



# STAT 453: Introduction to Deep Learning and Generative Models

---

Ben Lengerich

Lecture 19: Recurrent Neural Networks

November 10, 2025

Reading: See course homepage



# Our semester

Week	Lecture Dates	Topic
<b>Module 1: Introduction and Foundations</b>		
1	9/3	Course Introduction
2	9/8, 9/10	A Brief History of DL, Statistics / linear algebra / calculus review
3	9/15, 9/17	Single-layer networks Parameter Optimization and Gradient Descent
4	9/22, 9/24	Automatic differentiation with PyTorch, Cluster and cloud computing resources
<b>Module 2: Neural Networks</b>		
5	9/29, 10/1	Multinomial logistic regression, Multi-layer perceptrons and backpropagation
6	10/6, 10/8	Regularization Normalization / Initialization
7	10/13, 10/15	Optimization, Learning Rates CNNs
8	10/20, 10/22	Review, <b>Midterm Exam</b>

<b>Module 3: Intro to Generative Models</b>			
9	10/27, 10/29	A Linear Intro to Generative Models, Factor Analysis, Autoencoders, VAEs	
10	11/3, 11/5	Generative Adversarial Networks, Diffusion Models	Project Midway Report
<b>Module 4: Large Language Models</b>			
11	11/10, 11/12	Sequence Learning with RNNs Attention, Transformers	HW4
12	11/17, 11/19	GPT Architectures, Unsupervised Training of LLMs	
13	11/24, 11/26	Supervised Fine-tuning of LLMs, Prompts and In-context learning	HW5
14	12/1, 12/3	Foundation models, alignment, explainability Open directions in LLM research	
15	12/8, 12/10	<b>Project Presentations</b>	Project Final Report
16	12/17	<b>Final Exam</b>	Final Exam



# A quick vote...

Week	Lecture Dates	Topic
<b>Module 1: Introduction and Foundations</b>		
1	9/3	Course Introduction
2	9/8, 9/10	A Brief History of DL, Statistics / linear algebra / calculus review
3	9/15, 9/17	Single-layer networks Parameter Optimization and Gradient Descent
4	9/22, 9/24	Automatic differentiation with PyTorch, Cluster and cloud computing resources
<b>Module 2: Neural Networks</b>		
5	9/29, 10/1	Multinomial logistic regression, Multi-layer perceptrons and backpropagation
6	10/6, 10/8	Regularization Normalization / Initialization
7	10/13, 10/15	Optimization, Learning Rates CNNs
8	10/20, 10/22	Review, <b>Midterm Exam</b>

<b>Module 3: Intro to Generative Models</b>			
9	10/27, 10/29	A Linear Intro to Generative Models, Factor Analysis, Autoencoders, VAEs	
10	11/3, 11/5	Generative Adversarial Networks, Diffusion Models	Project Midway Report
<b>Module 4: Large Language Models</b>			
11	11/10, 11/12	Sequence Learning with RNNs Attention, Transformers	HW4
	11/17, 11/19	GPT Architectures, Unsupervised Training of LLMs	
13	11/24, 11/26	Supervised Fine-tuning of LLMs, Prompts and In-context learning	HW5
14	12/1,	Foundation models, alignment, explainability Open directions in LLM research	
15	12/8, 12/10	<b>Project Presentations</b>	Project Final Report
16	12/17	<b>Final Exam</b>	Final Exam

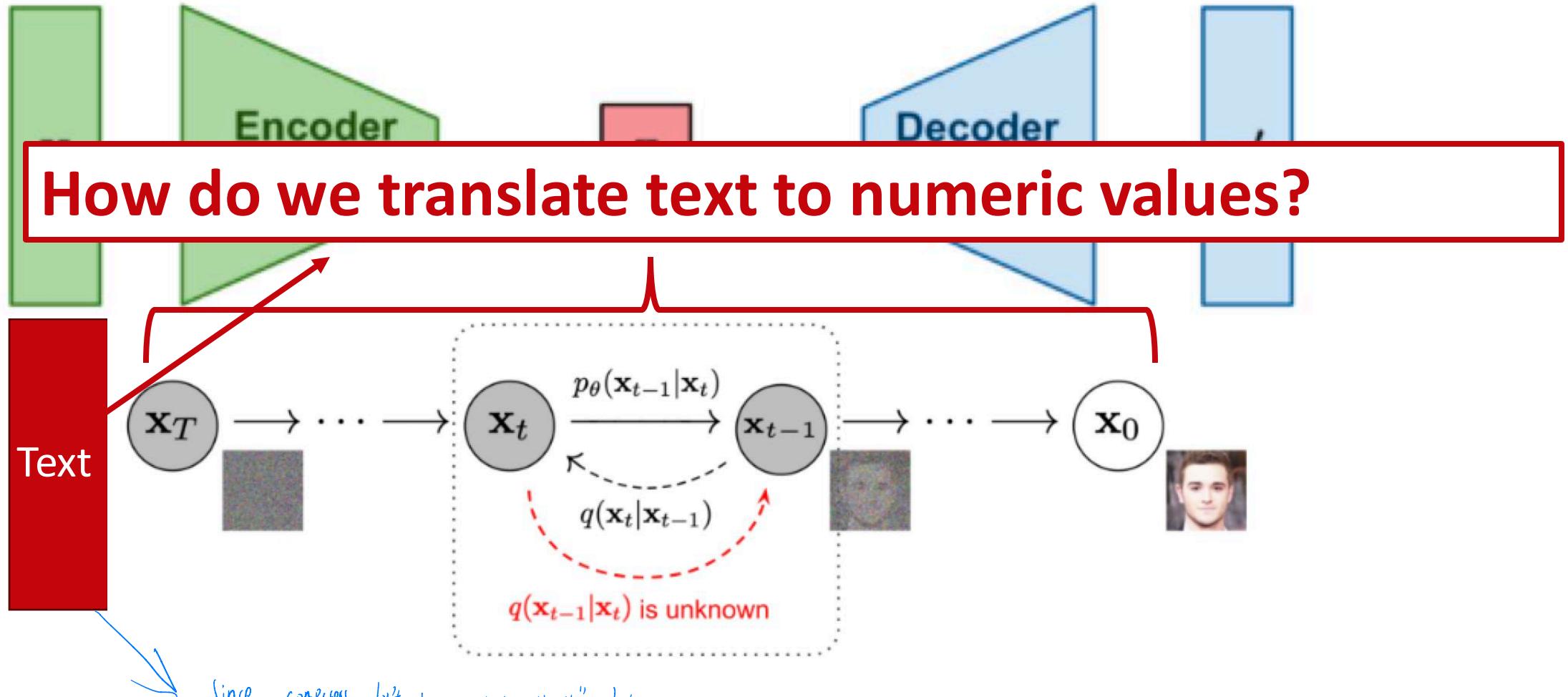


# HW4

---

- Released on [the website](#)
- Due next Friday
- Auto-encoder (4 parts) + bonus GAN
- We recommend Colab

# Recall “Conditioning on Text”...





# Challenges with text

- Variable length input
- Long-range dependencies *(Depends on previous words)*
- Want a generative model
- Scale
- Emergent properties of large language models

Module 4: Large Language Models			
11	11/10, 11/12	Sequence Learning with RNNs Attention, Transformers	HW4
12	11/17, 11/19	GPT Architectures, Unsupervised Training of LLMs	
13	11/24, 11/26	Supervised Fine-tuning of LLMs, Prompts and In-context learning	HW5
14	12/1, 12/3	Foundation models, alignment, explainability Open directions in LLM research	
15	12/8, 12/10	Project Presentations	Project Final Report
16	12/17	Final Exam	Final Exam

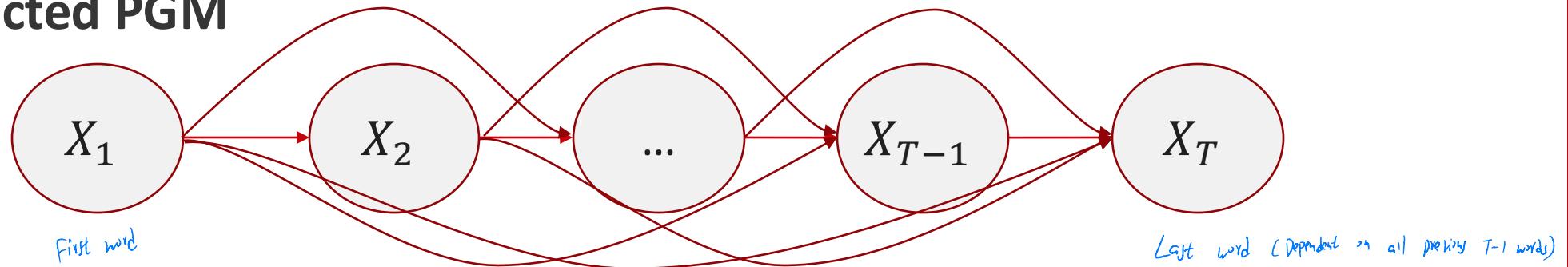
↓ (when the model is so big that it has the capability of doing something we did not train on)

# Where we're going

自回归 (自己预测自己的下一步，根据前面的话猜下一个字是什么)

GPT = Auto-regressive probabilistic model

- Directed PGM



$$P_\theta(X) = \prod_i \prod_t P_\theta(X_{i,t} | X_{i,<t})$$

↓  
 sentence  
 ↓  
 t's word

在 t 之前的的所有字

- Probabilistic objective: Max log-likelihood of observed seqs

$$\max_\theta \sum_i \sum_t \log P_\theta(X_{i,t} | X_{i,<t}) \rightarrow \text{让模型猜对正确单词的概率尽可能大}$$

[Radford et al., [Improving Language Understanding by Generative Pre-Training](#)]



# Today

---

- **Different Ways to Model Text**
- Sequence Modeling with RNNs
- Different Types of Sequence Modeling Tasks
- Backpropagation Through Time
- Long-Short Term Memory (LSTM)
- Many-to-one Word RNNs

Turn text into numeric values



# A classic approach: Bag-of-words

"Raw" training dataset

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

$x^{[2]}$  = "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 = {  
'and': 0,  
'is': 1  
'one': 2,  
'shining': 3,  
'sun': 4,  
'sweet': 5,  
'the': 6,  
'two': 7,  
'weather': 8,  
}

句子里的所有词拿出来按顺序列成一簇

Training set as design matrix

每句里9个词依次出现了几次

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

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

class labels

$X^{[3]}$  就被压缩成  
现在的这样的 9-D vector 3

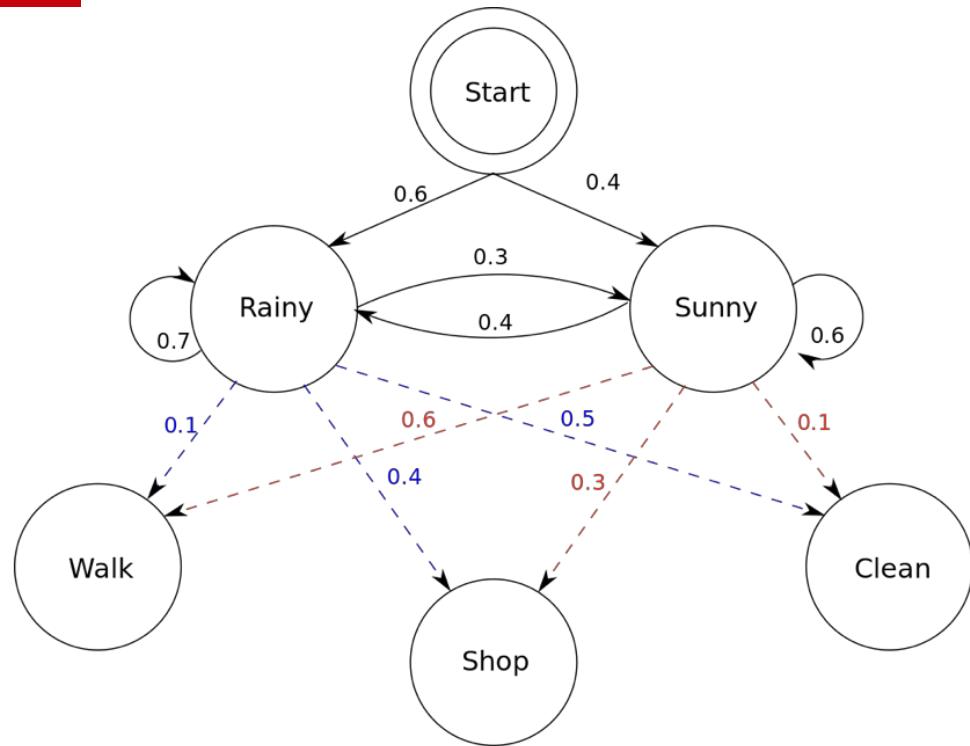
training

Classifier

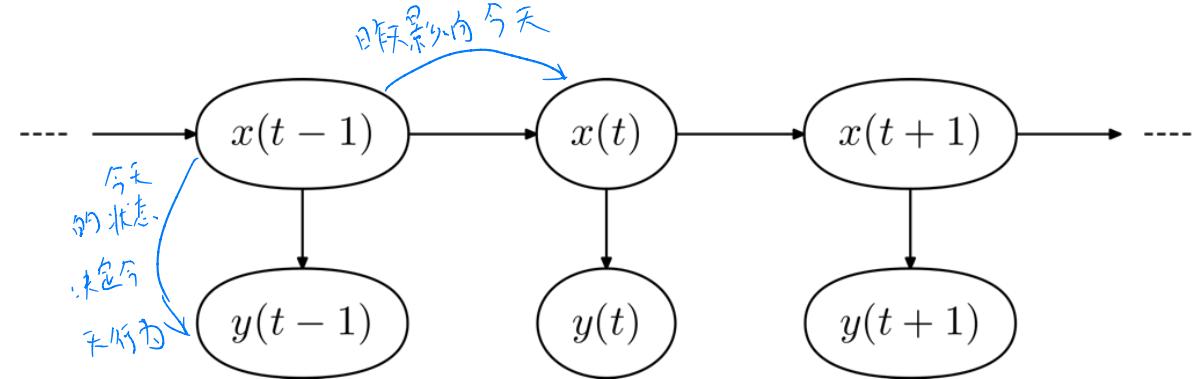
(e.g., logistic regression, MLP, ...)

(★ Drawback: Lose the ordering information of the words)

# Another classic approach: Hidden Markov Model



(引入了时间和因果)



Wikipedia example: each day, weather follows a Markov chain, and activities are observables

$$\mathbb{P}(Y_n = y | X_1 = x_1, \dots, X_n = x_n) = \mathbb{P}(Y_n = y | X_n = x_n)$$

$P(Y_n = y)$  取决于  $x_1, \dots, x_n$  的所有状态,

简化计算量

↳ 可以简化成当下的状态对当下的行为的影响

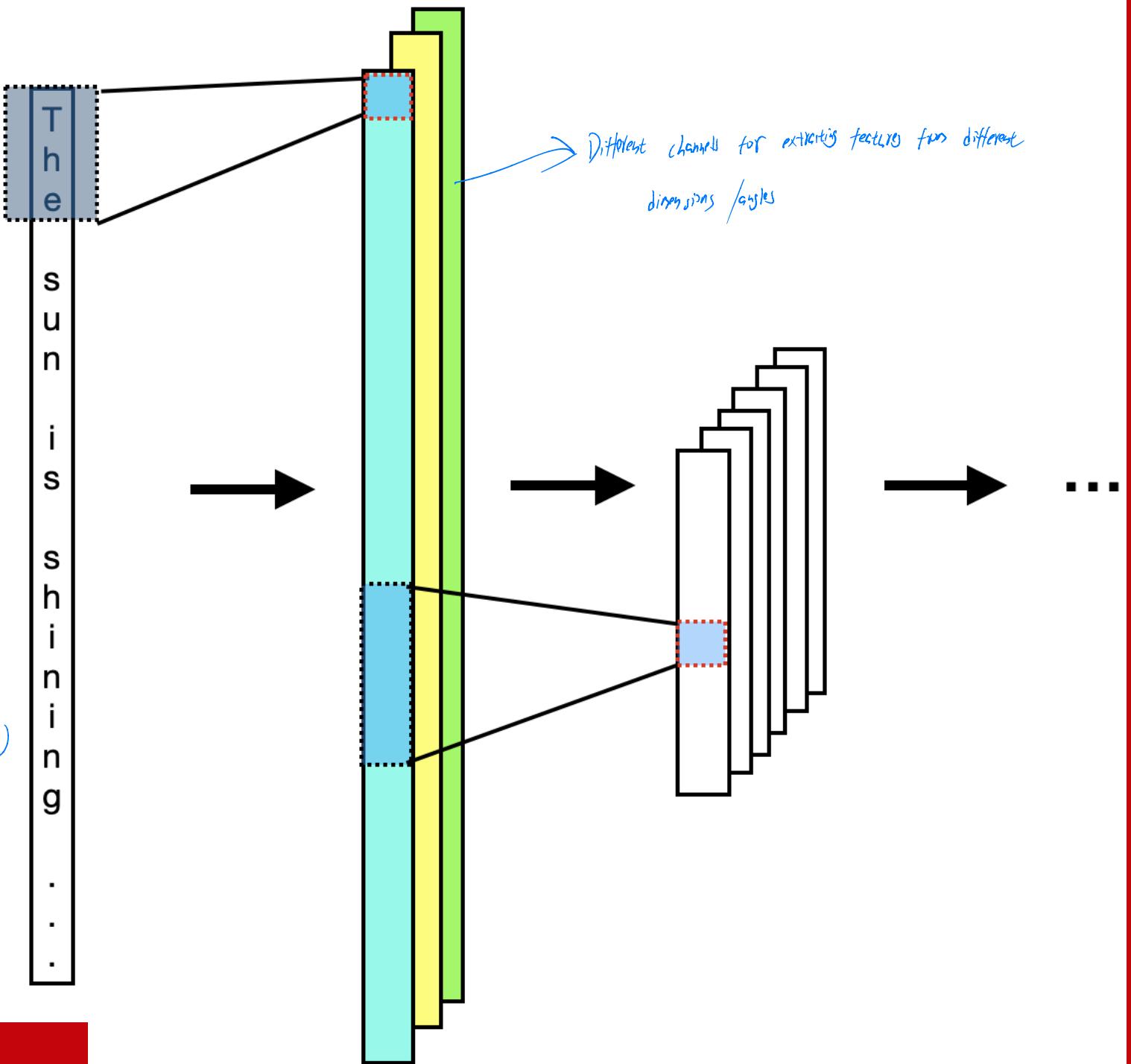
# Another approach:

## CNNs

Type text as a long image

Can't handle variable length input,  
→ need padding to max input length

CNN requires fixed input length (sentence can be short)  
we need to add 0s as padding to the end





# Today

---

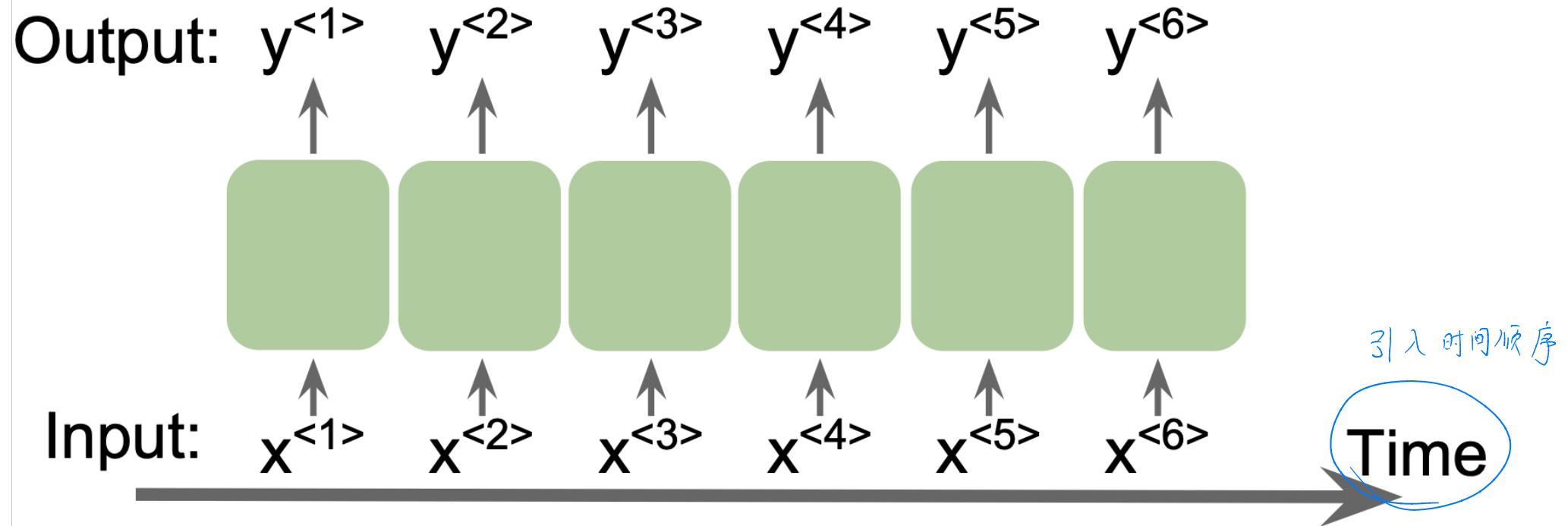
- Different Ways to Model Text
- **Sequence Modeling with RNNs**
- Different Types of Sequence Modeling Tasks
- Backpropagation Through Time
- Long-Short Term Memory (LSTM)
- Many-to-one Word RNNs



# Sequence data: order matters

The movie my friend has not seen is good

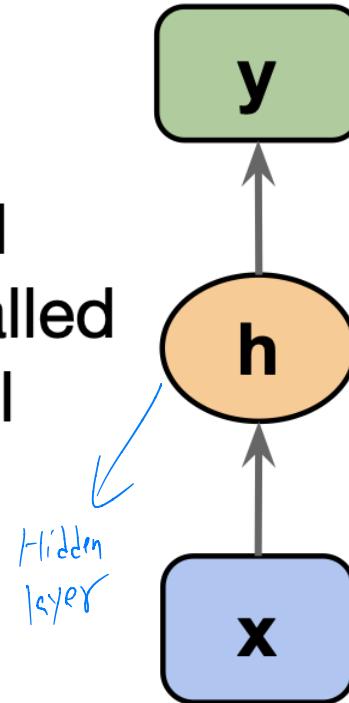
The movie my friend has seen is not good



# Recurrent Neural Networks (RNNs)

Networks we used previously: also called feedforward neural networks

↓  
Unable to process continuous stuff



Recurrent Neural Network (RNN)

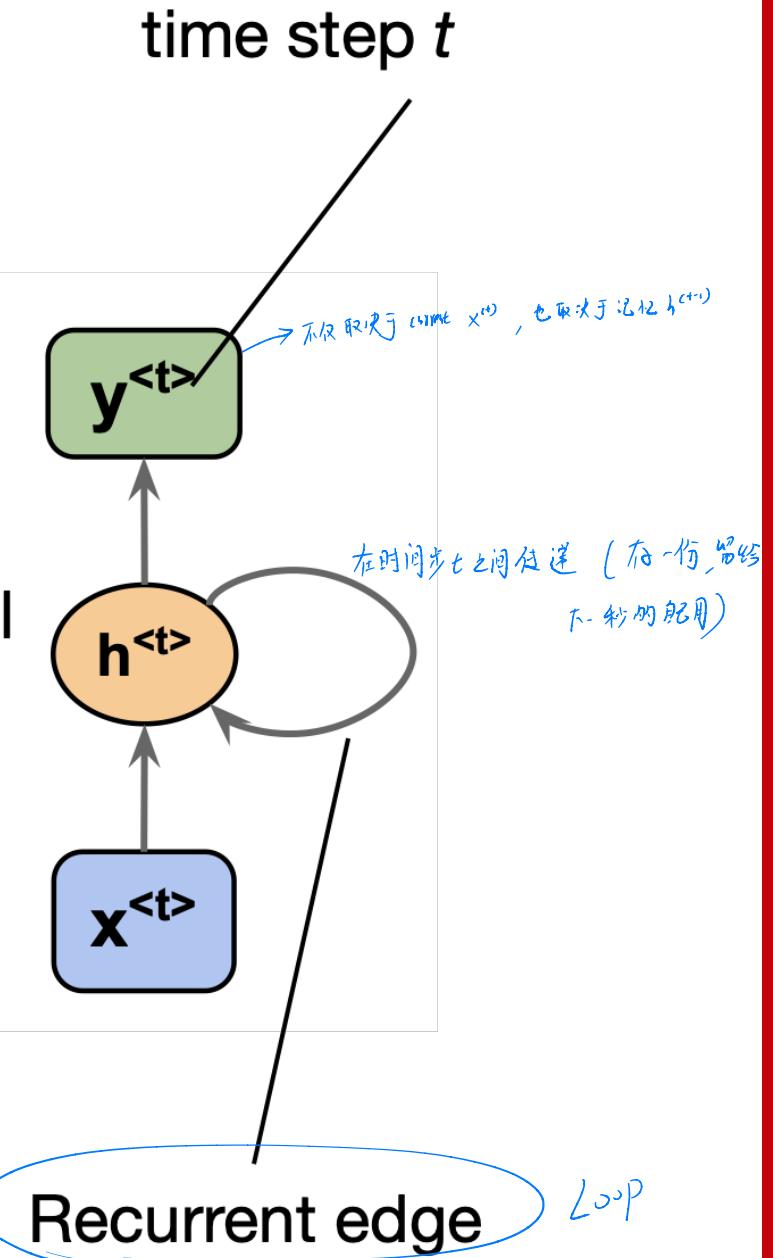


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

# Recurrent Neural Networks (RNNs)

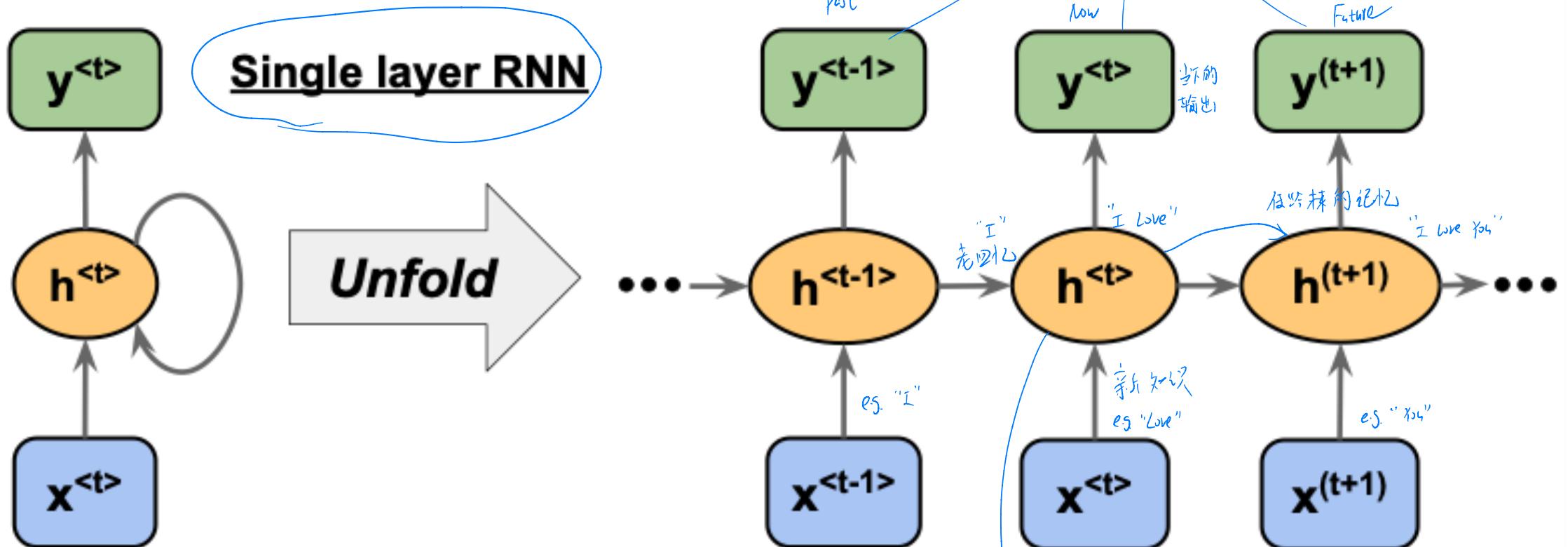


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

# Recurrent Neural Networks (RNNs)

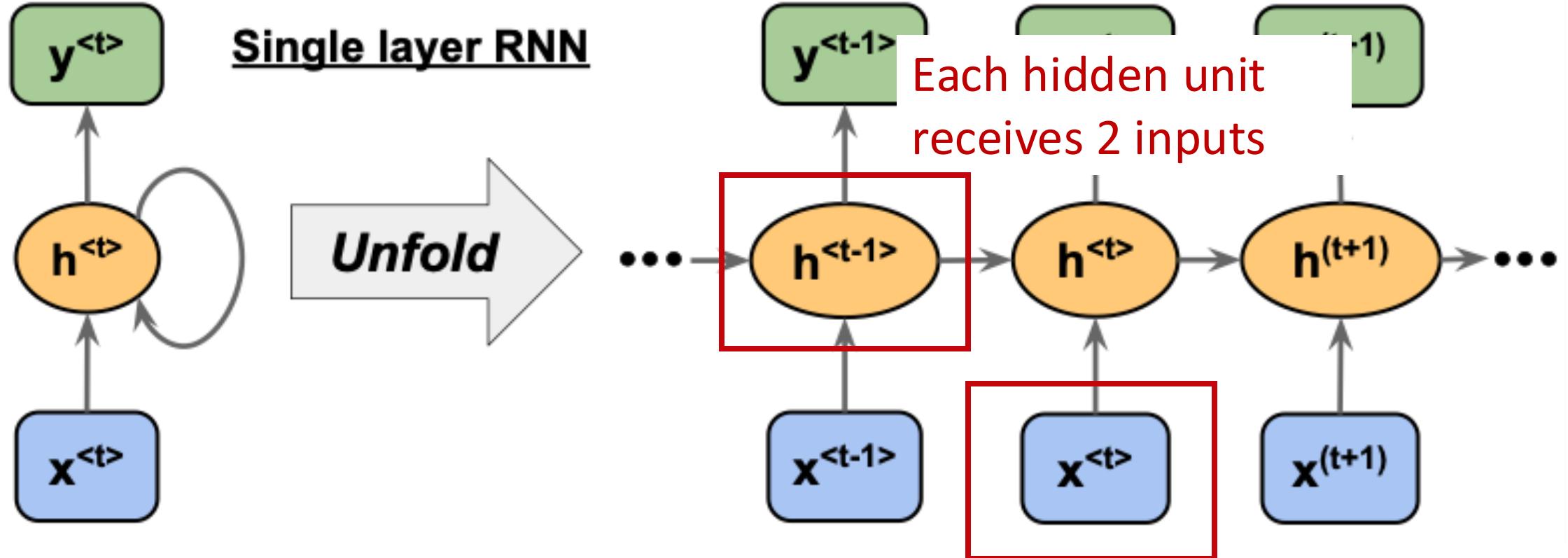


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

# Multilayer RNNs

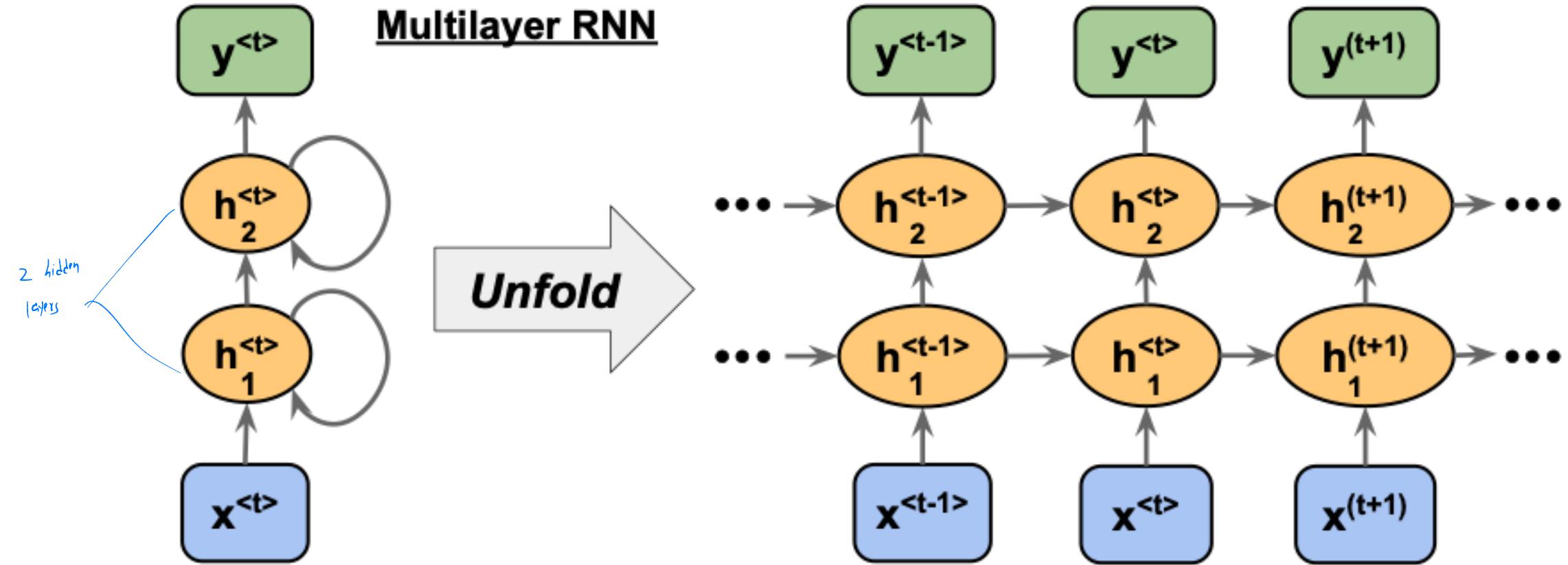


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

( ↗ From bottom to top: Extract features, understand more complex patterns  
 ↙ from left to right: Pass memories )

# Recurrence unlocks many types of sequence tasks

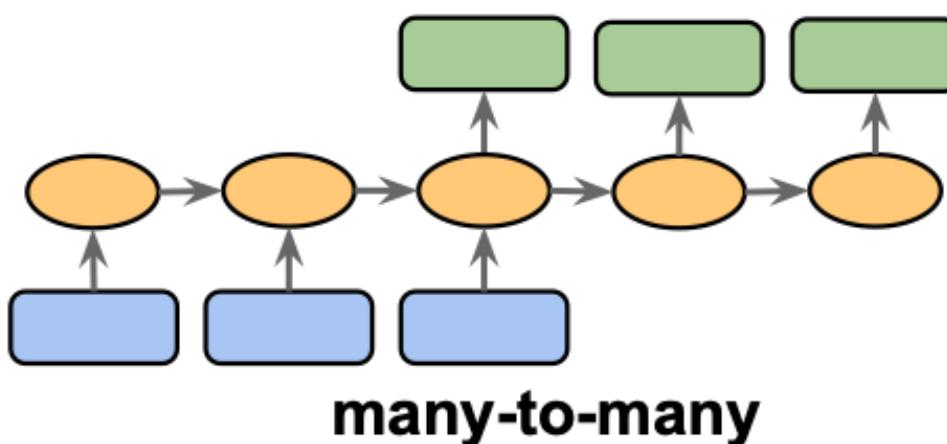
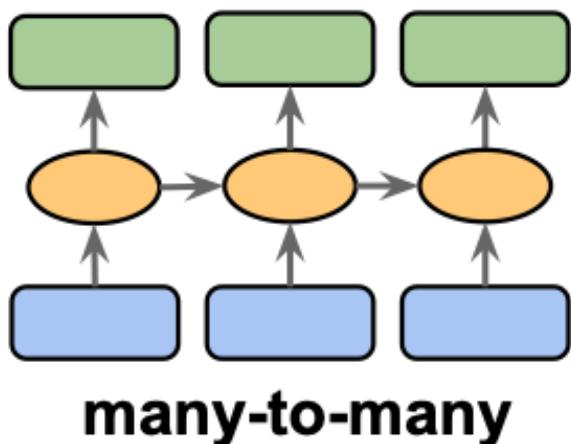
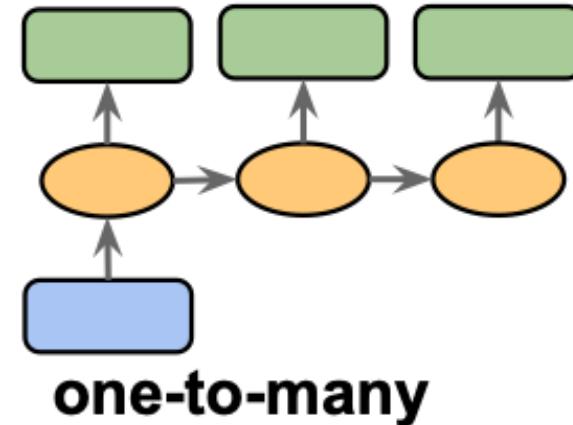
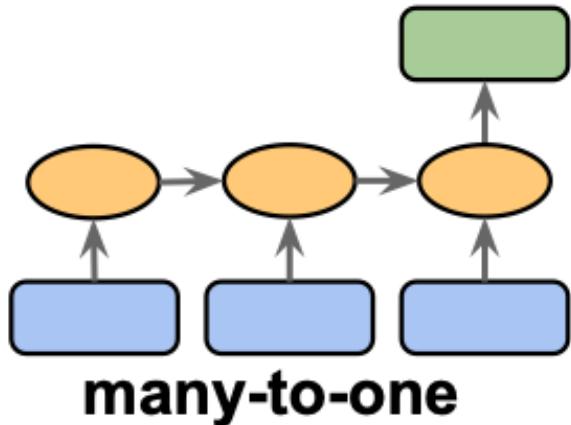


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019

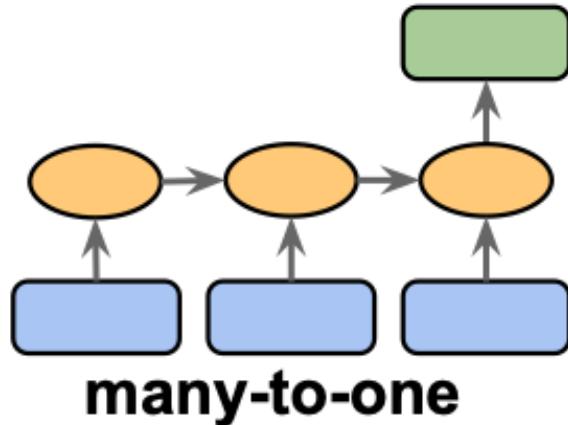


# Today

---

- Different Ways to Model Text
- Sequence Modeling with RNNs
- **Different Types of Sequence Modeling Tasks**
- Backpropagation Through Time
- Long-Short Term Memory (LSTM)
- Many-to-one Word RNNs

# Recurrence unlocks many types of sequence tasks

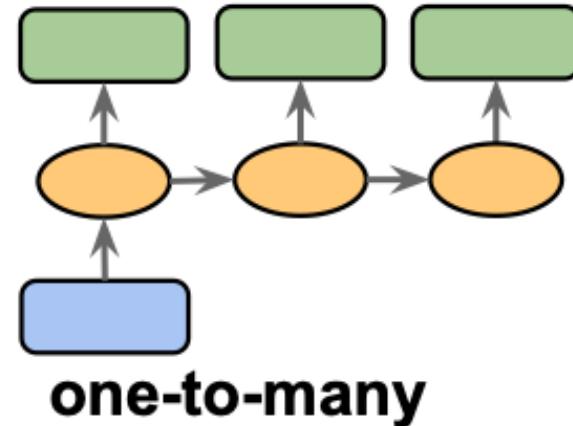


**Many-to-one:** The input data is a sequence, but the output is a fixed-size vector, not a sequence.

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

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

# Recurrence unlocks many types of sequence tasks



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

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

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

# Recurrence unlocks many types of sequence tasks

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

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

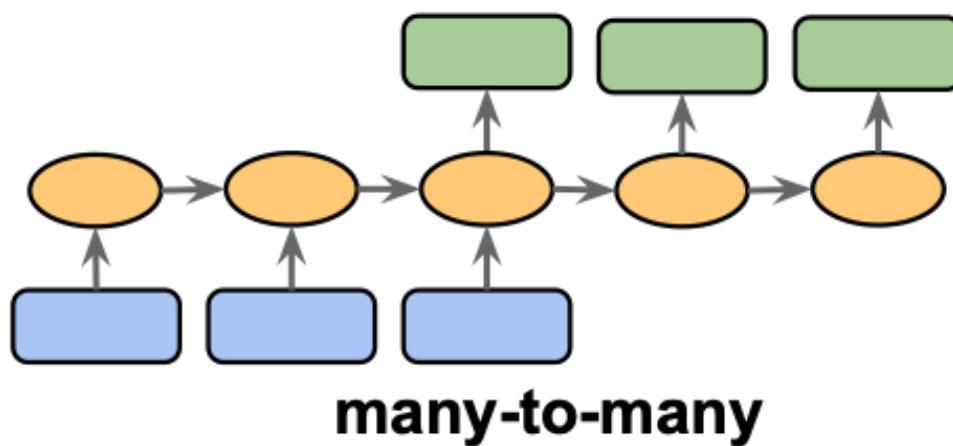
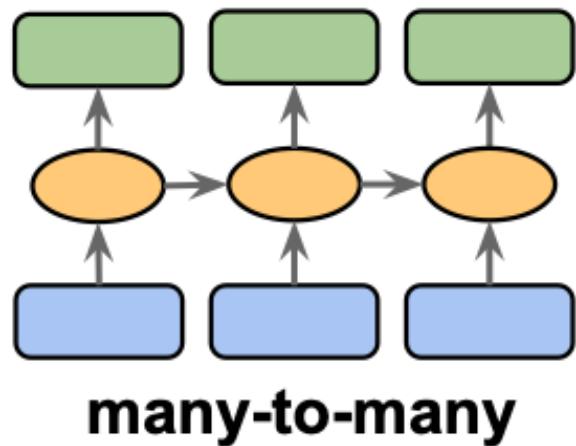


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019



# Today

---

- Different Ways to Model Text
- Sequence Modeling with RNNs
- Different Types of Sequence Modeling Tasks
- **Backpropagation Through Time**
- Long-Short Term Memory (LSTM)
- Many-to-one Word RNNs



# Under the hood: weight matrices in an RNN

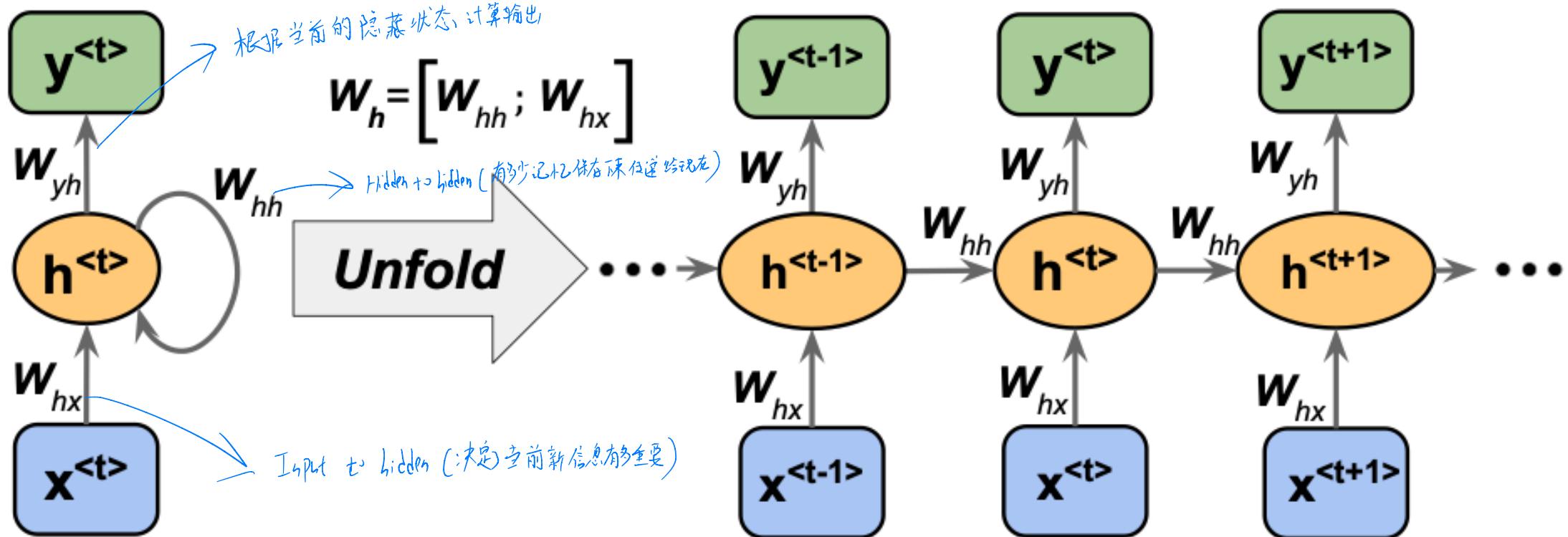
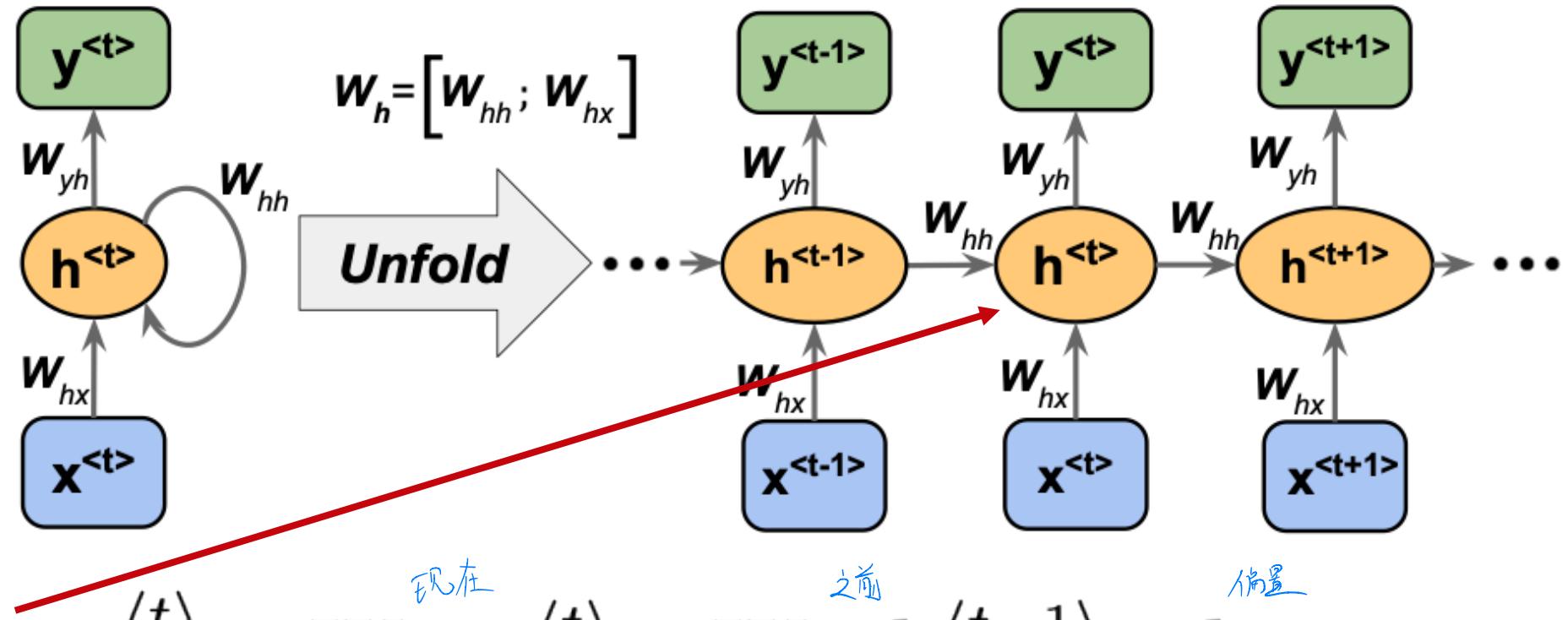


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019



# Under the hood: weight matrices in an RNN

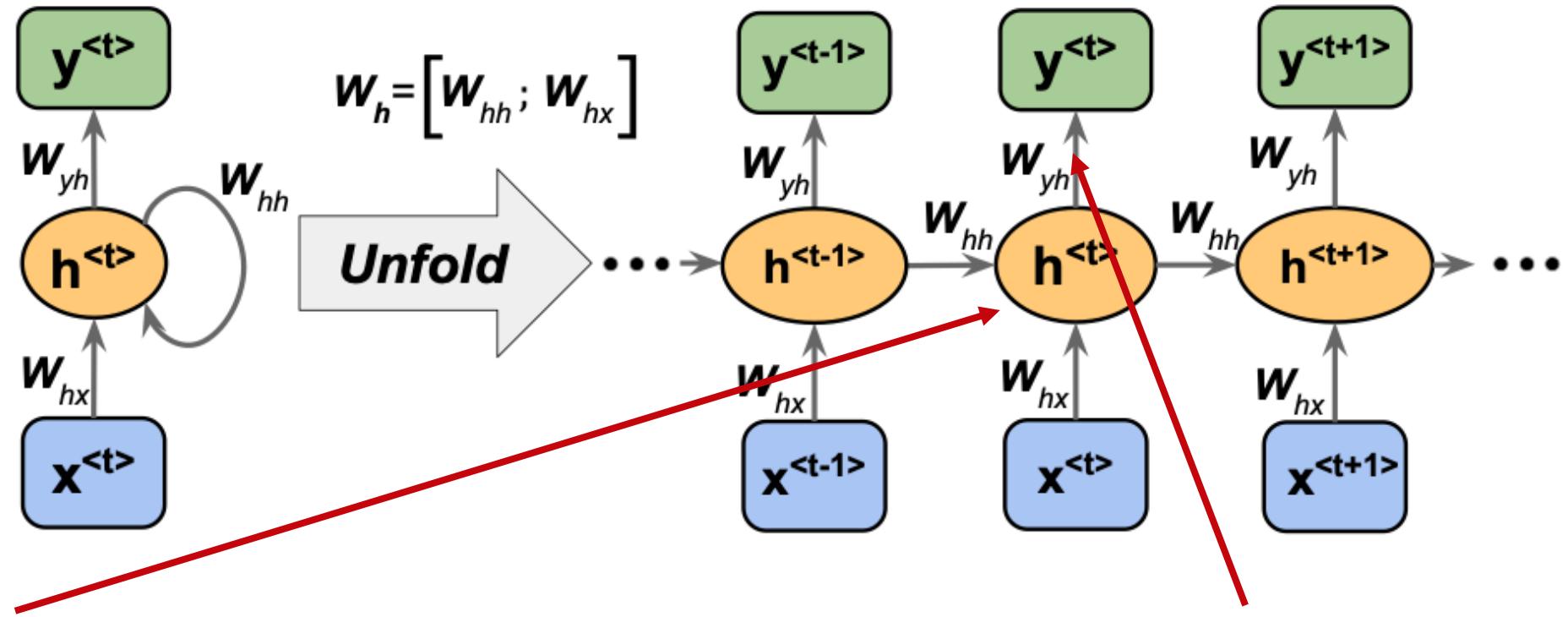


Net input:  $\mathbf{z}_h^{(t)} = \mathbf{W}_{hx} \mathbf{x}^{(t)} + \mathbf{W}_{hh} \mathbf{h}^{(t-1)} + \mathbf{b}_h$

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

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

# Under the hood: weight matrices in an RNN



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

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

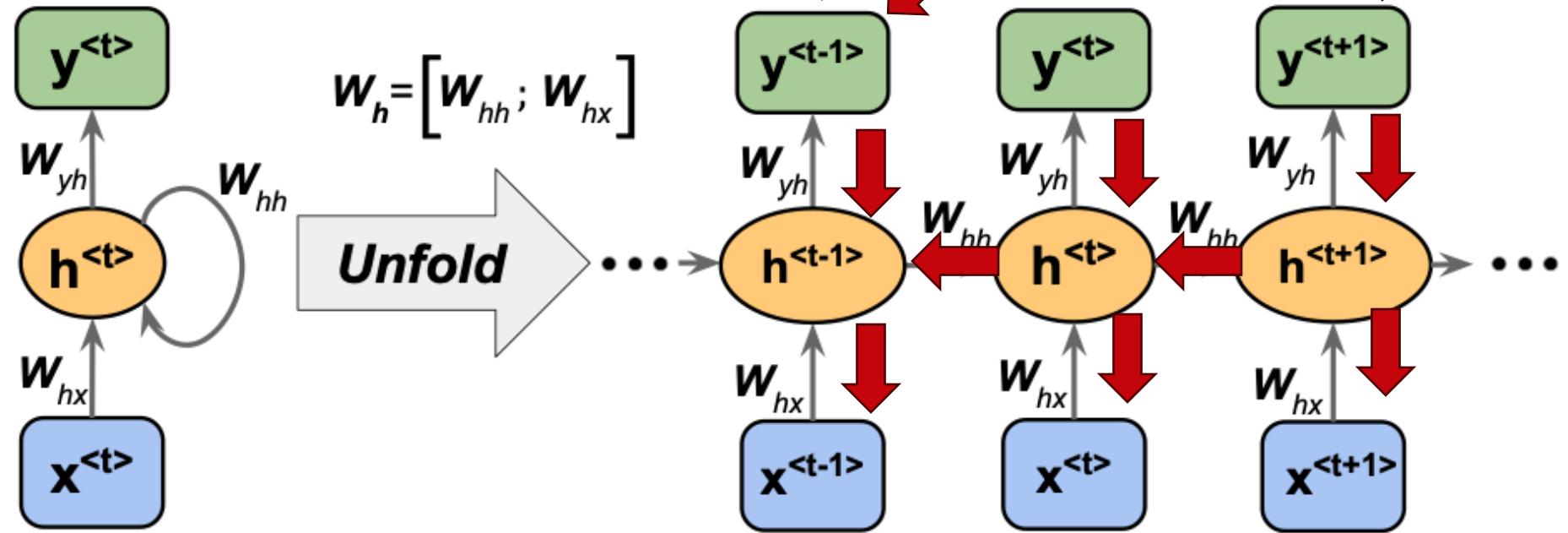
$$\text{Activation: } \mathbf{h}^{(t)} = \sigma_h(\mathbf{z}_h^{(t)})$$

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

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

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

(*每个时间步的  
错误总和*)



# Backpropagation through time

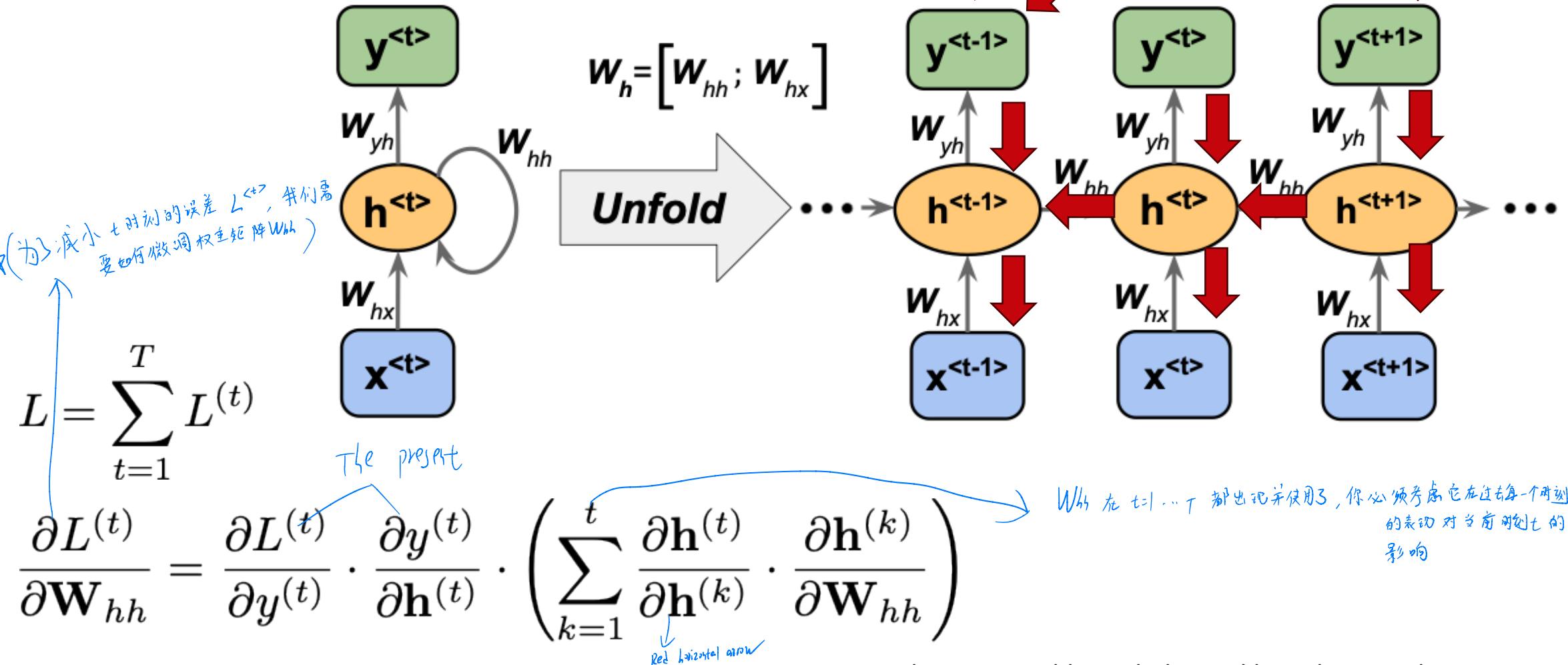
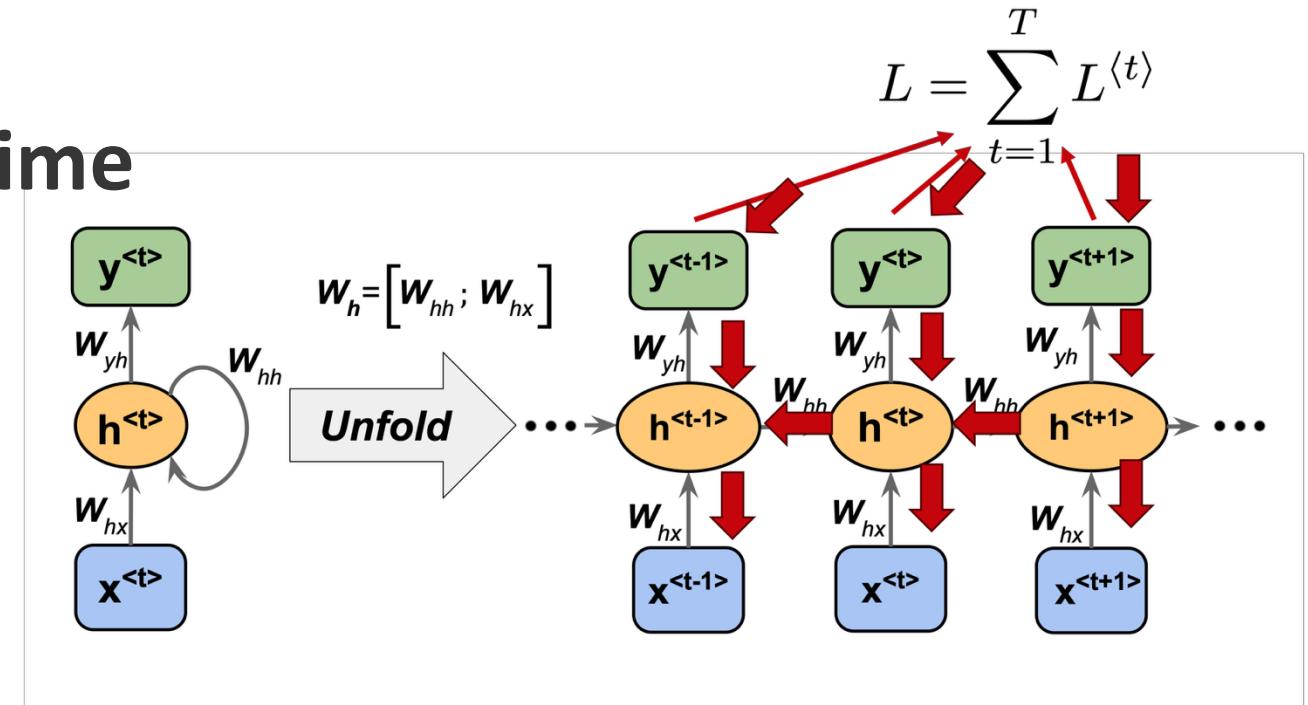


Image source: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Packt, 2019



# Backpropagation through time



Computed as a multiplication of adjacent time steps:

(每个相邻时间步的导数相乘)

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

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

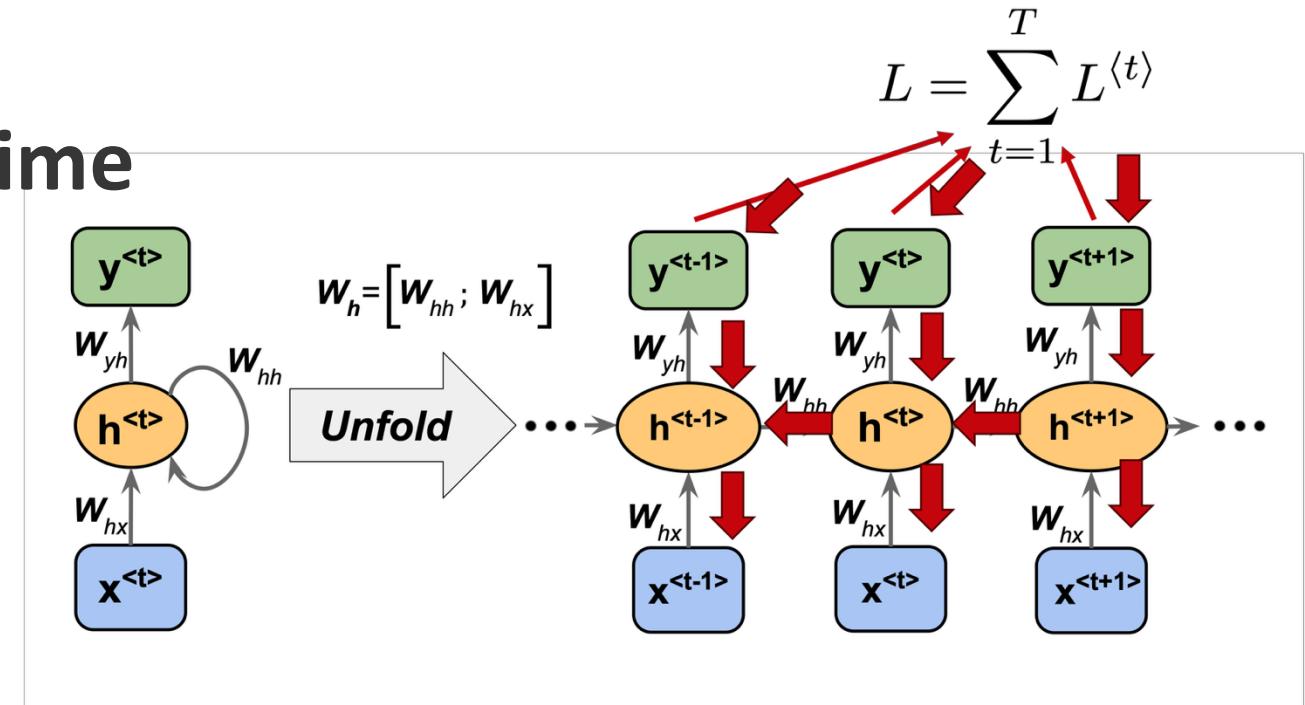
(第*k*秒的变化对第*t*秒造成的影响)

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

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

# Backpropagation through time

Straightforward, but problematic:  
vanishing / exploding gradients!



Computed as a multiplication of adjacent time steps:

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

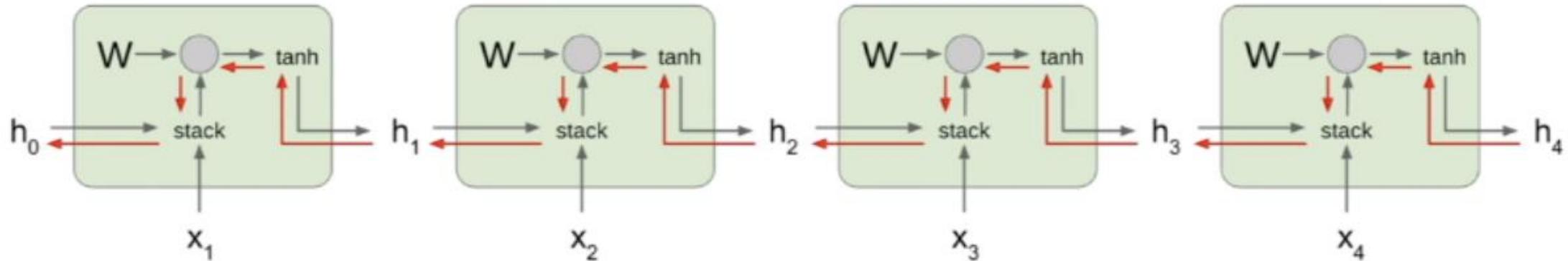
$$\frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left( \sum_{k=1}^t \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}}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

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

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

# A challenge: Vanishing / exploding gradients

$$\mathbf{h}_t = \tanh(W^{hh}\mathbf{h}_{t-1} + W^{hx}\mathbf{x}_t)$$



Computing gradient of  $\mathbf{h}_0$  involves many factors of  $\mathbf{W}$  (and repeated  $\tanh$ )

Largest singular value  $> 1$ :  
**Exploding gradients**

Largest singular value  $< 1$ :  
**Vanishing gradients**

**Gradient clipping:** Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```

*Bengio et al., 1994 “Learning long-term dependencies with gradient descent is difficult”*  
*Pascanu et al., 2013 “On the difficulty of training recurrent neural networks”*



# Solutions to Vanishing / Exploding Gradients

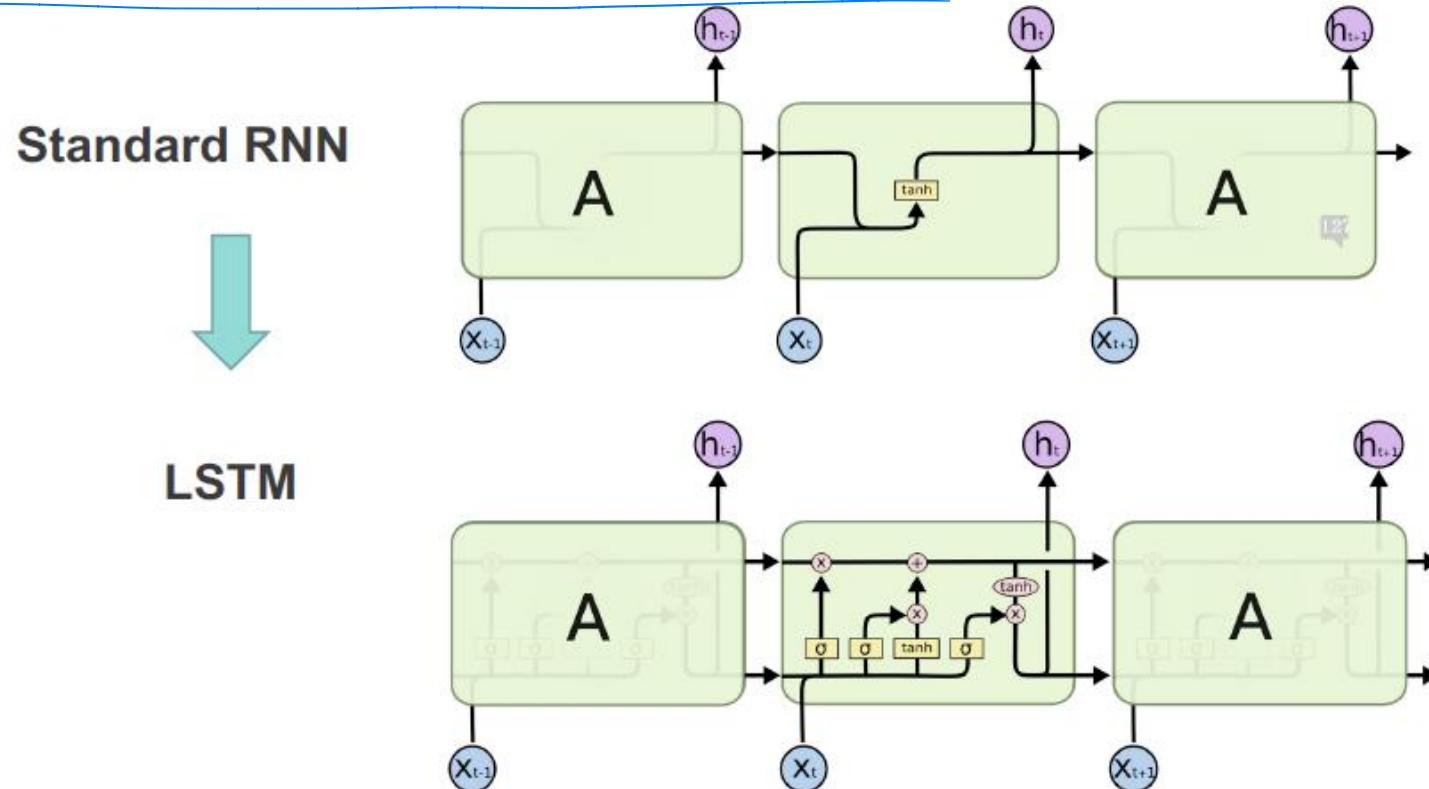
- **Gradient Clipping:** set a max value for gradients if they grow to large (solves only exploding gradient problem)
- **Truncated backpropagation through time (TBPTT):** limit the number of time steps the signal can backpropagate after each forward pass. E.g., even if the sequence has 100 elements/steps, we may only backpropagate through 20 or so.



人为规定双向传播时，每一步只把误差传回前 20 步，再往前直接截断

# Solutions to Vanishing / Exploding Gradients

**Long short-term memory (LSTM): uses a memory cell for modeling long-range dependencies and avoid vanishing gradient problems**





# Today

---

- Different Ways to Model Text
- Sequence Modeling with RNNs
- Different Types of Sequence Modeling Tasks
- Backpropagation Through Time
- **Long-Short Term Memory (LSTM)**
- Many-to-one Word RNNs

# Long-short term memory (LSTM)

- Not an oxymoron: **2 paths of memory**

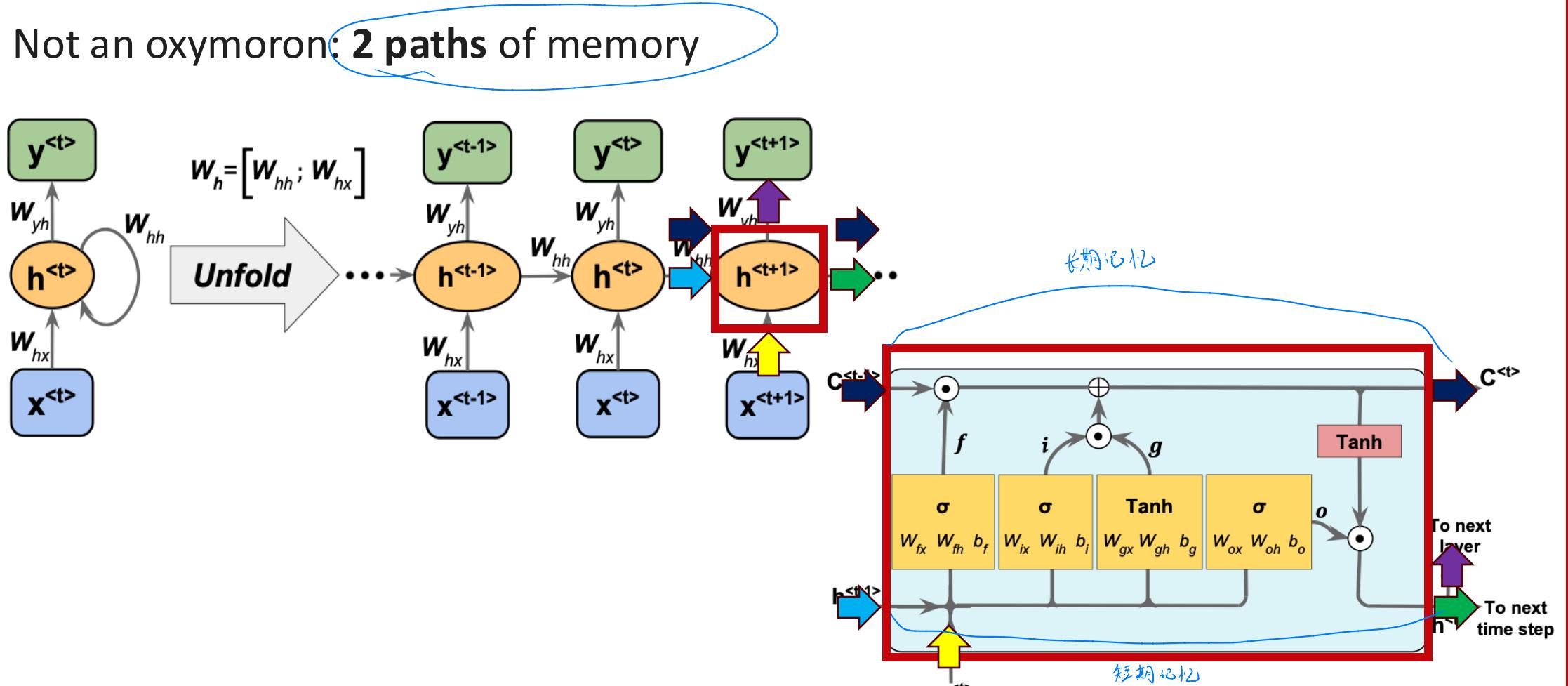


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Birmingham, UK: Packt Publishing, 2019.

# Long-short term memory (LSTM)

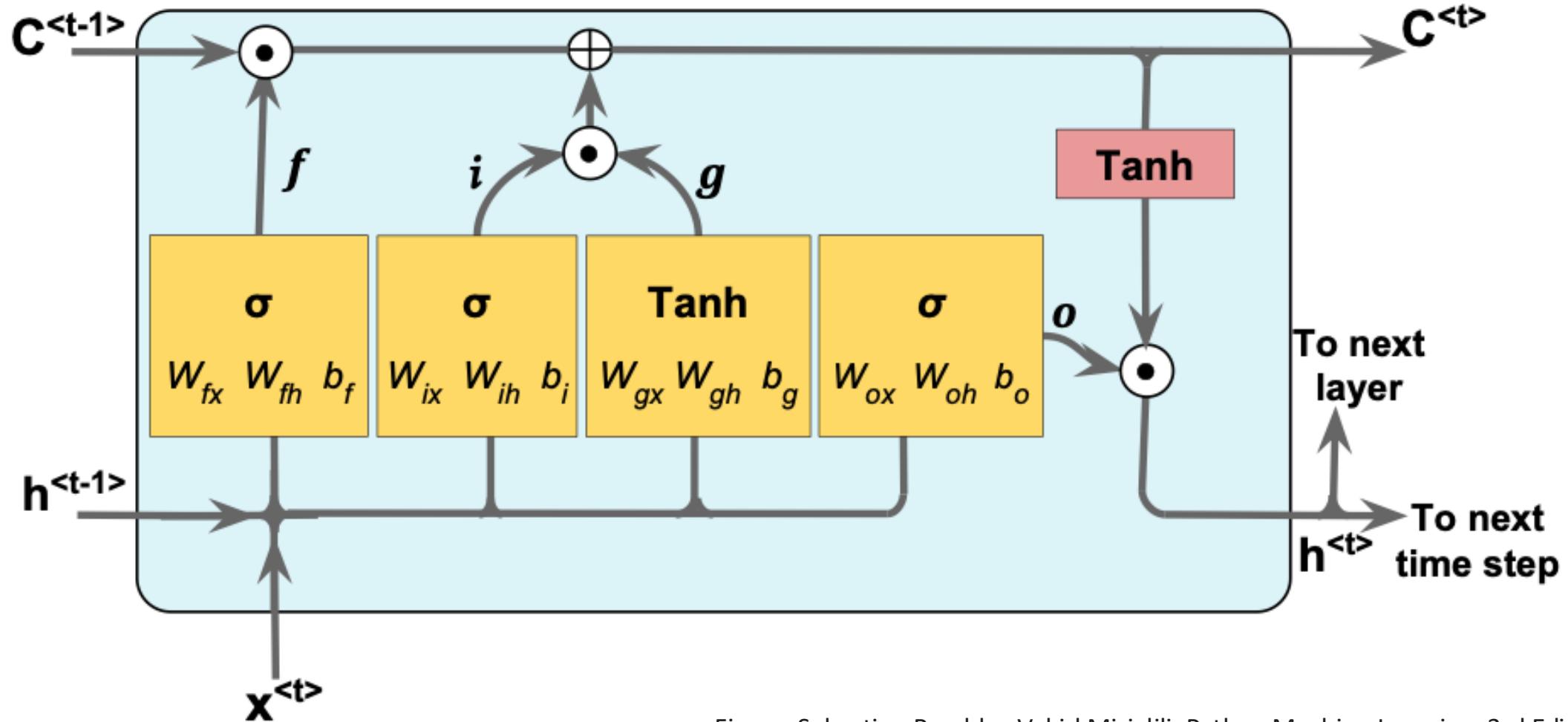


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019

# Inside LSTM

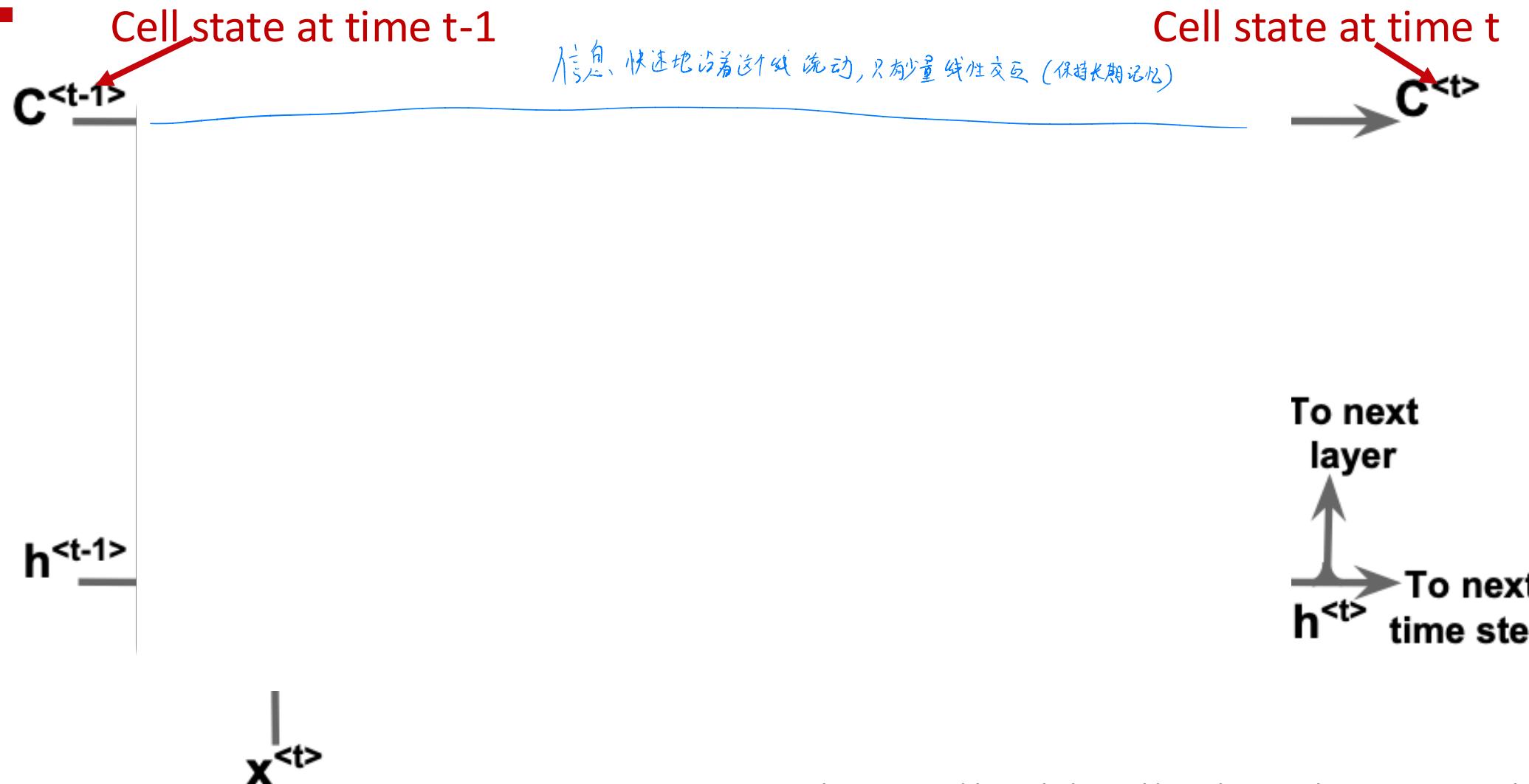


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019

# Inside LSTM

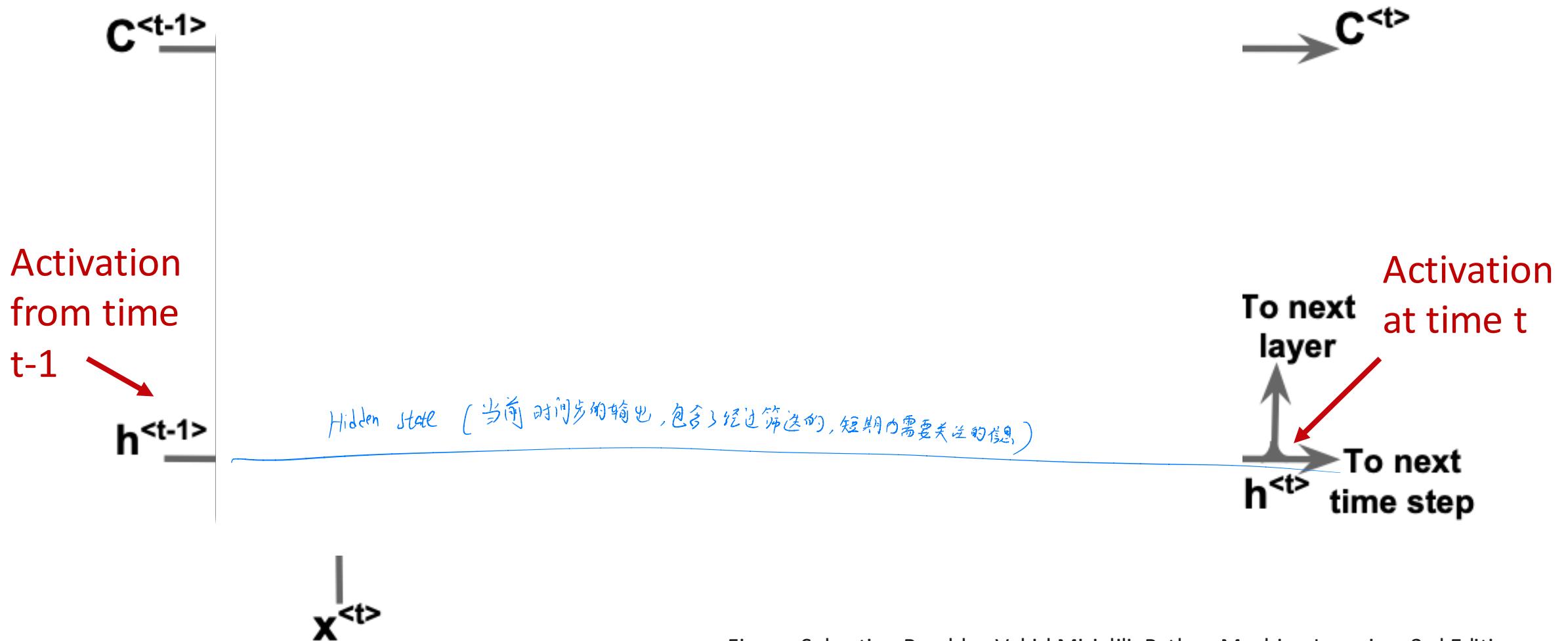
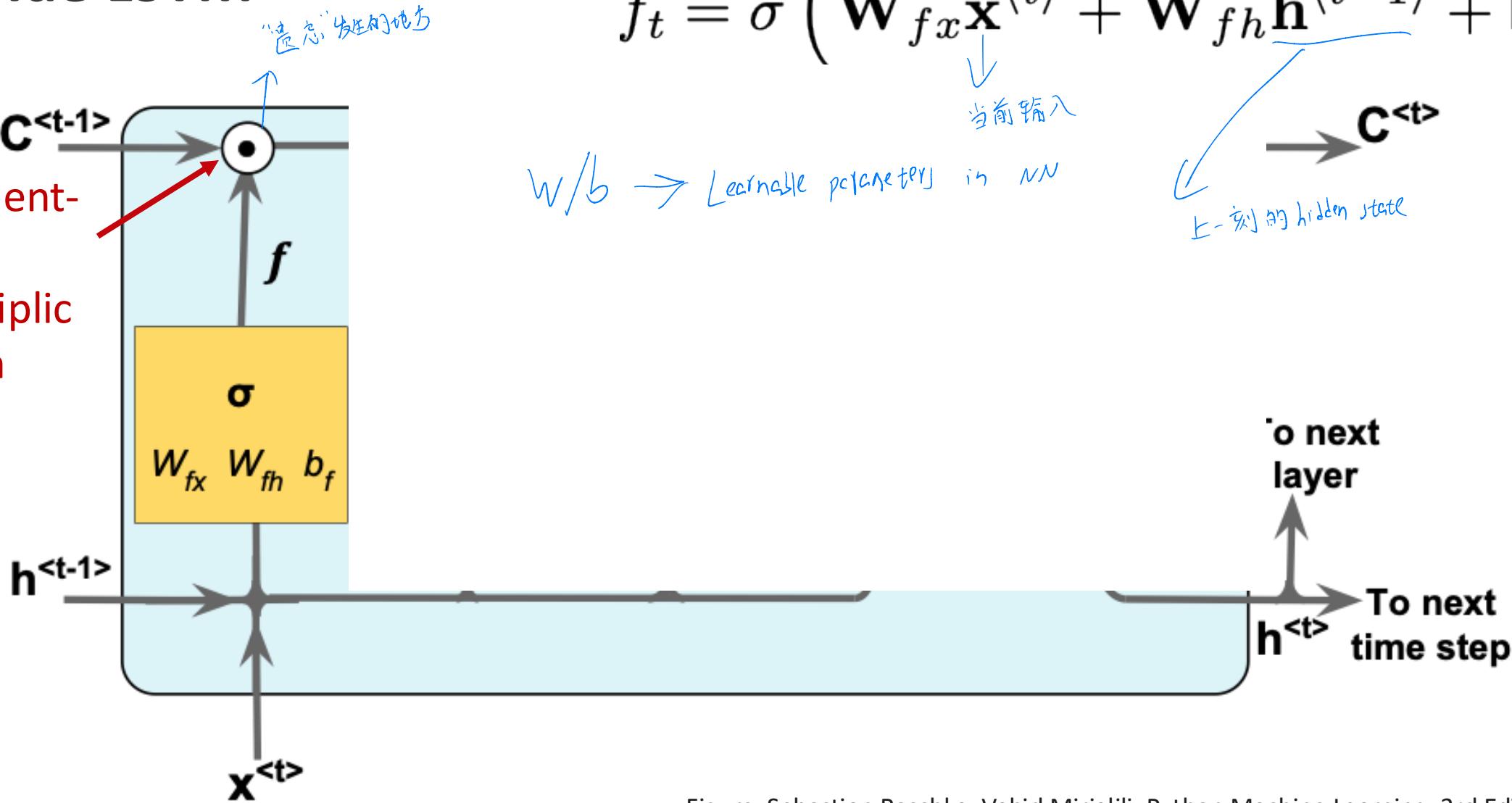


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019

# Inside LSTM

Element-wise multiplication



“Forget gate”: controls which information is remembered and which is forgotten

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

$W/b \rightarrow \text{Learnable parameters in NN}$



# Inside LSTM

决定信息是否需要更新？  
 “Input gate”:  $i_t = \sigma(\mathbf{W}_{ix}\mathbf{x}^{(t)} + \mathbf{W}_{ih}\mathbf{h}^{(t-1)} + b_i)$  → 输入门，决定更新哪里  
 “Input node”:  $g_t = \tanh(\mathbf{W}_{gx}\mathbf{x}^{(t)} + \mathbf{W}_{gh}\mathbf{h}^{(t-1)} + b_g)$  → 输入节点，准备新信息  
 新记忆内容

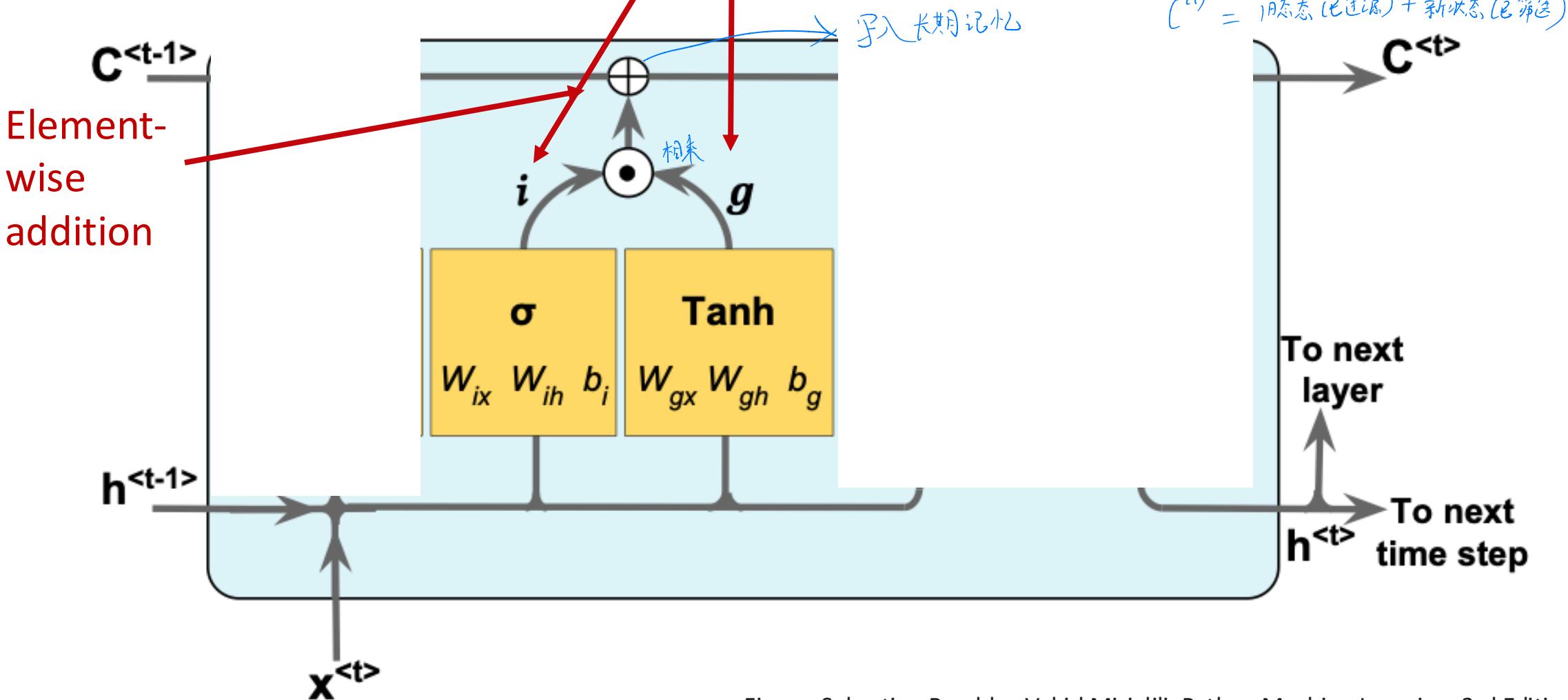


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019

# Inside LSTM

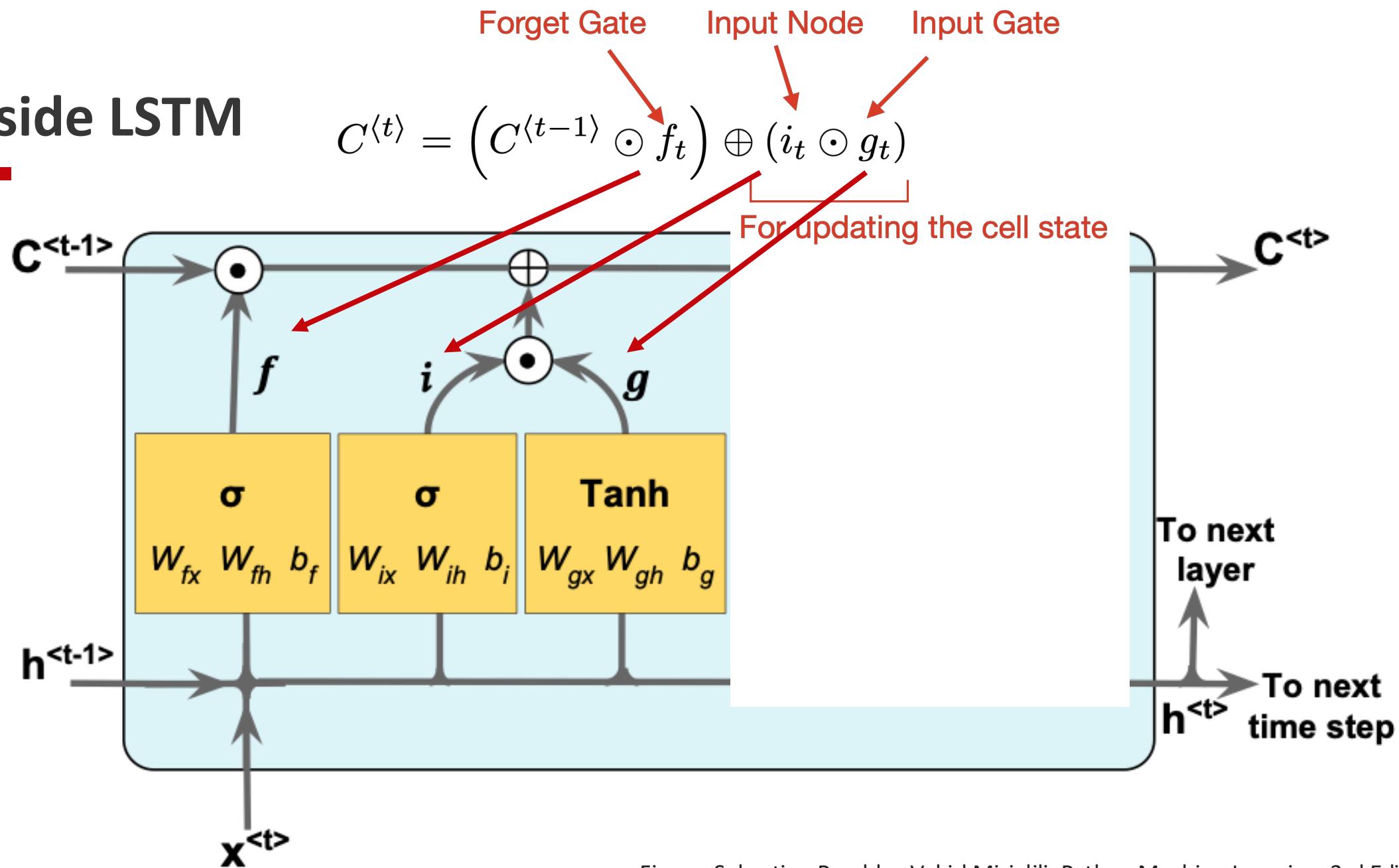


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019

# Inside LSTM

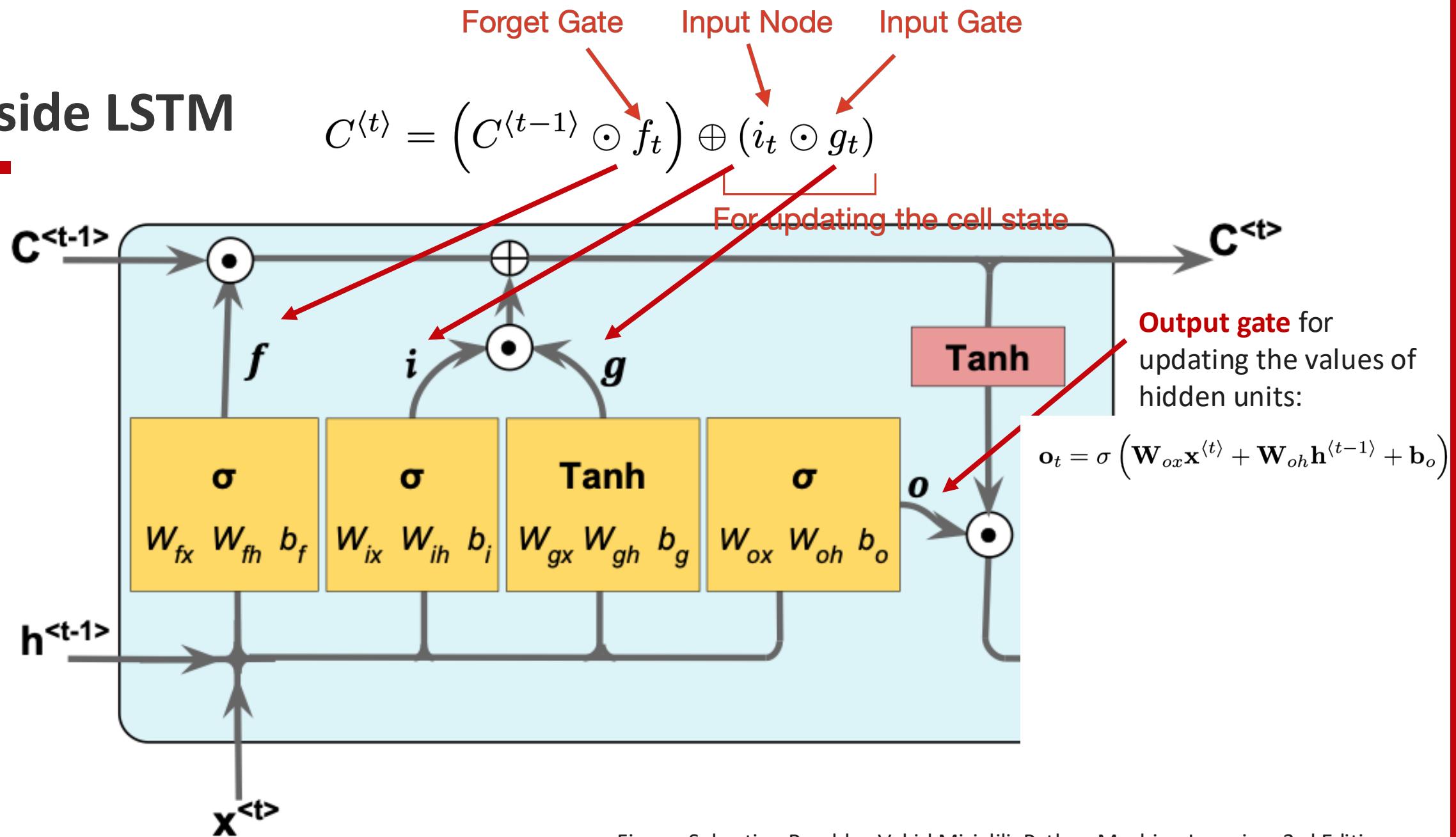


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019



# Inside LSTM

Forget Gate      Input Node      Input Gate

$$C^{(t)} = (C^{(t-1)} \odot f_t) \oplus (i_t \odot g_t)$$

(这是我们现在的全新的记忆)

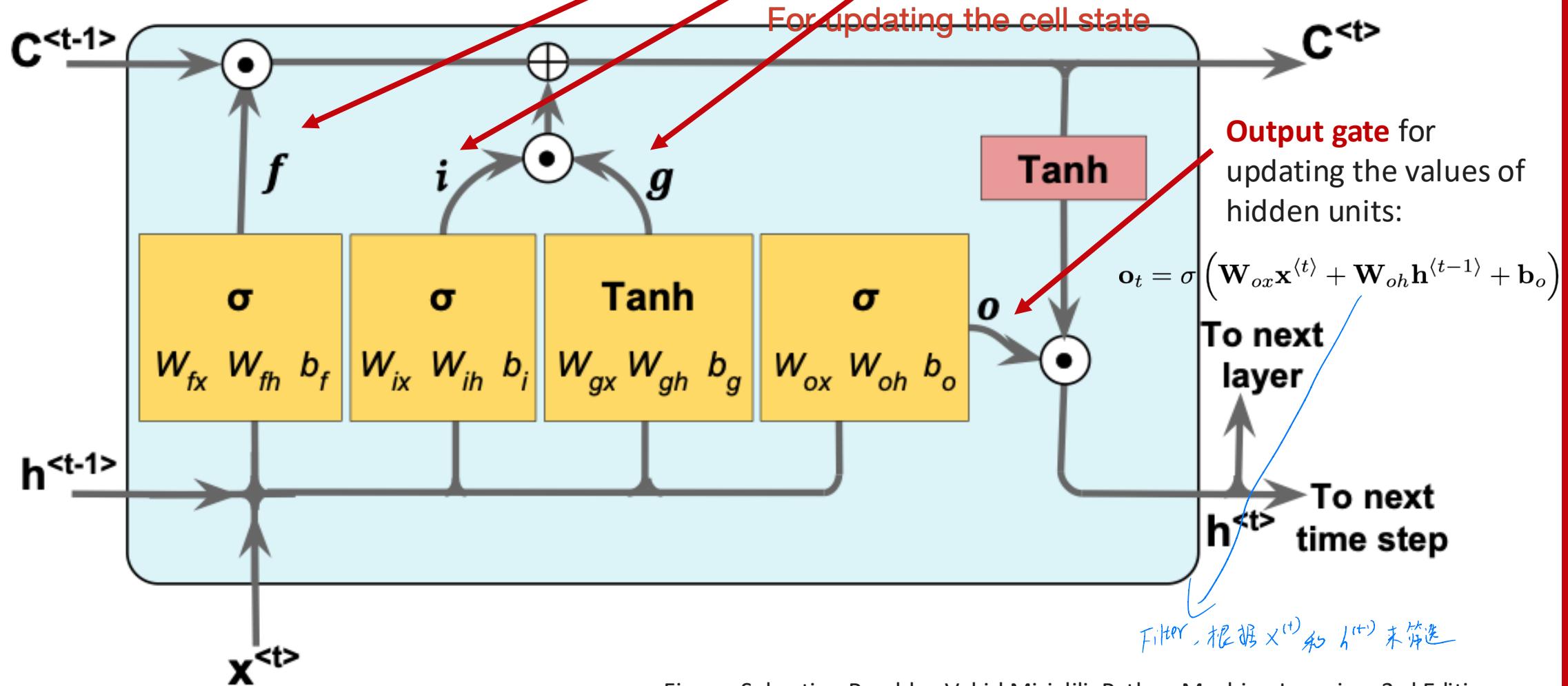


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition.  
Birmingham, UK: Packt Publishing, 2019



# LSTM Back Together

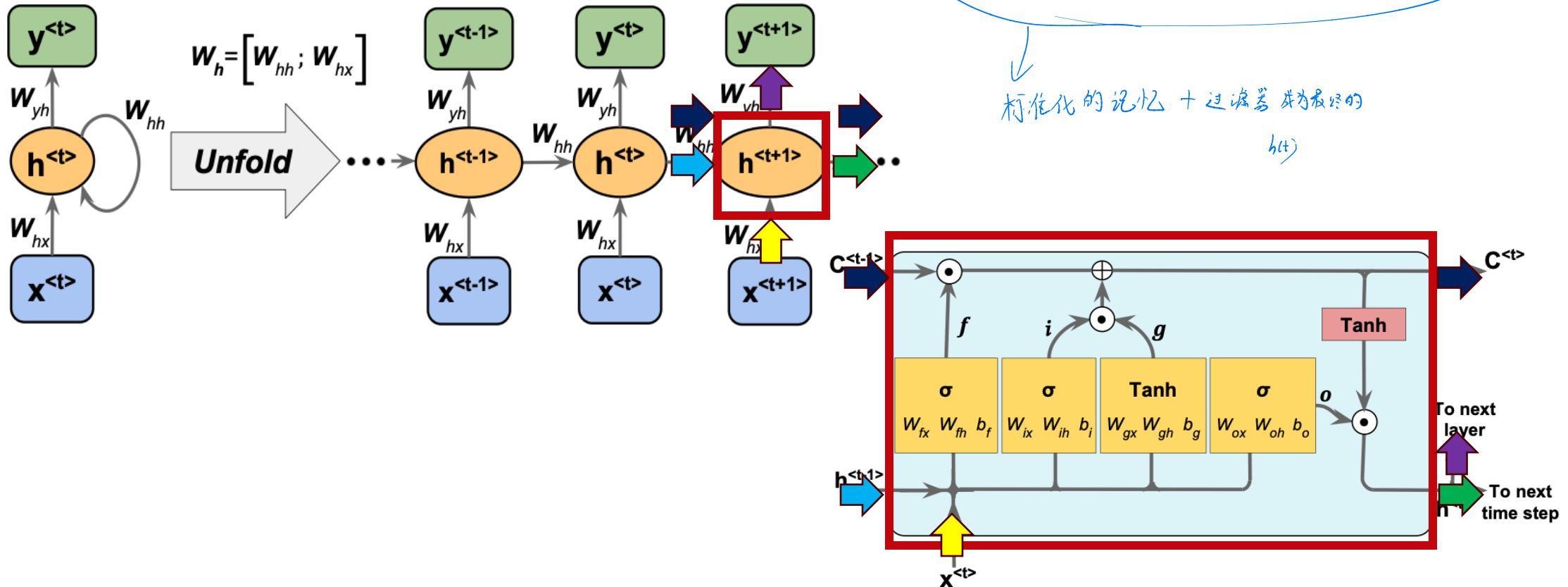


Figure: Sebastian Raschka, Vahid Mirjalili. Python Machine Learning. 3rd Edition. Birmingham, UK: Packt Publishing, 2019

## 0. 开局：两个信号进门

当前时间是  $t$ 。

- $x^{(t)}$  (新信息)：比如当前的单词 "dog"。
- $h^{(t-1)}$  (旧语境)：比如上一刻说出的“动词应该是单数”。
- (注意： $C^{(t-1)}$  此时正在顶部的传送带上等着被处理)

### 1. 遗忘门 (Forget Gate) —— “清理旧账”

- 动作：将  $x^{(t)}$  和  $h^{(t-1)}$  放入 Sigmoid 函数。
- 问题：“根据现在的单词 'dog'，以前记忆里的那些东西（比如关于 'he' 的信息）还有用吗？”
- 结果：生成一个 0 到 1 的向量  $f_t$ 。
- 执行： $C^{(t-1)} \odot f_t$ 。旧记忆里不重要的部分被擦除了。

### 2. 输入门 (Input Gate) —— “写入新知”

这里分两路并行：

- 筛选 ( $i_t$ )：Sigmoid 决定，“关于现在的 'dog'，哪些属性（如名词、单数、动物）值得记下来？”
- 创作 ( $g_t$ )：Tanh 生成全新的信息候选向量（包含了 'dog' 的所有属性）。
- 执行： $i_t \odot g_t$ 。只有被判断为“重要”的新属性被留下了。

### ★ 关键时刻：更新细胞状态 C

- 公式： $C^{(t)} = (\text{剩下的旧记忆}) \oplus (\text{筛选出的新记忆})$
- 意义：通过加法，我们将历史和现在融合。这就是为什么  $C$  能承载长期记忆的原因。现在， $C^{(t)}$  诞生了。

### 3. 输出门 (Output Gate) —— “整理发言”

现在  $C^{(t)}$  里存着极其丰富的信息（可能有这一句的，也有十句之前的）。但我们不需要把所有家底都亮出来。

- 准备 ( $o_t$ )：Sigmoid 再次观察  $x$  和  $h$ ，决定“为了预测下一个词，现在该关注记忆里的哪一部分？”
- 整理 ( $\tanh$ )：把庞大的  $C^{(t)}$  标准化到 -1 到 1 之间。
- 执行： $h^{(t)} = o_t \odot \tanh(C^{(t)})$ 。
- 结果： $h^{(t)}$  诞生了。它只包含当前任务最关心的那一点点长期记忆。

### 4. 结局：分道扬镳

- $h^{(t)}$ ：
  1. 输出给全连接层（做预测，比如预测下一个词是 "runs"）。
  2. 作为  $h^{(t-1)}$  传入下一个时间步。
- $C^{(t)}$ ：
  1. 作为  $C^{(t-1)}$  传入下一个时间步，继续它的长征。



# Today

---

- Different Ways to Model Text
- Sequence Modeling with RNNs
- Different Types of Sequence Modeling Tasks
- Backpropagation Through Time
- Long-Short Term Memory (LSTM)
- **Many-to-one Word RNNs**

# RNN Step 1: Build 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: Convert text 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: Convert 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 0 1 0 0 0 0 0 0 0 0]

[0 0 0 1 0 0 0 0 0 0 0]

...

[0 0 0 1 0 0 0 0 0 0 1]

[0 0 0 1 0 0 0 0 0 0 1]

[0 0 0 1 0 0 0 0 0 0 1]



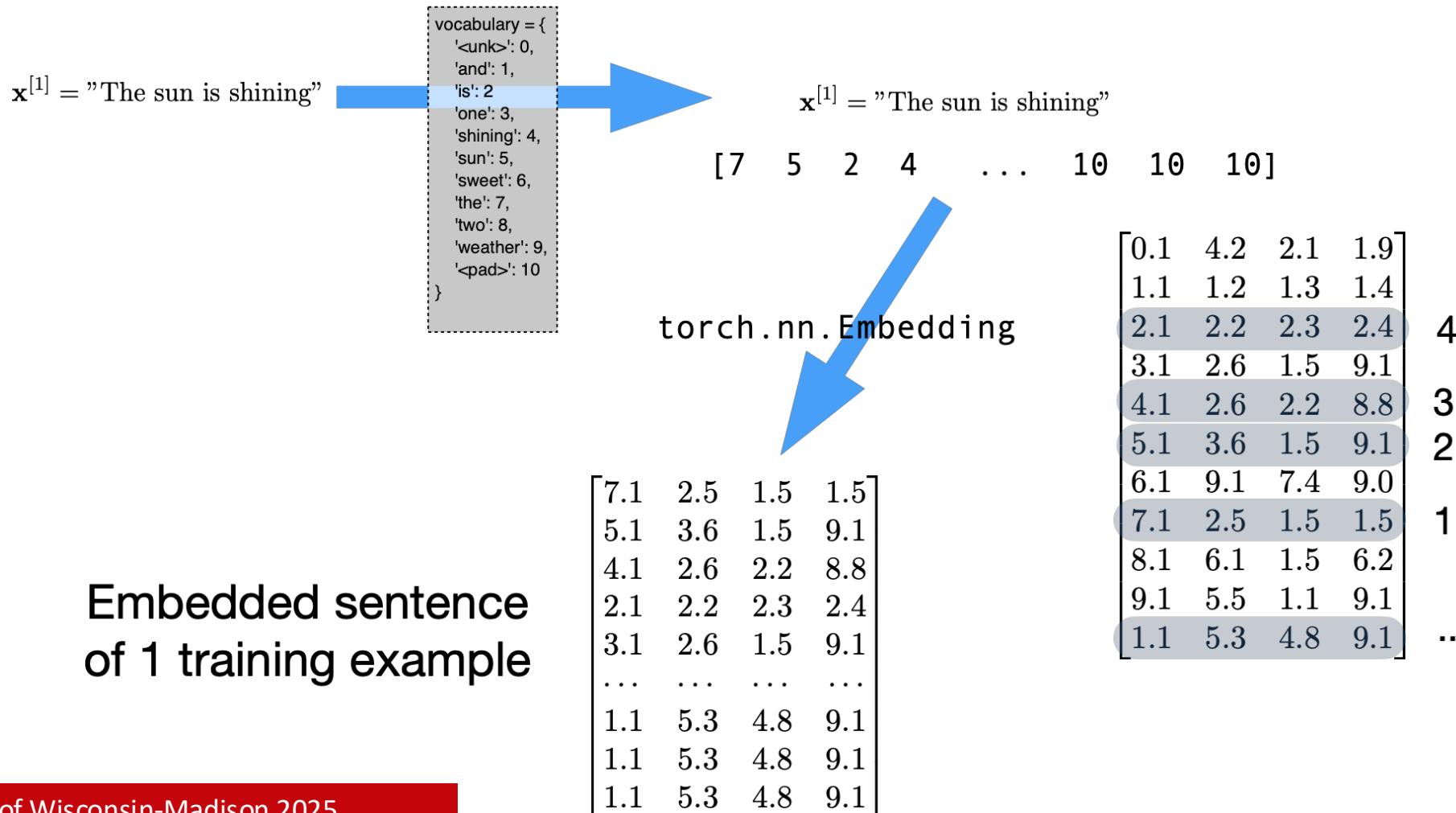
## RNN Step 4: Convert one-hot to embeddings

浪费空间 稀疏  
Embedding matrix

$$\begin{array}{l} \text{One-hot vector} \\ [0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0] \end{array} \times \begin{array}{c} \left[ \begin{matrix} 0.1 & 4.2 & 2.1 & 1.9 \\ 1.1 & 1.2 & 1.3 & 1.4 \\ 2.1 & 2.2 & 2.3 & 2.4 \\ 3.1 & 2.6 & 1.5 & 9.1 \\ 4.1 & 2.6 & 2.2 & 8.8 \\ 5.1 & 3.6 & 1.5 & 9.1 \\ 6.1 & 9.1 & 7.4 & 9.0 \\ 7.1 & 2.5 & 1.5 & 1.5 \\ 8.1 & 6.1 & 1.5 & 6.2 \\ 9.1 & 5.5 & 1.1 & 9.1 \\ 1.1 & 5.3 & 4.8 & 9.1 \end{matrix} \right] \\ = [7.1 \ 2.5 \ 1.5 \ 1.5] \end{array} = \text{Hidden layer output}$$

# PyTorch: Skip steps 3 and 4. Instead...

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





# LSTMs in PyTorch

<https://pytorch.org/docs/stable/generated/torch.nn.LSTM.html>

## Parameters

- **input\_size** – The number of expected features in the input  $x$
- **hidden\_size** – The number of features in the hidden state  $h$
- **num\_layers** – Number of recurrent layers. E.g., setting `num_layers=2` would mean stacking two LSTMs together to form a *stacked LSTM*, with the second LSTM taking in outputs of the first LSTM and computing the final results. Default: 1
- **bias** – If `False`, then the layer does not use bias weights  $b_{ih}$  and  $b_{hh}$ . Default: `True`
- **batch\_first** – If `True`, then the input and output tensors are provided as (batch, seq, feature). Default: `False`
- **dropout** – If non-zero, introduces a *Dropout* layer on the outputs of each LSTM layer except the last layer, with dropout probability equal to `dropout`. Default: 0
- **bidirectional** – If `True`, becomes a bidirectional LSTM. Default: `False`
- **proj\_size** – If  $> 0$ , will use LSTM with projections of corresponding size. Default: 0

## Examples:

```
>>> rnn = nn.LSTM(10, 20, 2)
>>> input = torch.randn(5, 3, 10)
>>> h0 = torch.randn(2, 3, 20)
>>> c0 = torch.randn(2, 3, 20)
>>> output, (hn, cn) = rnn(input, (h0, c0))
```



# Good reading

---

- [The Unreasonable Effectiveness of Recurrent Neural Networks](#) by Andrej Karpathy
- [On the difficulty of training recurrent neural networks](#) by Razvan Pascanu, Tomas Mikolov, Yoshua Bengio

Questions?

