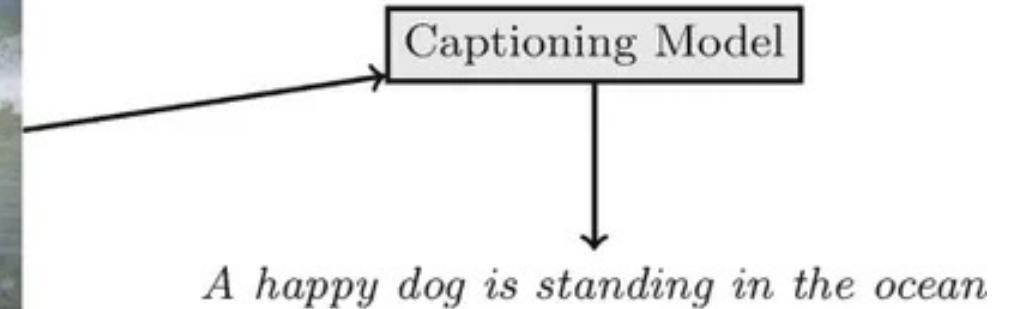
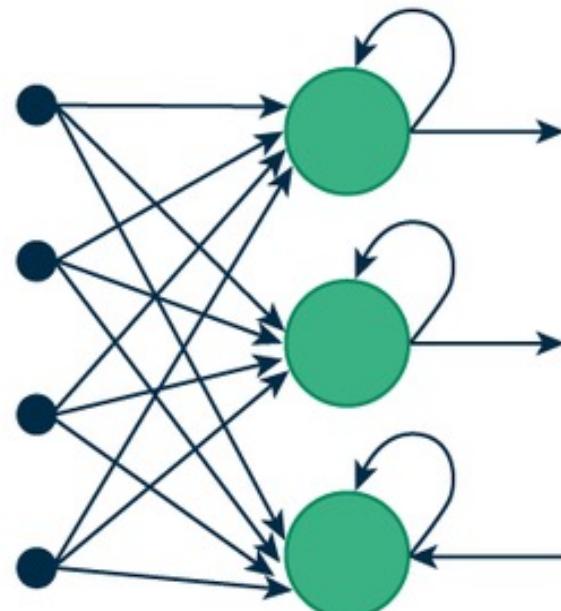


Recurrent Neural Networks

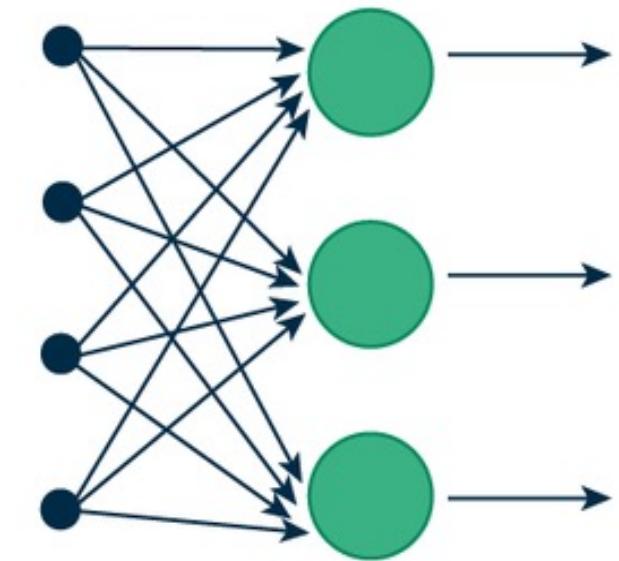
Deep Learning & Applications



Recurrent neural
networks use
sequential data to
solve common
temporal problems
seen in language
translation and
speech recognition



(a) Recurrent Neural Network



(b) Feed-Forward Neural Network

A Classic Approach for Text Classification: Bag-of-Words Model

"Raw" training dataset

$x^{[1]} = \text{"The sun is shining"}$

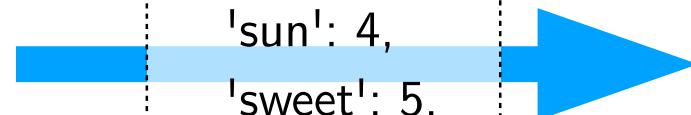
$x^{[2]} = \text{"The weather is sweet"}$

$x^{[3]} = \text{"The sun is shining,}$
 $\quad \text{the weather is sweet, and}$
 $\quad \text{one and one is two"}$

$y = [0, 1, 0]$

class labels

vocabulary = {
'and': 0,
'is': 1
'one': 2,
'shining': 3,
'sun': 4,
'sweet': 5,
'the': 6,
'two': 7,
'weather': 8,
}



Training set as design matrix

$$\mathbf{X} = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 1 \\ 2 & 3 & 2 & 1 & 1 & 1 & 2 & 1 & 1 \end{bmatrix}$$

$$\mathbf{y} = [0, 1, 0]$$

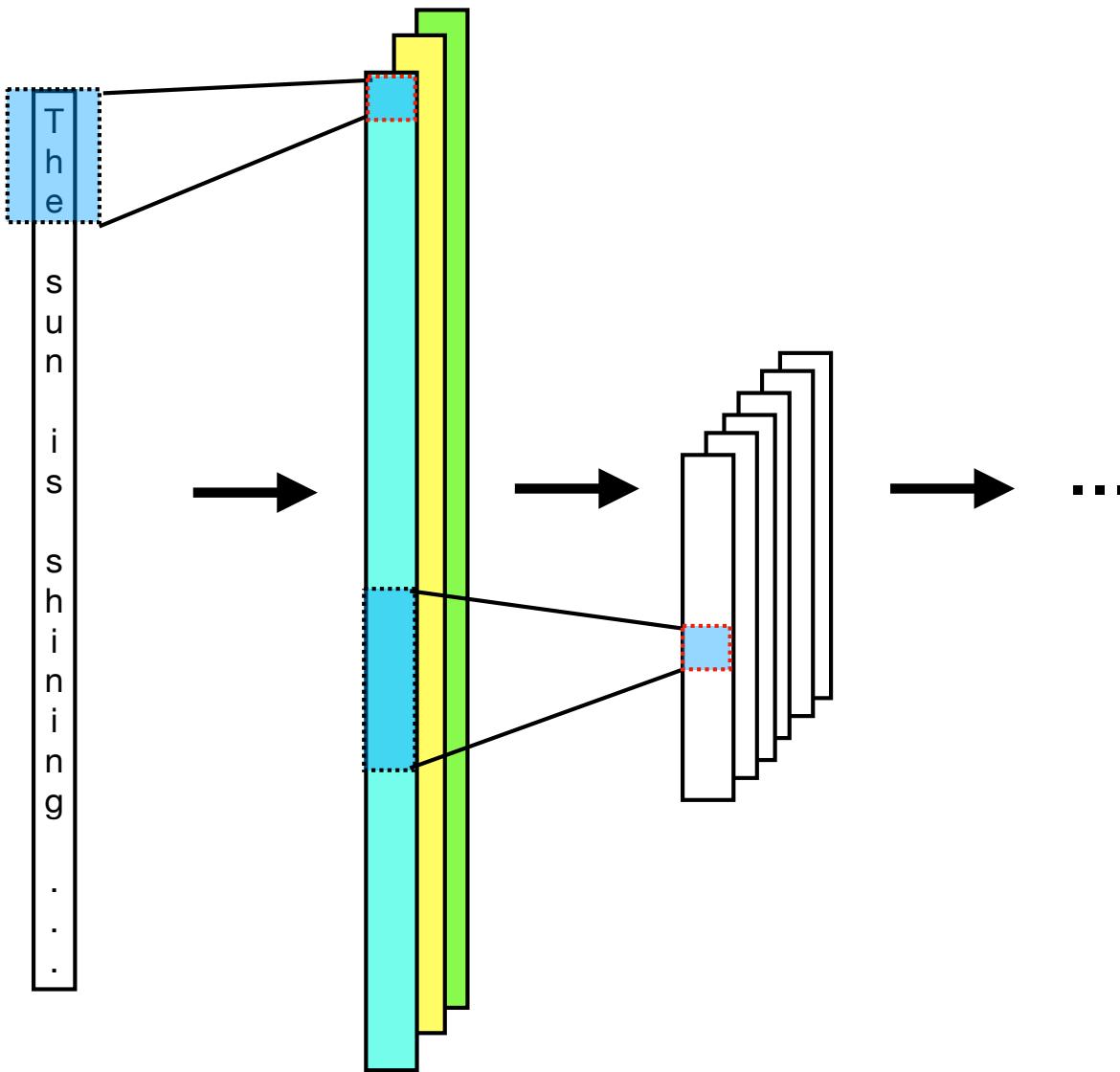
class labels

training

Classifier

(e.g., logistic regression, MLP, ...)

1D CNNs for text (and other sequence data)



Sequence data: order matters

The movie my friend has not seen is good

The movie my friend has seen is not good

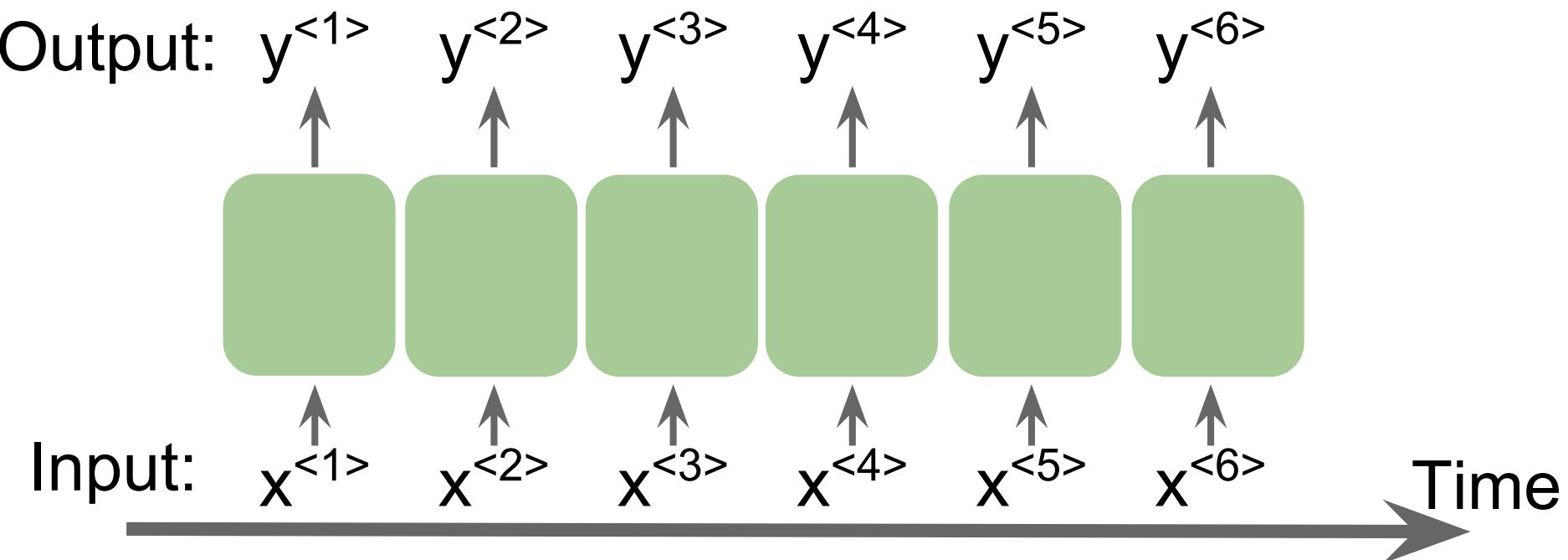
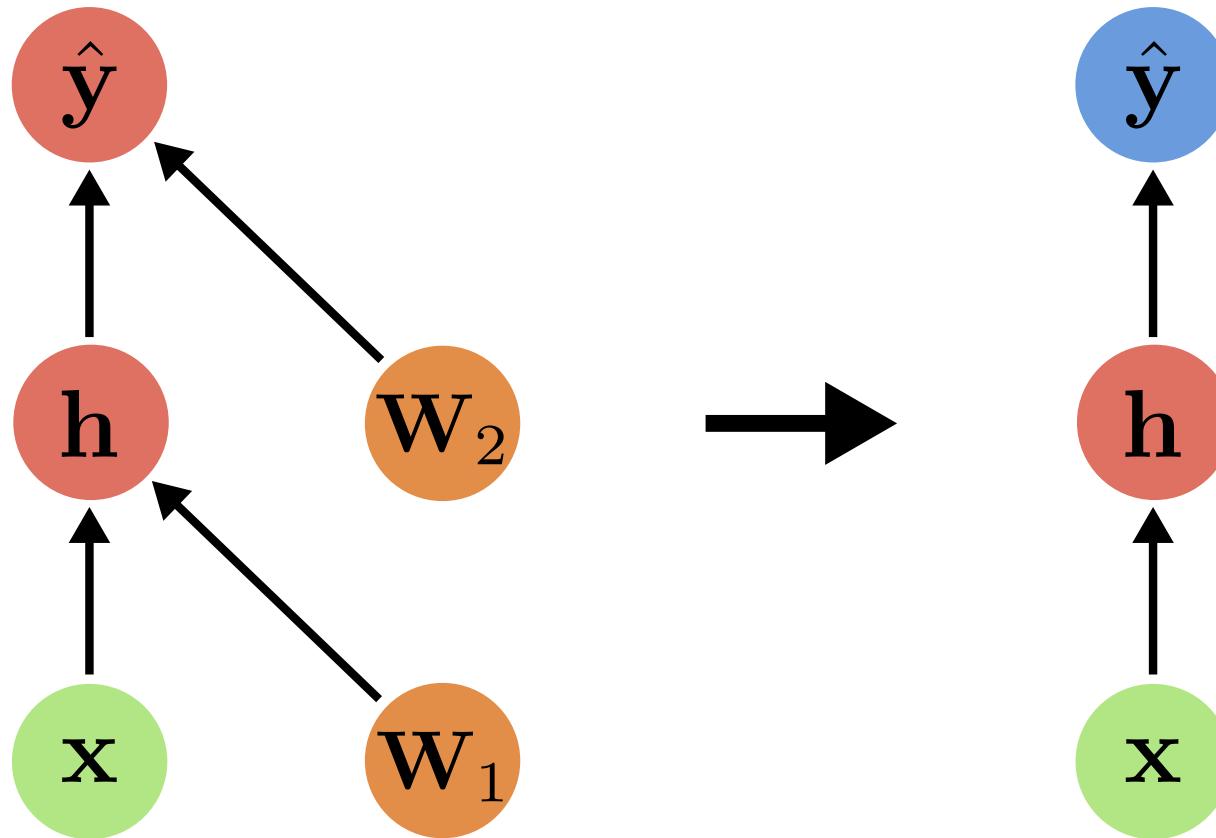


Image source: Sebastian Raschka, Vahid
Mirjalili. *Python Machine Learning*. 3rd
Edition. Packt, 2019

Recap: Computation Graphs



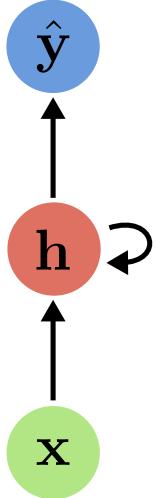
- ▶ Example of a **feedforward neural network** (here: MLP with 1 hidden layer)

Feedforward vs. Recurrent Neural Networks

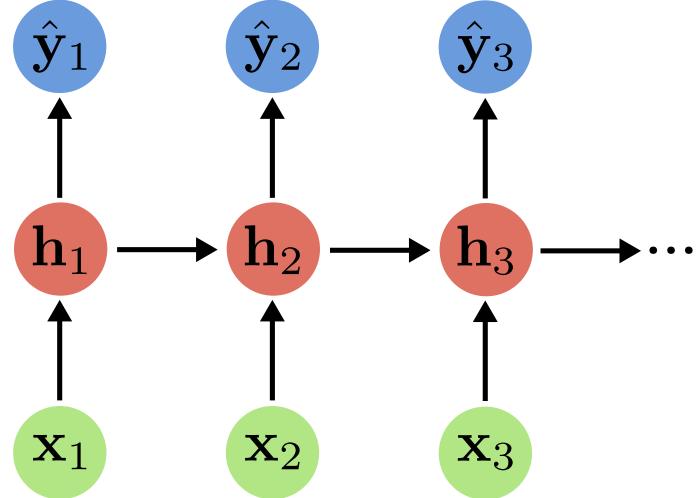
Feedforward
Neural Network



Recurrent Neural Network (RNN)
with feedback connection



Recurrent Neural Network (RNN)
unrolled over time (index = time t)

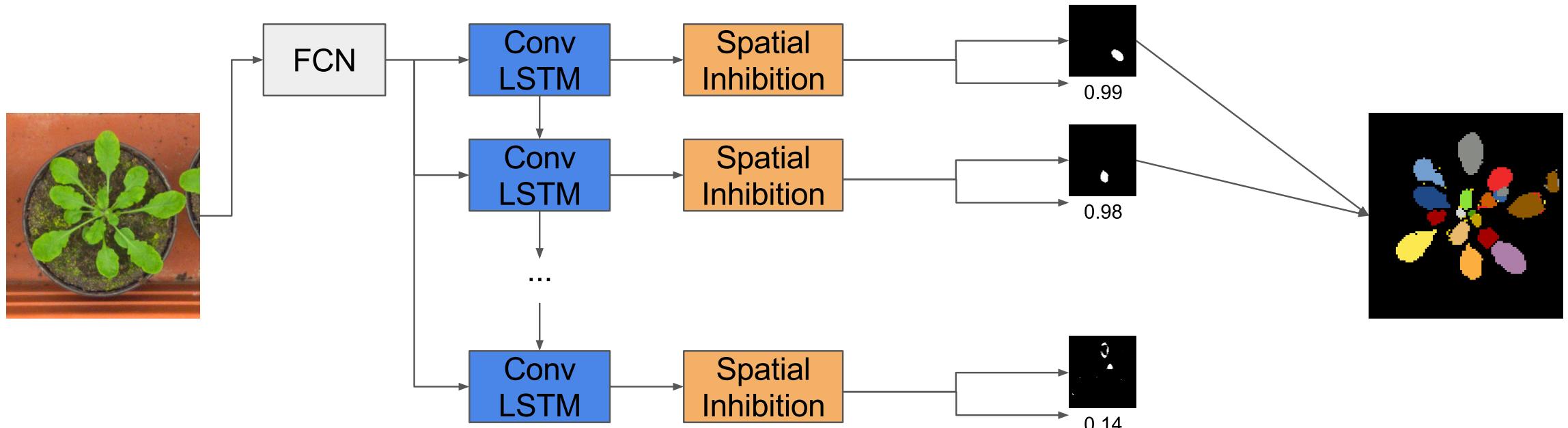


Recurrent Neural Networks (RNNs):

- ▶ Core idea: update **hidden state h** based on input and previous hidden state using same update rule (same/shared parameters) at each **time step**
- ▶ Allows for processing **sequences of variable length**, not only fixed-sized vectors
- ▶ **Infinite memory:** h is function of all previous inputs (long-term dependencies)

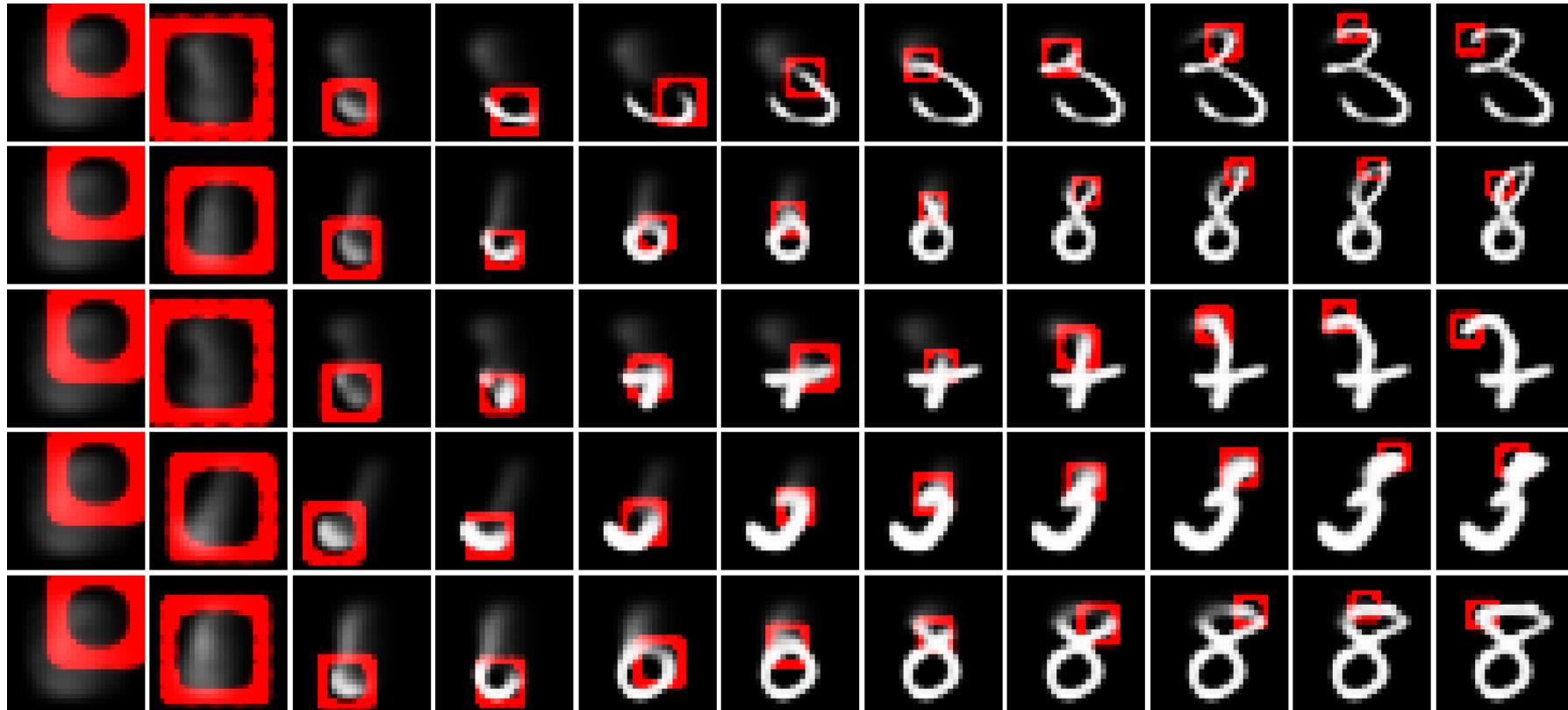
Recurrent Network Applications

Recurrent Instance Segmentation



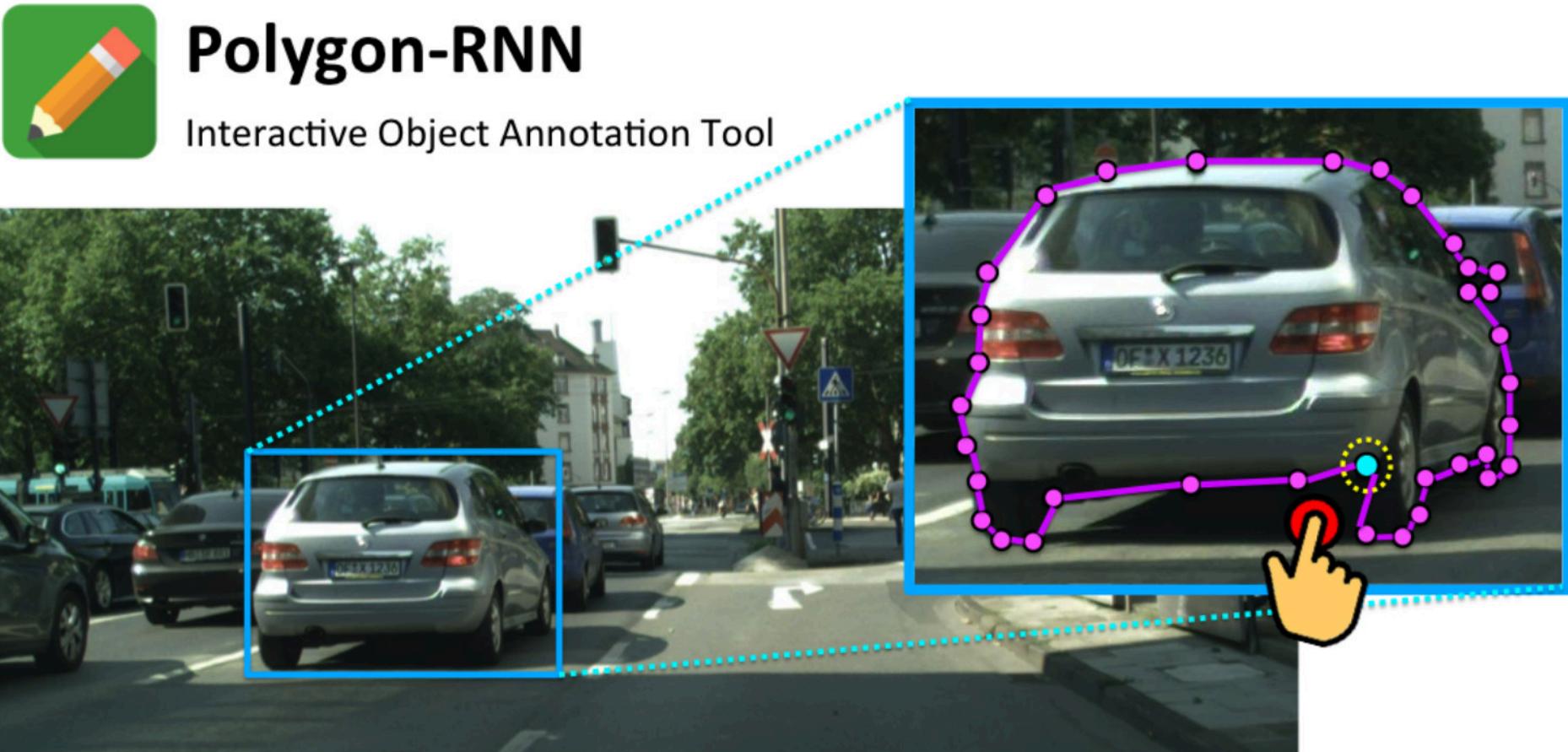
- At each time step, **segment** the next (not yet segmented) **part** of an object

Image Generation



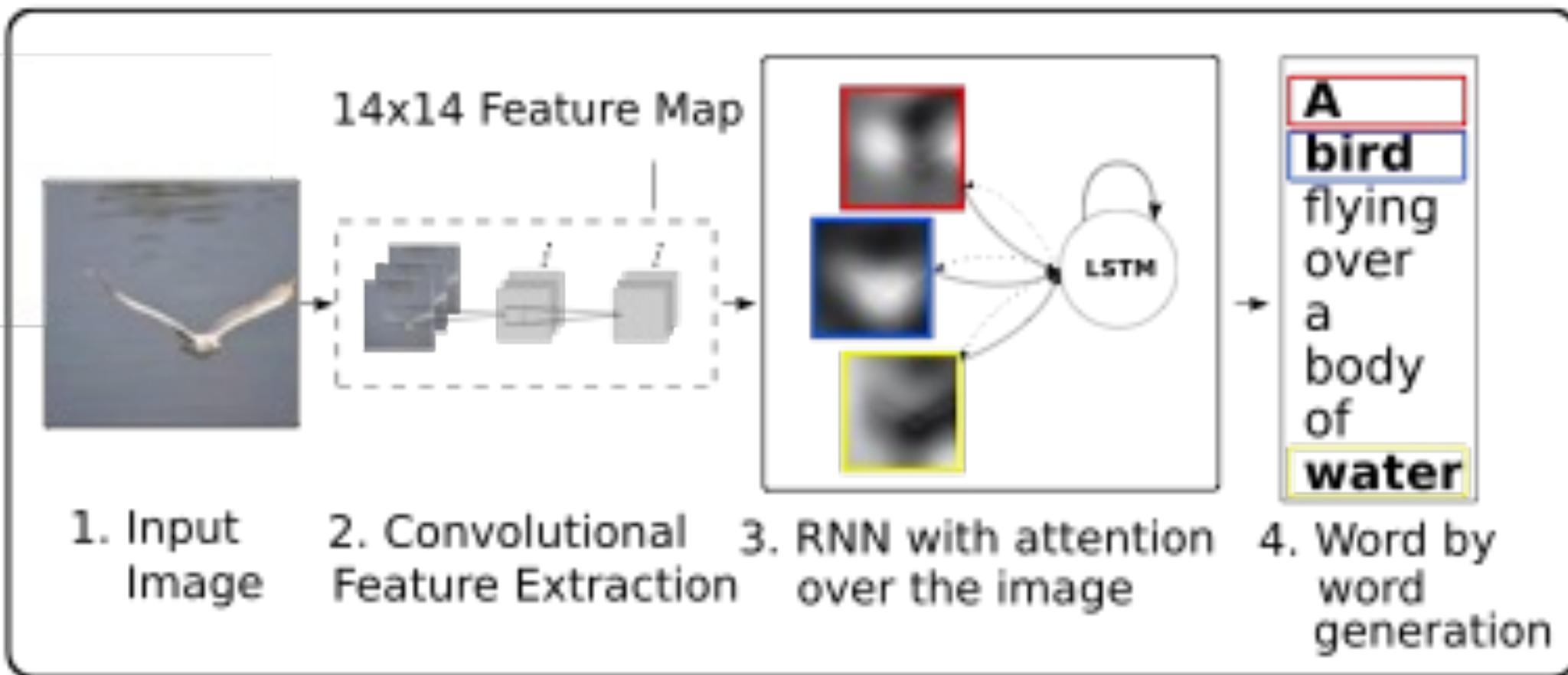
- ▶ Model for **sequential image generation** (red rectangle = attended region)
- ▶ Demonstrates that RNNs can also process/generate non-sequential data

Image Annotation



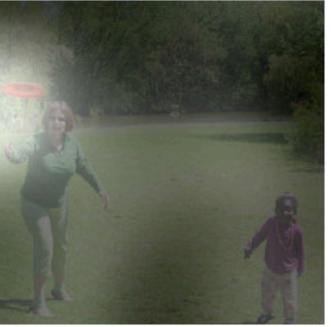
- ▶ Iteratively **annotate** the outline of 2D object instances in an image with **polygons**

Image Captioning



- Generate **image description** by sequentially looking at an image

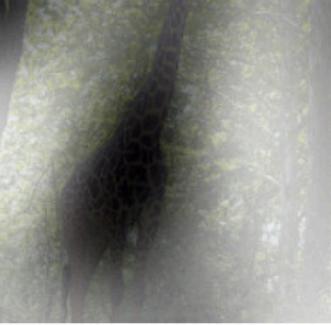
Image Captioning



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

► Successful caption generations

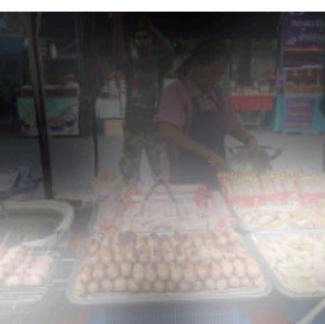
Image Captioning



A large white bird standing in a forest.

A woman holding a clock in her hand.

A man wearing a hat and
a hat on a skateboard.



A person is standing on a beach
with a surfboard.

A woman is sitting at a table
with a large pizza.

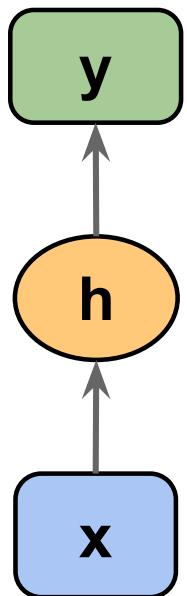


A man is talking on his cell phone
while another man watches.

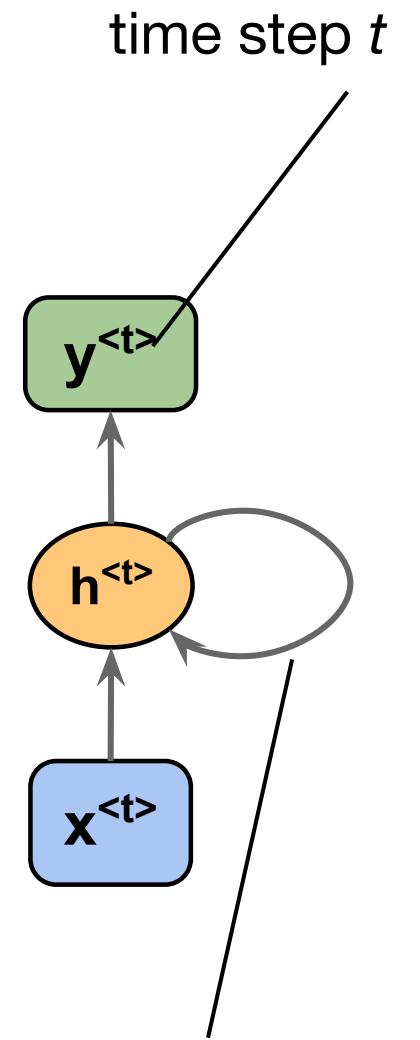
- ▶ Wrongly generated captions. Attention can reveal insights into what went wrong.

Overview

Networks we used previously: also called feedforward neural networks



Recurrent Neural Network (RNN)



Overview

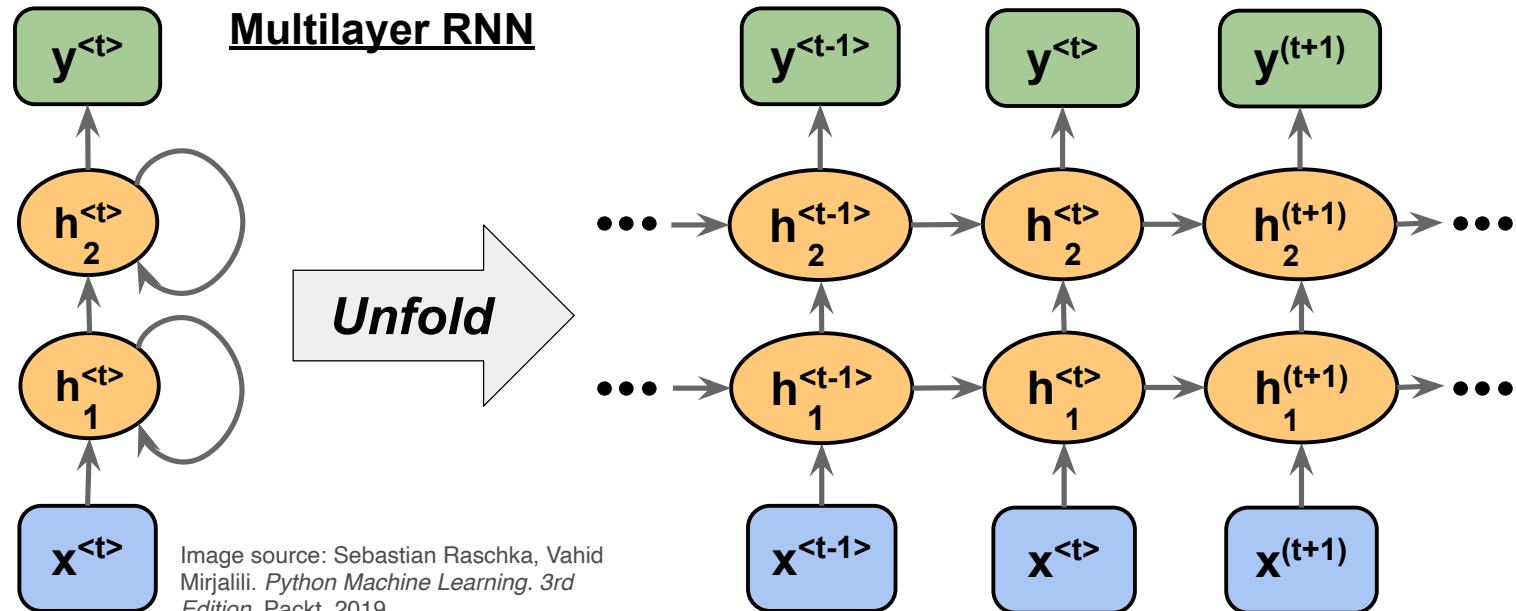
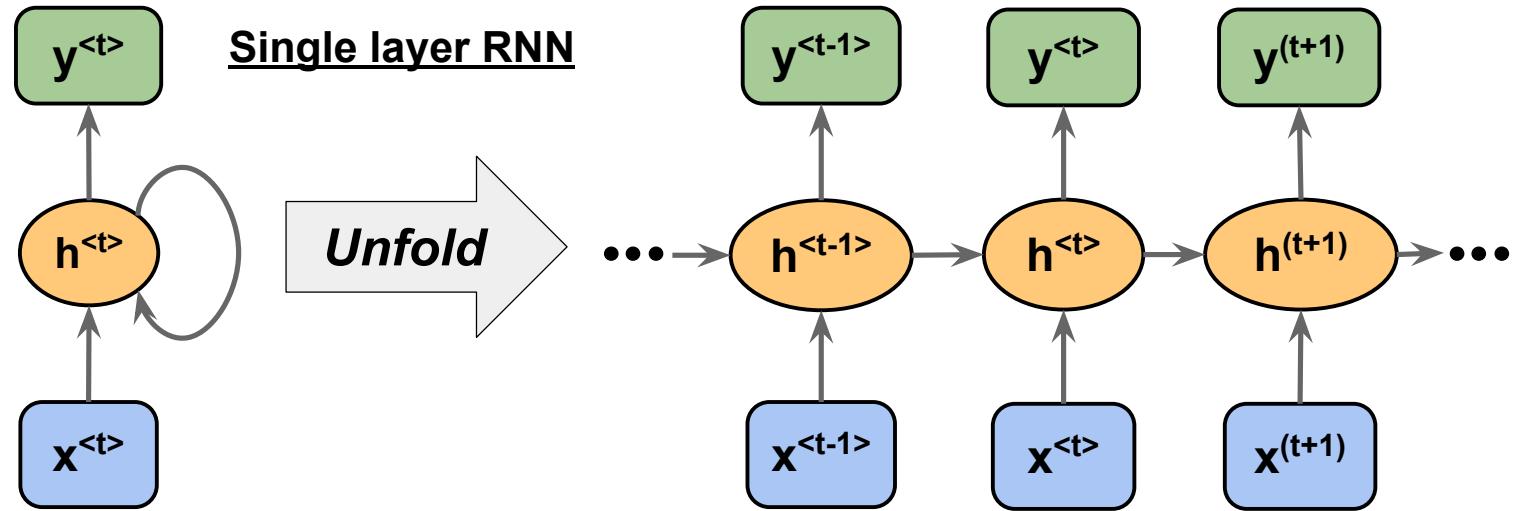
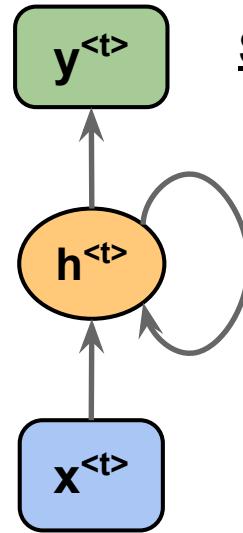
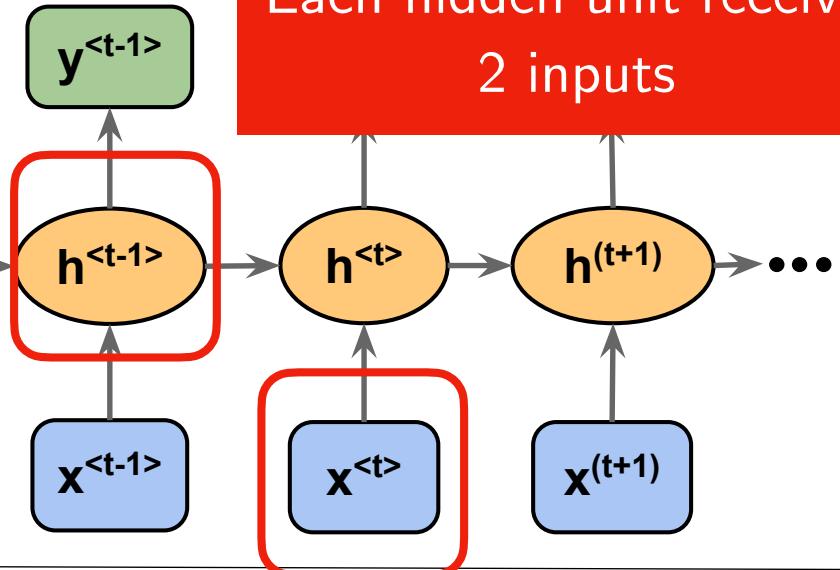


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

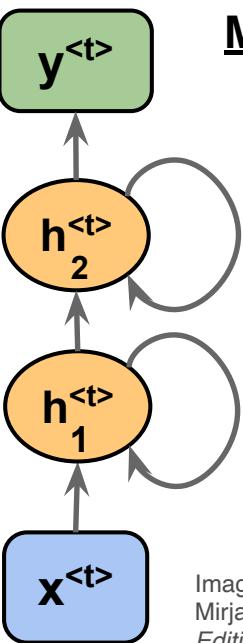


Single layer RNN

Unfold

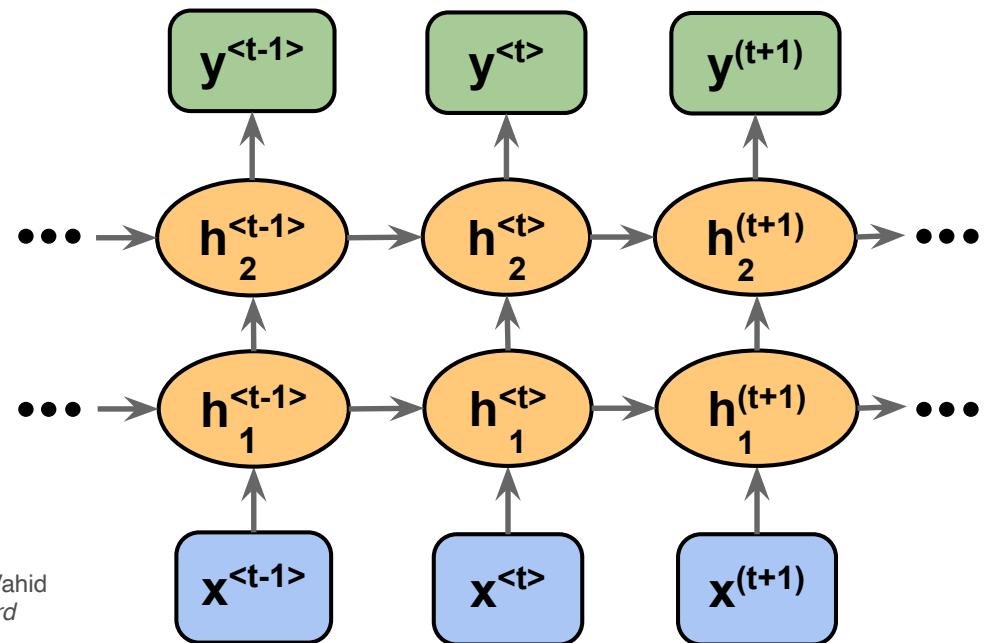


Each hidden unit receives
2 inputs

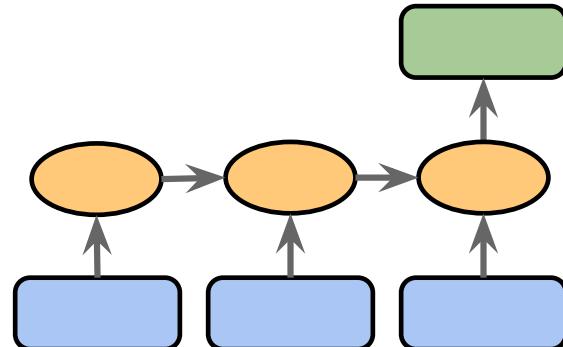


Multilayer RNN

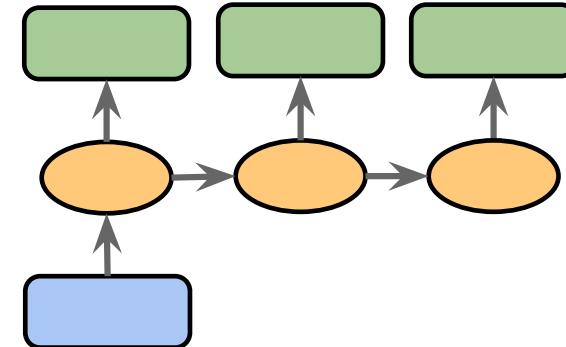
Unfold



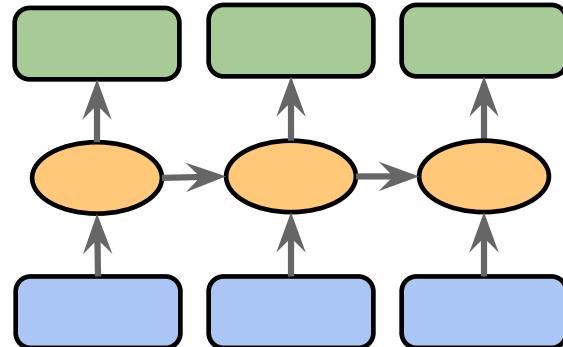
Different Types of Sequence Modeling Tasks



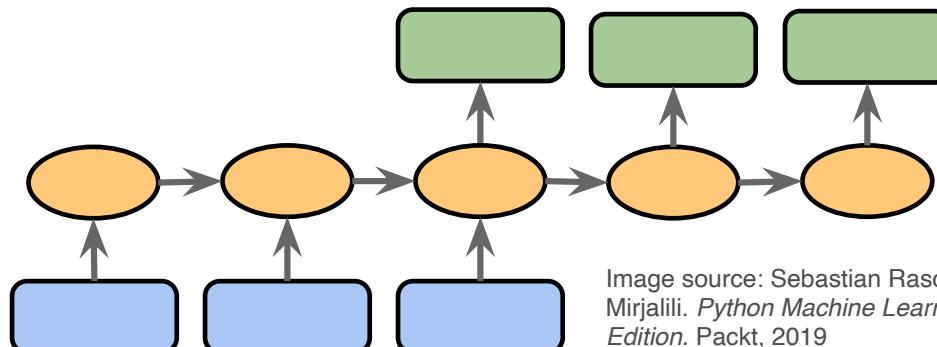
many-to-one



one-to-many



many-to-many



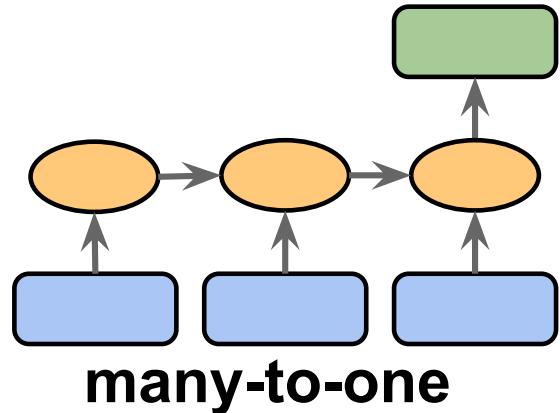
many-to-many

Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning, 3rd Edition*. Packt, 2019

Figure based on:

The Unreasonable Effectiveness of Recurrent Neural Networks by Andrej Karpathy (<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>)

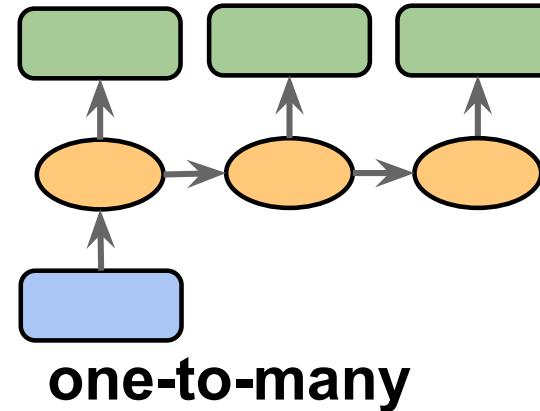
Different Types of Sequence Modeling Tasks



Many-to-one: The input data is a sequence, but the output is a fixed-size vector, not a sequence.

Ex.: sentiment analysis, the input is some text, and the output is a class label.

Different Types of Sequence Modeling Tasks



One-to-many: Input data is in a standard format (not a sequence), the output is a sequence.

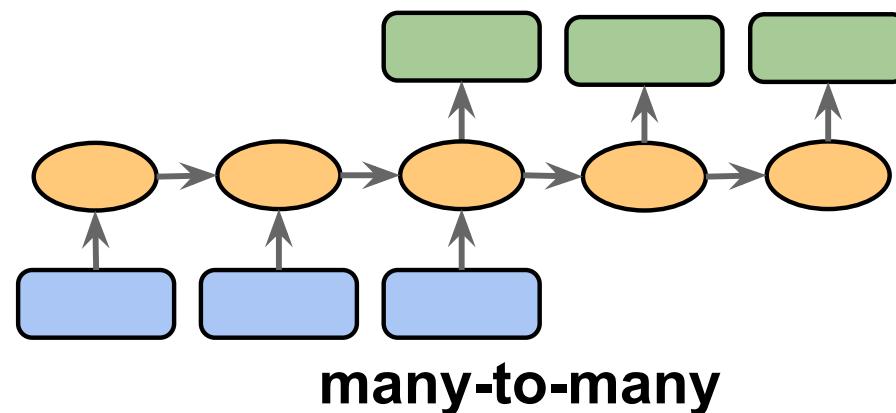
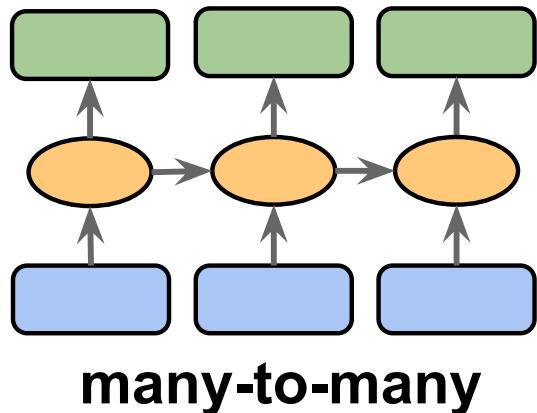
Ex.: Image captioning, where the input is an image, the output is a text description of that image

Different Types of Sequence Modeling Tasks

Many-to-many: Both inputs and outputs are sequences. Can be direct or delayed.

Ex.: Video-captioning, i.e., describing a sequence of images via text (direct).

Translating one language into another (delayed)



Weight matrices in a single-hidden layer RNN

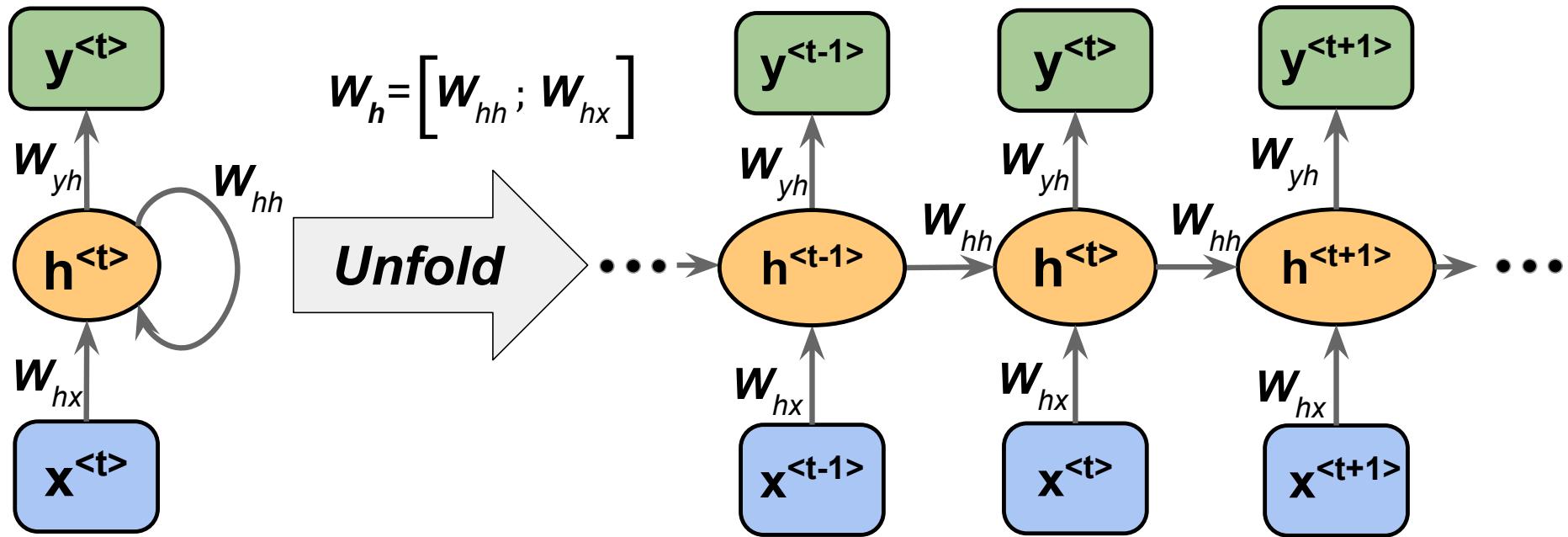
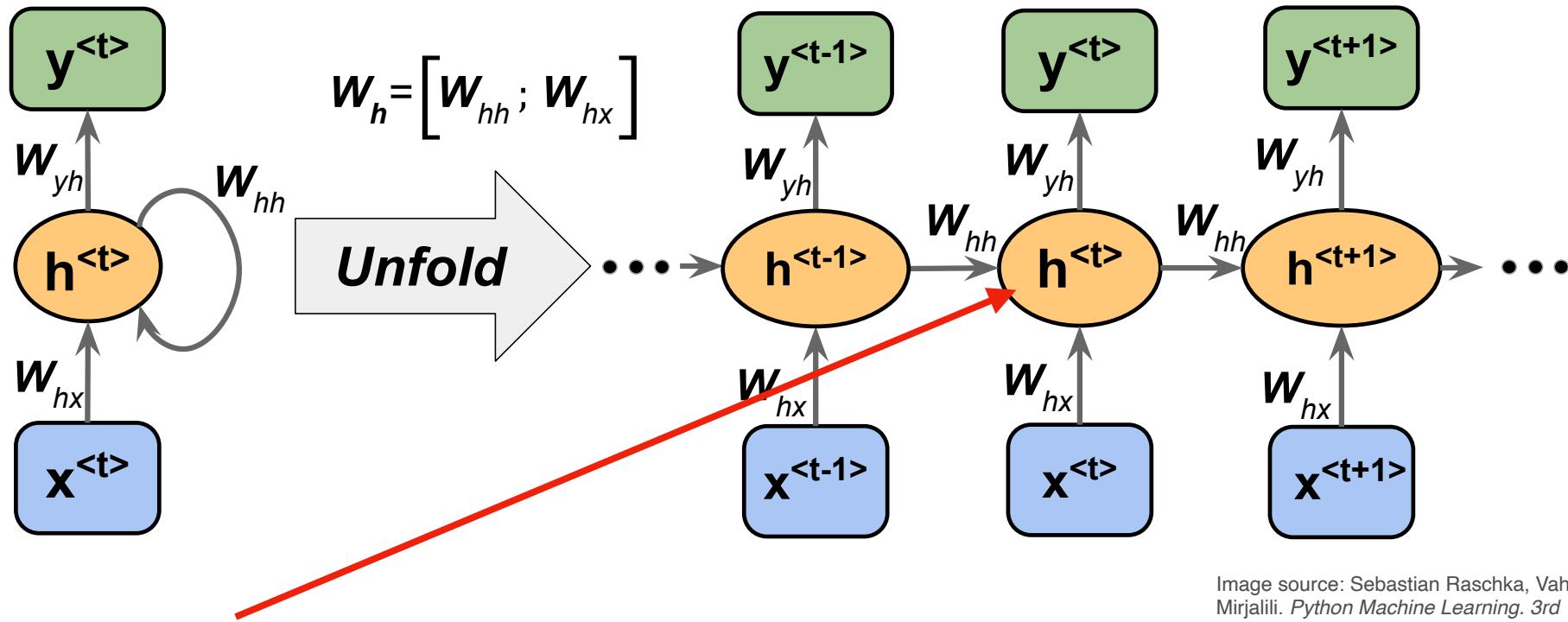


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Weight matrices in a single-hidden layer RNN



Net input:

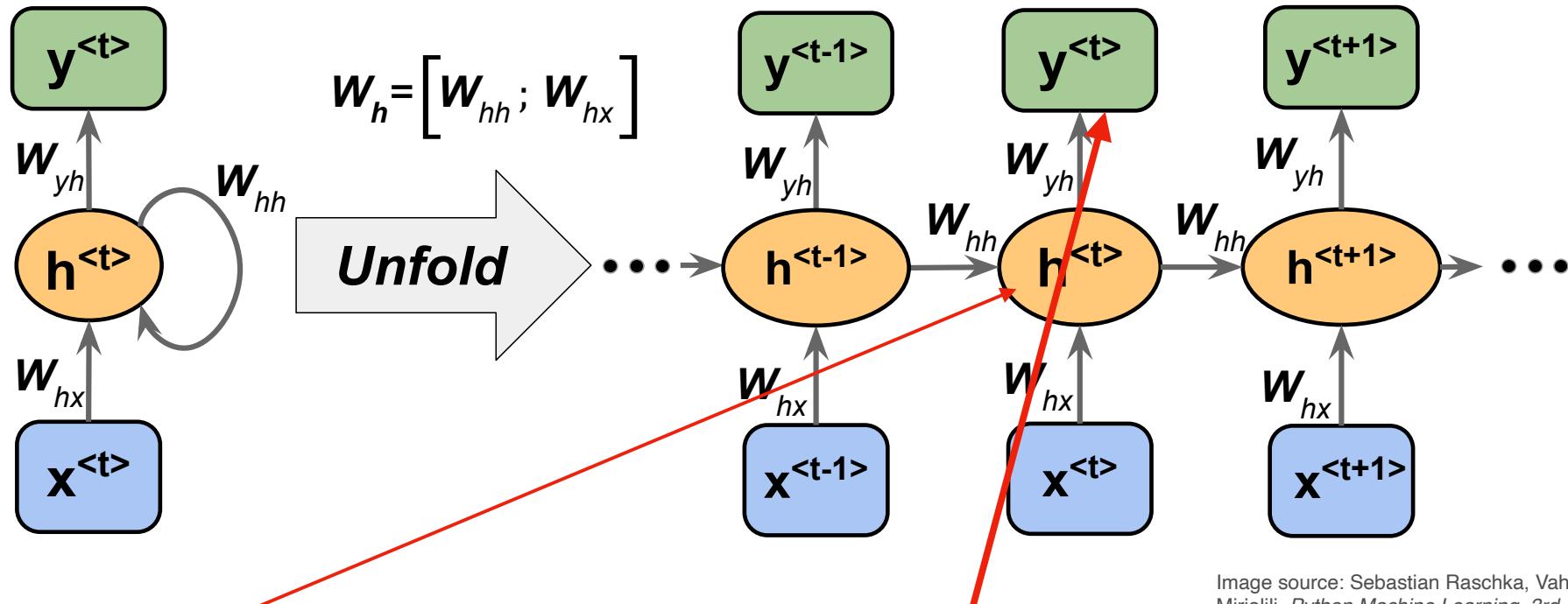
$$\mathbf{z}_h^{<t>} = \mathbf{W}_{hx} \mathbf{x}^{<t>} + \mathbf{W}_{hh} \mathbf{h}^{<t-1>} + \mathbf{b}_h$$

Activation:

$$\mathbf{h}^{<t>} = \sigma_h(\mathbf{z}_h^{<t>})$$

Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning, 3rd Edition*. Packt, 2019

Weight matrices in a single-hidden layer RNN



Net input:

$$\mathbf{z}_h^{(t)} = \mathbf{W}_{hx} \mathbf{x}^{(t)} + \mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{b}_h$$

Activation:

$$\mathbf{h}^{(t)} = \sigma_h(\mathbf{z}_h^{(t)})$$

Net input:

$$\mathbf{z}_y^{(t)} = \mathbf{W}_{yh} \mathbf{h}^{(t)} + \mathbf{b}_y$$

Output:

$$\mathbf{y}^{(t)} = \sigma_y(\mathbf{z}_y^{(t)})$$

Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Backpropagation through time

The overall loss can be computed as the sum over all time steps

Werbos, Paul J. "[Backpropagation through time: what it does and how to do it.](#)" *Proceedings of the IEEE* 78, no. 10 (1990): 1550-1560.

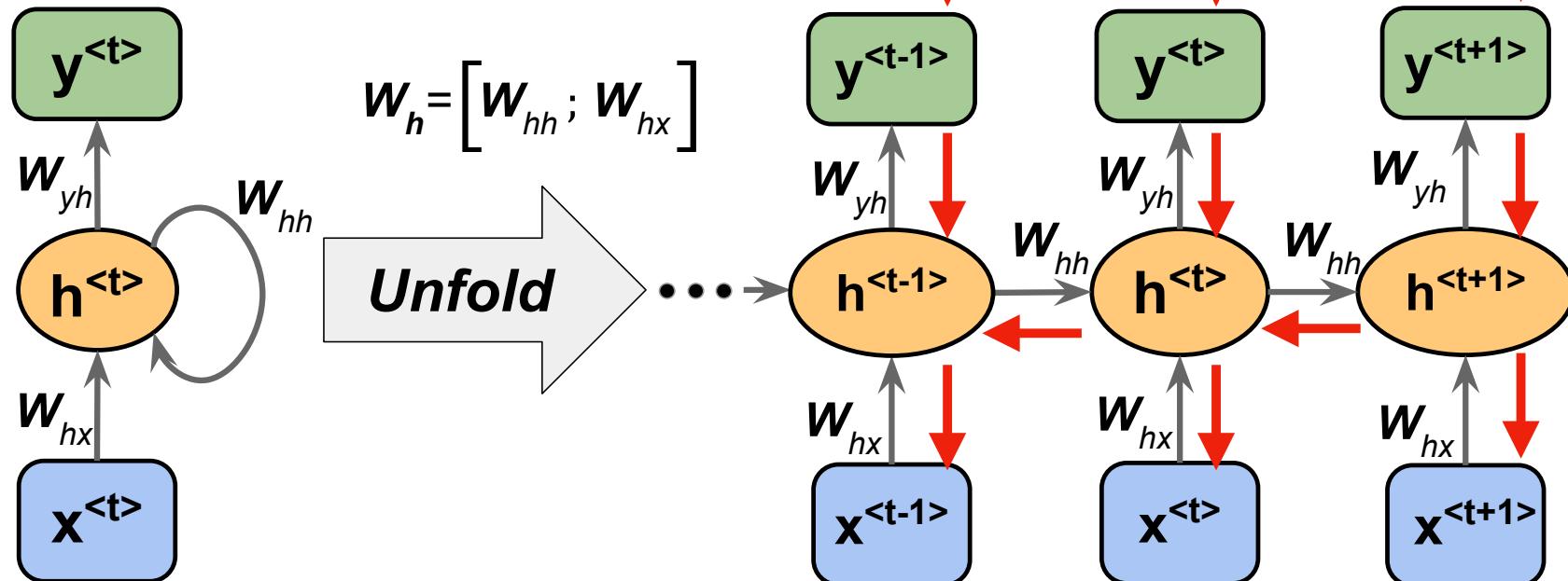
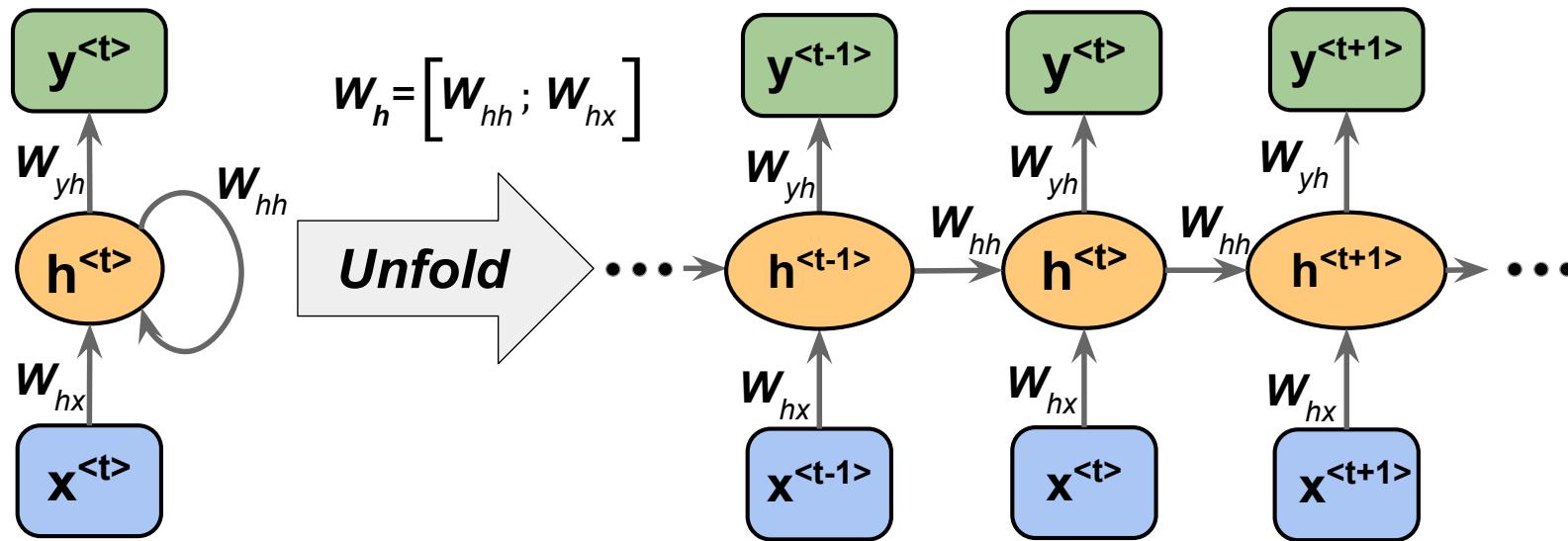


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Backpropagation through time

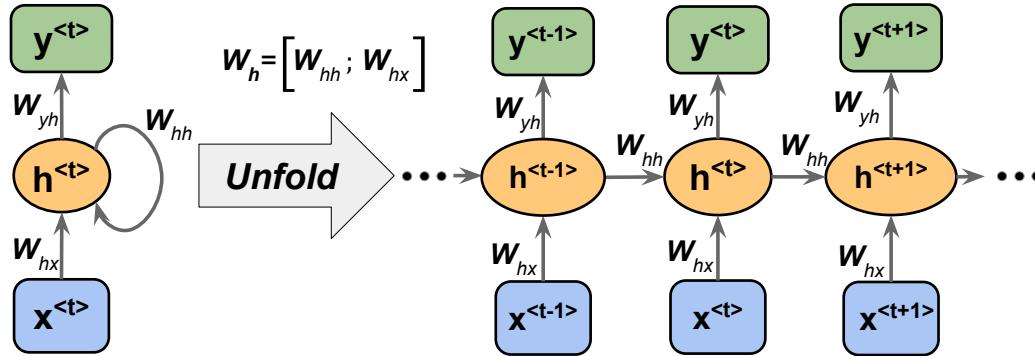


Werbos, Paul J. "[Backpropagation through time: what it does and how to do it.](#)" *Proceedings of the IEEE* 78, no. 10 (1990): 1550-1560.

$$L = \sum_{t=1}^T L^{(t)}$$

$$\frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}} \right)$$

Backpropagation through time



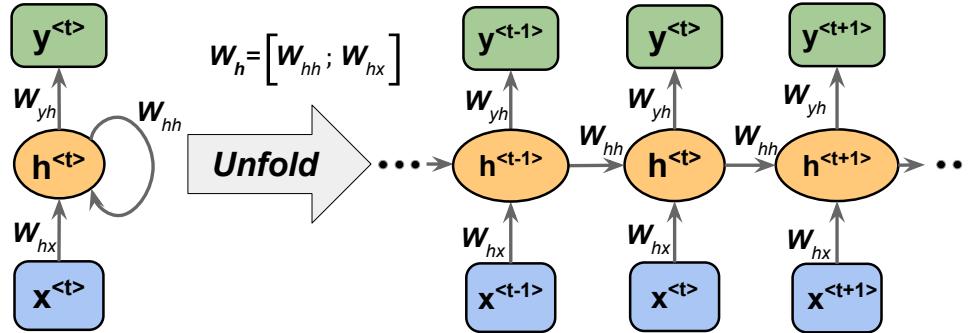
Werbos, Paul J. "[Backpropagation through time: what it does and how to do it.](#)" *Proceedings of the IEEE* 78, no. 10 (1990): 1550-1560.

$$\frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}} \right)$$

computed as a multiplication of adjacent time steps:

$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$$

Backpropagation through time



Werbos, Paul J. "[Backpropagation through time: what it does and how to do it.](#)" *Proceedings of the IEEE* 78, no. 10 (1990): 1550-1560.

$$L = \sum_{t=1}^T L^{(t)} \quad \frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \boxed{\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}} \right)$$

computed as a multiplication of adjacent time steps:

This is very problematic:
Vanishing/Exploding gradient problem!

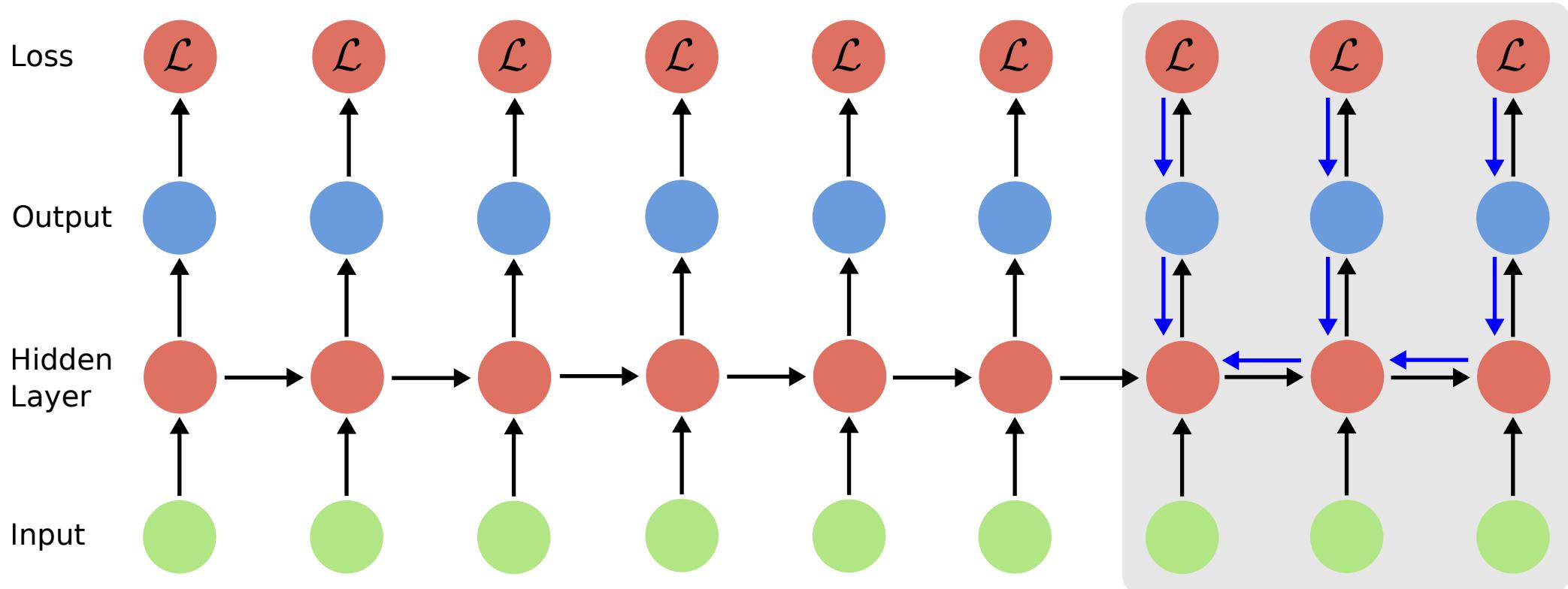
$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$$

Solutions to the vanishing/exploding gradient problems

- 1) Gradient Clipping: set a max value for gradients if they grow to large (solves only exploding gradient problem)
- 2) Truncated backpropagation through time (TBPTT)
simply limits the number of time steps the signal can backpropagate after each forward pass. E.g., even if the sequence has 100 elements/steps, we may only backpropagate through 20 or so
- 3) Long short-term memory (LSTM) -- uses a memory cell for modeling long-range dependencies and avoid vanishing gradient problems

Hochreiter, Sepp, and Jürgen Schmidhuber. "[Long short-term memory.](#)" *Neural computation* 9, no. 8 (1997): 1735-1780.

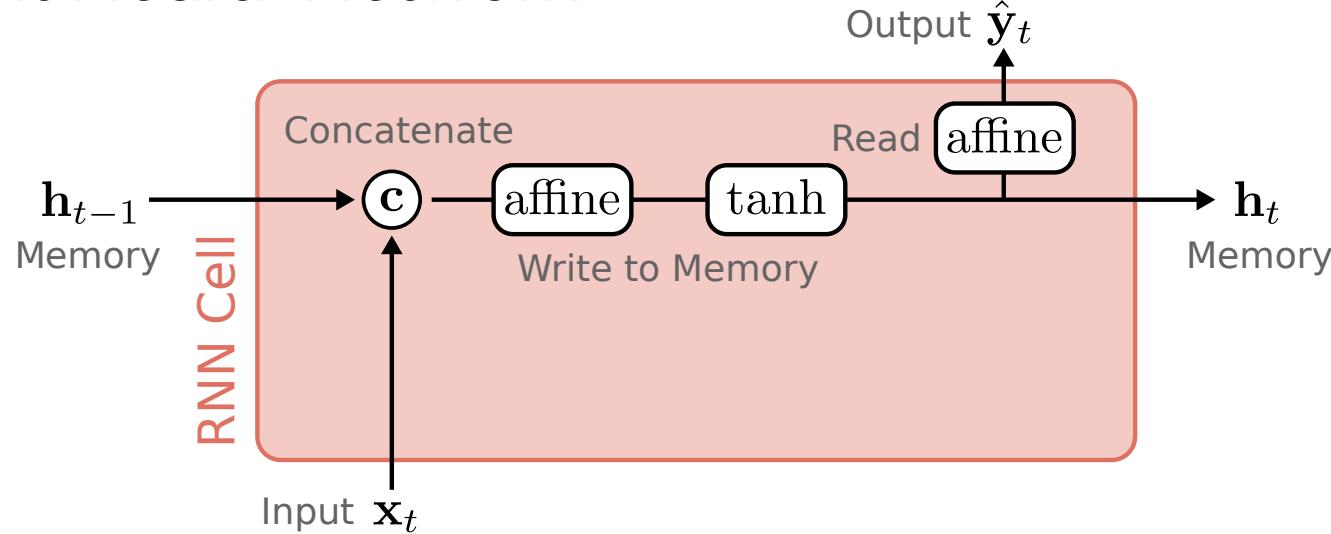
Truncated Backpropagation through Time



- ▶ Thus, one typically uses **truncated backpropagation through time** in practice
- ▶ Carry hidden states forward in time forever, but stop backpropagation earlier
- ▶ Total loss is sum of individual loss functions (= negative log likelihood)

Basic RNN Cell

Basic Recurrent Neural Network



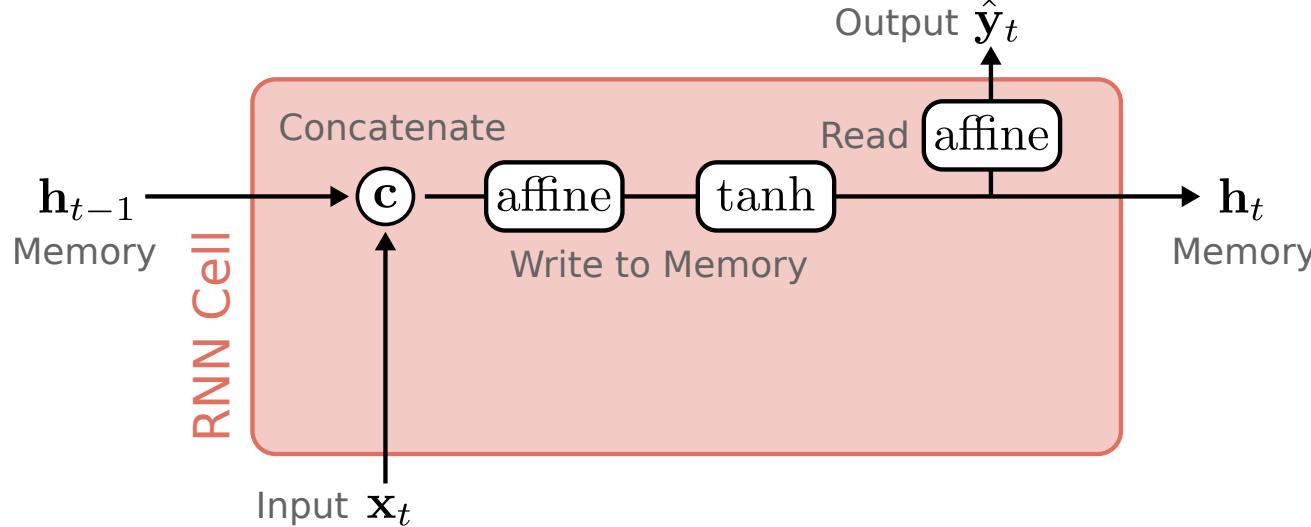
General Form:

$$\mathbf{h}_t = f_h(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

$$\hat{\mathbf{y}}_t = f_y(\mathbf{h}_t)$$

- We use t as the **time index** (in feedforward networks we used i as layer index)
- General form does not specify the form of the hidden and output mappings
- Important: f_h and f_y **do not change over time**, unlike in layers of feedforward net

Basic Recurrent Neural Network



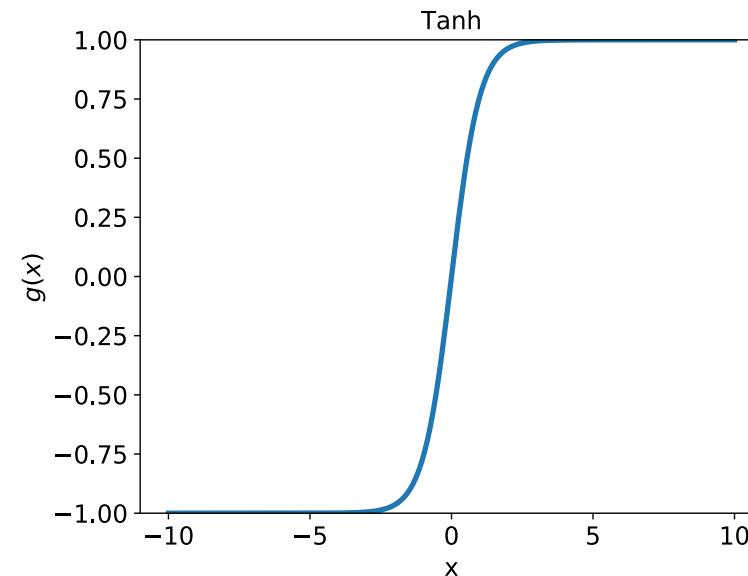
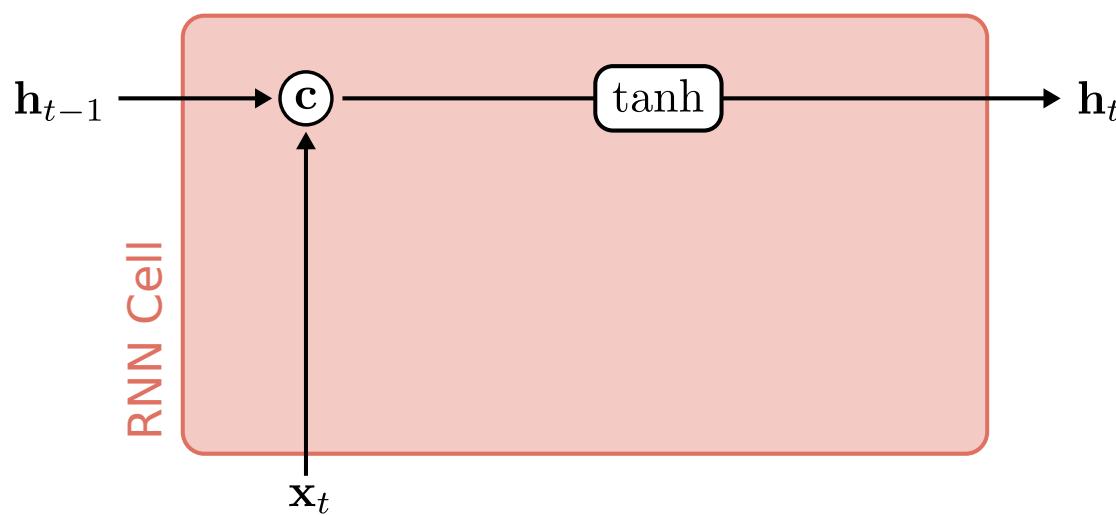
Single-Layer RNN:

$$\mathbf{h}_t = \tanh(\mathbf{A}_h \mathbf{h}_{t-1} + \mathbf{A}_x \mathbf{x}_t + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}_t = \mathbf{A}_y \mathbf{h}_t + \mathbf{b}_y$$

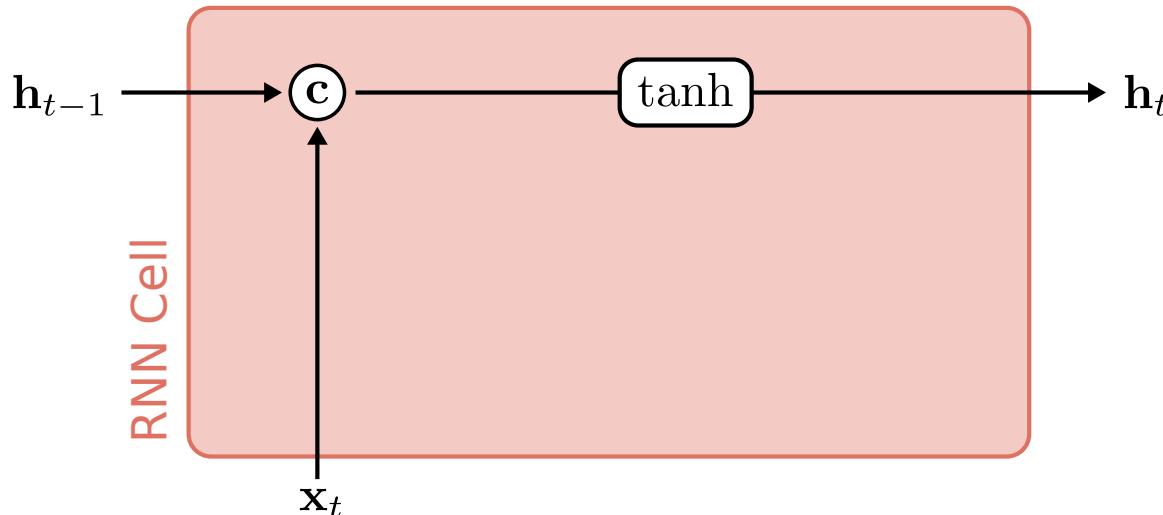
- ▶ Hidden state \mathbf{h}_t = linear combination of input \mathbf{x}_t and previous hidden state \mathbf{h}_{t-1}
- ▶ Output $\hat{\mathbf{y}}_t$ = linear prediction based on current hidden state \mathbf{h}_t
- ▶ $\tanh(\cdot)$ is commonly used as activation function (data is in the range $[-1, 1]$)
- ▶ Parameters $\mathbf{A}_h, \mathbf{A}_x, \mathbf{A}_y, \mathbf{b}_h, \mathbf{b}_y$ are constant over time

Recurrent Neural Network



- The state update H_t is modeled using a zero-centered $\tanh(\cdot)$
- $\tanh(\cdot)$ assumes that the processed data is in the range $[-1, 1]$
- Remark: we omit the affine transformations and the output layer for clarity

Vanishing / Exploding Gradients



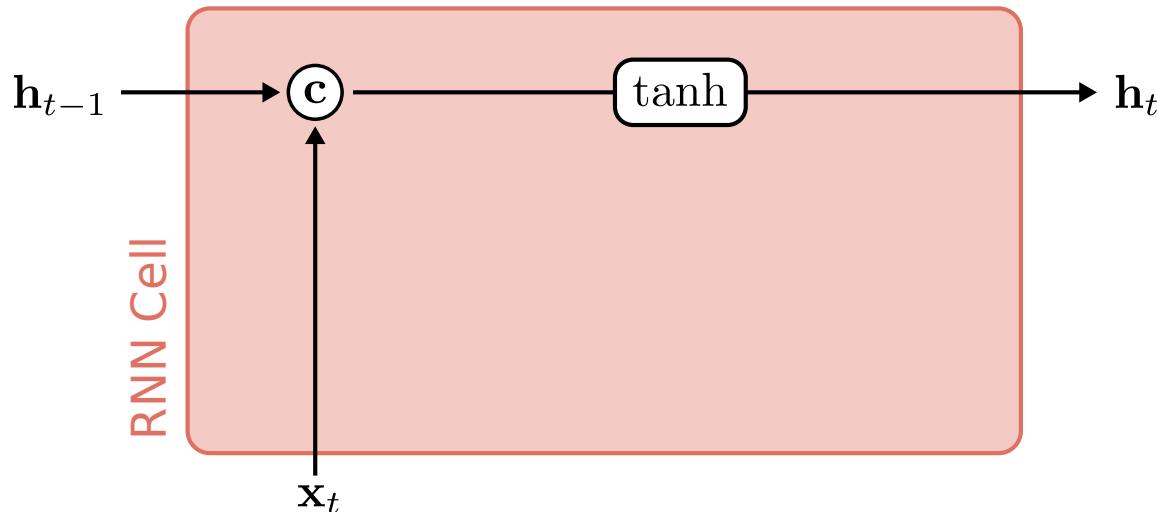
$$h_t = \tanh(a_h h_{t-1} + a_x x_t + b)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh' a_h \text{ with } \tanh' = \frac{\partial \tanh(x)}{\partial x}$$

What is the problem with vanilla RNNs?

- ▶ Let us consider an RNN with one dimensional hidden state $h_t \in \mathbb{R}$
- ▶ We have:
$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = \left(\prod_{i=t-k+1}^t \tanh'_i \right) a_h^k$$
- ▶ Thus, the gradient vanishes if $\tanh(\cdot)$ saturates as in feedforward networks
- ▶ RNNs require careful initialization to avoid saturating activation functions

Vanishing / Exploding Gradients



$$h_t = \tanh(a_h h_{t-1} + a_x x_t + b)$$

$$\approx a_h h_{t-1} + a_x x_t + b$$

$$\frac{\partial h_t}{\partial h_{t-1}} \approx a_h$$

However, gradients might still misbehave:

- ▶ Let us now assume that the RNN has been initialized well such that the activation functions are not saturated $\Rightarrow a_h h_{t-1} + a_x x_t + b \in [-1, 1] \Rightarrow \tanh(x) \approx x$
- ▶ We now have:

$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = a_h^k$$

Vanishing / Exploding Gradients

$$\frac{\partial h_t}{\partial h_{t-k}} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{t-k+1}}{\partial h_{t-k}} = a_h^k$$

For $a_h > 1$ gradients will **explode** (become very large, cause divergence):

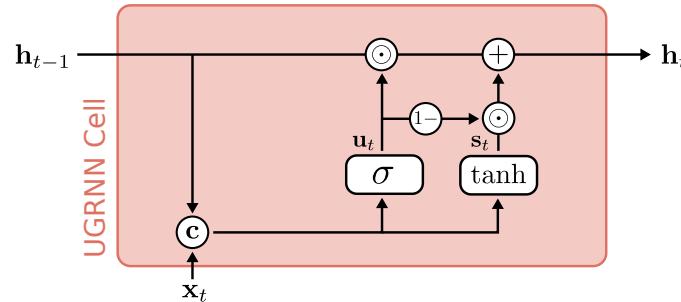
- ▶ Example: For $a_h = 1.1$ and $k = 100$ we have $\partial h_t / \partial h_{t-k} = a_h^k = 13781$
- ▶ This problem is often addressed in practice using **gradient clipping**
- ▶ Forward values do not explode due to bounded $\tanh(\cdot)$ activation function

For $a_h < 1$ gradients will **vanish** (no learning in earlier time steps):

- ▶ Example: For $a_h = 0.9$ and $k = 100$ we have $\partial h_t / \partial h_{t-k} = a_h^k = 0.0000266$
- ▶ Avoiding this problem requires an **architectural change**
- ▶ But residual connections do not work here as the parameters are shared across time and the input and desired output at each time step are different

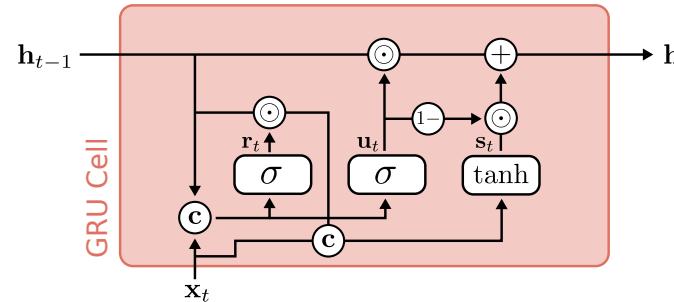
Gated Recurrent Networks

UGRNN



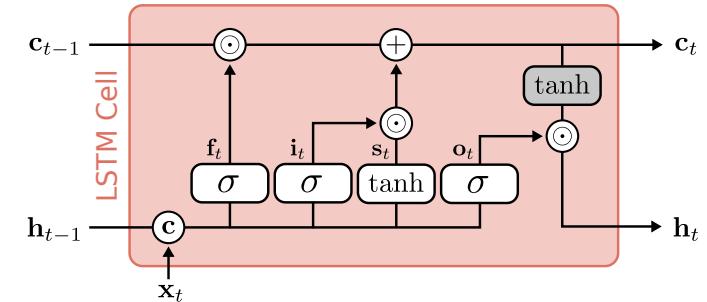
Collins, 2017

GRU



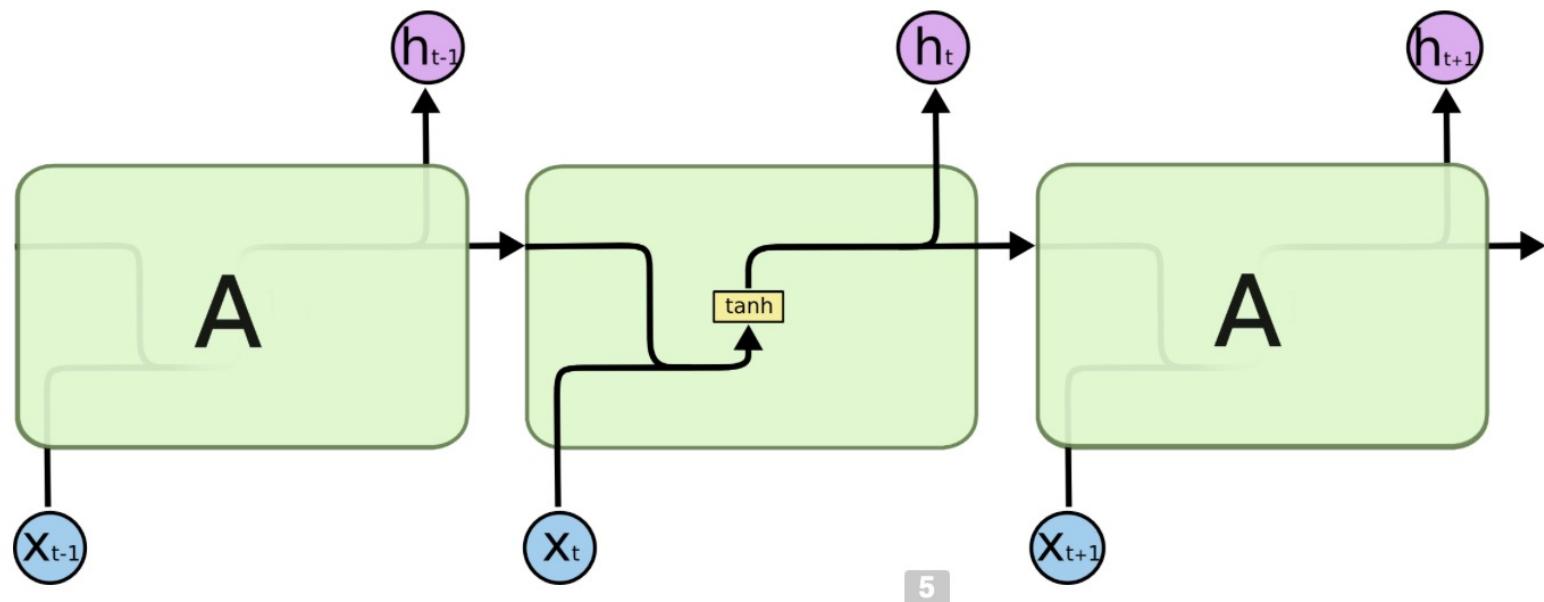
Cho, 2014

LSTM

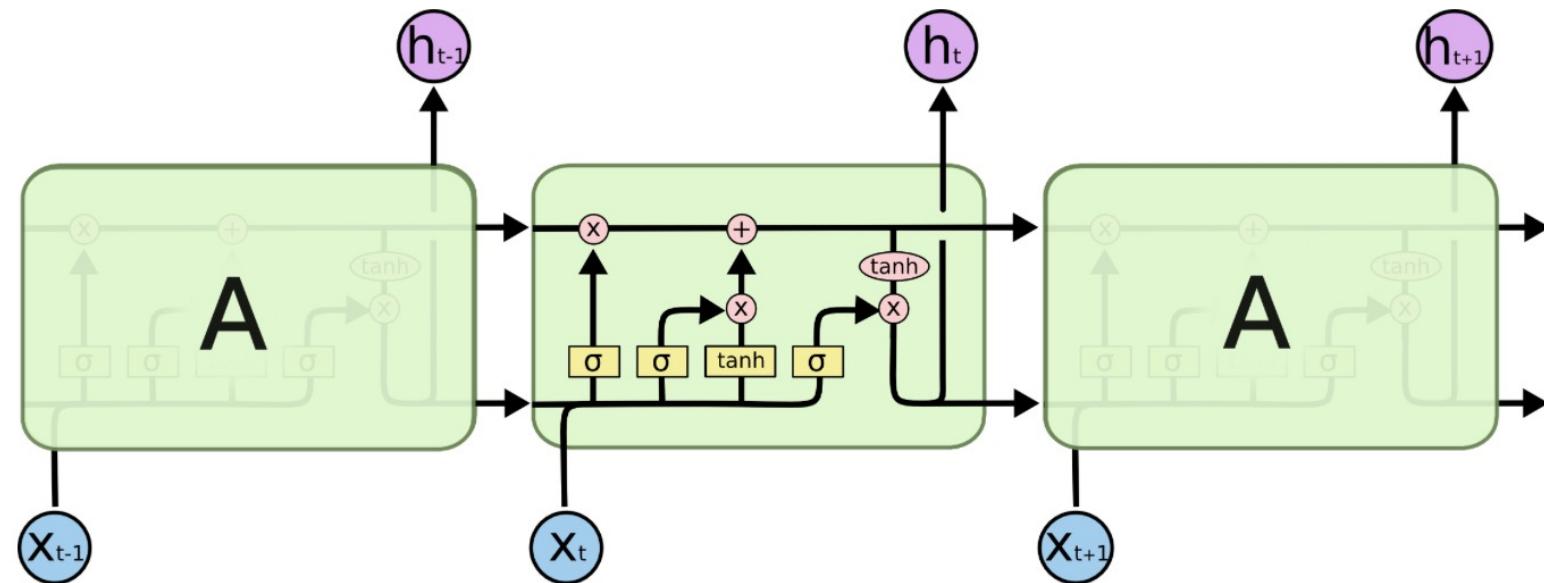


Hochreiter, 1997

- ▶ **UGRNN:** Update Gate Recurrent Neural Network
- ▶ **GRU:** Gated Recurrent Unit
- ▶ **LSTM:** Long Short-Term Memory
- ▶ LSTM was the first and most transformative (revolutionized NLP in 2015, e.g. at Google), but also most complex model. UGRNN and GRU work similarly well.
- ▶ Common to all architectures: **gates** for filtering information



5



Long-short term memory (LSTM)

LSTM cell:

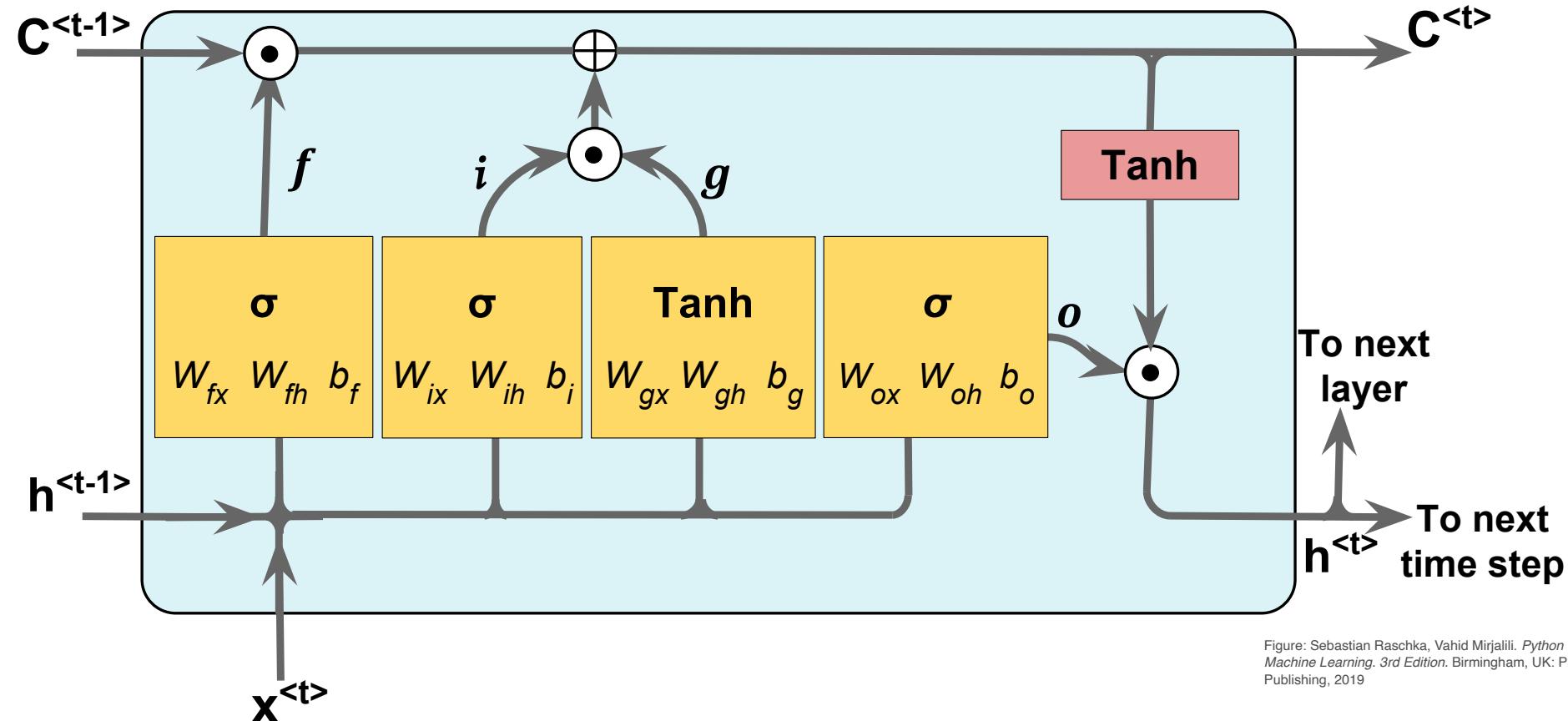
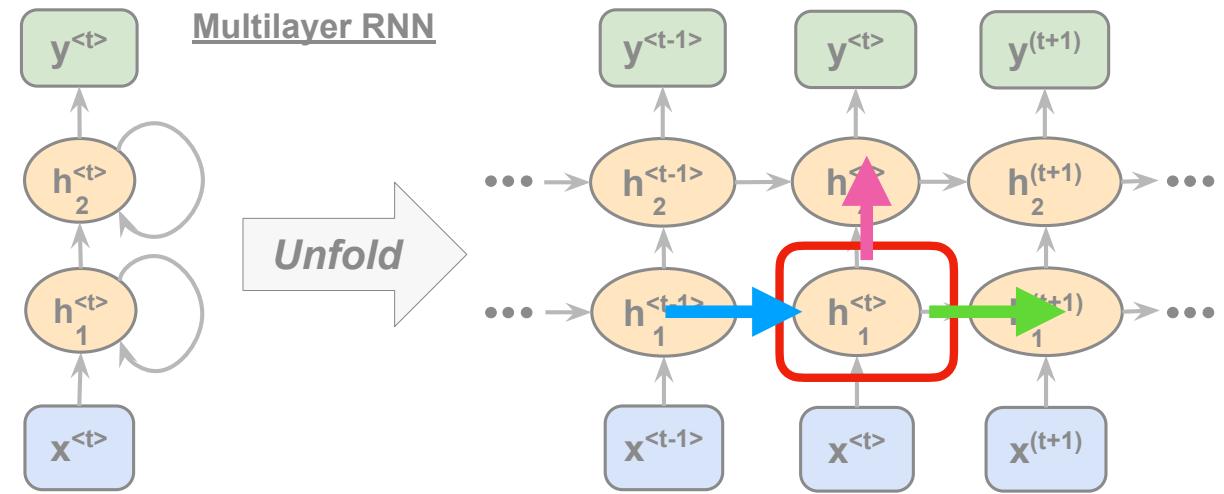
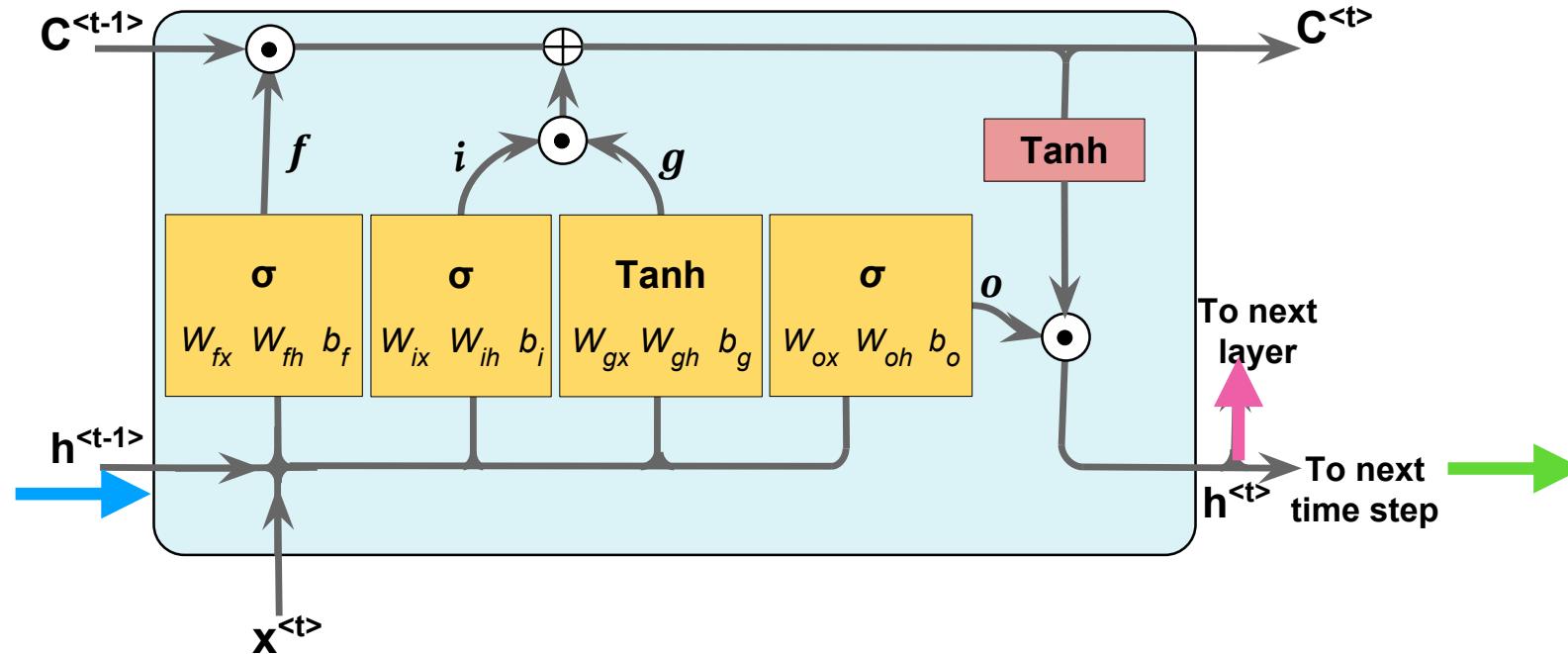


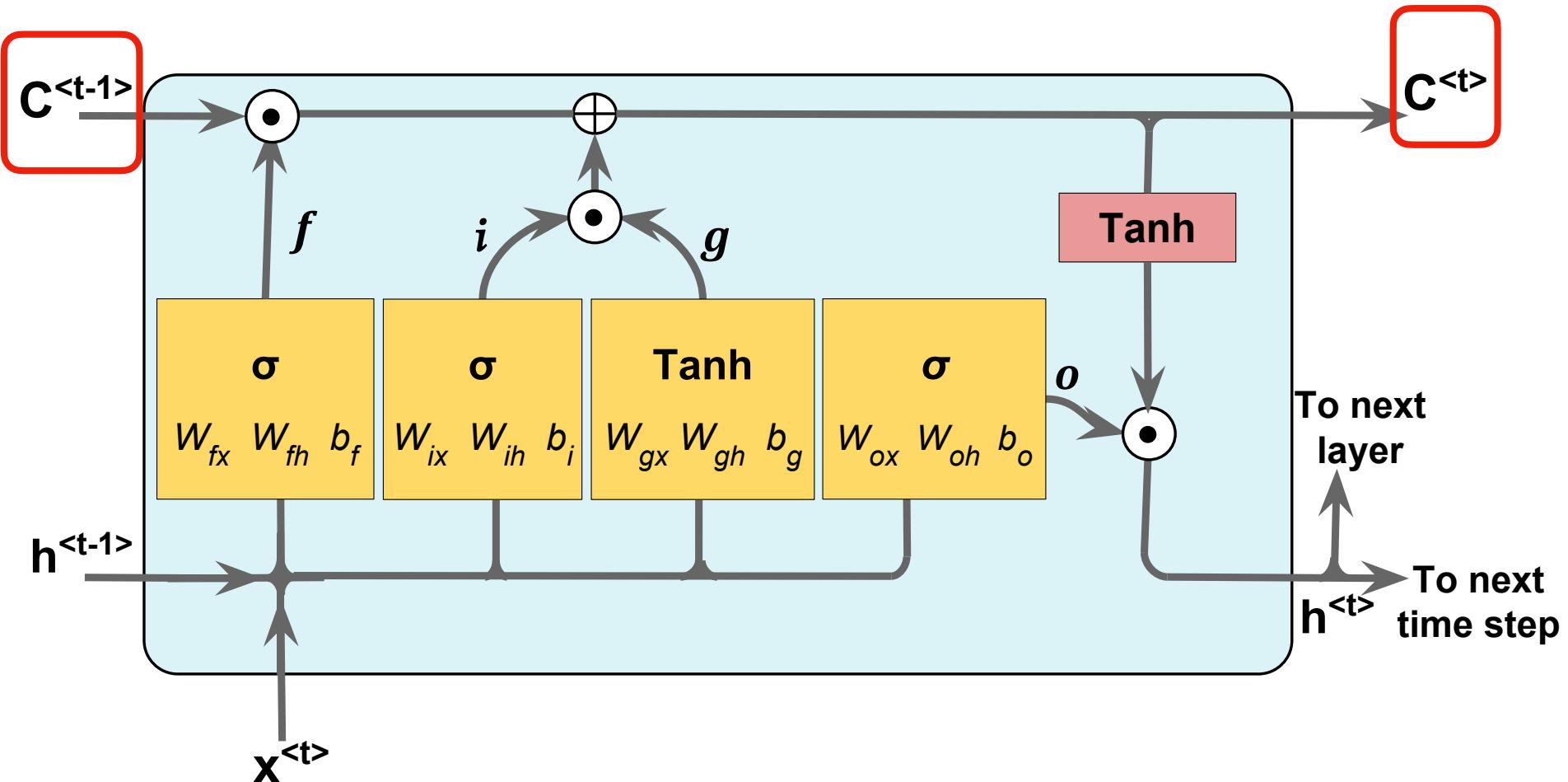
Figure: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Birmingham, UK: Packt Publishing, 2019



Long-short term memory (LSTM)

Cell state at previous time step

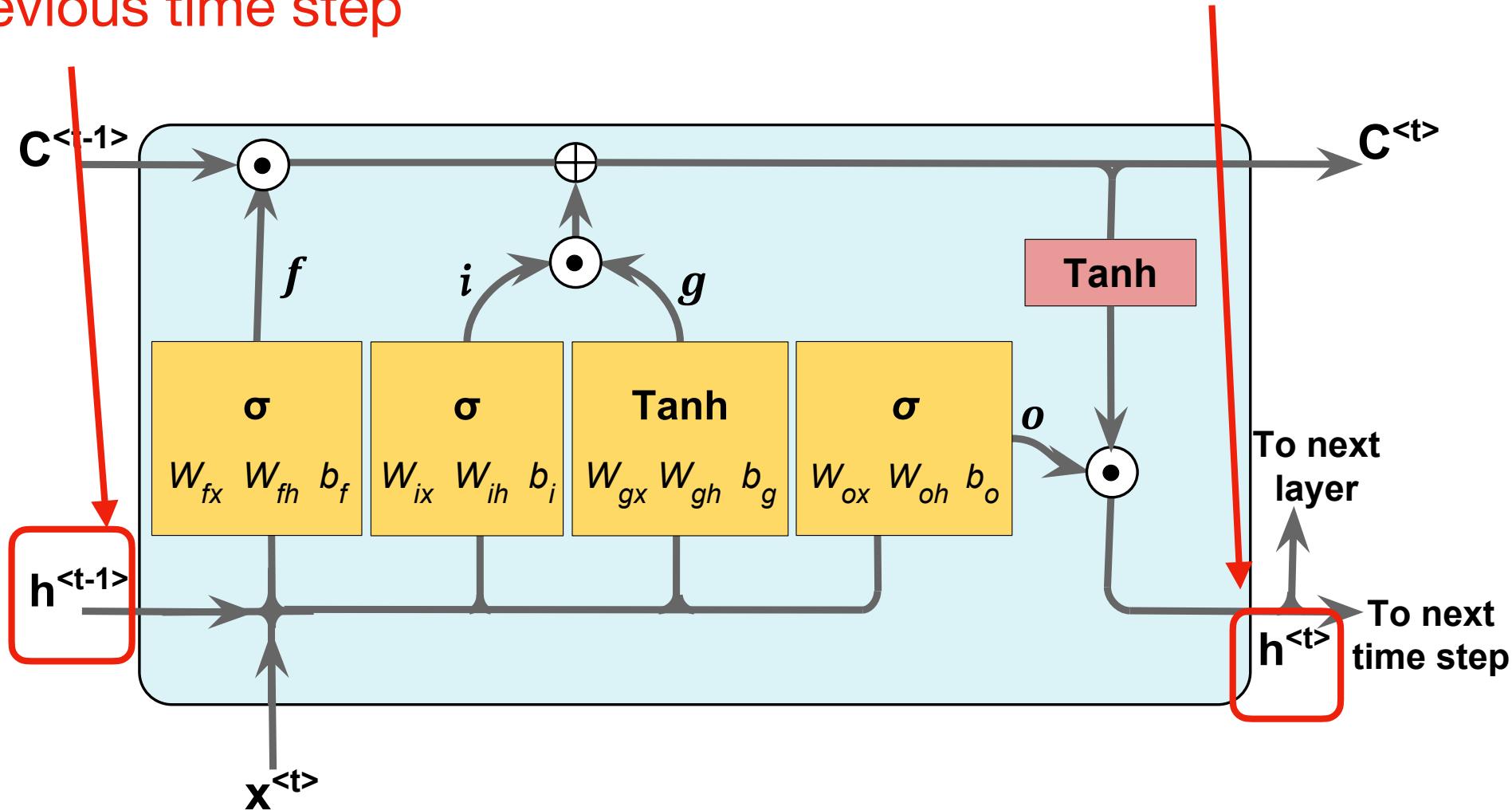
Cell state at current time step



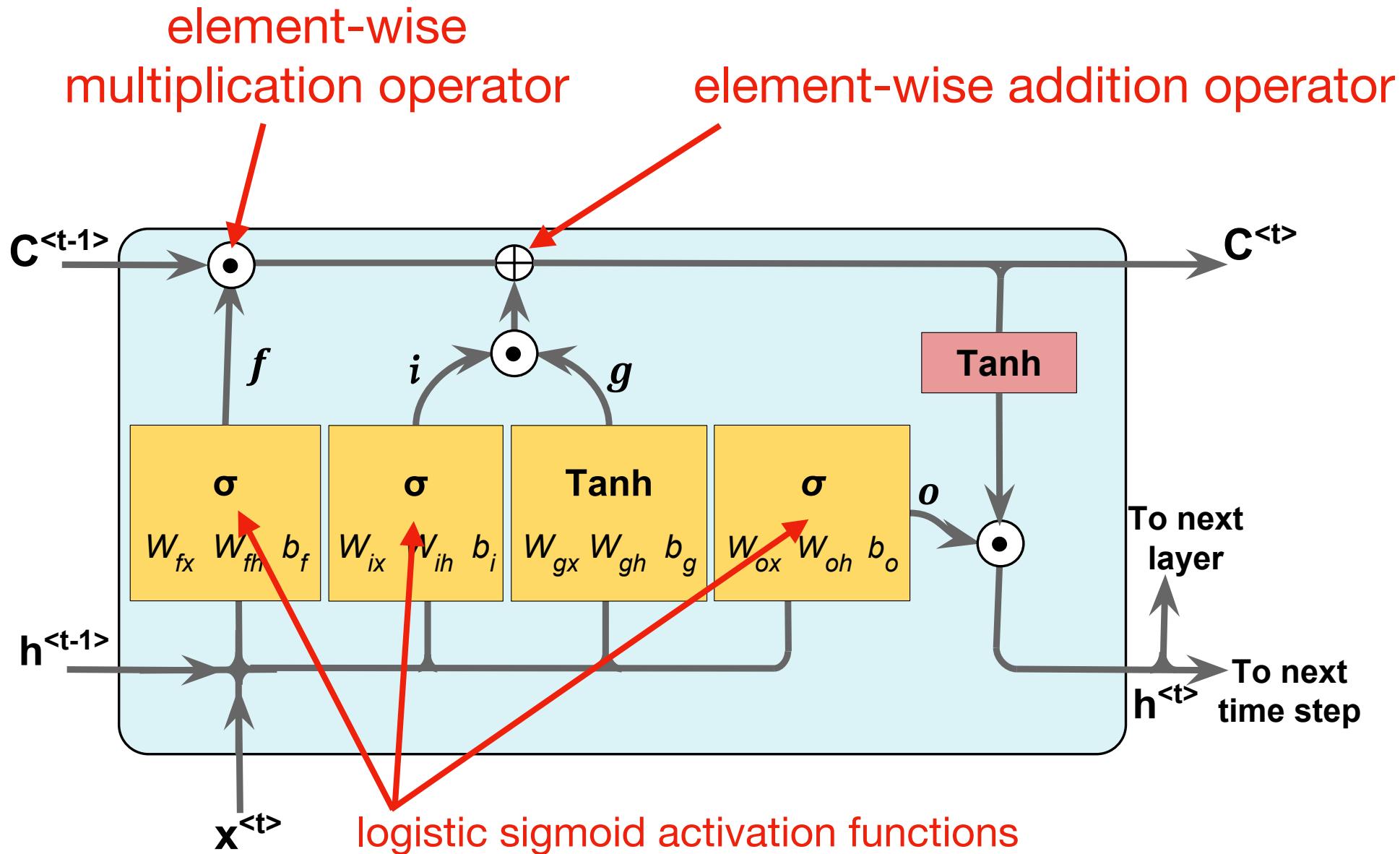
Long-short term memory (LSTM)

activation from
previous time step

activation for next time step



Long-short term memory (LSTM)

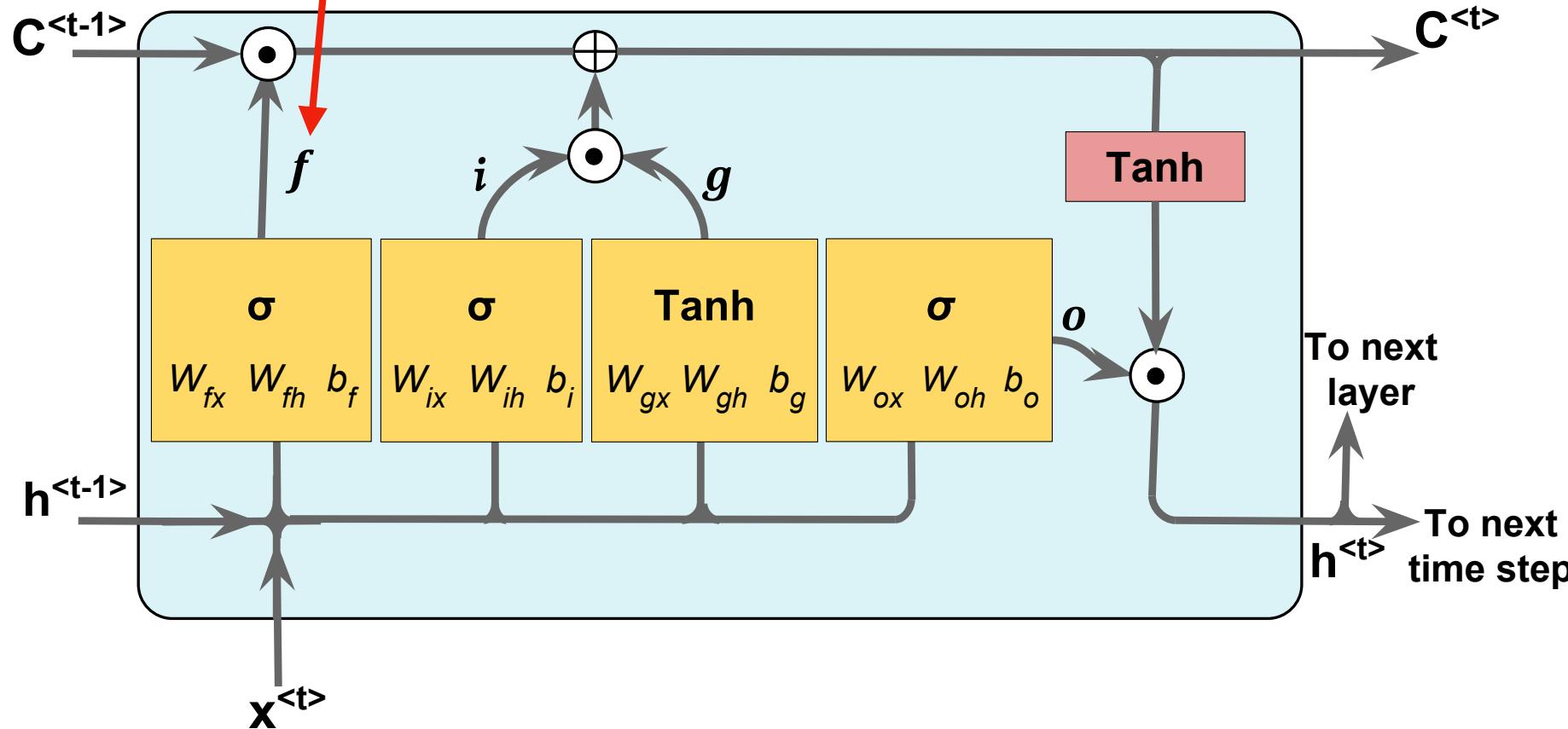


Long-short term memory (LSTM)

Gers, Felix A., Jürgen Schmidhuber, and Fred Cummins. "[Learning to forget: Continual prediction with LSTM](#)." (1999): 850-855.

"Forget Gate": controls which information is remembered, and which is forgotten;
can reset the cell state

$$f_t = \sigma \left(\mathbf{W}_{fx} \mathbf{x}^{(t)} + \mathbf{W}_{fh} \mathbf{h}^{(t-1)} + \mathbf{b}_f \right)$$

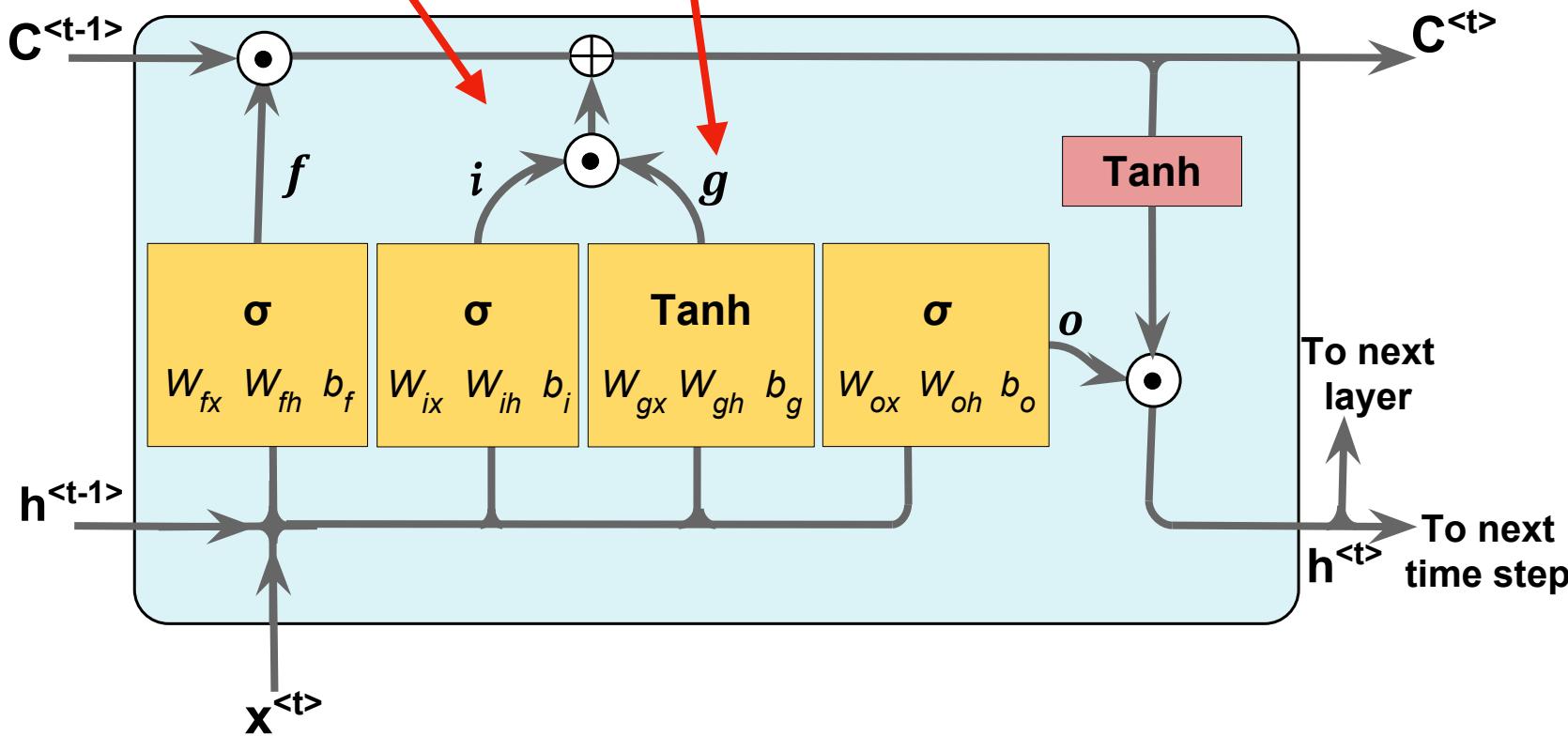


Long-short term memory (LSTM)

"Input Gate": $i_t = \sigma \left(\mathbf{W}_{ix} \mathbf{x}^{(t)} + \mathbf{W}_{ih} \mathbf{h}^{(t-1)} + \mathbf{b}_i \right)$

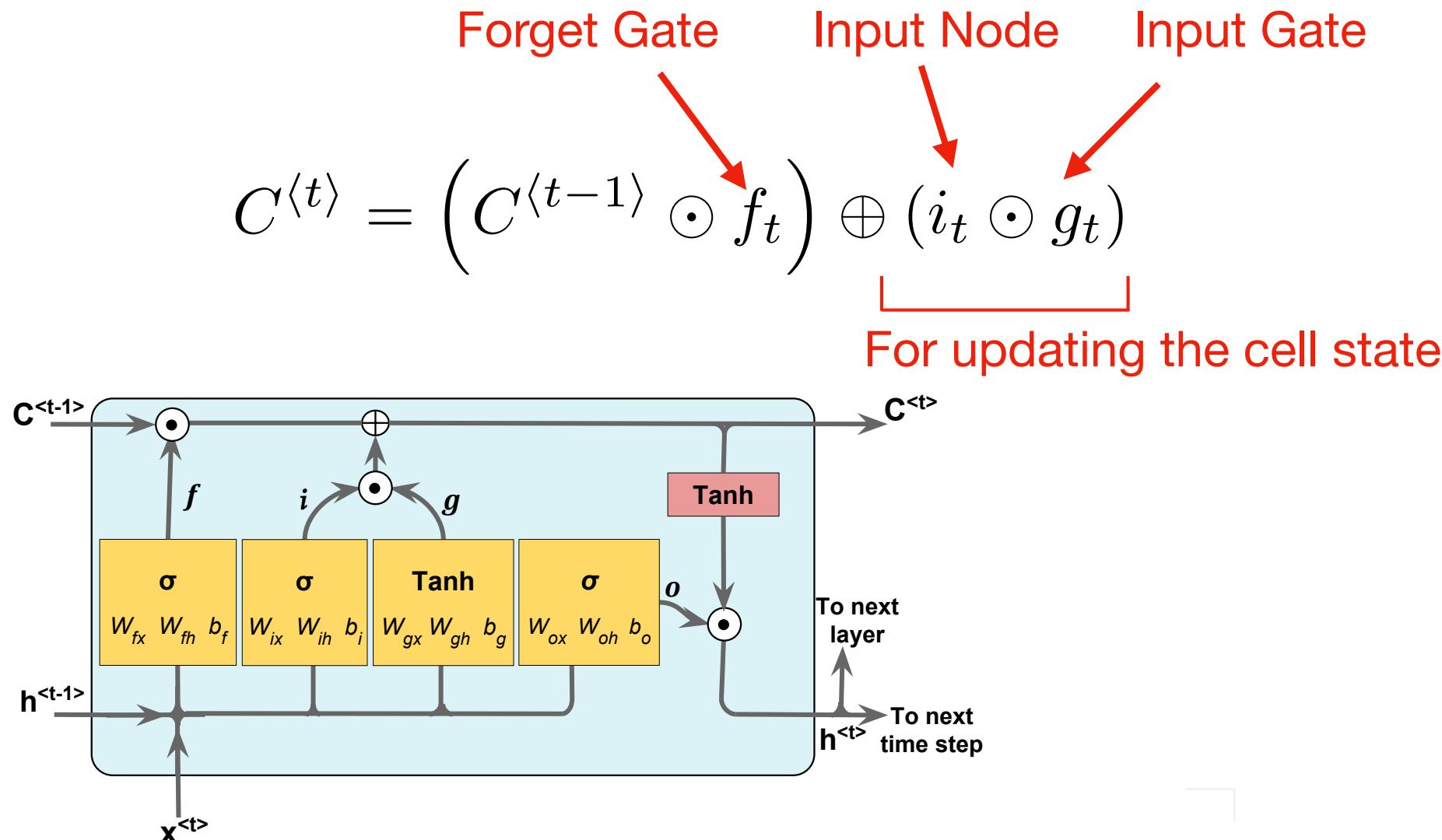
"Input Node":

$$\mathbf{g}_t = \tanh \left(\mathbf{W}_{gx} \mathbf{x}^{(t)} + \mathbf{W}_{gh} \mathbf{h}^{(t-1)} + \mathbf{b}_g \right)$$



Long-short term memory (LSTM)

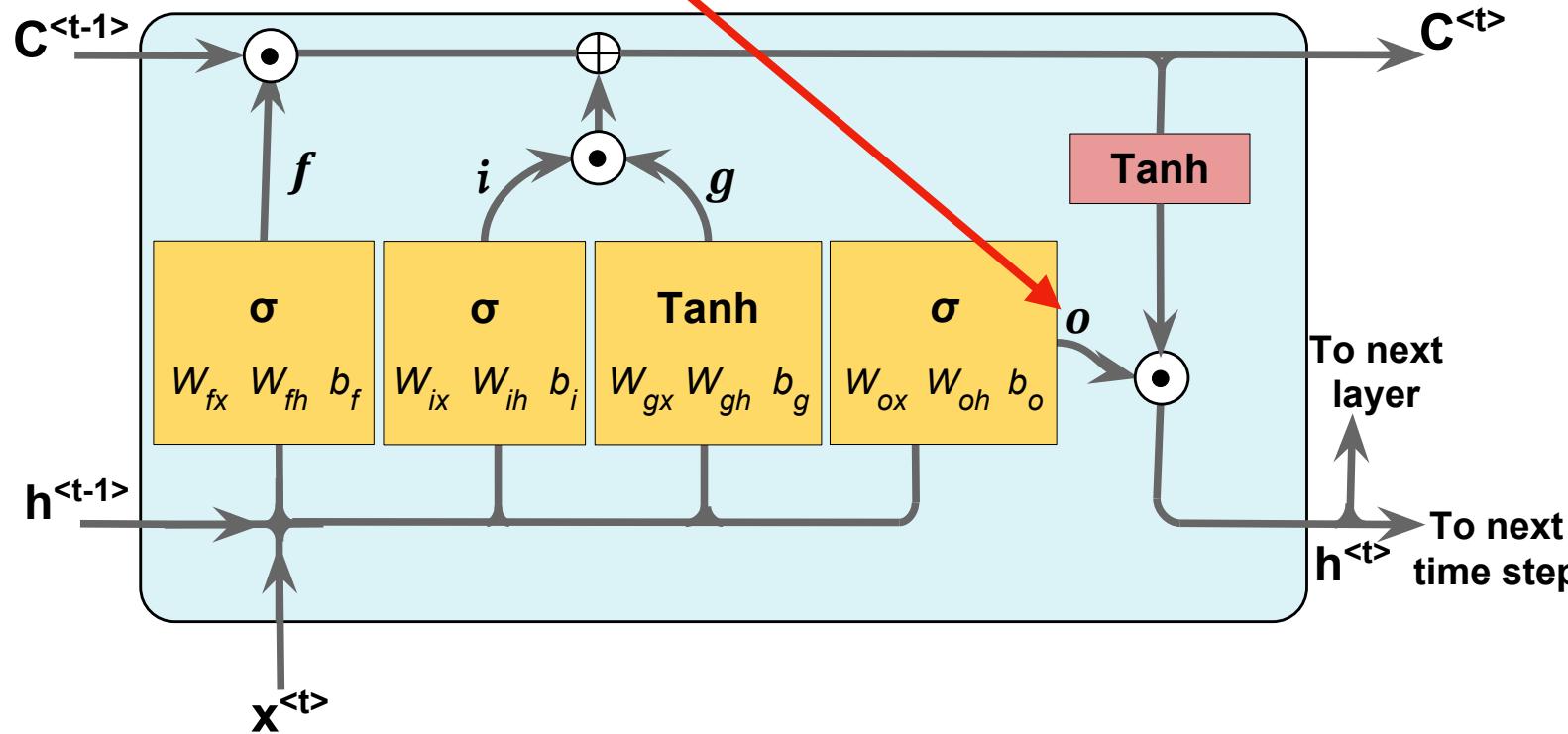
Brief summary of the gates so far ...



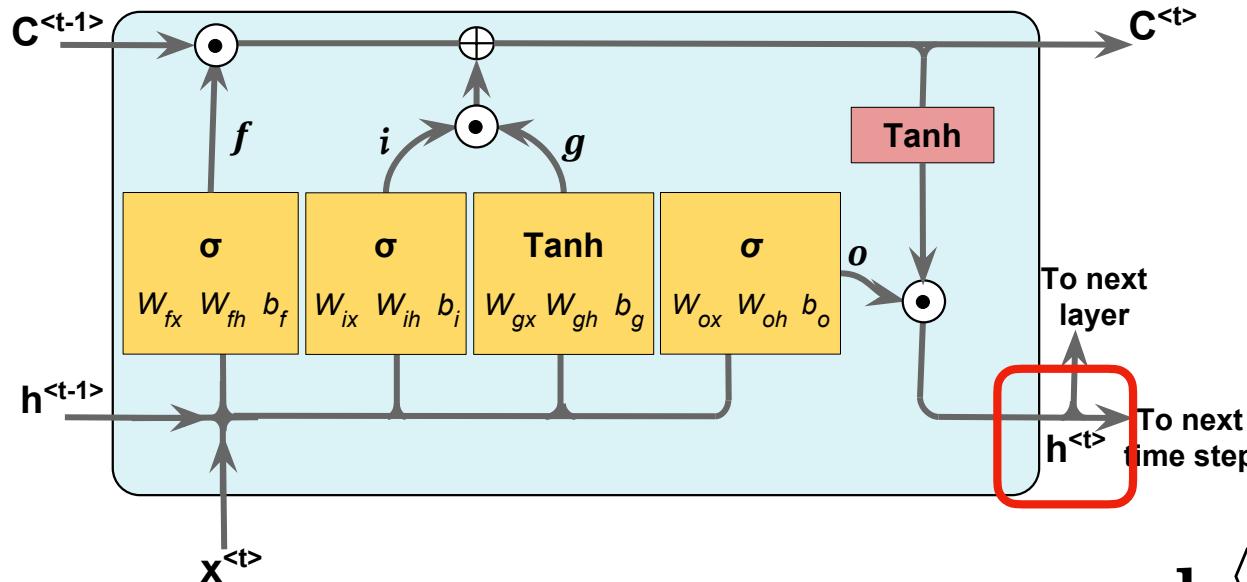
Long-short term memory (LSTM)

Output gate for updating the values of hidden units:

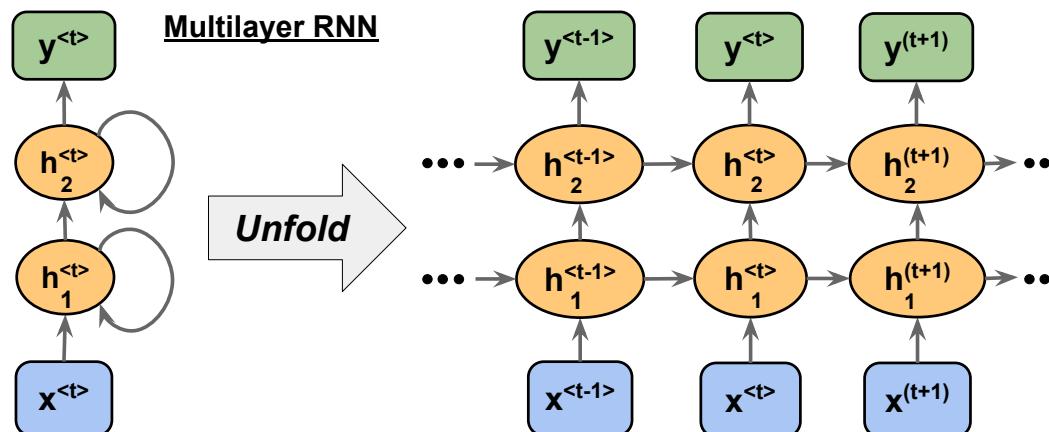
$$\mathbf{o}_t = \sigma \left(\mathbf{W}_{ox} \mathbf{x}^{(t)} + \mathbf{W}_{oh} \mathbf{h}^{(t-1)} + \mathbf{b}_o \right)$$



Long-short term memory (LSTM)

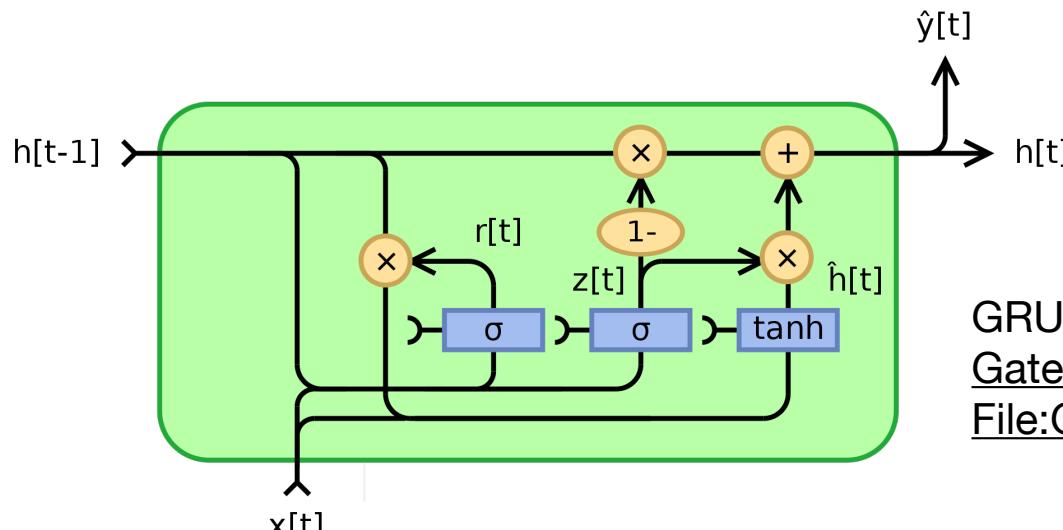


$$h^{(t)} = o_t \odot \tanh(C^{(t)})$$



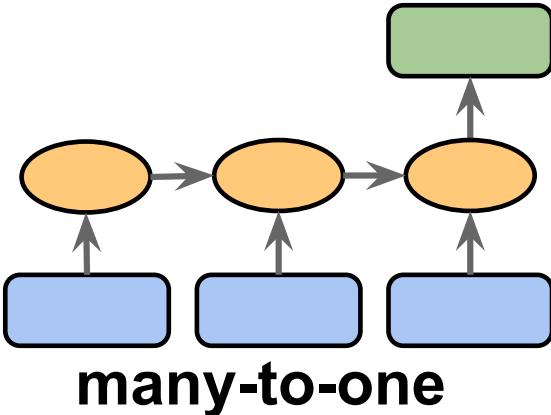
Long-short term memory (LSTM)

- Still popular and widely used today
- A recent, related approach is the Gated Recurrent Unit (GRU)
Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. "[Learning phrase representations using RNN encoder-decoder for statistical machine translation.](#)" *arXiv preprint arXiv:1406.1078* (2014).
- Nice article exploring LSTMs and comparing them to GRUs
Jozefowicz, Rafal, Wojciech Zaremba, and Ilya Sutskever.
["An empirical exploration of recurrent network architectures."](#) In *International Conference on Machine Learning*, pp. 2342-2350. 2015.

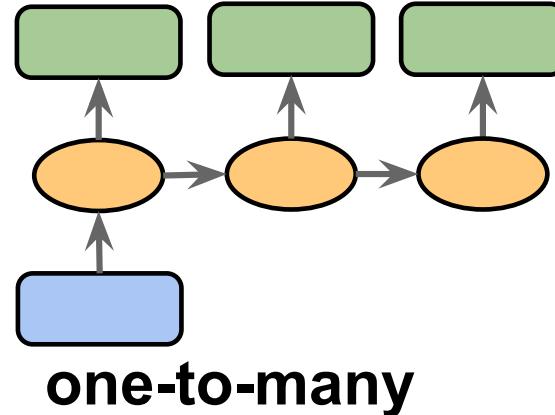


GRU image source: https://en.wikipedia.org/wiki/Gated_recurrent_unit#/media/File:Gated_Recurrent_Unit,_base_type.svg

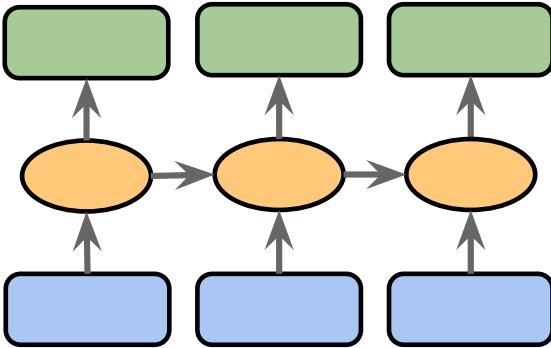
Different Types of Sequence Modeling Tasks



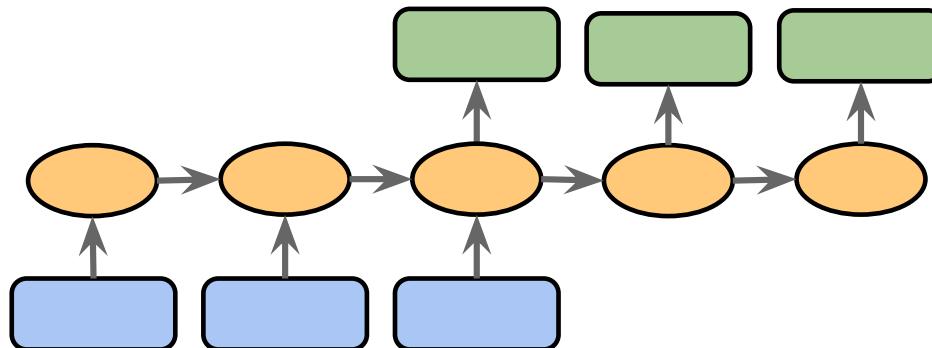
many-to-one



one-to-many



many-to-many



many-to-many

A Classic Approach for Text Classification: Bag-of-Words Model

"Raw" training dataset

$x^{[1]} = \text{"The sun is shining"}$

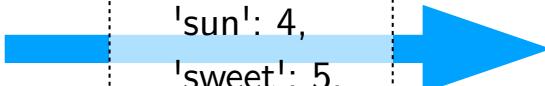
$x^{[2]} = \text{"The weather is sweet"}$

$x^{[3]} = \text{"The sun is shining,}$
 $\text{the weather is sweet, and}$
 $\text{one and one is two"}$

$$\mathbf{y} = [0, 1, 0]$$

class labels

```
vocabulary = {  
    'and': 0,  
    'is': 1  
    'one': 2,  
    'shining': 3,  
    'sun': 4,  
    'sweet': 5,  
    'the': 6,  
    'two': 7,  
    'weather': 8,  
}
```



Training set as design matrix

$$\mathbf{X} = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 1 \\ 2 & 3 & 2 & 1 & 1 & 1 & 2 & 1 & 1 \end{bmatrix}$$

$$\mathbf{y} = [0, 1, 0]$$

class labels

training

Classifier

(e.g., logistic regression, MLP, ...)

RNN Step 1): Building the Vocabulary

"Raw" training dataset

$x^{[1]}$ = "The sun is shining"

$x^{[2]}$ = "The weather is sweet"

$x^{[3]}$ = "The sun is shining,
the weather is sweet, and
one and one is two"

$y = [0, 1, 0]$

class labels



```
vocabulary = {  
    '<unk>': 0,  
    'and': 1,  
    'is': 2  
    'one': 3,  
    'shining': 4,  
    'sun': 5,  
    'sweet': 6,  
    'the': 7,  
    'two': 8,  
    'weather': 9,  
    '<pad>': 10  
}
```

RNN Step 2): Training Example Texts to Indices

"Raw" training dataset

$x^{[1]} = \text{"The sun is shining"}$

$x^{[2]} = \text{"The weather is sweet"}$

$x^{[3]} = \text{"The sun is shining,}$
the weather is sweet, and
one and one is two"

```
vocabulary = {  
    '<unk>': 0,  
    'and': 1,  
    'is': 2  
    'one': 3,  
    'shining': 4,  
    'sun': 5,  
    'sweet': 6,  
    'the': 7,  
    'two': 8,  
    'weather': 9,  
    '<pad>': 10  
}
```

$x^{[1]} = \text{"The sun is shining"}$

[7 5 2 4 ... 10 10 10]

$x^{[2]} = \text{"The weather is sweet"}$

[7 9 2 6 ... 10 10 10]

$x^{[3]} = \text{"The sun is shining,}$
the weather is sweet, and
one and one is two"

[7 5 2 4 ... 3 2 8]

RNN Step 3): Indices to One-Hot Representation

"Raw" training dataset

$x^{[1]} = \text{"The sun is shining"}$

$x^{[2]} = \text{"The weather is sweet"}$

$x^{[3]} = \text{"The sun is shining,}$
 $\text{the weather is sweet, and}$
 $\text{one and one is two"}$

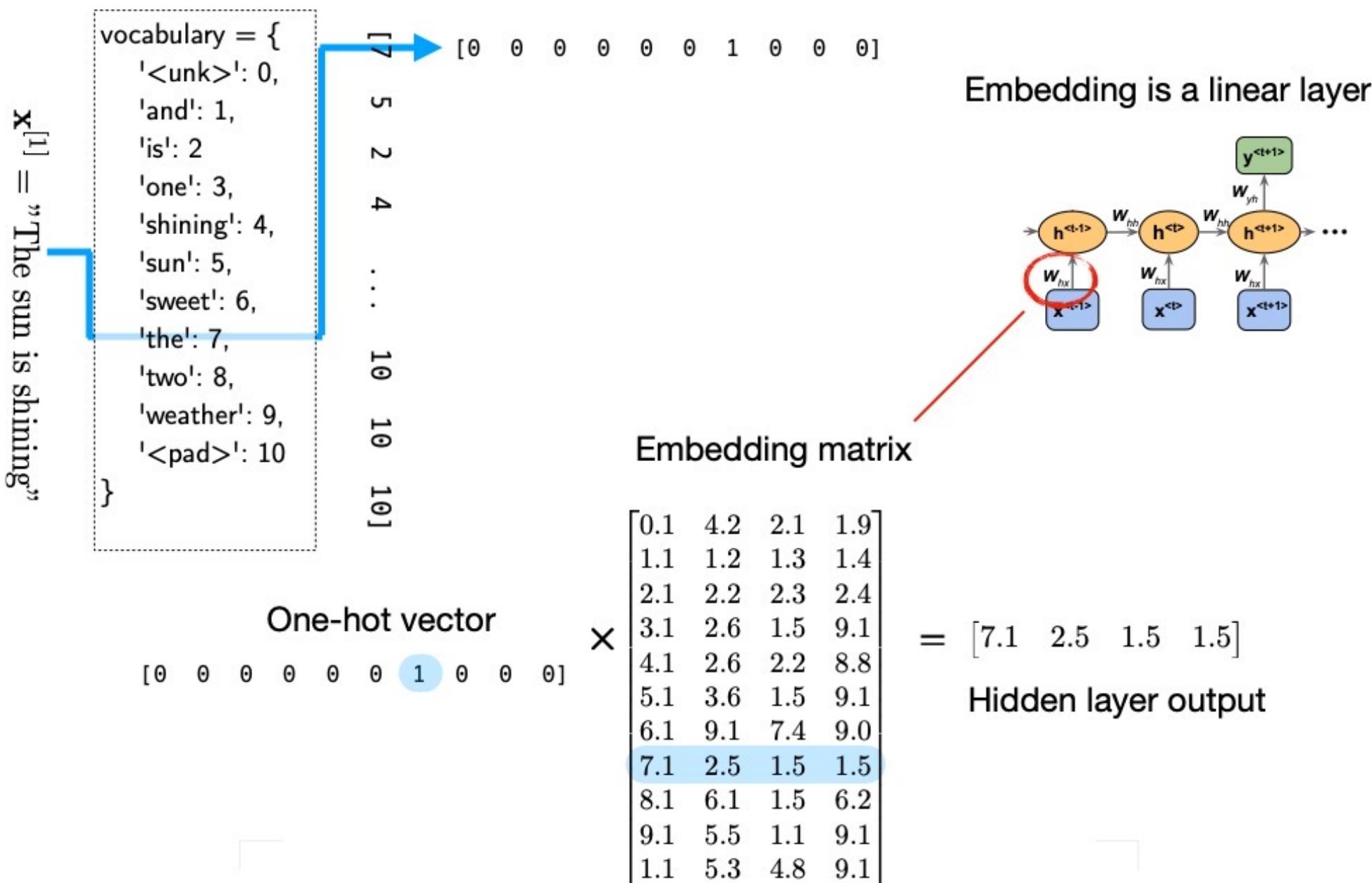
$x^{[1]} = \text{"The sun is shining"}$

```
vocabulary = {  
    '<unk>': 0,  
    'and': 1,  
    'is': 2  
    'one': 3,  
    'shining': 4,  
    'sun': 5,  
    'sweet': 6,  
    'the': 7,  
    'two': 8,  
    'weather': 9,  
    '<pad>': 10  
}
```

[7 5 2 4 . . 10 10 10]

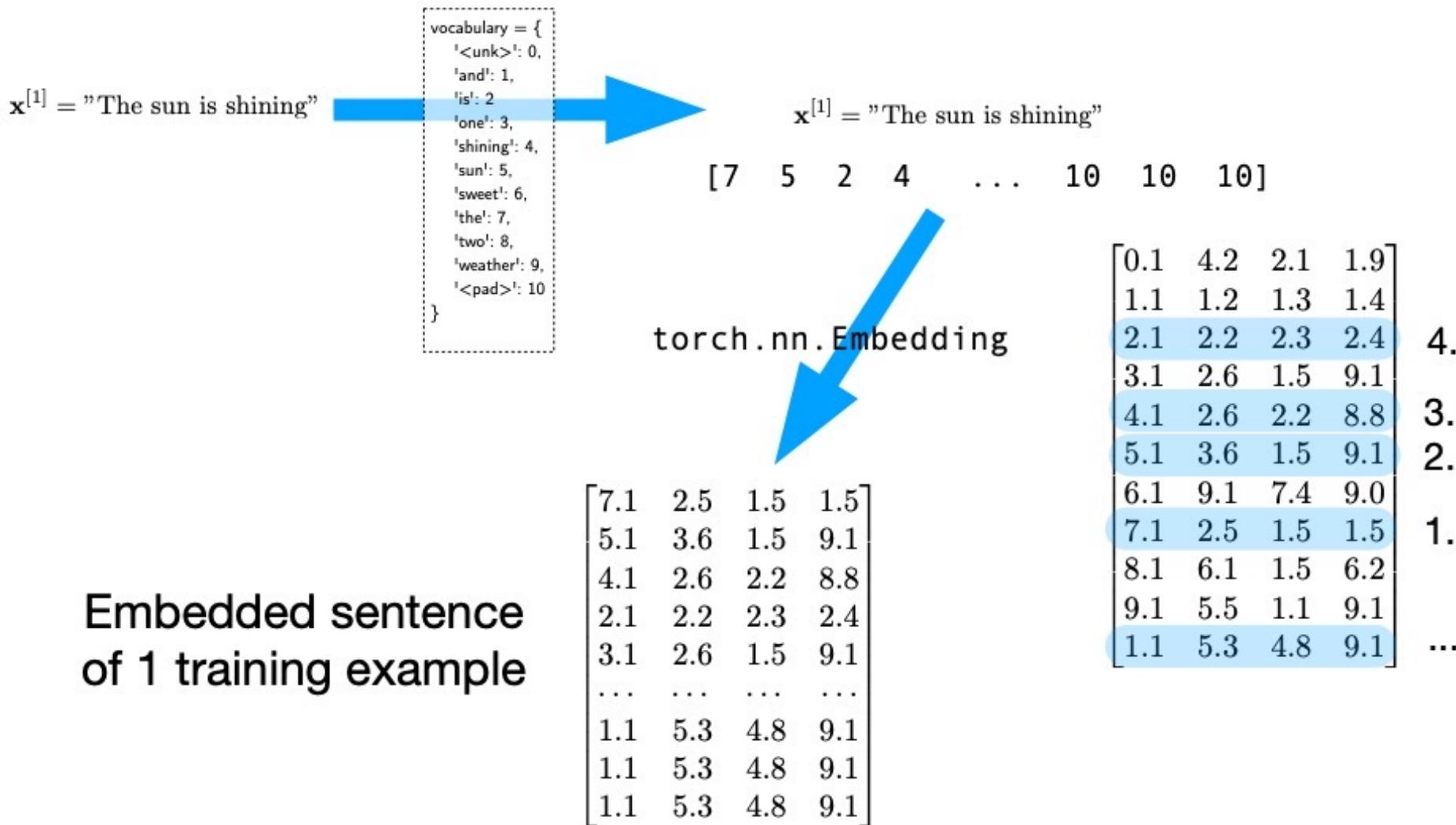
	[0 0 0 0 0 0 0 1 0 0 0]
	[0 0 0 0 0 1 0 0 0 0 0]
	[0 0 1 0 0 0 0 0 0 0 0]
	[0 0 0 1 0 0 0 0 0 0 0]
	[...]
10	[0 0 0 1 0 0 0 0 0 0 1]
10	[0 0 0 1 0 0 0 0 0 0 1]
10	[0 0 0 1 0 0 0 0 0 0 1]

RNN Step 4): One-Hot to Real via Embedding Matrix



In Practice, Skip Steps 3 & 4 And ...

use a lookup function (`torch.nn.Embedding`)



IMDB Dataset

- Dataset for binary sentiment classification
- 25,000 highly polar movie reviews for training, and 25,000 for testing

<https://ai.stanford.edu/~amaas/data/sentiment/>

	review	sentiment
0	In 1974, the teenager Martha Moxley (Maggie Gr...	1
1	OK... so... I really like Kris Kristofferson a...	0
2	***SPOILER*** Do not read this, if you think a...	0
3	hi for all the people who have seen this wonde...	1
4	I recently bought the DVD, forgetting just how...	0