# Deep Sequence Modeling

# Sequences in the Wild



Audio

# Sequences in the Wild

## Introduction to Deep Learning

Text

# A Sequence Modeling Problem:
## Predict the Next Word

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

given these words        predict the
                         next word

# Idea #1: Use a Fixed Window

"This morning I took my cat for a walk."

given these
two words

predict the
next word

One-hot feature encoding: tells us what each word is

$$[ 1\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0 ]$$

for        a

prediction

# Problem #1: Can't Model Long-Term Dependencies

"**France** is where I grew up, but I now live in Boston. I speak fluent ____."

We need information from **the distant past** to accurately predict the correct word.

# Idea #2: Use Entire Sequence as Set of Counts

"This morning I took my cat for a"

↓

"bag of words"

$$[\,0\;1\;0\;0\;1\;0\;0\;...\;0\;0\;1\;1\;0\;0\;0\;1\,]$$

↓

prediction

# Problem #2: Counts Don't Preserve Order

The food was good, not bad at all.

vs.

The food was bad, not good at all.

# Idea #3: Use a Really Big Fixed Window

"This morning I took my cat for a walk."

given these words · predict the next word

[ 1 0 0 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 0 1 0 ... ]

morning      I      took      this      cat

prediction

# Problem #3: No Parameter Sharing

[ 1 0 0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 0 1 0 ... ]

this      morning      took      the      cat

Each of these inputs has a **separate parameter:**

[ 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 1 ... ]
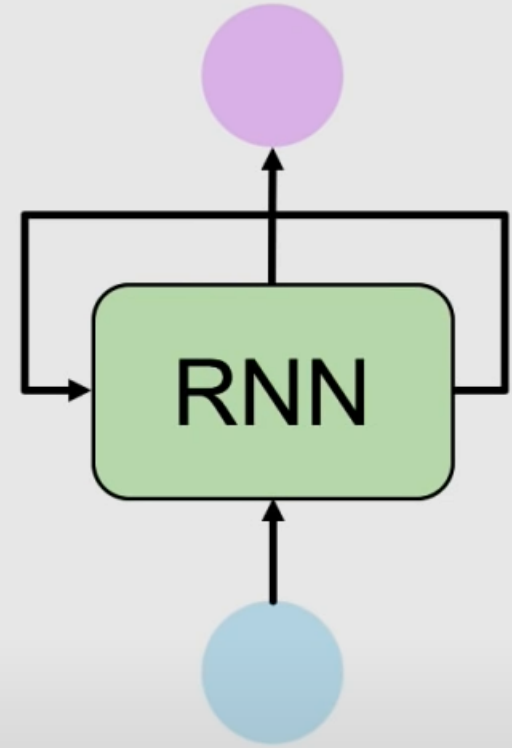
this      morning

Things we learn about the sequence **won't transfer** if
they appear **elsewhere** in the sequence.

# Sequence Modeling: Design Criteria

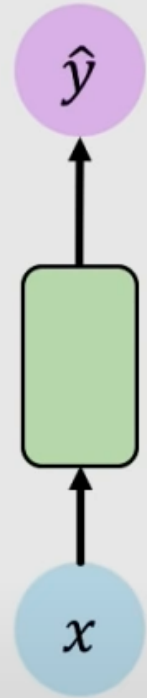To model sequences, we need to:

1.  Handle **variable-length** sequences

2.  Track **long-term** dependencies

3.  Maintain information about **order**
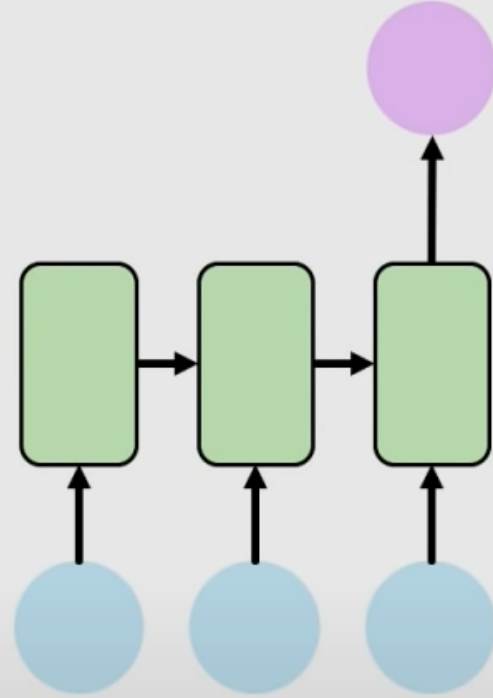
4.  **Share parameters** across the sequence

**Recurrent Neural Networks (RNNs)** as
an approach to sequence modeling problems
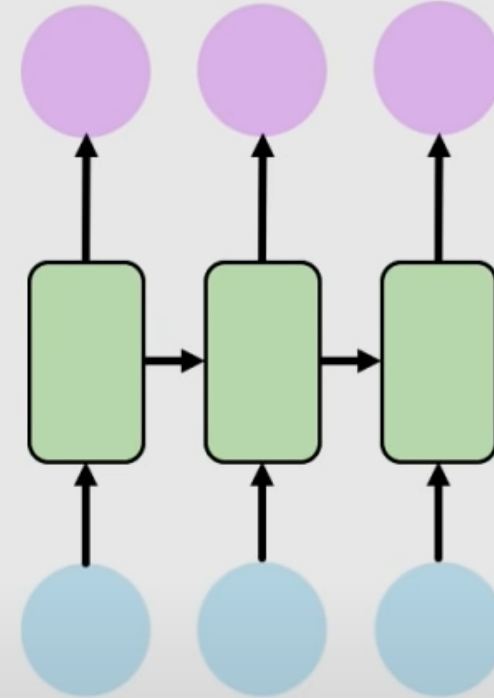
# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks for Sequence Modeling



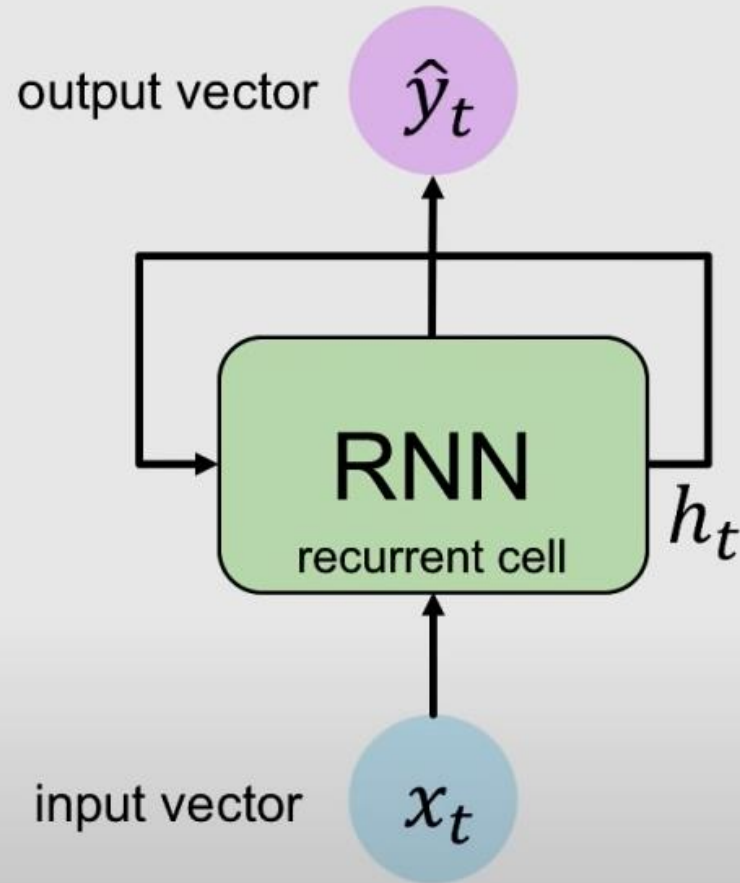One to One
"Vanilla" neural network

Many to One
*Sentiment Classification*

Many to Many
*Music Generation*

... and many other architectures and applications

# Recurrent Neural Network (RNN)

output vector $\hat{y}_t$

RNN
recurrent cell

$h_t$

input vector $x_t$

Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(h_{t-1}, x_t)$$

cell state — function parameterized by W — old state — input vector at time step $t$

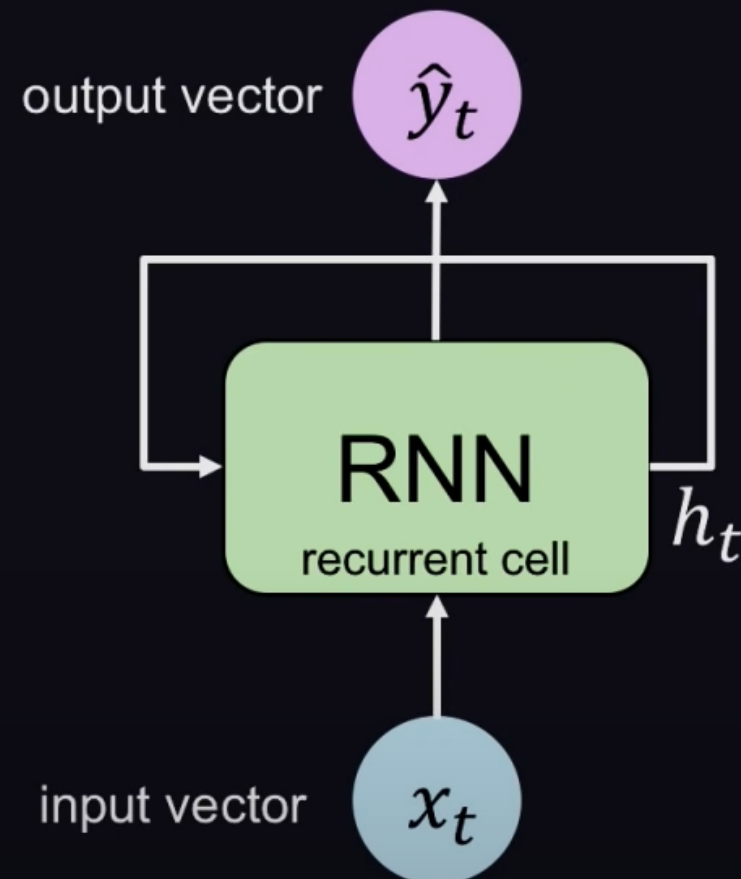Note: the same function and set of parameters are used at every time step

# RNN Intuition

```python
my_rnn = RNN()
hidden_state = [0, 0, 0, 0]


sentence = ["I", "love", "recurrent", "neural"]


for word in sentence:
    prediction, hidden_state = my_rnn(word, hidden_state)


next_word_prediction = prediction
# >>> "networks!"
```
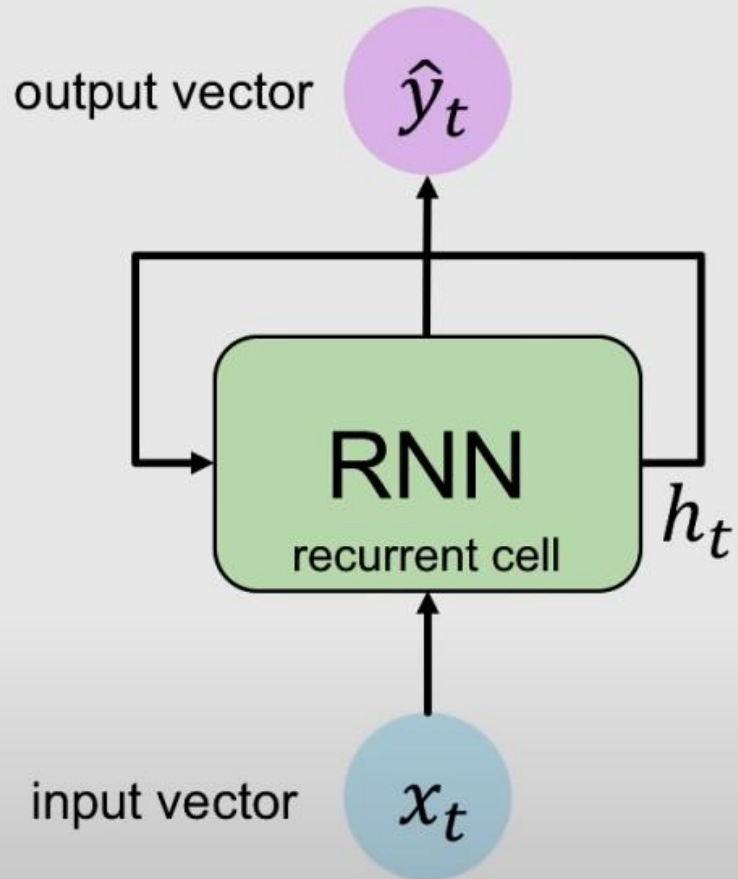
output vector $\hat{y}_t$

RNN
recurrent cell

$h_t$

input vector $x_t$

# RNN State Update and Output



output vector $\hat{y}_t$

**RNN** recurrent cell $h_t$

input vector $x_t$

**Output Vector**

$$\hat{y}_t = \boldsymbol{W_{hy}^T} h_t$$
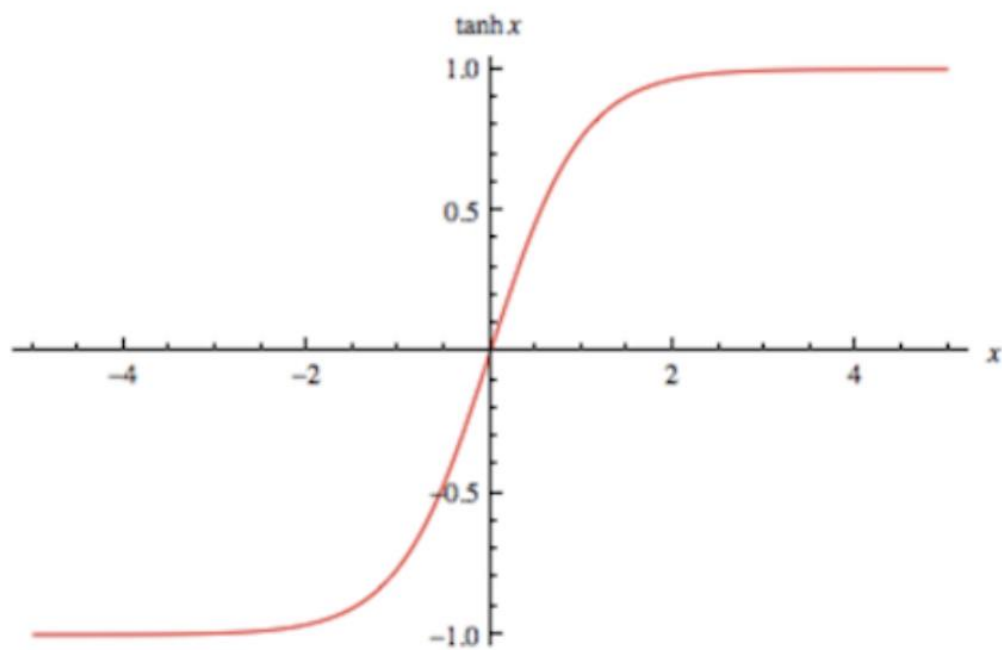
**Update Hidden State**

$$h_t = \tanh(\boldsymbol{W_{hh}^T} h_{t-1} + \boldsymbol{W_{xh}^T} x_t)$$
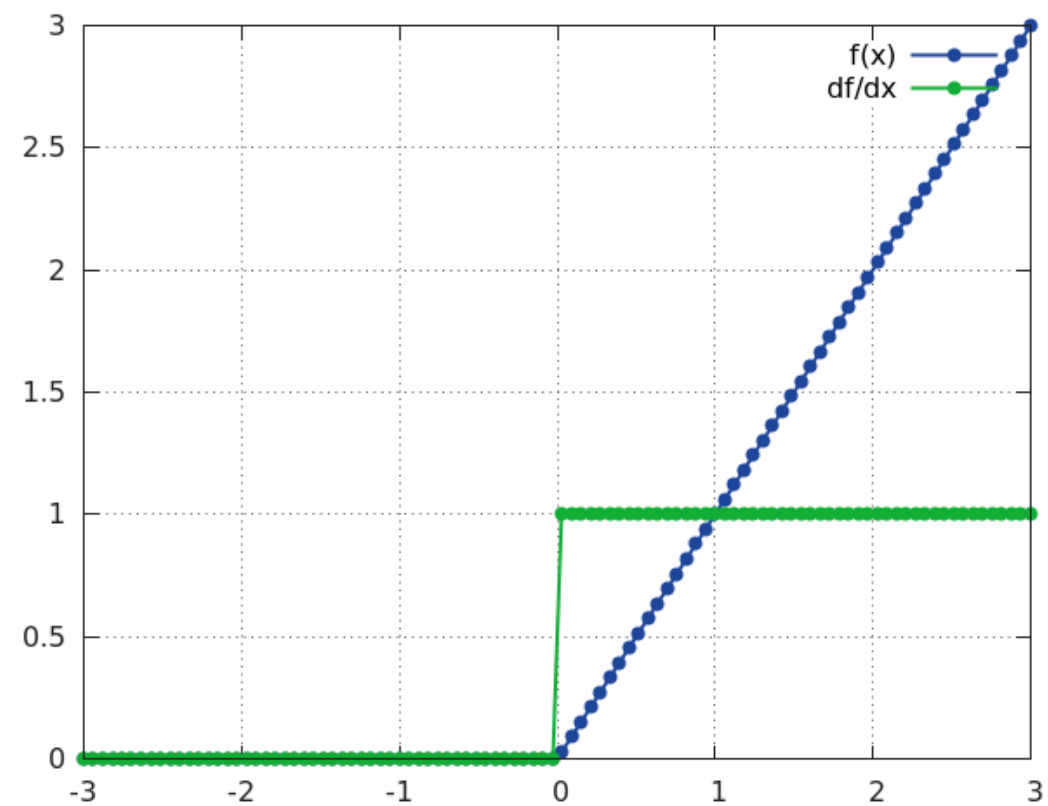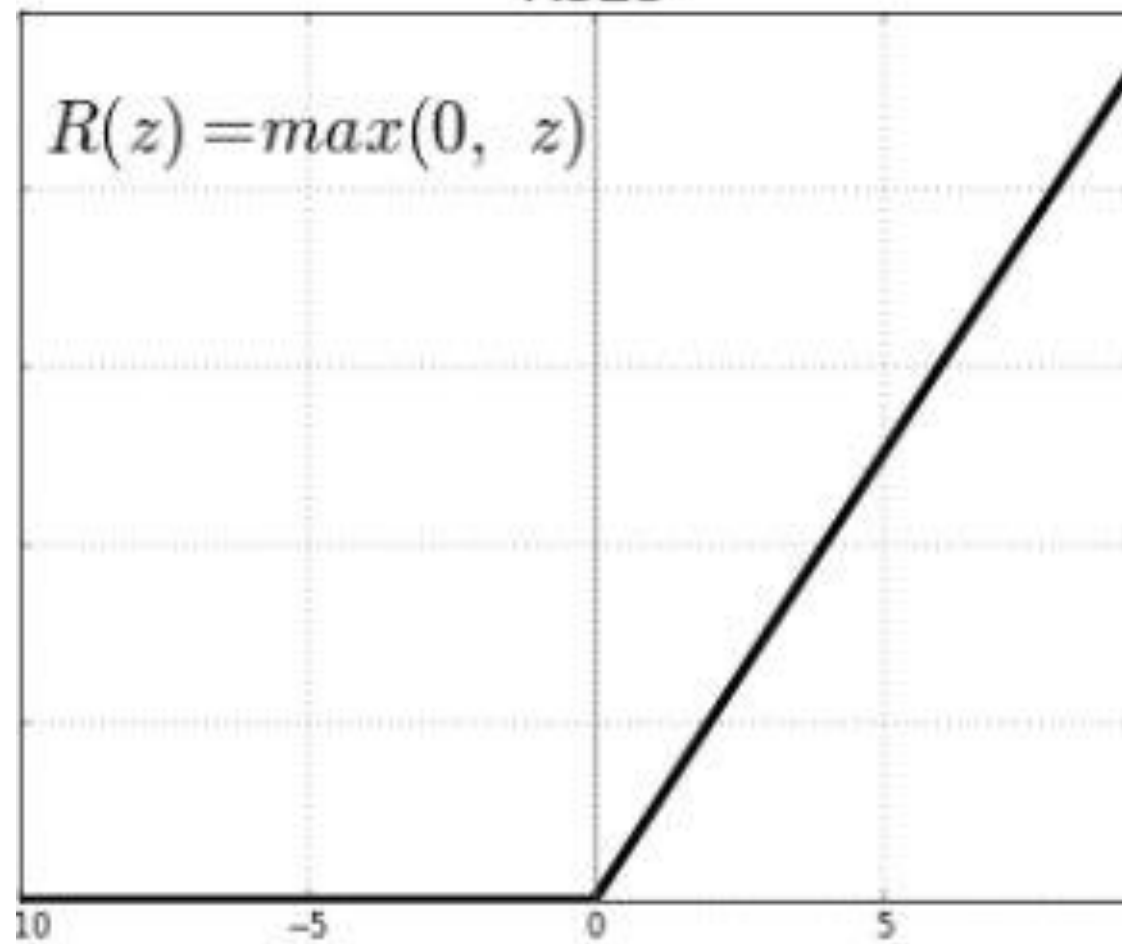
**Input Vector**

$$x_t$$

$$tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

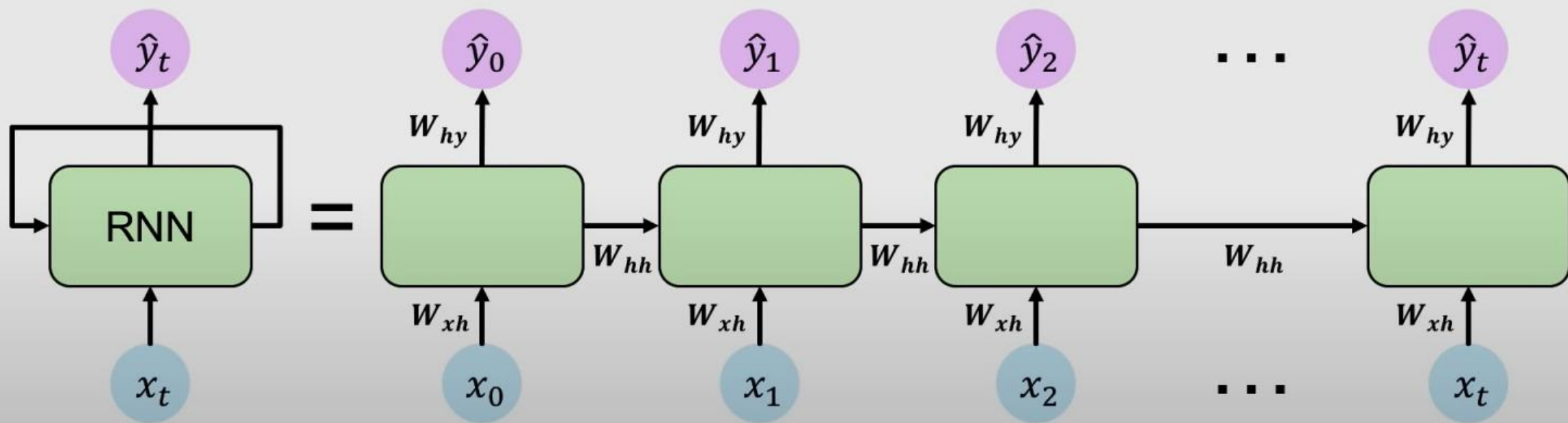Mathematical formula of the Tanh function

# ReLU

$$R(z) = max(0, \ z)$$

# RNNs from Scratch

```python
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_hy = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh( self.W_hh * self.h + self.W_xh * x )

        # Compute the output
        output = self.W_hy * self.h

        # Return the current output and hidden state
        return output, self.h
```
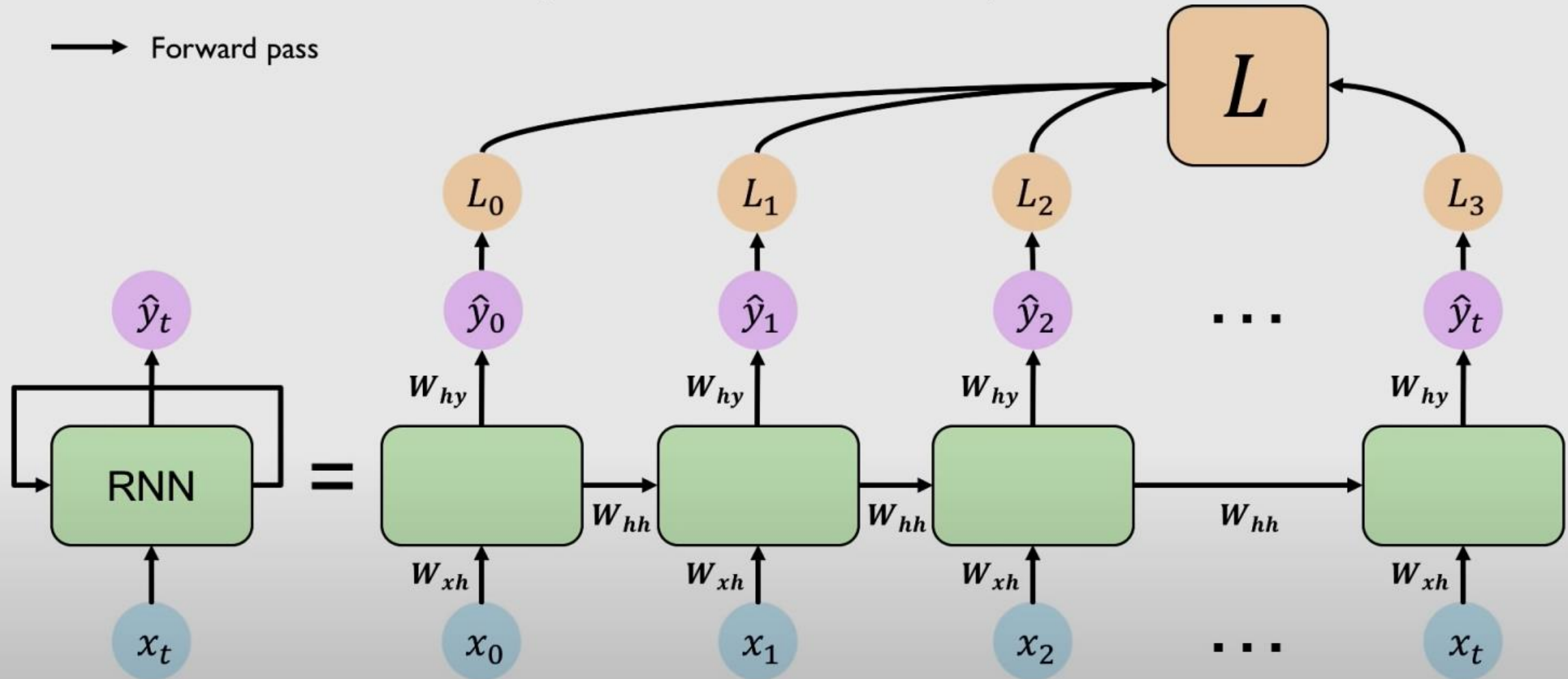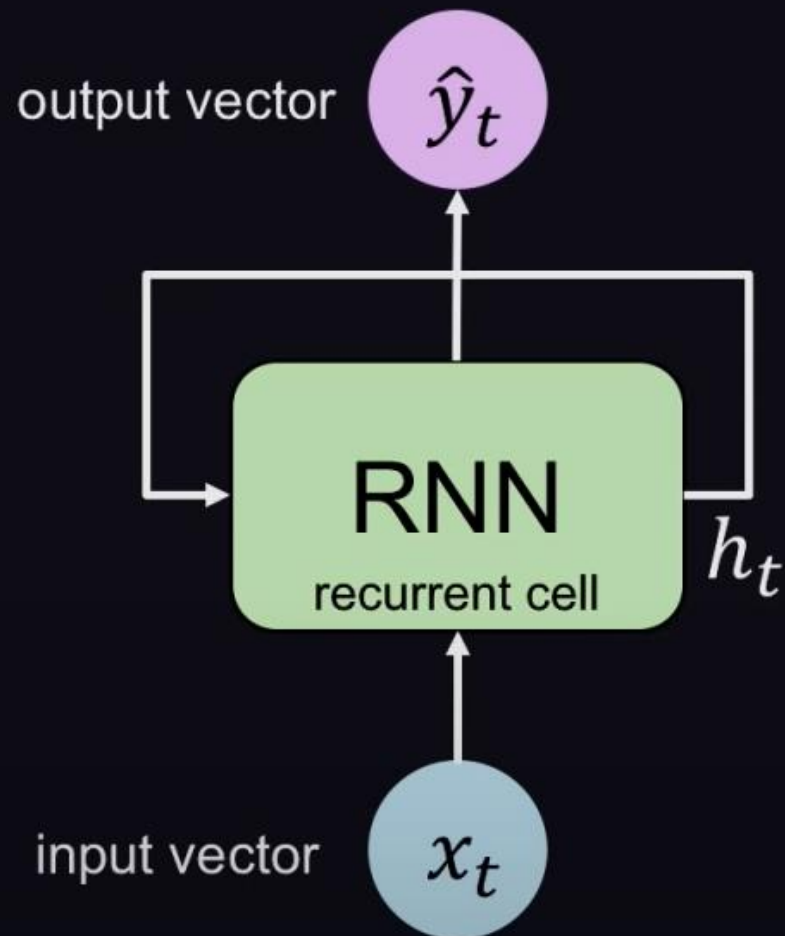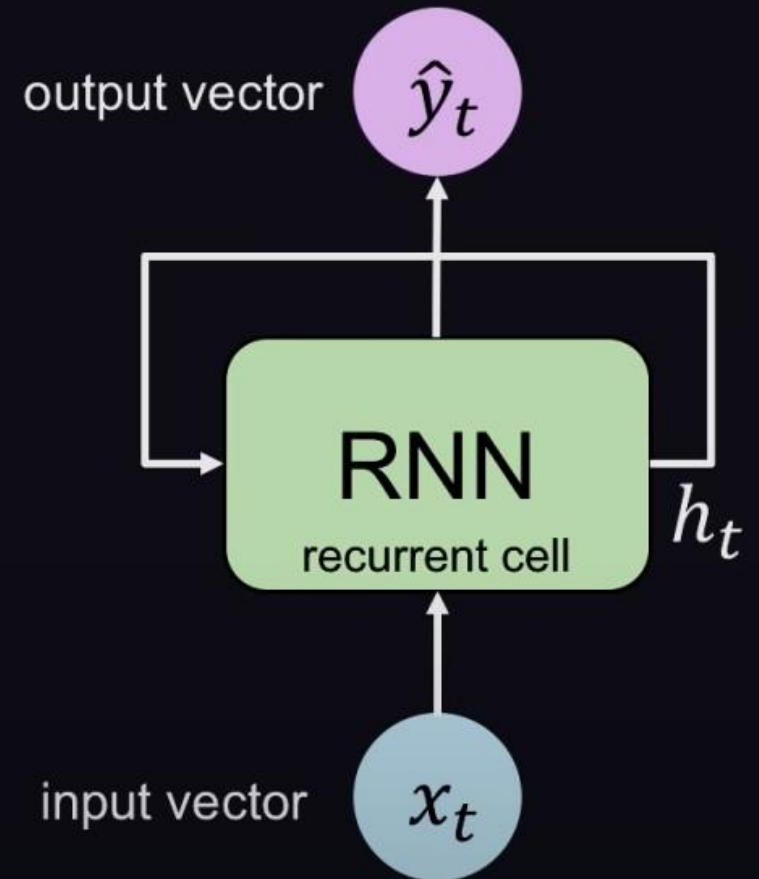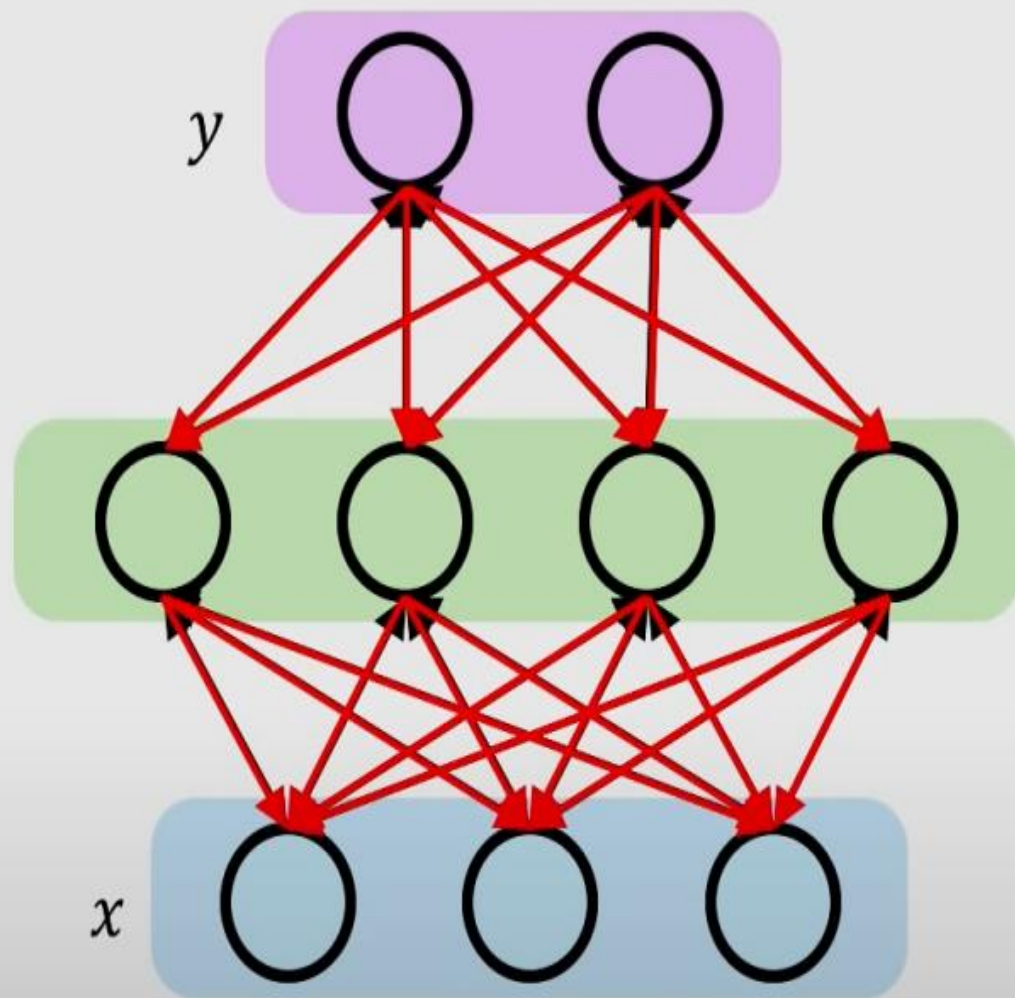
output vector $\hat{y}_t$

RNN
recurrent cell

$h_t$

input vector $x_t$

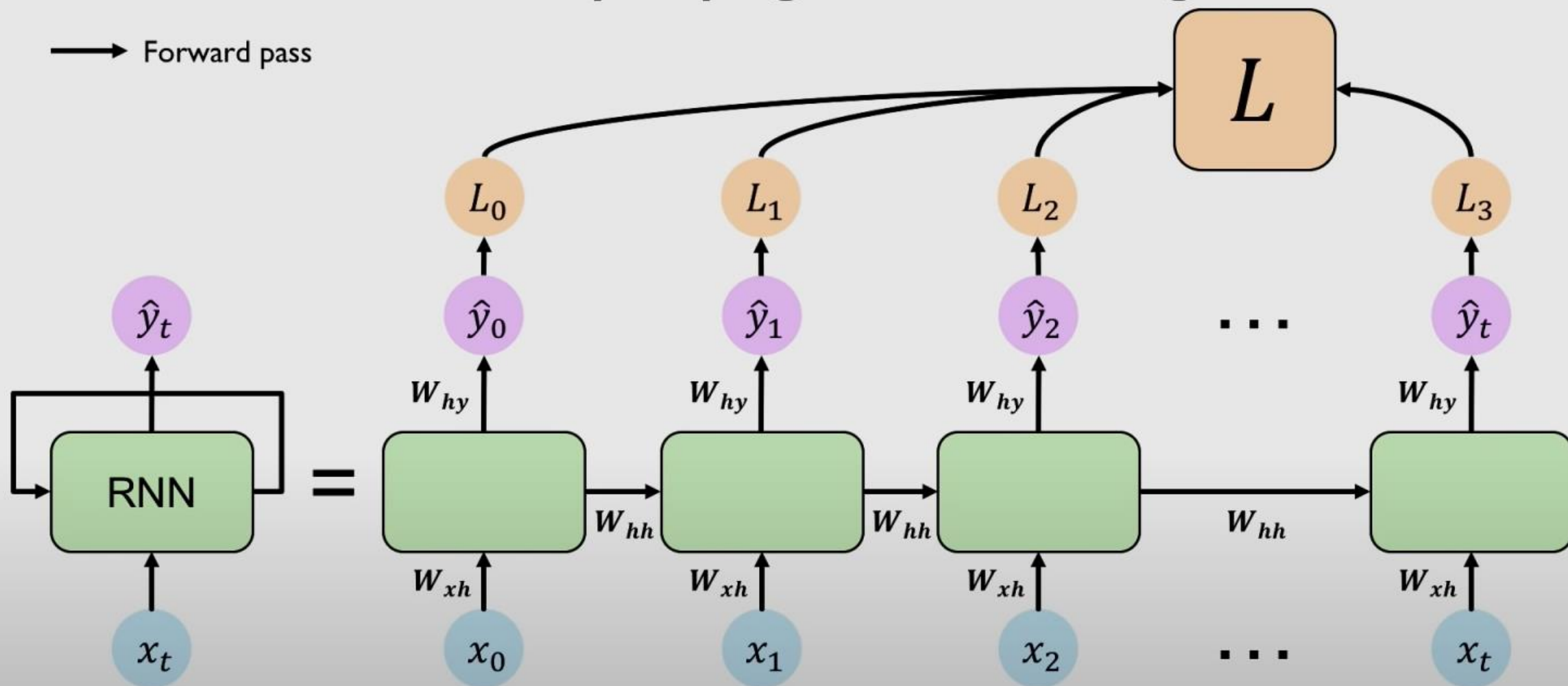Backpropagation Through Time (BPTT)

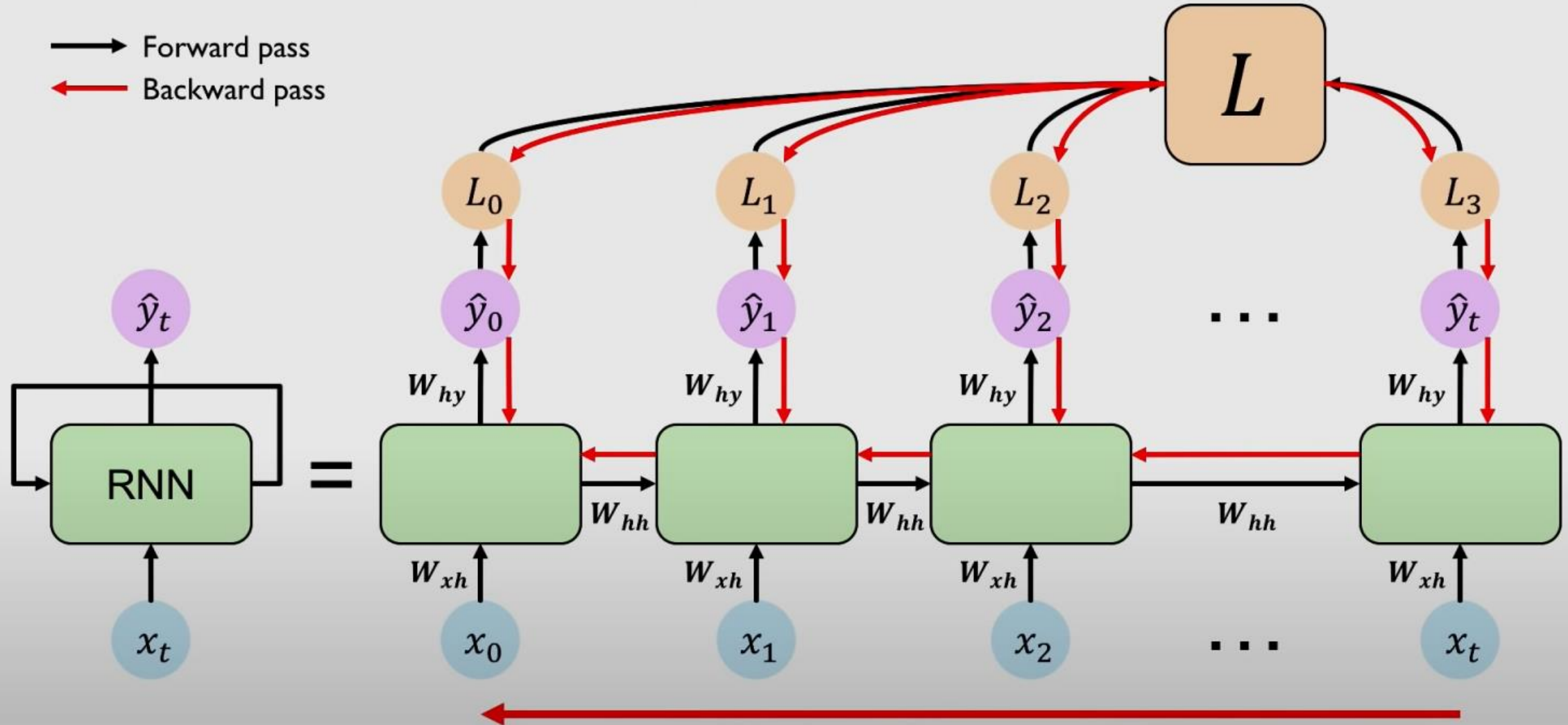# Recall: Backpropagation in Feed Forward Models



Backpropagation algorithm:

1.  Take the derivative (gradient) of the loss with respect to each parameter

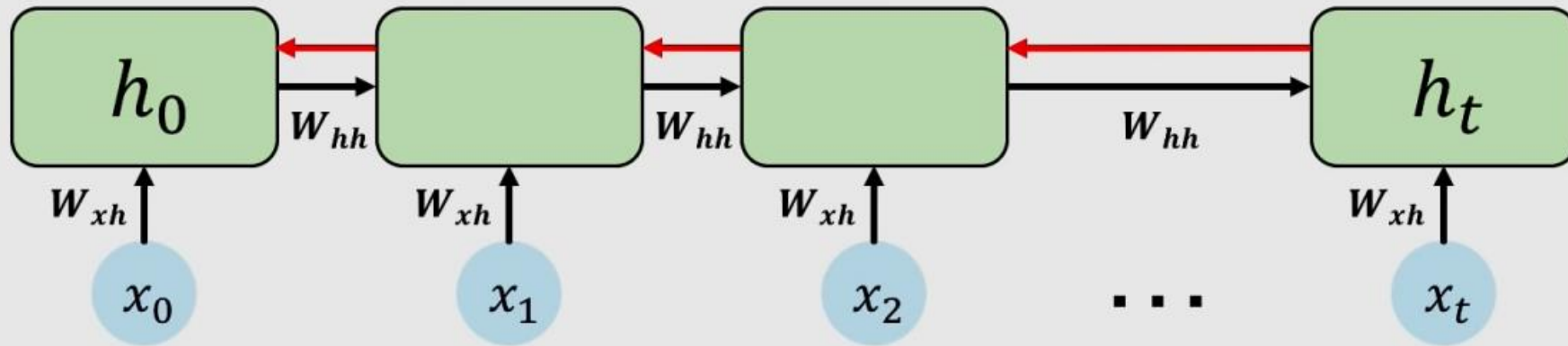2.  Shift parameters in order to minimize loss

RNNs: Backpropagation Through Time

# Standard RNN Gradient Flow: Exploding Gradients



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$ + repeated gradient computation!**

Many values > 1:
**exploding gradients**

**Gradient clipping** to scale big gradients

Many values < 1:
**vanishing gradients**

1. Activation function
2. Weight initialization
3. Network architecture

# The Problem of Long-Term Dependencies

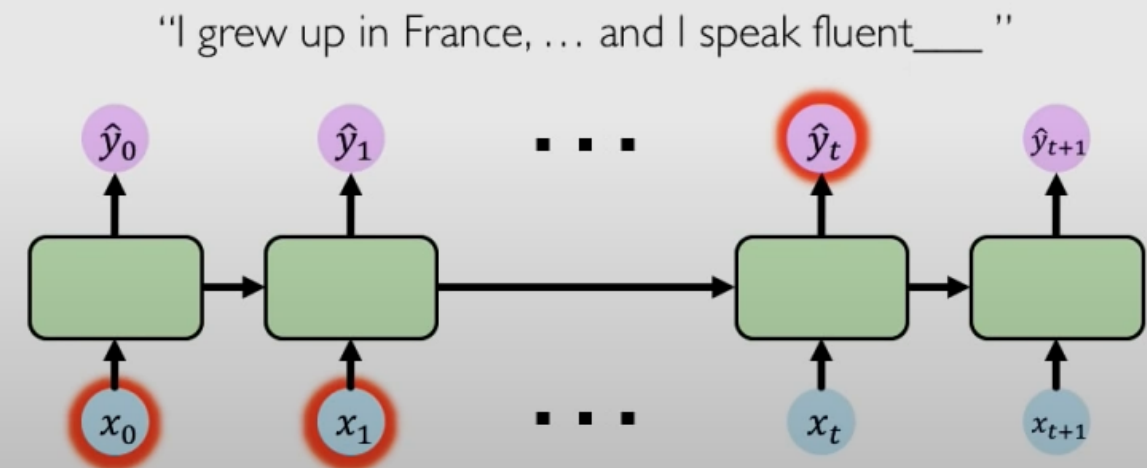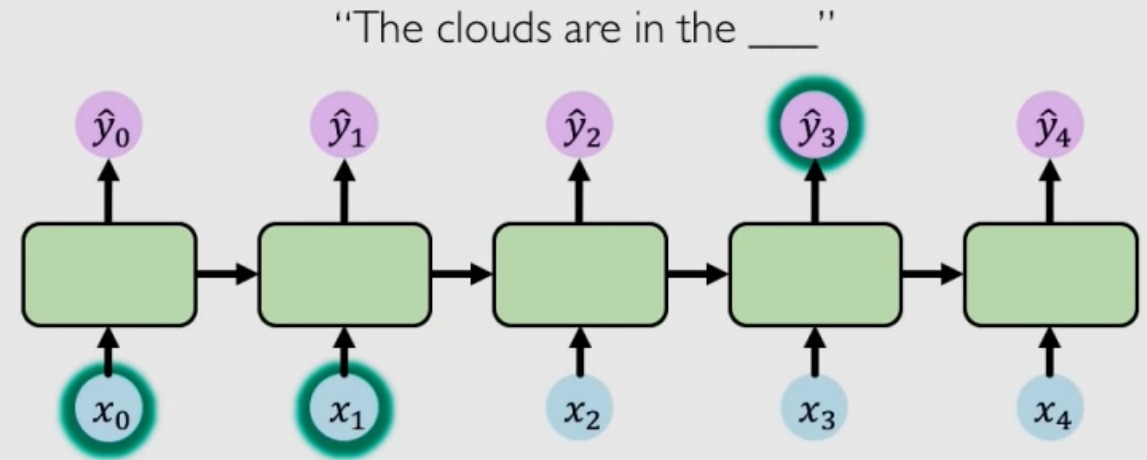**Why are vanishing gradients a problem?**
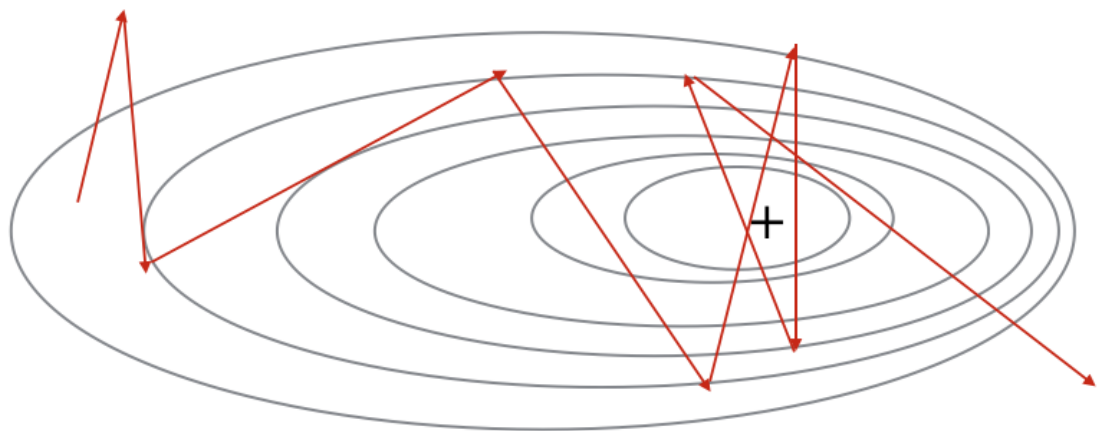
Multiply many **small numbers** together

↓

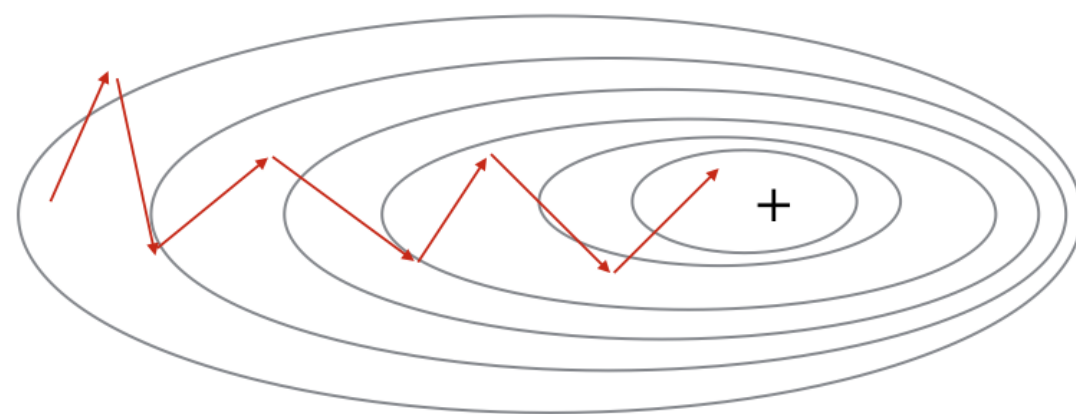Errors due to further back time steps have smaller and smaller gradients

↓

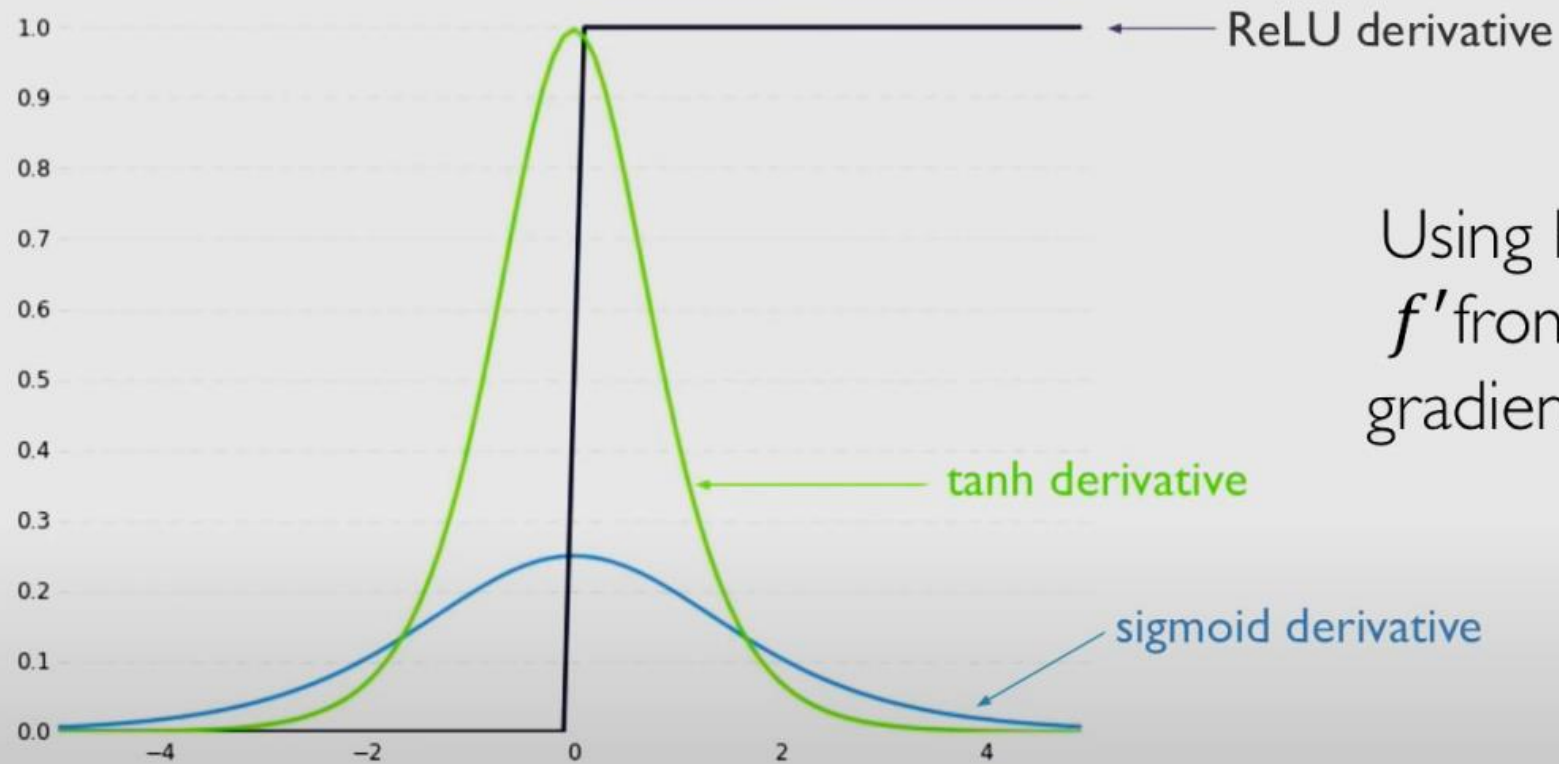Bias parameters to capture short-term dependencies

"The clouds are in the ___"



"I grew up in France, … and I speak fluent___ "

Without gradient clipping

With gradient clipping

# Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.