



# ECE 6504: Deep Learning for Perception

## Topics:

- Recurrent Neural Networks (RNNs)
- BackProp Through Time (BPTT)
- Vanishing / Exploding Gradients
- [Abhishek:] Lua / Torch Tutorial

Dhruv Batra  
Virginia Tech

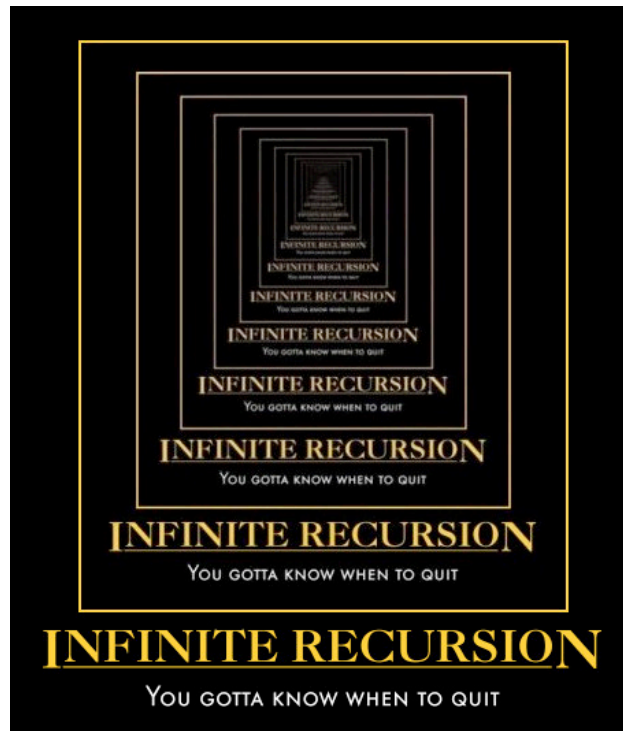
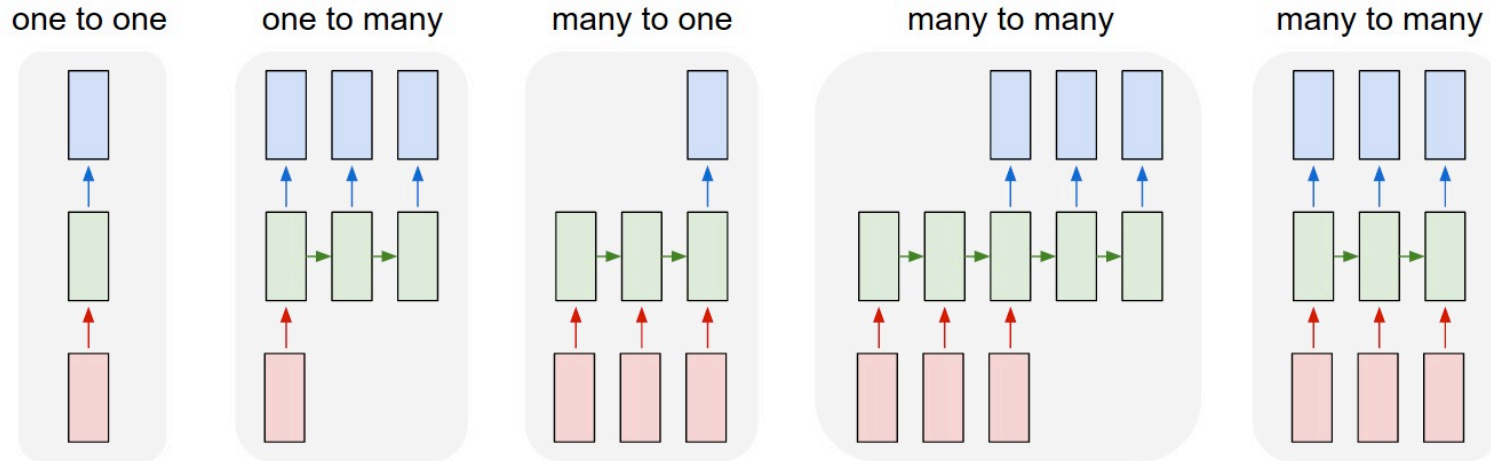
# Administrativa

- HW3
  - Out today
  - Due in 2 weeks
  - Please please please please please start early
  - <https://computing.ece.vt.edu/~f15ece6504/homework3/>

# Plan for Today

- Model
  - Recurrent Neural Networks (RNNs)
- Learning
  - BackProp Through Time (BPTT)
  - Vanishing / Exploding Gradients
- [Abhishek:] Lua / Torch Tutorial

# New Topic: RNNs

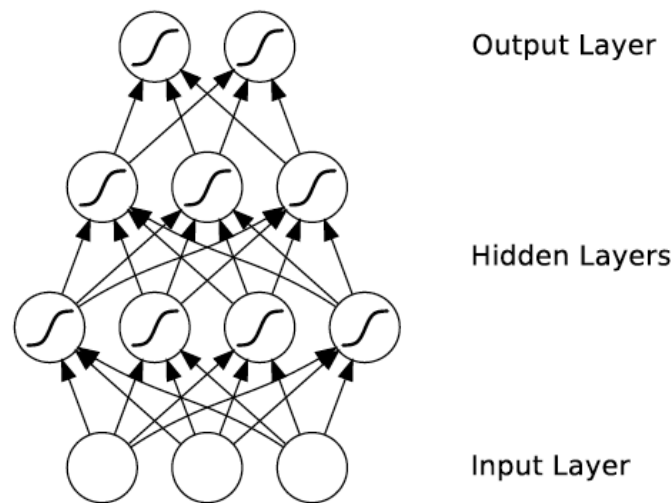


# Synonyms

- Recurrent Neural Networks (RNNs)
- Recursive Neural Networks
  - General family; think graphs instead of chains
- Types:
  - Long Short Term Memory (LSTMs)
  - Gated Recurrent Units (GRUs)
  - Hopfield network
  - Elman networks
  - ...
- Algorithms
  - BackProp Through Time (BPTT)
  - BackProp Through Structure (BPTS)

# What's wrong with MLPs?

- Problem 1: Can't model sequences
  - Fixed-sized Inputs & Outputs
  - No temporal structure
- Problem 2: Pure feed-forward processing
  - No “memory”, no feedback



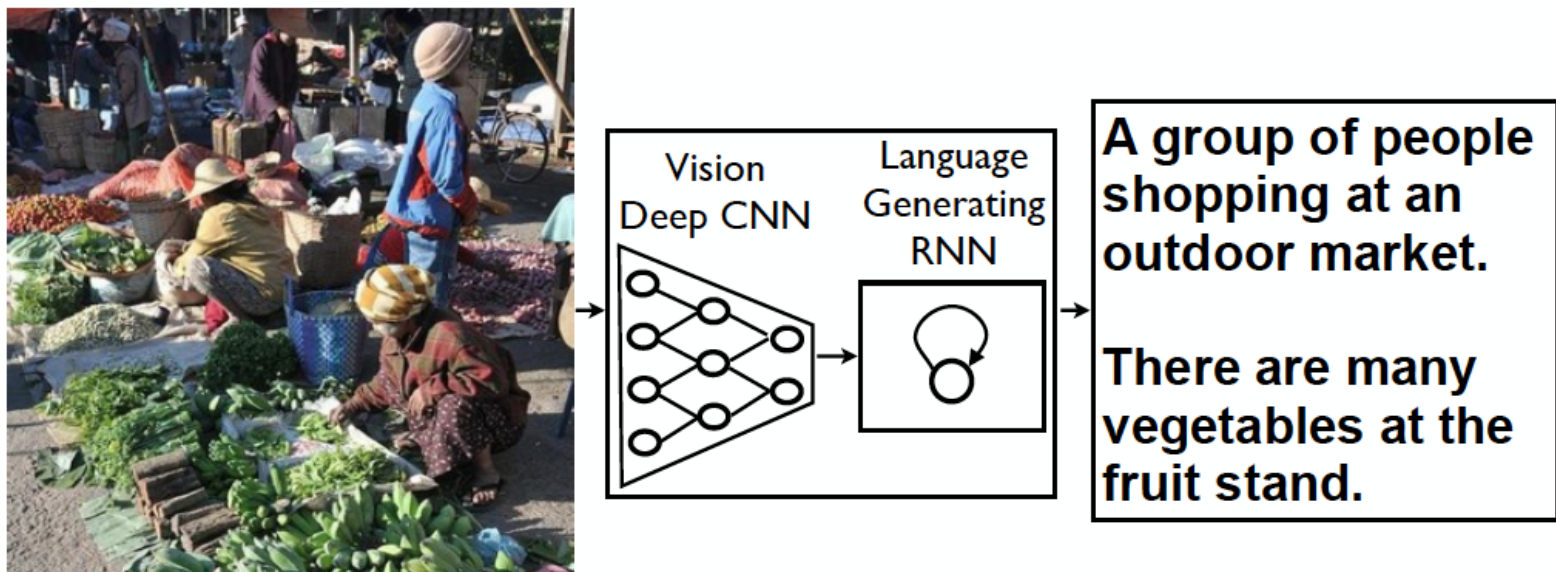
# Sequences are everywhere...

*Foreign Minister.* → FOREIGN MINISTER.

 → THE SOUND OF

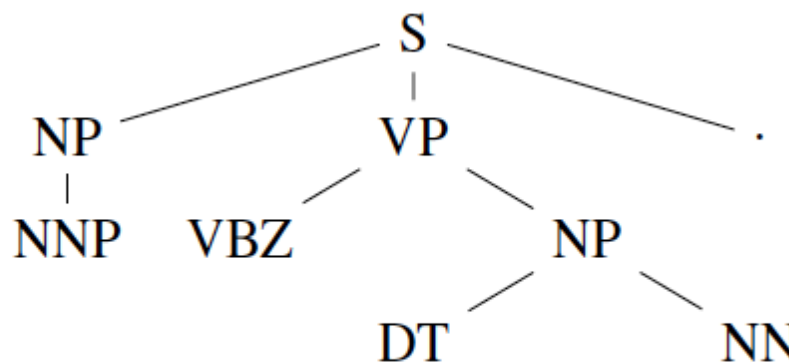
$a_1=2$     $a_2=0$     $a_3=1$     $a_4=3$     $a_5=4$     $a_6=2$     $a_7=5$   
 $x =$  bringen   sie   bitte   das   auto   zurück   .  
↙  
 $y =$  please   return   the   car   .

# Even where you might not expect a sequence...



John has a dog .

→



John has a dog .

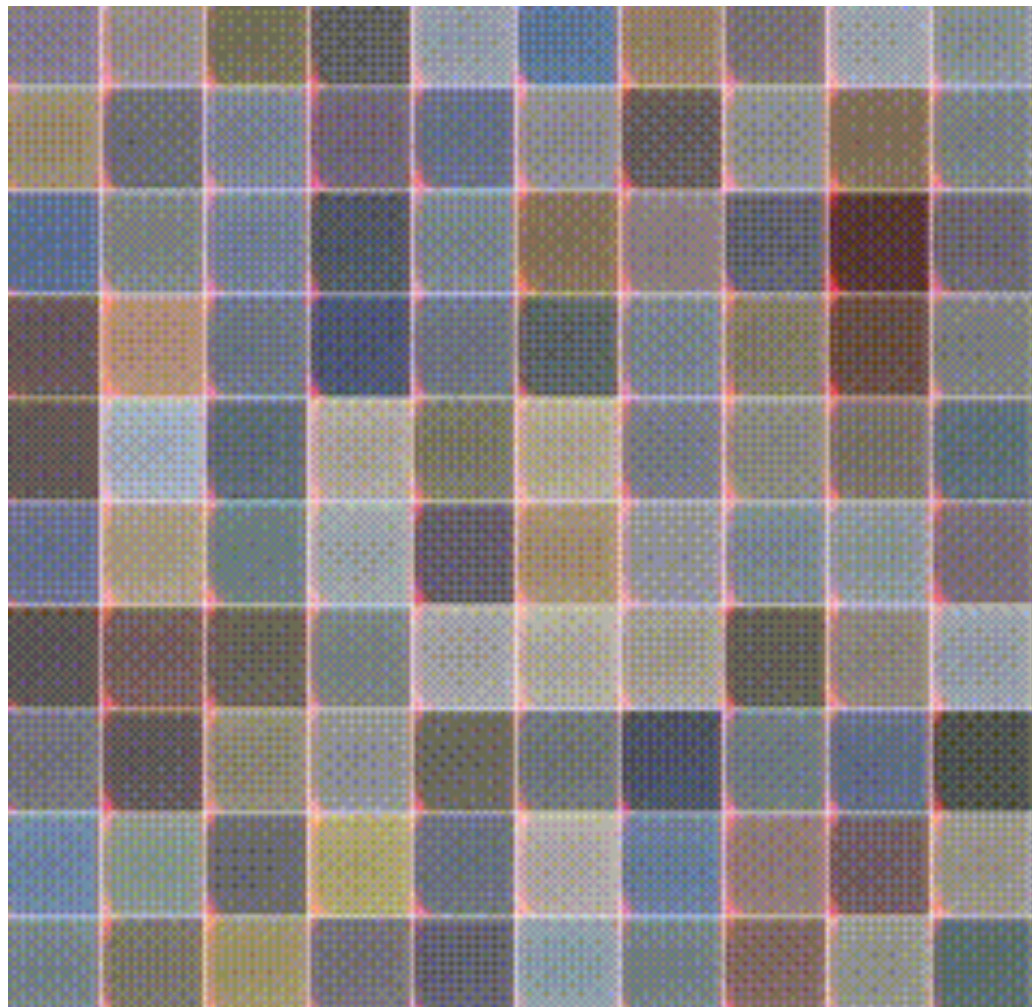
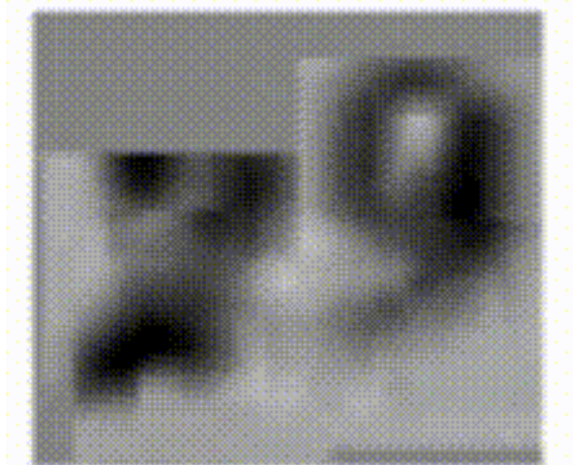
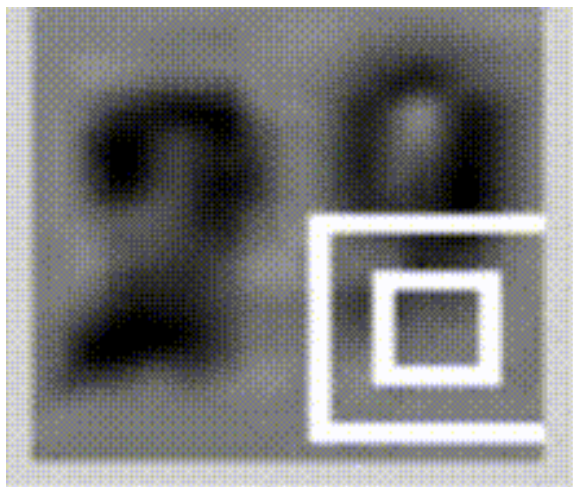
→

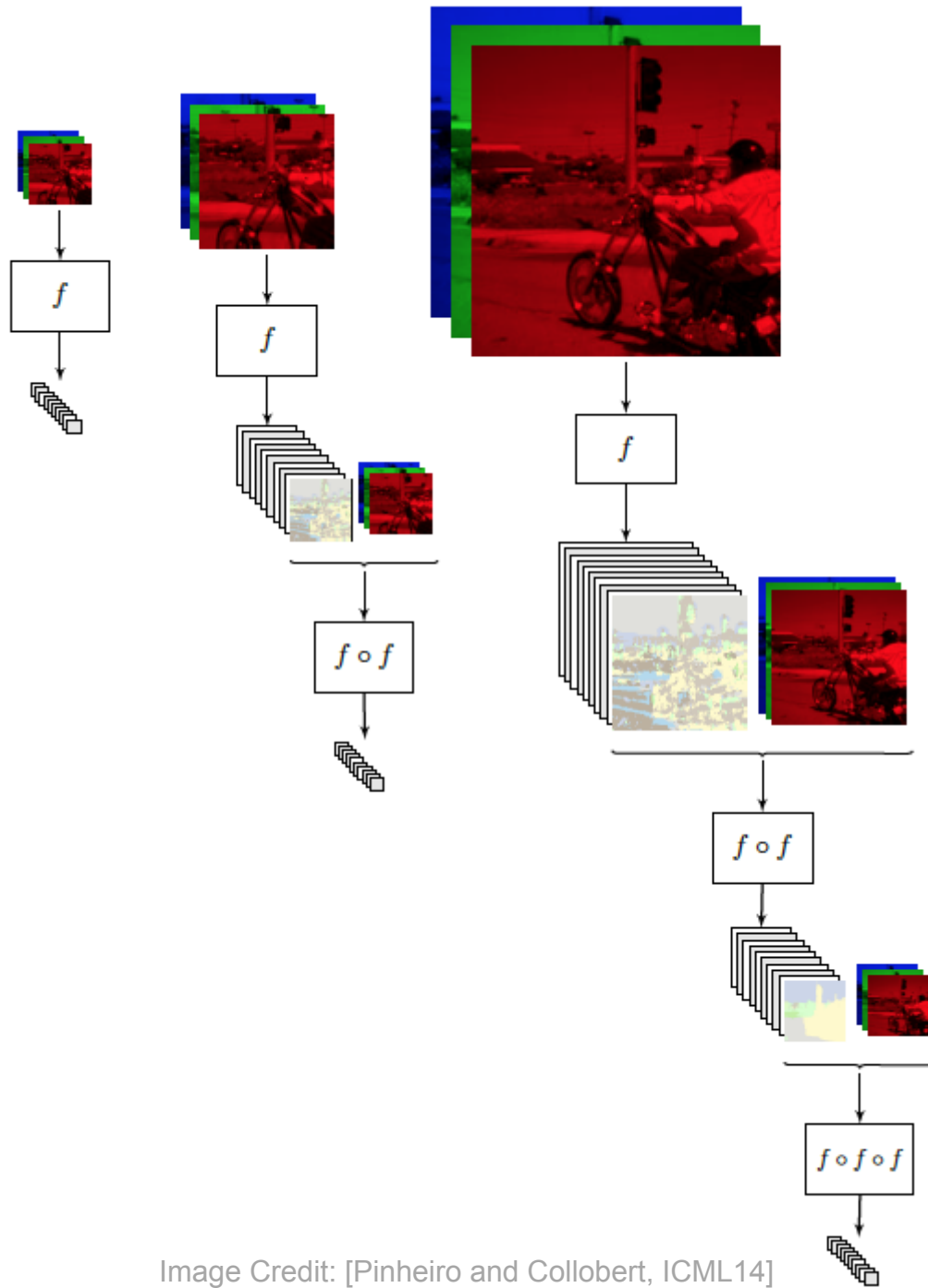
$(S (NP NNP)_{NP} (VP VBZ (NP DT NN)_{NP})_{VP} .)_{S}$



# Even where you might not expect a sequence...

- Input ordering = sequence





# Why model sequences?

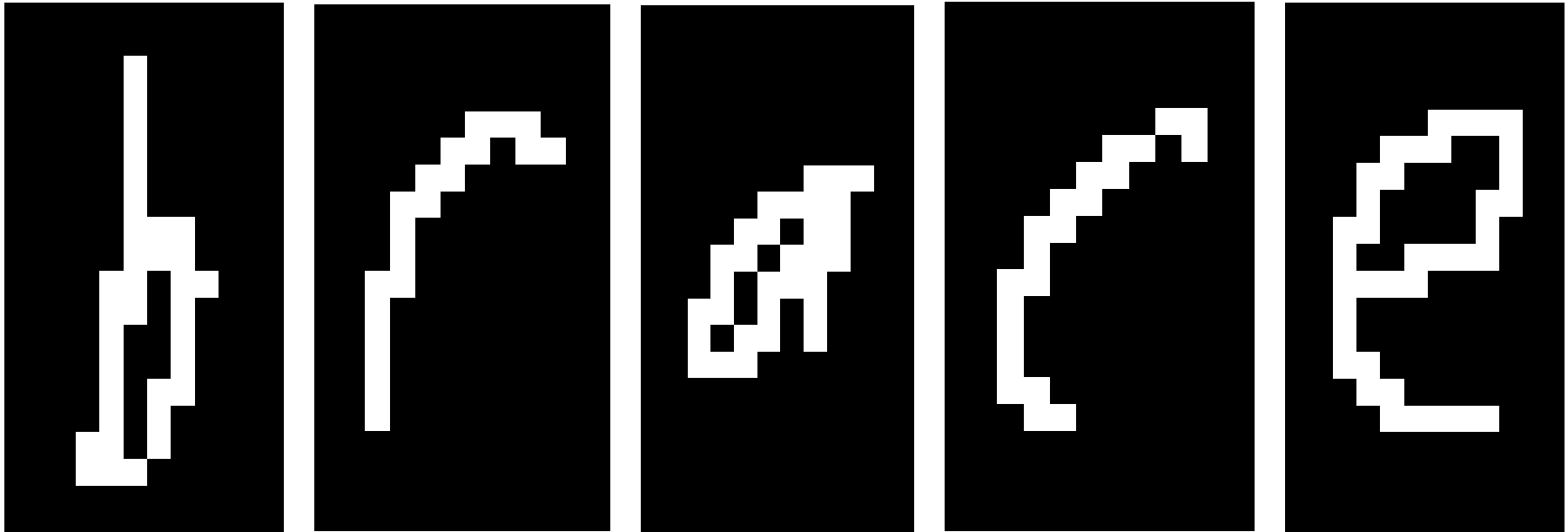
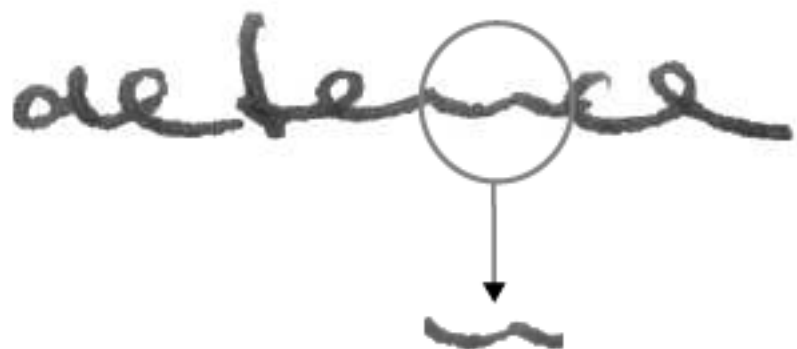
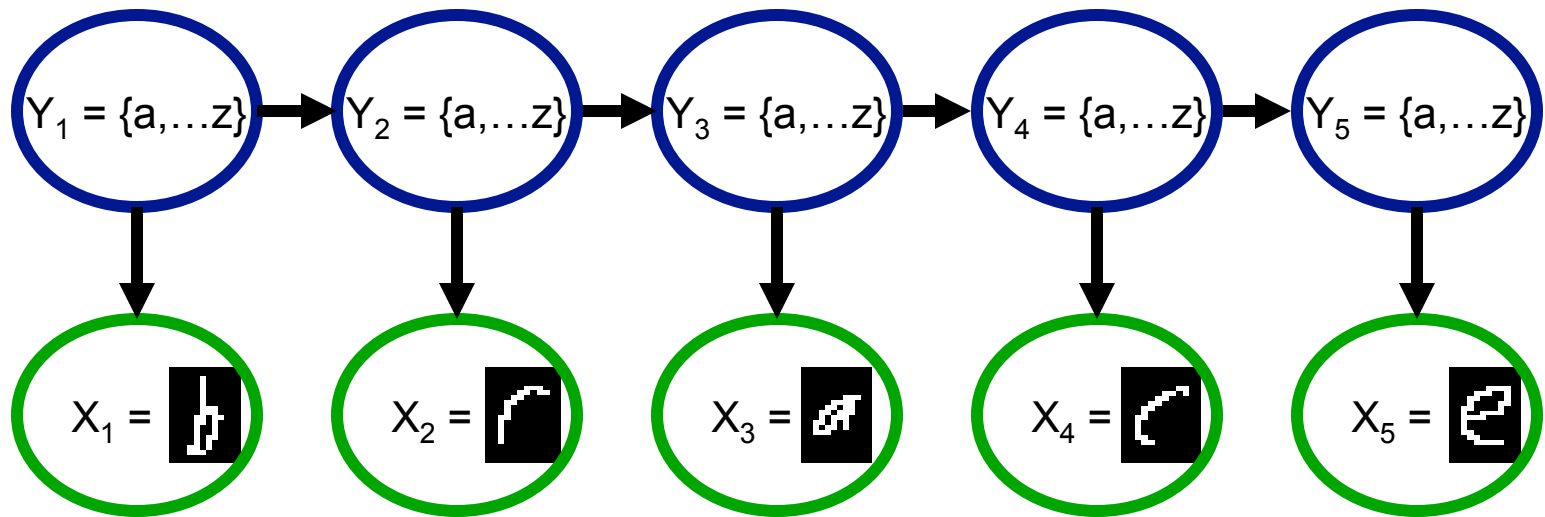


Figure Credit: Carlos Guestrin

# Why model sequences?



# Name that model

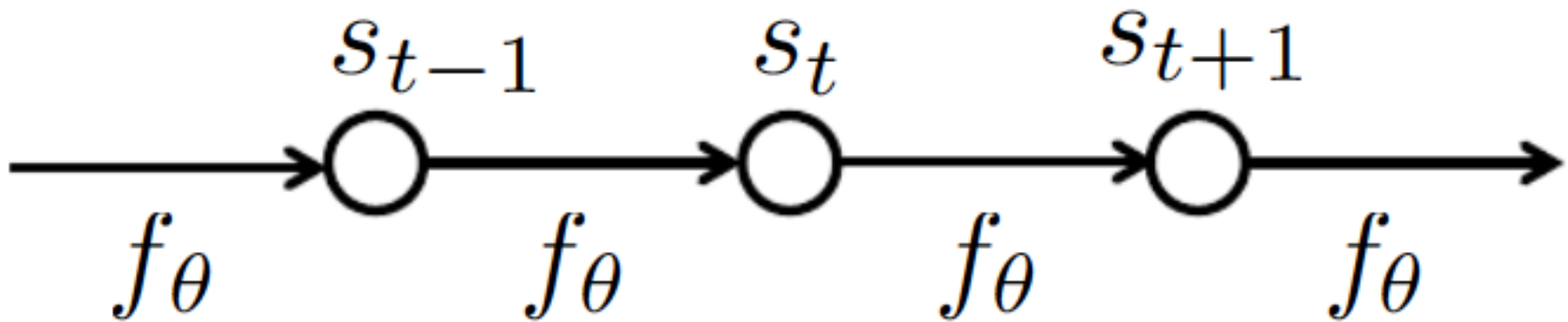


*Hidden Markov Model (HMM)*

# How do we model sequences?

- No input

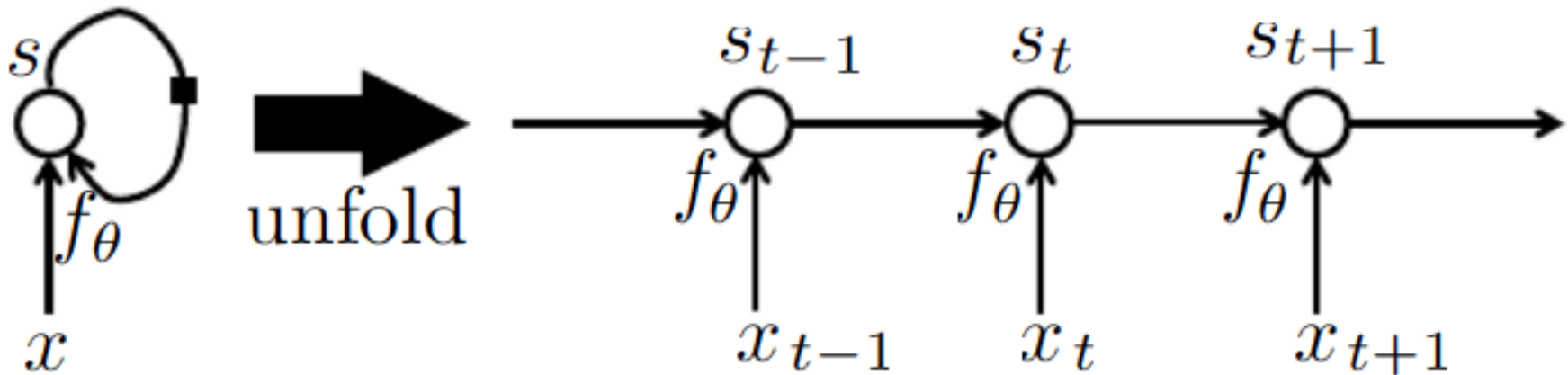
$$s_t = f_\theta(s_{t-1})$$



# How do we model sequences?

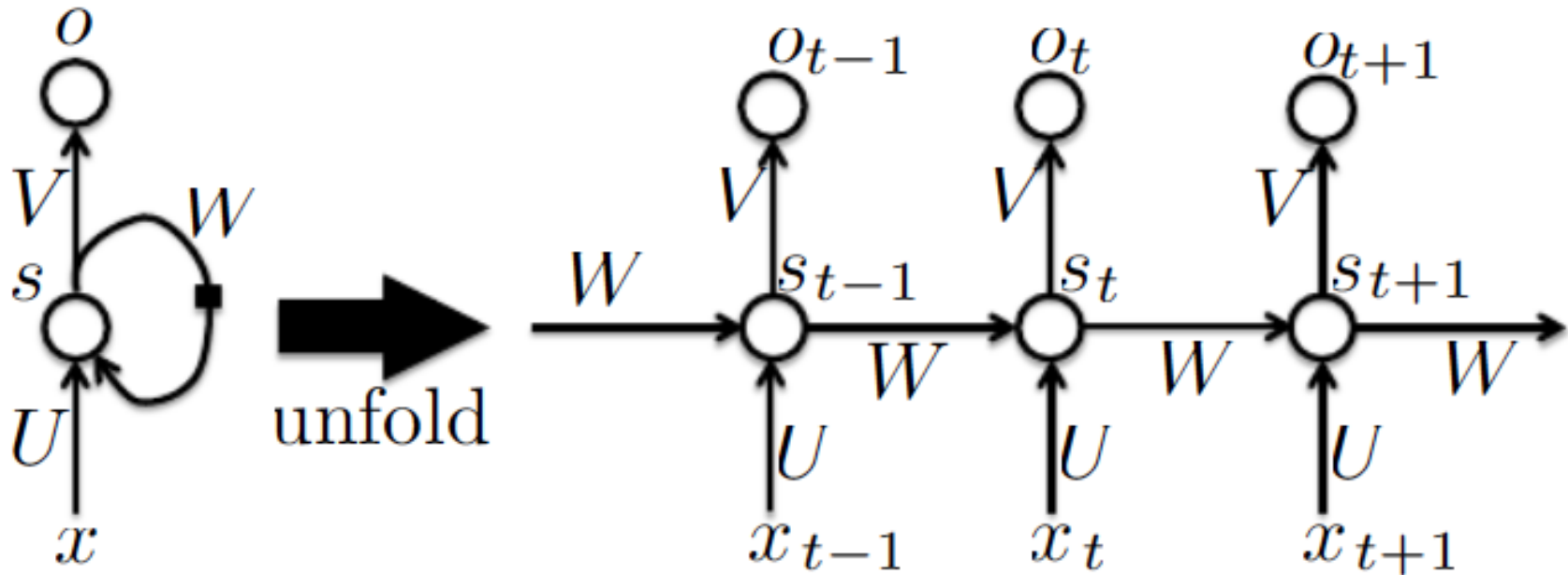
- With inputs

$$s_t = f_{\theta}(s_{t-1}, x_t)$$



# How do we model sequences?

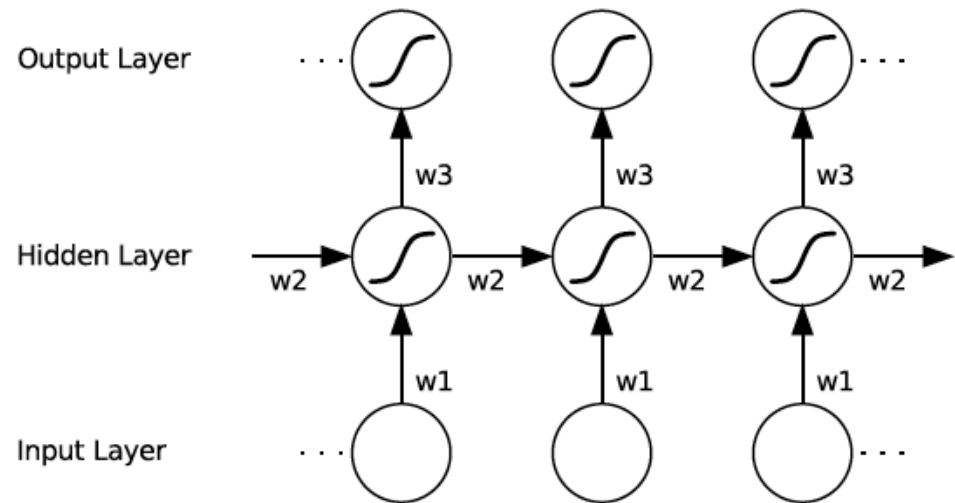
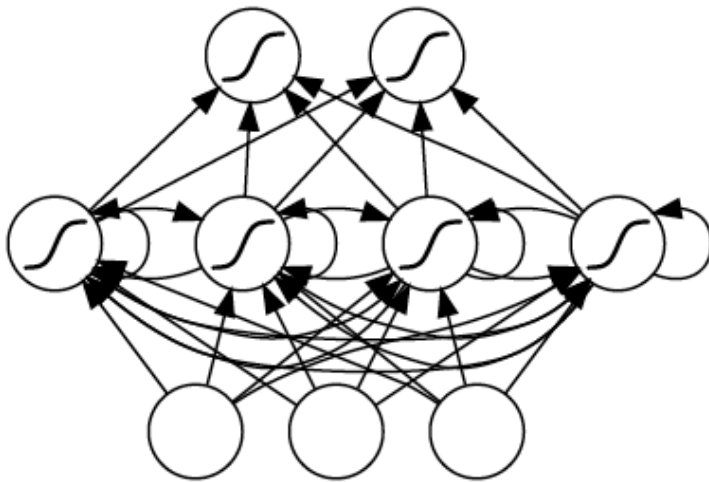
- With inputs and outputs





# How do we model sequences?

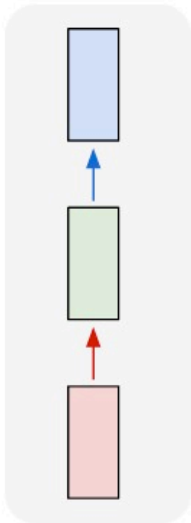
- With Neural Nets



# How do we model sequences?

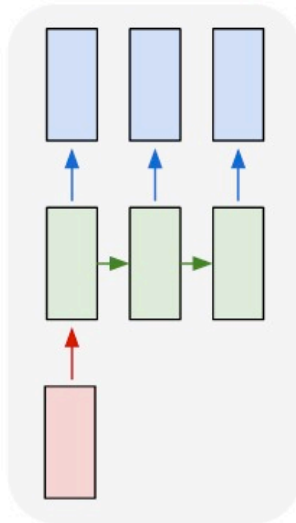
- It's a spectrum...

one to one



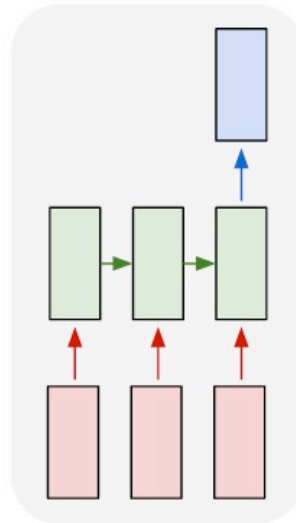
Input: No sequence  
Output: No sequence  
Example: "standard" classification / regression problems

one to many



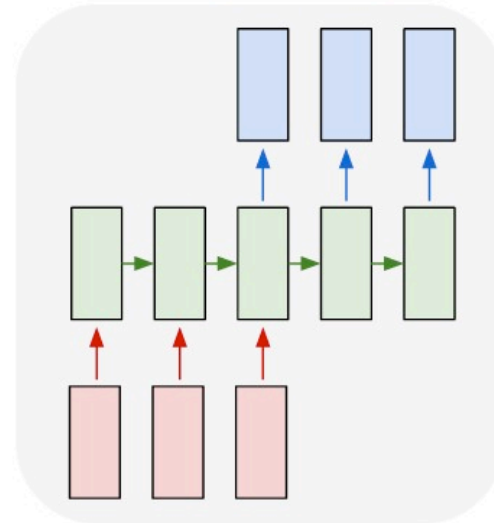
Input: No sequence  
Output: Sequence  
Example: Im2Caption

many to one



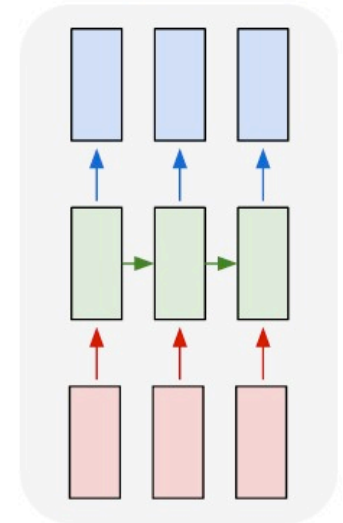
Input: Sequence  
Output: No sequence  
Example: sentence classification, multiple-choice question answering

many to many

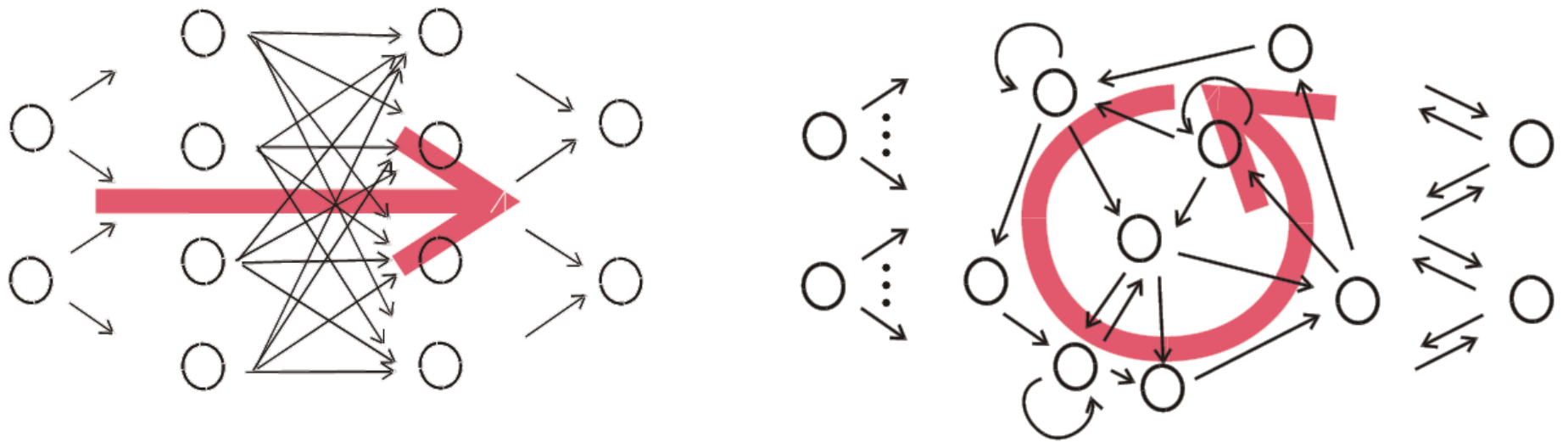


Input: Sequence  
Output: Sequence  
Example: machine translation, video captioning, open-ended question answering, video question answering

many to many



# Things can get arbitrarily complex



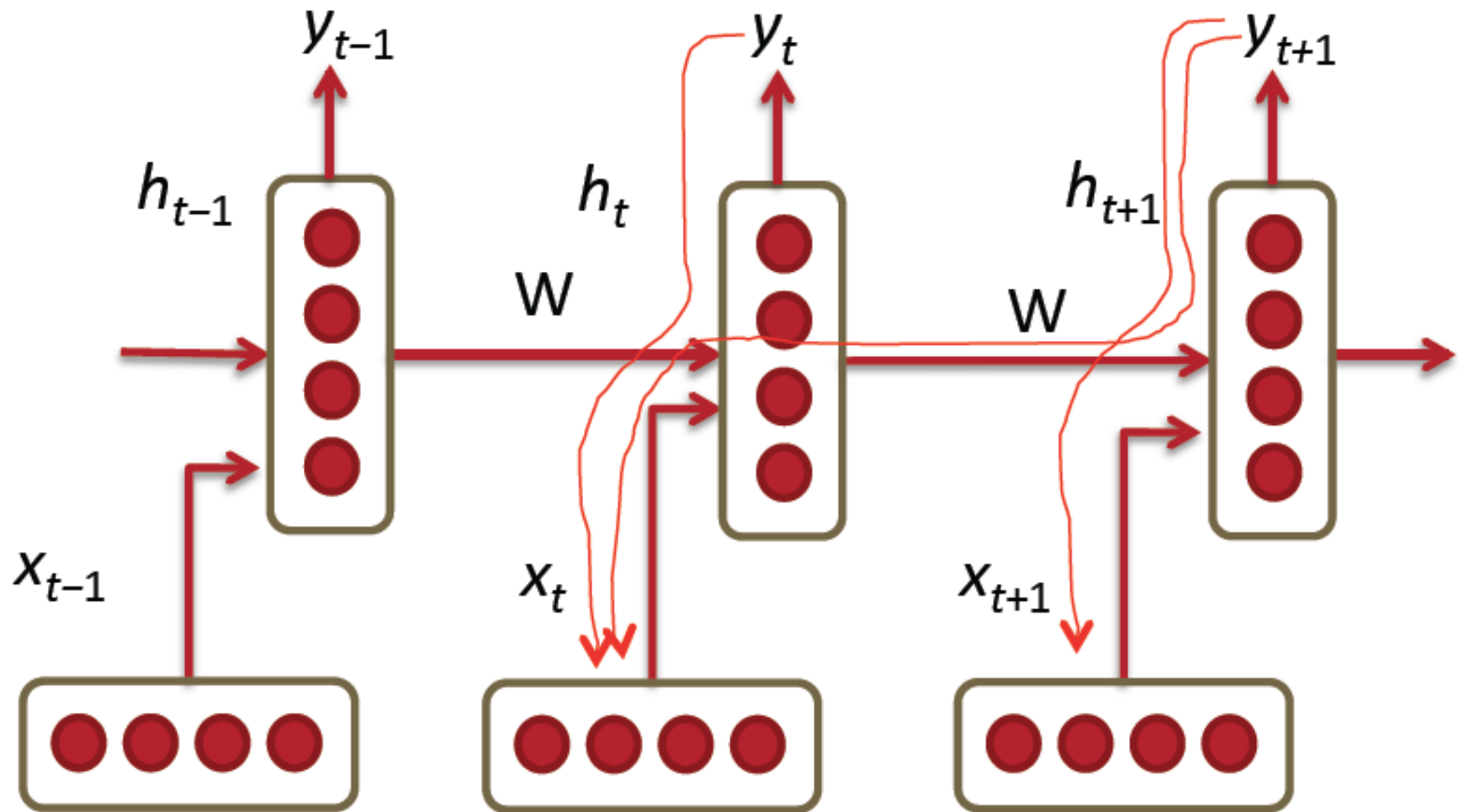
# Key Ideas

- Parameter Sharing + Unrolling
  - Keeps numbers of parameters in check
  - Allows arbitrary sequence lengths!
- “Depth”
  - Measured in the usual sense of layers
  - Not unrolled timesteps
- Learning
  - Is tricky even for “shallow” models due to unrolling

# Plan for Today

- Model
  - Recurrent Neural Networks (RNNs)
- Learning
  - BackProp Through Time (BPTT)
  - Vanishing / Exploding Gradients
- [Abhishek:] Lua / Torch Tutorial

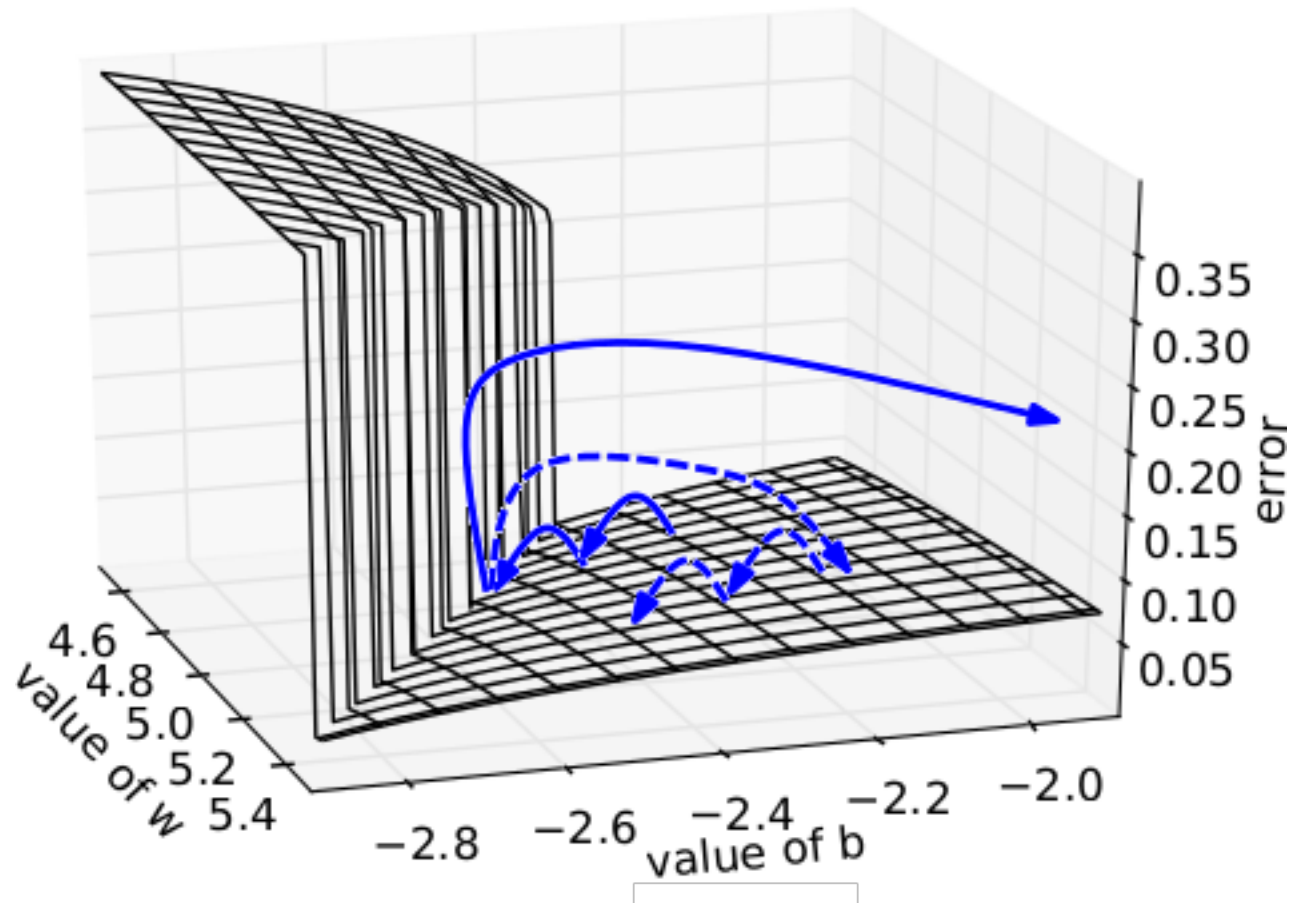
# BPTT



# Illustration [Pascanu et al]

- Intuition

- Error surface of a single hidden unit RNN; High curvature walls
- Solid lines: standard gradient descent trajectories
- Dashed lines: gradient rescaled to fix problem



# Fix #1

- Pseudocode

---

**Algorithm 1** Pseudo-code for norm clipping the gradients whenever they explode

---

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$

**if**  $\|\hat{\mathbf{g}}\| \geq \textit{threshold}$  **then**

$$\hat{\mathbf{g}} \leftarrow \frac{\textit{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$

**end if**

---



# Fix #2

- Smart Initialization and ReLus
  - [Socher et al 2013]
  - *A Simple Way to Initialize Recurrent Networks of Rectified Linear Units*, Le et al. 2015

