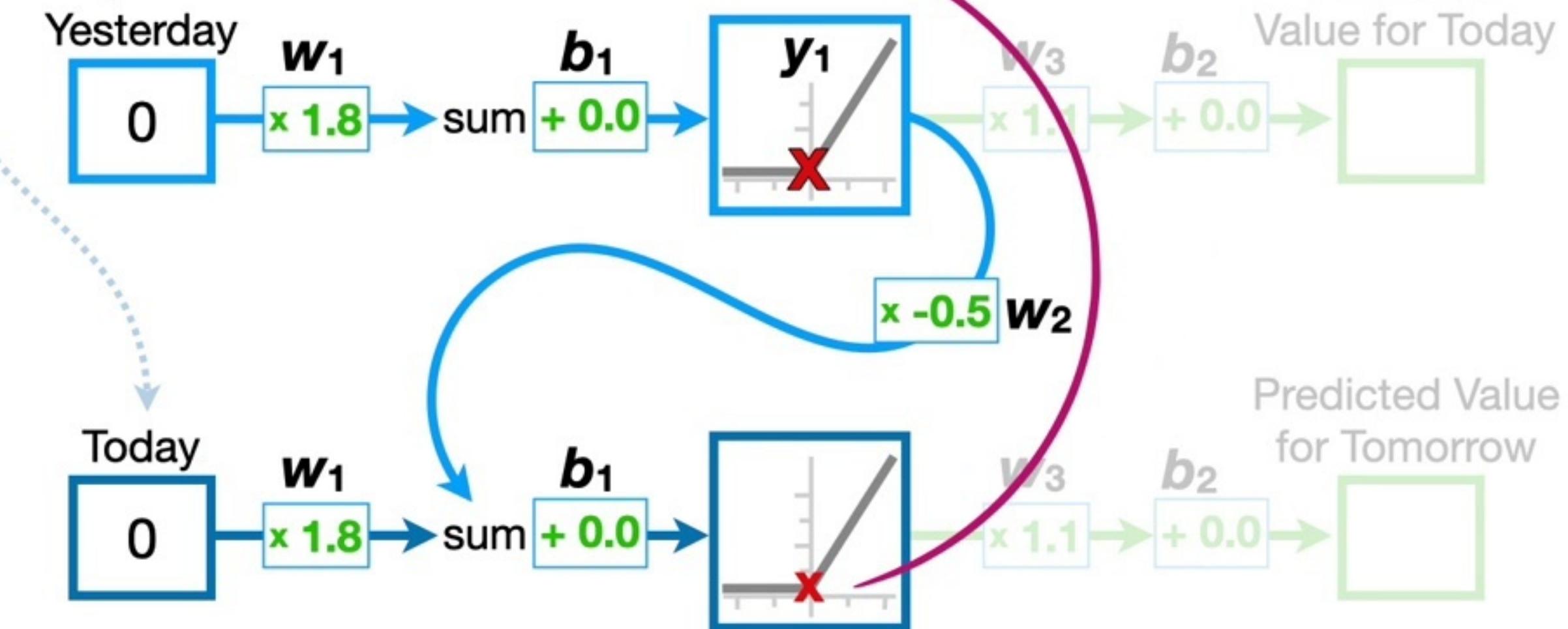


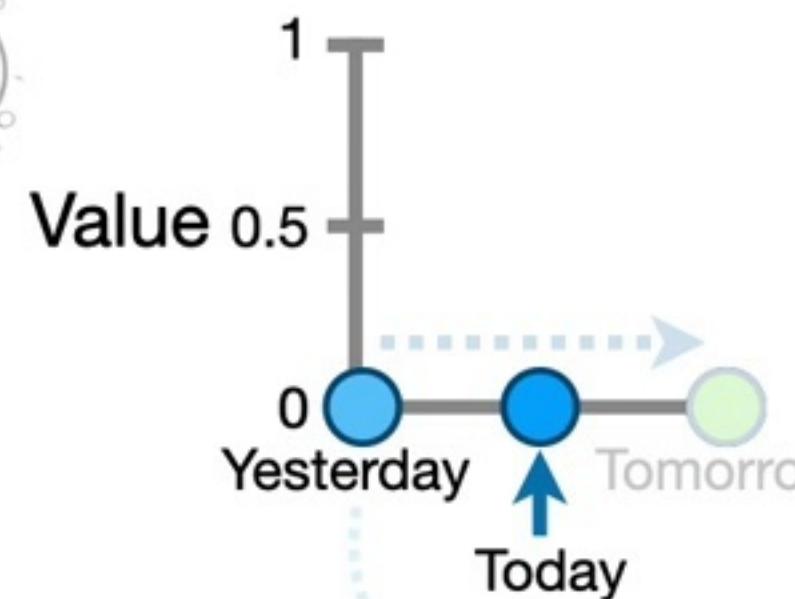
...and keep
doing the math...

$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = \text{x-axis coordinate}$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

$f(x) = \max(0, x) = \text{y-axis coordinate}$



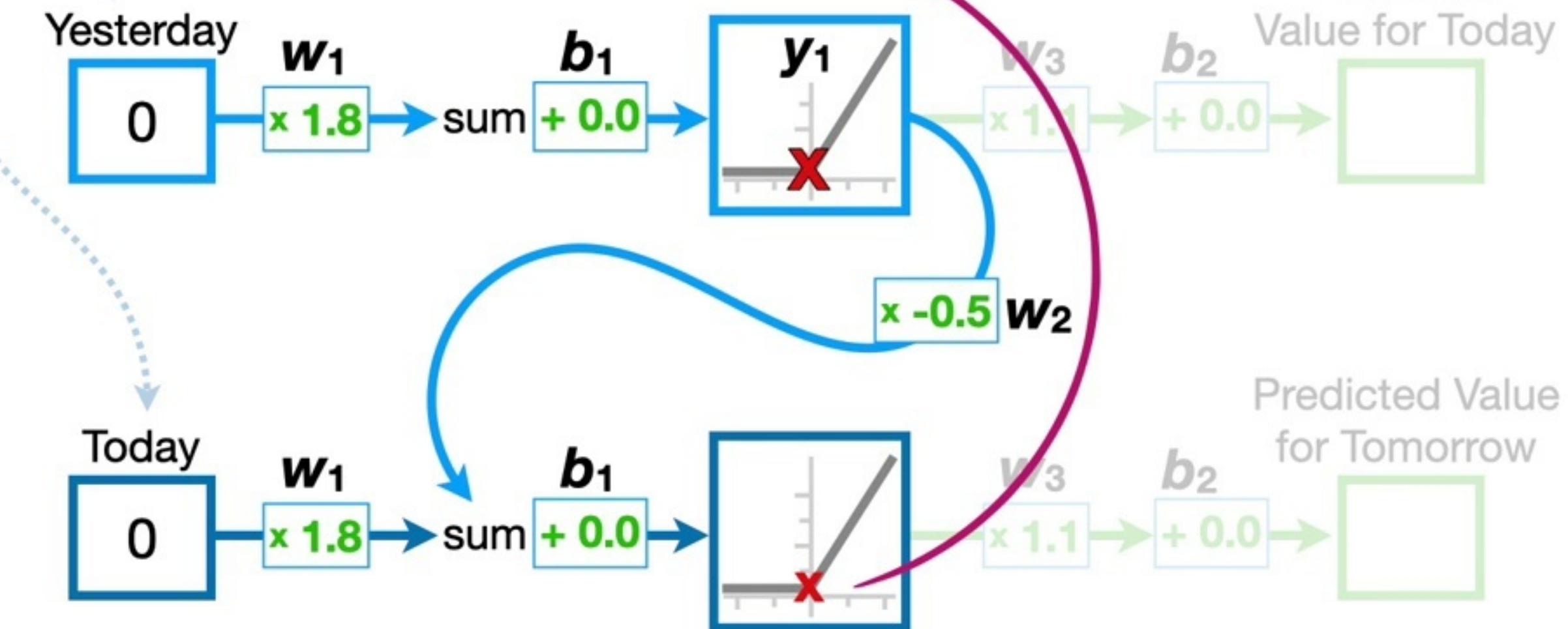


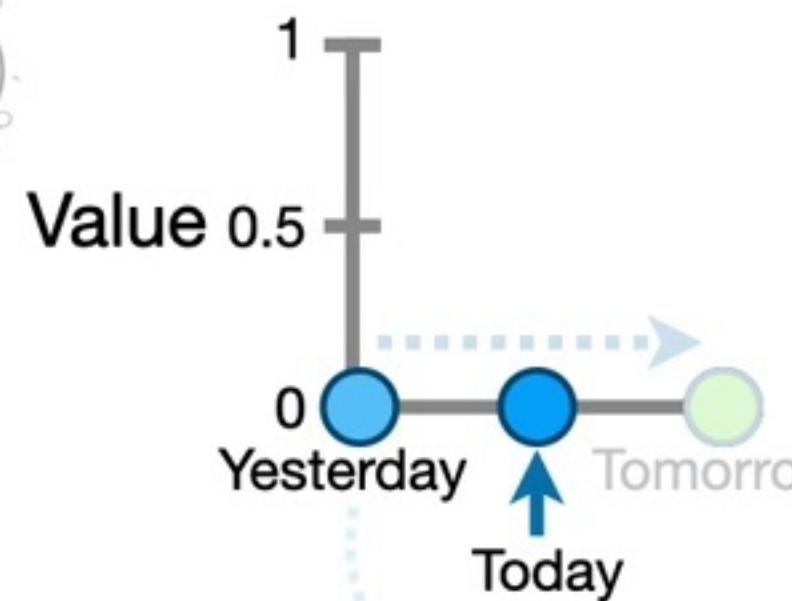
...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = x\text{-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

$$f(0) = \max(0, x) = y\text{-axis coordinate}$$



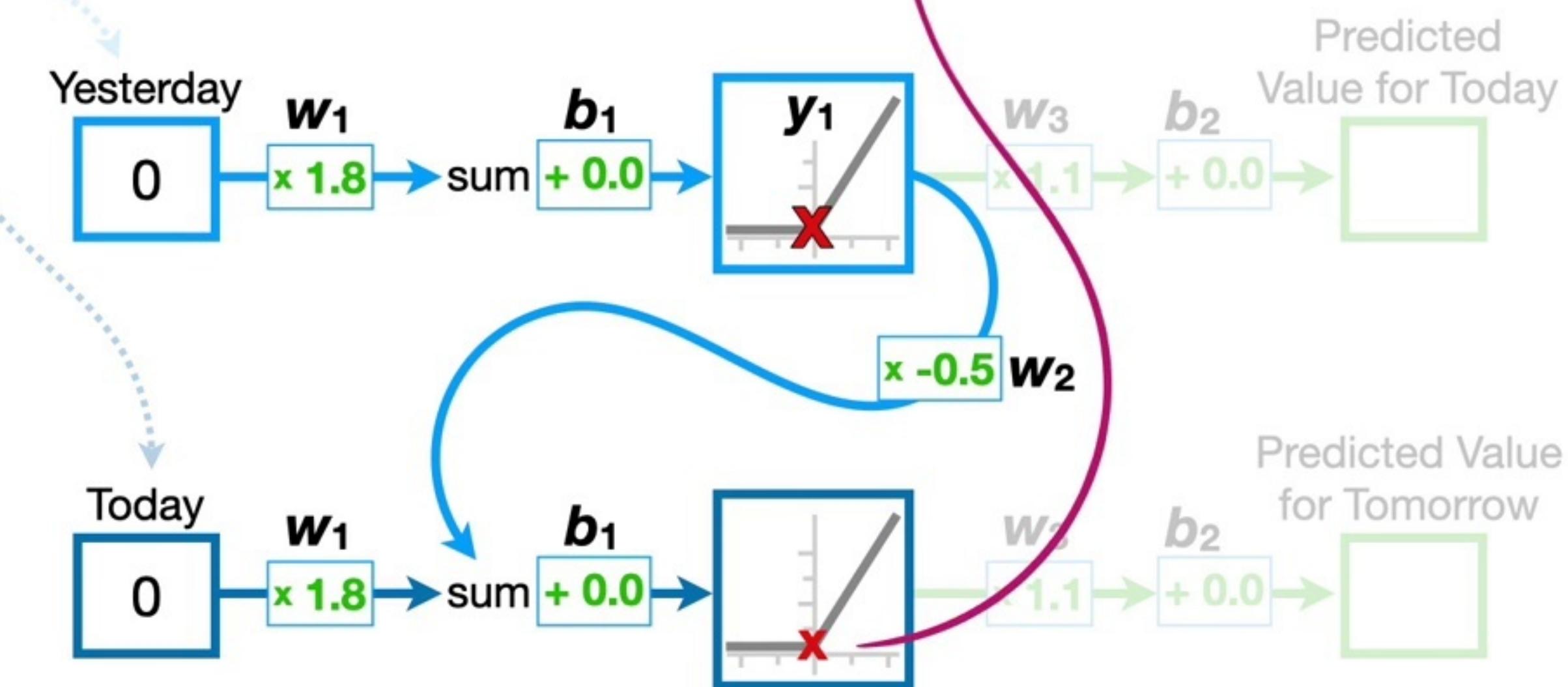


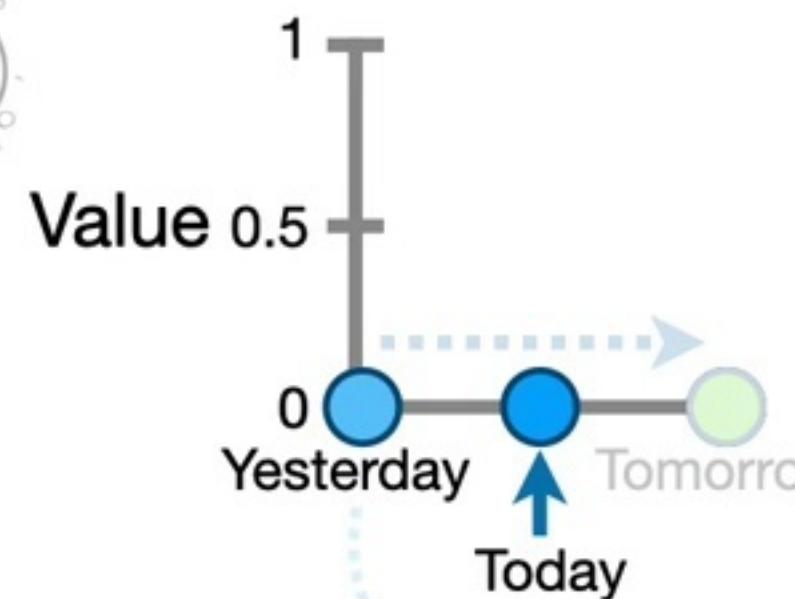
...and keep
doing the math..

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = \text{x-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

$f(0) = \max(0, x)$ = y-axis coordinate



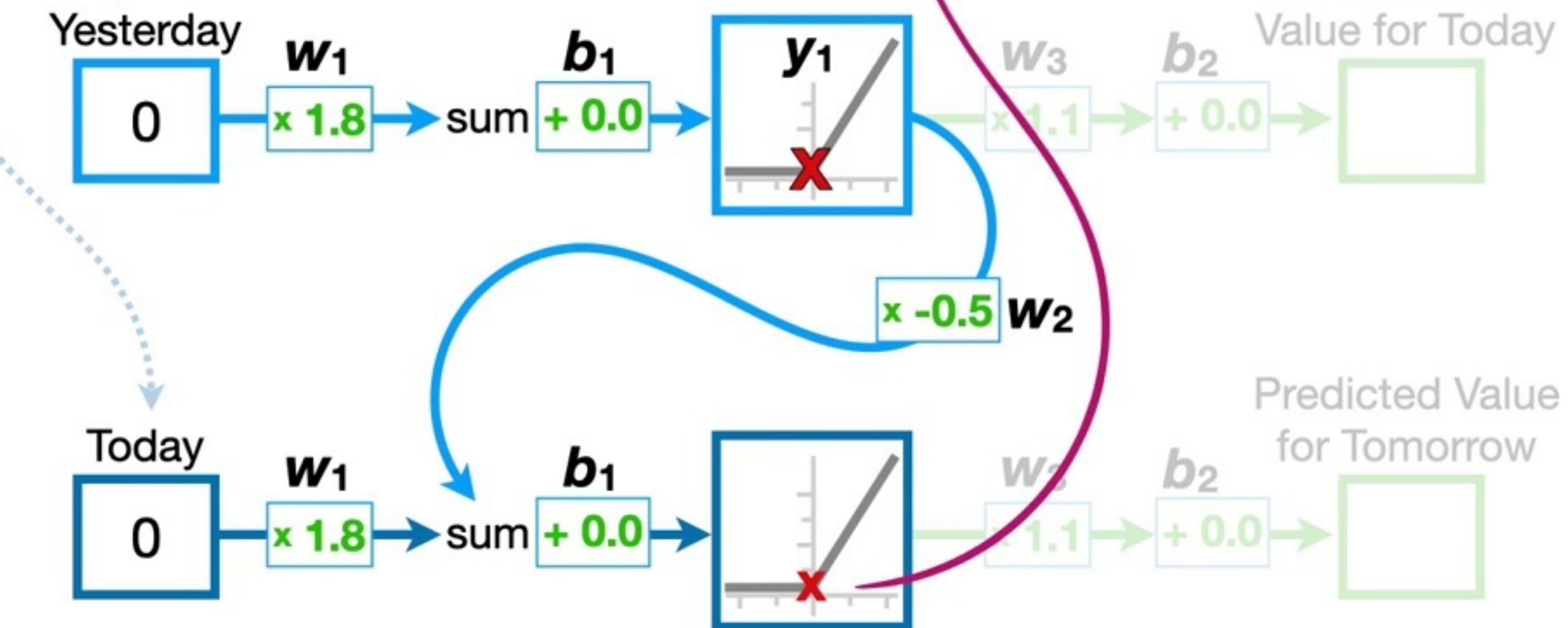


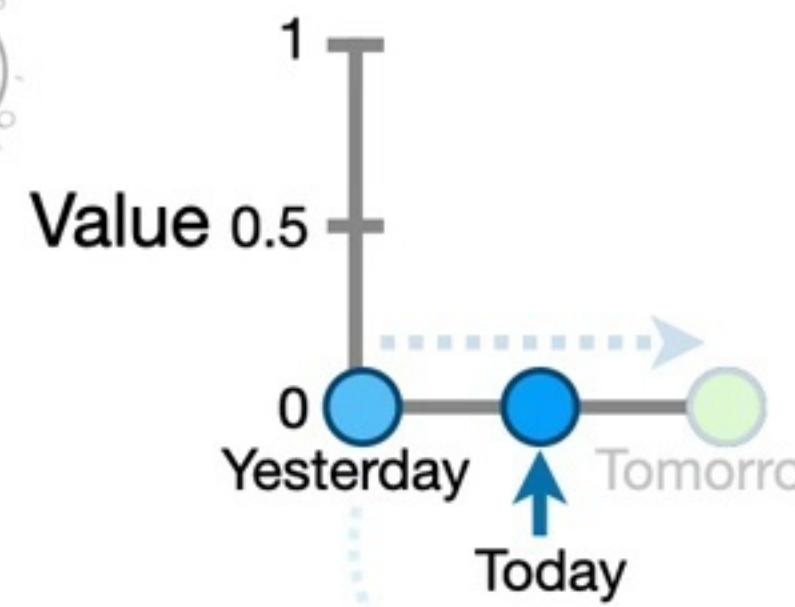
...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = \text{x-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

$$f(0) = \max(0, 0) = \text{y-axis coordinate}$$



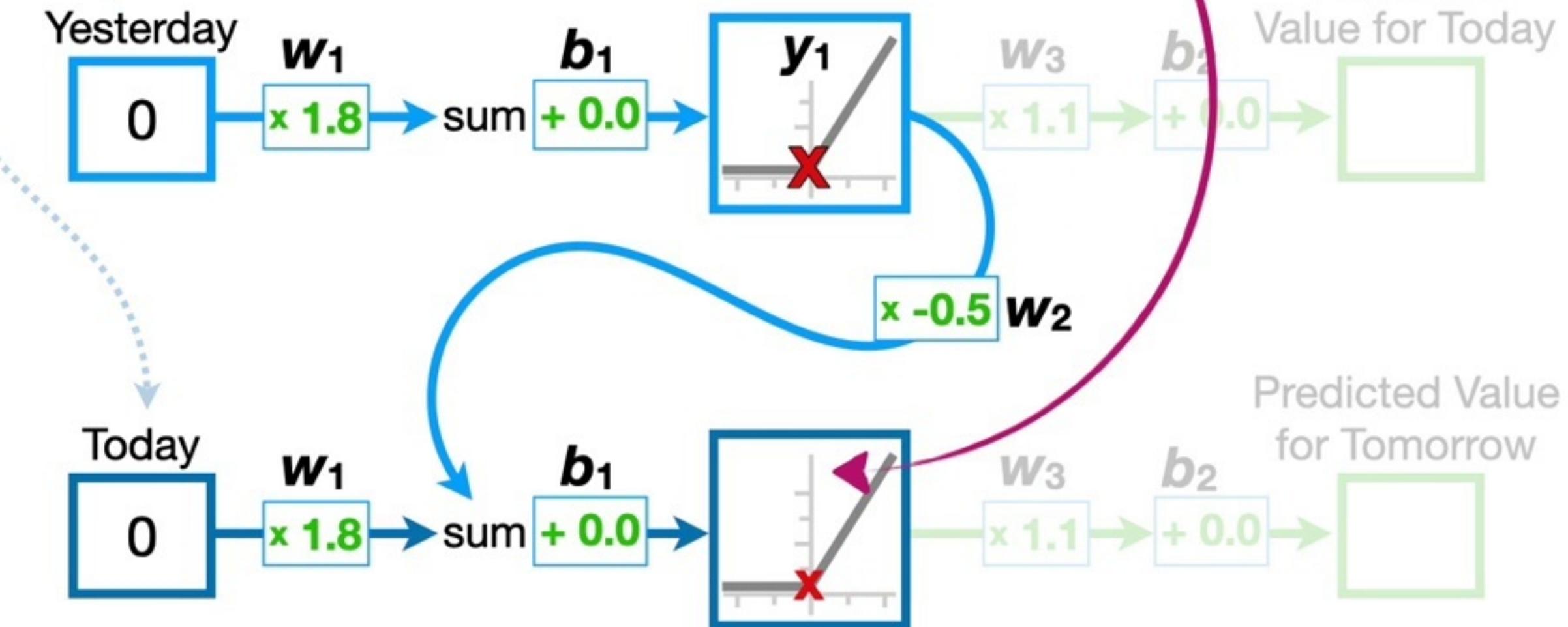


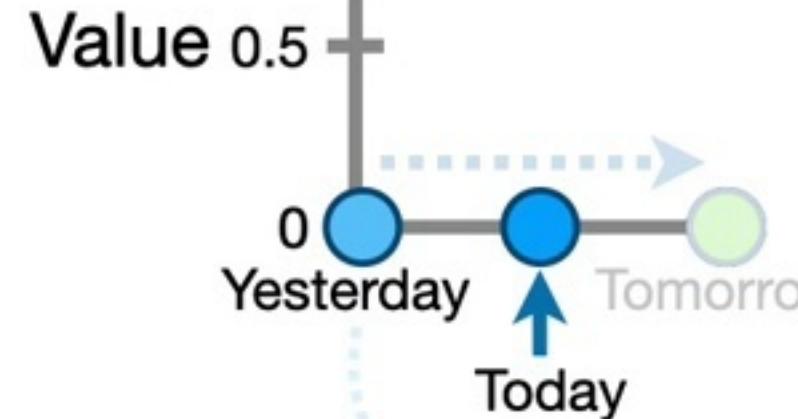
...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = \text{x-axis coordinate}$$

$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

$$f(0) = \max(0, 0) = \text{y-axis coordinate}$$



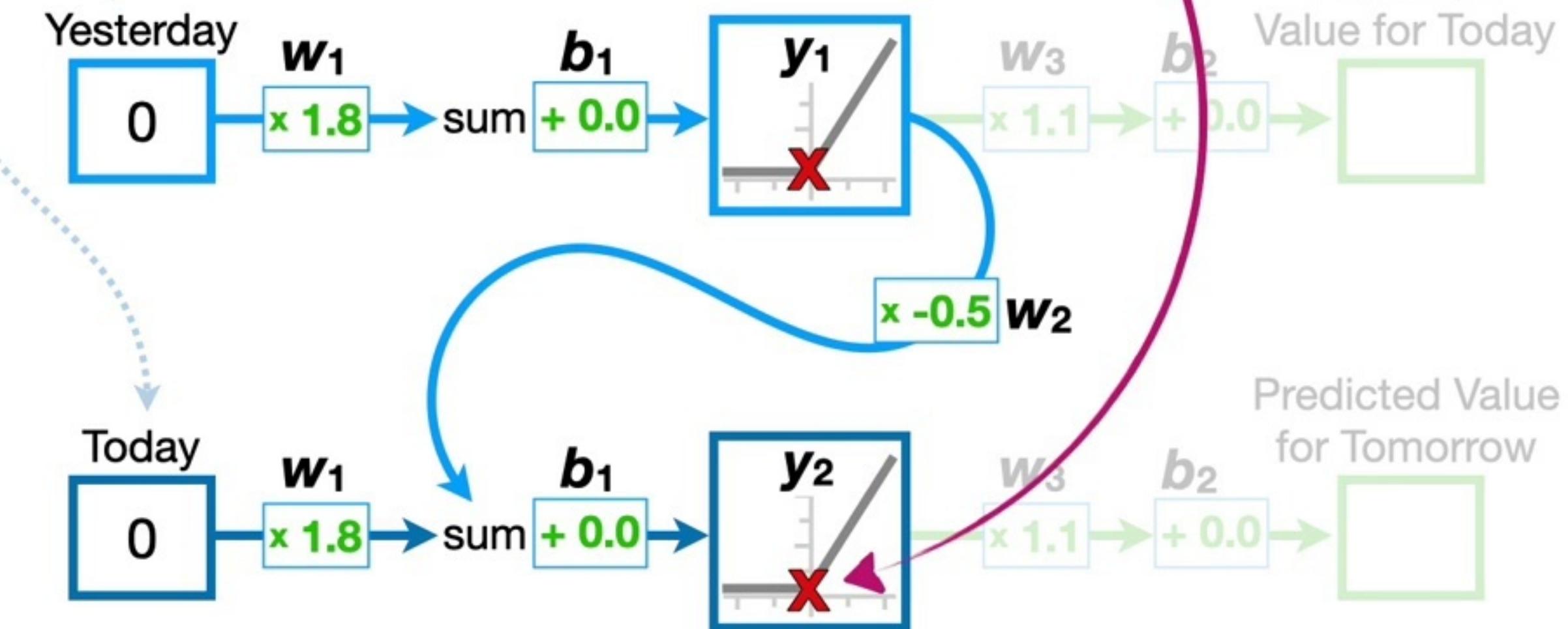


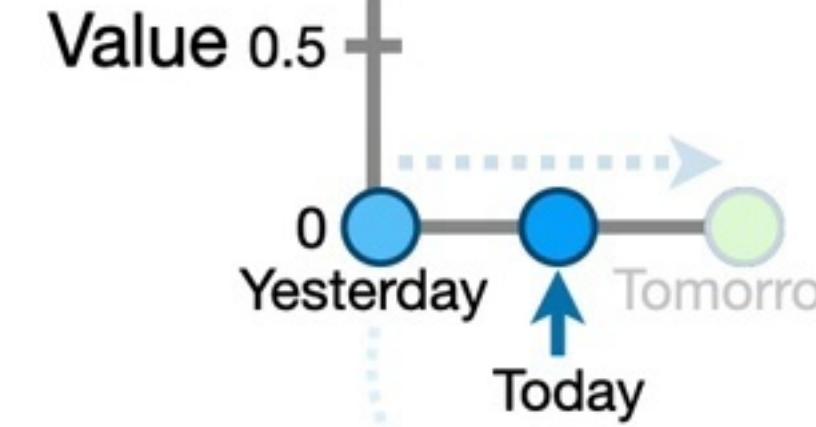
...and keep
doing the math...

$$(\text{Today} \times w_1) + (y_1 \times w_2) + b_1 = x\text{-axis coordinate}$$

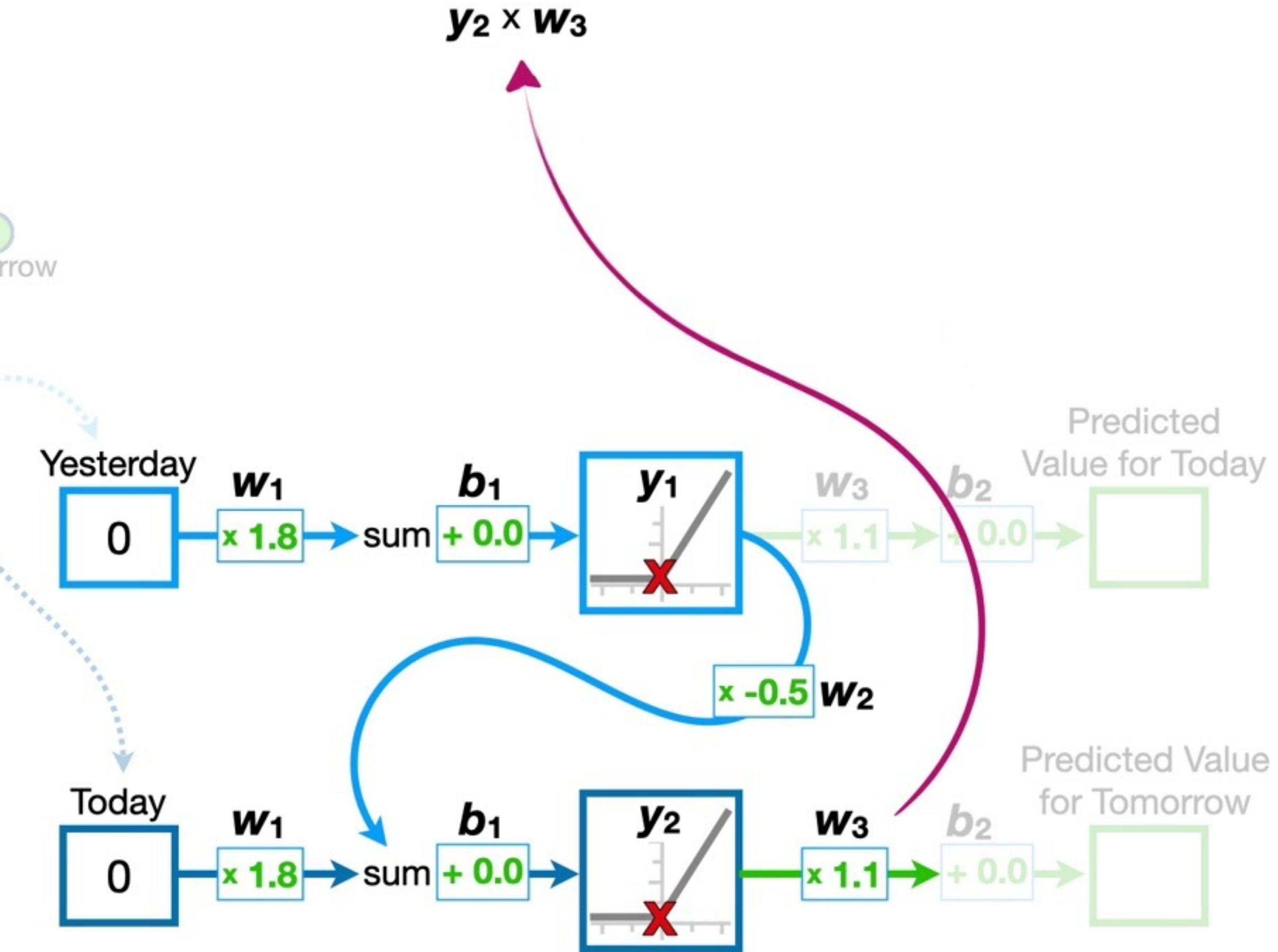
$$0 \times 1.8 + (0 \times -0.5) + 0.0 = 0$$

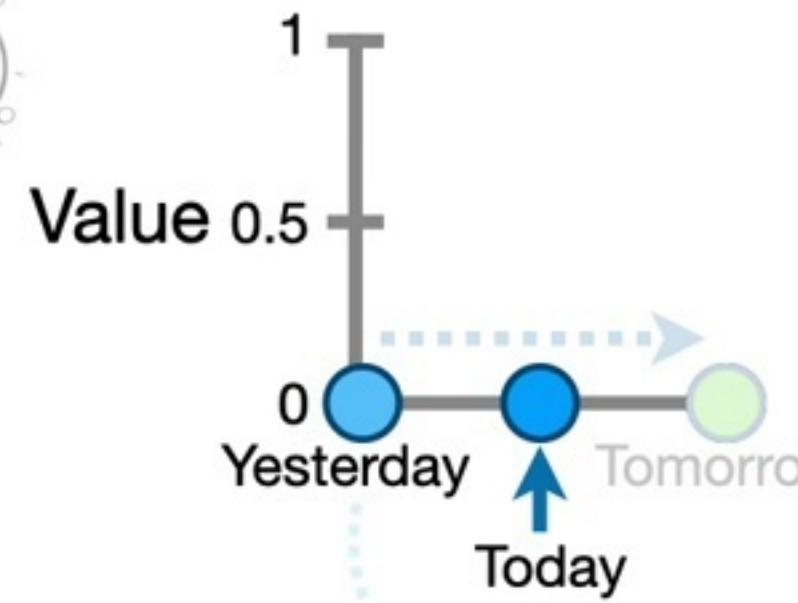
$$f(0) = \max(0, 0) = 0 = y_2$$



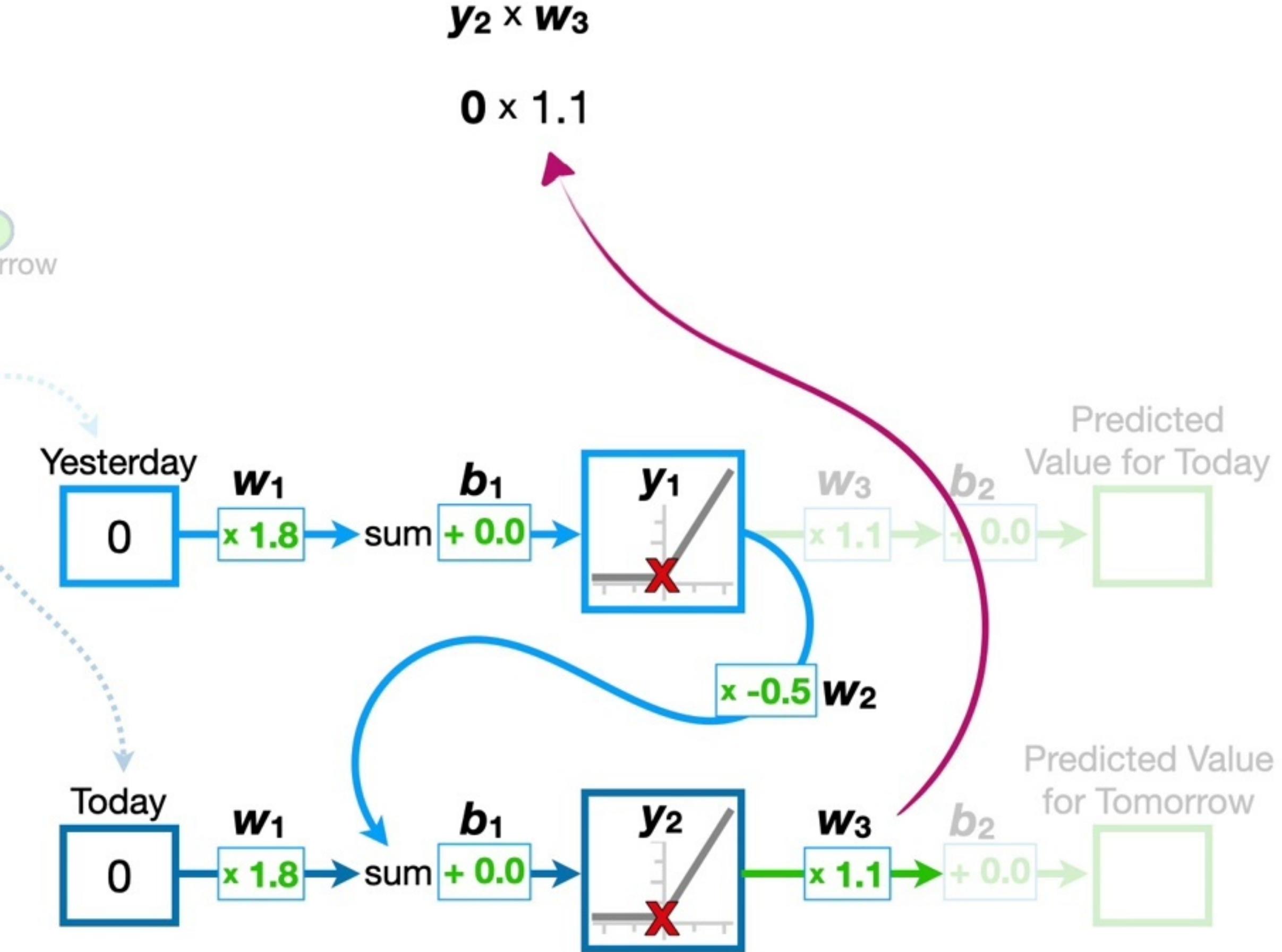


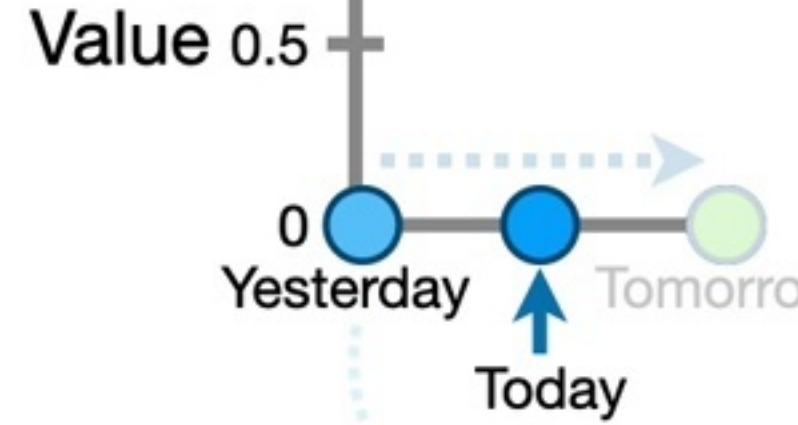
...and keep
doing the math...



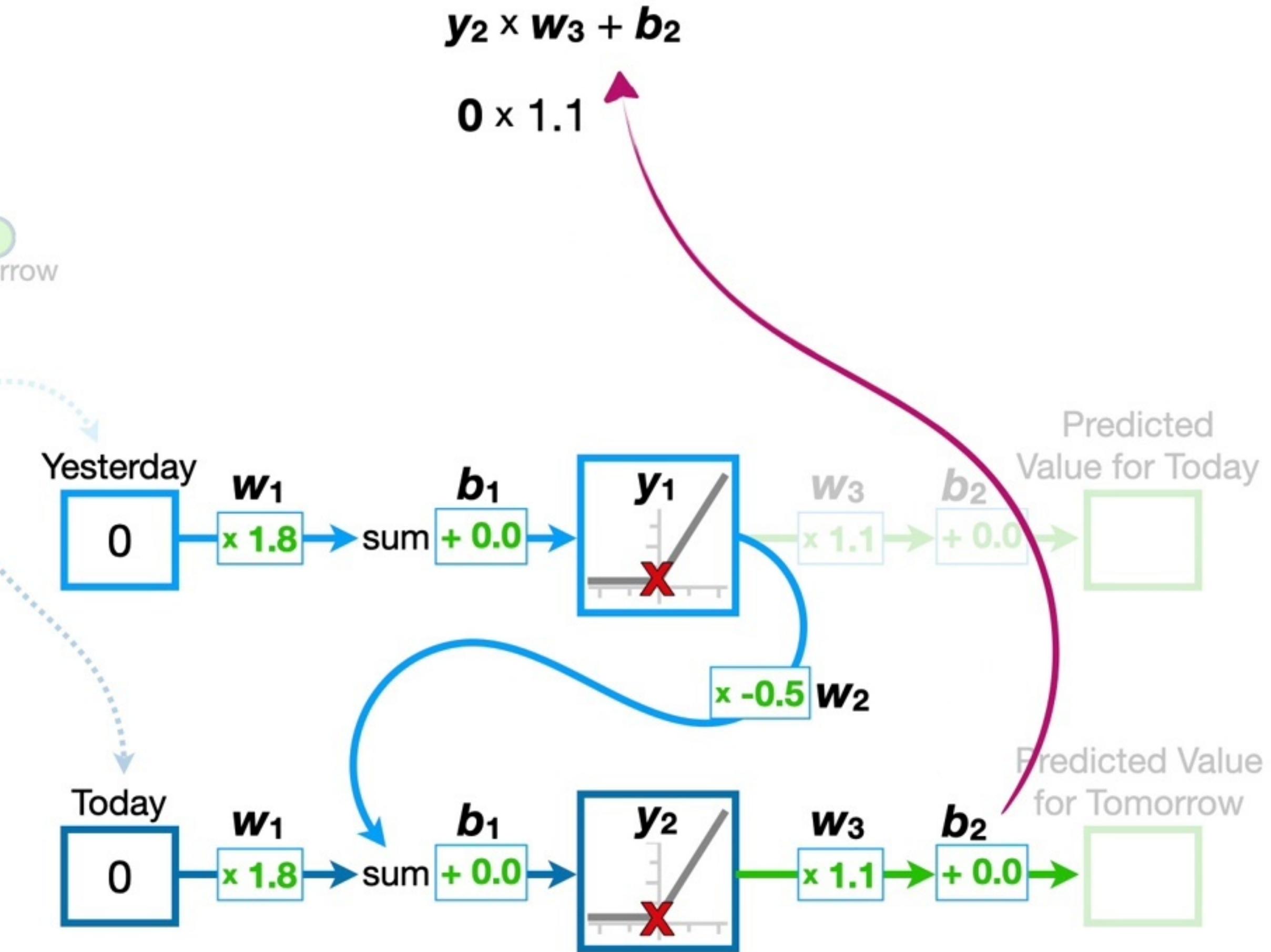


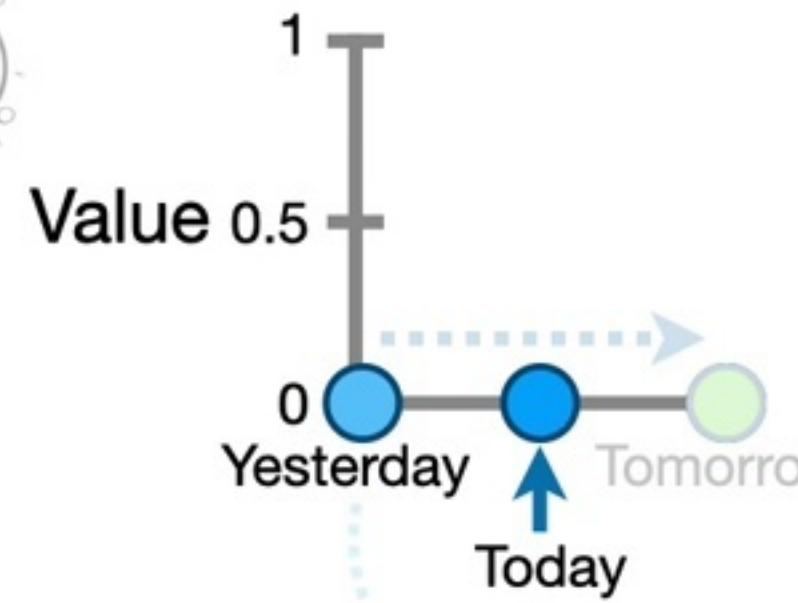
...and keep
doing the math...



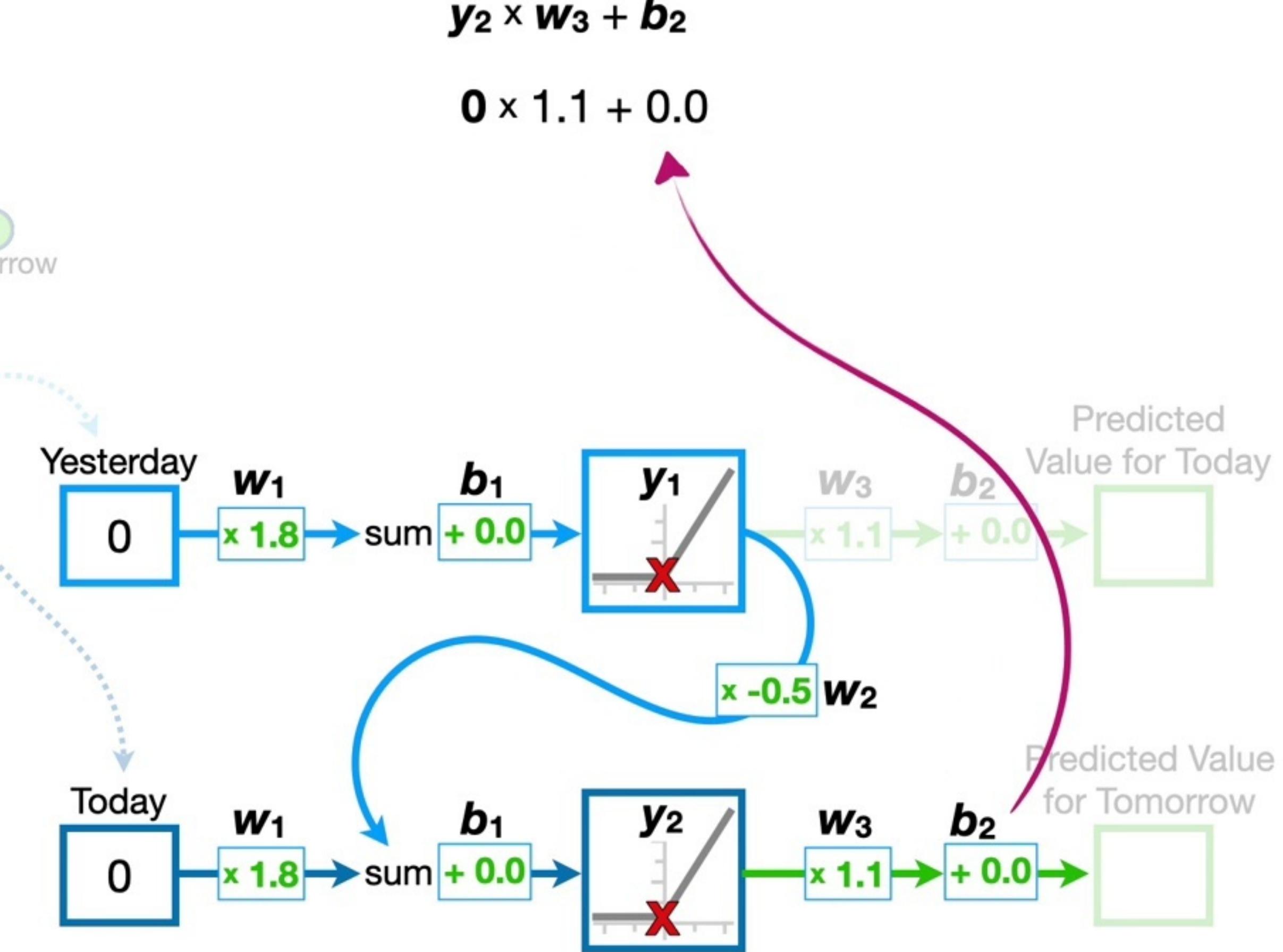


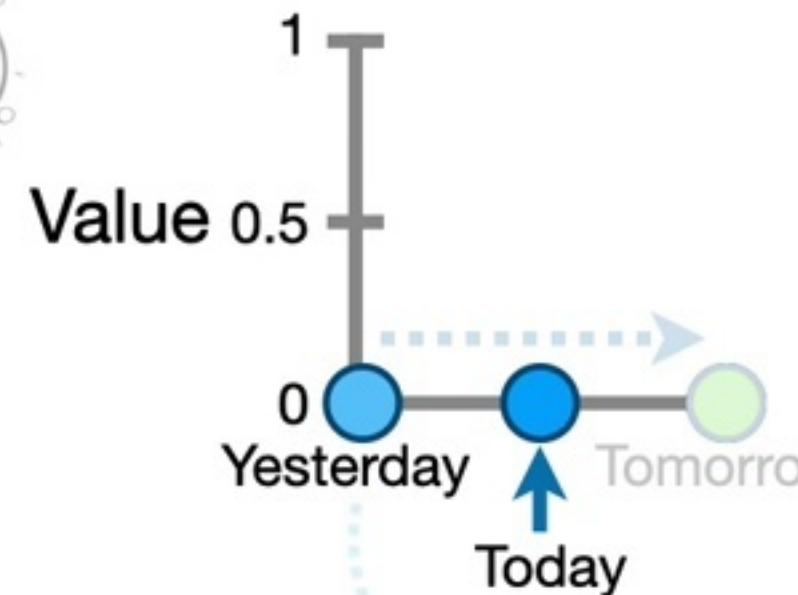
...and keep
doing the math...





...and keep
doing the math...

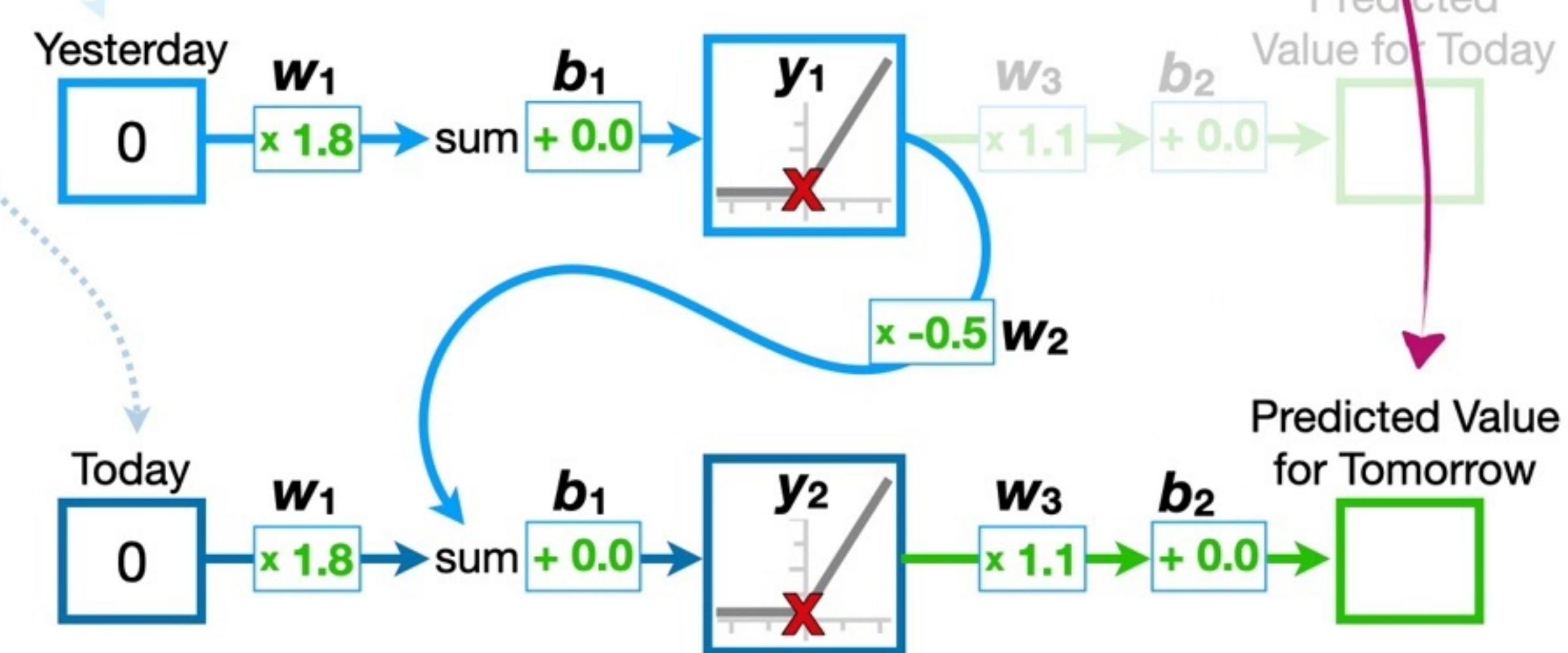


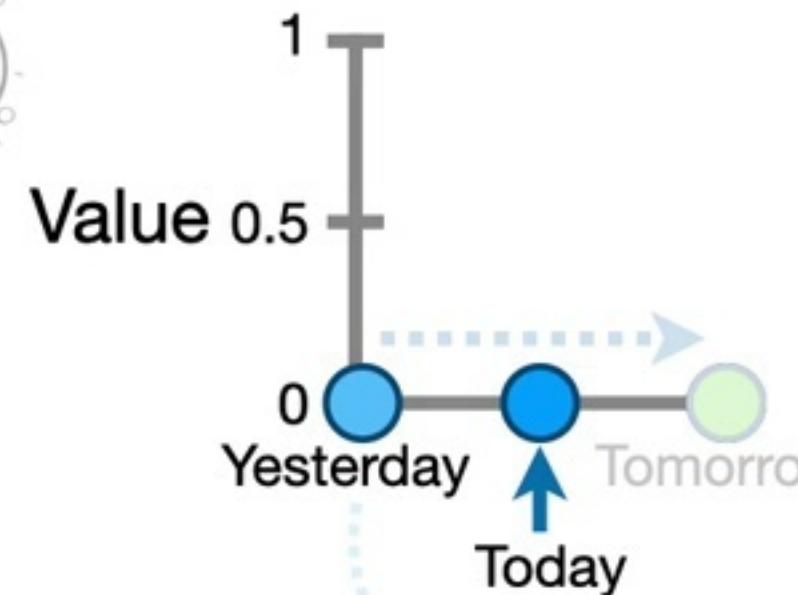


...and that gives us
a predicted value
for **tomorrow**, 0...

$y_2 \times w_3 + b_2 = \text{Prediction for Tomorrow}$

$$0 \times 1.1 + 0.0$$





...and that gives us
a predicted value
for **tomorrow**, 0...

Yesterday



w₁

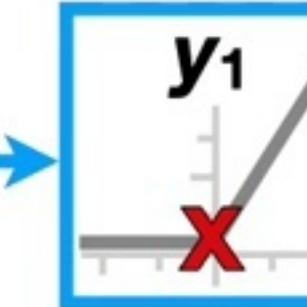
$$\times 1.8$$

sum

$$+ 0.0$$

b₁

$$+ 0.0$$



$$\times -0.5 \quad \mathbf{w_2}$$

Today



w₁

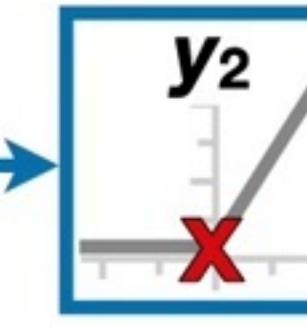
$$\times 1.8$$

sum

$$+ 0.0$$

b₁

$$+ 0.0$$



w₃

$$\times 1.1$$

b₂

$$+ 0.0$$

Predicted Value
for Tomorrow

$$0$$

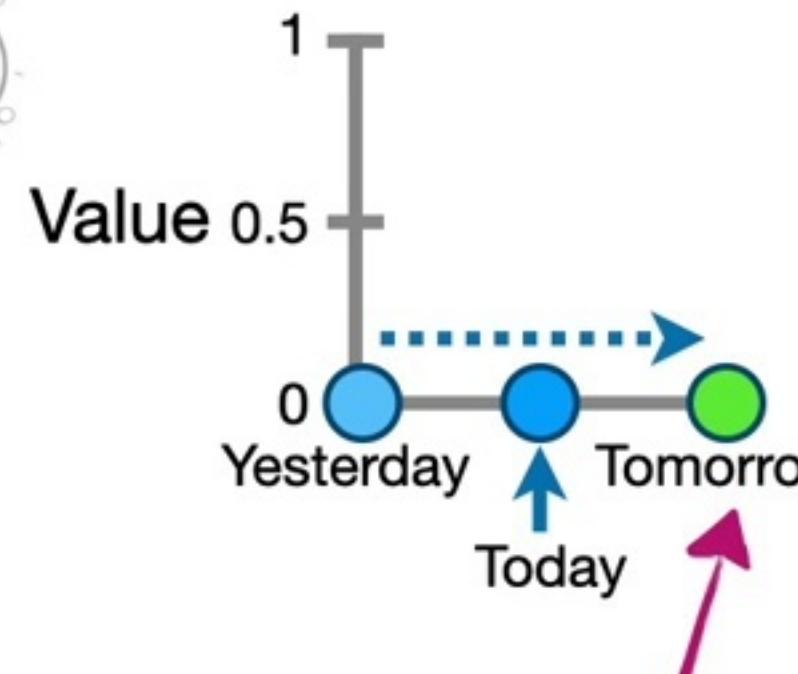
y₂ × w₃ + b₂ = Prediction for Tomorrow

$$0 \times 1.1 + 0.0 = 0$$

Predicted
Value for Today

$$0$$

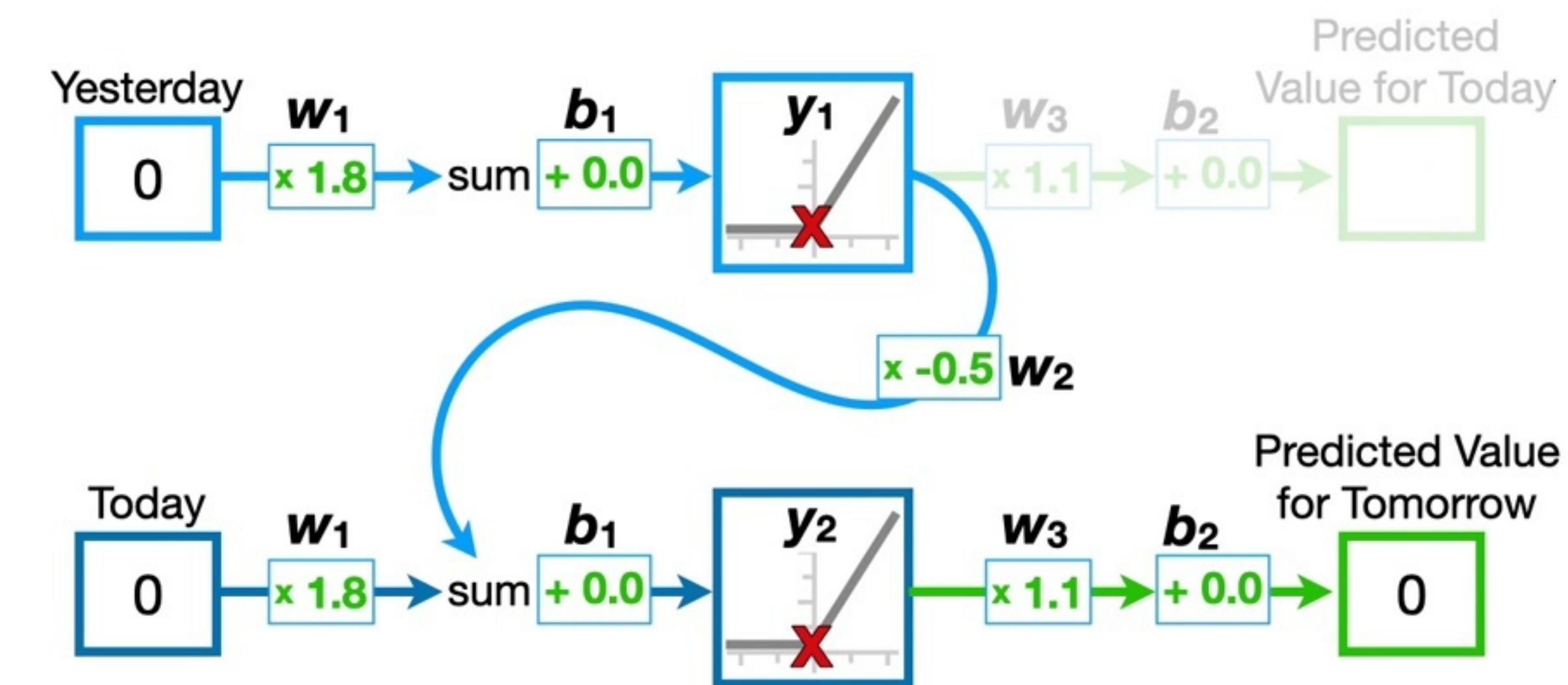
$$0$$

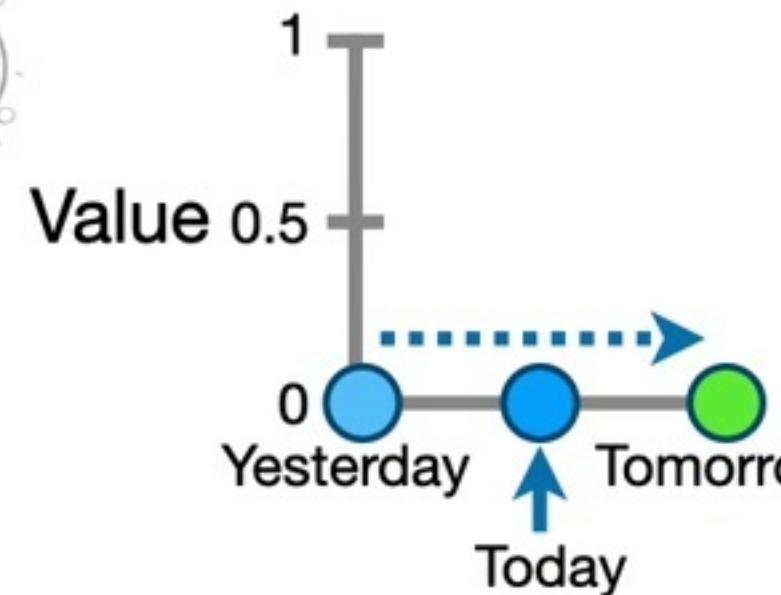


...which is
consistent with the
original observation.

$$y_2 \times w_3 + b_2 = \text{Prediction for Tomorrow}$$

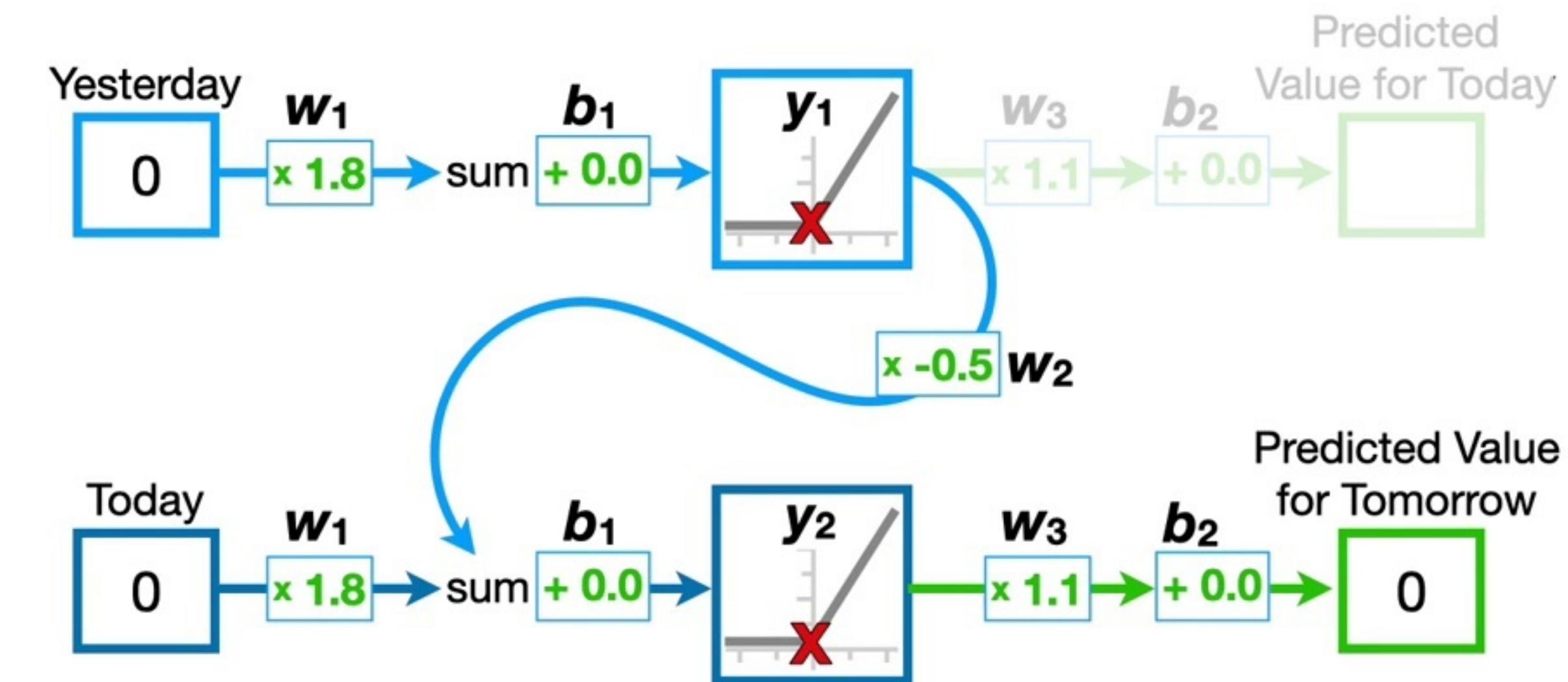
$$0 \times 1.1 + 0.0 = 0$$

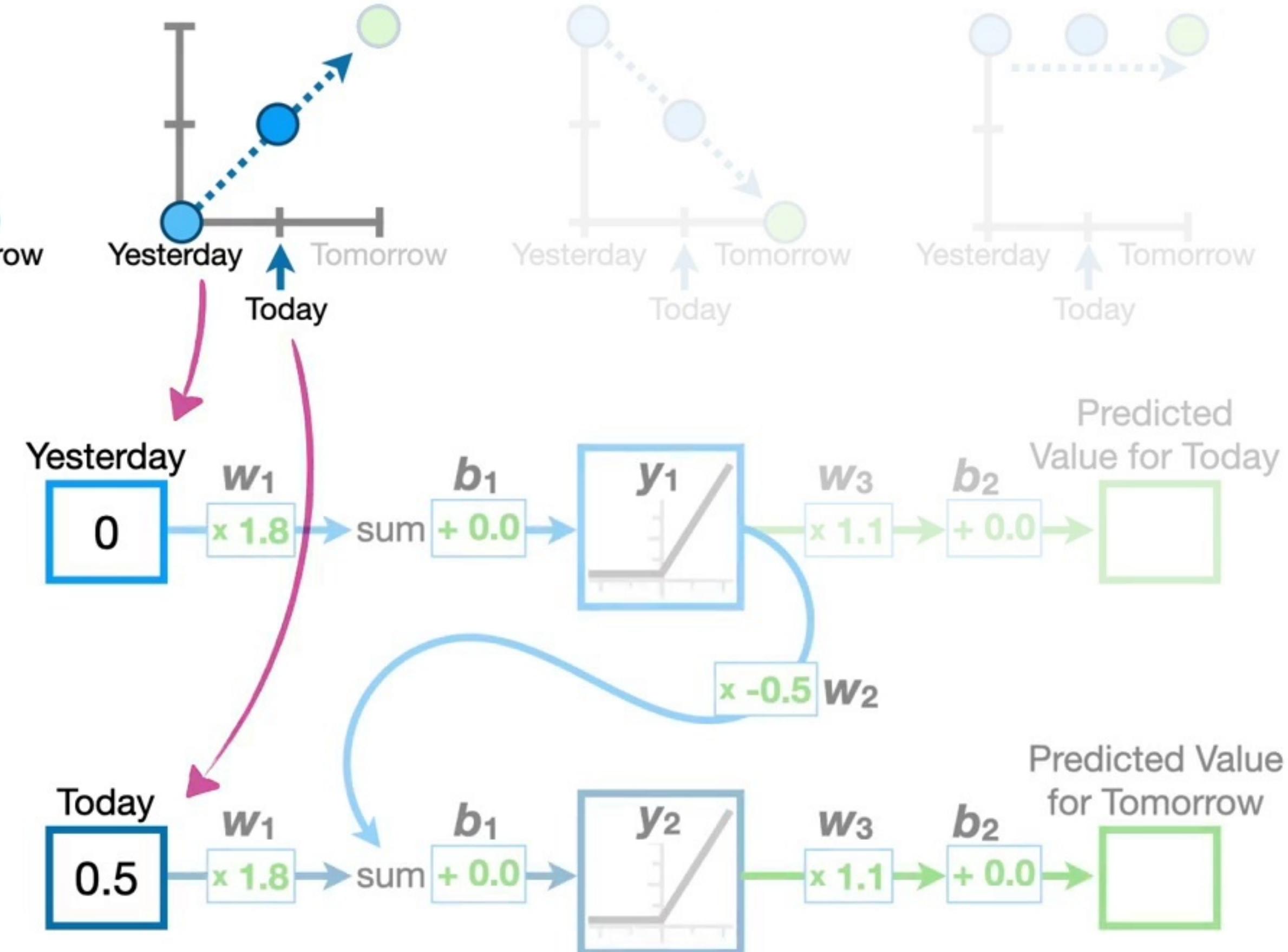
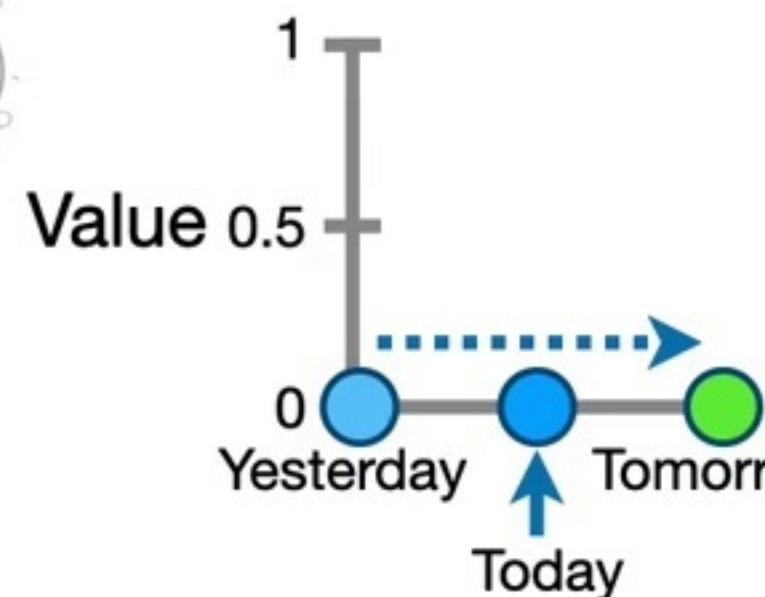




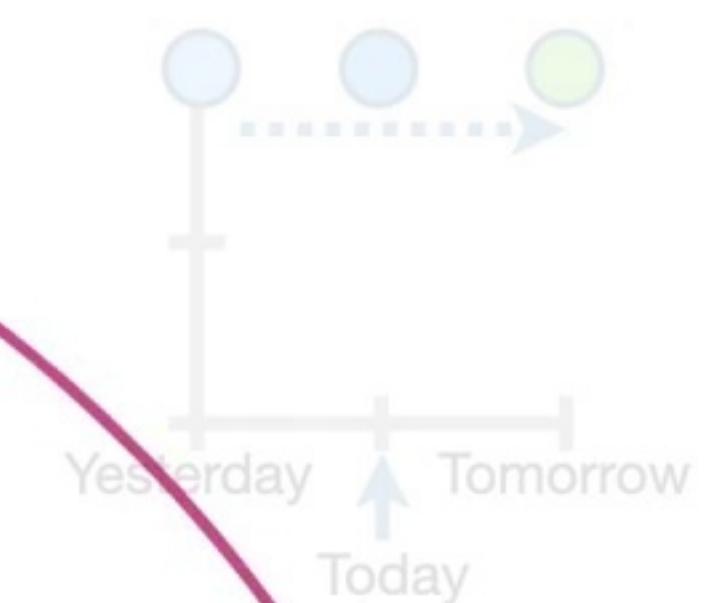
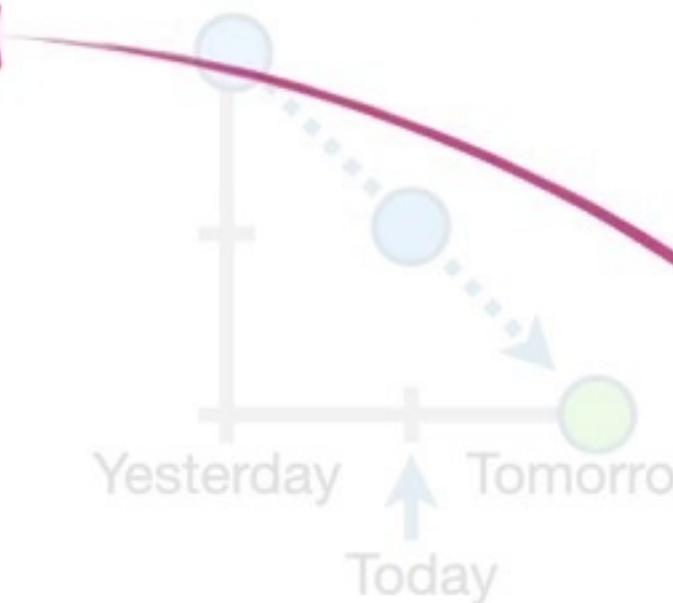
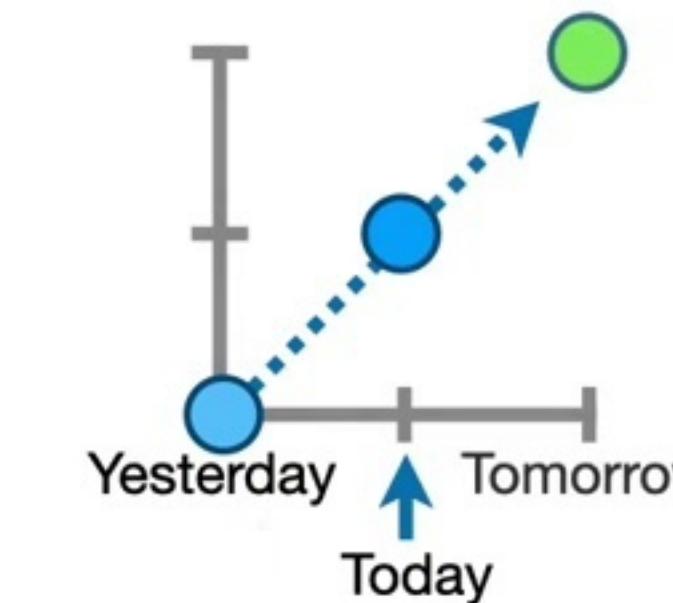
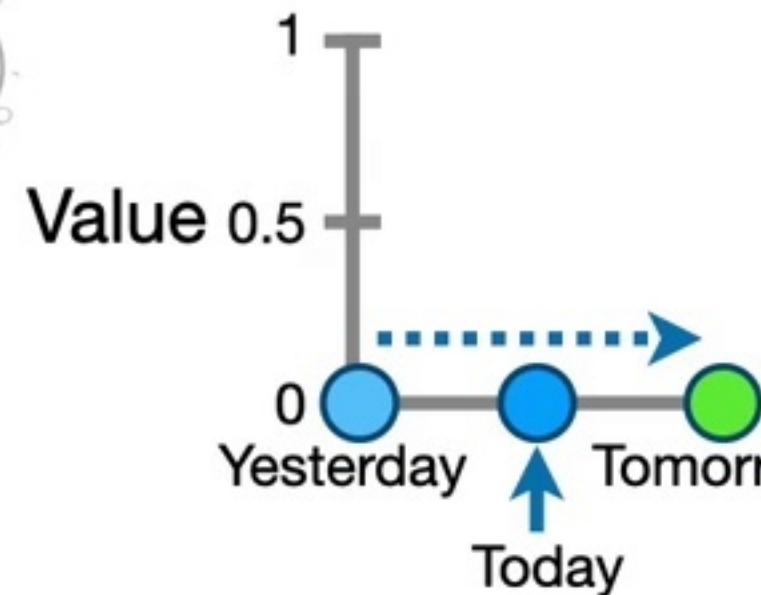
$$y_2 \times w_3 + b_2 = \text{Prediction for Tomorrow}$$
$$0 \times 1.1 + 0.0 = 0$$

In other words, the recurrent neural network correctly predicted tomorrow's value.

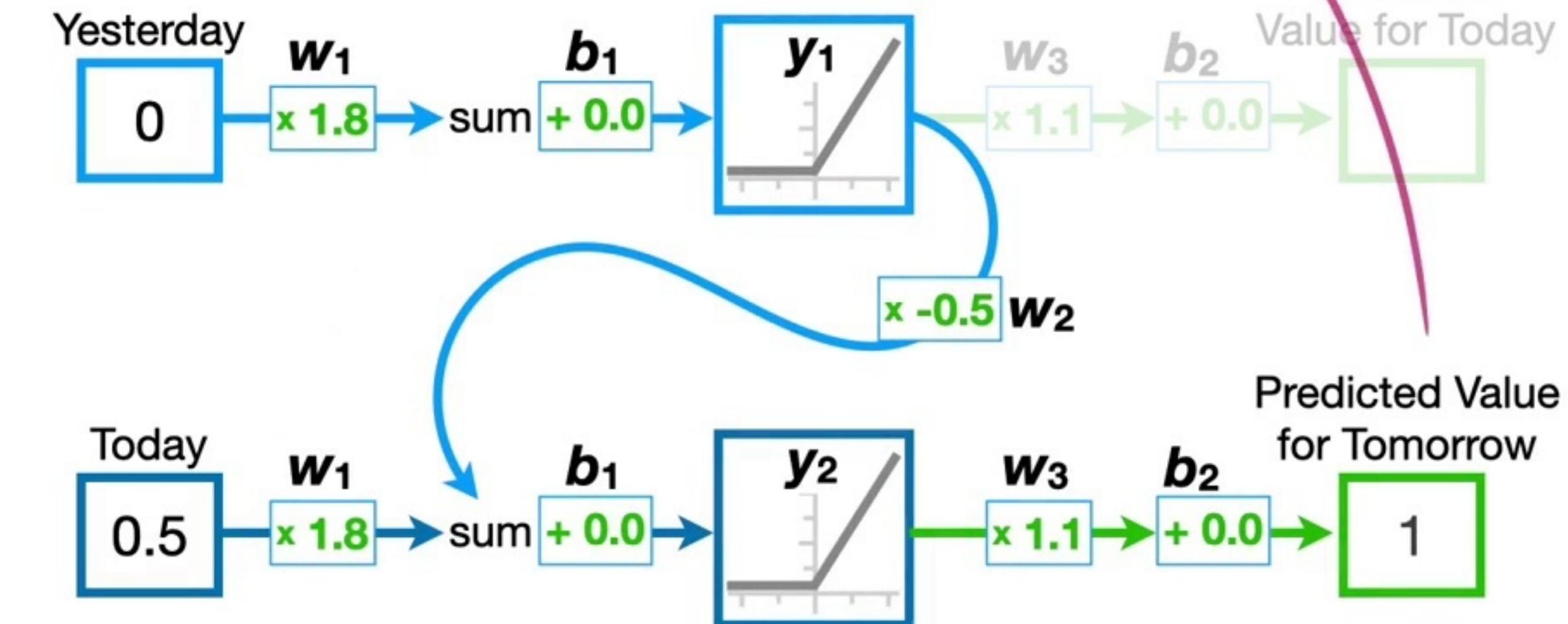


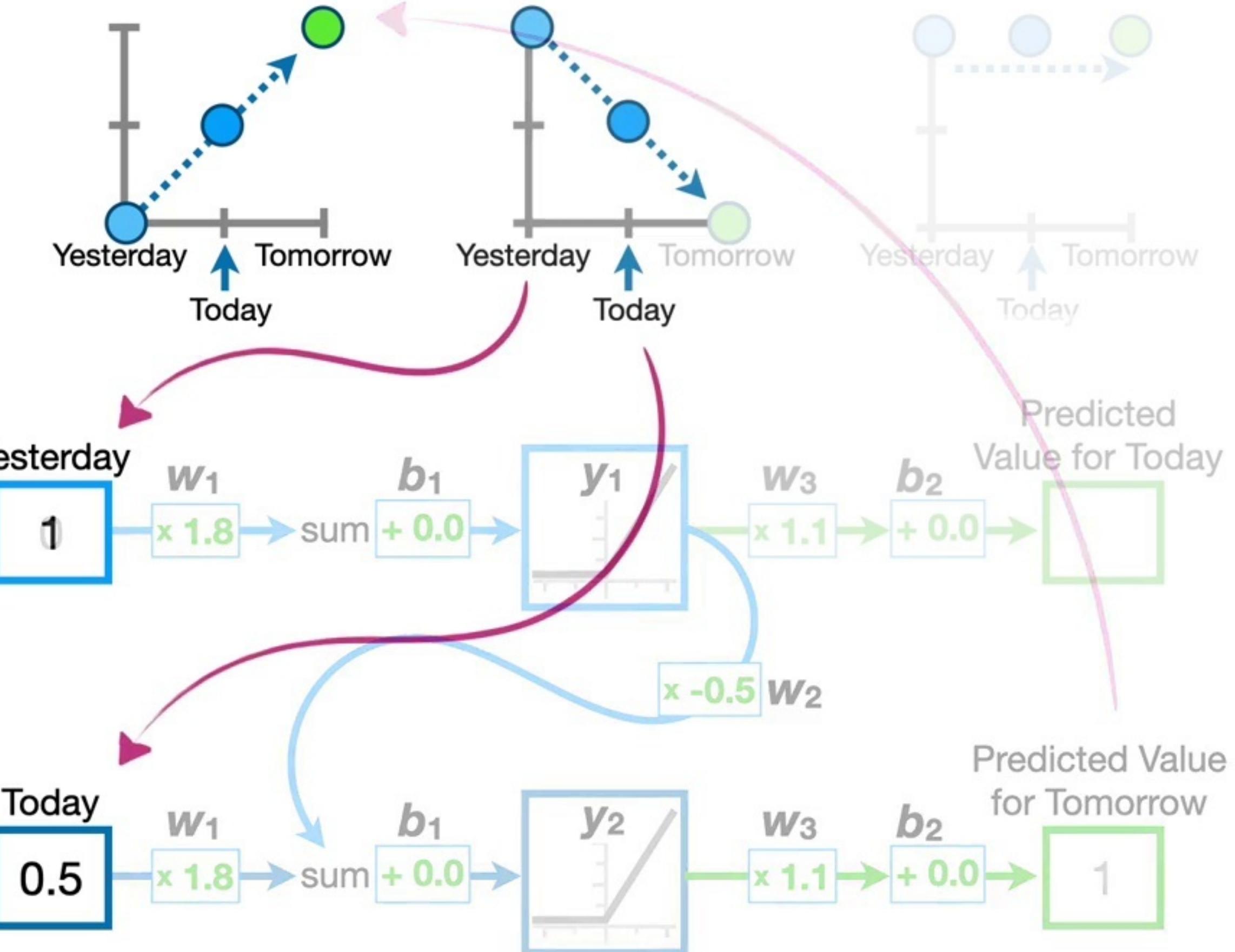
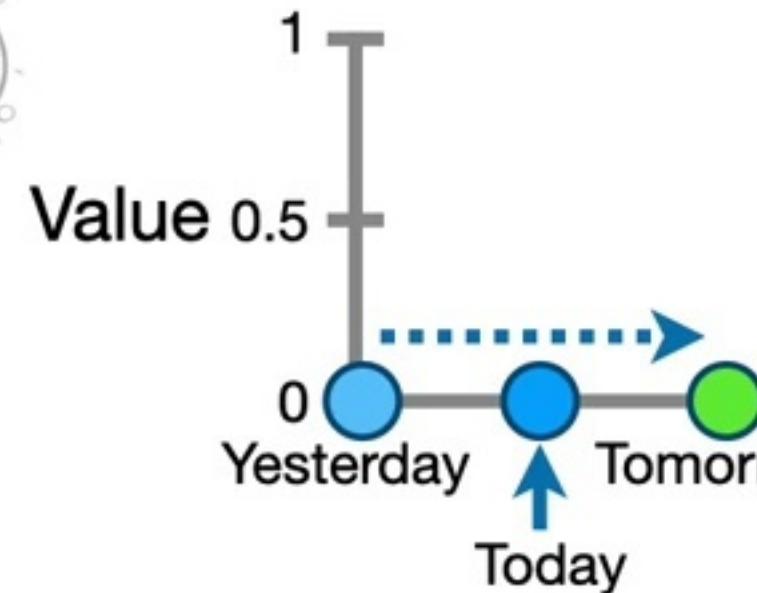


Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.

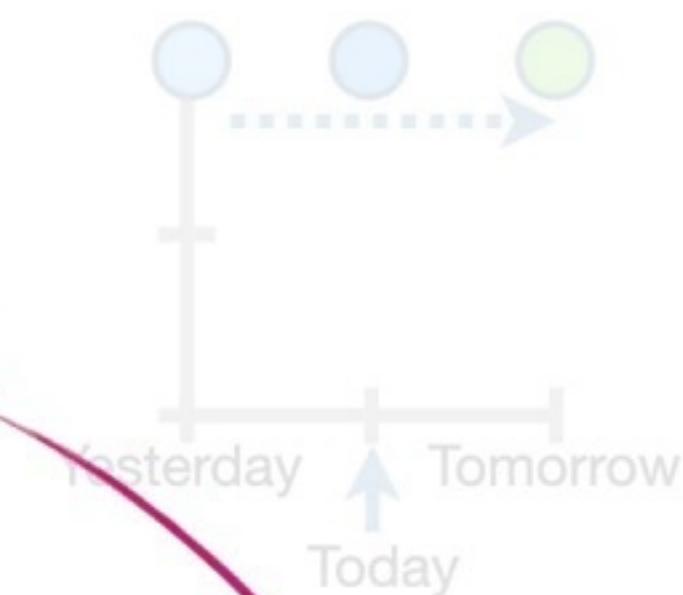
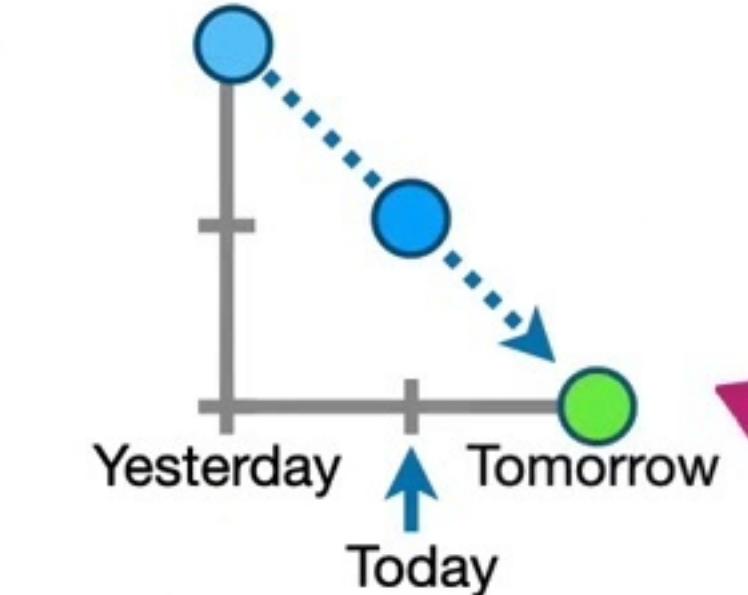
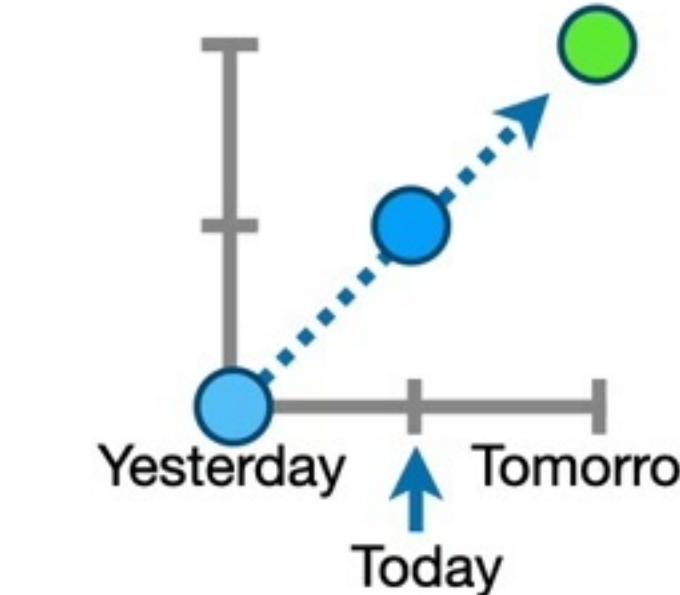
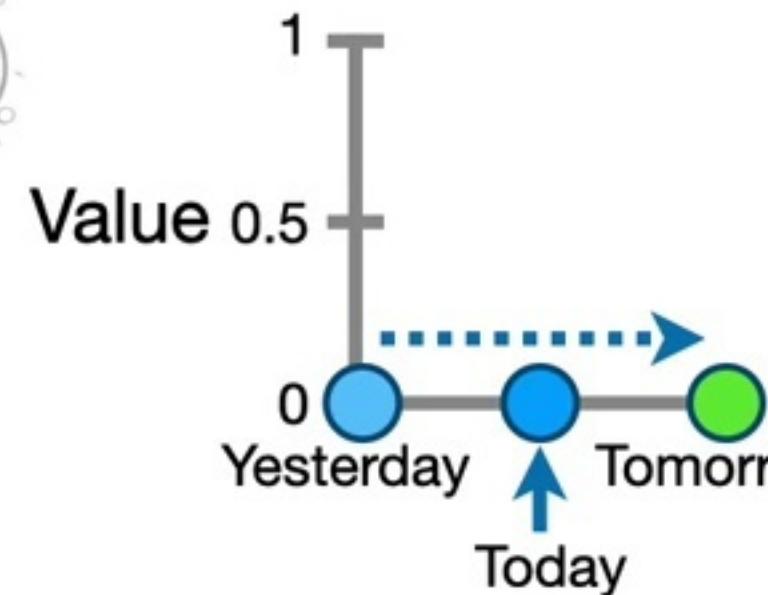


Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.

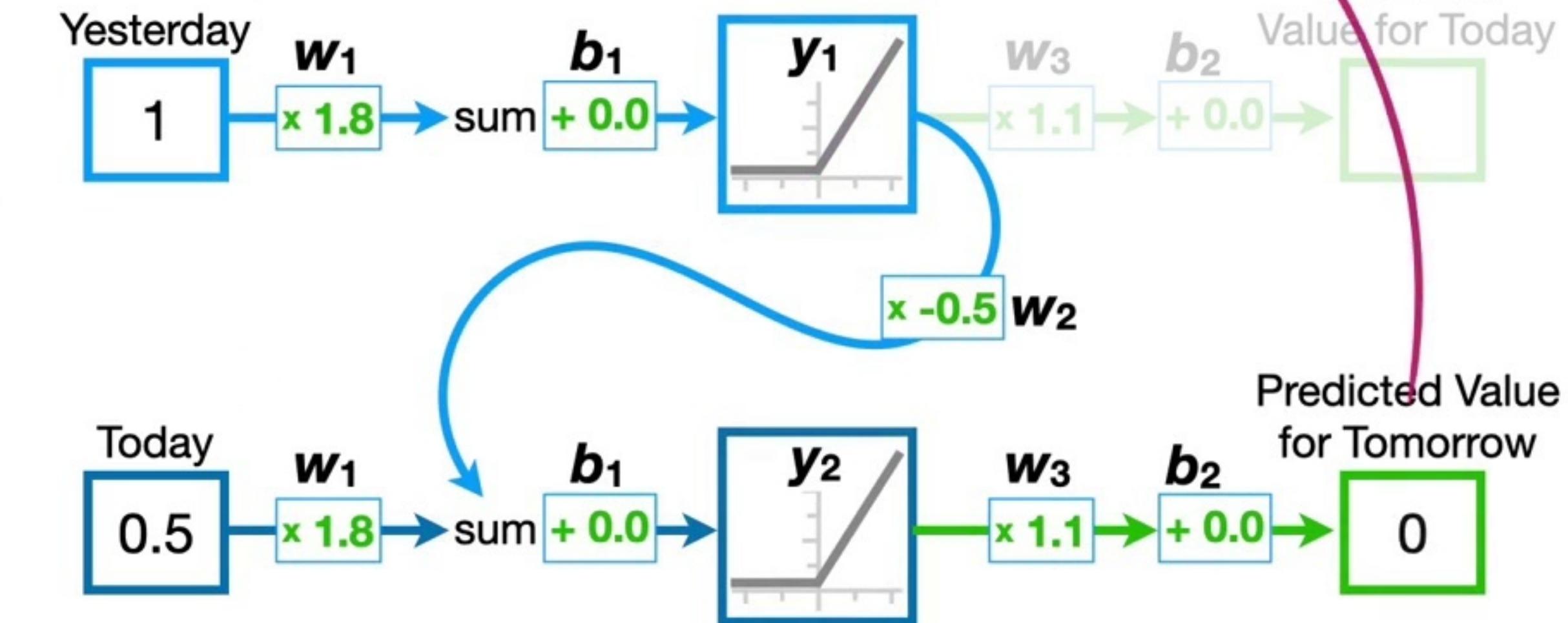


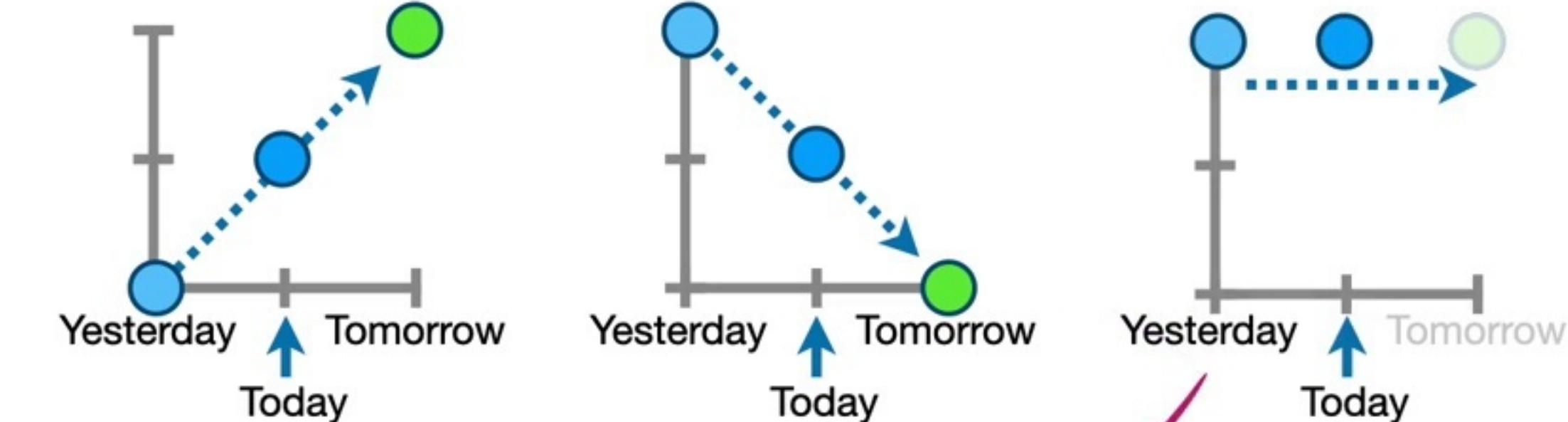
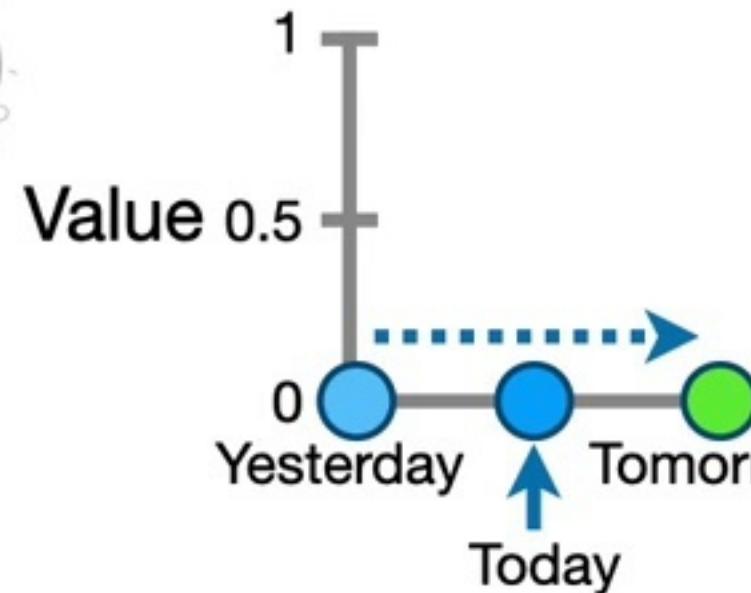


Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.

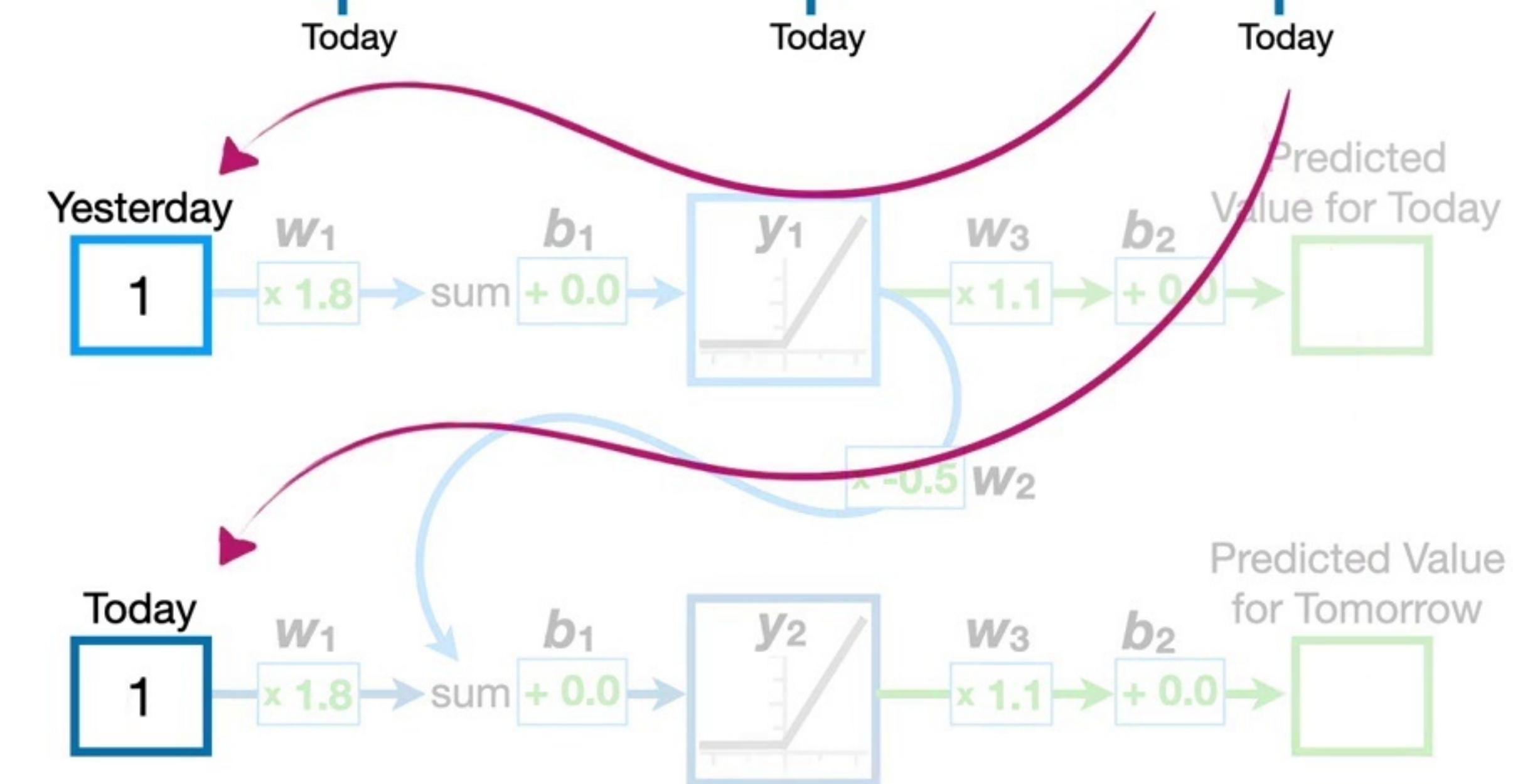


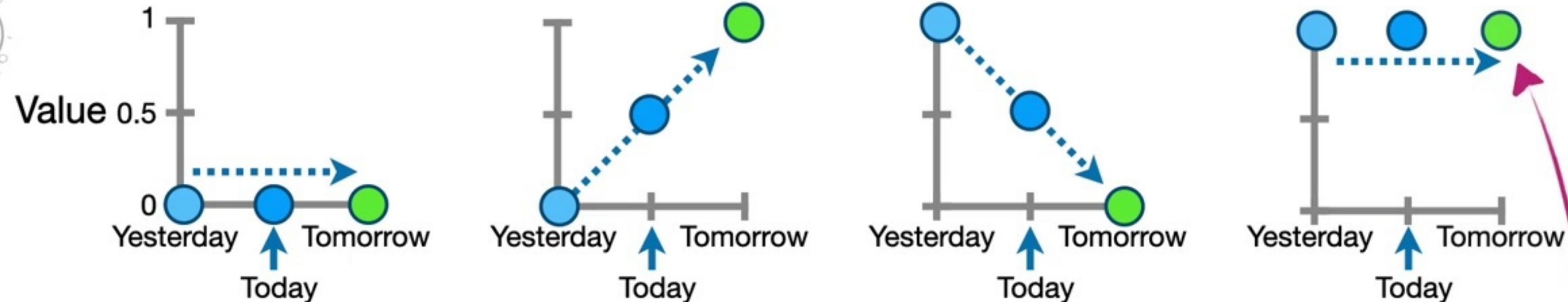
Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.



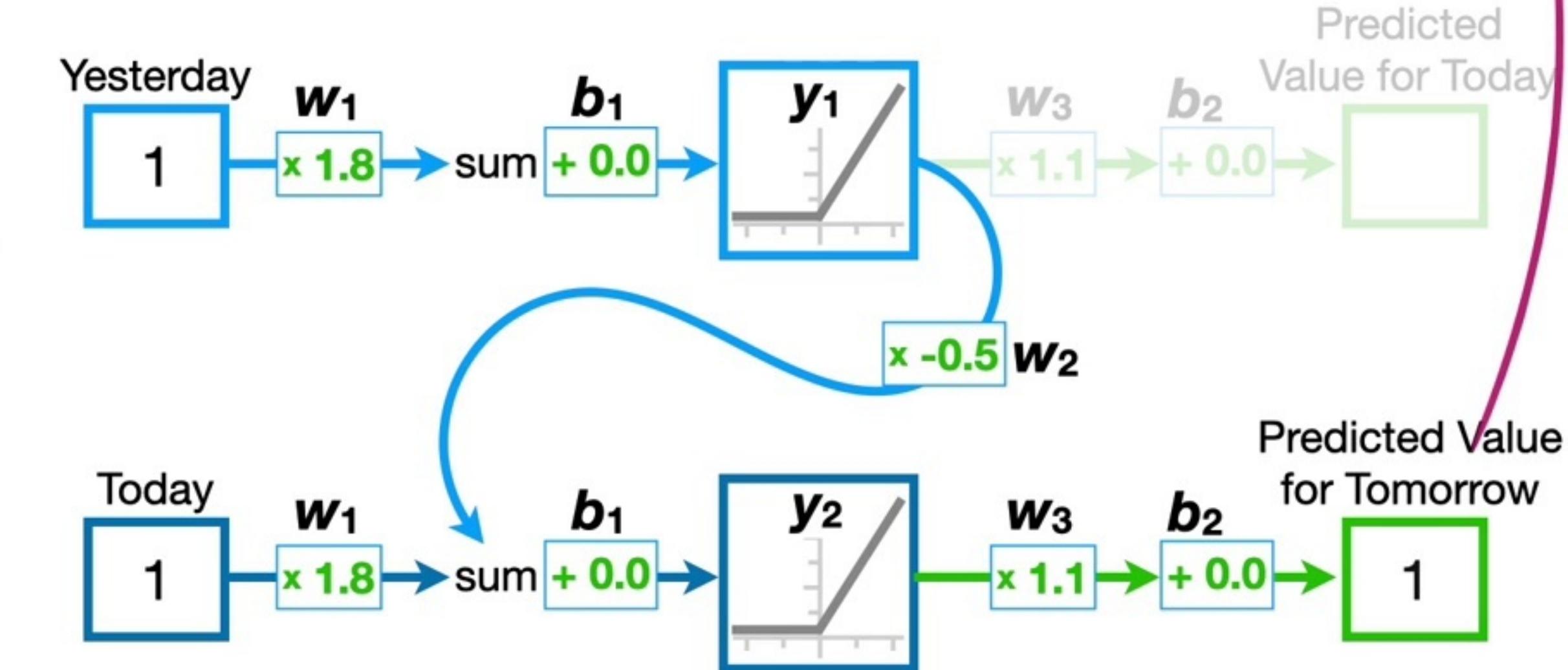


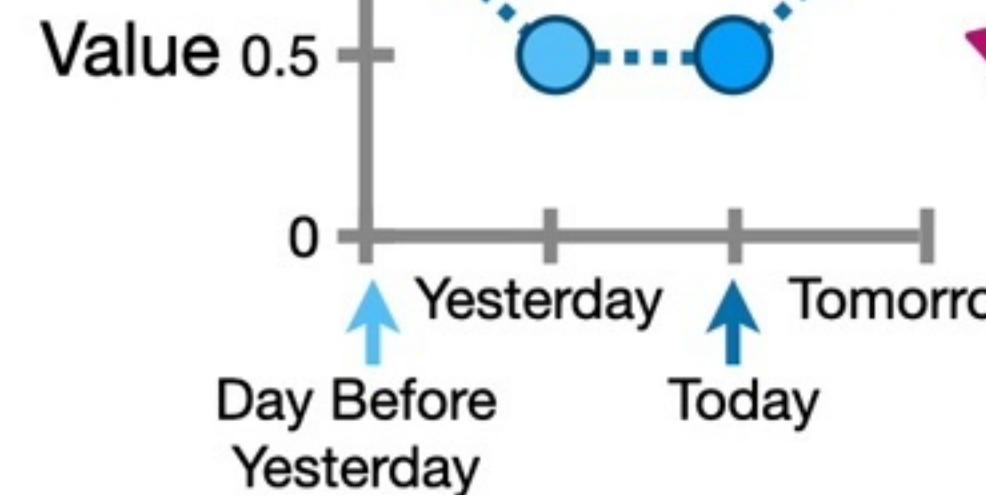
Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.



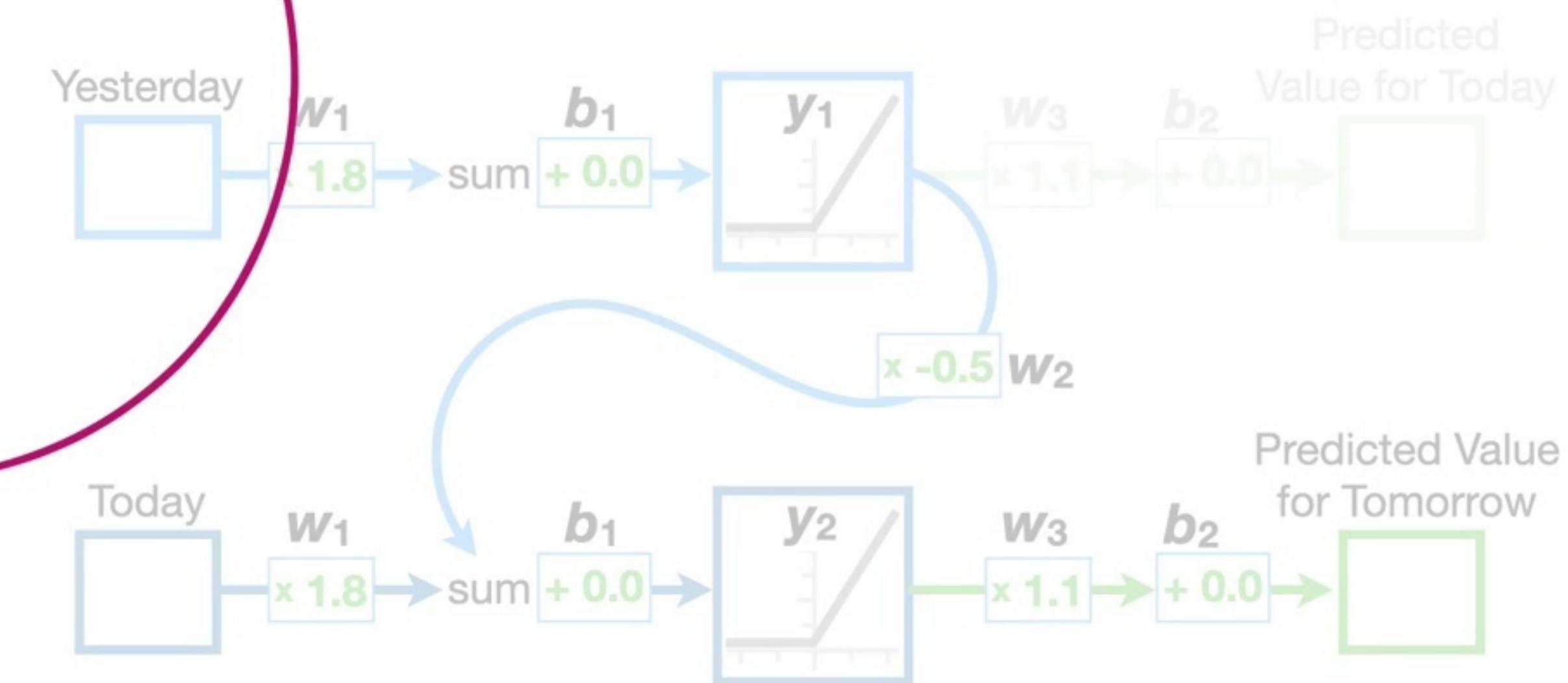


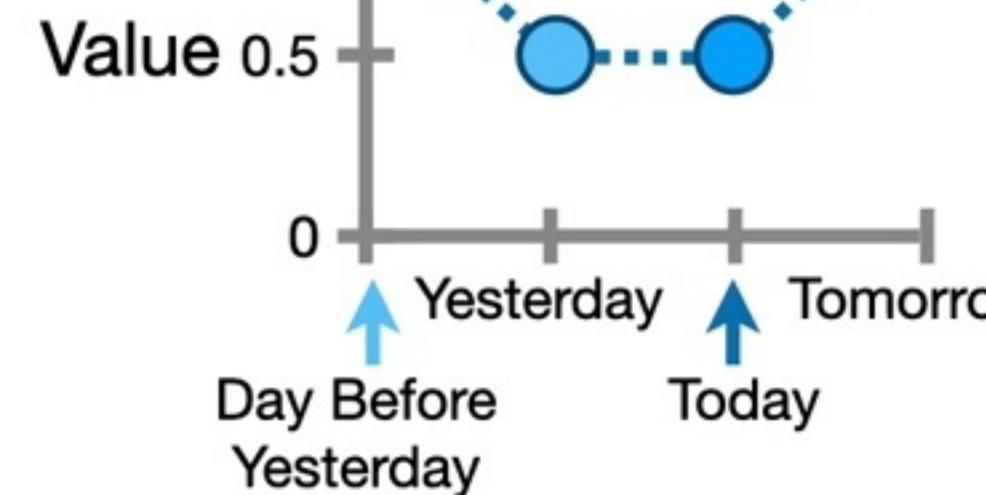
Likewise, when we run **yesterday** and **today's** values for the other scenarios through the recurrent neural network we predict the correct values for **tomorrow**.



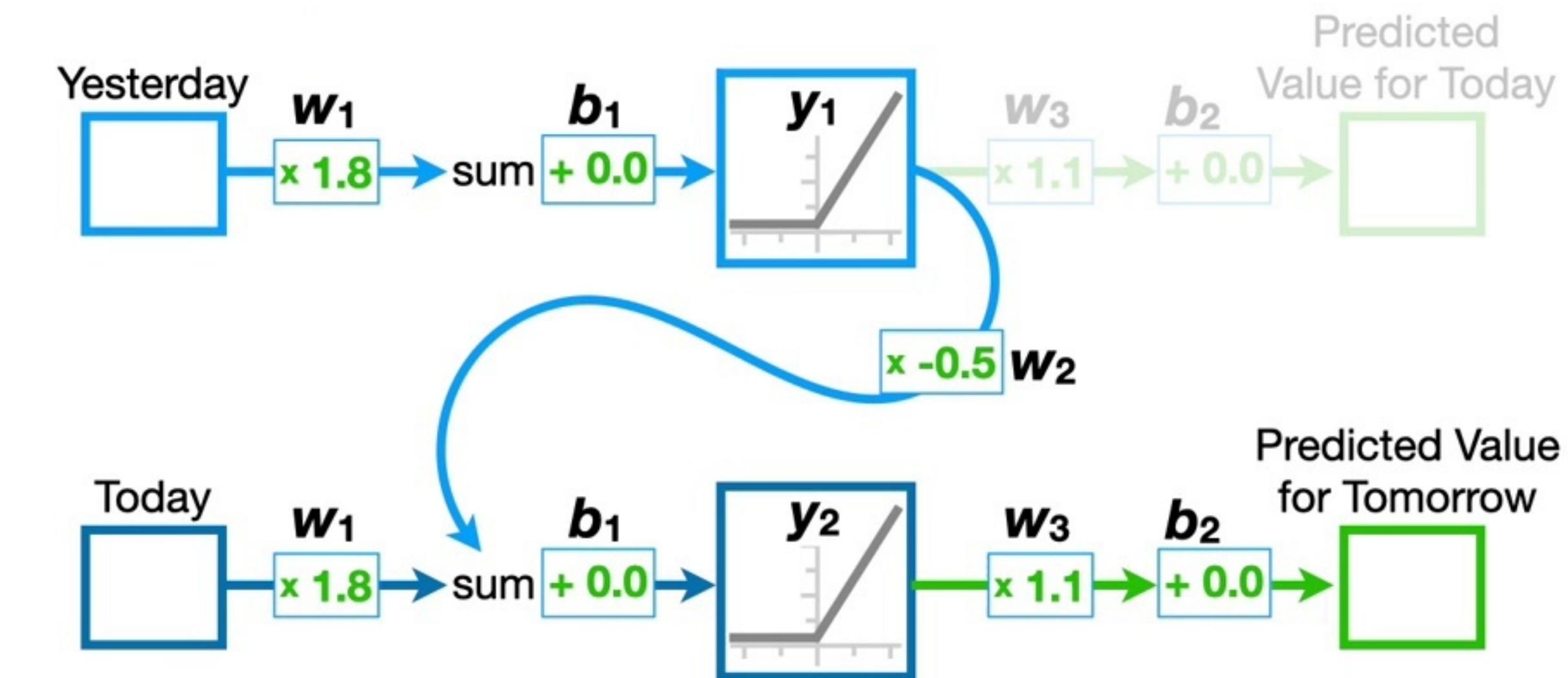


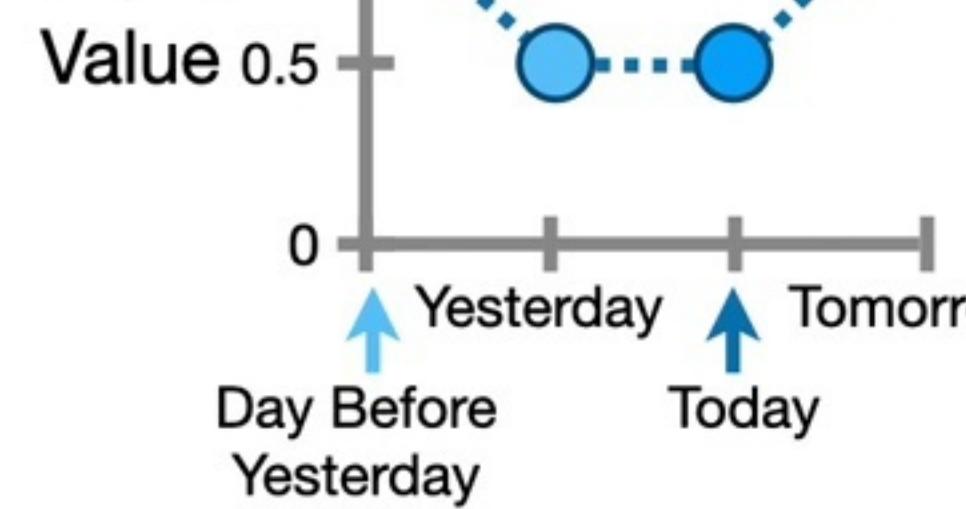
When we want to
use **3** days of data to
make a prediction
about **tomorrow's**
price, like this...



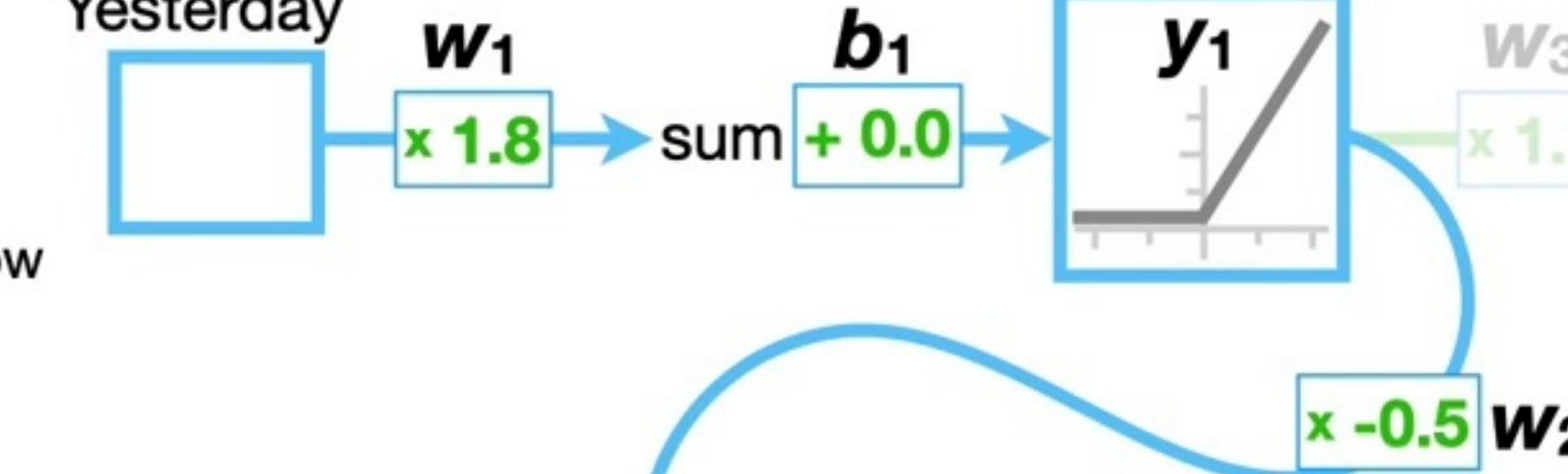


...then we just keep
unrolling the
recurrent neural
network until we
have an input for
each day of data.

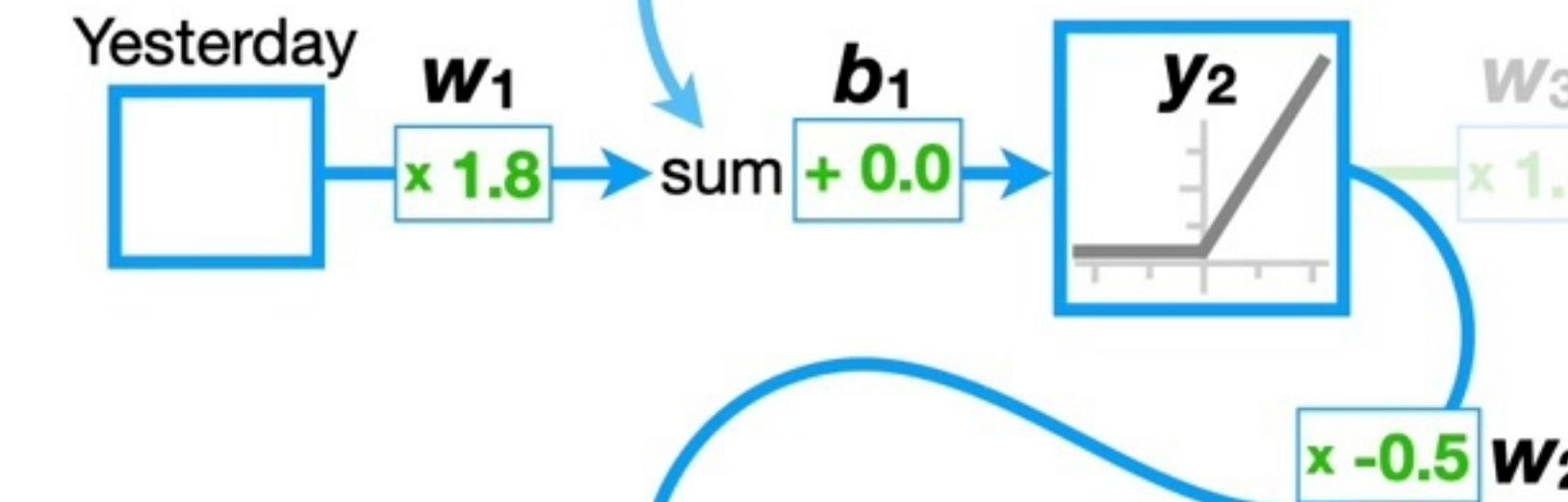




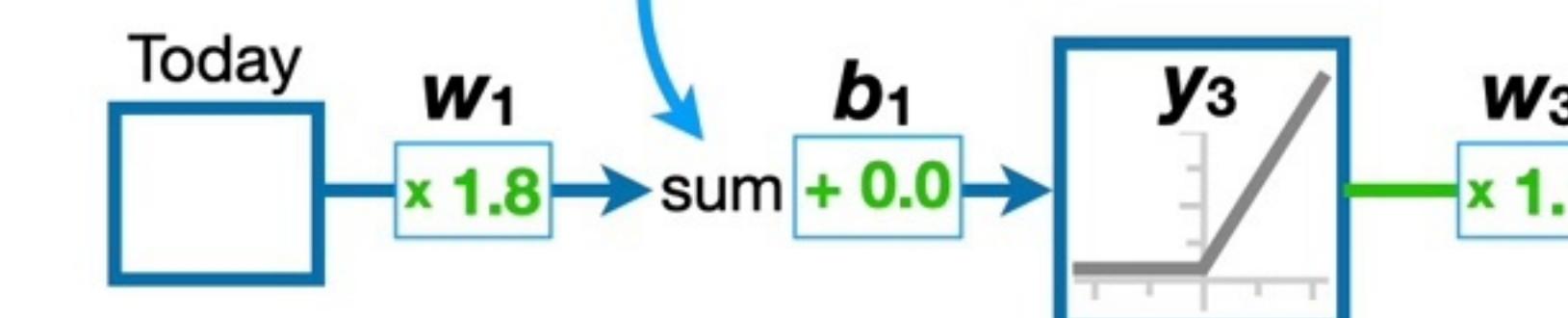
Day Before
Yesterday



Predicted Value
for Yesterday

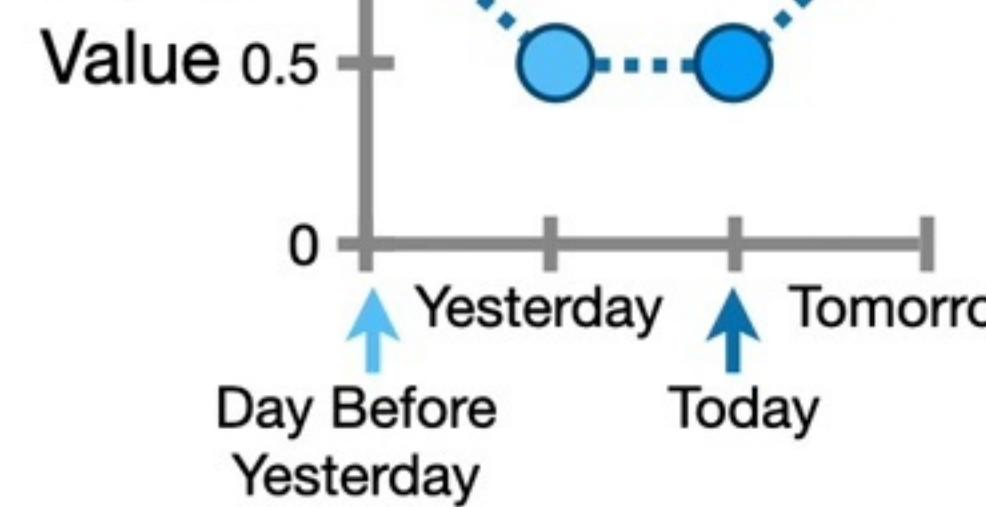


Predicted
Value for Today



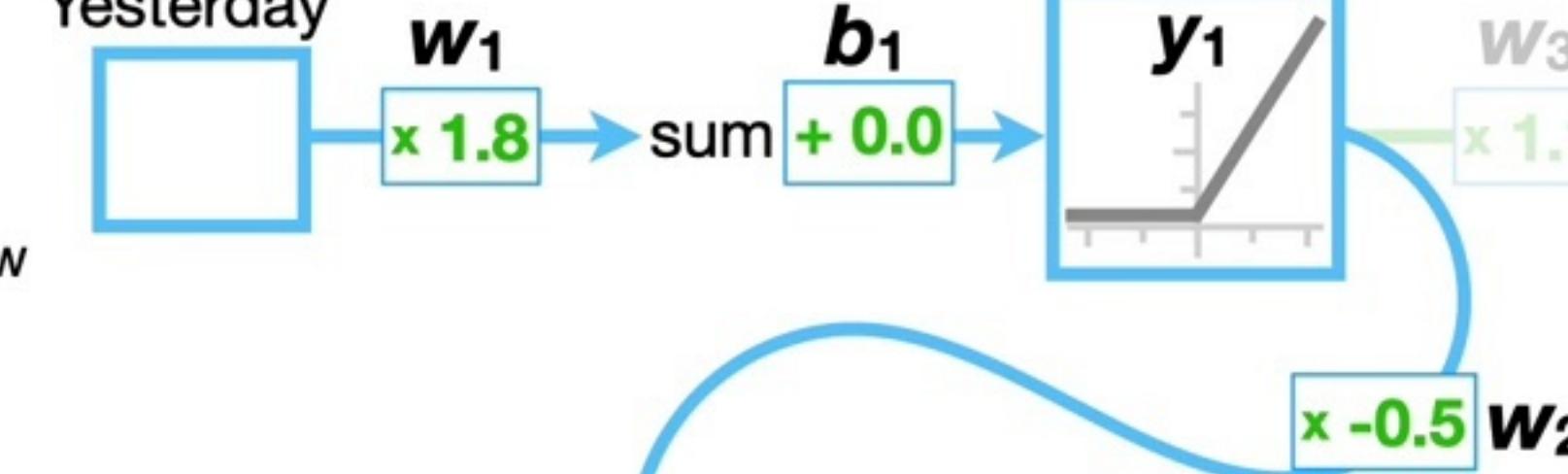
Predicted Value
for Tomorrow

...then we just keep
unrolling the
recurrent neural
network until we
have an input for
each day of data.

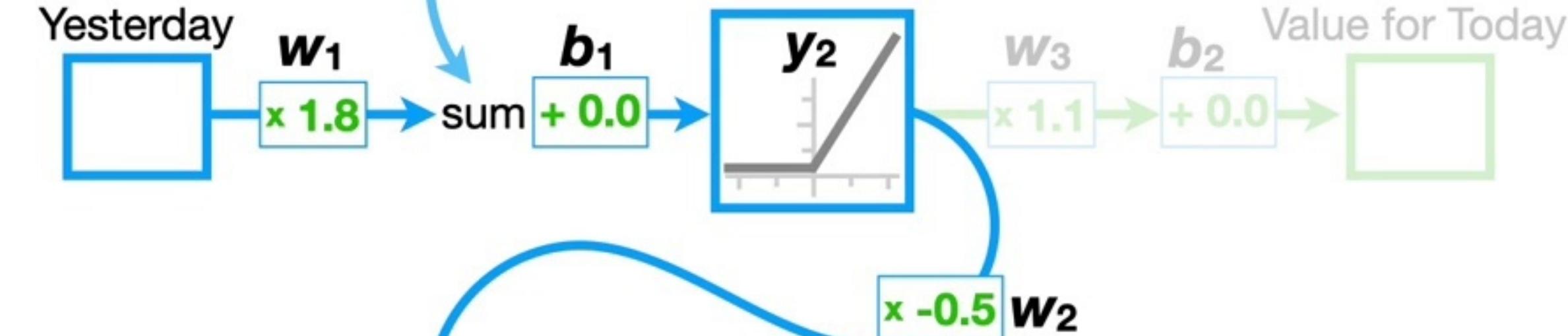


Then we plug the values into the inputs, always from oldest to newest.

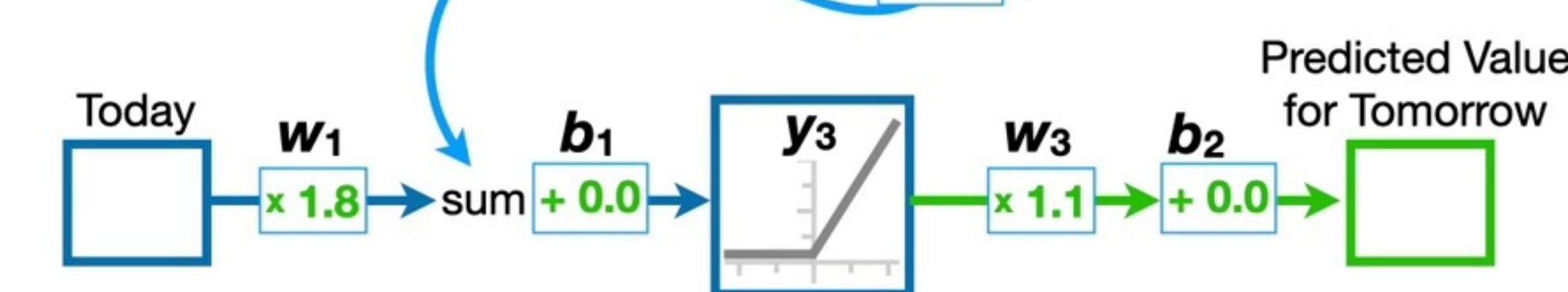
Day Before
Yesterday



Yesterday



Today





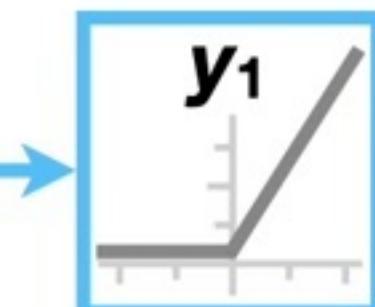
In this case, that means we start by plugging in the value for the **day before yesterday**...

Day Before
Yesterday



sum

+ 0.0



$x -0.5 w_2$

Predicted Value
for Yesterday



Predicted
Value

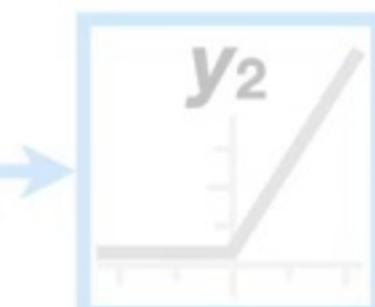
Value for Today

Yesterday



sum

+ 0.0



$x -0.5 w_2$

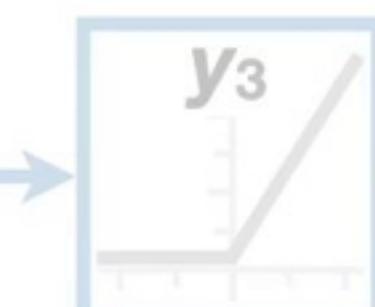
Predicted Value
for Tomorrow

Today



sum

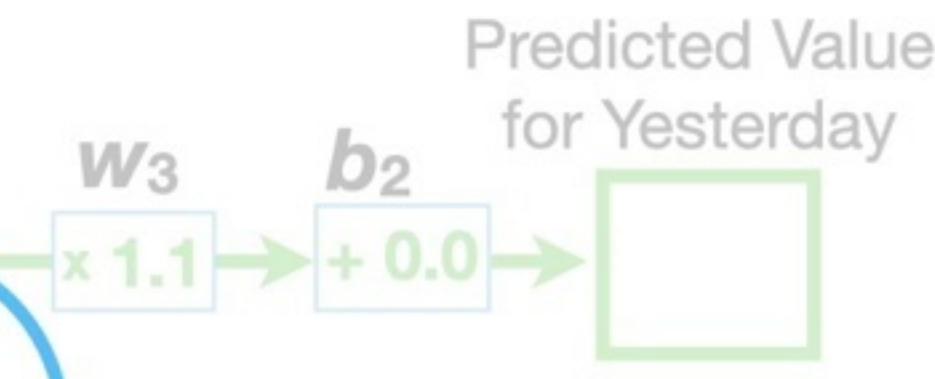
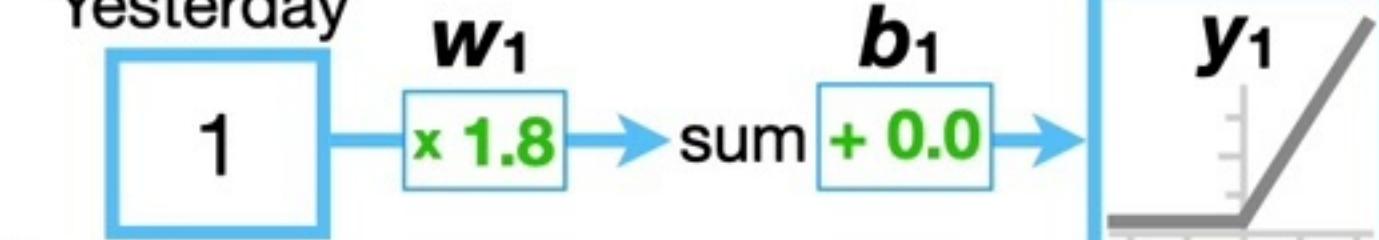
+ 0.0



$x -0.5 w_2$

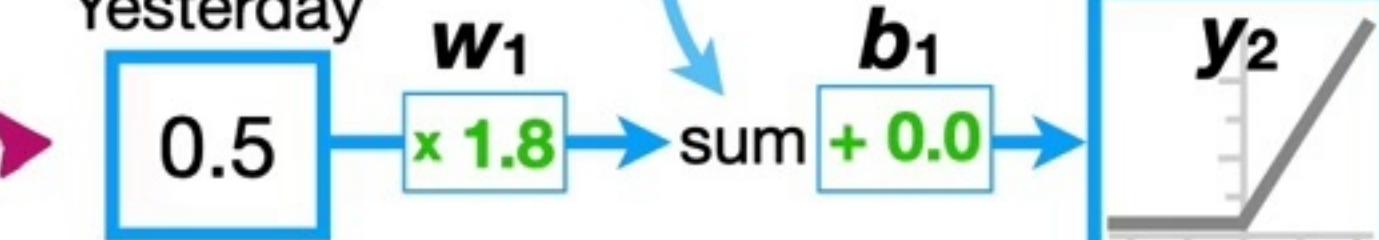


Day Before
Yesterday



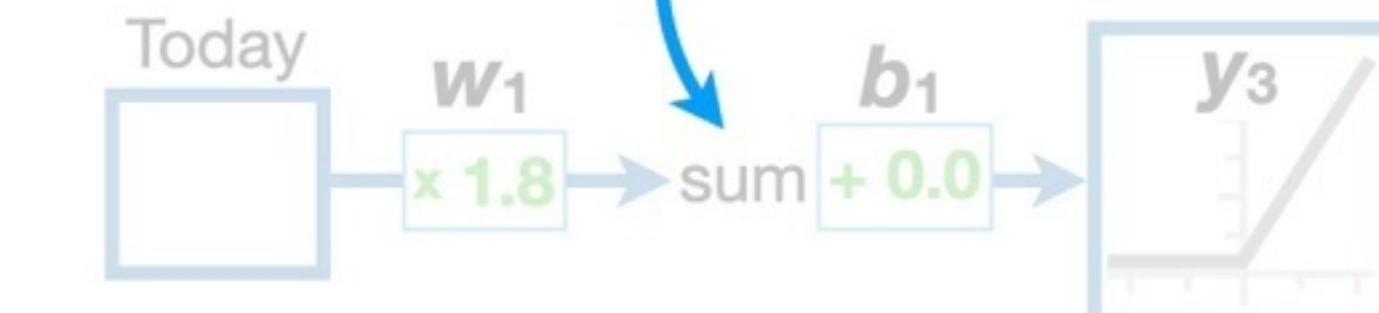
Predicted Value
for Yesterday

Yesterday

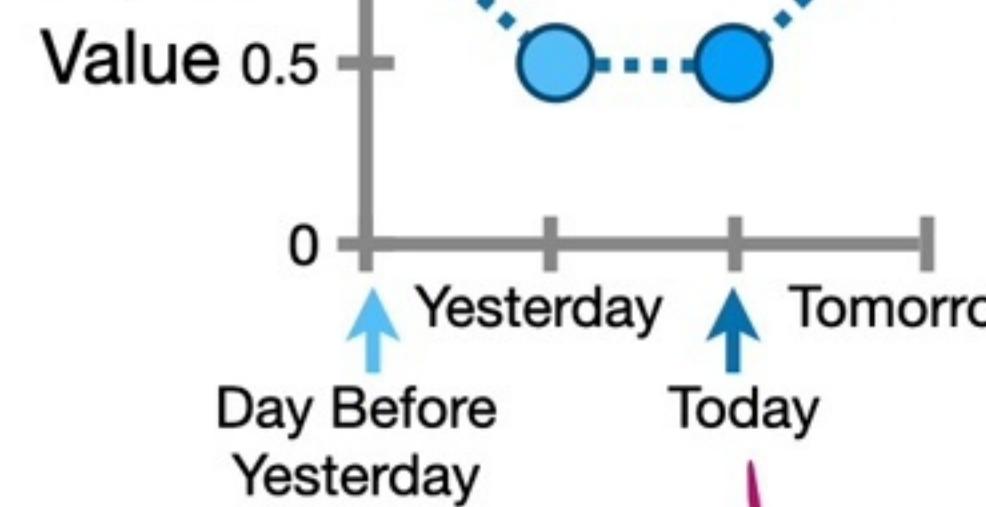


Predicted
Value for Today

...then we plug in
yesterday's value...



Predicted Value
for Tomorrow



Day Before
Yesterday

$$w_1 \times 1.8 + b_1 + 0.0 = y_1$$

$$w_3 \times 1.1 + b_2 + 0.0 = \text{Predicted Value for Yesterday}$$

$$w_1 \times 1.8 + b_1 + 0.0 = y_2$$

$$w_3 \times 1.1 + b_2 + 0.0 = \text{Predicted Value for Today}$$

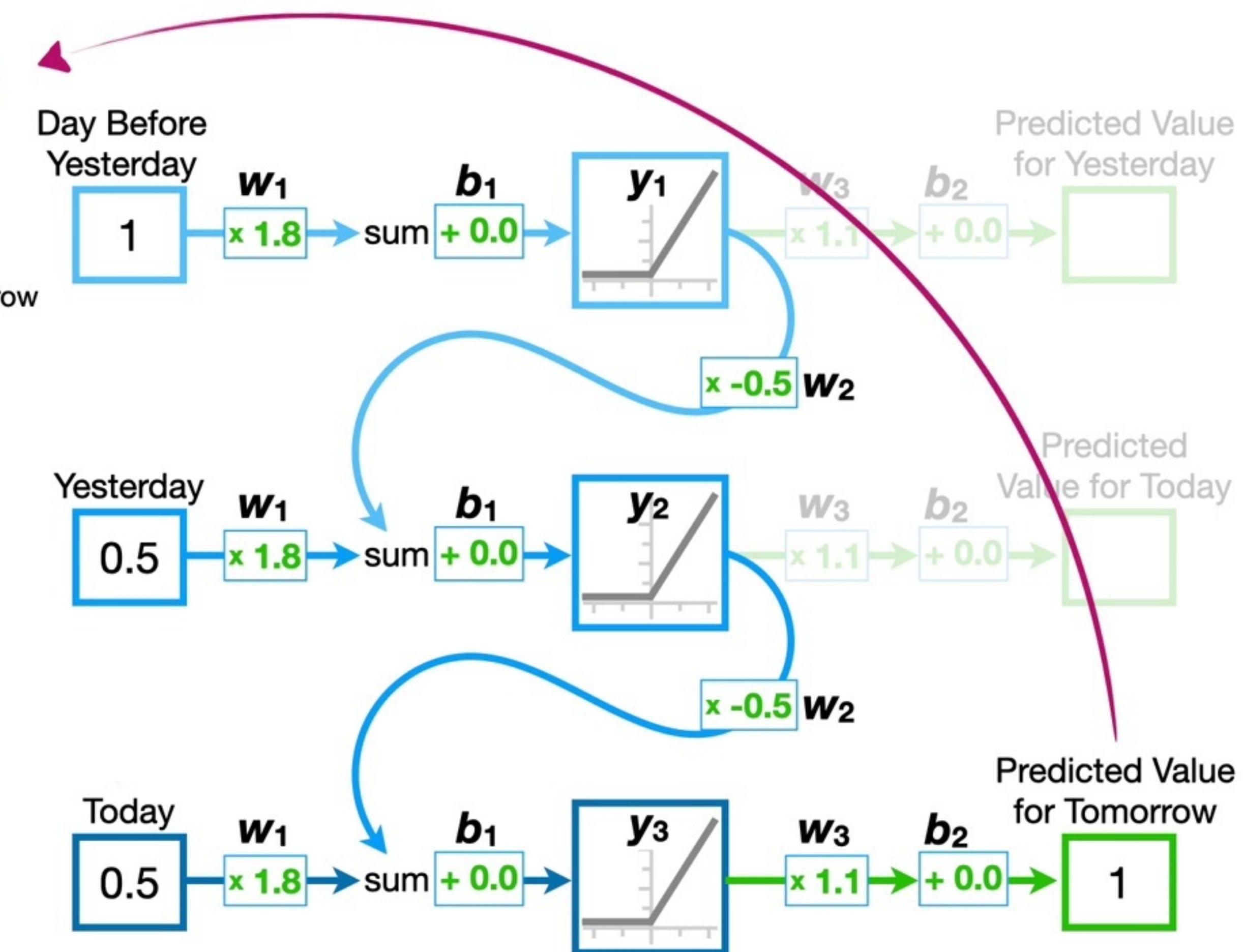
$$w_1 \times 1.8 + b_1 + 0.0 = y_3$$

$$w_3 \times 1.1 + b_2 + 0.0 = \text{Predicted Value for Tomorrow}$$

...and then we plug
in **today's** value.

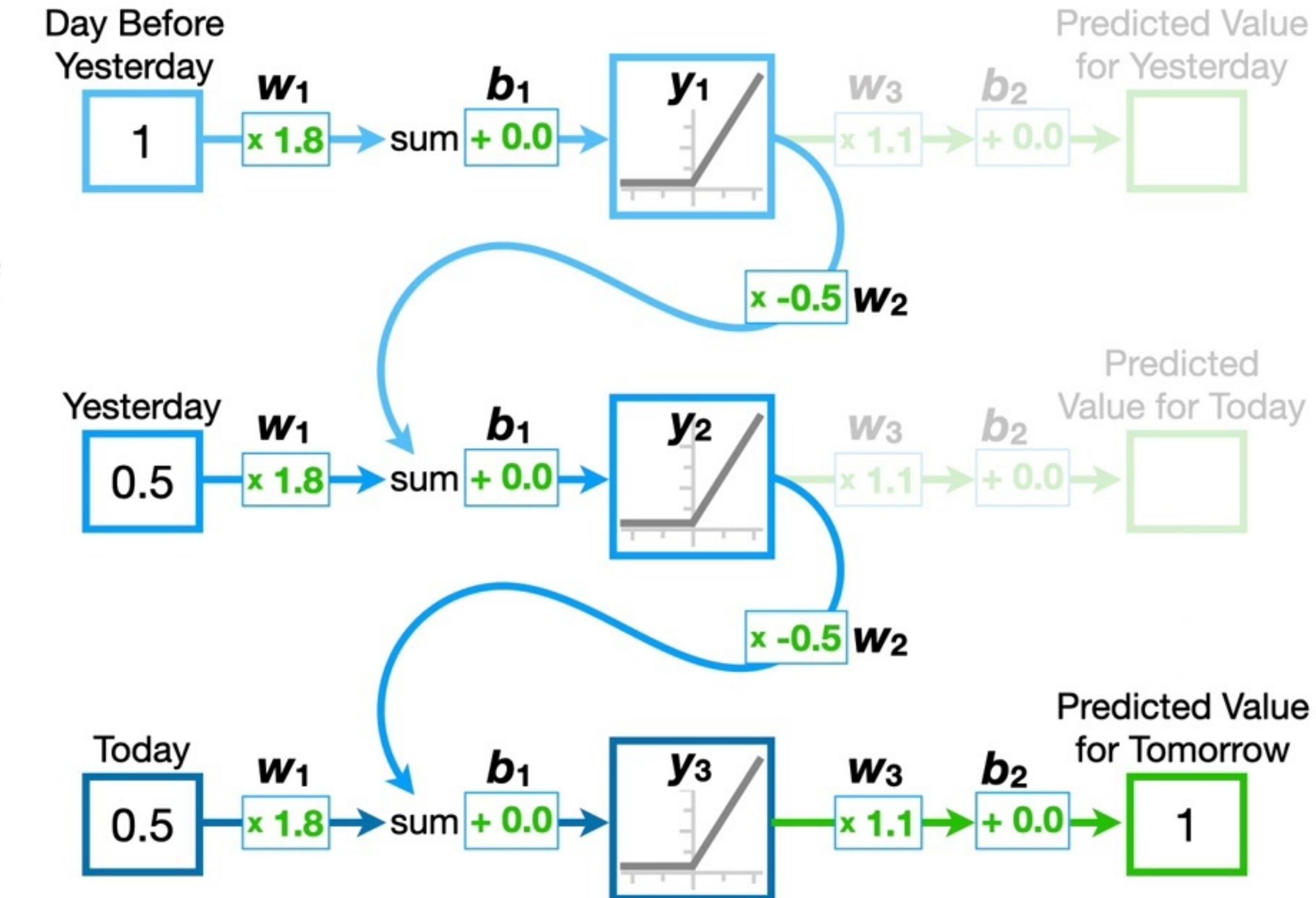


And when we do the math, the last output gives us the prediction for **tomorrow**.



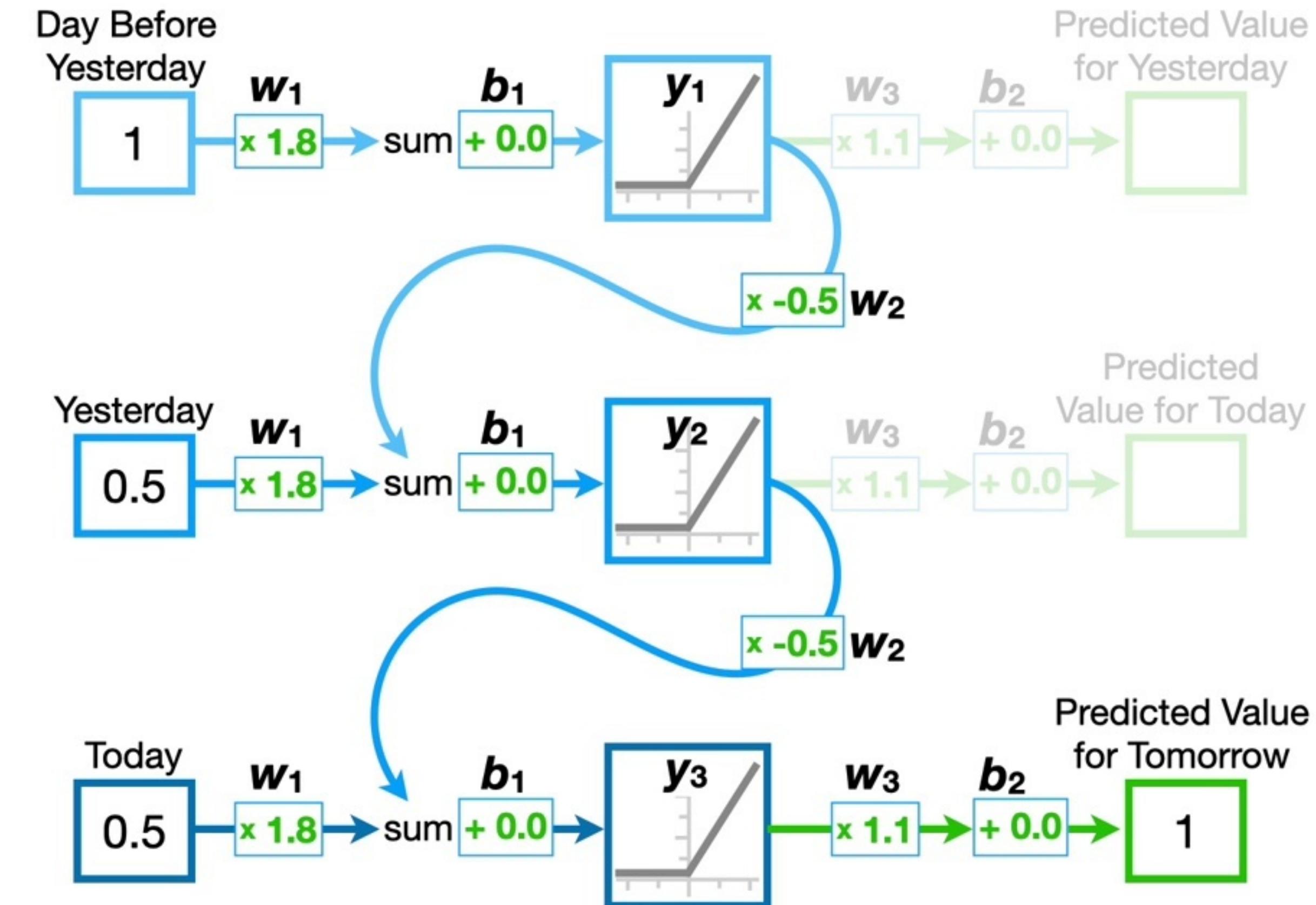


NOTE: Regardless of how many times we **unroll** a recurrent neural network, the **weights** and **biases** are shared across every input.



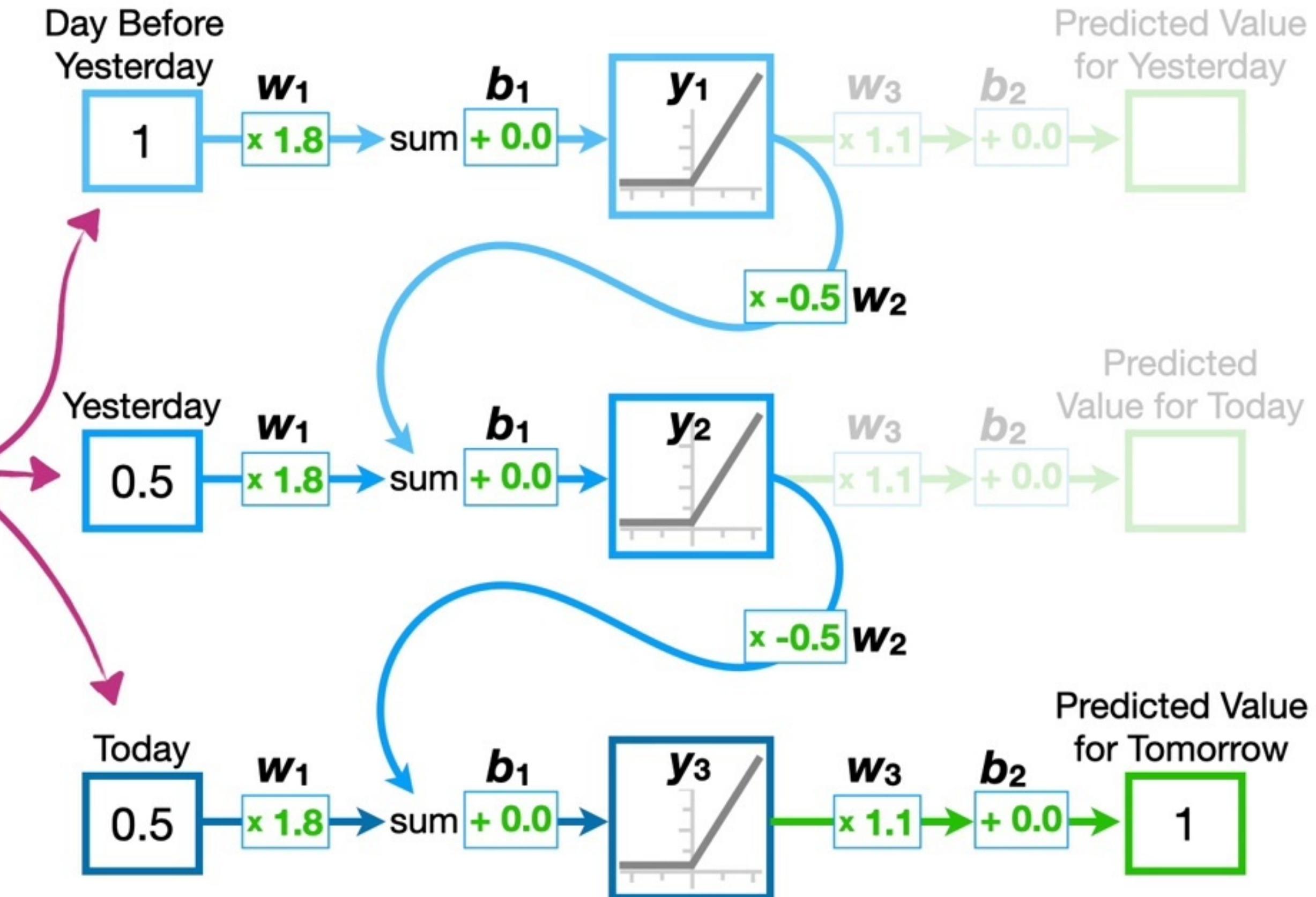


In other words, even though this **unrolled** network has **3** inputs...





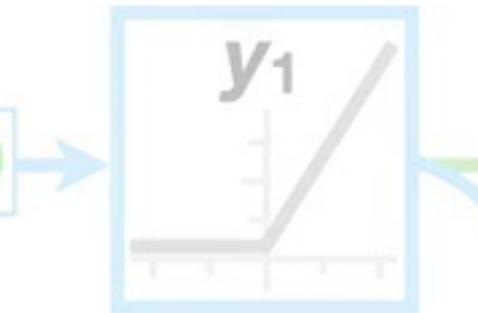
In other words, even though this **unrolled** network has **3** inputs...





...the **weight**, w_1 ,
is the same for all
3 inputs.

Day Before
Yesterday

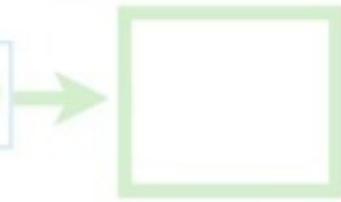
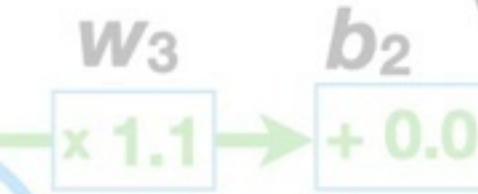
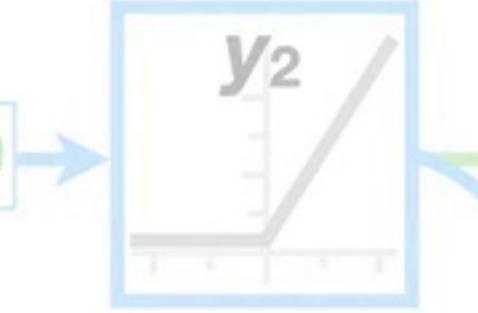


Predicted Value
for Yesterday

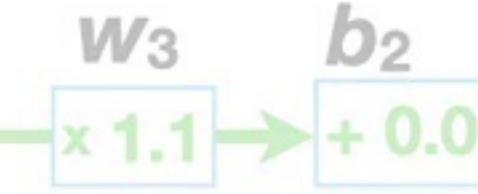
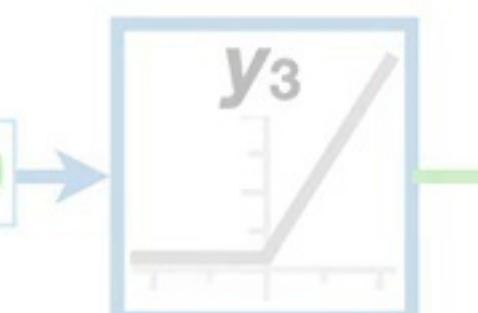
Predicted
Value for Today

Predicted Value
for Tomorrow

Yesterday



Today

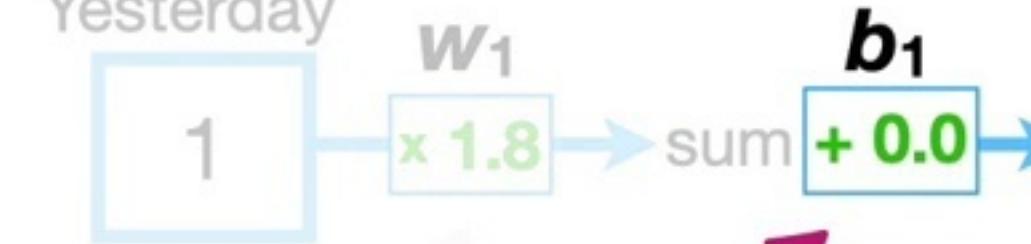


1



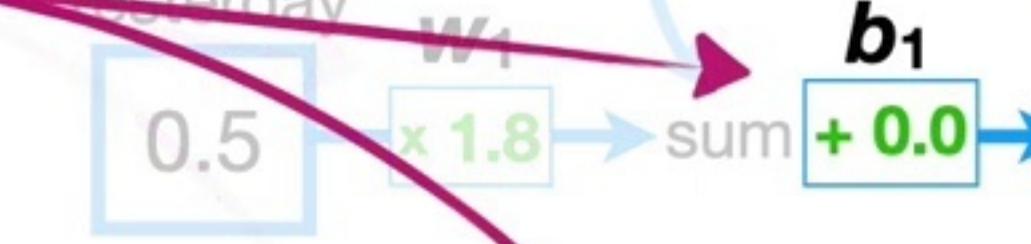
And the **bias**, b_1 ,
is also the same
for all **3** inputs.

Day Before
Yesterday



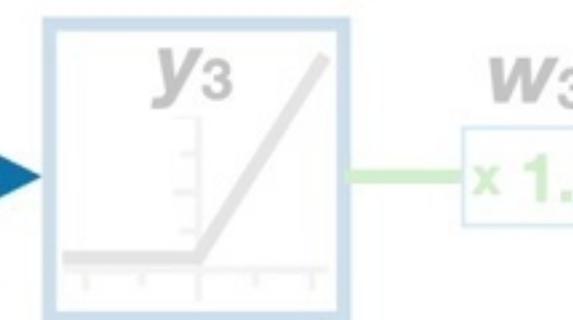
Predicted Value
for Yesterday

Yesterday



Predicted
Value for Today

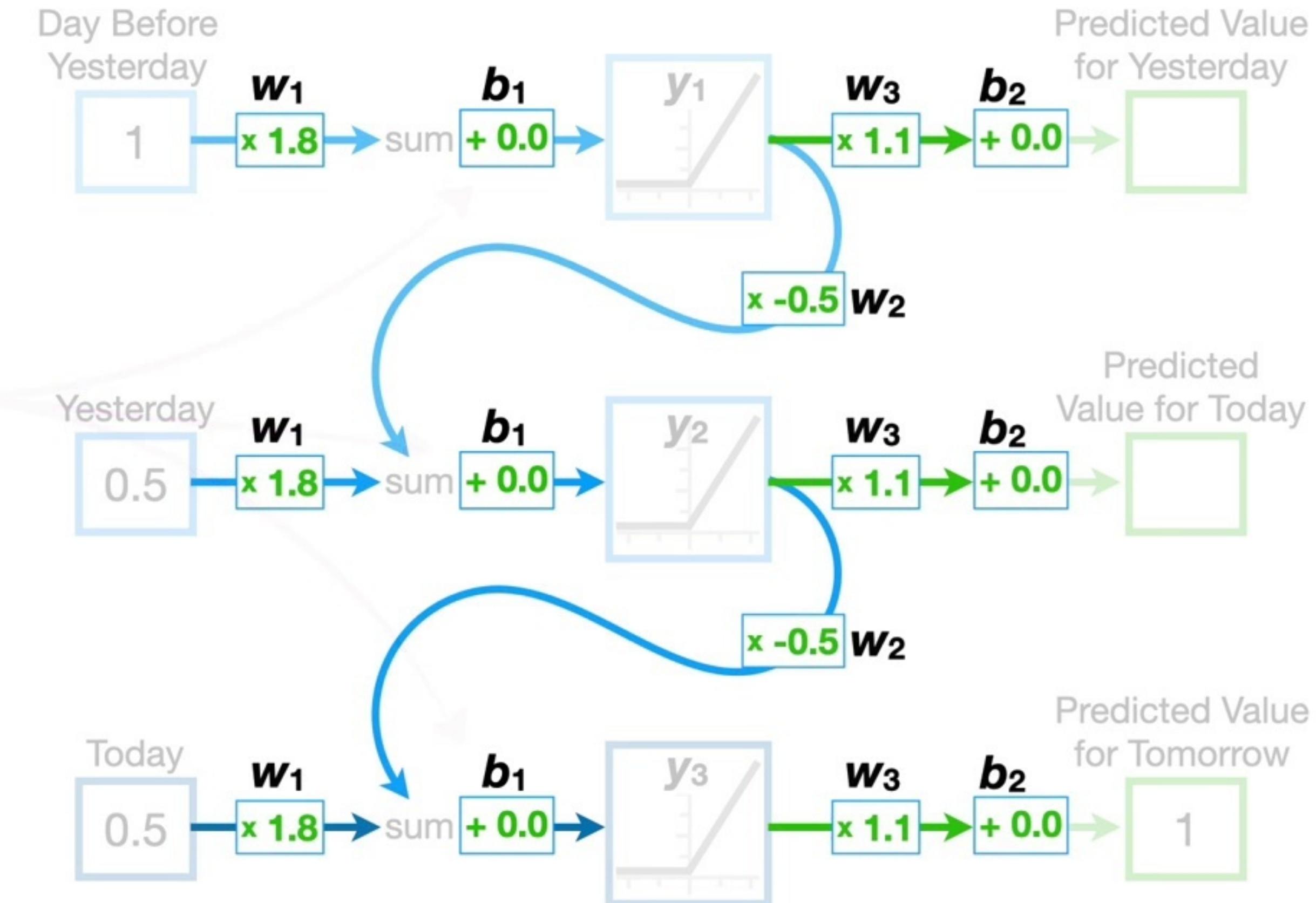
Today



Predicted Value
for Tomorrow

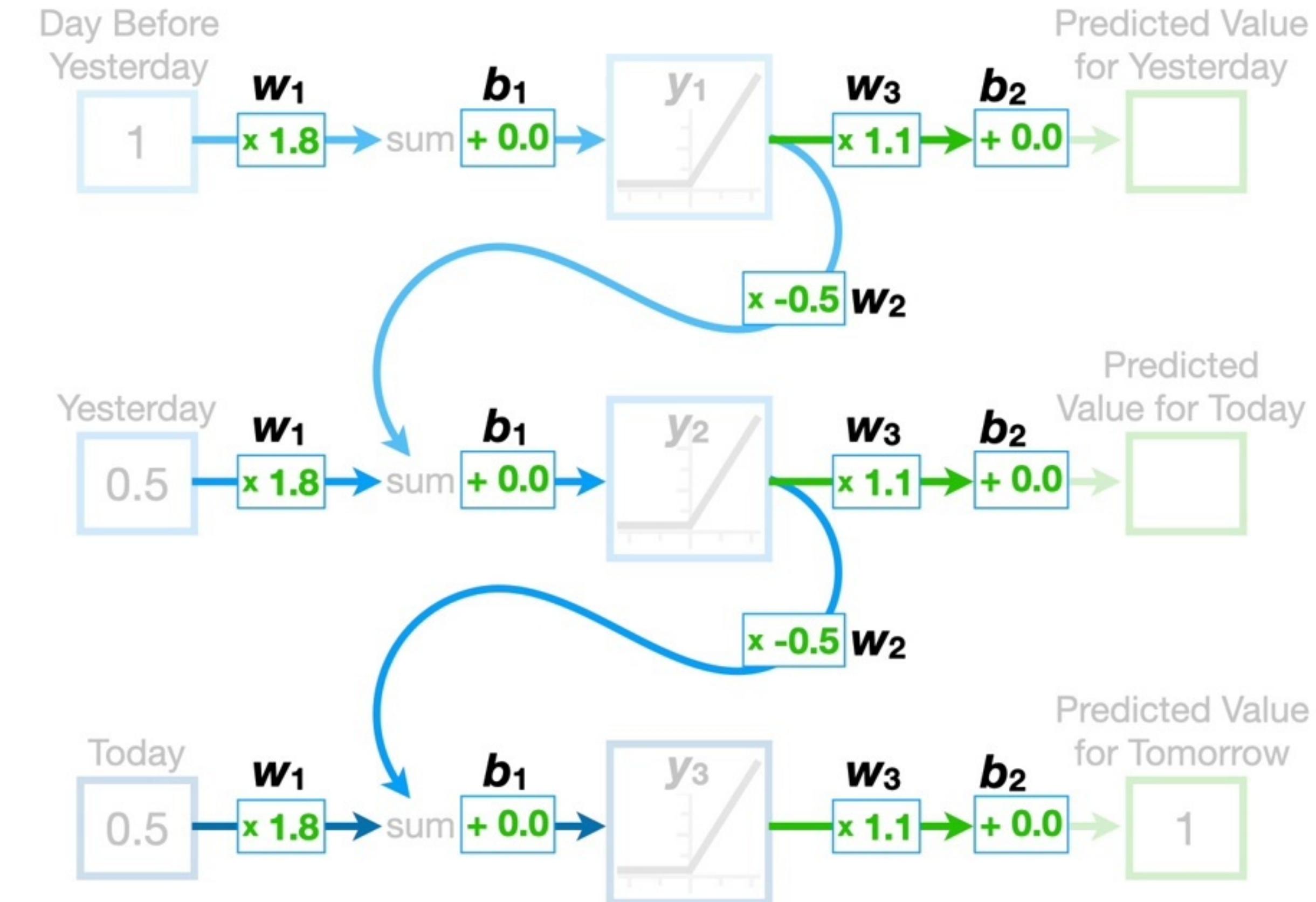


Likewise, all of the other **weights** and **biases** are shared.



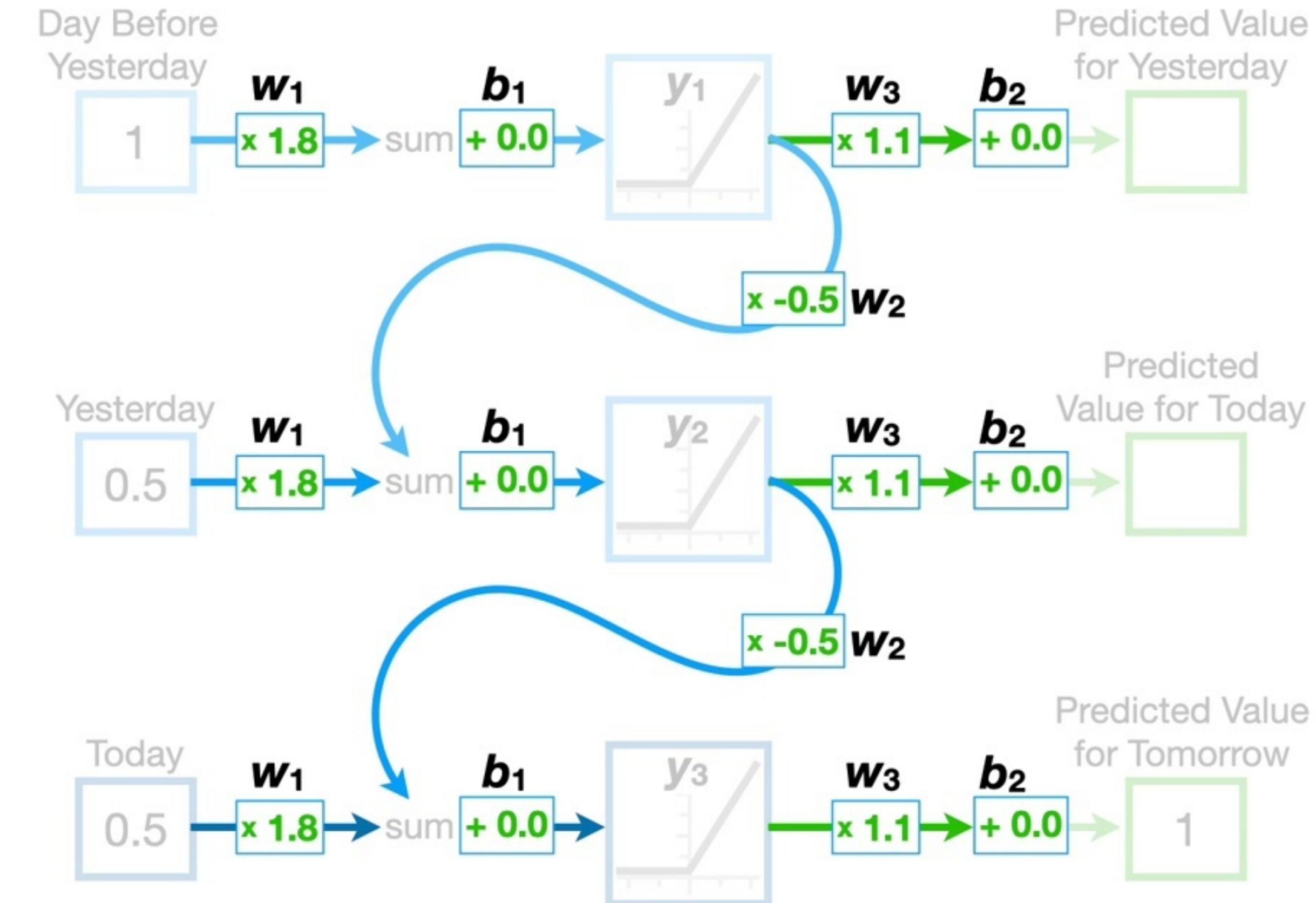


So, no matter how many times we **unroll** a recurrent neural network, we never increase the number of **weights** and **biases** that we have to train.



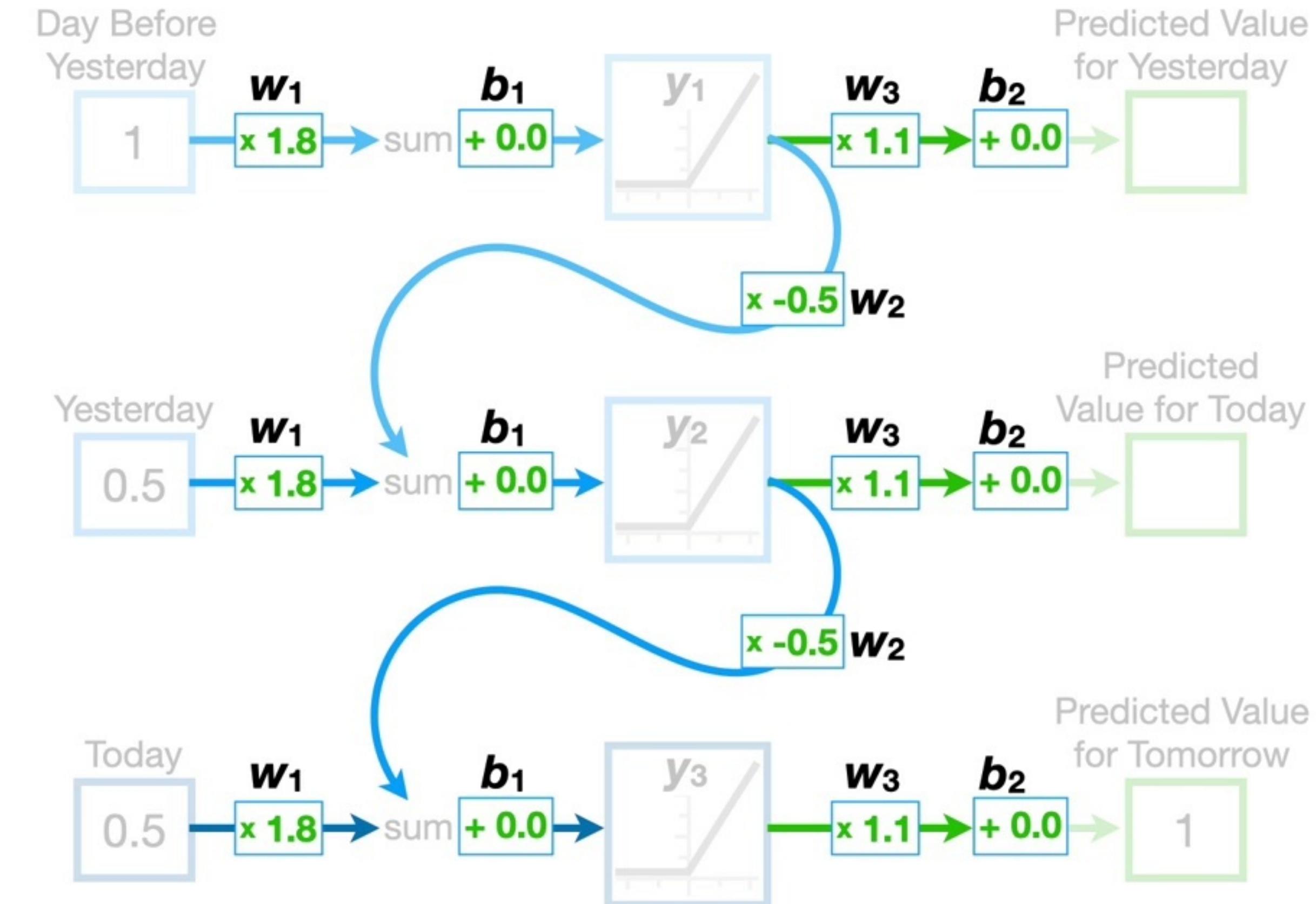


OK, now that we've talked about what makes **basic** recurrent neural networks so cool...



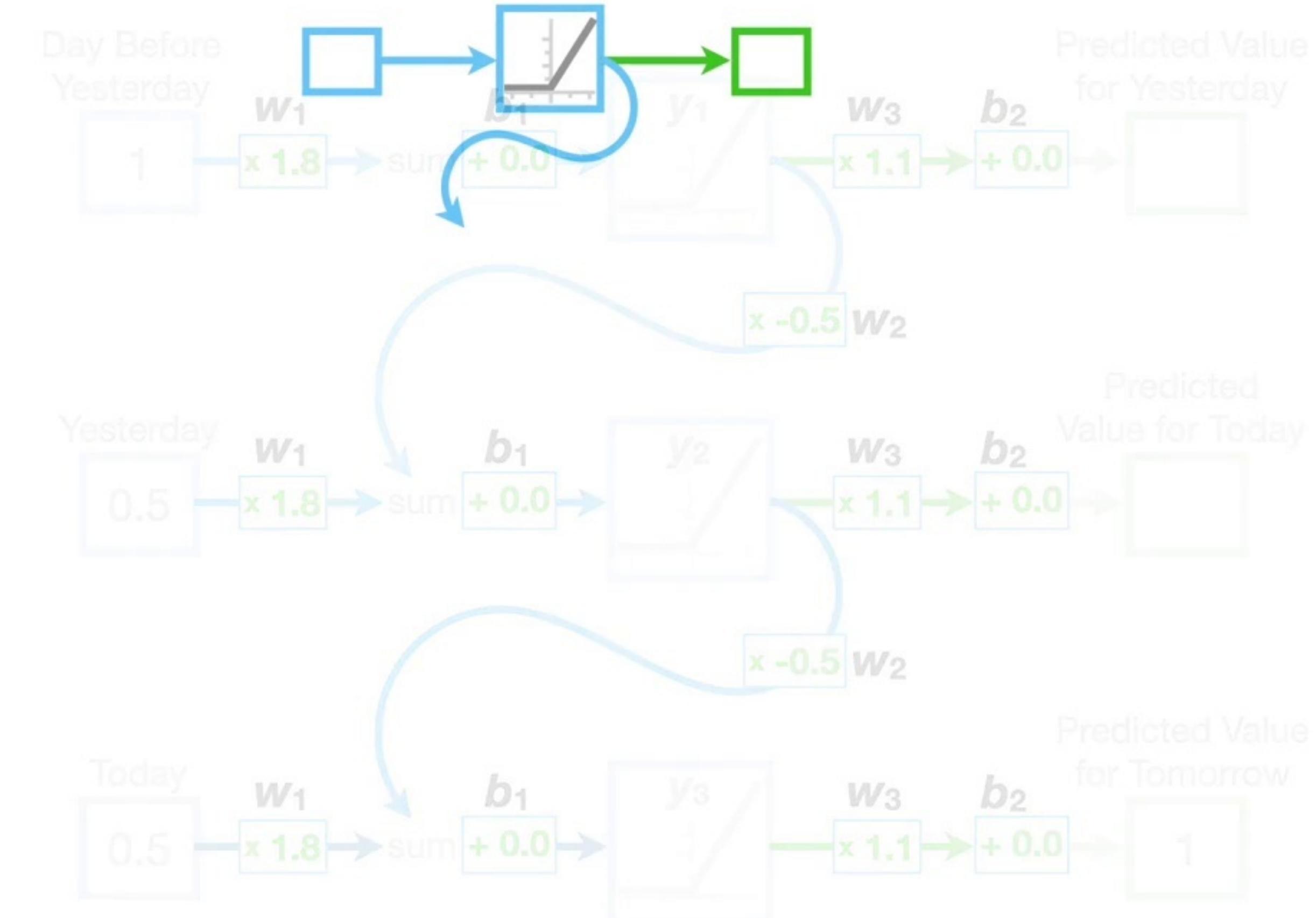


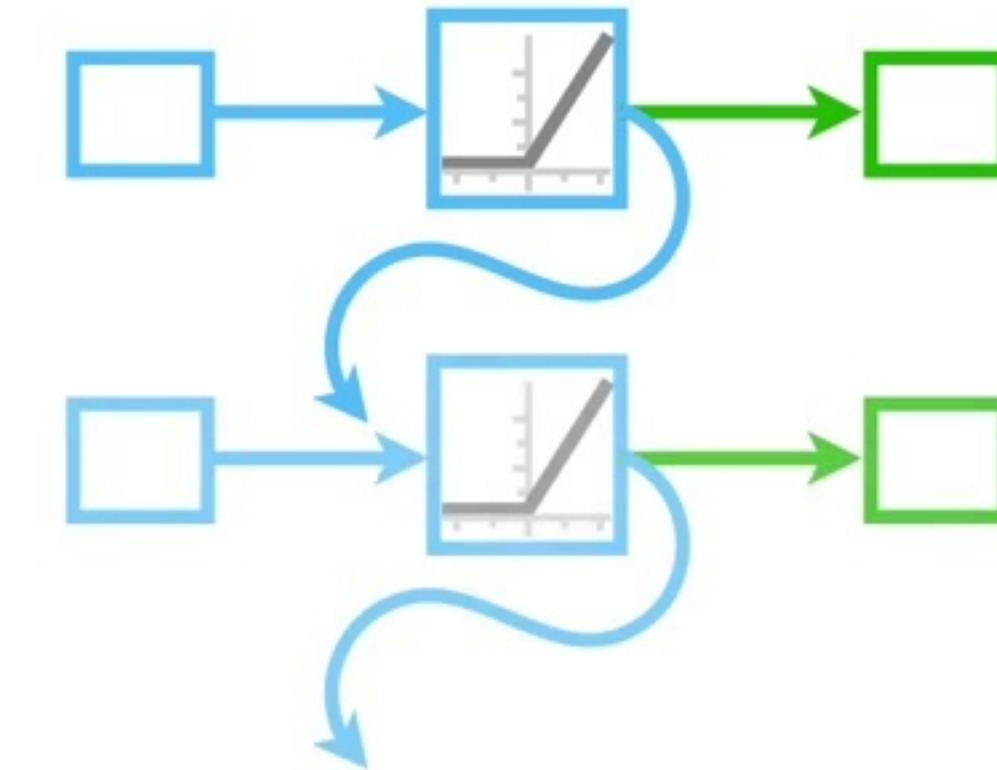
...let's briefly talk
about why they are
not used very often.





One big problem is
that the more we
unroll a recurrent
neural network, the
harder it is to train.

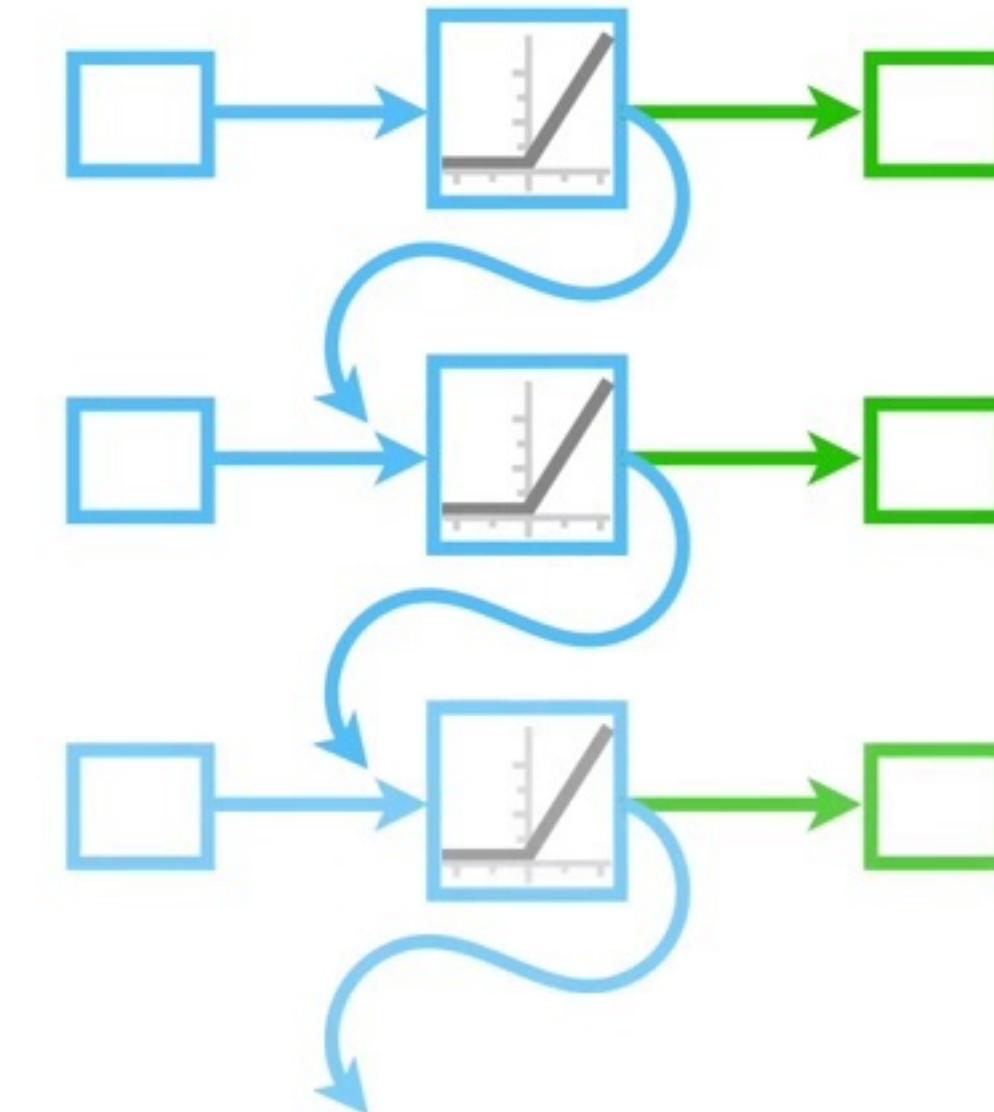




One big problem is
that the more we
unroll a recurrent
neural network, the
harder it is to train.

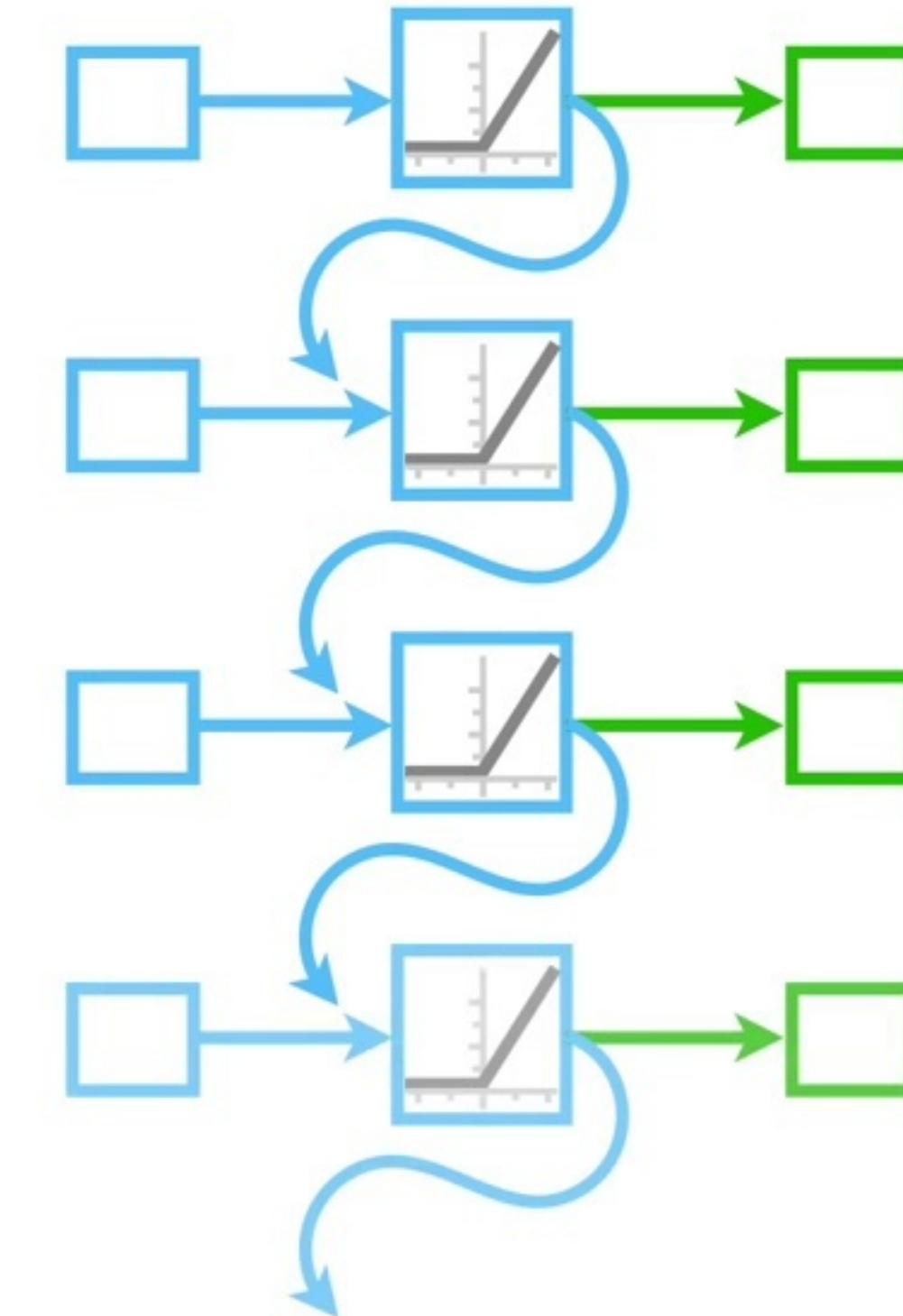


One big problem is
that the more we
unroll a recurrent
neural network, the
harder it is to train.



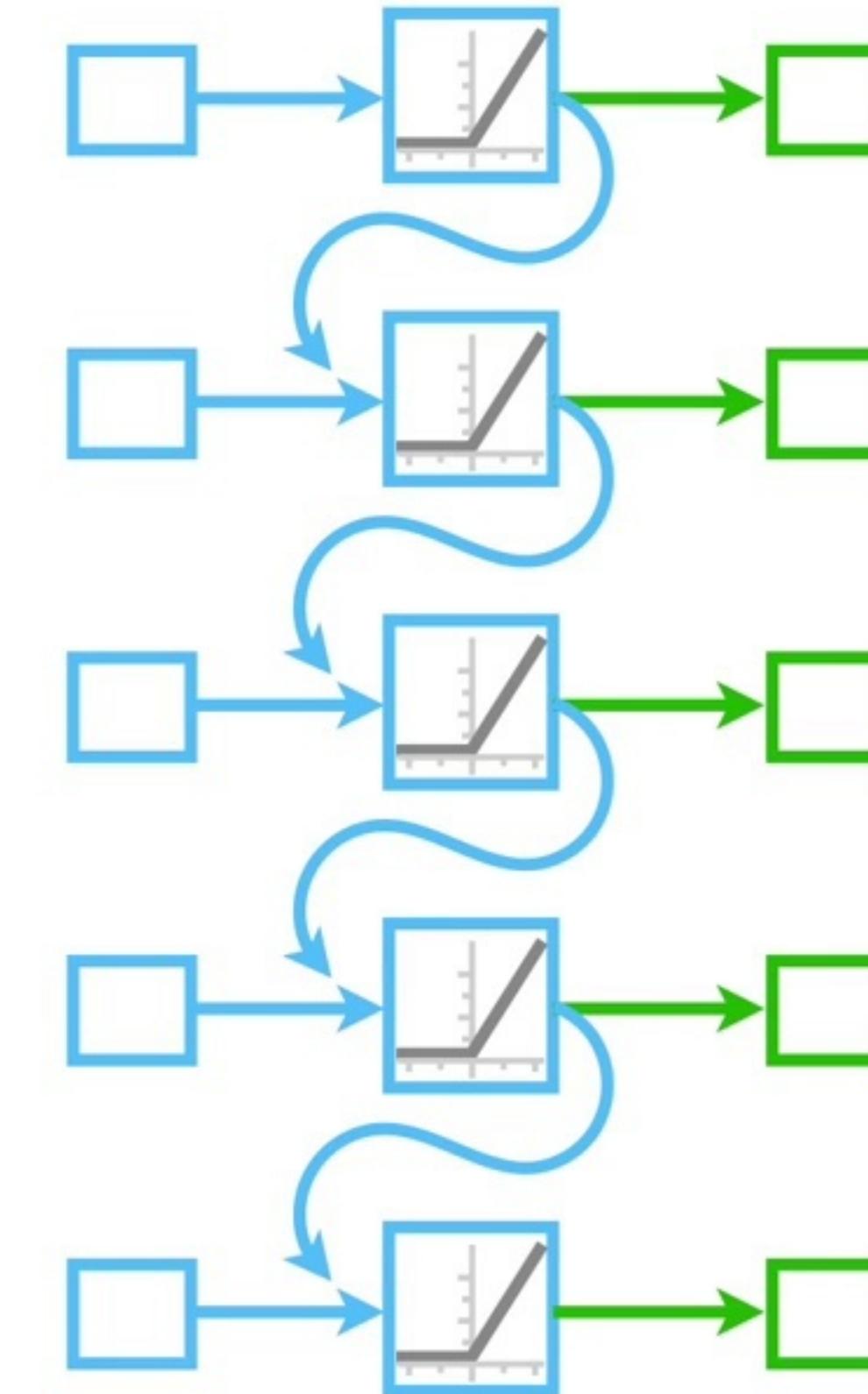


One big problem is
that the more we
unroll a recurrent
neural network, the
harder it is to train.



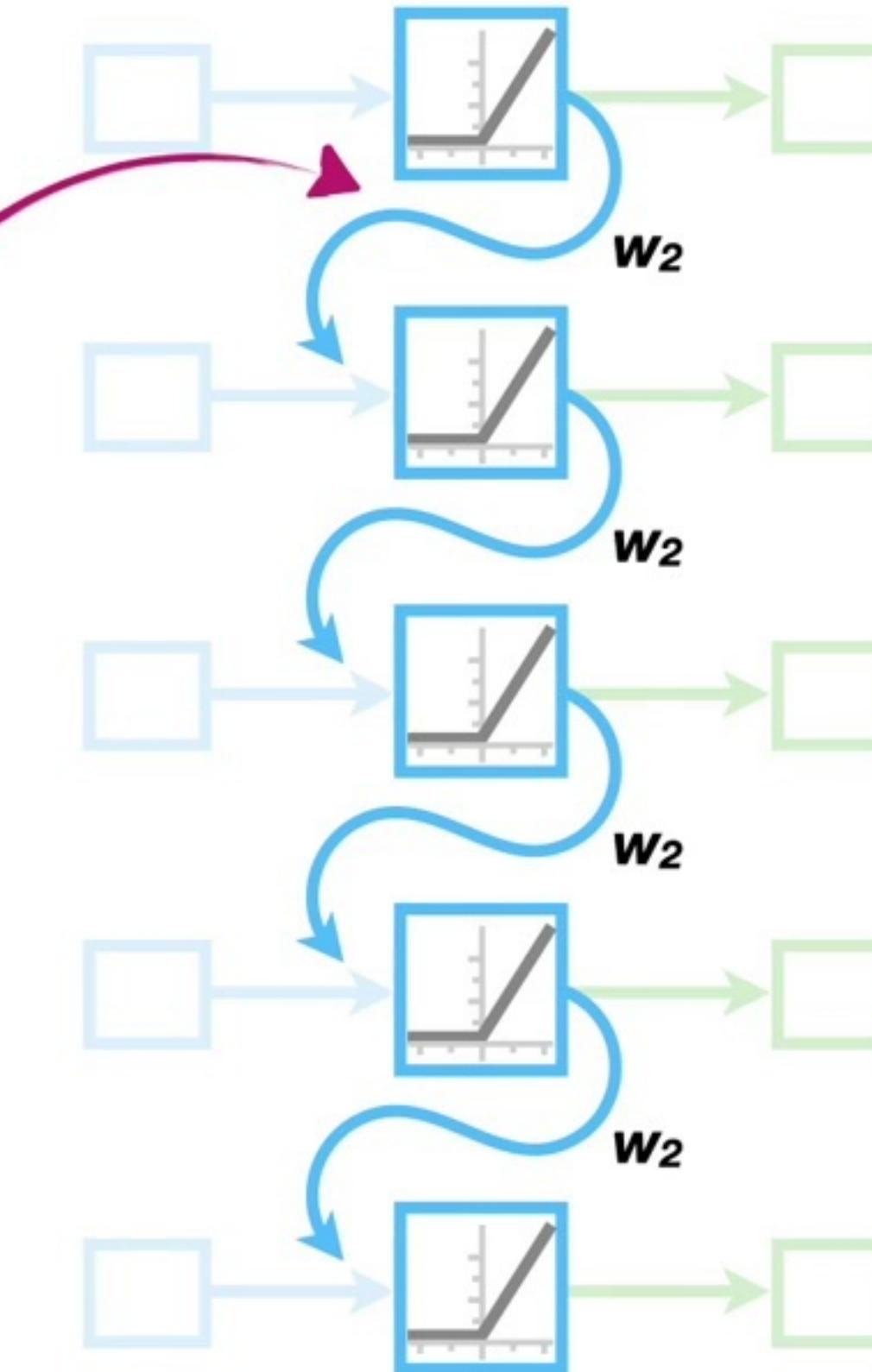


This problem is
called **The**
Vanishing/Exploding
Gradient Problem.





In our example, **The Vanishing/Exploding Gradient** problem has to do with the **weight** along the squiggle that we copy each time we **unroll** the network.



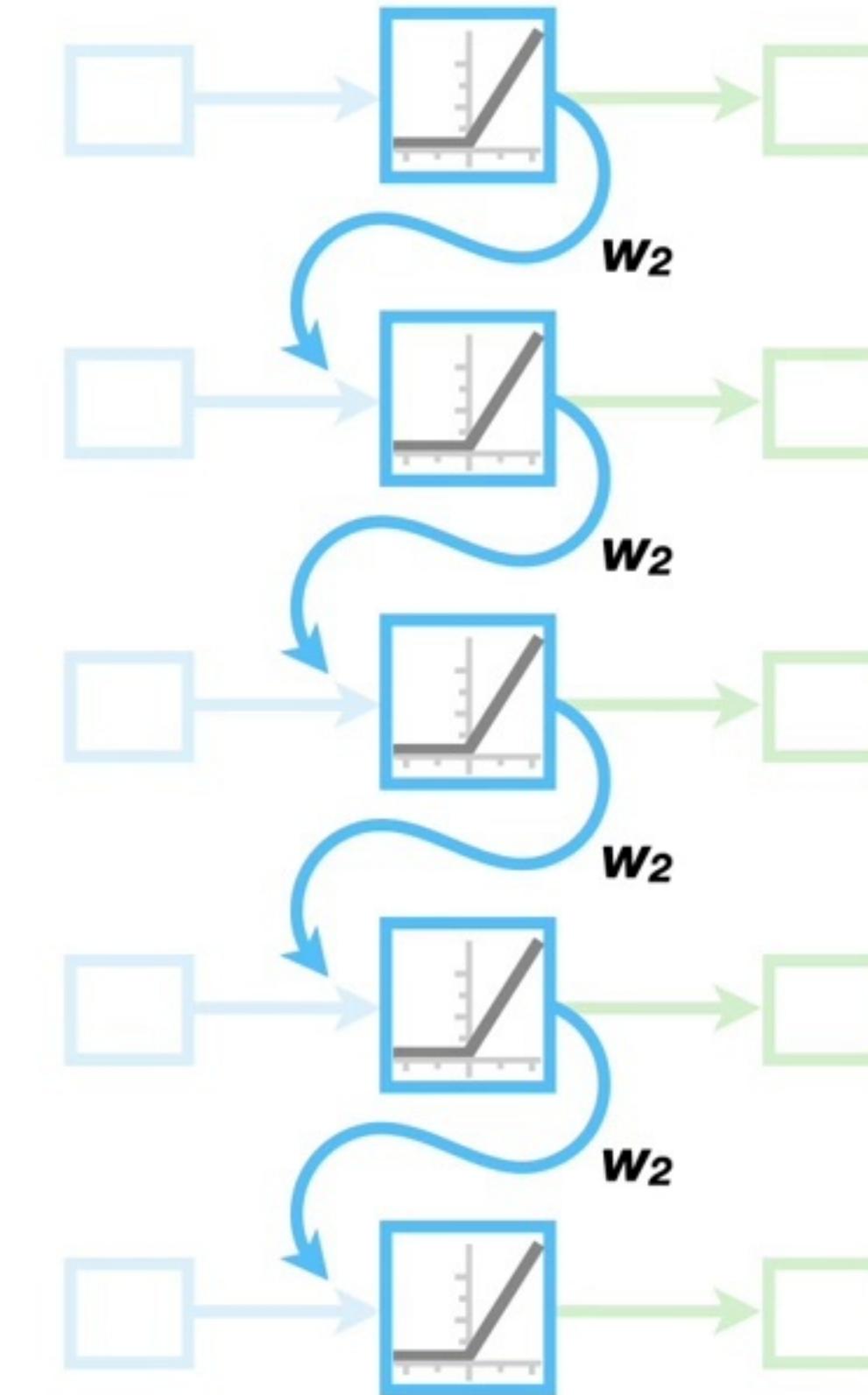


NOTE: To make it easier to understand **The Vanishing/Exploding Gradient Problem**, we're going to ignore the other **weights and biases** in this network and just focus on w_2 .





Also, just to remind you, when we optimize neural networks with **Backpropagation**, we first find the derivatives, or **Gradients**, for each parameter.





$$\frac{d \text{ } SSR}{d \text{ } w_1} = \frac{d \text{ } SSR}{d \text{ } \text{Predicted}} \times \frac{d \text{ } \text{Predicted}}{d \text{ } y_2} \times \frac{d \text{ } y_2}{d \text{ } x_2} \times \frac{d \text{ } x_2}{d \text{ } w_1} + \dots$$

Also, just to remind you, when we optimize neural networks with **Backpropagation**, we first find the derivatives, or **Gradients**, for each parameter.



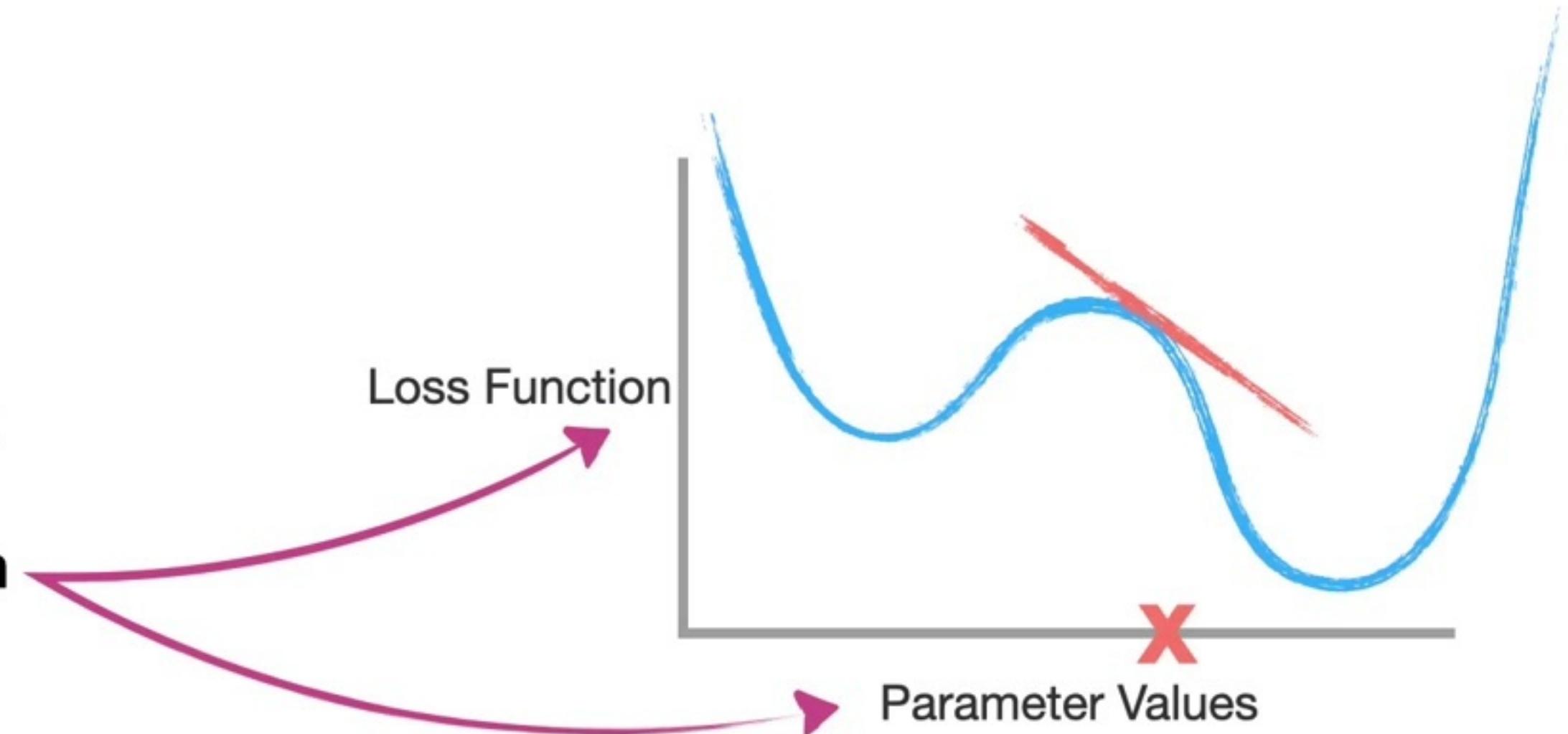
$$\frac{d \text{ } SSR}{d \text{ } w_1} = \frac{d \text{ } SSR}{d \text{ } \text{Predicted}} \times \frac{d \text{ } \text{Predicted}}{d \text{ } y_2} \times \frac{d \text{ } y_2}{d \text{ } x_2} \times \frac{d \text{ } x_2}{d \text{ } w_1} + \dots$$

We then plug those **Gradients** into the **Gradient Descent** algorithm to find the parameter values that minimize a **Loss Function**, like the **sum of the squared residuals**.



$$\frac{d \text{ SSR}}{d w_1} = \frac{d \text{ SSR}}{d \text{ Predicted}} \times \frac{d \text{ Predicted}}{d y_2} \times \frac{d y_2}{d x_2} \times \frac{d x_2}{d w_1} + \dots$$

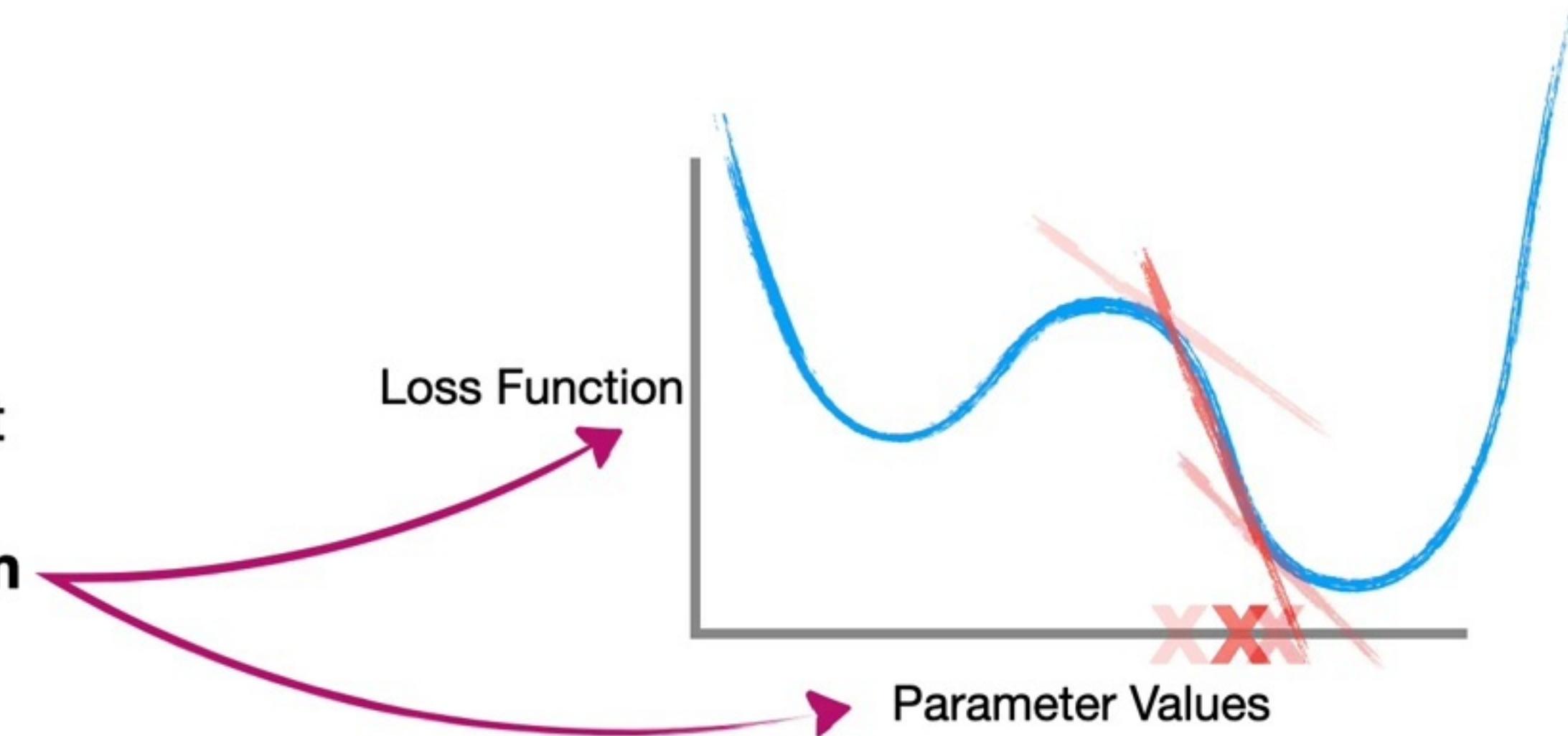
We then plug those **Gradients** into the **Gradient Descent** algorithm to find the parameter values that minimize a **Loss Function**, like the **sum of the squared residuals**.





$$\frac{d \text{ SSR}}{d w_1} = \frac{d \text{ SSR}}{d \text{ Predicted}} \times \frac{d \text{ Predicted}}{d y_2} \times \frac{d y_2}{d x_2} \times \frac{d x_2}{d w_1} + \dots$$

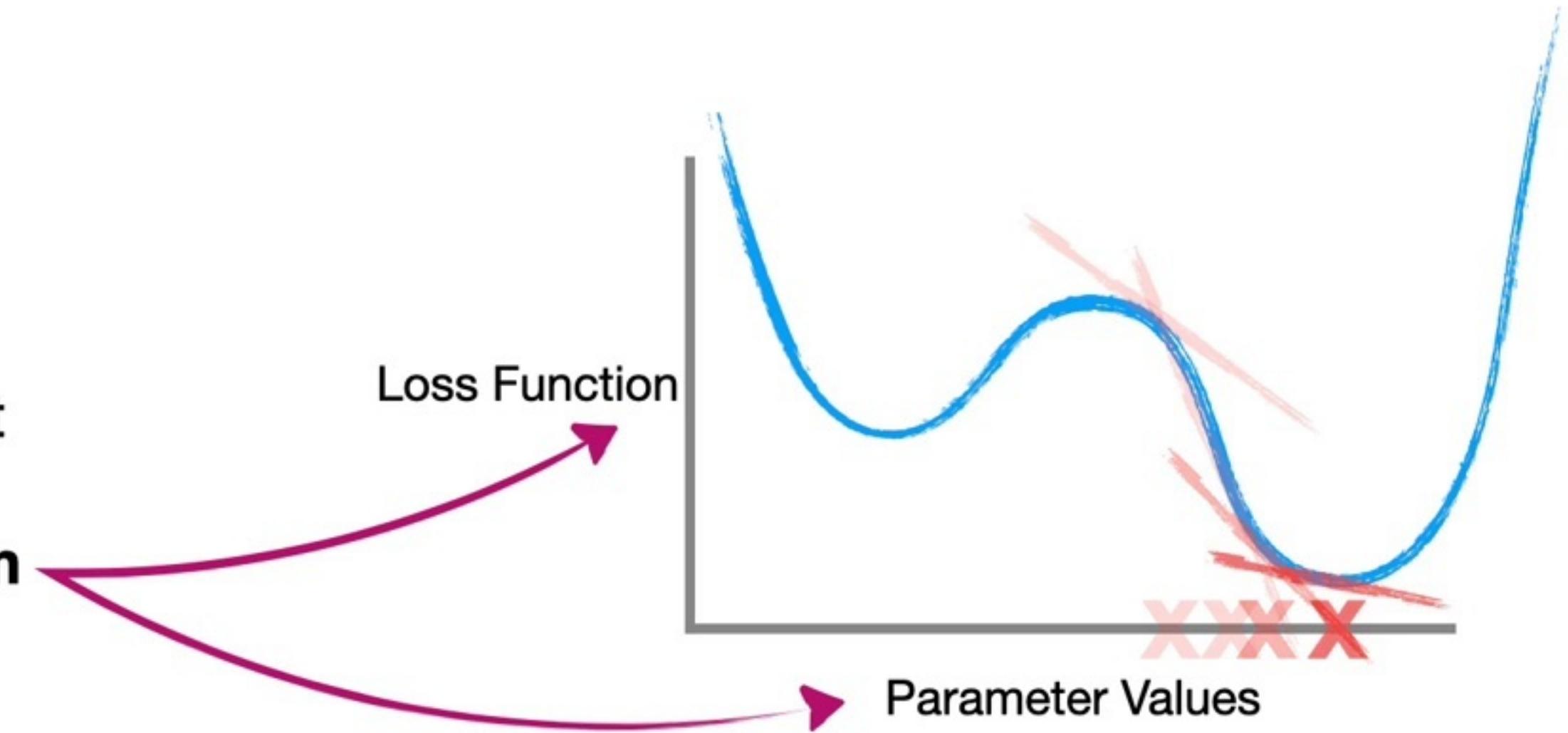
We then plug those **Gradients** into the **Gradient Descent** algorithm to find the parameter values that minimize a **Loss Function**, like the **sum of the squared residuals**.





$$\frac{d \text{ SSR}}{d w_1} = \frac{d \text{ SSR}}{d \text{ Predicted}} \times \frac{d \text{ Predicted}}{d y_2} \times \frac{d y_2}{d x_2} \times \frac{d x_2}{d w_1} + \dots$$

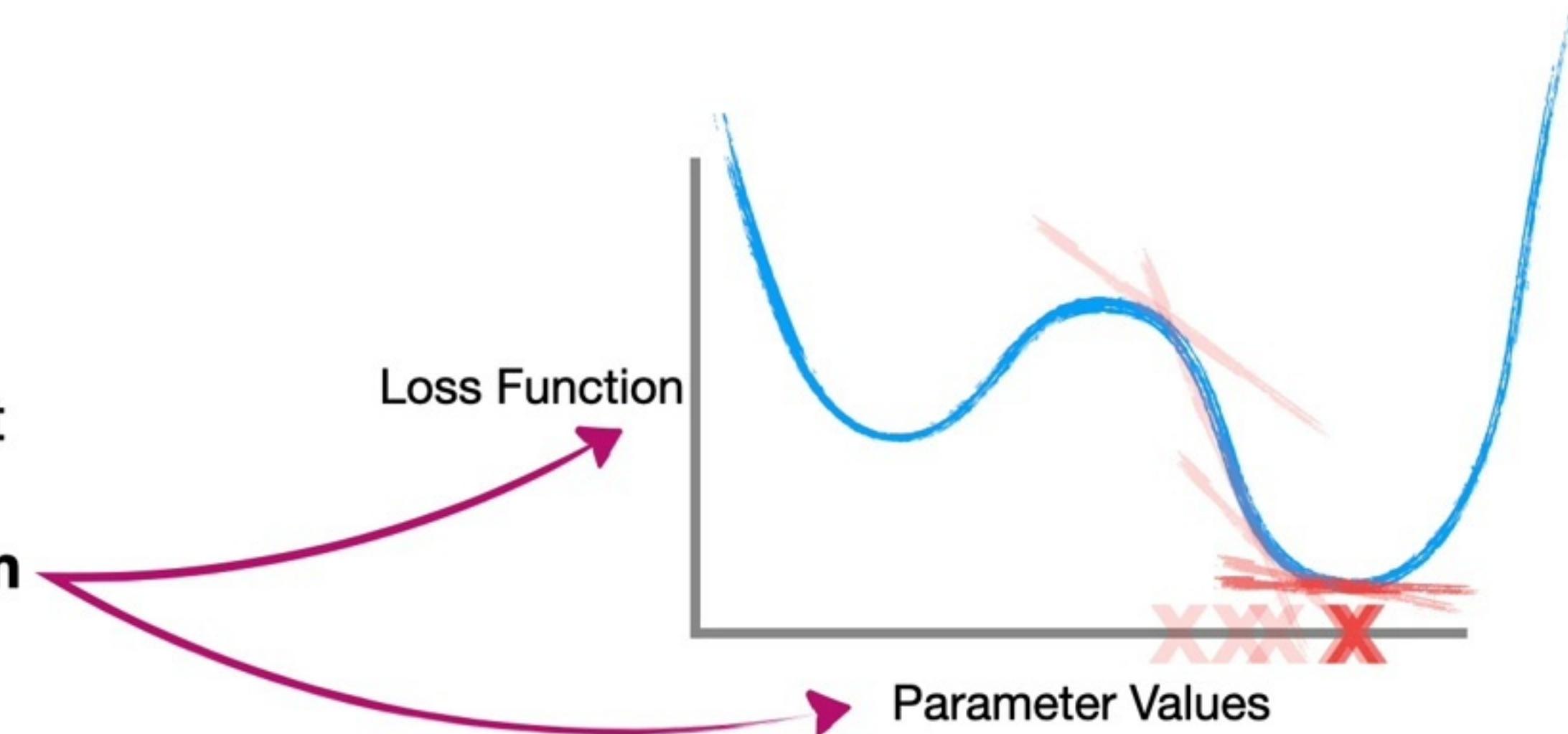
We then plug those **Gradients** into the **Gradient Descent** algorithm to find the parameter values that minimize a **Loss Function**, like the **sum of the squared residuals**.





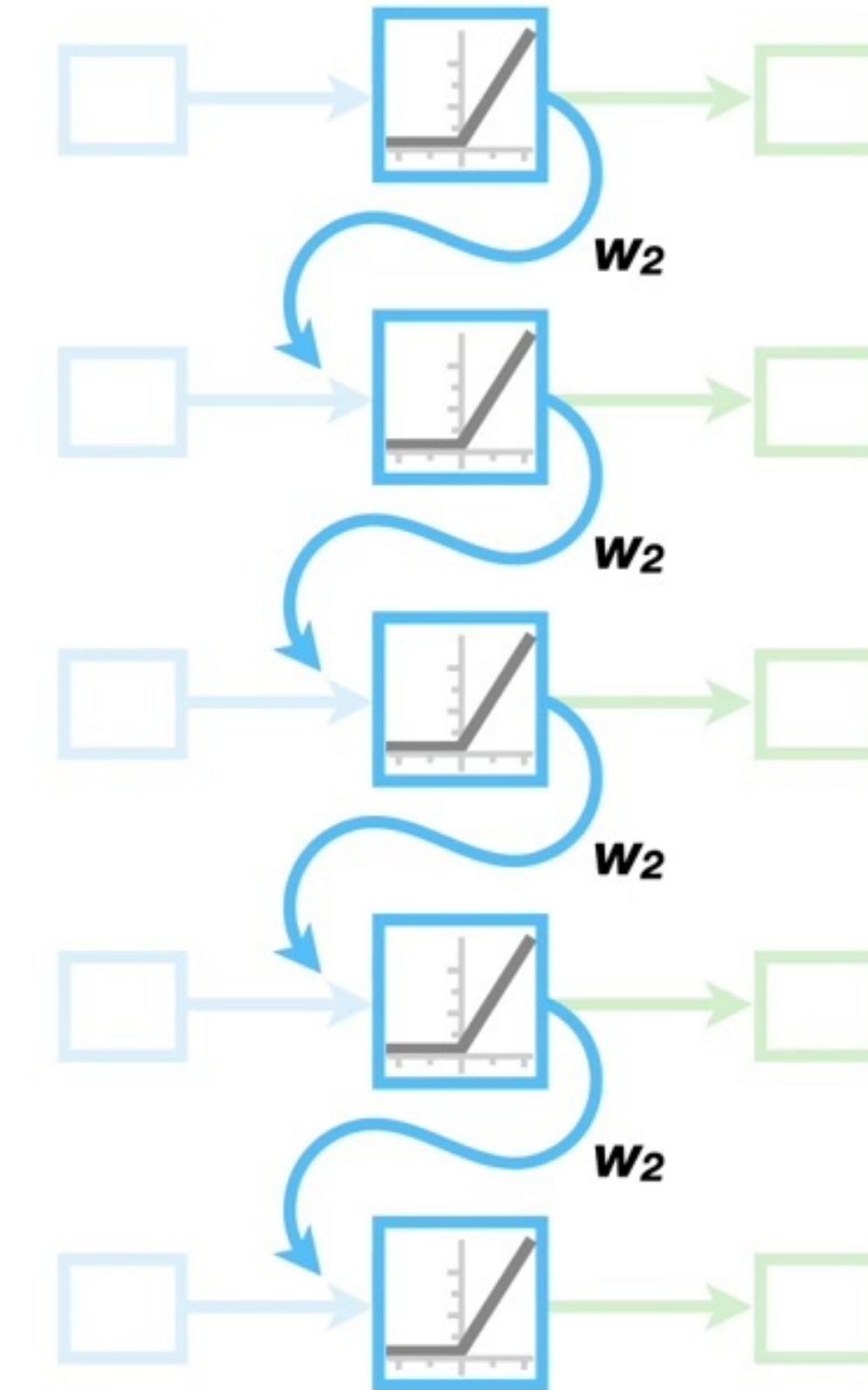
$$\frac{d \text{ SSR}}{d w_1} = \frac{d \text{ SSR}}{d \text{ Predicted}} \times \frac{d \text{ Predicted}}{d y_2} \times \frac{d y_2}{d x_2} \times \frac{d x_2}{d w_1} + \dots$$

We then plug those **Gradients** into the **Gradient Descent** algorithm to find the parameter values that minimize a **Loss Function**, like the **sum of the squared residuals**.





...we're going to start
by showing how a
gradient can
Explode.





In our example, the gradient will **explode**, when we set w_2 to any value larger than 1.





So let's set
 $w_2 = 2$.





Now, the first input value, **Input₁**...





...will be multiplied
by **2** on the first
squiggle...

$\text{Input}_1 \times 2$





...and then
multiplied by **2** by
the next squiggle...

$\text{Input}_1 \times 2 \times 2$





...and again by the
next squiggle...

Input₁ $\times 2 \times 2 \times 2$





...and again by the
last squiggle.

Input₁ $\times 2 \times 2 \times 2 \times 2$





In other words, since
we **unrolled** the
recurrent neural
network **4 times**...

Input₁ $\times 2 \times 2 \times 2 \times 2$





...we multiply the input value by w_2 , which is 2, raised to the number of times we **unrolled**, which is 4.

$$\text{Input}_1 \times 2 \times 2 \times 2 \times 2$$

$$= \text{Input}_1 \times 2^4$$

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





And that means the first input value is amplified **16** times before it gets to the final copy of the network.

$$\mathbf{Input}_1 \times 2 \times 2 \times 2 \times 2$$

$$= \mathbf{Input}_1 \times 2^4$$

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





Now, if we had **50** sequential days of stock market data, which, to be honest, really isn't that much data...

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

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





...then we would **unroll**
the network **50** times...

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

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





...and 2^{50} is
A HUGE NUMBER.

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

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





And this **HUGE NUMBER**
is why they call this an
Exploding Gradient
Problem.

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

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





$$\frac{d \text{SSR}}{d w_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d w_1} + \dots$$

If we tried to train this recurrent neural network with backpropagation, this **HUGE NUMBER** would find its way into some of the gradients...

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

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





$$\frac{d \text{SSR}}{d w_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d w_1} + \dots$$

If we tried to train this recurrent neural network with backpropagation, this **HUGE NUMBER** would find its way into some of the gradients...

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

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

$$\frac{d x_1}{d w_1} = \frac{d}{d w_1} \text{Input}_1 \times w_1 \times w_2^{50} \dots$$



$$\frac{d \text{SSR}}{d w_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d w_1} + \dots$$

If we tried to train this recurrent neural network with backpropagation, this **HUGE NUMBER** would find its way into some of the gradients...

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

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

$$\begin{aligned}\frac{d x_1}{d w_1} &= \frac{d}{d w_1} \text{Input}_1 \times w_1 \times w_2^{50} \dots \\ &= \frac{d}{d w_1} \text{Input}_1 \times w_1 \times \text{A HUGE NUMBER} \dots\end{aligned}$$

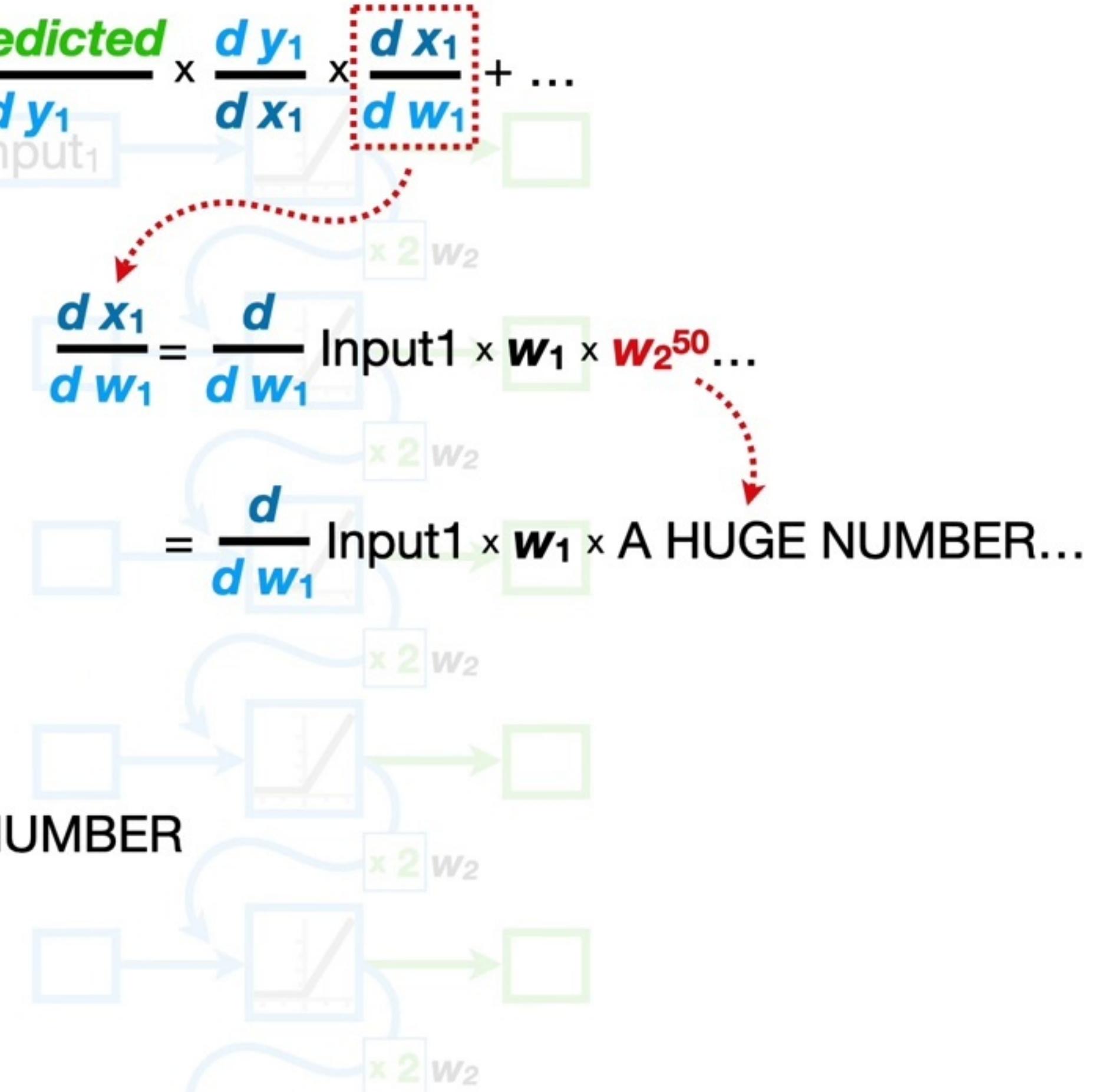


$$\frac{d \text{SSR}}{d w_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d w_1} + \dots$$

...and that would make it hard to take small steps to find the optimal **weights and biases**.

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

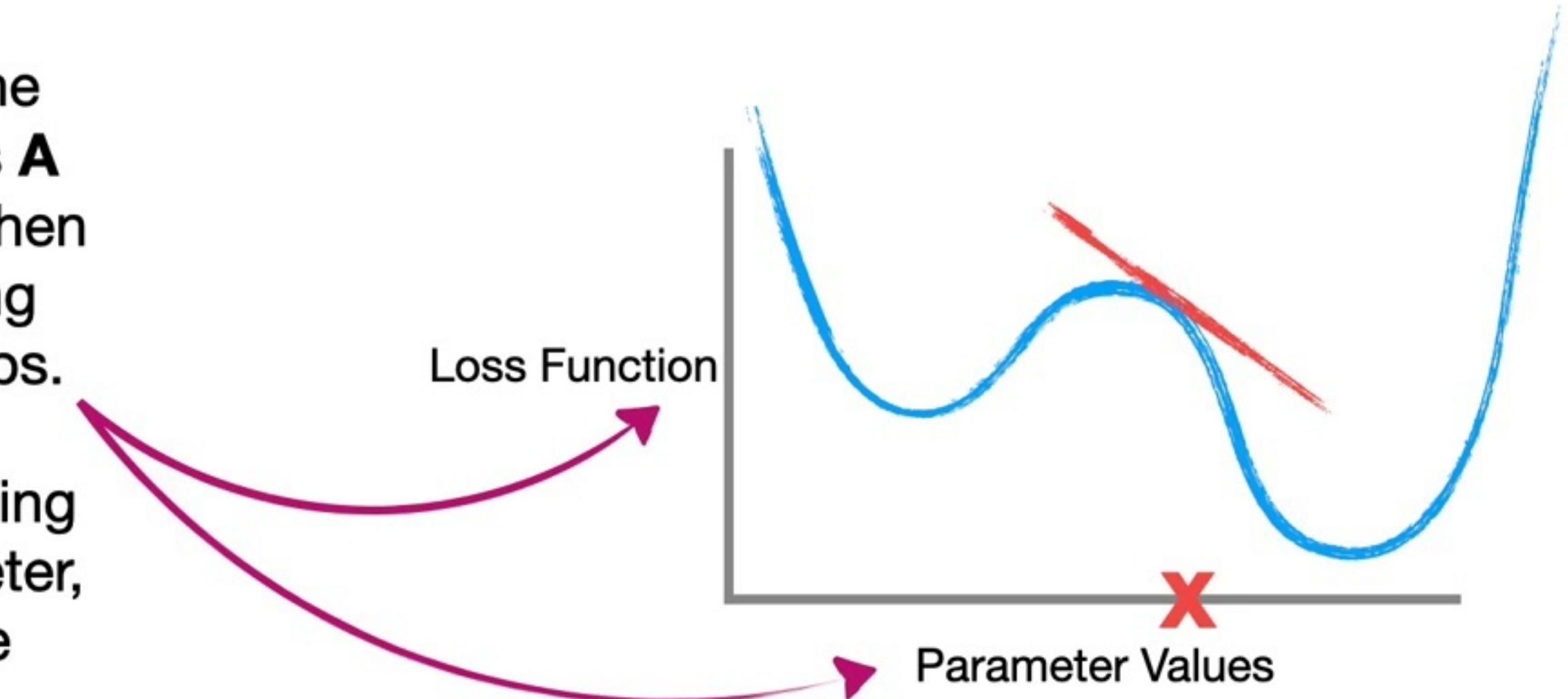
$$= \text{Input}_1 \times w_2^{\text{Num. Unroll}}$$





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

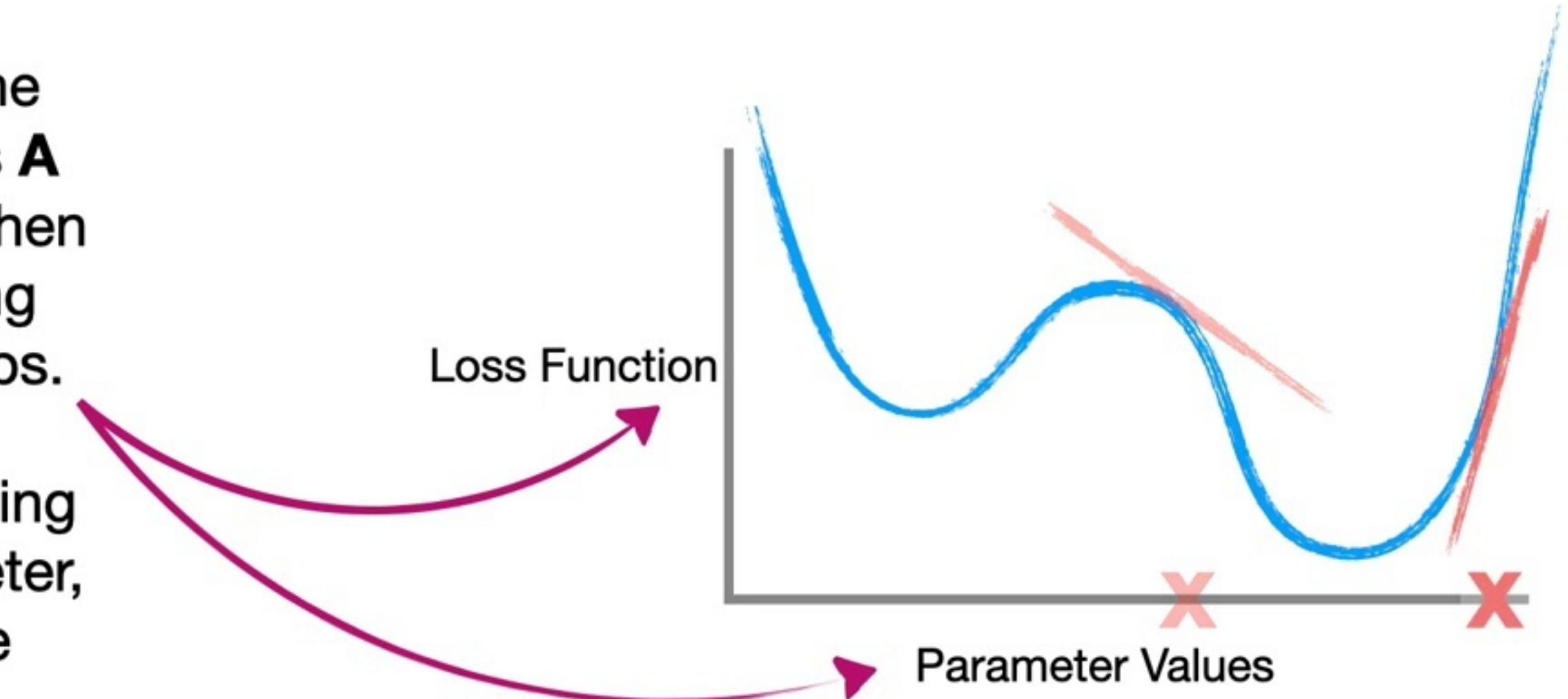
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

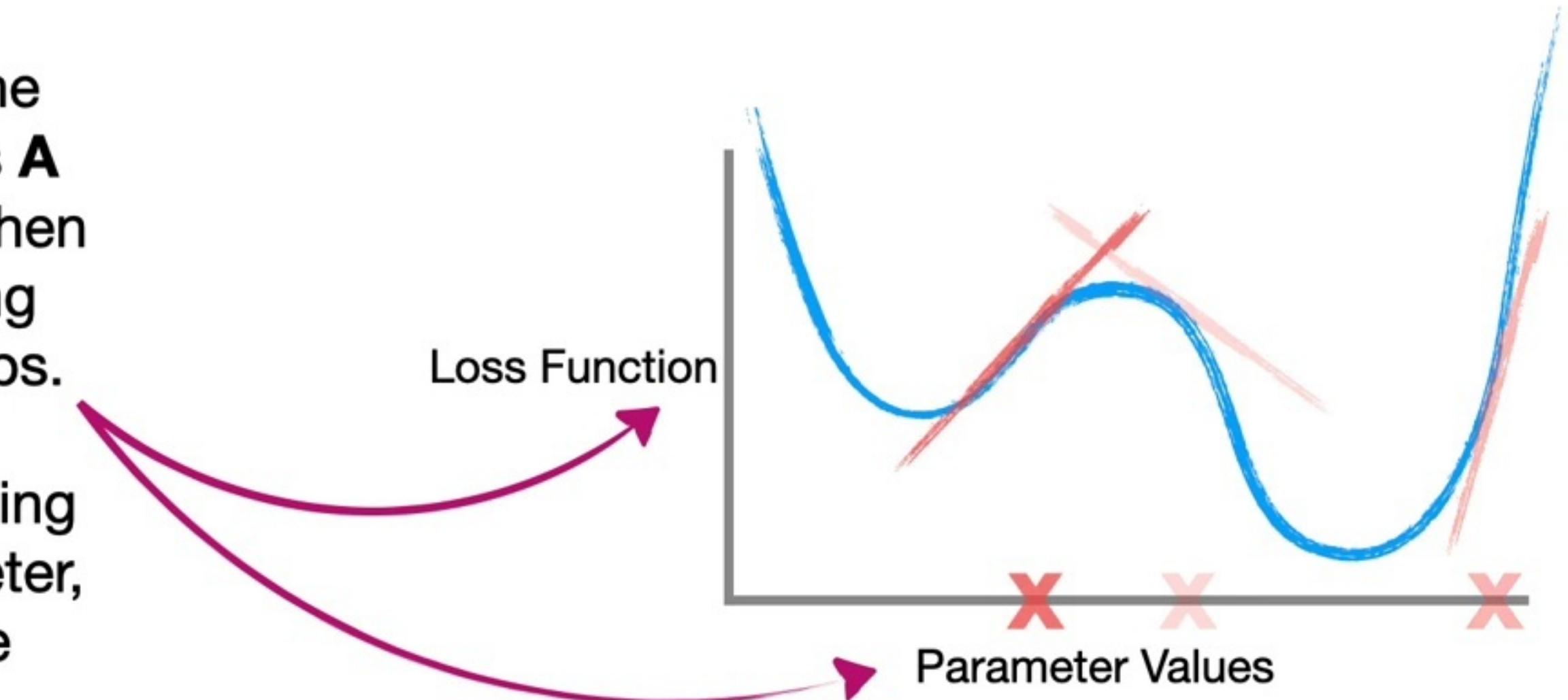
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

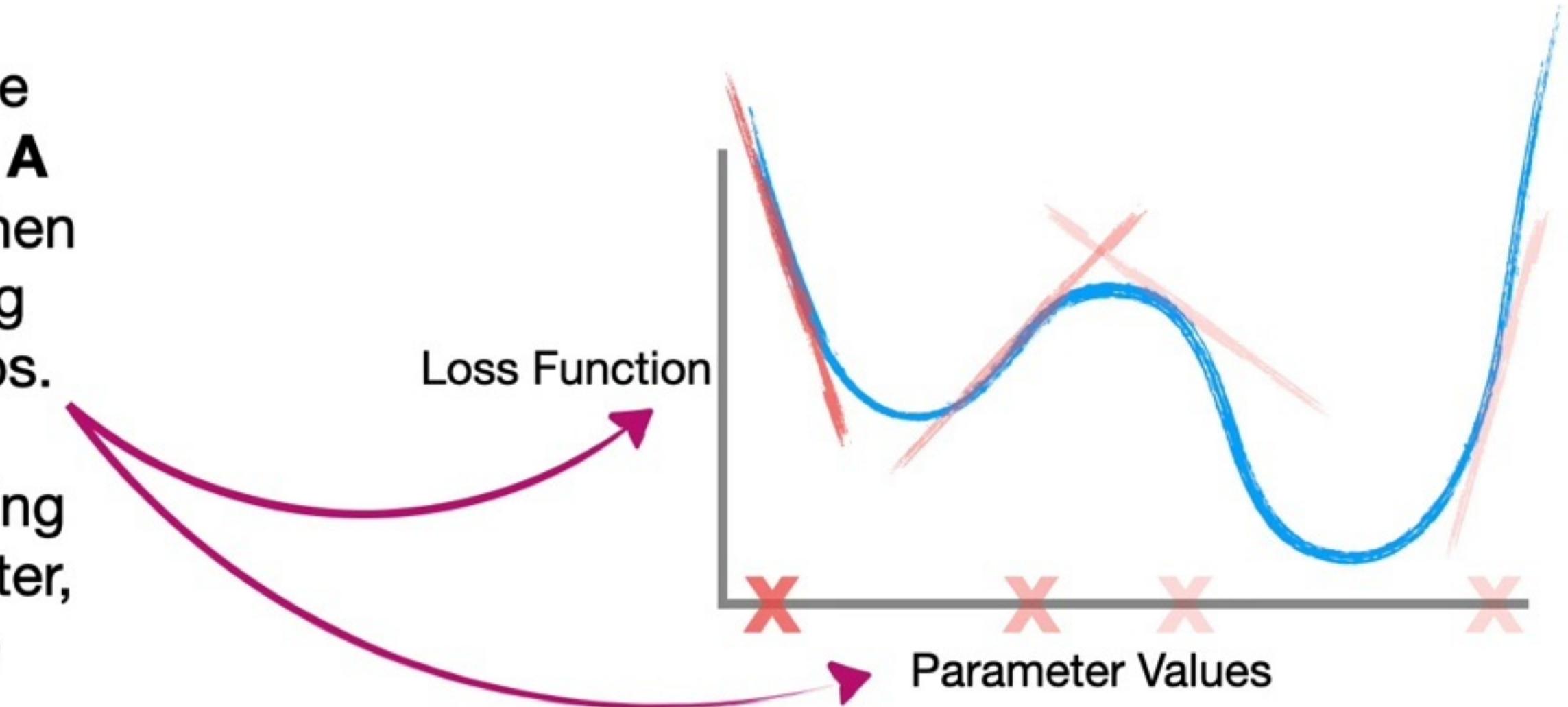
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

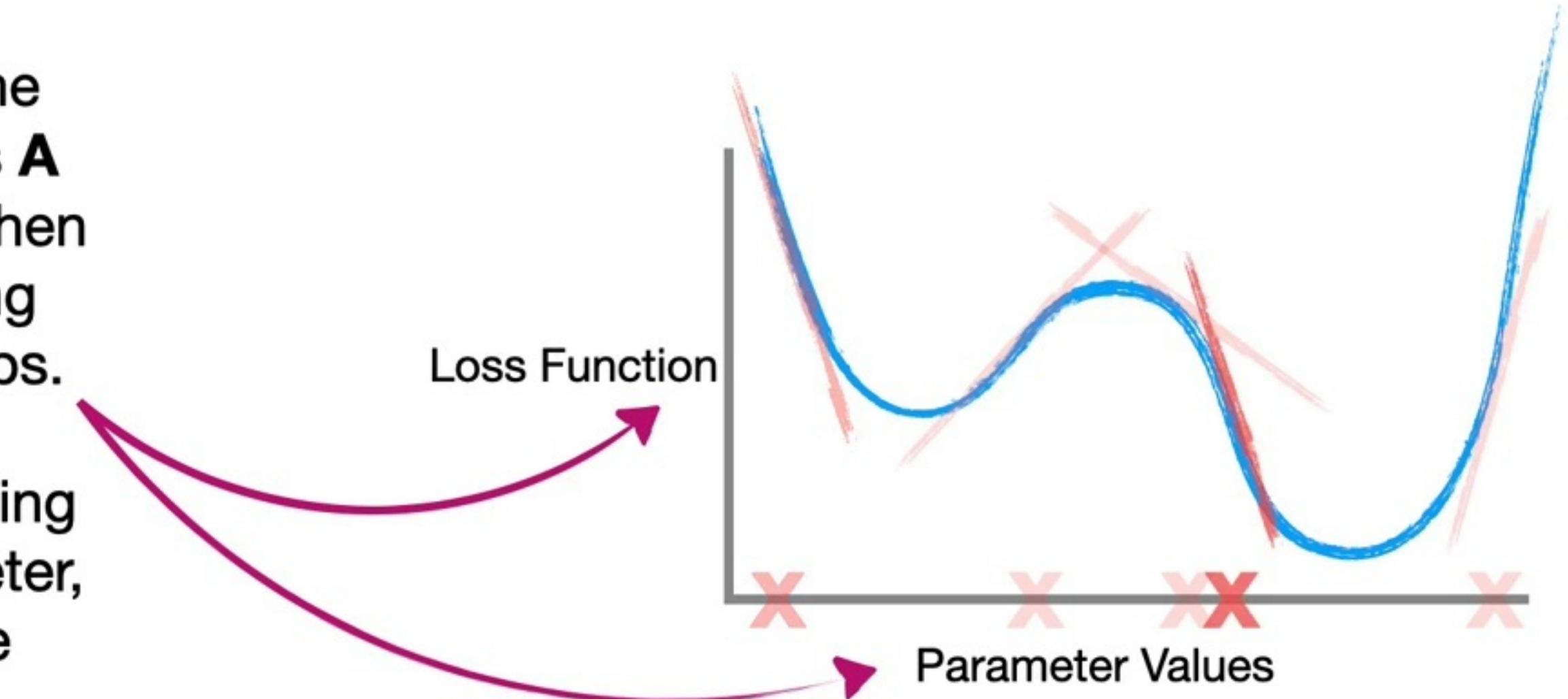
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

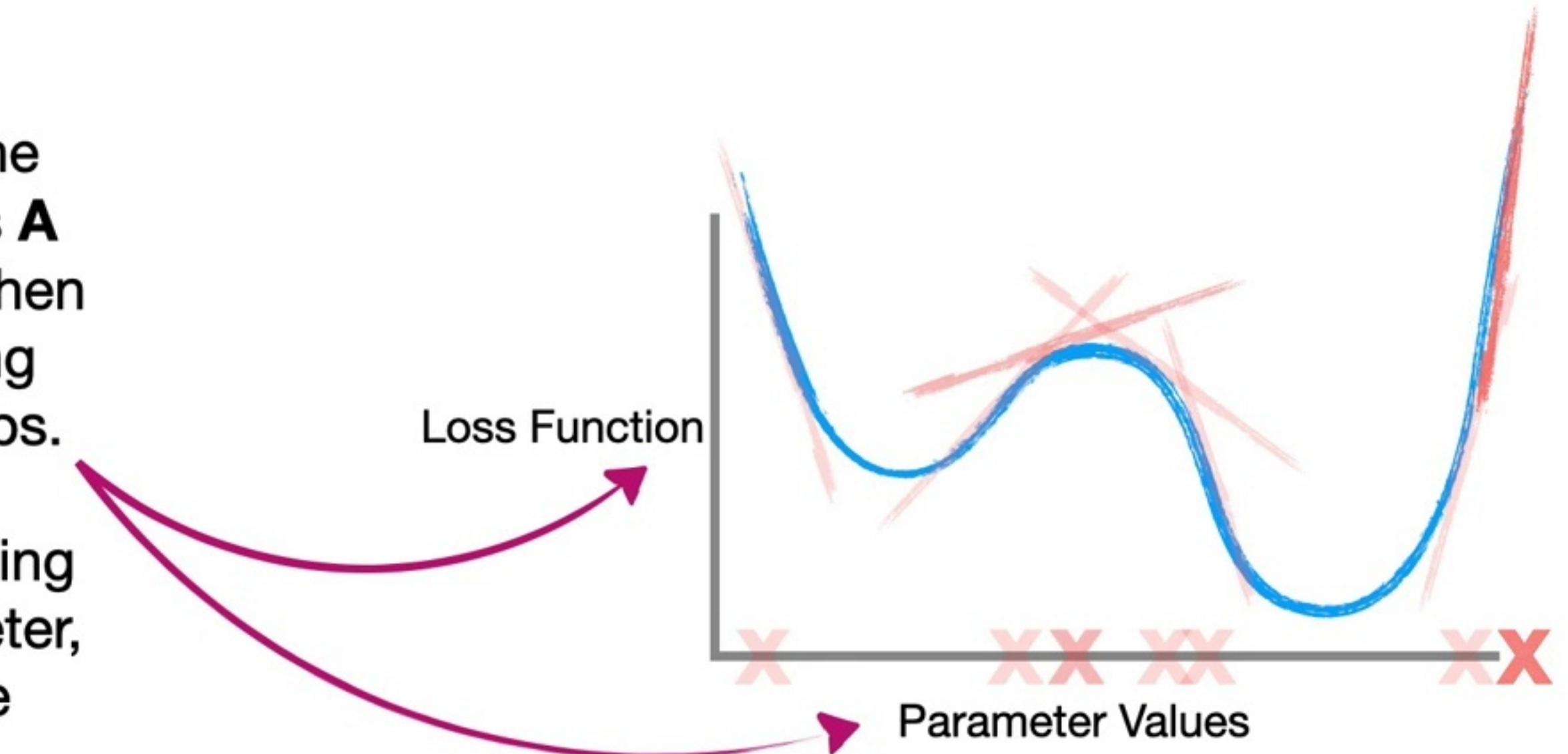
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

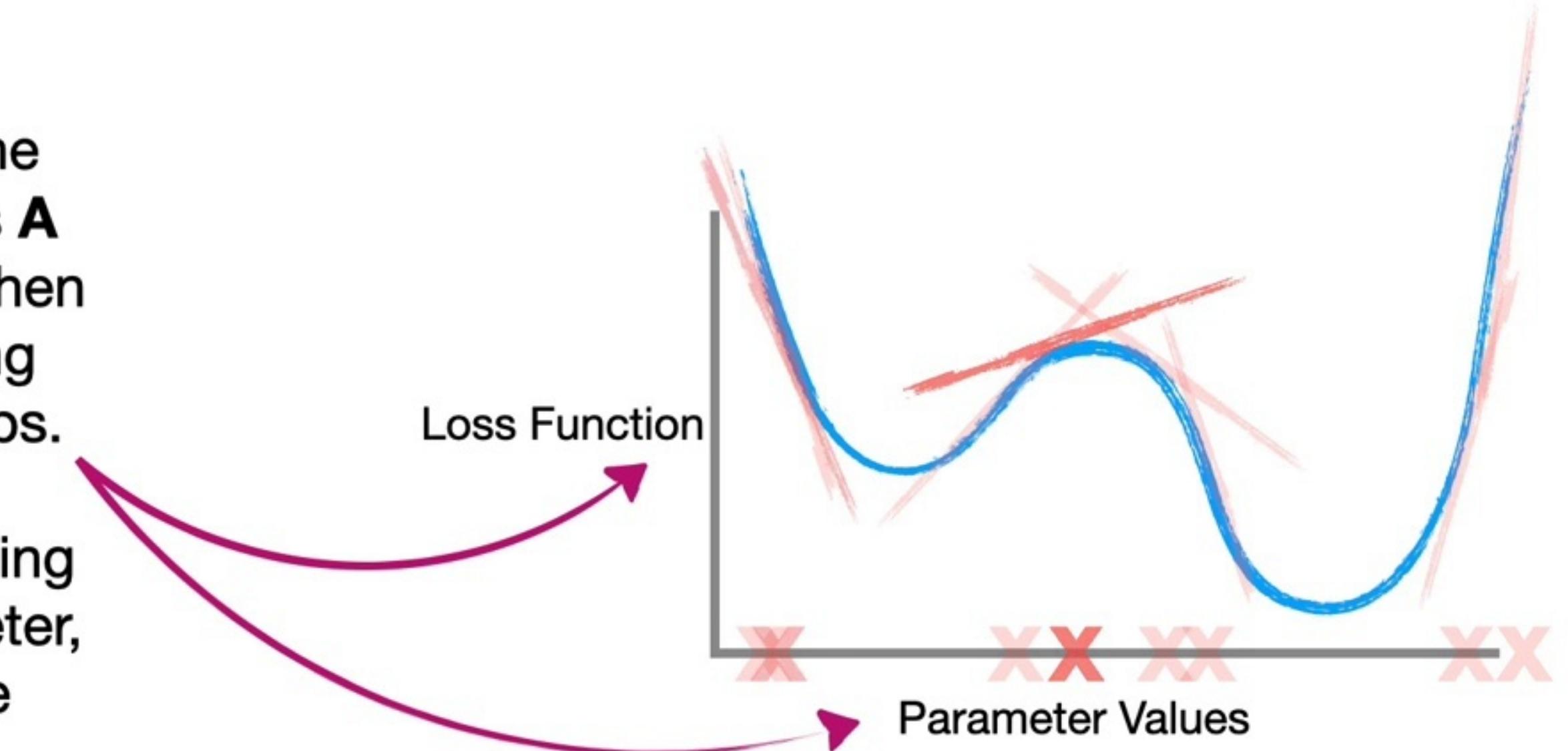
And instead of finding the optimal parameter, we'll just bounce around a lot.





However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

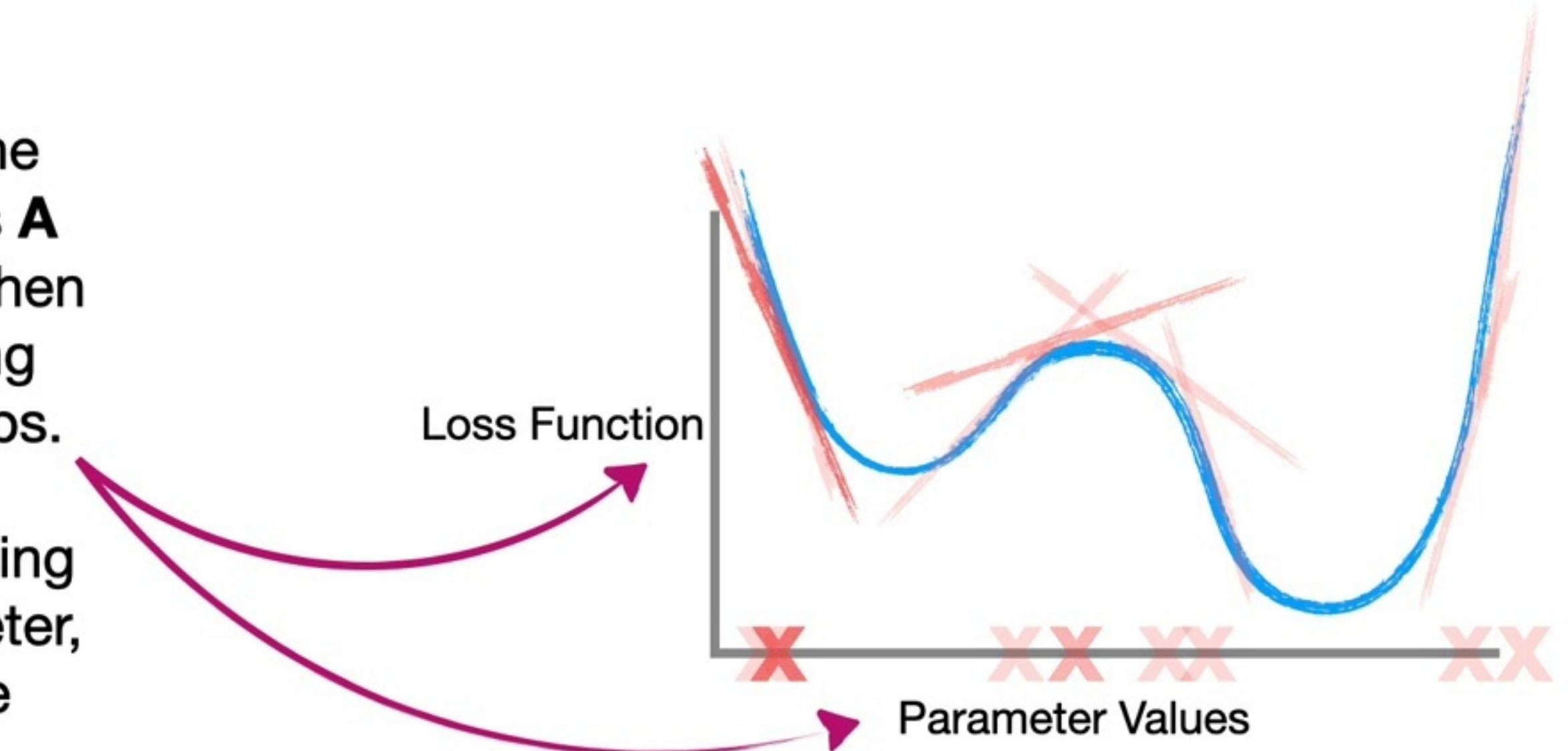
And instead of finding the optimal parameter, we'll just bounce around a lot.





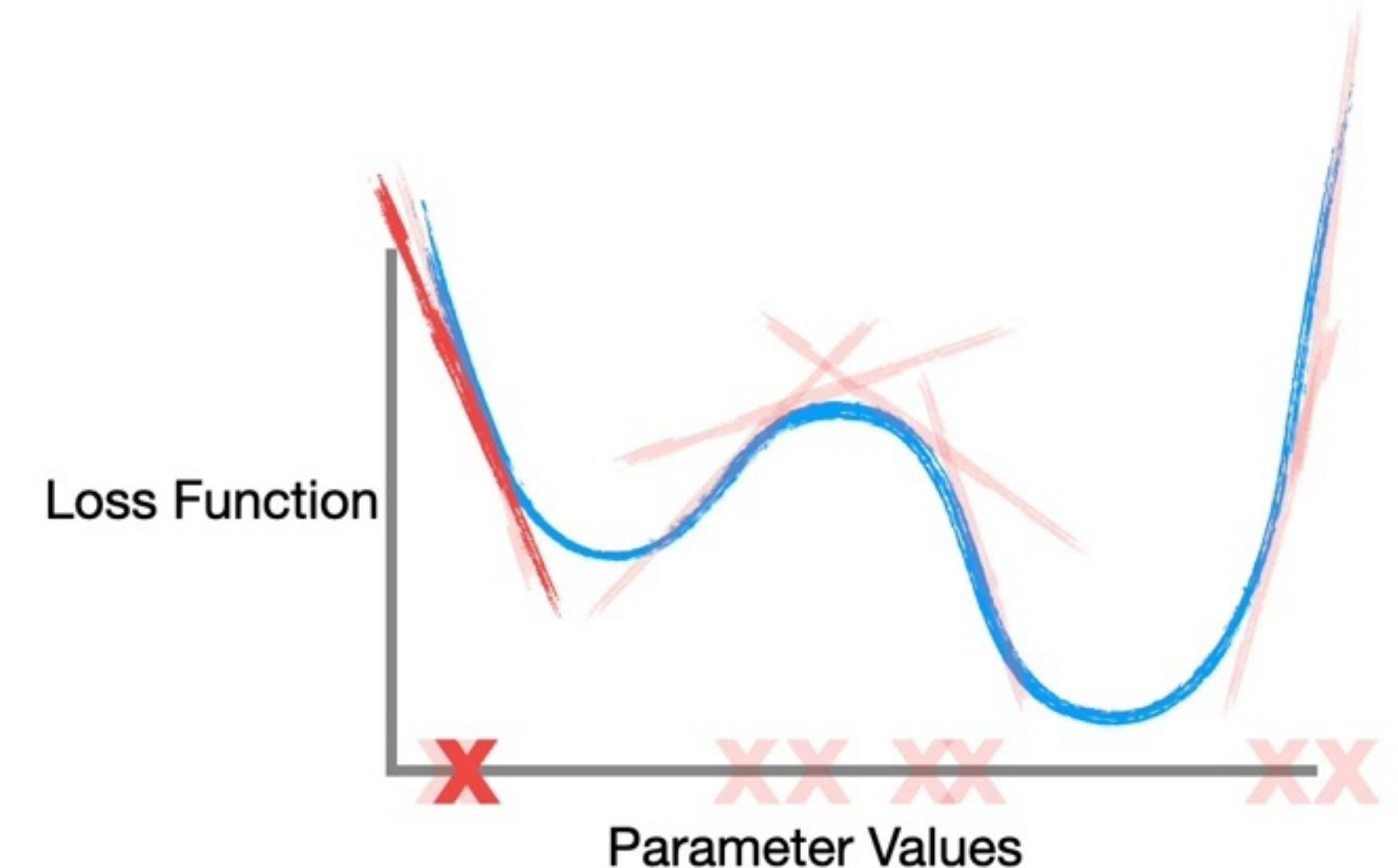
However, when the **Gradient** contains **A HUGE NUMBER**, then we'll end up taking relatively large steps.

And instead of finding the optimal parameter, we'll just bounce around a lot.





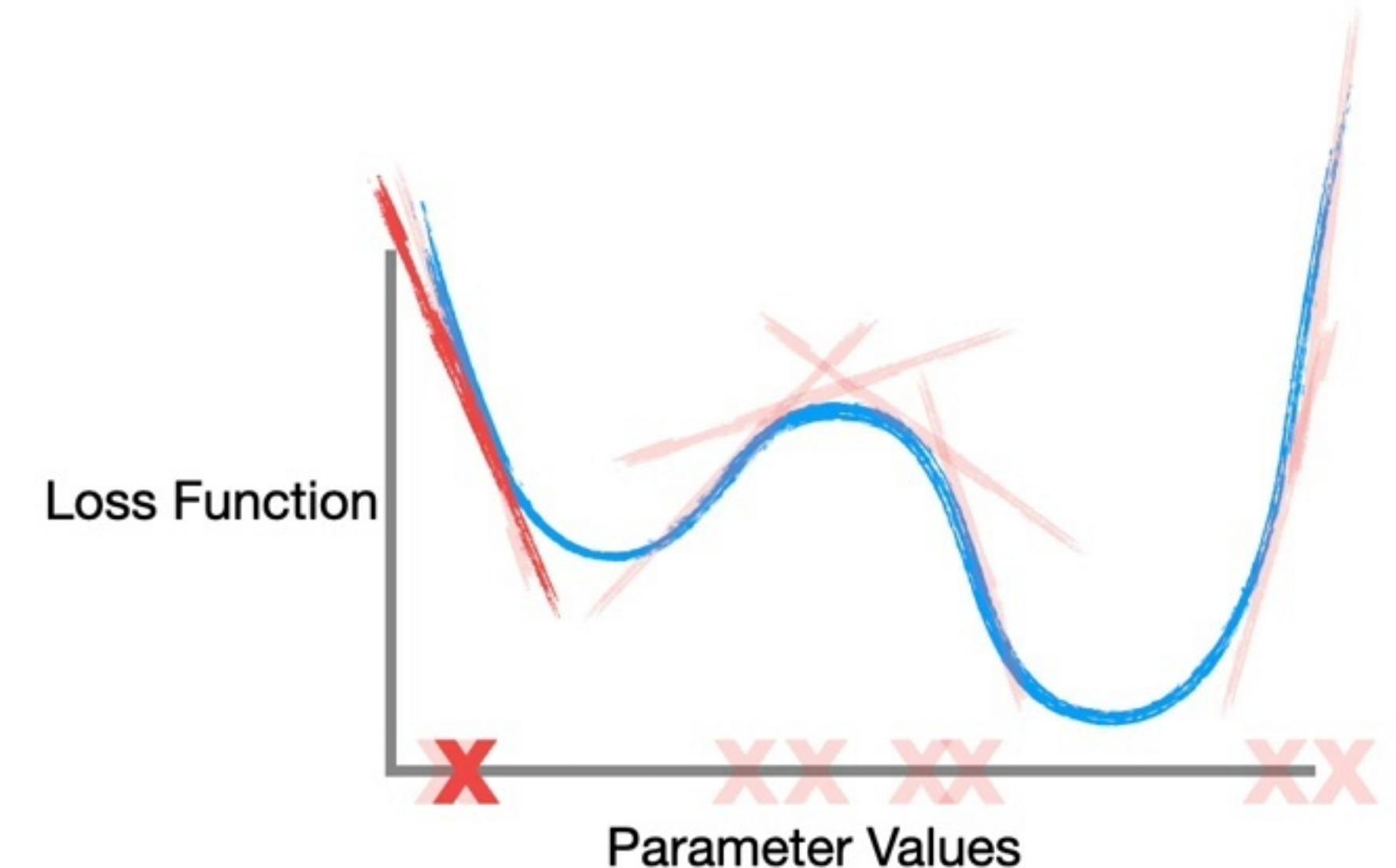
One way to prevent
**The Exploding
Gradient Problem**
would be to limit w_2
to values < 1 .





One way to prevent
**The Exploding
Gradient Problem**
would be to limit w_2
to values < 1 .

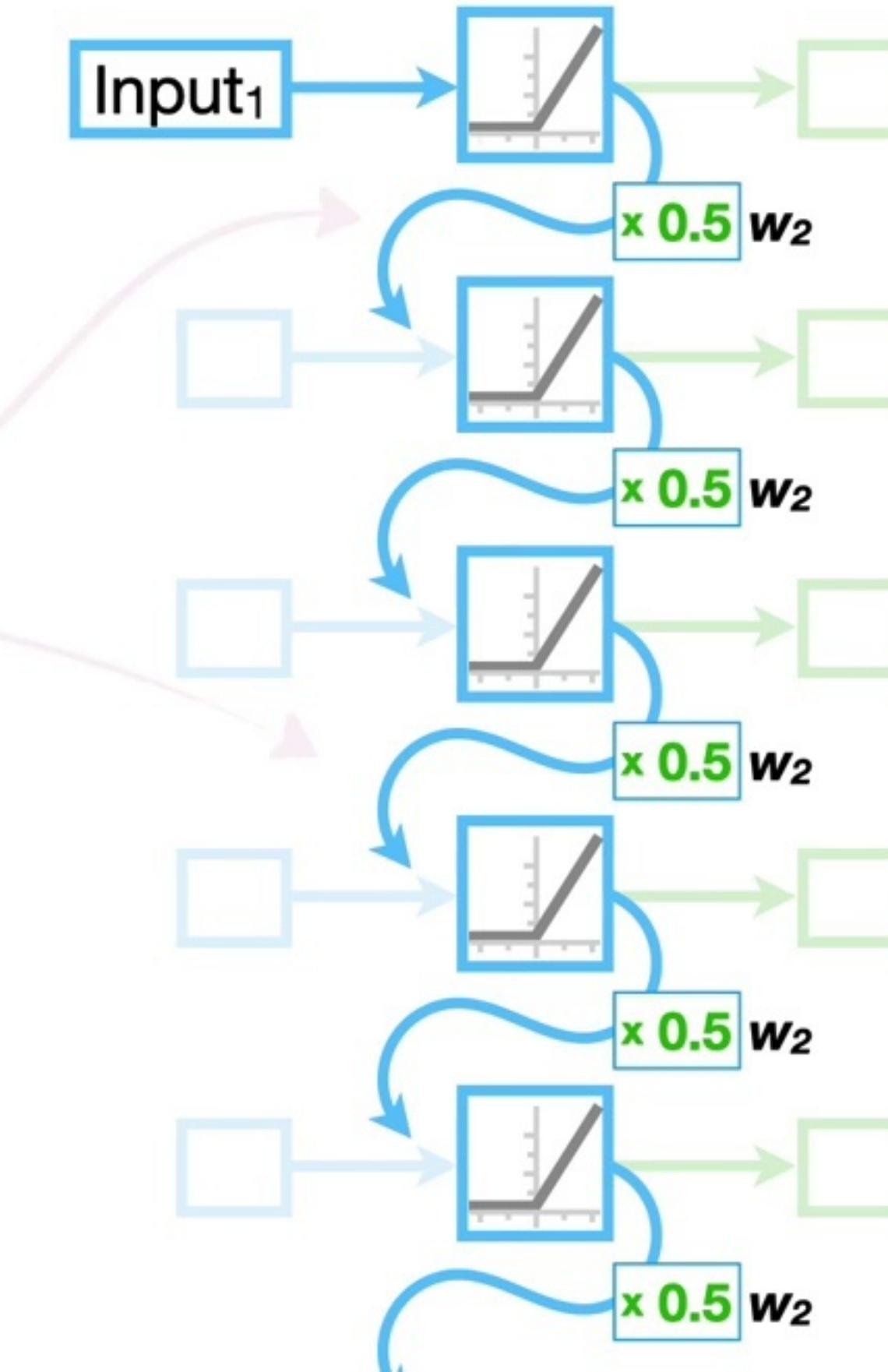
However, this results
in **The Vanishing
Gradient Problem**.





Now, just like before, we multiply the first input by w_2 raised to the number of times we **unroll** the network.

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





So if we have **50** sequential input values, that means multiplying **Input₁** by **0.5⁵⁰**...

Input₁ × 0.5⁵⁰

Input₁ × **w₂** Num. Unroll





**...and 0.5^{50} is a number
super close to 0.**

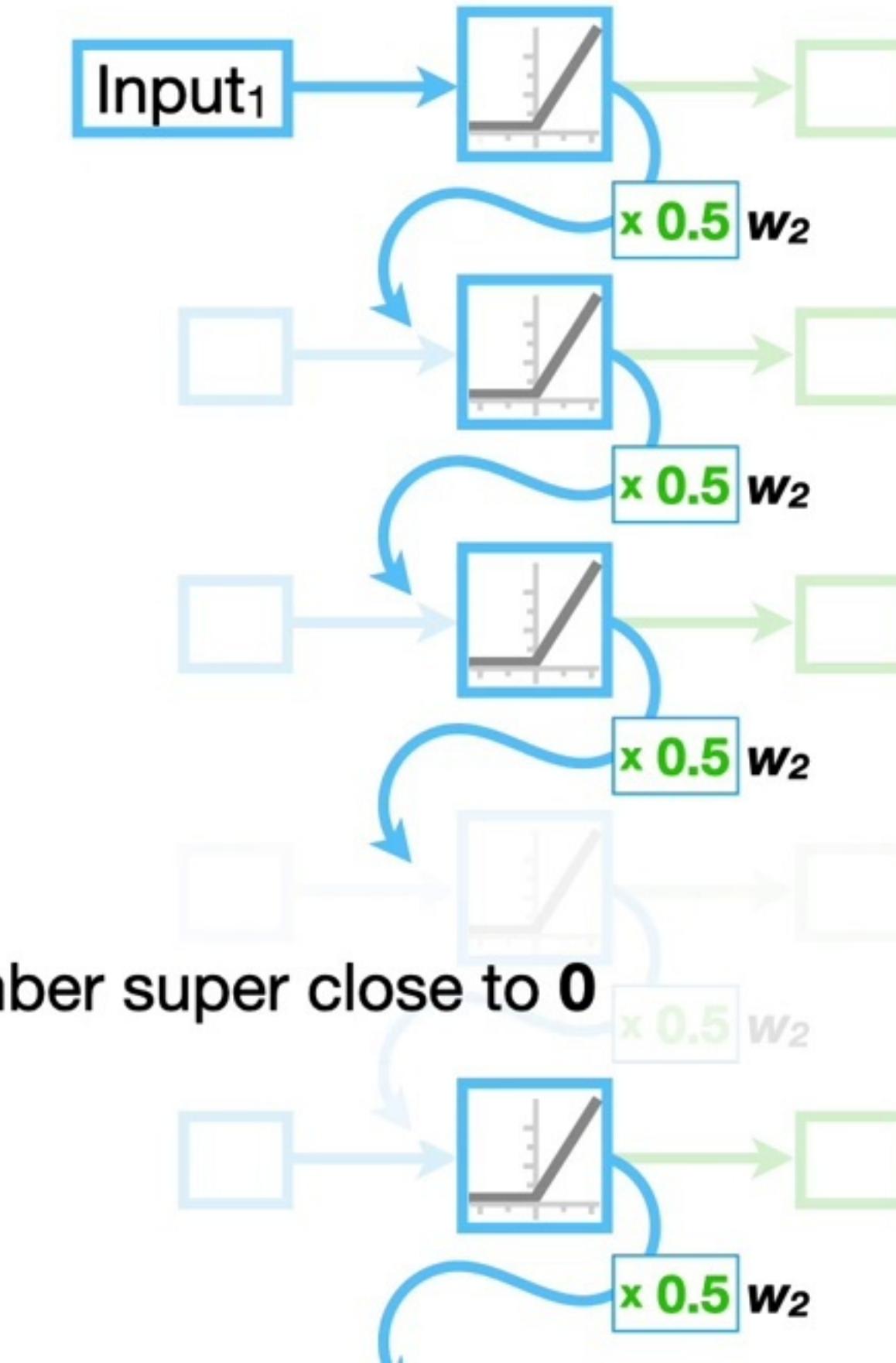
$$\mathbf{Input}_1 \times 0.5^{50} = \mathbf{Input}_1 \times \text{a number super close to 0}$$

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





Because this number is
super close to 0, this is
called the **Vanishing
Gradient Problem**.

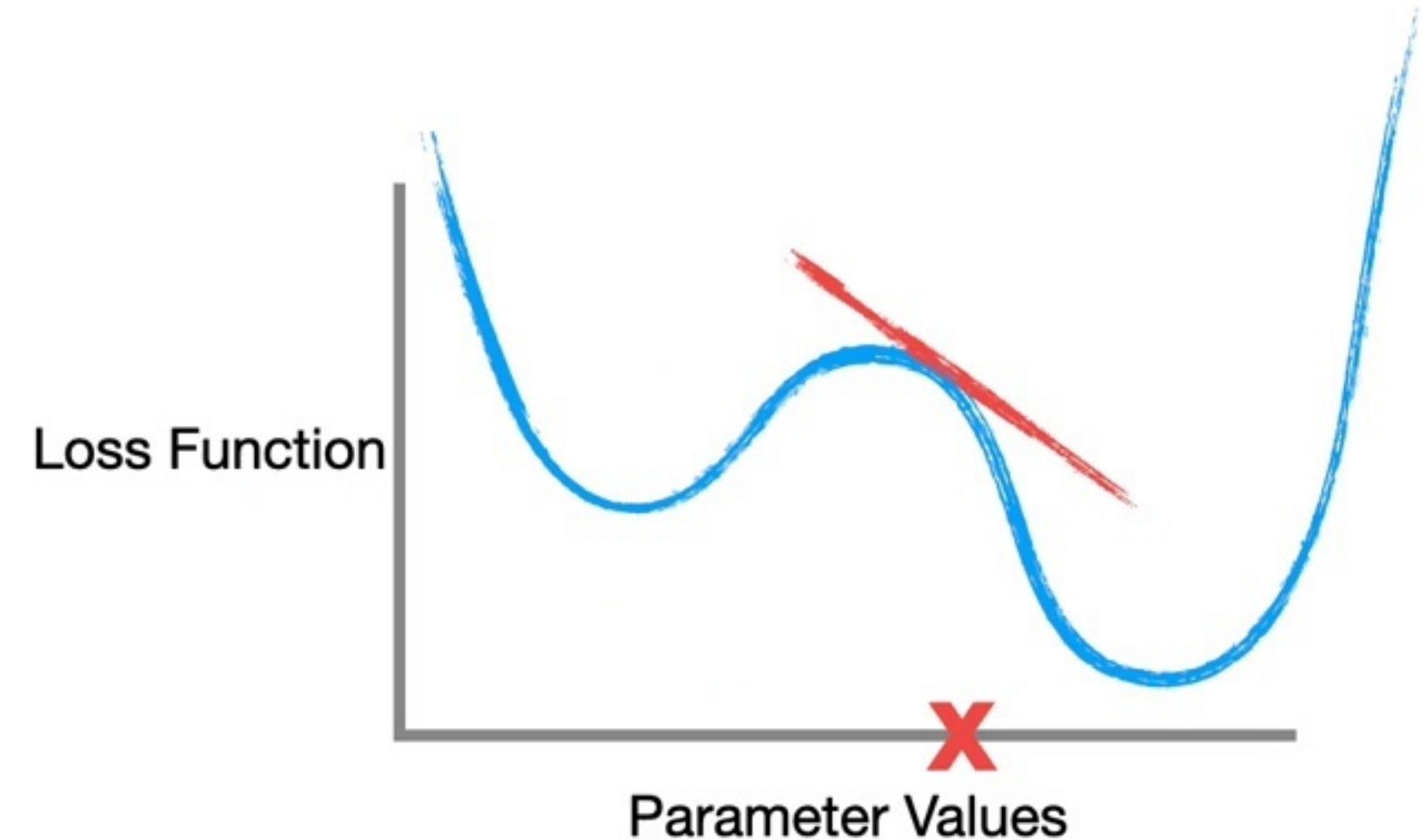


$$\text{Input}_1 \times 0.5^{50} = \text{Input}_1 \times \text{a number super close to } 0$$

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

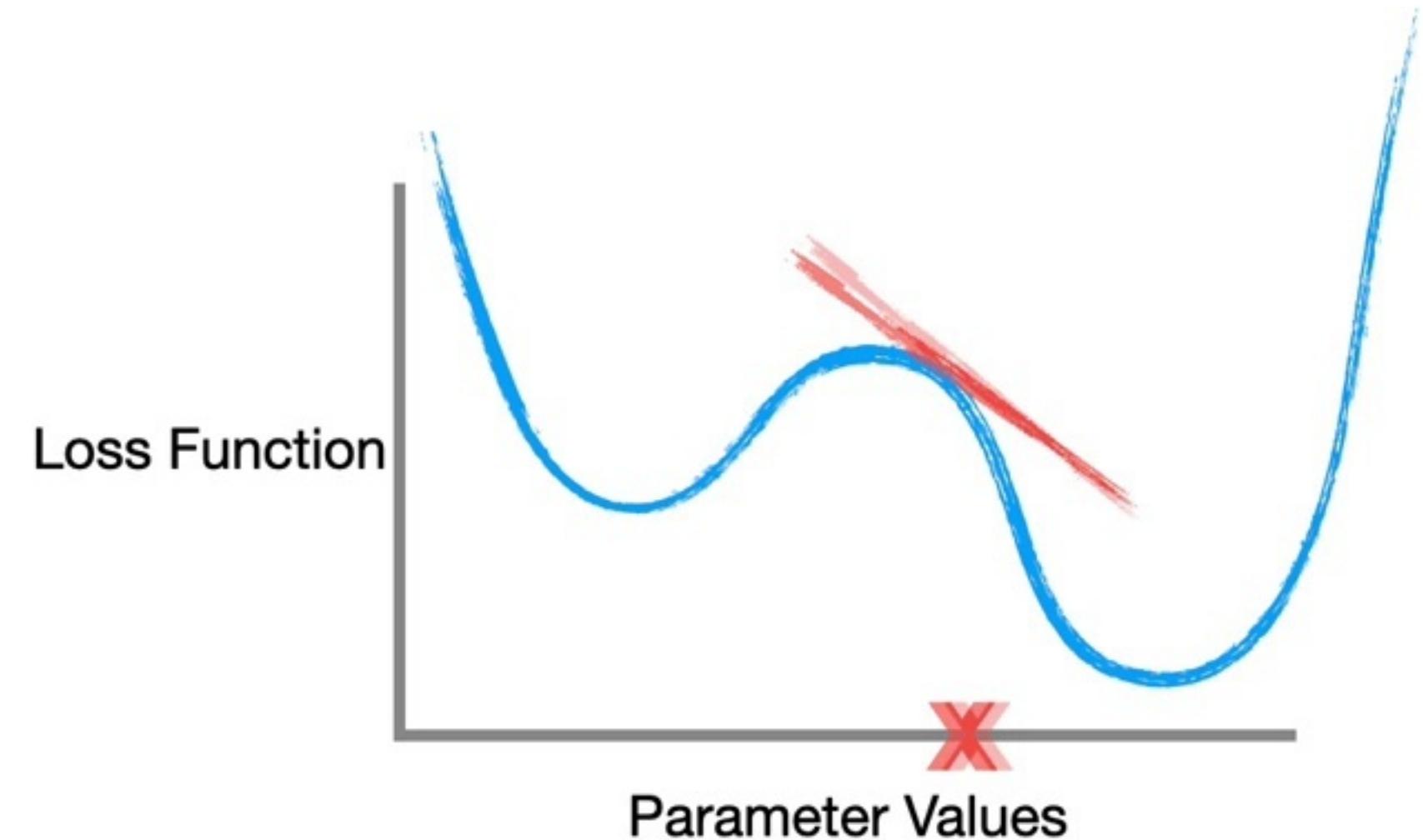


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



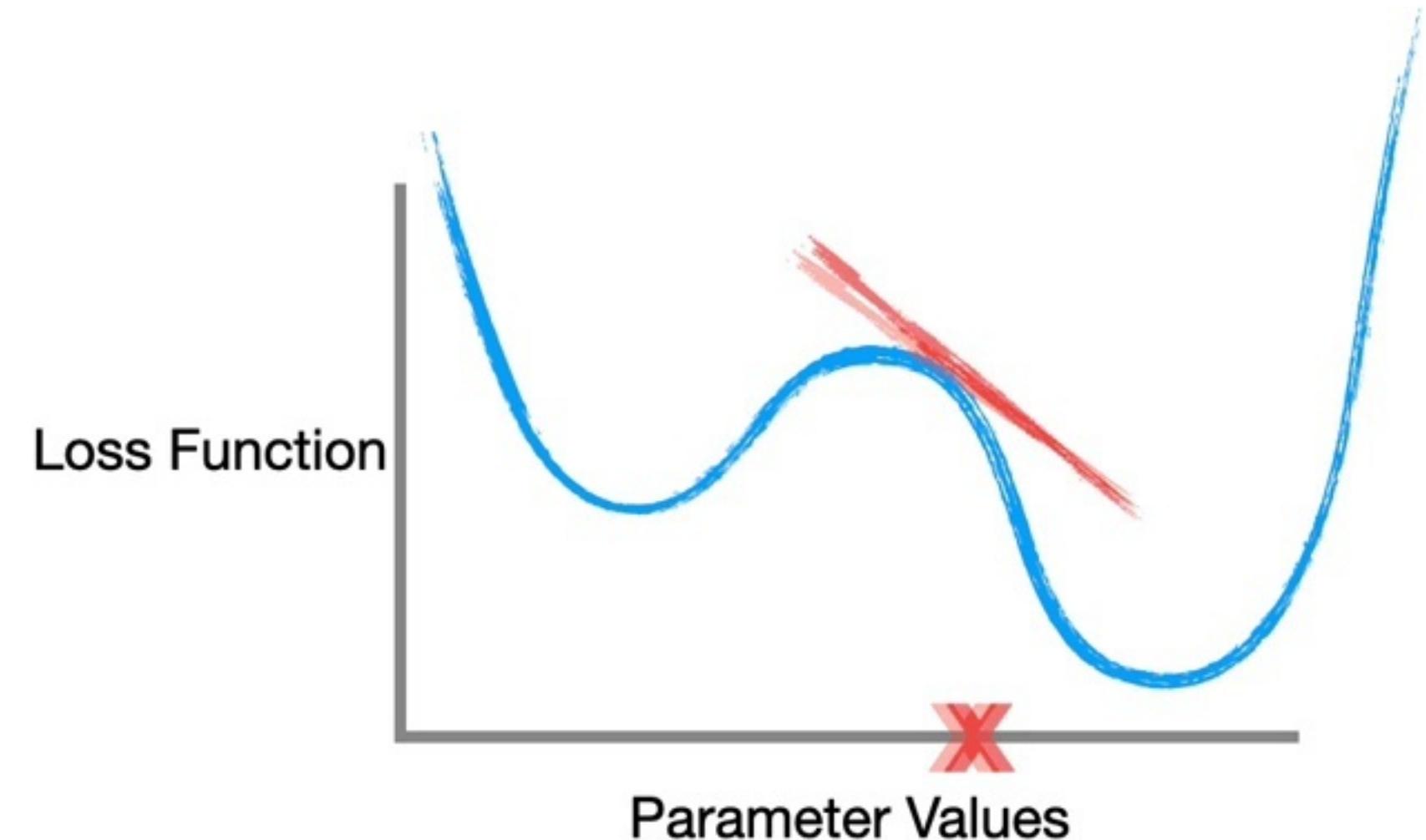


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



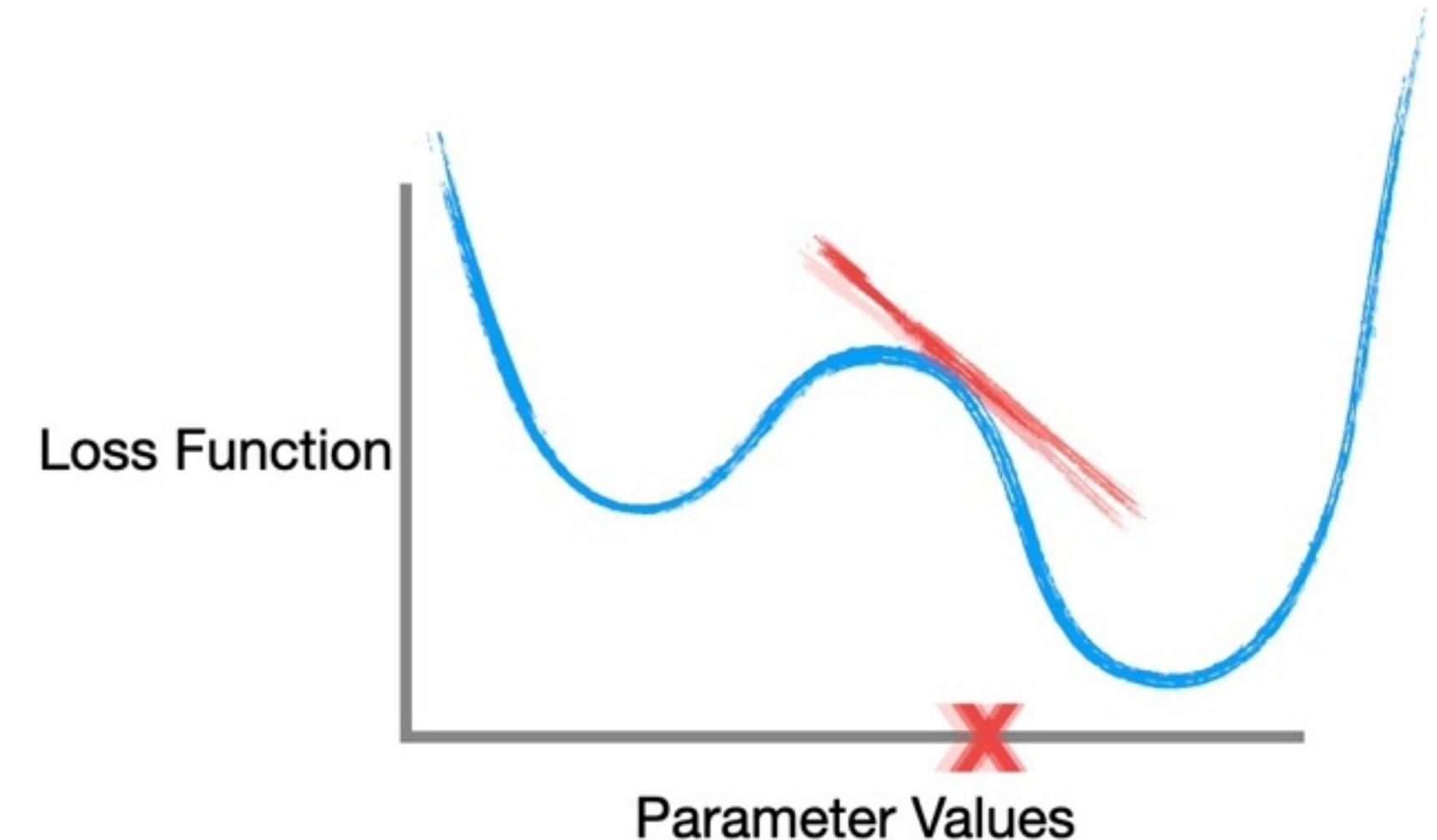


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



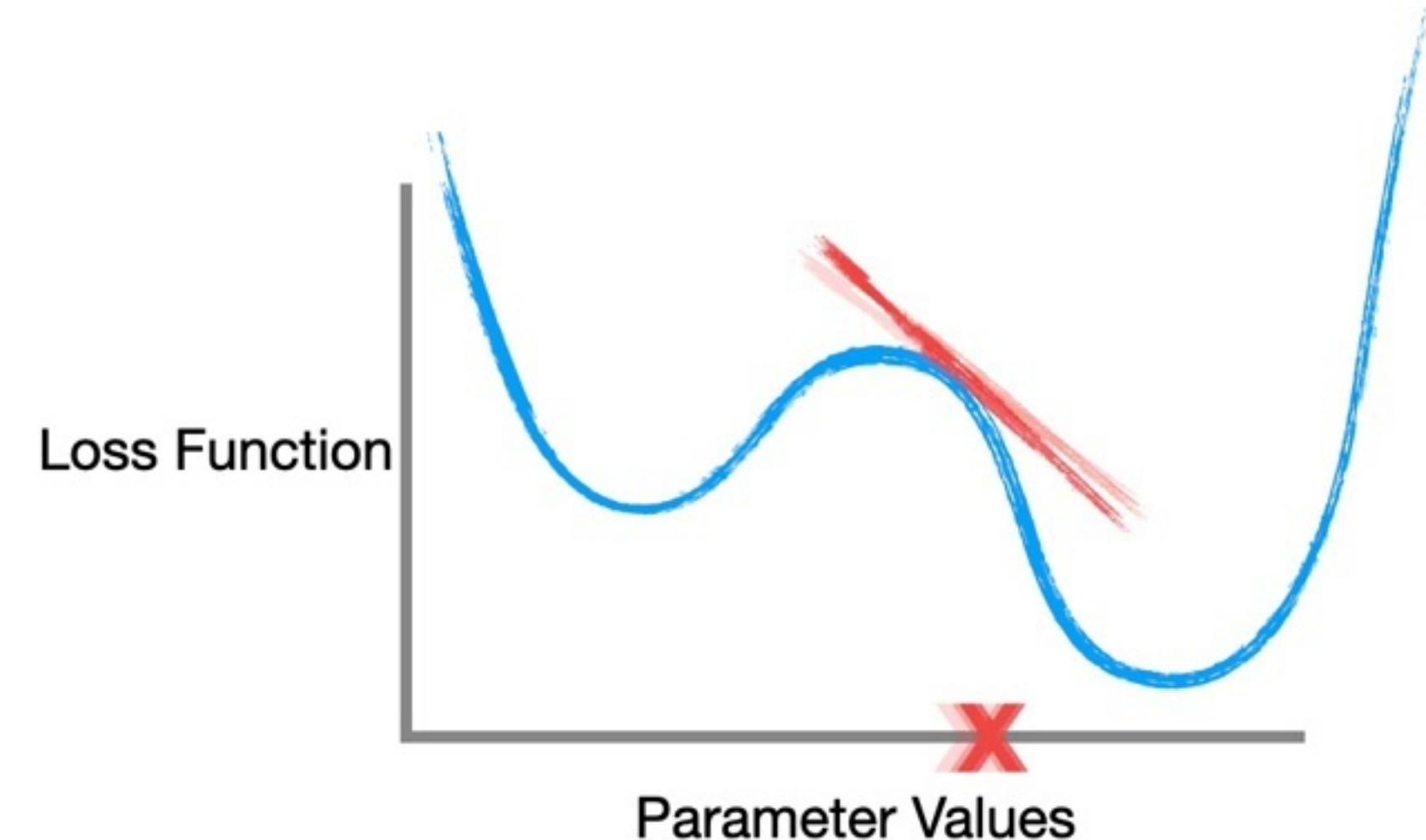


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



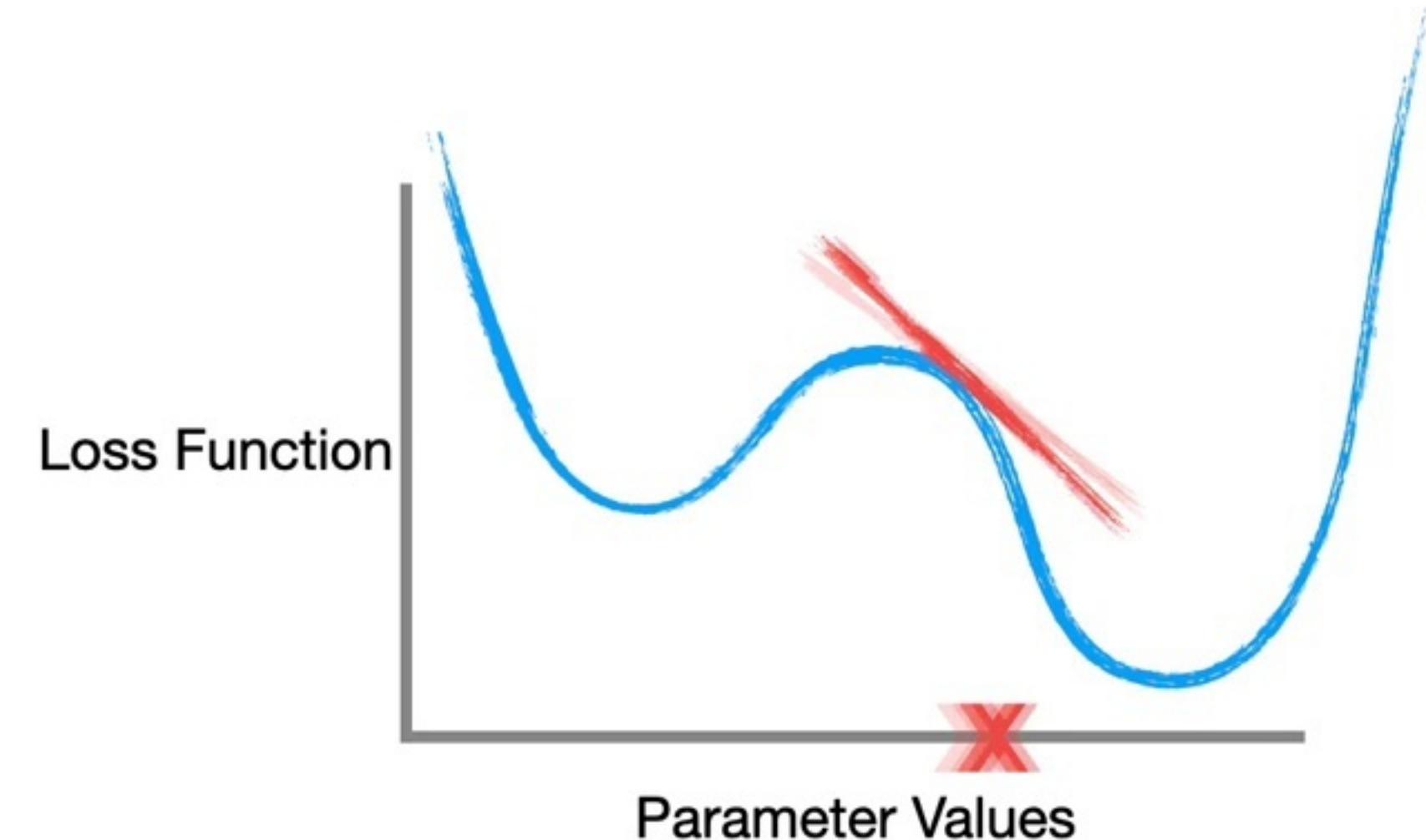


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



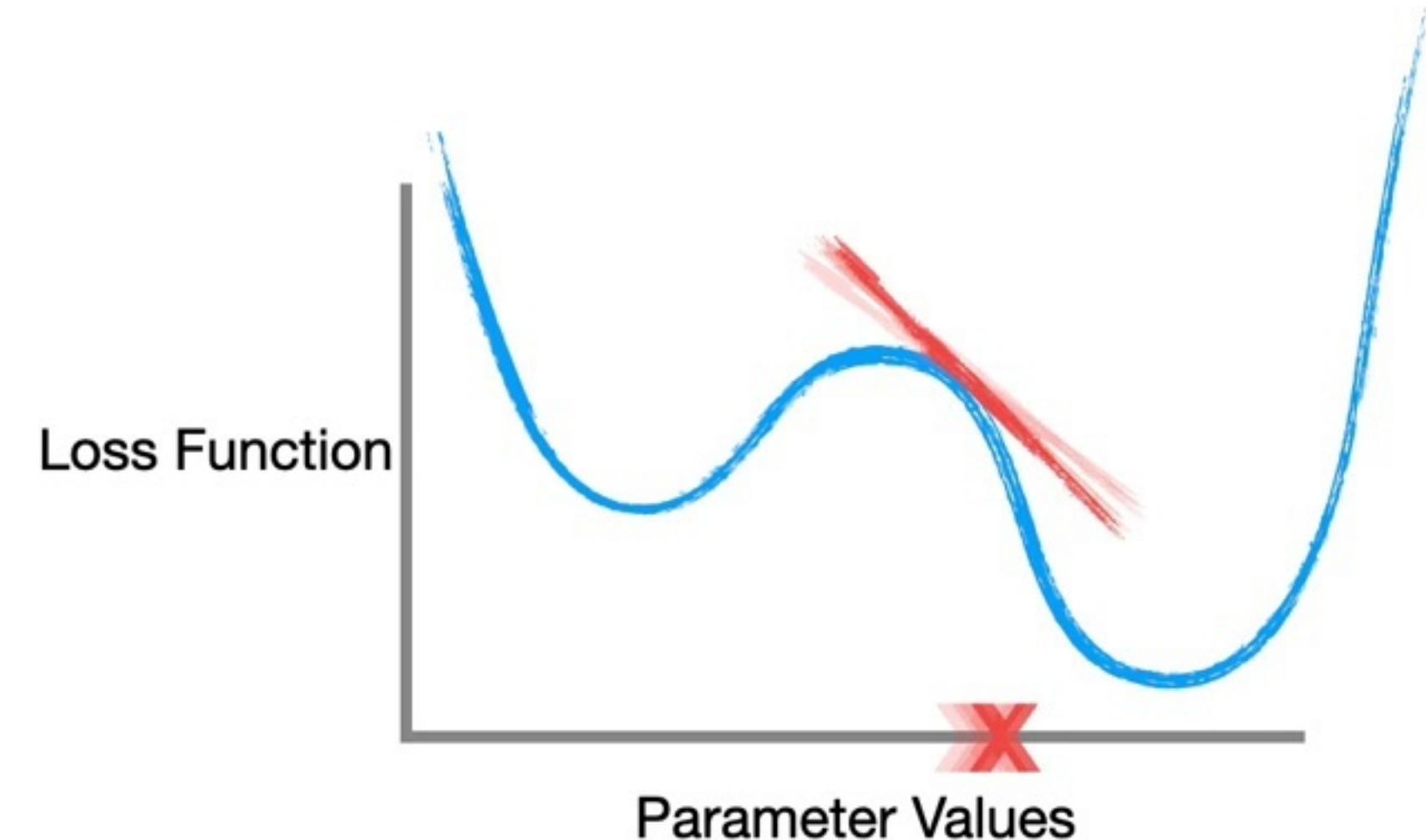


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



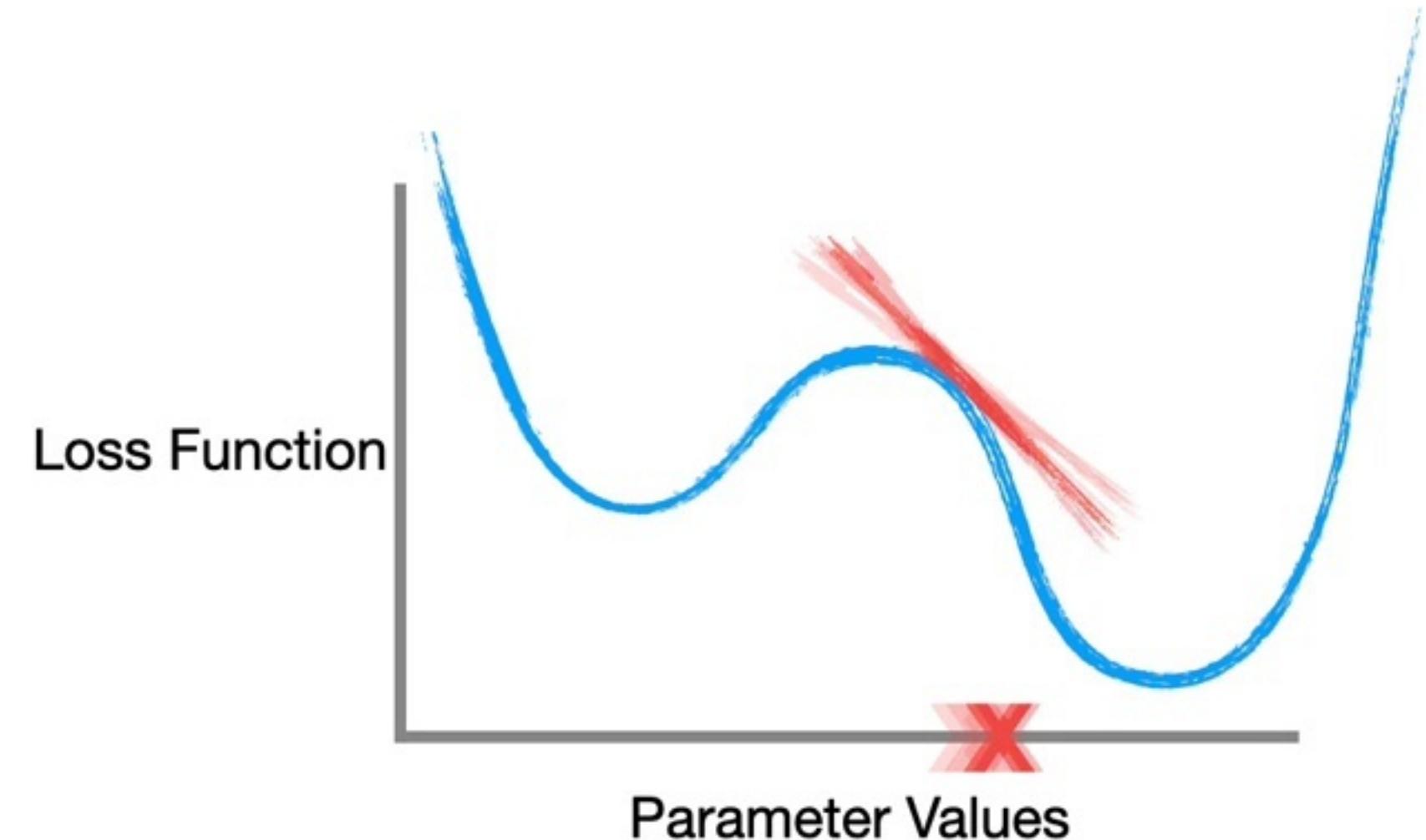


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



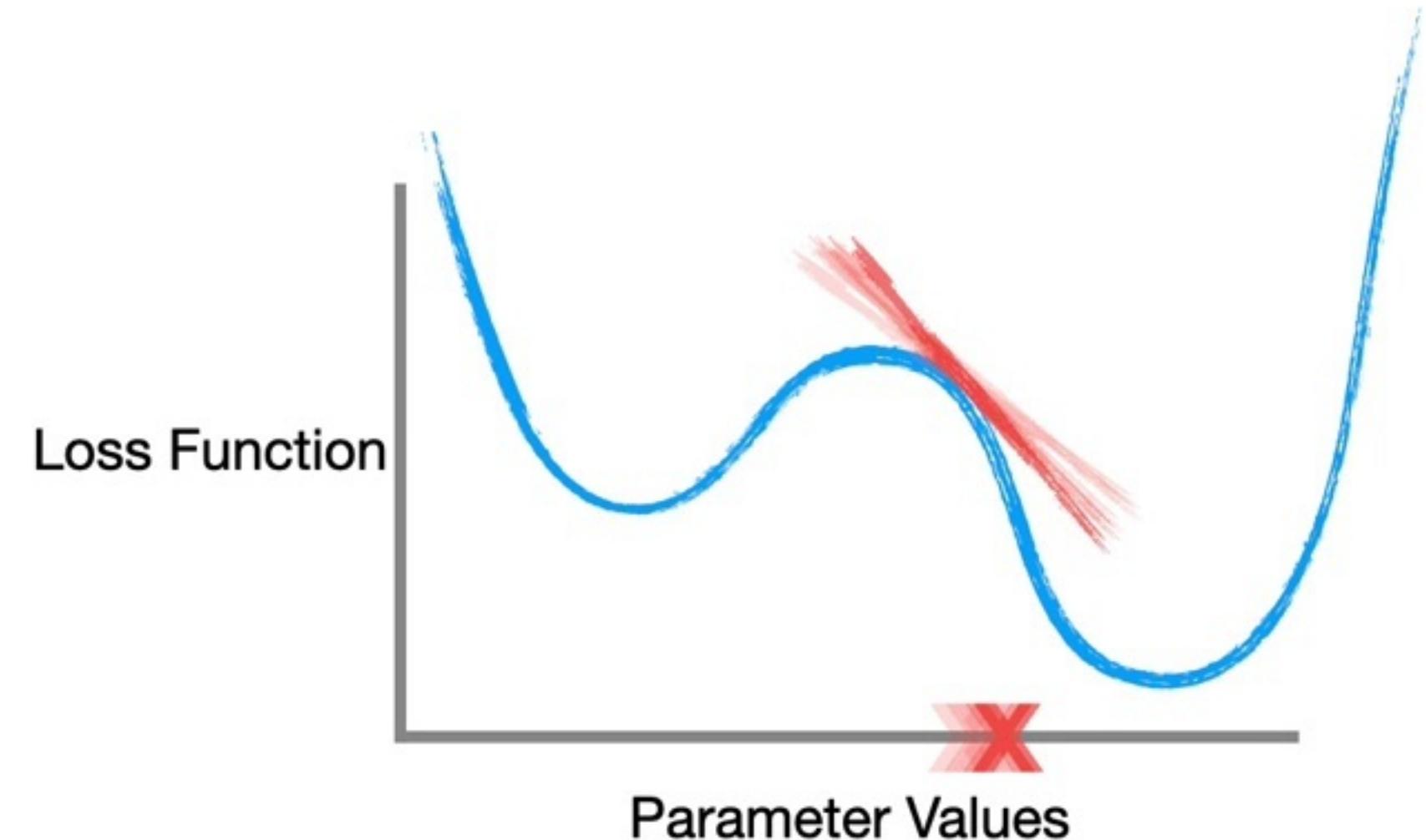


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



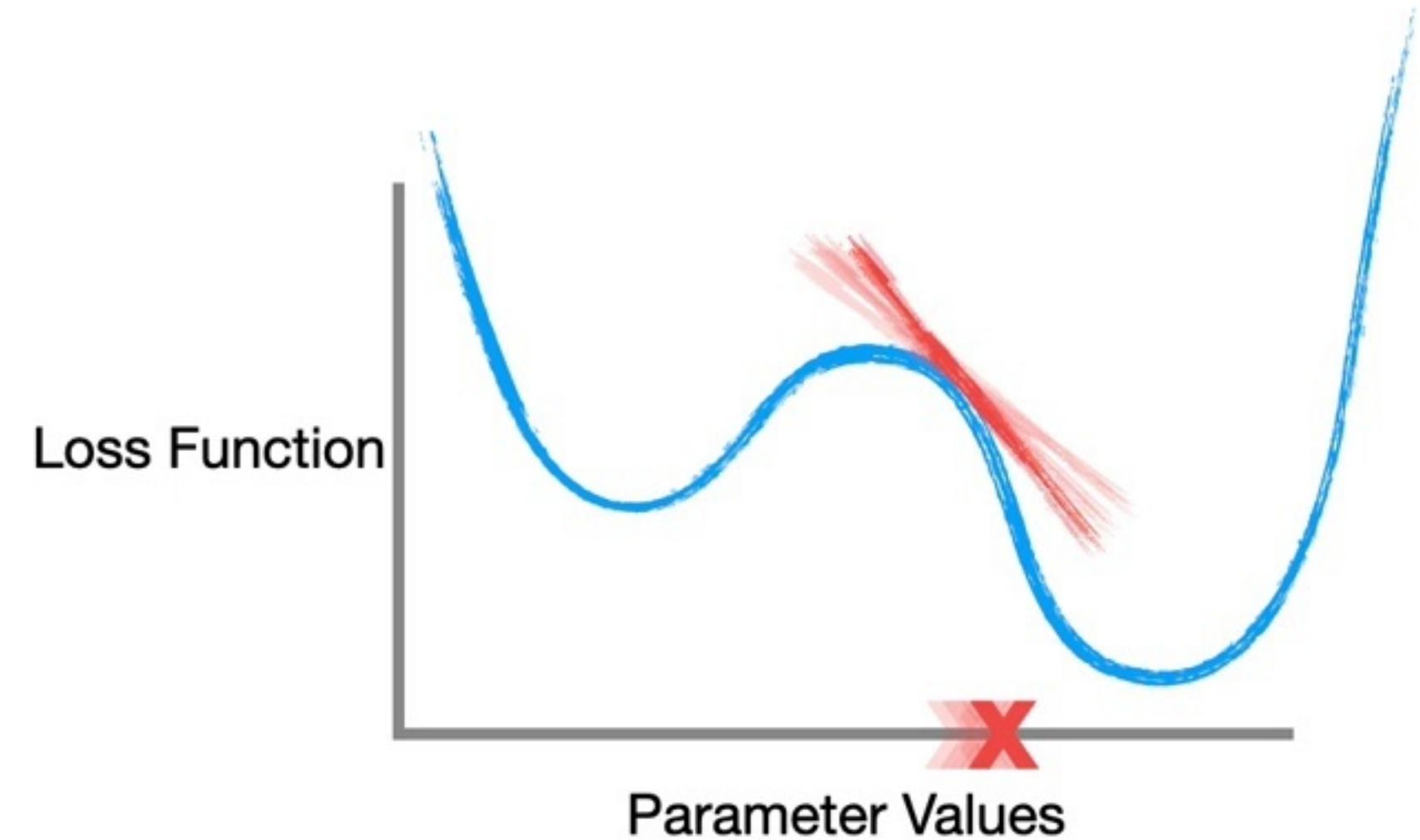


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.



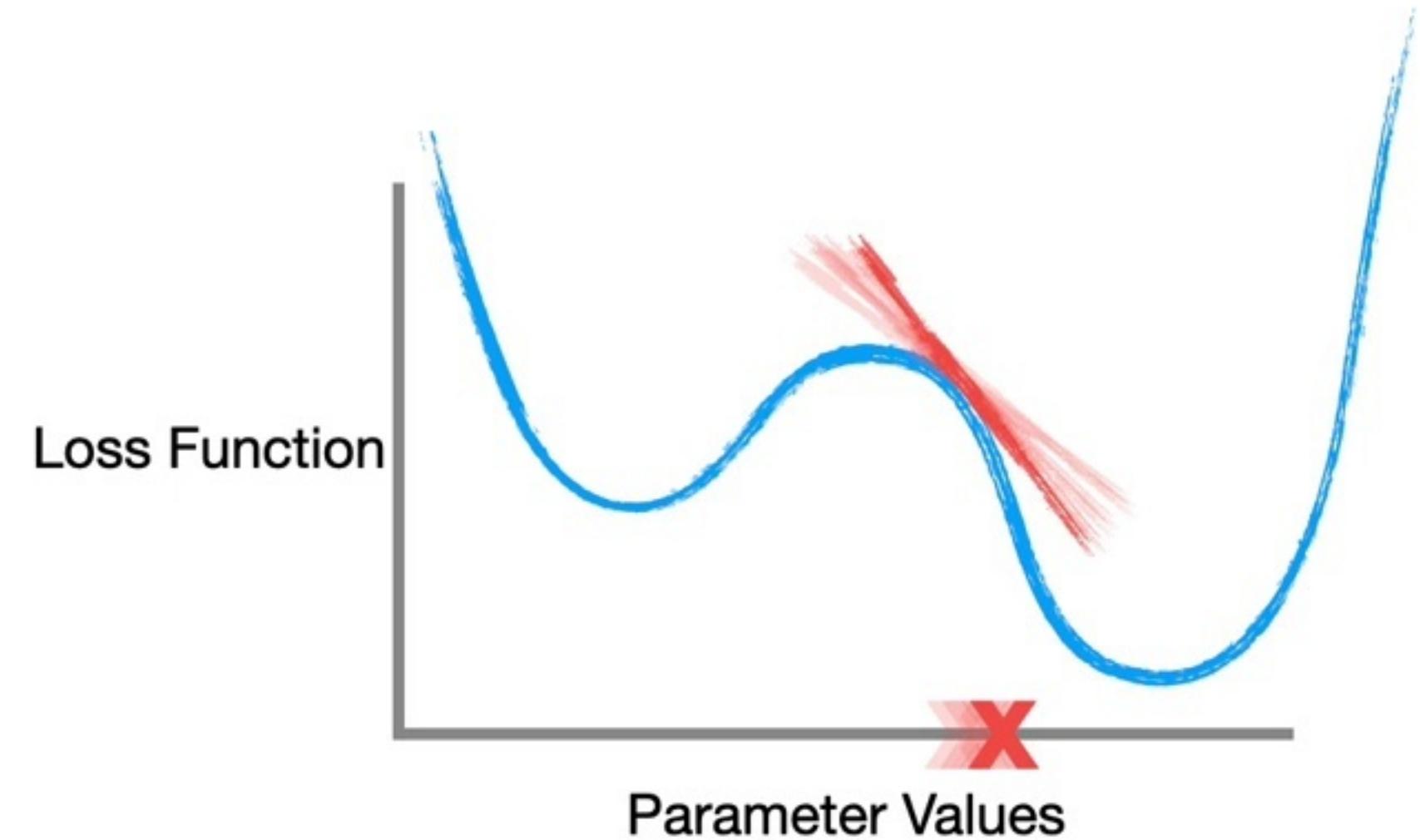


Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.





Now, when optimizing a parameter, instead of taking steps that are too large, we end up taking steps that are too small.





And as a result, we end up hitting the maximum number of steps we are allowed to take before we find the optimal value.

