

CS 4803 / 7643: Deep Learning

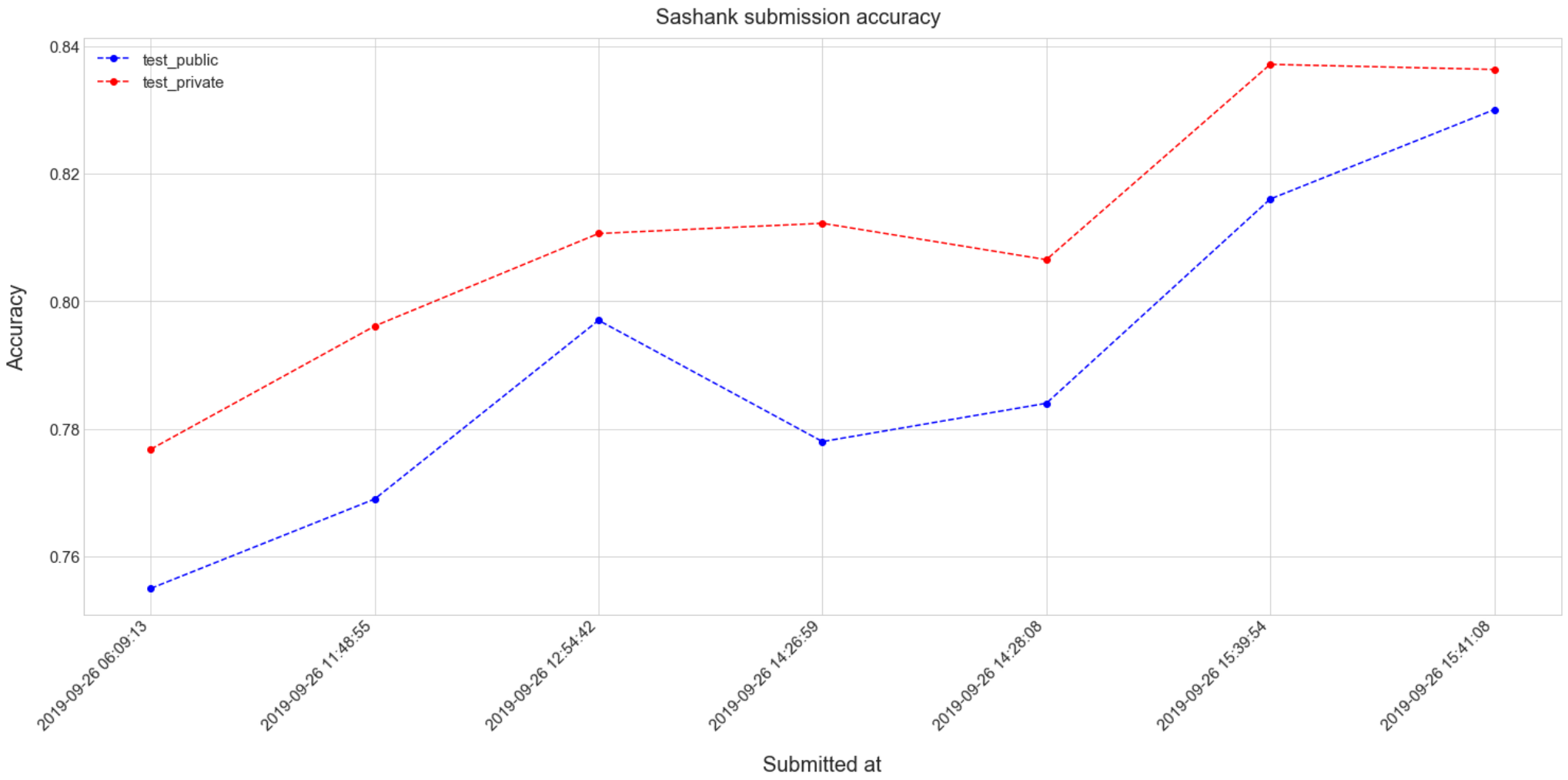
Topics:

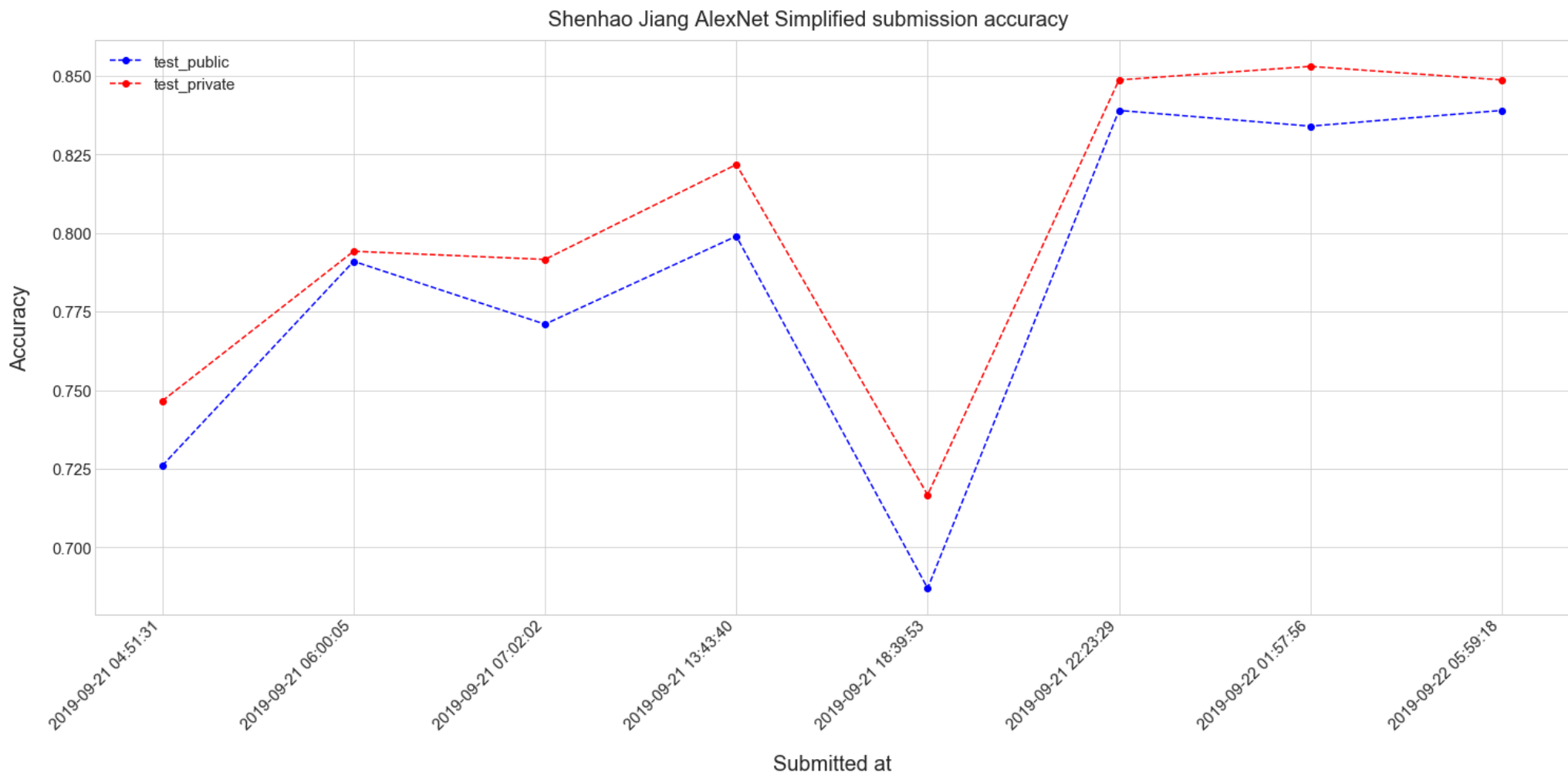
- (Finish) Convolutional Neural Networks
 - Transposed convolutions
- Recurrent Neural Networks (RNNs)]

Dhruv Batra
Georgia Tech

Administrativa

- HW1 Challenge Final Analysis
 - <https://evalai.cloudcv.org/web/challenges/challenge-page/431/leaderboard/1200>
 - Qualitative Trends
- HW2 Reminder
 - Due: 10/10, 11:55pm
 - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/assets/hw2.pdf





Plan for Today

- (Finish) Convolutional Neural Networks
 - Transposed convolutions
- Recurrent Neural Networks (RNNs)
 - A new model class
 - Learning: BackProp Through Time (BPTT)

Other Computer Vision Tasks

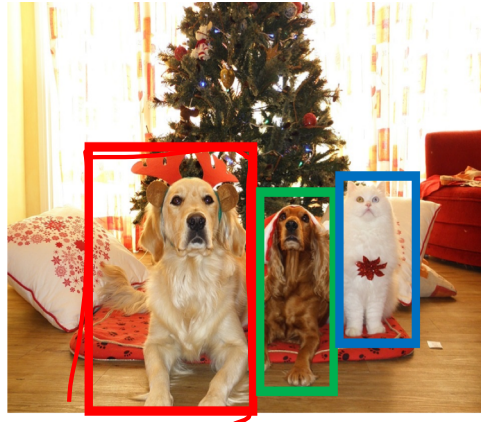
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

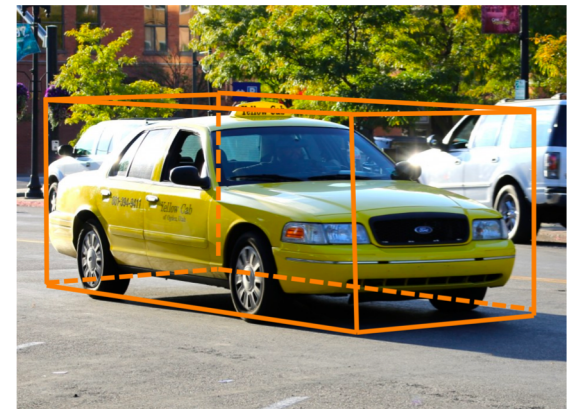
2D Object Detection



DOG, DOG, CAT

Object categories +
2D bounding boxes

3D Object Detection



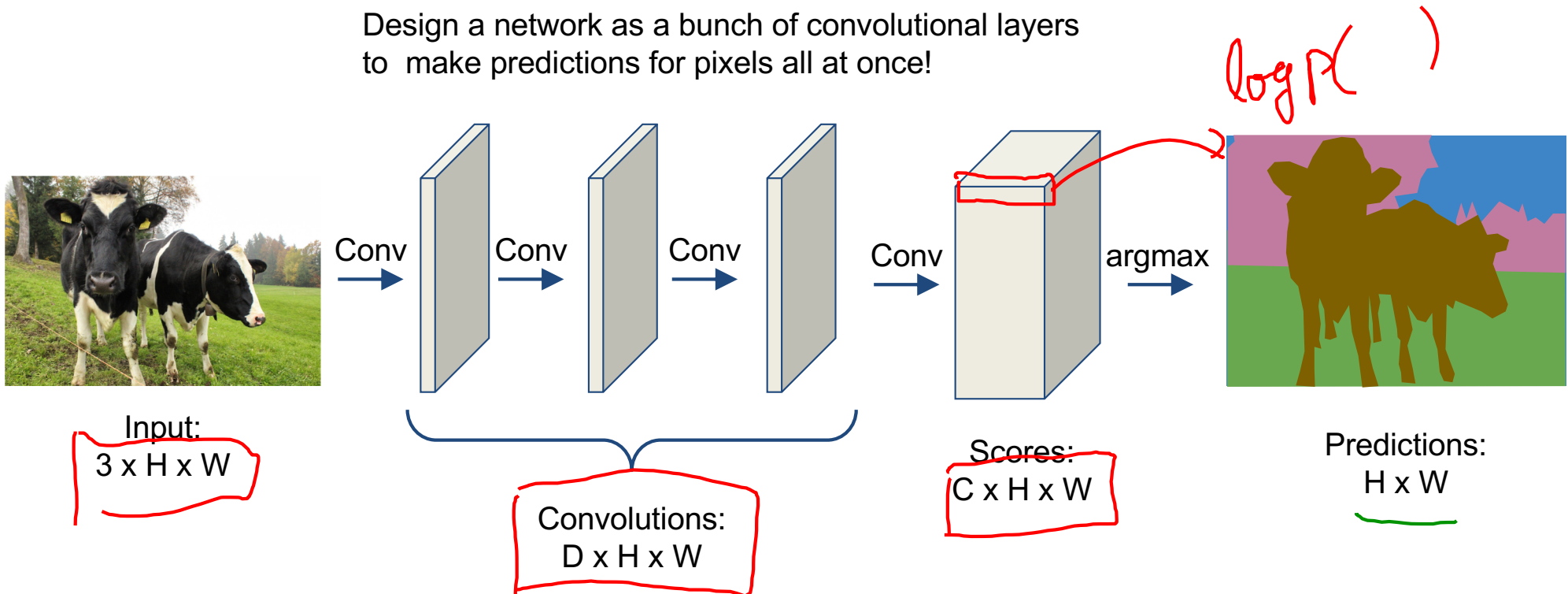
Car

Object categories +
3D bounding boxes

[This image](#) is [CC0 public domain](#)

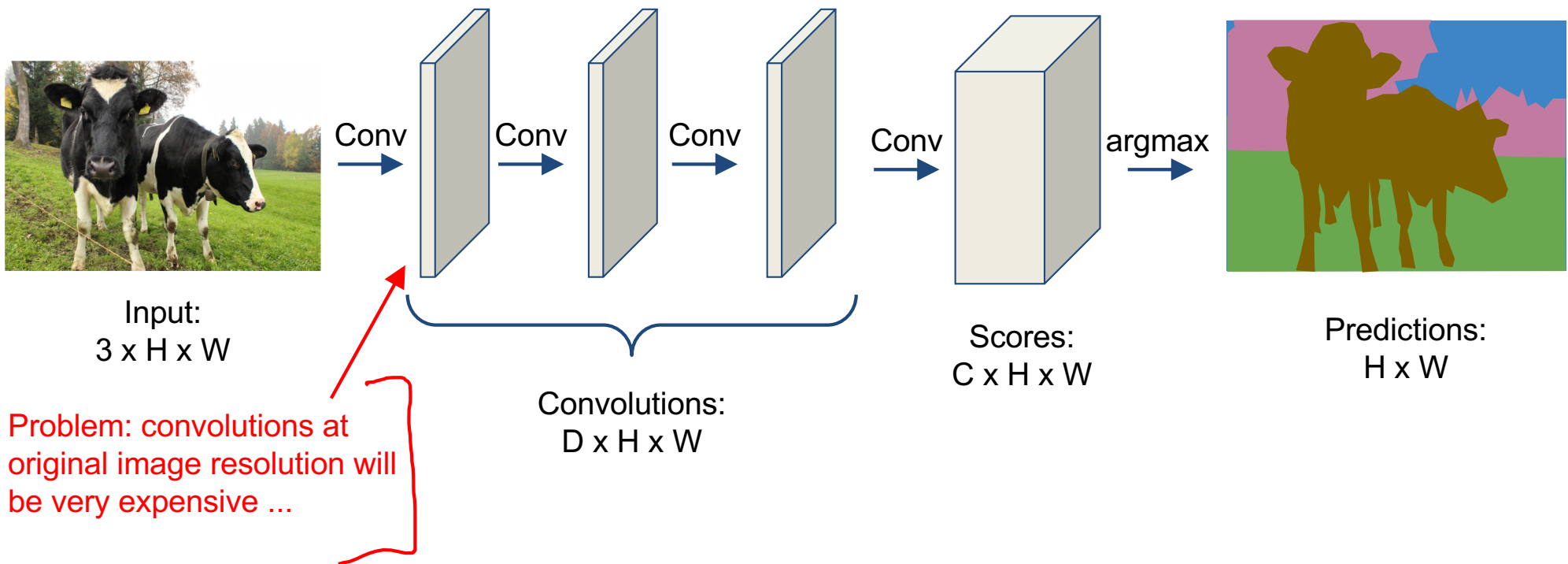
Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

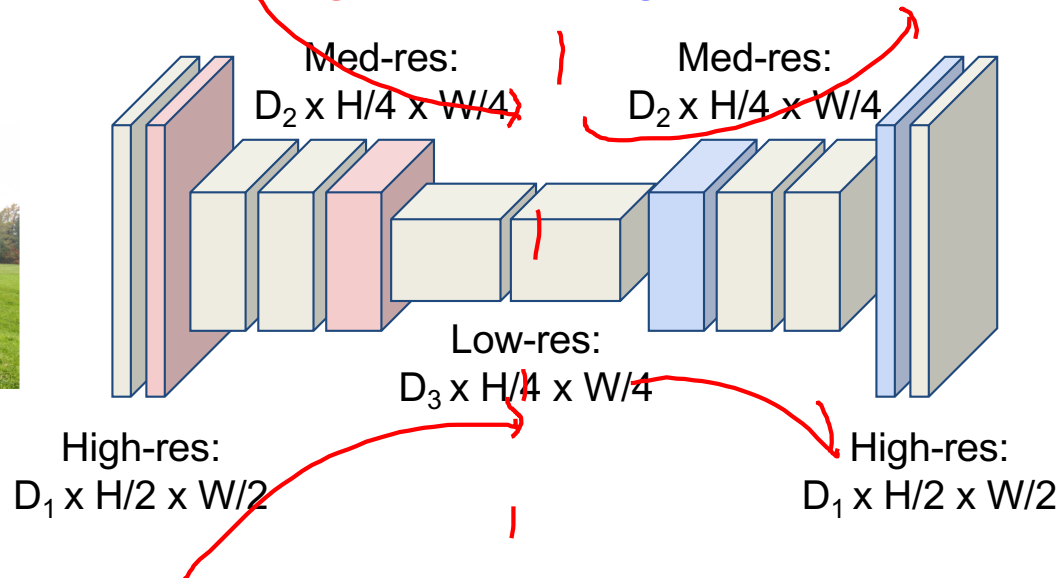


Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

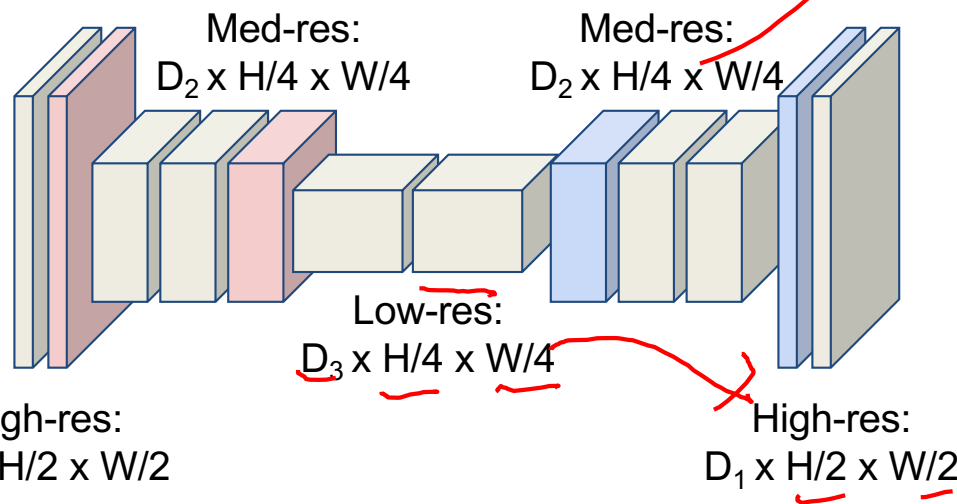
Semantic Segmentation Idea: Fully Convolutional

Downsampling:
Pooling, strided convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Upsampling:
???



Predictions:
 $H \times W$

In-Network upsampling: “Unpooling”

Nearest Neighbor

1	2
3	4

Input: 2 x 2



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output: 4 x 4

“Bed of Nails”

<u>1</u>	2
3	4

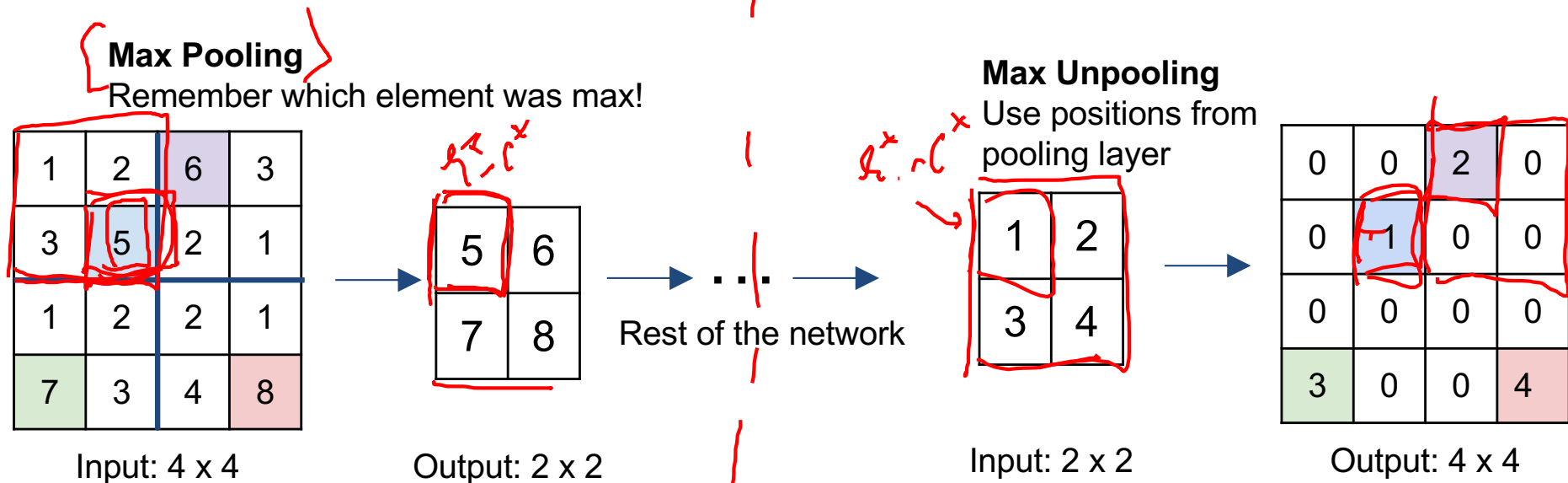
Input: 2 x 2



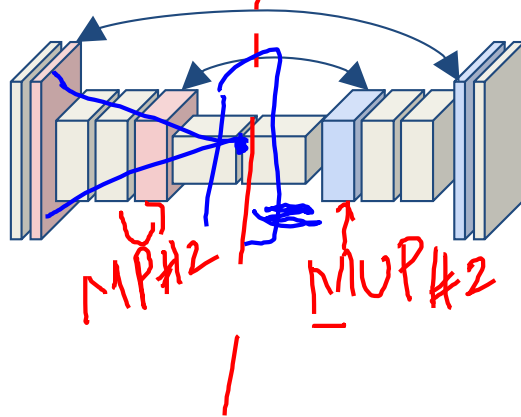
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output: 4 x 4

In-Network upsampling: "Max Unpooling"



Corresponding pairs of downsampling and upsampling layers

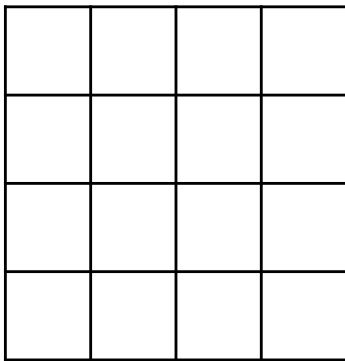


→ Transposed Convolutions

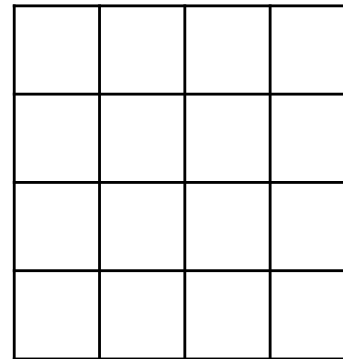
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

Learnable Upsampling: Transpose Convolution

Recall: Typical 3 x 3 convolution, stride 1 pad 1



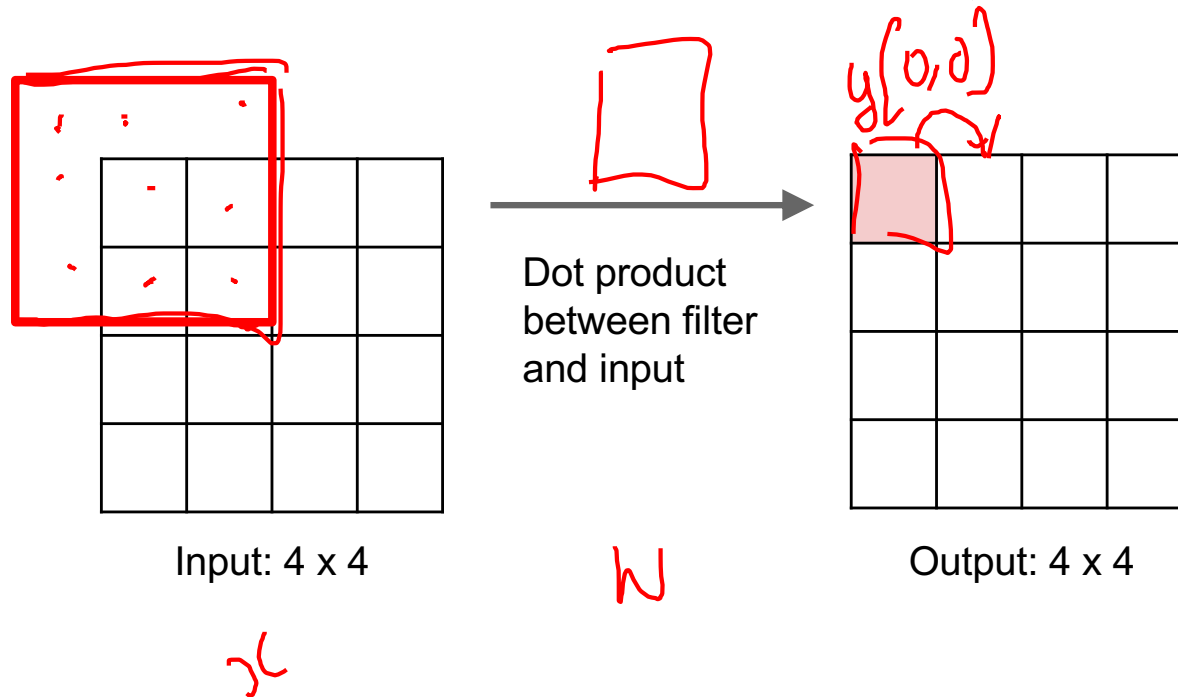
Input: 4 x 4



Output: 4 x 4

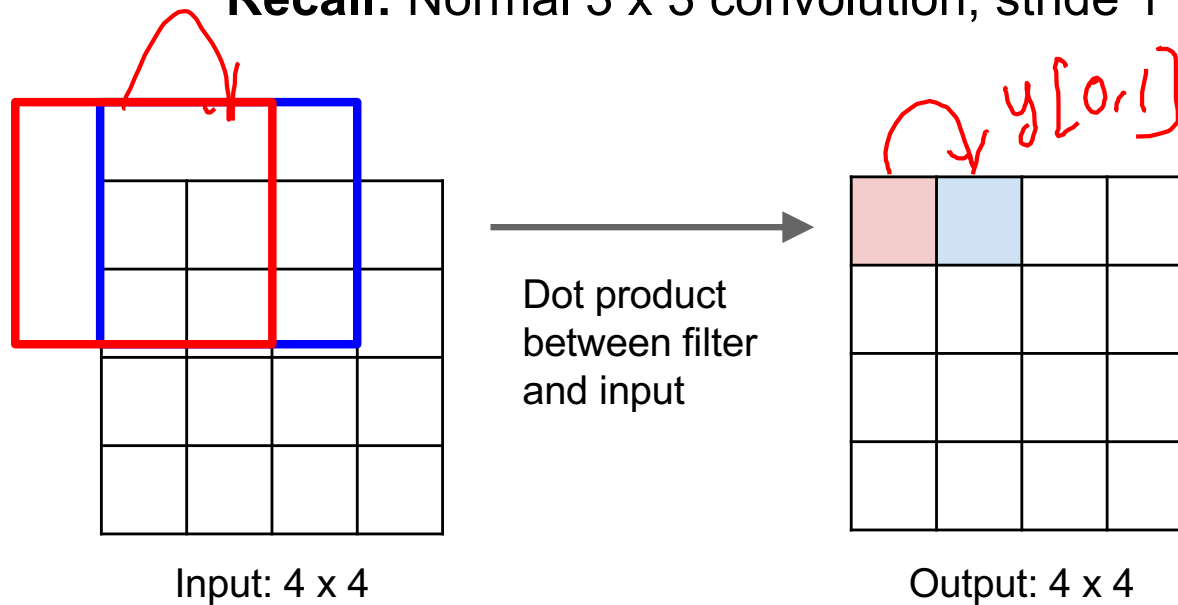
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1



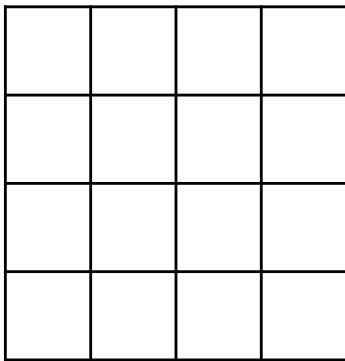
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1

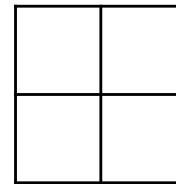


Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



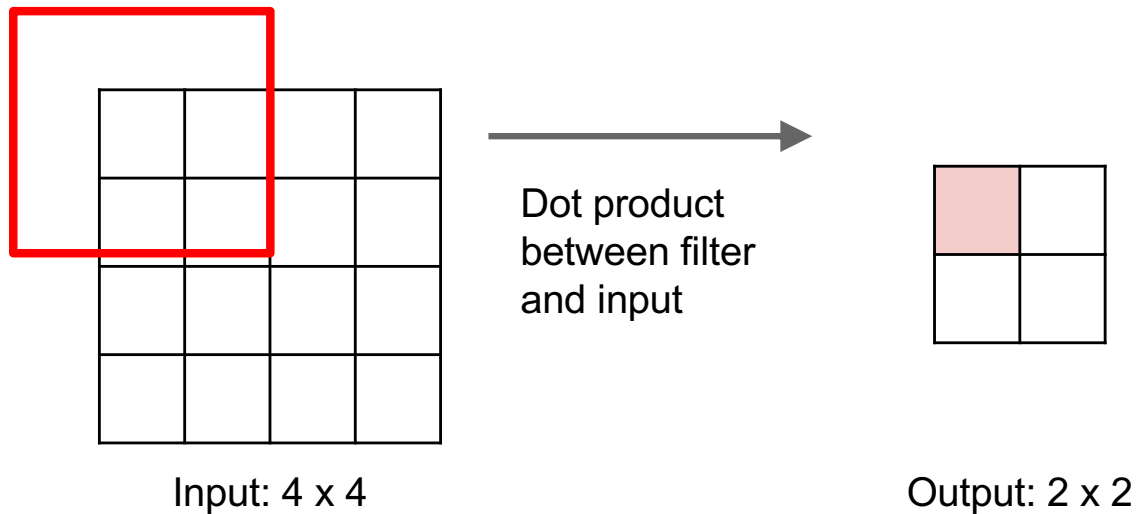
Input: 4 x 4



Output: 2 x 2

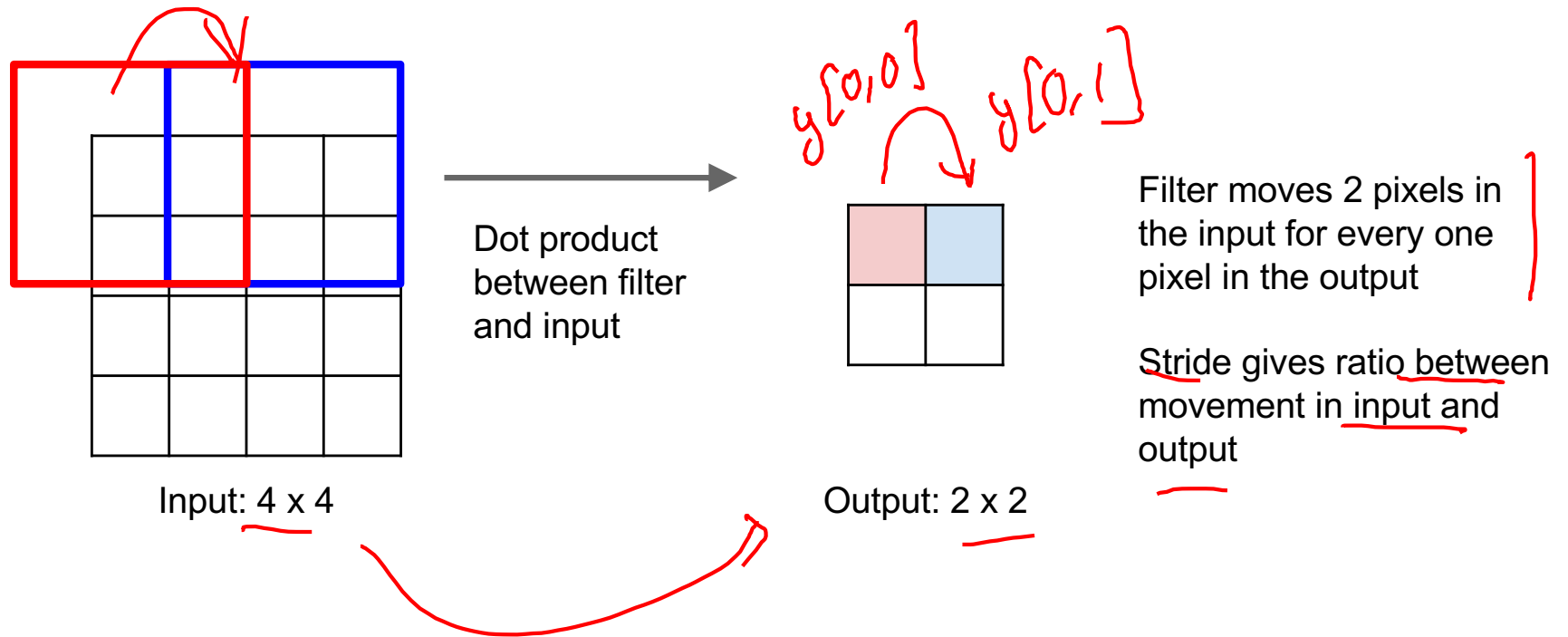
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



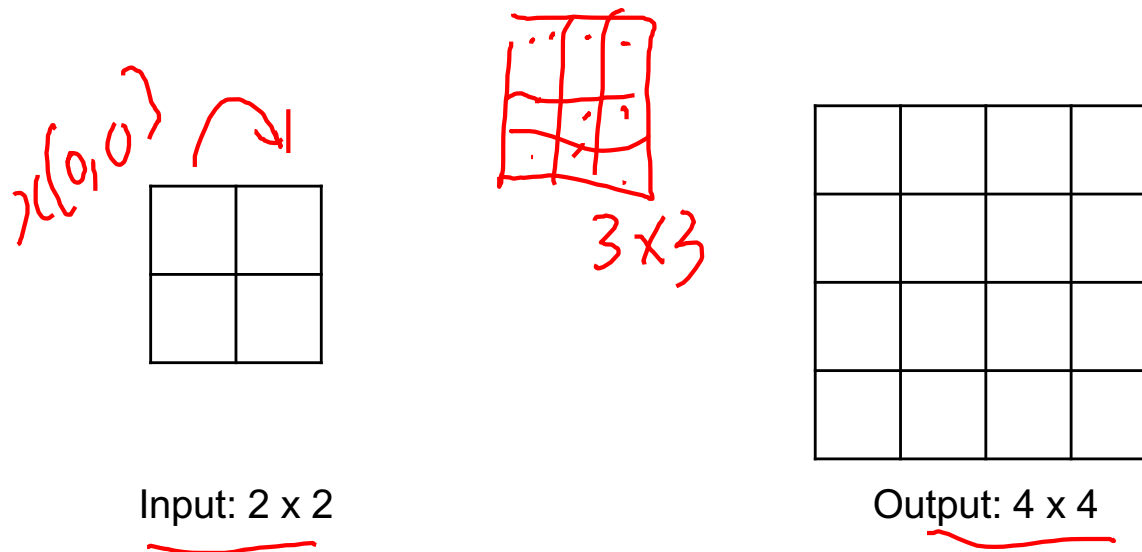
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



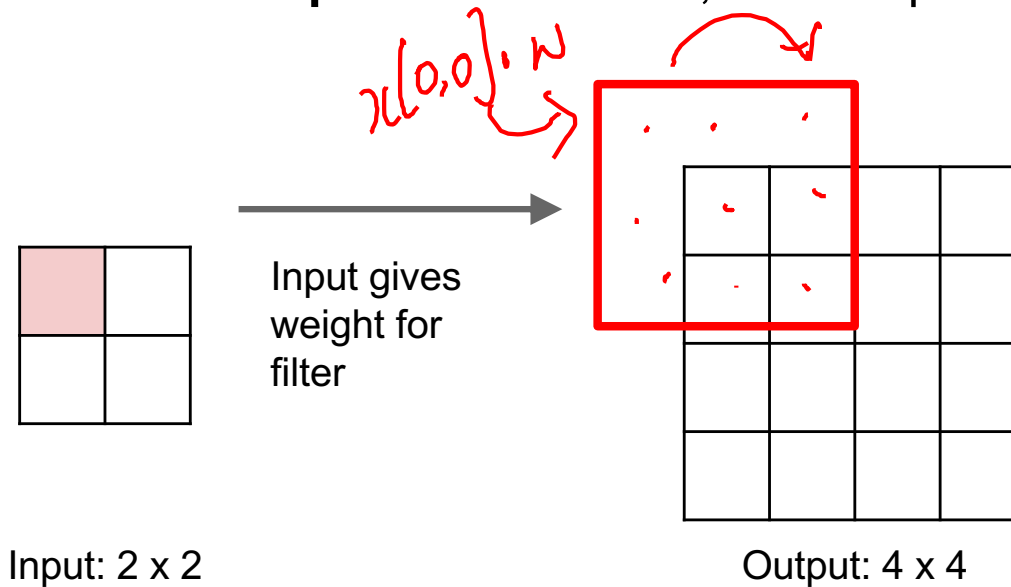
Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

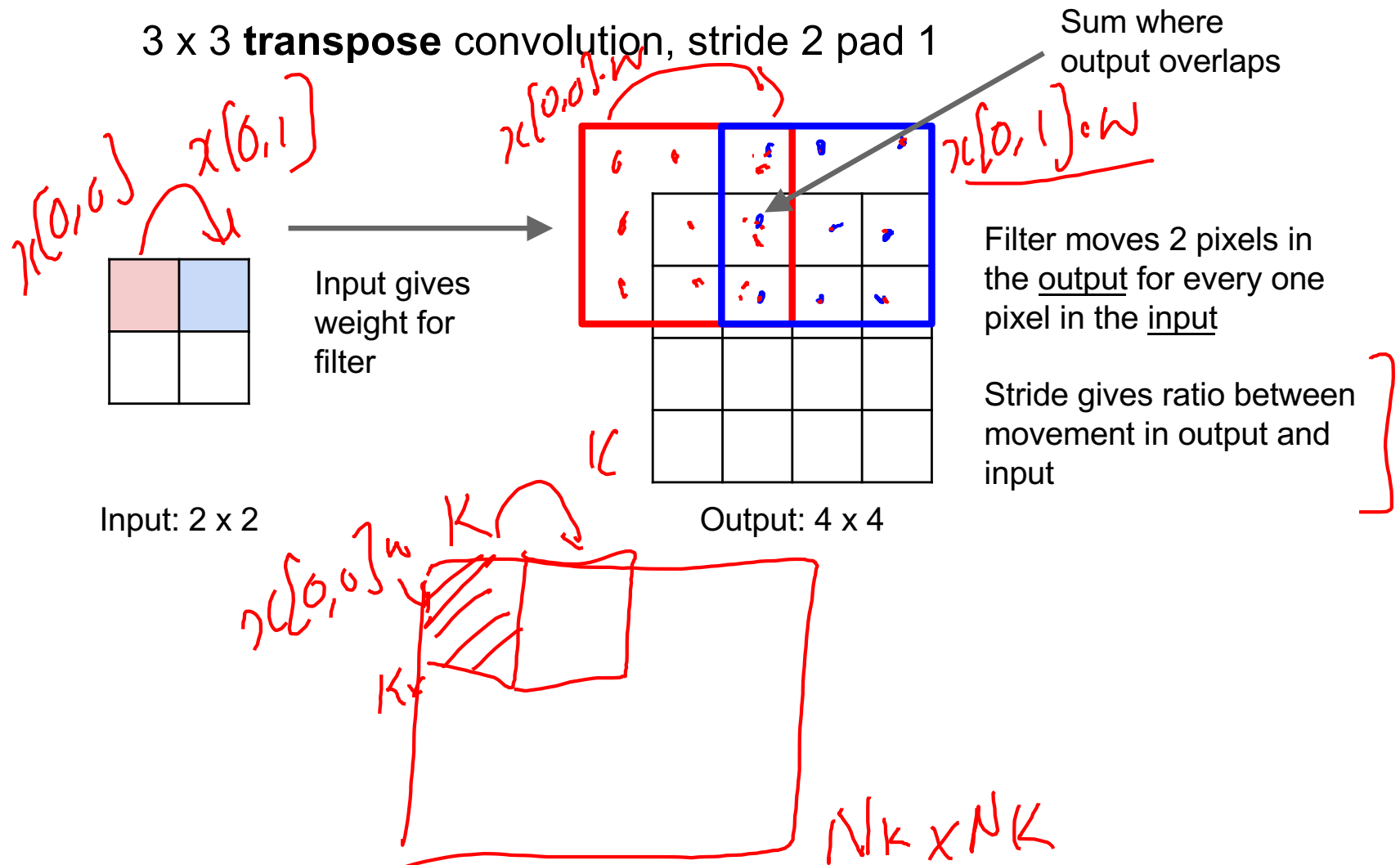


Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1



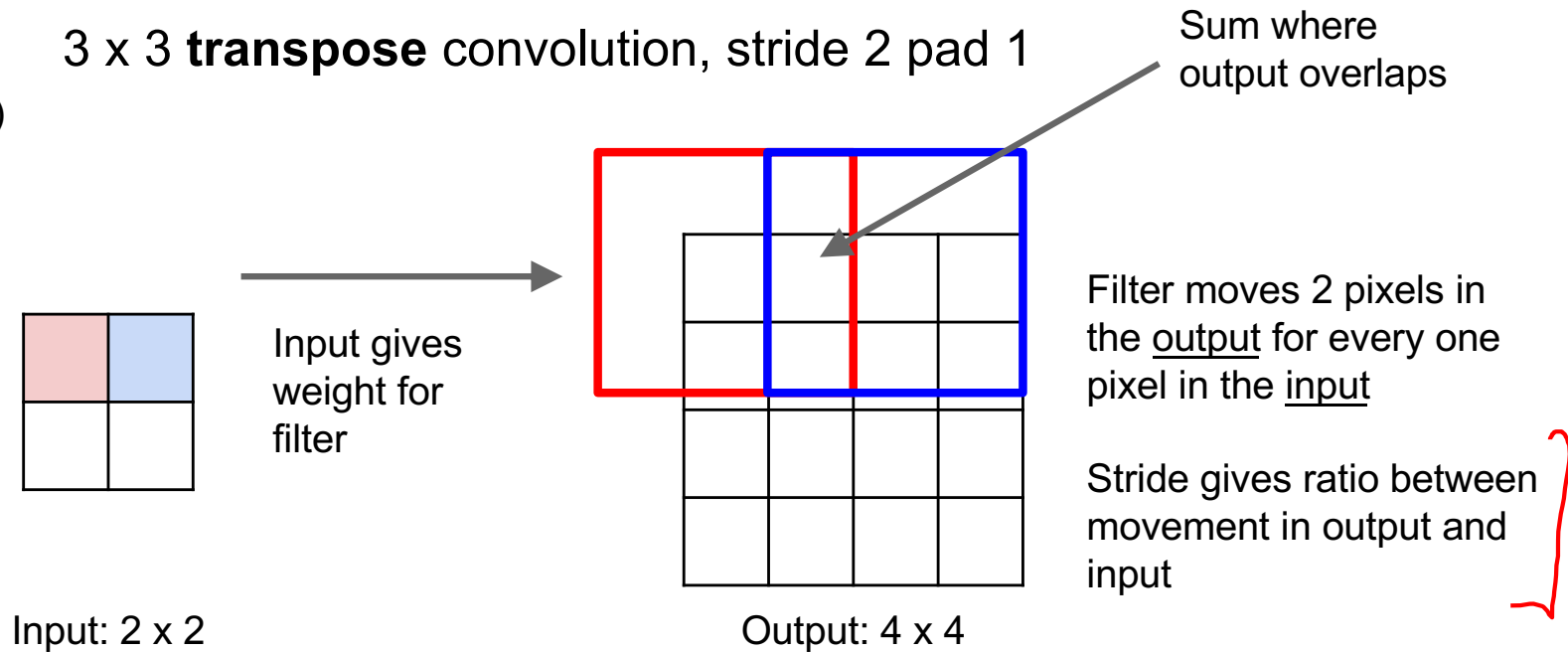
Learnable Upsampling: Transpose Convolution



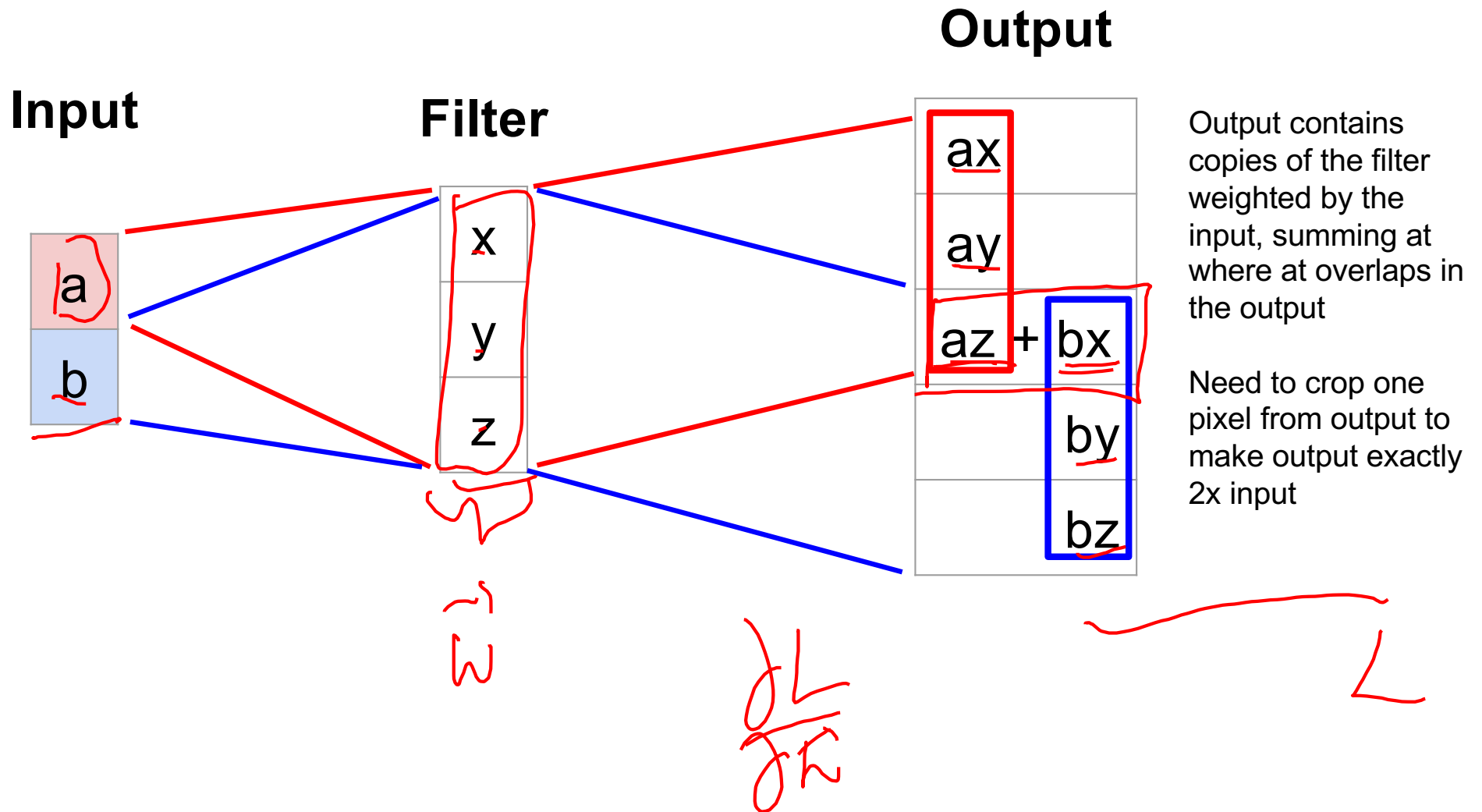
Learnable Upsampling: Transpose Convolution

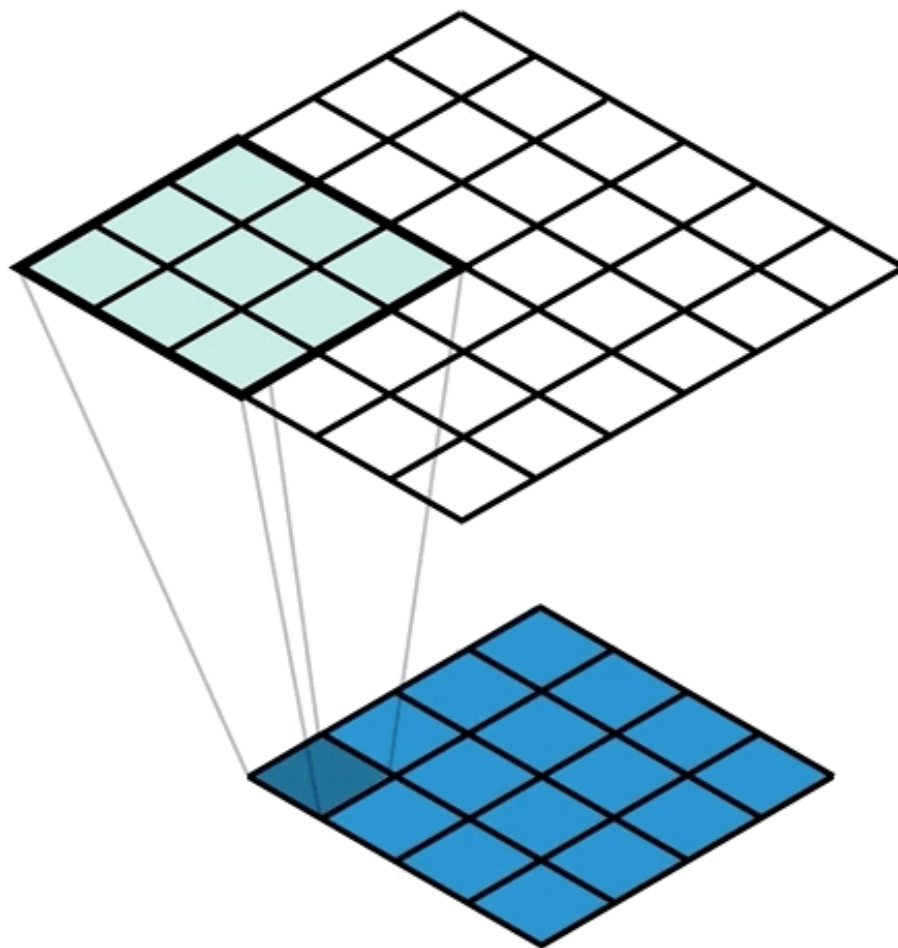
Other names:

- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution



Transpose Convolution: 1D Example

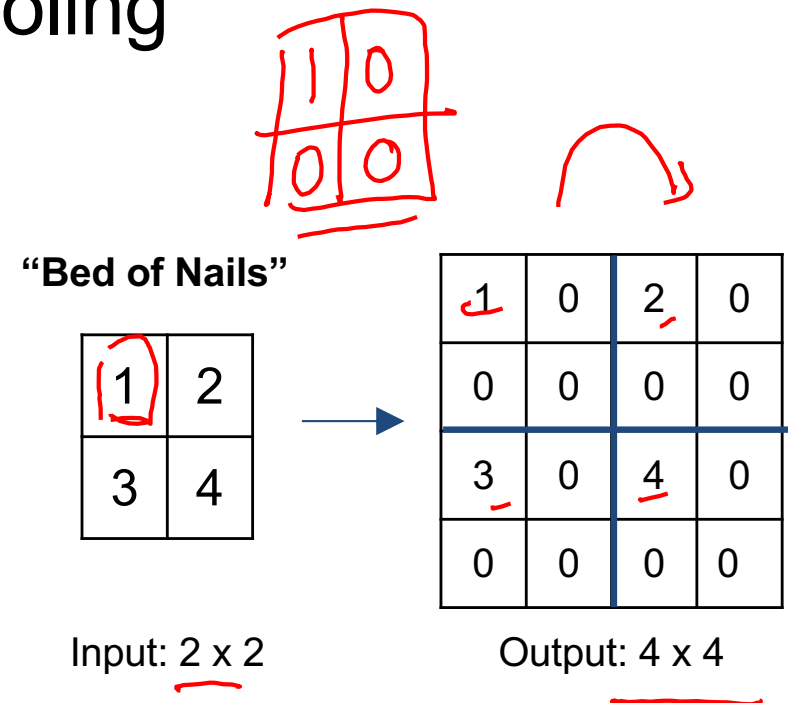
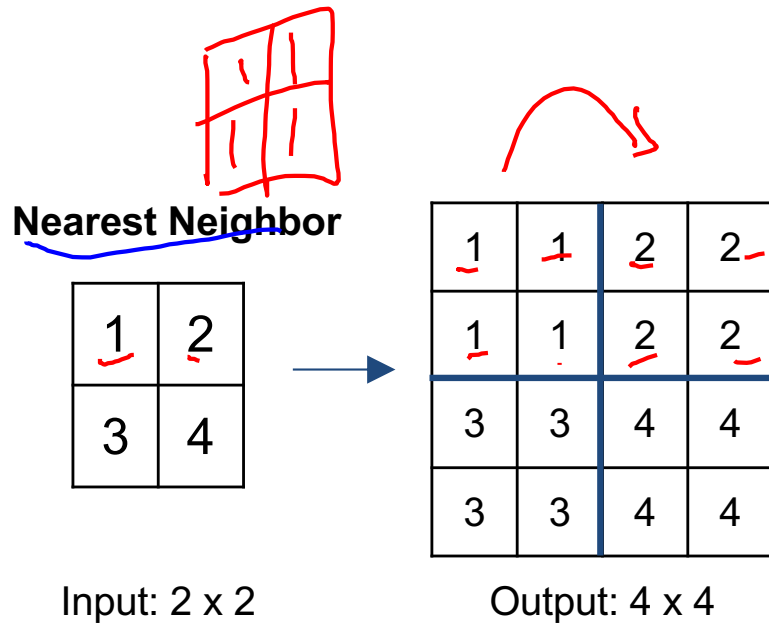




Transposed Convolution

- <https://distill.pub/2016/deconv-checkerboard/>

In-Network upsampling: “Unpooling”



Why this operation?

What is deconvolution?

- (Non-blind) Deconvolution

(conv) $y = x * w$

Deconv → Blind : Given y , produce x, w

Deconv ↘ Non-blind : Given y & w , find x

What is deconvolution?

- (Non-blind) Deconvolution

$$\vec{y} = W \vec{x} \iff x = W^T y$$

If W is orthonormal

$$\vec{w} = \begin{bmatrix} -1 & 0 & +1 \end{bmatrix}$$

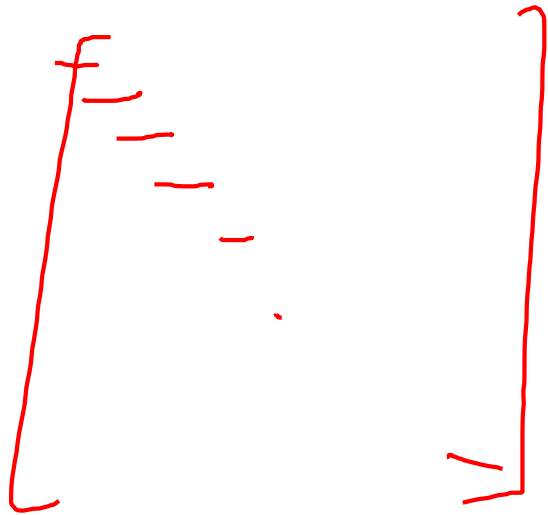
$$\begin{bmatrix} +1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} * \begin{bmatrix} 0 \\ +1 \\ 0 \\ -1 \\ \vdots \end{bmatrix}$$

$$y = w * x$$

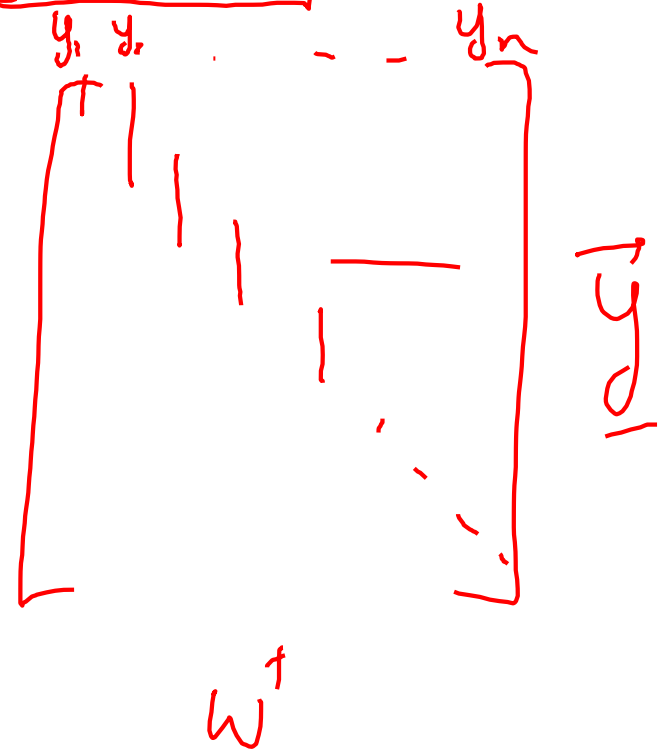
$$\begin{bmatrix} w_k & 0 & \dots & 0 & 0 \\ w_{k-1} & w_k & \dots & 0 & 0 \\ w_{k-2} & w_{k-1} & +1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ w_1 & \dots & \dots & w_k & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & w_1 & \dots & w_{k-1} & w_k \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \vdots & w_1 & w_2 \\ \vdots & \vdots & \vdots & 0 & w_1 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

What does “deconvolution” have to do with “transposed convolution”?



$\vec{x} =$



“transposed convolution” is a convolution!

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} \underline{x} & \underline{y} & z & 0 & 0 & 0 \\ 0 & \underline{x} & \underline{y} & z & 0 & 0 \\ 0 & 0 & \underline{x} & \underline{y} & z & 0 \\ 0 & 0 & 0 & \underline{x} & \underline{y} & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} \underline{ay} + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

“transposed convolution” is a convolution!

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

$$[x \ y \ z]$$

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

$$[z \ y \ x]$$

“transposed convolution” is a convolution!

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

When stride=1, convolution transpose is just a regular convolution (with different padding rules)

But not always

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel
size=3, stride=2, padding=1

But not always

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

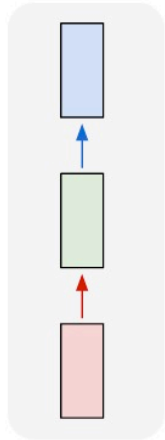
When stride>1, convolution transpose is no longer a normal convolution!

Plan for Today

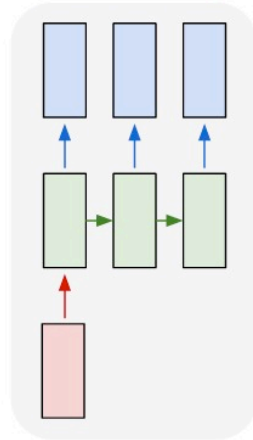
- (Finish) Convolutional Neural Networks
 - Transposed convolutions
- Recurrent Neural Networks (RNNs)
 - A new model class
 - Learning: BackProp Through Time (BPTT)

New Topic: RNNs

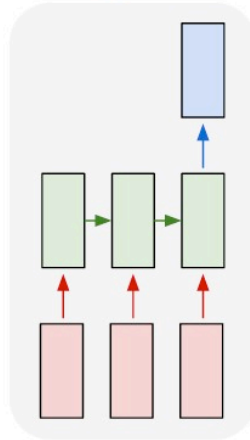
one to one



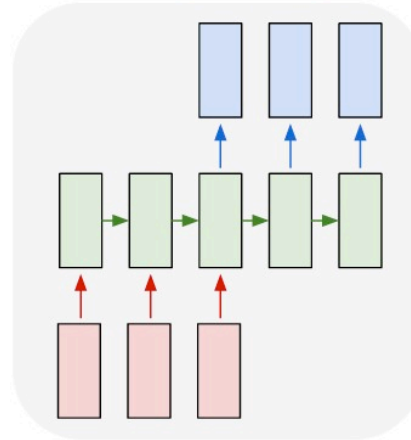
one to many



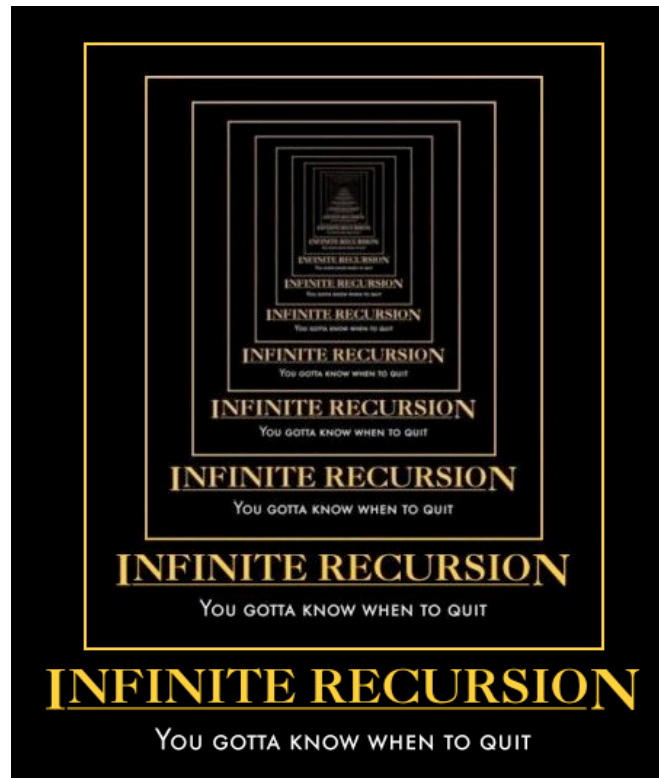
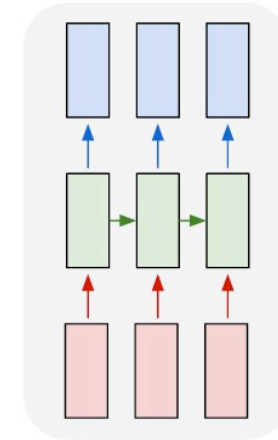
many to one



many to many



many to many

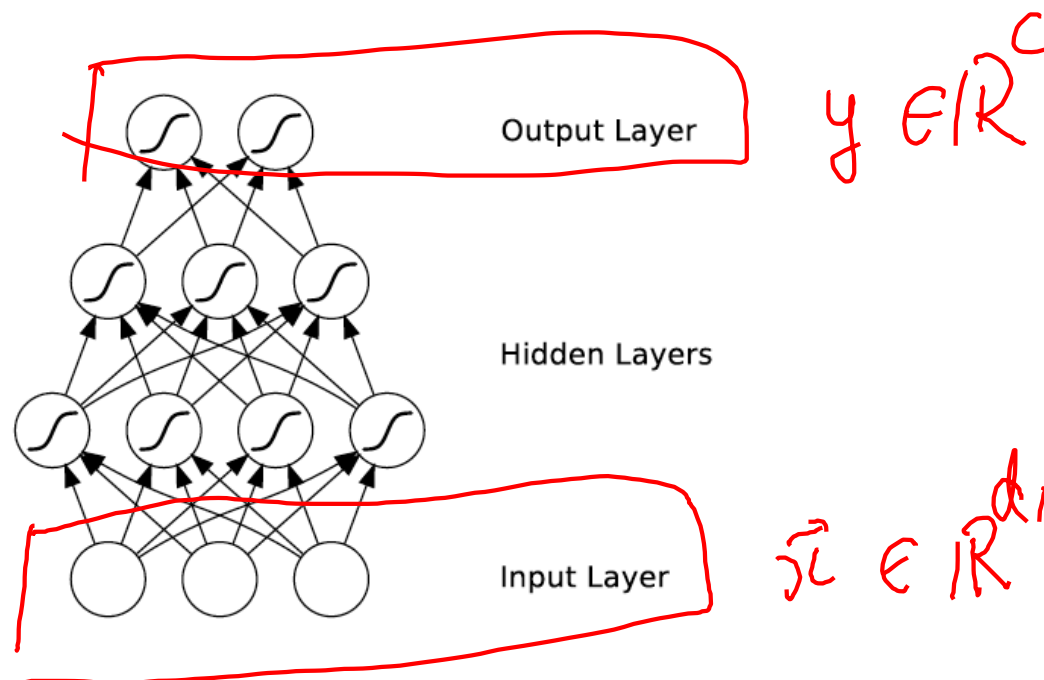


New Words

- 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)

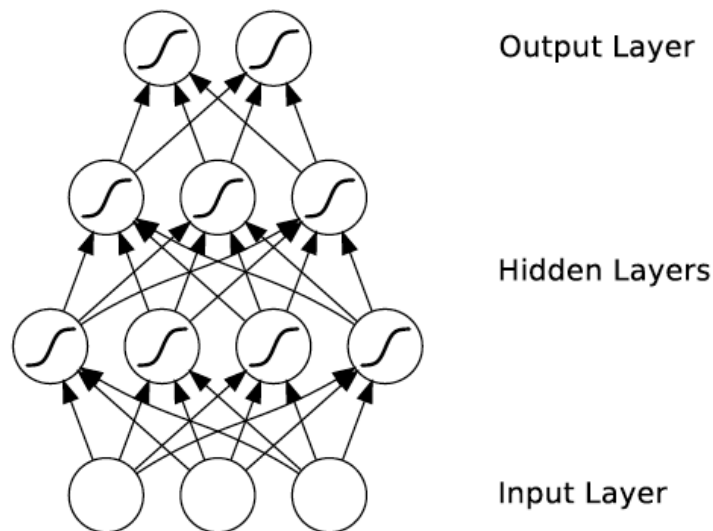
What's wrong with MLPs?

- 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?

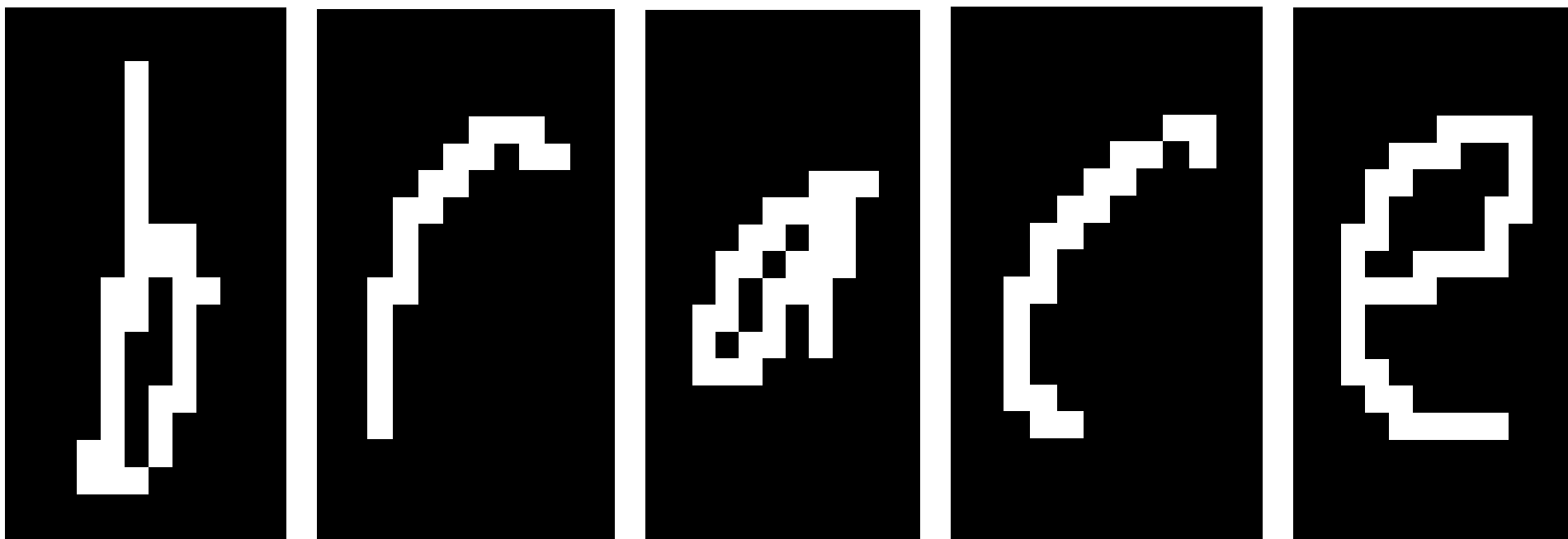
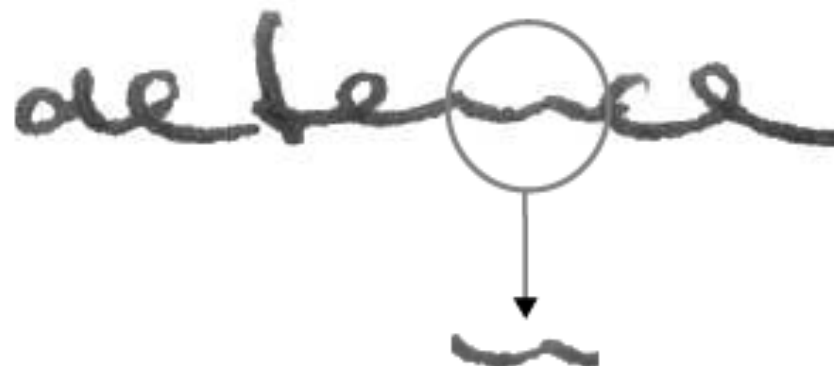


Figure Credit: Carlos Guestrin

Why model sequences?



Sequences are everywhere...

Foreign Minister. → FOREIGN MINISTER.

[Waveform] → THE SOUND OF

$x =$ bringen sie bitte das auto zurück .
 $a_1=2$ $a_2=0$ $a_3=1$ $a_4=3$ $a_5=4$ $a_6=2$ $a_7=5$

$y =$ please return the car .

↙

Even where you might not expect a sequence...

Classify images by taking a series of “glimpses”



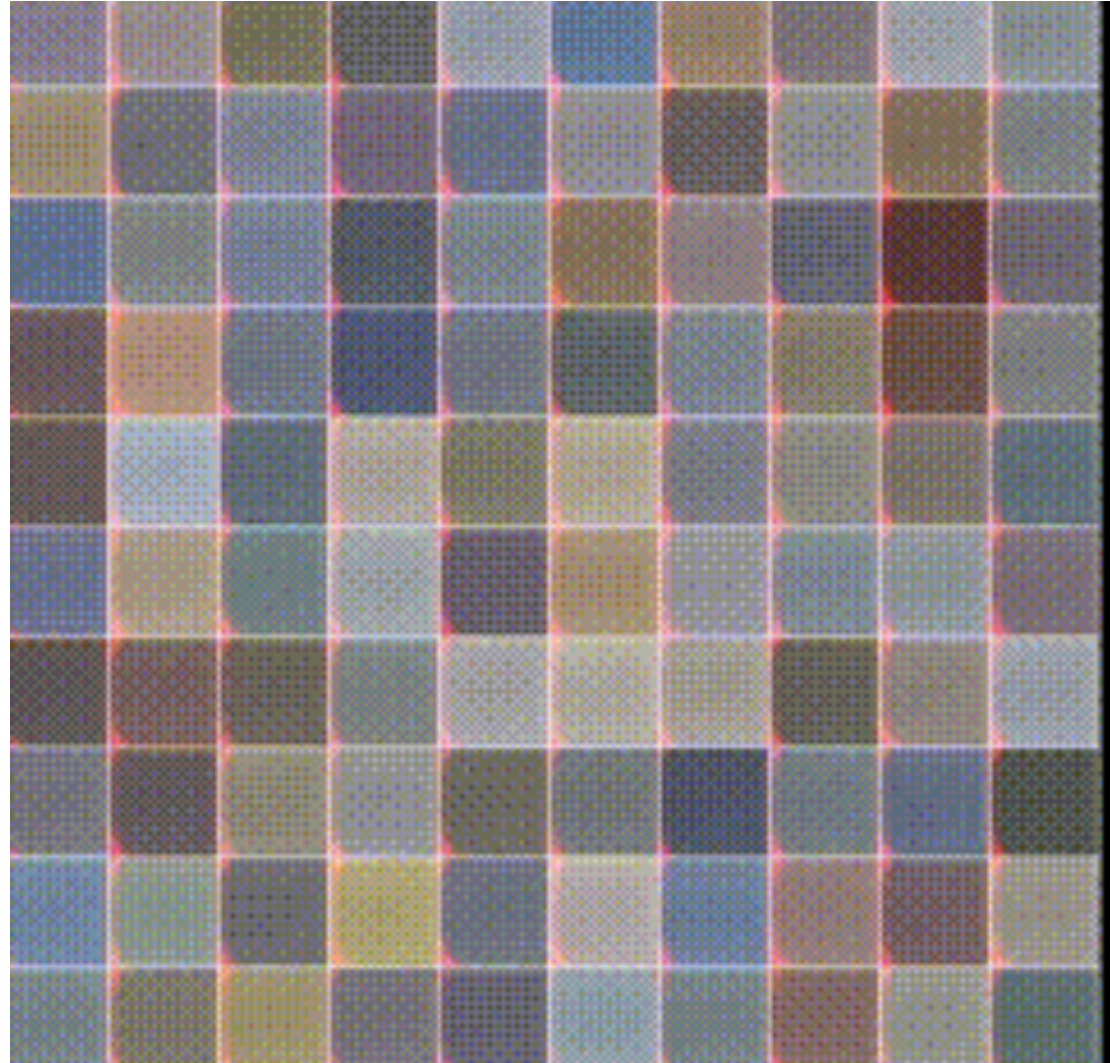
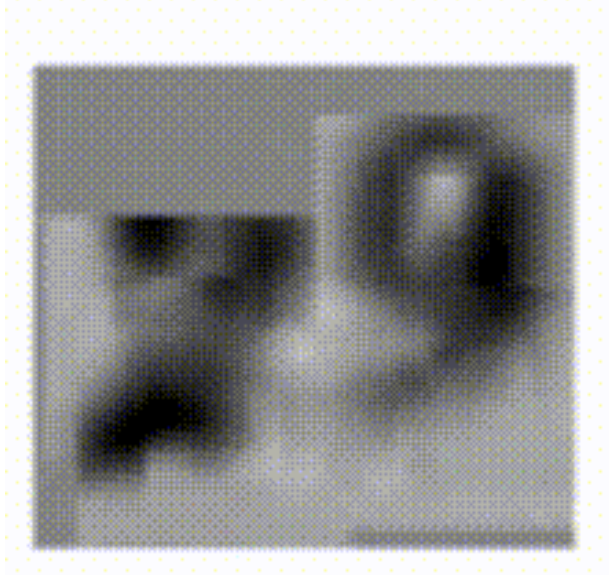
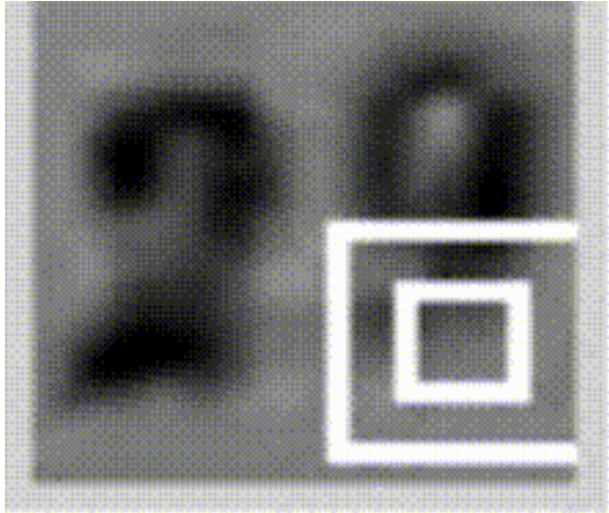
Ba, Mnih, and Kavukcuoglu, “Multiple Object Recognition with Visual Attention”, ICLR 2015.

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...

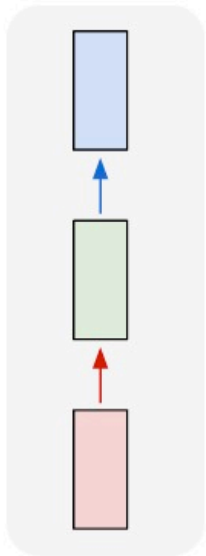
- Output ordering = sequence



Sequences in Input or Output?

- It's a spectrum...

one to one



Input: No
sequence

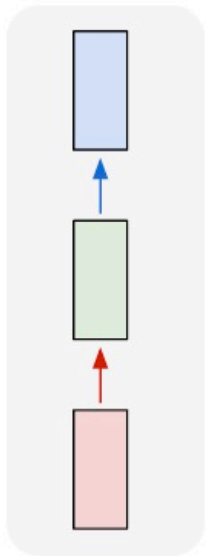
Output: No
sequence

Example:
"standard"
classification
regression
problems

Sequences in Input or Output?

- It's a spectrum...

one to one

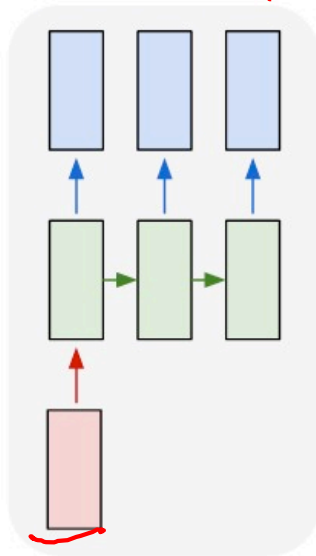


Input: No sequence

Output: No sequence

Example:
"standard"
classification /
regression
problems

one to many



Input: No sequence

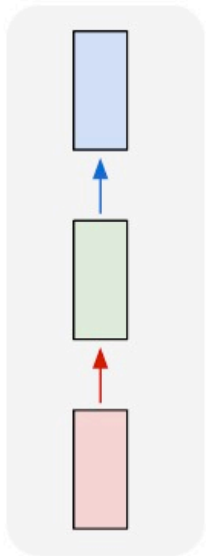
Output: Sequence

Example:
Im2Caption

Sequences in Input or Output?

- It's a spectrum...

one to one

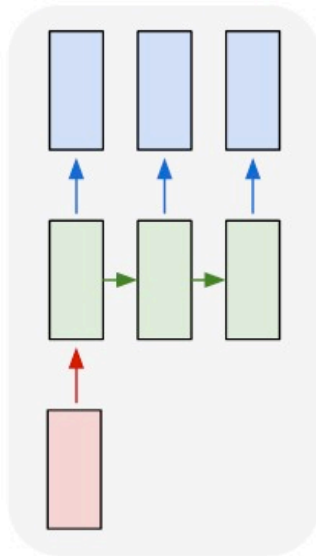


Input: No sequence

Output: No sequence

Example: "standard" classification / regression problems

one to many

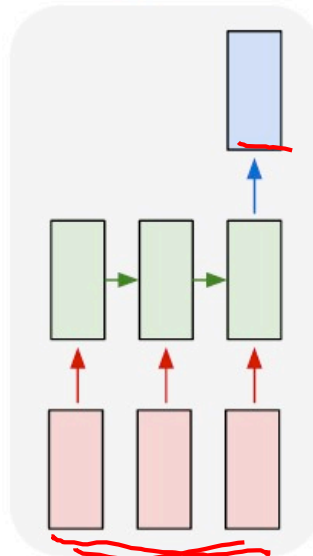


Input: No sequence

Output: Sequence

Example: Im2Caption

many to one



Input: Sequence

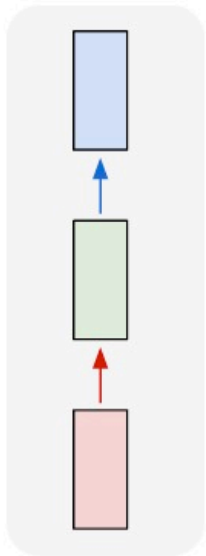
Output: No sequence

Example: sentence classification, multiple-choice question answering

Sequences in Input or Output?

- It's a spectrum...

one to one

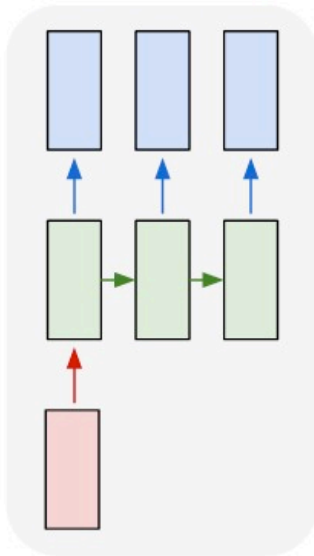


Input: No sequence

Output: No sequence

Example: "standard" classification / regression problems

one to many

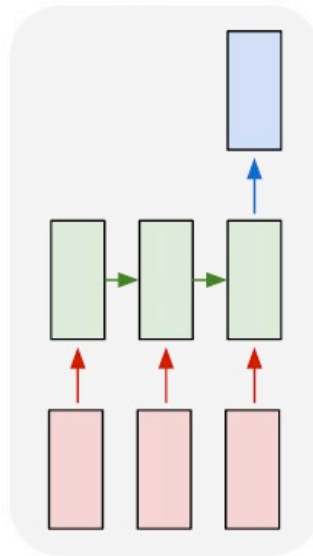


Input: No sequence

Output: Sequence

Example: Im2Caption

many to one

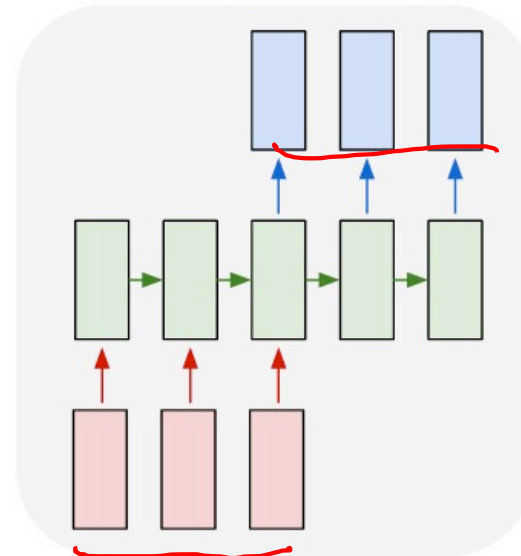


Input: Sequence

Output: No sequence

Example: sentence classification, multiple-choice question answering

many to many

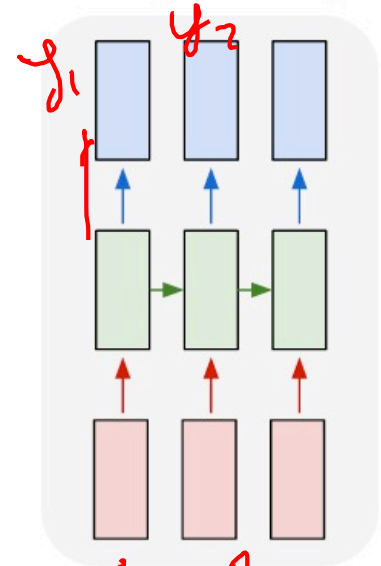


Input: Sequence

Output: Sequence

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

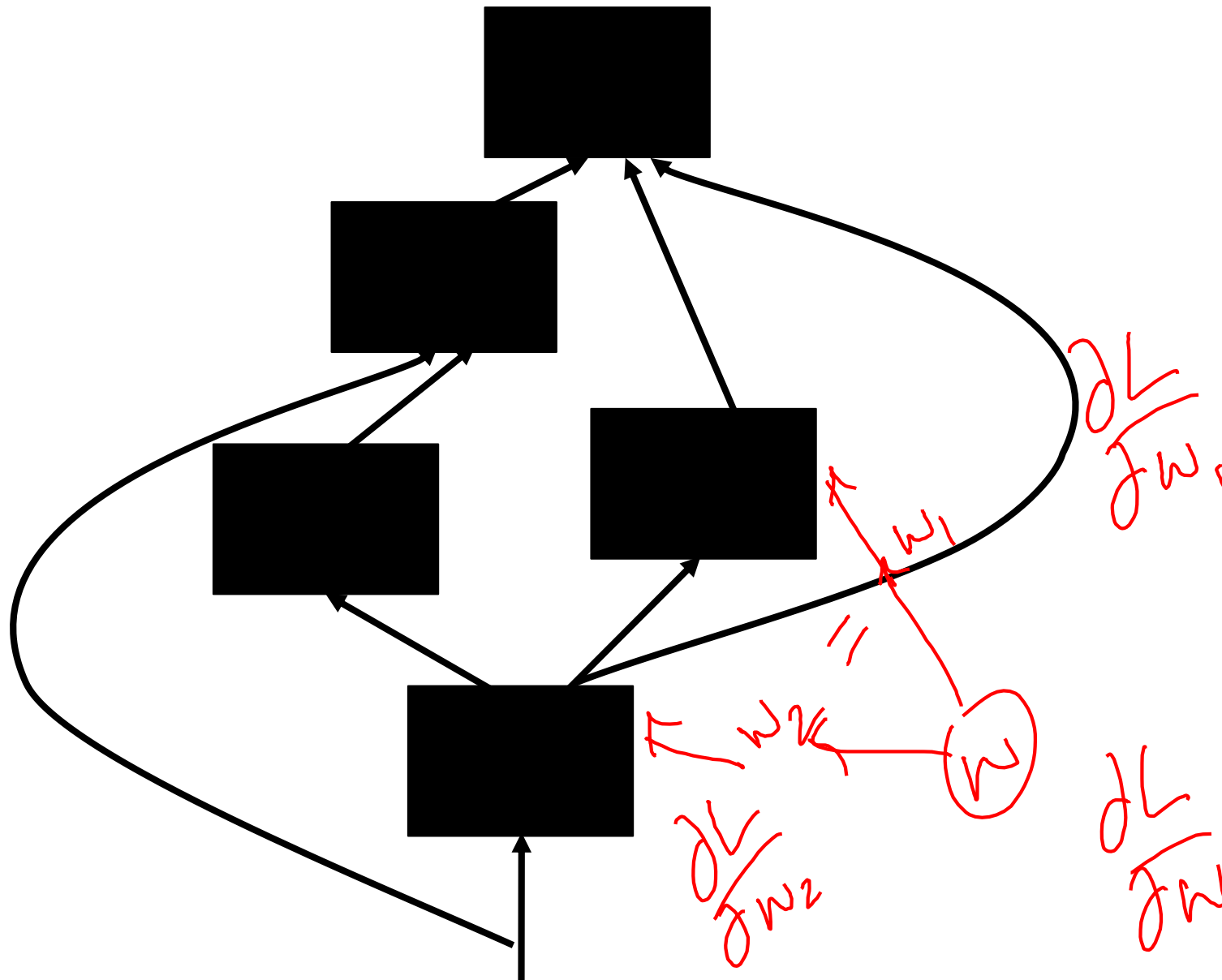
many to many



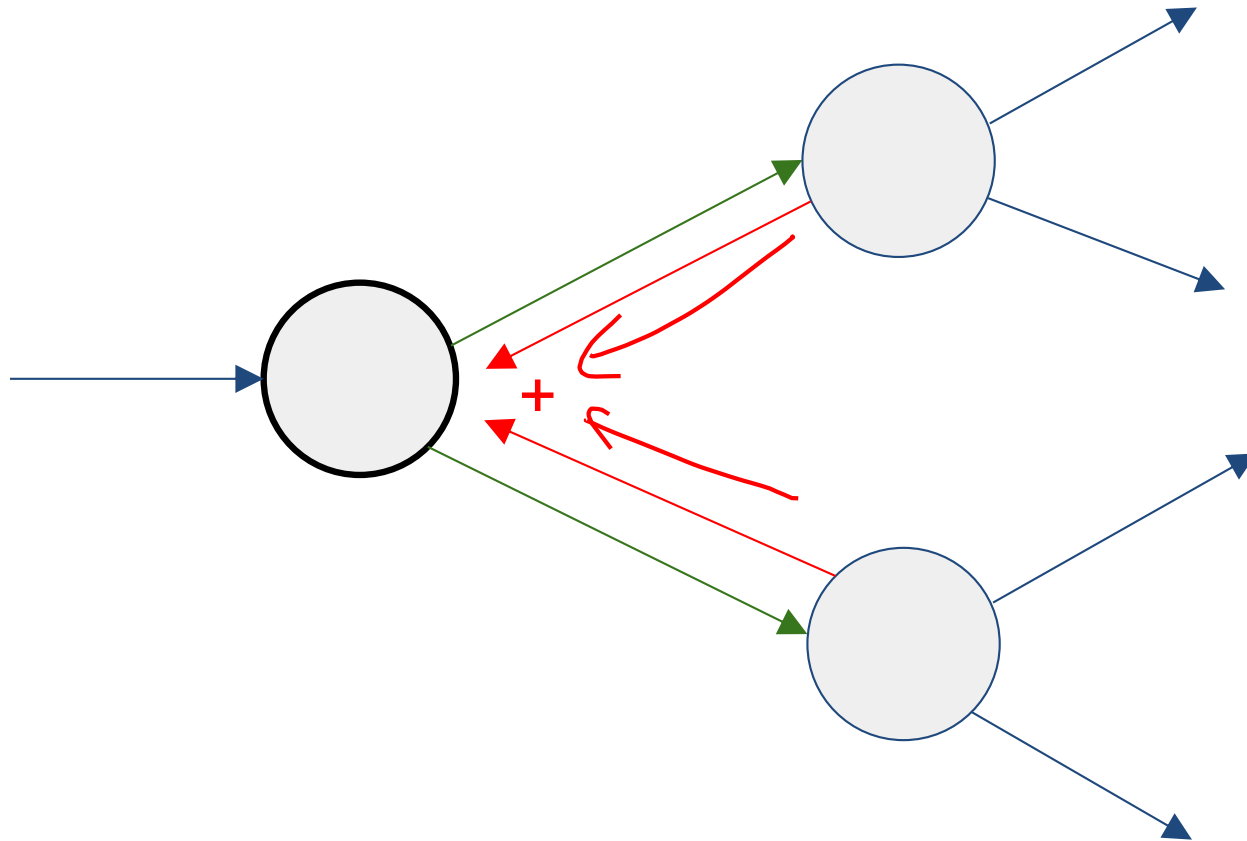
2 Key Ideas

- Parameter Sharing
 - in computation graphs = adding gradients

Computational Graph



Gradients add at branches



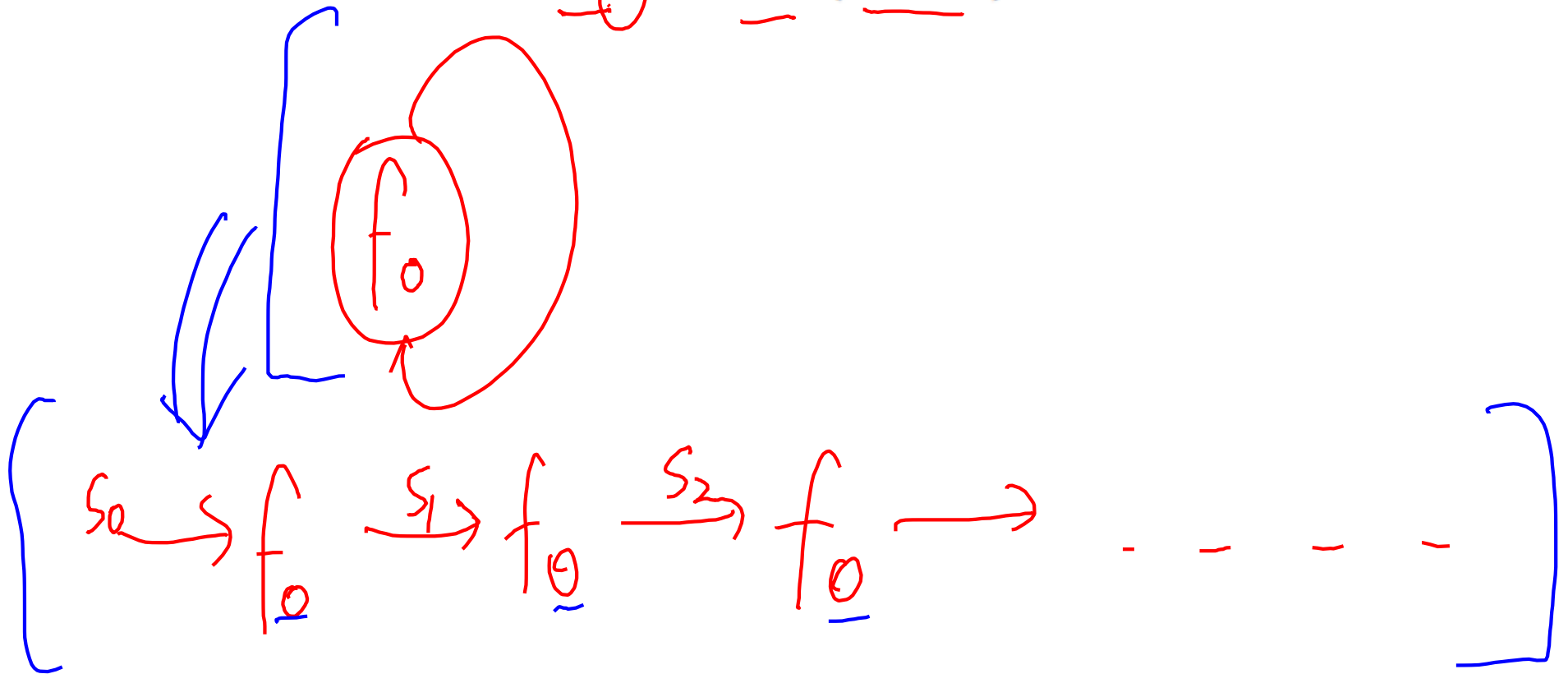
2 Key Ideas

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

How do we model sequences?

- No input

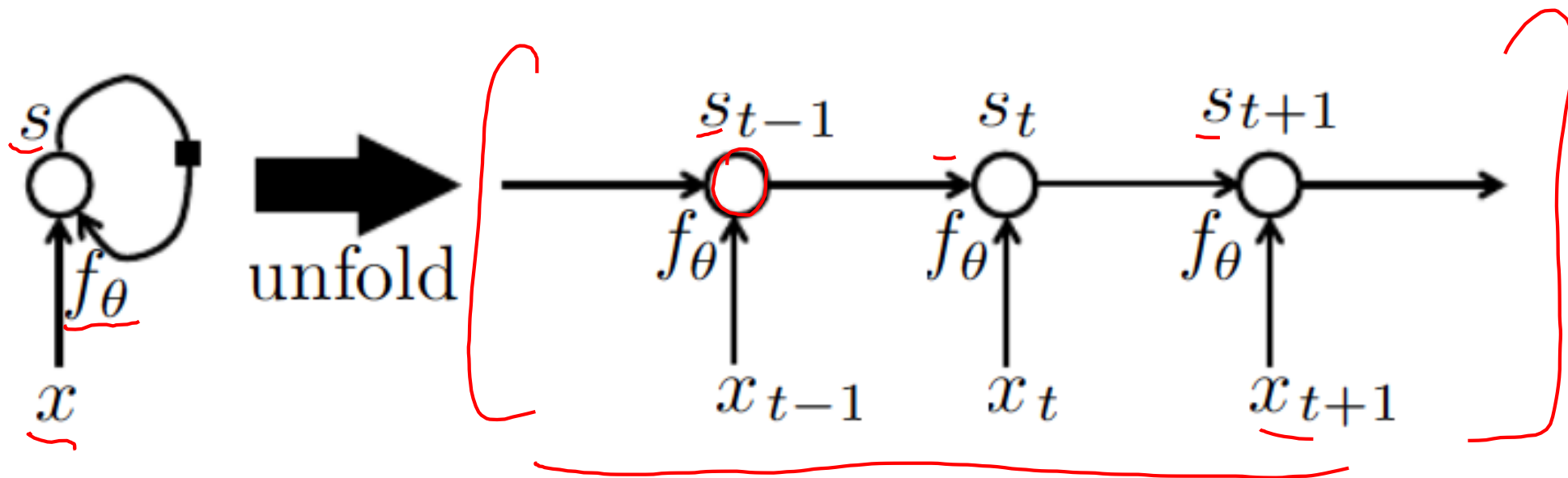
$$s_t = f_\theta(s_{t-1})$$



How do we model sequences?

- With inputs

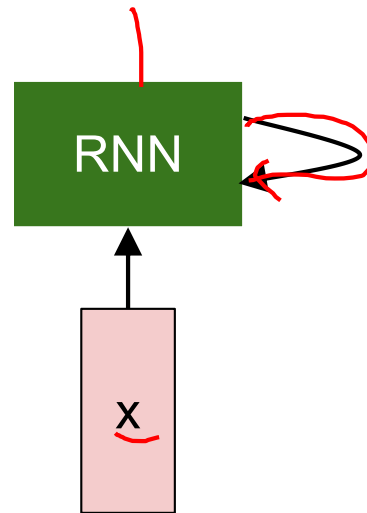
$$s_t = f_\theta(s_{t-1}, x_t)$$



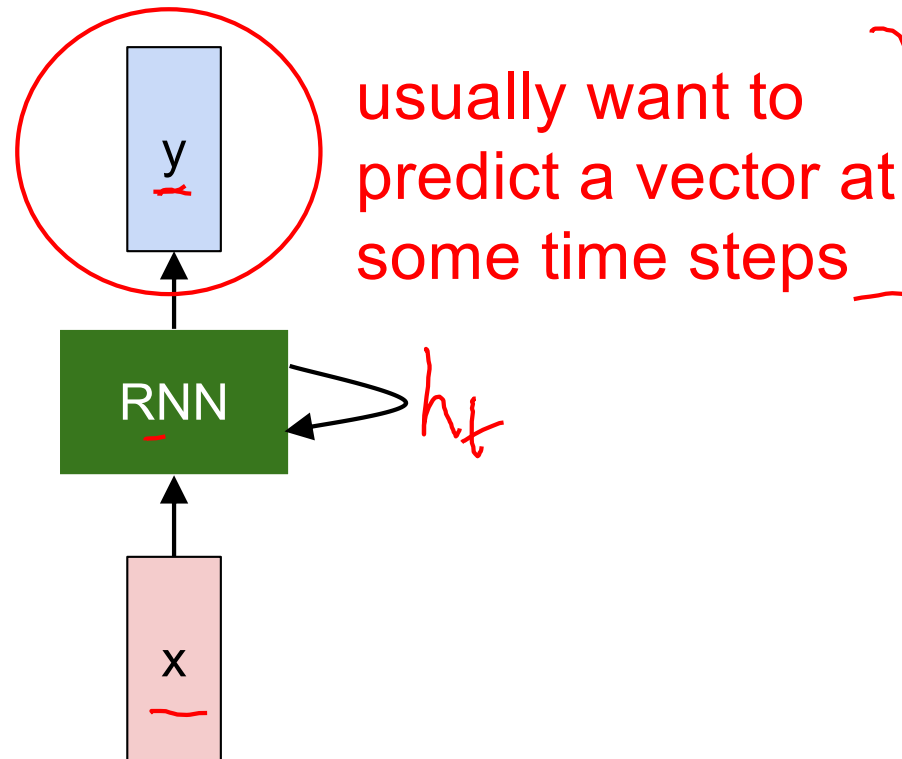
2 Key Ideas

- 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



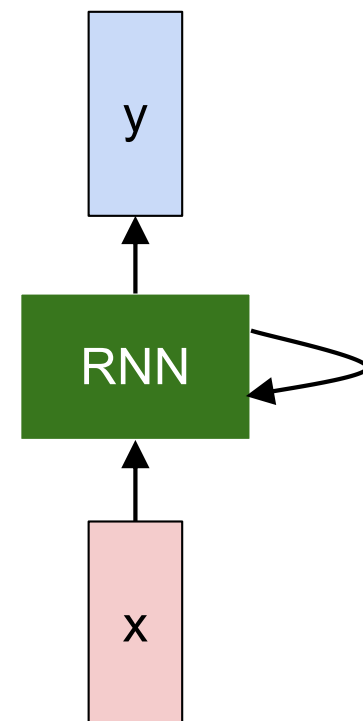
Recurrent Neural Network

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

$$[y_t = f_{W_2}(h_t)]$$

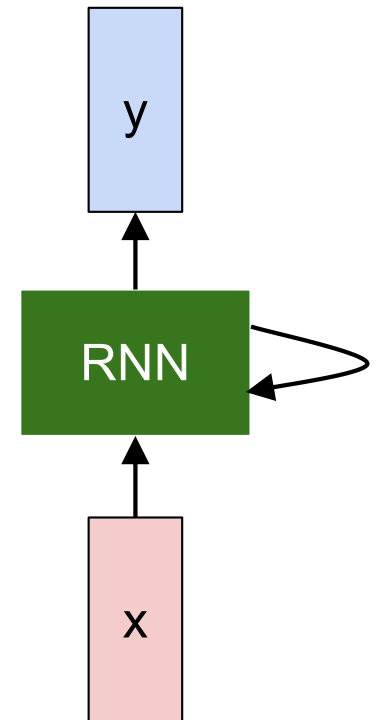


Recurrent Neural Network

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

$$h_t = f_{\underline{W}}(h_{t-1}, x_t)$$

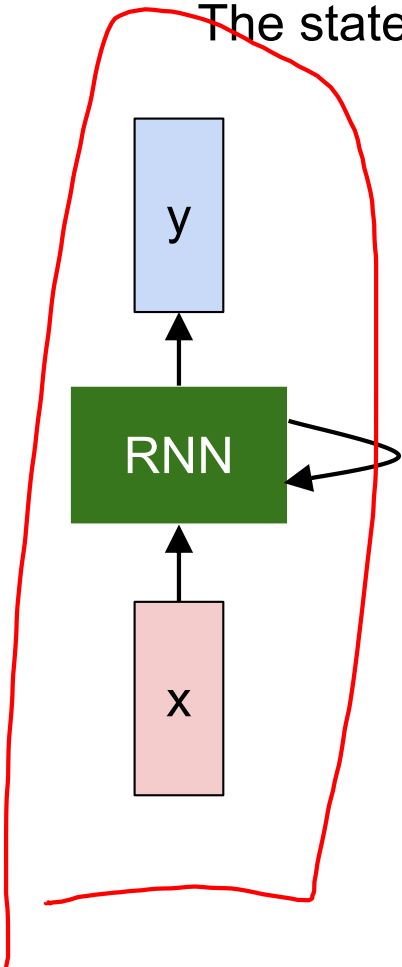
Notice: the same function and the same set of parameters are used at every time step.



(Vanilla) Recurrent Neural Network

$h_t \in \mathbb{R}^{d_2}$
 $x_t \in \mathbb{R}^{d_1}$

The state consists of a single "hidden" vector h :



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

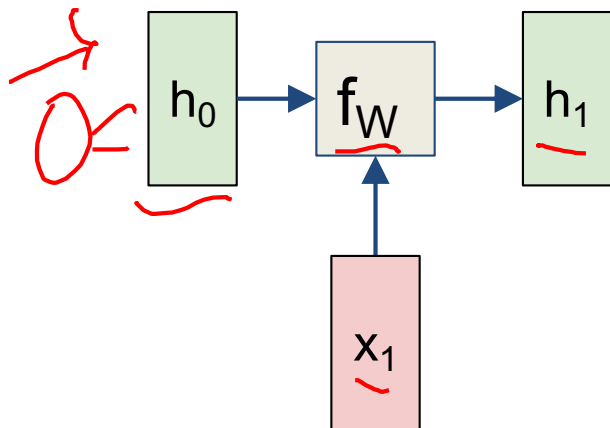
$$\underline{h_t} = \underline{f_W}(\underline{h_{t-1}}, \underline{x_t})$$



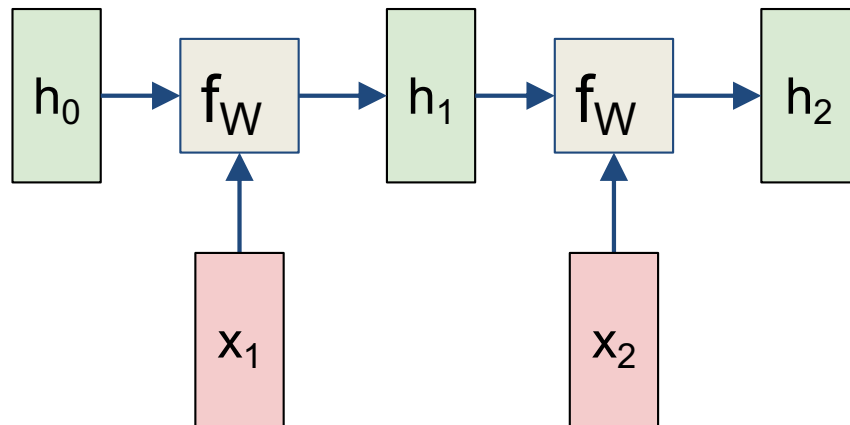
$$\underline{h_t} = \underline{\tanh}(W_{hh} \underline{h_{t-1}} + W_{xh} \underline{x_t} + \underline{b_h})$$



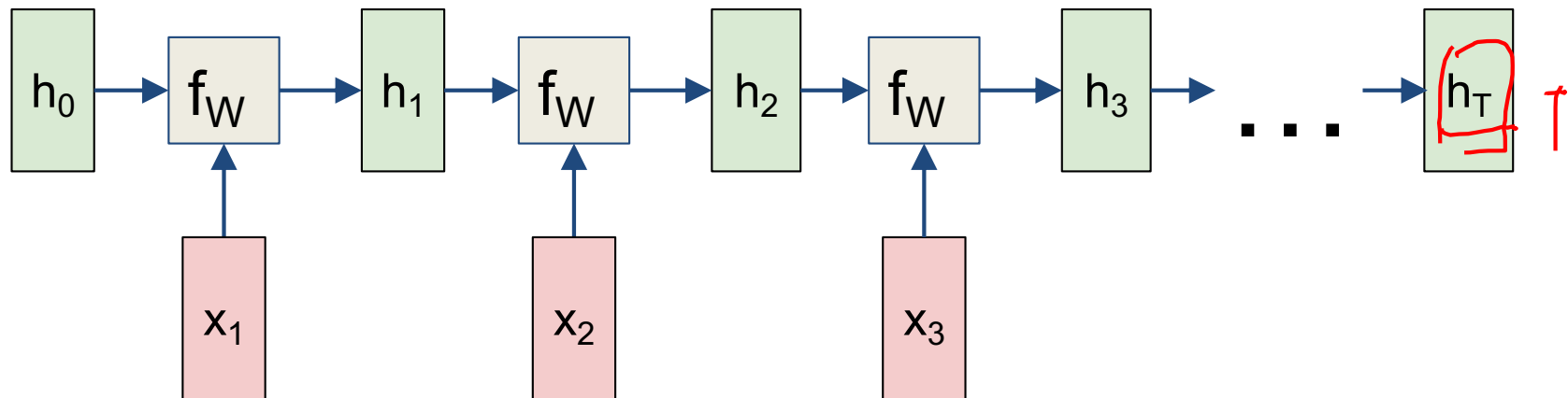
RNN: Computational Graph



RNN: Computational Graph

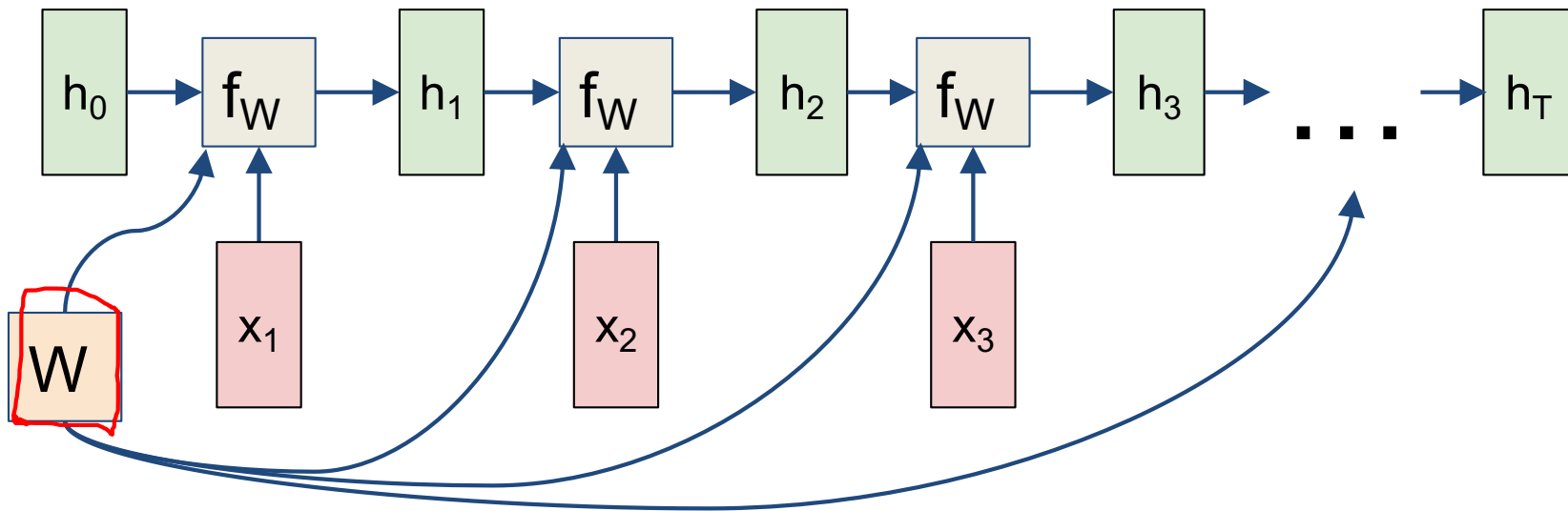


RNN: Computational Graph

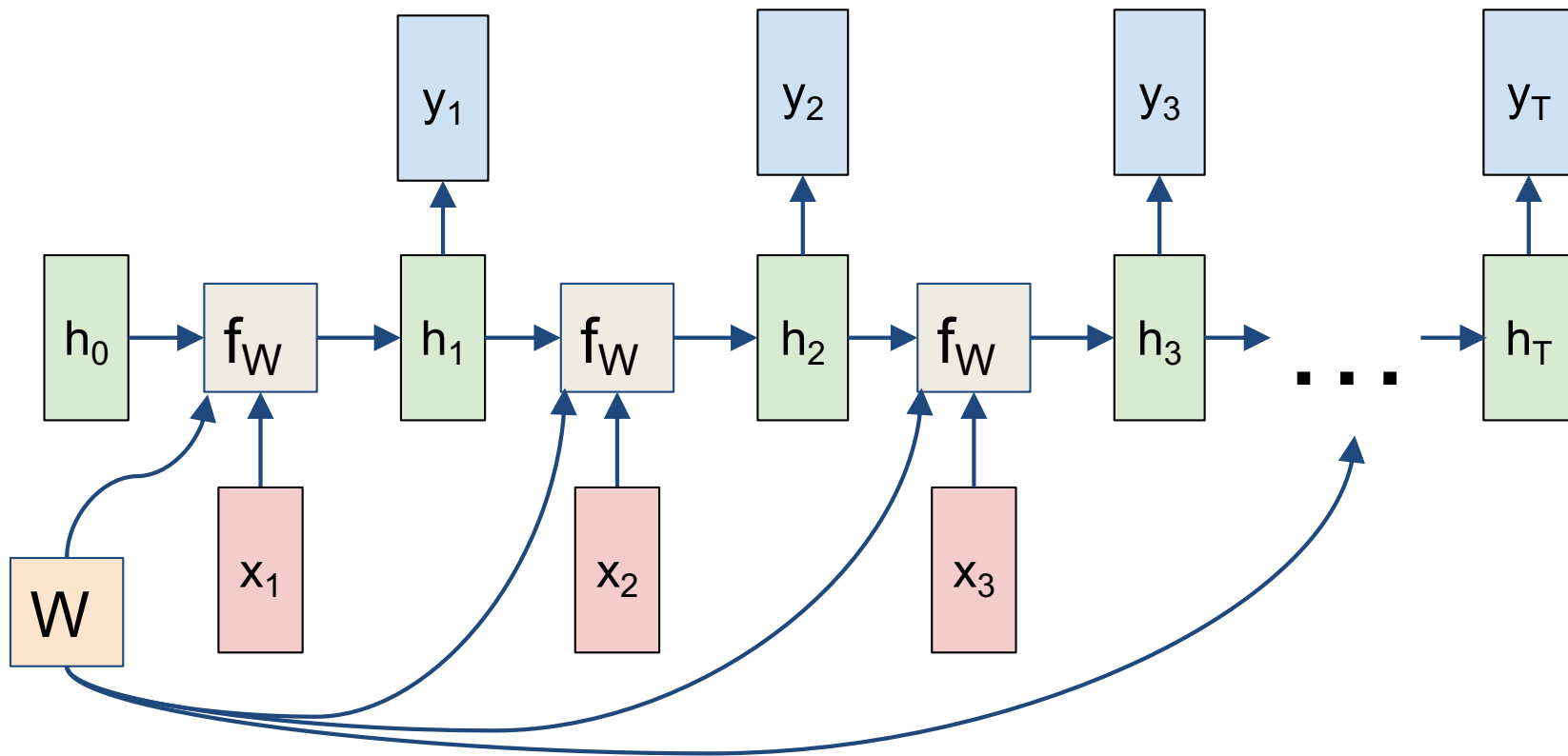


RNN: Computational Graph

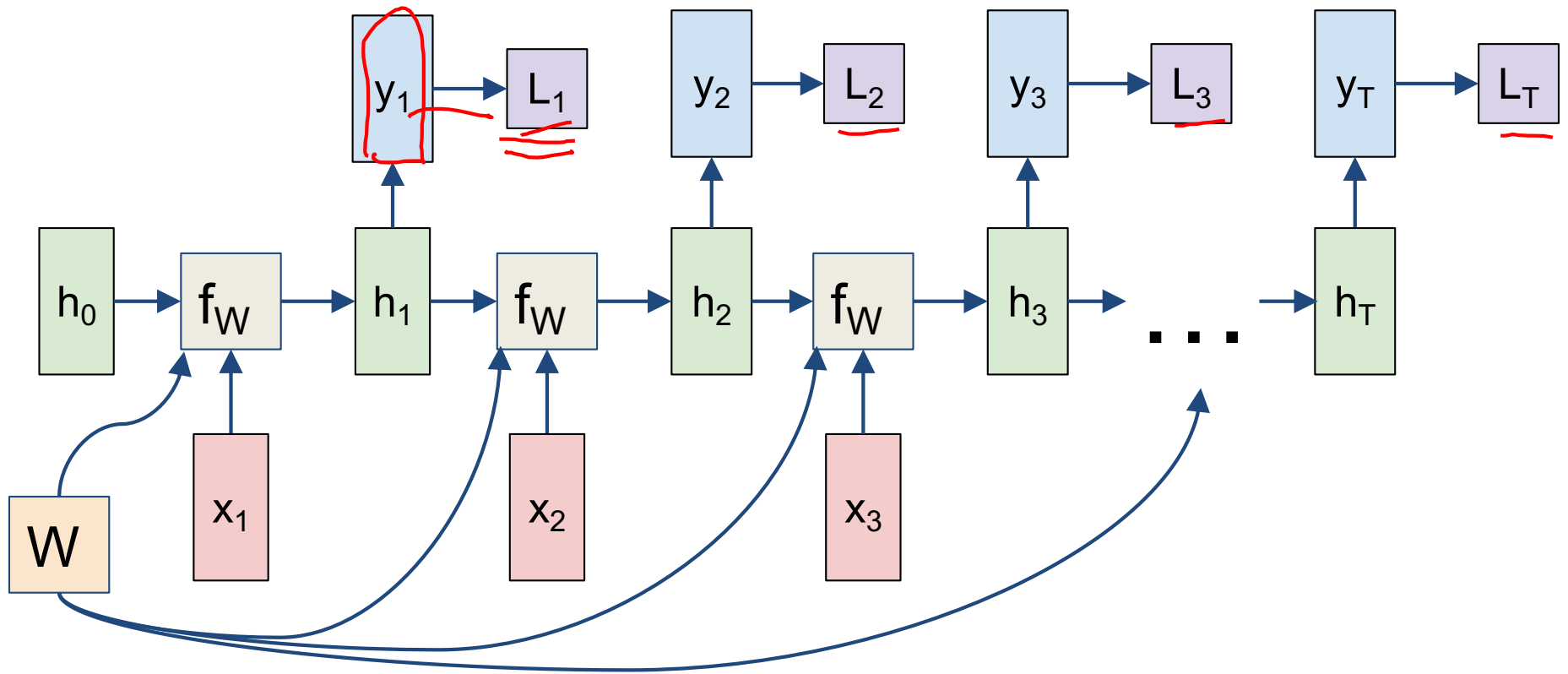
Re-use the same weight matrix at every time-step



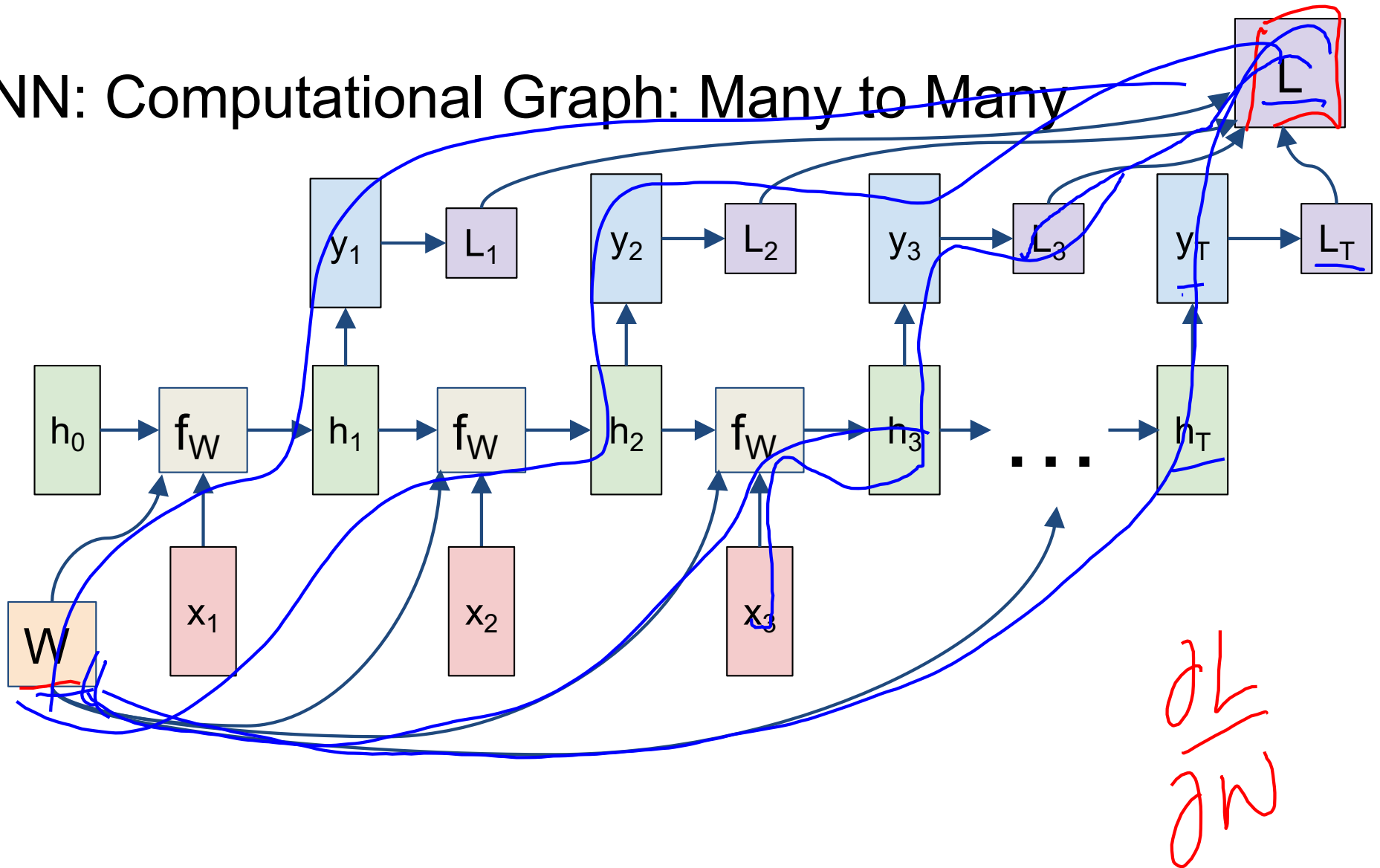
RNN: Computational Graph: Many to Many



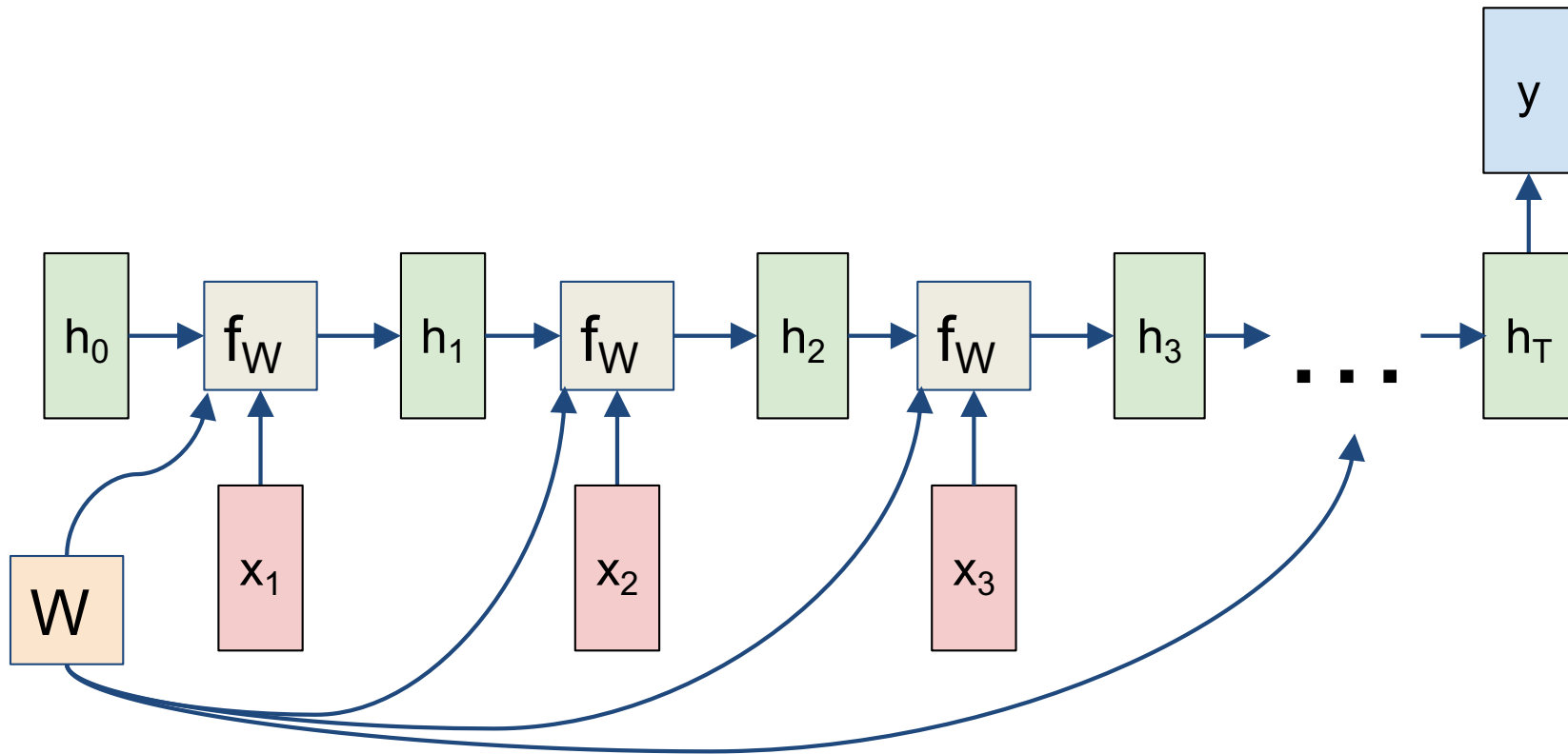
RNN: Computational Graph: Many to Many



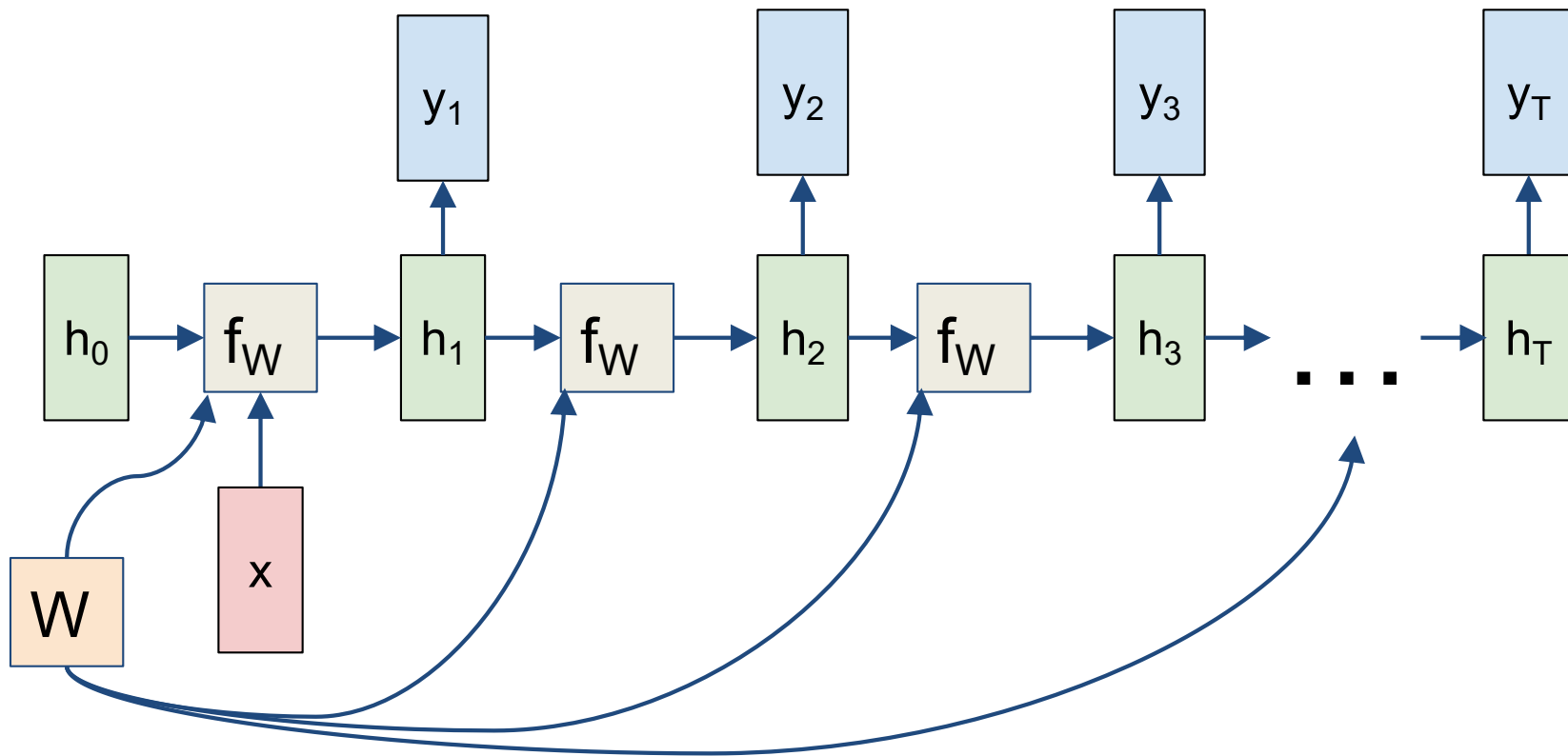
RNN: Computational Graph: Many to Many



RNN: Computational Graph: Many to One

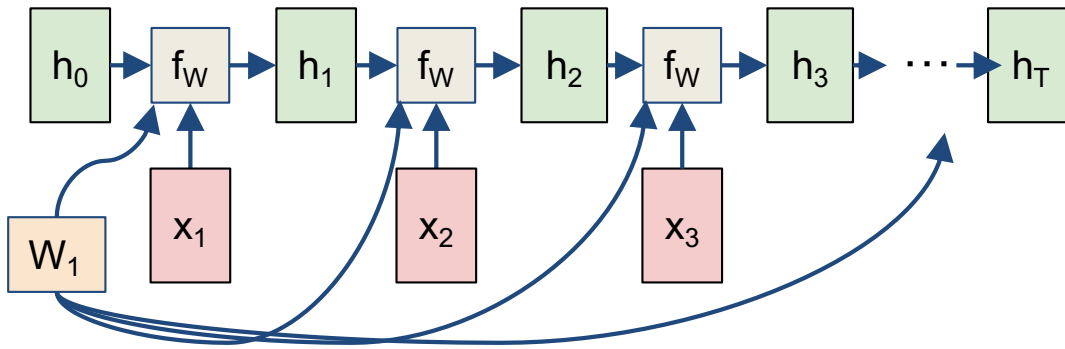


RNN: Computational Graph: One to Many



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

