# INF721

2024/2

# Deep Learning

## L13: Recurrent Neural Networks

UFV

# Logistics

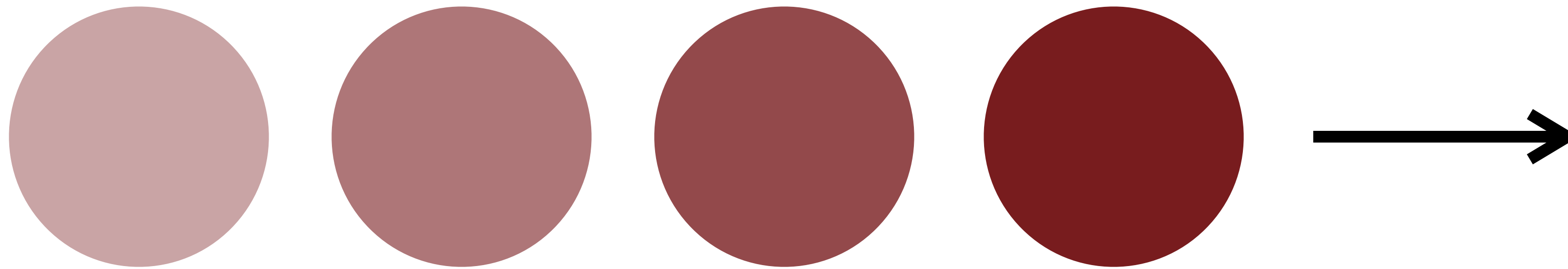**Announcements**

‣ PA3 is due this Wednesday, 11:59pm

**Last Lecture**

‣ Input normalization

‣ Batch normalization

‣ Layer normalization

UFV

# Lecture Outline

▶ Sequential Problems

▶ Recurrent Neural Networks (RNNs)

▶ Type of RNNs

▶ Backpropagation Through Time

▶ Language Models

▶ Exploding/Vanishing Gradients

UFV

# 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!

UFV

# Sequential Problems in Artificial Intelligence

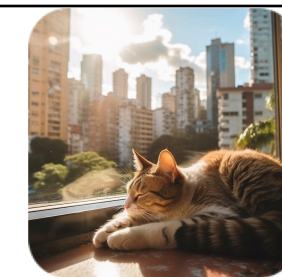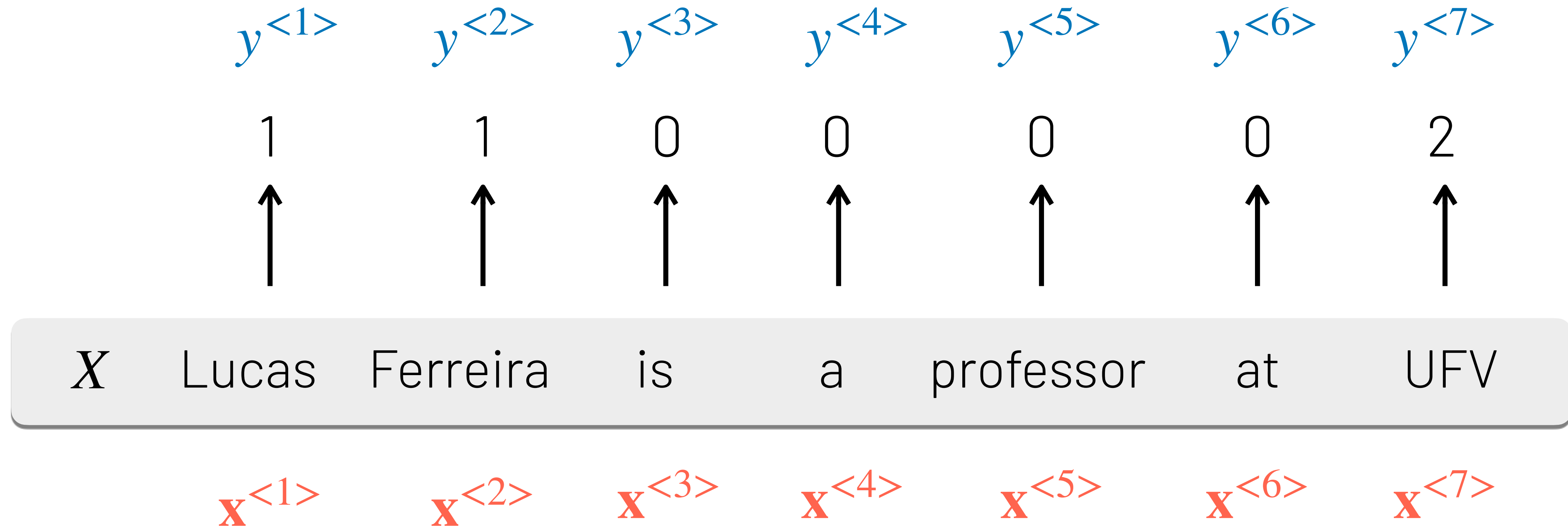| | Input | Output |
|---|---|---|
| **Speech Recognition** |  | "Alexa, play The Beatles on Spotify" |
| **Sentiment Analysis** | "This is a terrible product." | ★☆☆☆☆ |
| **Machine Translation** | "The book is on the table." | "O livro está em cima da mesa." |
| **Image Captioning** |  | "A cat lying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

# 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

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

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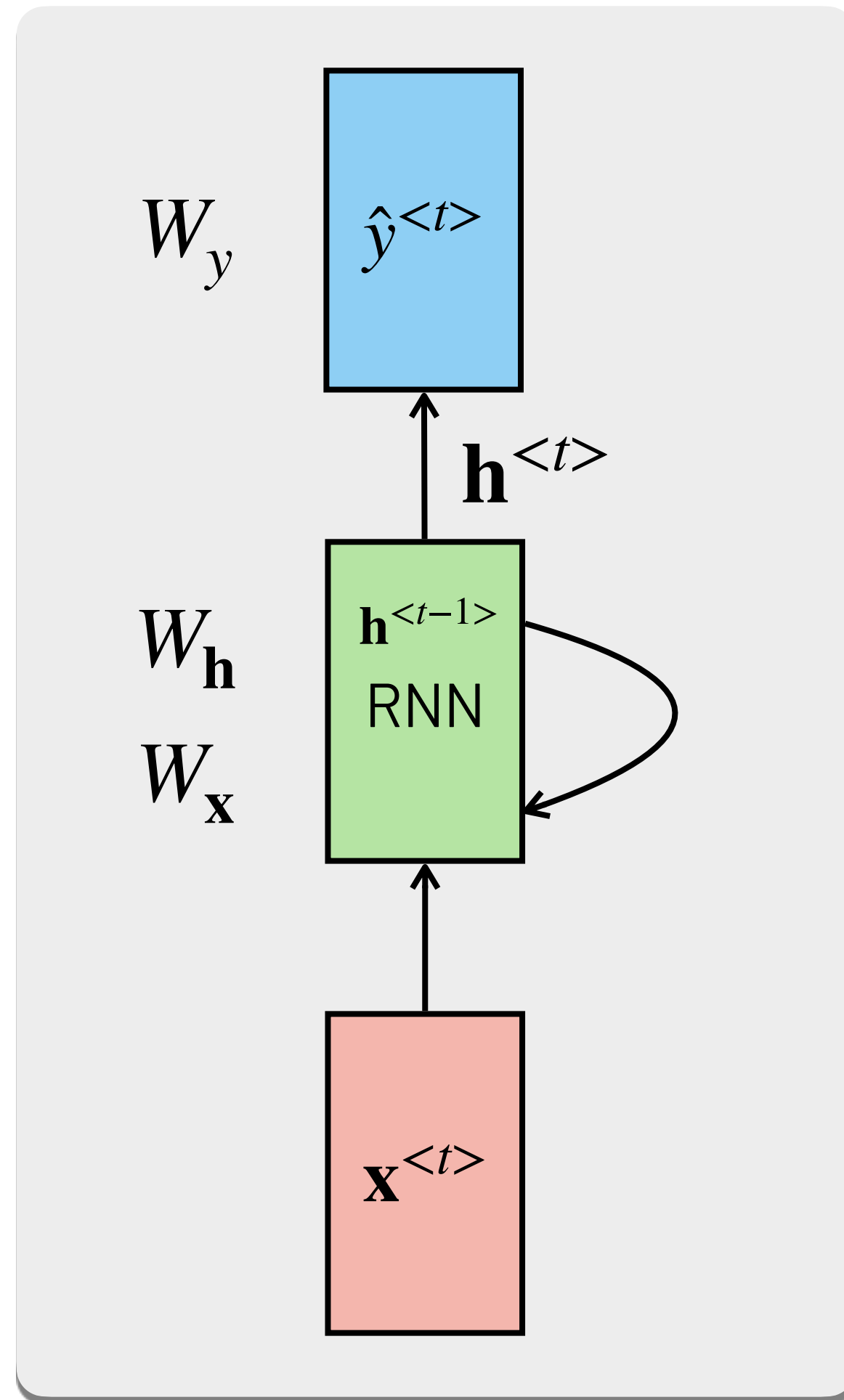
UFV

# 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)



$W_y$

$\hat{y}^{<t>}$

$\mathbf{h}^{<t>}$

$W_\mathbf{h}$
$W_\mathbf{x}$
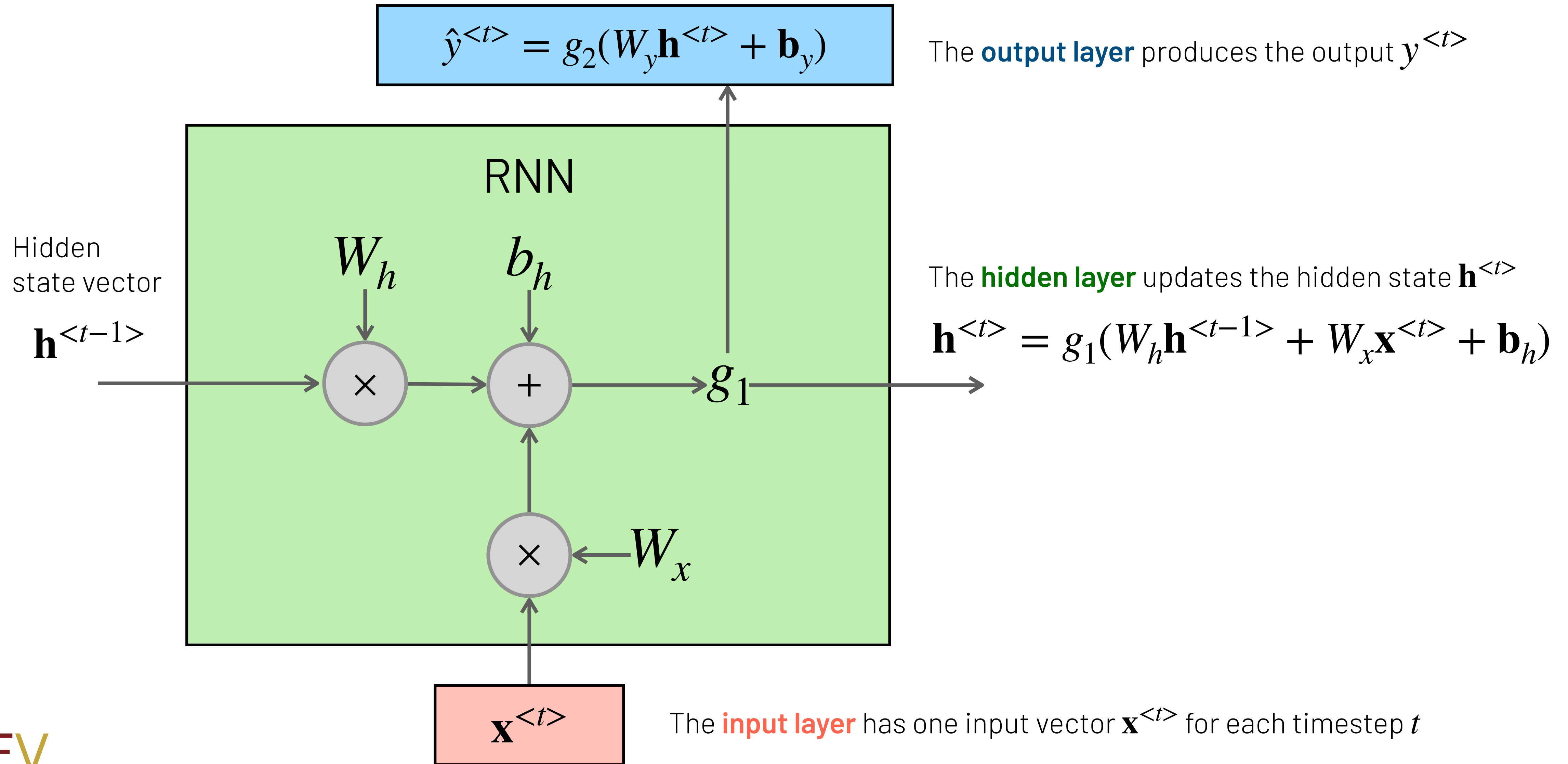
$\mathbf{h}^{<t-1>}$

RNN

$\mathbf{x}^{<t>}$

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 $y^{<t>}$

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

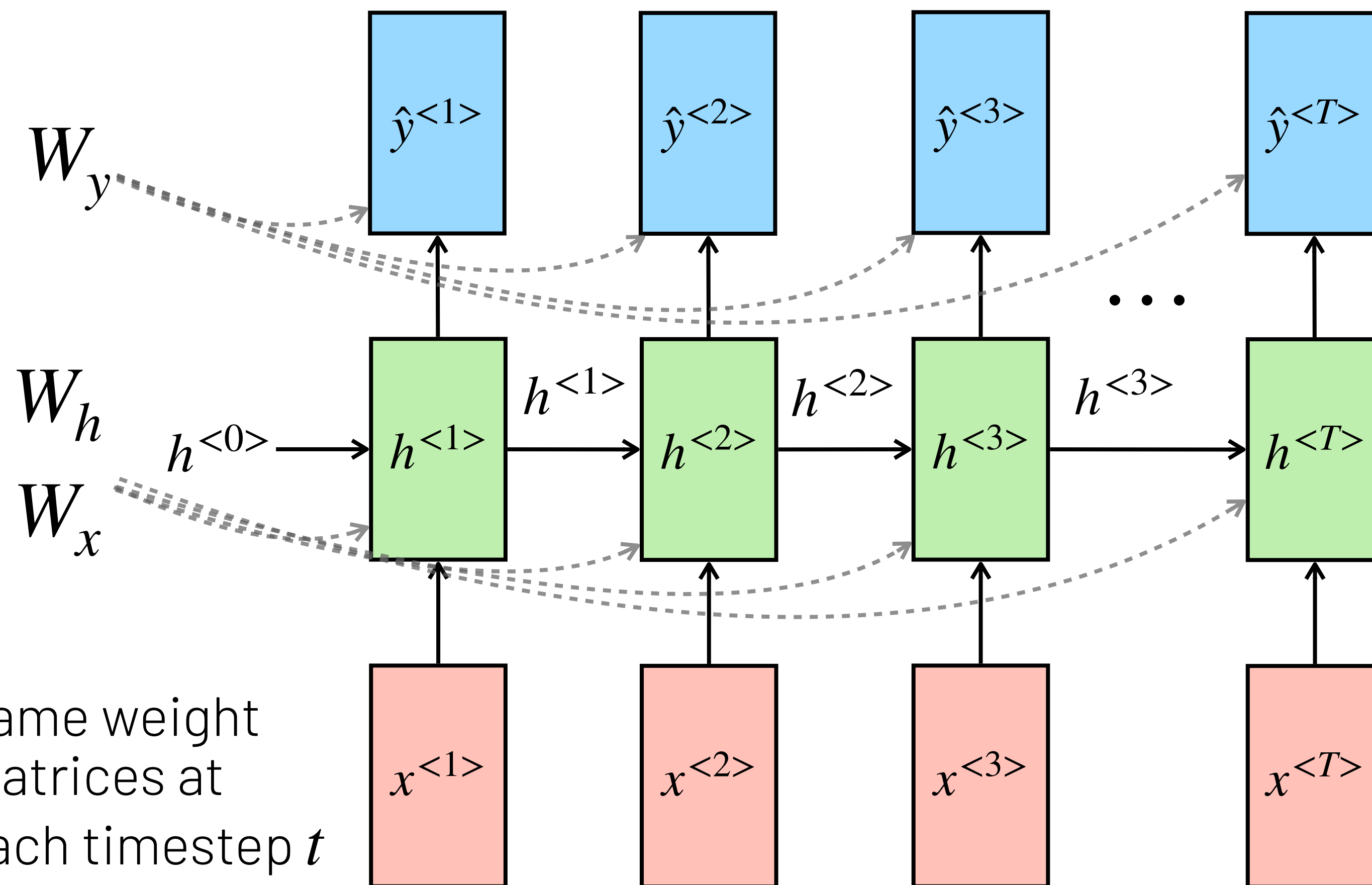$$\hat{y}^{<t>} = g_2(W_y\mathbf{h}^{<t>} + \mathbf{b}_y)$$

▸ $g_1$: hidden layer activation function (tanh/relu)

▸ $g_2$: output layer activation function (sigmoid/softmax)

UFV

# Recurrent Neural Networks (RNNs)

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

The **output layer** produces the output $y^{<t>}$

RNN

$W_h$     $b_h$

Hidden
state vector

The **hidden layer** updates the hidden state $\mathbf{h}^{<t>}$

$\mathbf{h}^{<t-1>}$

$\times$   $+$   $g_1$

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

$\times$   $W_x$

$\mathbf{x}^{<t>}$

The **input layer** has one input vector $\mathbf{x}^{<t>}$ for each timestep $t$

UFV

# Recurrent Neural Networks (RNNs)

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



$$h^{<1>} = g_1(W_h h^{<0>} + W_x x^{<1>} + b_h)$$
$$\hat{y}^{<1>} = g_2(W_y h^{<1>} + b_y)$$

$$h^{<2>} = g_1(W_h h^{<1>} + W_x x^{<2>} + b_h)$$
$$\hat{y}^{<2>} = g_2(W_y h^{<2>} + b_y)$$

$$h^{<3>} = g_1(W_h h^{<2>} + W_x x^{<3>} + b_h)$$
$$\hat{y}^{<3>} = g_2(W_y h^{<3>} + b_y)$$

$$h^{<T>} = g_1(W_h h^{<T-1>} + W_x x^{<T>} + b_h)$$
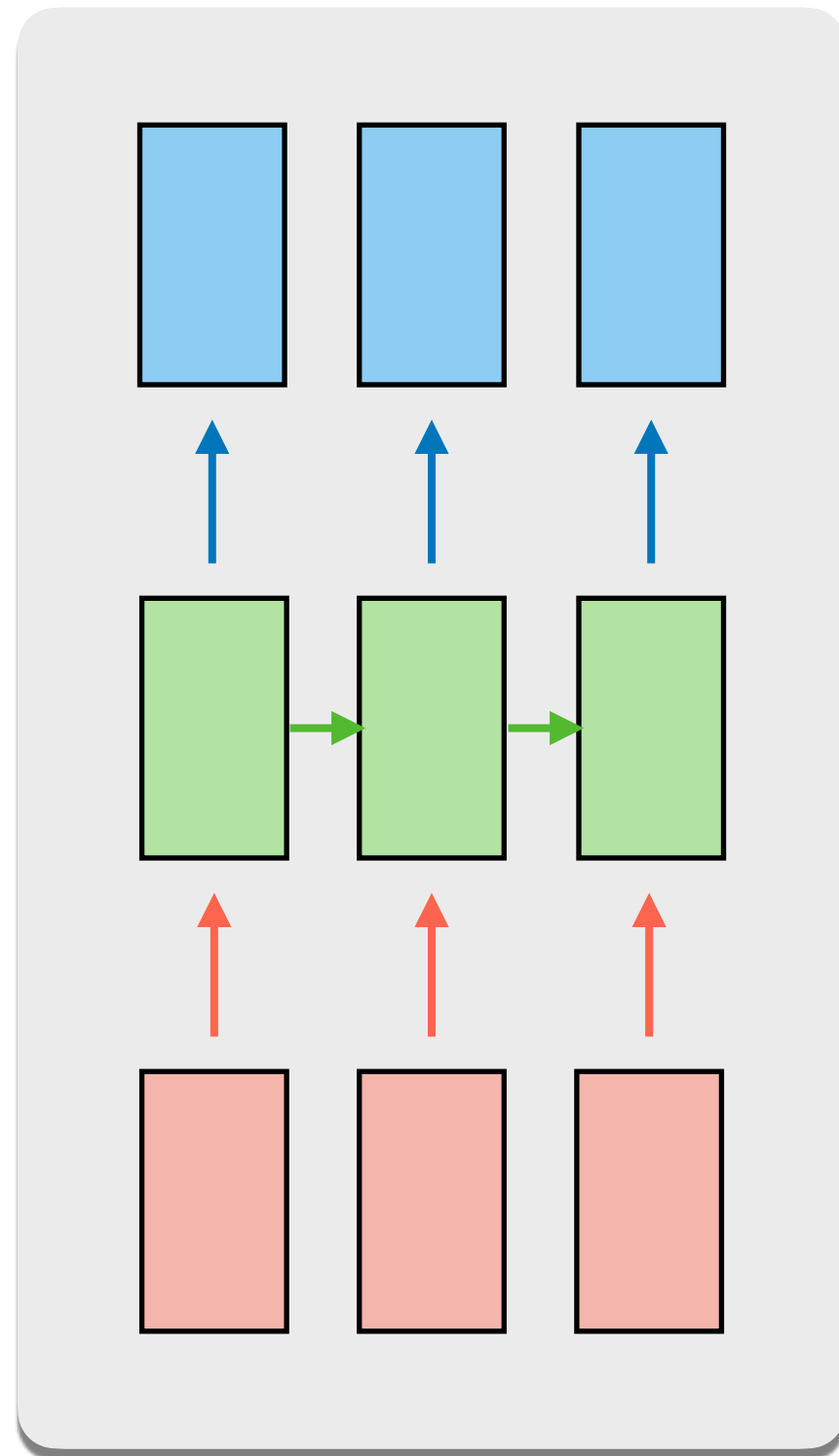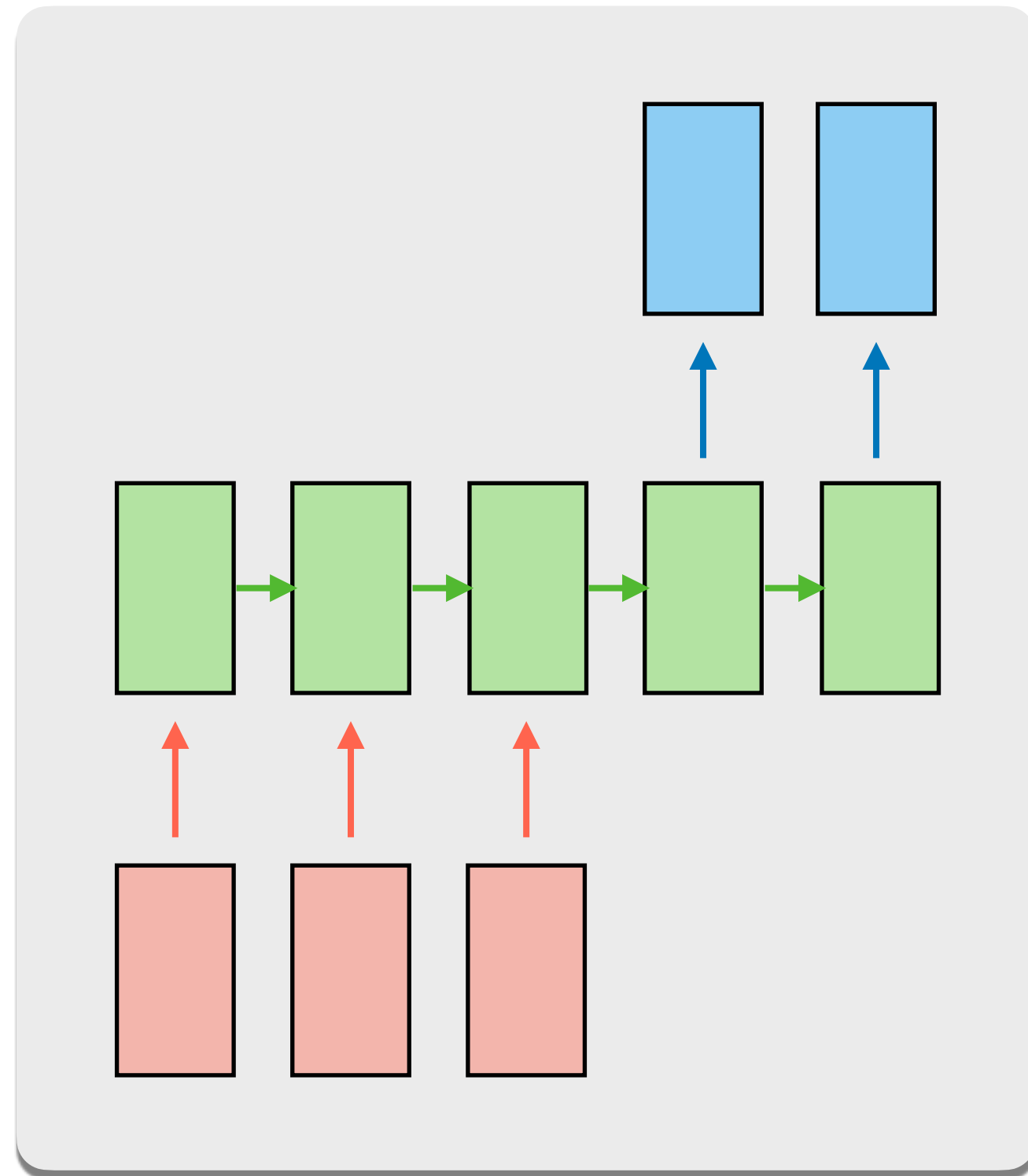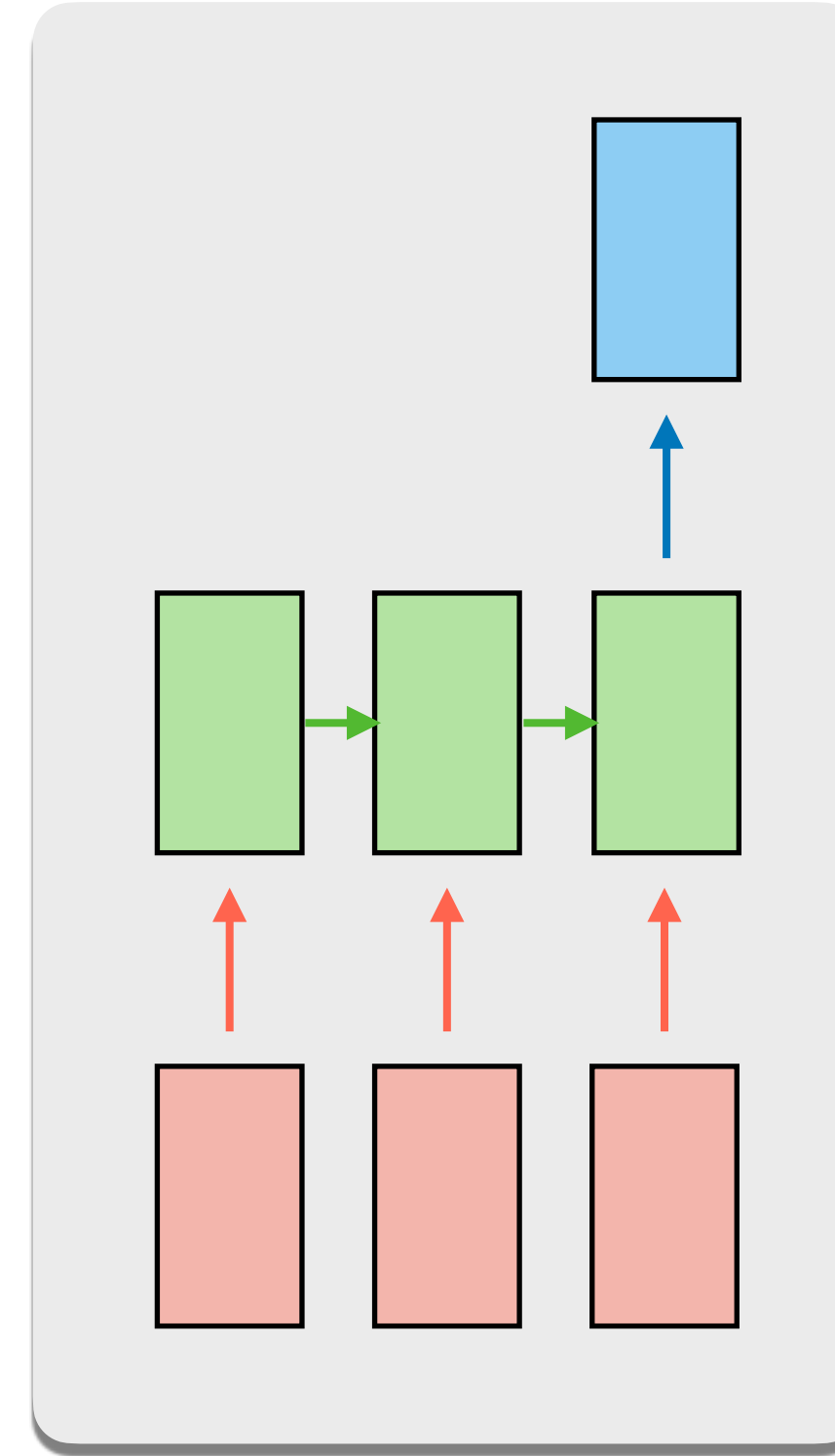$$\hat{y}^{<T>} = g_2(W_y h^{<T>} + b_y)$$

$W_y$

$W_h$

$W_x$

$h^{<0>}$

$h^{<1>}$    $h^{<2>}$    $h^{<3>}$

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$    $\hat{y}^{<3>}$    $\hat{y}^{<T>}$

$h^{<1>}$   $h^{<2>}$   $h^{<3>}$   $h^{<T>}$

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$    $x^{<T>}$

Same weight matrices at each timestep $t$

UFV

# Types of RNNs



**Many** to **Many**

**Example**
Named Entity Recognition

**Many** to **Many (Seq2Seq)**

**Example**
Machine Translation

**Many** to **one**

**Example**
Sentiment Analysis

**One** to **many**

**Example**
Image Description

**one** to **one**

MLP

UFV

12

# Many to Many

An output is produced for each timestep state $t$

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

$$y^{<t>} = g_2(W_y h^{<t>})$$

$$h^{<t>} = g_1(W_h h^{<t-1>} + W_x x^{<t>} + b_h)$$

$L^{<1>}$    $L^{<2>}$    $L^{<3>}$    $L^{<T>}$

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$    $\hat{y}^{<3>}$    $\hat{y}^{<T>}$
1    0    0    2

$h^{<0>}$    $h^{<1>}$    $h^{<2>}$    $h^{<3>}$    $h^{<T>}$

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$    $x^{<4>}$
Lucas    Is    a    UFV

UFV

13

# Many to One

$$L = L^{<t>}$$

$L^{<T>}$

$\hat{y}^{<T>}$
1.05

$$y^{<T>} = g_2(W_y h^{<T>})$$

The network output is produced
only once at the last time step $t$

$h^{<0>}$    $h^{<1>}$    $h^{<2>}$    $h^{<3>}$    $\cdots$    $h^{<T>}$

$$h^{<t>} = g_1(W_h h^{<t-1>} + W_x x^{<t>} + b_h)$$

$x^{<1>}$
It's

$x^{<2>}$
a

$x^{<3>}$
terrible

$\cdots$

$x^{<T>}$
product

UFV

# One to Many

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

An output is produced for each time step $t$

$$y^{<t>} = g_2(W_{yh} h^{<t>})$$

$$h^{<t>} = g_1(W_{hh} h^{<t-1>} + W_{hx} x^{<t>} + b_h)$$

Encoder

$$x^{<t>} = \hat{y}^{<t-1>}$$

Use **x** to initialize $h^{<0>}$

$L^{<1>}$ $L^{<2>}$ $L^{<3>}$ $L^{<T>}$

$\hat{y}^{<1>}$ A $\hat{y}^{<2>}$ cat $\hat{y}^{<3>}$ lying $\hat{y}^{<T>}$ \<EOS\>

$h^{<0>}$ $h^{<1>}$ $h^{<2>}$ $h^{<3>}$ $h^{<T>}$

$y^{<0>}$ \<SOS\> $\hat{y}^{<1>}$ A $\hat{y}^{<2>}$ cat $\hat{y}^{<T-1>}$ window

UFV

15

# Backpropagation Through Time



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

Process all $T$ elements of the sequence to calculate the loss

Backpropagate through the entire sequence to calculate the gradient
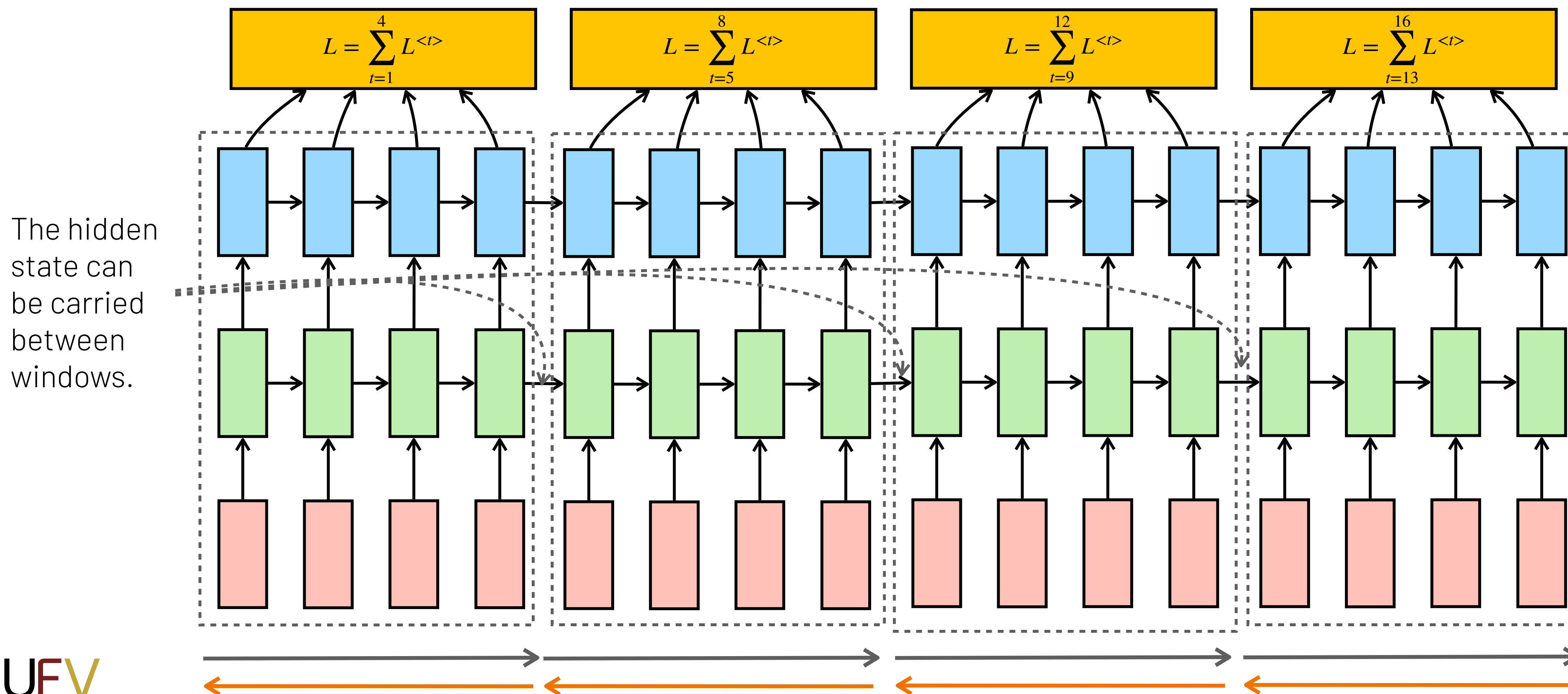
# Truncated Backpropagation Through Time

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)



The hidden state can be carried between windows.

$$L = \sum_{t=1}^{4} L^{<t>}$$

$$L = \sum_{t=5}^{8} L^{<t>}$$

$$L = \sum_{t=9}^{12} L^{<t>}$$

$$L = \sum_{t=13}^{16} L^{<t>}$$

# Next Lecture

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

GRUs and LSTMs for processing with very long sequences.

UFV