
Neural Networks for sequential data

29 January, 2020

Alexey Zaytsev, head of Laboratory LARSS, PhD

ISP course on sequential data processing

-
- **Why bother about Neural Networks for sequential data?**
 - How to do Neural Networks for sequential data?

Complex sequence data examples

Speech recognition



"The quick brown fox jumped over the lazy dog."

Sentiment classification

"There is nothing to like in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter
avec moi?



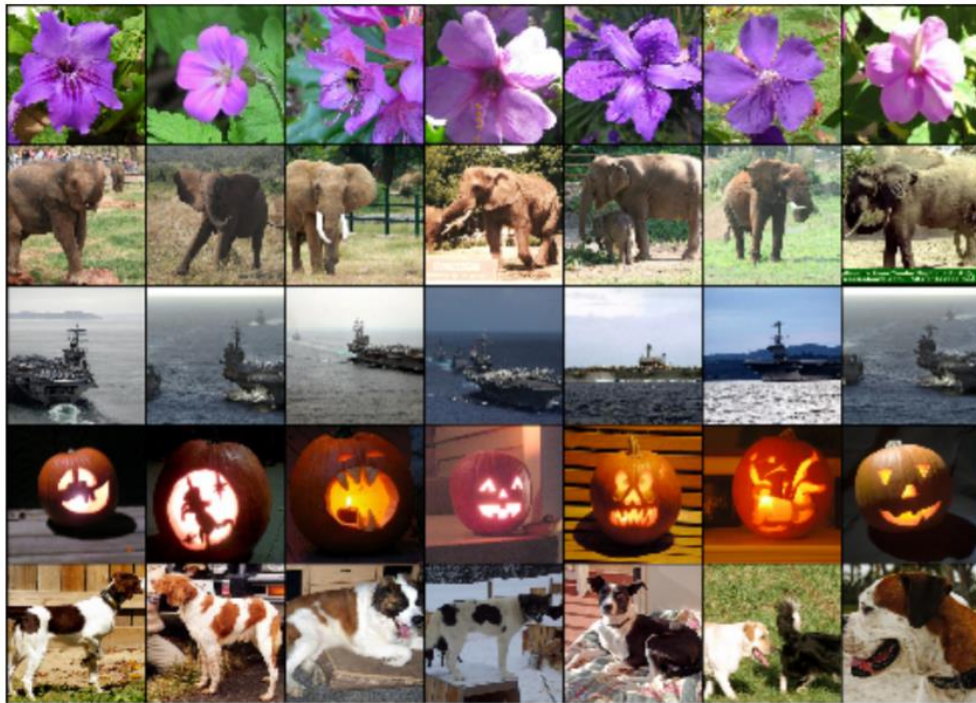
Do you want to sing with me?

Video activity
recognition



Running

Why bother with deep learning: Deep learning is better, then a human



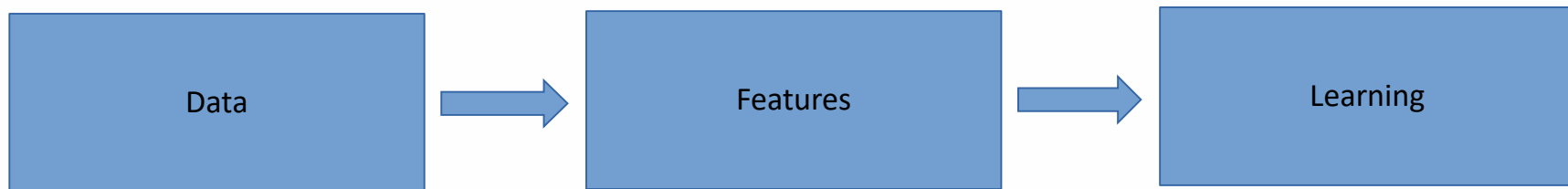
Human error percentage is

5.1%

DL error percentage is

3.57%

Classic approach to Machine Learning



Image



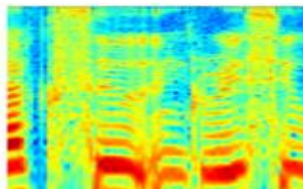
Features



Recognition



Speech

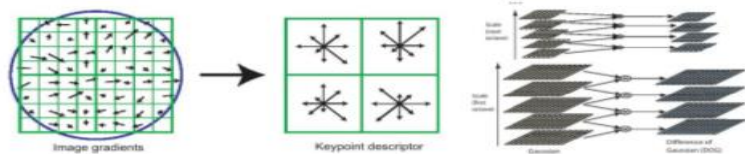


Features

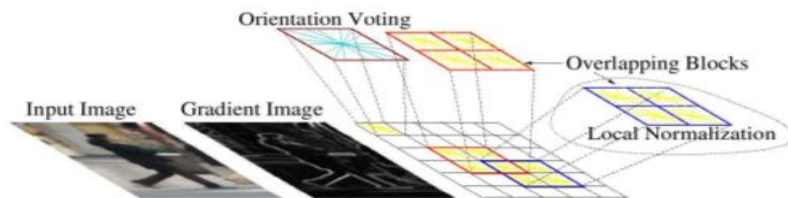


Recognition

An art of feature construction

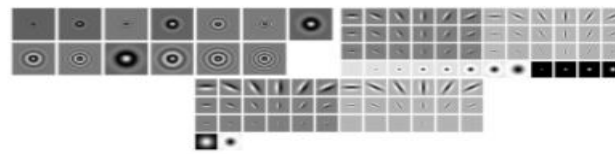
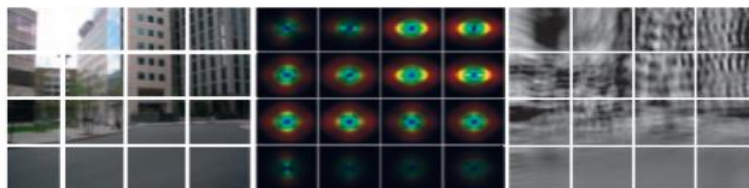


SIFT

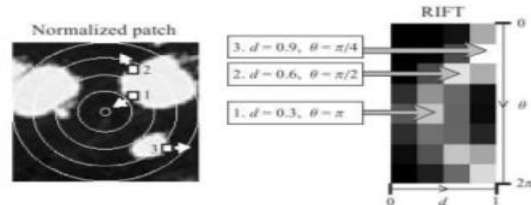


HoG

GIST

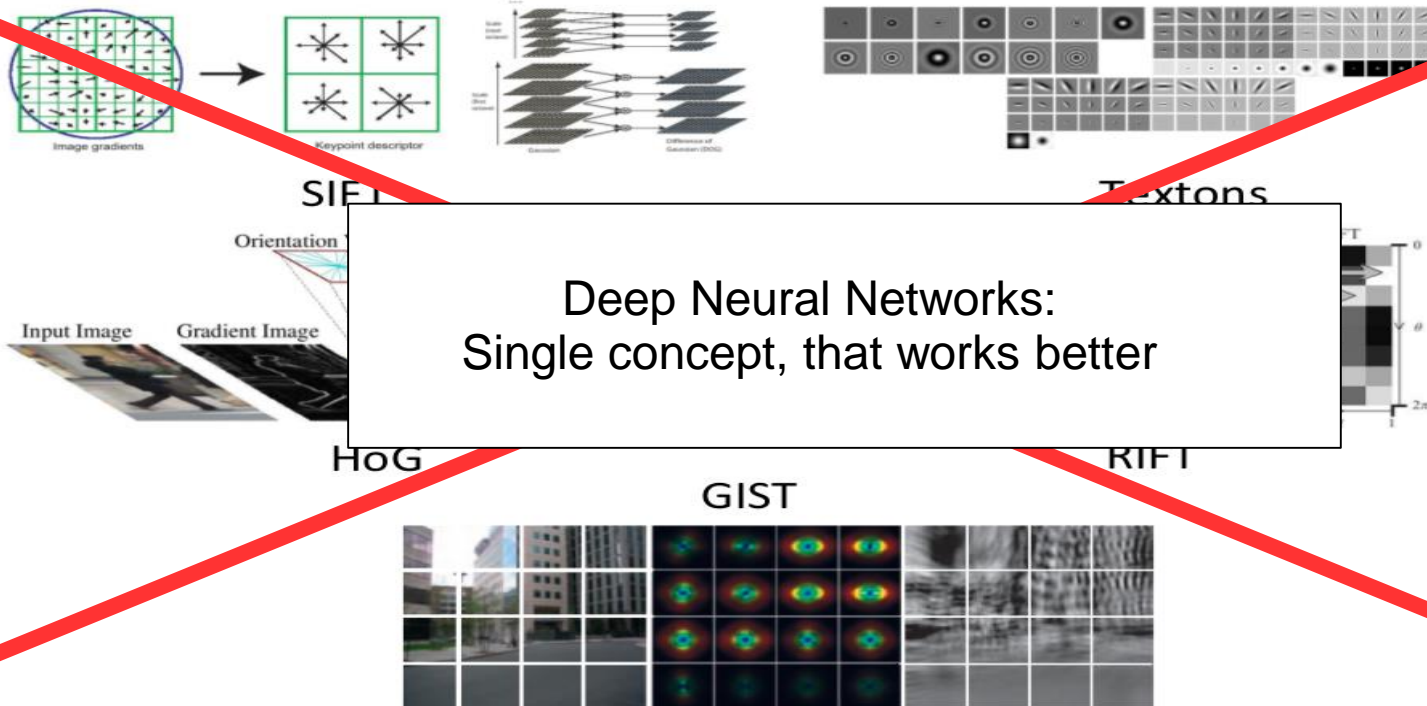


Textons



RIFT

No art now, just engineering



Why and how Deep learning works

- Availability of Graphical processing units (GPU)
- Large scale “big” data
- Open-source libraries
- Complex structured data



PYTORCH



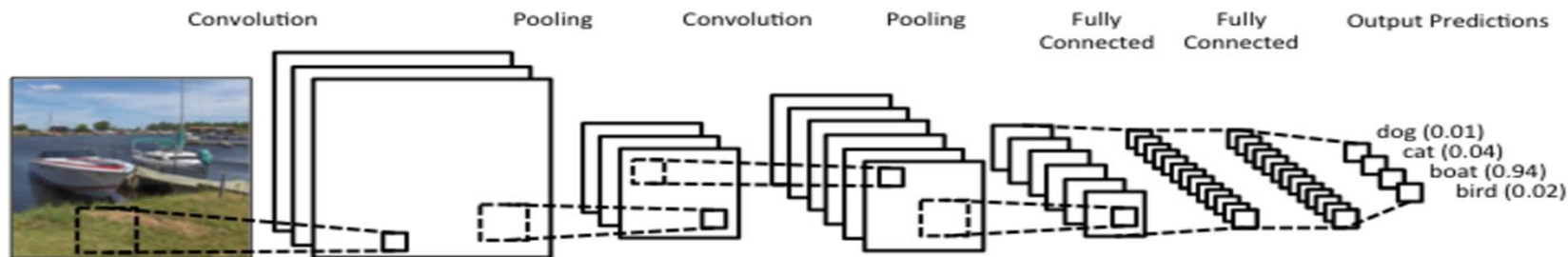
Complex sequence data examples

		Large data set	Structured data
Speech recognition		+	+
Sentiment classification	"There is nothing to like in this movie."	+	+
DNA sequence analysis	AGCCCCTGTGAGGAACTAG	+	+
Machine translation	Voulez-vous chanter avec moi?	+	+
Video activity recognition		+	+

- Why bother about Neural Networks for sequential data?
- **How to do Neural Networks for sequential data?**

Convolutional Neural Networks

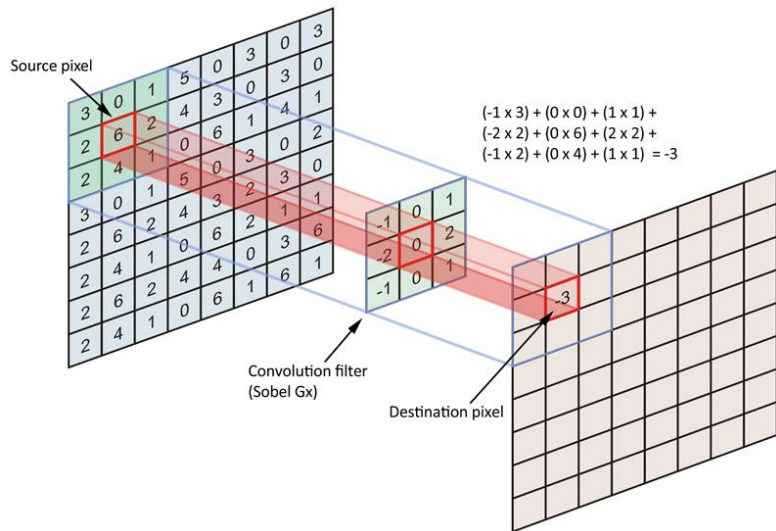
- Mathematical definition: combination of simple transformations
- There can be a lot of layers
- Convolutional layers help



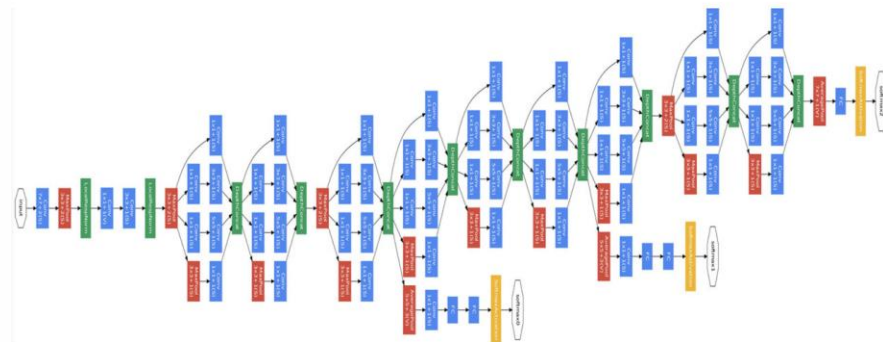
<https://playground.tensorflow.org/>

Convolutional Neural Networks

- Mathematical definition: combination of simple transformations
- There can be a lot of layers
- Convolutional layers help

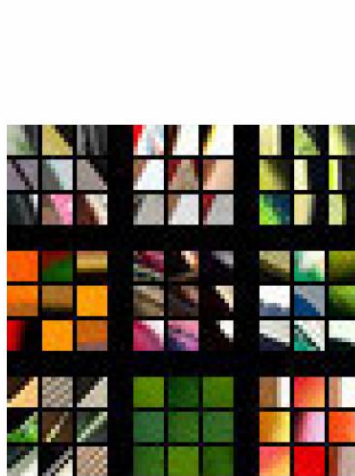
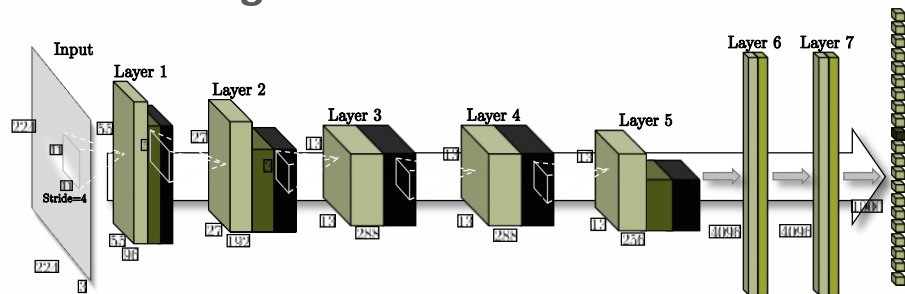


Single convolutional layer



GoogLeNet architecture

Flow of data through Convolutional Neural Network



Layer 1

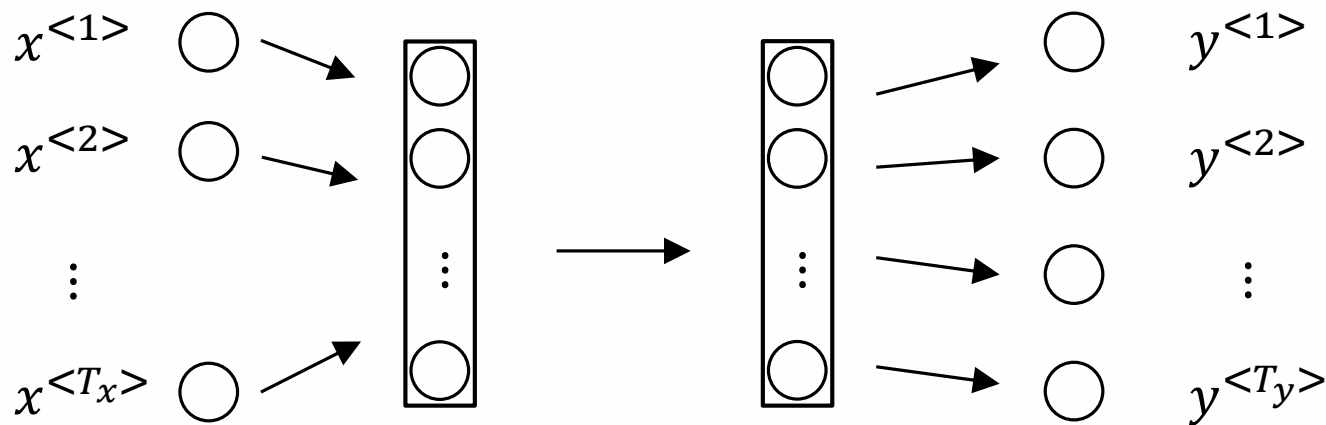


Layer 2



Layer 5

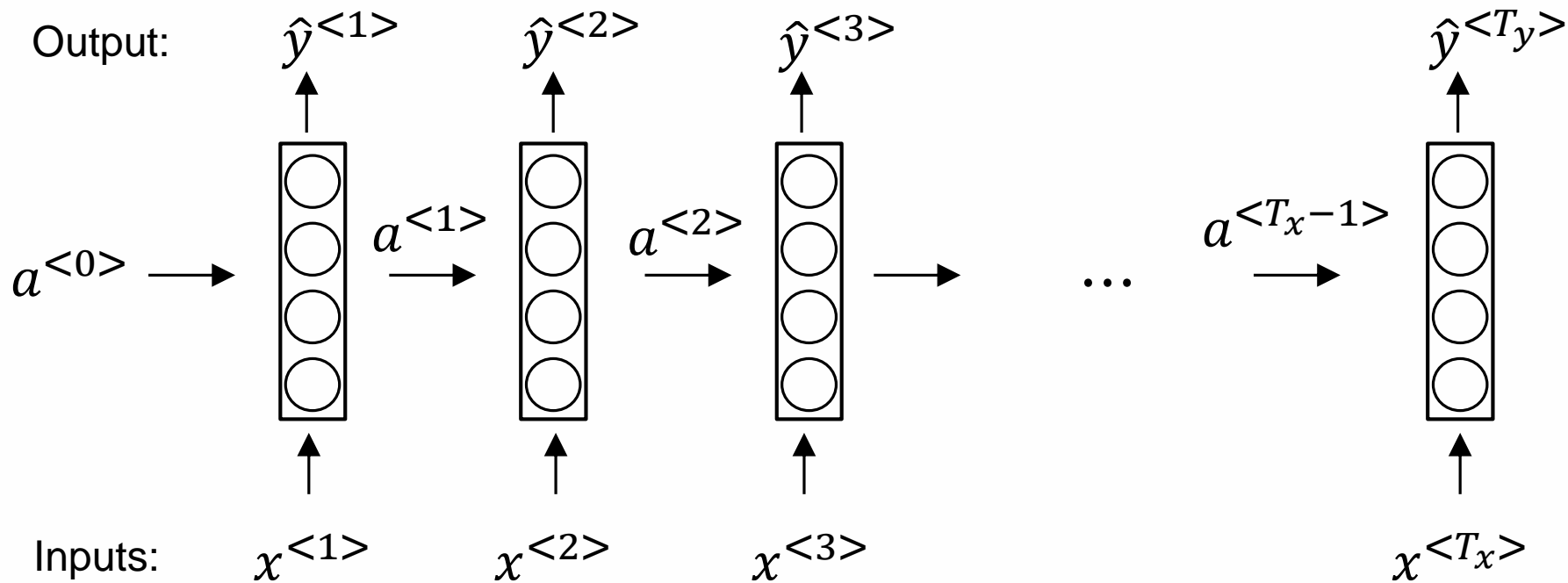
Why not a standard fully-connected or convolutional network?



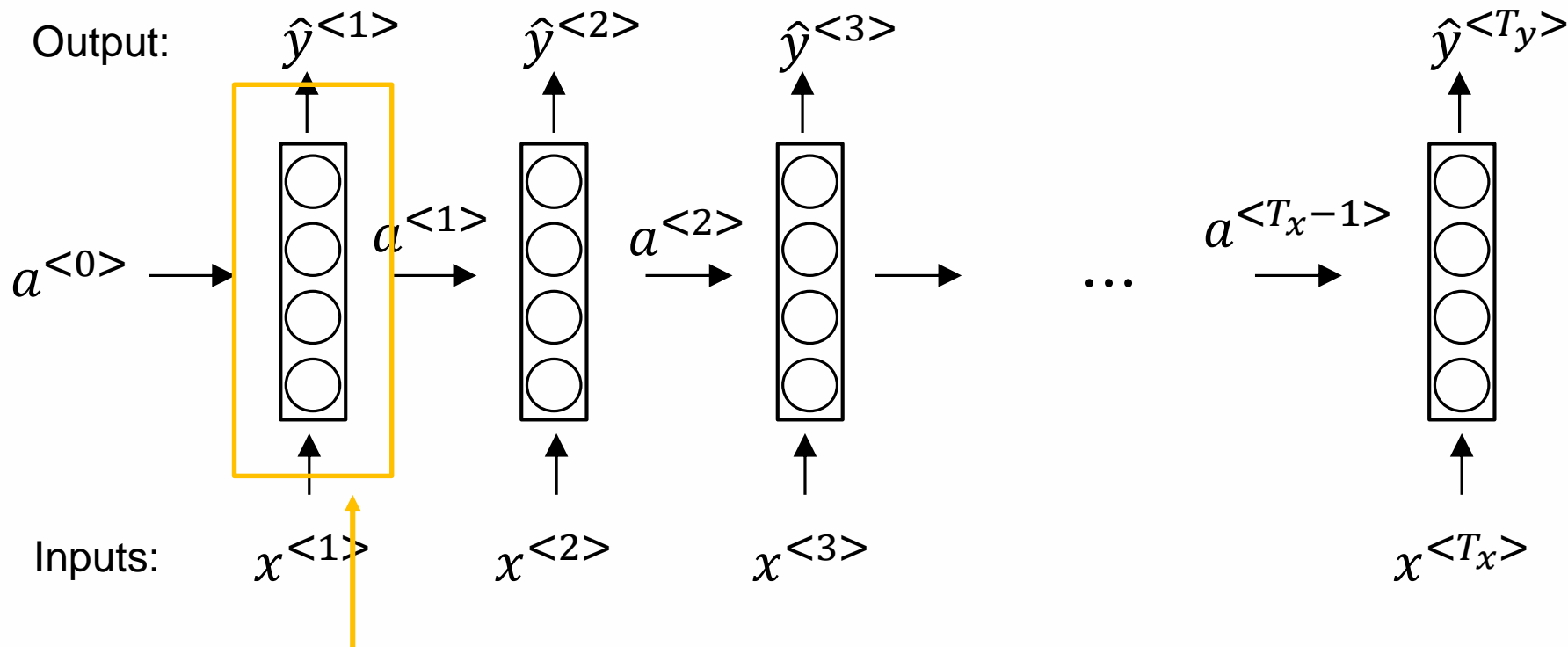
Problems:

- Inputs, outputs can be different lengths in different examples
- Doesn't share features learned across different positions of text

Forward propagation through Recurrent Neural Network

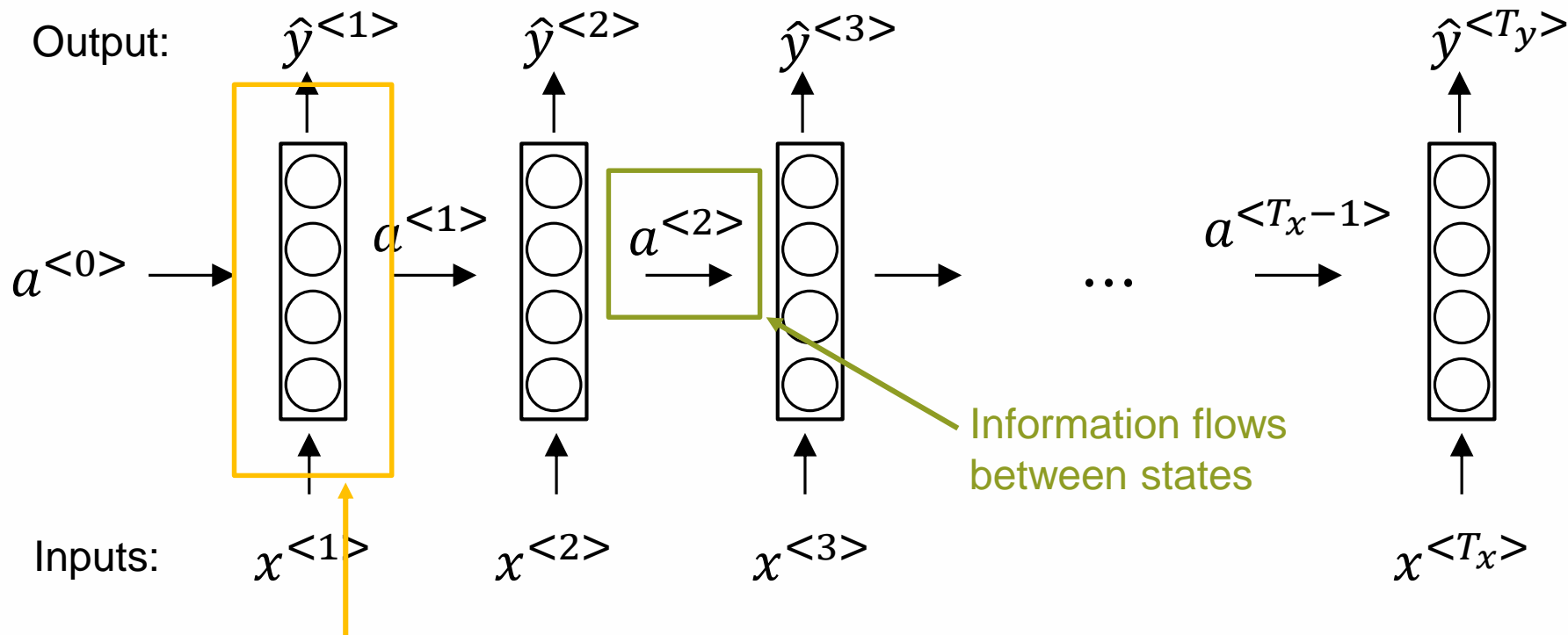


Forward propagation through Recurrent Neural Network



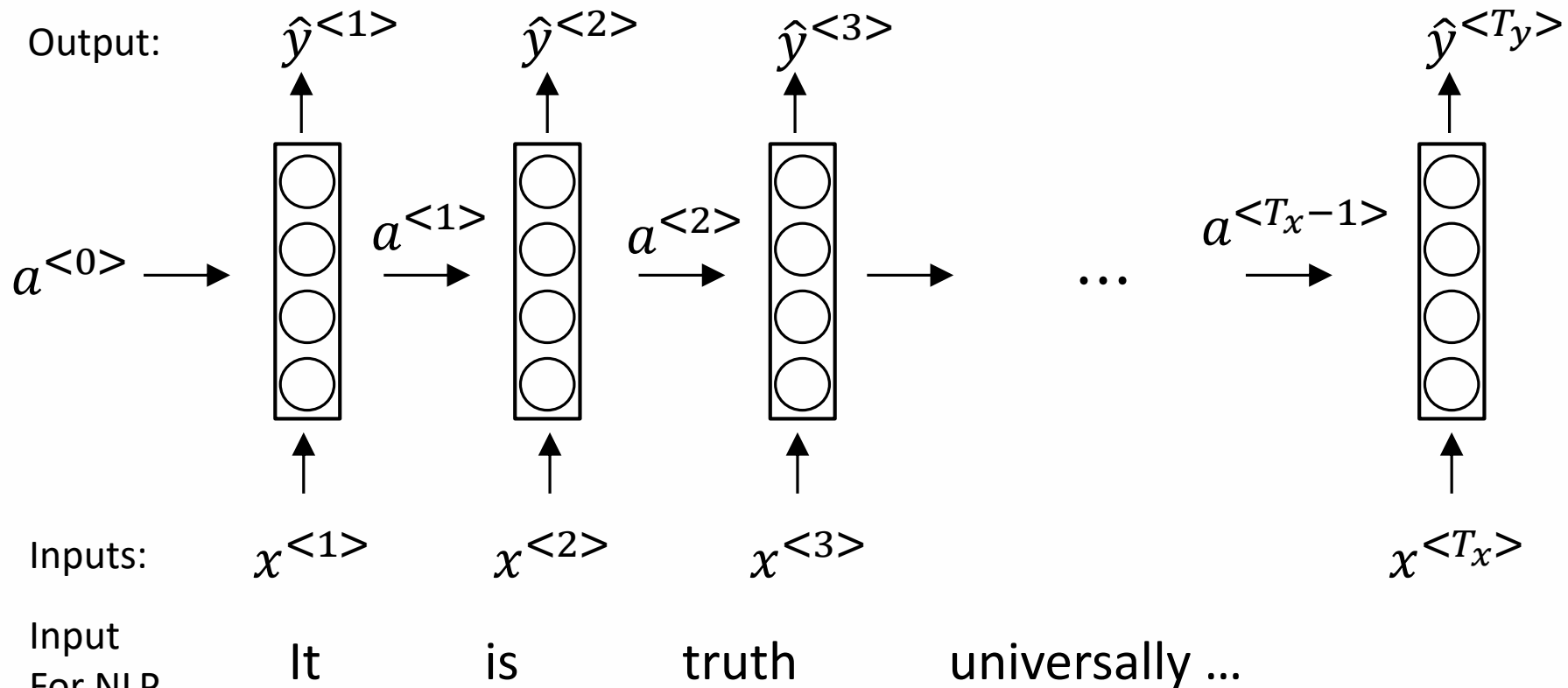
Processing unit is the same for all time moments.
Units has parameters we want to learn!

Forward propagation through Recurrent Neural Network



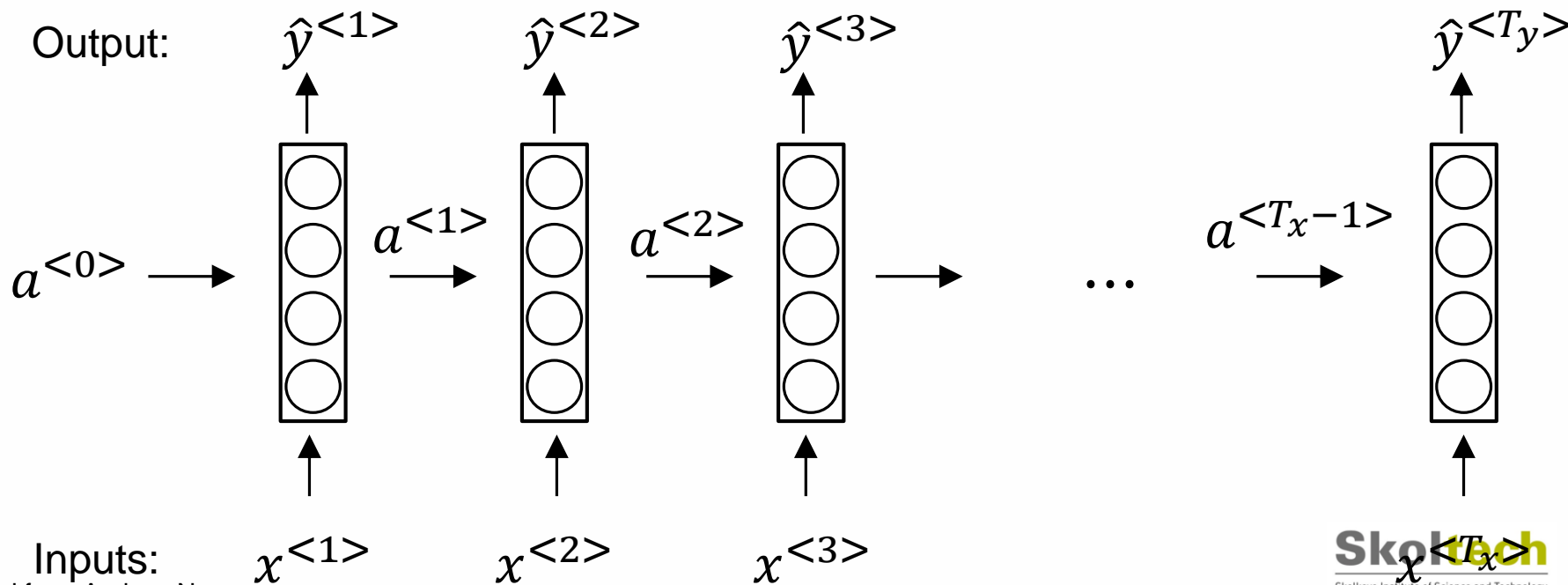
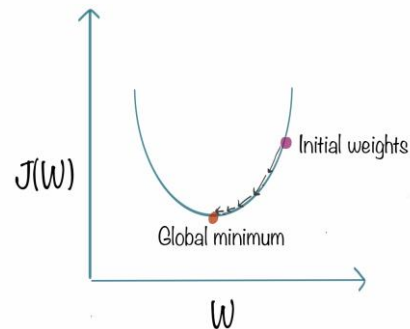
Processing unit is the same for all time moments.
Units has parameters we want to learn!

Forward propagation through Recurrent Neural Network



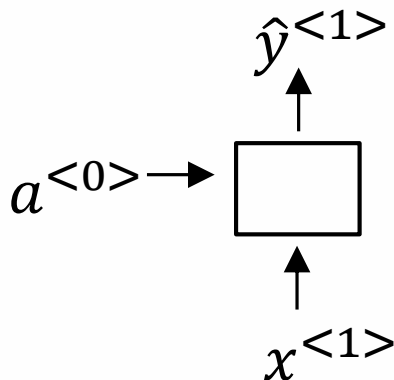
Backward propagation through Recurrent Neural Network

- Our loss function is $\mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$
- We want to estimate parameters, so we need derivatives!



Zoo of RNN (Recurrent Neural Network) models

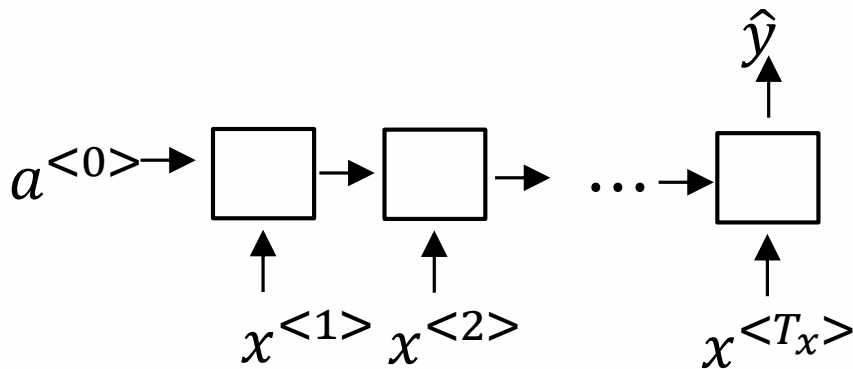
One to one



We have a separate input & output each time

- Image classification from cameras

Many to one

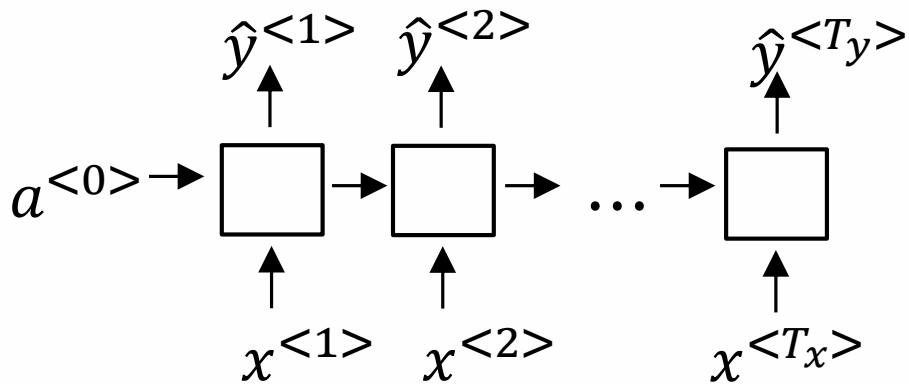


We have a single output for a sequence

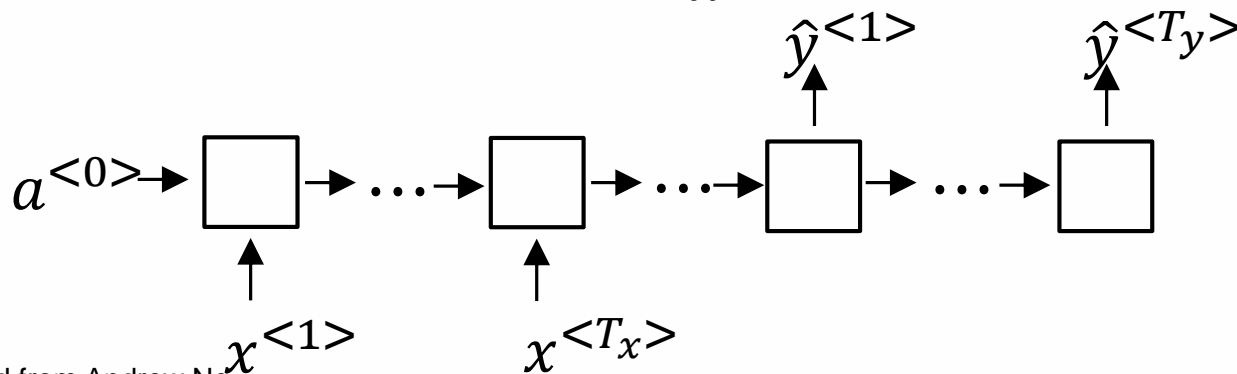
- Sentiment of a sentence: good or bad review?

Zoo of RNN (Recurrent Neural Network) models

Many to many



- Time series prediction
- Anomaly detection



- Translation

Forward propagation through Recurrent Neural Network

Initial sequence:

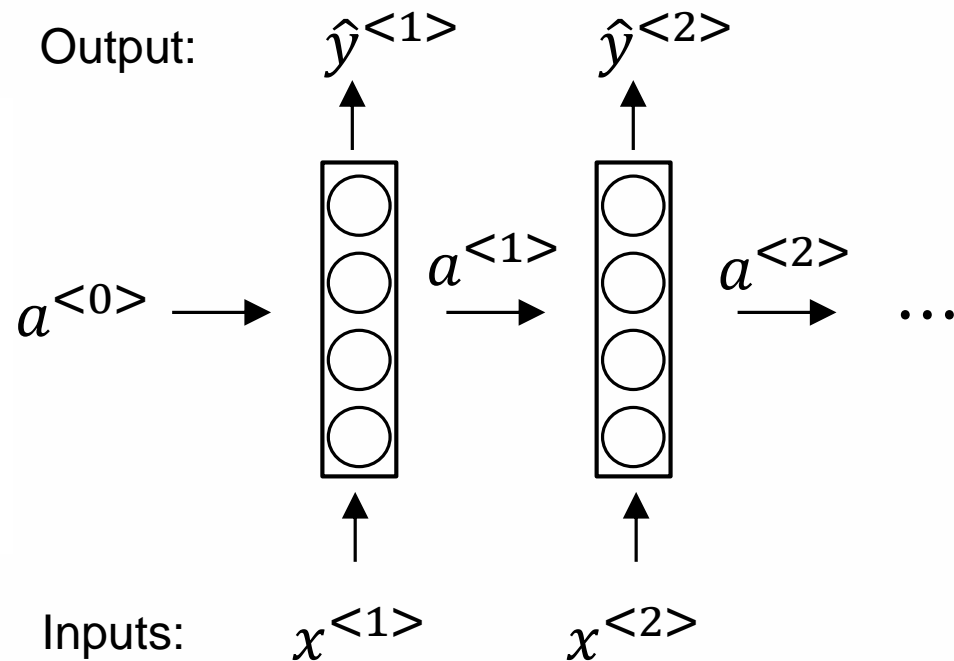
$$x^{<1:n>} = x^{<1>}, x^{<2>}, \dots, x^{<T>}, x_i \in \mathbb{R}^{d_{in}}$$

For each input $x_{1:i}$ we get an output y_i :

$$y^{<i>} = RNN(x^{<1:i>}), y^{<i>} \in \mathbb{R}^{d_{out}}$$

For the whole sequence $x^{<1:n>}$:

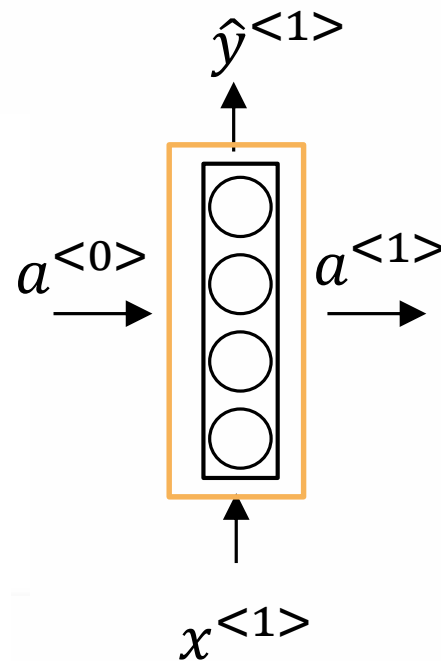
$$y^{<1:n>} = RNN^*(x^{<1:n>}), y^{<i>} \in \mathbb{R}^{d_{out}}$$



Selection of RNN architecture

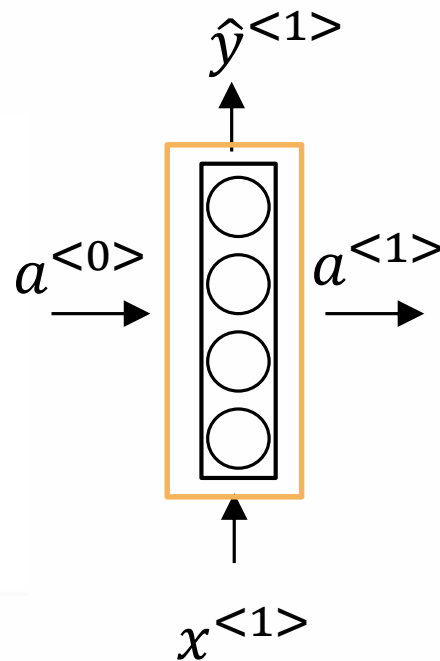
Simplest RNN unit: what is going on inside **yellow** rectangle?

- $RNN^*(x^{<1:n>}, a^{<0>}) = y^{<1:n>}$
- $y^{<i>} = g(W^{out}[a^{<i>}, x^{<i>}] + b)$



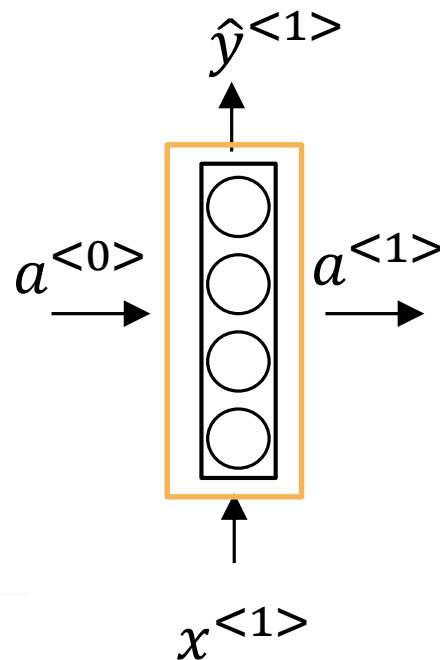
Simplest RNN unit: what is going on inside yellow rectangle?

- $RNN^*(x^{<1:n>}, a^{<0>}) = y^{<1:n>}$
- $y^{<i>} = g(W^{out}[a^{<i>}, x^{<i>}] + b)$
- R is a recursive activation function. It depends on inputs $x^{<t>}$ and output of the previous state $a^{<t-1>}$ (vector of the previous state)
- $a^{<i>} = R(a^{<i-1>}, x^{<i>})$
- $R(a^{<i-1>}, x^{<i>}) = g(W^{hid}[a^{<i-1>}, x^{<i>}] + b)$, $[a^{<i>}, x^{<i>}]$ is the concatenation of $a^{<i>}$ and $x^{<i>}$
- $x^{<i>} \in \mathbb{R}^{d_{in}}, y^{<i>} \in \mathbb{R}^{d_{out}}, a^{<i>} \in \mathbb{R}^{d_{hid}}$

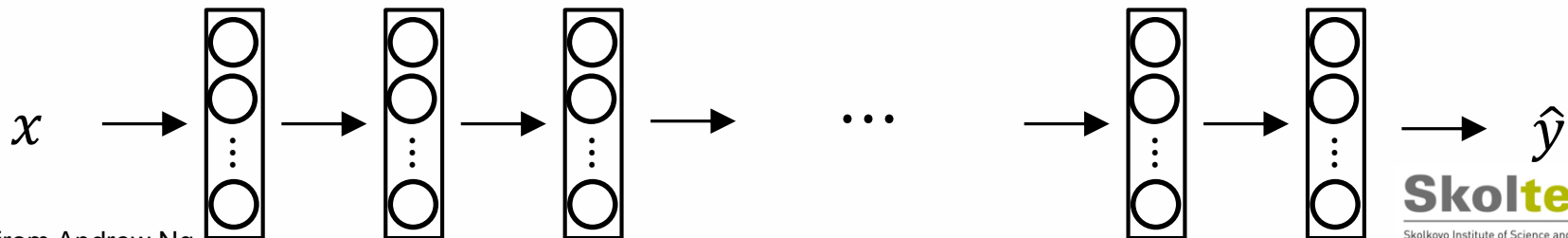
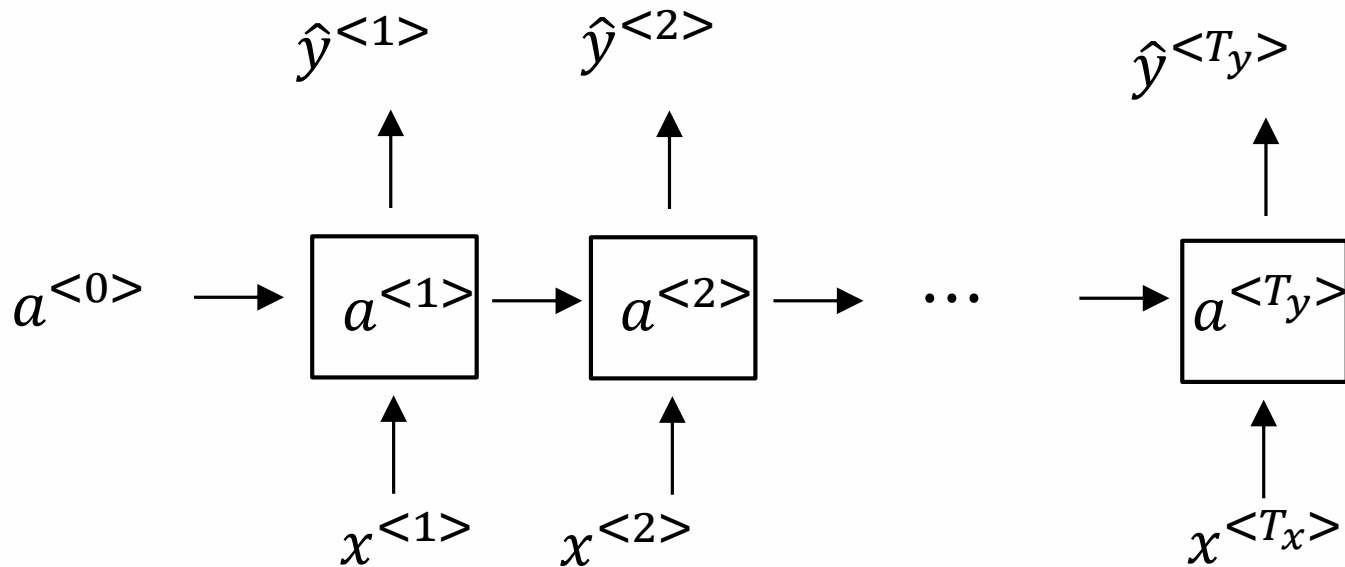


Simplest RNN unit: what is going on inside yellow rectangle?

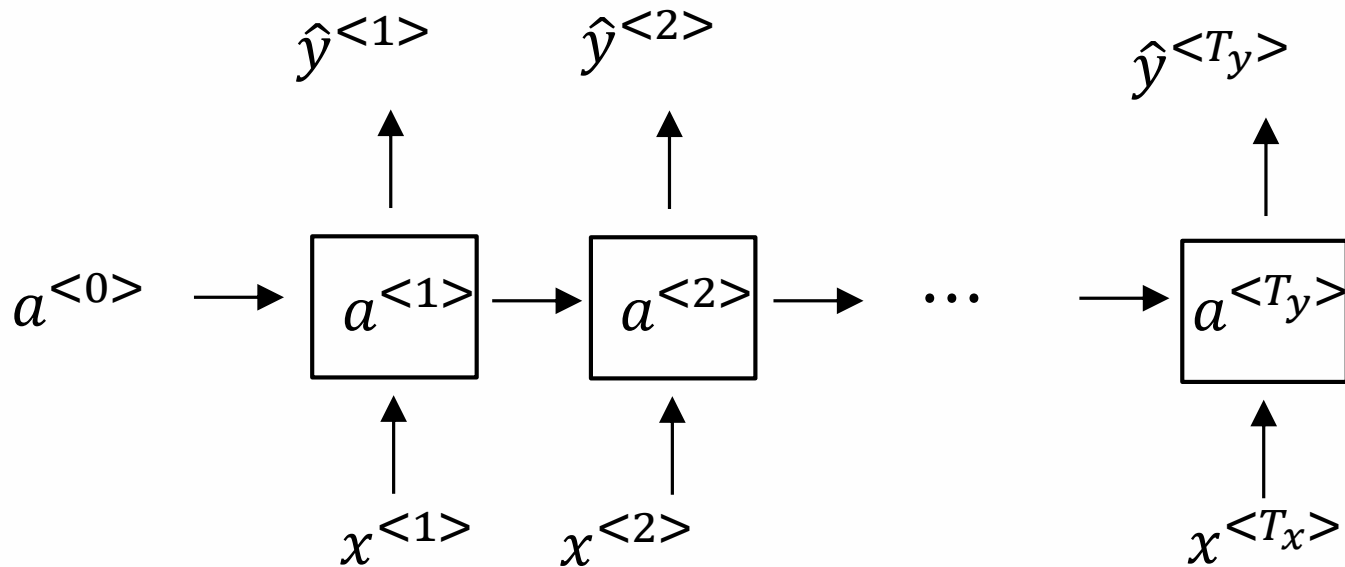
- $RNN^*(x^{<1:n>}, a^{<0>}) = y^{<1:n>}$
- $y^{<i>} = g(W^{out}[a^{<i>}, x^{<i>}] + b)$
- R is a recursive activation function. It depends on inputs $x^{<t>}$ and output of the previous state $a^{<t-1>}$ (vector of the previous state)
- $a^{<i>} = R(a^{<i-1>}, x^{<i>})$
- $R(a^{<i-1>}, x^{<i>}) = g(W^{hid}[a^{<i-1>}, x^{<i>}] + b)$, $[a^{<i>}, x^{<i>}]$ is the concatenation of $a^{<i>}$ and $x^{<i>}$
- $x^{<i>} \in \mathbb{R}^{d_{in}}, y^{<i>} \in \mathbb{R}^{d_{out}}, a^{<i>} \in \mathbb{R}^{d_{hid}}$
- Parameters of Neural Network are $W^{hid} \in \mathbb{R}^{(d_{in}+d_{out}) \times d_{hid}}, W^{out} \in \mathbb{R}^{d_{hid} \times d_{out}}$



Gradient vanish and/or explode



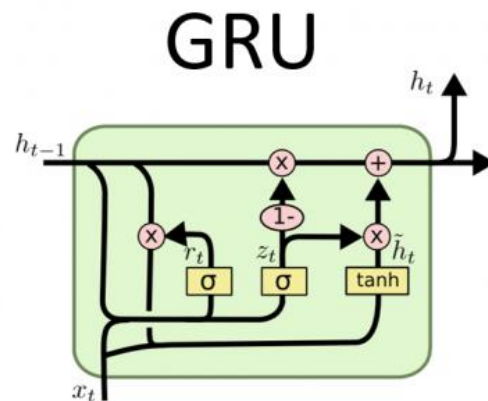
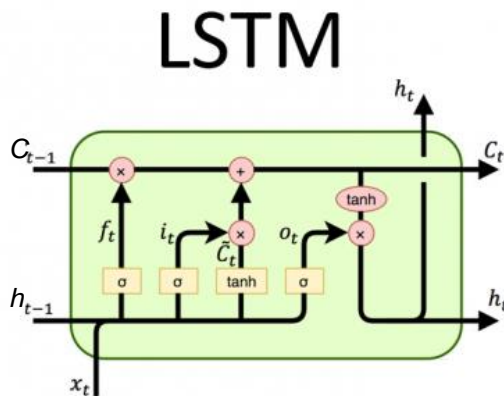
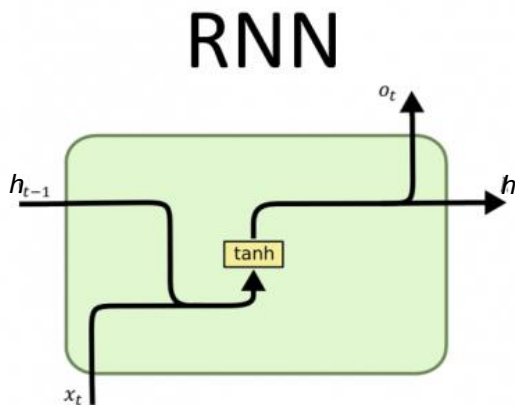
Gradient vanish and/or explode



Neural networks forget fast, and it is hard to learn long-term dependencies

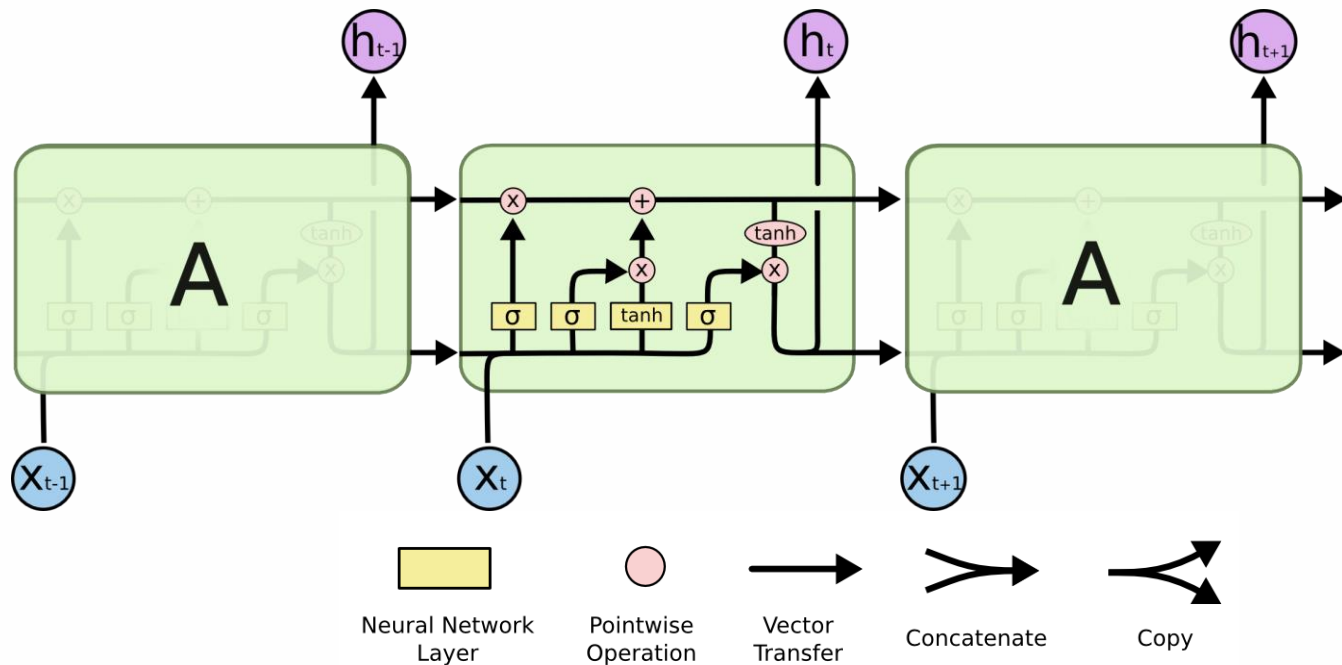
Better RNN units: LSTM and GRU

- LSTM: long short term memory
- GRU: Gated recurrent unit

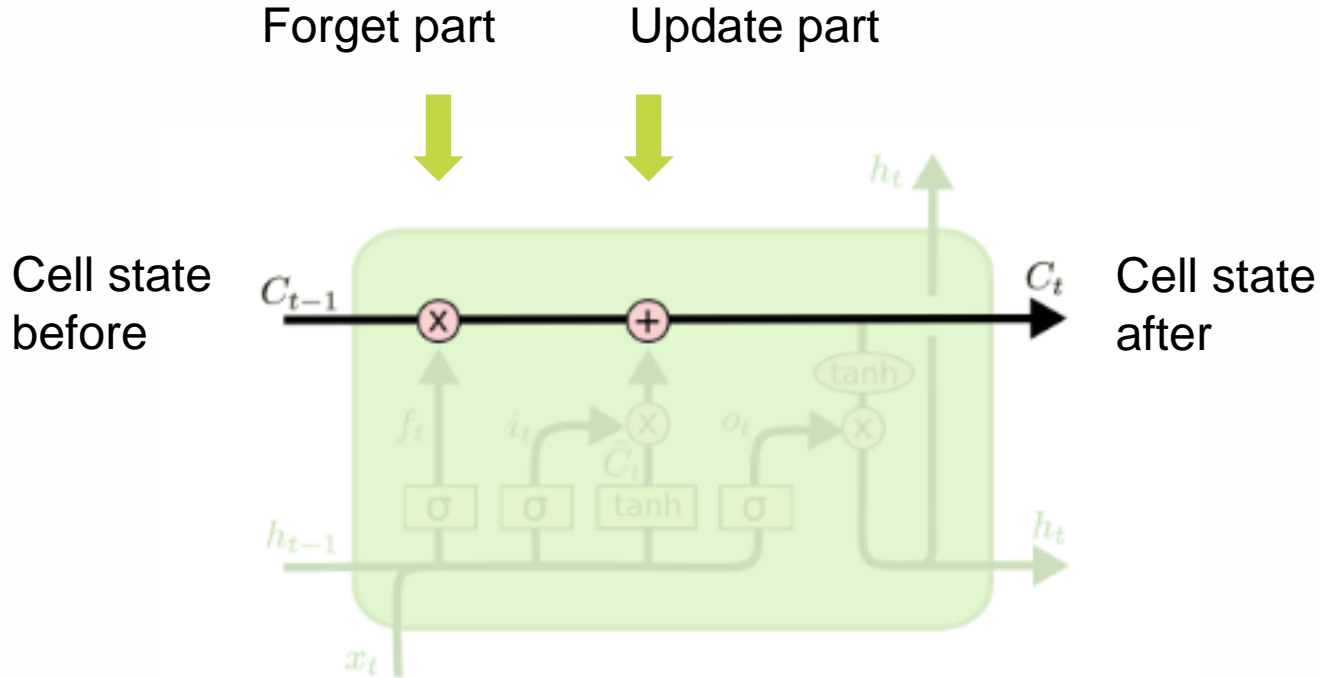


Details on how LSTM works

Remembering information for long periods of time is practically the default behavior of LSTM

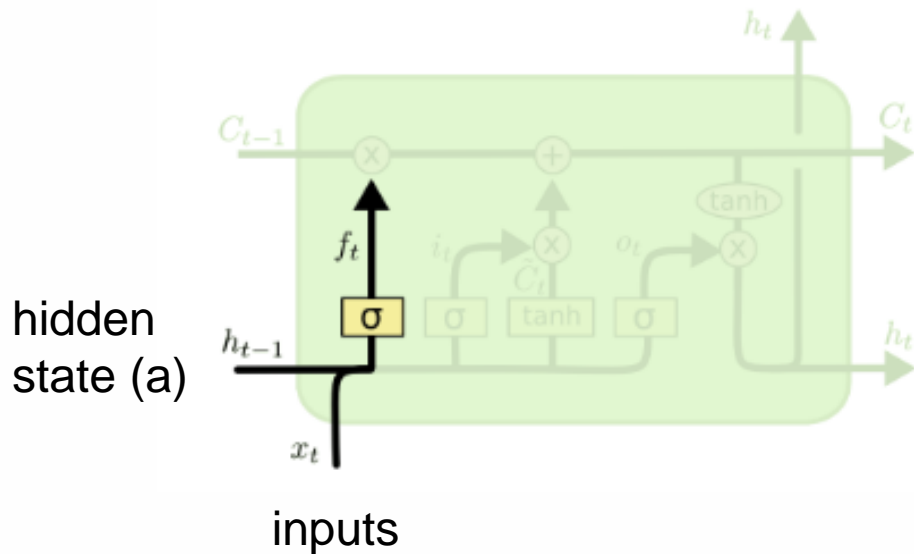


Long term memory part – Cell state



Forget part

Identify how much should we forget:
sigmoid returns value between 0 and 1

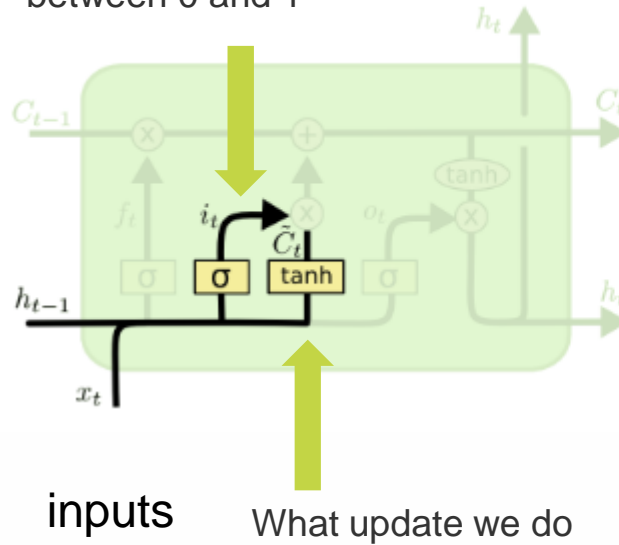


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Update part

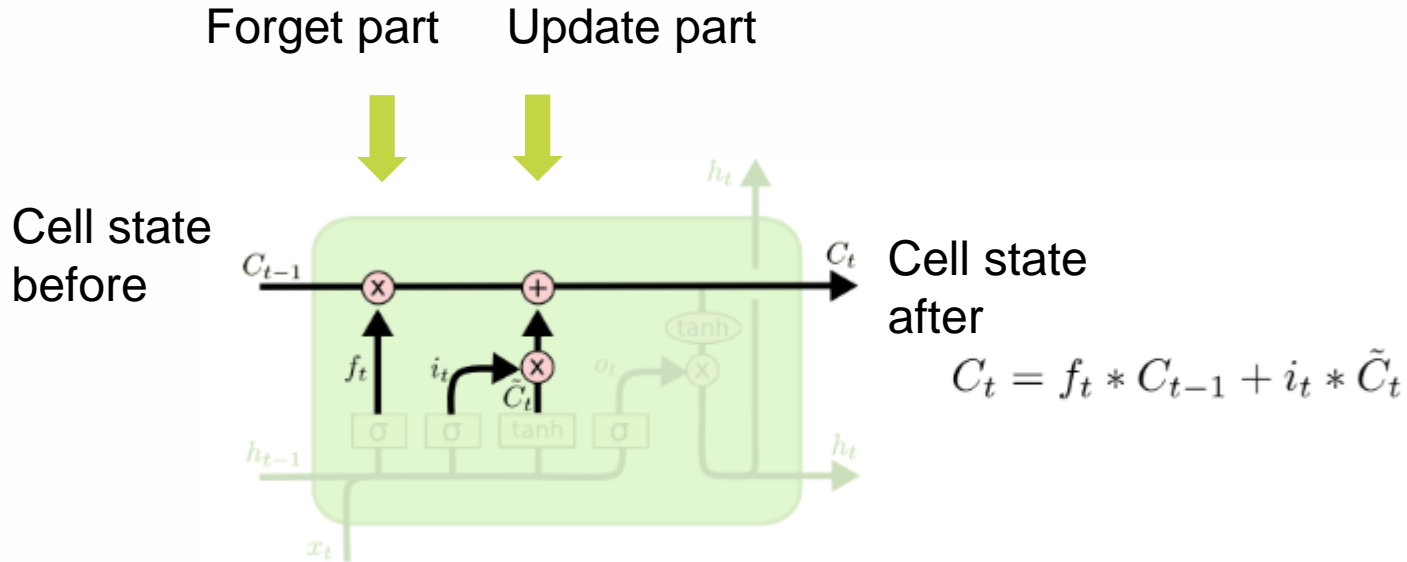
Identify how much should we update: sigmoid returns value between 0 and 1

hidden state (a)

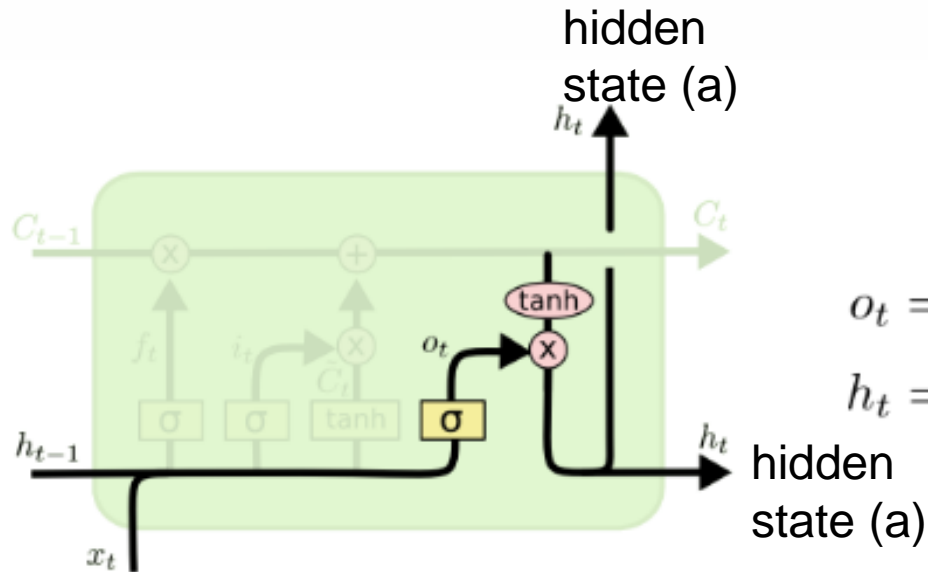


$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Long term memory part – Cell state



Update everything else

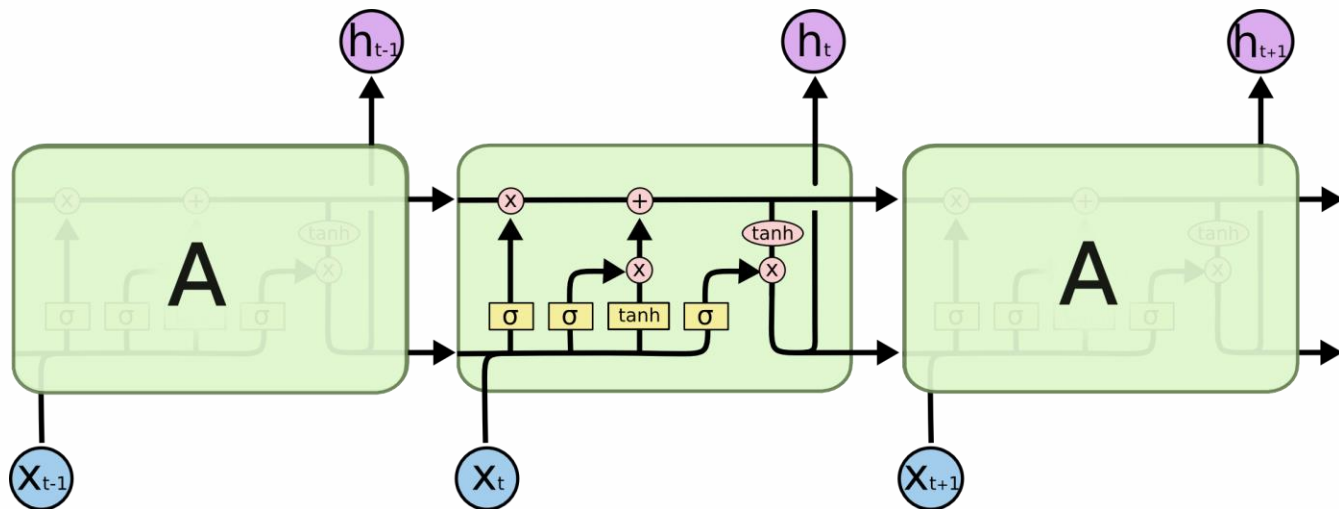


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

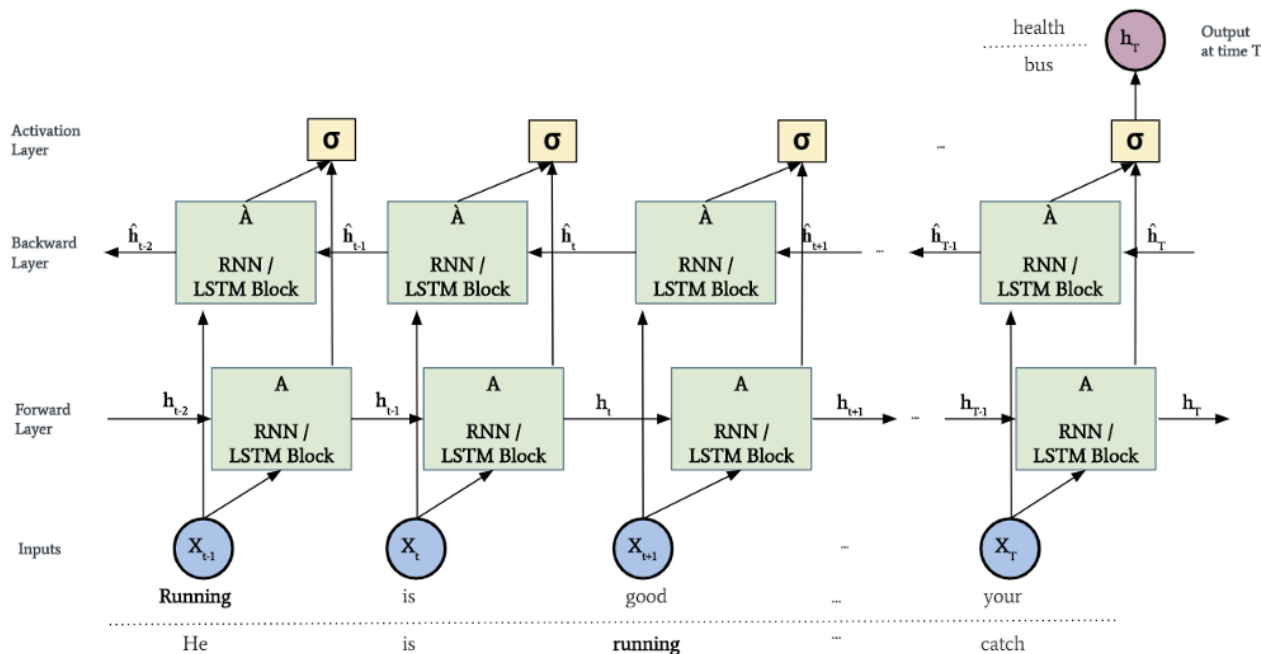
$$h_t = o_t * \tanh (C_t)$$

Details on how LSTM works

- There are cell and hidden (activation) states
- LSTM block forgets and updates cell state during processing at one block



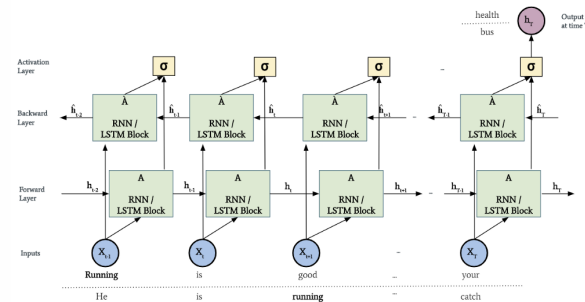
Other architectures: bidirectional LSTM



Other architectures: bidirectional LSTM

Bidirectional LSTM are useful when we benefit from the future data:

- handwriting recognition
- speech recognition
- protein Structure Prediction (bioinformatics)



Take-home messages

- For some types of data classic methods fail:
we need to learn a representation i.e. extract features automatically
- Neural Networks provide enough flexibility for this problem for various data types
- The basic architecture is Recurrent Neural Network RNN
- But we can do better in terms of keeping the necessary information with LSTM and GRU blocks/architectures