

**INF721**

2024/2

**UFV**

# Deep Learning

## L13: Recurrent Neural Networks

# Logistics

## Announcements

- ▶ PA3 is due this Wednesday, 11:59pm

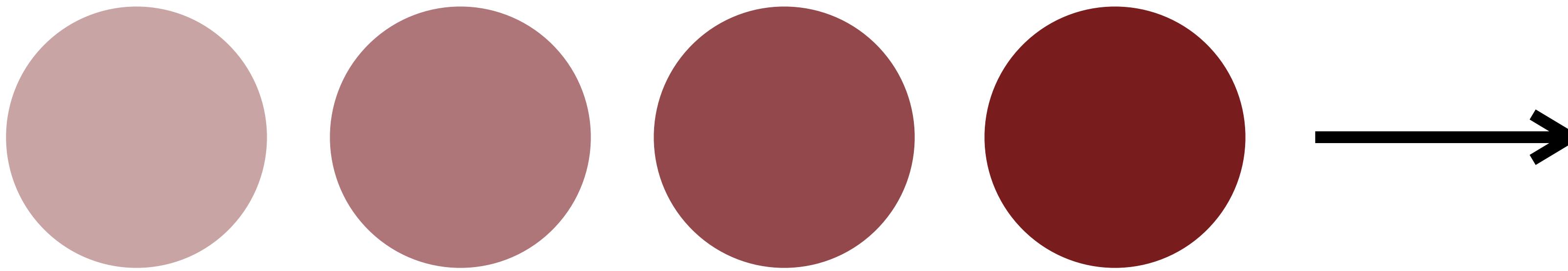
## Last Lecture

- ▶ Input normalization
- ▶ Batch normalization
- ▶ Layer normalization

# Lecture Outline

- ▶ Sequential Problems
- ▶ Recurrent Neural Networks(RNNs)
- ▶ Type of RNNs
- ▶ Backpropagation Through Time
- ▶ Language Models
- ▶ Exploding/Vanishing Gradients

# What is the next position of this ball?



Recurrent Neural Networks are used for classification, regression or generation of sequential data!

# What letter comes after T in the alphabet?

R S T U

Recurrent Neural Networks are used for classification, regression or generation of sequential data!

# Sequential Problems in Artificial Intelligence

	Input	Output
<b>Speech Recognition</b>		"Alexa, play The Beatles on Spotify"
<b>Sentiment Analysis</b>	"This is a terrible product."	
<b>Machine Translation</b>	"The book is on the table."	"O livro está em cima da mesa."
<b>Image Captioning</b>		"A cat lying by the window."
<b>Music Generation</b>	None	
<b>Named Entity Recognition</b>	"Lucas Ferreira is a professor at UFV"	"Lucas Ferreira is a professor at <b>UFV</b> "

# Example: Named Entity Recognition

Locate and classify named entities mentioned in unstructured text:

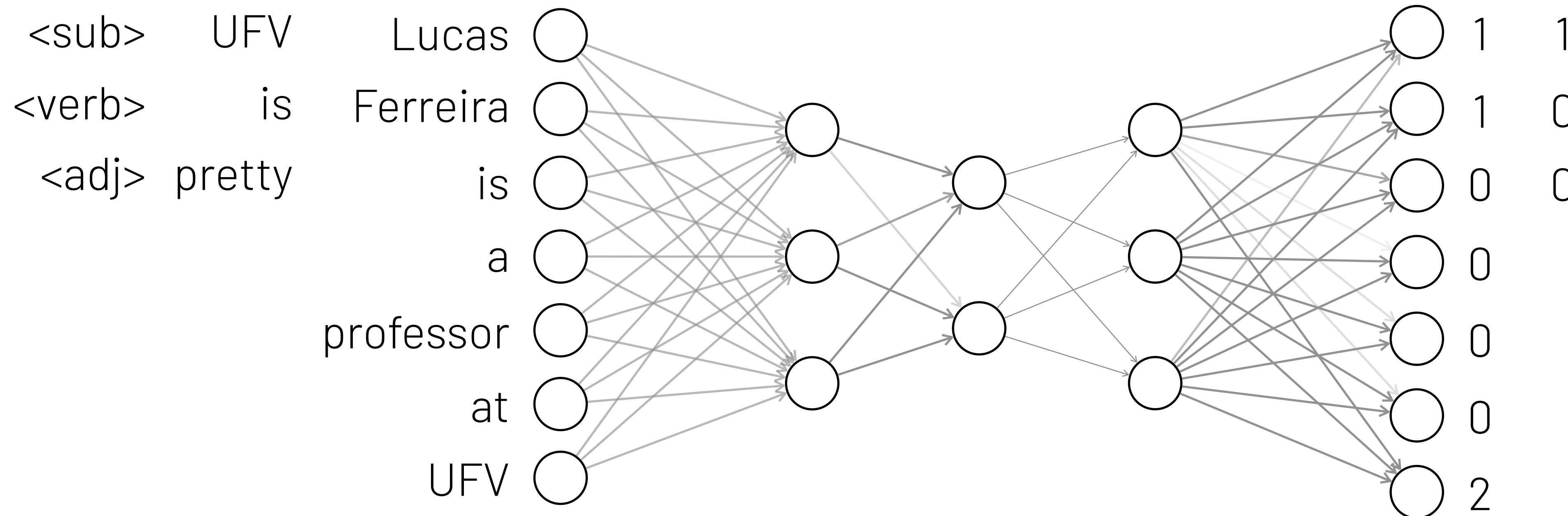
$y^{<1>}$	$y^{<2>}$	$y^{<3>}$	$y^{<4>}$	$y^{<5>}$	$y^{<6>}$	$y^{<7>}$
1	1	0	0	0	0	2

$X$  Lucas Ferreira is a professor at UFV

$x^{<1>}$   $x^{<2>}$   $x^{<3>}$   $x^{<4>}$   $x^{<5>}$   $x^{<6>}$   $x^{<7>}$

In sequential problems, each input  $x^{<t>}$  can have an associated output  $y^{<t>}$

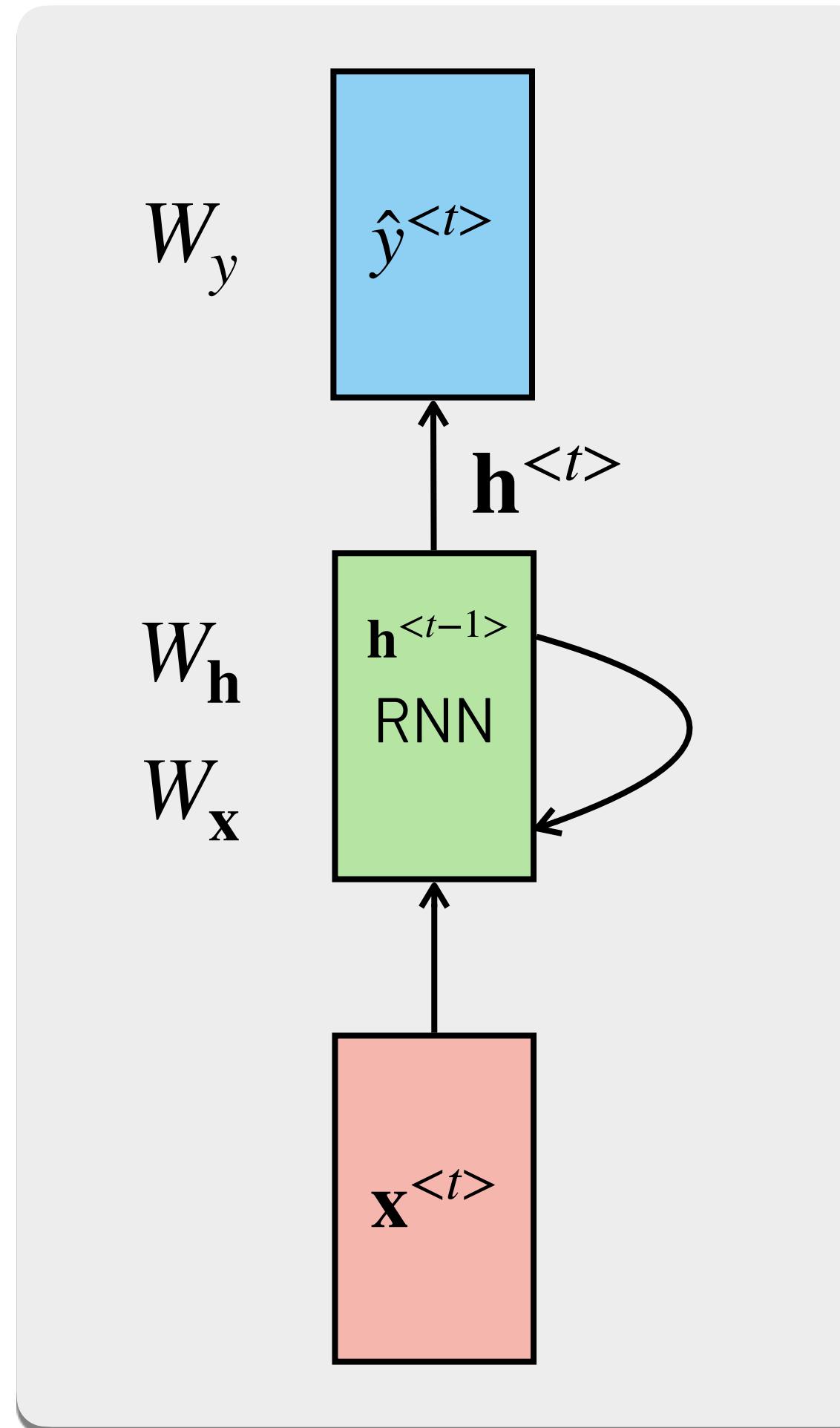
# Why not MLPs for sequential problems?



**Problem 1:** Inputs and outputs may have different sizes in different examples.

**Problem 2:** MLPs do not capture temporal dependencies between elements of a sequence.

# Recurrent Neural Networks (RNNs)



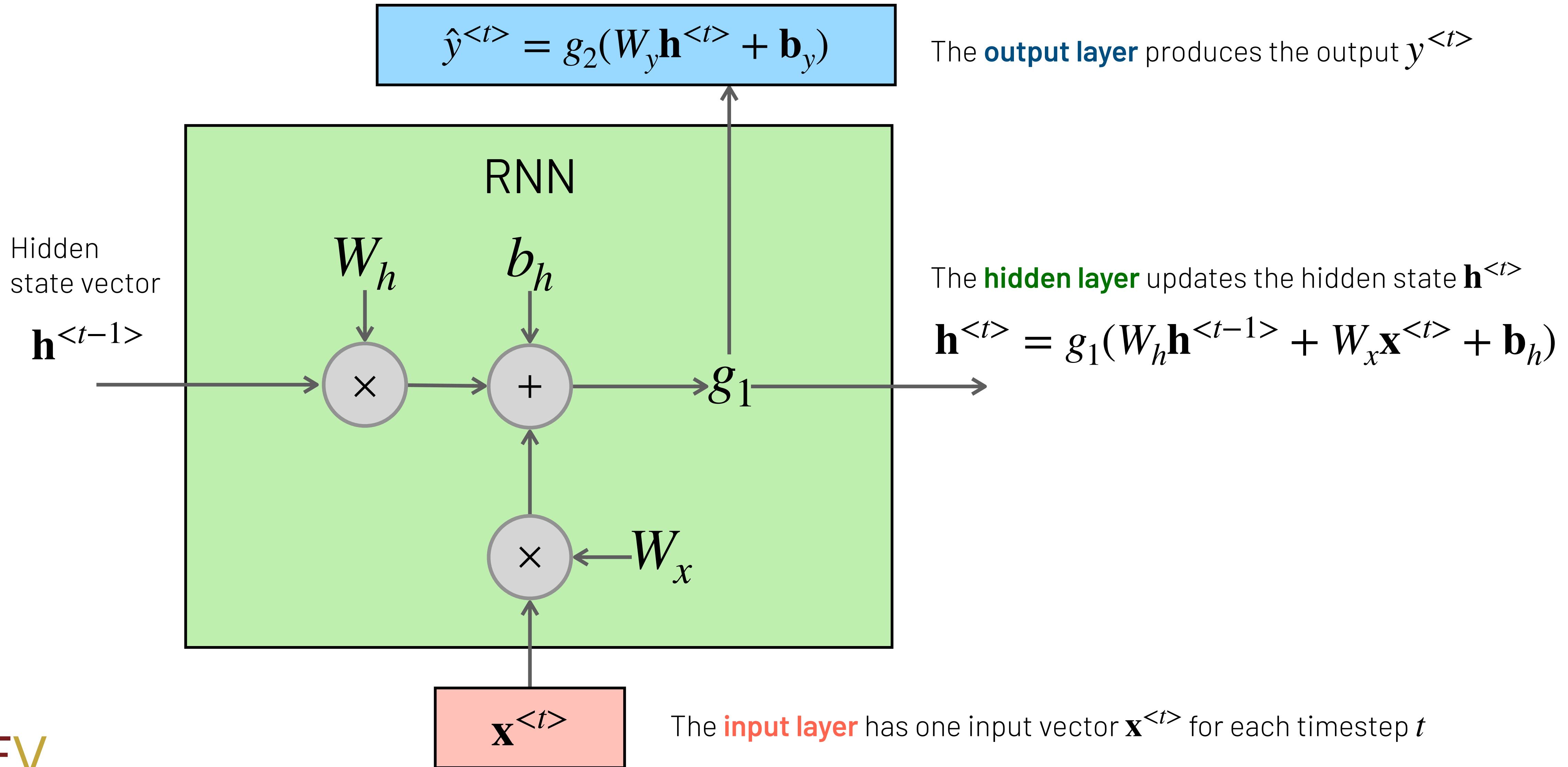
RNNs process each input element  $\mathbf{x}^{<t>}$  at a time, keeping a state (vector)  $\mathbf{h}^{<t>}$  that is updated at each time step  $t$  to produce the output  $\hat{y}^{<t>}$

$$\mathbf{h}^{<t>} = g_1(W_h \mathbf{h}^{<t-1>} + W_x \mathbf{x}^{<t>} + \mathbf{b}_h)$$

$$\hat{y}^{<t>} = g_2(W_y \mathbf{h}^{<t>} + \mathbf{b}_y)$$

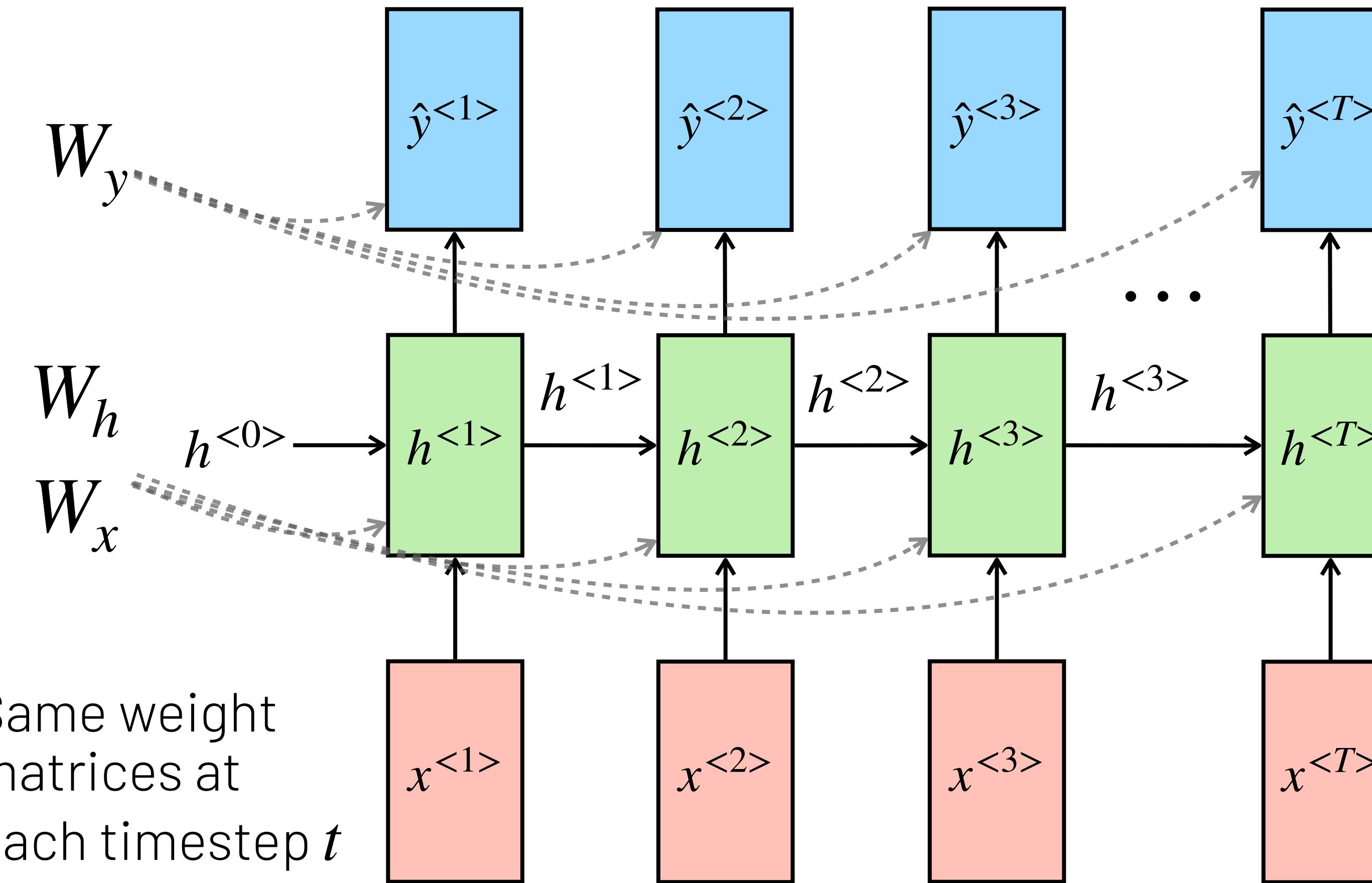
- ▶  $g_1$ : hidden layer activation function (tanh/relu)
- ▶  $g_2$ : output layer activation function (sigmoid/softmax)

# Recurrent Neural Networks (RNNs)



# Recurrent Neural Networks (RNNs)

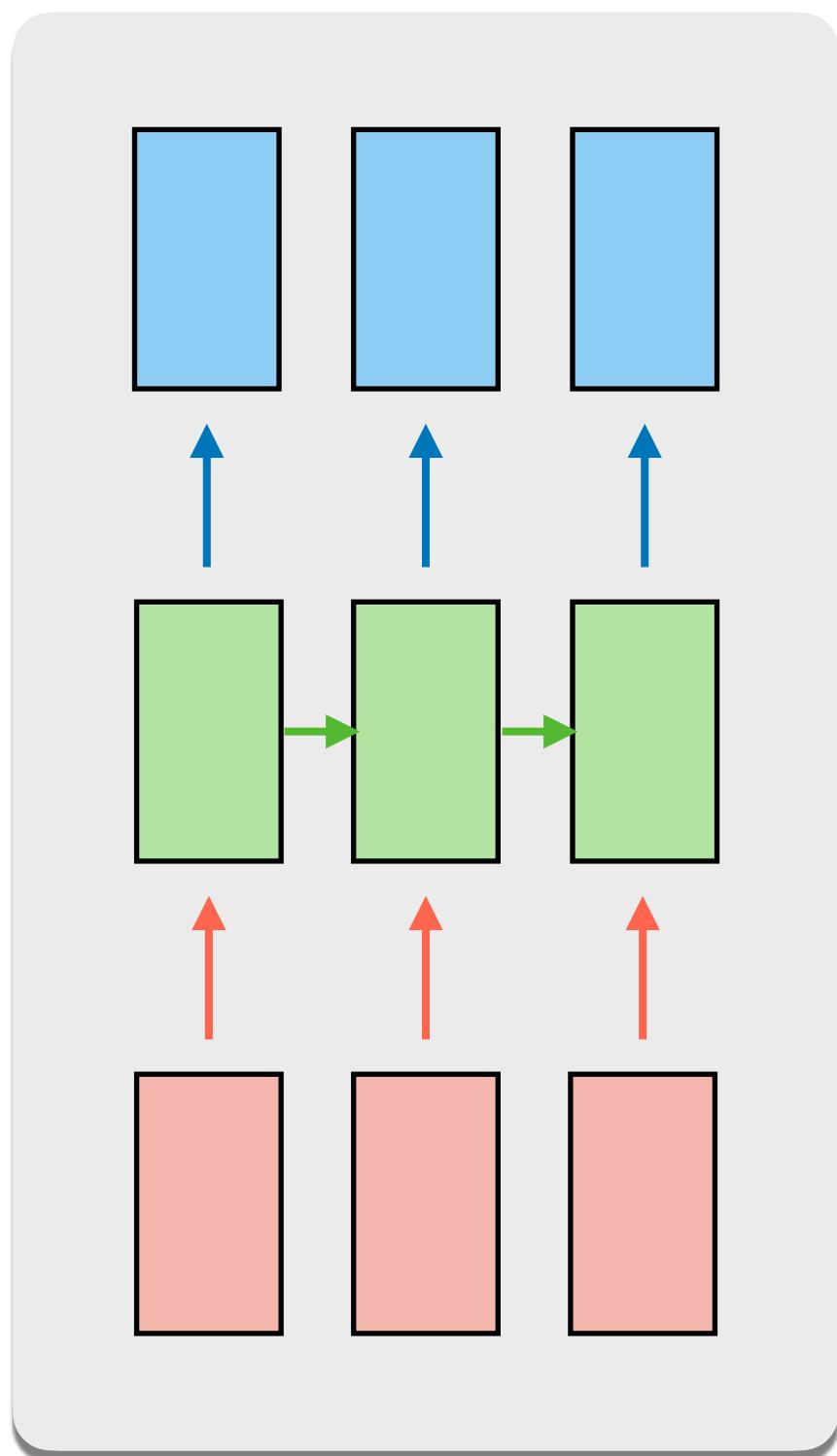
RNNs can be seen unrolled over a fixed number of timesteps  $T$



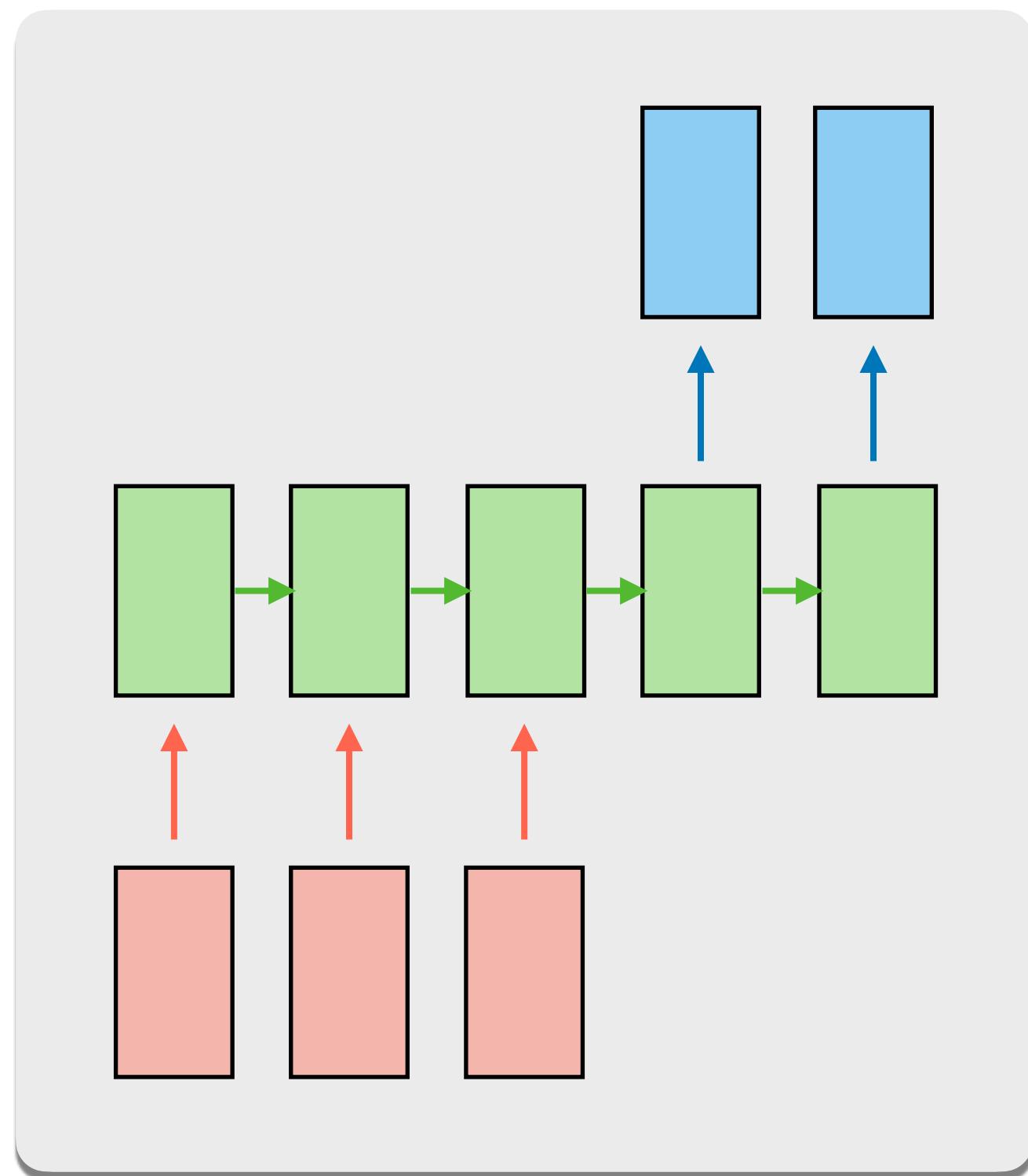
$$\begin{aligned} h^{<1>} &= g_1(W_h h^{<0>} + W_x x^{<1>} + b_h) \\ \hat{y}^{<1>} &= g_2(W_y h^{<1>} + b_y) \\ h^{<2>} &= g_1(W_h h^{<1>} + W_x x^{<2>} + b_h) \\ \hat{y}^{<2>} &= g_2(W_y h^{<2>} + b_y) \\ h^{<3>} &= g_1(W_h h^{<2>} + W_x x^{<3>} + b_h) \\ \hat{y}^{<3>} &= g_2(W_y h^{<3>} + b_y) \\ &\vdots \\ h^{<T>} &= g_1(W_h h^{<T-1>} + W_x x^{<T>} + b_h) \\ \hat{y}^{<T>} &= g_2(W_y h^{<T>} + b_y) \end{aligned}$$

# Types of RNNs

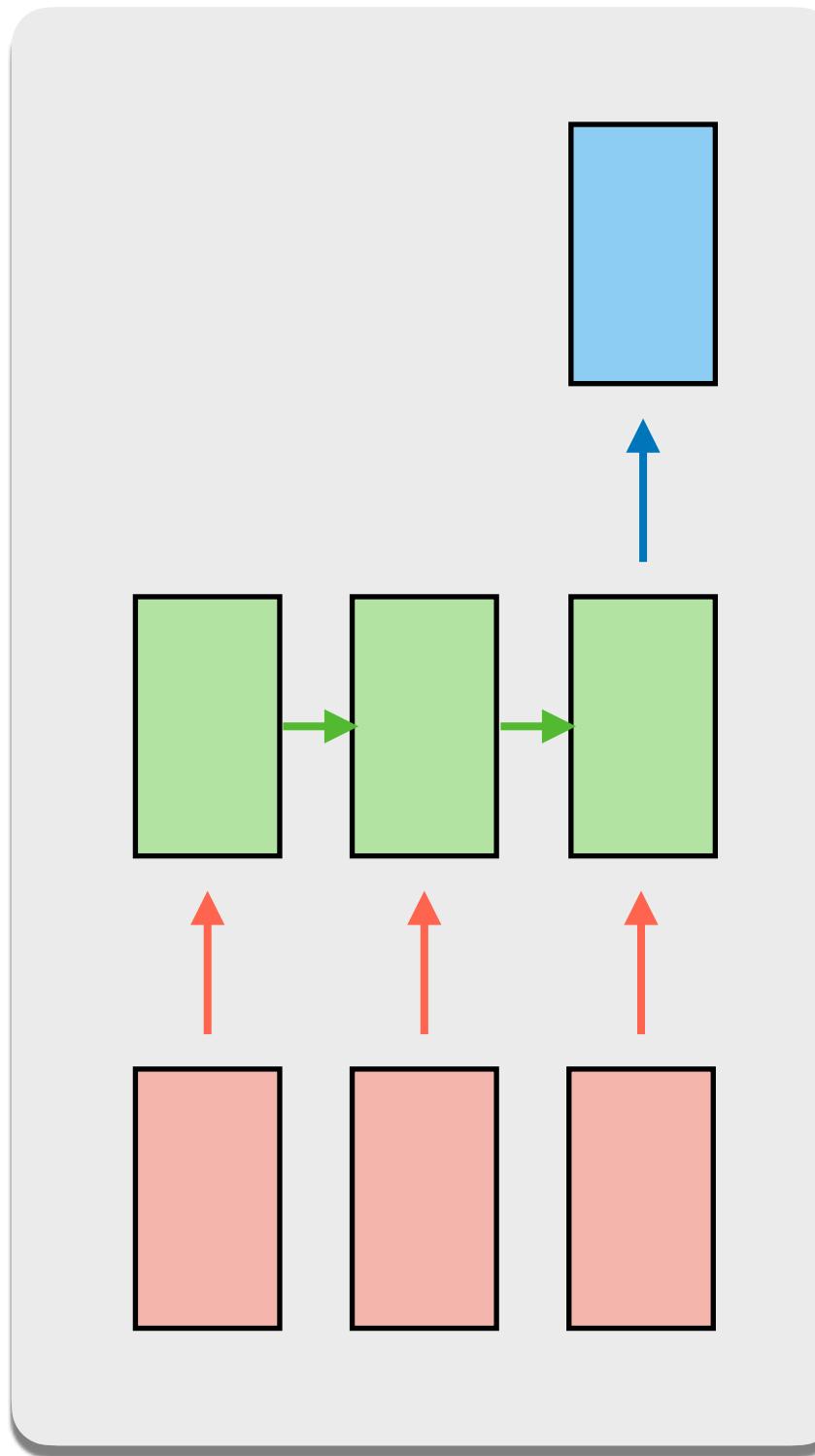
Many to Many



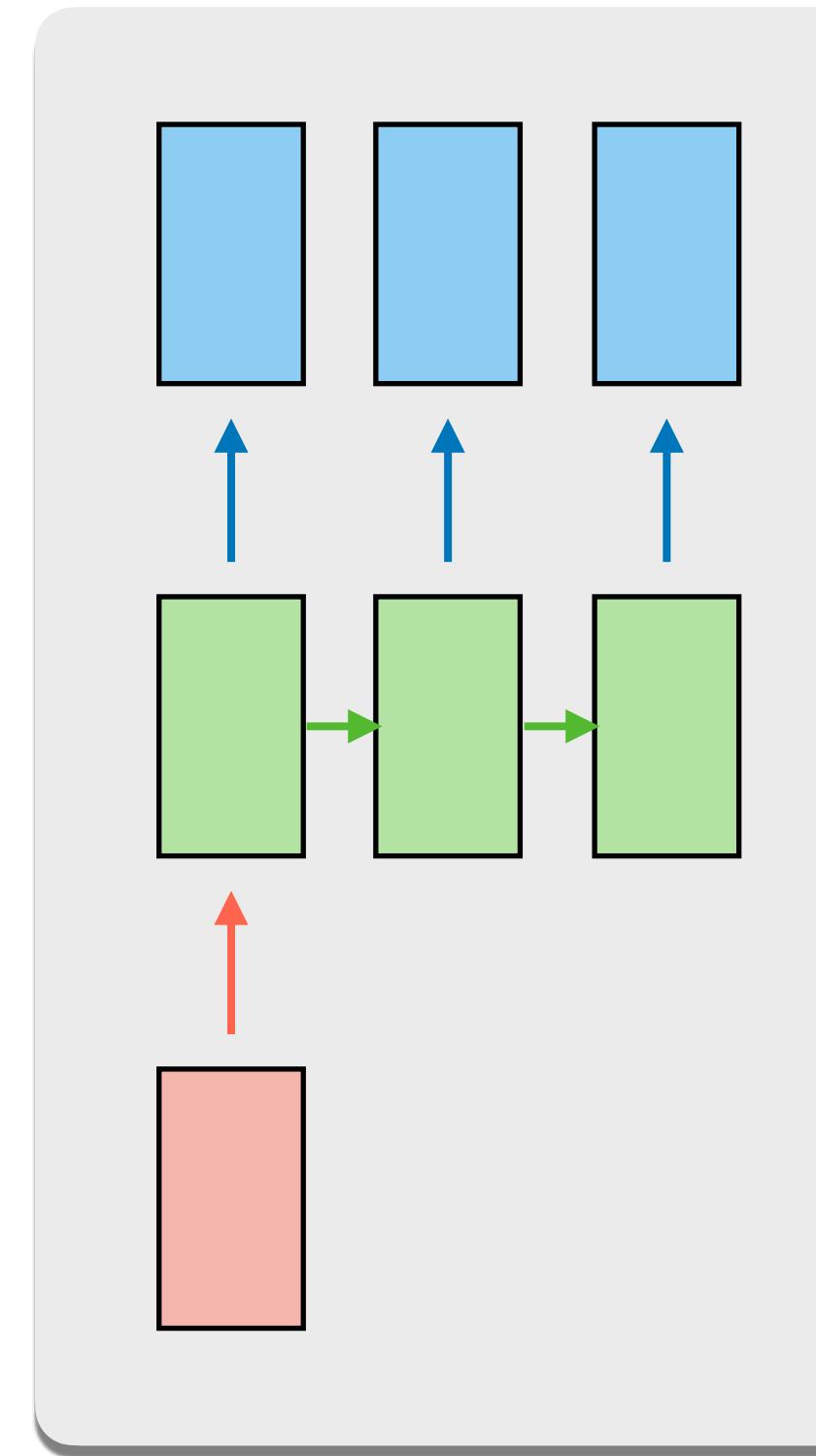
Many to Many (Seq2Seq)



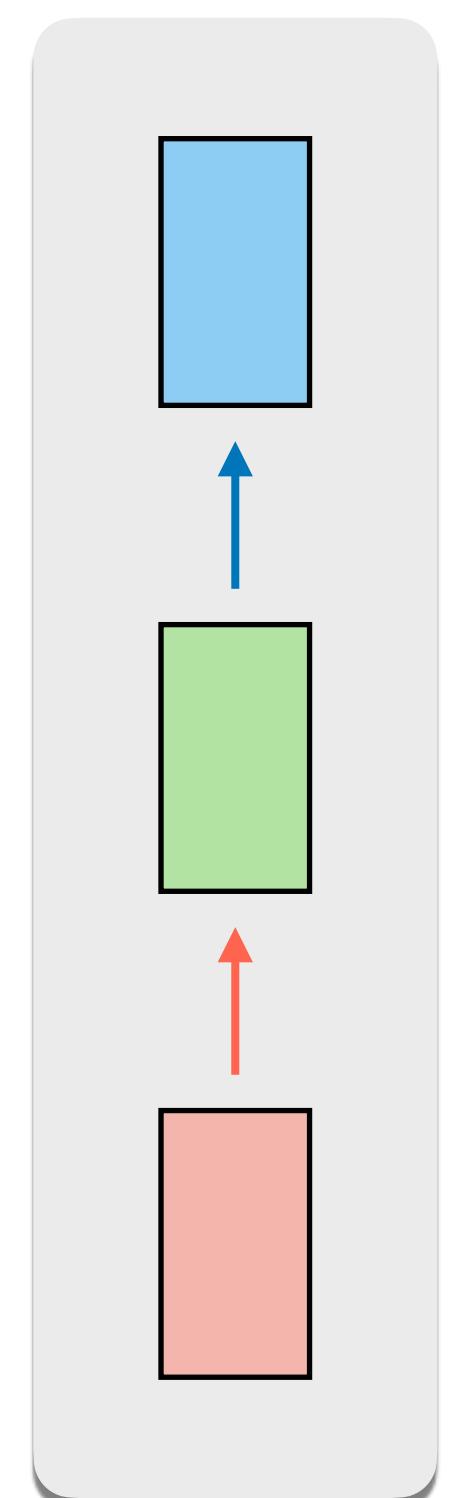
Many to one



One to many



one to one



**Example**

Named Entity Recognition

**Example**

Machine Translation

**Example**

Sentiment Analysis

**Example**

Image Description

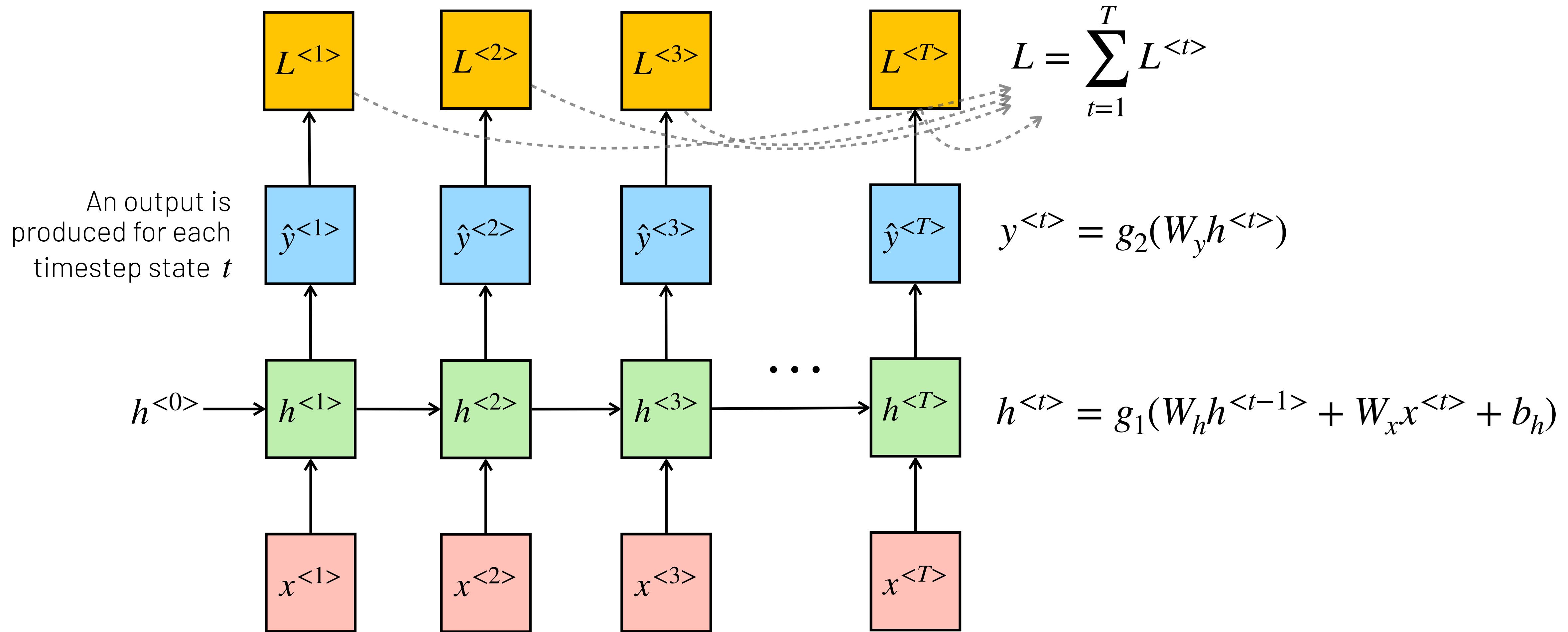
MLP

# Many to Many

Named Entity Recognition

"Lucas Ferreiis a professor at UFV"

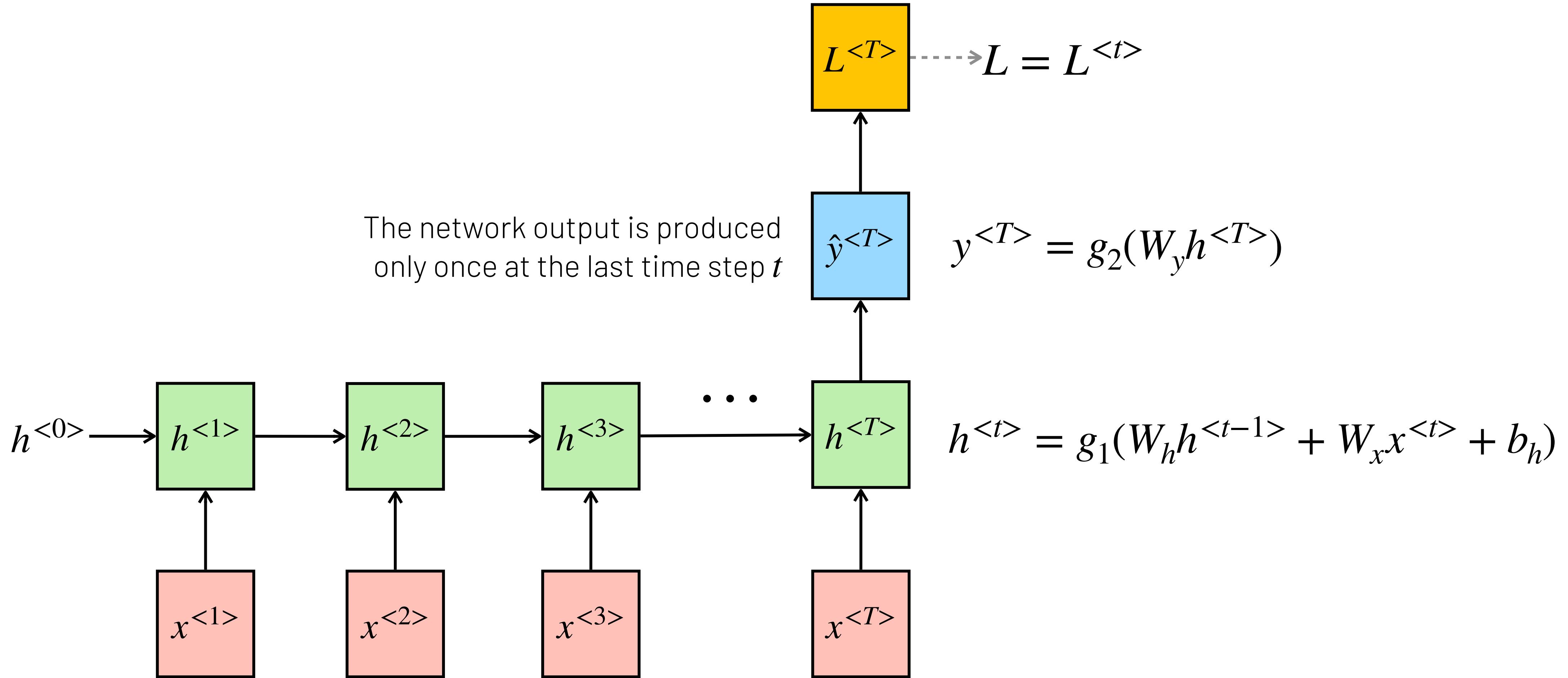
1100002



# Many to One

## Sentiment Analysis

"This is a terrible product."

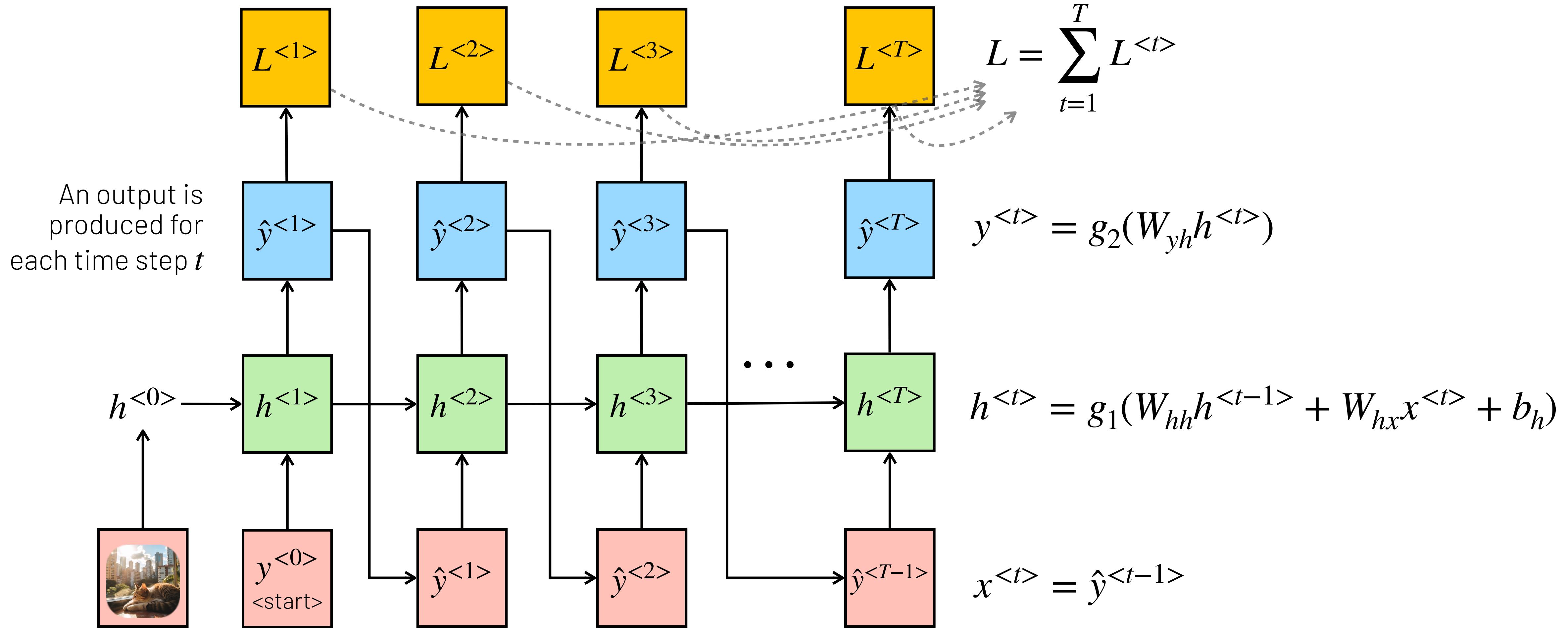


# One to Many

Image Captioning



"A cat lying by the window."

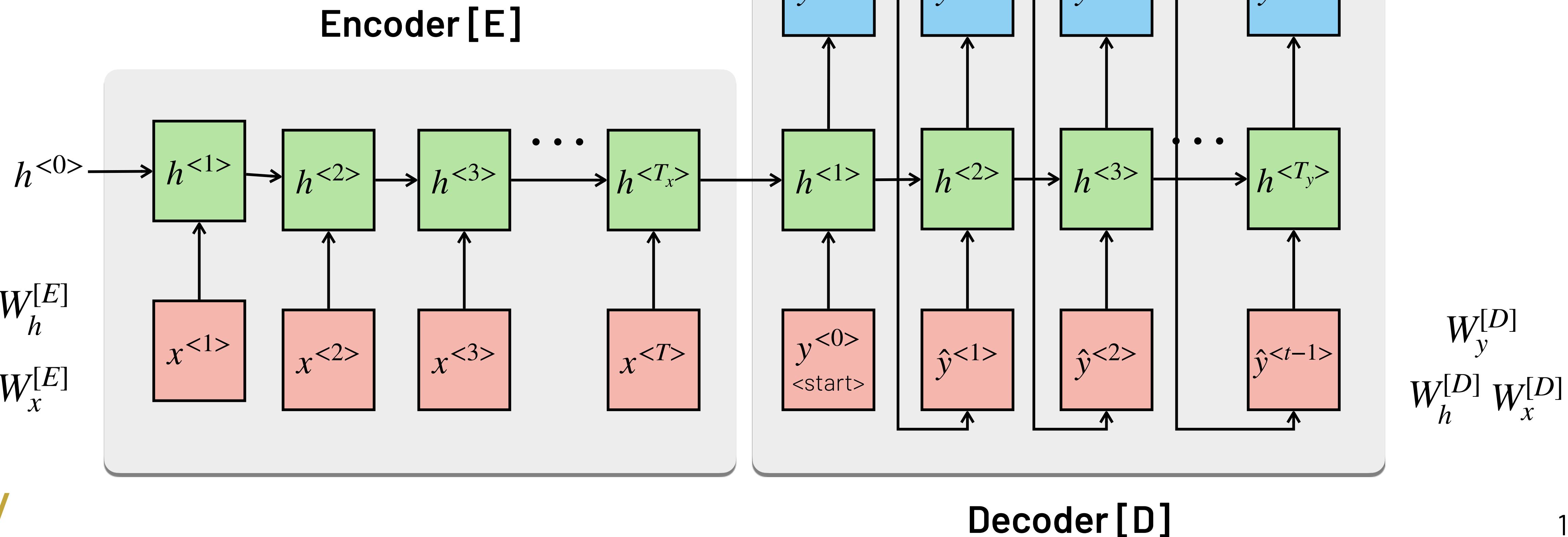


# Seq2Seq

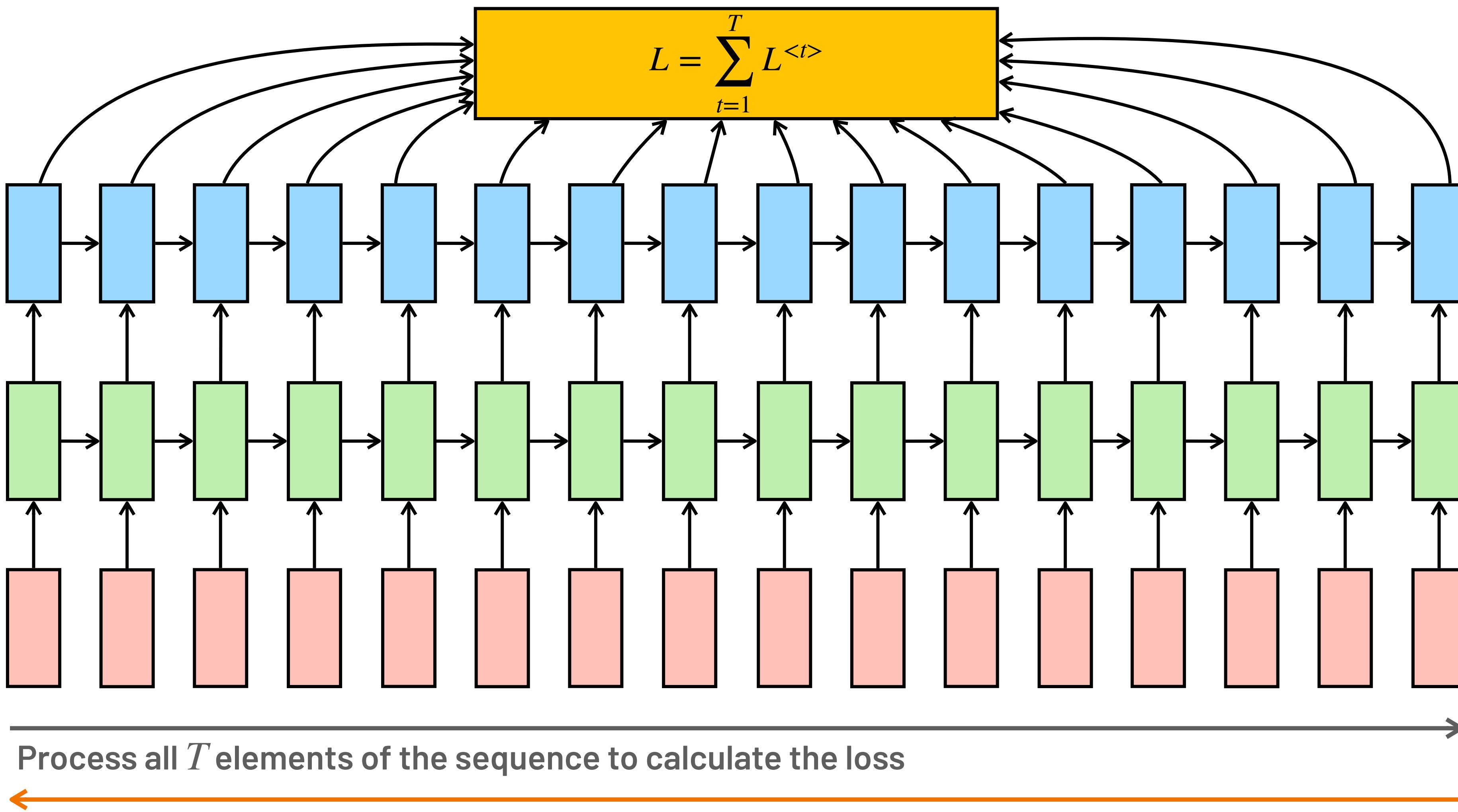
**Machine Translation** “The book is on the table.” “O livro está em cima da mesa.”

The input  $x$  is processed with an **encoder** network and its final hidden state  $h^{<T_x>}$  is used to initialize the hidden state of another **decoder** network, which produces an output for each time step  $t_y$ .

$$L = \sum_{t=1}^{T_y} L^{<t>}$$

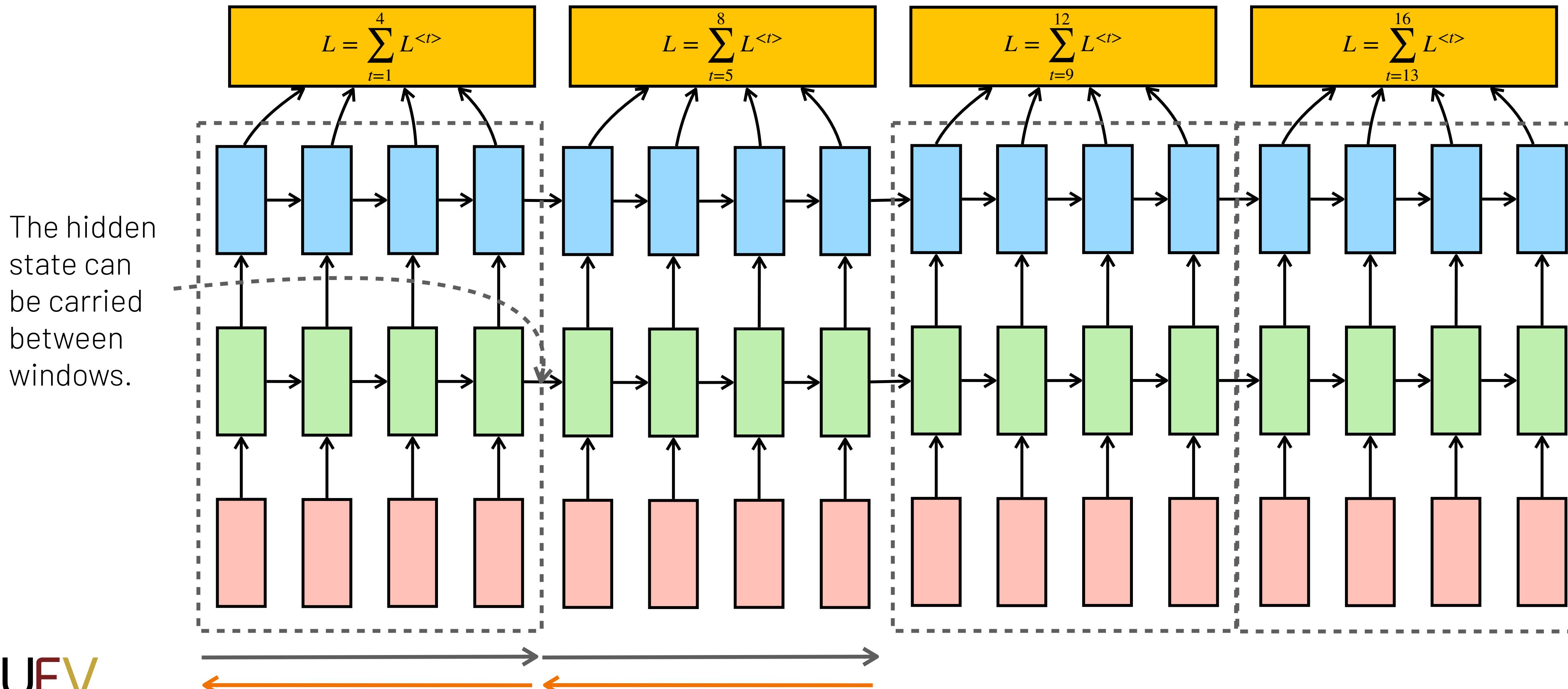


# Backpropagation Through Time



# Truncated Backpropagation Through Time

If the size of the sequence to be processed is very large or infinite (e.g., time series), perform propagation and backpropagation in windows of size  $j$  (e.g., 4)



# Next Lecture

**L14:** Recurrent Neural Networks (Part II)

GRUs and LSTMs for processing with very long sequences.