



# CS 4803 / 7643: Deep Learning

## Topics:

- Recurrent Neural Networks (RNNs)
- BackProp Through Time (BPTT)

Dhruv Batra  
Georgia Tech

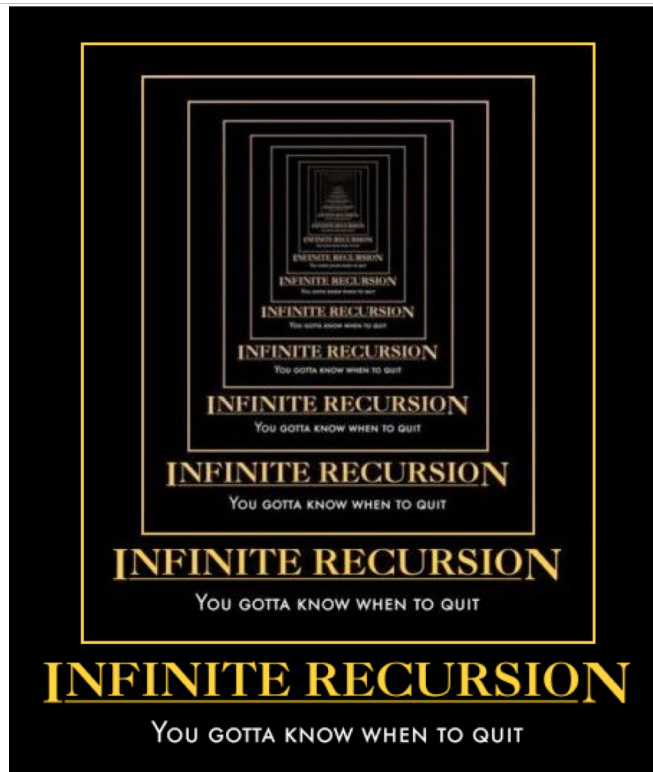
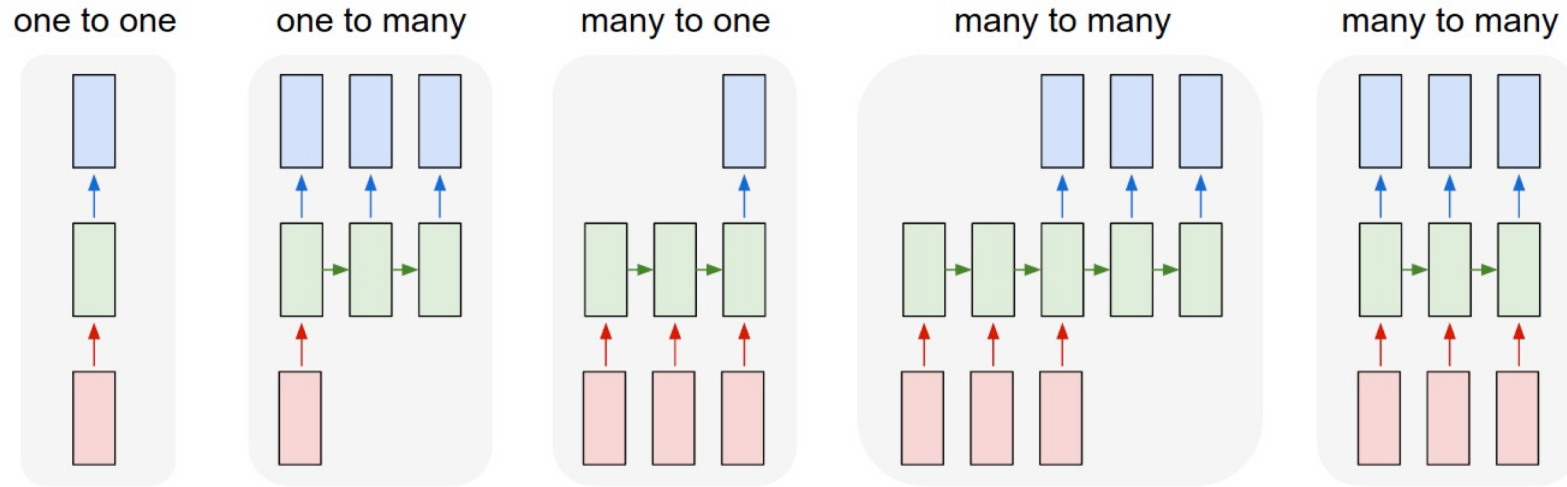
# Administrativa

- HW3 Released
  - Due: 11/06, 11:55pm
  - Last HW
  - Focus on projects after this
  - [https://www.cc.gatech.edu/classes/AY2019/cs7643\\_fall/assets/hw3.pdf](https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/assets/hw3.pdf)

# Plan for Today

- Model
  - Recurrent Neural Networks (RNNs)
- Learning
  - BackProp Through Time (BPTT)

# New Topic: RNNs



# New Words

- **Recurrent Neural Networks (RNNs)**

- **Recursive Neural Networks**

- General family; think graphs instead of chains

- **Types:**

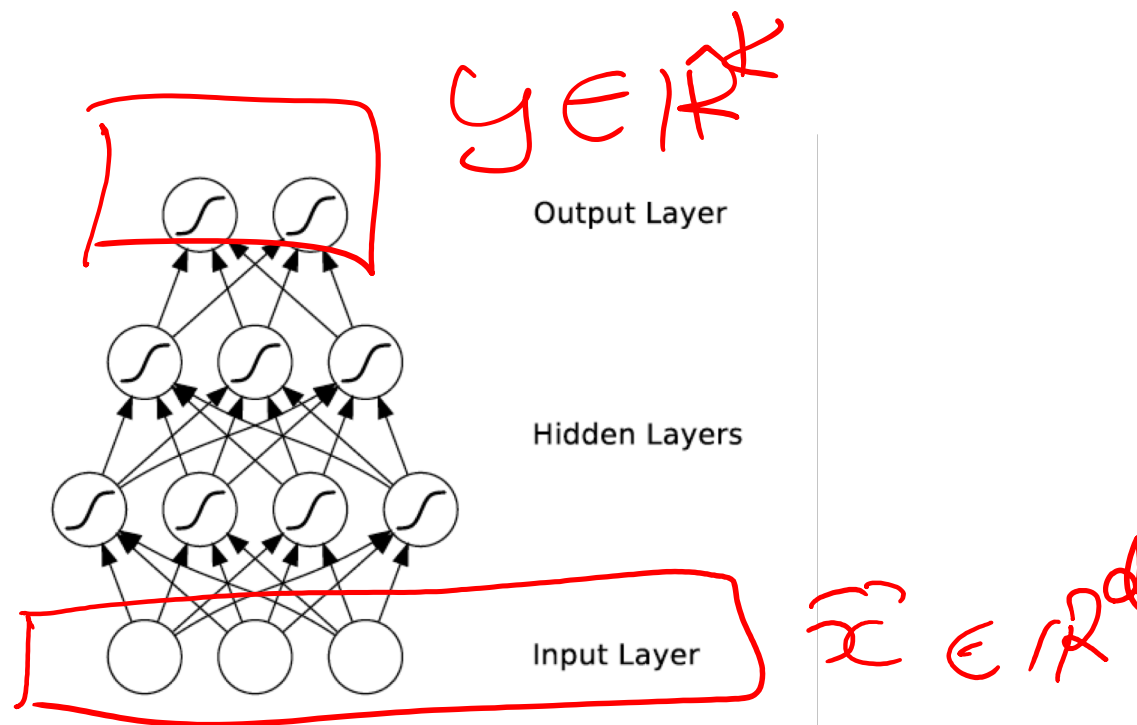
- **“Vanilla” RNNs (Elman Networks)**
- Long Short Term Memory (LSTMs)
- Gated Recurrent Units (GRUs)
- ...

- **Algorithms**

- BackProp Through Time (BPTT)
- BackProp Through Structure (BPTS)

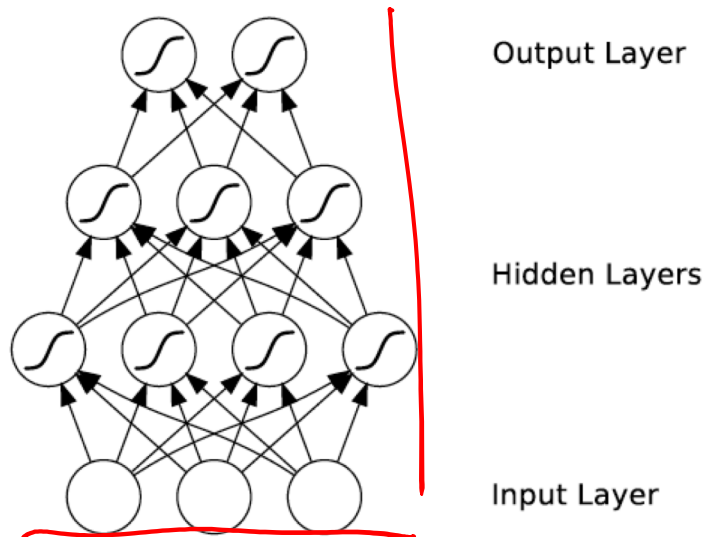
# What's wrong with MLPs?

- Problem 1: Can't model sequences
  - Fixed-sized Inputs & Outputs
  - No temporal structure

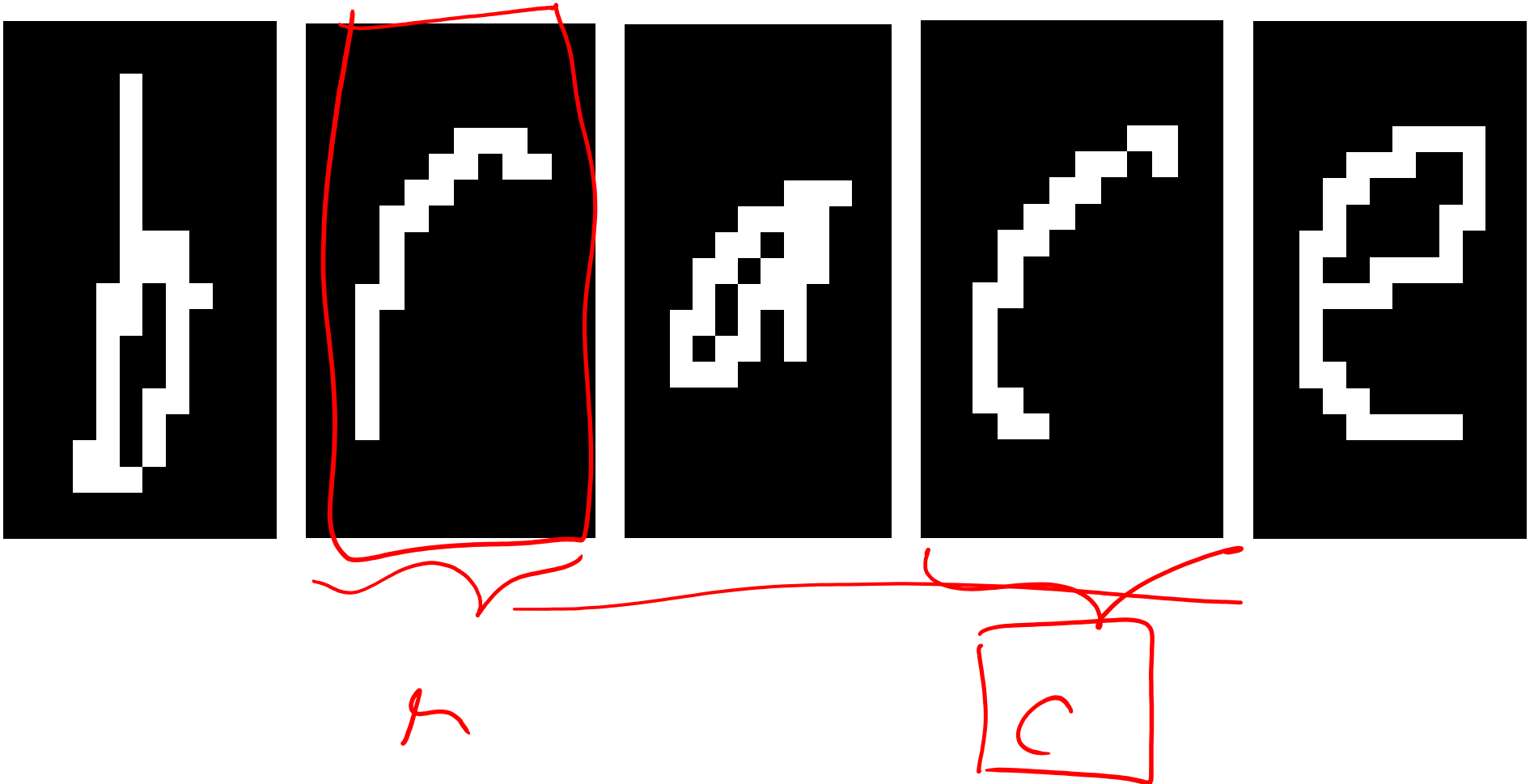


# What's wrong with MLPs?

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

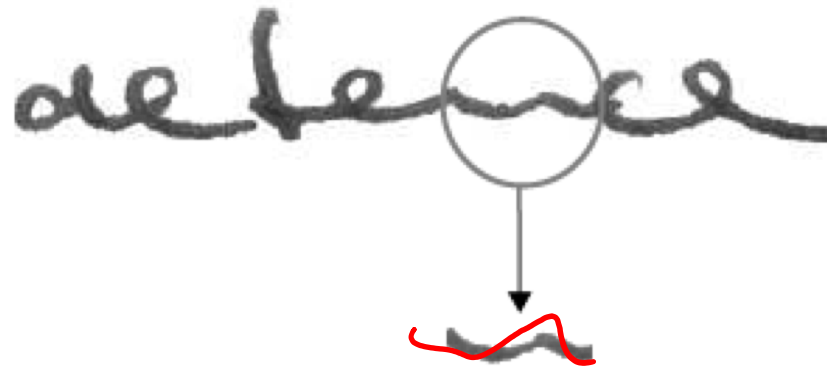


# Why model sequences?





# Why model sequences?



# Sequences are everywhere...

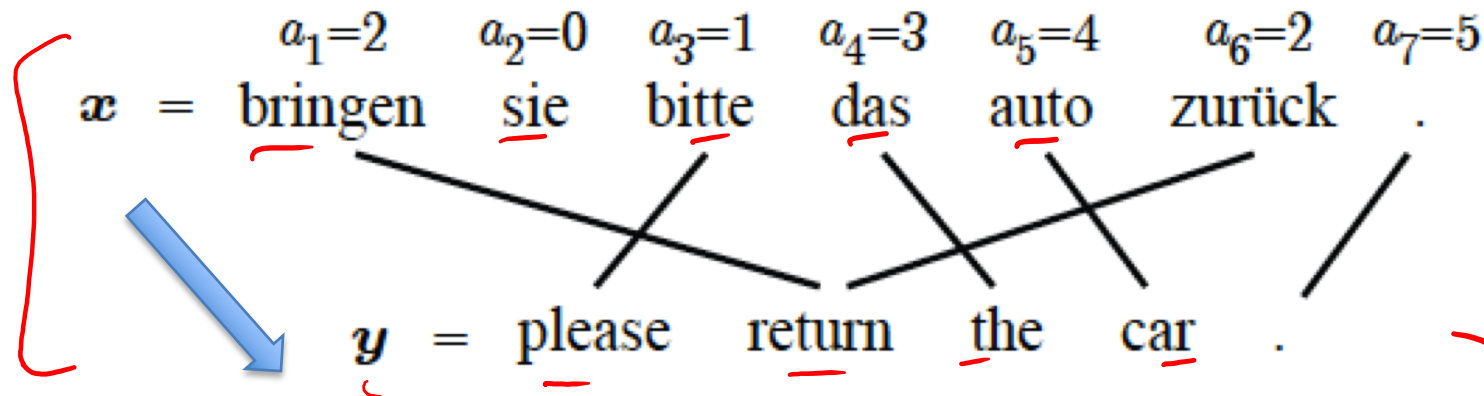
Foreign Minister.



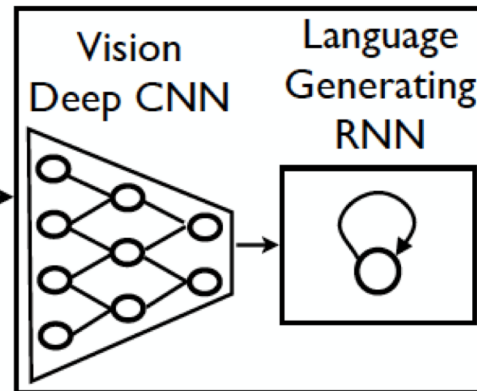
FOREIGN MINISTER.



THE SOUND OF



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

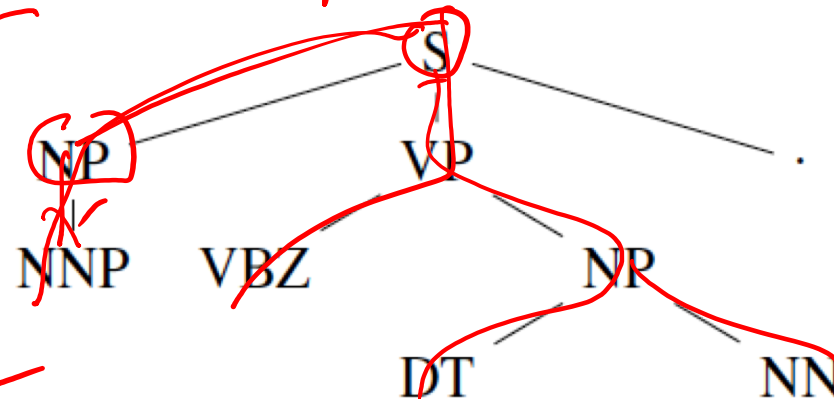


A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

John has a dog .

→



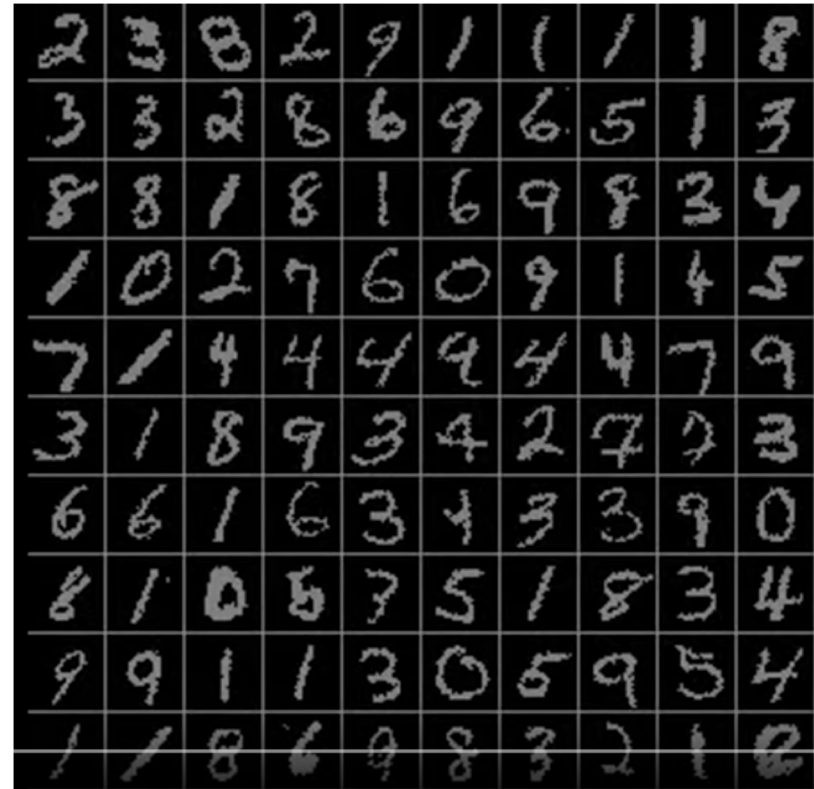
John has a dog .

→

(S (NP NNP )<sub>NP</sub> (VP VBZ (NP DT NN )<sub>NP</sub> )<sub>VP</sub> . )<sub>S</sub>

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

Classify images by taking a series of “glimpses”



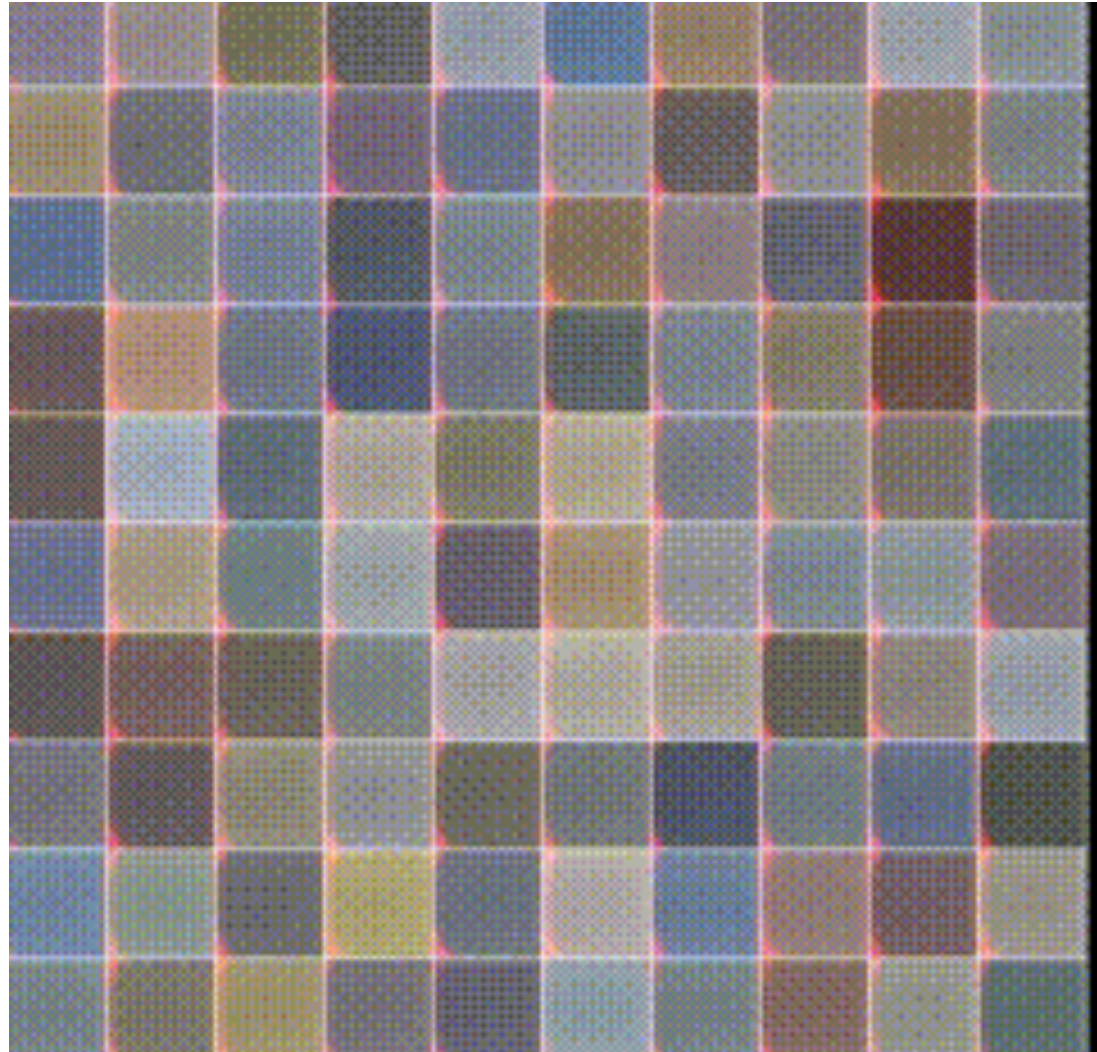
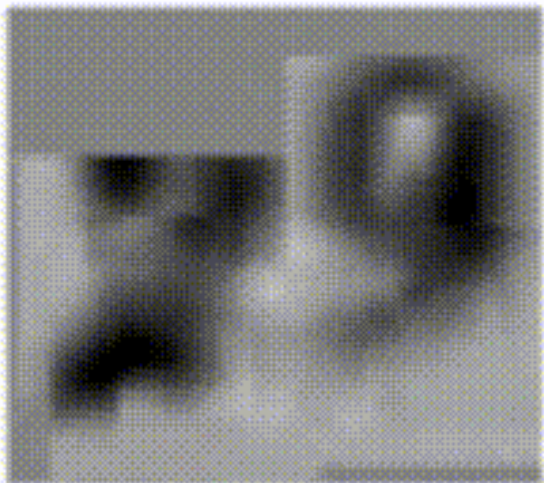
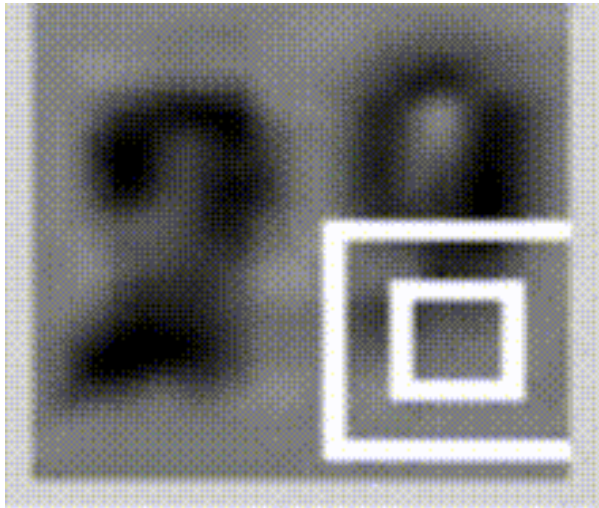
Ba, Mnih, and Kavukcuoglu, "Multiple Object Recognition with Visual Attention", ICLR 2015.

Gregor et al, "DRAW: A Recurrent Neural Network For Image Generation", ICML 2015

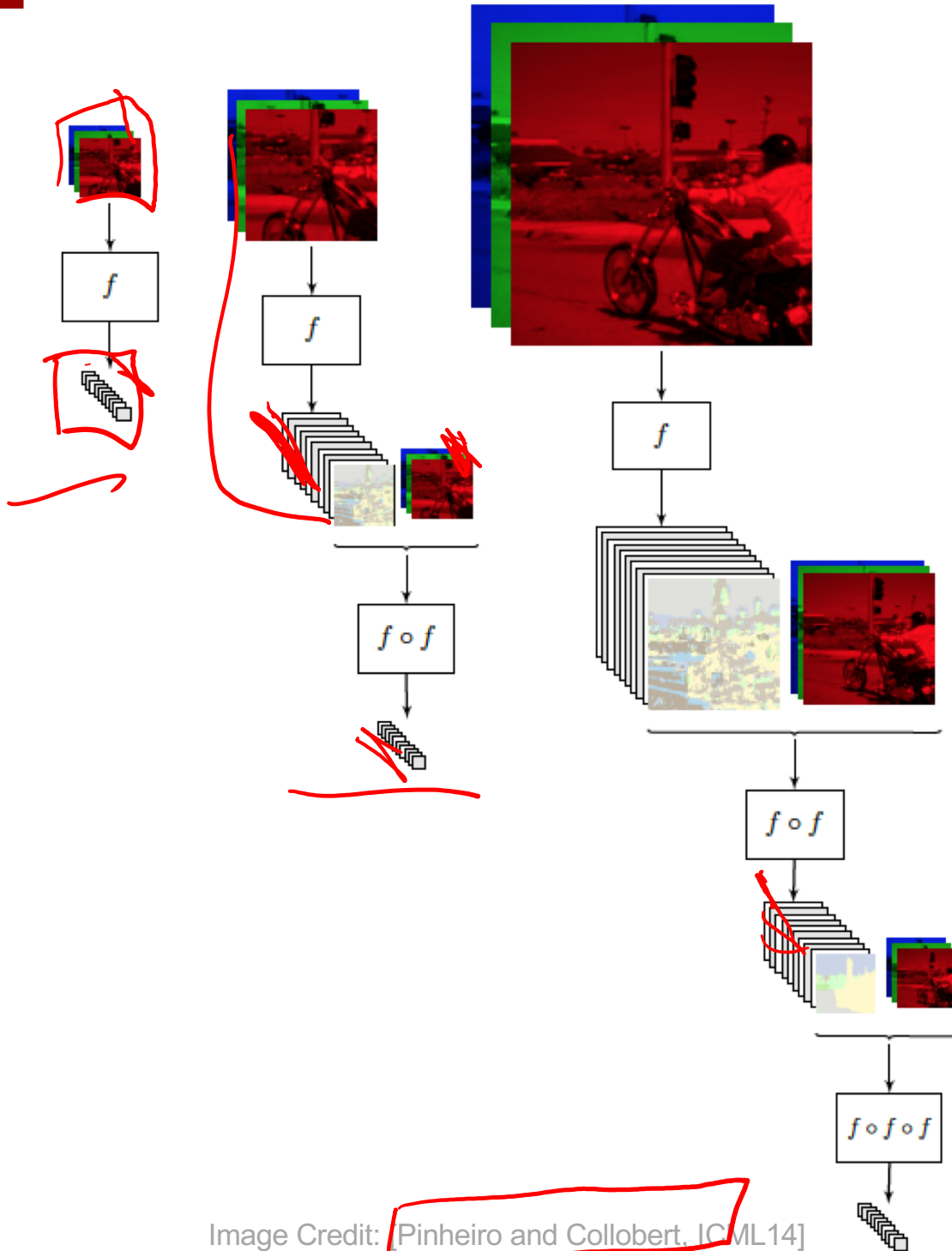
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission.

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

- Output ordering = sequence



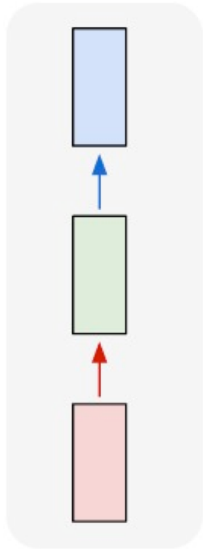
$7 \times 7 + K$



# Sequences in Input or Output?

- It's a spectrum...

one to one



Input: No  
sequence

Output: No  
sequence

Example:  
"standard"  
classification /  
regression  
problems

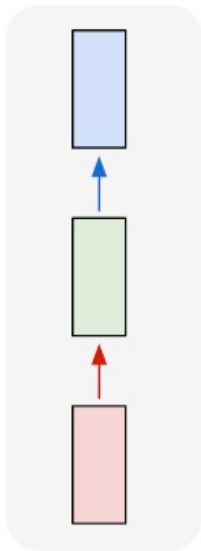
(C) Dhruv Batra

Image Credit: Andrej Karpathy

# Sequences in Input or Output?

- It's a spectrum...

one to one



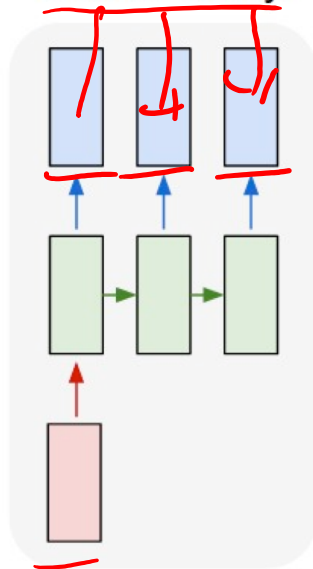
Input: No  
sequence

Output: No  
sequence

Example:  
"standard"  
classification /  
regression  
problems

(C) Dhruv Batra

one to many



Input: No sequence

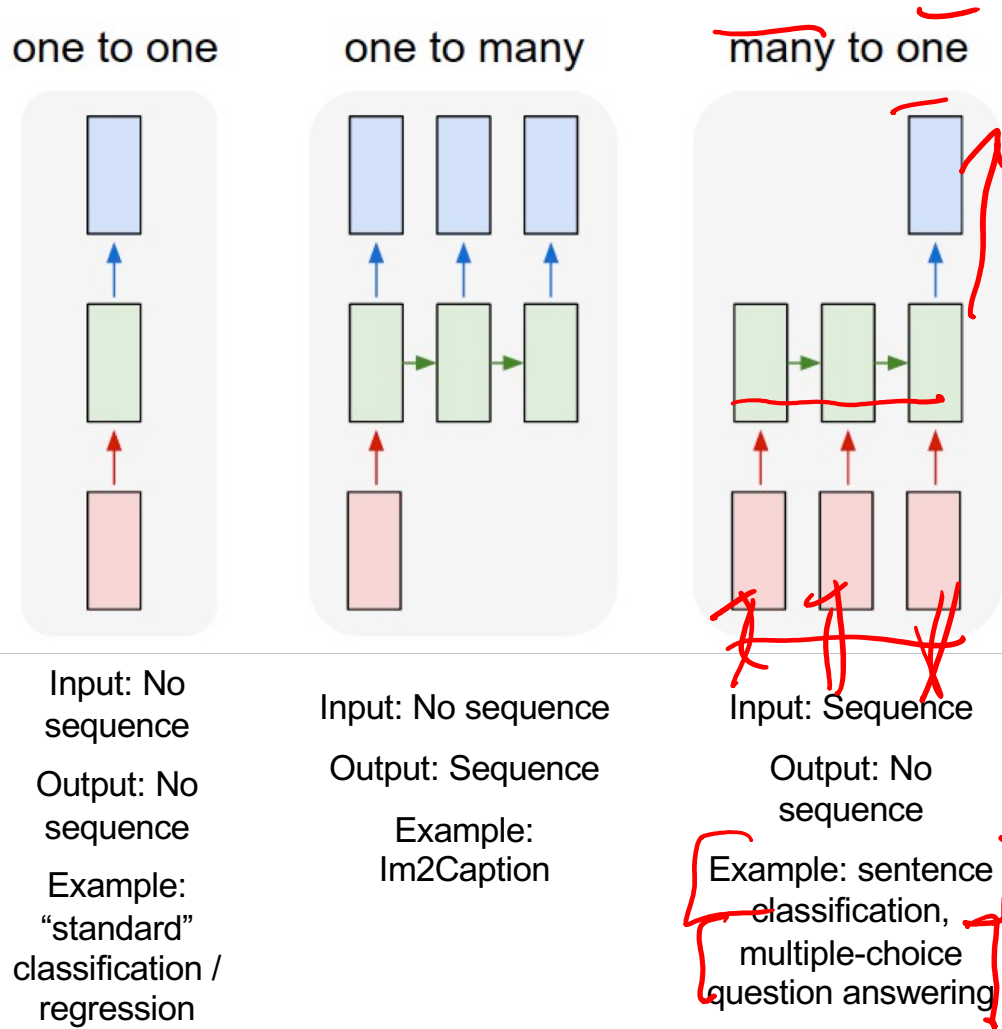
Output: Sequence

Example:  
Im2Caption



# Sequences in Input or Output?

- It's a spectrum...

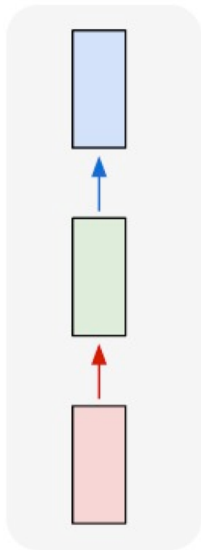


# Sequences in Input or Output?

o

- It's a spectrum...

one to one

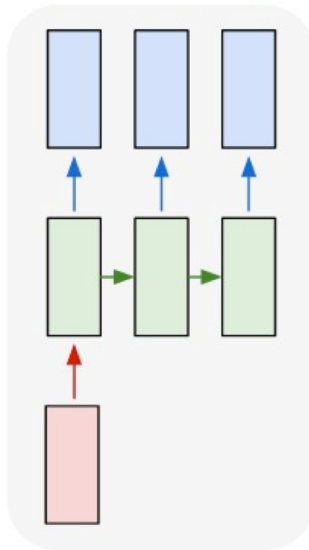


Input: No sequence

Output: No sequence

Example: "standard" classification / regression problems

one to many

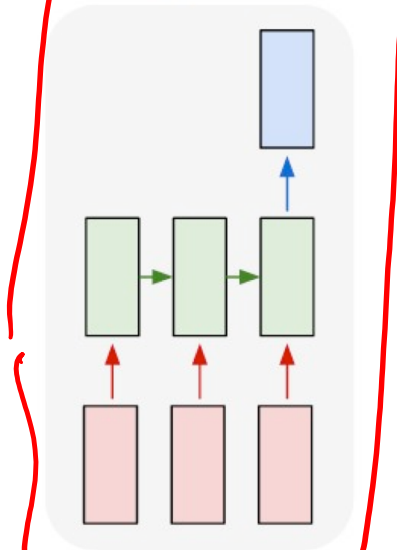


Input: No sequence

Output: Sequence

Example: Im2Caption

many to one

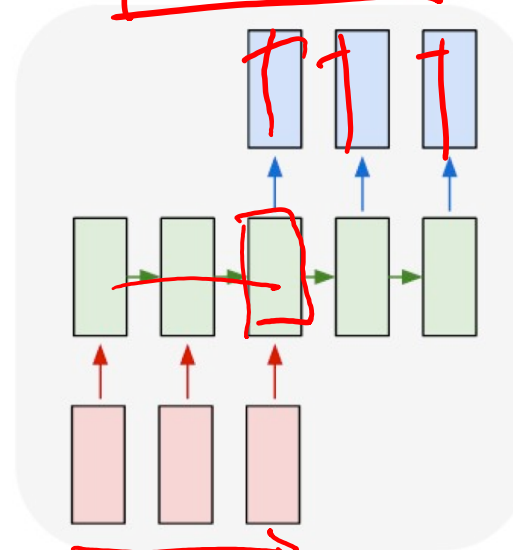


Input: Sequence

Output: No sequence

Example: sentence classification, multiple-choice question answering

many to many

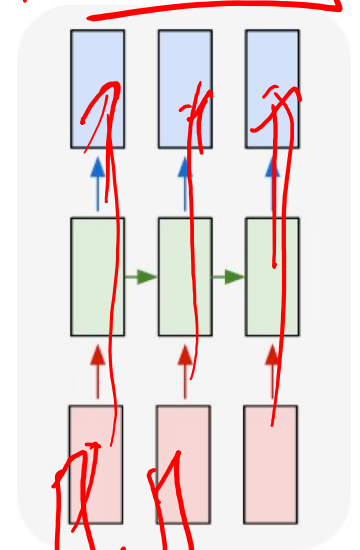


Input: Sequence

Output: Sequence

Example: machine translation, video classification, video captioning, open-ended question answering

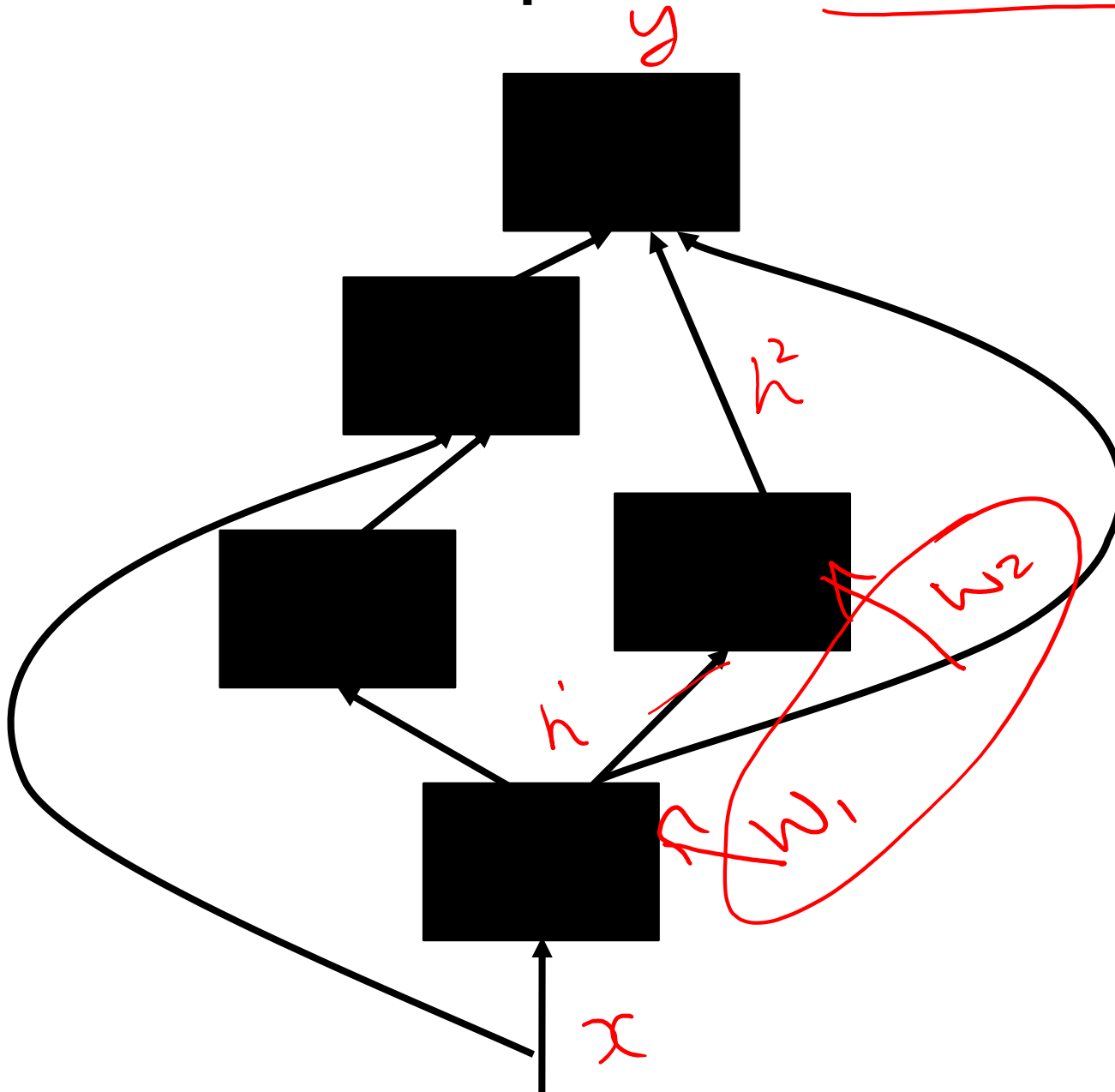
many to many



# 2 Key Ideas

- Parameter Sharing
  - in computation graphs = adding gradients

# Computational Graph



$$f(w_1, w_2)$$

$$w_1(t) \quad w_2(t)$$

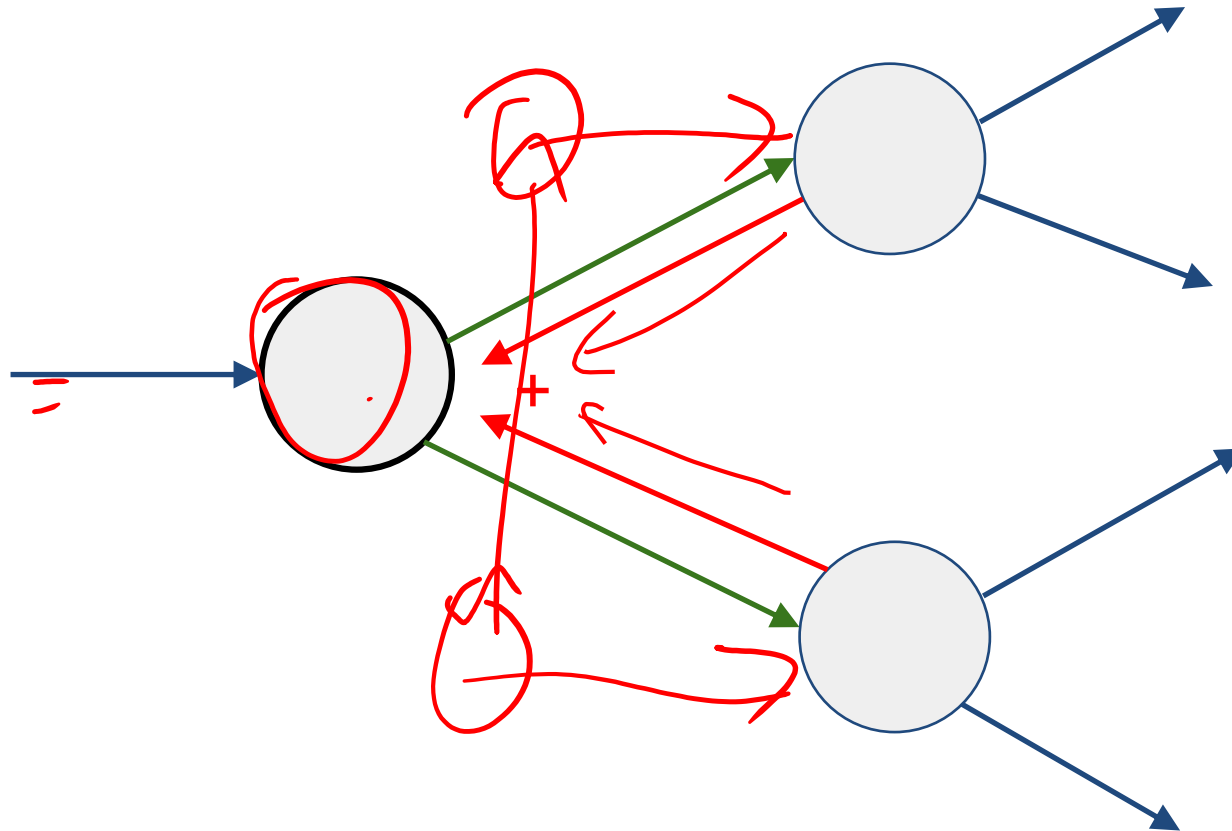
$$\frac{\partial f}{\partial w_1} = \lim_{\delta \rightarrow 0}$$

$$\frac{f(w_1 + \delta, w_2) - f(w_1, w_2)}{\delta}$$

$$\boxed{\frac{df}{dt}} = \boxed{\frac{\partial f}{\partial w_1}} \frac{\partial w_1}{\partial t} + \boxed{\frac{\partial f}{\partial w_2}} \frac{\partial w_2}{\partial t}$$

$$w_1 = t = w_2$$

# Gradients add at branches



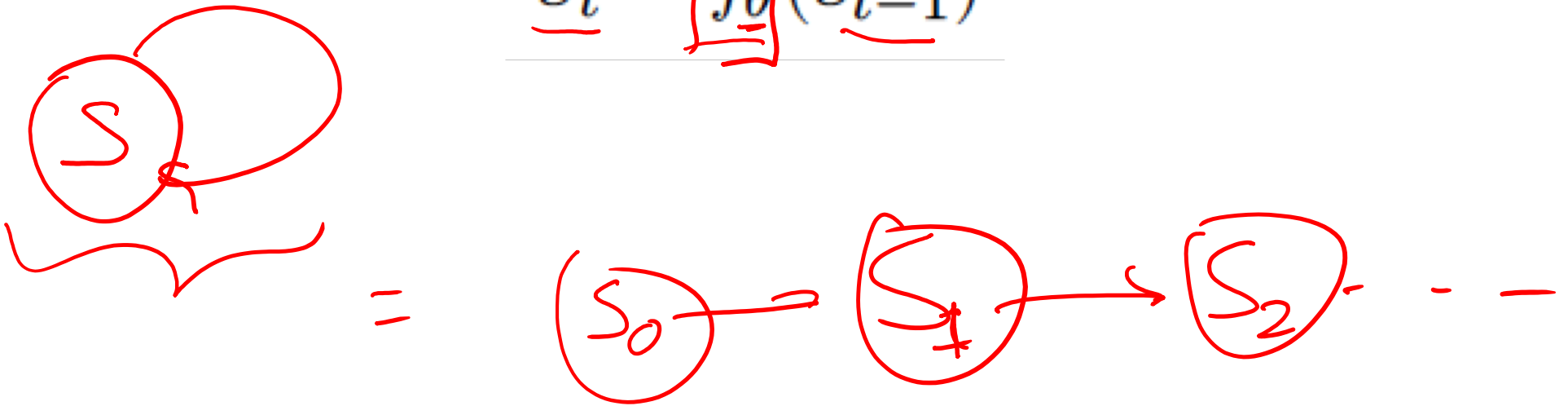
# 2 Key Ideas

- Parameter Sharing
  - in computation graphs = adding gradients
- “Unrolling”
  - in computation graphs with parameter sharing

# How do we model sequences?

- No input

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

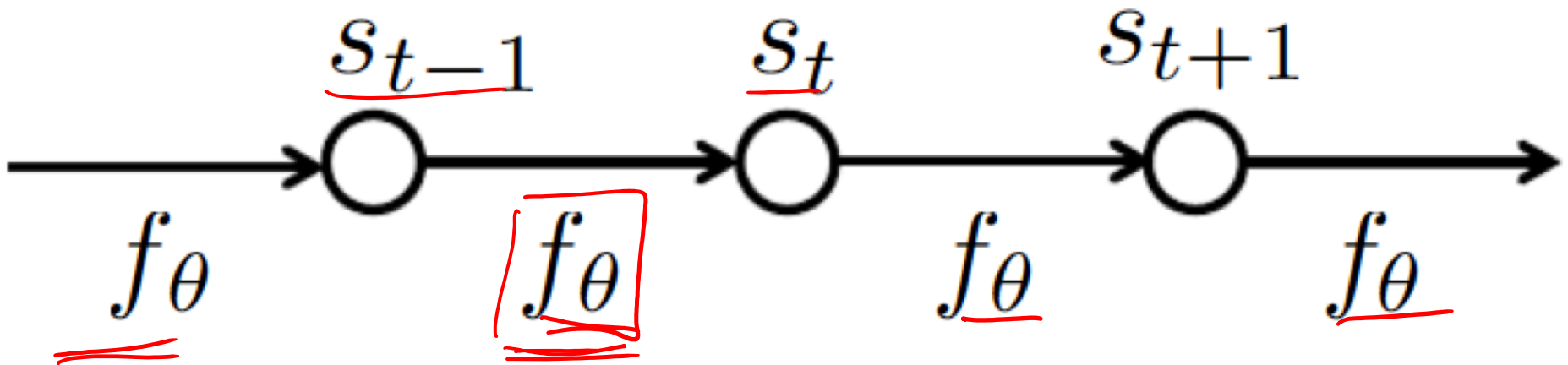




# How do we model sequences?

- No input

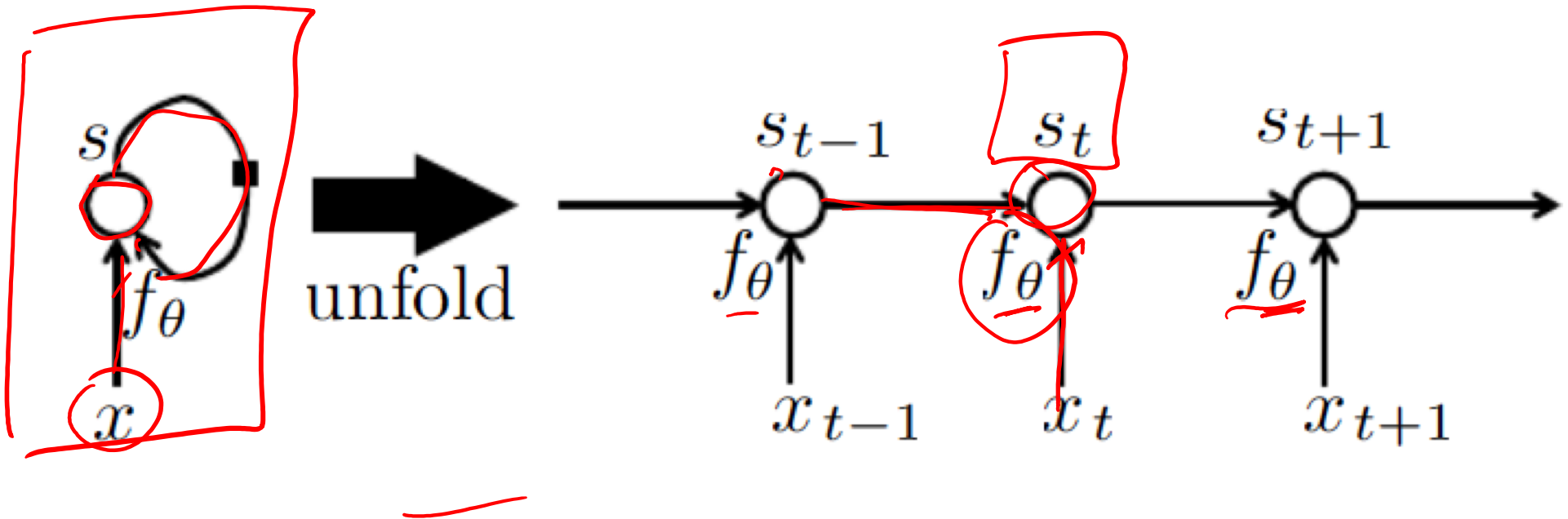
$$s_t = f_{\theta}(s_{t-1})$$



# How do we model sequences?

- With inputs

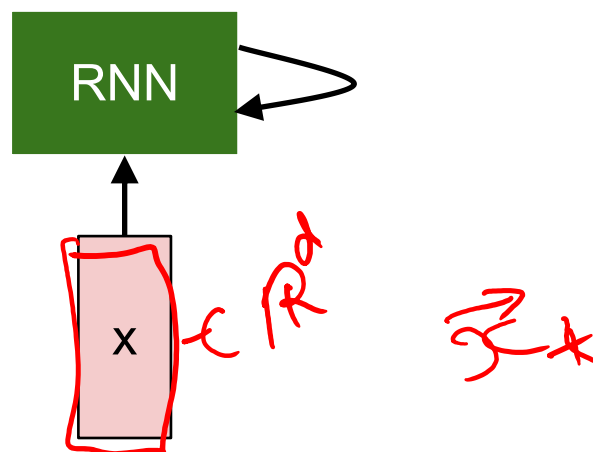
$$s_t = f_{\theta}(s_{t-1}, x_t)$$



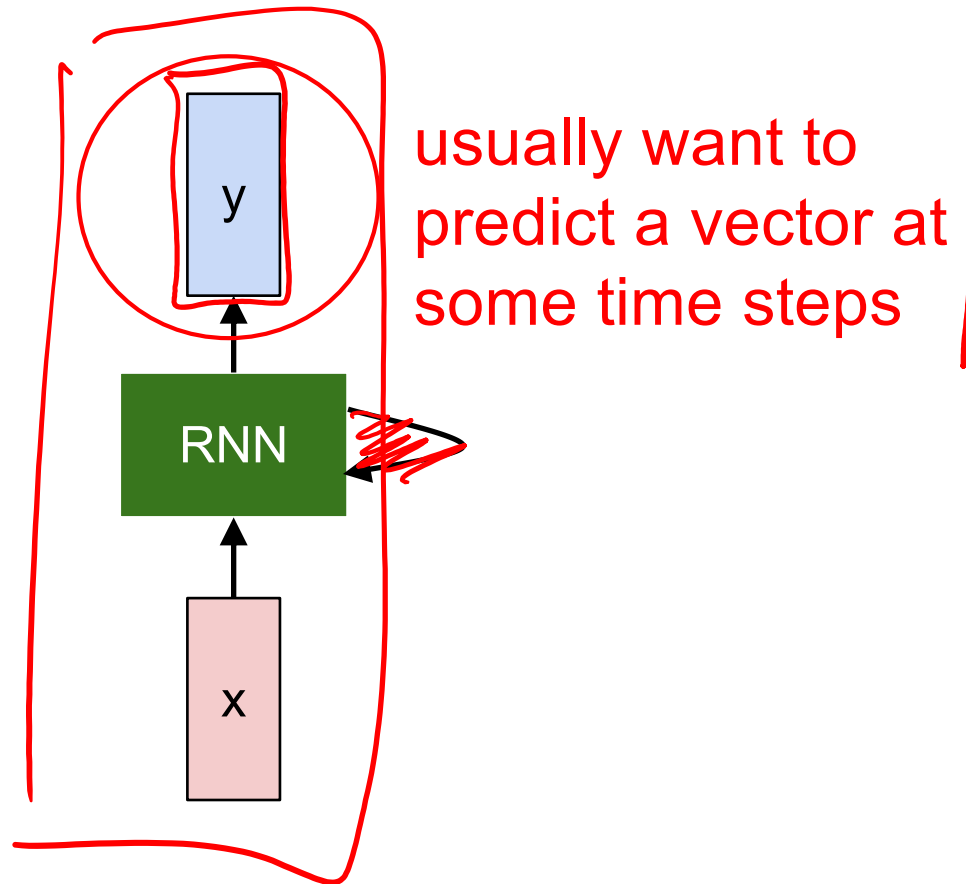
# 2 Key Ideas

- Parameter Sharing
  - in computation graphs = adding gradients
- “Unrolling”
  - in computation graphs with parameter sharing
- Parameter sharing + Unrolling
  - Allows modeling arbitrary sequence lengths!
  - Keeps numbers of parameters in check

# Recurrent Neural Network



# Recurrent Neural Network



# Recurrent Neural Network

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

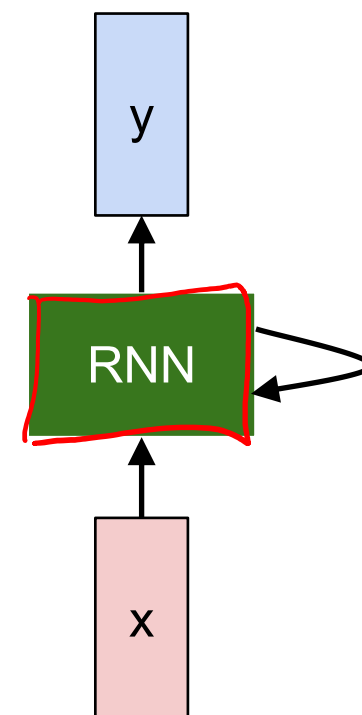
$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state

some function with parameters  $W$

old state

input vector at some time step

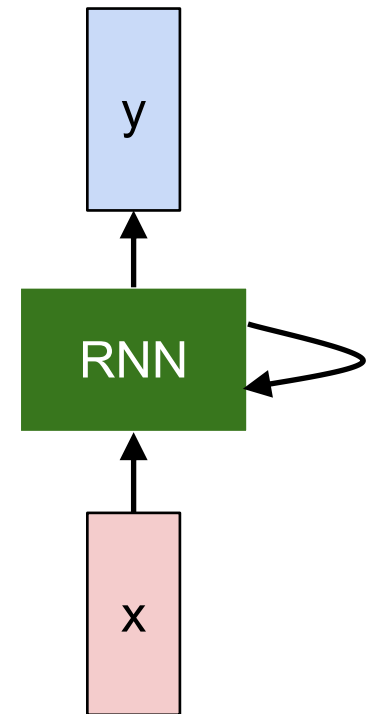


# Recurrent Neural Network

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

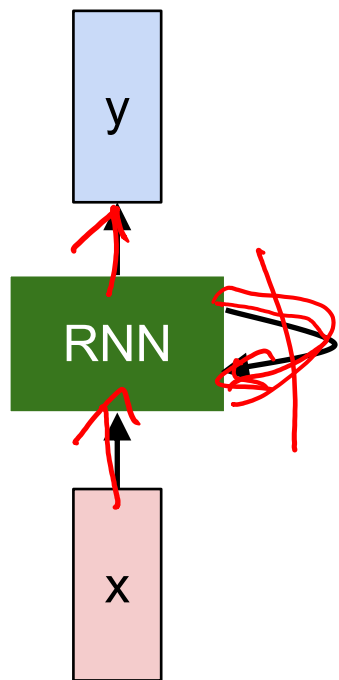
$$h_t = f_W(h_{t-1}, x_t)$$

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  $h$ :



$$y_t = W_{hy} h_t + b_y$$

cell

$$h_t = f_W(\vec{h}_{t-1}, \vec{x}_t)$$

$$W \begin{matrix} d_2 \times (d_1 + d_2) \\ \hline + d_2 \end{matrix}$$

$$[W_{hh} \quad W_{xh}]$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} \vec{x}_t + \vec{b}_h)$$

$$\vec{x}_t \in \mathbb{R}^{d_1}$$

$$\vec{h}_t \in \mathbb{R}^{d_2}$$

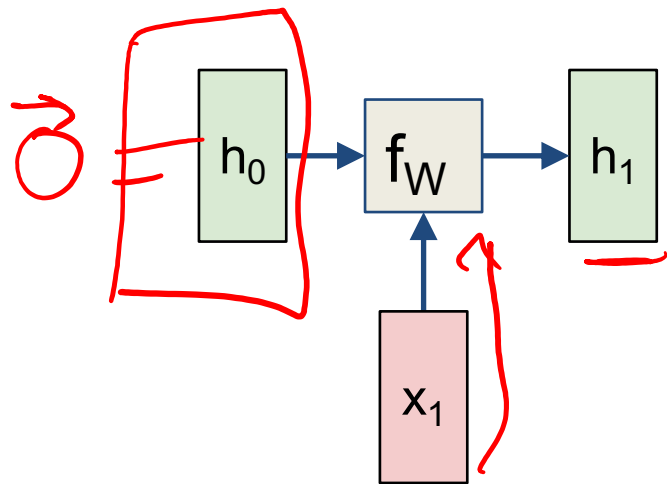
$$h_t = \tanh \left( \begin{matrix} \vec{h}_{t-1} \\ \vec{x}_t \end{matrix} + \vec{b} \right)$$

Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

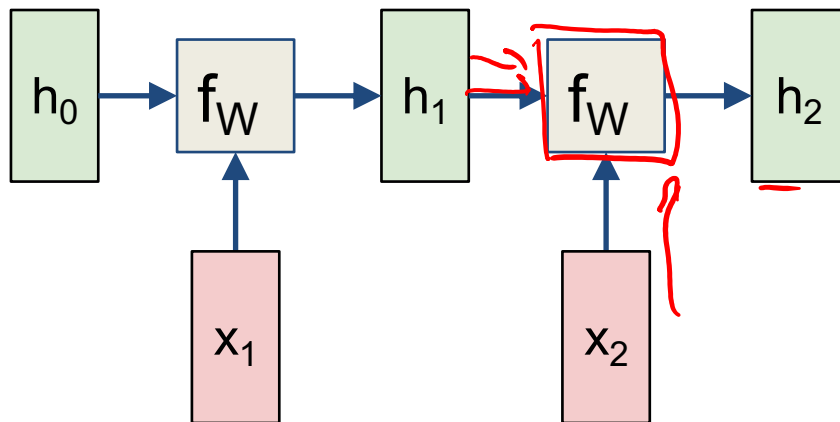
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



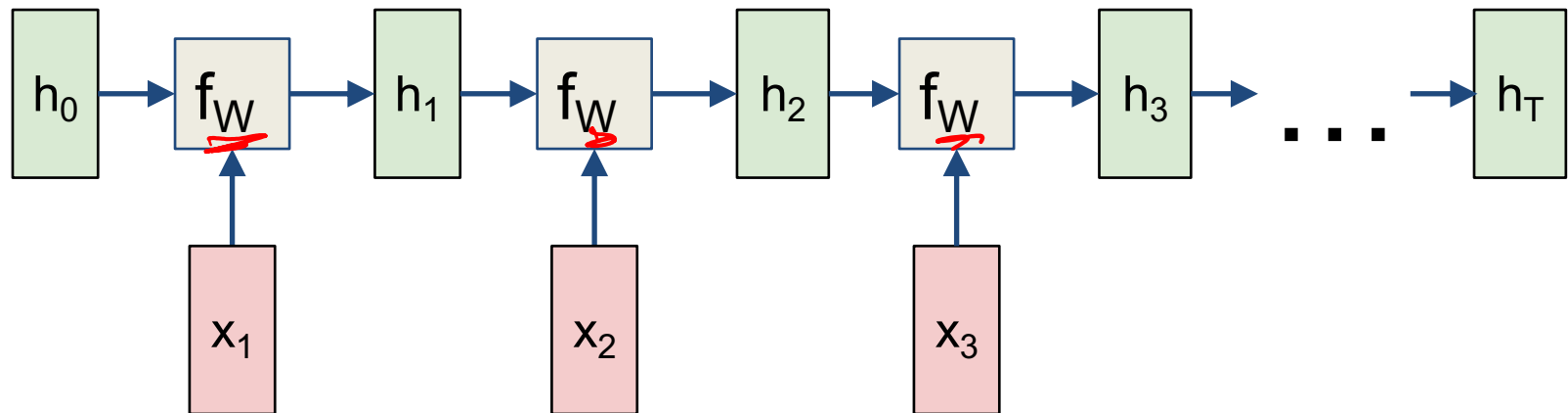
# RNN: Computational Graph



# RNN: Computational Graph

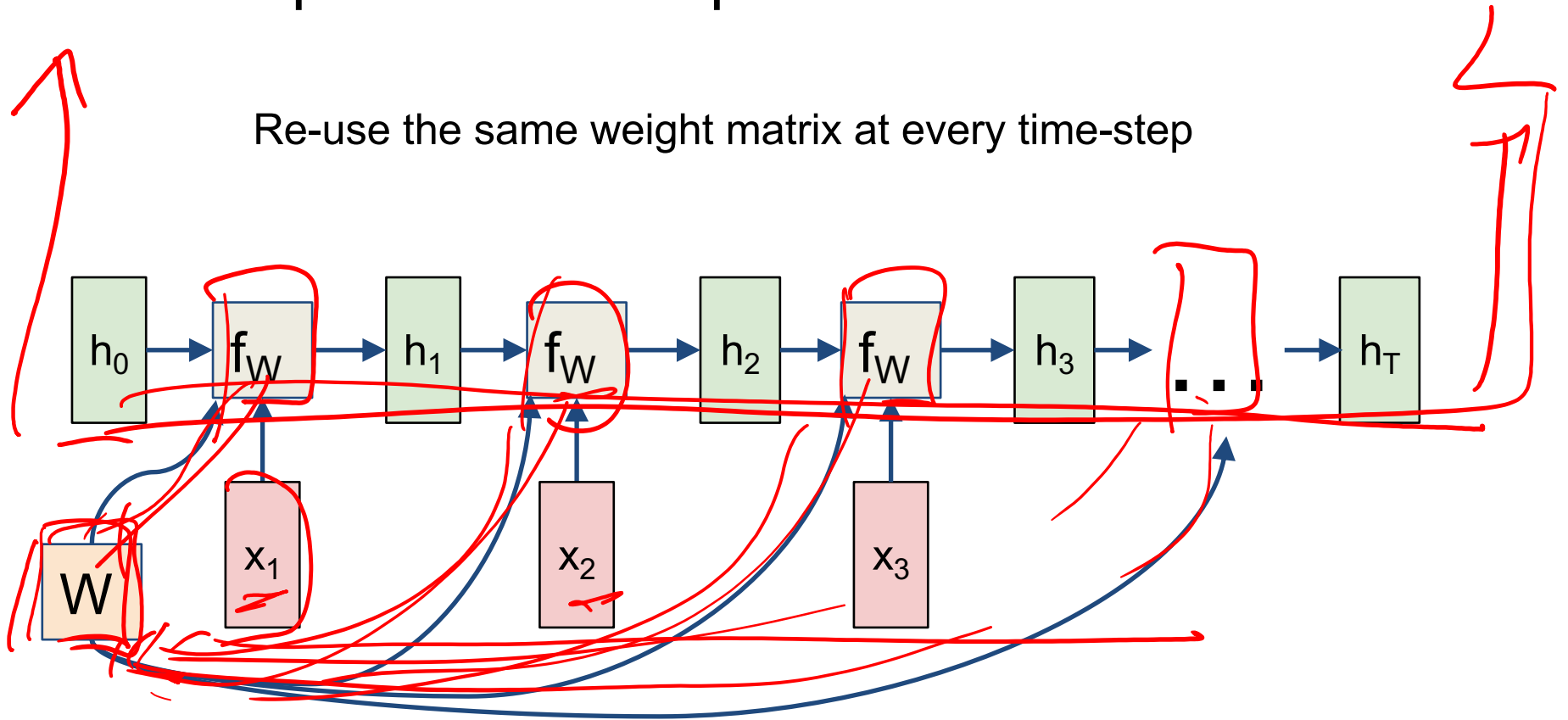


# RNN: Computational Graph

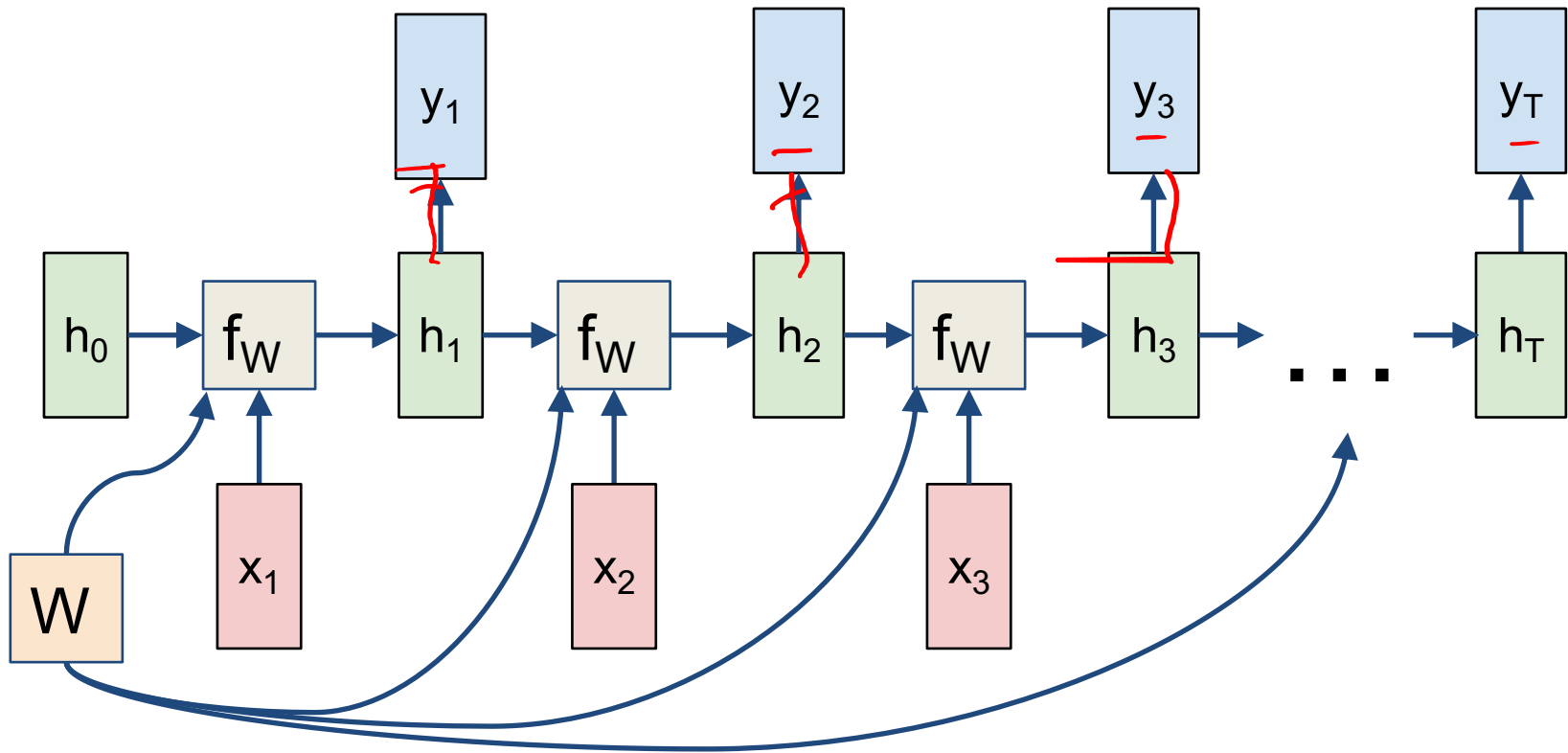


# RNN: Computational Graph

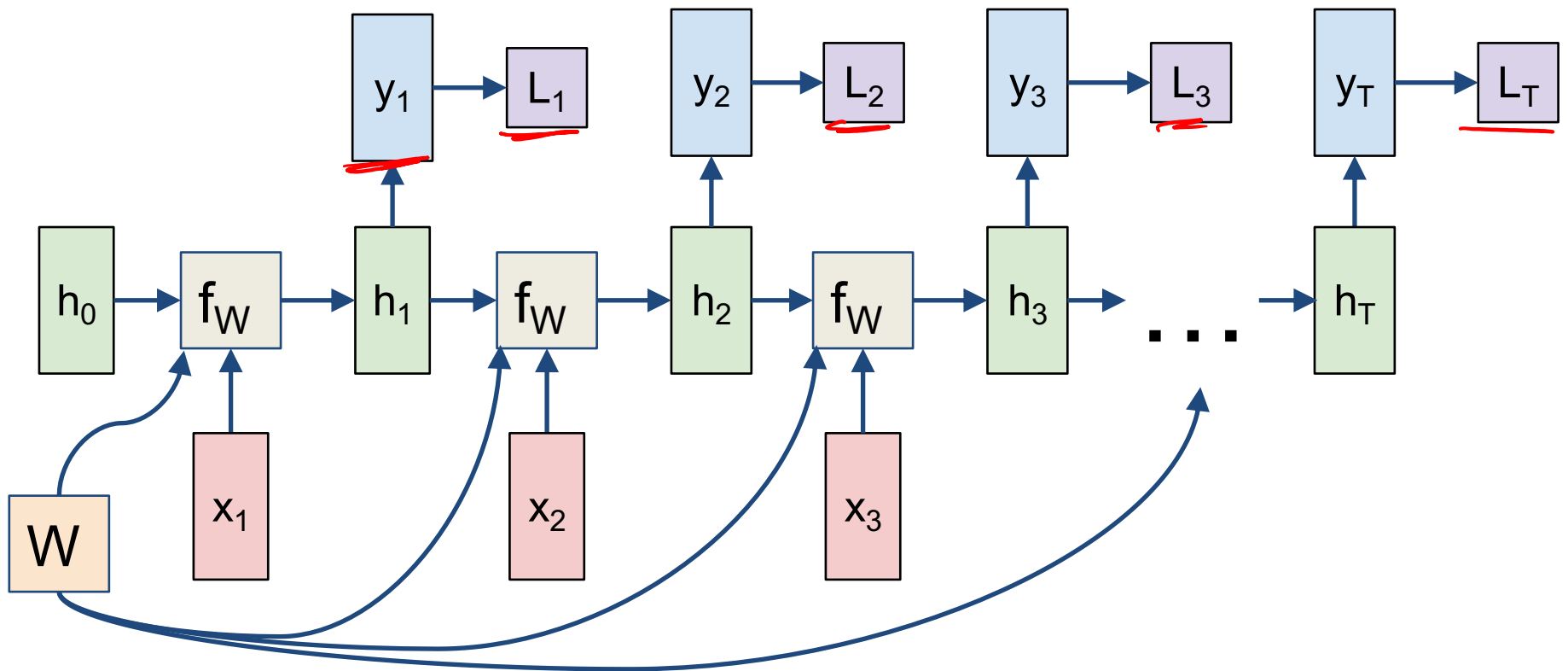
Re-use the same weight matrix at every time-step



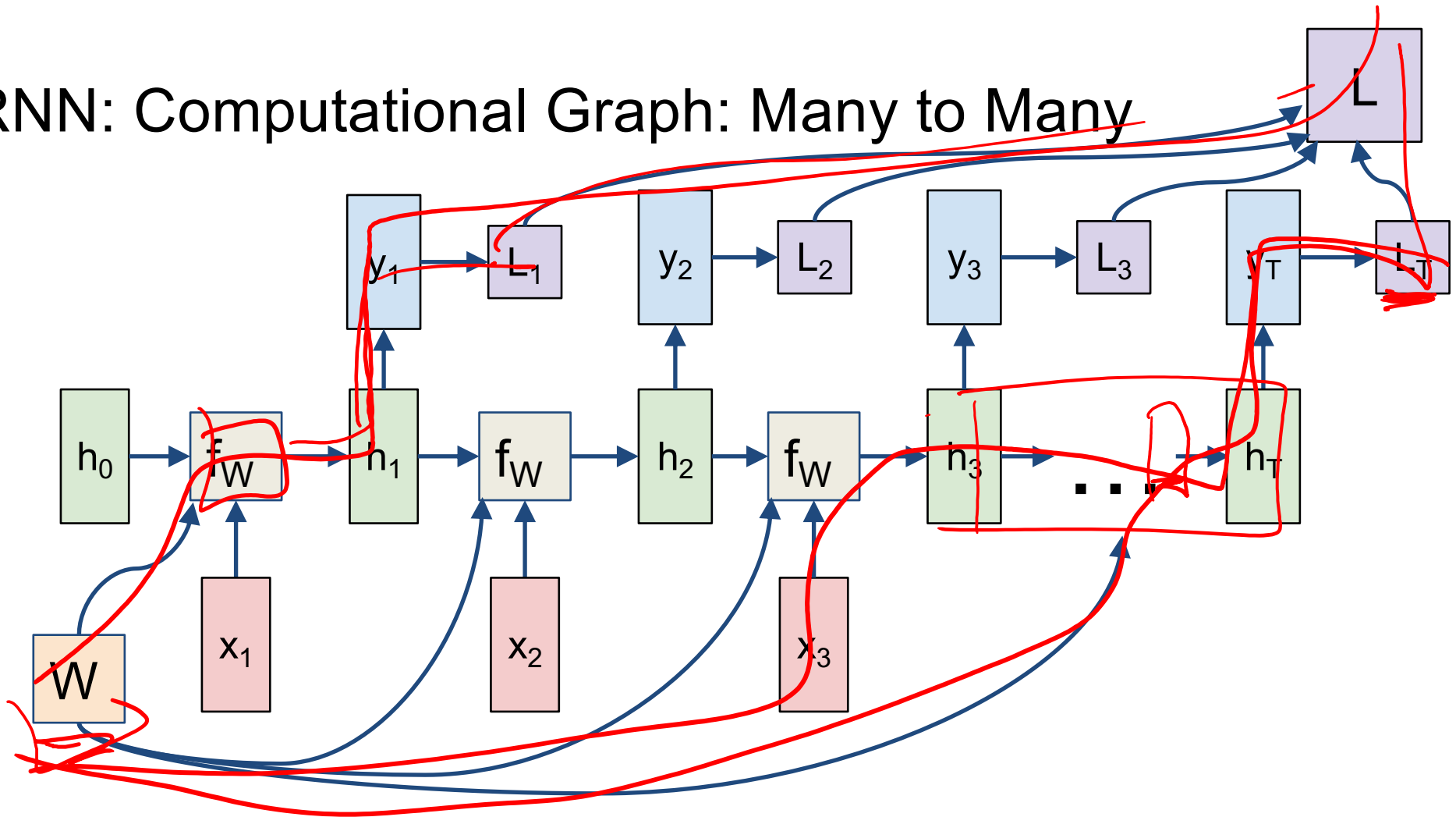
# RNN: Computational Graph: Many to Many



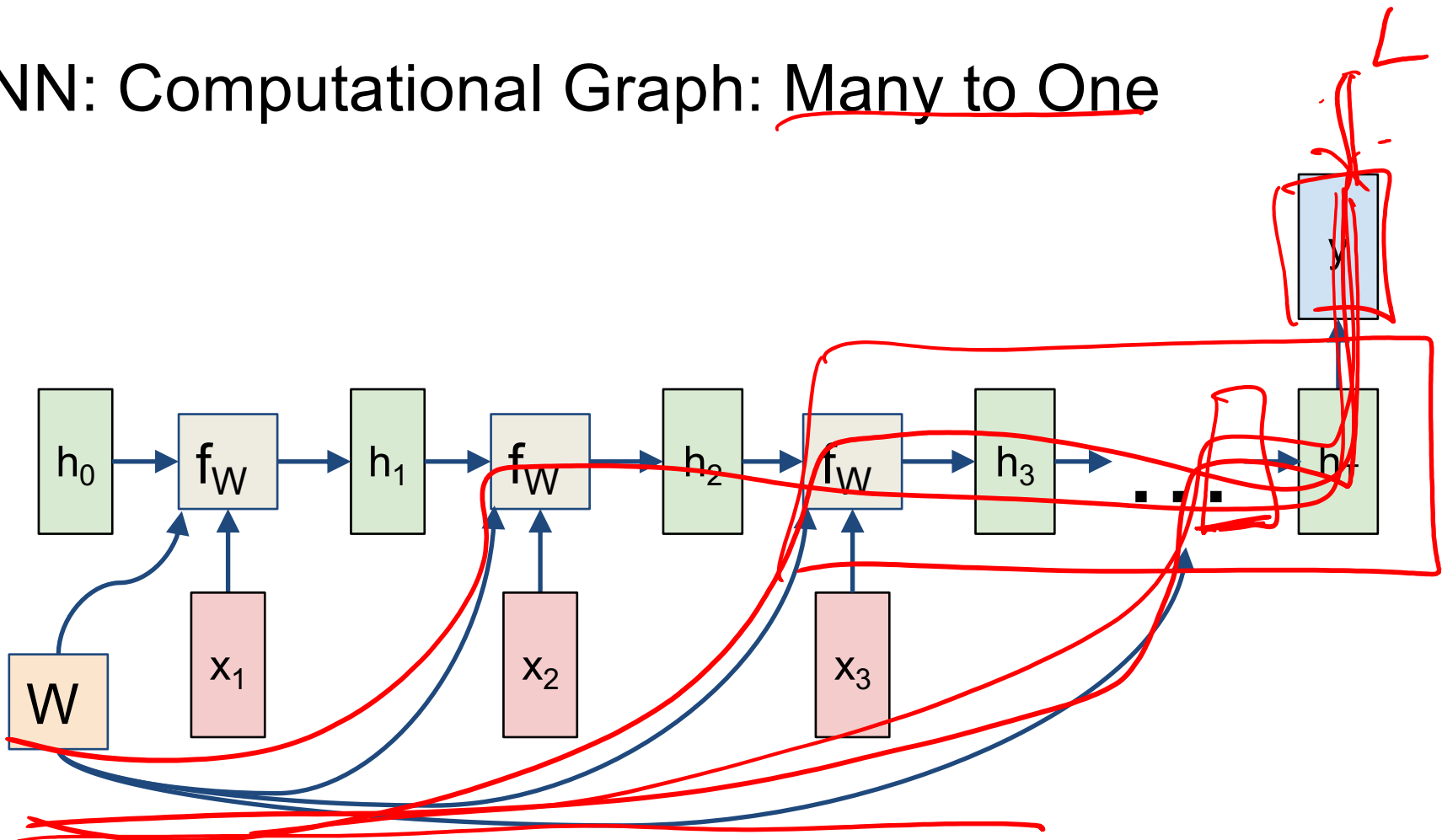
# RNN: Computational Graph: Many to Many



# RNN: Computational Graph: Many to Many

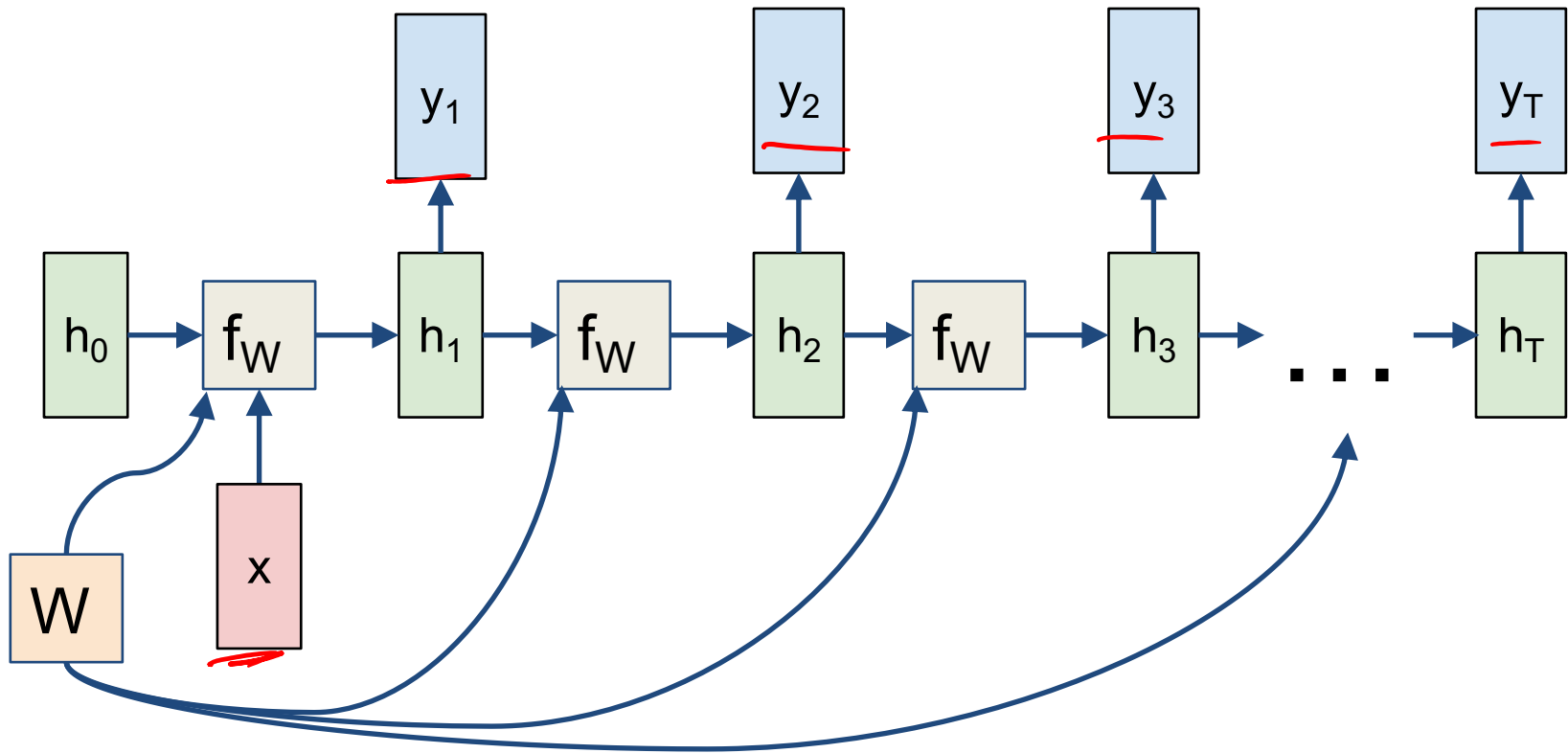


# RNN: Computational Graph: Many to One



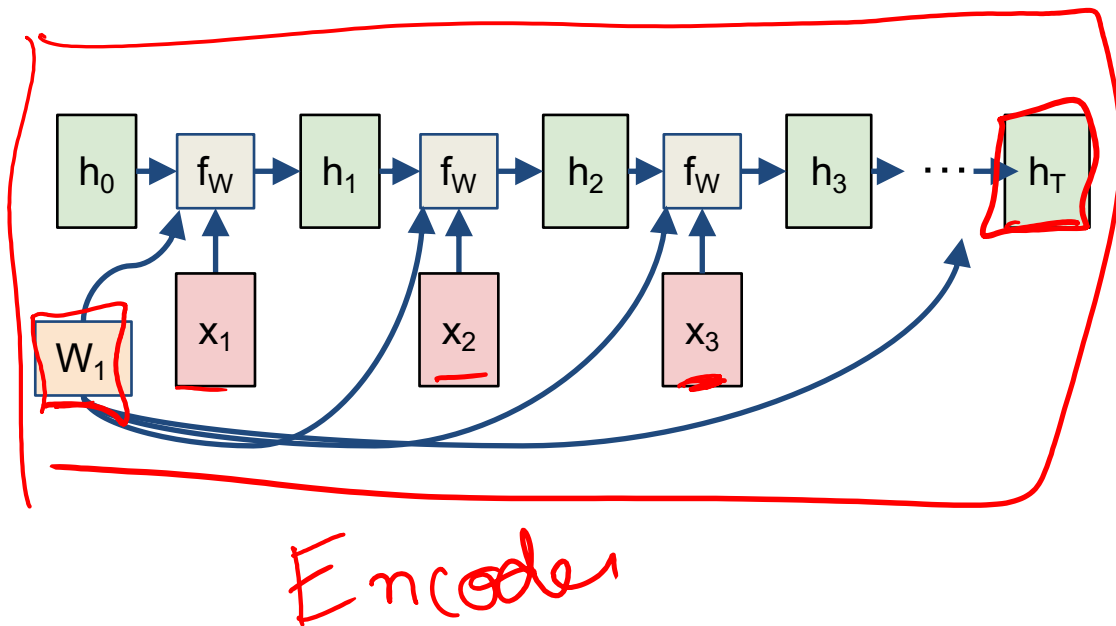


# RNN: Computational Graph: One to Many



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

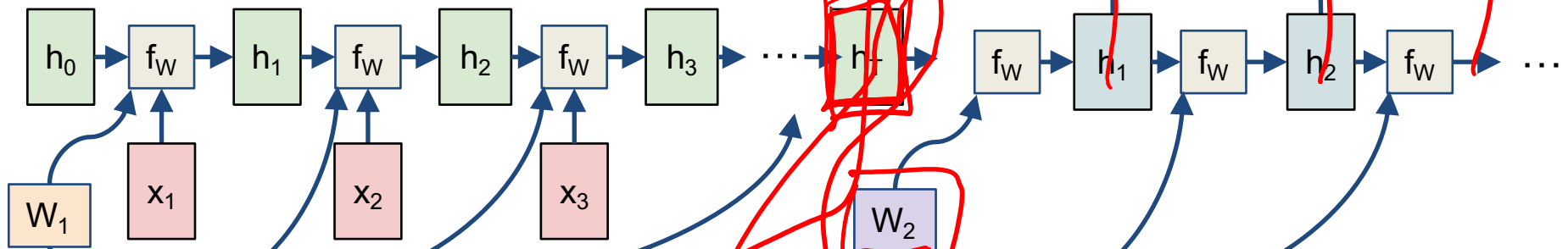
**Many to one:** Encode input sequence in a single vector



# 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



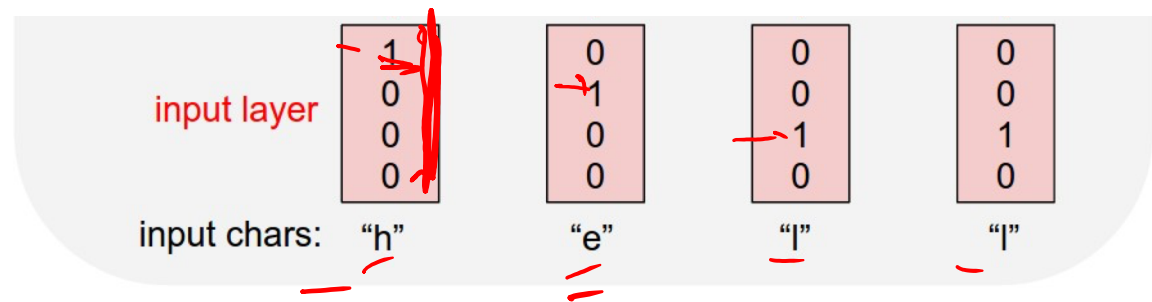
Decode

$$P_{\theta}(x_1, \dots, x_T) = P(x_1)P(x_2|x_1) \dots P(x_k|x_1, \dots, x_{k-1}) \dots$$

**Example:**  
**Character-level**  
**Language Model**

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
"hello"

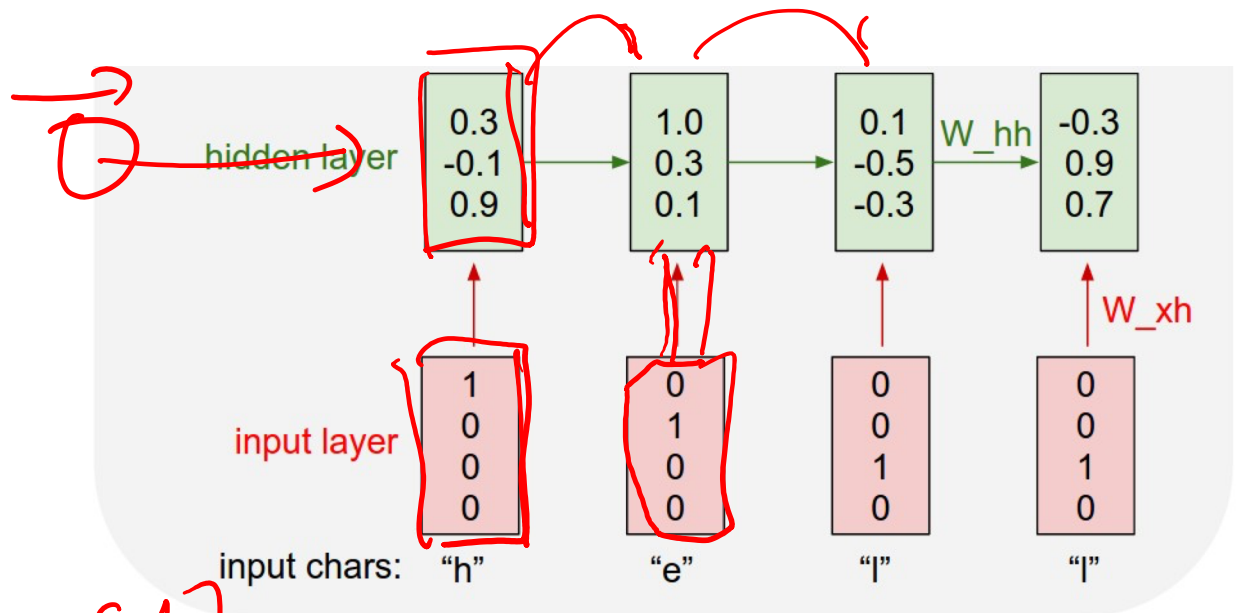


# Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
"hello"

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



$$\begin{bmatrix} | & | & | & | \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} | \\ \times \\ | \\ \times \end{matrix} \begin{matrix} | \\ | \\ | \\ | \end{matrix}$$

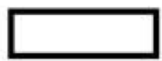
$3 \times 4$

# Distributed Representations Toy Example

- Local vs Distributed

(a)

no pattern



# Distributed Representations Toy Example

- Can we interpret each dimension?

(a)

no pattern



(b)

no pattern



*vertical*

*horizontal*


*rectangle*

*ellipse*




# Power of distributed representations!

Local


$$\bullet \bullet \circ \bullet = VR + HR + HE = ?$$

Distributed


$$\bullet \bullet \circ \bullet = V + H + E \approx \bigcirc$$

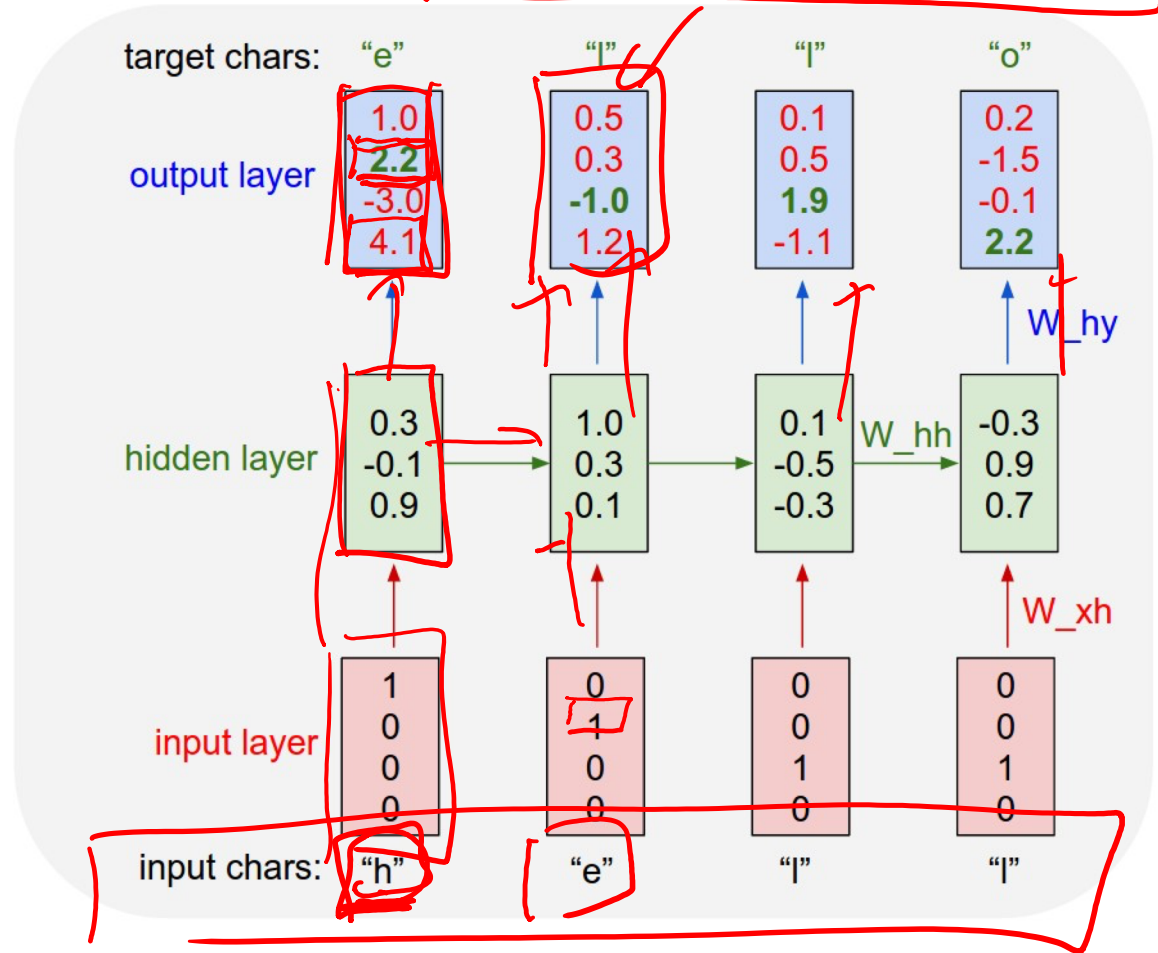


$$\max_w \log P(x_1 \dots x_t | w) = \sum_t \log P(x_t | x_{1:t-1}, w)$$

# Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
"hello"

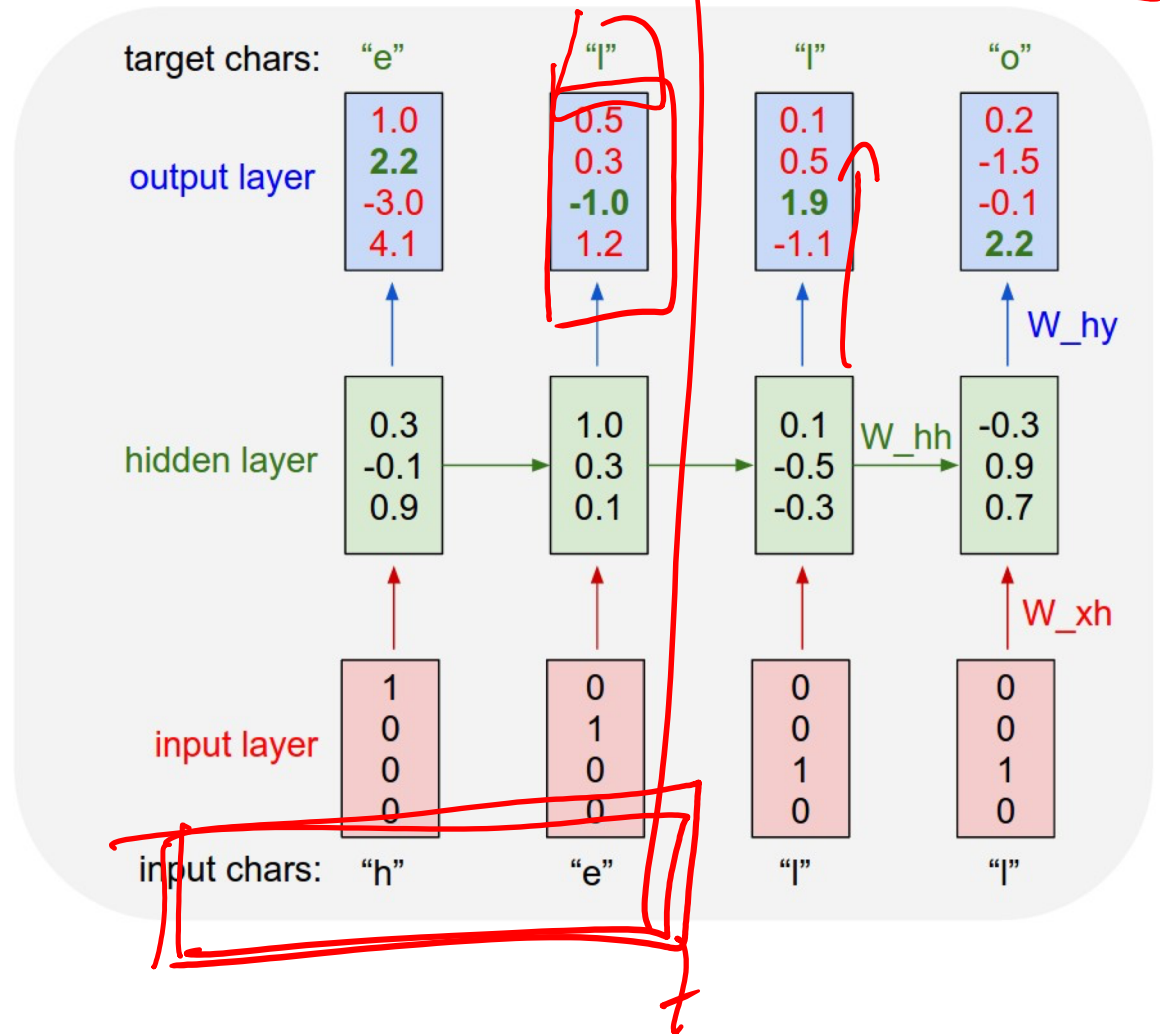


# Training Time: MLE / "Teacher Forcing"

## Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
"hello"

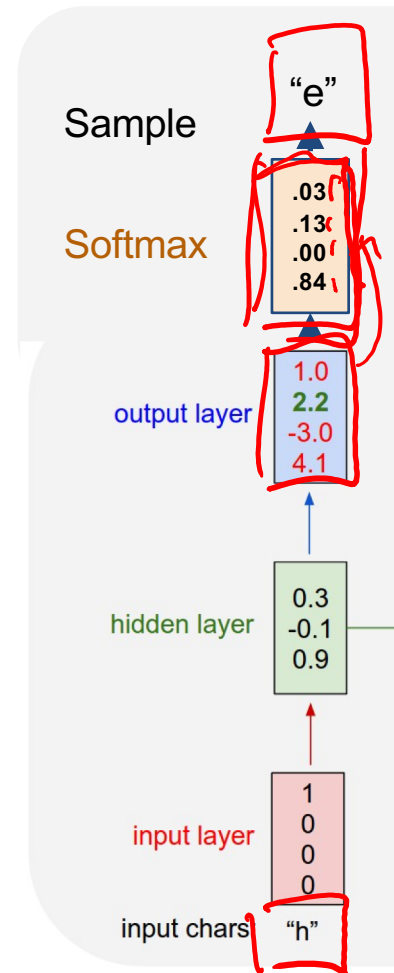


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

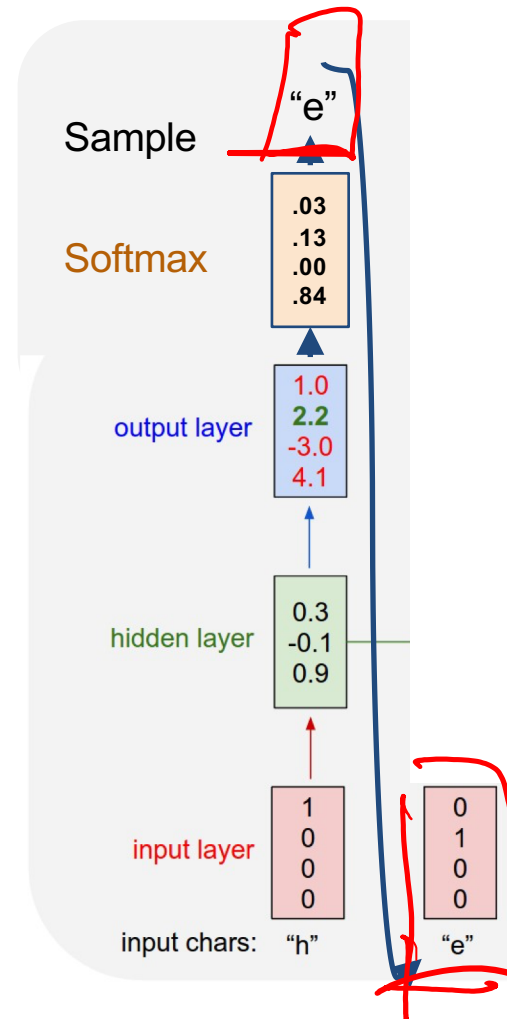


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

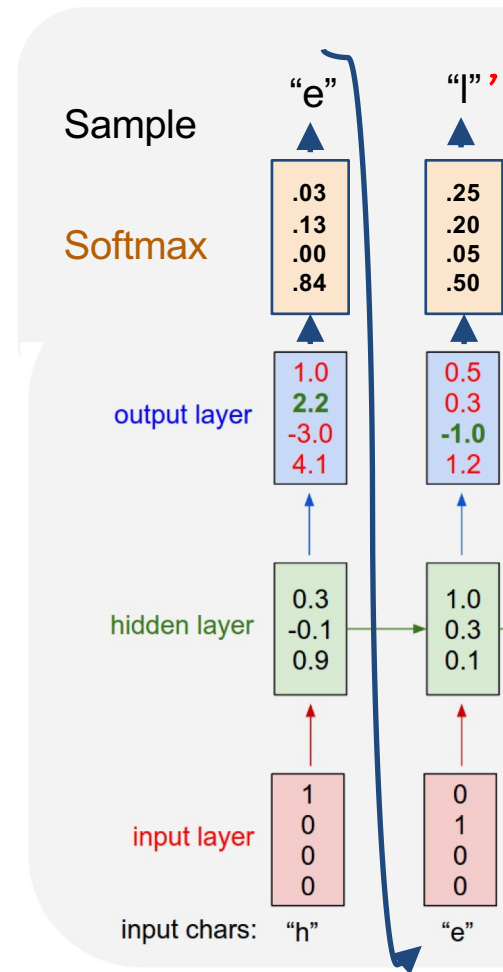


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model



# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

