



# Deep Learning

**CNNs & RNNs**

Tuesday 19<sup>th</sup> November

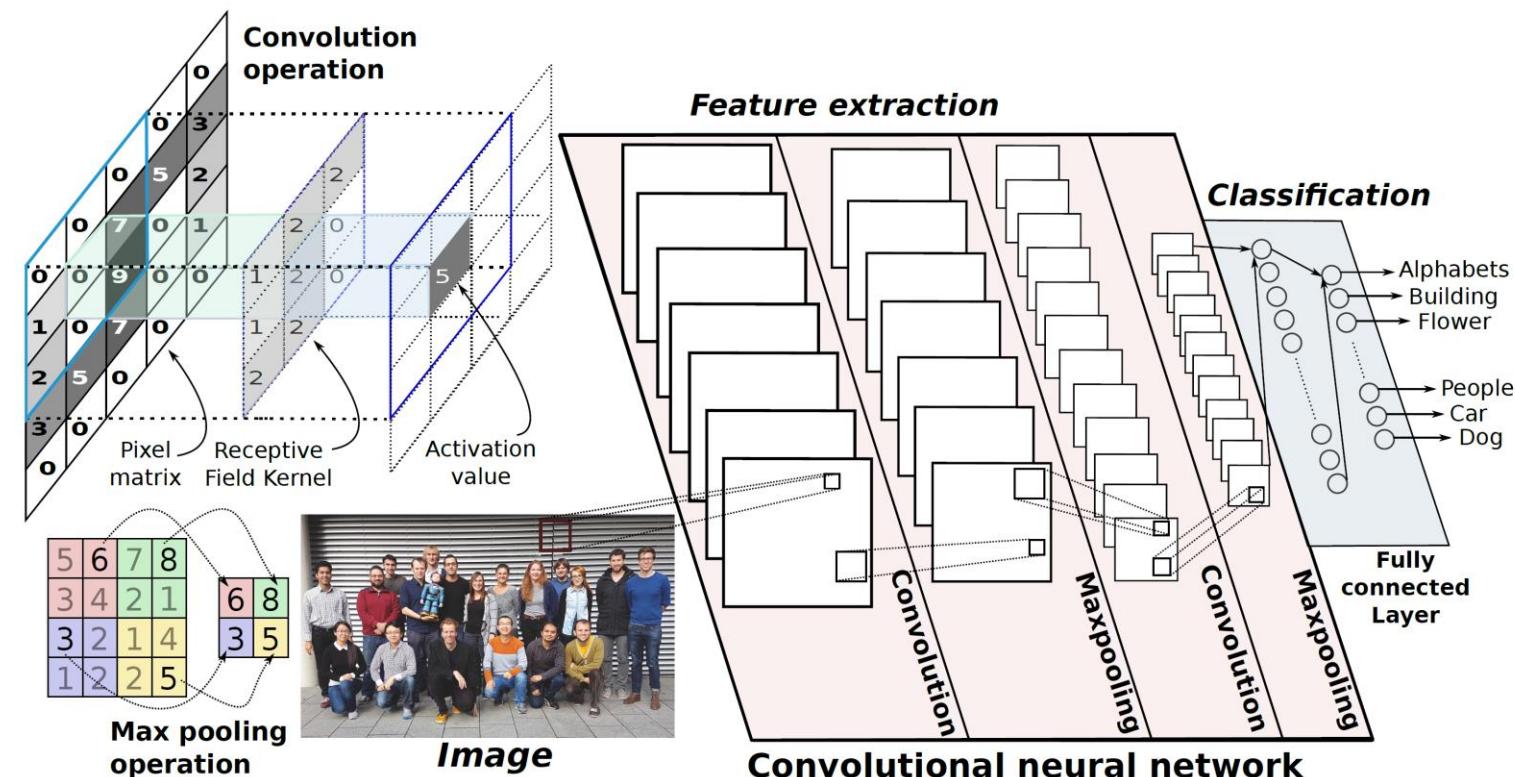
Dr. Nicholas Cummins

# Convolutional Neural Network (CNNs)

Improve generalisation by reusing the same weights to  
detect the same feature in multiple places

- **What is a Convolutional Neural Network?**
  - A CNN is a neural network with **convolution operations** instead of matrix multiplications in at least one of the layers
  - **Special form of feed-forward network**
    - Reuse the same neurons for repetitive convolution tasks
      - Convolution  $\approx$  Filtering
    - Convolutional kernels recognise patterns in a signal
    - Reuses weights to detect the same patterns in multiple places
    - Reduces overfitting and leads to much more accurate models

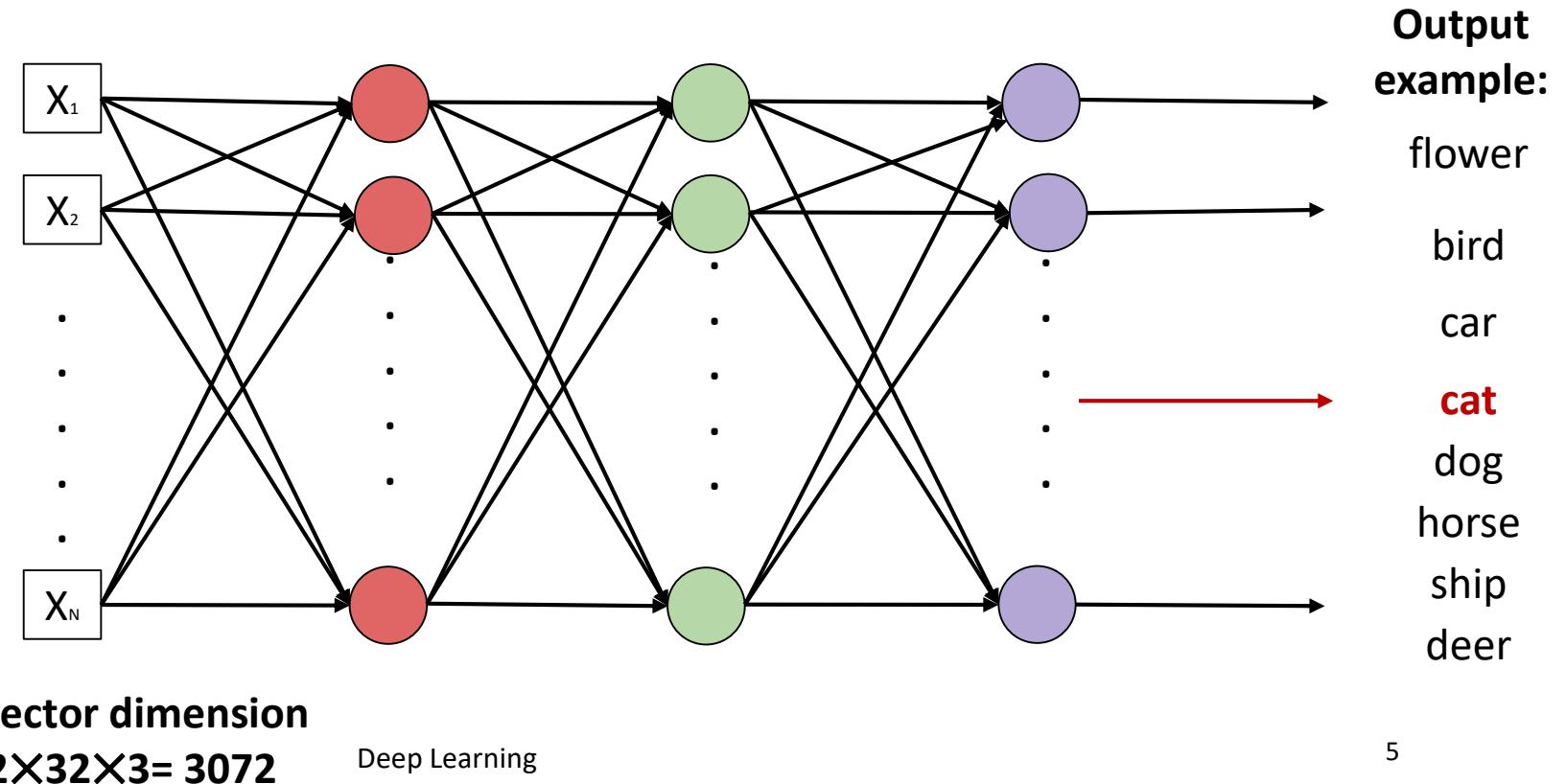
- **Convolutional Neural Network (CNN)**
  - Convolutional kernels perform feature extraction



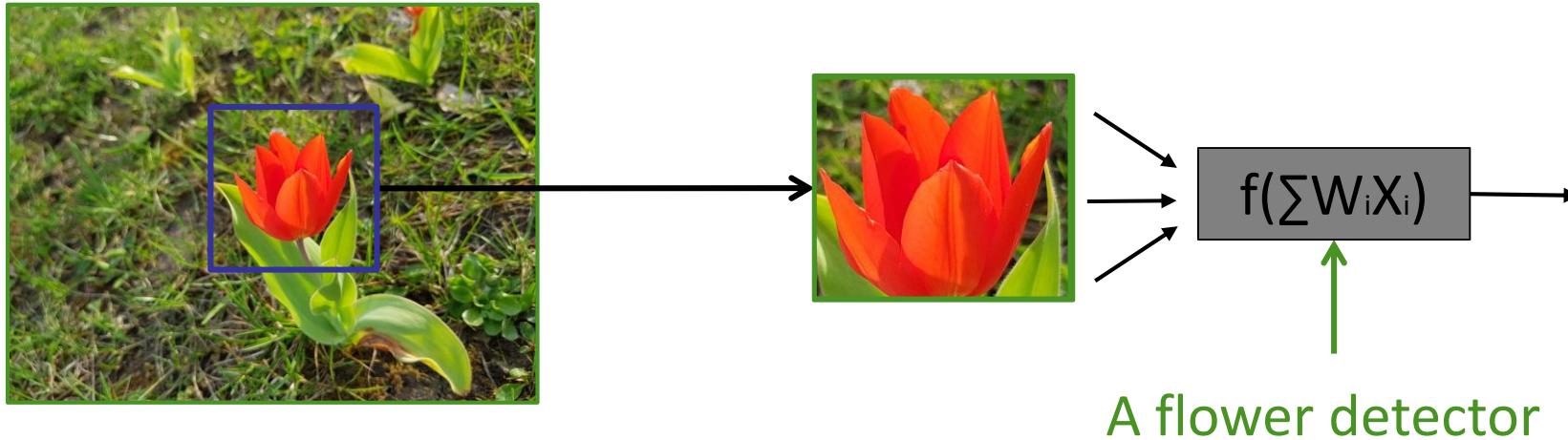
- When applying a fully connected feedforward neural network
  - Many thousands of weights and connections
  - Feedforward networks can be simplified by considering properties of input signal



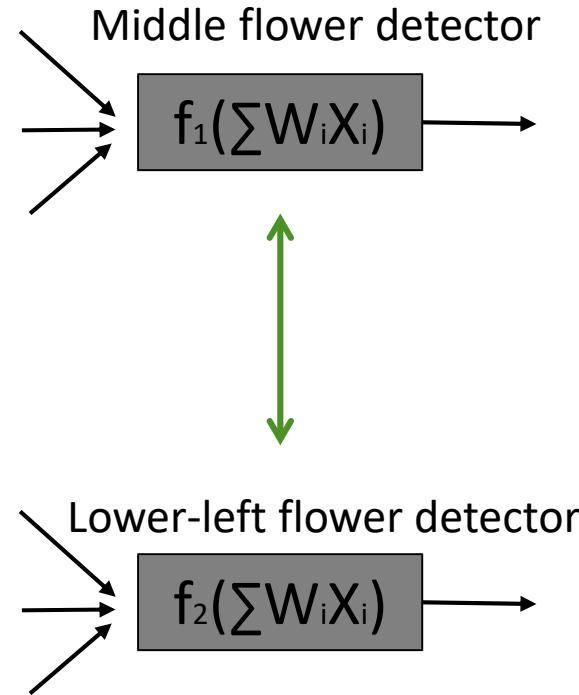
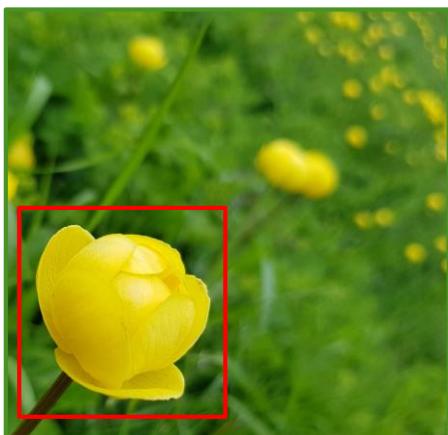
An RGB image can be represented as pixels  
 $(32 \times 32 \times 3)$



- Key patterns can be much smaller than the whole signal
  - Therefore a network does not have to see the whole signal to discover the pattern
  - Only need to be connected to small region ( $X_i$ )
  - Fewer parameters ( $W_i$ ) are required



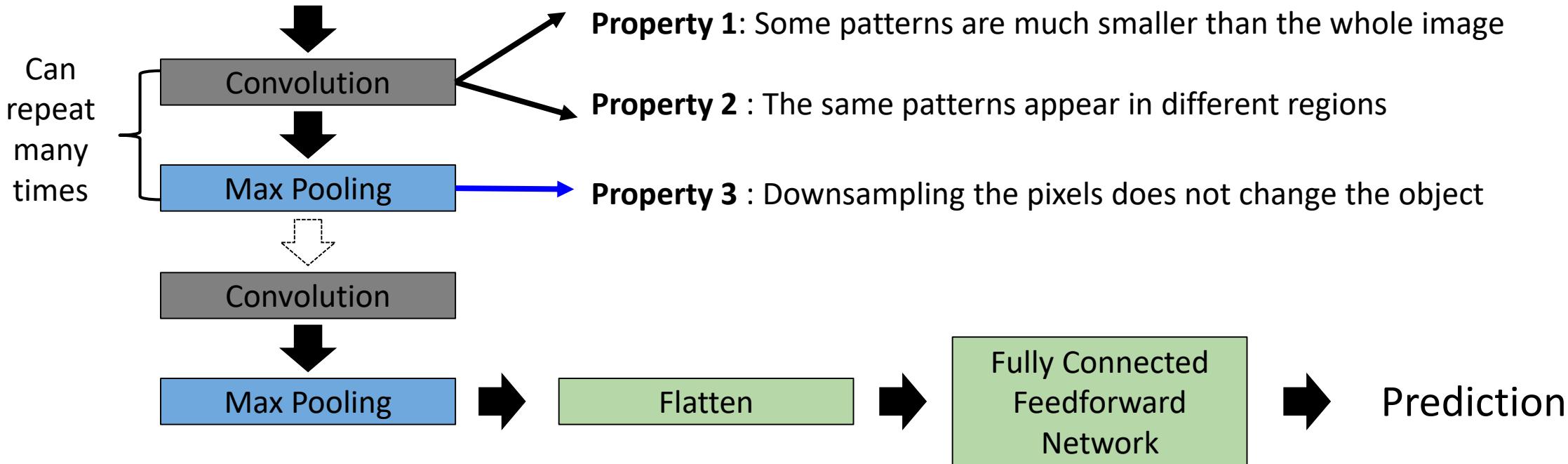
- Similar patterns appear in different regions



- Both networks are flower detectors
- They detect almost **the same** thing
- Therefore, they can use **the same** set of parameters

- **Downsampling does not change the key properties**
  - E.g., images are made **smaller** by downsampling
  - Less parameters for the network to process





## Convolution

- Convolution is the manipulation of two *signals* to form a third
- *Notation:* Convolution is denoted by the star  $*$  operator

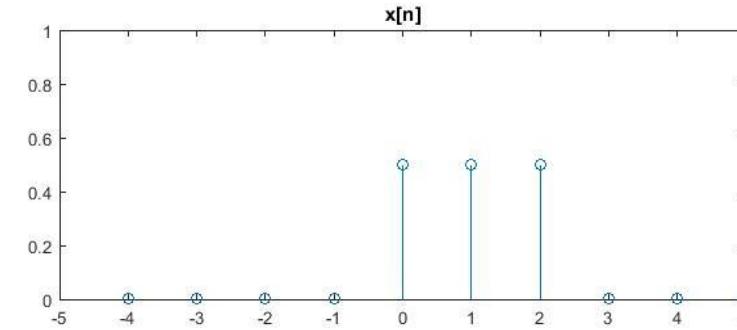
$$y[n] = x[n] * h[n] = \sum_{i=-\infty}^{\infty} x[i]h[n-i]$$

$$y[n] = x[m, n] * h[m, n] = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} x[i, j]h[m - i, n - j]$$

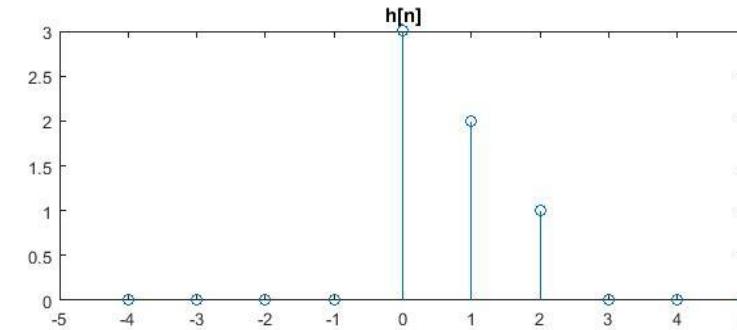
## Convolution Worked Example:

- Calculate and plot  $y[n] = x[n] * h[n]$  when:

$$x[n] = \begin{cases} 0.5 & \text{for } n = 0, 1, 2 \\ 0 & \text{otherwise} \end{cases}$$



$$h[n] = \begin{cases} 3, 2, 1 & \text{for } n = 0, 1, 2 \\ 0 & \text{otherwise} \end{cases}$$



## Convolution Worked Example:

- Convolution can be expressed as an inner product

$$\begin{aligned}y[n] &= x[n] * h[n] \\y[n] &= \sum_{k=-\infty}^{\infty} x[k] h[n-k] \\y[n] &= \sum_{k=-\infty}^{\infty} x[k] \tilde{h}[k-n] \\y[n] &= \langle x, \bar{h}_n \rangle\end{aligned}$$

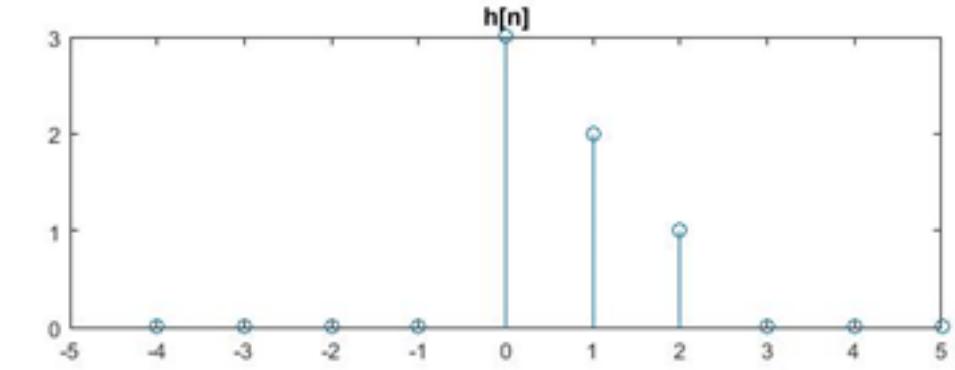
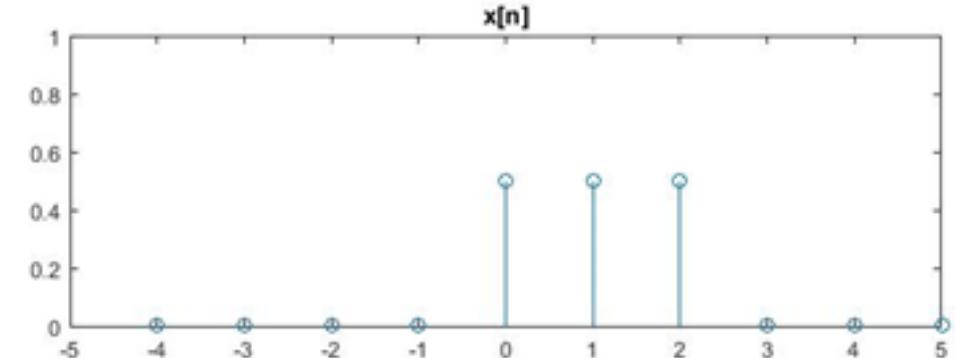
- $\bar{h} \equiv h[-i]$  is a **time reversed** version  $h \equiv h[i]$  (**flip**)
- $n^{th}$  output obtained by *delaying*  $\bar{h}$   $n$  times (**shift**) and taking *inner product* (**multiple**)

## Convolution Worked Example:

- $x = \{0.5, 0.5, 0.5\}$ ,  $h = \{3, 2, 1\}$

### Flip, shift, multiple:

- $y[0] = \langle x, \bar{h}_0 \rangle$  where  $\bar{h}_0 = \{3, 0, 0\}$
- $y[0] = 1.5$
- $y[1] = \langle x, \bar{h}_1 \rangle$  where  $\bar{h}_1 = \{2, 3, 0\}$
- $y[1] = 2.5$
- $y[2] = \langle x, \bar{h}_2 \rangle$  where  $\bar{h}_2 = \{1, 2, 3\}$
- $y[2] = 3.0$

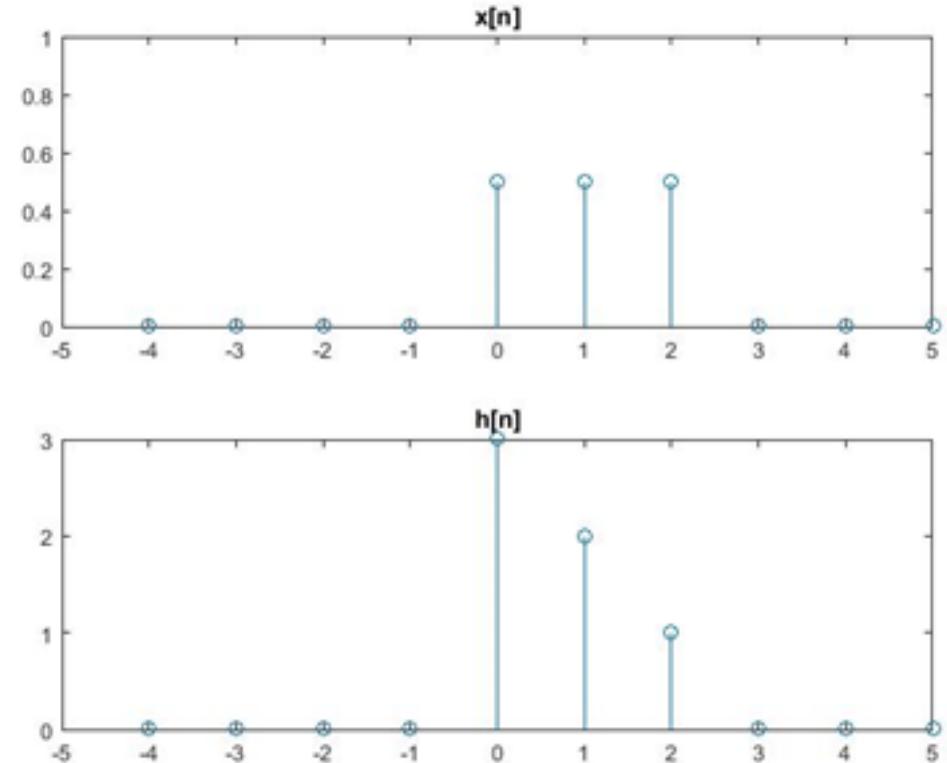


## Convolution Worked Example:

- $x = \{0.5, 0.5, 0.5\}$ ,  $h = \{3, 2, 1\}$

### Flip, shift, multiple:

- $y[3] = \langle x, \bar{h}_3 \rangle$  where  $\bar{h}_0 = \{0, 1, 2\}$
- $y[3] = 1.5$
- $y[4] = \langle x, \bar{h}_1 \rangle$  where  $\bar{h}_1 = \{0, 0, 1\}$
- $y[4] = 0.5$



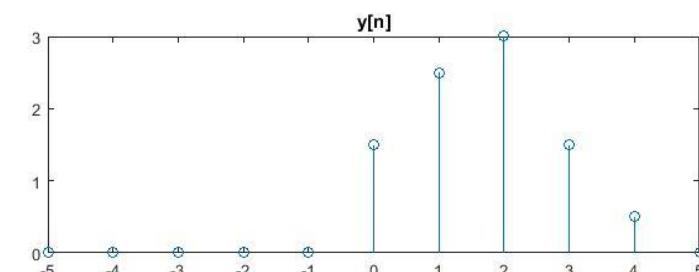
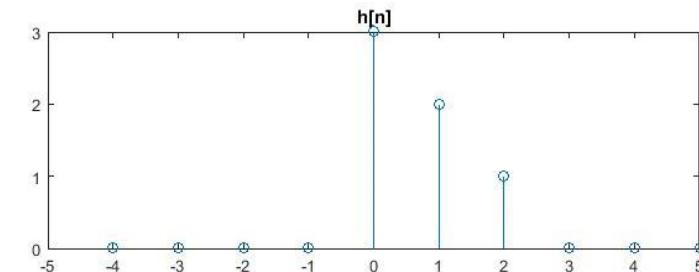
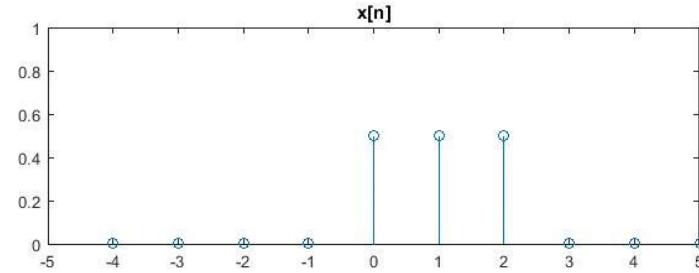
Therefore  $y = x * h = \{1.5, 2.5, 3, 1.5, 0.5\}$

## Convolution Worked Example:

$$x[n] = \begin{cases} 0.5 & \text{for } n = 0, 1, 2 \\ 0 & \text{otherwise} \end{cases}$$

$$h[n] = \begin{cases} 3, 2, 1 & \text{for } n = 0, 1, 2 \\ 0 & \text{otherwise} \end{cases}$$

$$y[n] = x[n] * h[n]$$

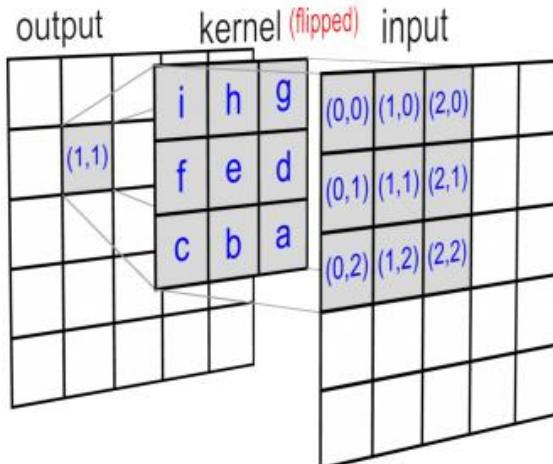


- Convolution in two dimensions

m	-1	0	1
-1	a	b	c
0	d	e	f
1	g	h	i

### Example 2D Kernel

- The origin is located at the middle of kernel



### 2D Convolution

- The kernel is flipped in both the horizontal and vertical directions

- **Convolution in two dimensions**

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline -1 & -2 & -1 \\ \hline 0 & 0 & 0 \\ \hline 1 & 2 & 1 \\ \hline \end{array}$$

=

$$\begin{array}{|c|c|c|} \hline -13 & -20 & -17 \\ \hline -18 & -24 & -18 \\ \hline 13 & 20 & 17 \\ \hline \end{array}$$

$$\begin{aligned}y[1,1] &= \sum_j \sum_i x[i,j] h[1-i, 1-j] \\&= x[0,0]h[1,1] + x[1,0]h[0,1] + x[2,0]h[-1,1] \\&\quad + x[0,1]h[1,0] + x[1,1]h[0,0] + x[2,1]h[-1,0] \\&\quad + x[0,2]h[1,-1] + x[1,0]h[0,-1] + x[2,2]h[-1,-1] \\&= 1 \cdot 1 + 2 \cdot 2 + 3 \cdot 1 \\&= 4 \cdot 0 + 5 \cdot 0 + 6 \cdot 0 \\&= 7 \cdot -1 + 8 \cdot -2 + 9 \cdot -1 \\&= -24\end{aligned}$$

- **Convolution and Correlation**

Input

$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$

Filter

$F_{11}$	$F_{12}$
$F_{21}$	$F_{22}$

**Correlation**

$$o_{11} = F_{11}x_{11} + F_{12}x_{12} + F_{21}x_{21} + F_{22}x_{22}$$

$$o_{12} = F_{11}x_{12} + F_{12}x_{13} + F_{21}x_{22} + F_{22}x_{23}$$

$$o_{21} = F_{11}x_{21} + F_{12}x_{22} + F_{21}x_{31} + F_{22}x_{32}$$

$$o_{22} = F_{11}x_{22} + F_{12}x_{23} + F_{21}x_{32} + F_{22}x_{33}$$

- **Convolution and Correlation**

$O_{11}$		
$F_{11}x_{11}$	$F_{12}x_{12}$	$x_{13}$
$F_{21}x_{21}$	$F_{22}x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$



$O_{12}$		
$x_{11}$	$F_{11}x_{12}$	$F_{12}x_{13}$
$x_{21}$	$F_{21}x_{22}$	$F_{22}x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$



$O_{21}$		
$x_{11}$	$x_{12}$	$x_{13}$
$F_{11}x_{21}$	$F_{12}x_{22}$	$x_{23}$
$F_{21}x_{31}$	$F_{22}x_{32}$	$x_{33}$



$O_{22}$		
$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$F_{11}x_{22}$	$F_{12}x_{23}$
$x_{31}$	$F_{21}x_{32}$	$F_{22}x_{33}$



- **Convolution and Correlation**

Input

$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$

Filter

$F_{11}$	$F_{12}$
$F_{21}$	$F_{22}$

## Convolution

$$o_{11} = F_{22}x_{11} + F_{21}x_{12} + F_{12}x_{21} + F_{11}x_{22}$$

$$o_{12} = F_{22}x_{12} + F_{21}x_{13} + F_{12}x_{22} + F_{11}x_{23}$$

$$o_{21} = F_{22}x_{21} + F_{21}x_{22} + F_{12}x_{31} + F_{11}x_{32}$$

$$o_{22} = F_{22}x_{22} + F_{21}x_{23} + F_{12}x_{32} + F_{11}x_{33}$$

- **Convolution and Correlation**

- Convolution is correlation with the filter rotated 180 degrees

## Correlation

Input

$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$

## Filter

$F_{11}$	$F_{12}$
$F_{21}$	$F_{22}$

## Convolution

Input

$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$

## Filter

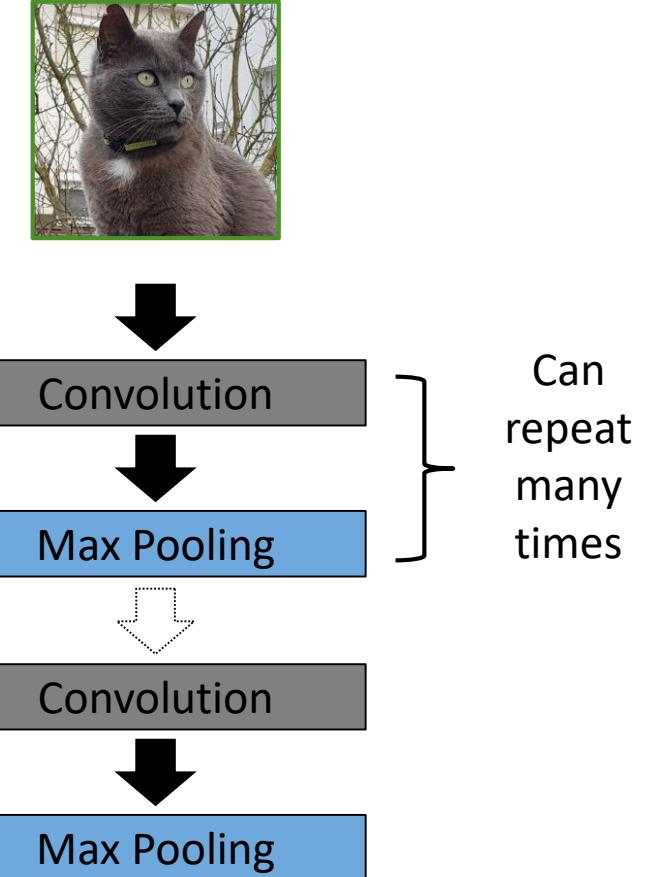
$F_{22}$	$F_{21}$
$F_{12}$	$F_{11}$

- **Convolutional Layers**

- **Property 1:** Some patterns are much smaller than the whole image
- **Property 2 :** The same patterns appear in different regions

- **Pooling Layers**

- **Property 3 :** Downsampling the pixels does not change the object



0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
<b>1</b>	0	0	0	0	<b>1</b>
0	<b>1</b>	0	0	<b>1</b>	0
0	0	<b>1</b>	<b>1</b>	0	0

6 x 6 image

Apply **small filters** to detect small patterns

Each filter has a size of **3 x 3**

-1	<b>1</b>	-1
-1	<b>1</b>	-1
-1	<b>1</b>	-1

Filter 1

<b>1</b>	-1	-1
-1	<b>1</b>	-1
-1	-1	<b>1</b>

Filter 2

⋮ ⋮ ⋮ ⋮ ⋮

**Note:** Only the size of the filters is specified; the weights are initialised to arbitrary values before the start of training.

The weights of the **filters are learned** through the CNN training process

- **Key Parameters**

- **Filter size** – defines the height and width of the filter kernel
  - E.g., a filter kernel of size three would have nine weights
- **Stride** – determines the number of steps to move in each spatial direction while performing convolution.
- **Padding** – appends zeroes to the boundary of an image to control the size of the output of convolution
  - When we convolve an image of a specific size by a filter, the resulting image is generally smaller than the original image

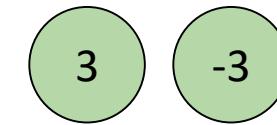
stride = 1

0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

6 x 6 image

-1	1	-1
-1	1	-1
-1	1	-1

Filter 1



Compute the **dot product** between the filter and a small 3 x 3 chunk of the image

**Note: We are technically performing correlation not convolution**

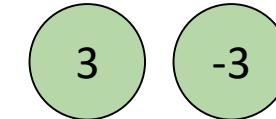
stride = 2

0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

6 x 6 image

-1	1	-1
-1	1	-1
-1	1	-1

Filter 1



We set stride = 1 below

Compute the **dot product** between the filter and a small 3 x 3 chunk of the image

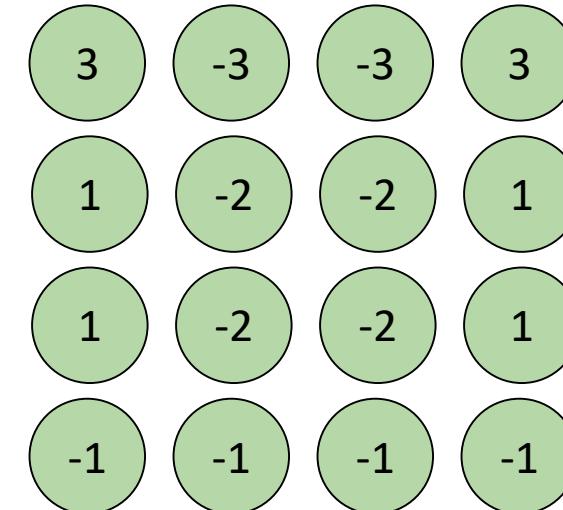
stride = 1

0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
<b>1</b>	0	0	0	0	<b>1</b>
0	<b>1</b>	0	0	<b>1</b>	0
0	0	<b>1</b>	<b>1</b>	0	0

6 x 6 image

-1	<b>1</b>	-1
-1	<b>1</b>	-1
-1	<b>1</b>	-1

Filter 1

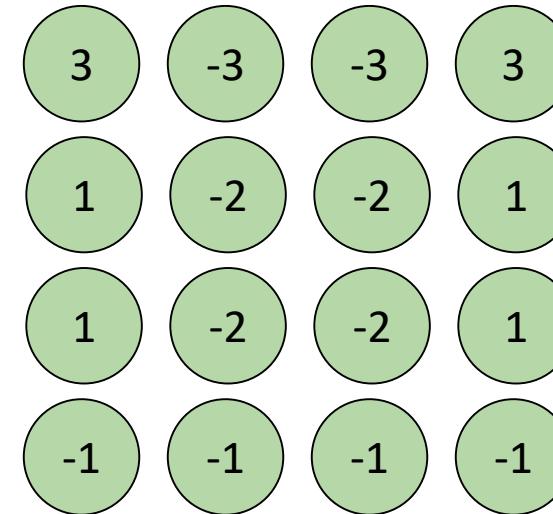


4 x 4 image

stride = 1, filter size = 3

0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

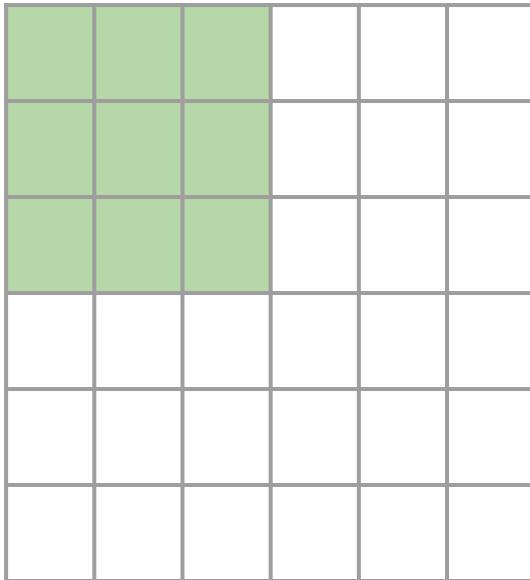
6 x 6 image



4 x 4 image

$$\text{output size: } (6 - 3) / 1 + 1 = 4$$

filter size =  $F$



$N \times N$  image

output size:  $(N - F) / \text{stride} + 1$

for example:  $N = 6, F = 3$

$\text{stride} = 1 \rightarrow (6-3)/1 + 1 = 4$

$\text{stride} = 2 \rightarrow (6-3)/2 + 1 = 2.5 : \backslash$

$\text{stride} = 3 \rightarrow (6-3)/3 + 1 = 2$

**zero-padding to the border**

0	0	0	0	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

**N x N** image

For example:  $N = 6$ ,  $F = 3$ , stride = 1

*Without 0-padding:*

output size is  $(6-3)/1 + 1 = 4$

*With 0-padding with 1 pixel border:*

output size is  $(6-3+2 \times 1)/1 + 1 = 6$

The output size is then the same as the input!

## zero-padding to the border

0	0	0	0	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

$N \times N$  image

In general, stride=1, filters of size  $F \times F$ ,  
then zero-padding with  $(F-1)/2$ ,  
to preserve size spatially

e.g.  $F = 3 \rightarrow$  zero pad with 1 pixel to the border  
 $F = 5 \rightarrow$  zero pad with 2 pixels to the border  
 $F = 7 \rightarrow$  zero pad with 3 pixels to the border

stride = 1

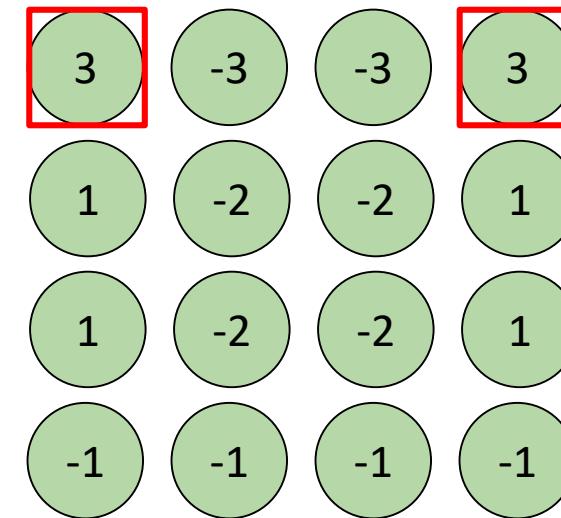
0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

6 x 6 image

-1	<b>1</b>	-1
-1	<b>1</b>	-1
-1	<b>1</b>	-1

Filter 1

detect a vertical line

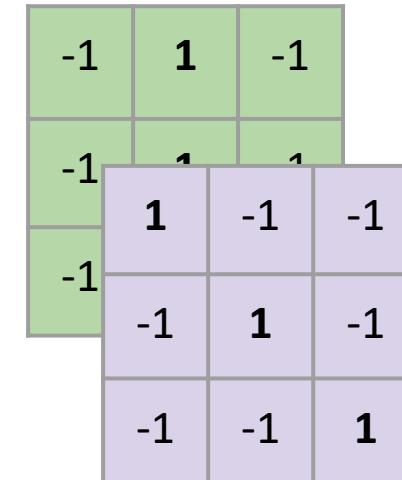


The same pattern in **different locations** are detected with the same filter

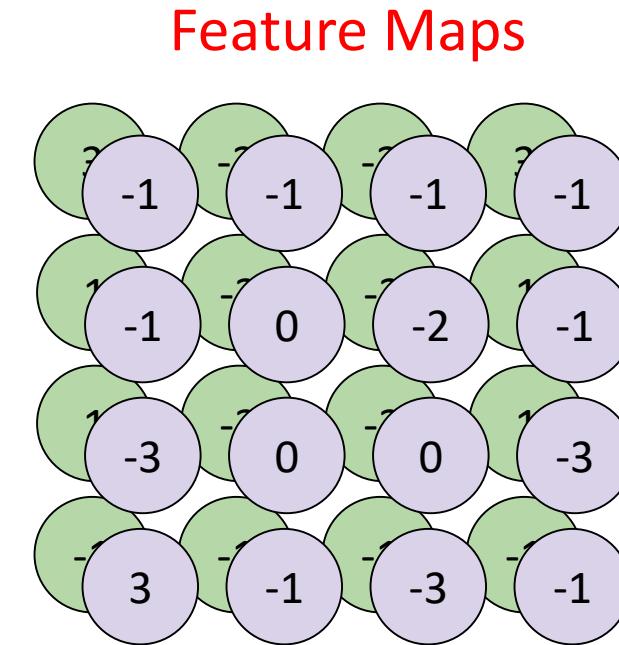
stride = 1

0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

6 x 6 image



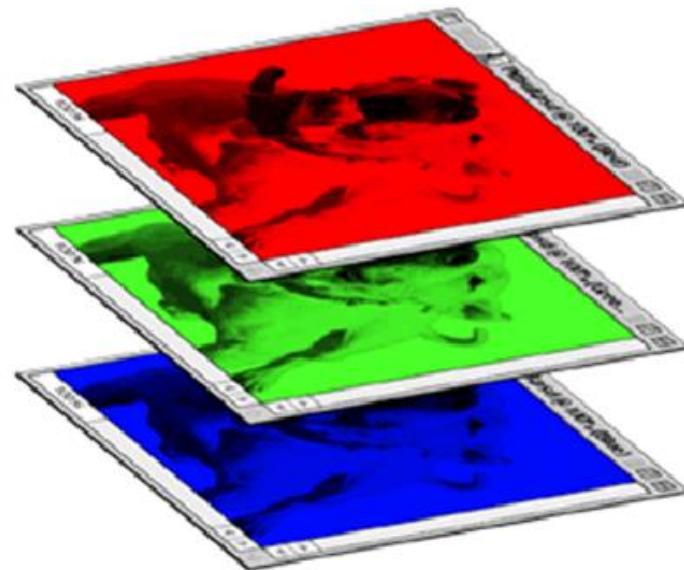
Filter 2



4 x 4 x (#filters)

Do the same process for every filter

RGB images



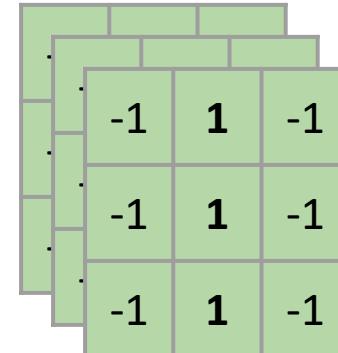
3 channels



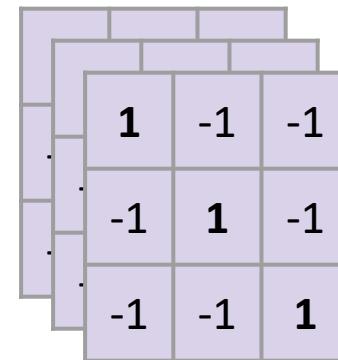
0	1	0	0	1	0	
0	0	1	0	0	1	0
0	0	1	0	0	1	0
1	0	1	0	0	1	0
0	1	0	0	0	0	1
0	0	1	0	0	1	0
0	0	1	1	0	0	

$6 \times 6 \times 3$

Filters always extend the full depth of the input volume



Filter 1

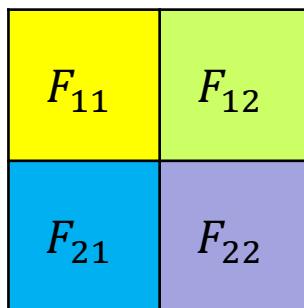


Filter 2

- **Forward pass:**

- Input  $X$  of size  $3 \times 3$  using a single filter  $F$  of size  $2 \times 2$ , no padding and  $stride = 1$ , generating an output  $O$  of size  $2 \times 2$

$X_{11}$	$X_{12}$	$X_{13}$
$X_{21}$	$X_{22}$	$X_{23}$
$X_{31}$	$X_{32}$	$X_{33}$



$O_{11}$	$O_{12}$
$O_{21}$	$O_{22}$

$$O_{11} = F_{11}x_{11} + F_{12}x_{12} + F_{21}x_{21} + F_{22}x_{22}$$

$$O_{12} = F_{11}x_{12} + F_{12}x_{13} + F_{21}x_{22} + F_{22}x_{23}$$

$$O_{21} = F_{11}x_{21} + F_{12}x_{22} + F_{21}x_{31} + F_{22}x_{32}$$

$$O_{22} = F_{11}x_{22} + F_{12}x_{23} + F_{21}x_{32} + F_{22}x_{33}$$

Image source: <https://medium.com>

- **Backward pass:**

- $\frac{\partial L}{\partial o}$  the loss gradient from the next layer
- $\frac{\partial L}{\partial X}$  the loss gradient to be passed back to the proceeding layer
- $\frac{\partial L}{\partial F}$  the loss gradient to update the filter

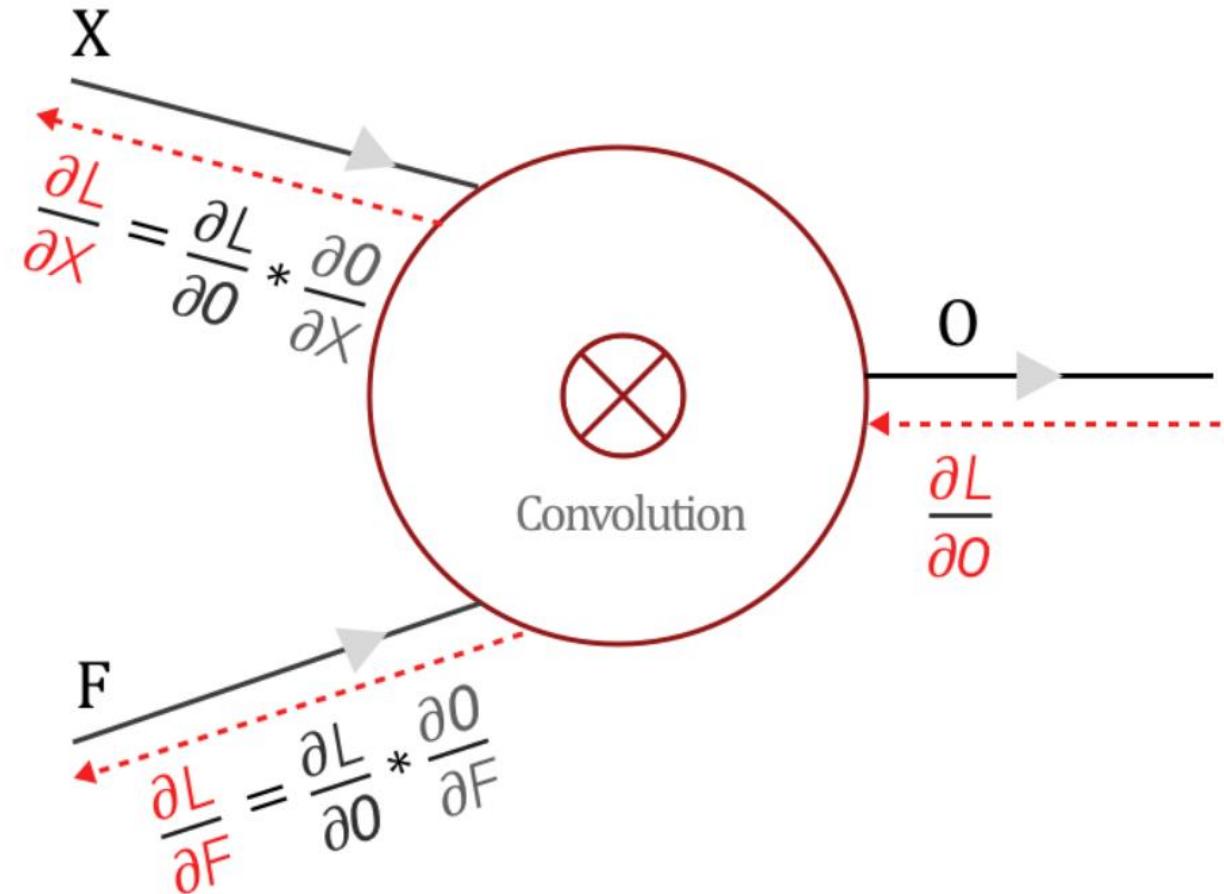


Image source: <https://medium.com>

- $\frac{\partial L}{\partial X}$  the loss gradient to be passed back

- Convolution between a  $F$  rotated by 180° and the loss gradient  $\frac{\partial L}{\partial O}$ 
  - I.e., “Proper” Convolution

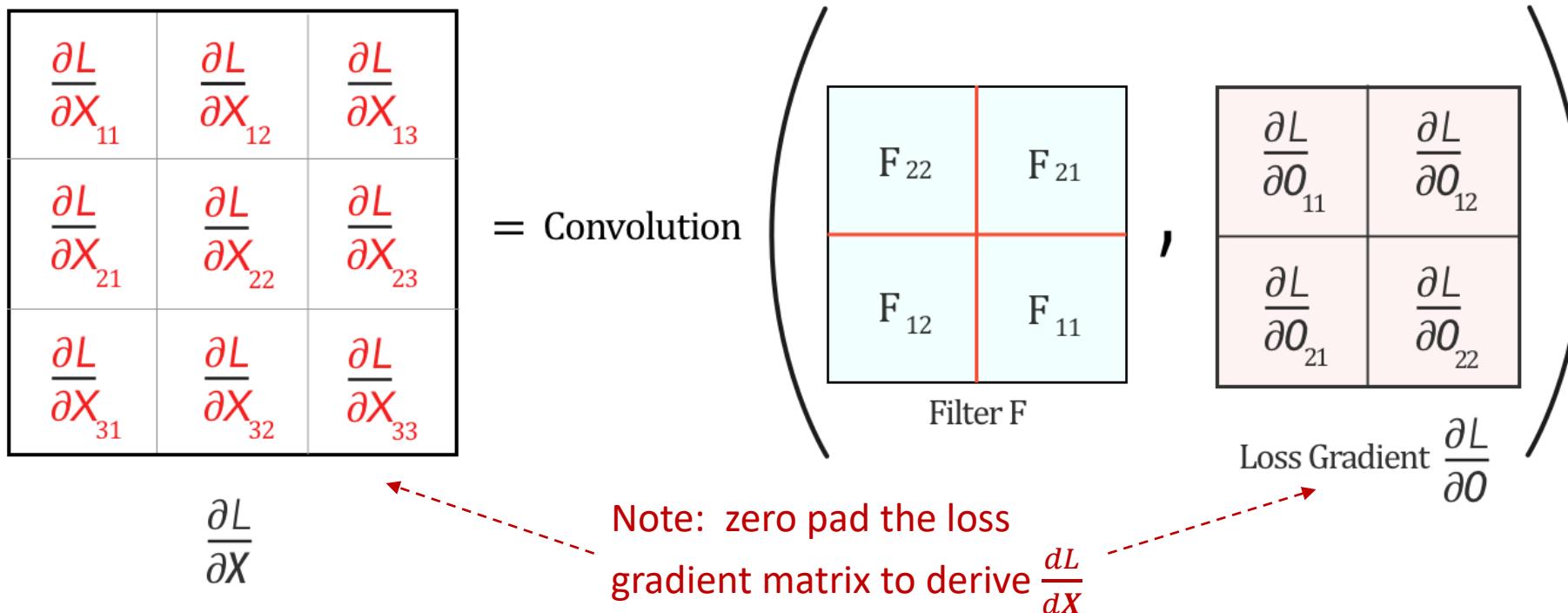


Image source: <https://medium.com>

- $\partial L / \partial F$  the loss gradient to update the filter
  - Convolution of input matrix  $X$  and loss gradient  $\frac{\partial L}{\partial o}$

$\frac{\partial L}{\partial F_{11}}$	$\frac{\partial L}{\partial F_{12}}$
$\frac{\partial L}{\partial F_{21}}$	$\frac{\partial L}{\partial F_{22}}$

= Convolution

$X_{11}$	$X_{12}$	$X_{13}$
$X_{21}$	$X_{22}$	$X_{23}$
$X_{31}$	$X_{32}$	$X_{33}$

$\frac{\partial L}{\partial o_{11}}$	$\frac{\partial L}{\partial o_{12}}$
$\frac{\partial L}{\partial o_{21}}$	$\frac{\partial L}{\partial o_{22}}$

## Key Parameters:

- Accepts an input of size  $W_1 \times H_1 \times D_1$
- Requires 4 hyperparameters:
  - Number of filters  $K$
  - Size of the filters  $F$
  - The stride  $S$
  - The amount of zero padding  $P$
- Produce an output of size  $W_2 \times H_2 \times D_2$ , where
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$
  - $D_2 = K$
- With parameter sharing, it introduces  $F \times F \times D_1$  weights per filter, for a total of  $(F \times F \times D_1) \times K$  weights and  $K$  biases.

## Common settings:

$K$ : powers of 2, such as 32, 64, 128, 512

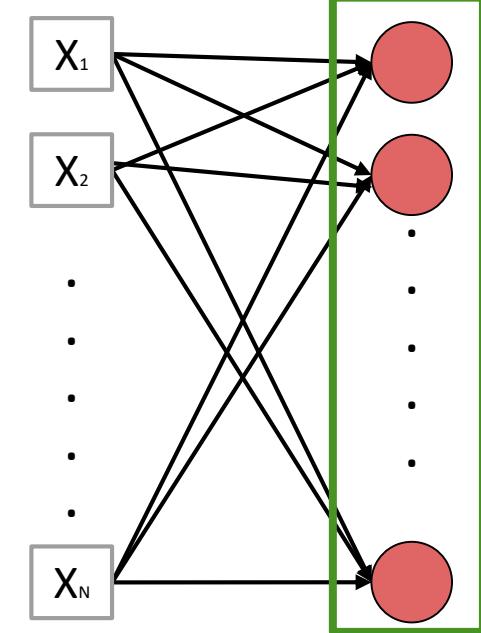
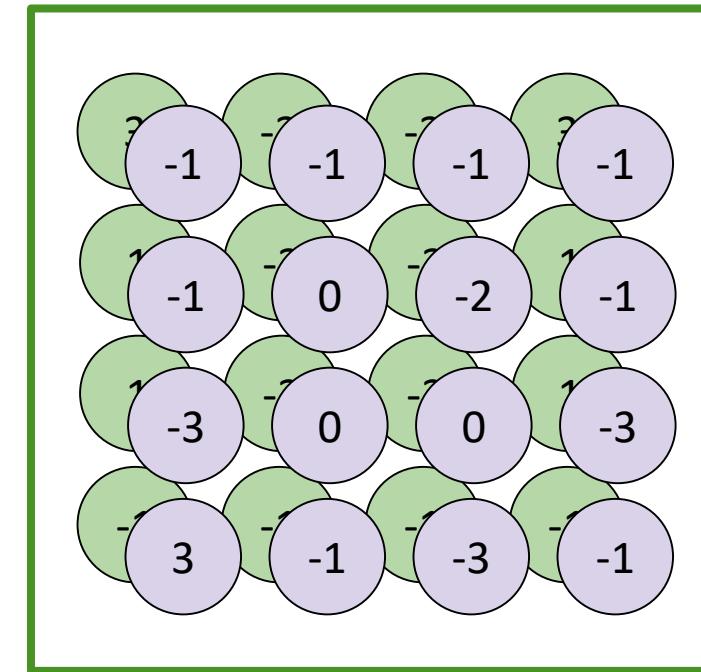
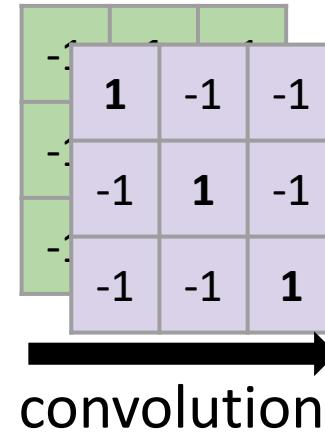
$F = 3, S=1, P=2$

$F = 5, S=1, P=2$

... ...

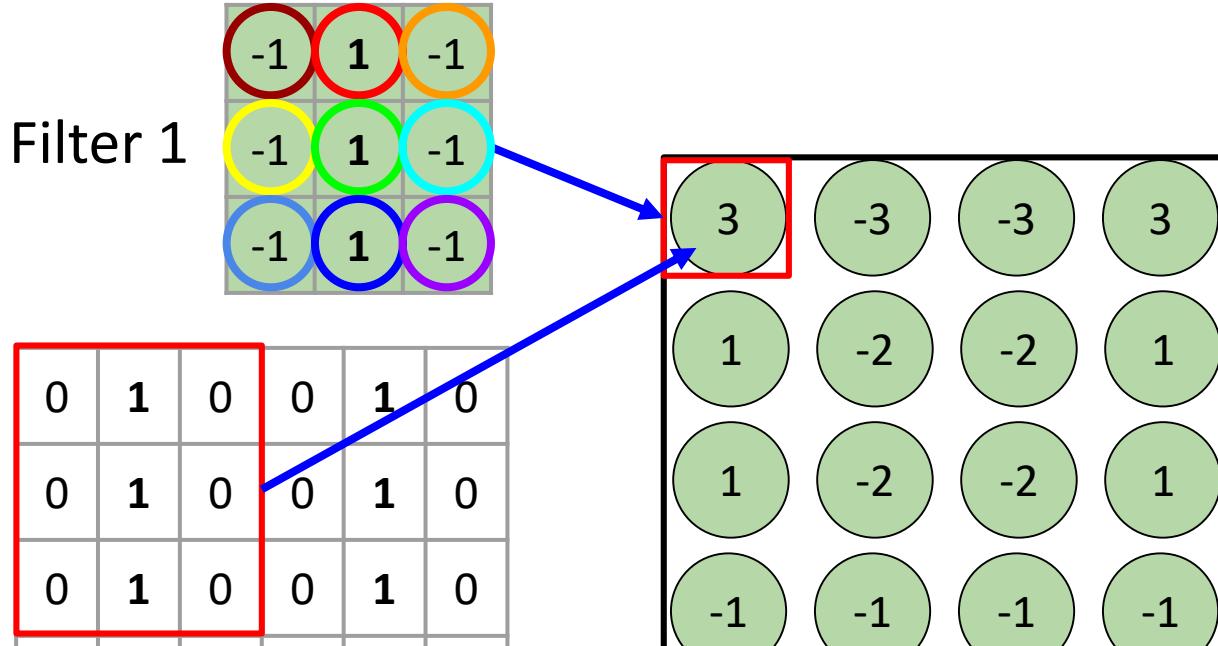
# Convolution vs Fully Connected

0	1	0	0	1	0
0	1	0	0	1	0
0	1	0	0	1	0
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

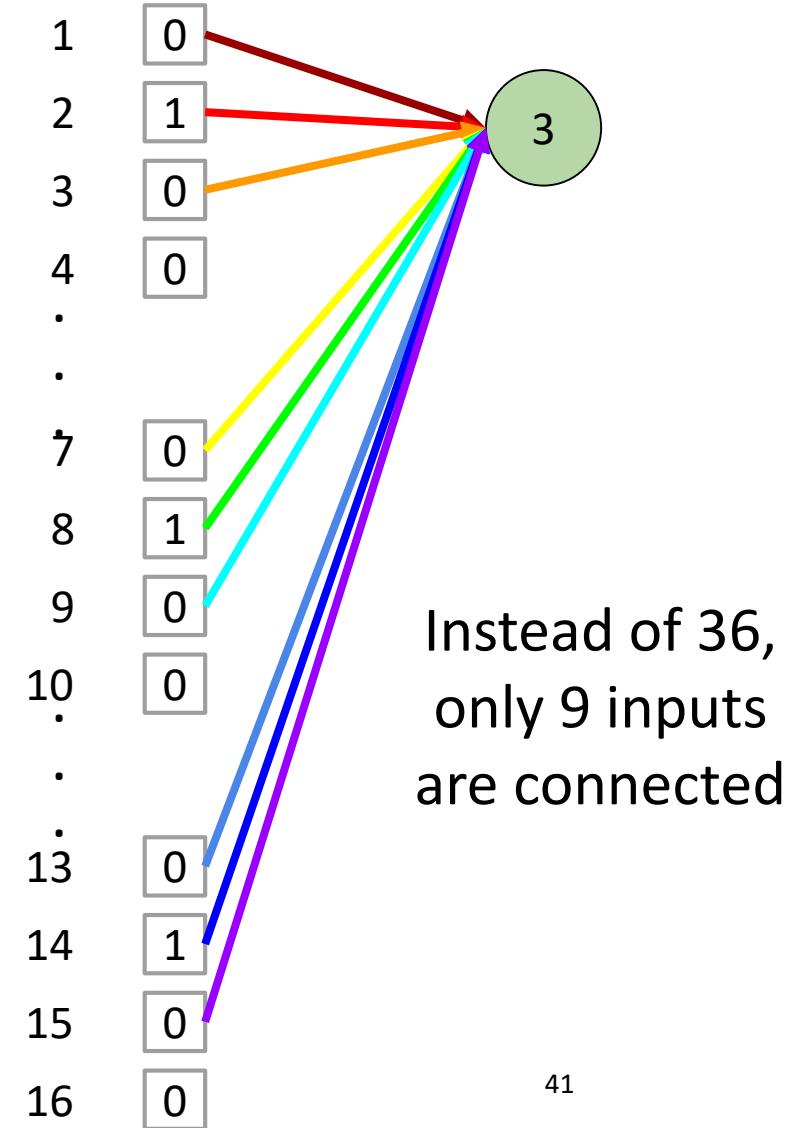


fully connected

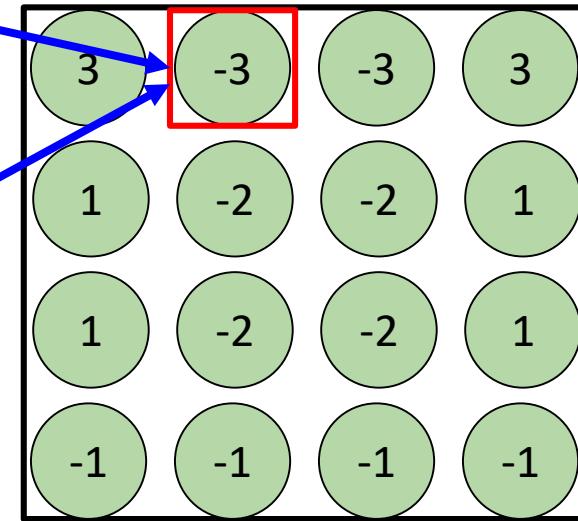
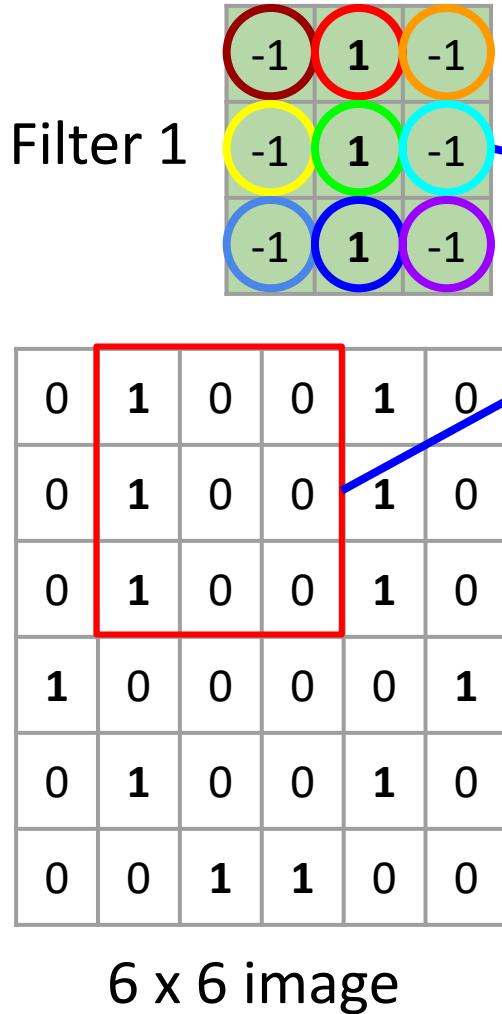
# Convolution vs Fully Connected



Less parameters to learn!

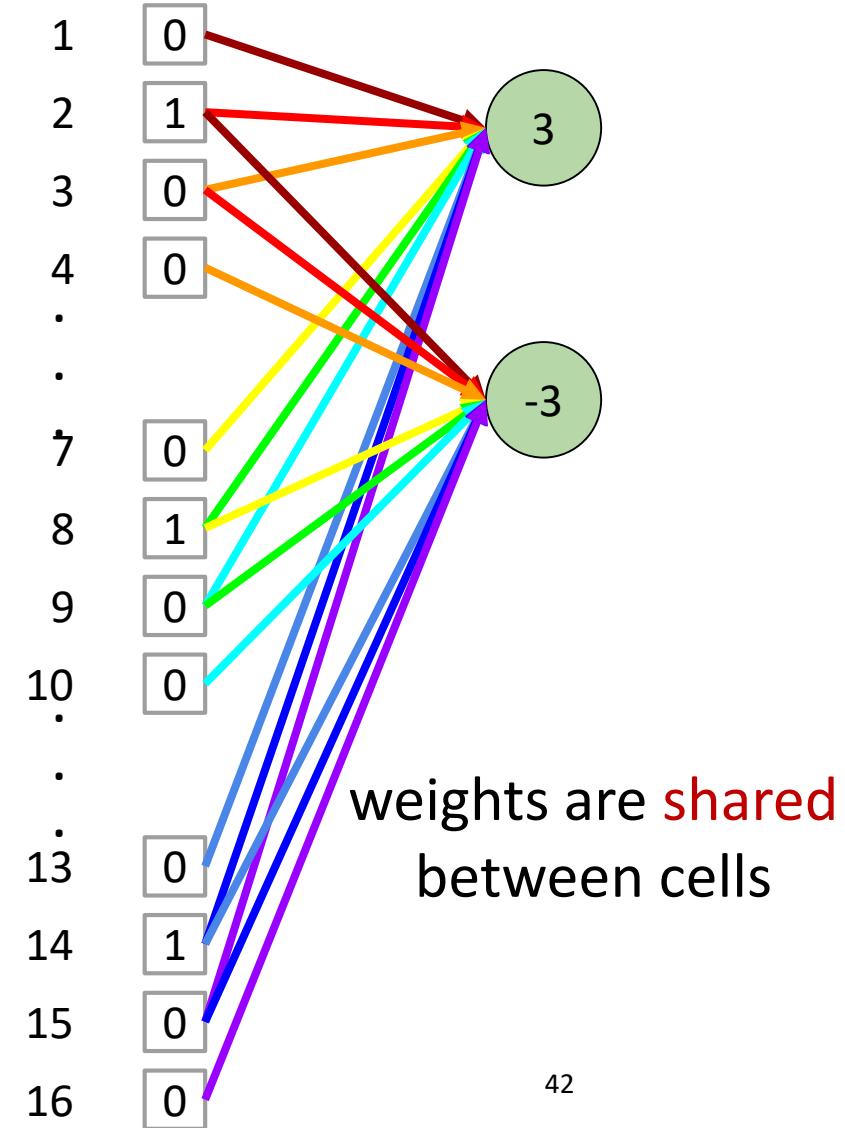


# Convolution vs Fully Connected



Less parameters to learn!

Even less parameters to learn!



- **The core idea**
  - Instead of having a large, dense linear layer with a connection from every input to every output, we have lots of small convolutional layers
    - Convolutional layers usually with fewer inputs and a single output
  - The result is a smaller subset of kernel predictions, which are used as input to the next layer.
  - Convolutional layers usually have many kernels

- **The core idea**

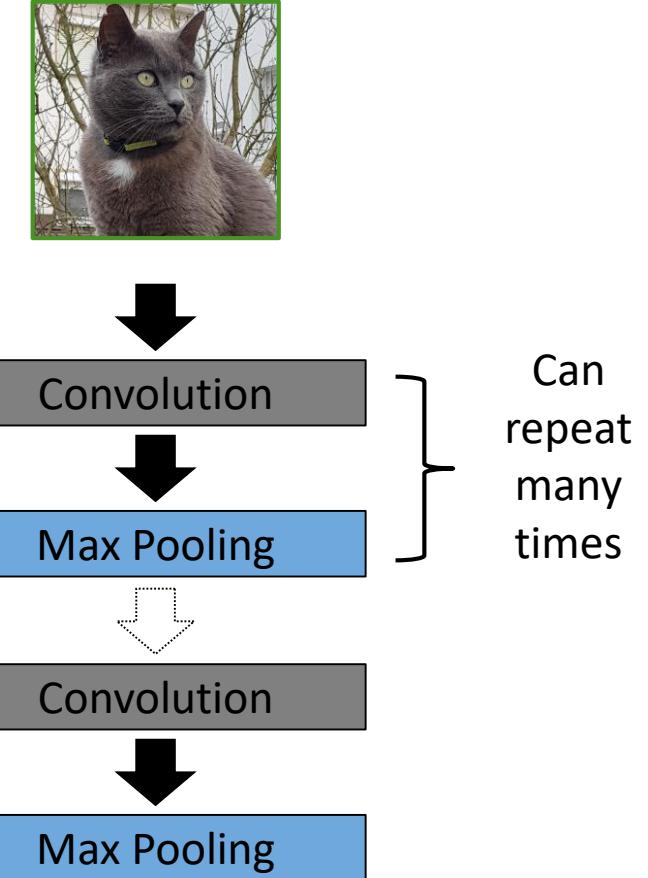
- Each kernel to learn a particular pattern and then search for the existence of that pattern somewhere in the image
- A single, small set of weights can train over a much larger set of training examples
- This changes the ratio of weights to datapoints on which those weights are being trained
- This has a powerful impact on the network, drastically reducing its ability to overfit to training data and increasing its ability to generalise

- **Convolutional Layers**

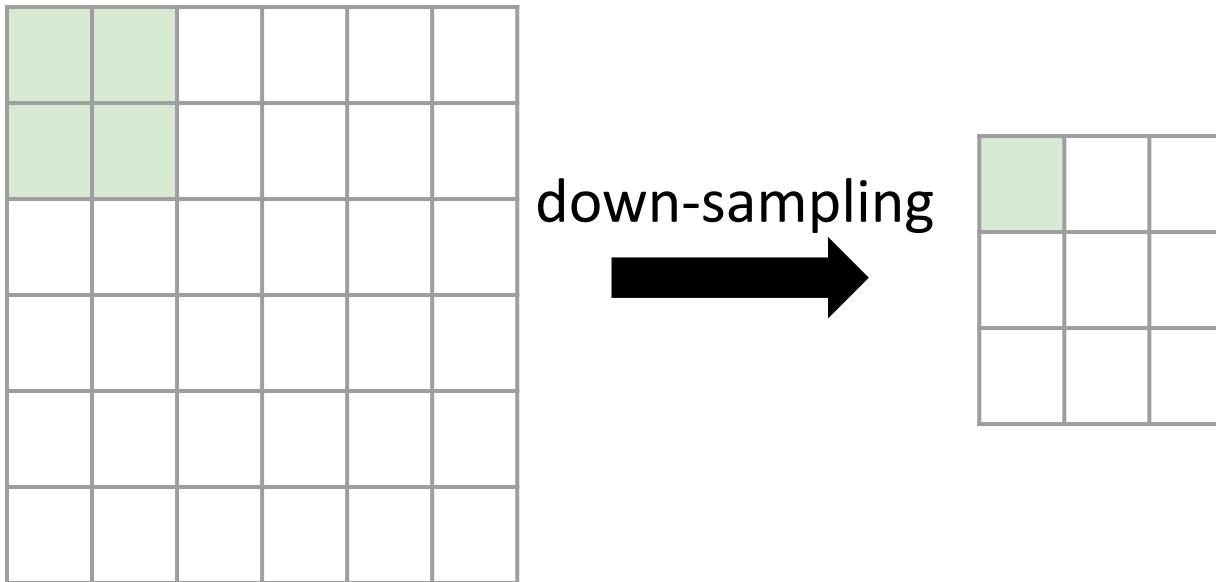
- **Property 1:** Some patterns are much smaller than the whole image
- **Property 2 :** The same patterns appear in different regions

- **Pooling Layers**

- **Property 3 :** Downsampling the pixels does not change the object



Pooling layers are usually present after a convolutional layer.  
They provide a **down-sampled** version of the convolution output.



In this example, a 2x2 region is used as input of the pooling.  
There are different types of pooling, the most used is **max pooling**.

# Max Pooling

Filter 1

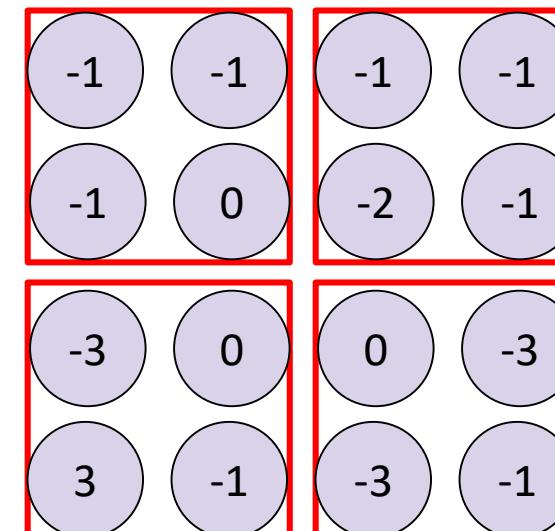
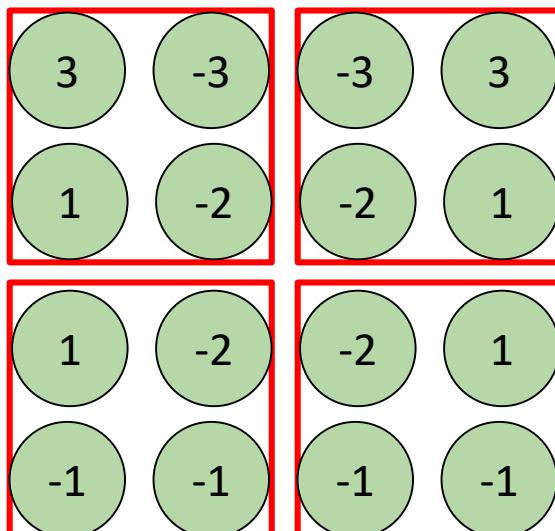
-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

1	-1	-1
-1	1	-1
-1	-1	1

Pooling size =  $2 \times 2$

Stride = 2



Operates over each  
feature map  
independently

Invariant to small  
differences in the input

0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
0	<b>1</b>	0	0	<b>1</b>	0
<b>1</b>	0	0	0	0	<b>1</b>
0	<b>1</b>	0	0	<b>1</b>	0
0	0	<b>1</b>	<b>1</b>	0	0

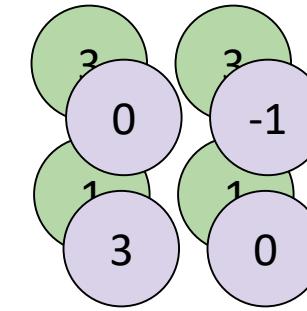
6 x 6 image



Convolution



Max Pooling



each filter is a channel

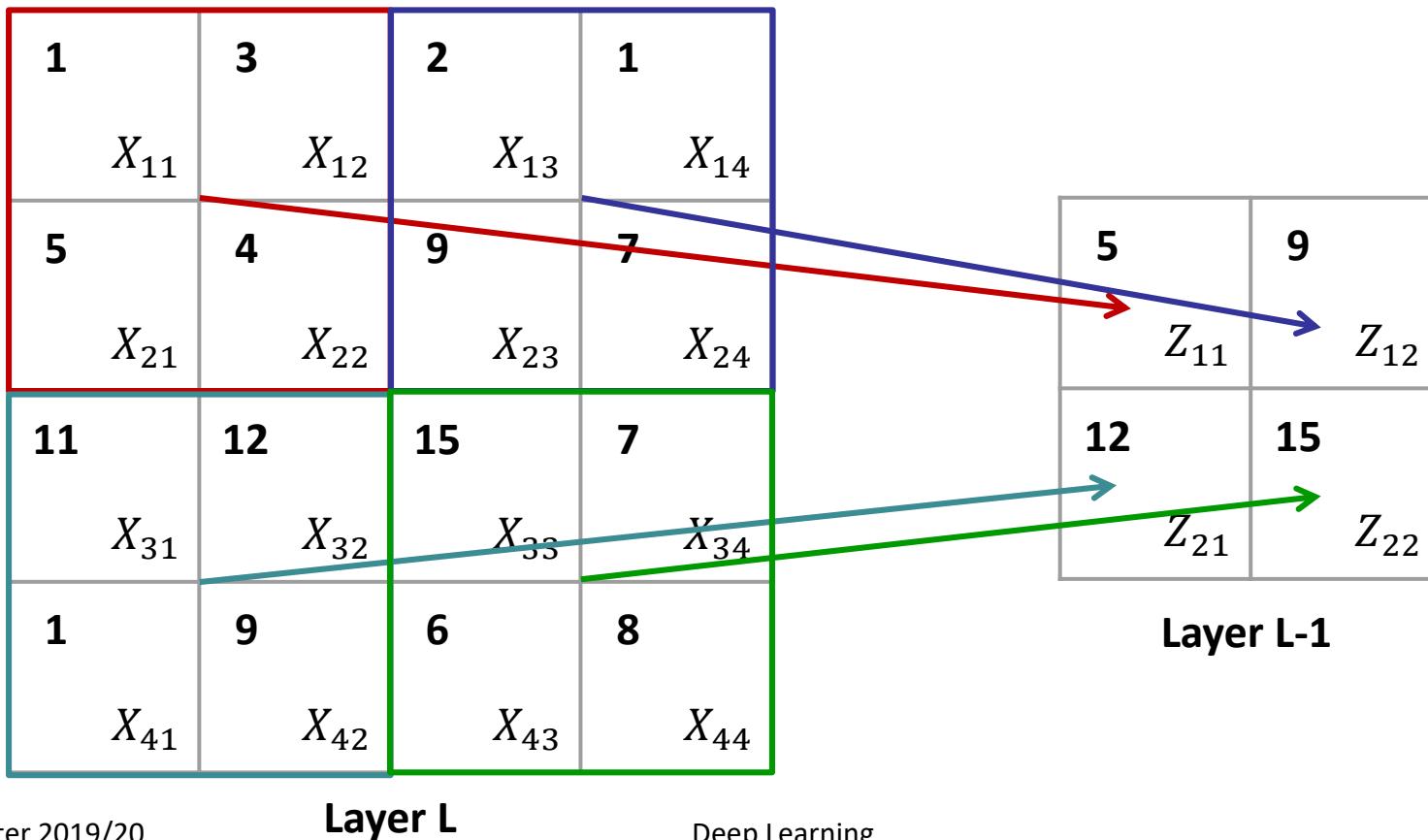
2 feature maps

each of size **2 x 2**

**Smaller and more manageable**

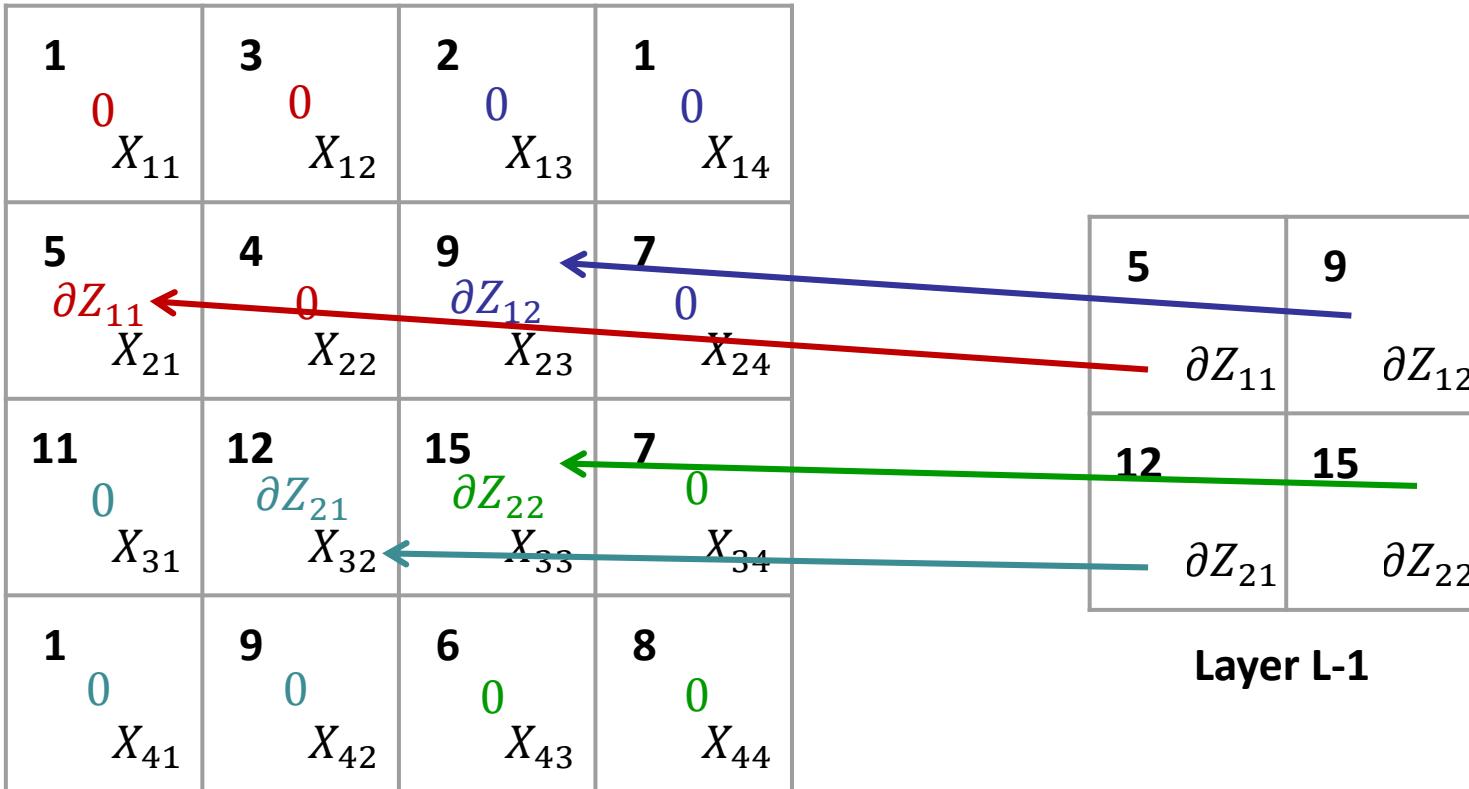
- **Worked Example**

- Receptive field for max pooling is  $2 \times 2$ , and the stride size is 2



- **Worked Example**

- Gradient gets passed to the cell it came from, all other cells get 0



## Key Parameters:

- Accepts an input of size  $W_1 \times H_1 \times D_1$
- Requires 2 hyperparameters:
  - Size of the filters  $F$
  - The stride  $S$
- Produce an output of size  $W_2 \times H_2 \times D_2$ , where
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- It introduces **zero** learnable parameters since it computes a fixed function of the input.

## Common settings:

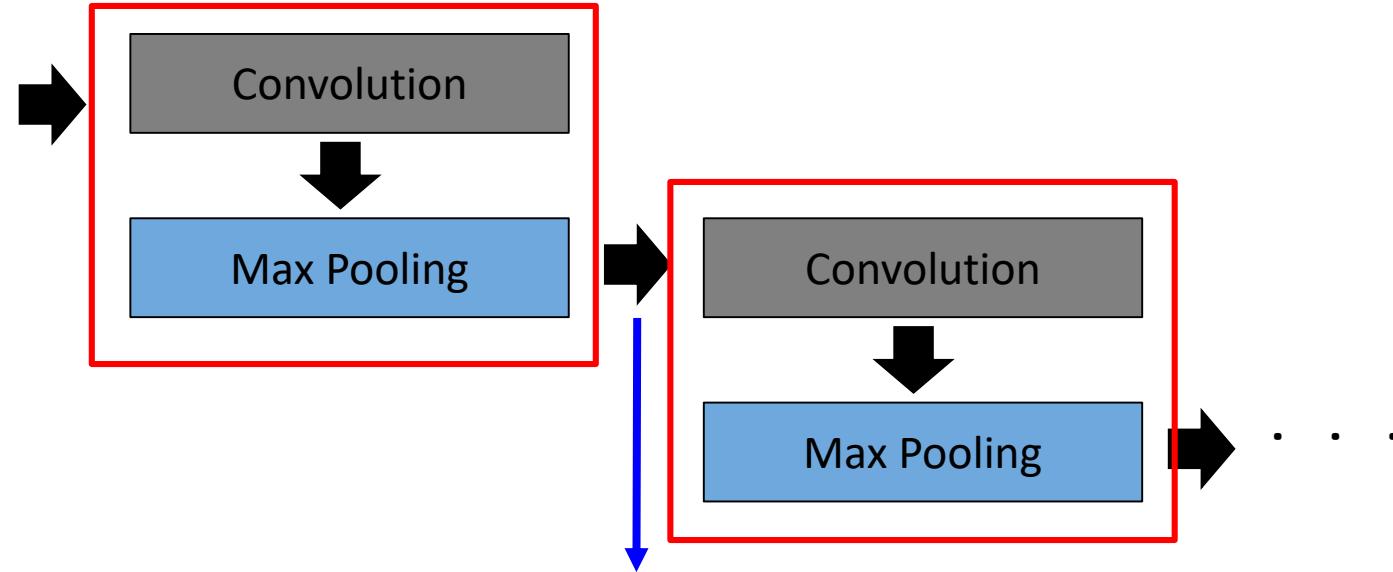
$F = 2, S=2$

$F = 3, S=2$

... ...

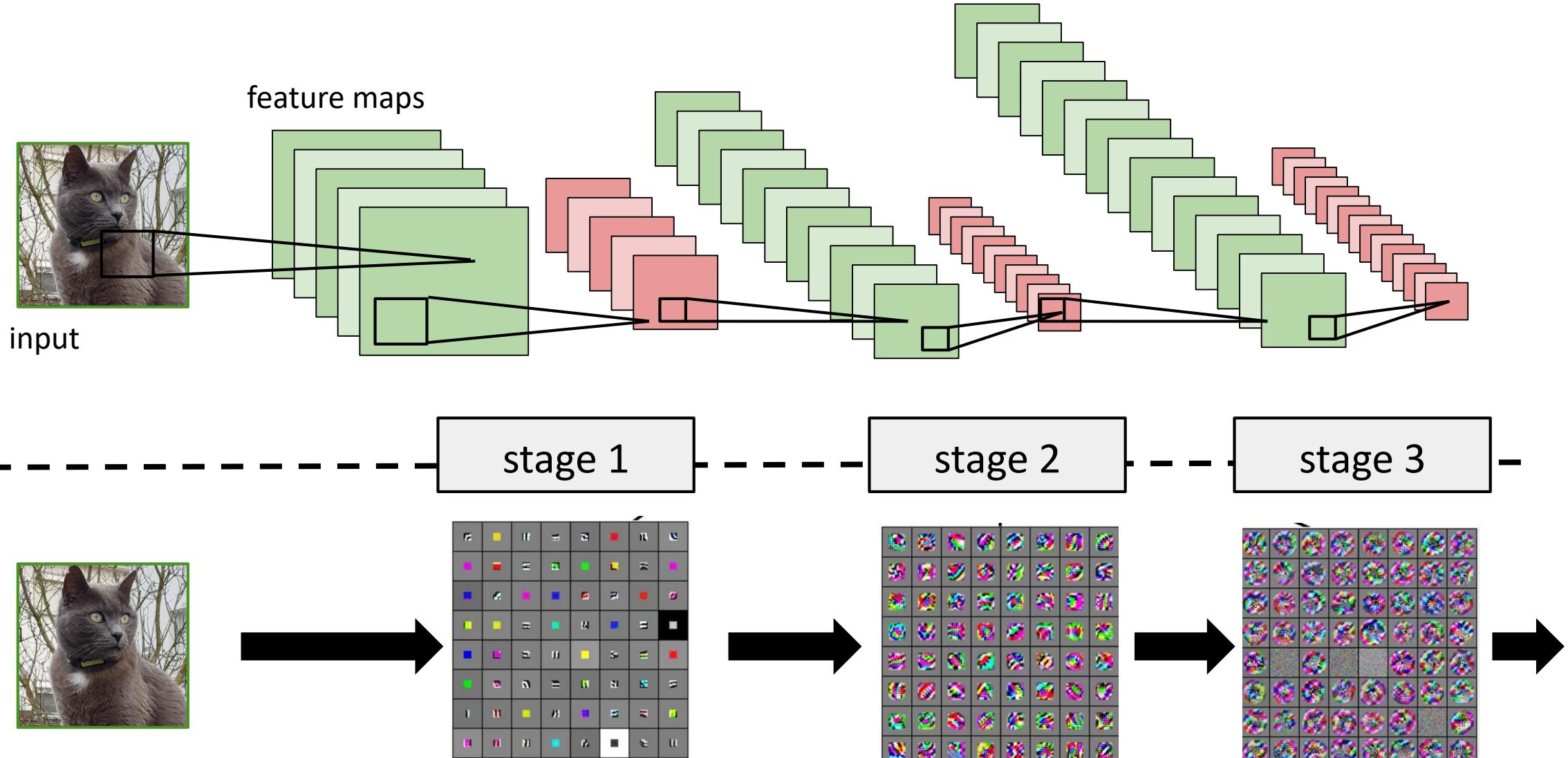


**Can be repeat many times**



Output can be regarded as new images:

- Smaller than the original images
- The depth of new images is the number of filters



- **Flatten the feature maps**

- Takes the output of the previous layers, and “flattens” them and turns them into a single vector
- This representation is then feed into a standard feedforward network to perform classification

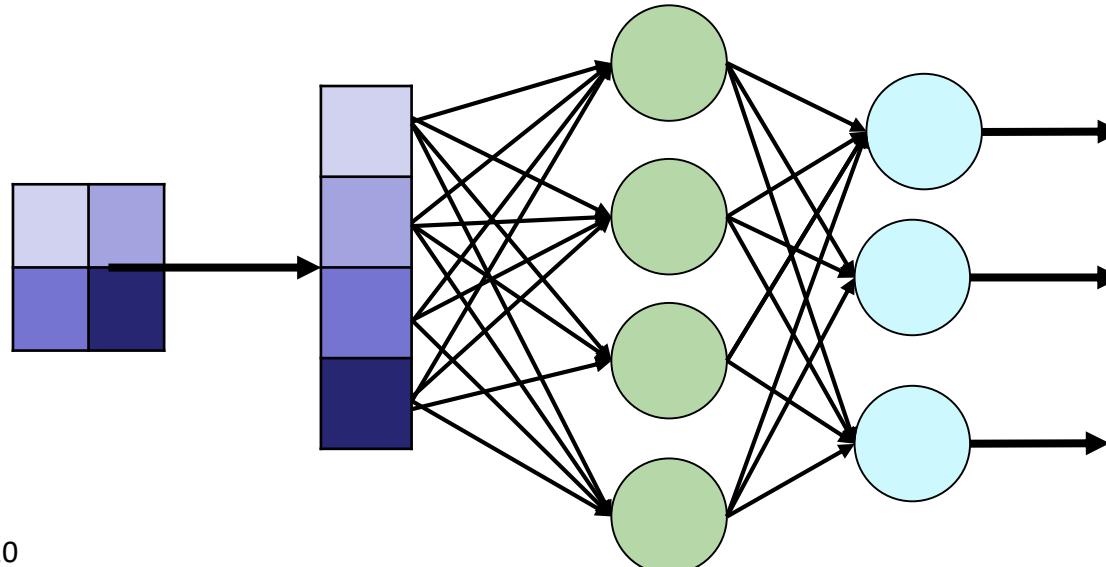
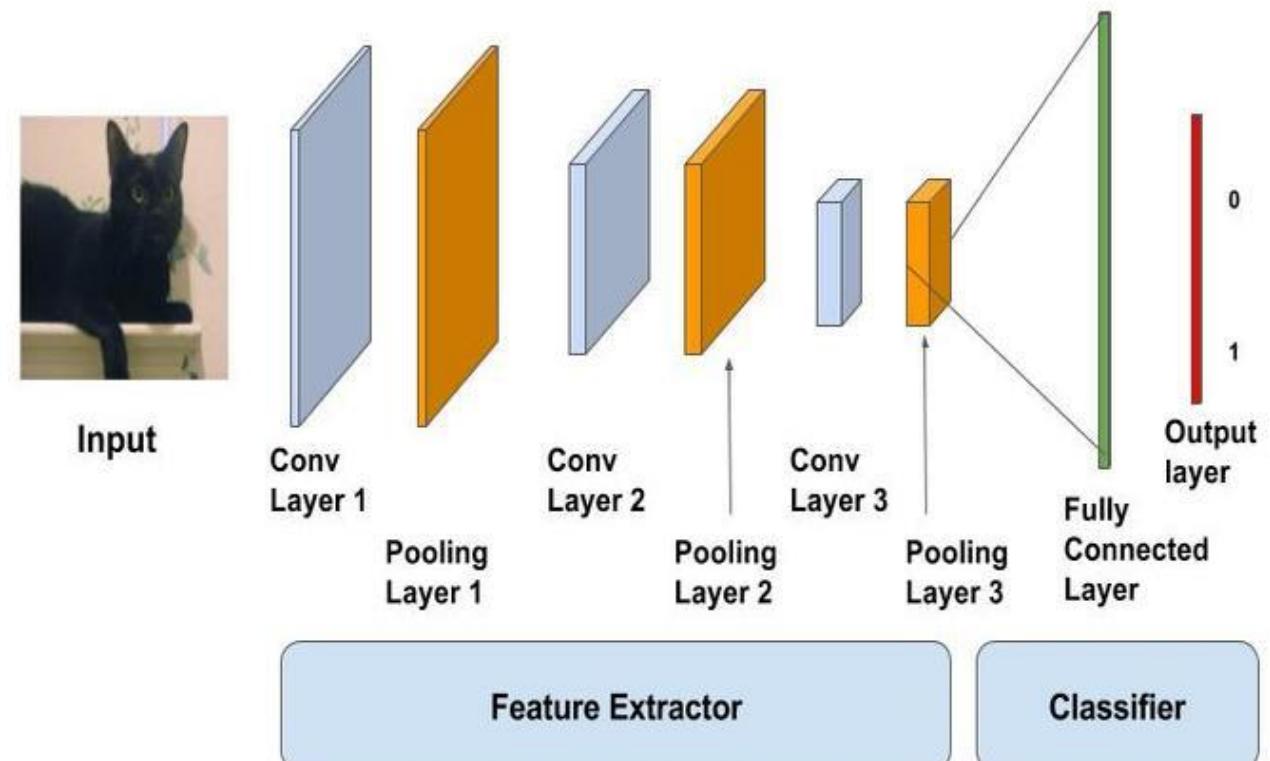


Image source: <https://www.learnopencv.com/>

- **Flatten the feature maps**
  - **First fully connected layer:**
    - Takes the inputs from the feature analysis and applies weights to predict the correct label
  - **Fully connected output layer:**
    - Gives the final probabilities for each label

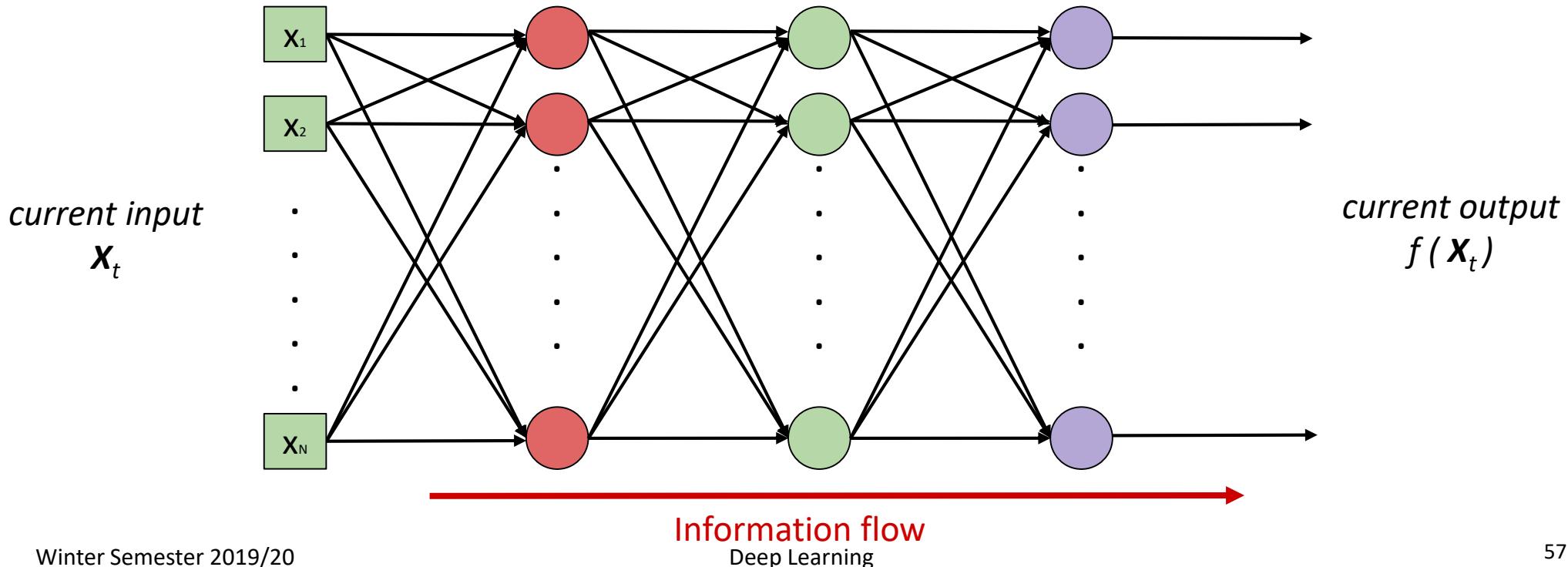


# Recurrent Neural Network (RNNs)

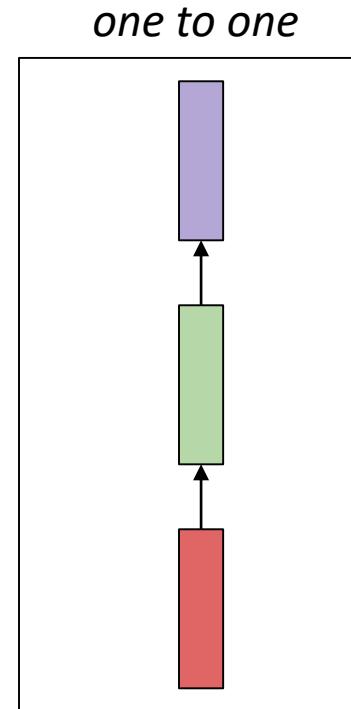
Perform predict over arbitrary-length sequences

## “Vanilla” Neural Networks

- Information only propagates forward through the network
  - Previous, and future, results do not affect current prediction



## “Vanilla” Neural Networks



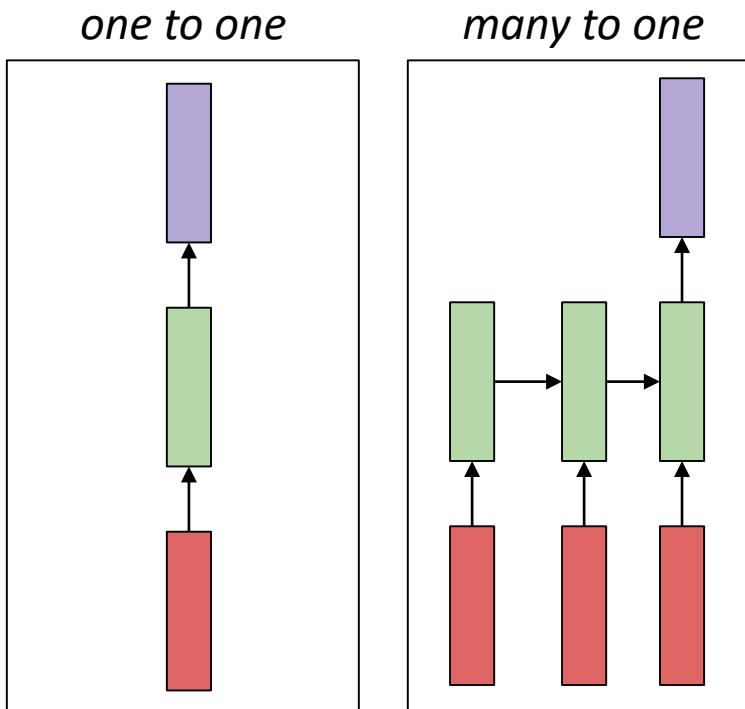
**Feedforward Neural Network:**

- *one input ↗ one output*
- *classifier or regressor*
- ***no context*** information

Neural Network needs **memory!**

Can the information/knowledge from **previous inputs** be included?

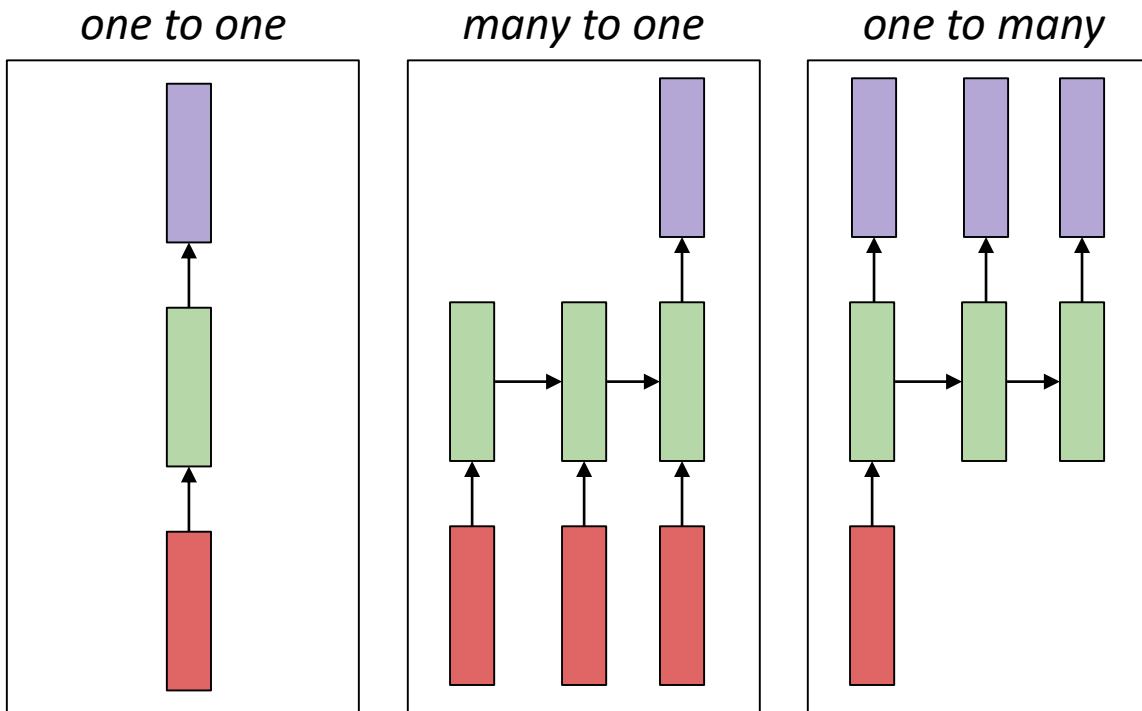
## Processing sequential inputs/outputs



Flexible to handle sequential inputs of varying length

- *speech recognition: speech waves ↗ one word*
- *sentiment classification: words ↗ positive/negative*
- *video classification: images ↗ one category*

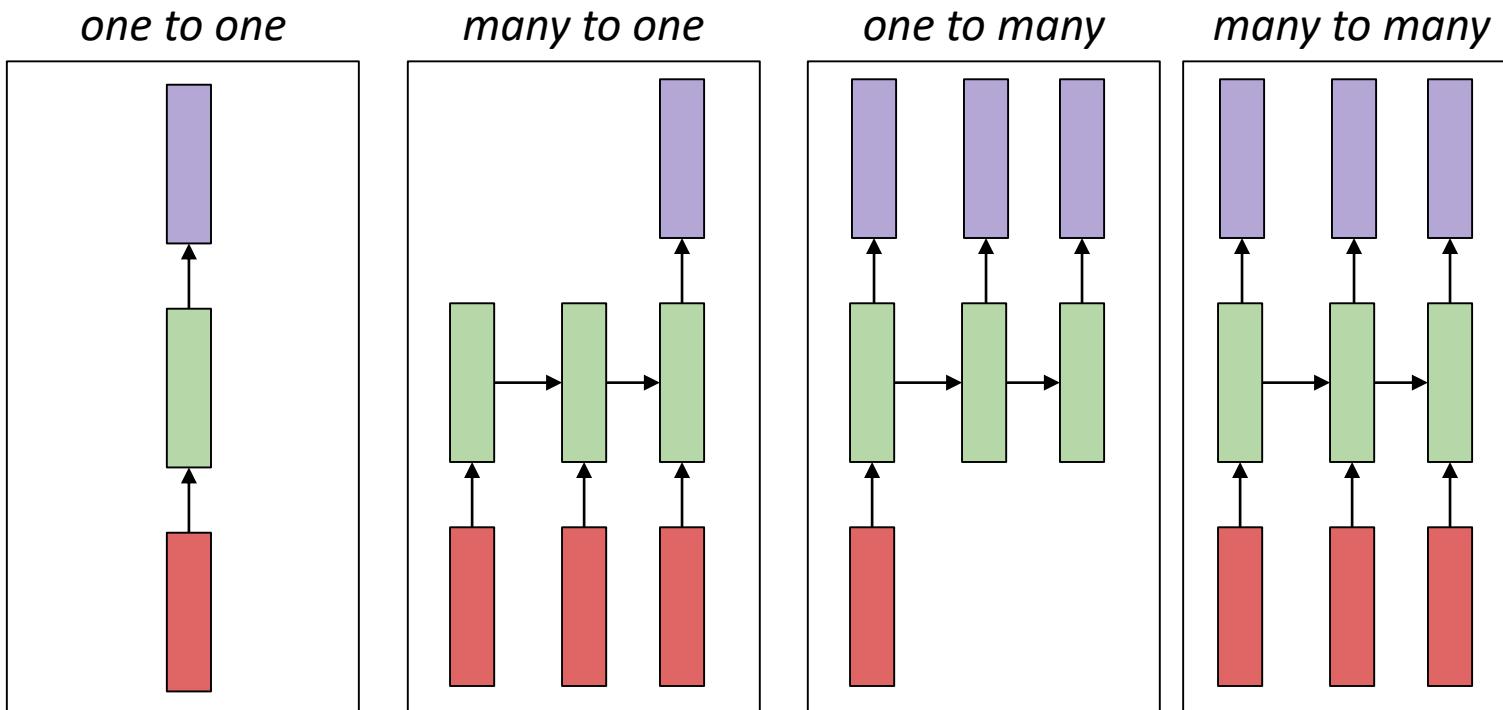
## Processing sequential inputs/outputs



To generate sequential outputs

- *music composition*
- *image captioning*
- *natural text generation*

## Processing sequential inputs/outputs

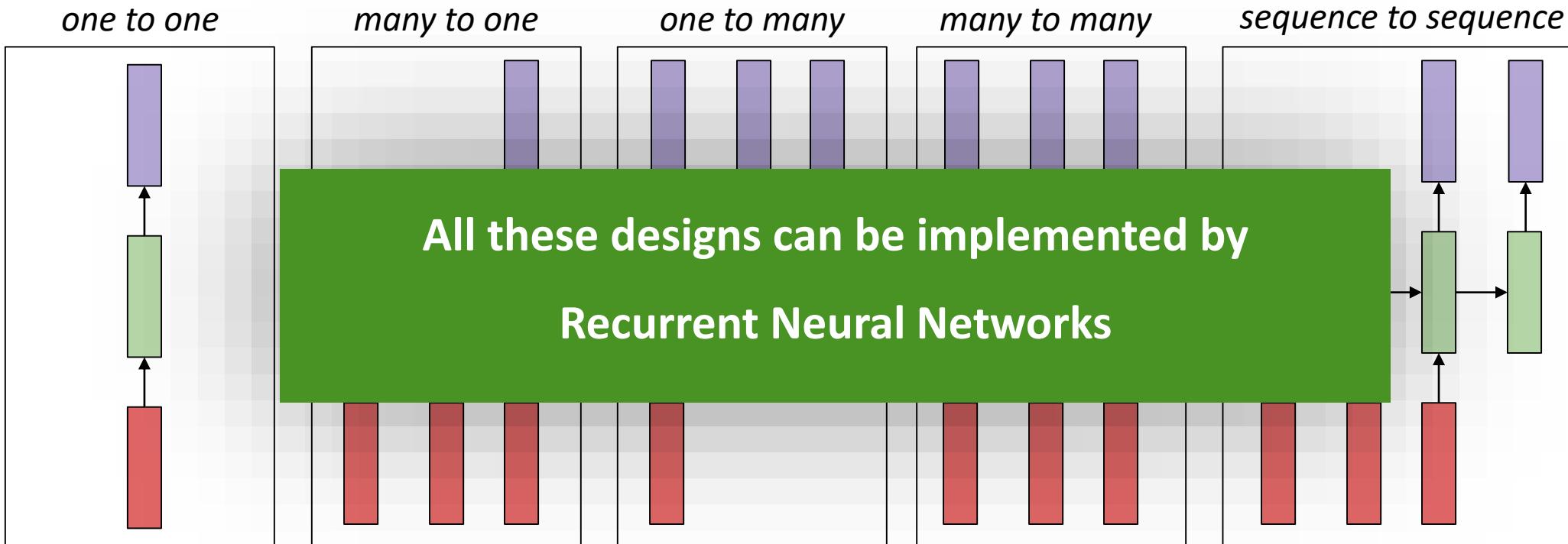


Inputs and outputs of same length

- *continuous emotion prediction*
- *speech enhancement*
- *video classification on frame level*

## Processing sequential inputs/outputs

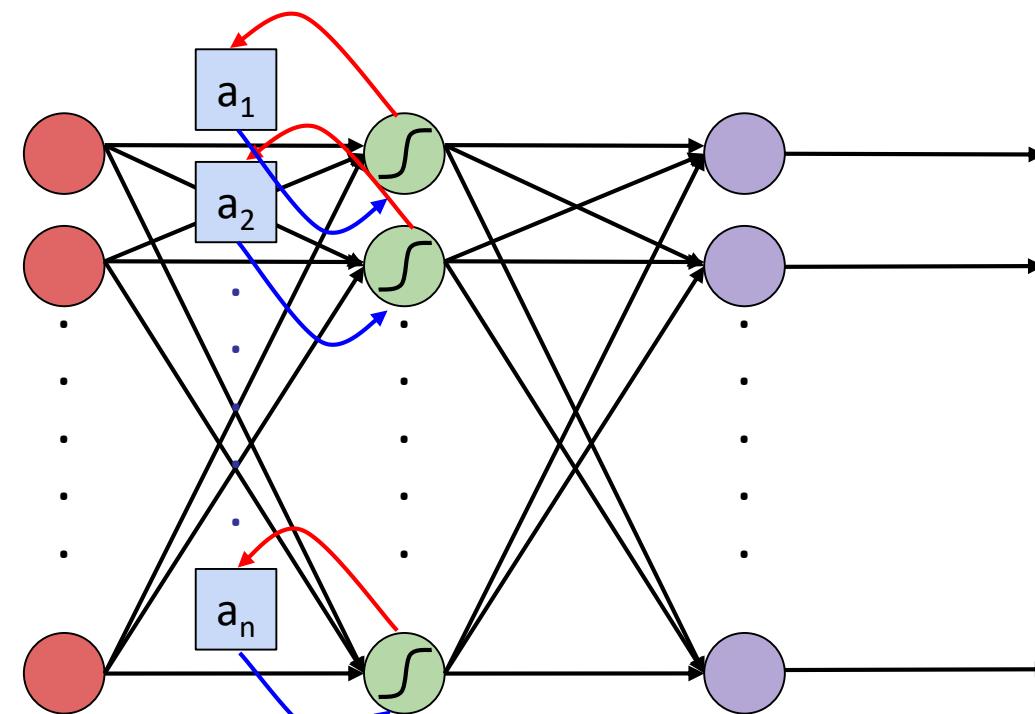
Different lengths:  
➤ *machine translation*  
➤ *question answering*



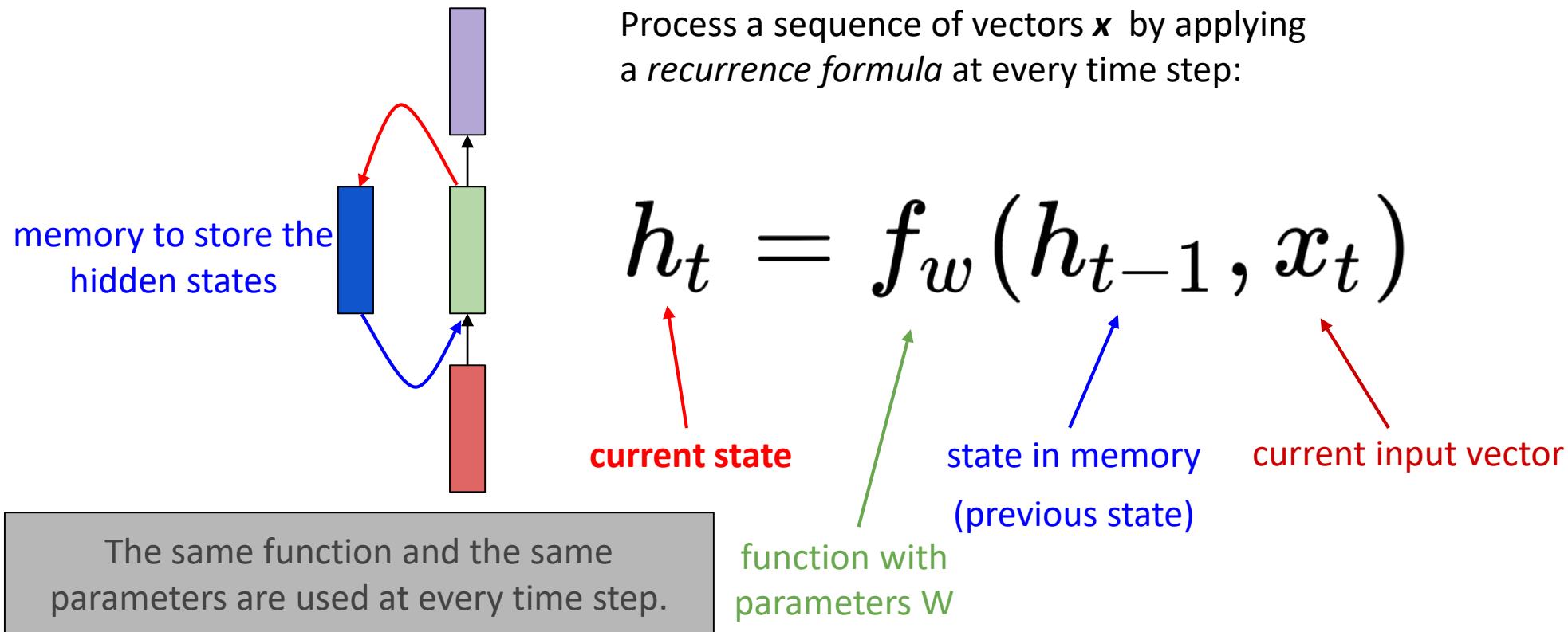
## Inclusion of feedback into the network structure

Output of hidden layer are stored in the memory

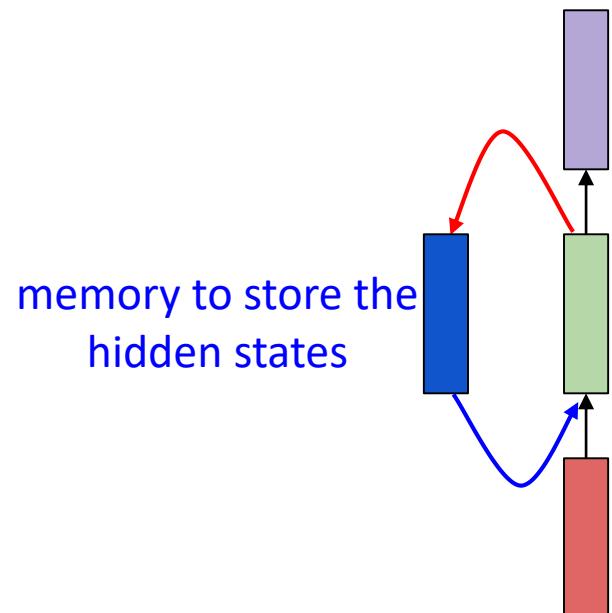
Values in the memory are considered as additional input in the next time step



## Inclusion of feedback into the network structure



## Inclusion of feedback into the network structure



Process a sequence of vectors  $x$  by applying a *recurrence formula* at every time step:

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

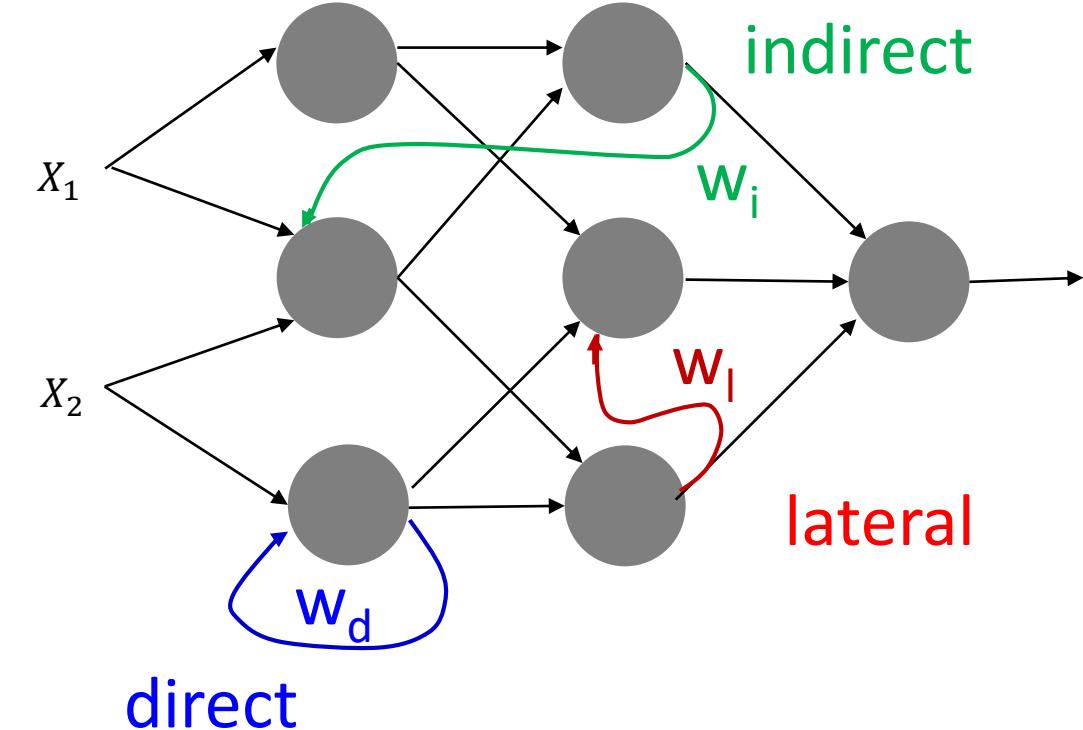


$$h_t = \text{tahn}(W_{hh} h_{t-1} + W_{xh} x_t)$$

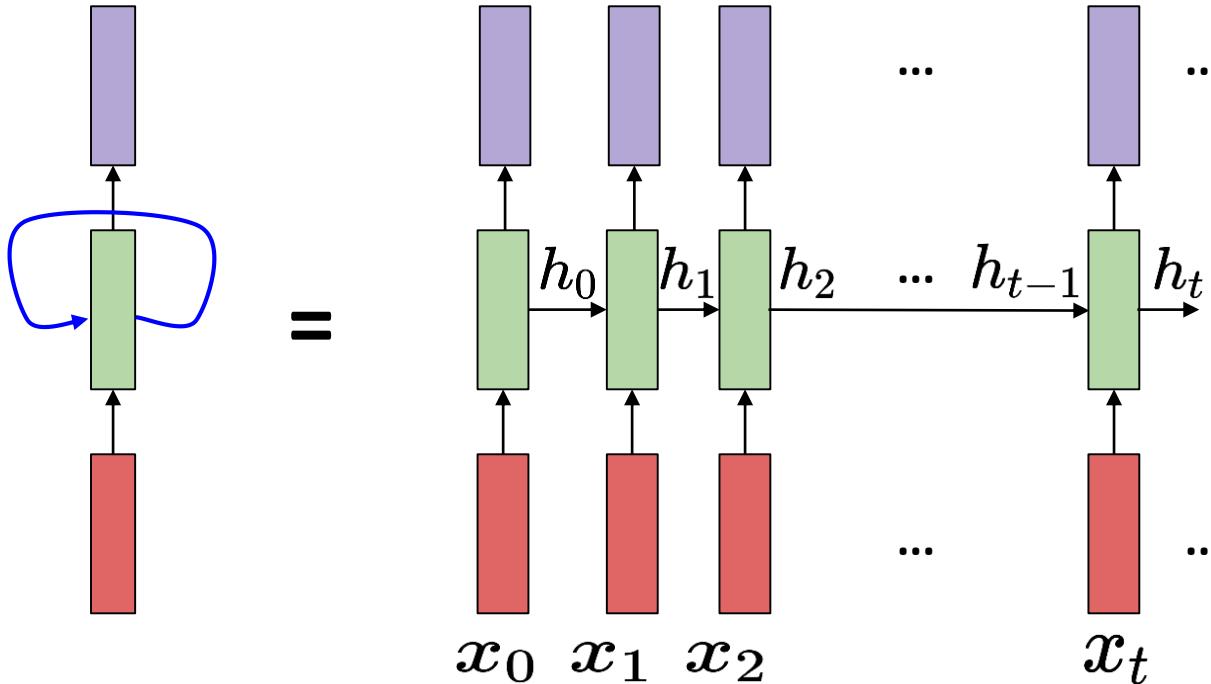
$$y_t = W_{hy} h_t$$

## Recurrent Connections

- **Direct feedback**
  - There exists a connection from a neuron  $j$  to itself with weight  $w_{j,j}$
- **Indirect feedback**
  - Activity is sent back to the previous layer in the neural network
- **Lateral feedback:**
  - Connect neurons within a layer

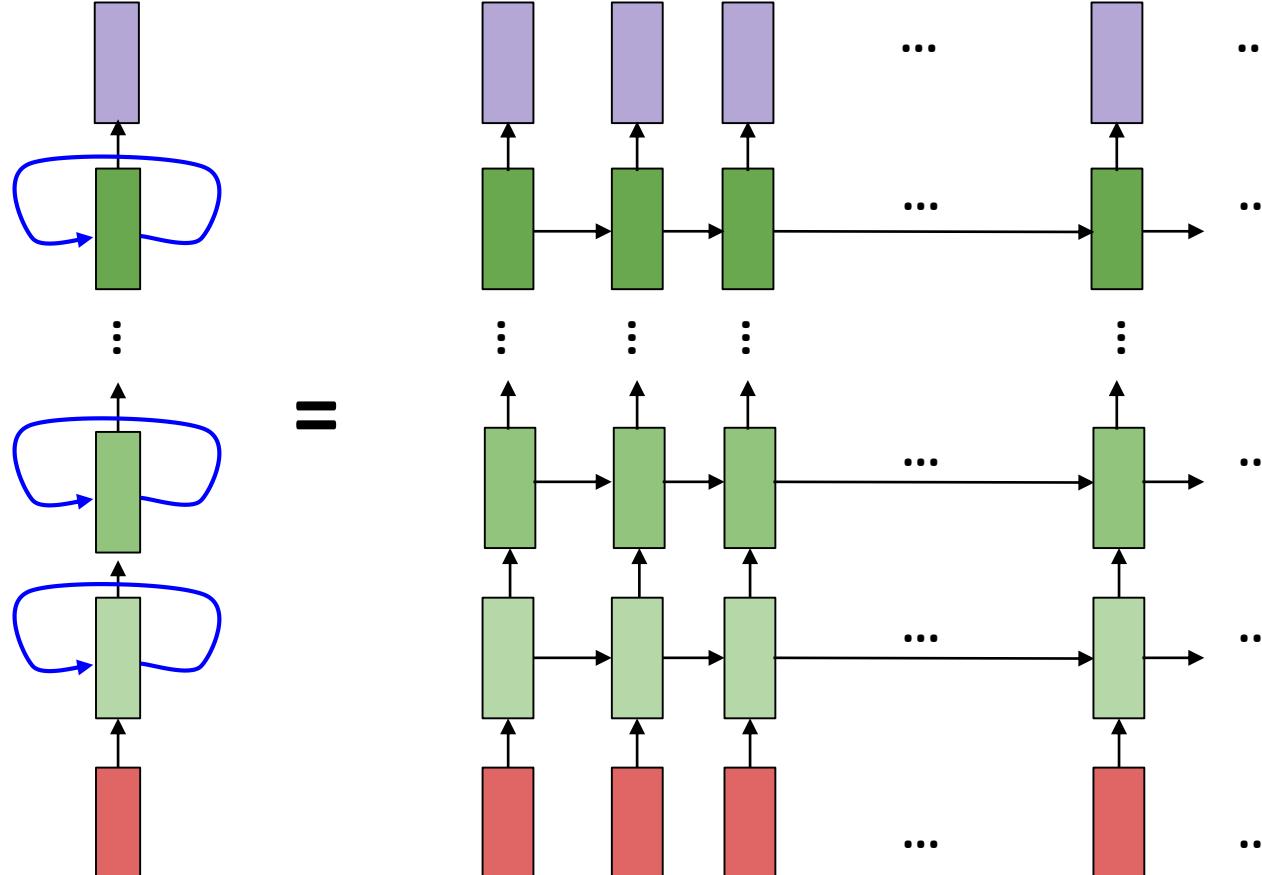


## Unrolled RNN



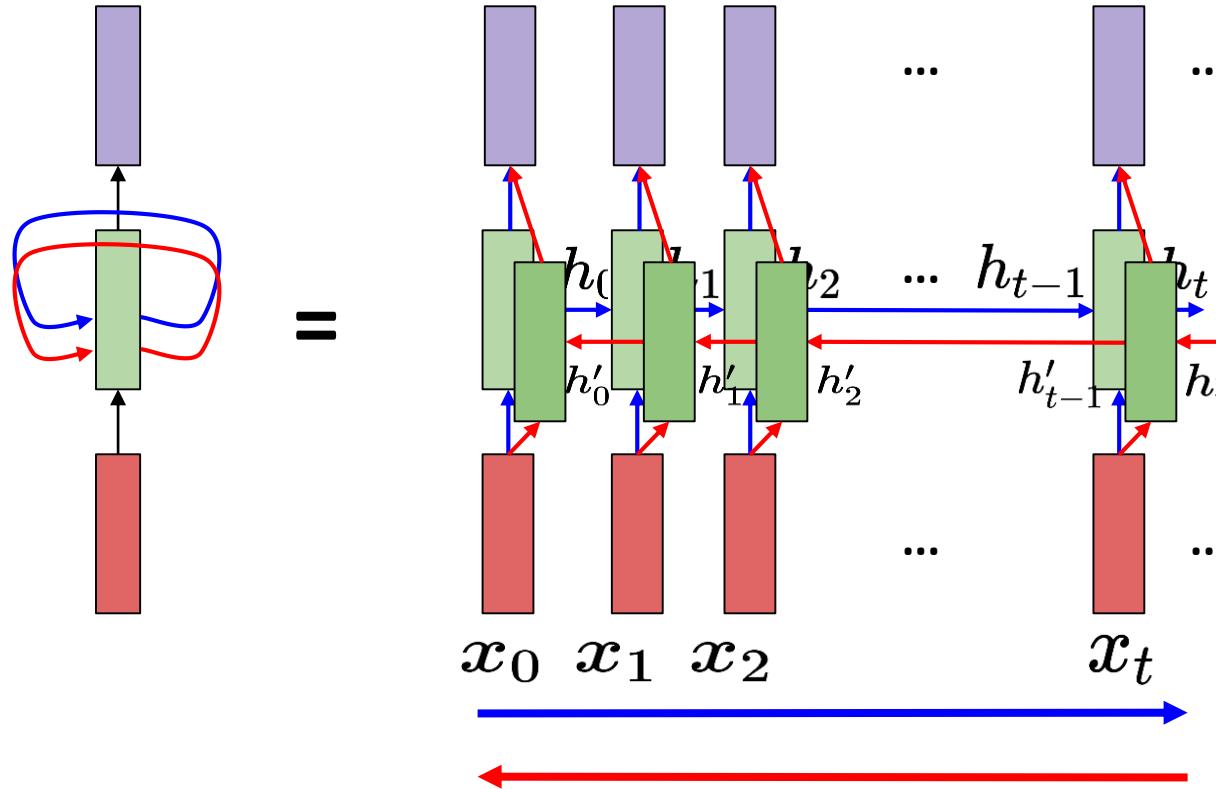
- Reuse the same weight matrix at every time step
- Makes the network easier to train

## Deep RNN



- Not easy to train
- Regularisation helps
  - *Batch normalisation*
  - *Dropout*
  - *These topics will be covered in future lecture*

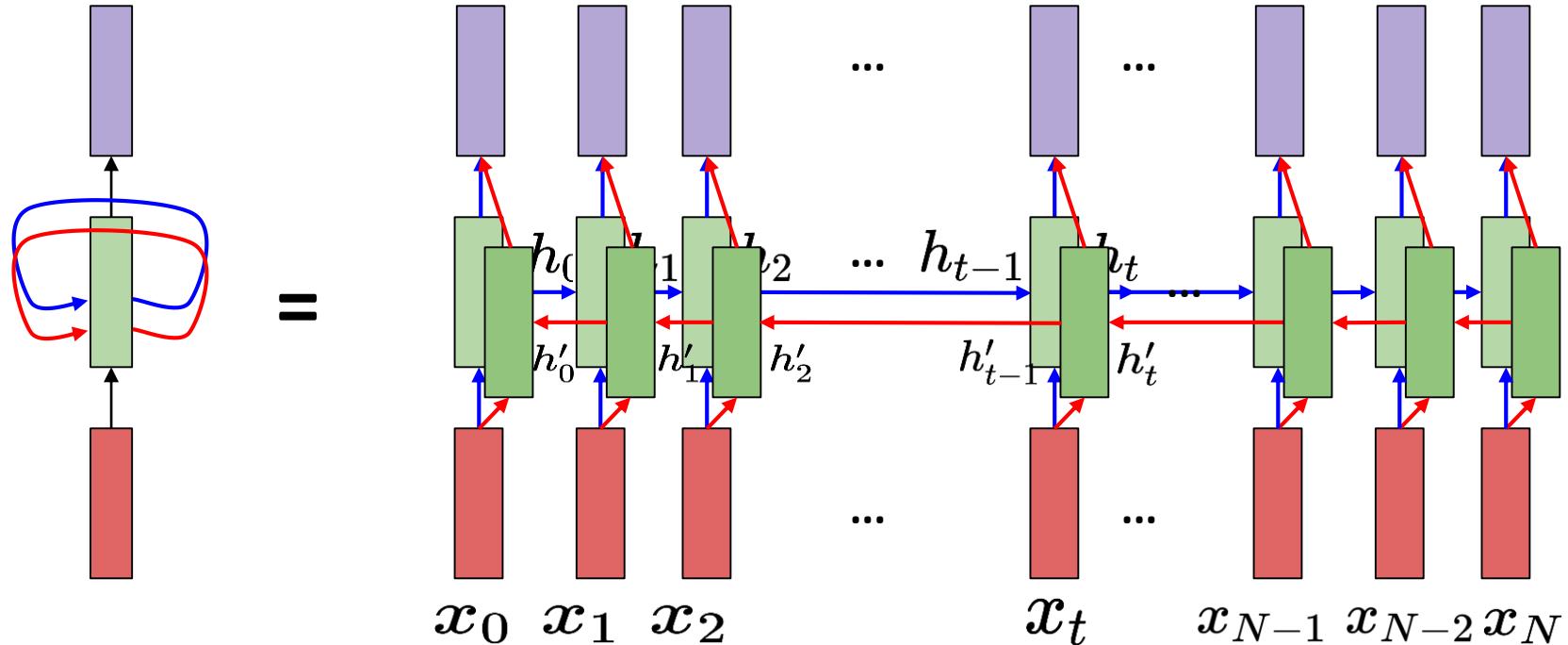
## Bidirectional RNN



Two hidden layers:

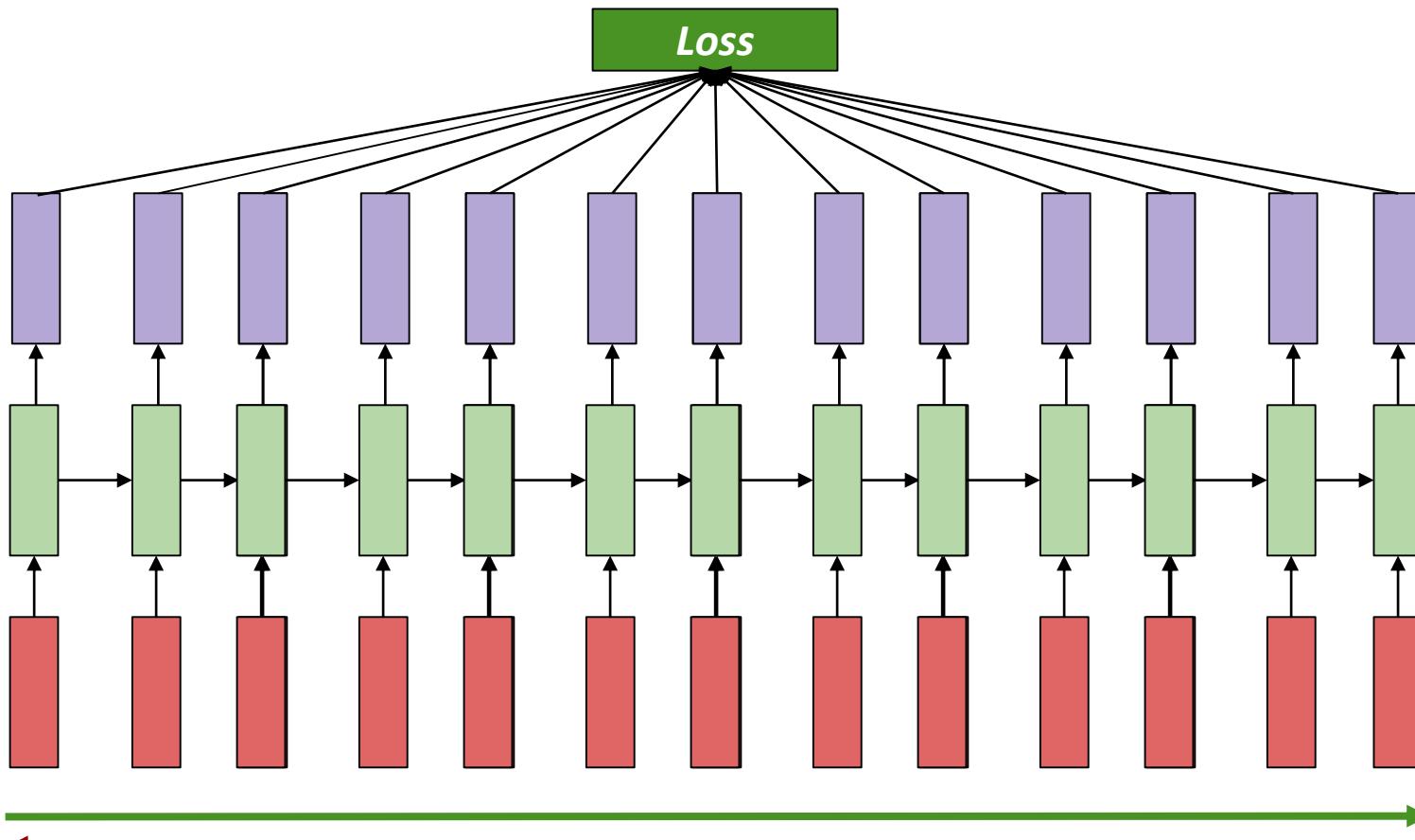
- Scan the input sequences separately
- One *forward* and one *backward*
- Connect to the same output layer
- Learn both *past* and *future* context

## Bidirectional RNN



To determine  $y_t$  you have to  
see the whole sequence

## Backpropagation through time (BPTT)



**Forward**

Run through entire sequence to compute the **loss**

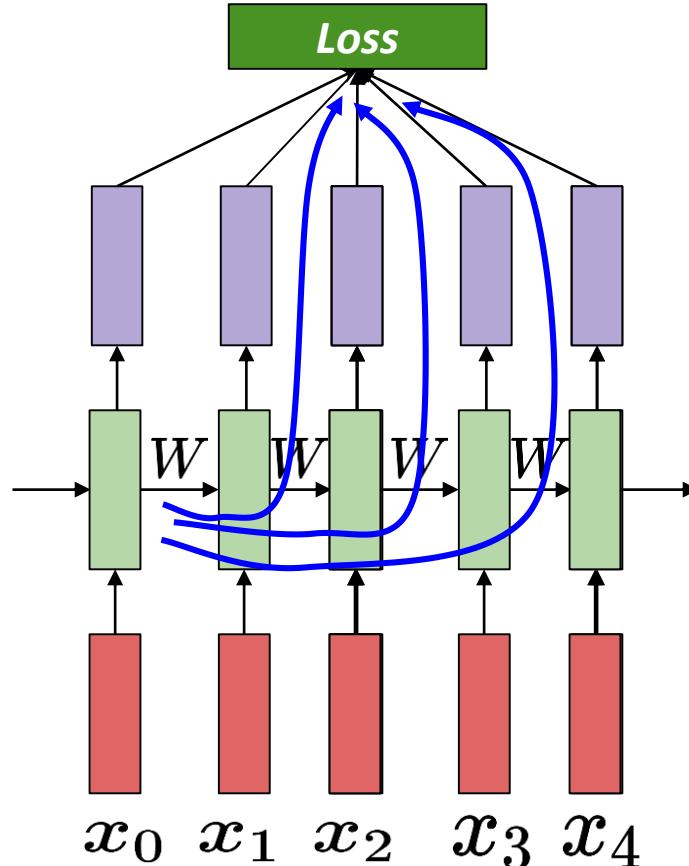


**Backward**

Run through entire sequence to compute the **gradient** and update the weight matrix:

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

## Gradient flow in simple RNNs



$$w \leftarrow w - \eta \partial L / \partial w$$

Issue: **W** occurs each timestep

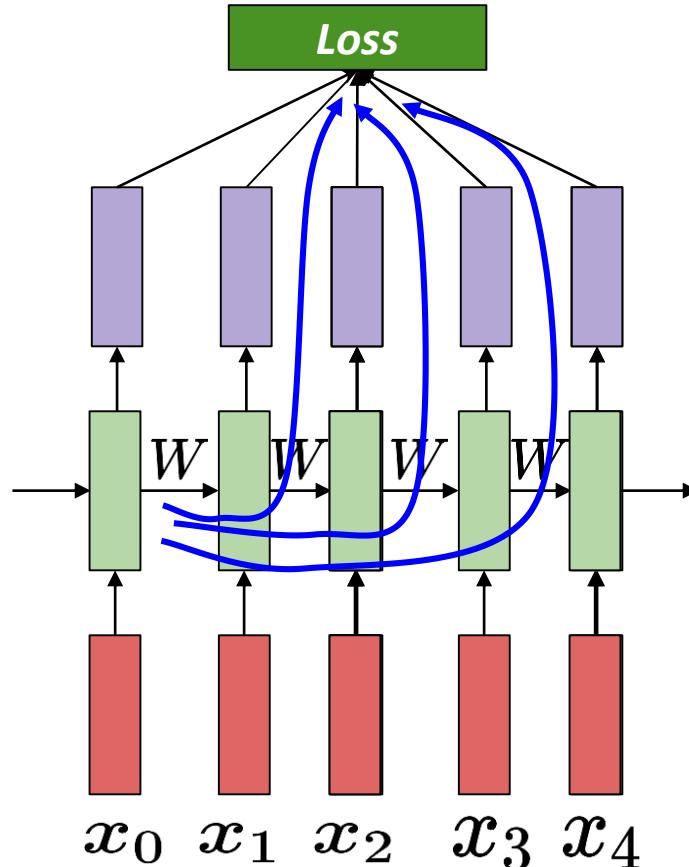
Every path from **W** to **Loss** is one dependency

All paths from **W** to **Loss** need to be involved

We need to sum up the gradients at each time step

$$\frac{\partial L}{\partial w} = \sum_{j=0}^T \frac{\partial L_j}{\partial w}$$

## Gradient flow in simple RNNs



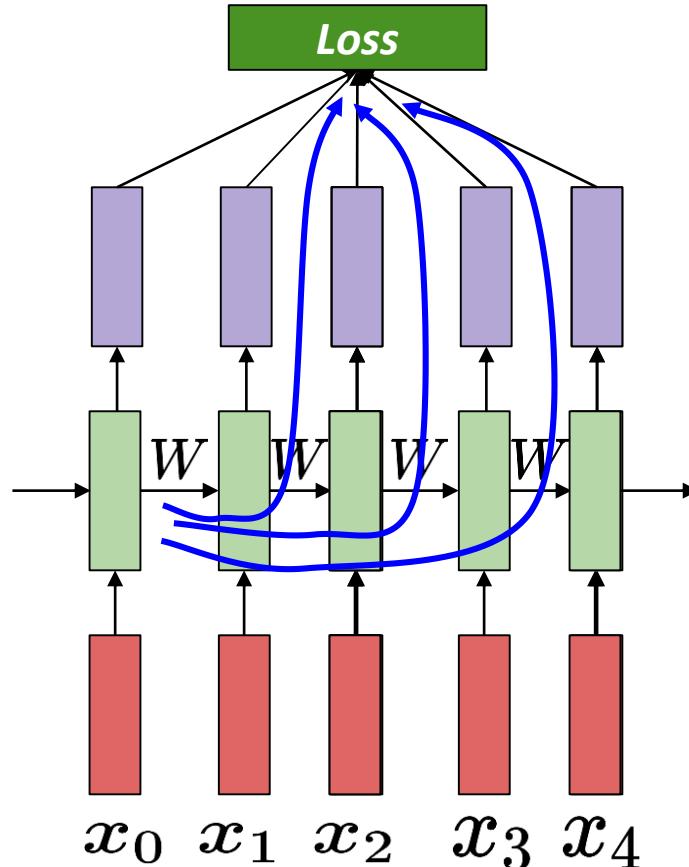
We need to **sum up the gradients** at each time step

$$\frac{\partial L}{\partial w} = \sum_{j=0}^T \frac{\partial L_j}{\partial w}$$

To calculate these gradients we use **multiple iterations** of chain rule of differentiation

$$\frac{\partial L}{\partial w} = \sum_{j=0}^T \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{t=k+1}^j \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_k}{\partial w}$$

## Gradient flow in simple RNNs



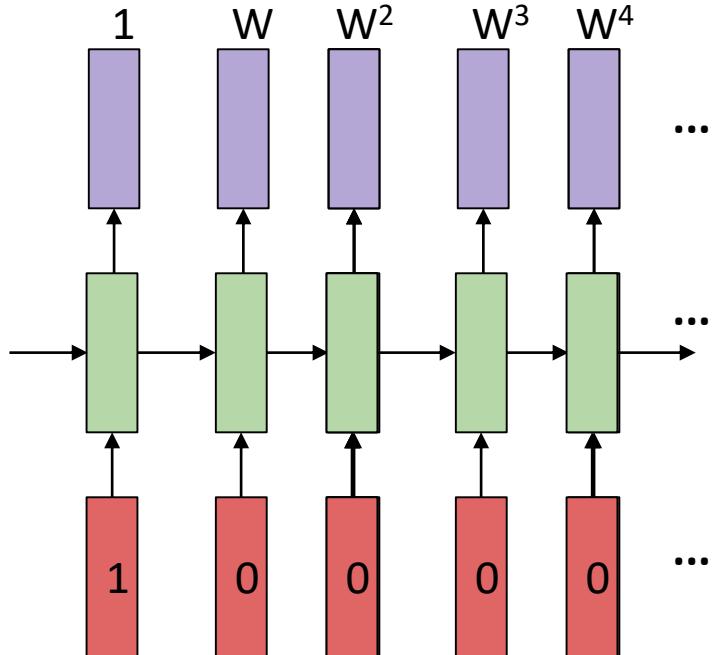
$$\frac{\partial L}{\partial w} = \sum_{j=0}^T \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{t=k+1}^j \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_k}{\partial w}$$

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

computing gradient of  $h_{t-1}$  involves  $W_{hh}$

Repeated matrix multiplications leads to **vanishing** and **exploding** gradients.

## Gradient flow in simple RNNs



*Toy Example:*

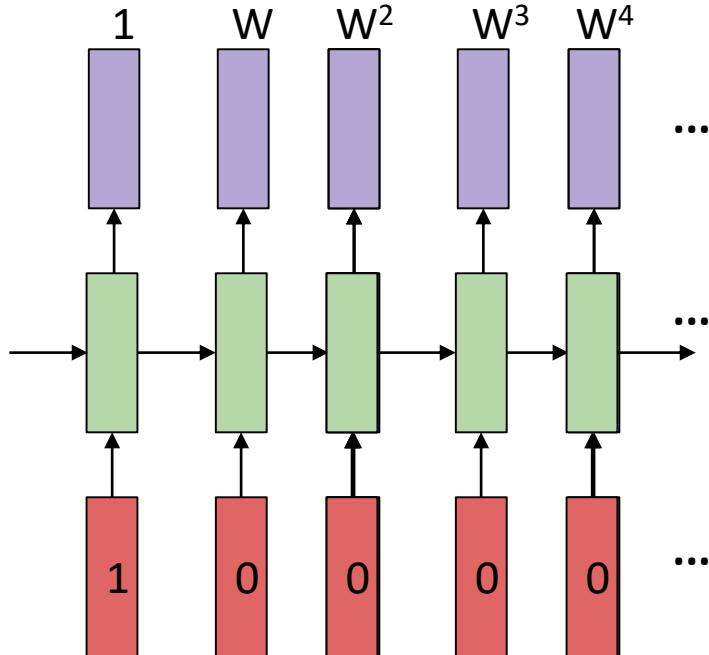
- 1-d input,  $x_0=1$ , and  $x_t=0$  ( $t>0$ )
- Sequence length is **1000**
- $W_{xh}=W_{hy}=1$
- Linear activation function

$$h_t = W_{hh} h_{t-1} + x_t$$

$$y_t = h_t$$

$$y_{1000} = W_{hh}^{999}$$

## Gradient flow in simple RNNs



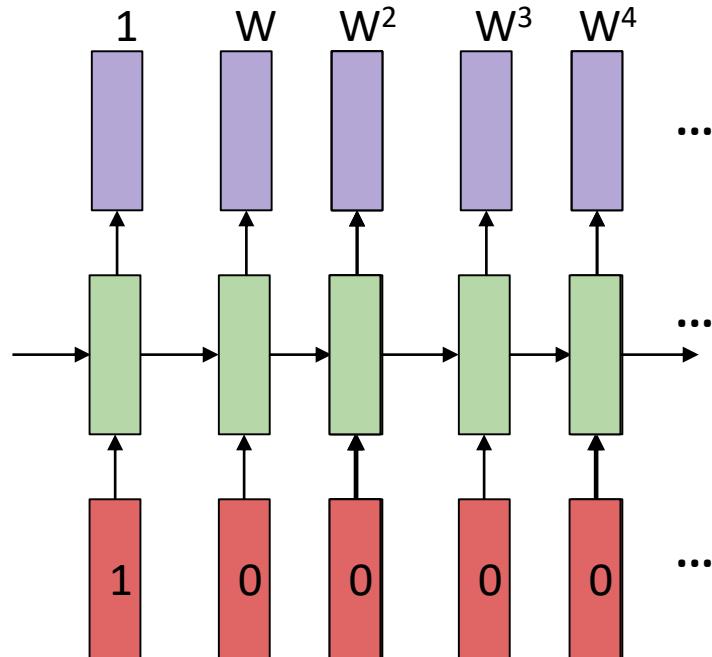
$\partial L / \partial w$  is too large

→ Exploding gradients

- Regulate via gradient clipping
  - Scale gradient if value gets to big

$$\text{if } \|g\| > \text{threshold}, g \leftarrow \frac{\text{threshold} \times g}{\|g\|}$$

## Gradient flow in simple RNNs



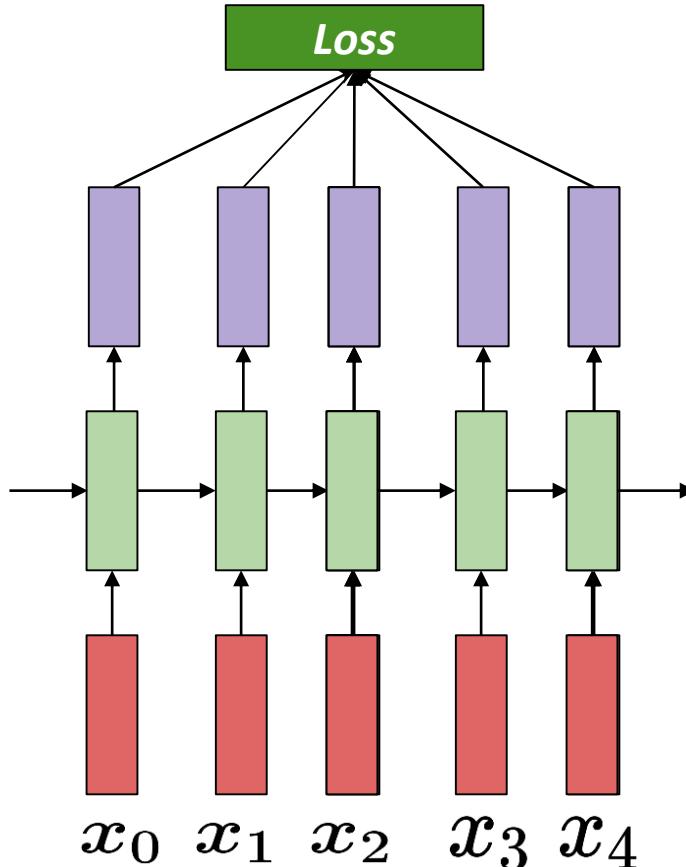
$$\begin{aligned} w_{hh} = 0.99 &\rightarrow y_{1000} \approx 0 \\ w_{hh} = 0.01 &\rightarrow y_{1000} \approx 0 \end{aligned}$$

$\partial L / \partial w$  is too small

→ **Vanishing gradients**

- Need to change RNN architecture
- **Gated RNNs**

## Identity-RNN



- Network weights are initialised to the identity matrix
- Activation functions are all set to ReLU
- This encourages the network computations to stay close to the identity function
- **Gradient does not decay** as derivatives stay either 0 or 1
- Error is propagated all the way back

$$\frac{\partial L}{\partial w} = \sum_{j=0}^T \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{t=k+1}^j \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_k}{\partial w}$$

$$h_t = h_{t-1} + f(x_t) \rightarrow \left( \frac{\partial h_t}{\partial h_{t-1}} \right) = 1$$

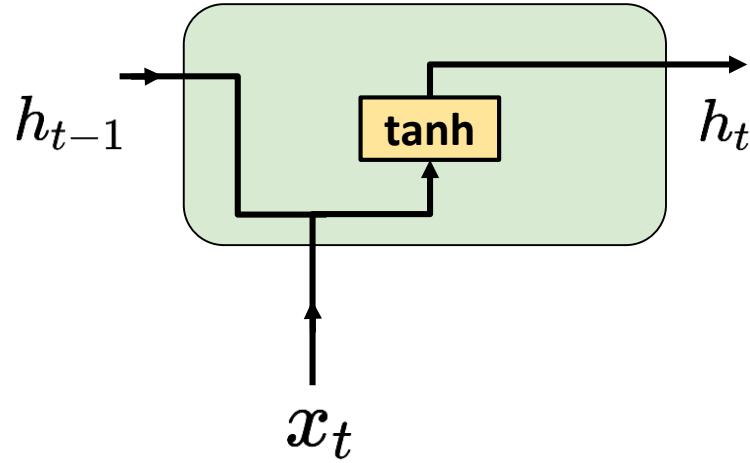
## Highly effective sequence models

- Gated RNNs are based on the idea of creating paths that have derivatives that neither vanish nor explode
- Gated RNNs have connection weights that may change at each time step
- Gated RNNs also allow a network to forget an old state
- Instead of manually deciding when to clear the state, the network to learn to decide when to do it.

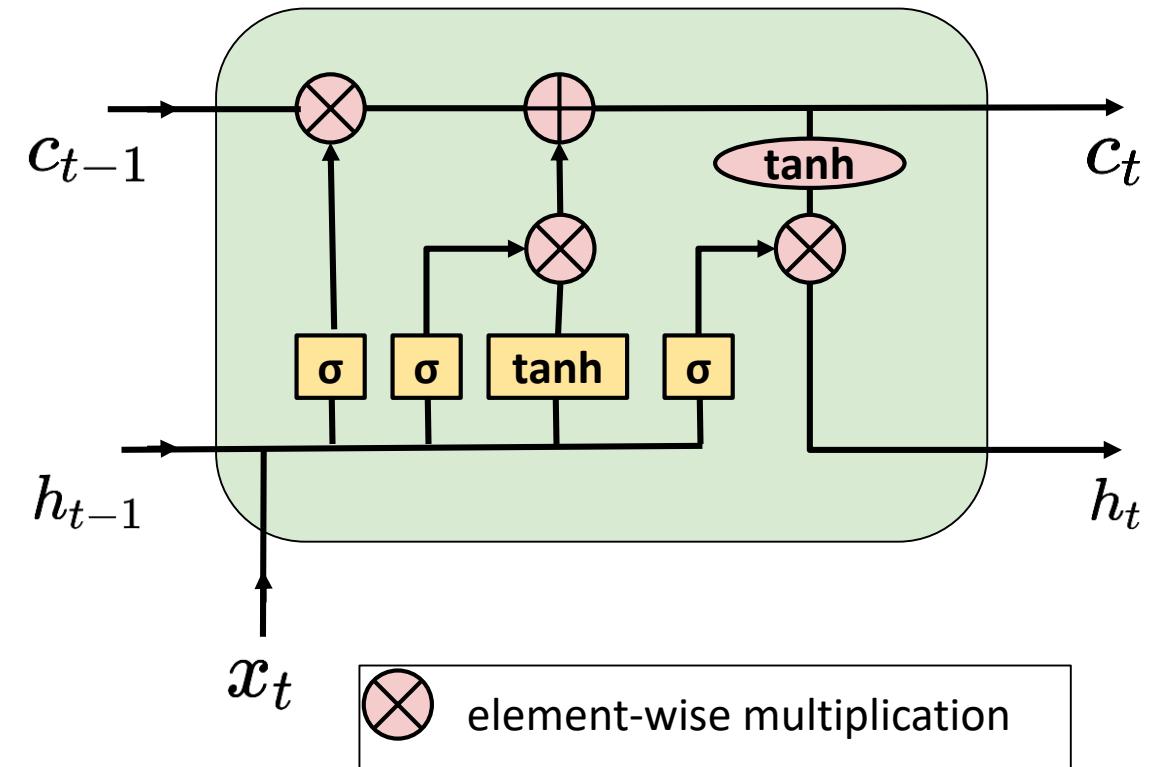
- **The LSTM Recurrent Unit**

- The internal structure of an LSTM unit allows it to keep or forget information over time
- **Cell State:** the memory of the network, carries relevant information throughout the processing of the sequence
- Information gets added or removed to the cell state via gates
  - **Forget gate:** what is relevant to keep from prior step
  - **Input gate:** what information is relevant to add from current step
  - **Output gate:** what the next hidden state should be
- The gates learn to determine what is relevant during training

## Simple RNN unit

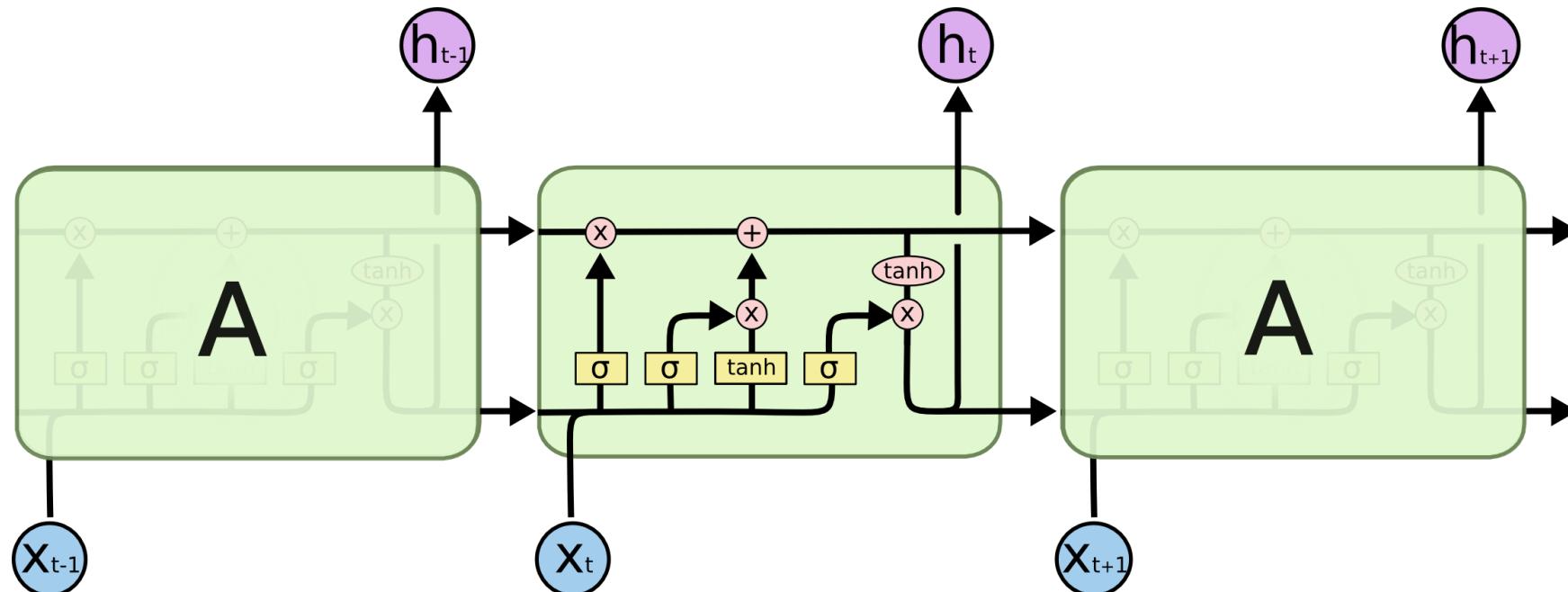


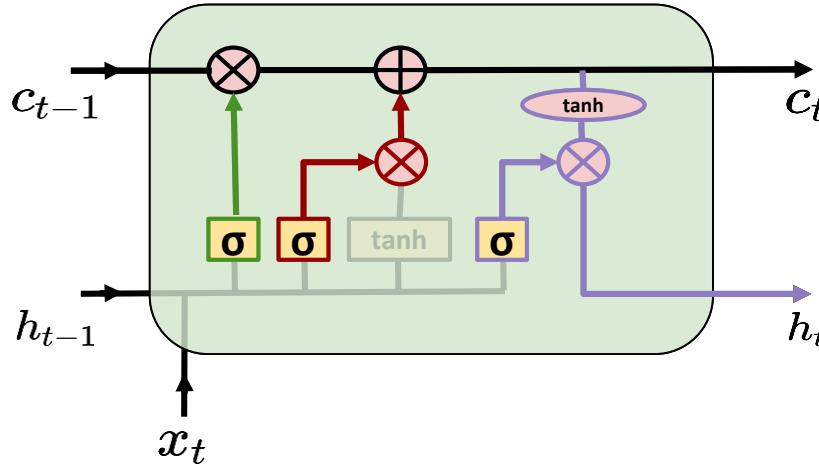
## Gated LSTM unit



- **Unrolled LSTM network**

- There are four interacting networks: **cell state**, **input gate**, **forget gate**, **output gate**





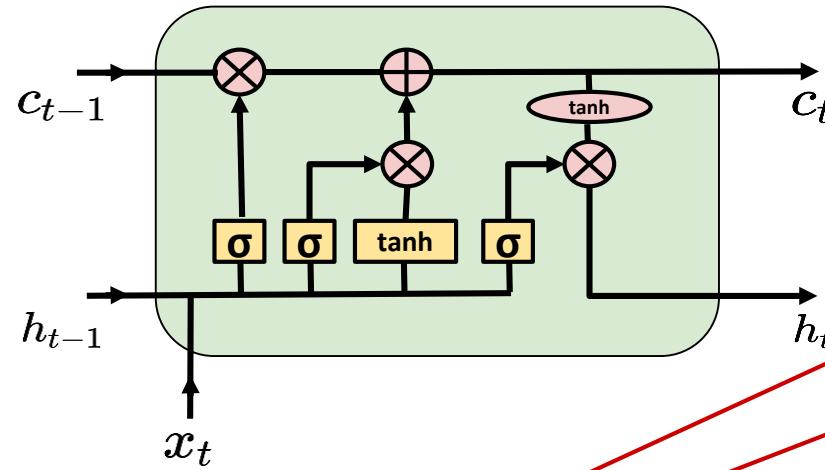
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

## Cell state $c$ :

- Core of LSTM unit
- Information is removed or added to the cell
- Controlled by three gates:
  - *Input gate*: how much to write to cell
  - *Forget gate*: how much to erase/keep in the cell
  - *Output gate*: how much to output from the cell



$$g_t = \tanh(W_g \cdot [h_{t-1}, x_t] + b_g)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

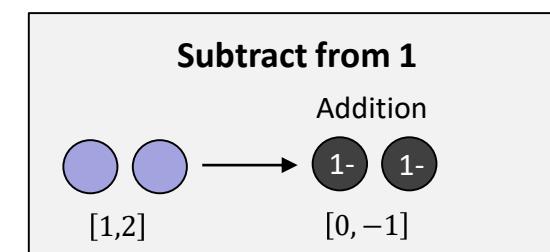
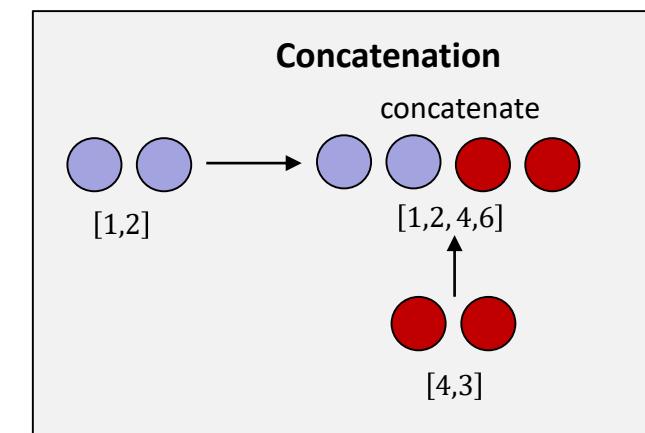
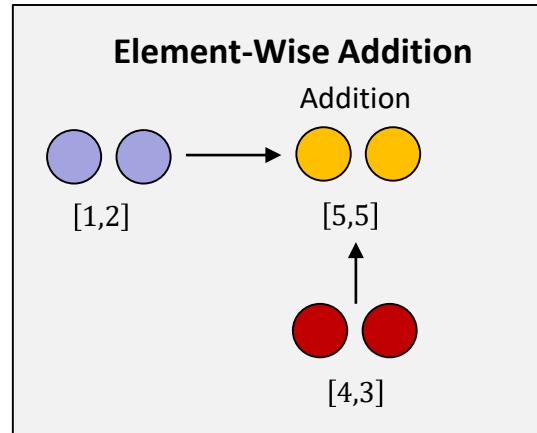
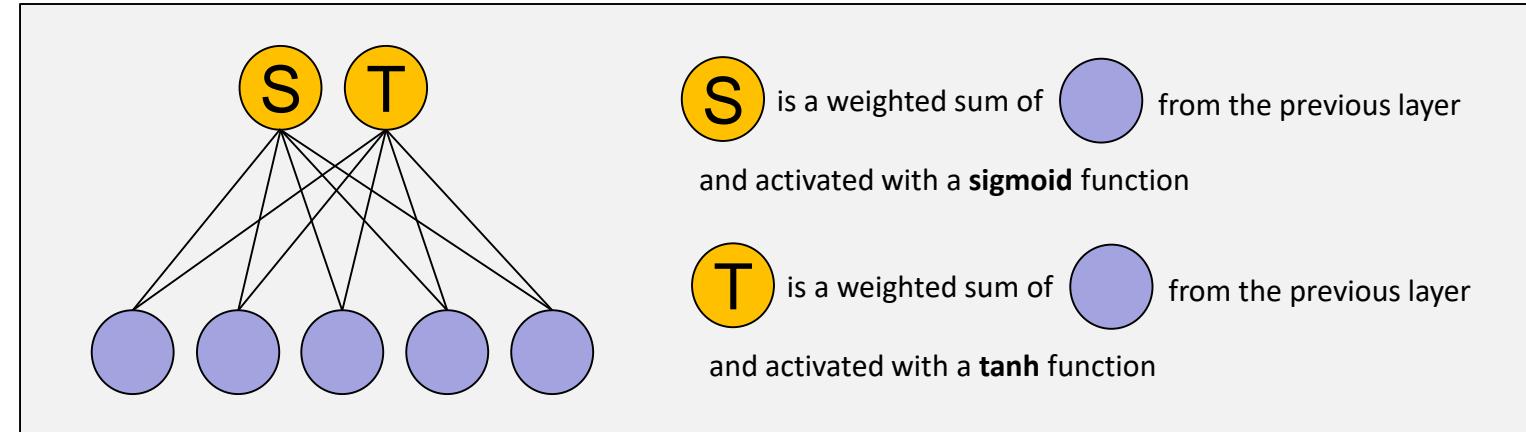
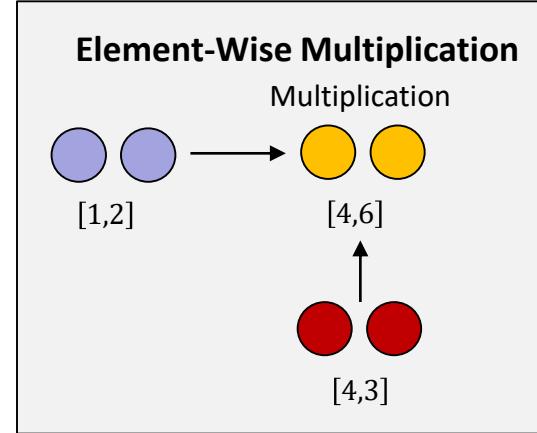
$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$

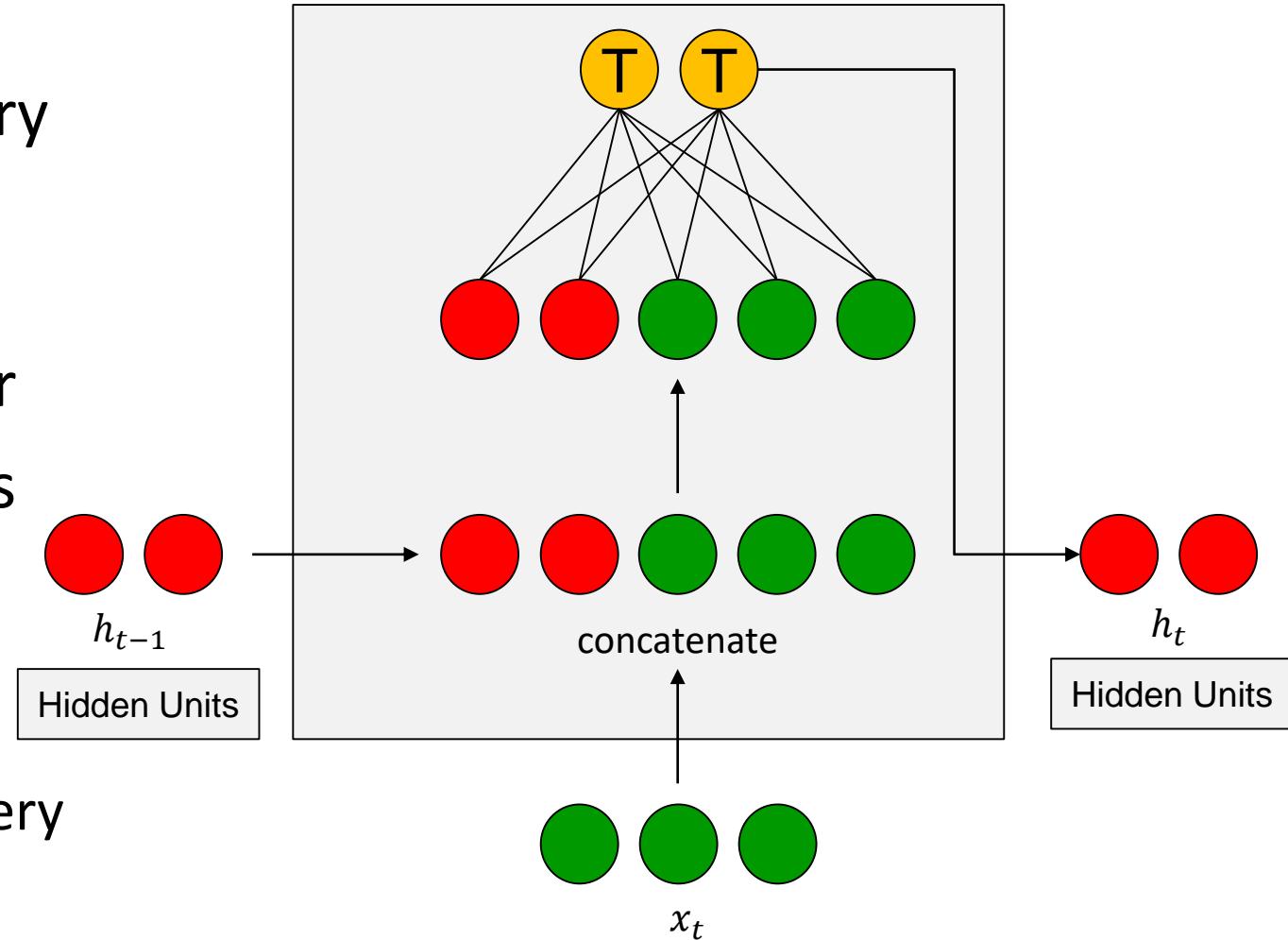
More weights to learn

As in the *simple RNN*  
We reuse the same weight  
matrix at every time step

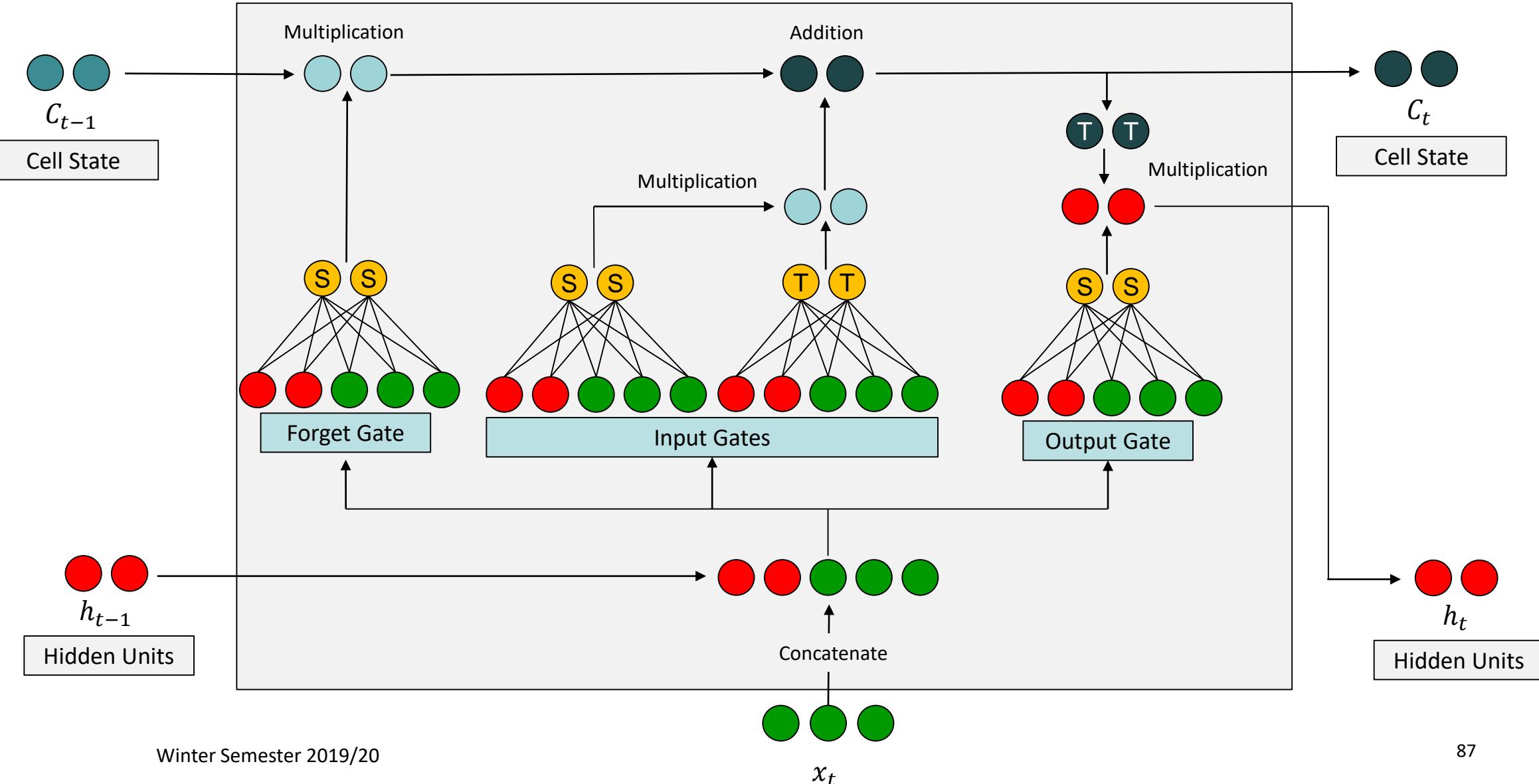


- **Vanilla RNNs**

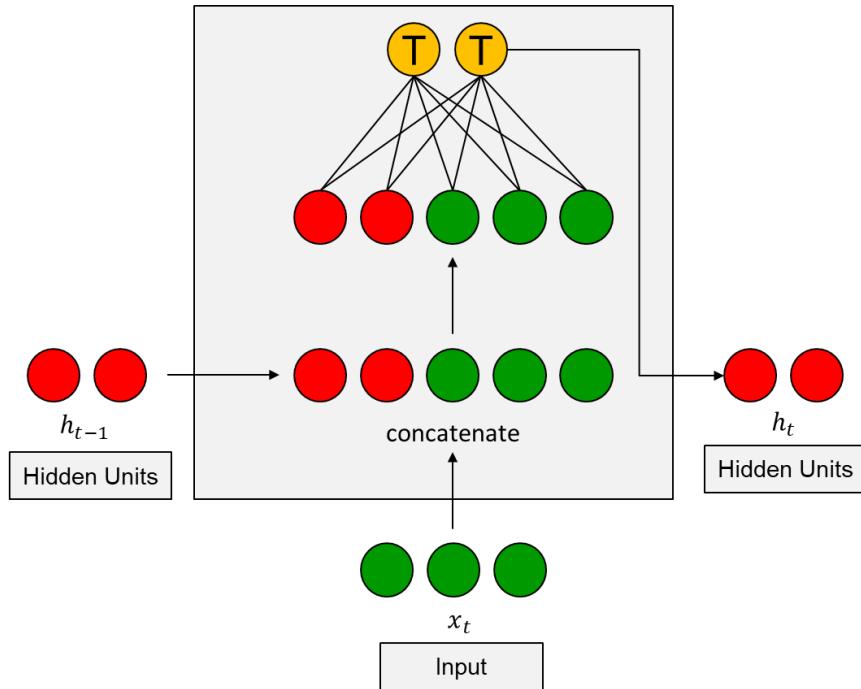
- Suffer from short-term memory
- If a sequence is long enough they cannot carry information from earlier time steps to later
- Vanishing/Exploding Gradients
  - If a gradient value becomes extremely small, it doesn't contribute too much learning
  - Large error gradients result in very large updates during training.



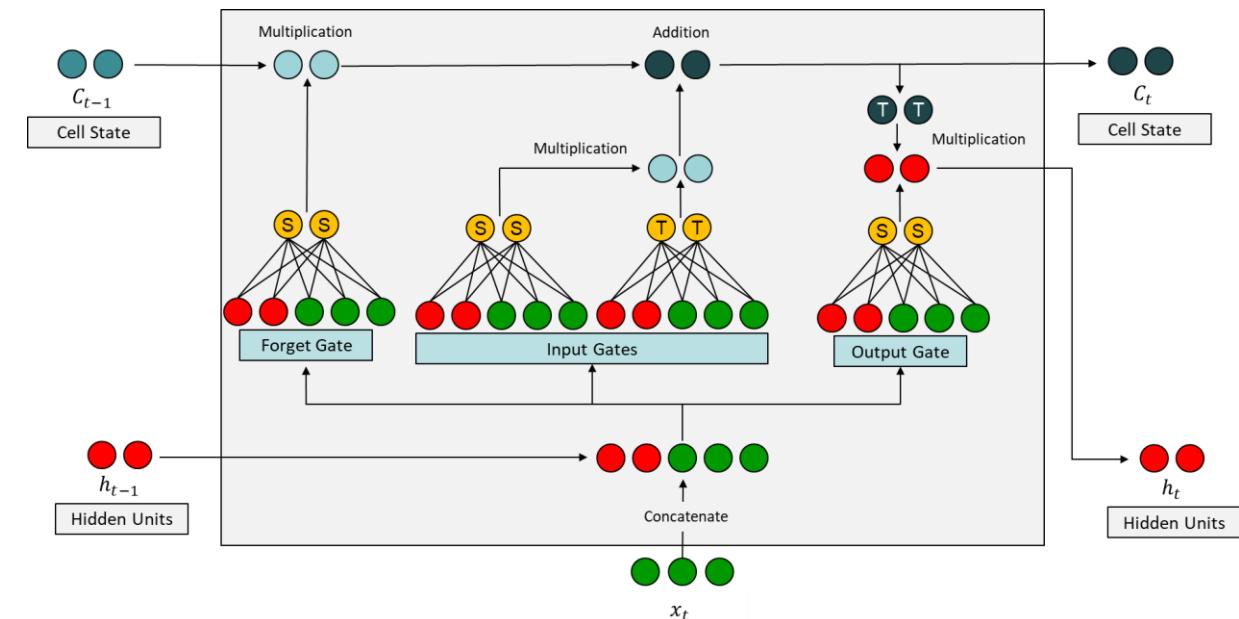
# Long Short-term Memory (LSTM)

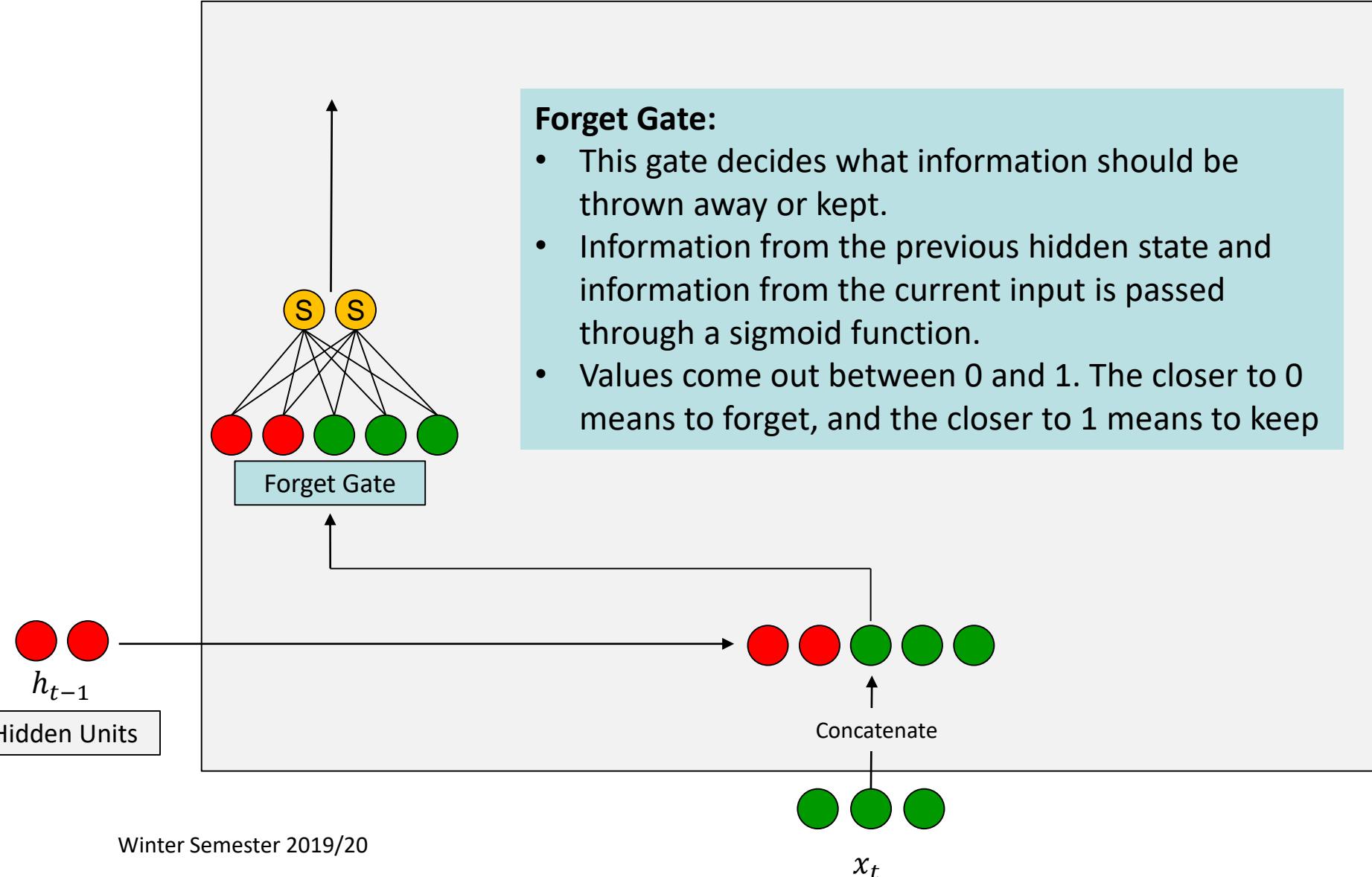


## Simple RNN unit



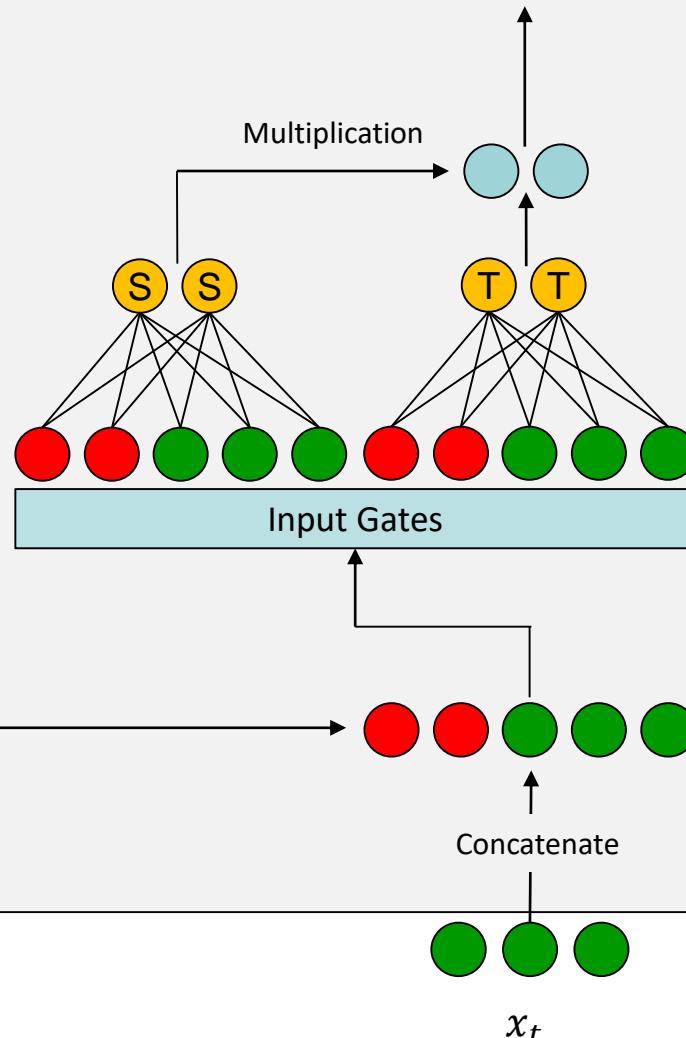
## Gated LSTM unit





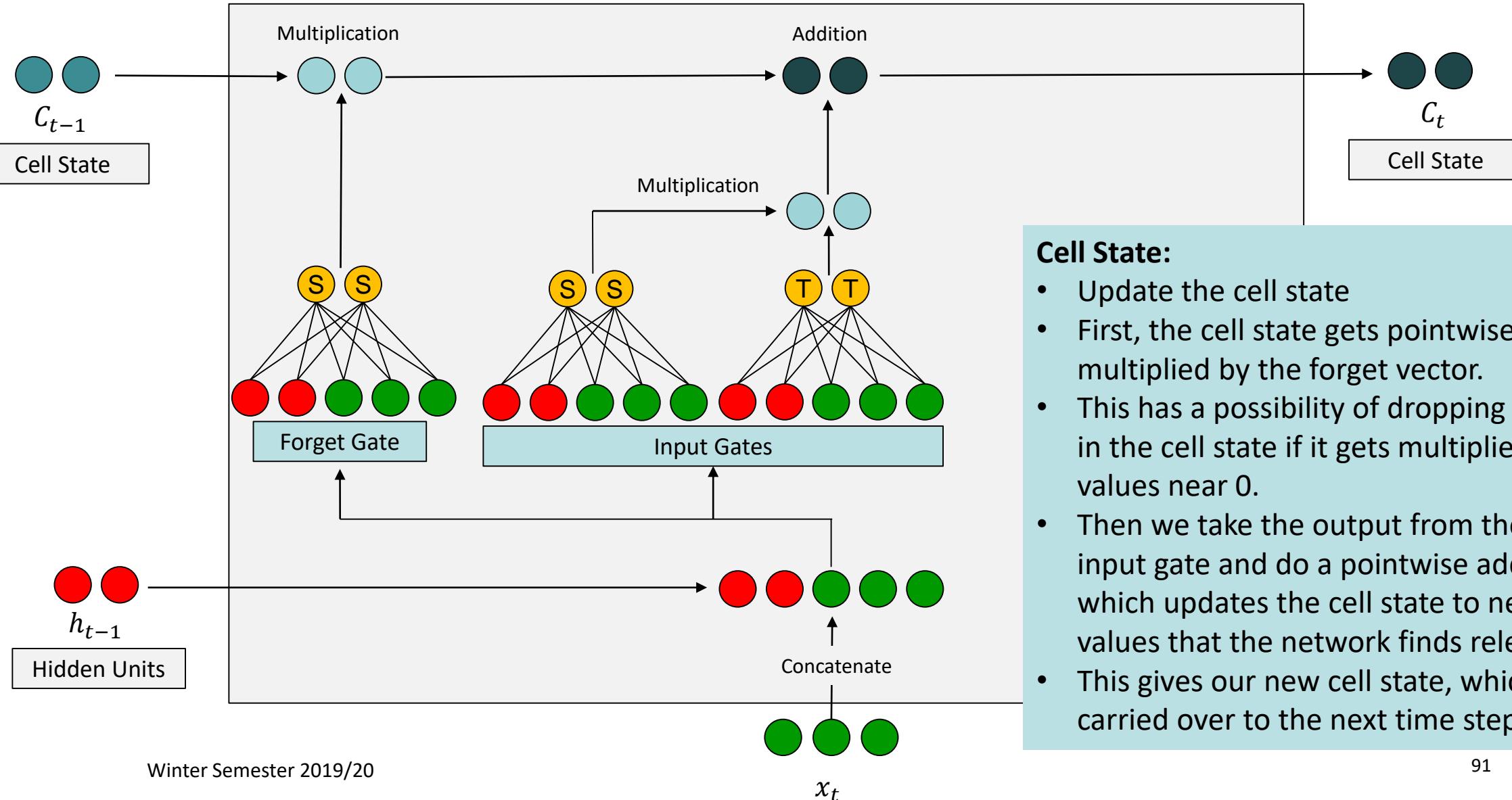
## Input Gate:

- Updates the cell state
- First, we pass the previous hidden state and current input into a sigmoid function.
- This decides which values will be updated by transforming the values to be between 0 and 1
- 0 means not important, and 1 means important



## Input Gate:

- You also pass the hidden state and current input into the tanh function to squish values between -1 and 1
- This helps to regulate the network.
- Then you multiply the tanh output with the sigmoid output.
- The sigmoid output will decide which information is important to keep from the tanh output



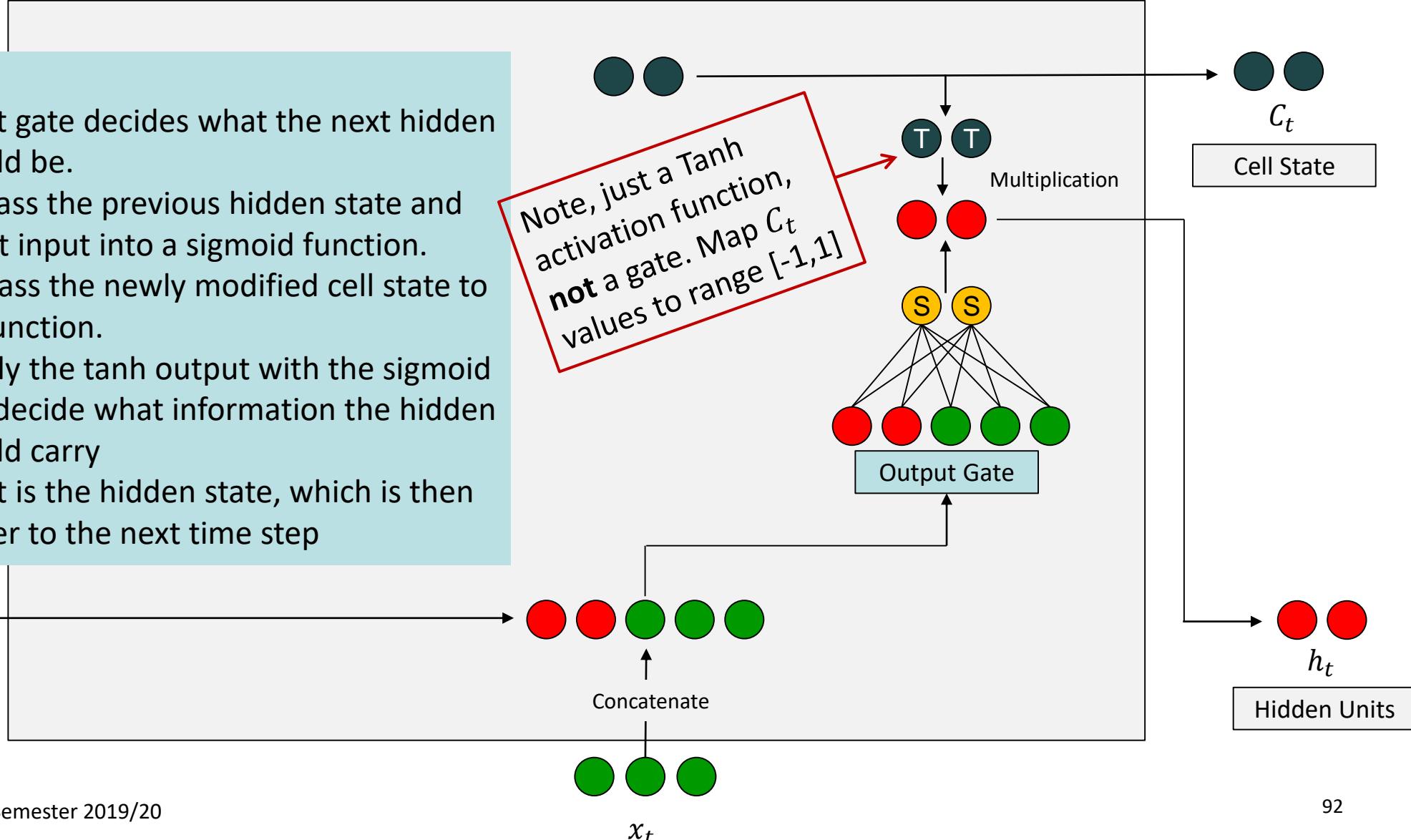
### Cell State:

- Update the cell state
- First, the cell state gets pointwise multiplied by the forget vector.
- This has a possibility of dropping values in the cell state if it gets multiplied by values near 0.
- Then we take the output from the input gate and do a pointwise addition which updates the cell state to new values that the network finds relevant
- This gives our new cell state, which is carried over to the next time step

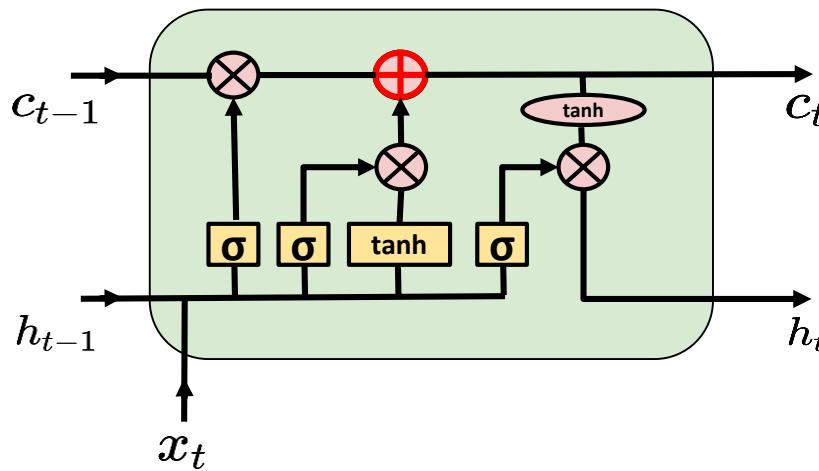
## Output Gate

- The output gate decides what the next hidden state should be.
- First, we pass the previous hidden state and the current input into a sigmoid function.
- Then we pass the newly modified cell state to the tanh function.
- We multiply the tanh output with the sigmoid output to decide what information the hidden state should carry
- The output is the hidden state, which is then carried over to the next time step

  
 $h_{t-1}$   
Hidden Units



## LSTM can deal with the gradient vanishing:

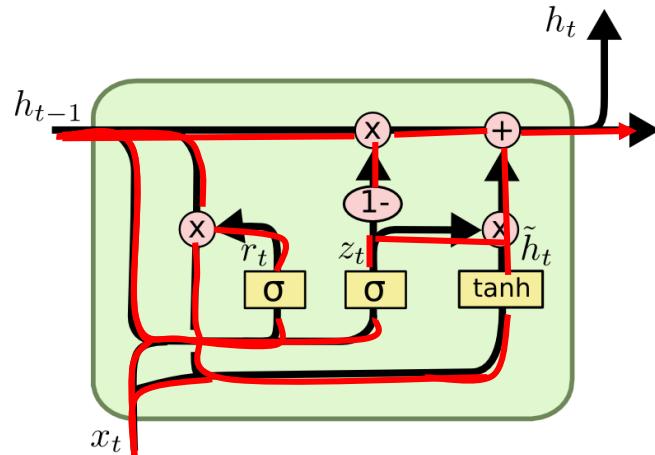


- Previous cell state and input are **added together**
- The influence never disappears unless the forget gate is completely closed
- Instead of computing new state as a matrix product with the old one, it rather computes the additive interactions
- Gradient flow is therefore improved
- Store/retrieve information over longer time periods

- **The Gated Recurrent Unit (GRU)**
  - A simplified variant of the LSTM unit
  - LSTM has 3 gates: input, output and forget
  - The GRU we only have two gates an update and a reset
    - **Update gate:** Decides how much of previous memory to keep
    - **Reset gate:** How to combine new input with previous value
  - Unlike the LSTM, the GRU has no persistent cell state

## Gated Recurrent Unit

- Single gating unit simultaneously controls the forgetting factor and the decision to update the state unit



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

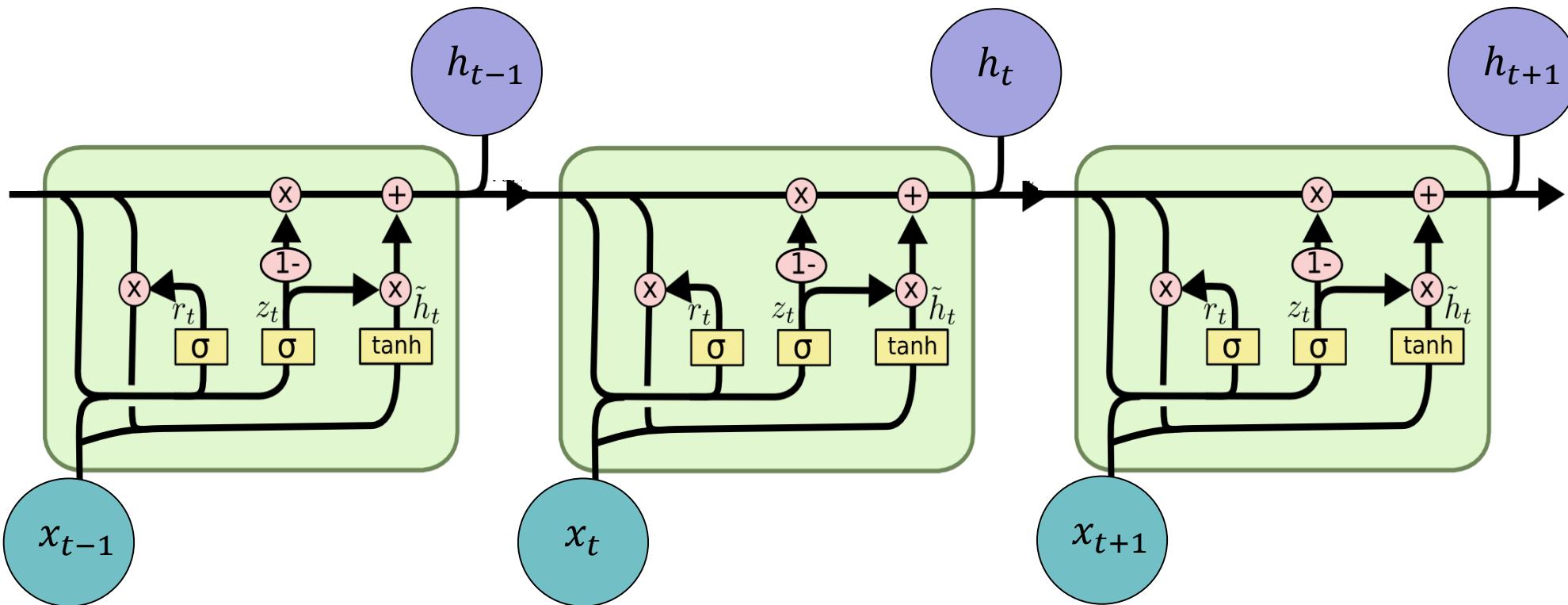
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

$$g_t = \tanh(W_g \cdot [r_t \odot h_{t-1}, x_t] + b_g)$$

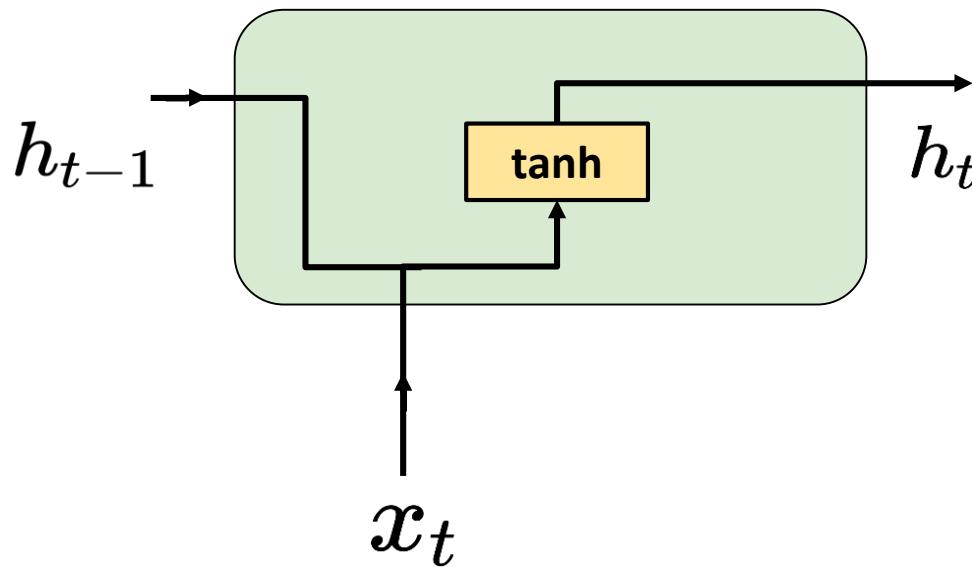
$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot g_t$$

- **Unrolled GRU network**

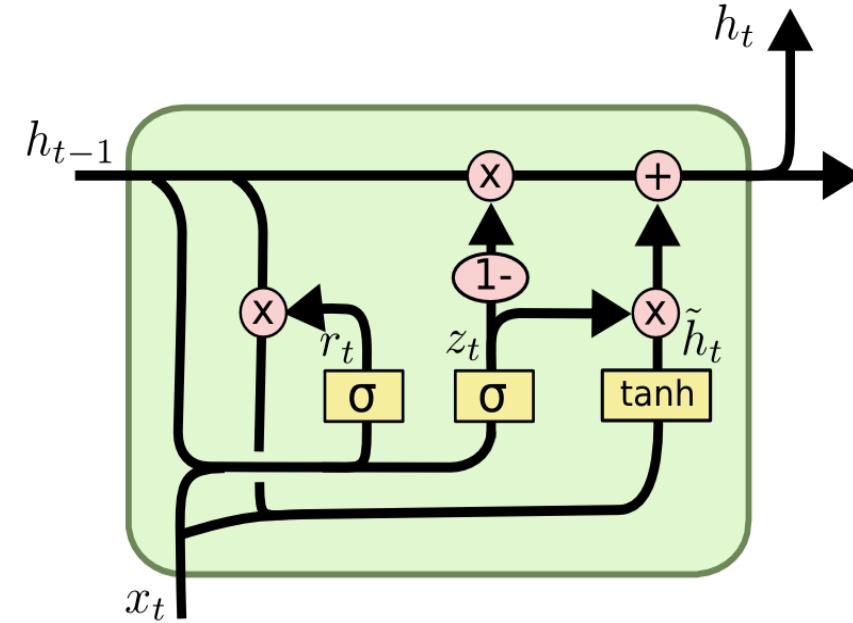
- There are three interacting networks: **reset gate**, **update gate**, **output gate**



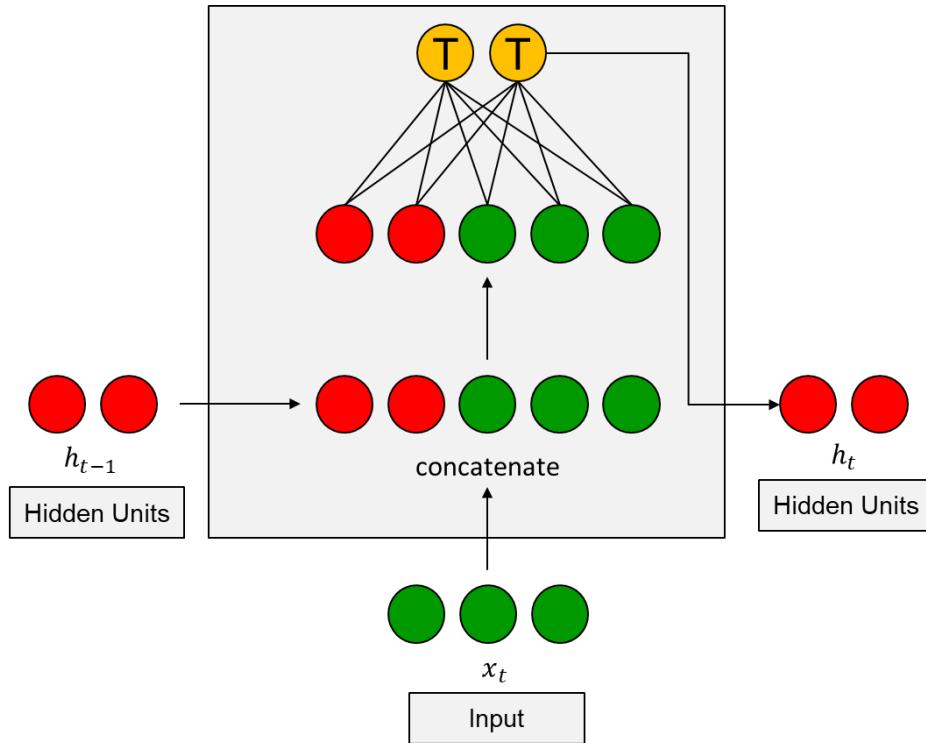
Simple RNN unit



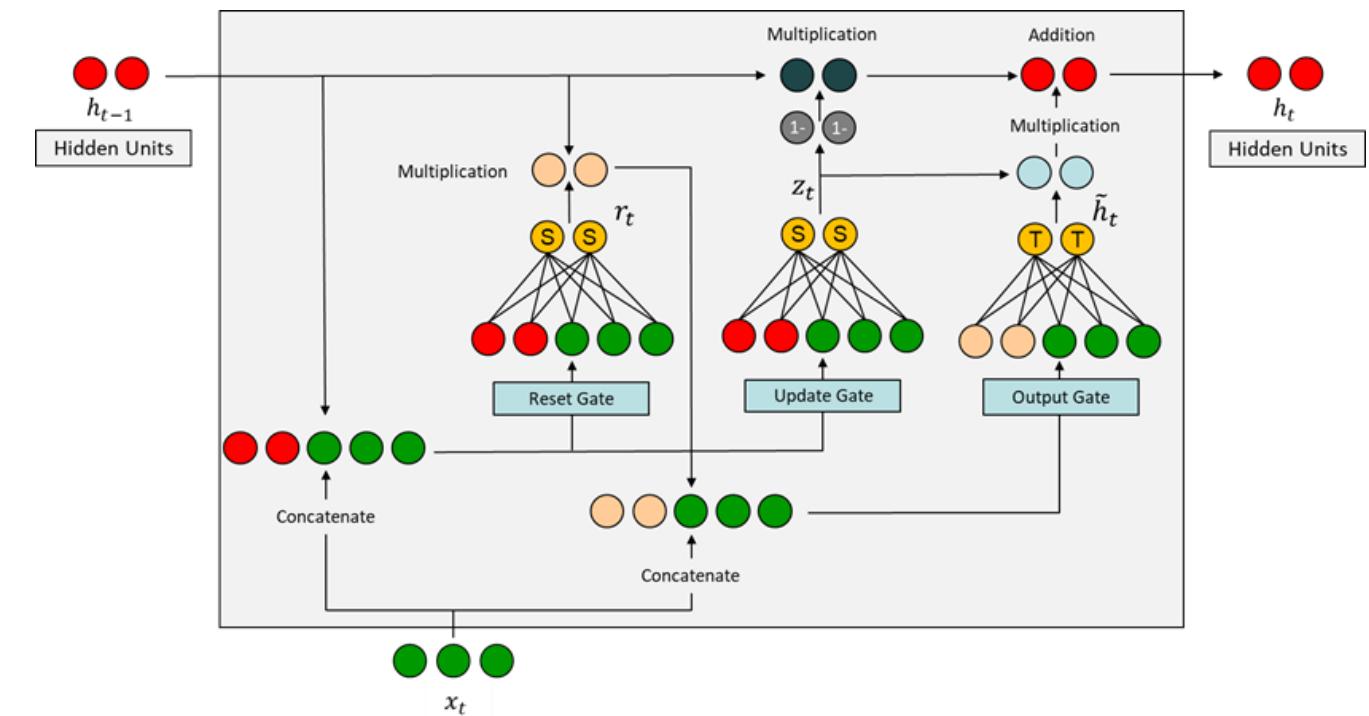
Gated Recurrent Unit

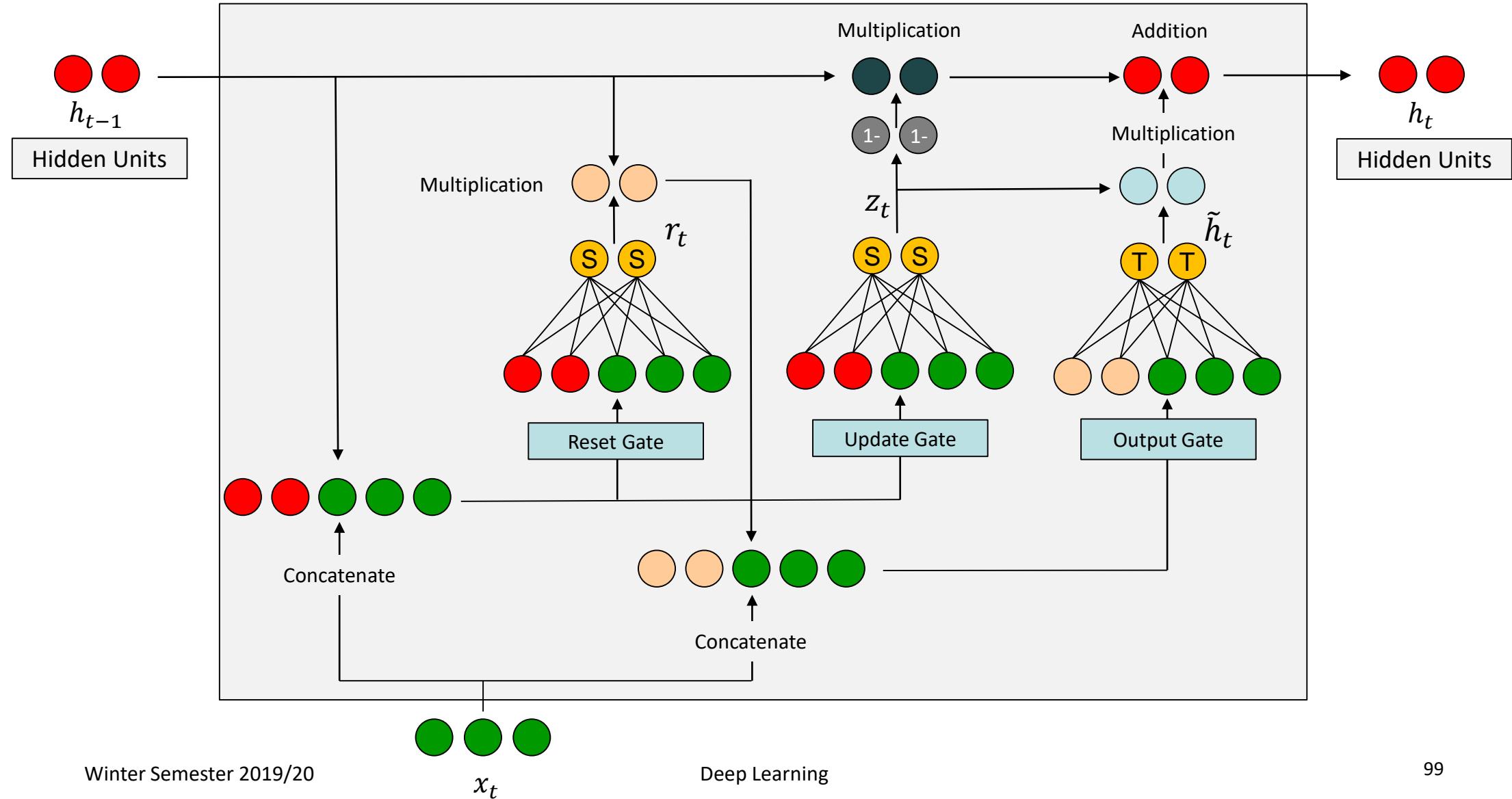


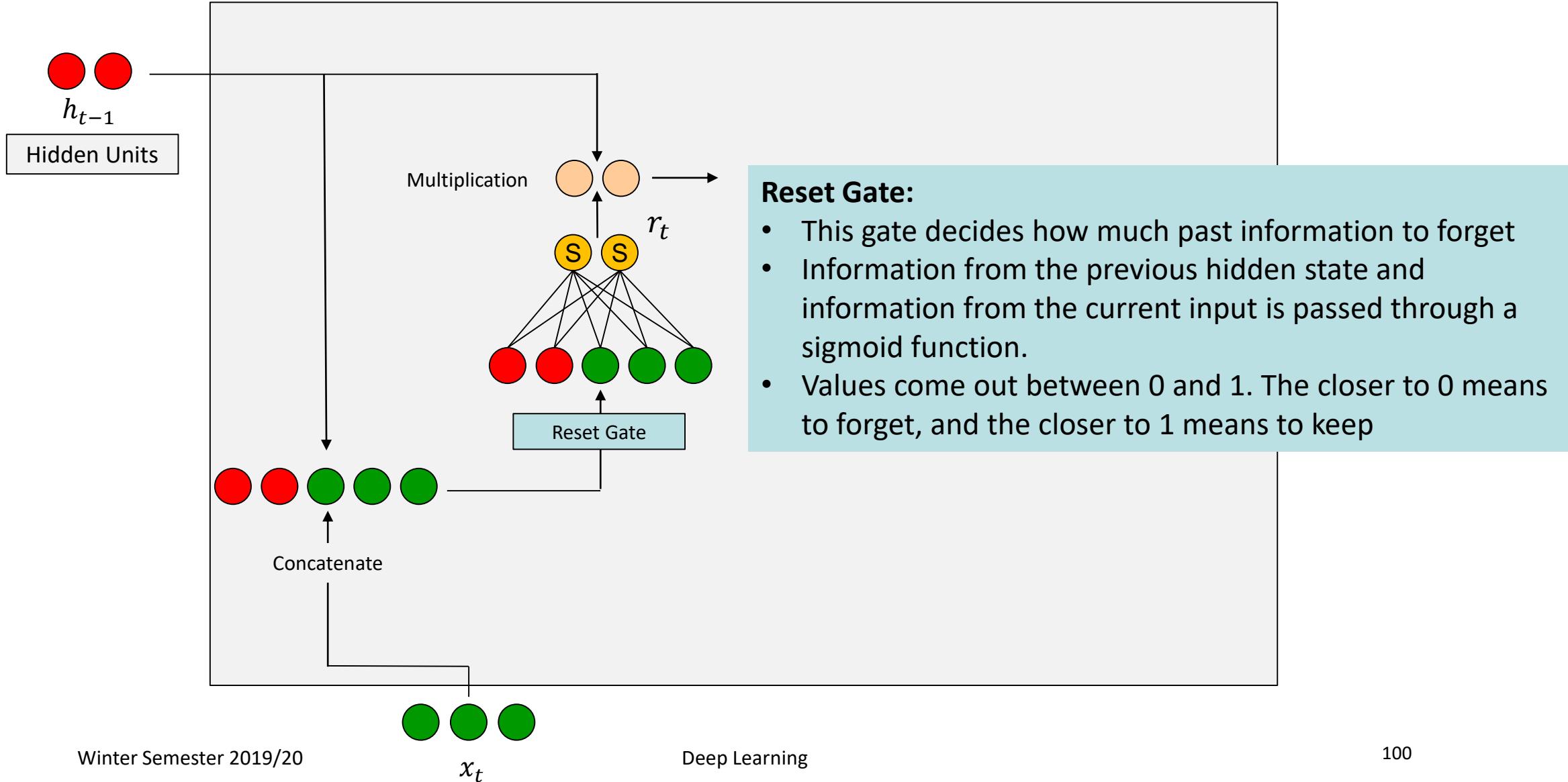
## Simple RNN unit

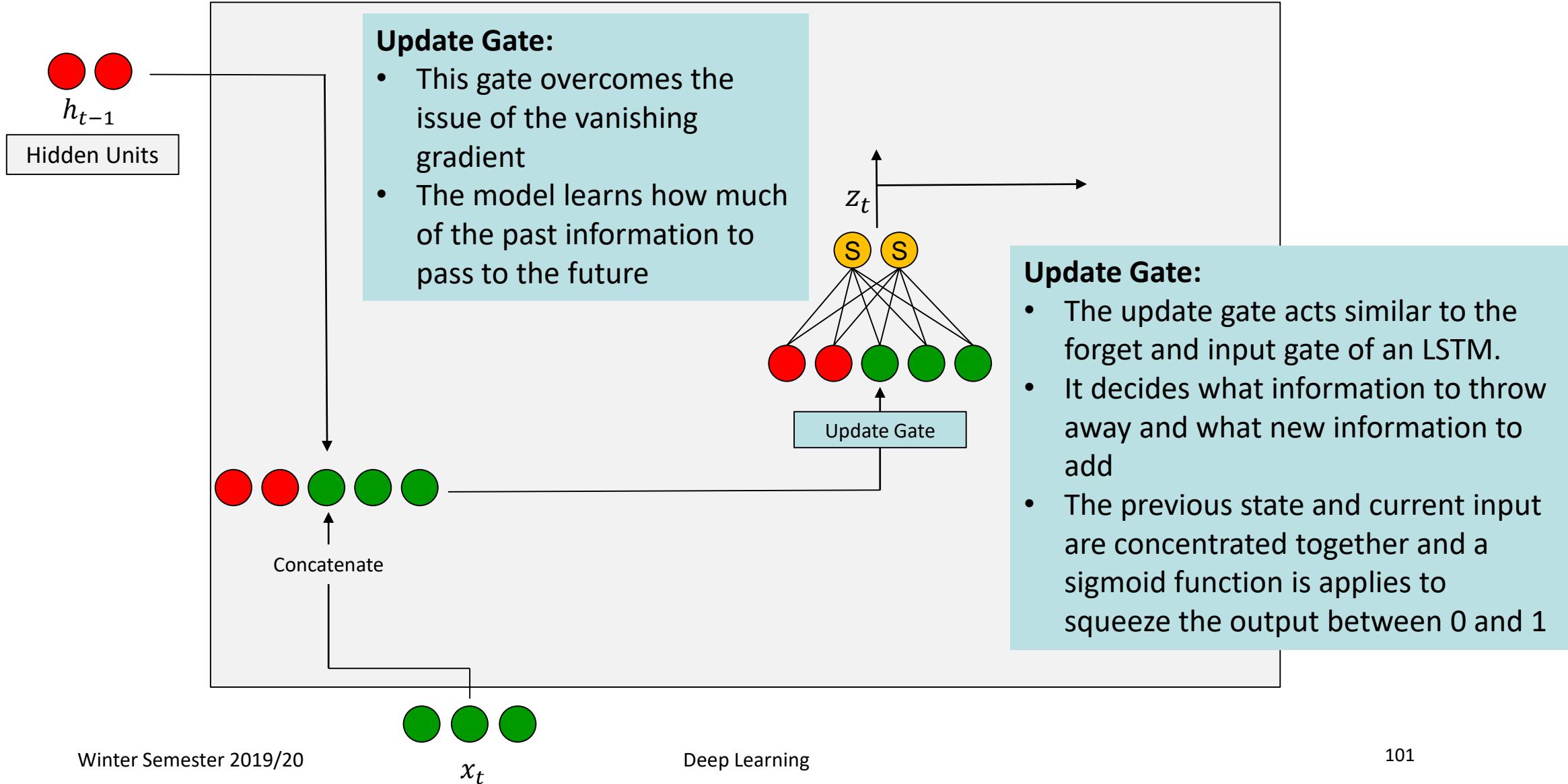


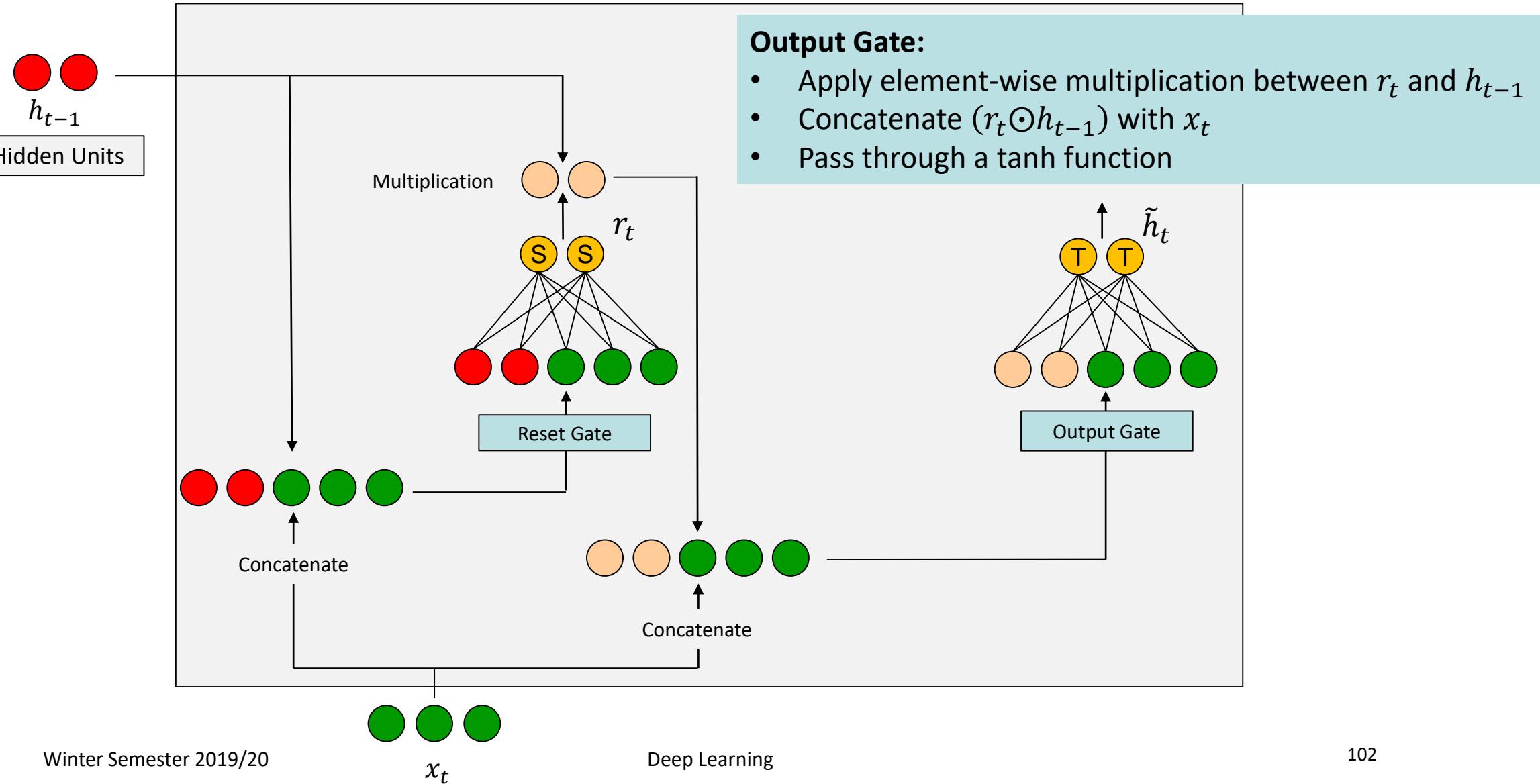
## Gated Recurrent Unit

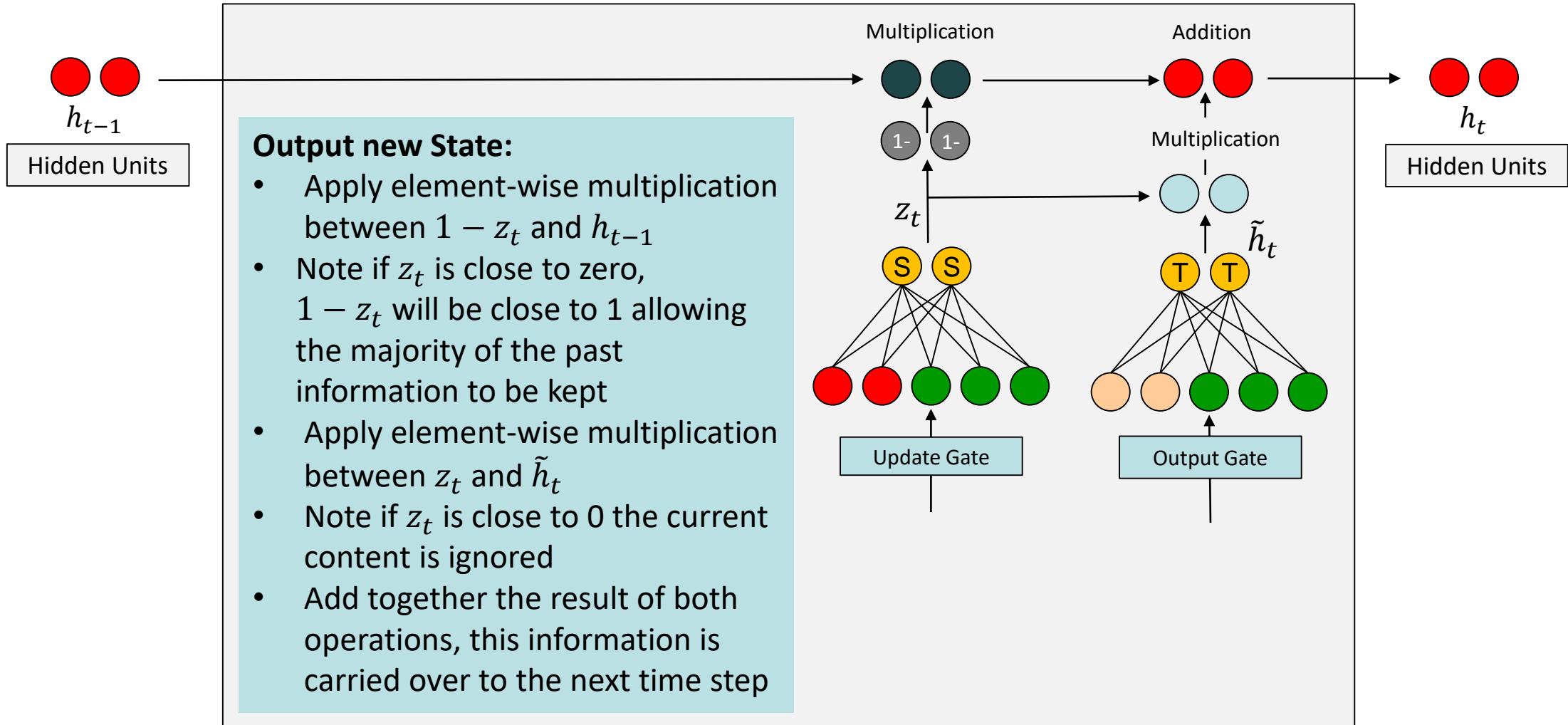






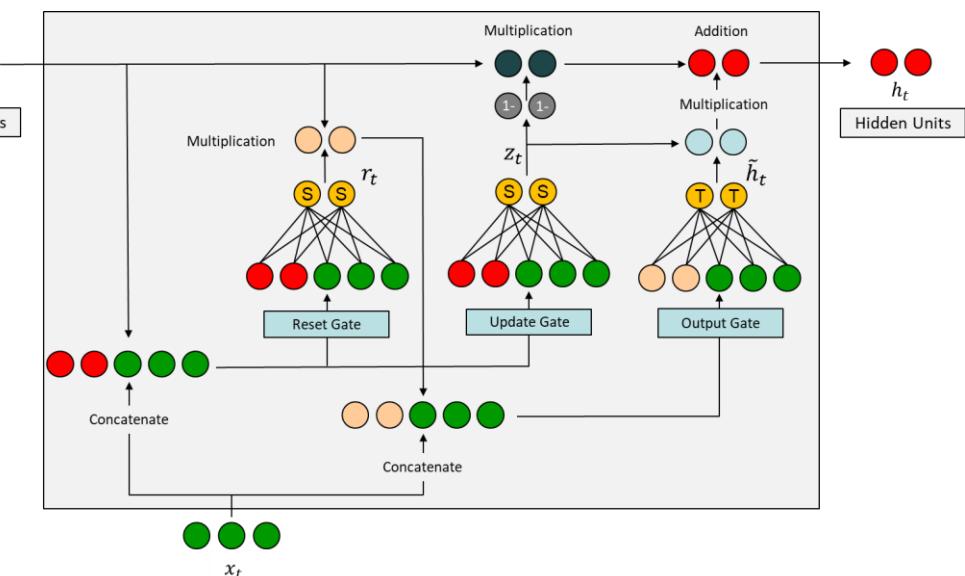
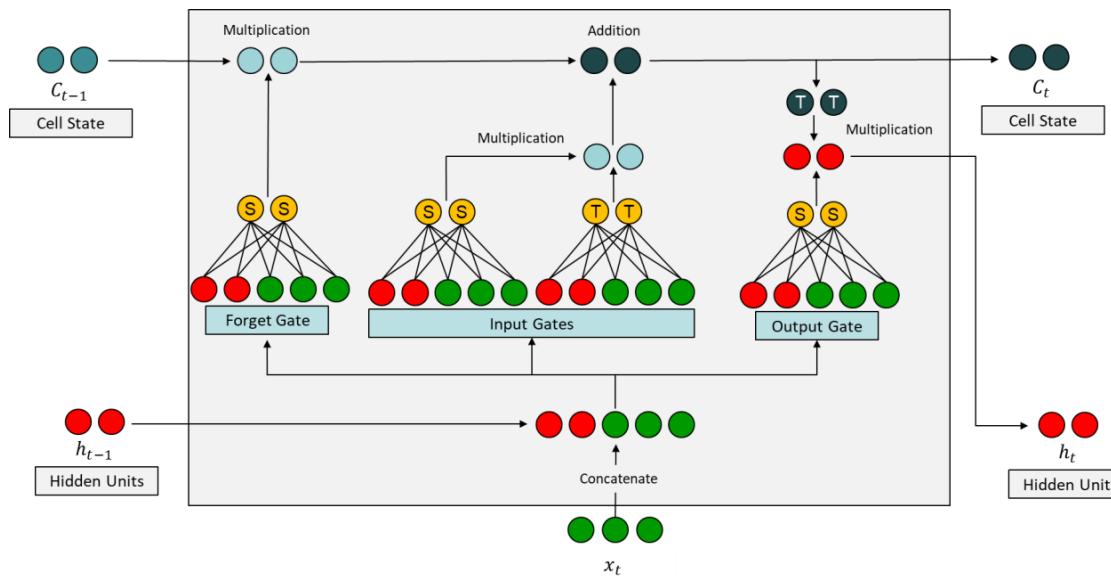






- **Comparing GRU and LSTM**

- Both GRU and LSTM are better than a simple RNN
- GRU is faster to train with fewer parameters
- GRU can perform better than LSTM on certain tasks



- RNN's incorporation of feedback into the network structure to allow for processing sequence data
- Vanilla RNNs suffers from short-term memory
  - Vanishing Gradient issues
- LSTM's and GRU's were created as a method to mitigate short-term memory using gating mechanisms
  - Gates are just neural networks that regulate the flow of information flowing through the sequence chain

# End-to-End Network

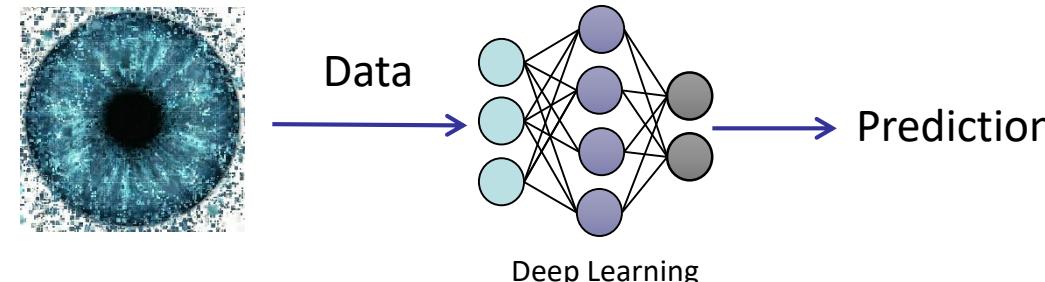
A network that accepts ‘raw’ input and produces output.  
The learning that optimizes the network weights by  
considering the inputs and outputs directly

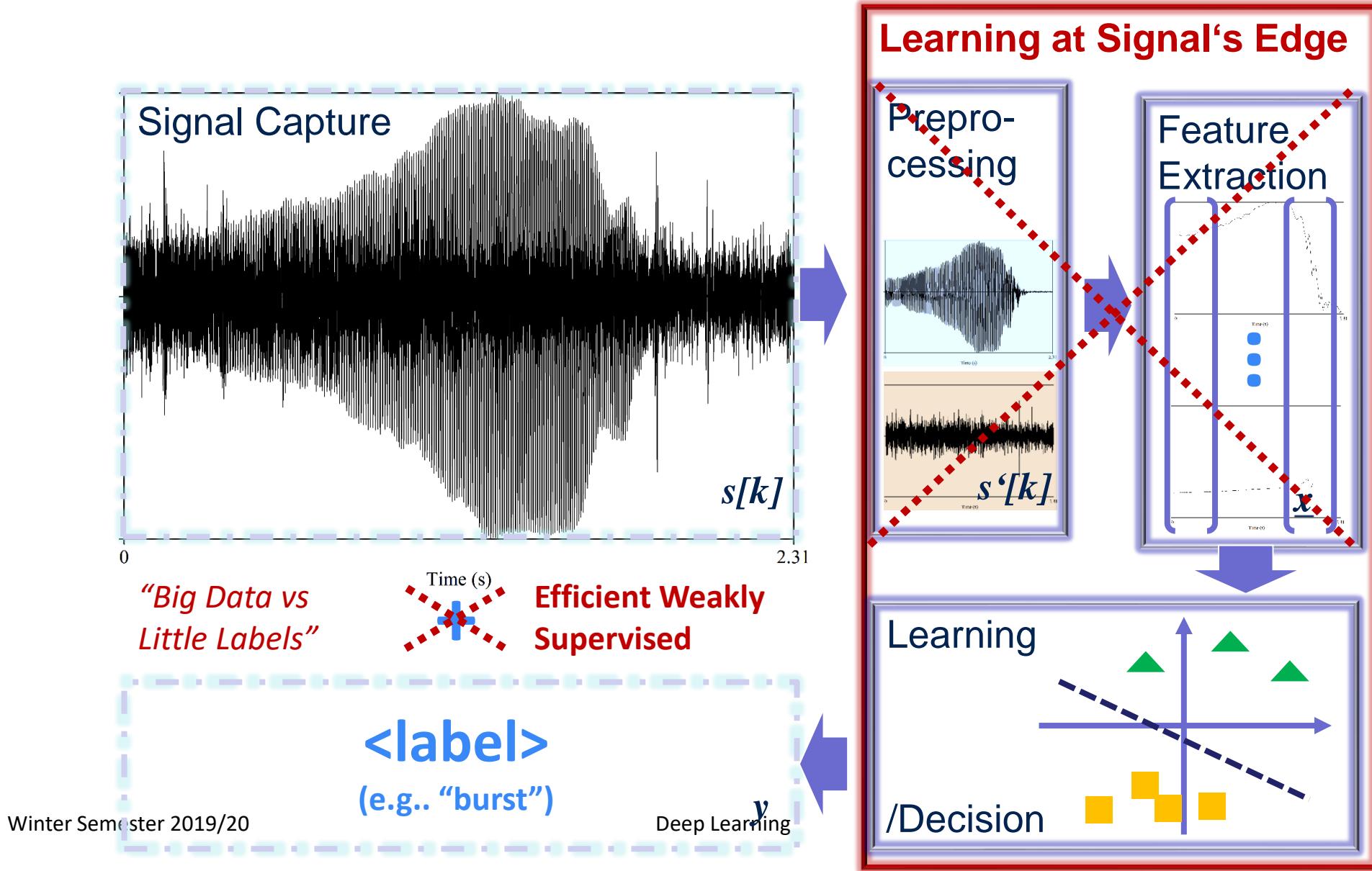
## End to End Processing

- Conventional Machine learning relies of features to reduce data complexity & make patterns more visible



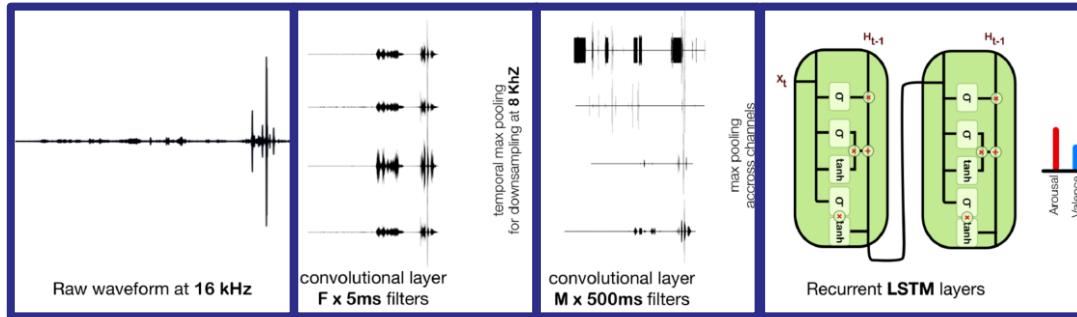
- Using CNNs and RNNs high-level features can be learnt from data in an incremental manner.
  - Eliminates the need of domain expertise and feature extraction.





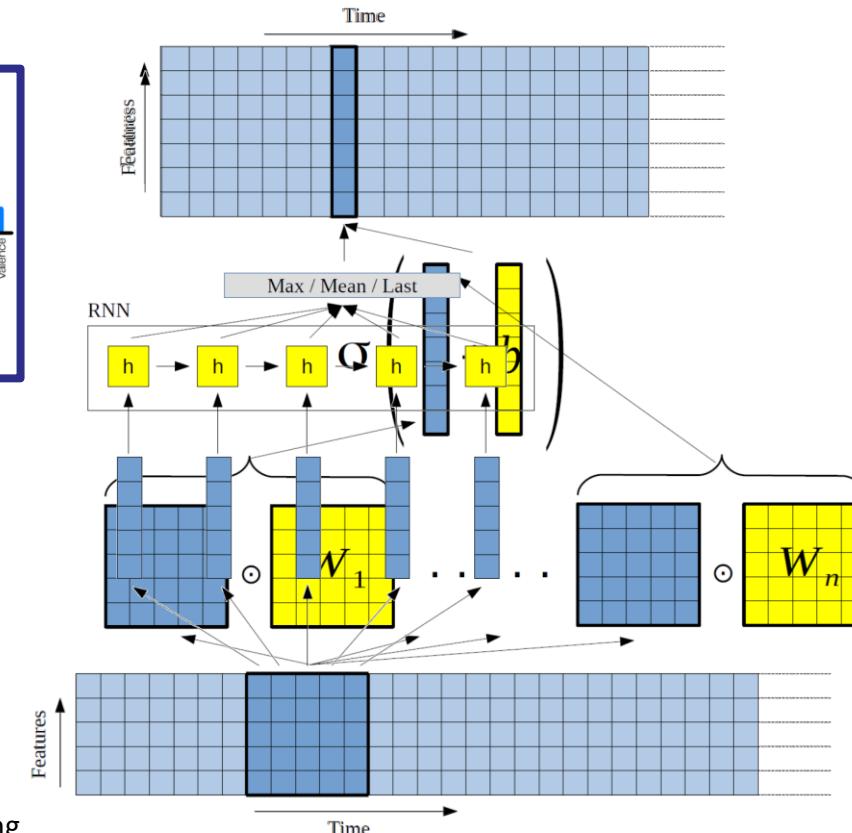
## Pattern Recognition 2.0?

- CNN + LSTM → CLSTM ?



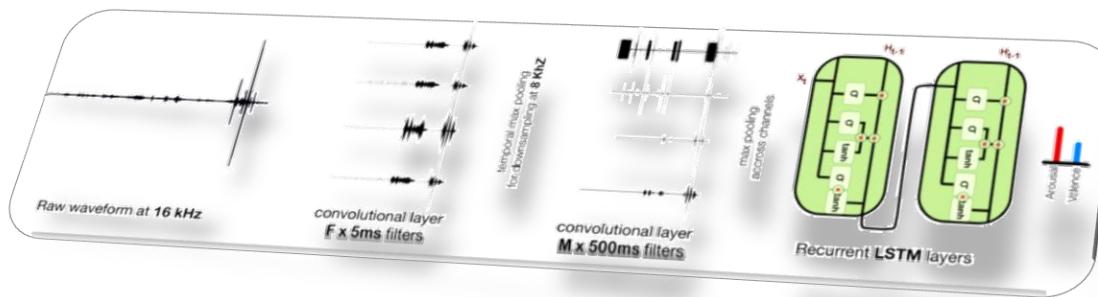
G. Trigeorgis, F. Ringeval, R. Bruckner, E. Marchi, M. Nicolaou, B. Schuller, and S. Zafeiriou, "Adieu Features? End-to-End Speech Emotion Recognition using a Deep Convolutional Recurrent Network," in Proceedings 41st IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2016, (Shanghai, P. R. China), pp. 5200–5204, IEEE, March 2016.

Arousal	CCC
Baseline	.366
e2e	.686



## End-to-End – a black box?

- CNN activations correlate with standard speech features



energy range (.77)  
loudness (.73)  
F0 mean (.71)

