



UNIVERSIDADE
CATÓLICA
PORTUGUESA

BRAGA

Deep Learning

Session 14

Introduction to Recurrent Neural Networks (RNNs)

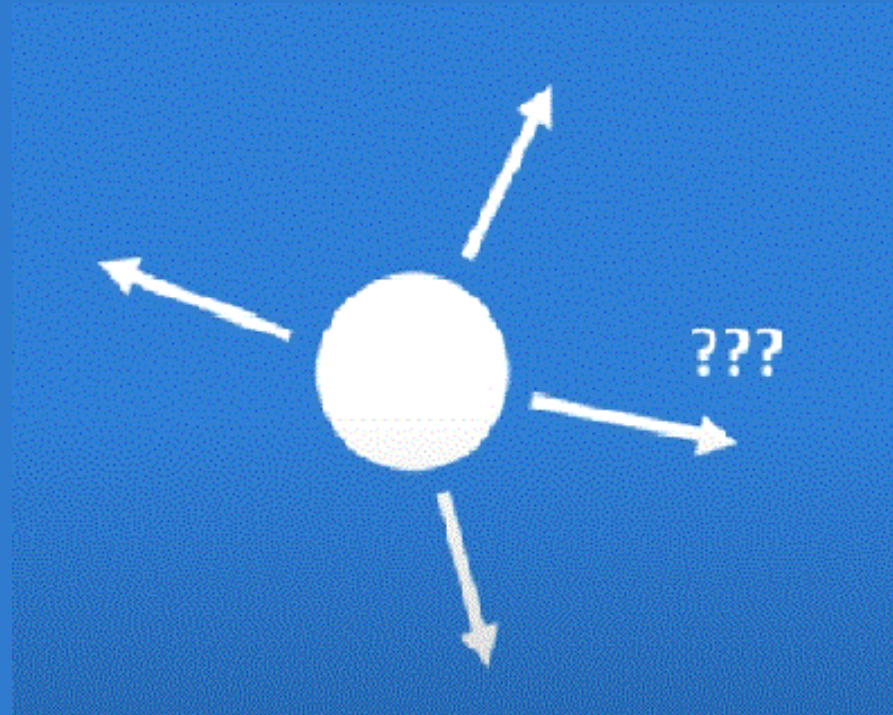
Applied Data Science

2024/2025

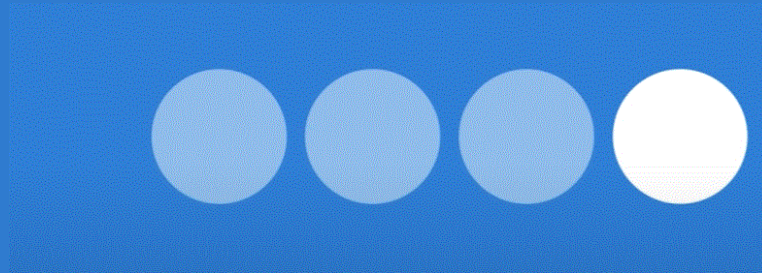
Given an image of a ball can you predict where it will go next?



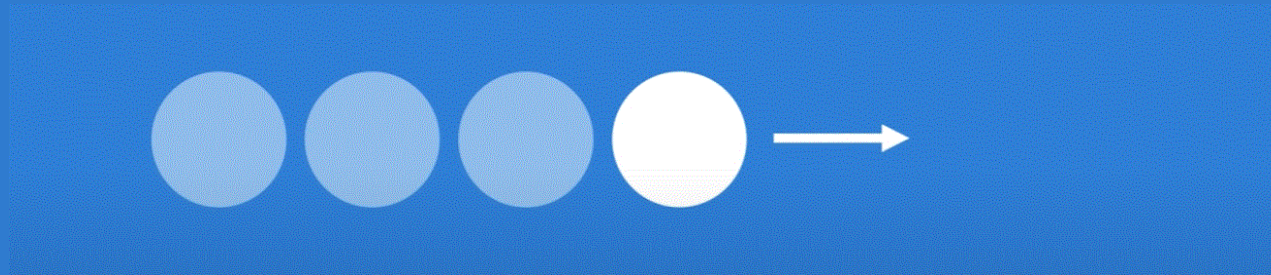
UNIVERSIDADE
CATOLICA
PORTUGUESA
BRAGA



Given an image of a ball can you predict where it will go next?

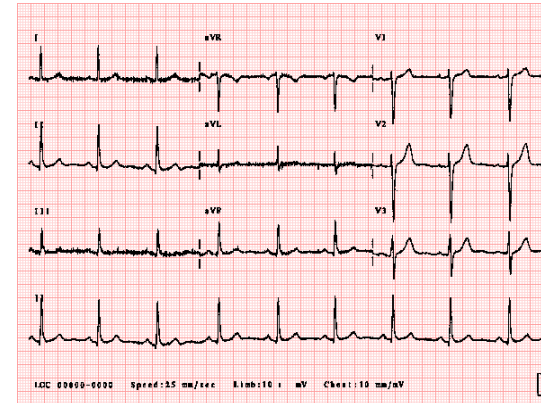
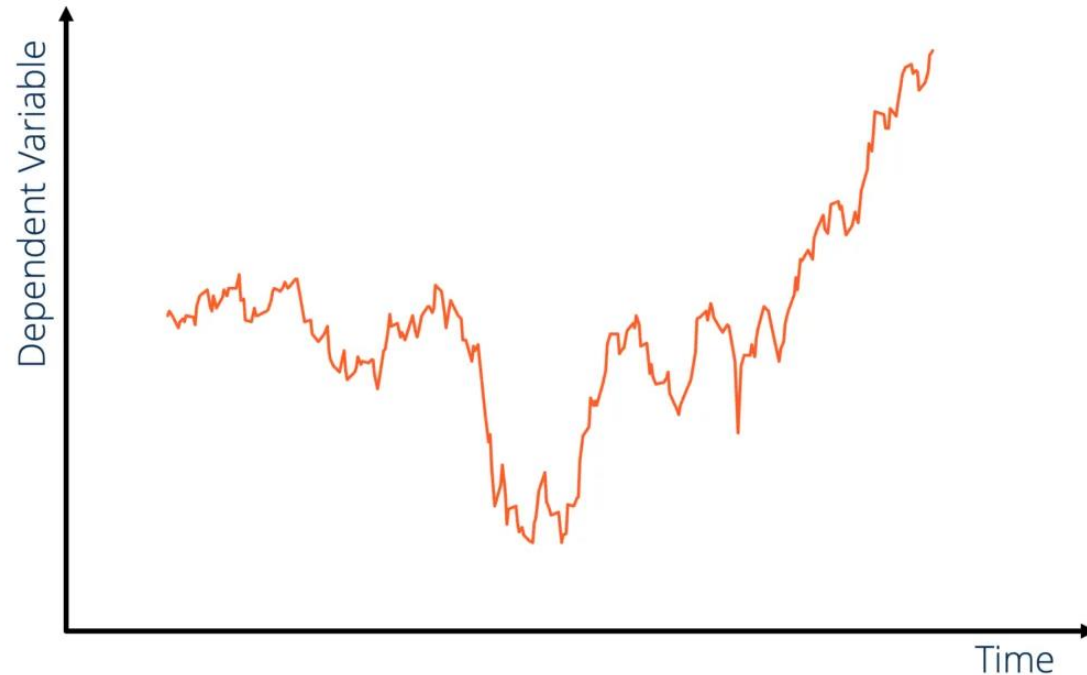


Given an image of a ball can you predict where it will go next?



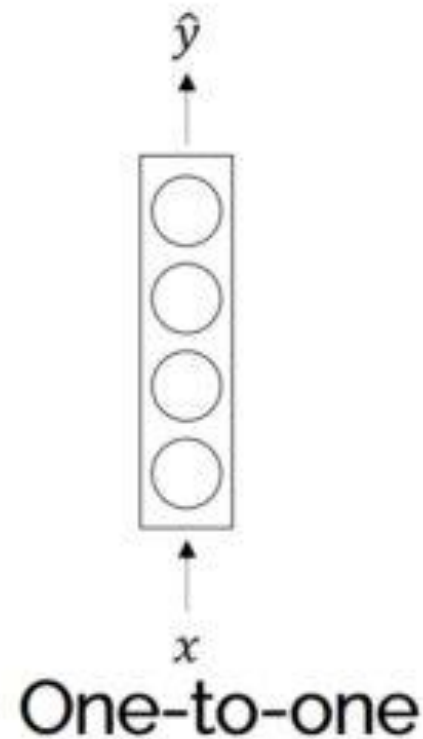
Sequential Data

- Elements in a sequence occur in a certain order
- Elements depend on each other



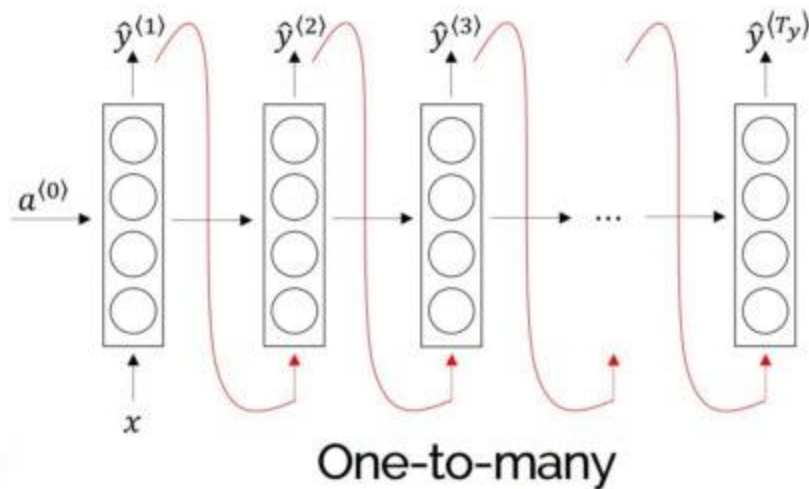
Sequence Applications: One-to-One

- **Description:** Processes a **single input** to produce a **single output**. Typically used in traditional feedforward networks, **not truly recurrent**.



Sequence Applications: One-to-Many

- **Description:** Takes **one input** and generates a **sequence of outputs**, often used for generating text or music.



A person is walking along a beach with a big dog



A black and white dog carries a tennis ball in its mouth



A soccer player takes a soccer ball in the grass



A man is doing a trick on a snowboard



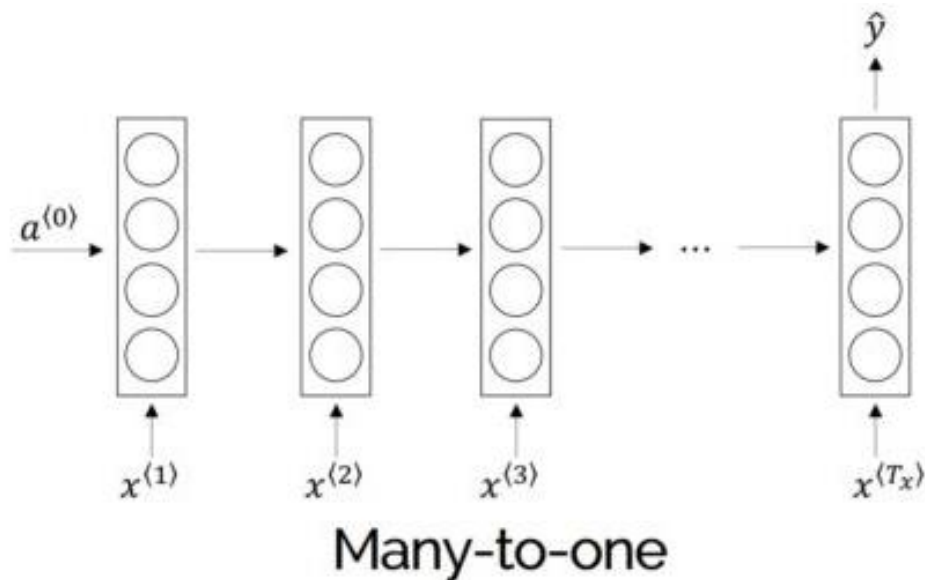
A surfer dives into the ocean



A black and white dog leaps to catch a Frisbee


Sequence Applications: Many-to-One

- **Description:** Takes a **sequence of inputs** and produces a **single output**, ideal for analyzing sequences.




CRITIC REVIEWS FOR *STAR WARS: THE LAST JEDI*


All Critics (371) | Top Critics (51) | Fresh (336) | Rotten (35)

 What's most interesting to me about The Last Jedi is Luke's return as the mentor rather than the student, grappling with his failure in this new role, and later aspiring to be the wise and patient teacher.

December 26, 2017 | Rating: 3/4 | [Full Review...](#)


 **Leah Pickett**
Chicago Reader
★ Top Critic

 POSITIVE

 Fanatics will love it; for the rest of us, it's a tolerably good time.

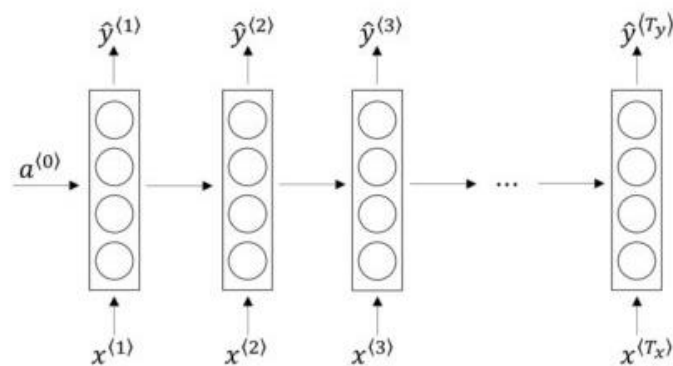
December 15, 2017 | Rating: B | [Full Review...](#)

 **Peter Rainer**
Christian Science Monitor
★ Top Critic

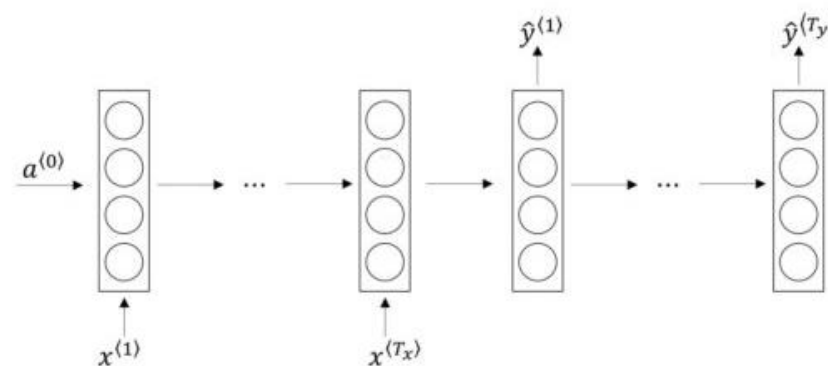
 NEUTRAL

Sequence Applications: Many-to-Many

- **Description:** Processes a **sequence of inputs** to generate a **sequence of outputs**, useful for tasks where both input and output are sequential.



Many-to-many



Many-to-many

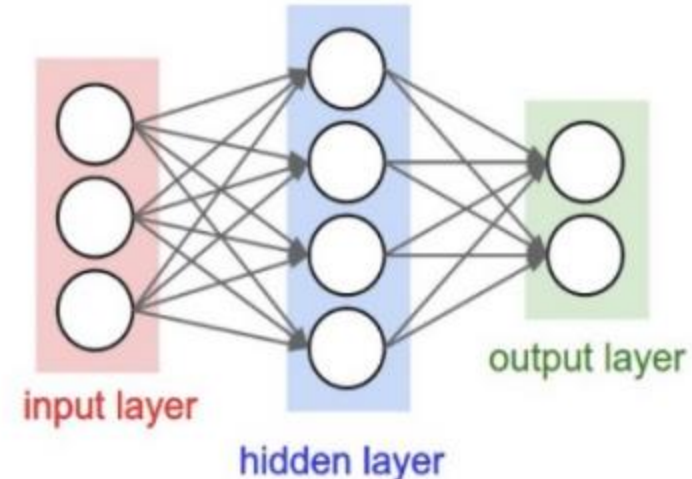
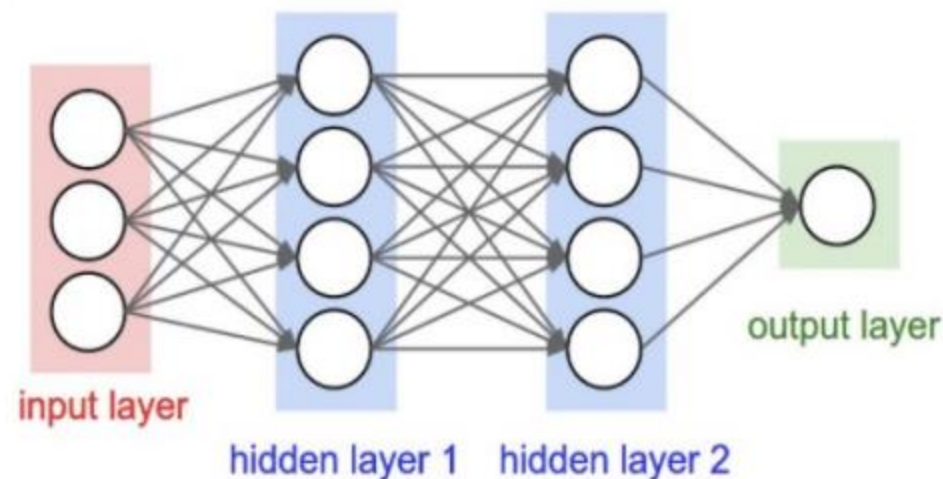
Sequence Applications: Many-to-Many

- **Description:** Processes a **sequence of inputs** to generate a **sequence of outputs**, useful for tasks where both input and output are sequential.



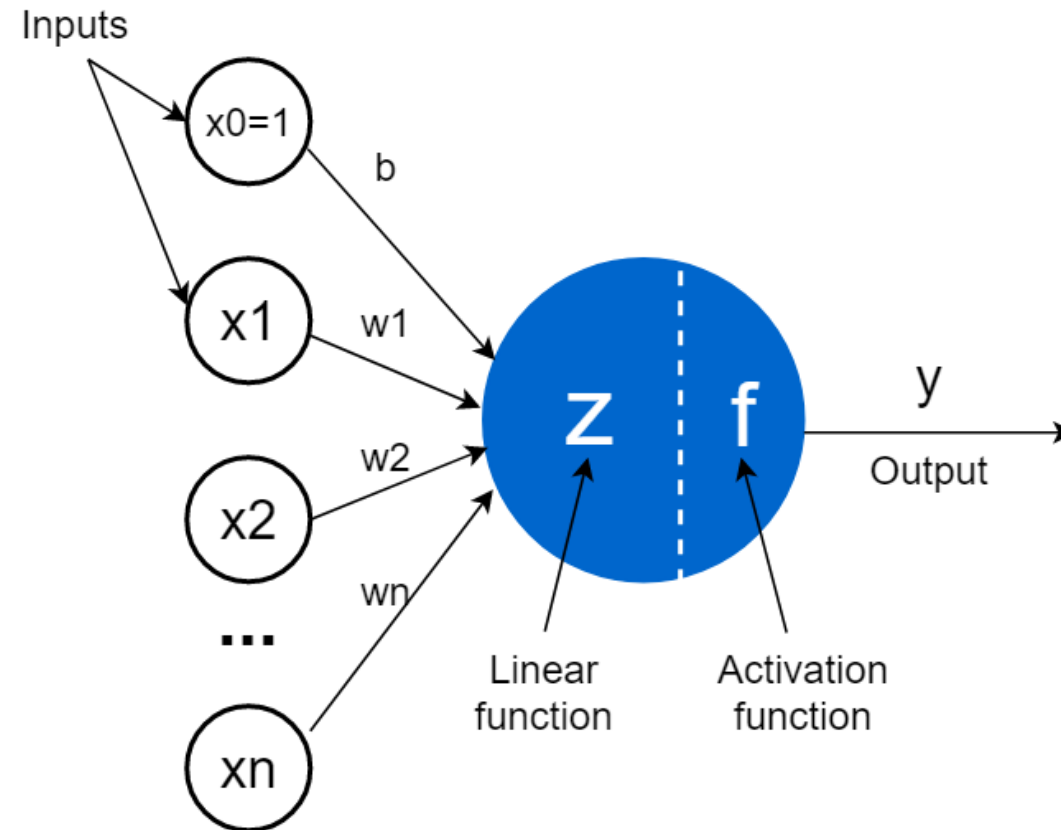
Recall: Feedforward Neural Networks

- Problems:
 - Many model parameters;
 - Input / output sizes are fixed;
 - No memory of past since weights are learned independently.



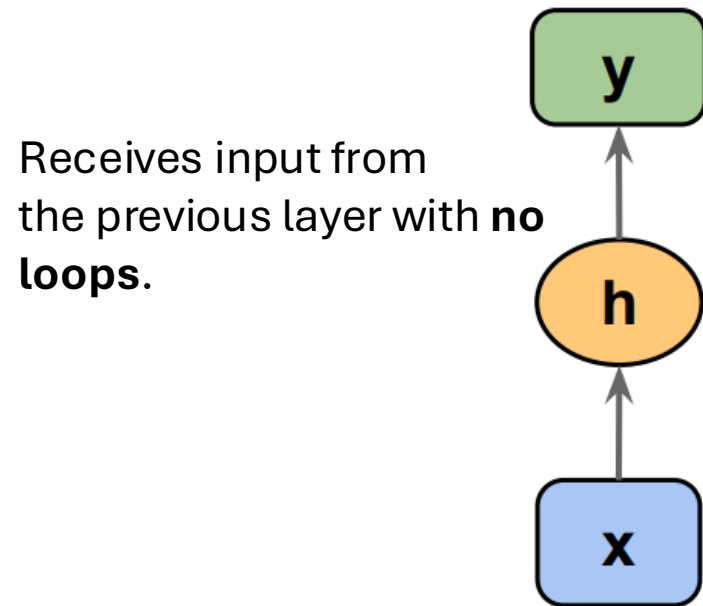
- Each layer serves as input to the next layer with no loops

Recall: The Perceptron



Recurrent Neural Networks (RNNs)

- **Main idea:** use hidden state to capture information about the past



Recurrent Neural Network (RNN)

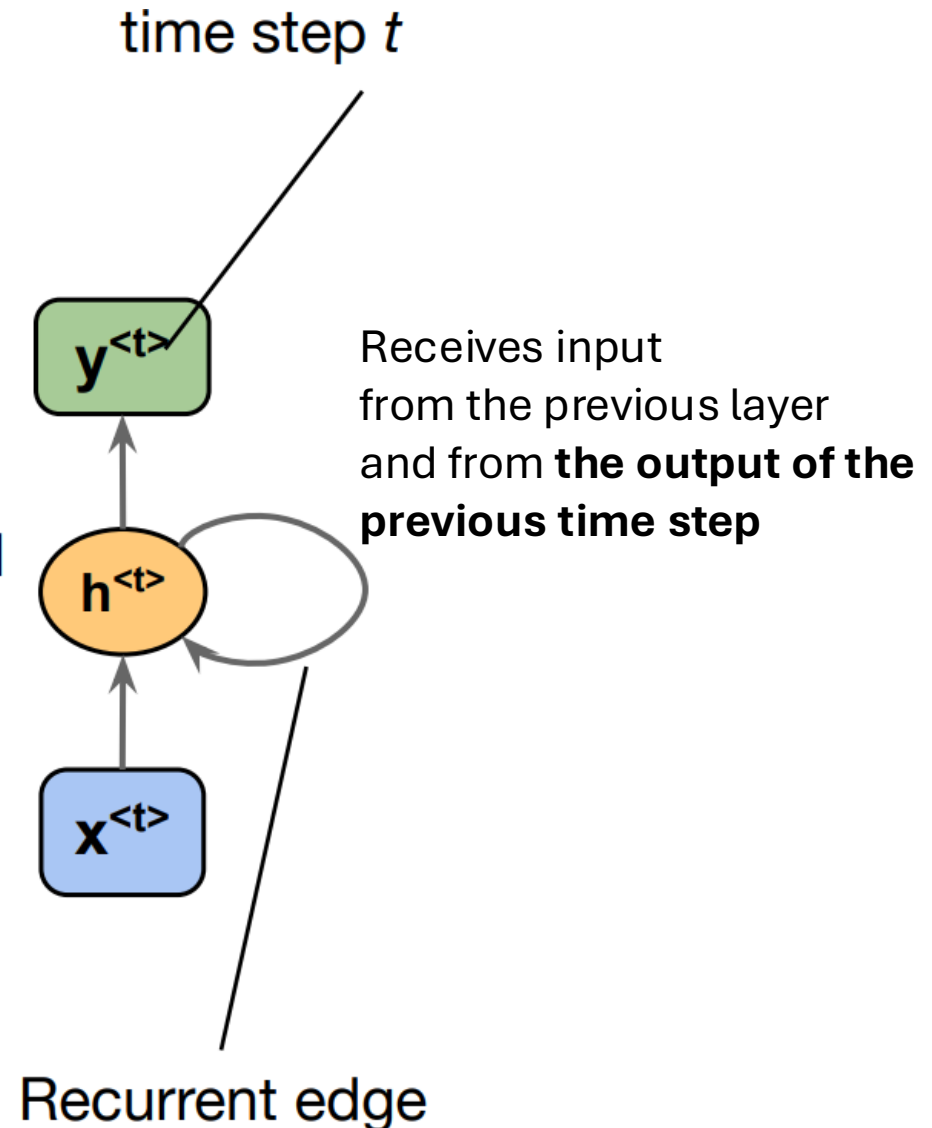


Image source: Sebastian Raschka, Vahid Mirjalili, *Python Machine Learning*, 3rd Edition, Packt, 2019

Recurrent Neural Networks (RNNs)

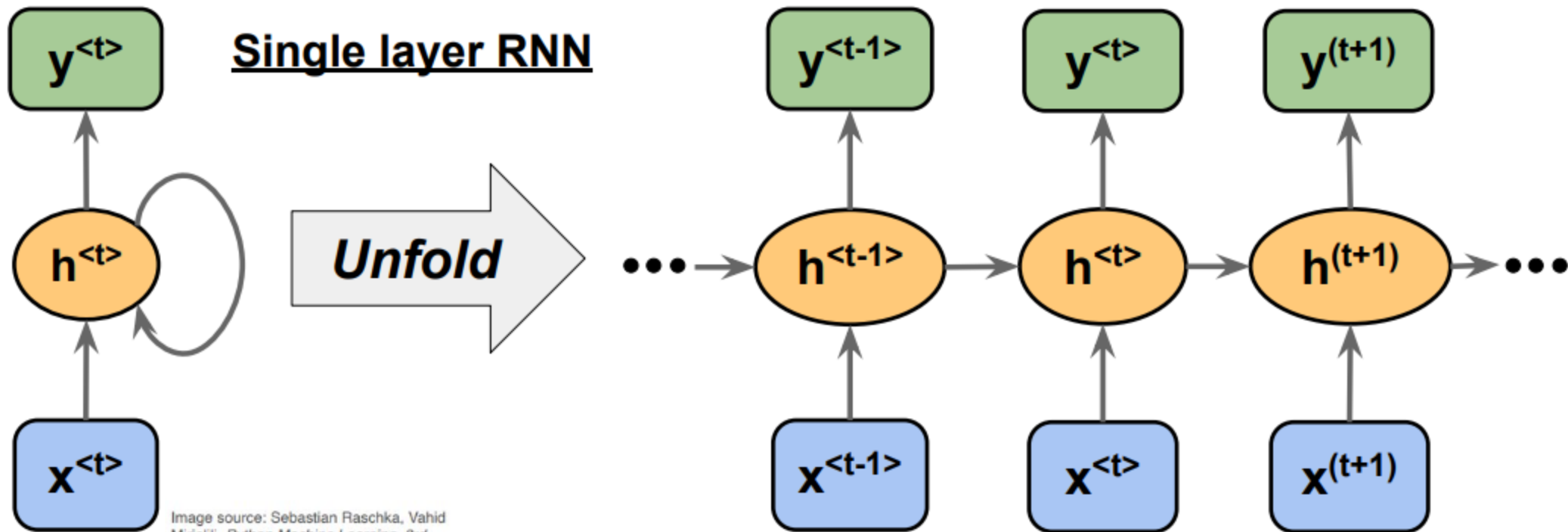


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Recurrent Neural Networks (RNNs)

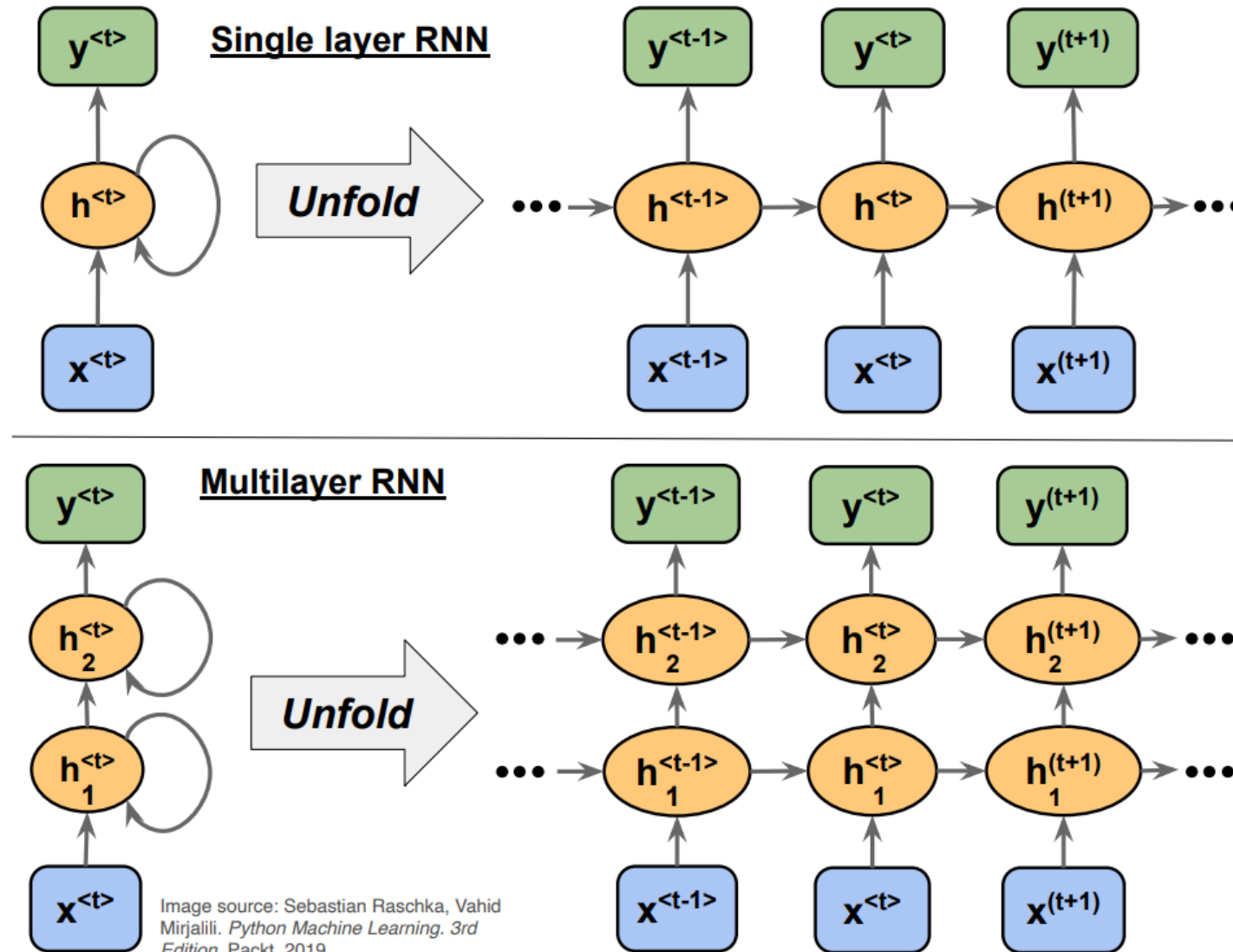
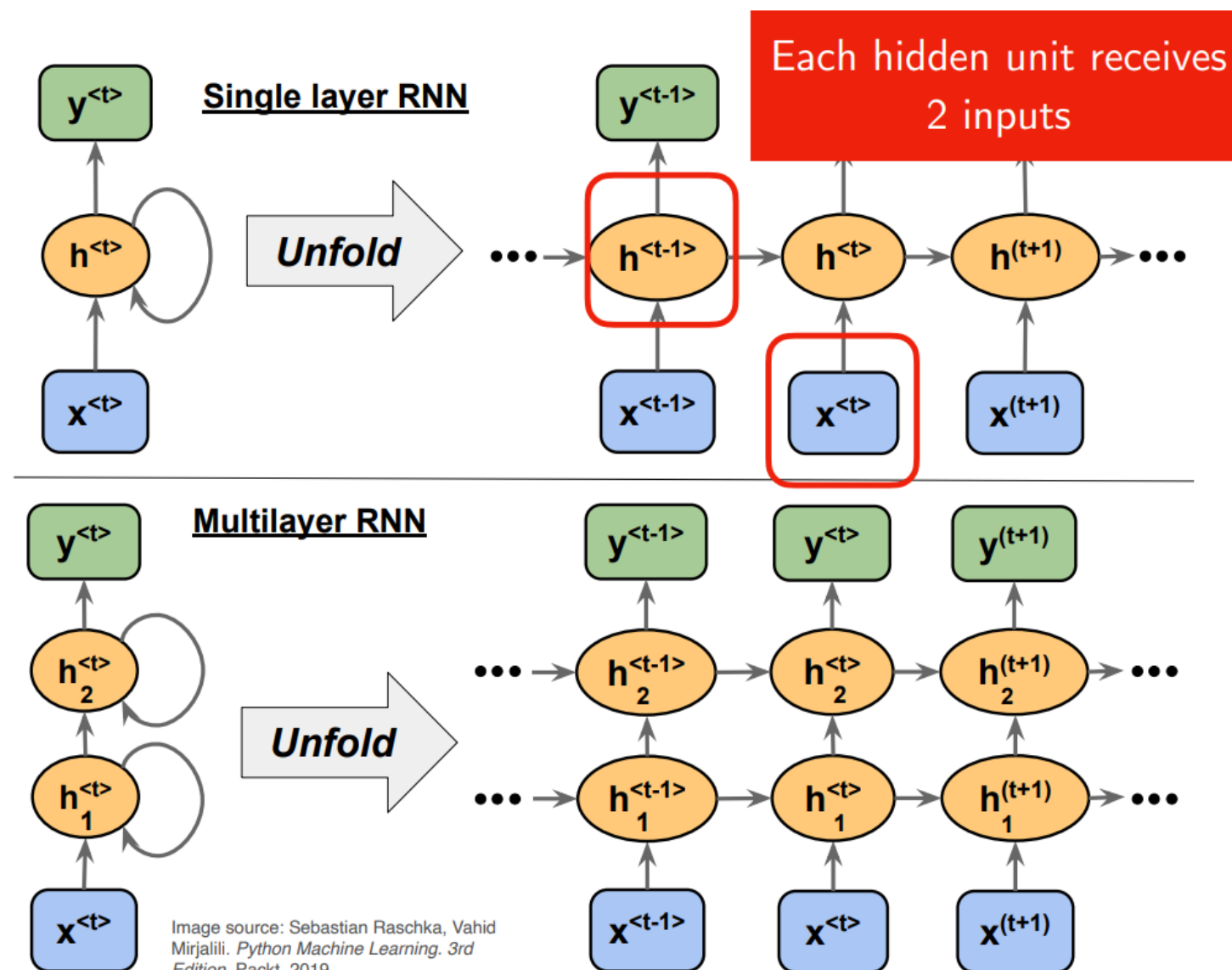


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Recurrent Neural Networks (RNNs)



$$\hat{y}_t = f(\underbrace{x_t}_{\text{input}}, \underbrace{h_{t-1}}_{\text{past memory}})$$

output

Backpropagation Through Time

- Weight matrices in a single-hidden layer RNN

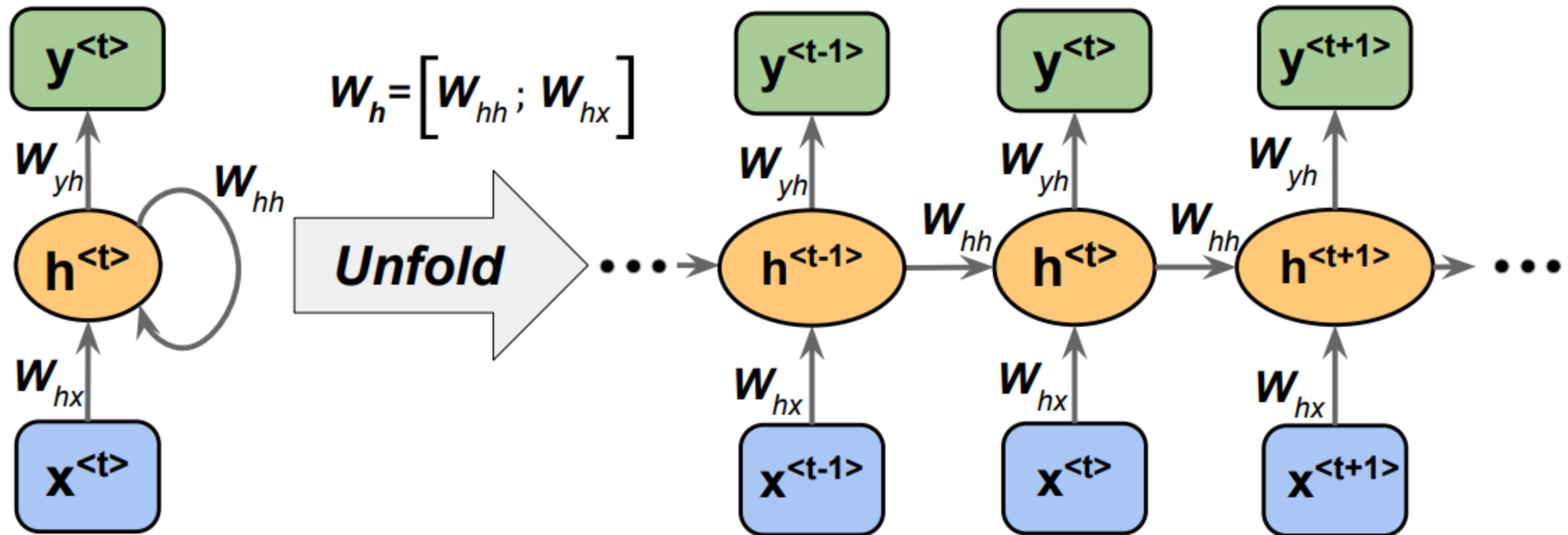


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Backpropagation Through Time

- Weight matrices in a single-hidden layer RNN

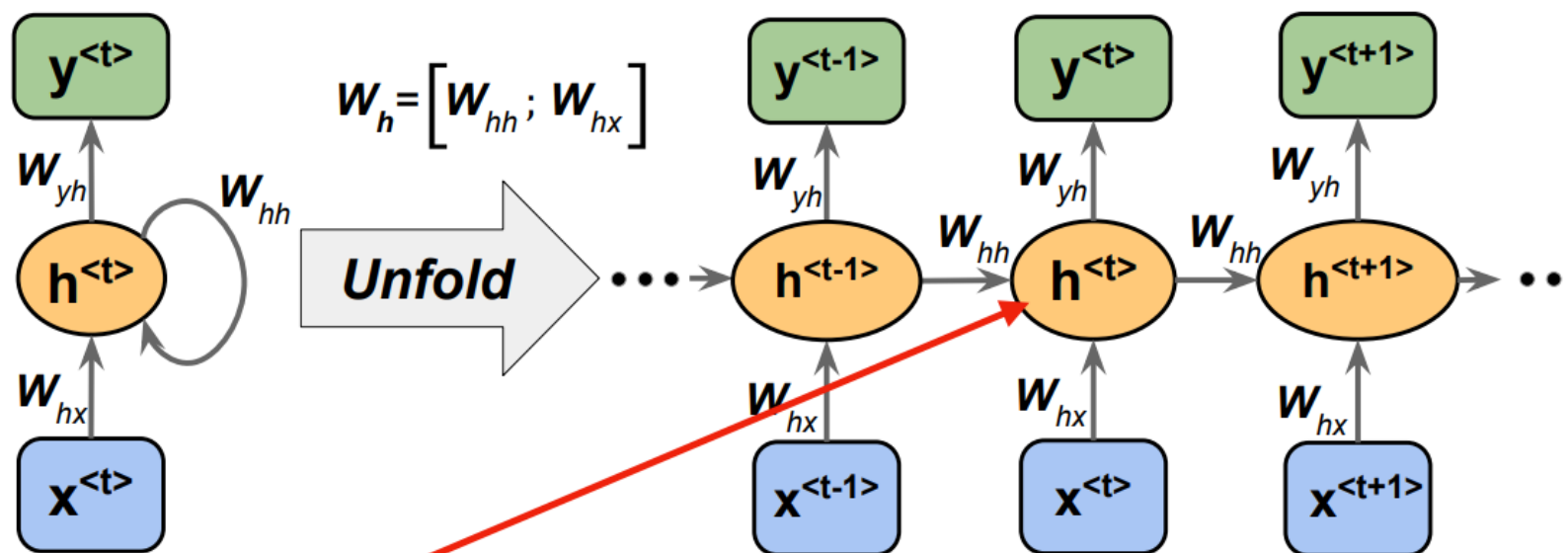


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Net input:

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

Activation:

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

Backpropagation Through Time

- Weight matrices in a single-hidden layer RNN

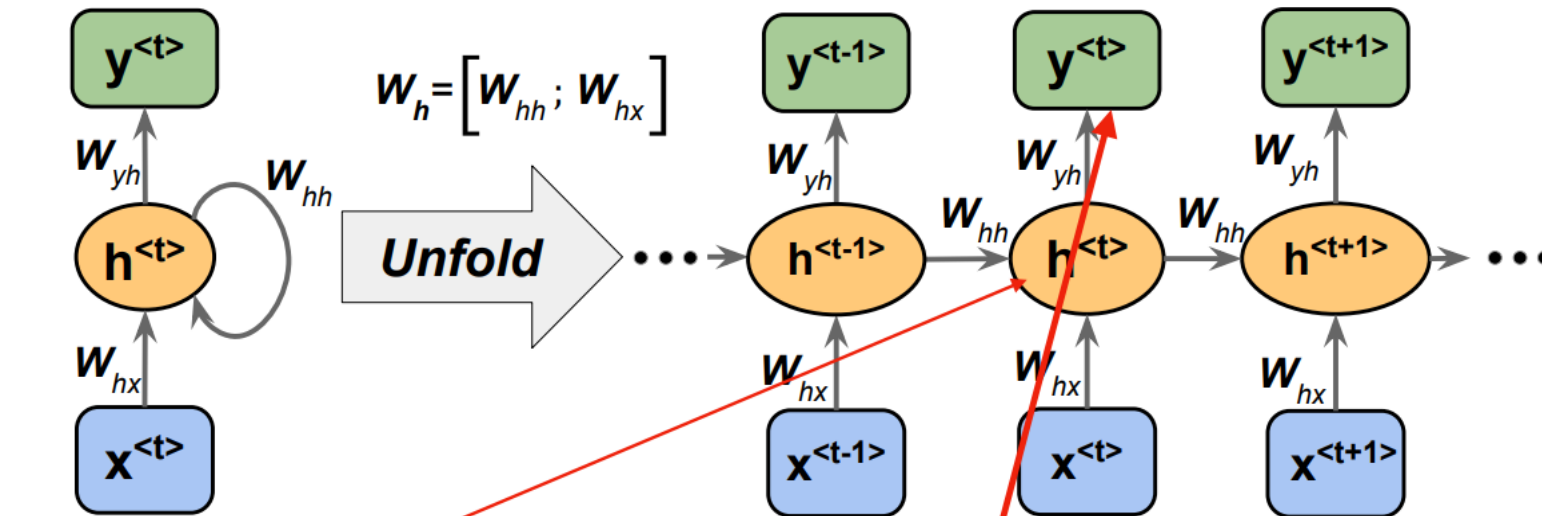


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

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:

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

Backpropagation Through Time

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

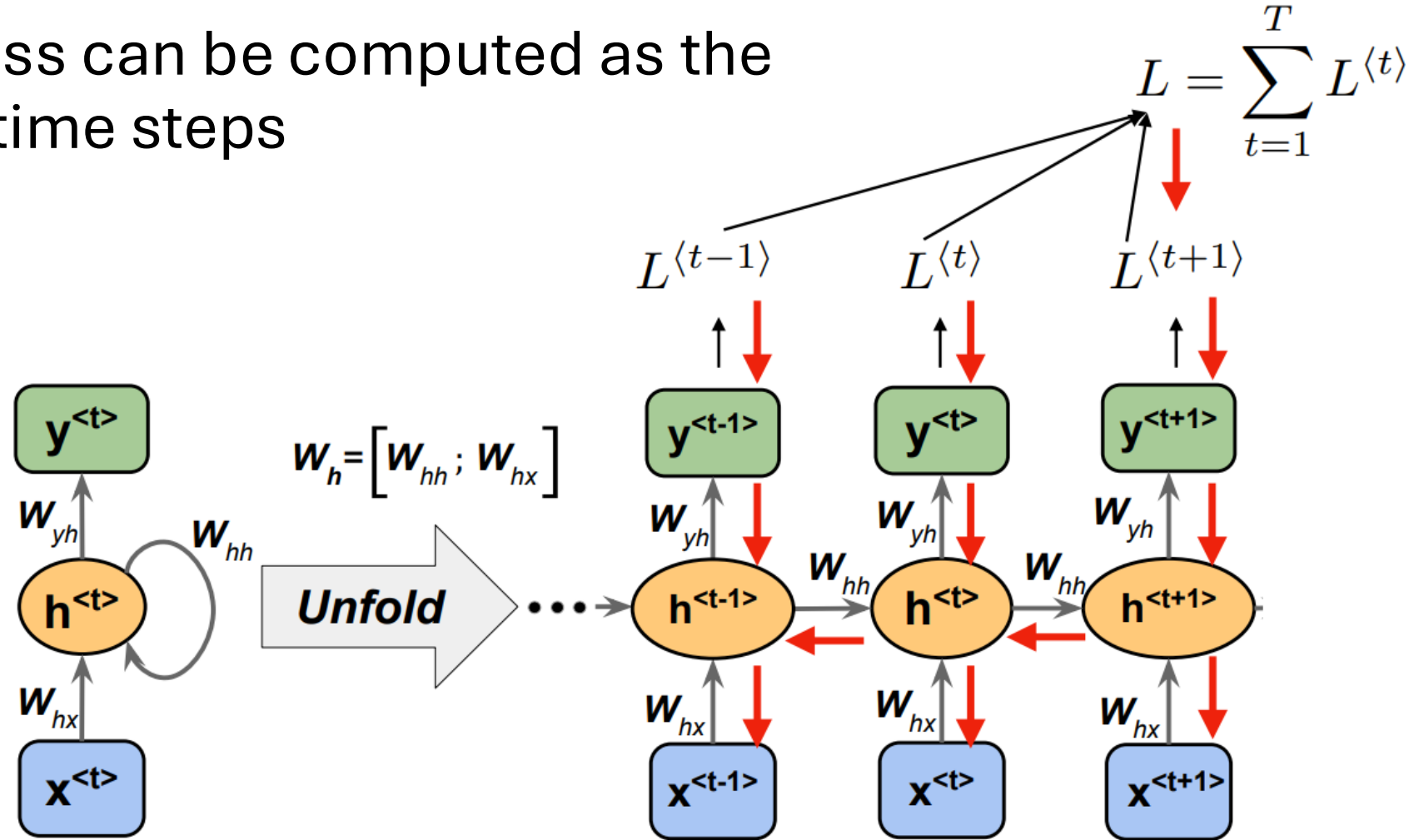
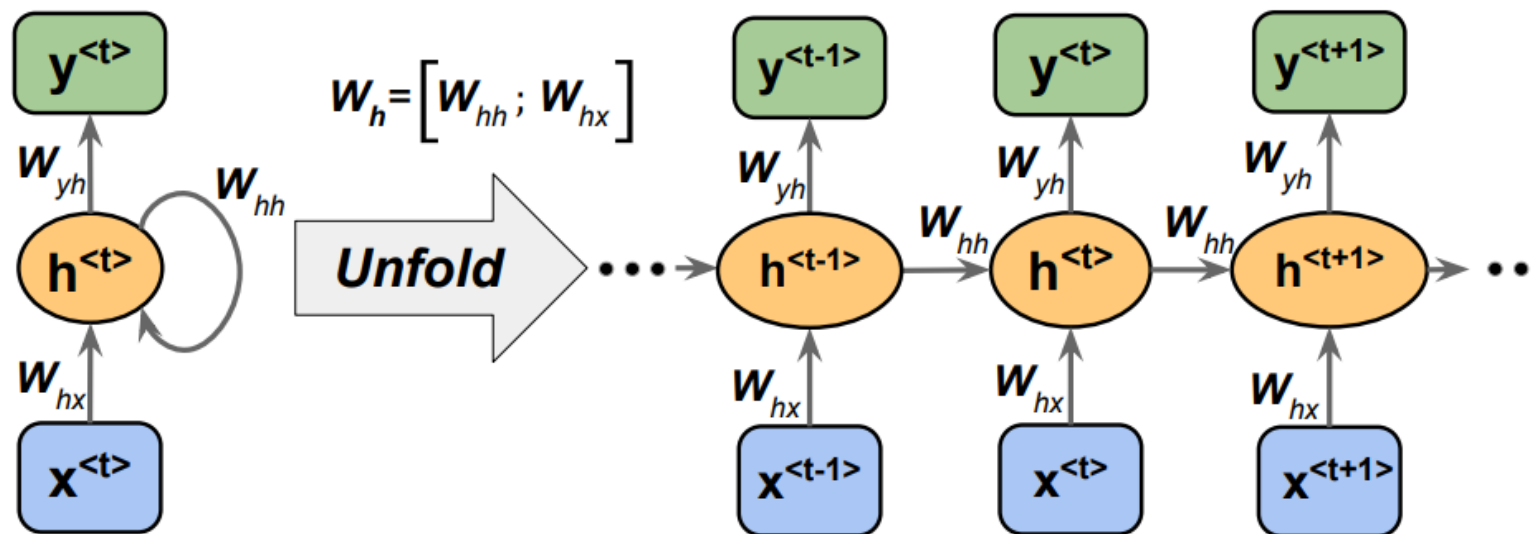


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

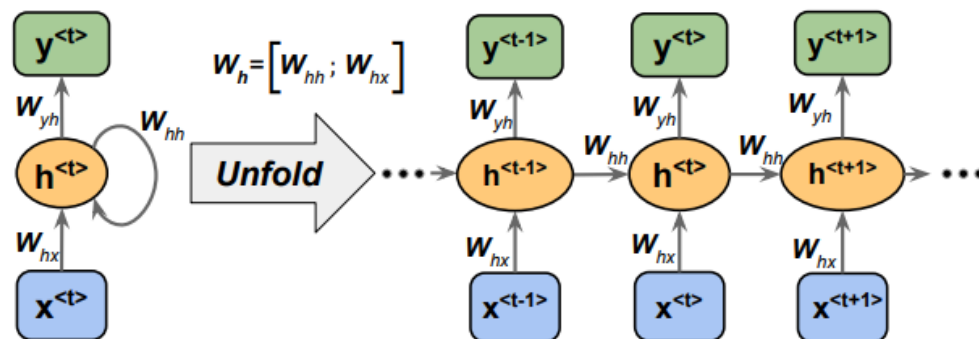
Backpropagation Through Time



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

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

Backpropagation Through Time

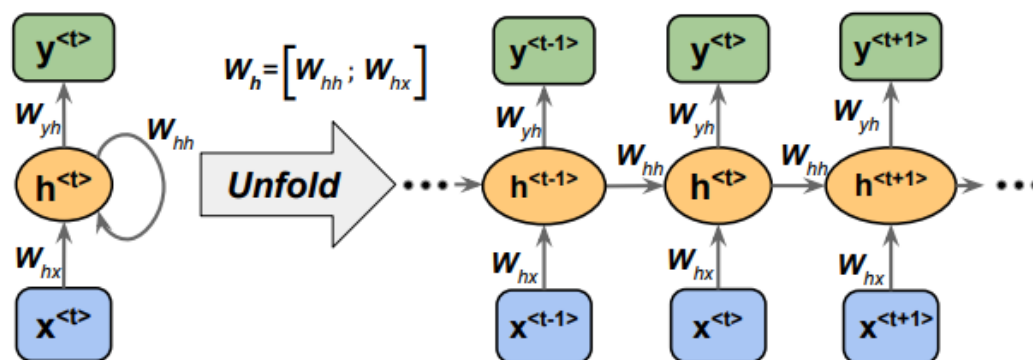


$$\frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \boxed{\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}}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

computed as a multiplication of adjacent time steps:

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

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.

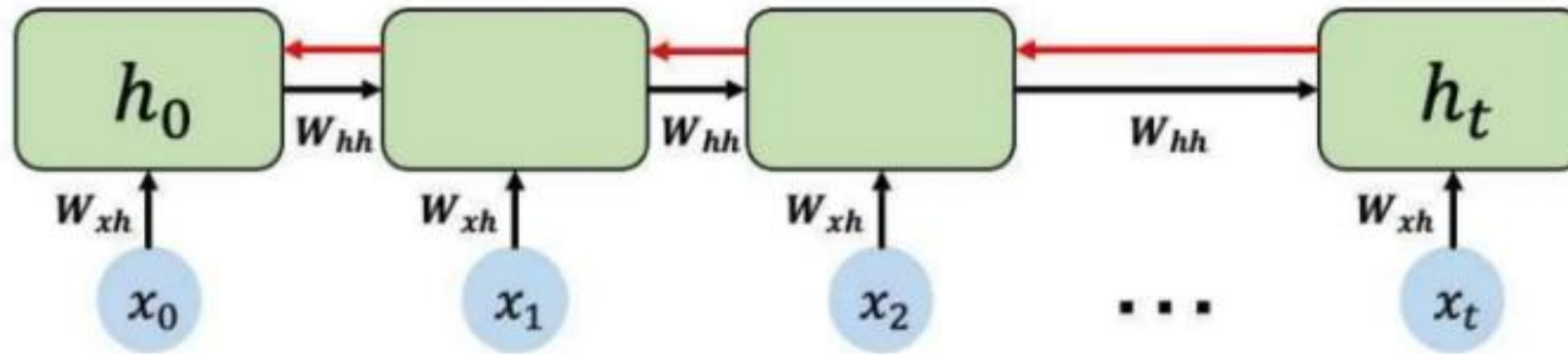
$$L = \sum_{t=1}^T L^{(t)} \quad \frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial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

computed as a multiplication of adjacent time steps:

This is very problematic:
Vanishing/Exploding gradient problem!

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

Standard RNN Gradient Flow: Exploding Gradients

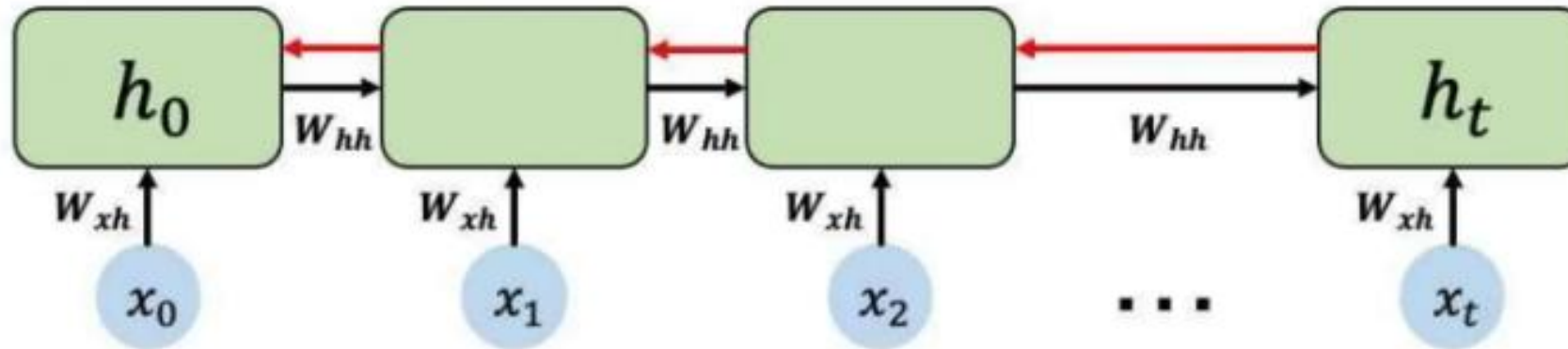


Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Standard RNN Gradient Flow: Vanishing Gradients



Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1 :
vanishing gradients

1. Activation function
2. Weight initialization
3. Network architecture

The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies

The Problem of Long-Term Dependencies

"The clouds are in the ____"

Why are vanishing gradients a problem?

Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies

The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

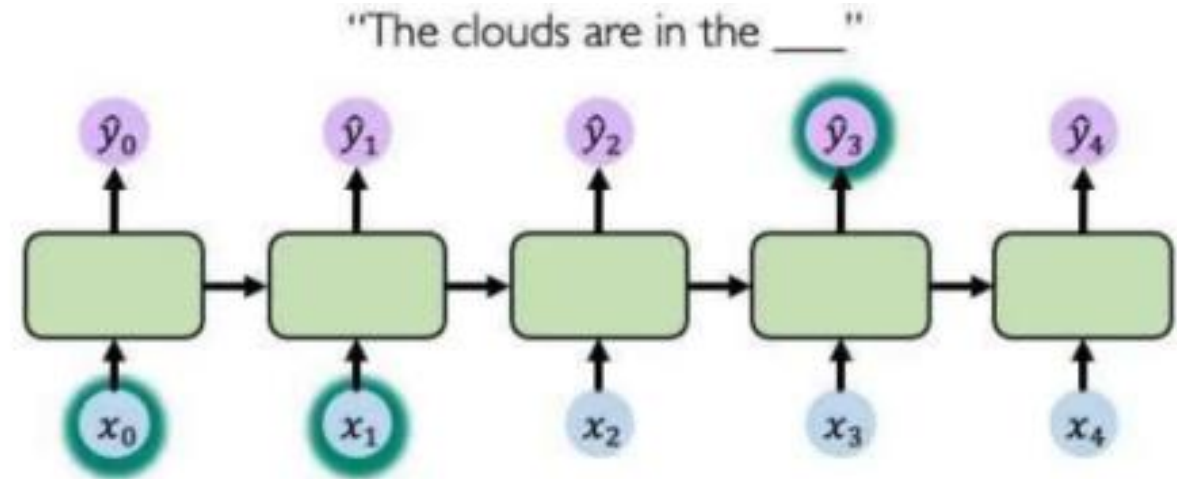
Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

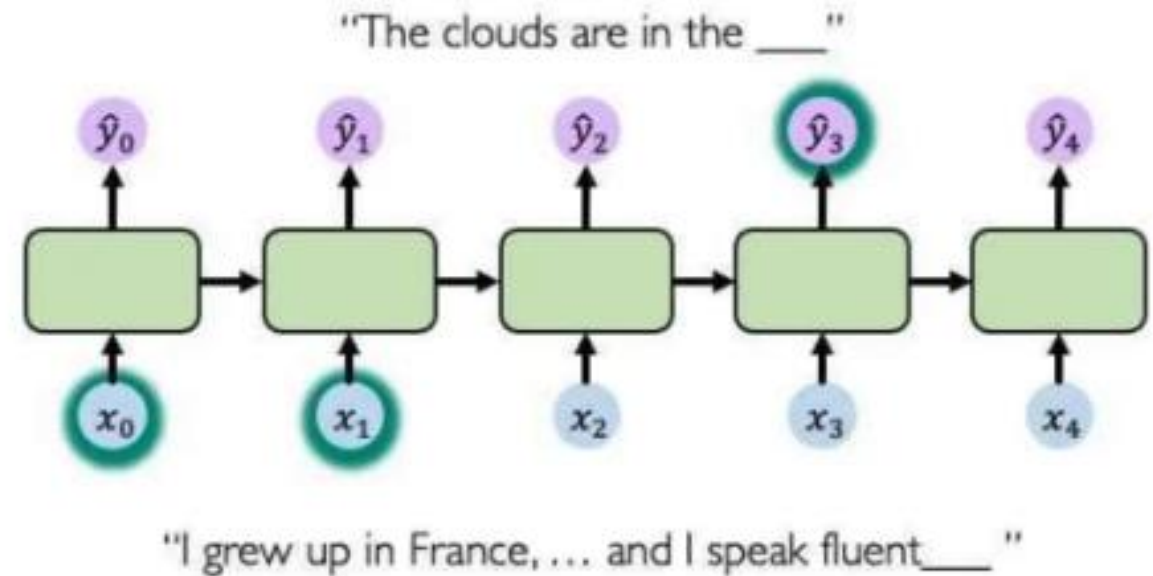
Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

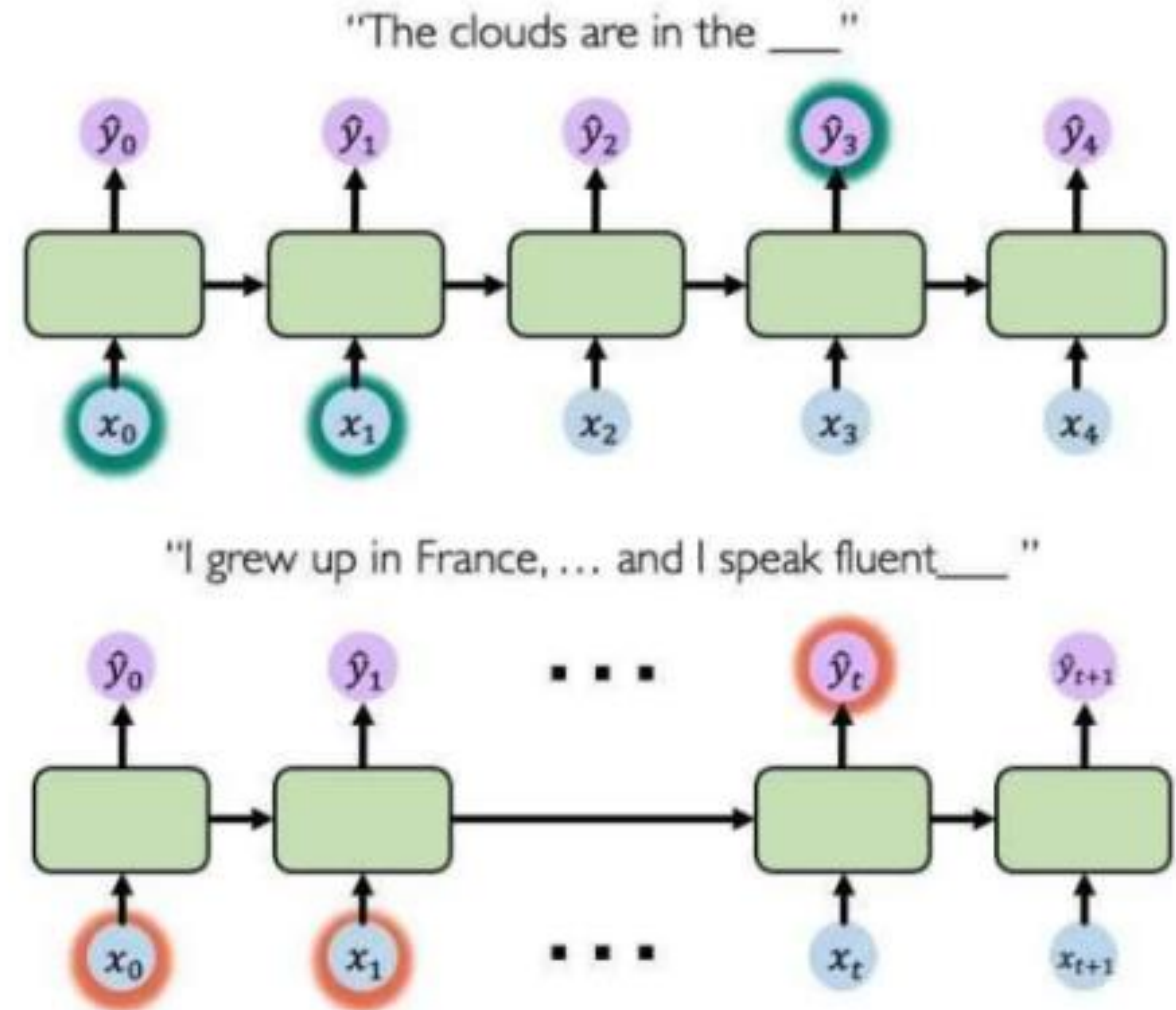
Multiply many **small numbers** together



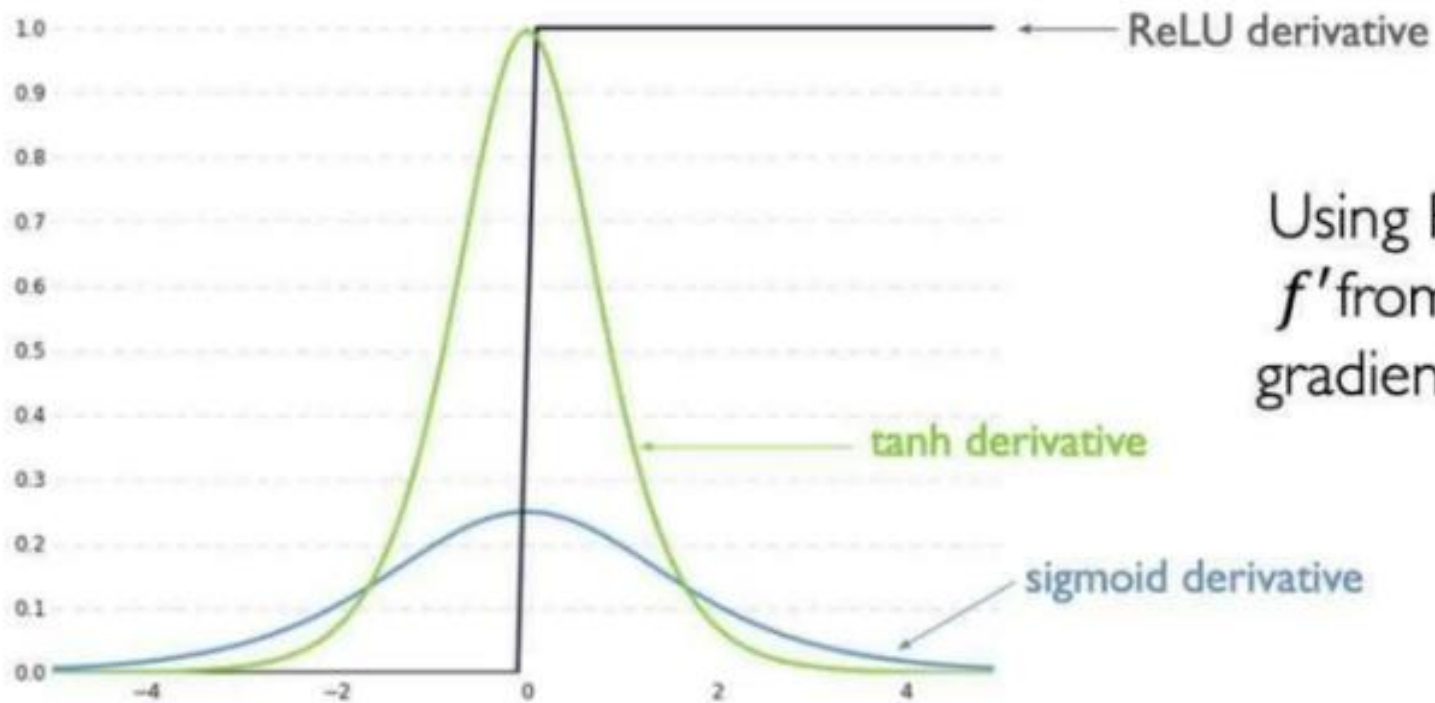
Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



Solutions to the Vanishing/Exploding Gradients



Using ReLU prevents f' from shrinking the gradients when $x > 0$

Solutions to the Vanishing/Exploding Gradients



Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

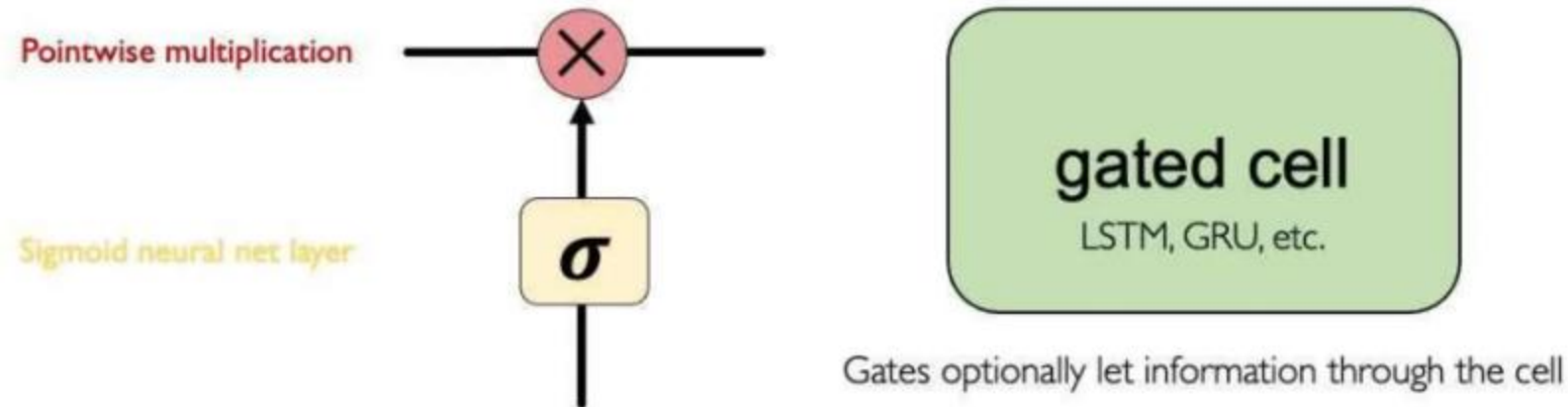
Solutions to the Vanishing/Exploding Gradients



- **Gradient Clipping:** Set a maximum threshold for gradients to prevent them from becoming excessively large, effectively addressing the exploding gradient problem.
- **Truncated Backpropagation Through Time (TBPTT):** Limits the number of time steps for which gradients are backpropagated after each forward pass. For example, with a sequence of 100 steps, TBPTT may only backpropagate through the most recent 20 steps.
- **Gated Cells:** use gates to selectively add or remove information within each recurrent unit

Gated Cells

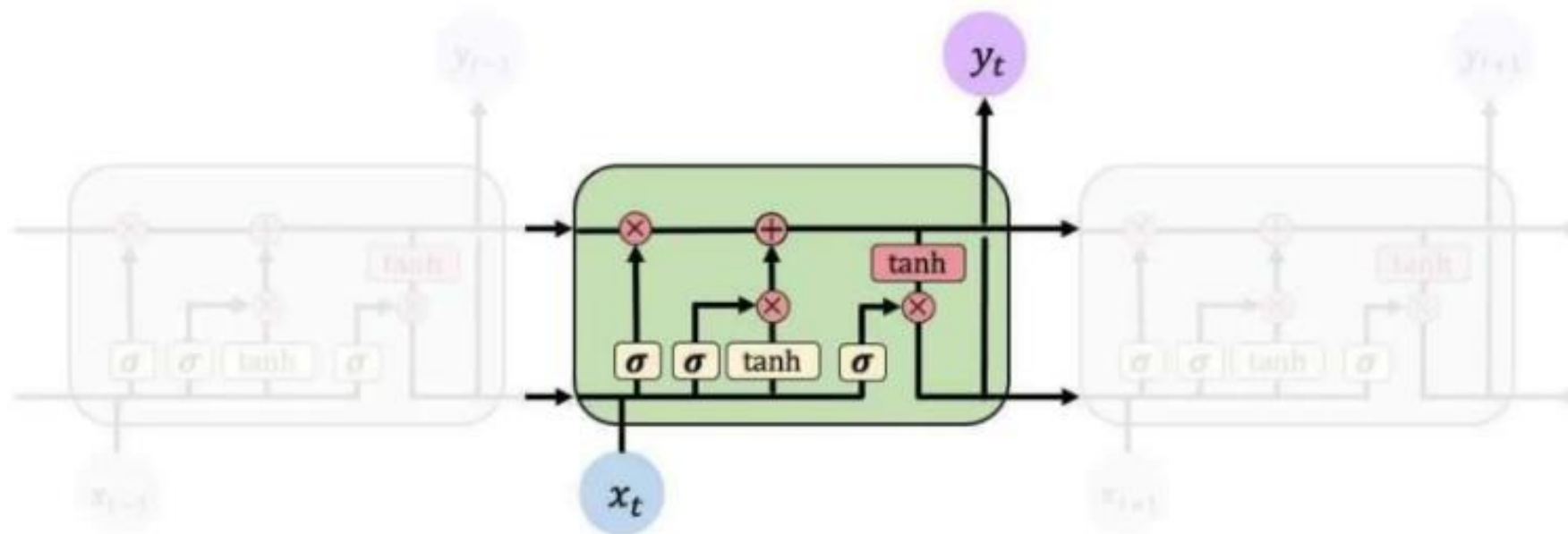
- **Idea:** use gates to selectively add or remove information within each recurrent unit



- **Long Short Term Memory (LSTM)** networks rely on gated cells to track information throughout many time steps.

Long-Short Term Memory (LSTM)

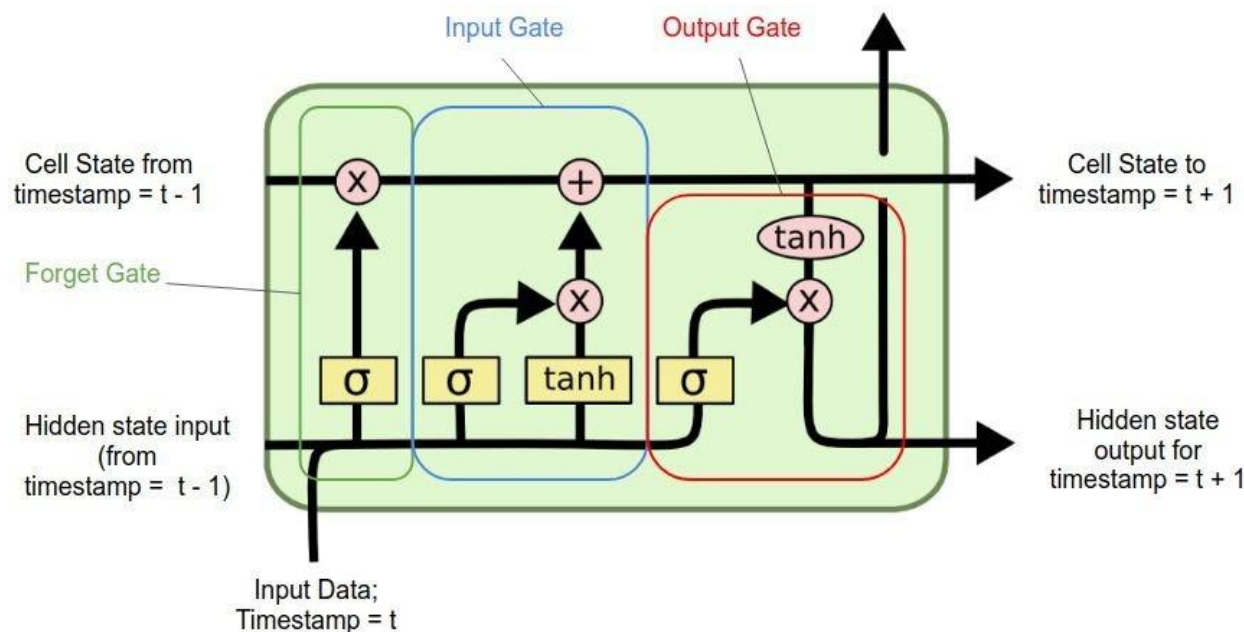
- Gated LSTM cells control information flow:
1) Forget 2) Store 3) update 4) Output



- LSTM cells are able to track information throughout many timesteps

Long-Short Term Memory (LSTM)

- Gated LSTM cells control information flow:
 - 1) Forget
 - 2) Store
 - 3) update
 - 4) Output



- LSTM cells are able to track information throughout many timesteps

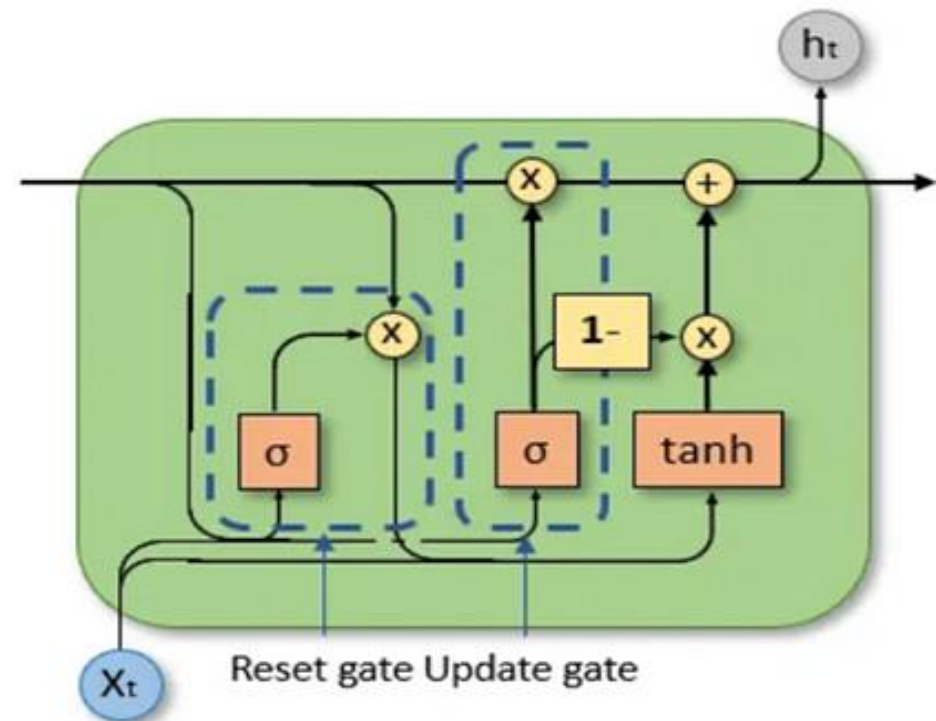


LSTMs Key Concepts

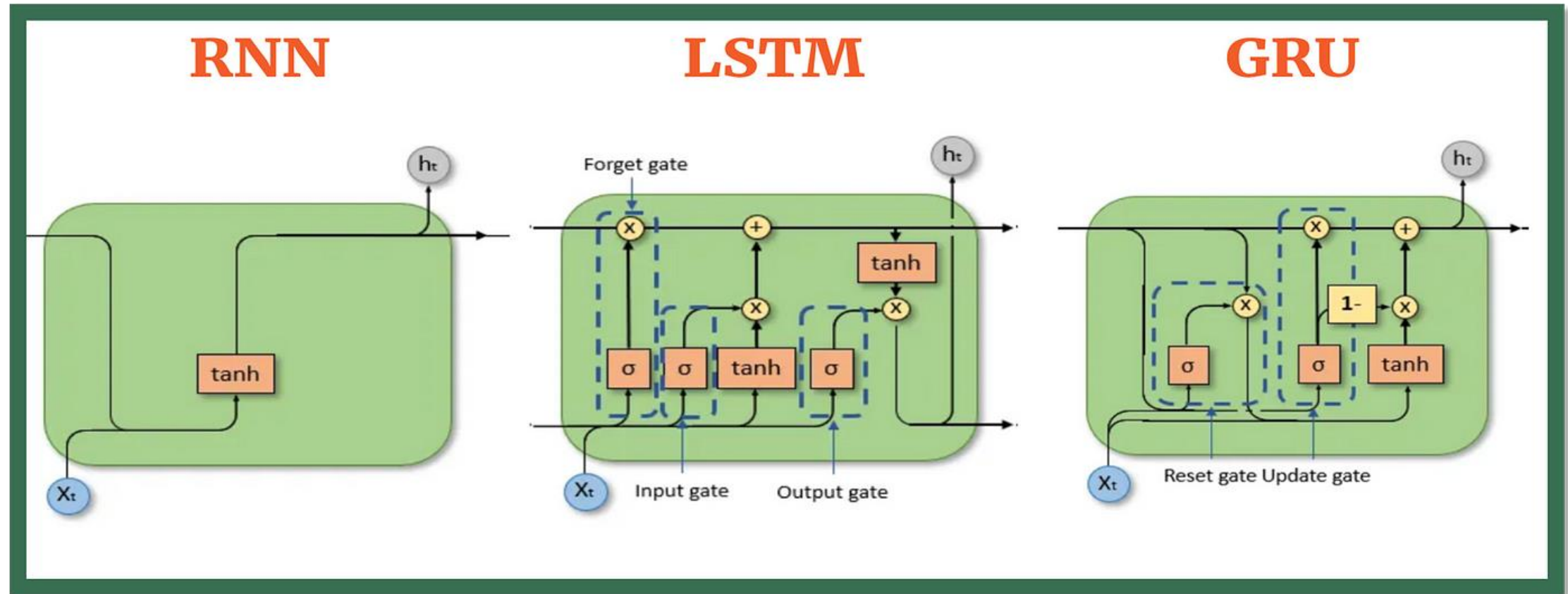
- Maintain a **cell state**
- Use **gates** to control the **flow of information**
 - **Forget** gate gets rid of irrelevant information
 - **Store** relevant information from current input
 - Selectively **update** the cell state
 - **Output** gate returns a filtered version of the cell state
- Backpropagation through time with partially **uninterrupted gradient flow**

Gated Recurrent Unit (GRU)

- Simplifies LSTM by merging: (1) cell and hidden states and (2) forget and input gates



Simple RNN vs LSTM vs GRU



A Sequence Modeling Problem: Predict the Next Word



"This morning I took my cat for a walk."

given these words

predict the
next word

A Sequence Modeling Problem: Predict the Next Word

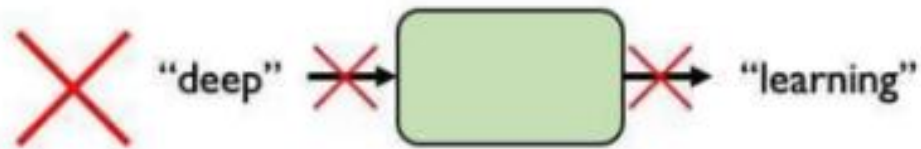
"This morning I took my cat for a walk."

given these words

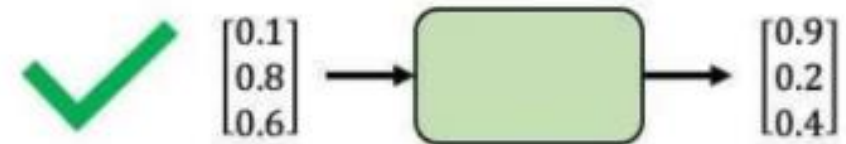
predict the
next word



Representing Language to a Neural Network

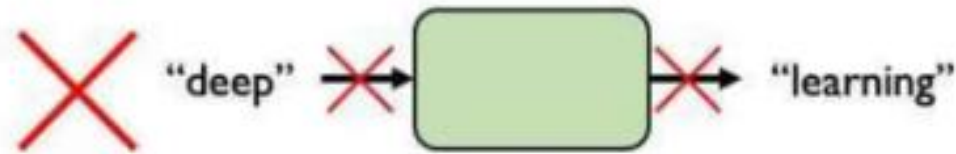


Neural networks cannot interpret words

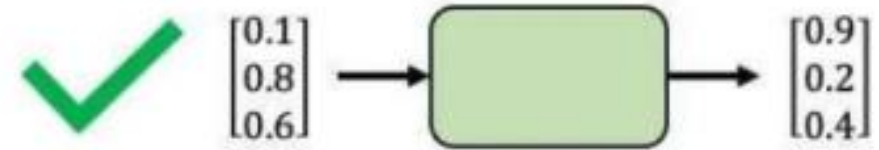


Neural networks require numerical inputs

A Sequence Modeling Problem: Predict the Next Word



Neural networks cannot interpret words



Neural networks require numerical inputs

Embedding: transform indexes into a vector of fixed size.

this cat for
my took
a | walk
morning

1. Vocabulary:
Corpus of words

a → 1
cat → 2
...
walk → N

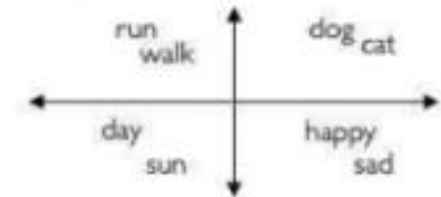
2. Indexing:
Word to index

One-hot embedding

"cat" = $[0, 1, 0, 0, 0, 0]$

i -th index

Learned embedding



3. Embedding:
Index to fixed-sized vector

Handle Variable Sequence Lengths



The food was great

vs.

We visited a restaurant for lunch

vs.

We were hungry but cleaned the house before eating

Capture Differences in Sequence Order



The food was good, not bad at all.

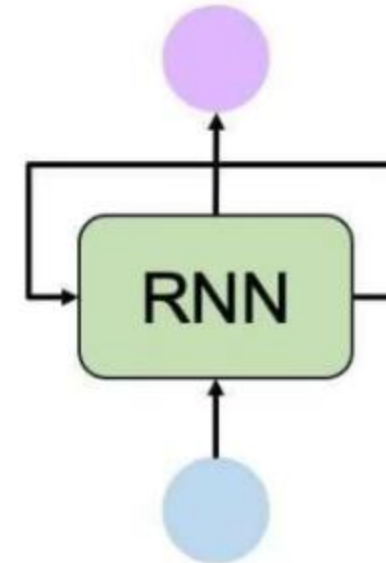
vs.

The food was bad, not good at all.



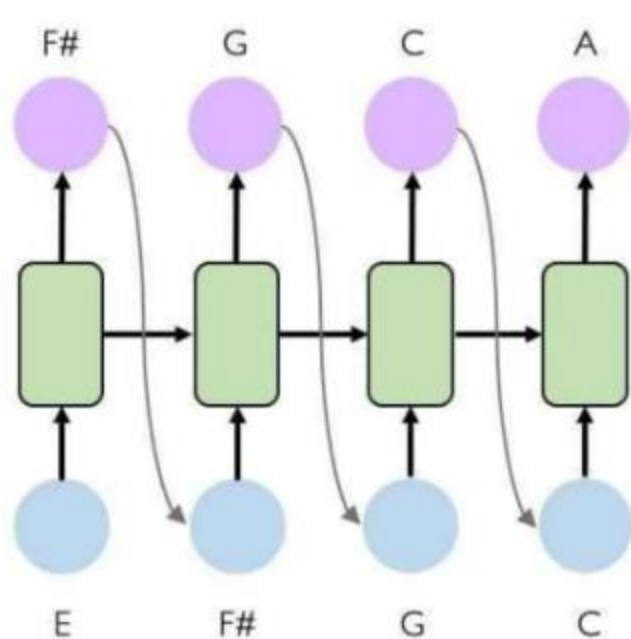
Sequence Modeling: Design Criteria

- To model sequences, we need:
 1. Handle variable-length sequences
 2. Track long-term dependencies
 3. Maintain information about order
 4. Share parameters across the sequence
- RNNs meet these sequence modeling design criteria



RNNs Applications and Limitations

- Example Task: Music Generation

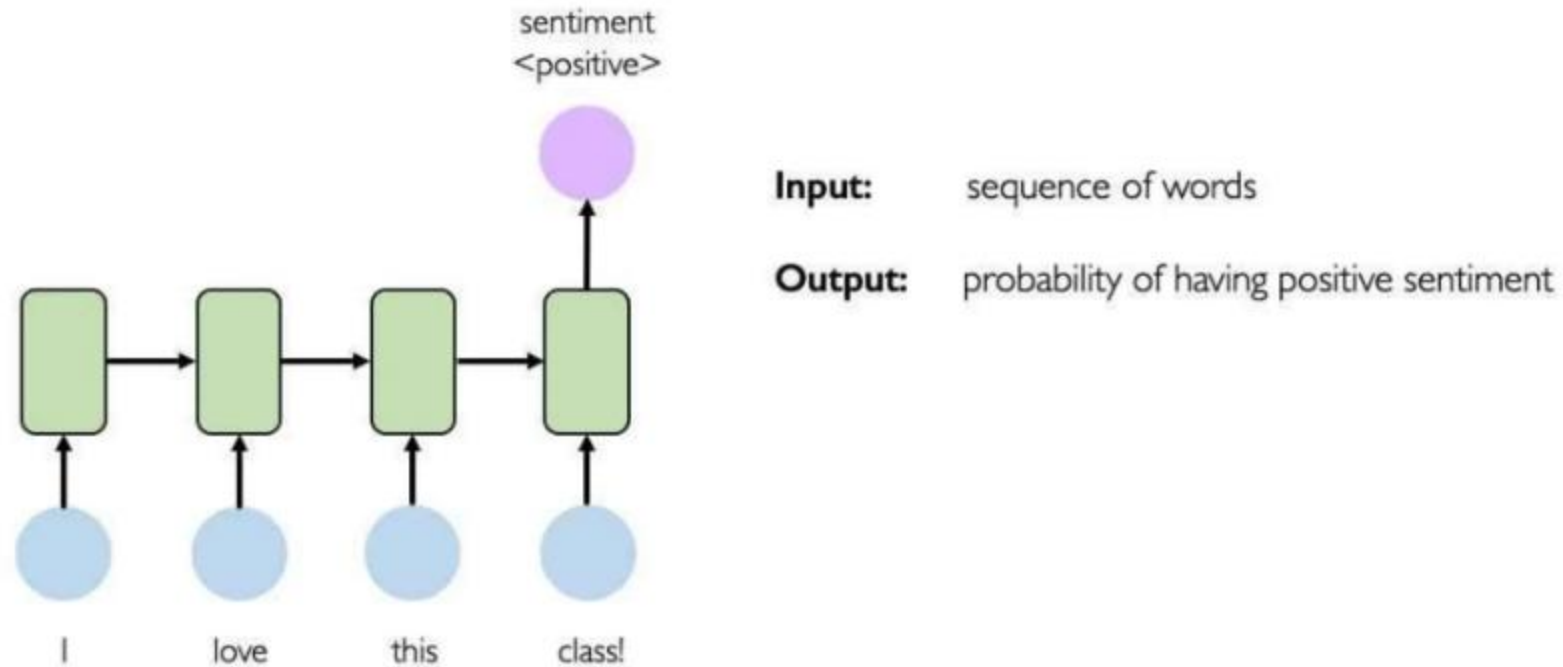


Input: sheet music
Output: next character in sheet music



RNNs Applications and Limitations

- Example Task: Sentiment Classification





RNNs Applications and Limitations

- Limitations of RNNs:
 - Encoding bottleneck
 - Slow, no parallelization
 - Not long memory