# Intro. to Sequence Modelling

Rina BUOY

# Why Sequence Modelling?



Given an image of a ball,
can you predict where it will go next?

# Why Sequence Modelling?

# Why Sequence Modelling?

**Previous positions**
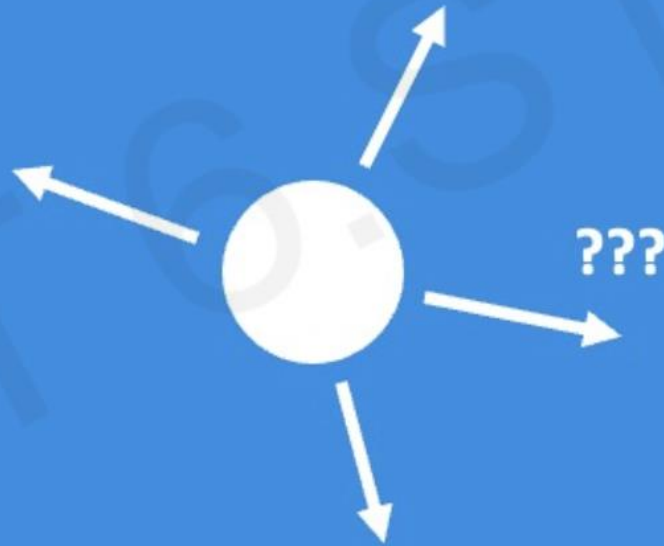


Given an image of a ball,
can you predict where it will go next?

# Why Sequence Modelling?

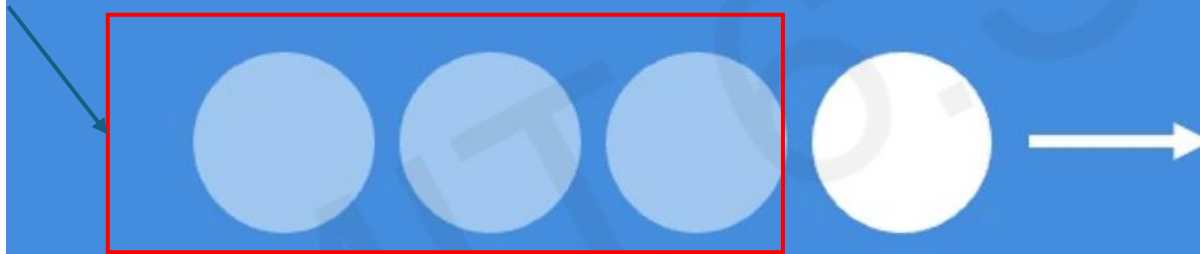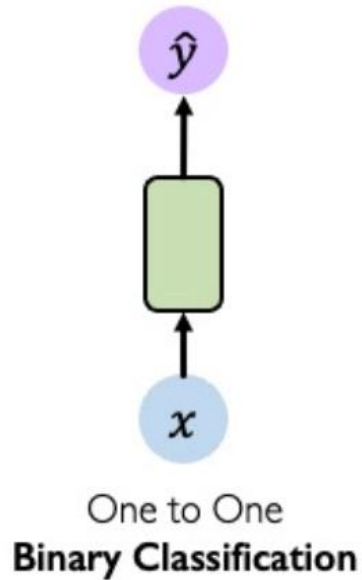**Given an image of a ball, can you predict where it will go next?**
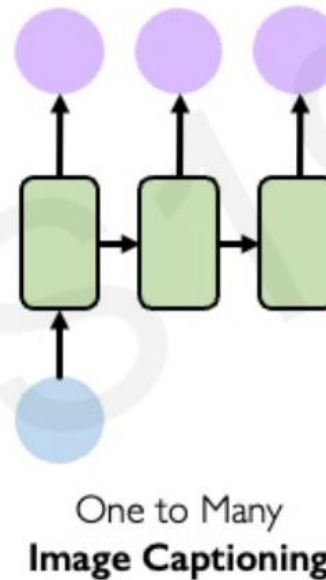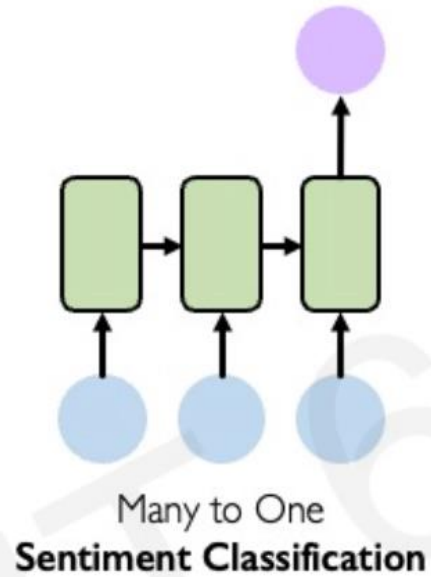
**Previous positions**

# Other sequential data



Audio

# Sequences in Applications



One to One
**Binary Classification**

"Will I pass this class?"
Student → Pass?

Many to One
**Sentiment Classification**

Ivar Hagendoorn
@IvarHagendoorn
Follow

The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online
introtodeeplearning.com
12:45 PM - 12 Feb 2018

One to Many
**Image Captioning**

"A baseball player throws a ball."

Many to Many
**Machine Translation**

# Perceptron

# Simple Neural Networks



$$\boldsymbol{x} \in \mathbb{R}^m \qquad \hat{\boldsymbol{y}} \in \mathbb{R}^n$$

# Simple Neural Networks



$$\boldsymbol{x_t} \in \mathbb{R}^m \qquad\qquad \boldsymbol{\hat{y}_t} \in \mathbb{R}^n$$

# Handling Multiple Timesteps



output vector    $\hat{y}_t$        $\hat{y}_0$        $\hat{y}_1$        $\hat{y}_2$

input vector    $x_t$        $x_0$        $x_1$        $x_2$

$$\hat{y}_t = f(x_t)$$

# Adding Context



output vector — $\hat{y}_t$

input vector — $x_t$

$\hat{y}_0$   $\hat{y}_1$   $\hat{y}_2$

$h_0$   $h_1$

$x_0$   $x_1$   $x_2$

$$\hat{y}_t = f(x_t,\ h_{t-1})$$

output   input   past memory

# Recurrent Neurons



output vector   $\hat{y}_t$

recurrent cell   $h_t$

input vector   $x_t$

$\hat{y}_0$   $\hat{y}_1$   $\hat{y}_2$

$h_0$   $h_1$

$x_0$   $x_1$   $x_2$

$$\hat{y}_t = f(x_t, \; h_{t-1})$$

output   input   past memory

# Recurrent Neural Networks (RNN)

output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

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

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

cell state     function    input     old state
              with weights
                  W

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

RNNs have a **state**, $h_t$, that is updated **at each time step** as a sequence is processed
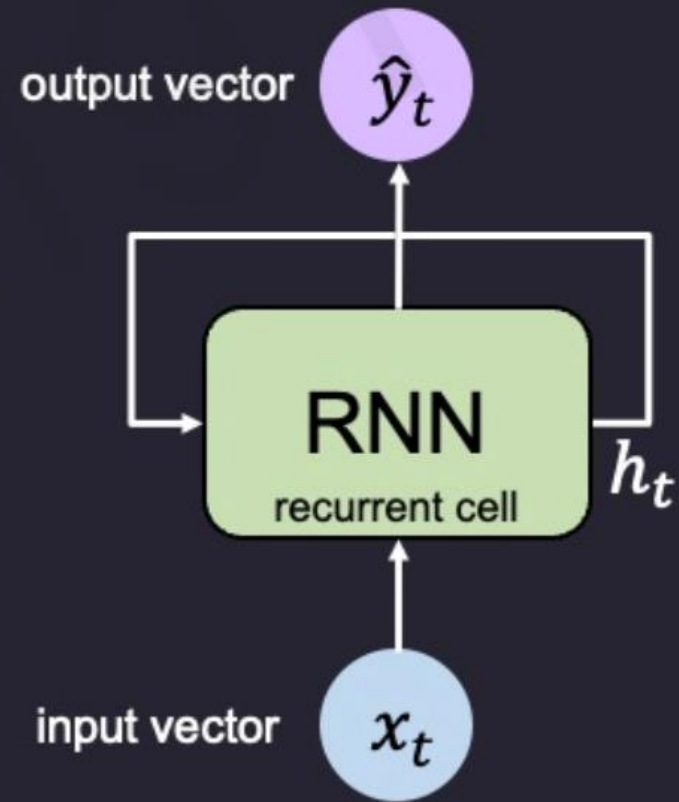
# Recurrent Neural Networks



```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!"
```
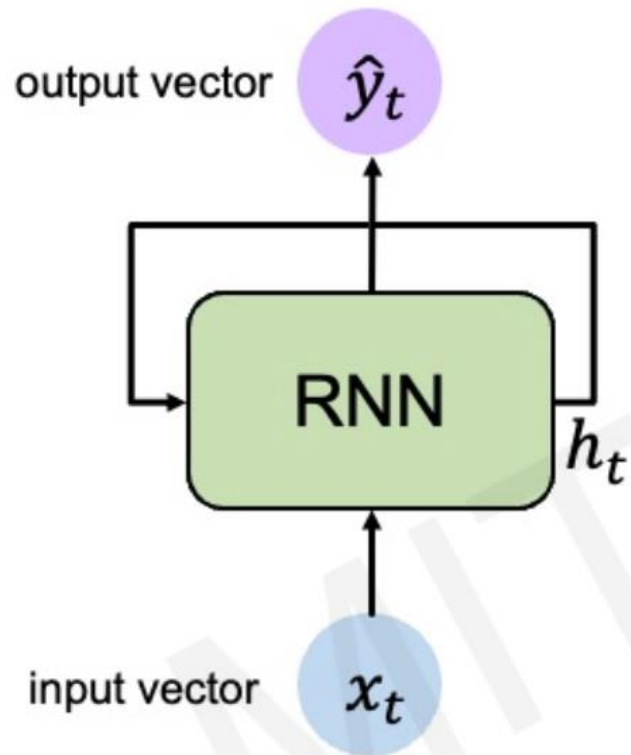
output vector $\hat{y}_t$

RNN
recurrent cell

$h_t$

input vector $x_t$

# Recurrent Neural Networks



output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$
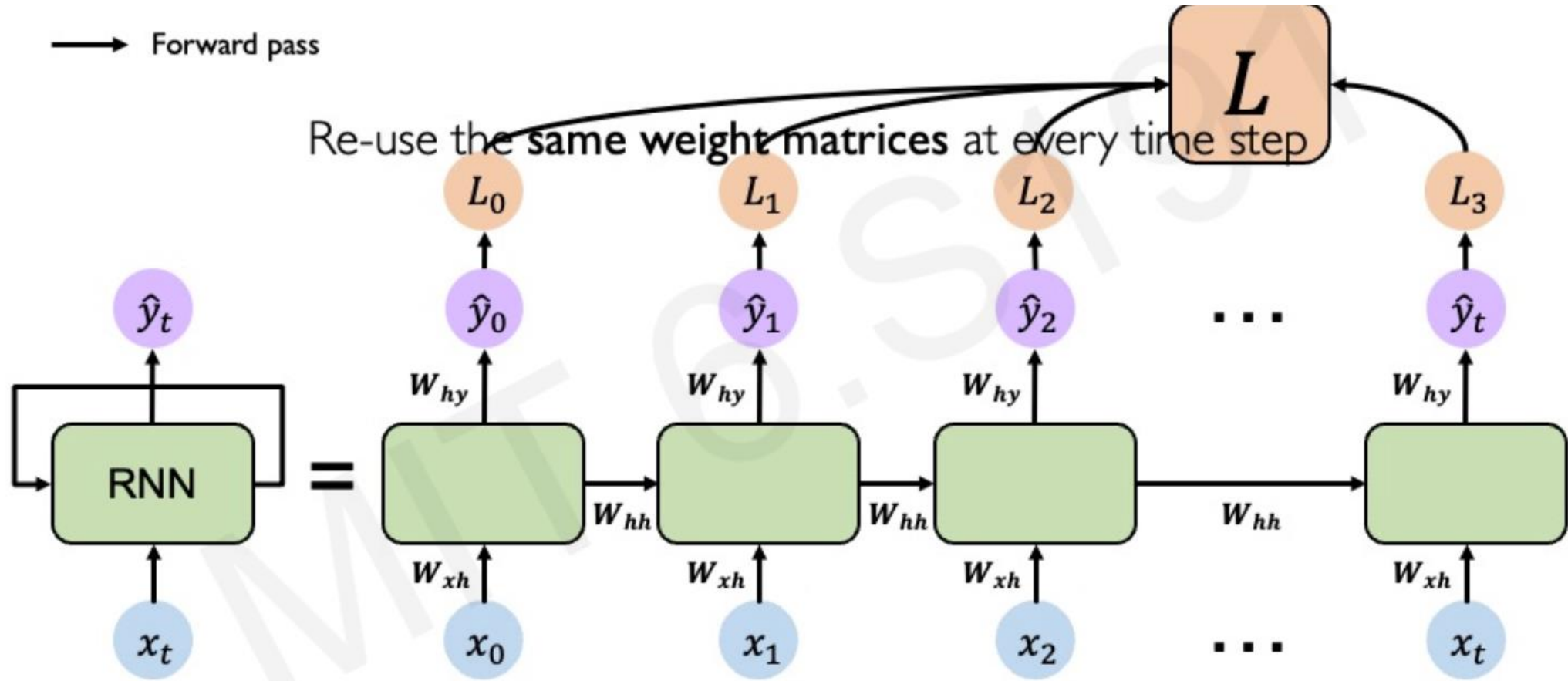
**Output Vector**
$$\hat{y}_t = W_{hy}^T h_t$$

**Update Hidden State**
$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$
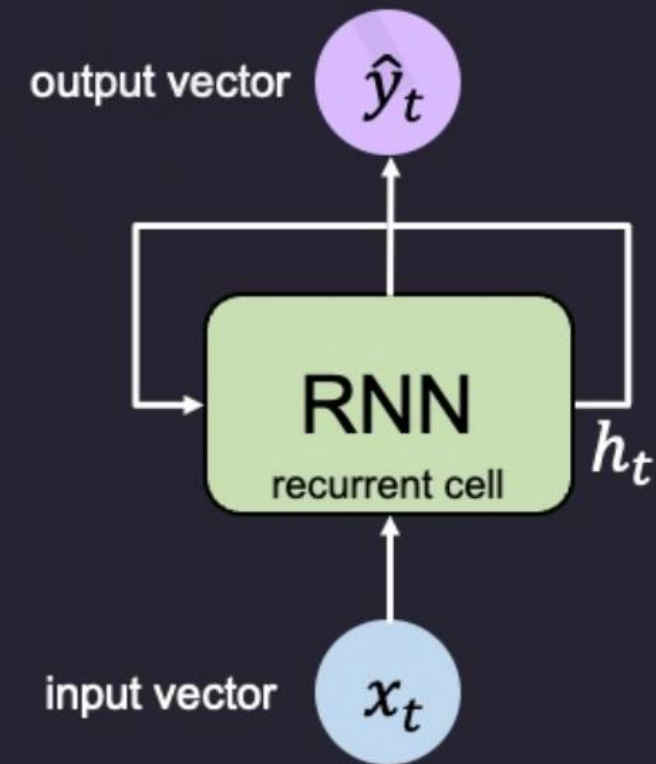
**Input Vector**
$$x_t$$
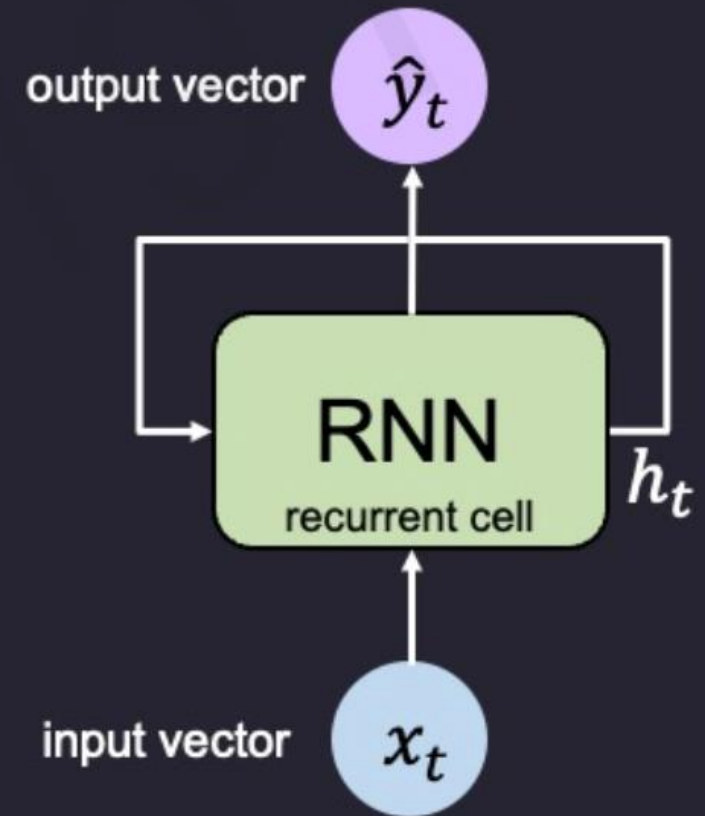
# Recurrent Neural Networks



17

# Recurrent Neural Networks

```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
```
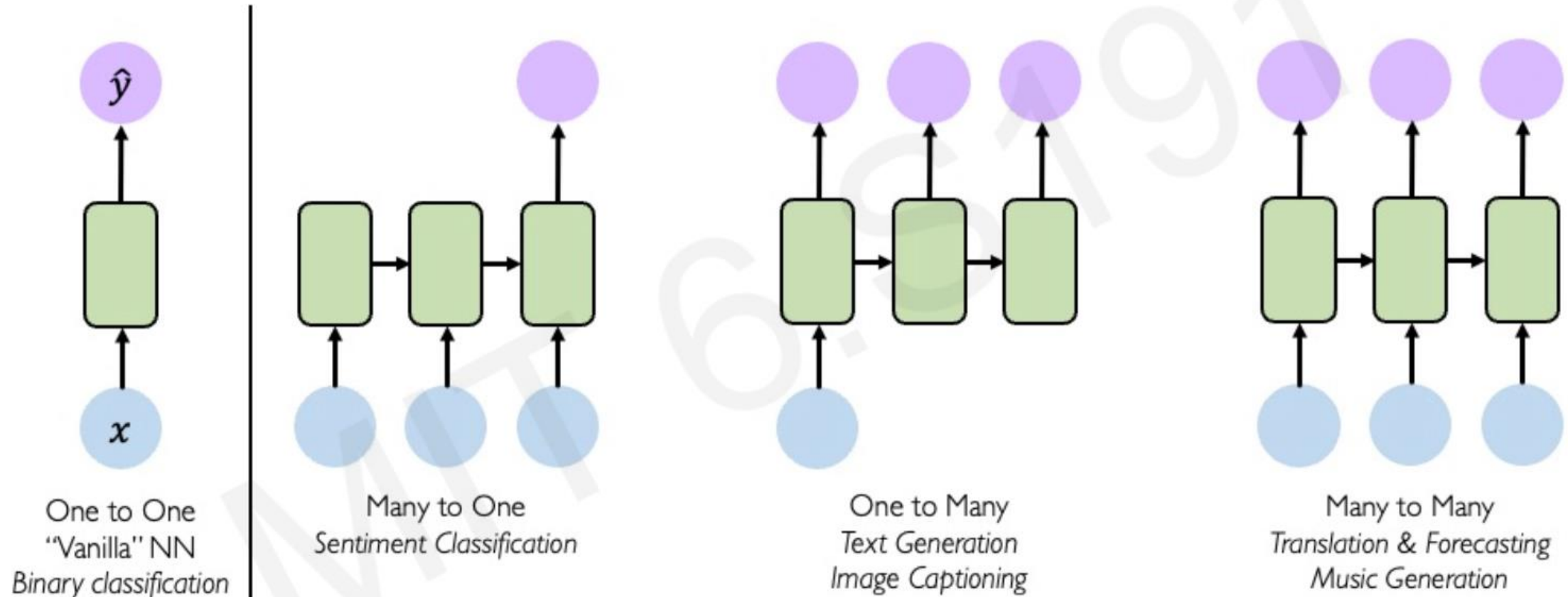
output vector $\hat{y}_t$

RNN
recurrent cell

$h_t$

input vector $x_t$

# Recurrent Neural Networks

# Recurrent Neural Networks



One to One
"Vanilla" NN
Binary classification

Many to One
Sentiment Classification

One to Many
Text Generation
Image Captioning

Many to Many
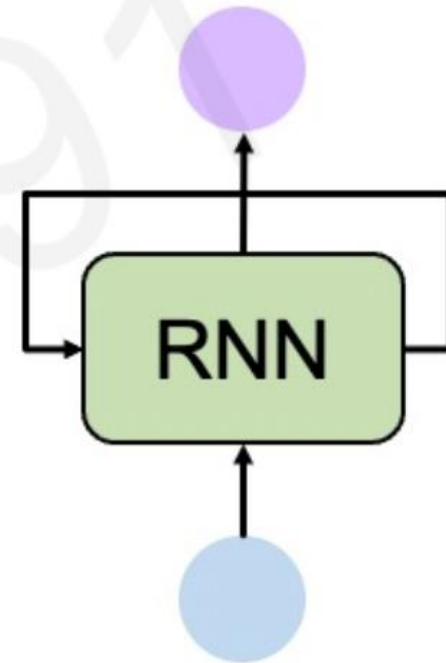Translation & Forecasting
Music Generation

# Design Factors

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)** meet
these sequence modeling design criteria

# Next Word Prediction – GPT Task
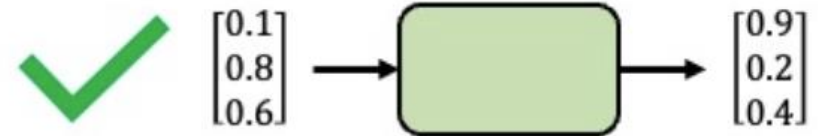
"This morning I took my cat for a walk."

given these words          predict the
                           next word

## Representing Language to a Neural Network



"deep" → ☐ → "learning"

Neural networks cannot interpret words

$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix}$ → ☐ → $\begin{bmatrix} 0.9 \\ 0.2 \\ 0.4 \end{bmatrix}$

Neural networks require numerical inputs
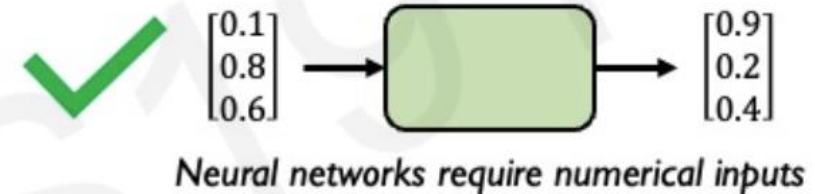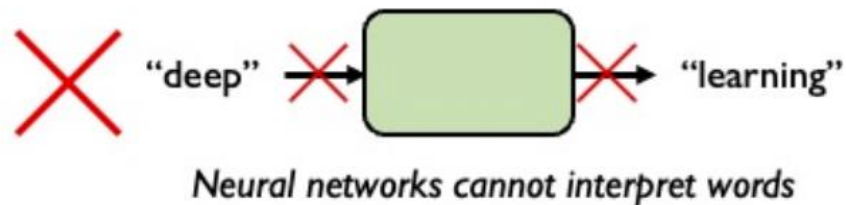
# Next Word Prediction – GPT Task



*Neural networks cannot interpret words*

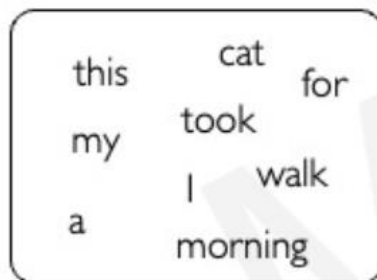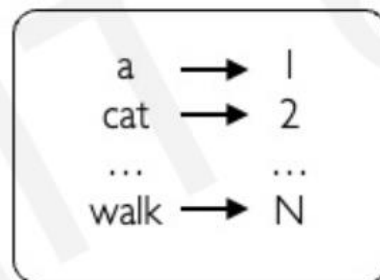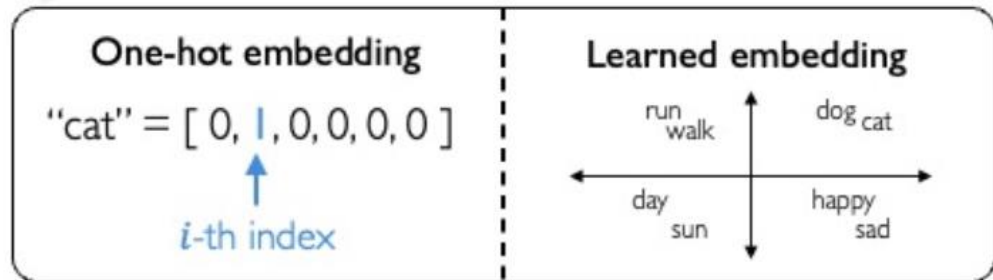*Neural networks require numerical inputs*

Embedding: transform indexes into a vector of fixed size.

**1. Vocabulary:**
Corpus of words

**2. Indexing:**
Word to index

**3. Embedding:**
Index to fixed-sized vector

One-hot embedding

"cat" = [ 0, 1, 0, 0, 0, 0 ]

$i$-th index

Learned embedding

# Variable Length

The food was great

vs.

We visited a restaurant for lunch

vs.

We were hungry but cleaned the house before eating

# Long-term Dependencies

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

*J'aime 6.S191!*

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

# Word Order



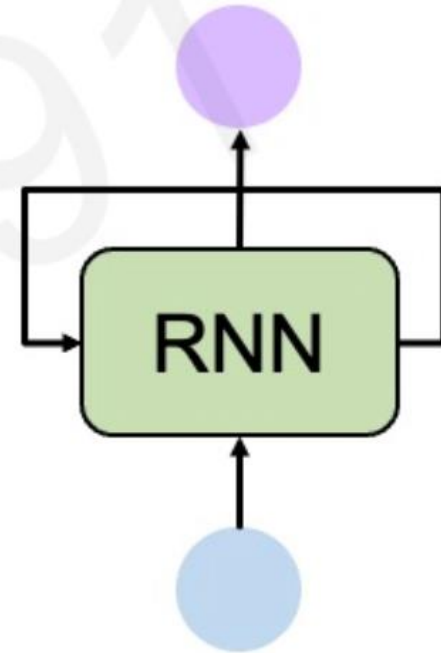The food was good, not bad at all.

vs.
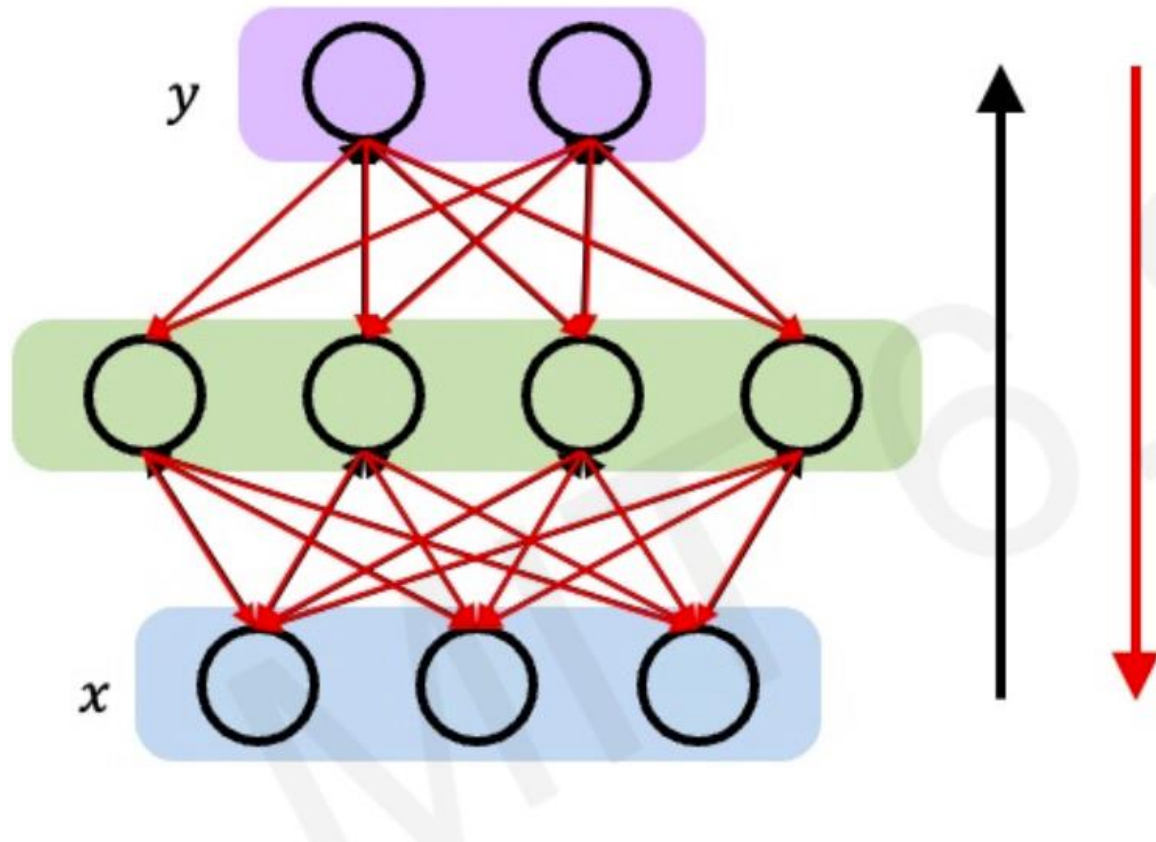
The food was bad, not good at all.

# RNN

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)** meet
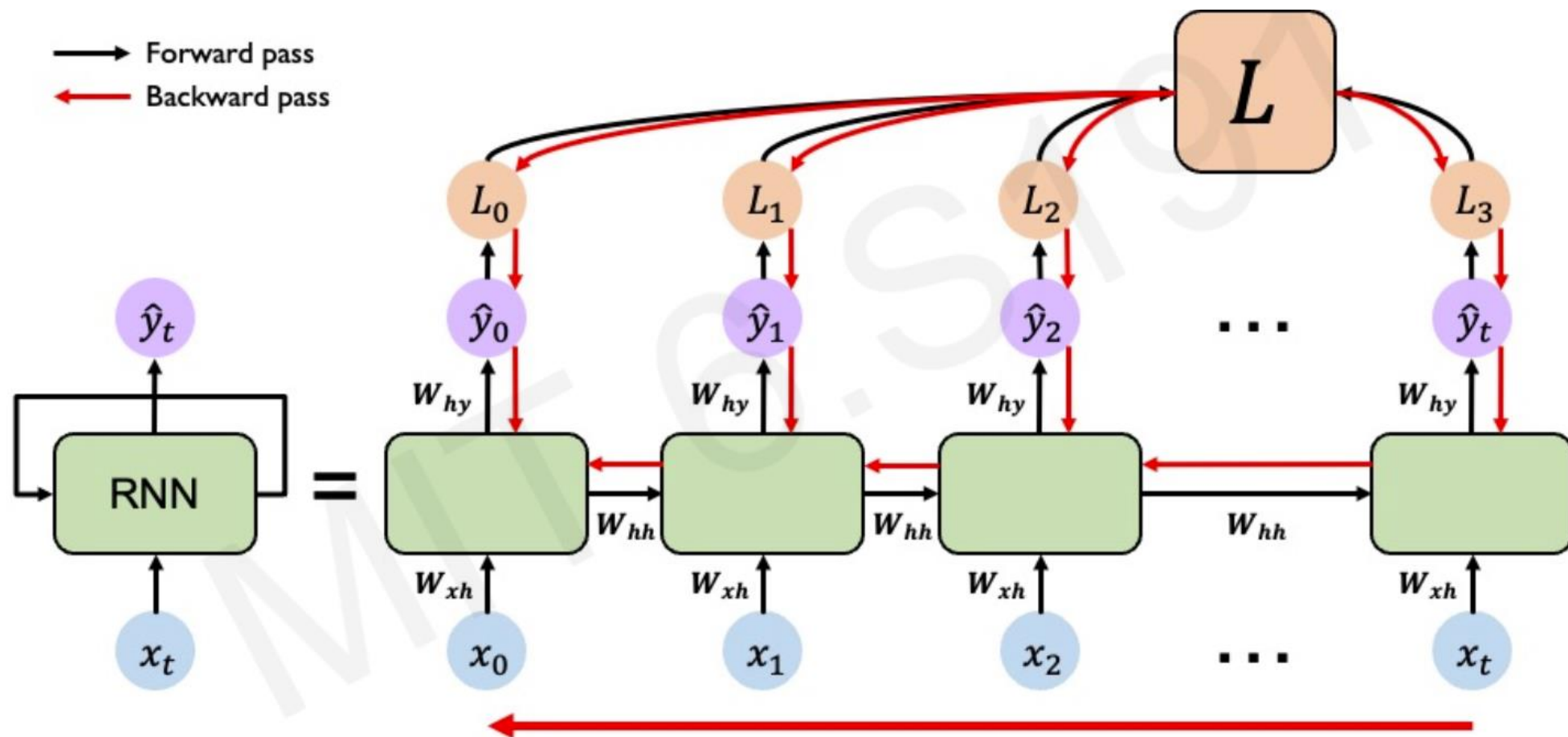these sequence modeling design criteria
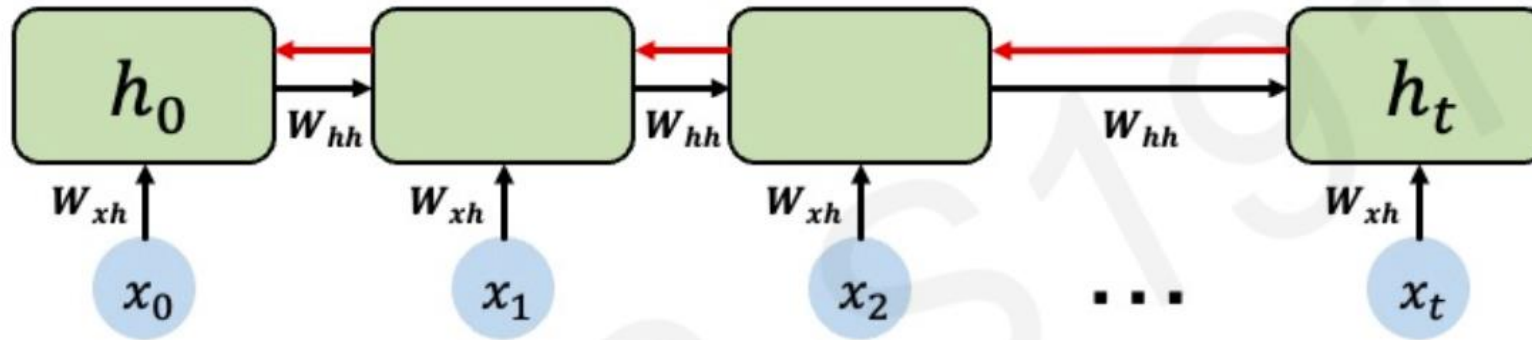
# Gradient Backprop.



**Backpropagation algorithm:**

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

2. Shift parameters in order to minimize loss
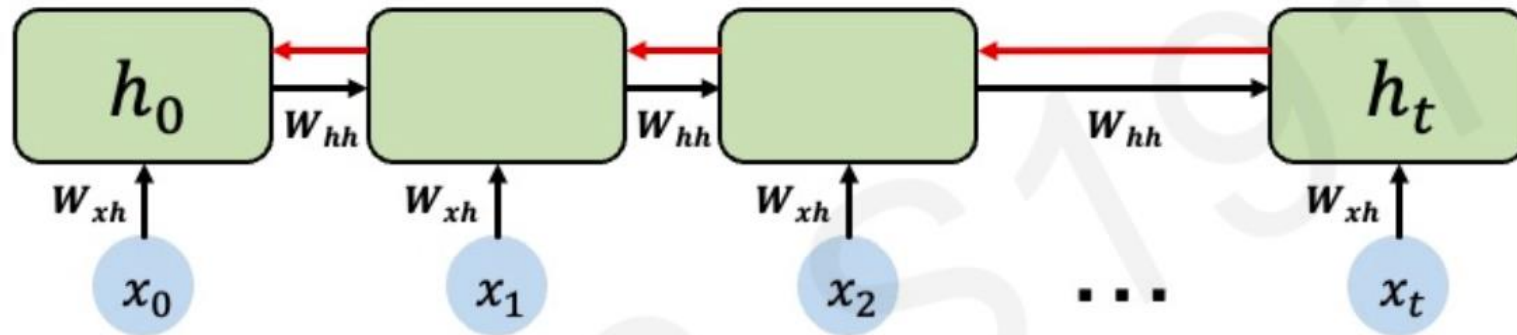
# Time Backprop.

# Exploding Gradient



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

# Vanishing Gradient



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

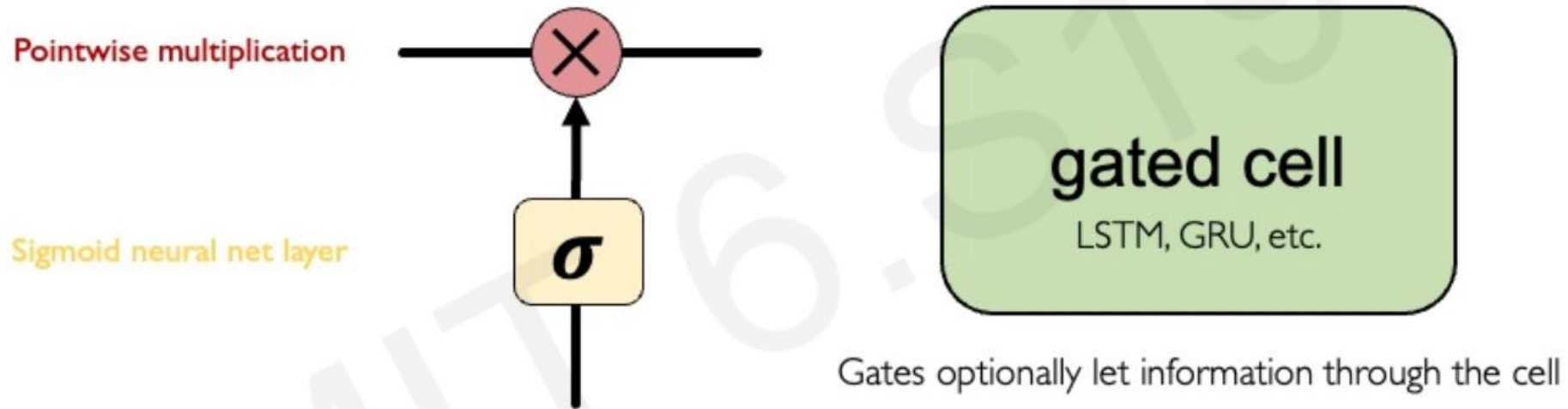Many values > 1:
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1:
**vanishing gradients**

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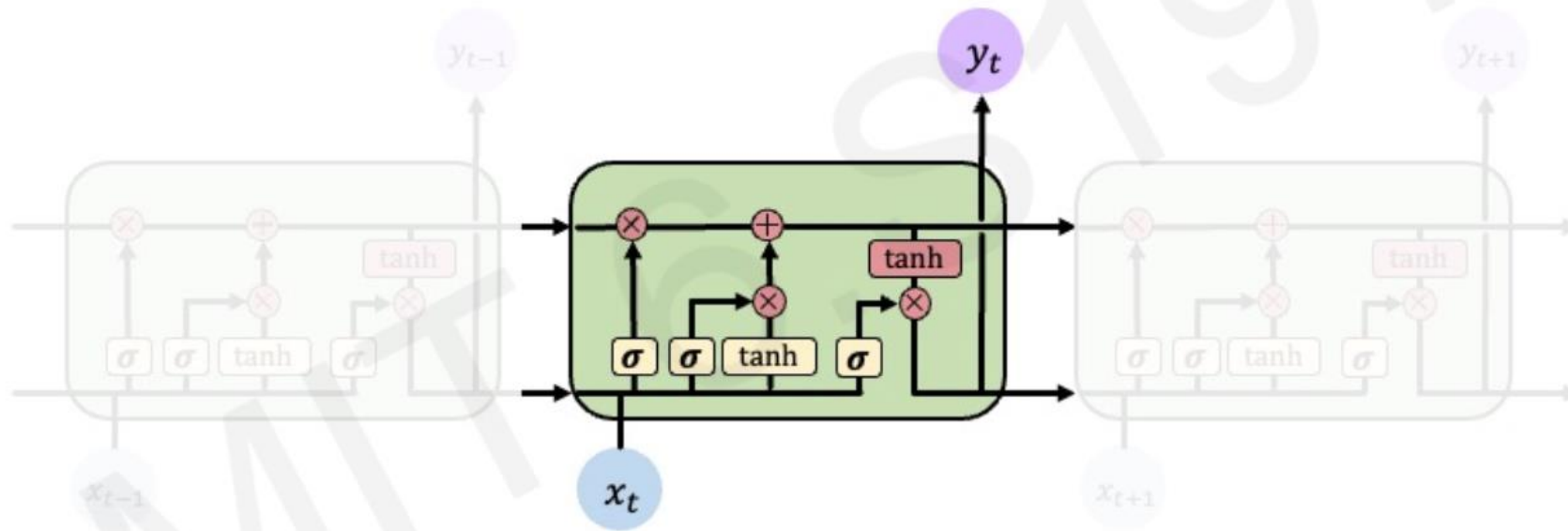2. Weight initialization
3. Network architecture

# LSTM

Idea: use **gates** to selectively **add** or **remove** information within **each recurrent unit with**

Pointwise multiplication

Sigmoid neural net layer

$\sigma$

gated cell

LSTM, GRU, etc.

Gates optionally let information through the cell

**Long Short Term Memory (LSTMs)** networks rely on a gated cell to track information throughout many time steps.

# Long-term Dependencies



Gated LSTM cells **control information flow:**
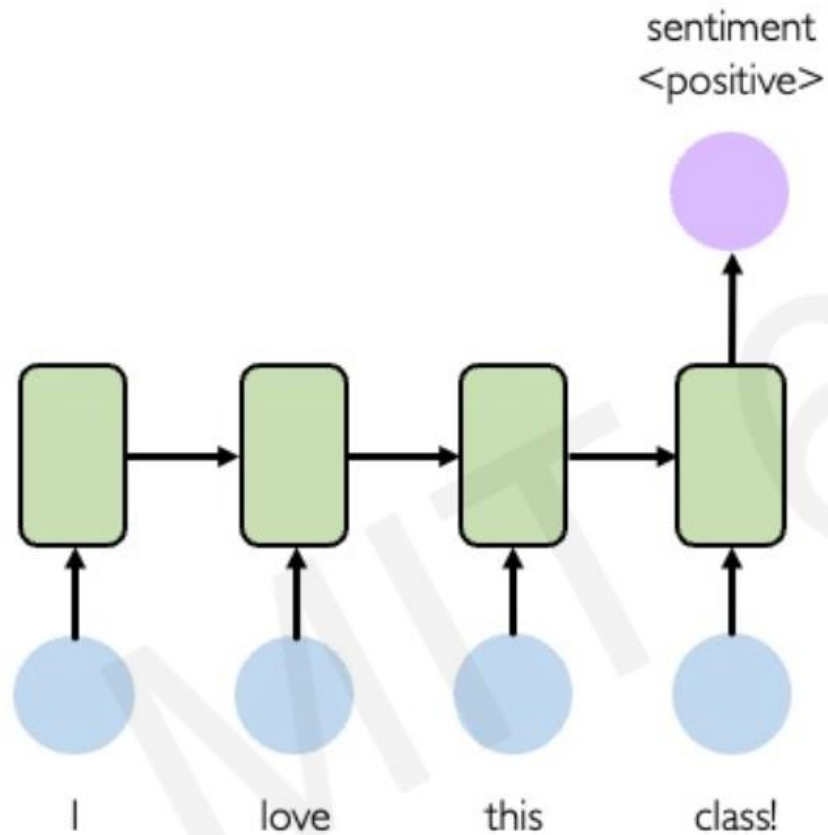1) Forget   2) Store   3) Update   4) Output

LSTM cells are able to track information throughout many timesteps

`tf.keras.layers.LSTM(num_units)`

# LSTM

1. Maintain a **cell state**

2. Use **gates** to control the **flow of information**
   - **Forget** gate gets rid of irrelevant information
   - **Store** relevant information from current input
   - Selectively **update** cell state
   - **Output** gate returns a filtered version of the cell state

3. Backpropagation through time with partially **uninterrupted gradient flow**
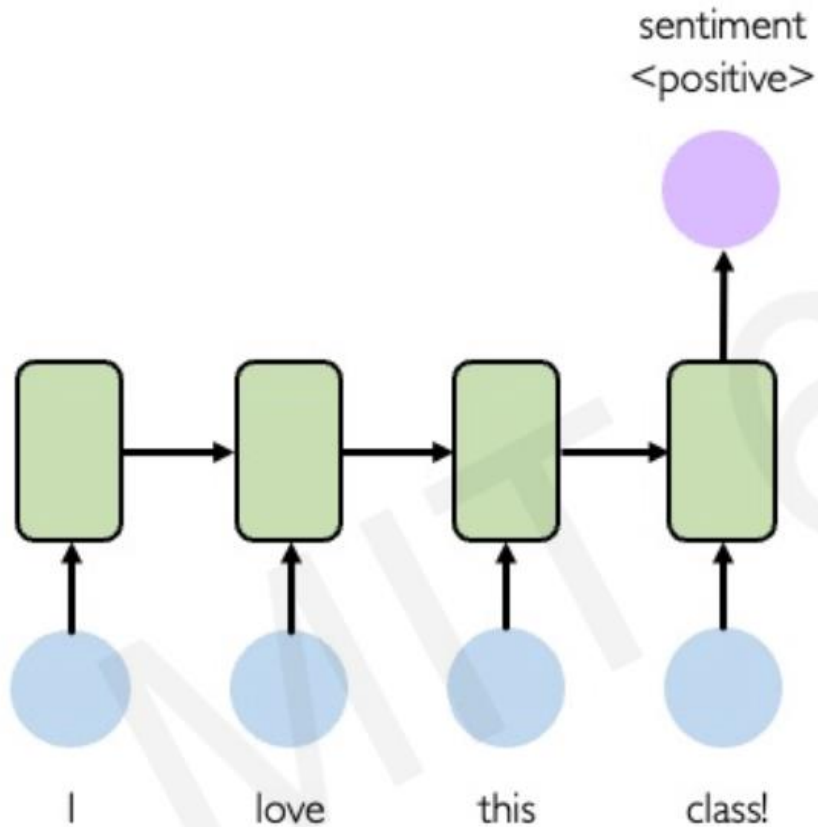
# Sentiment Analysis



sentiment
<positive>

**Input:** sequence of words

**Output:** probability of having positive sentiment

```
loss = tf.nn.softmax_cross_entropy_with_logits(y, predicted)
```

I love this class!

# Challenges with RNN Models



sentiment
<positive>

I    love    this    class!

**Limitations of RNNs**

Encoding bottleneck

Slow, no parallelization

Not long memory

# Attention is all you need!