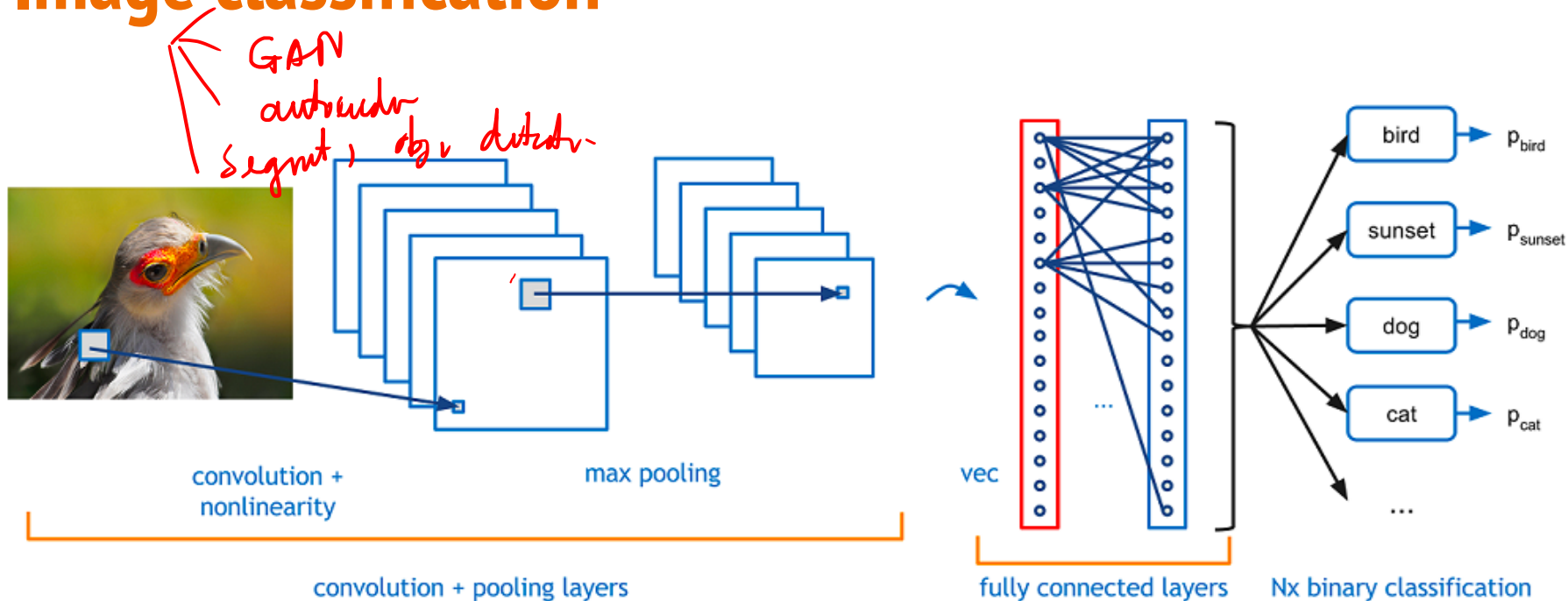

Recurrent Neural Network (RNN)

— Tuan Nguyen - AI4E —

Outline

- Motivation for RNN
- Introduction to RNN
- The structure RNN
- Deep RNN
- RNN application

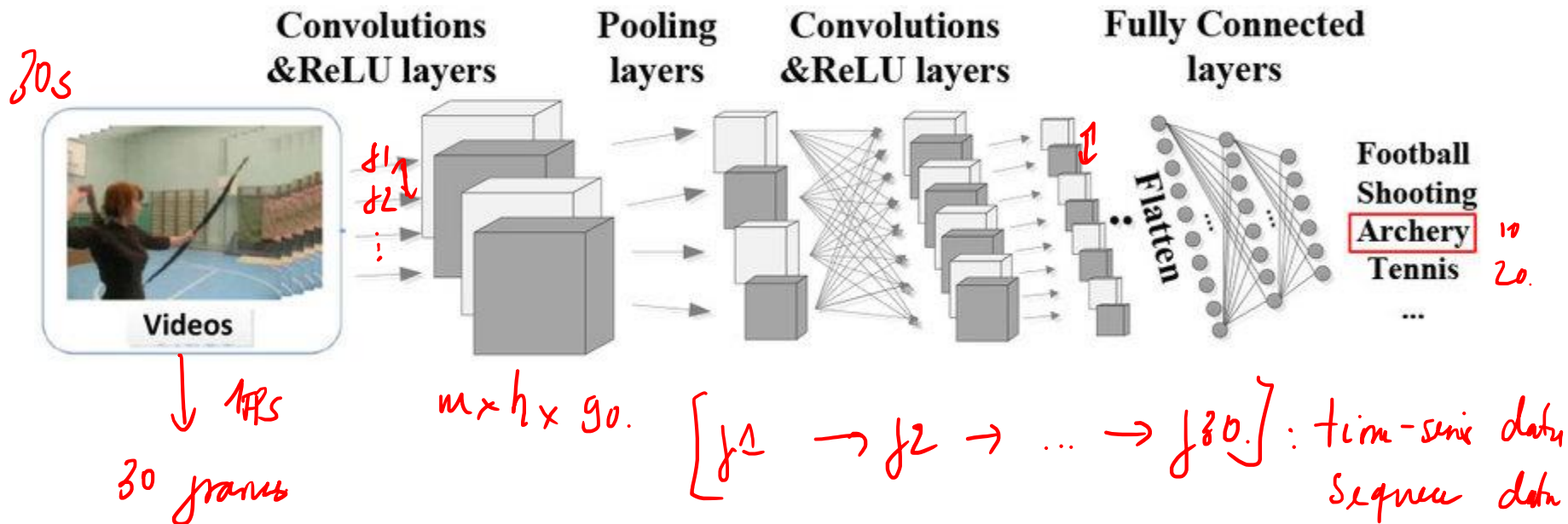
Image classification



Video classification

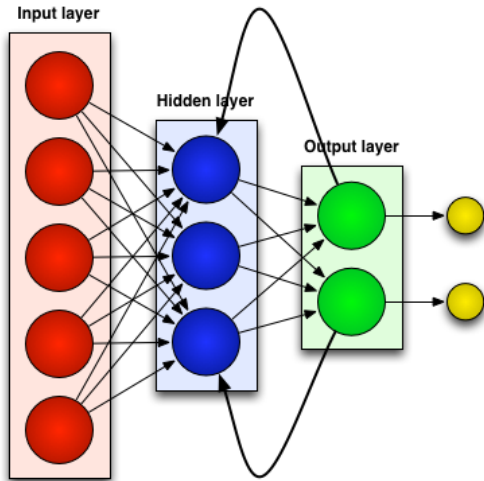
Trần Vũ Việt Nam

https://www.youtube.com/watch?v=BmoVg3BCwF0&feature=emb_title

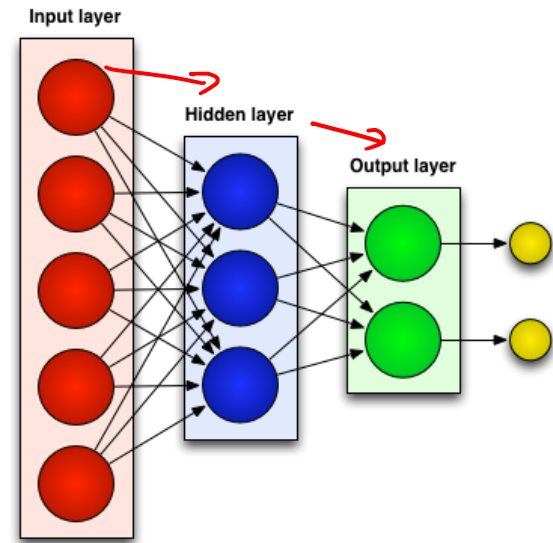


Neural network

Recurrent Neural Network (RNN)

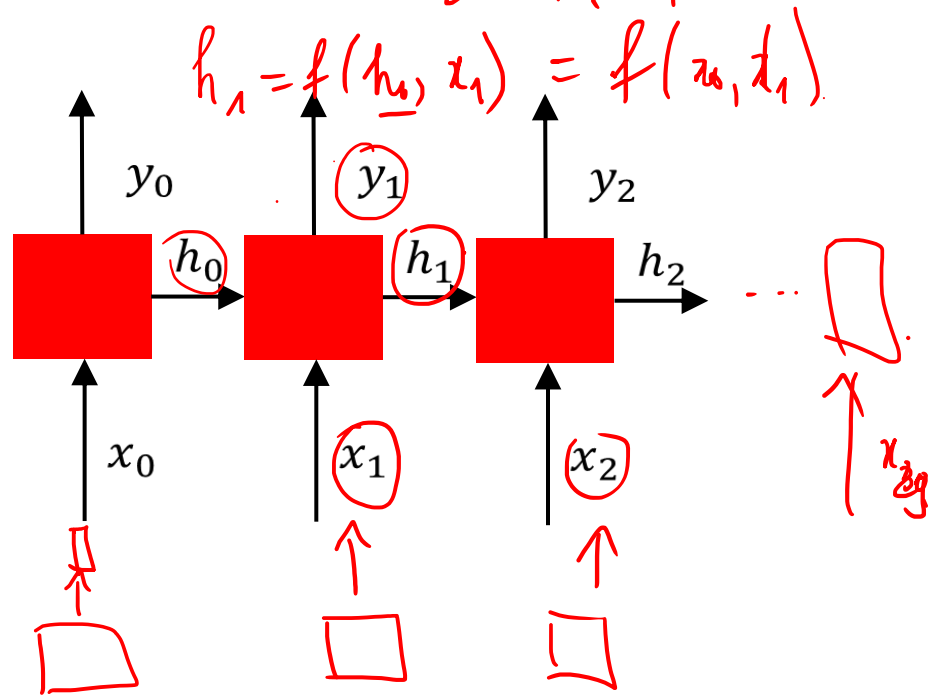
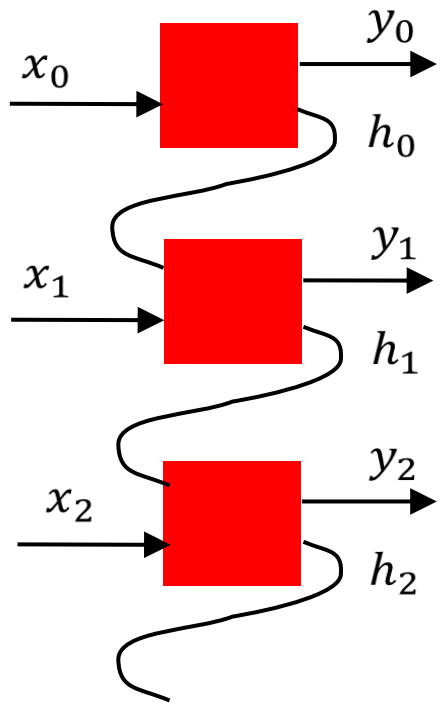


Feed forward neural network



Recurrent Neural Network

Usually drawn as:

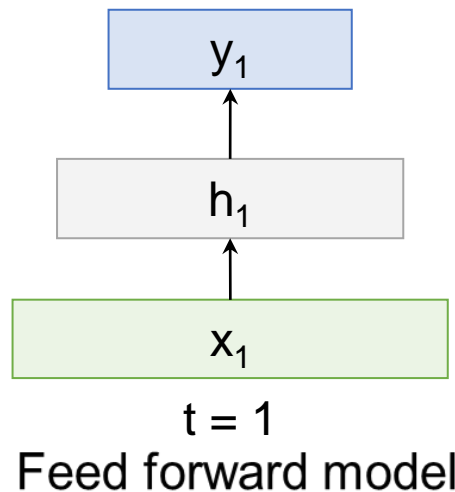
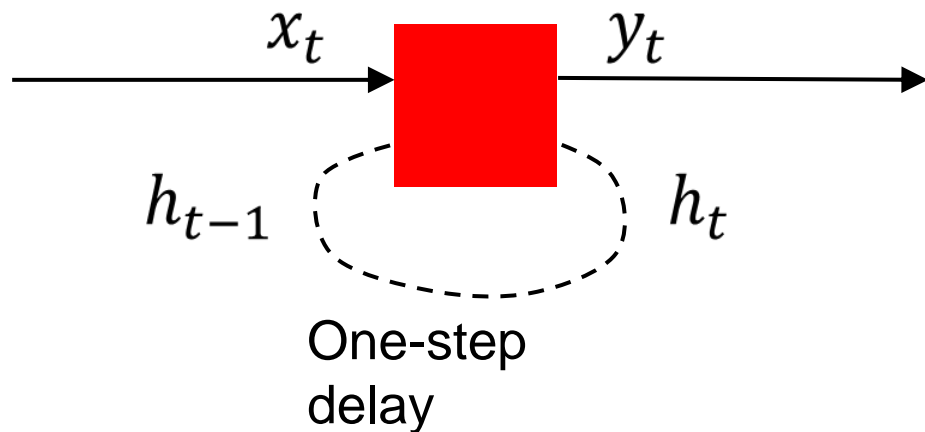


- Input can vector

$$h_0 = f(a_0)$$
$$h_1 = f(h_0, x_1)$$
$$h_2 = f(h_1, x_2)$$
$$h_1 = f(h_0, x_1) = f(x_0, x_1)$$

Recurrent neural network

Recurrent networks introduce cycles and a notion of time.



- They are designed to process sequences of data x_1, \dots, x_n and can produce sequences of outputs y_1, \dots, y_m .

RNN Formula

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

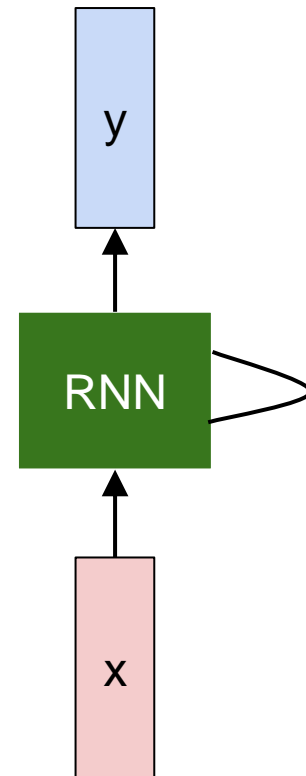
$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state

some function with parameters W

old state

input vector at some time step

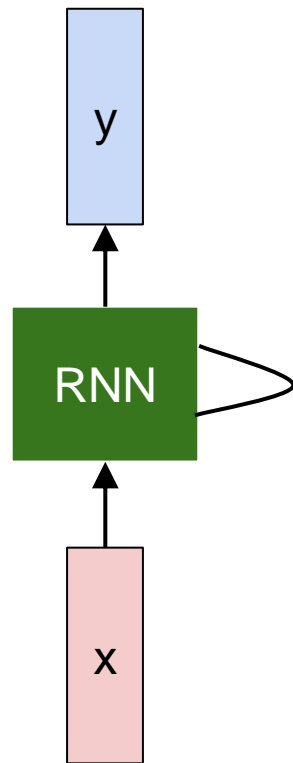


RNN Formula

We can process a sequence of vectors \mathbf{x} by applying a recurrence formula at every time step:

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

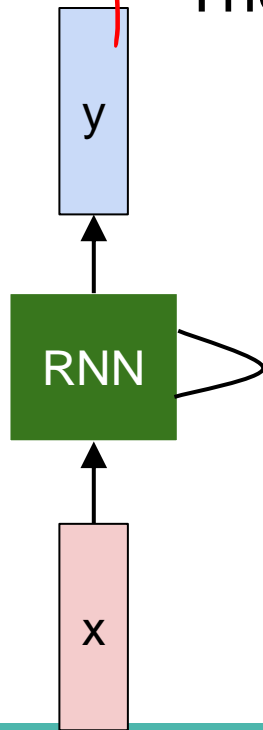
Notice: the same function and the same set of parameters are used at every time step.



RNN Formula

· sigmoid: $(0 \rightarrow 1)$

· relu.



$(-1, 1)$

The state consists of a single “hidden” vector \mathbf{h} :

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

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

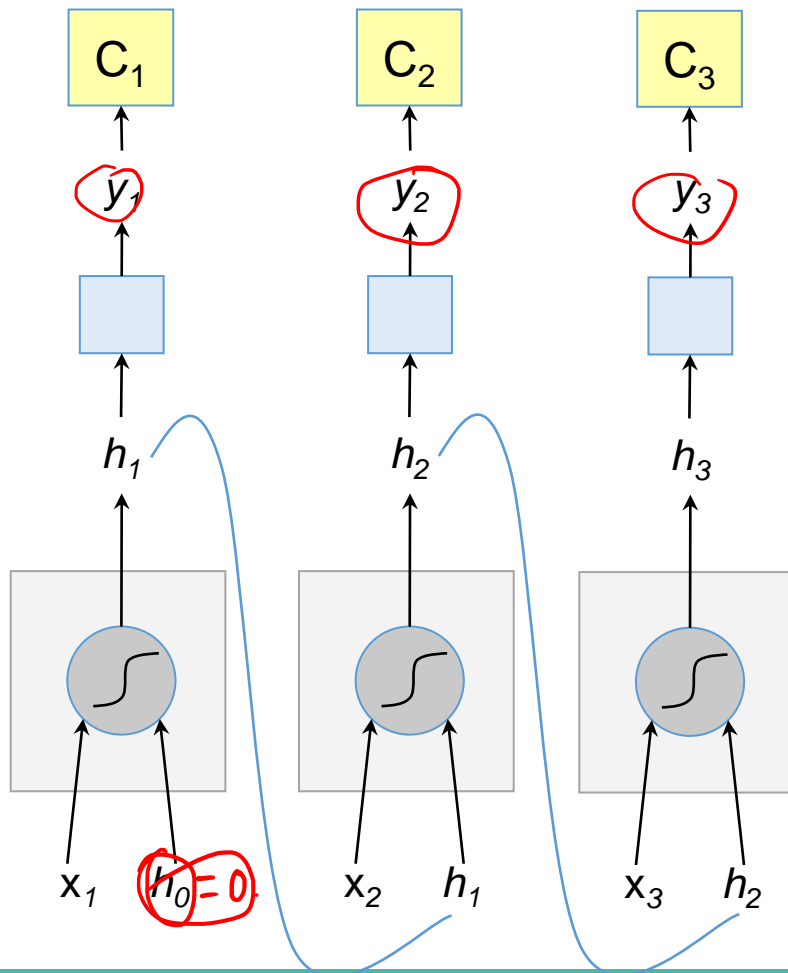
$$\begin{aligned} x &\in \mathbb{R}^{m \times 1} \\ h &\in \mathbb{R}^{n \times 1} \\ y &\in \mathbb{R}^{k \times 1} \end{aligned}$$

$$N(0, 1)$$

Batch-norm

LSTM
GRU

Forward

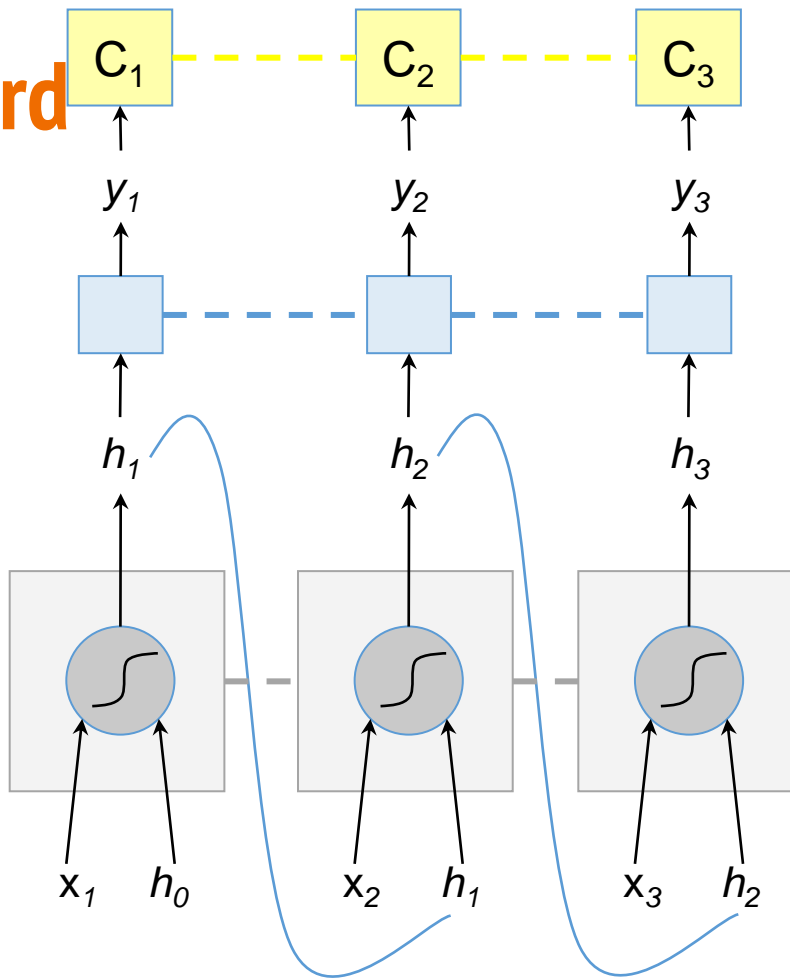


$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = F(h_t)$$

$$C_t = \text{Loss}(y_t, \text{GT}_t)$$

Forward



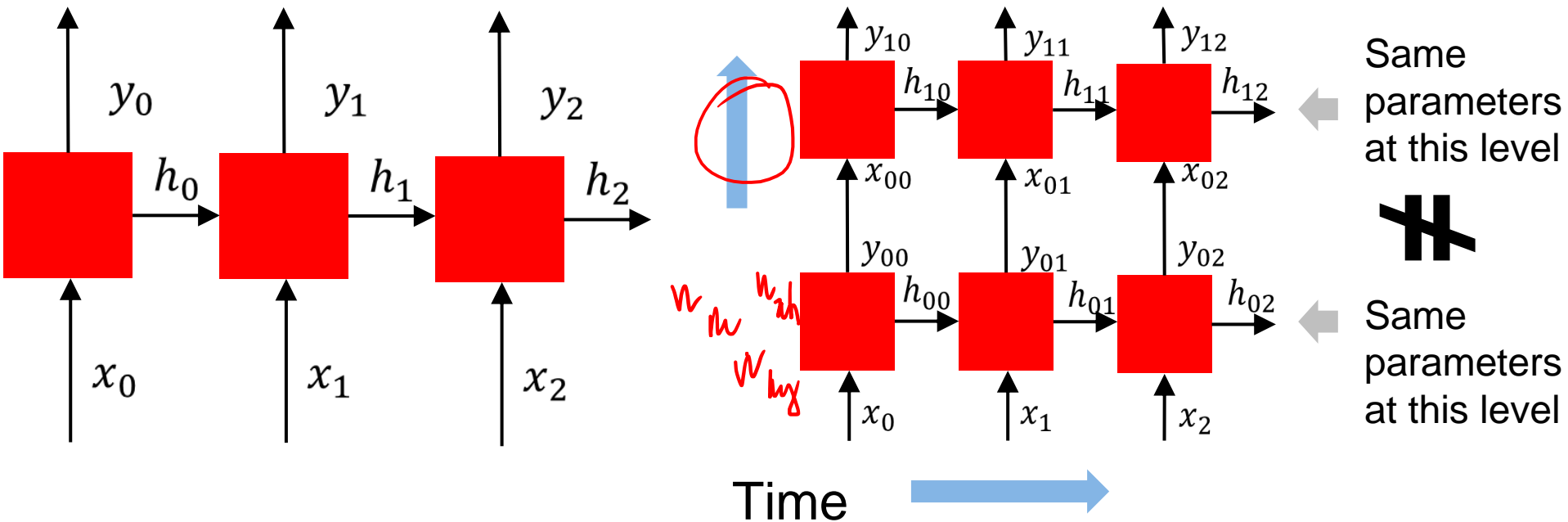
$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = F(h_t)$$

$$C_t = \text{Loss}(y_t, \text{GT}_t)$$

--- indicates shared weights

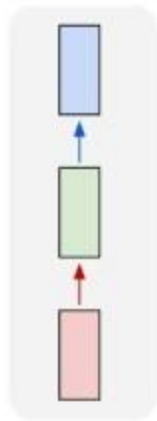
Deep RNN



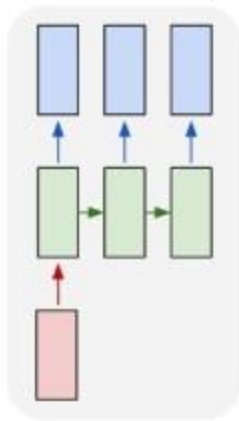
Recurrent neural network problem

one: 1 ảnh, 1 từ.
many: video, câu.

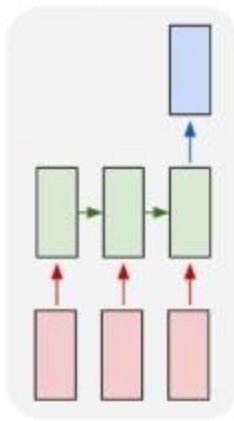
one to one



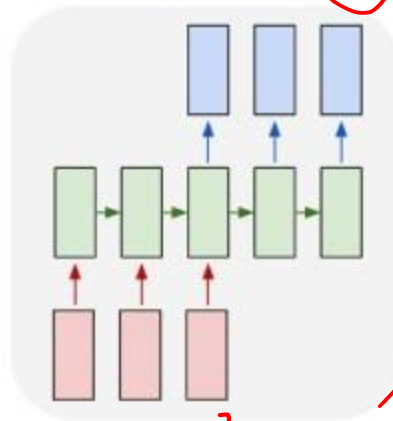
one to many



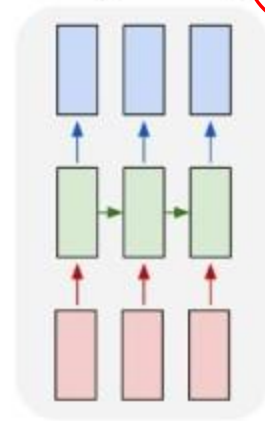
many to one



many to many (1)



many to many (2)



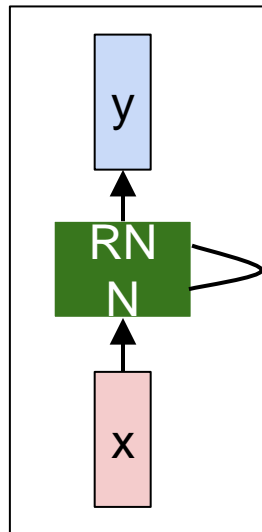
deep fake :
ảnh → video
image captioning :
ảnh → Caption

- mô tả video.
- dịch.

Character-level language model example

Vocabulary:
[h,e,l,o]

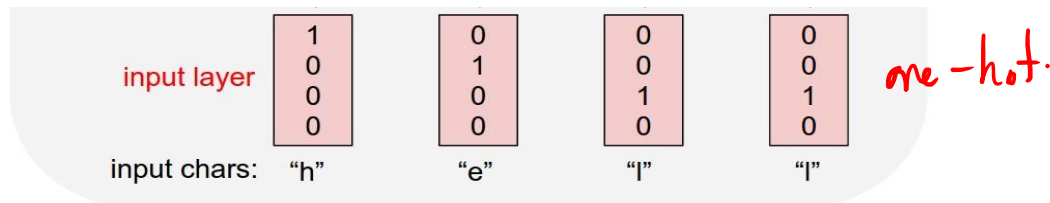
Example training
sequence:
“hello”



Character-level language model example

Vocabulary:
[h,e,l,o]

Example training
sequence:
“hello”



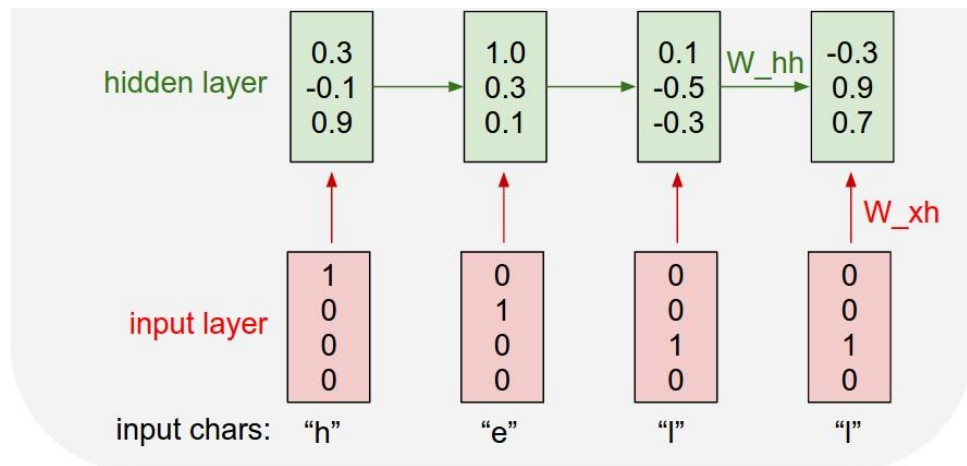
Character-level language model example

Vocabulary:

[h,e,l,o]

Example training
sequence:
“hello”

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



Character-level language model example

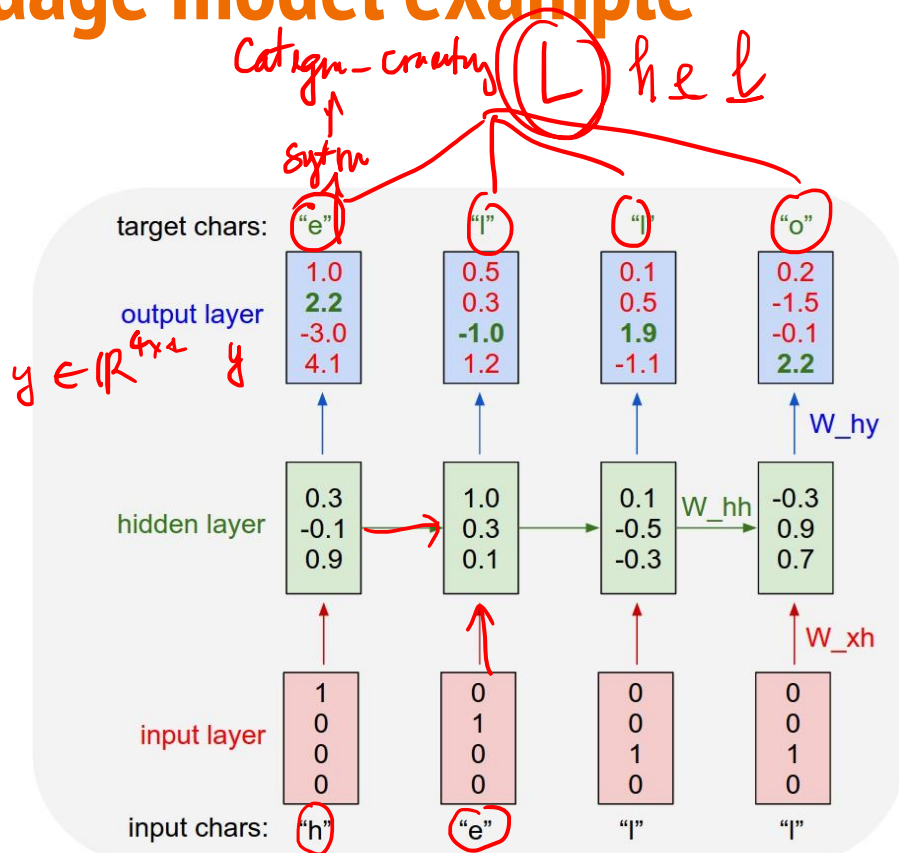
Vocabulary:

[h,e,l,o]

phân loại

Example training
sequence:

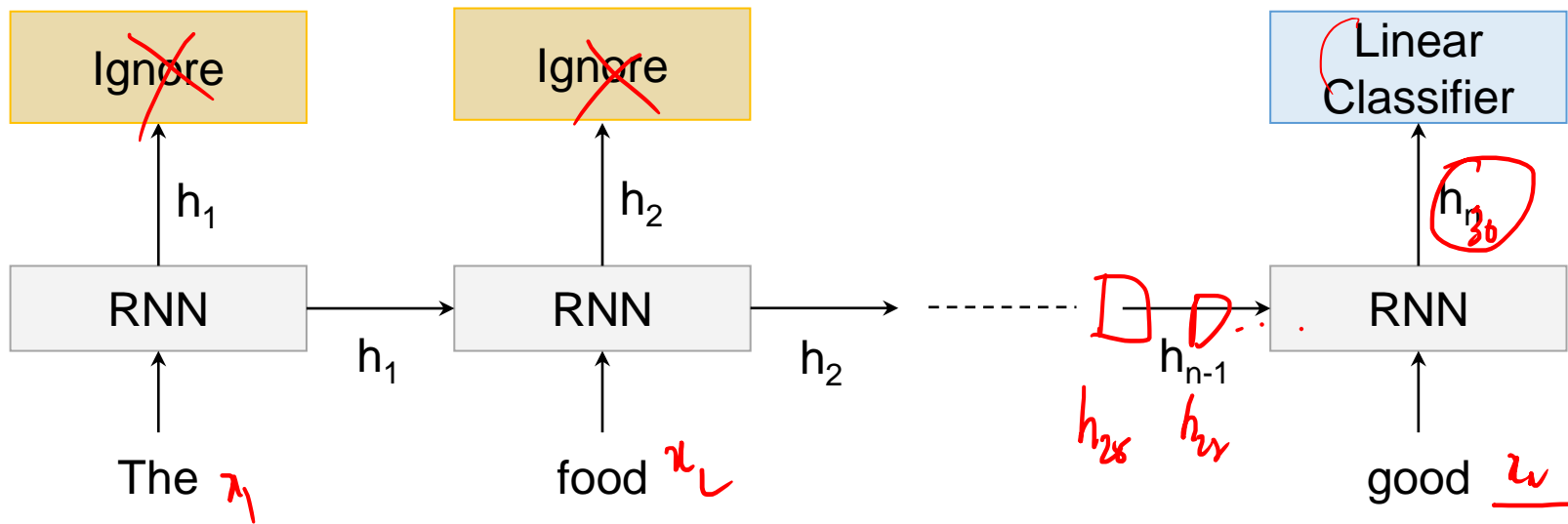
“hello”



Sentiment Classification

~~RNN~~ \leftarrow LSTM
GRU.

$$h_1 = f(x_1)$$
$$h_2 = f(x_2, h_1) = f(x_2, f(x_1)) = f(x_1, x_2)$$



Sentiment Classification

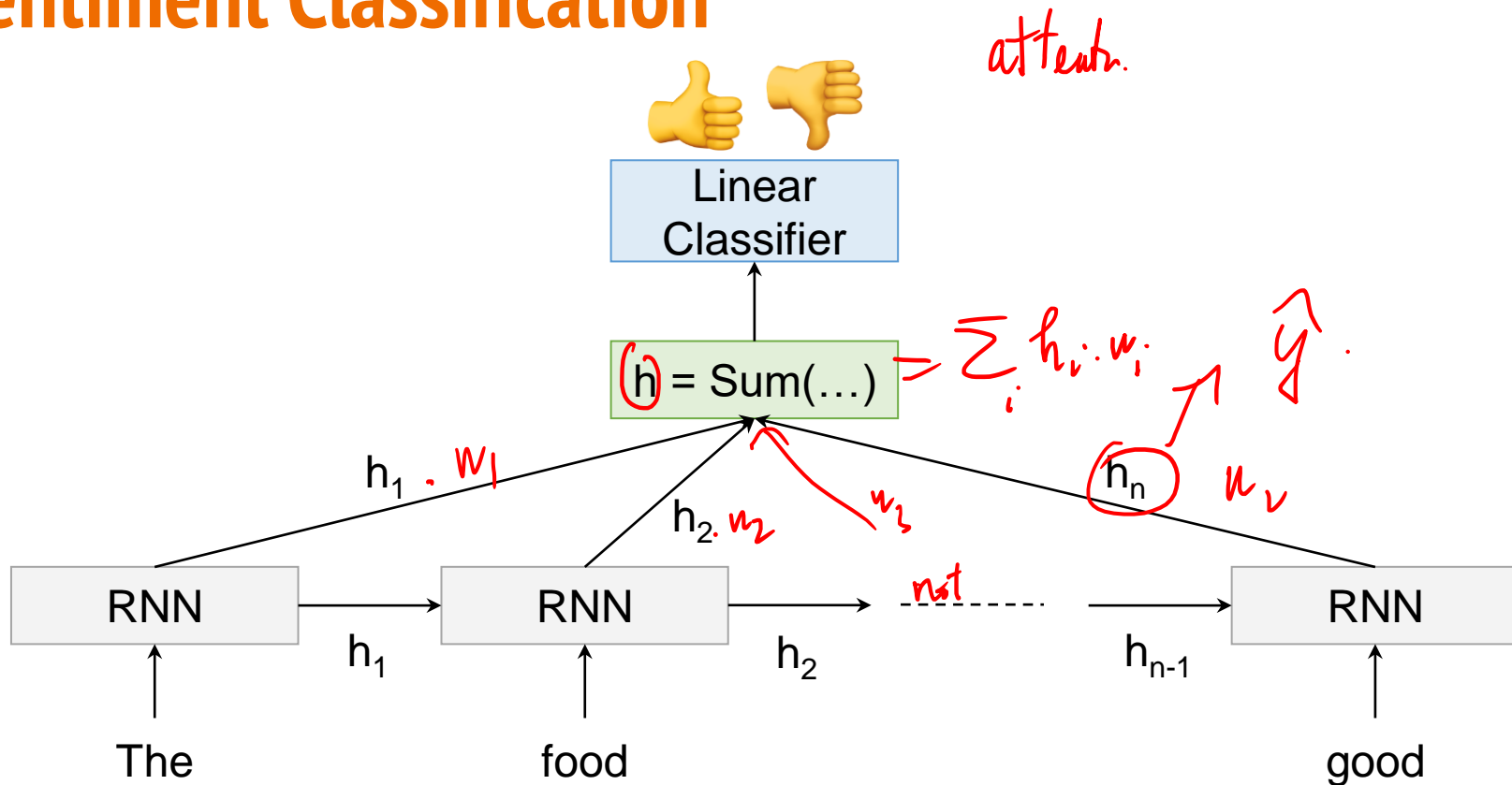


Image Captioning

- Given an image, produce a sentence describing its contents
- Inputs: Image feature (from a CNN)
- Outputs: Multiple words (let's consider one word)



: The man with a straw hat

Image captioning

hello

max = 30

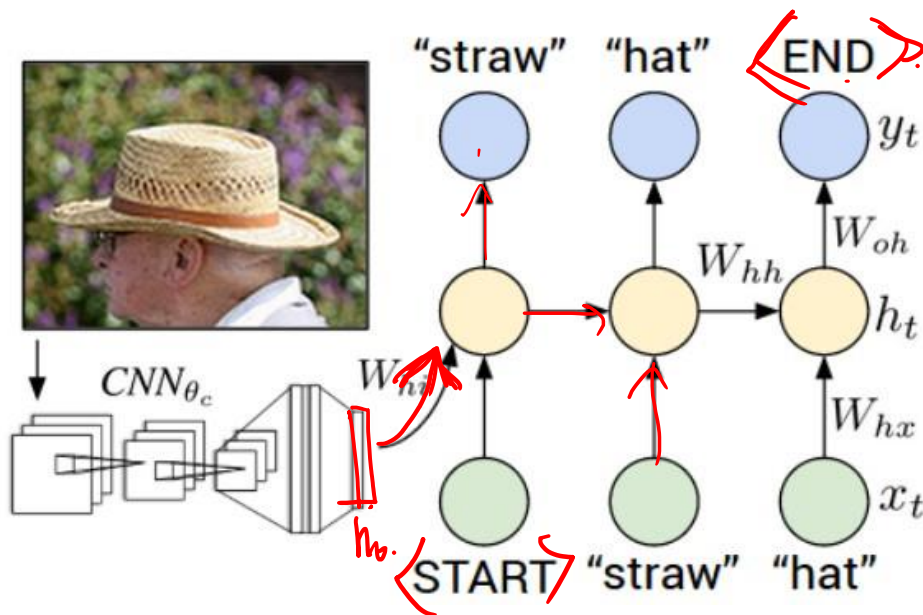
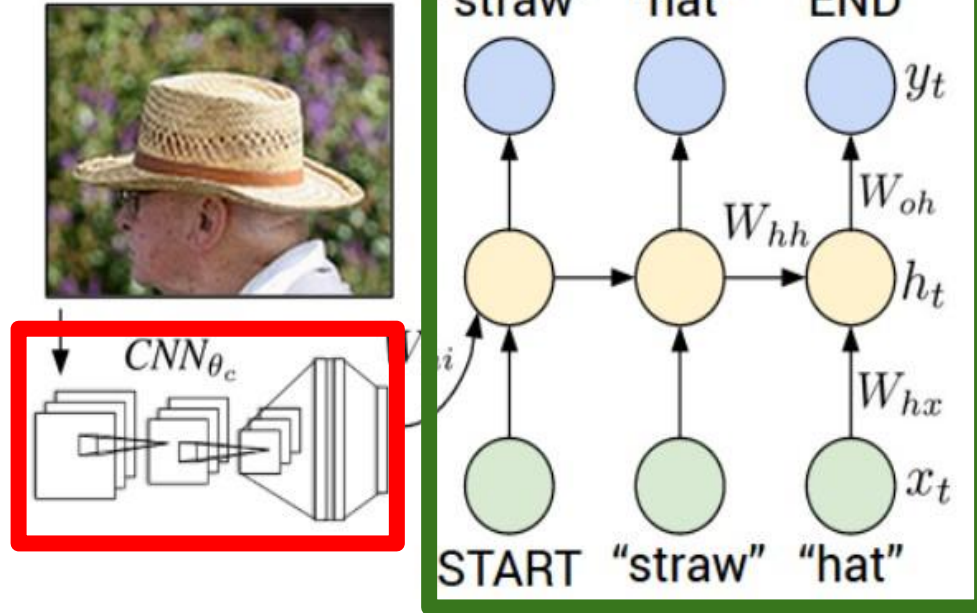


Image captioning

Recurrent Neural Network



Convolutional Neural Network

Image captioning

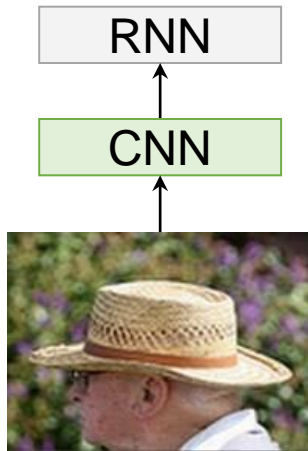


Image captioning

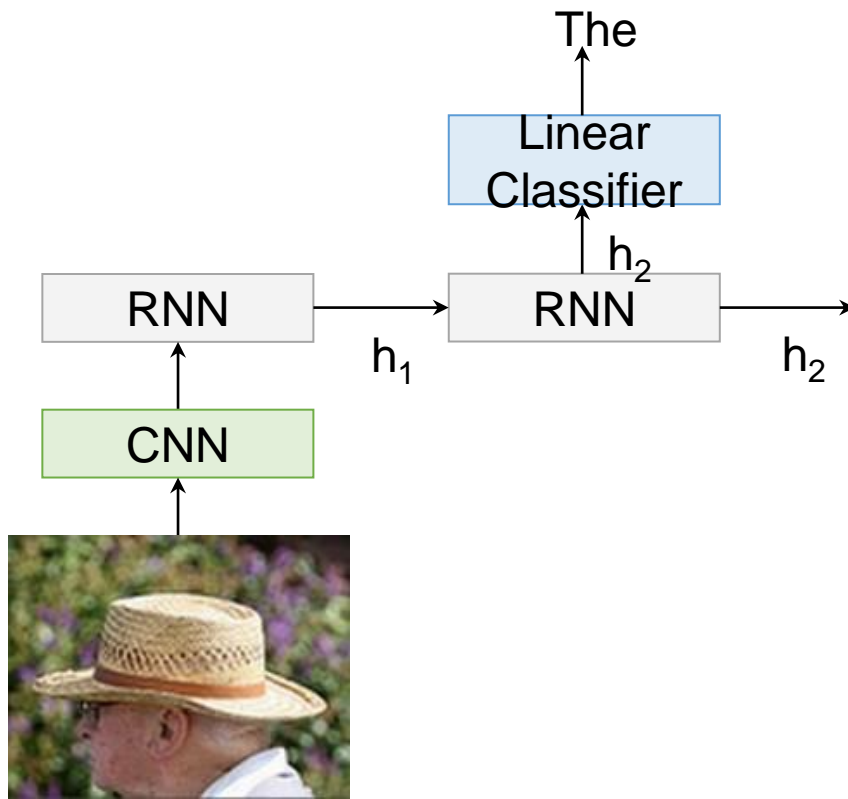
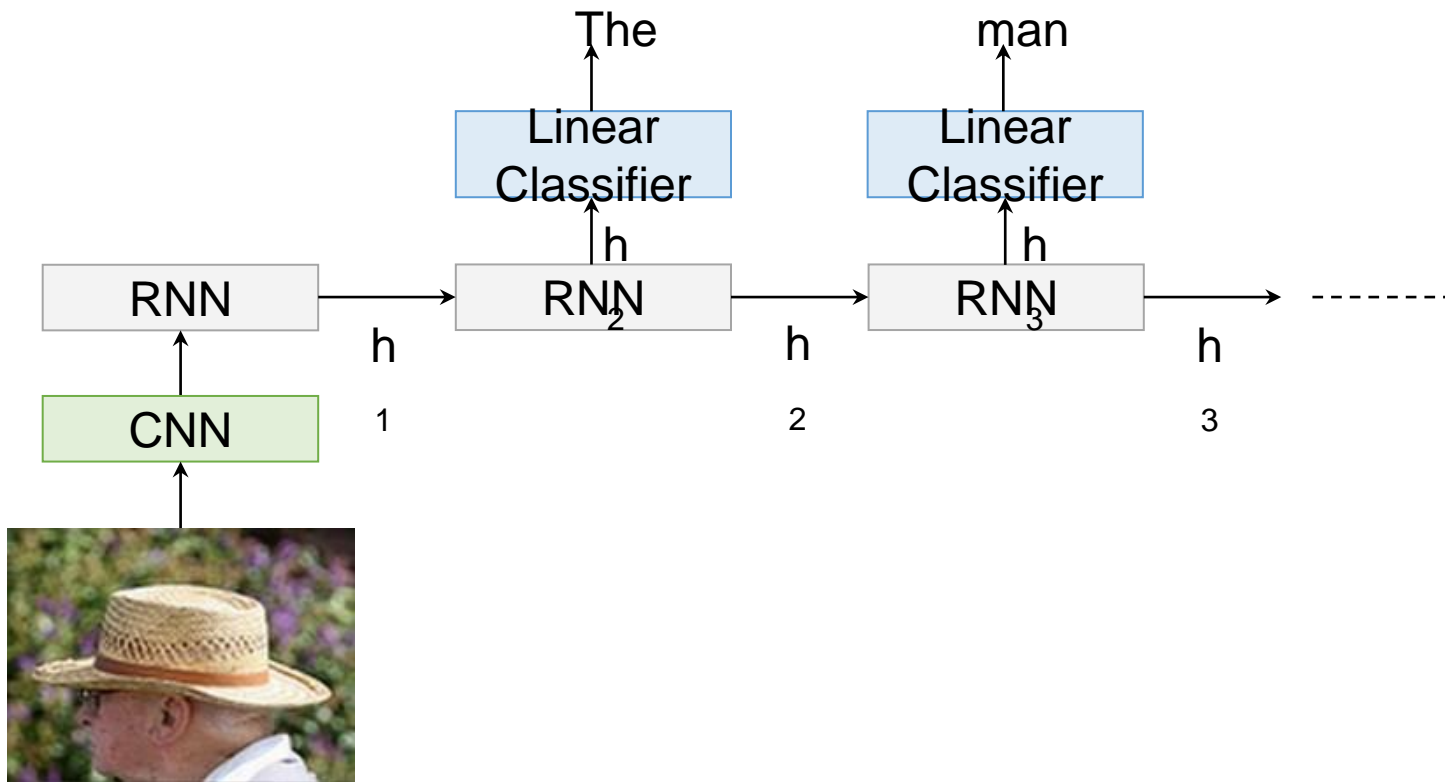


Image captioning





test image

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

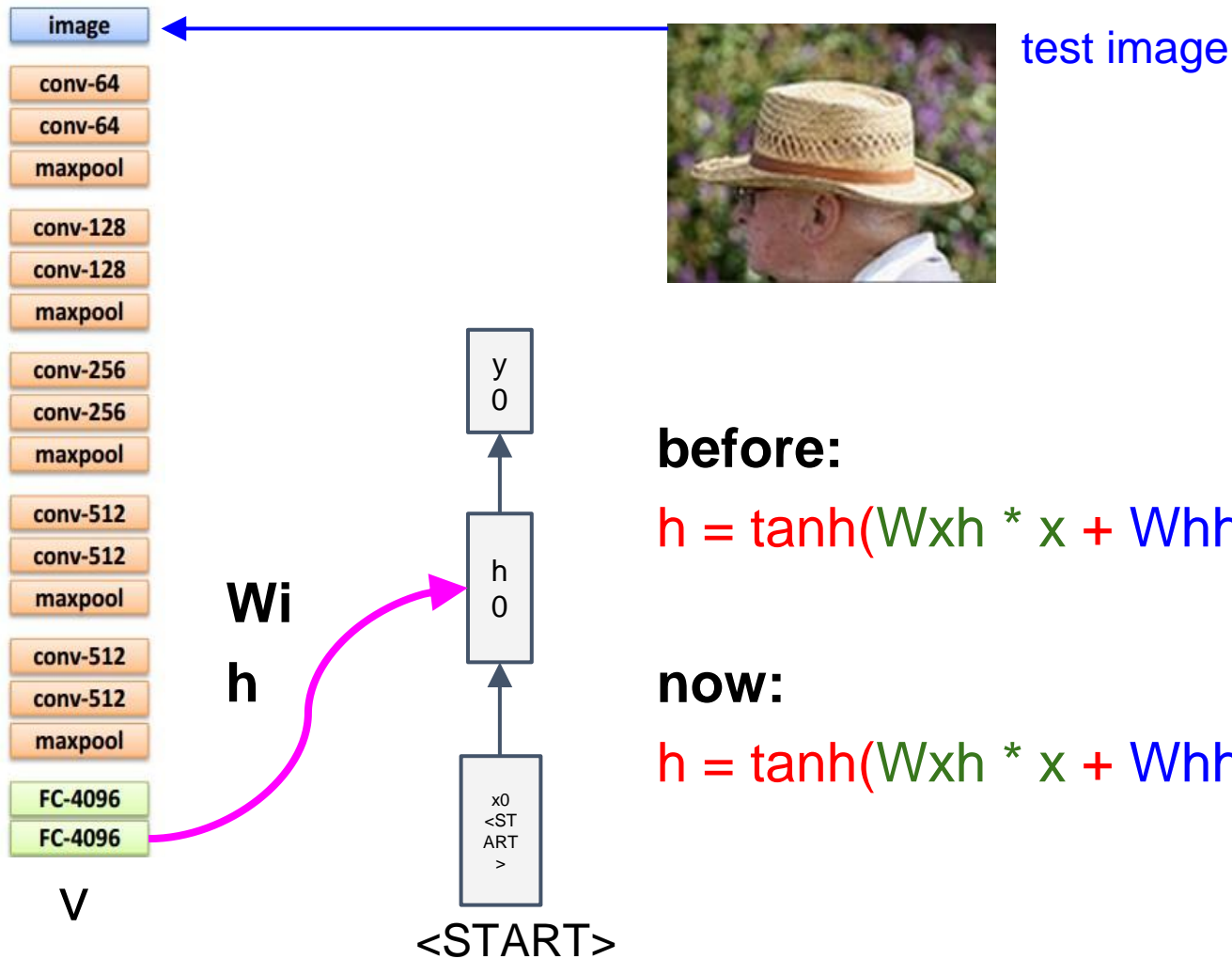
FC-4096

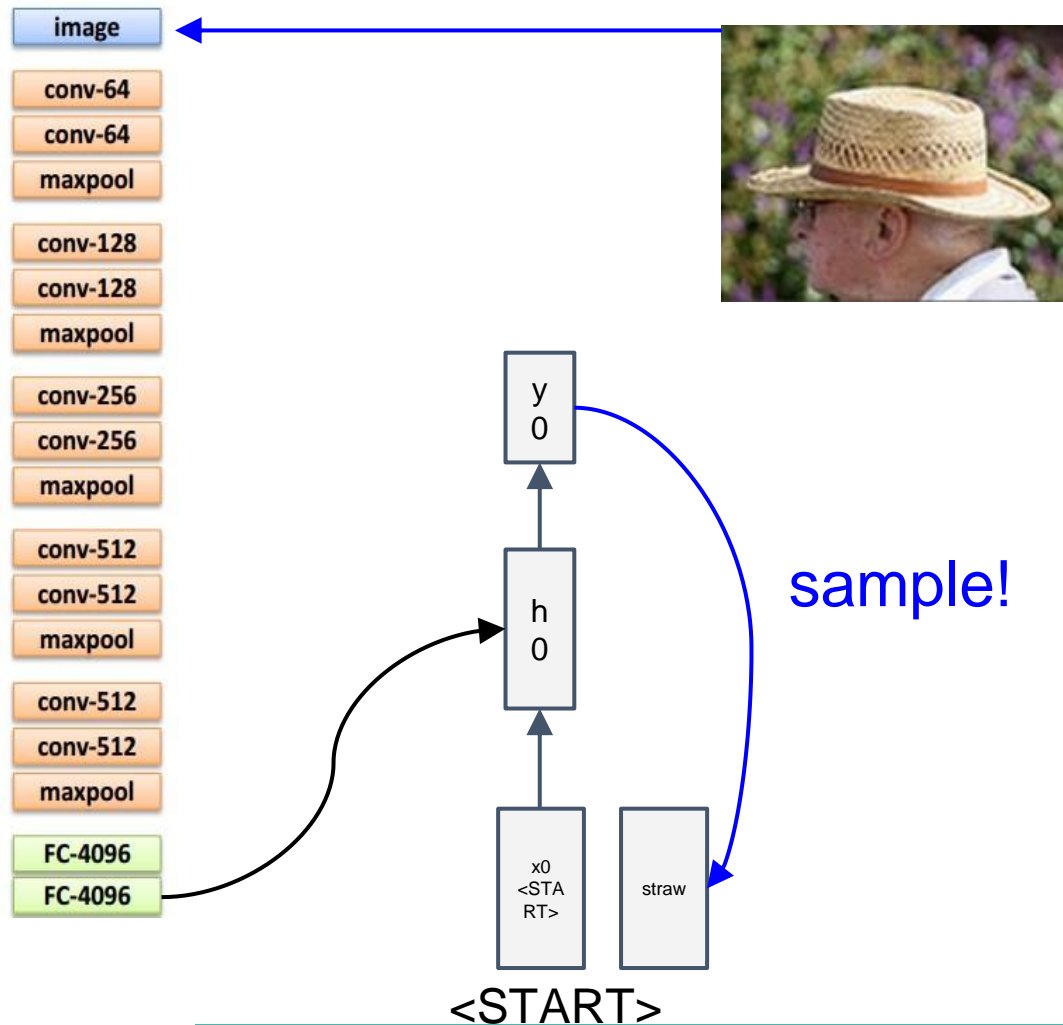


test image

x0
<ST
ART
>

<START>





test image

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

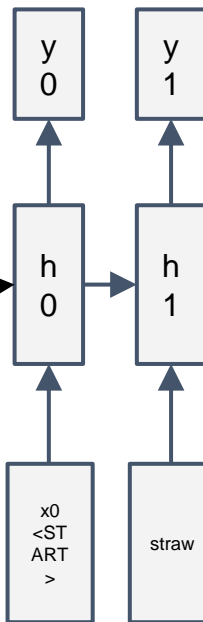
maxpool

FC-4096

FC-4096



test image



<START>

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096



test image

y
0

y
1

h
0

h
1

x0
<ST
ART
>

stra
w

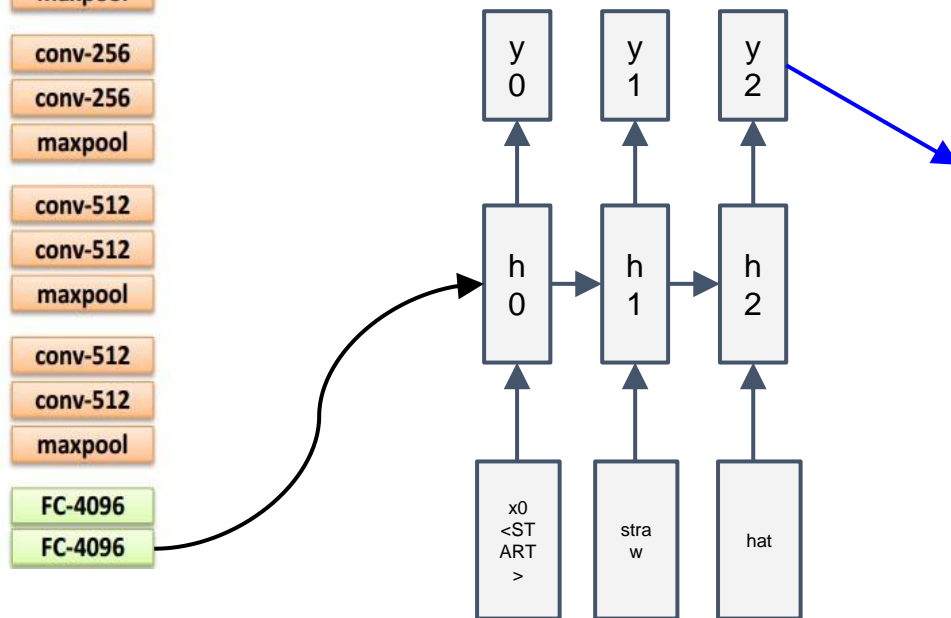
hat

sample!

<START>



test image



sample
<END> token
=> finish.

<START>

Dataset

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO

*[Tsung-Yi Lin et al.
2014]*

mscoco.org

currently:

~120K images

~5 sentences each

Prediction



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



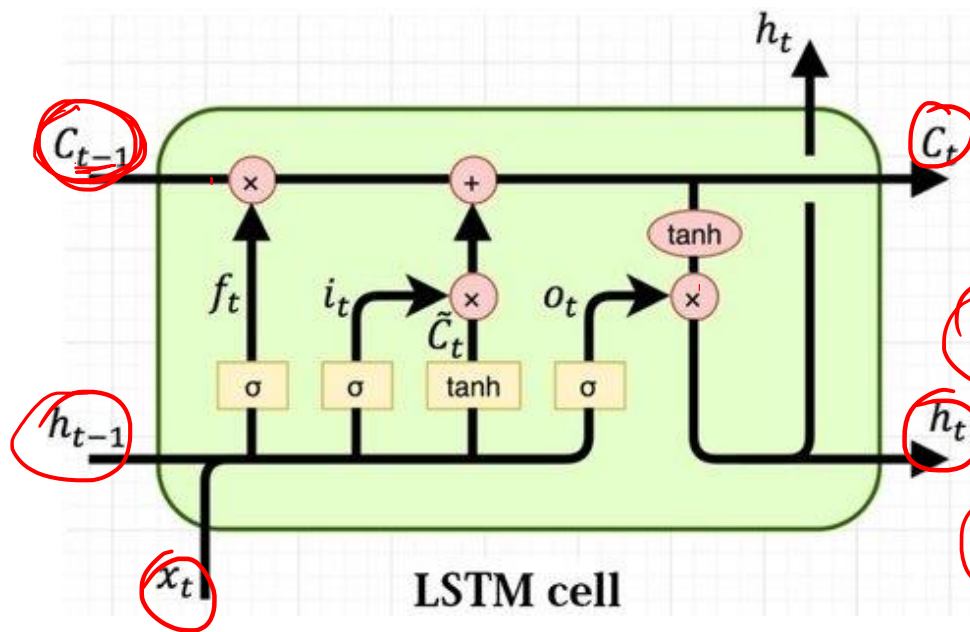
"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

LSTM

~~RNN~~ short term memory
 $0 \rightarrow 0 \rightarrow 0 \rightarrow 0 \rightarrow 0$



$$\underline{i_t} = \underline{\sigma}(x_t U^i + h_{t-1} W^i)$$

$$\underline{f_t} = \underline{\sigma}(x_t U^f + h_{t-1} W^f)$$

$$\underline{o_t} = \underline{\sigma}(x_t U^o + h_{t-1} W^o)$$

$$\underline{\tilde{C}_t} = \tanh(x_t U^g + h_{t-1} W^g)$$

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

$$\underline{h_t} = \tanh(C_t) * \underline{o_t}$$

Q&A

