

# Introducción a Redes Neuronales Artificiales

Mauricio A. Álvarez PhD,  
H.F. Garcia   C. Guarnizo (TA)



Universidad Tecnológica de Pereira, Pereira, Colombia

## 1 Introduction

## 2 Recurrent Neural Networks

- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications



## 1 Introduction

## 2 Recurrent Neural Networks

- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications



# Sequence Learning

- 1 Sequence learning is the study of machine learning algorithms designed for sequential data.
- 2 Discrete dynamical systems based on delay steps.
- 3 Language model is one of the most interesting topics that use sequence labeling.
  - Language modeling.
  - Machine Translation.



# Types of RNNs

one to one

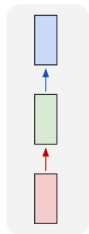
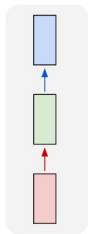


Image Classification  
Vanilla RNN



# Types of RNNs

one to one



one to many

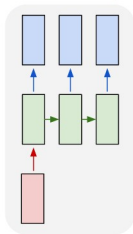
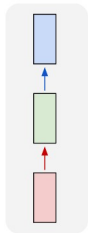


Image Captioning

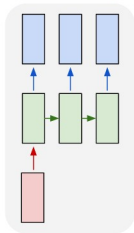


# Types of RNNs

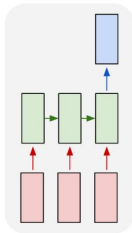
one to one



one to many



many to one

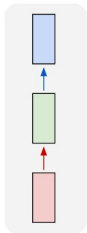


Sentiment Classification

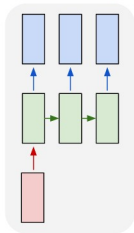


# Types of RNNs

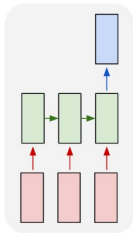
one to one



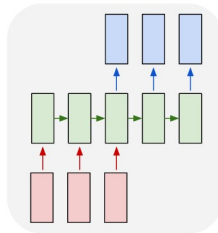
one to many



many to one



many to many



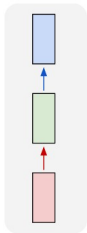
Machine Translation



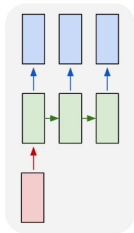


# Types of RNNs

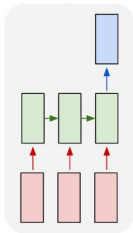
one to one



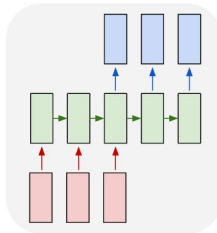
one to many



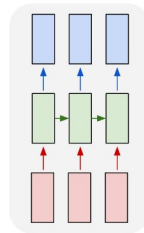
many to one



many to many



many to many



Video Labeling (frame)



## 1 Introduction

## 2 Recurrent Neural Networks

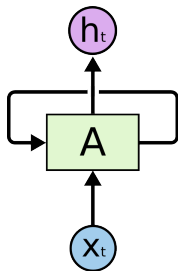
- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications



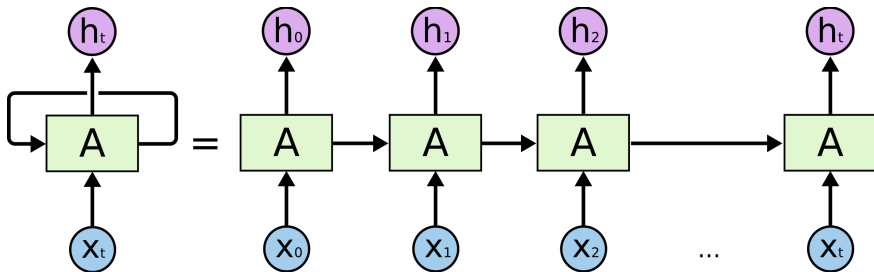
# Recurrent Neural Network - RNN

$$h_t \in \mathbb{R}^N, x_t \in \mathbb{R}^M \rightarrow h_t = f_W(h_{t-1}, x_t), y_t = g_{W_o}(h_t)$$



# Recurrent Neural Network - RNN

$$h_t \in \mathbb{R}^N, x_t \in \mathbb{R}^M \rightarrow h_t = f_W(h_{t-1}, x_t), y_t = g_{W_o}(h_t)$$



## 1 Introduction

## 2 Recurrent Neural Networks

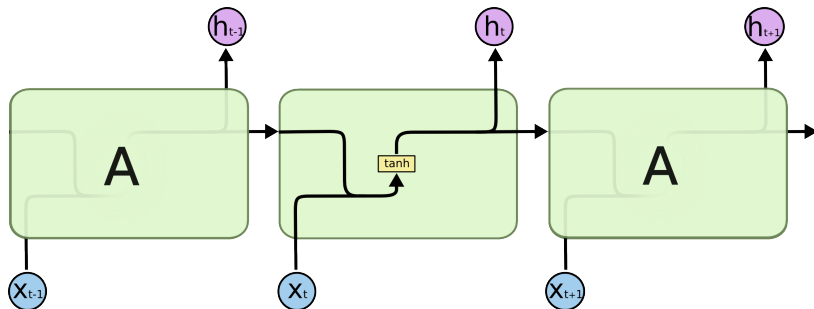
- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications

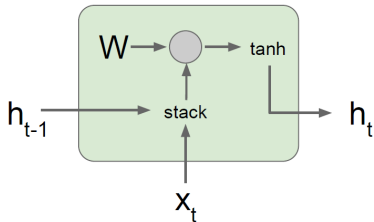


# Vanilla RNN

$$h_t \in \mathbb{R}^N, x_t \in \mathbb{R}^M \rightarrow h_t = \tanh(W_{hh}h_{t-1} + W_{hx}x_t), y_t = \mathbf{W}_o h_t$$

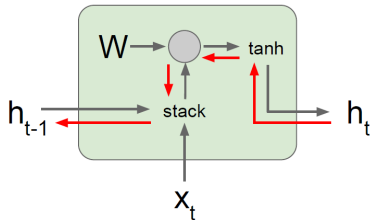


# Vanilla RNN



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

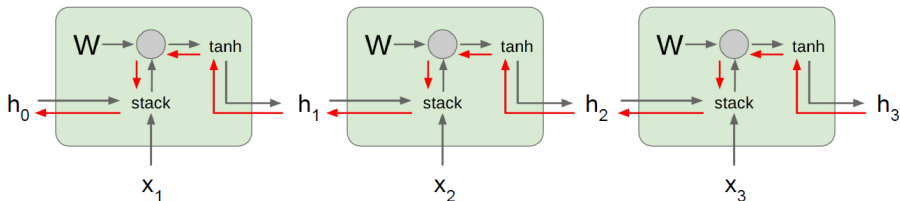
# Vanilla RNN



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$



# Vanilla RNN



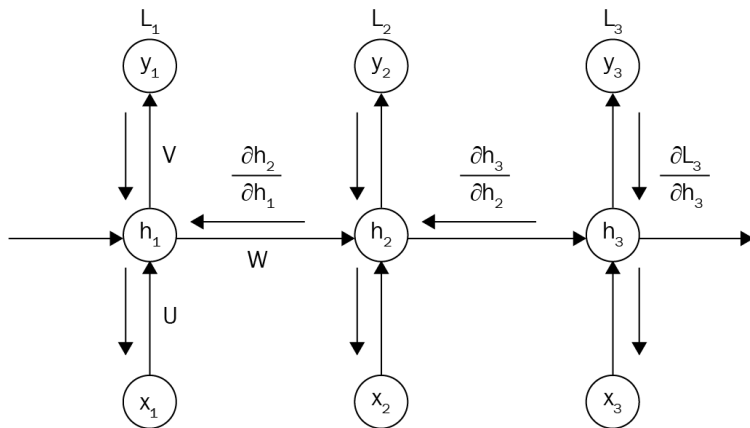
$$h_1 = \tanh(W_{hh}h_0 + W_{hx}x_1)$$

$$h_2 = \tanh(W_{hh}(\tanh(W_{hh}h_0 + W_{hx}x_1) + W_{hx}x_2)$$

$$h_3 = \tanh(\textcolor{red}{W}_{hh}(\tanh(\textcolor{red}{W}_{hh} \tanh(\textcolor{red}{W}_{hh}h_0 + W_{hx}x_1) + W_{hx}x_2) + W_{hx}x_3)$$



# Vanilla RNN: Gradients



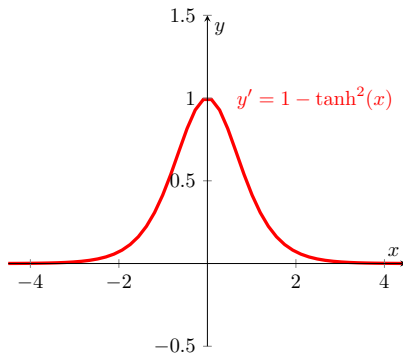
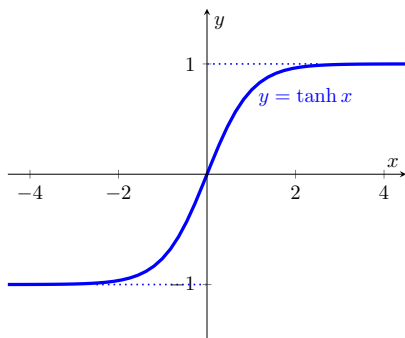
# Vanilla RNN: Gradient problems

$$\begin{aligned}\frac{\partial E}{\partial \theta} &= \sum_{1 \leq t \leq T} \frac{\partial E_t}{\partial \theta} \\ \frac{\partial E_t}{\partial \theta} &= \sum_{1 \leq k \leq t} \left( \frac{\partial E_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial \theta} \right) \\ \frac{\partial h_t}{\partial h_k} &= \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{t \geq i > k} W_{hh}^\top \text{diag}(\sigma'(h_{i-1}))\end{aligned}$$

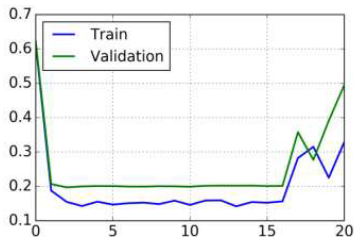
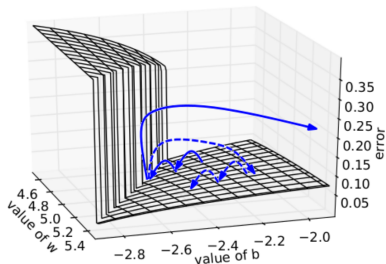
R. Pascanu - On the difficulty of training recurrent neural networks, 2013.



# Vanilla RNN: Vanishing Gradients



# Vanilla RNN: Exploding Gradients



$$\frac{\partial h_t}{\partial h_k} = \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{t \geq i > k} W_{hh}^\top \text{diag}(\sigma'(h_{i-1}))$$

## **1** Introduction

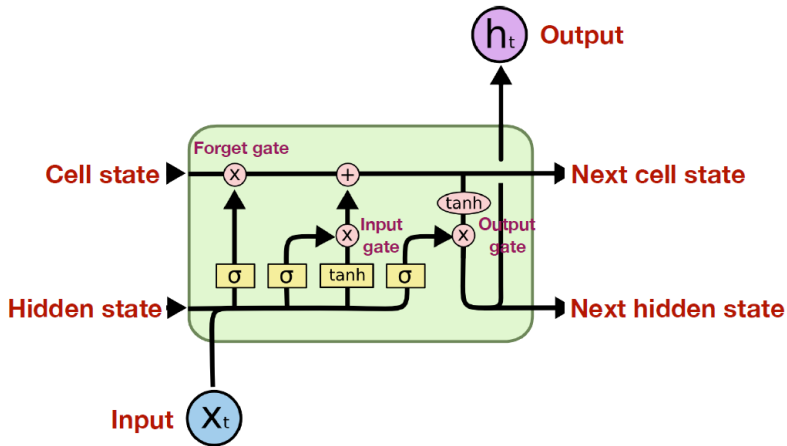
## **2** Recurrent Neural Networks

- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## **3** Applications



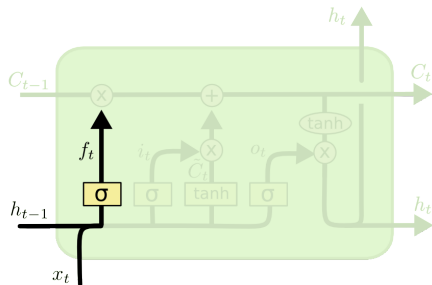
# LSTM - Long-Short Term Memory



<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>



# LSTM - Forget and Keep gate

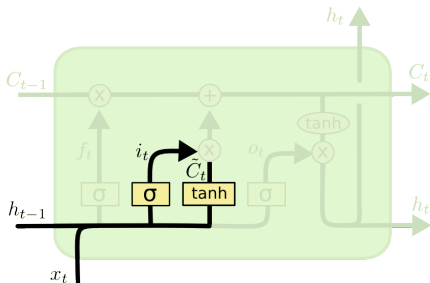


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Decide which information to throw away from the cell state.



# LSTM - Input Gate



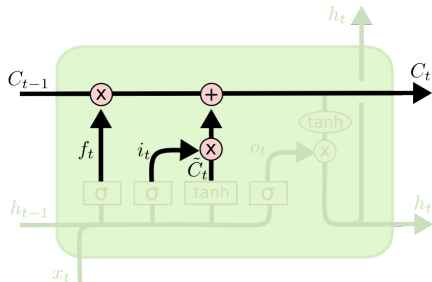
$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Decide which information to store to the cell state.



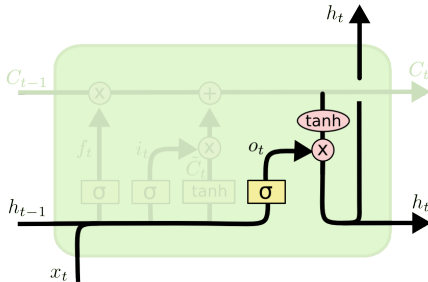
# LSTM - Update Cell State



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Update the cell state scaled by input and forget gates.

# LSTM - Output Gate



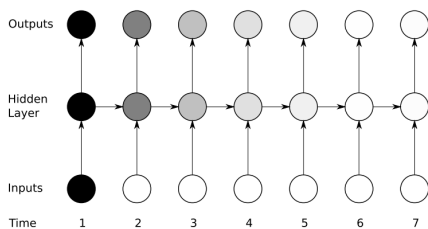
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

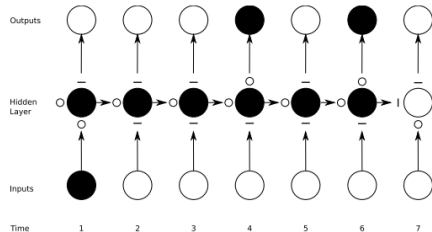
Output based on the updated cell state.



# LSTM vs. Vanilla RNN



Vanilla RNN



LSTM



## 1 Introduction

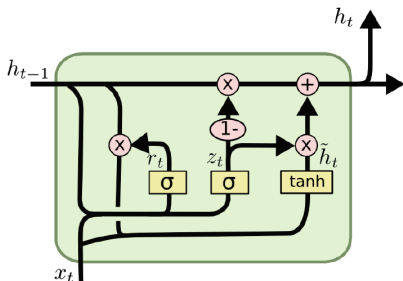
## 2 Recurrent Neural Networks

- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications



# Gated Recurrent Unit - GRU



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

<https://towardsdatascience.com/understanding-gru-networks-2ef37df6c9be>



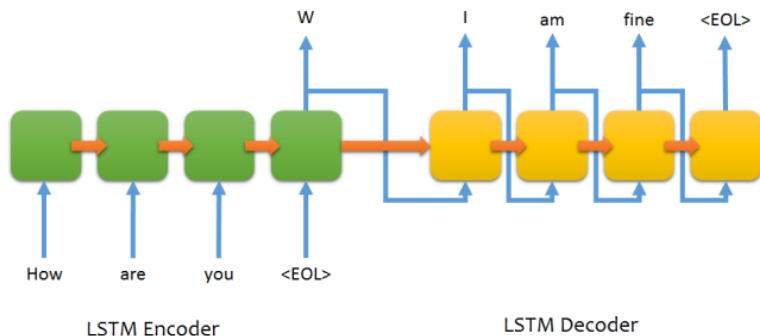
## 1 Introduction

## 2 Recurrent Neural Networks

- Vanilla RNN
- LSTM - Long-Short Term Memory
- Gated Recurrent Unit - GRU

## 3 Applications



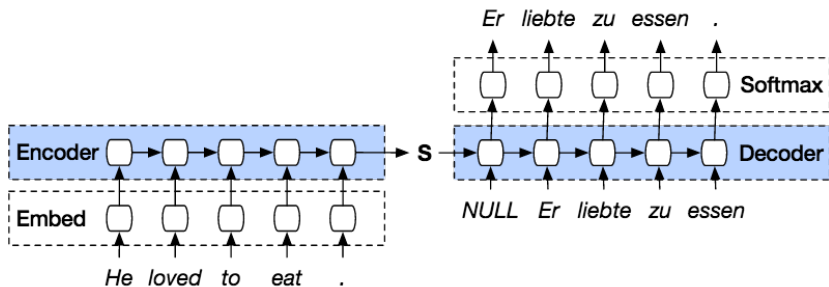


<https://github.com/farizrahman4u/seq2seq>





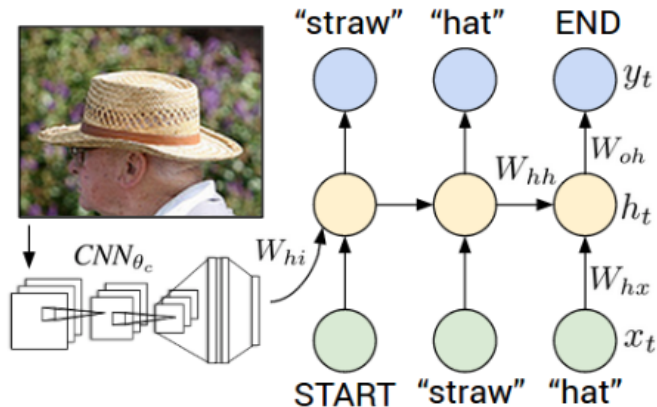
# Machine Translation



<https://github.com/farizrahman4u/seq2seq>



# Image Captioning



<http://bit.ly/neuraltalkdemo>



- 1 Vanilla RNNs presents two main issues during its learning based on backprop, i.e. *Vanishing* and *Exploding* gradients.
- 2 LSTM and GRU networks are able to overcome the above problems, and are the most used RNNs.

