

# Recurrent Neural Networks

M. Soleymani  
Sharif University of Technology  
Spring 2025

Most slides have been adopted from Fei Fei Li and colleagues lectures, cs231n, Stanford  
and some slides from Bhiksha Raj, 11-785, CMU

# Sequences in the wild

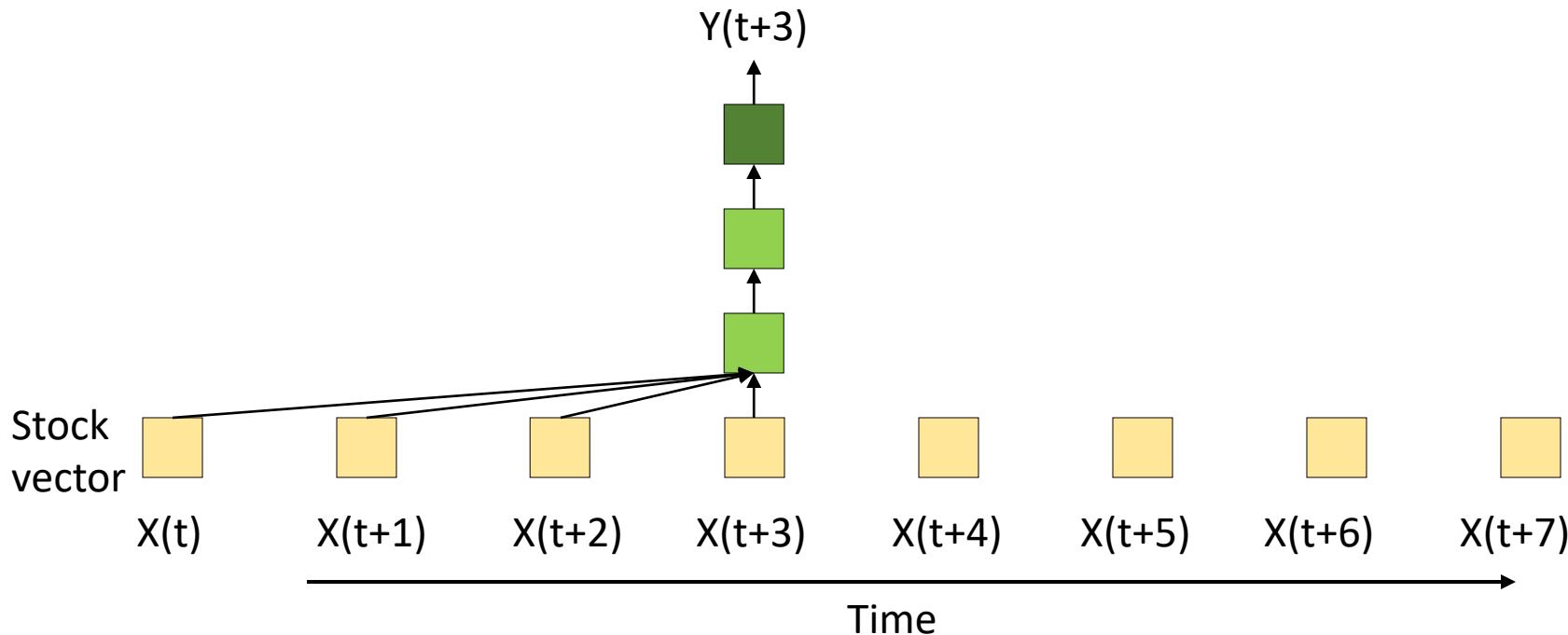


Source: Ava Amini & Alex Amini, Deep learning course, MIT

# Modelling Series

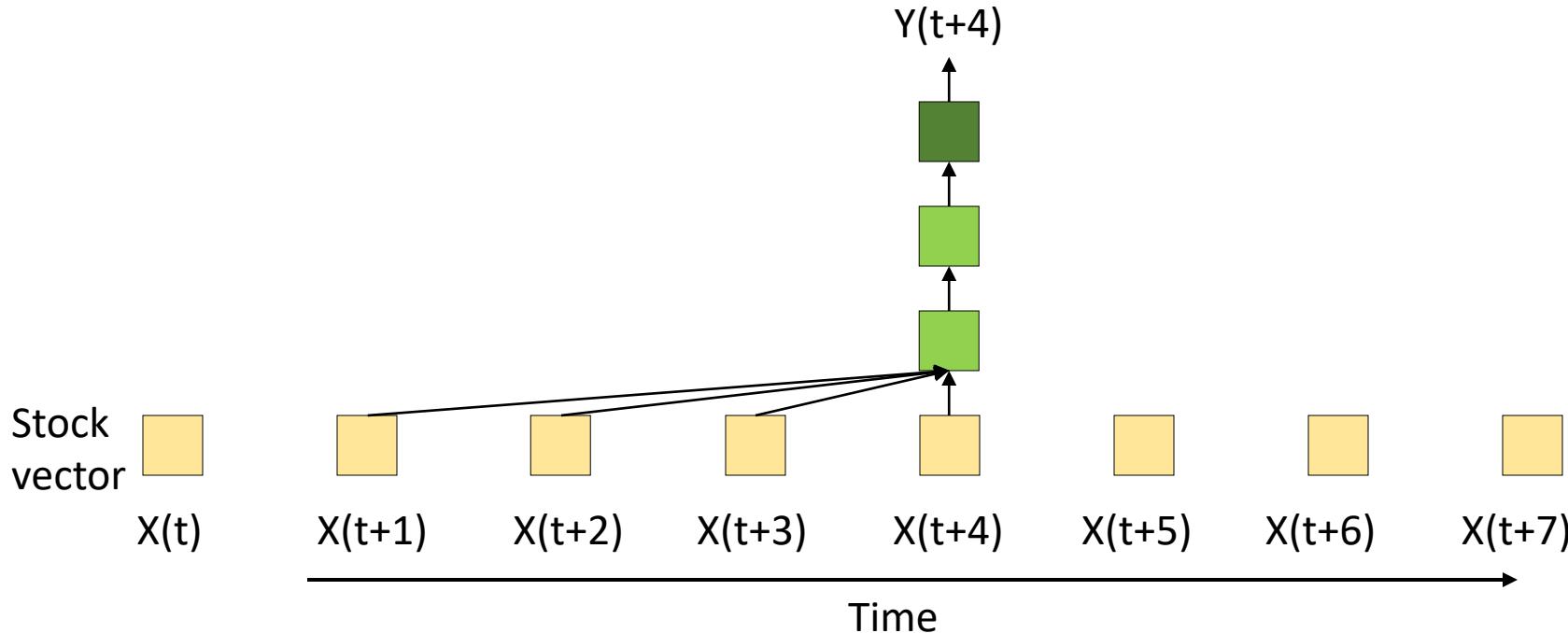
- In many situations one must consider a *series* of inputs to produce an output
  - Outputs too may be a series
- Examples: ..

# The stock predictor



- The sliding predictor
  - Look at the last few days
  - This is just a convolutional neural net applied to series data
    - Also called a *Time-Delay Neural Network (TDNN)*

# The stock predictor



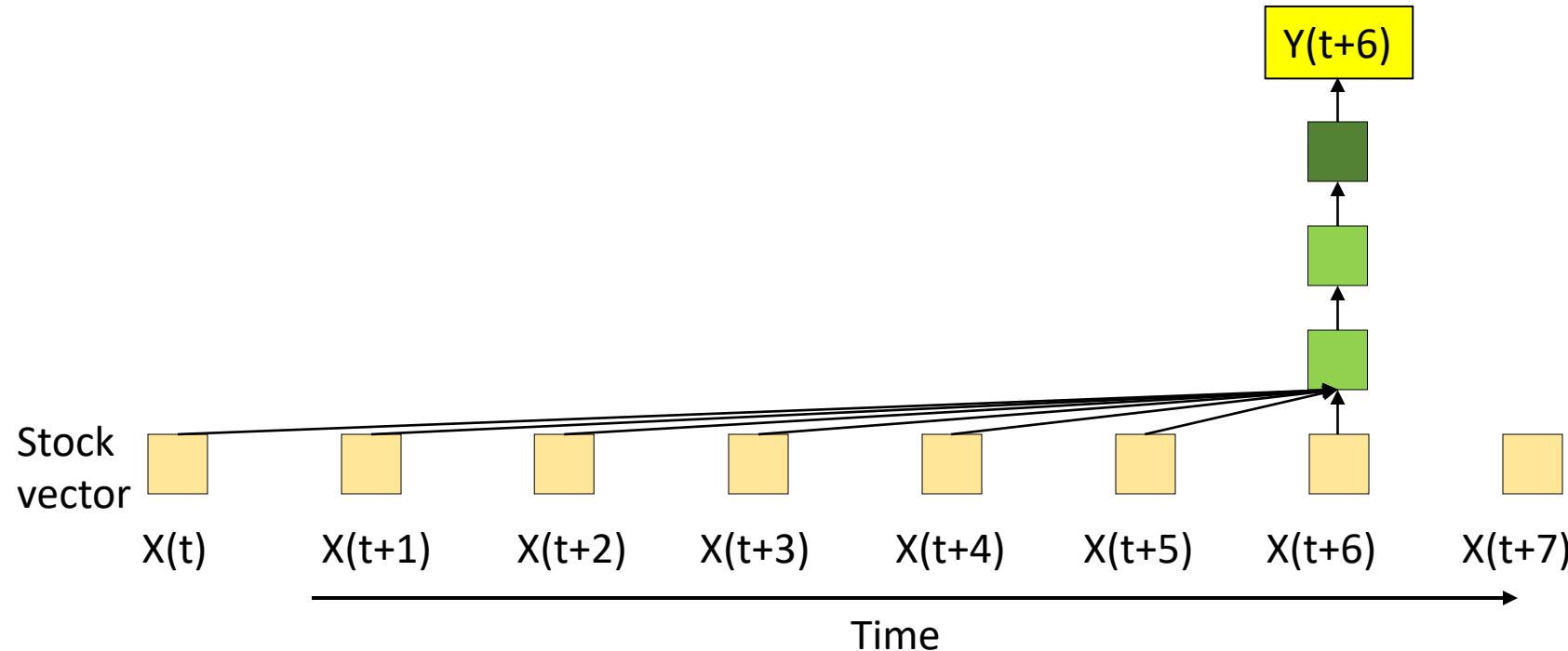
- The sliding predictor
  - Look at the last few days
  - This is just a convolutional neural net applied to series data
    - Also called a *Time-Delay Neural Network (TDNN)*

# Finite-response model

- This is a *finite response* system
  - Something that happens *today* only affects the output of the system for  $N$  days into the future
    - $N$  is the *width* of the system

$$Y_t = f(X_t, X_{t-1}, \dots, X_{t-N})$$

# Finite-response



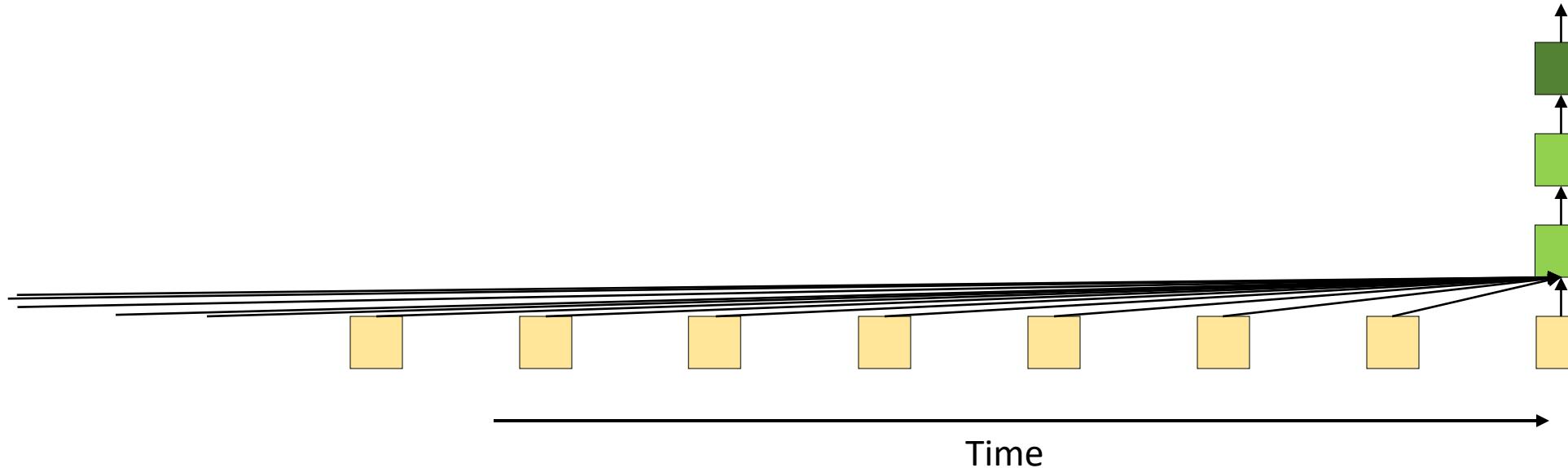
- Problem: Increasing the “history” makes the network more complex
  - No worries, we have the CPU and memory
    - Or do we?

# Systems often have long-term dependencies



- Longer-term trends
  - Weekly trends in the market
  - Monthly trends in the market
  - Annual trends
  - Though longer historic tends to affect us less than more recent events..

# We want *infinite* memory



- Required: *Infinite* response systems
  - What happens today can continue to affect the output forever
    - Possibly with weaker and weaker influence

$$Y_t = f(X_t, X_{t-1}, \dots, X_{t-\infty})$$

# An alternate model for infinite response: state-space model

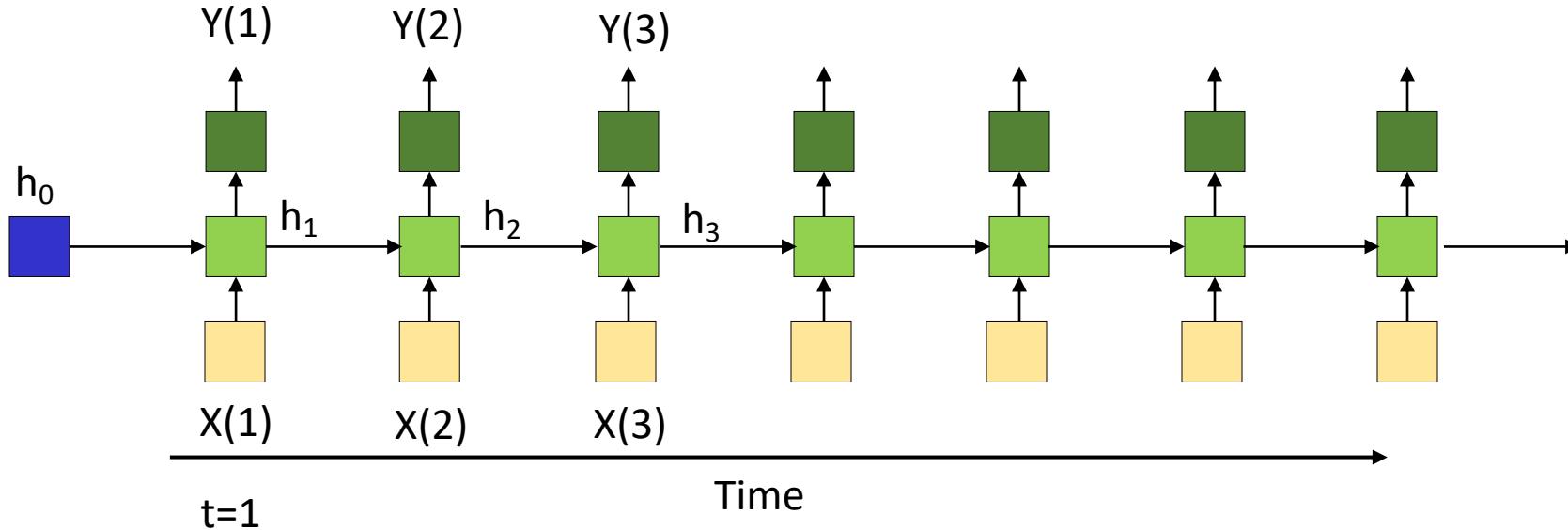
- State-space model:

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

$$y_t = g(h_t)$$

- $h_t$  is the *state* of the network
  - Model directly embeds the memory in the state
- Need to define initial state  $h_0$
- This is a *fully recurrent* neural network
  - Or simply a *recurrent neural network*
- *State* summarizes information about the entire past

# The simple state-space model



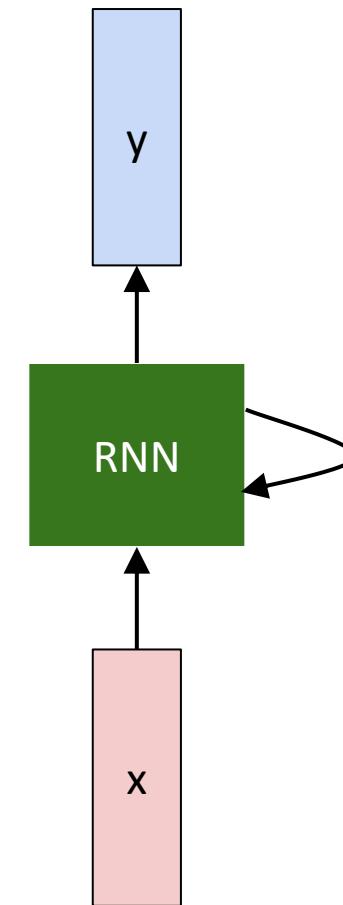
- The state (green) at any time is determined by the input at that time, and the state at the previous time
- Also known as a recurrent neural net

# Recurrent Neural Network

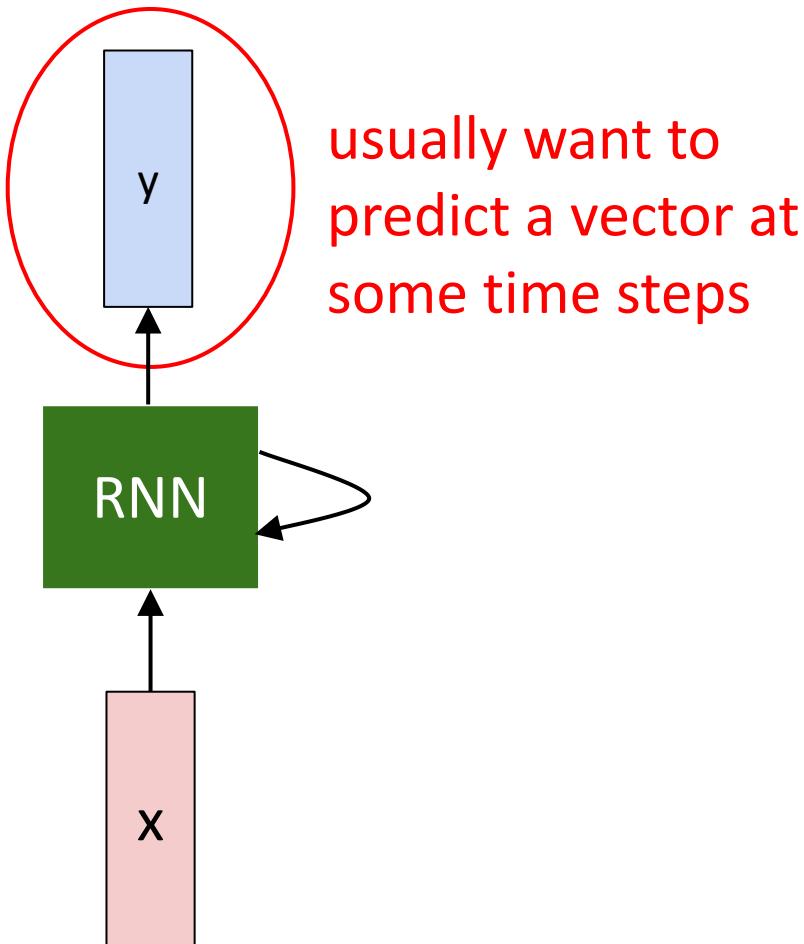
We can process a sequence of vectors  $\mathbf{x}$  by applying a recurrence formula at every time step:

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

new state      old state      input vector at some time step  
some function with parameters W



# Recurrent Neural Network

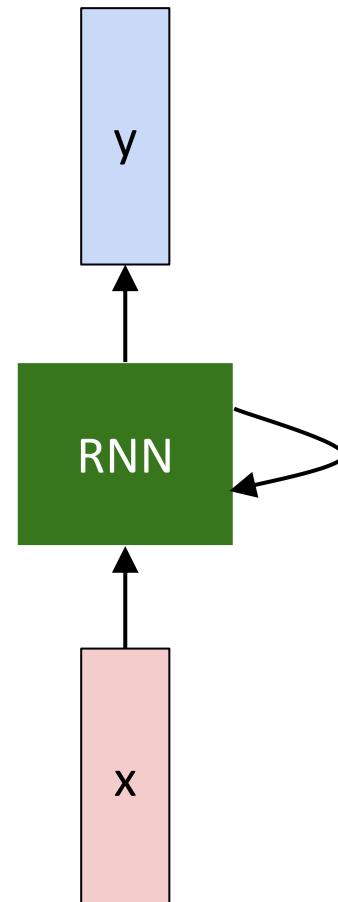


# Recurrent Neural Network

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

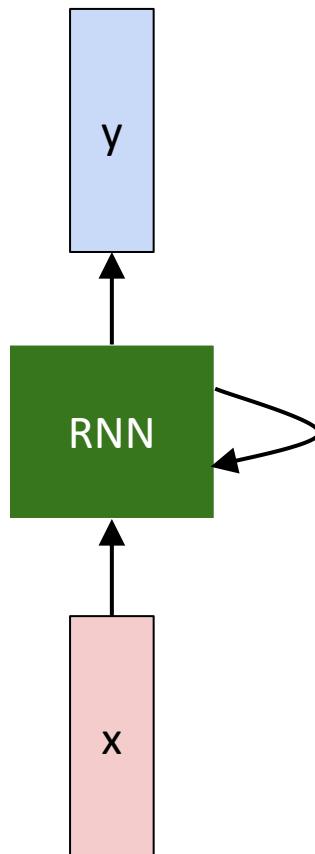
$$h_t = f_W(h_{t-1}, x_t)$$

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



# (Vanilla) Recurrent Neural Network

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



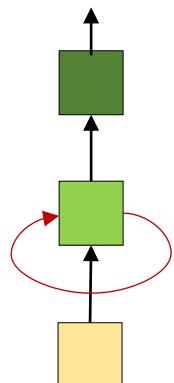
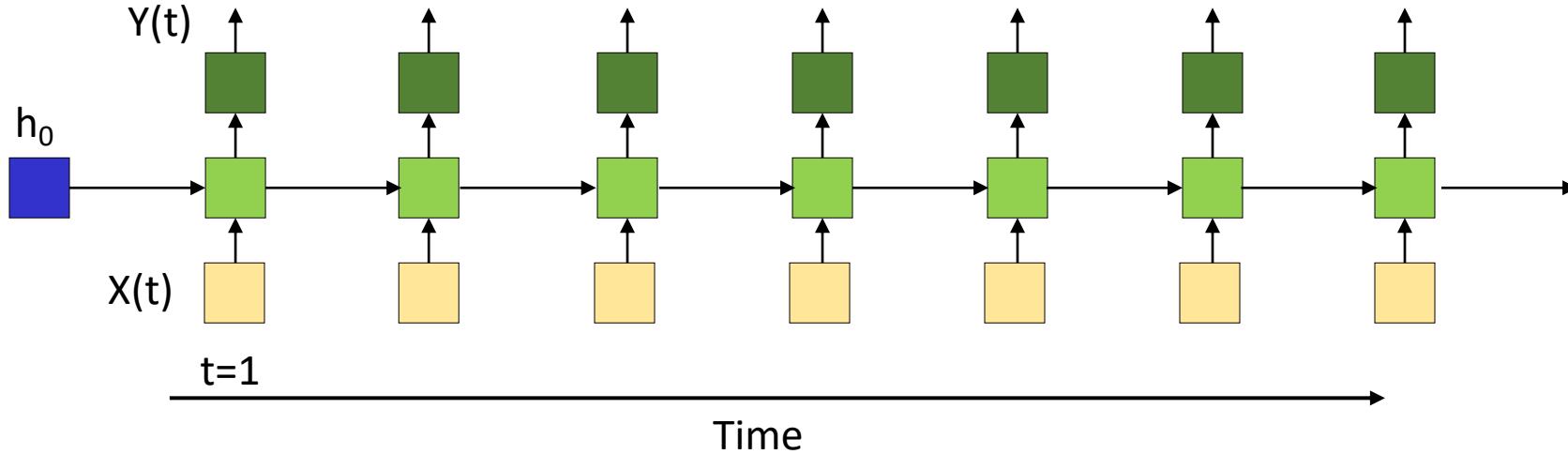
$$h_t = f_W(h_{t-1}, x_t)$$



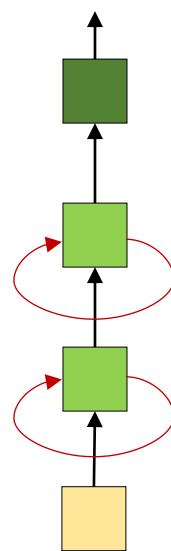
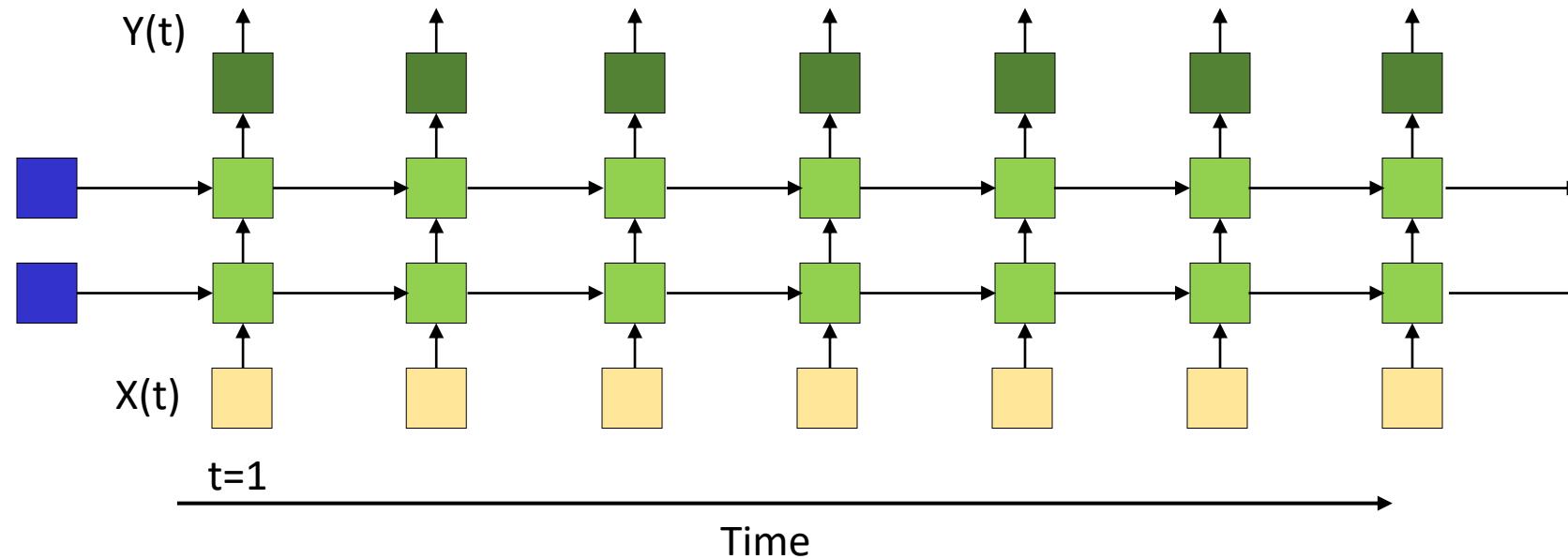
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

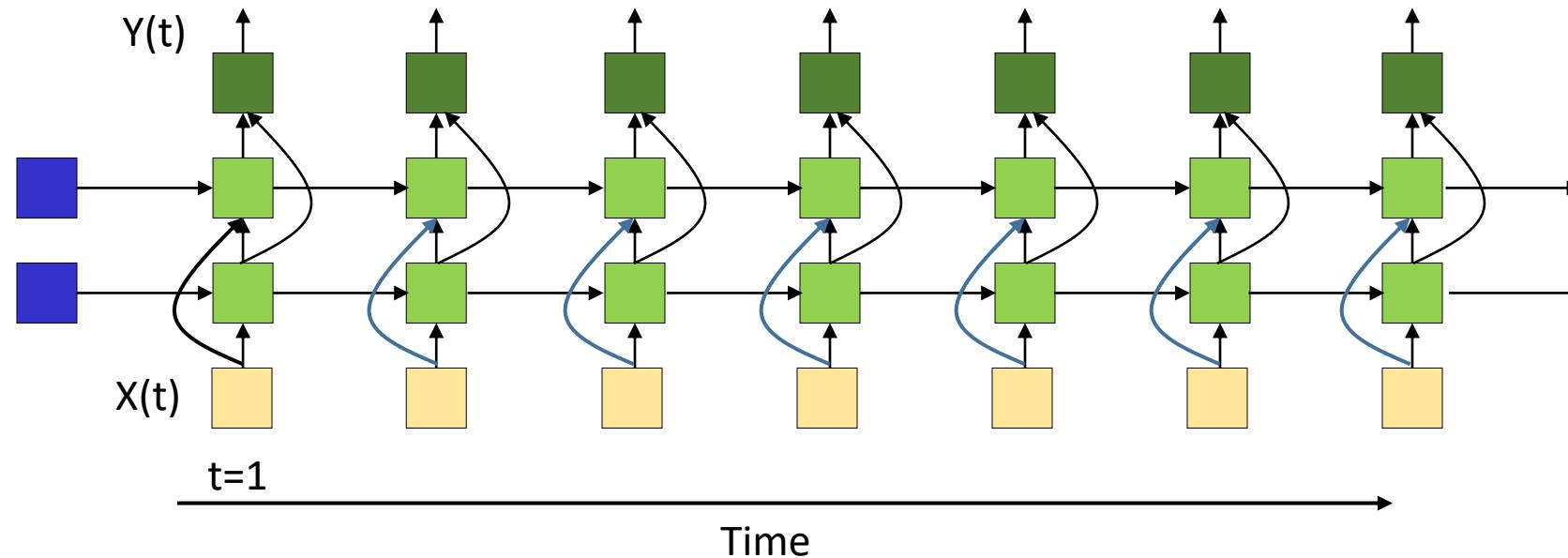
# Single hidden layer RNN



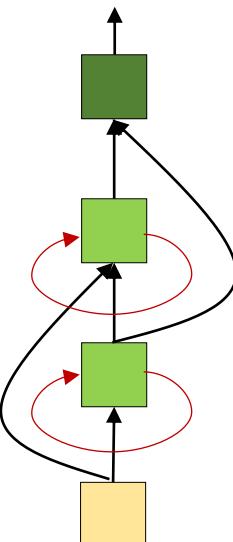
# Multiple recurrent layer RNN



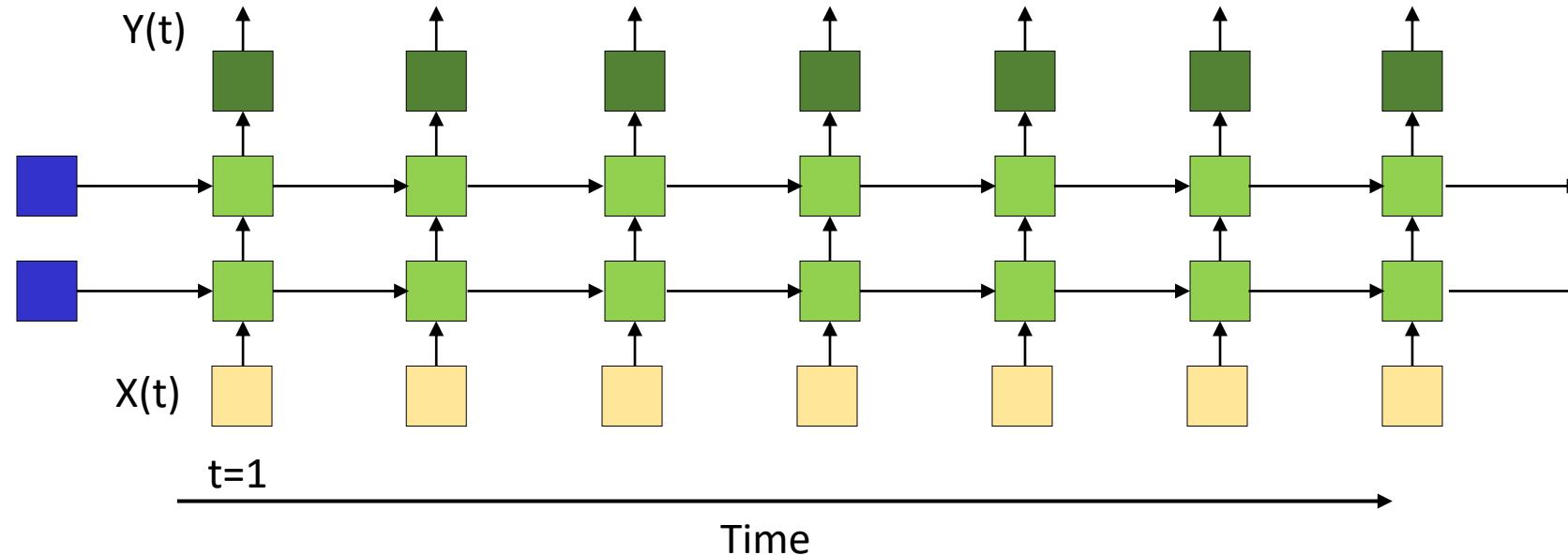
# Multiple recurrent layer RNN



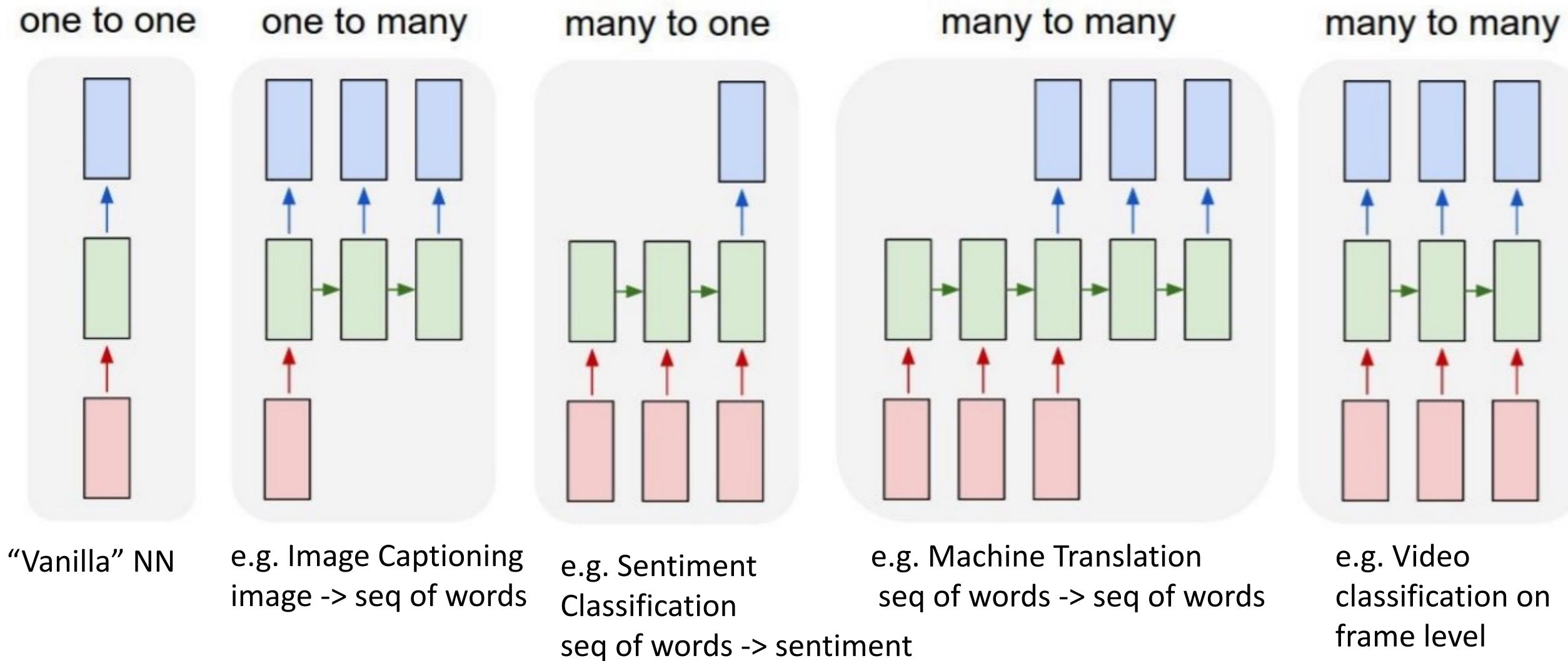
- We can also have skips..



# The simplest structures are most popular

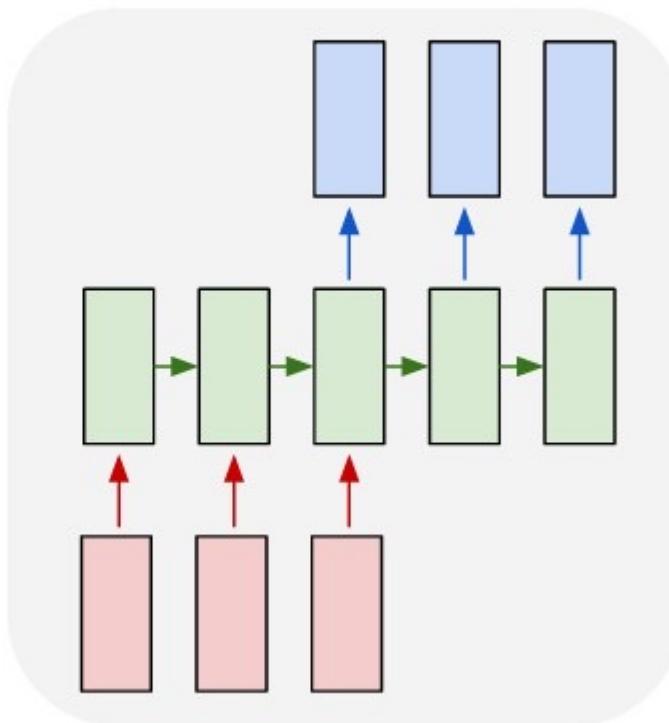


# Recurrent Neural Networks: Process Sequences

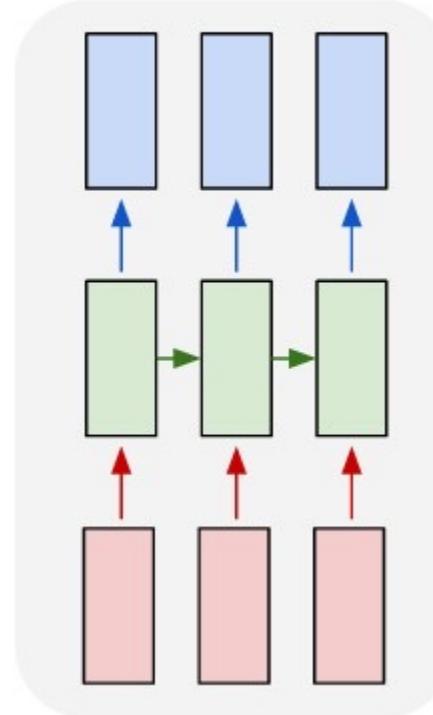


# Variants

many to many



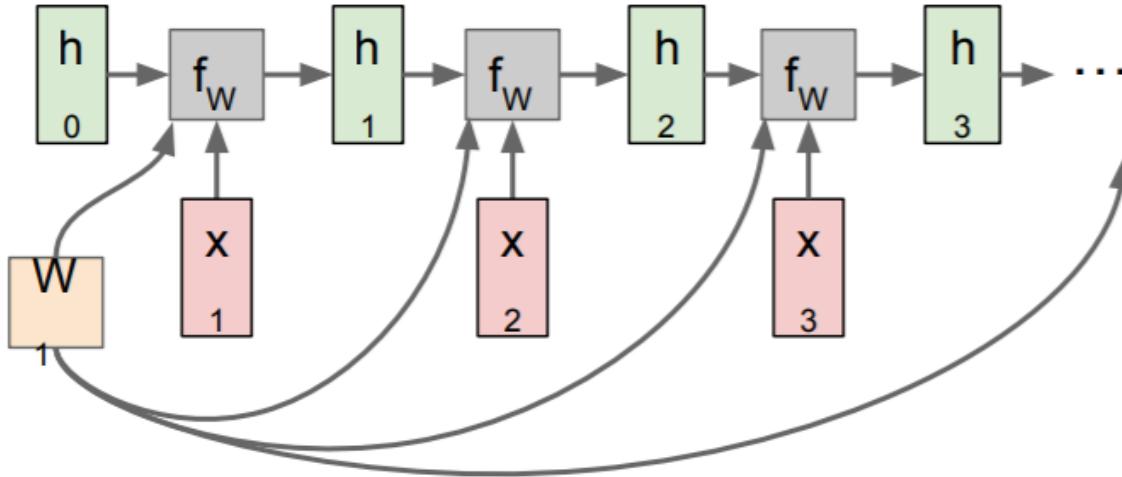
many to many



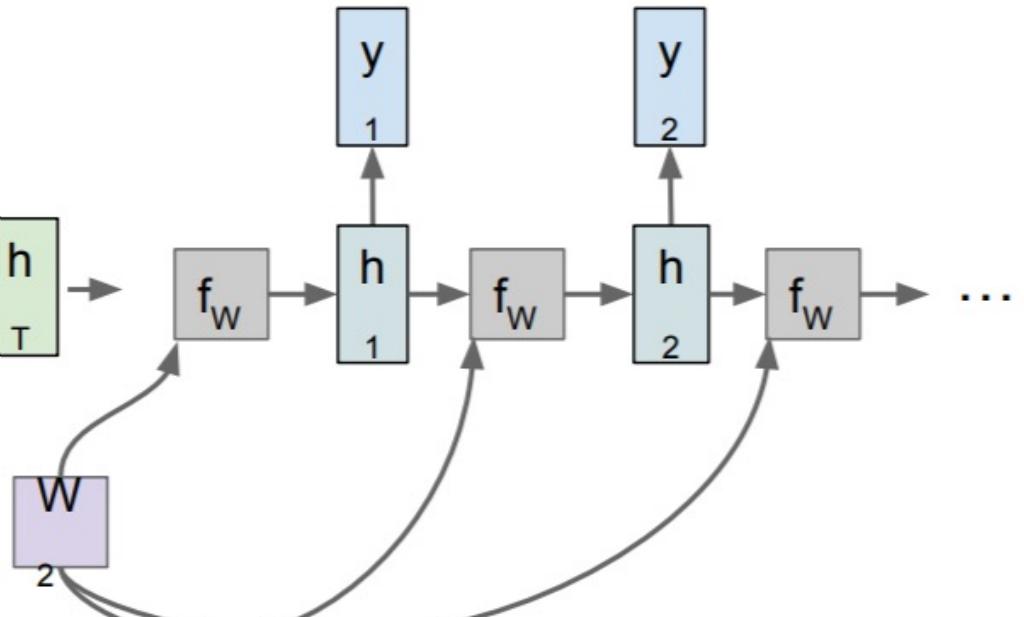
- 1: *Delayed* sequence to sequence
- 2: Sequence to sequence, e.g. stock problem, label prediction
- Etc...

# Sequence to Sequence: Many-to-one + one-to-many

**Many to one:** Encode input sequence in a single vector



**One to many:** Produce output sequence from single input vector

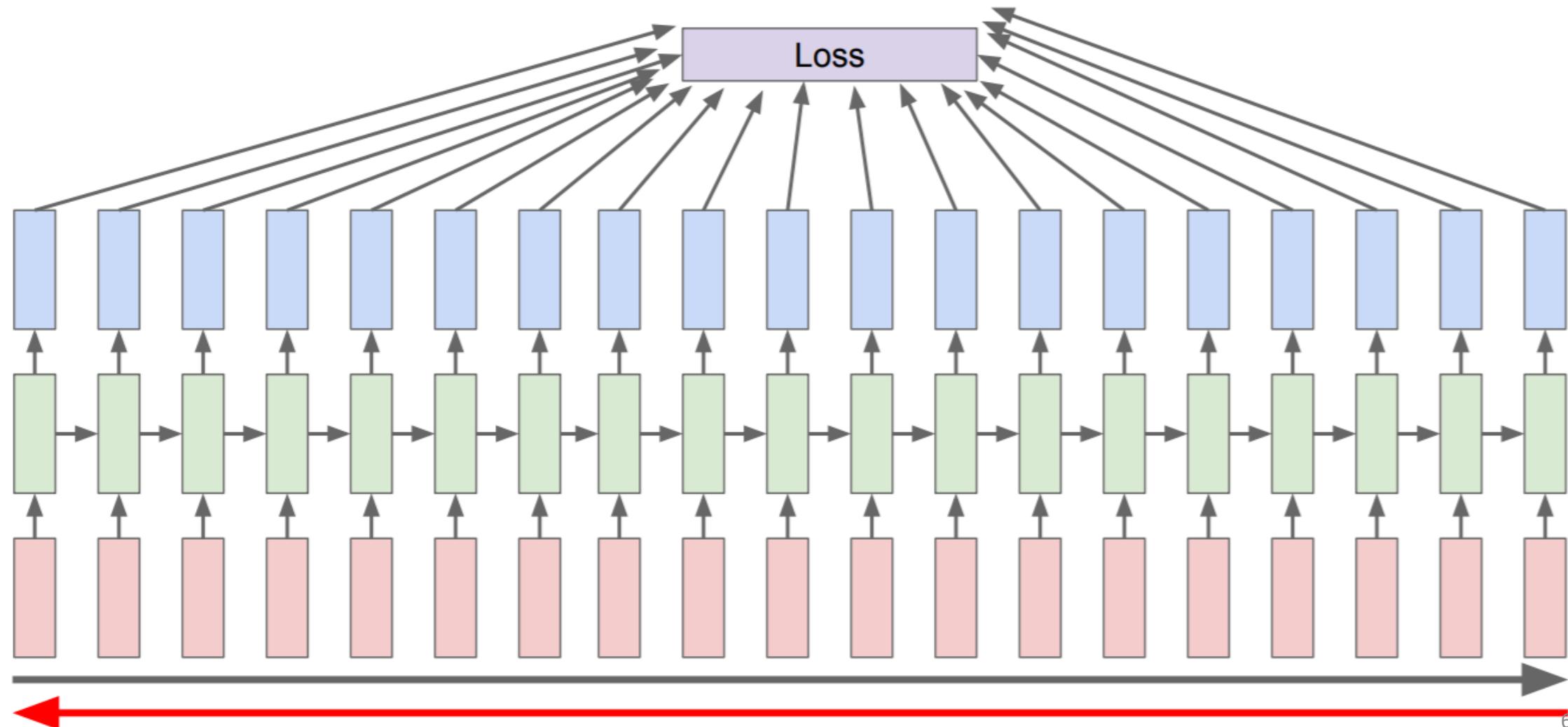


# Story so far

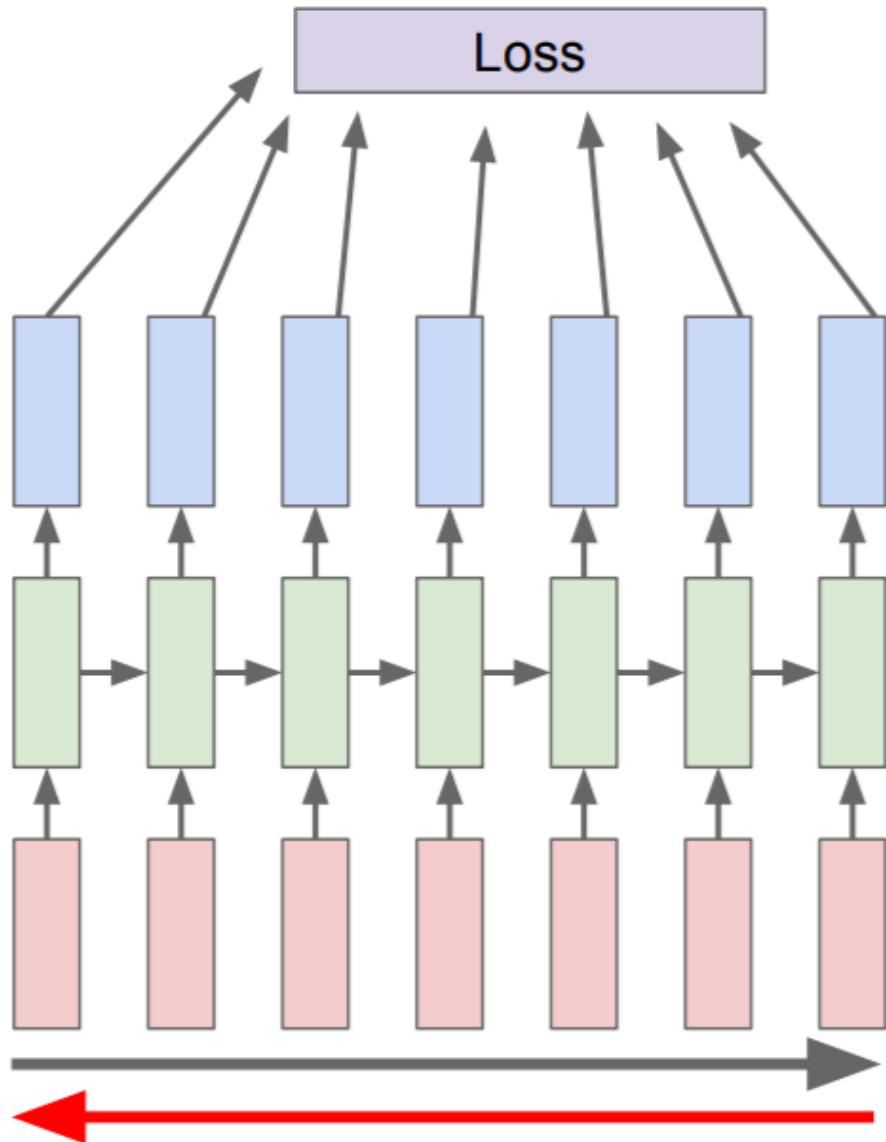
- Time series analysis must consider past inputs along with current input
- Looking into the infinite past requires recursion
- State-space models retain information about the past through recurrent hidden states
  - These are “**fully recurrent**” networks
  - The initial values of the hidden states are generally learnable parameters as well

# Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

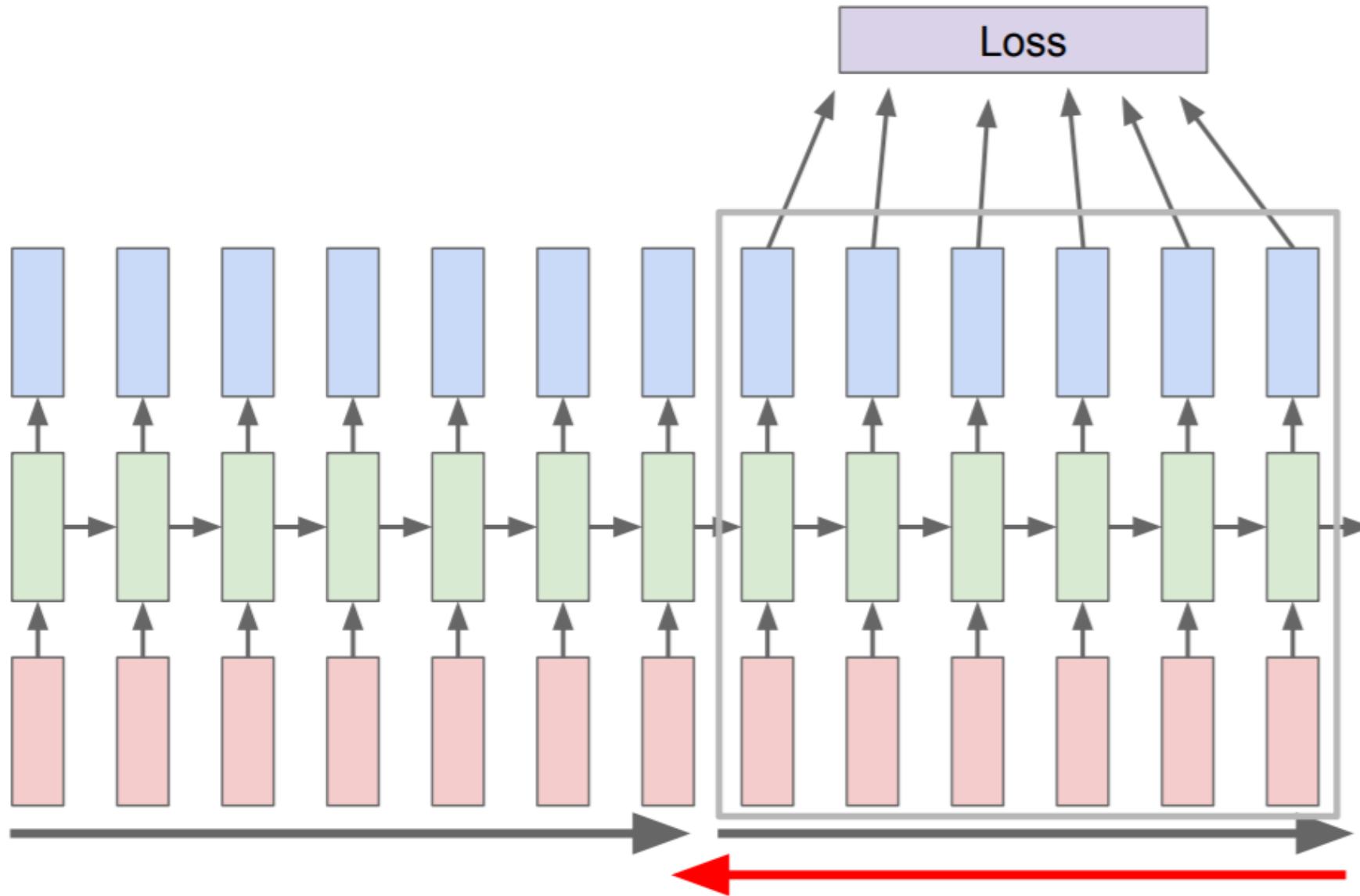


# Truncated Backpropagation through time



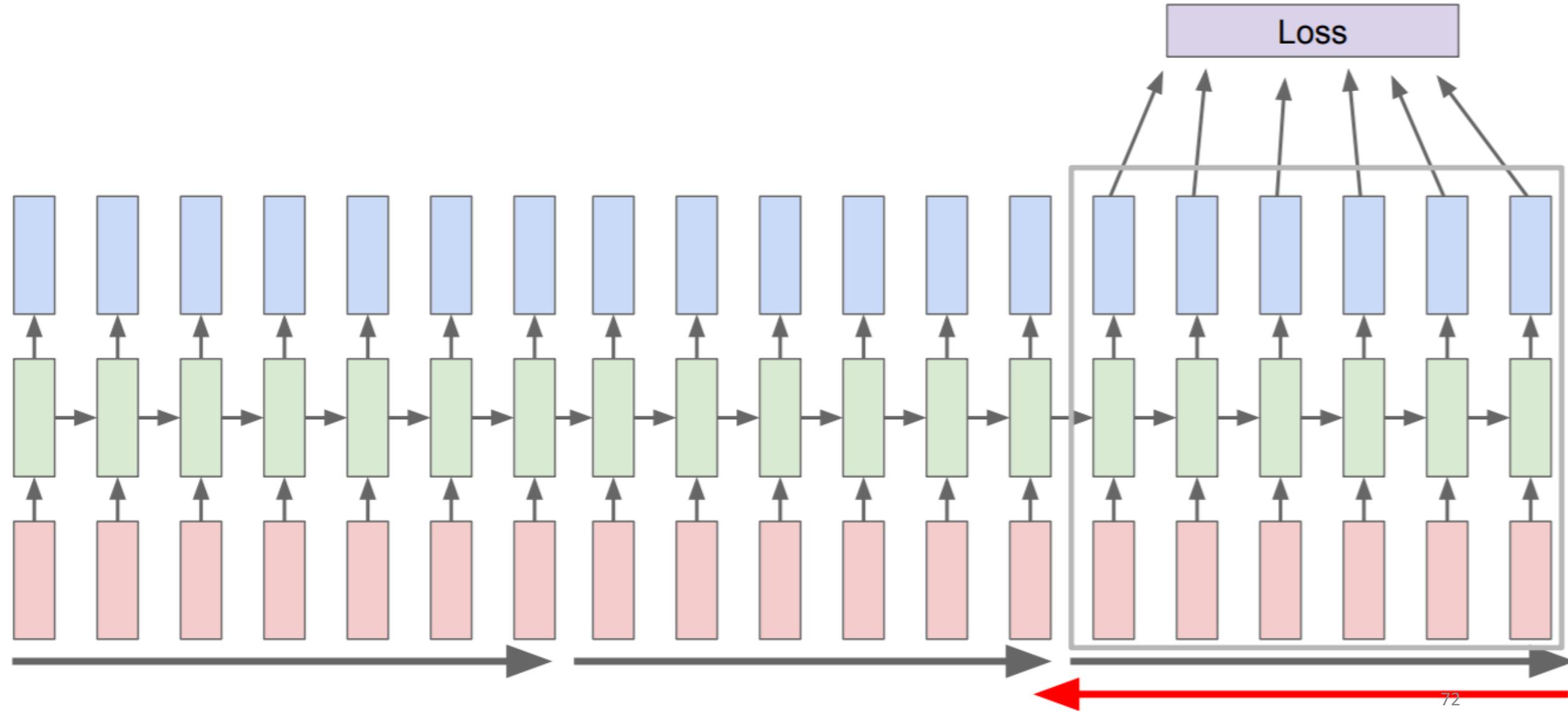
Run forward and backward through  
chunks of the sequence instead of whole  
sequence

# Truncated Backpropagation through time

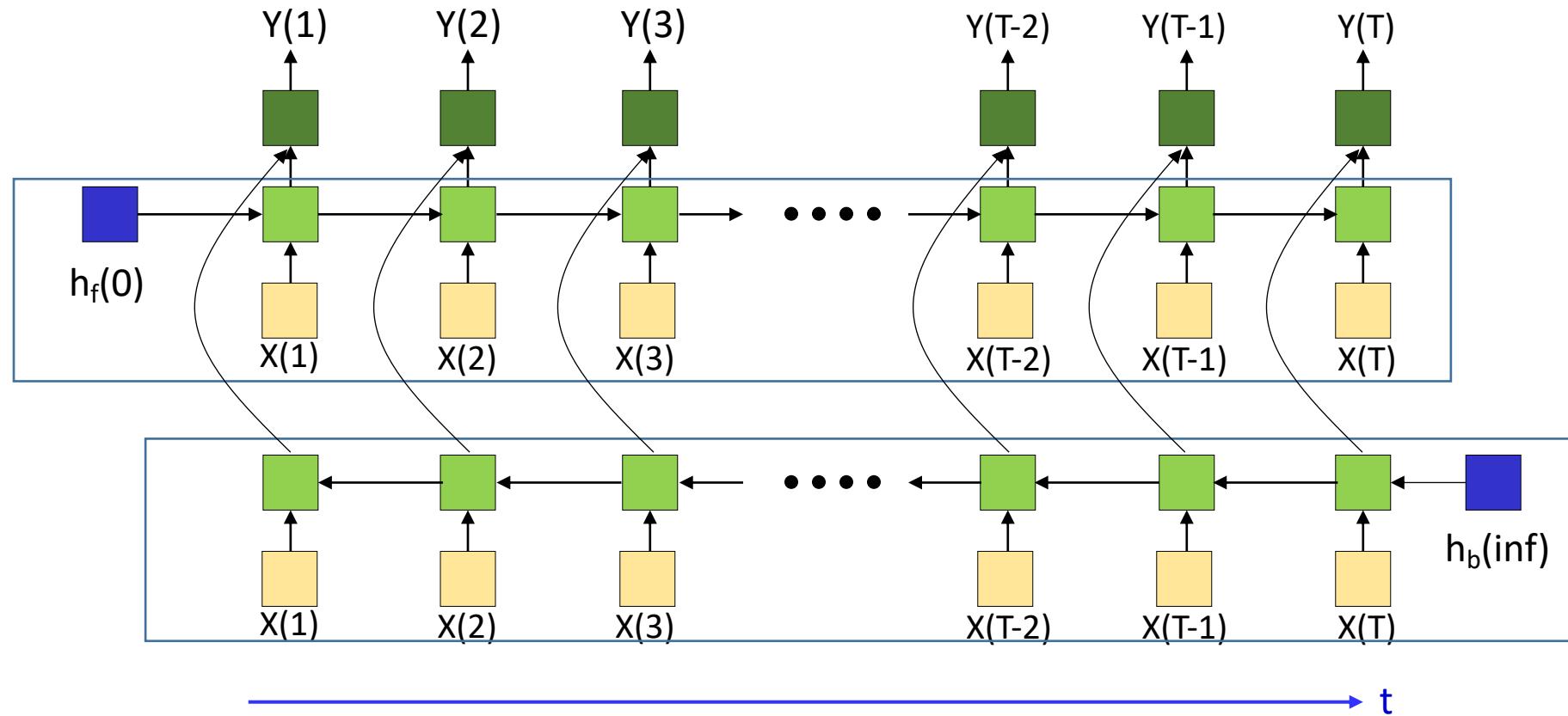


Carry hidden states  
forward in time forever,  
but only backpropagate  
for some smaller  
number of steps

# Truncated Backpropagation through time

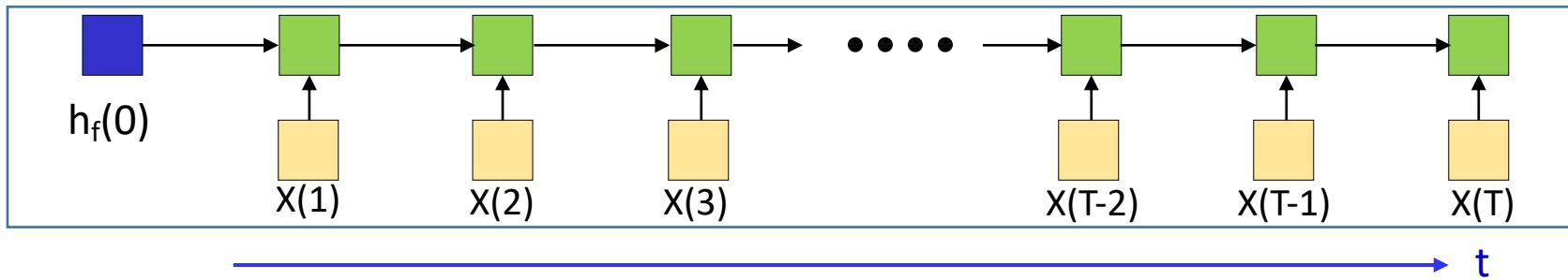


# Bidirectional RNN



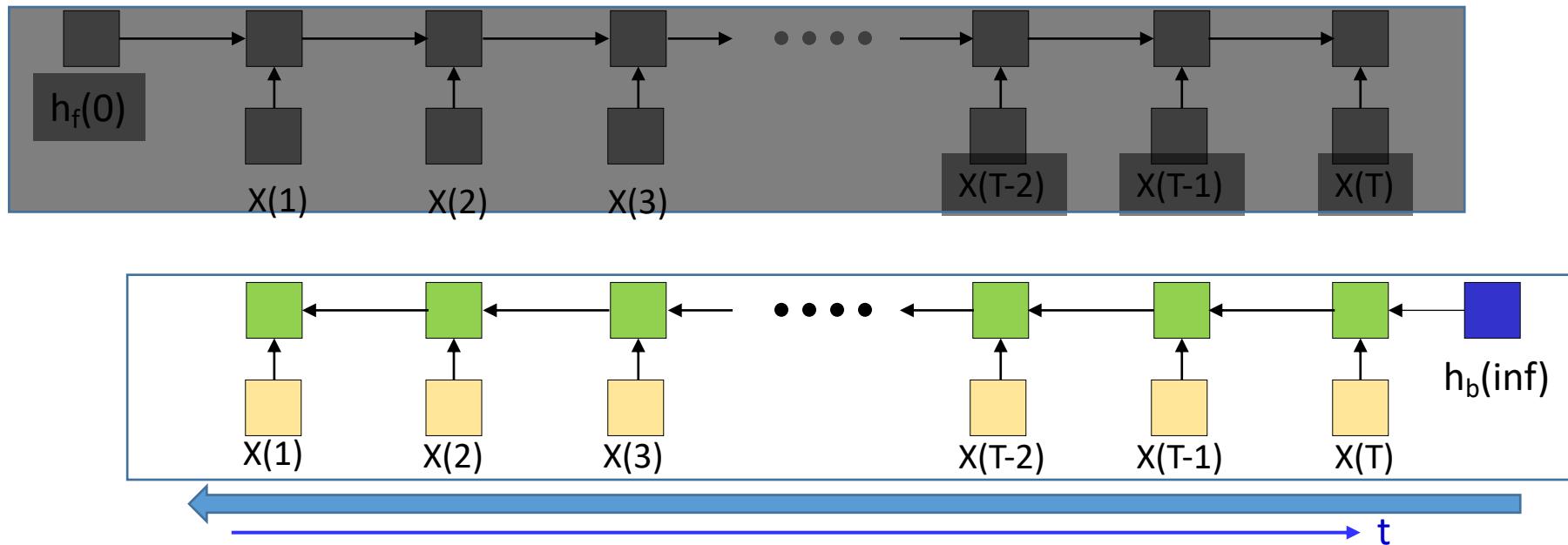
- A forward net process the data from  $t=1$  to  $t=T$
- A backward net processes it backward from  $t=T$  down to  $t=1$

# Bidirectional RNN: Processing an input string



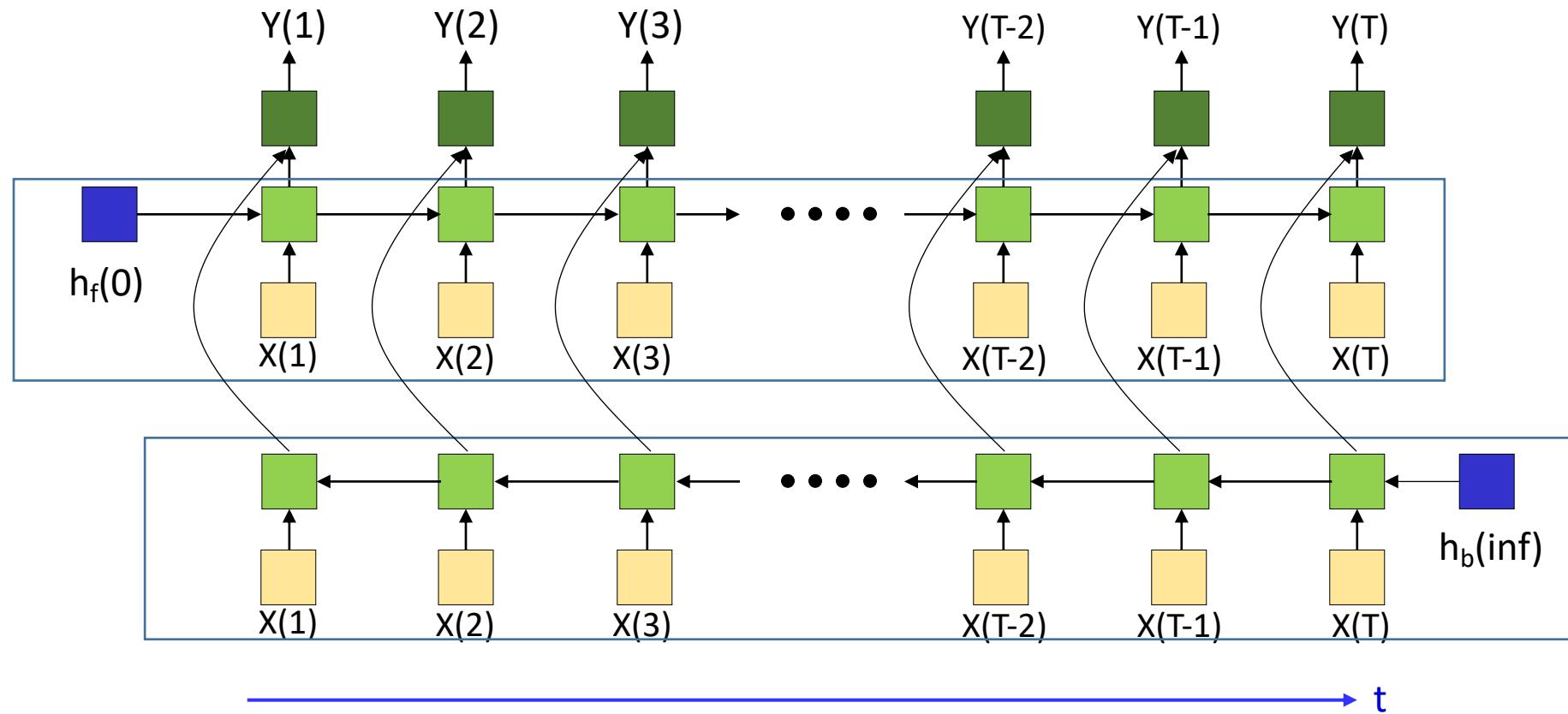
- The forward net process the data from  $t=1$  to  $t=T$ 
  - Only computing the hidden states, initially

# Bidirectional RNN: Processing an input string



- The backward nets processes the input data in *reverse time*, end to beginning
  - Initially only the hidden state values are computed
    - Clearly, this is not an online process and requires the *entire* input data
  - Note: *This is not the backward pass of backprop.*

# Bidirectional RNN: Processing an input string



- The computed states of both networks are used to compute the final output at each time

# Story so far

- RNNs for sequence modeling
  - Handle inputs with variable lengths
  - In theory, can track long term history
  - Share parameters across the sequence
  - Consider the order of inputs

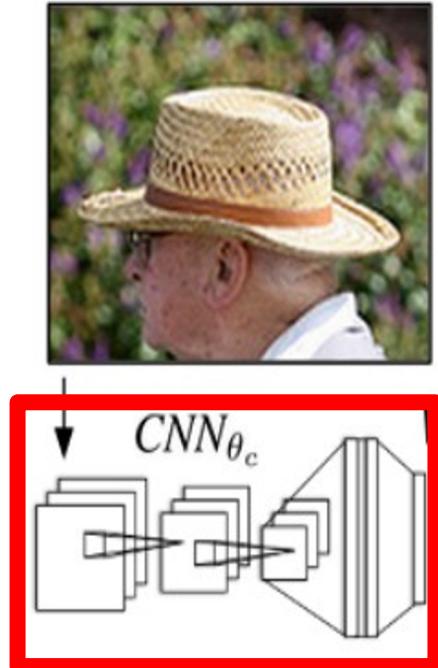
# RNNs..

- Excellent models for time-series analysis tasks
  - Time-series prediction
  - Time-series classification
  - Sequence prediction..

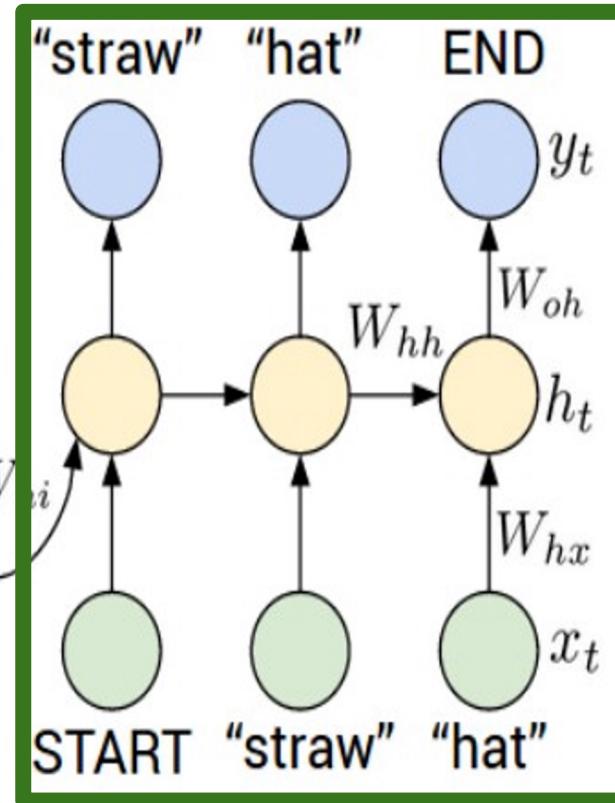
# RNN tradeoffs

- RNN Advantages:
  - Can process any length input
  - Computation for step  $t$  can (in theory) use information from many steps back
  - Model size doesn't increase for longer input
  - Same weights applied on every timestep, so there is symmetry in how inputs are processed.
- RNN Disadvantages:
  - Recurrent computation is slow
  - In practice, difficult to access information from many steps back

# Image Captioning



## Recurrent Neural Network



## Convolutional Neural Network

test image



image



test 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

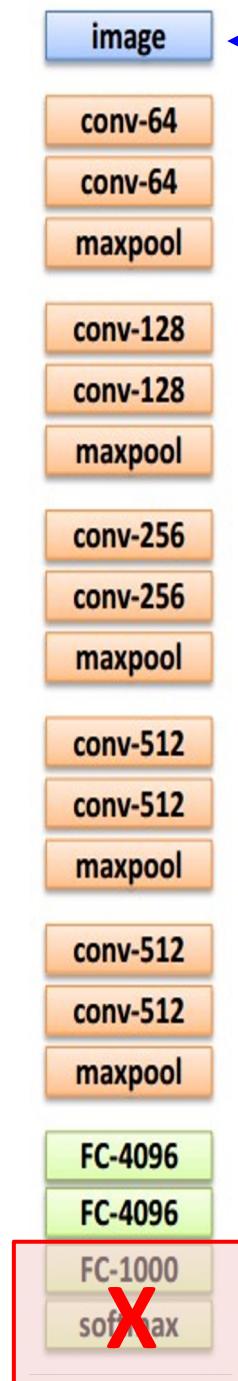
maxpool

FC-4096

FC-4096

FC-1000

softmax



test image

image



test 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



<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

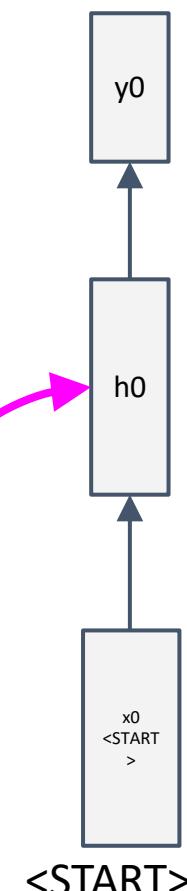
FC-4096

FC-4096

v



test image



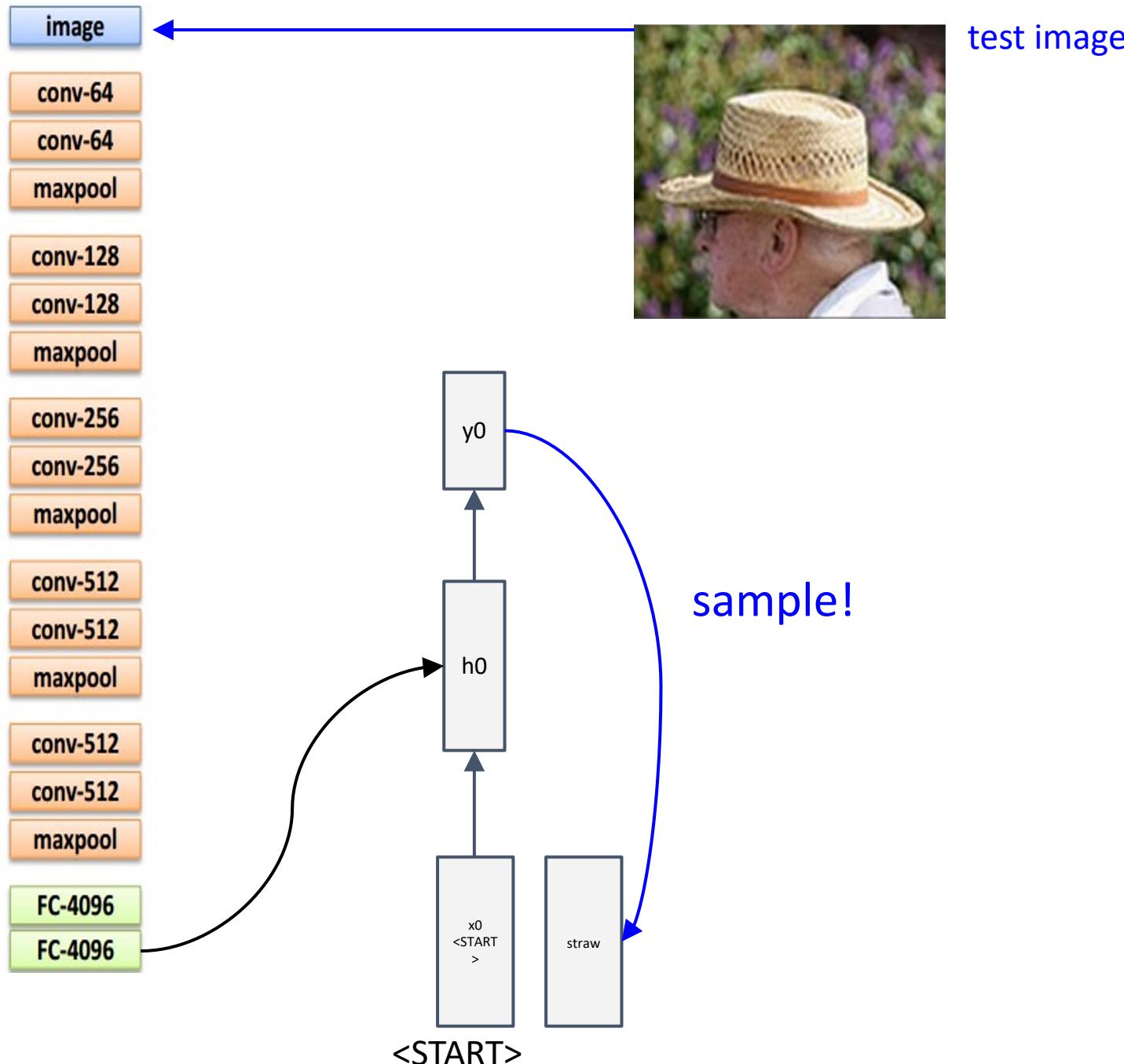
Why didn't we use "v" as the  $h_0$ ?

**before:**

$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

**now:**

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * v)$$



image



test image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

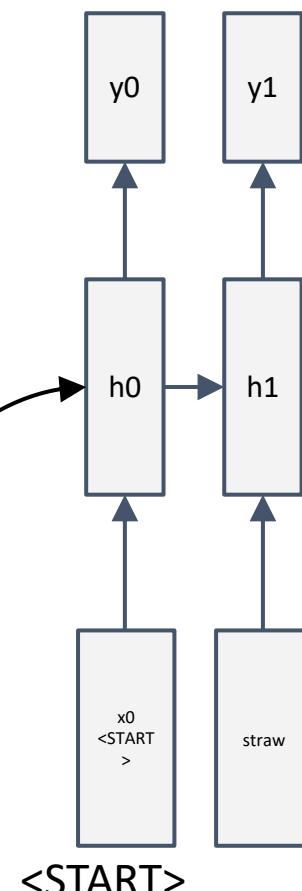
conv-512

conv-512

maxpool

FC-4096

FC-4096



image

← test image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

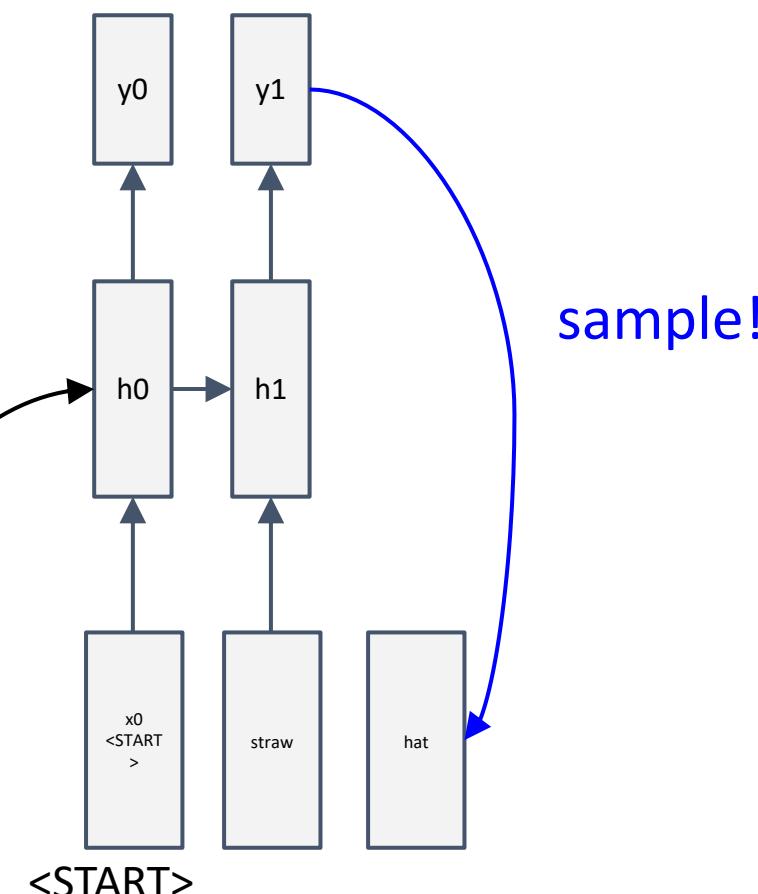
conv-512

conv-512

maxpool

FC-4096

FC-4096



image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

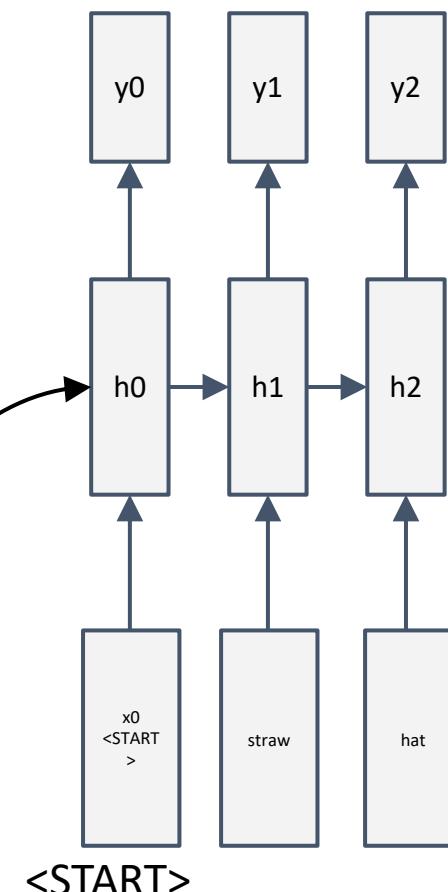
conv-512

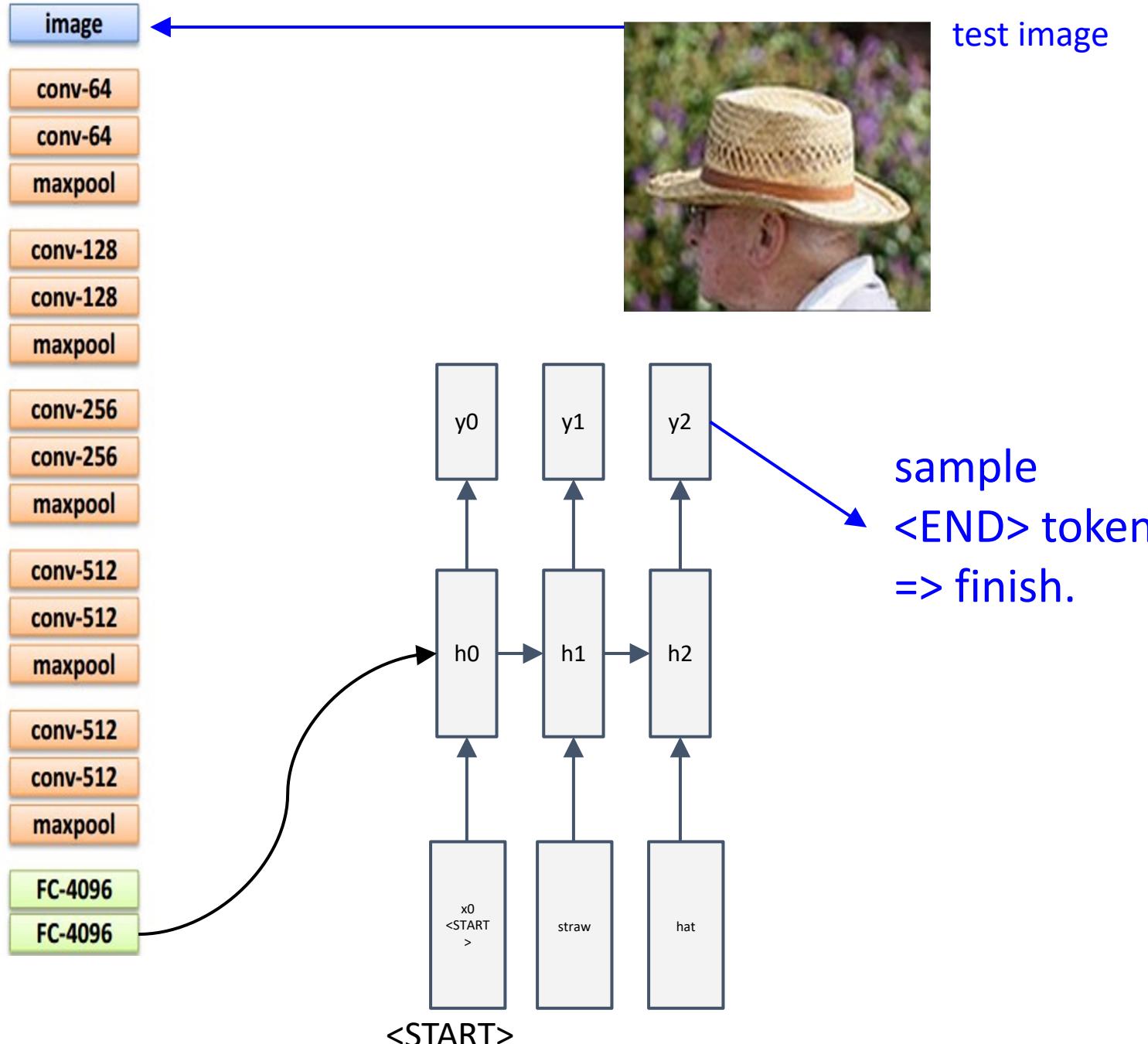
conv-512

maxpool

FC-4096

FC-4096





# Image Sentence Datasets

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](http://mscoco.org)

currently:  
~120K images  
~5 sentences each

# Image Captioning: Example Results



*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*

Captions generated using neuraltalk2  
All images are CC0 Public domain:  
[cat suitcase](#), [cat tree](#), [dog](#), [bear](#),  
[surfers](#), [tennis](#), [giraffe](#), [motorcycle](#)

# Image Captioning: Failure Cases

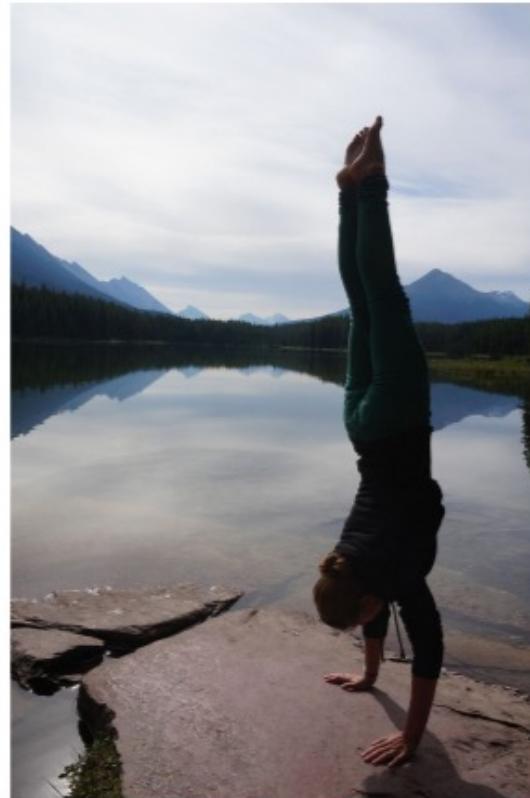
Captions generated using neuraltalk2  
All images are CC0 Public domain:  
[cat suitcase](#), [cat tree](#), [dog](#), [bear](#),  
[surfers](#), [tennis](#), [giraffe](#), [motorcycle](#)



*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*

# Visual Question Answering (VQA)



**Q: What endangered animal is featured on the truck?**

- A: A bald eagle.
- A: A sparrow.
- A: A humming bird.
- A: A raven.



**Q: Where will the driver go if turning right?**

- A: Onto 24 1/4 Rd.
- A: Onto 25 1/4 Rd.
- A: Onto 23 1/4 Rd.
- A: Onto Main Street.



**Q: When was the picture taken?**

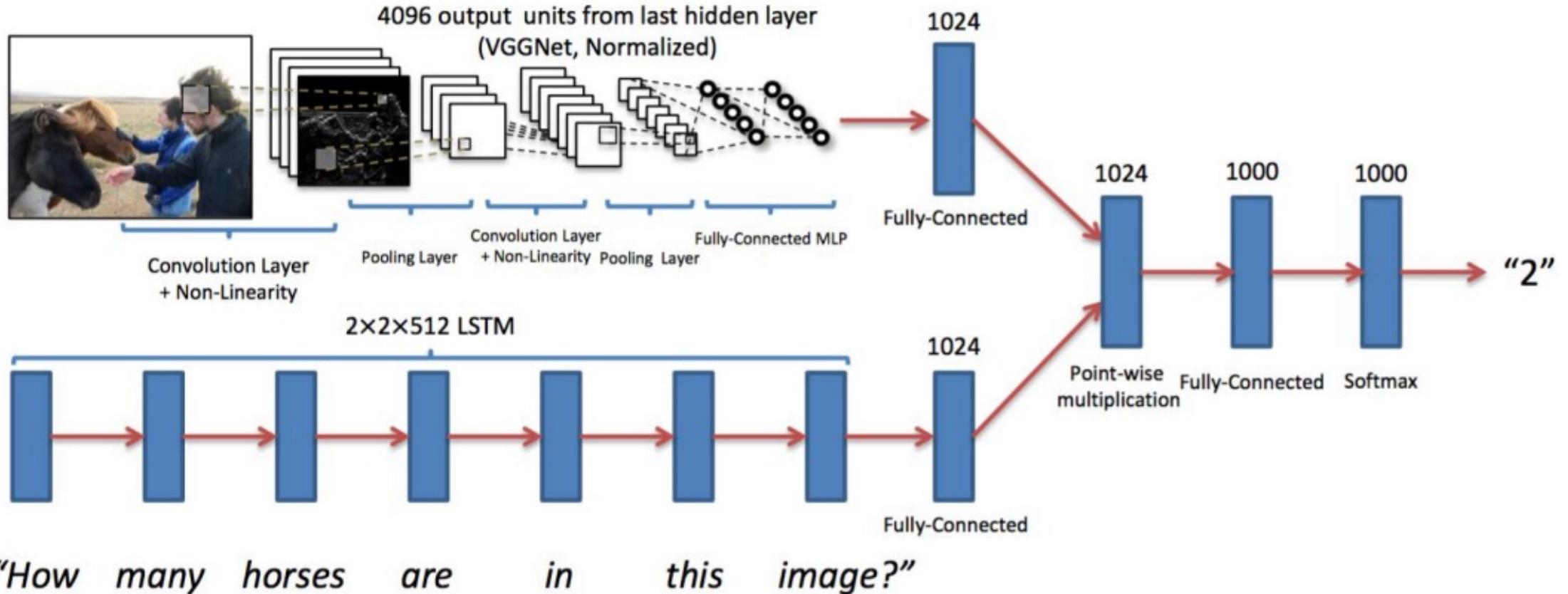
- A: During a wedding.
- A: During a bar mitzvah.
- A: During a funeral.
- A: During a Sunday church service.



**Q: Who is under the umbrella?**

- A: Two women.
- A: A child.
- A: An old man.
- A: A husband and a wife.

# Visual Question Answering (VQA)



Agrawal et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2015

Figures from Agrawal et al, copyright IEEE 2015. Reproduced for educational purposes.

## Sources:

- Sharif University of Technology, 40719 (DL Course), Spring 2025