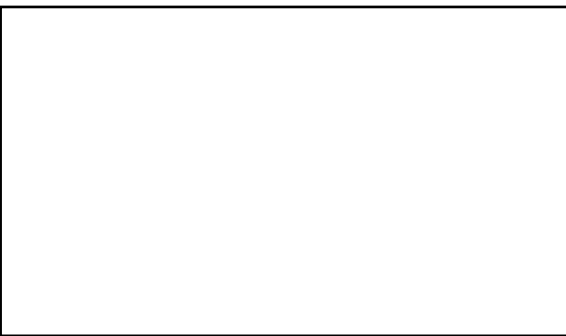




Recurrent Neural Networks

Marcin Walkowski 2022



Plan for today

- Sequence models motivation
- RNNs
- LSTMs
- **Kahoot time!**



Resources

- [Andrew Ng - Sequence Models](#)
- [Leo Dirac - LSTM is dead. Long Live Transformers!](#)
- [Christopher Olah - Understanding LSTM Networks](#)

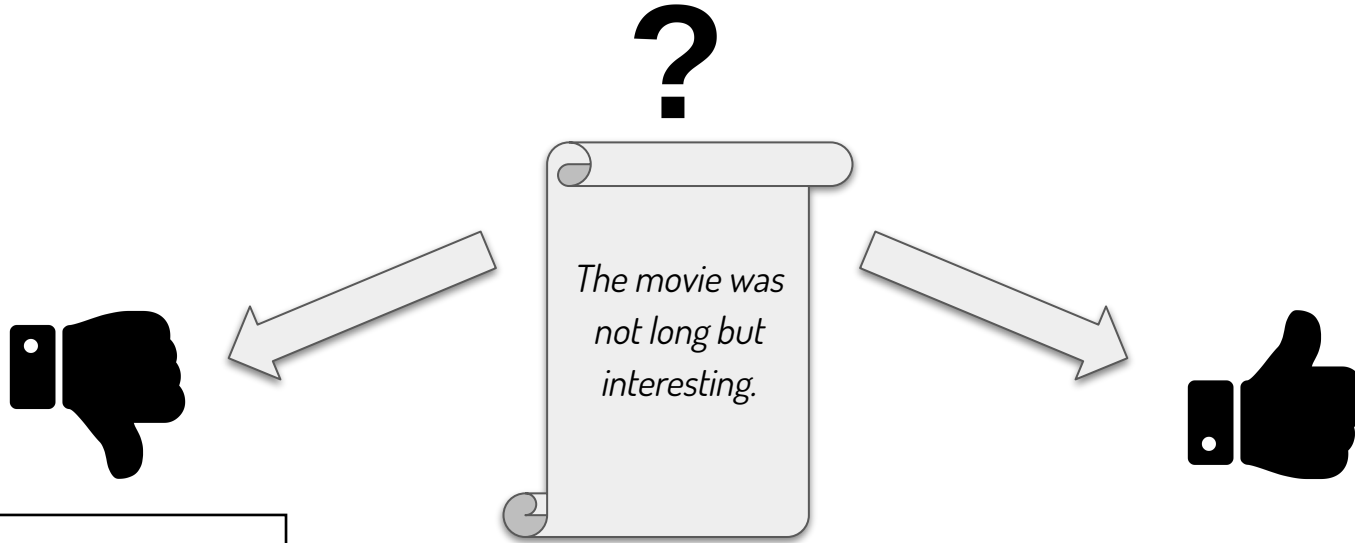


Motivation - Sequence data tasks

- Sentiment Classification variable \mapsto fixed size
- ECG Classification variable \mapsto fixed size
- Machine Translation variable \mapsto variable size
- ASR and TTS variable \mapsto variable size
- Music Generation fixed \mapsto variable size



Motivation - Sentiment Classification



Motivation – Bag-of-words

“The movie was not long but interesting.”

{'the': 1, 'movie': 1, 'song': 0, 'was': 1, 'is': 0, 'not': 1, 'long': 1, 'but': 1, 'interesting': 1}

“I like to sleep and I like to eat”

{'i': 2, 'like': 2, 'to': 2, 'sleep': 1, 'and': 1, 'eat': 1}



Motivation - Bag-of-words

- Fixed-length vector equal to dictionary size
- Information about the frequency
- No information about the order



Motivation - Bag-of-words

- Fixed-length vector equal to dictionary size
- Information about the frequency
- No information about the order

“The movie was not long but interesting.”

“The movie was long but not interesting.”

{'the': 1, 'movie': 1, 'song': 0, 'was': 1, 'is': 0, 'not': 1, 'long': 1, 'but': 1, 'interesting': 1}

Order matters!



Motivation - Bag-of-words

- n-grams are potential solution
- Vector size grows exponentially

The PWN spelling dictionary has approx. 140 thousands words.
By using bigrams, we get **19.6 billion** values.



RNN



Recurrent **N**eural **N**etwork



RNN

$$f(x^{<1>}, x^{<2>}, \dots, x^{<T_x>})?$$

- How to compute a function for variable-length data?



RNN

$$f(x^{<1>}, x^{<2>}, \dots, x^{<T_x>})?$$

- How to compute a function for variable-length data?
- Recursively. Final activation is the final output



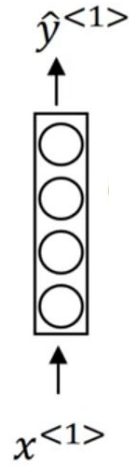
RNN - Word Representation

the
movie
was
not
long
but
interesting

was \mapsto $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
 $x^{<t>}$



RNN - Diagram

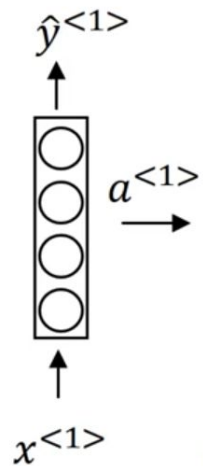


Source of diagrams in the next sections:

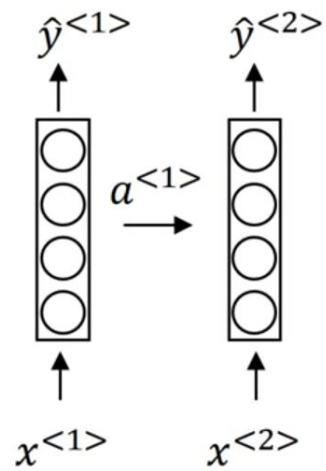
<https://www.coursera.org/learn/nlp-sequence-models>



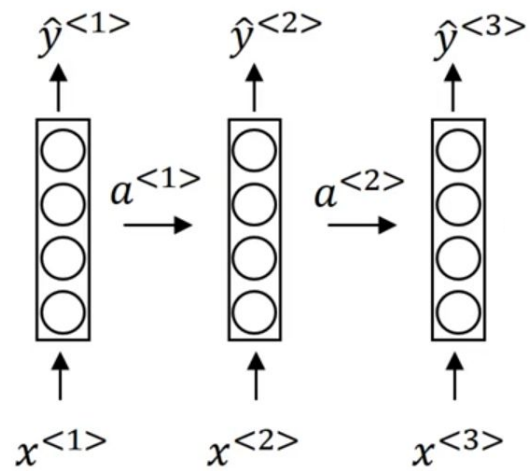
RNN - Diagram



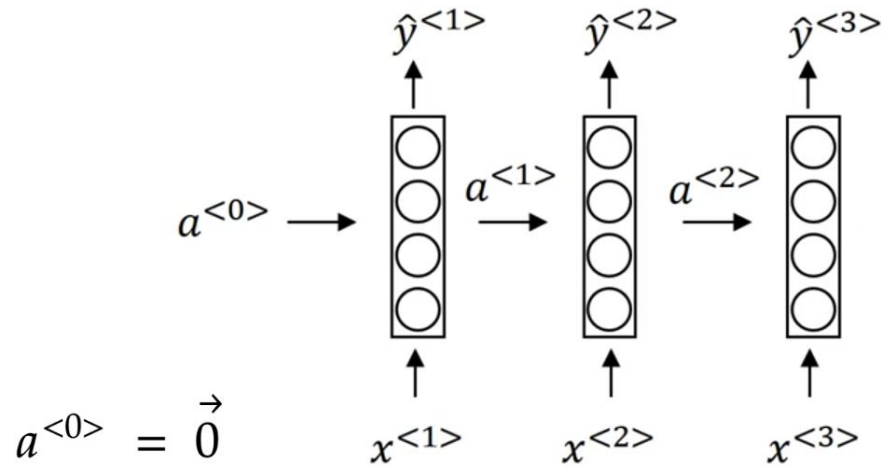
RNN - Diagram



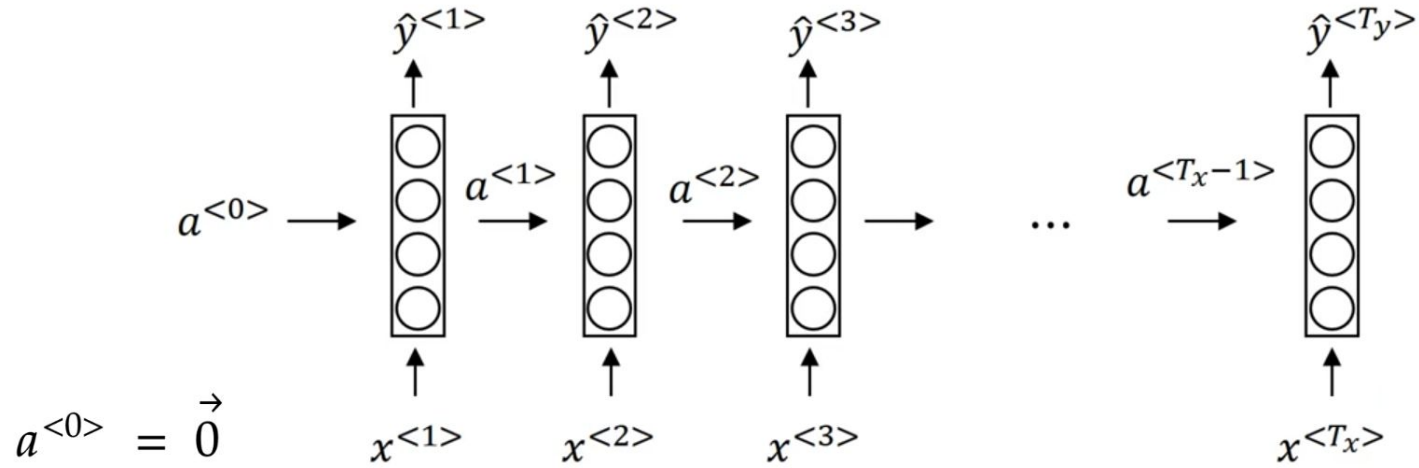
RNN - Diagram



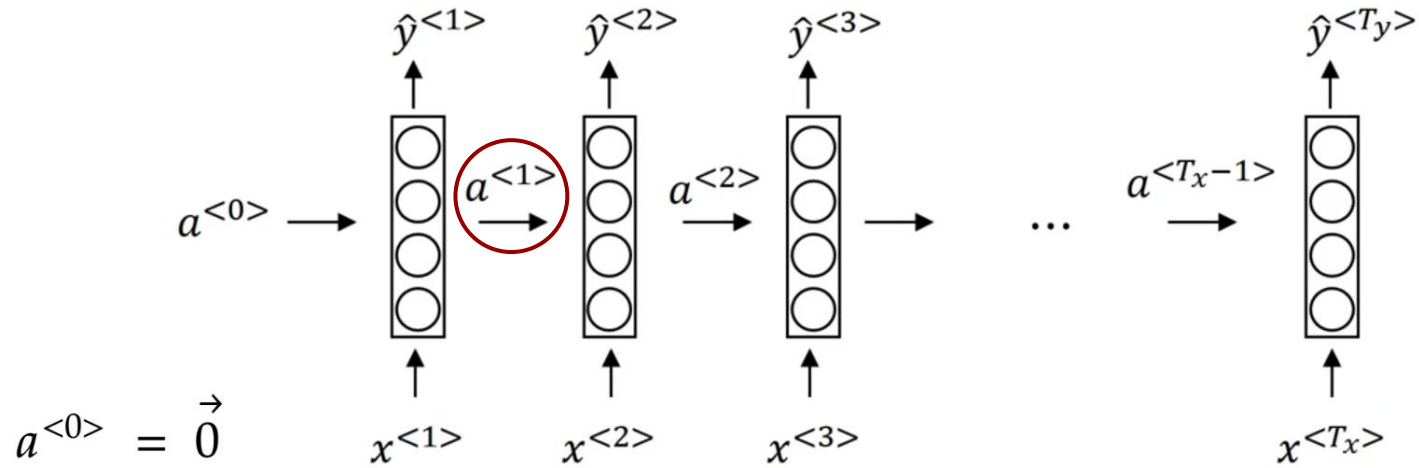
RNN - Diagram



RNN - Diagram



RNN - Diagram

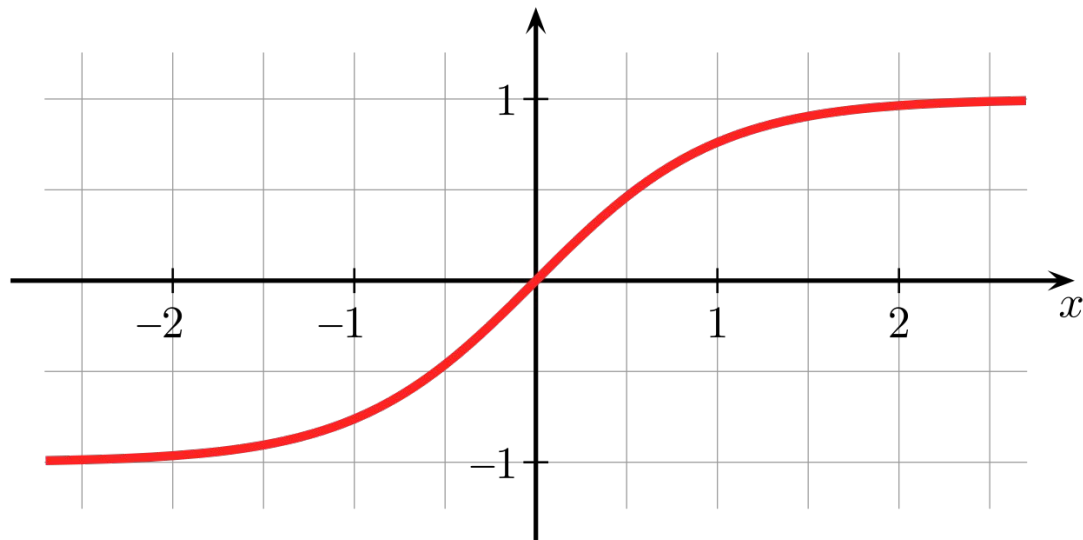


$$a^{<1>} = g(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$



Hyperbolic Tangent

$$g(x) = \tanh(x)$$

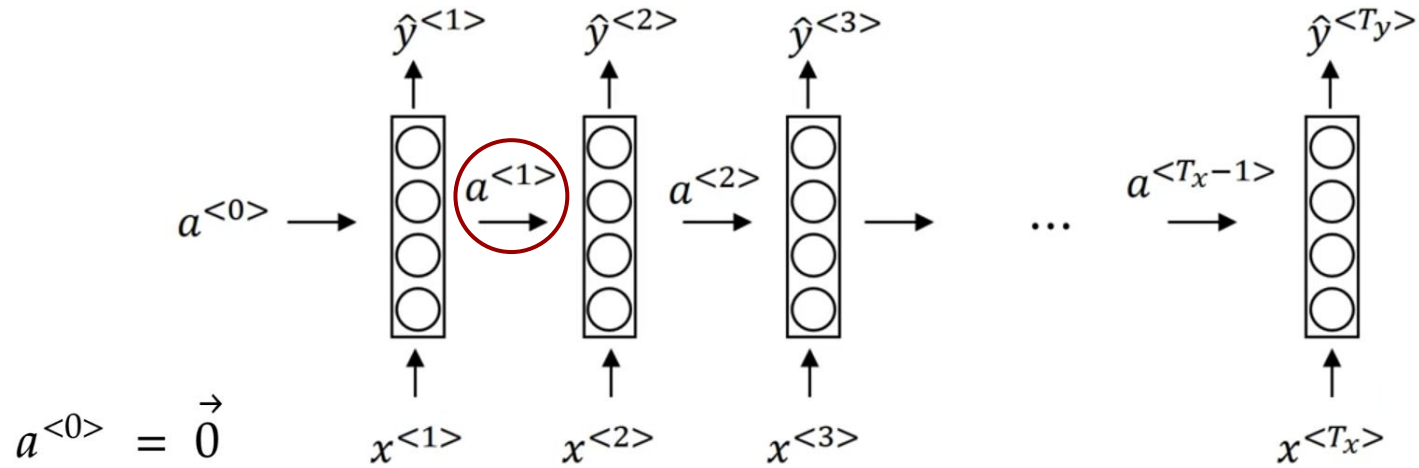


Source:

https://en.wikipedia.org/wiki/Hyperbolic_functions



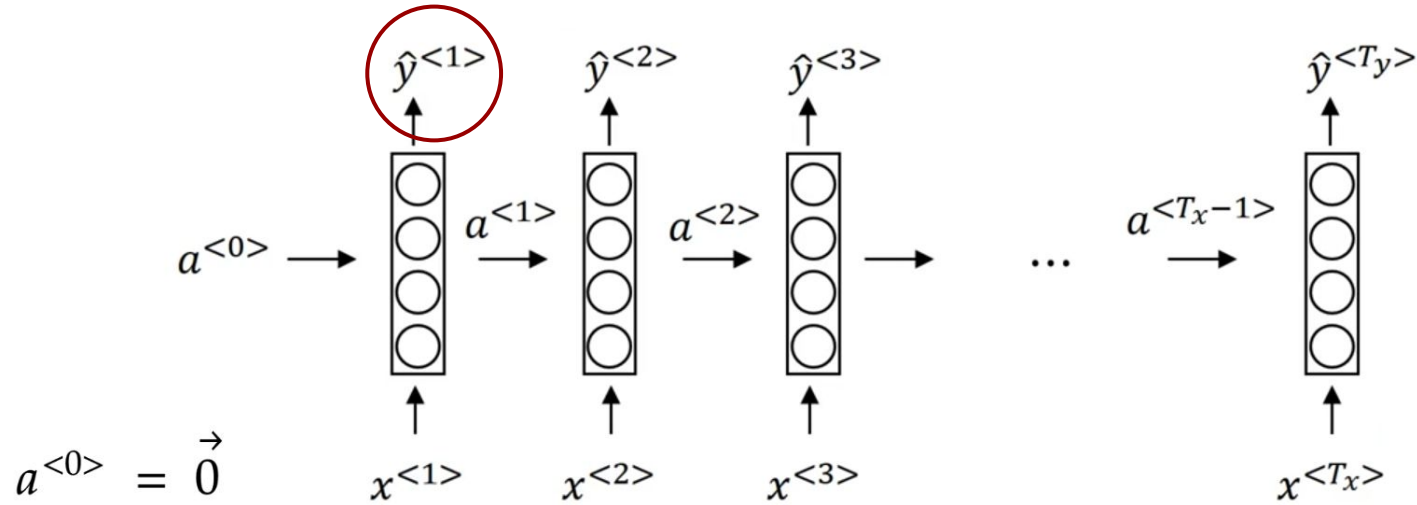
RNN - Diagram



$$a^{<1>} = g(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$



RNN - Diagram

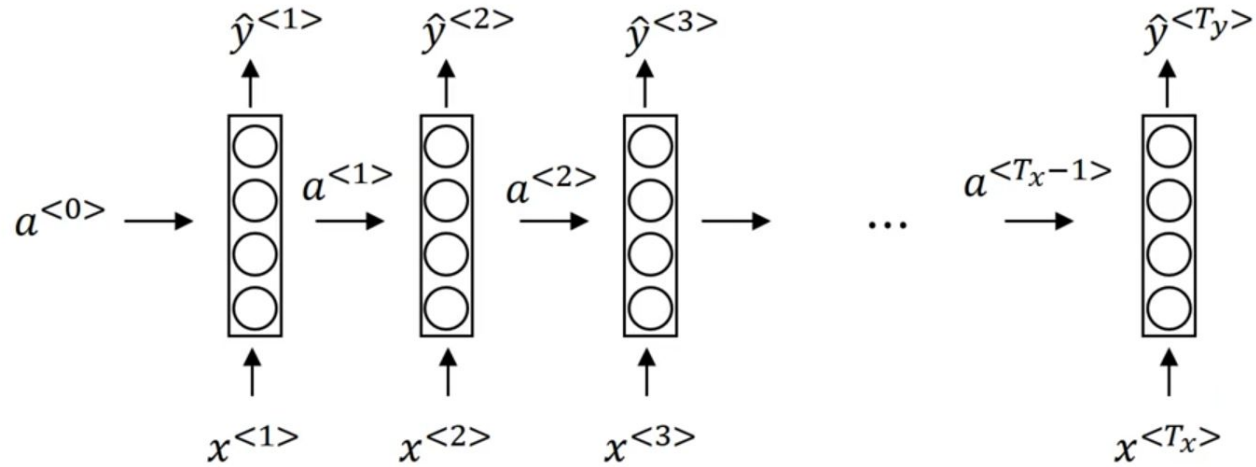


$$a^{<1>} = g(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g(W_{ya}a^{<1>} + b_y)$$



RNN - Diagram

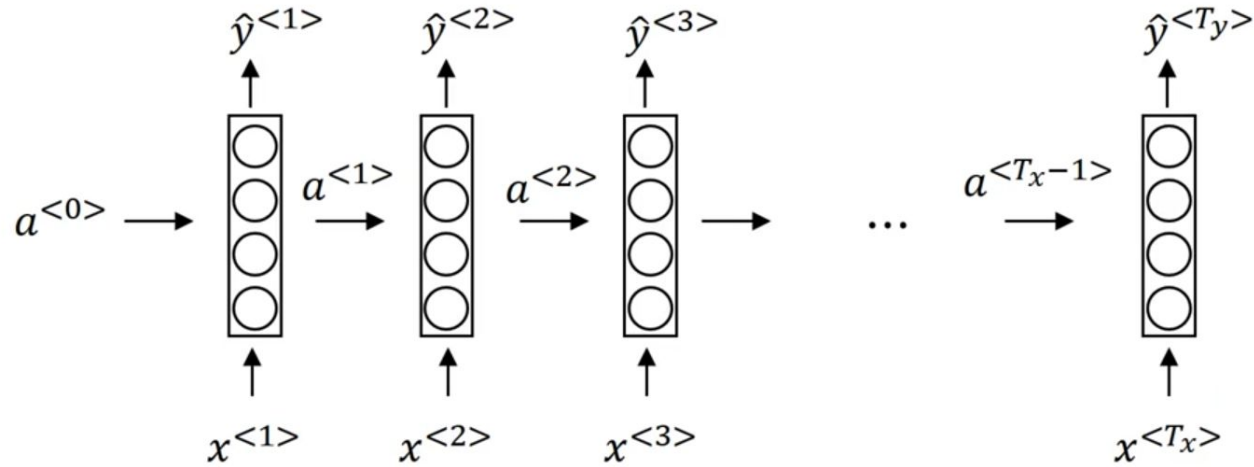


$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$



RNN - Diagram

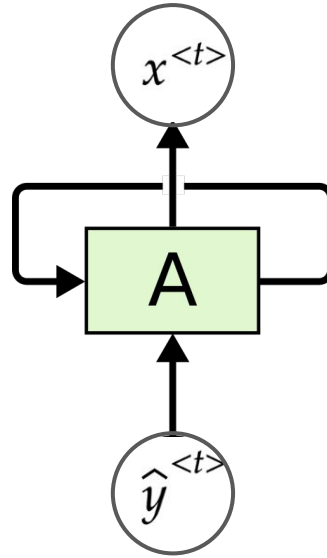


$$a^{<t>} = g(\underline{W_{aa}}a^{<t-1>} + \underline{W_{ax}}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(\underline{W_{ya}}a^{<t>} + b_y)$$



RNN - Diagram

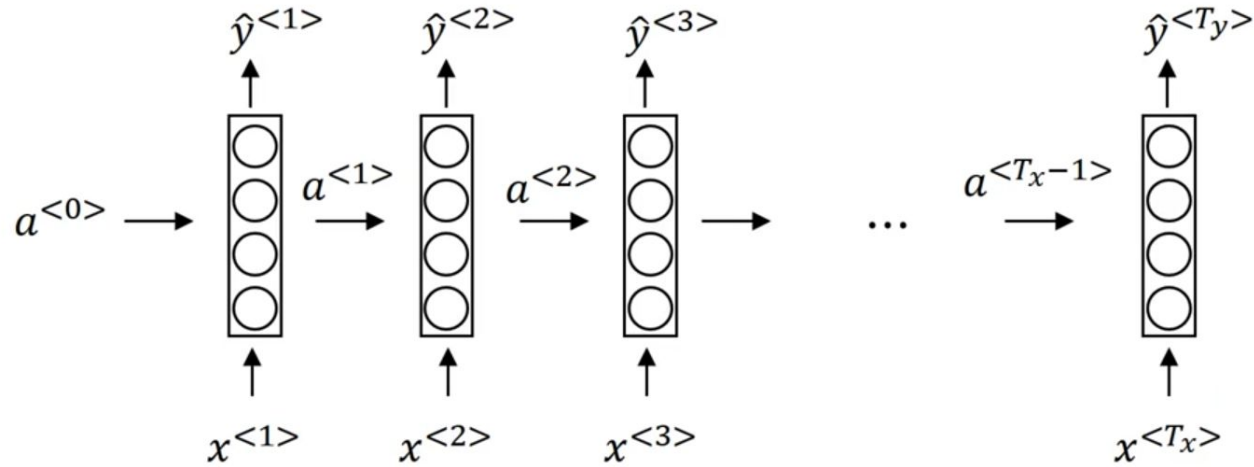


$$a^{<t>} = g(\underline{W_{aa}}a^{<t-1>} + \underline{W_{ax}}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(\underline{W_{ya}}a^{<t>} + b_y)$$



RNN - Diagram



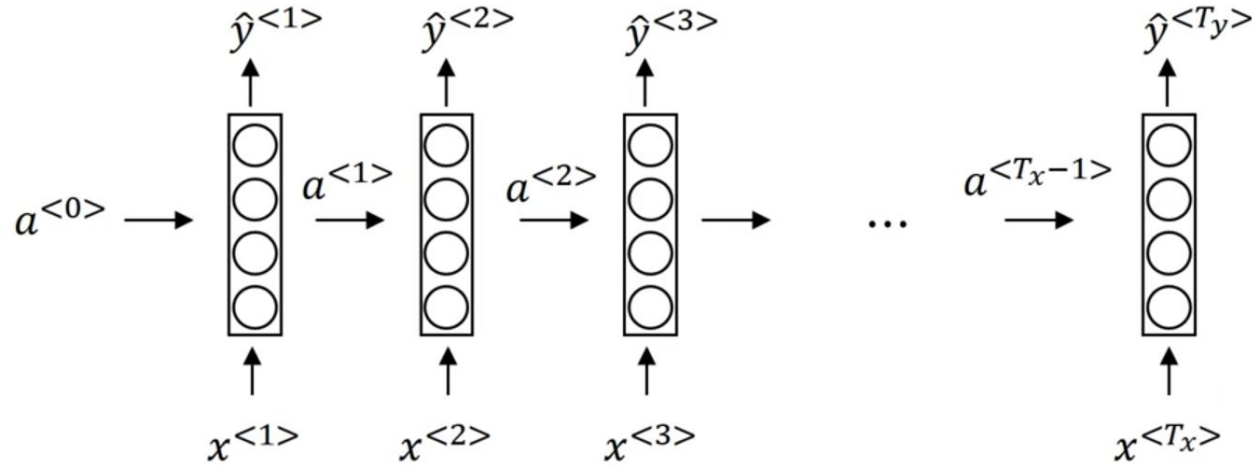
$$a^{<t>} = g(\underline{W_{aa}}a^{<t-1>} + \underline{W_{ax}}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(\underline{W_{ya}}a^{<t>} + b_y)$$



RNN - Diagram

$$\text{Loss}(\hat{y}^{<T_y>}, y^{<T_y>})$$

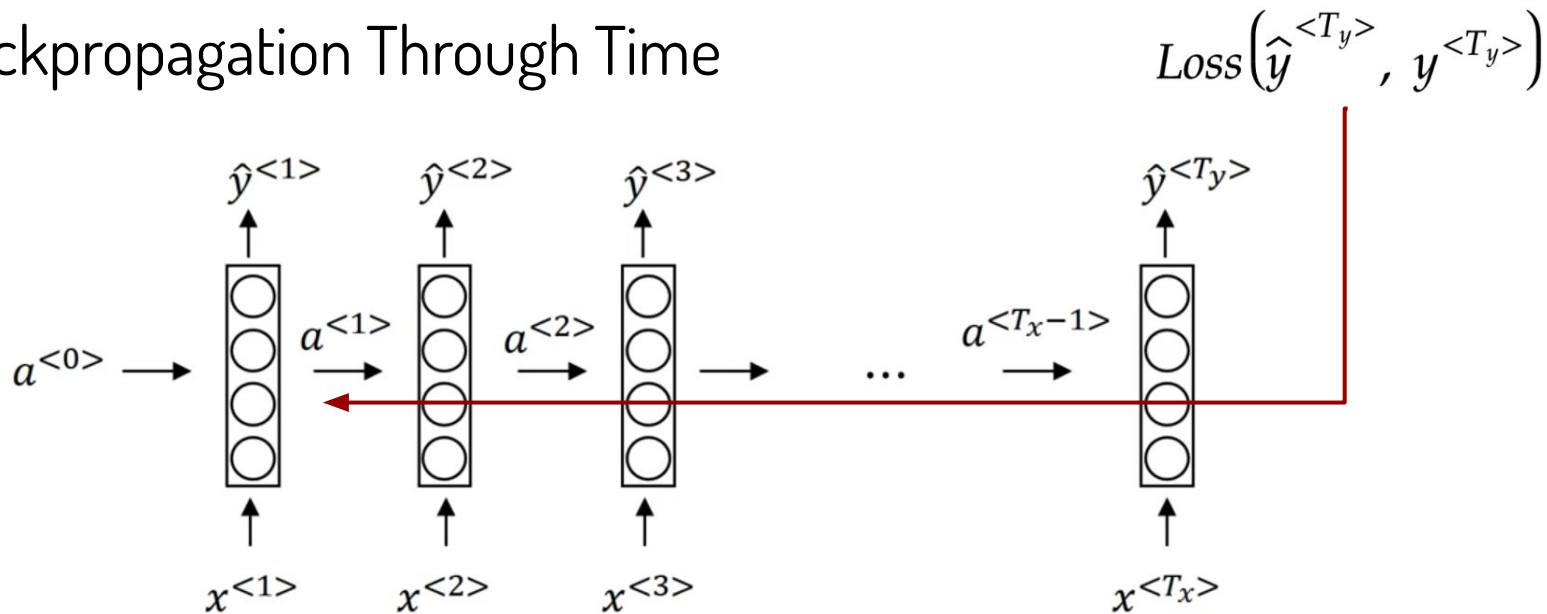


$$a^{<t>} = g(\underline{W_{aa}}a^{<t-1>} + \underline{W_{ax}}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(\underline{W_{ya}}a^{<t>} + b_y)$$



RNN - Backpropagation Through Time

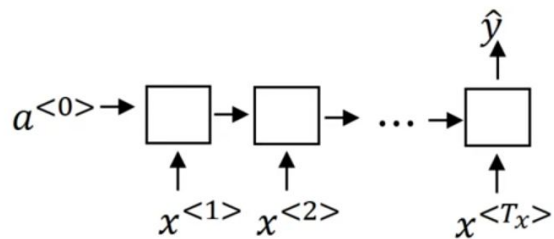


$$a^{<t>} = g(\underline{W_{aa}} a^{<t-1>} + \underline{W_{ax}} x^{<t>} + b_a)$$

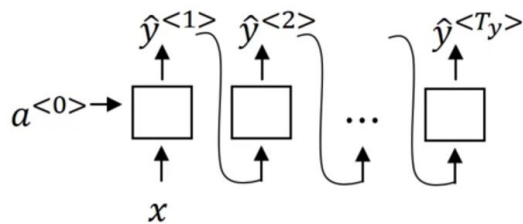
$$\hat{y}^{<t>} = g(\underline{W_{ya}} a^{<t>} + b_y)$$



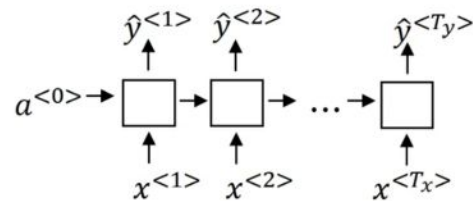
RNN Types & Variations



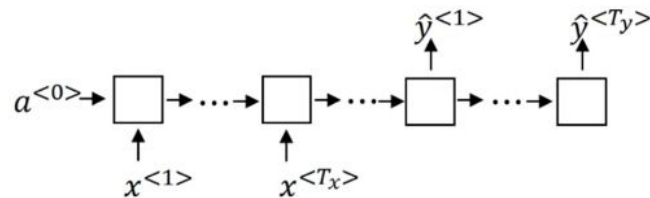
Many to one



One to many



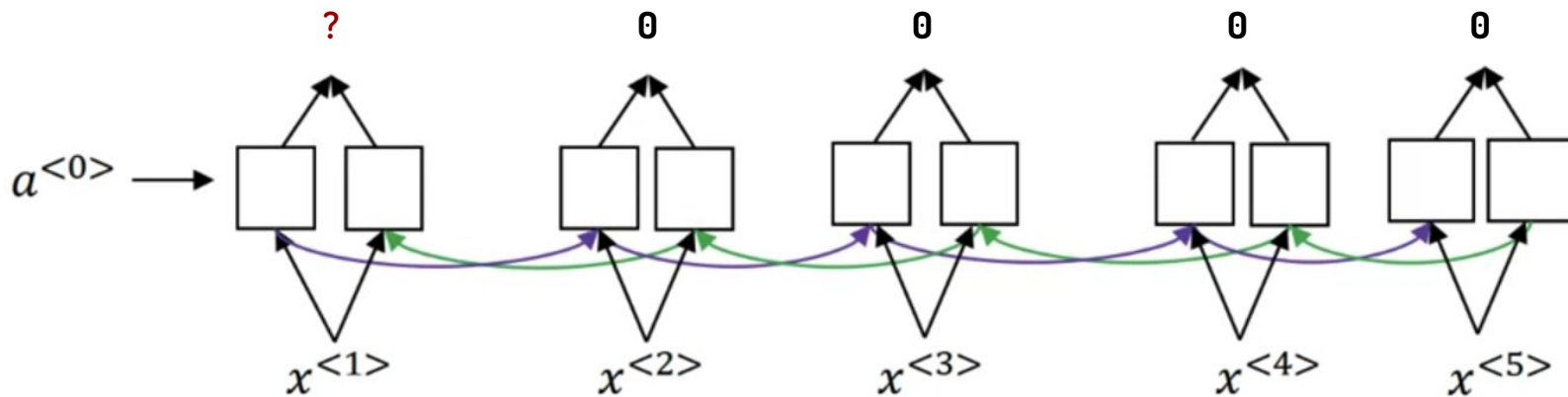
Many to many



Many to many



RNN Types & Variations - Bidirectional



$x^{<1>}$

$x^{<2>}$

$x^{<3>}$

$x^{<4>}$

$x^{<5>}$

1. **Teddy**
2. **Teddy**

bears
Roosevelt

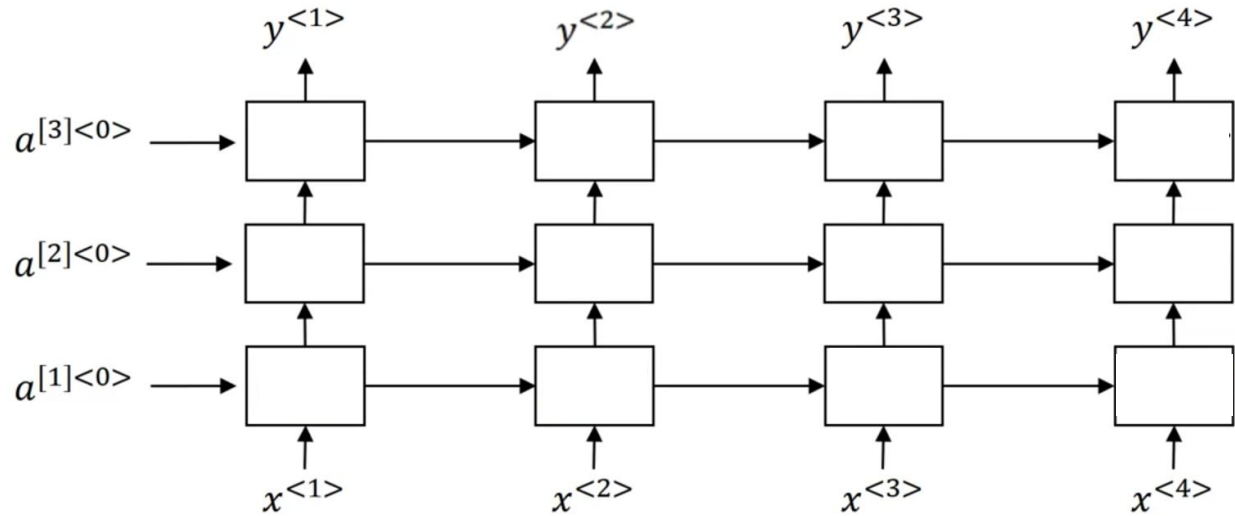
are
was

on
a

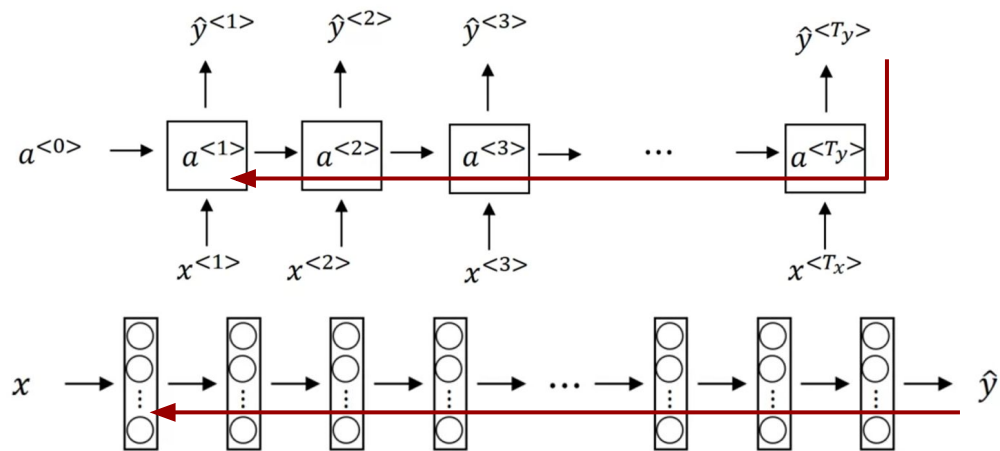
sale
president



RNN Types & Variations - Deep



RNN Problems



RNN Problems - Vanishing and Exploding Gradients

- $0.9^{**} 100 = 2.6561e-05$
- $0.9^{**} 300 = 1.8739e-14$
- Gradient Clipping - solution to explosion

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$



LSTM



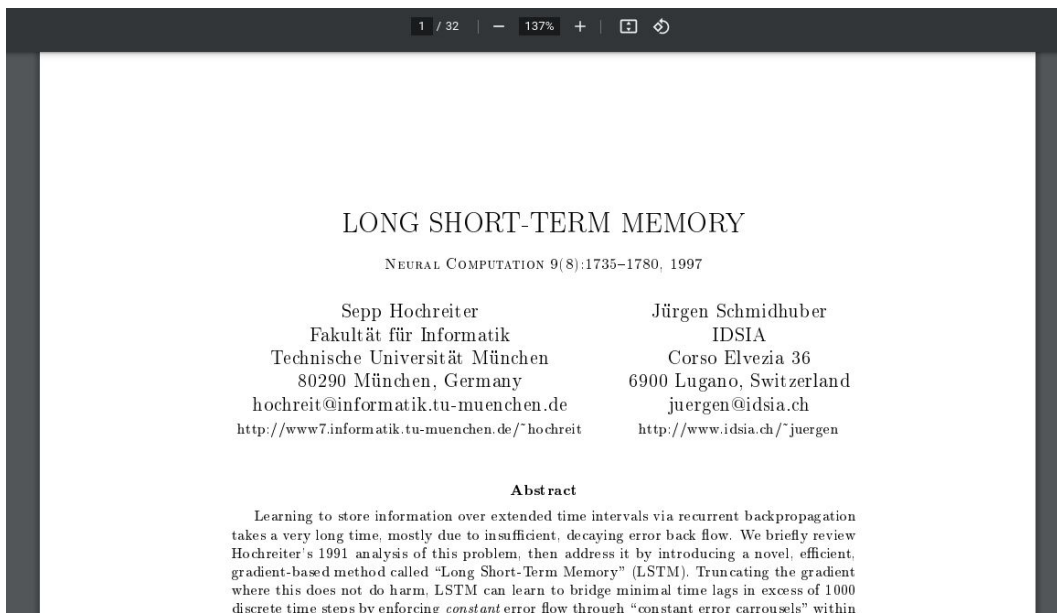
Long **S**hort-**T**erm **M**emory



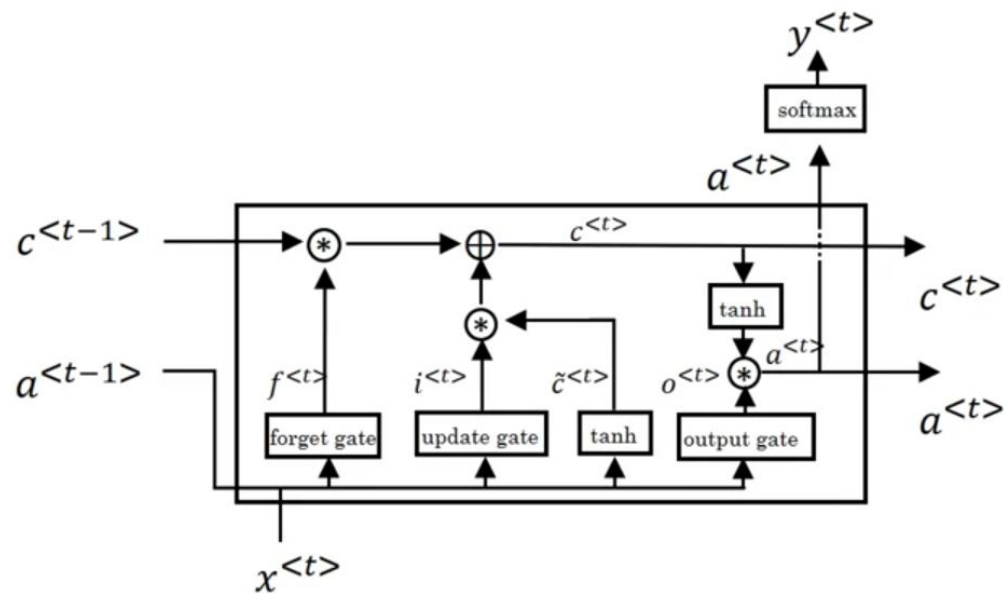
LSTM

- A kind of RNN
- Addressing vanishing gradients

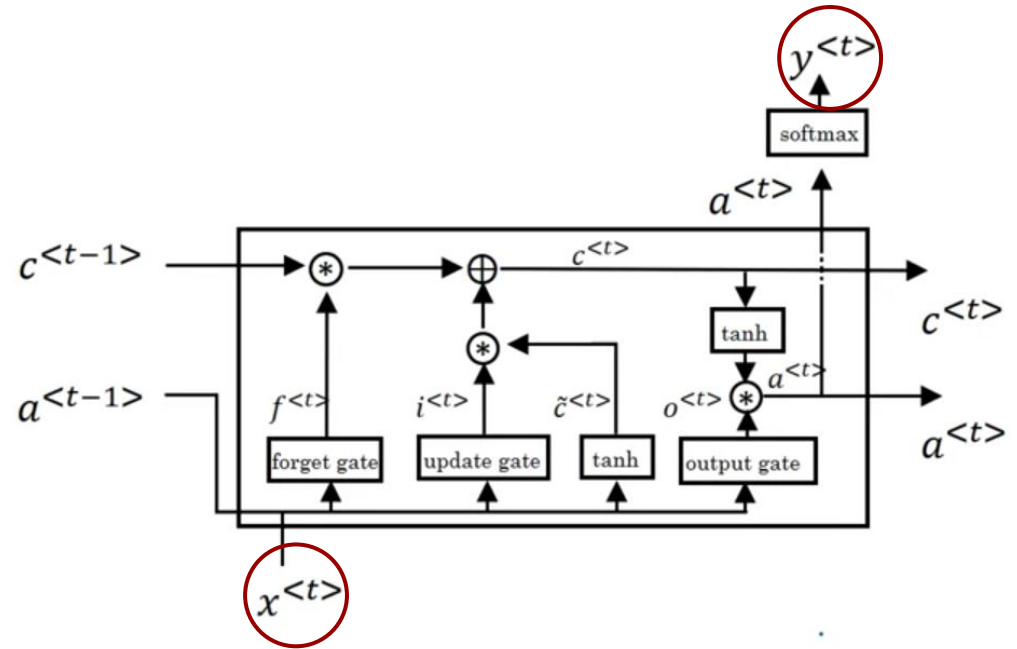
Source: https://www.researchgate.net/publication/13853244_Long_Short-Term_Memory



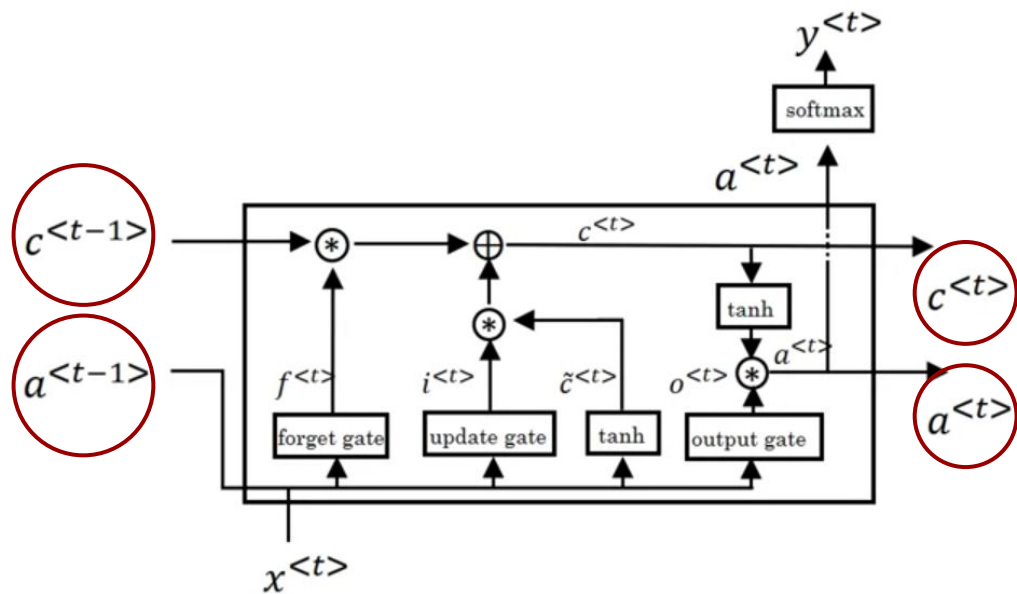
LSTM - Diagram



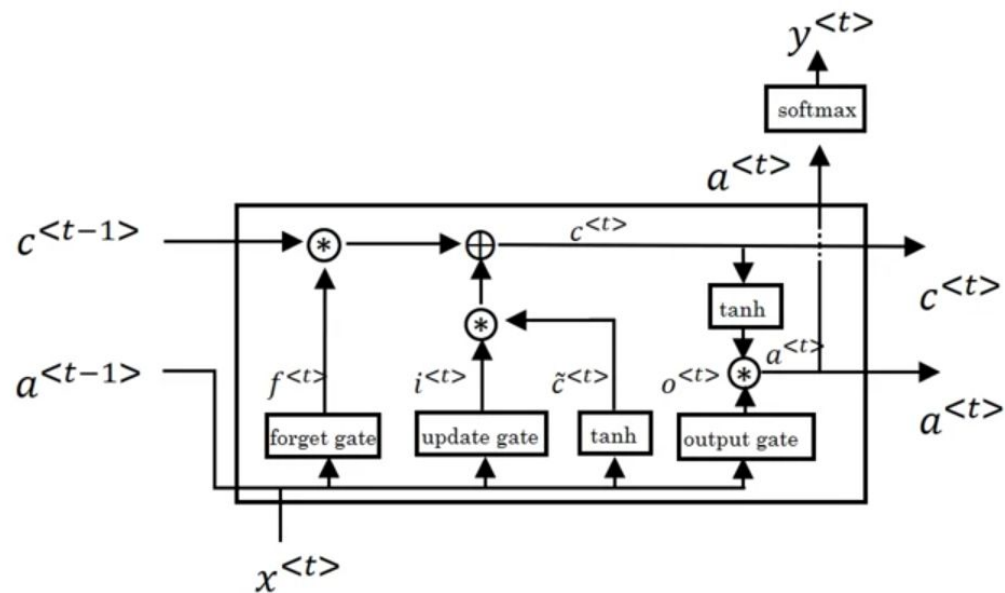
LSTM - Diagram



LSTM - Diagram

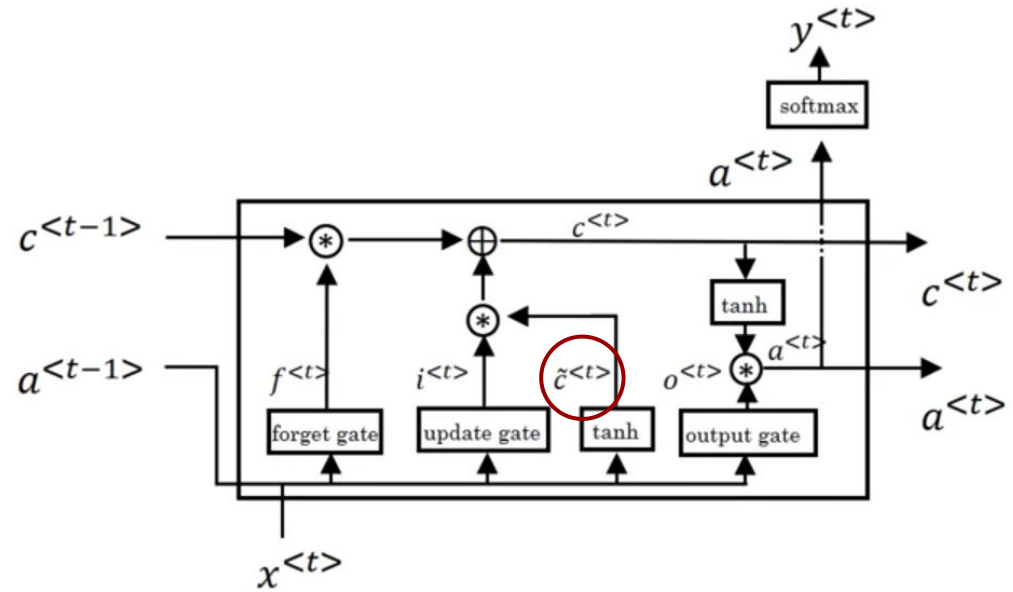


LSTM - Diagram



LSTM - Diagram

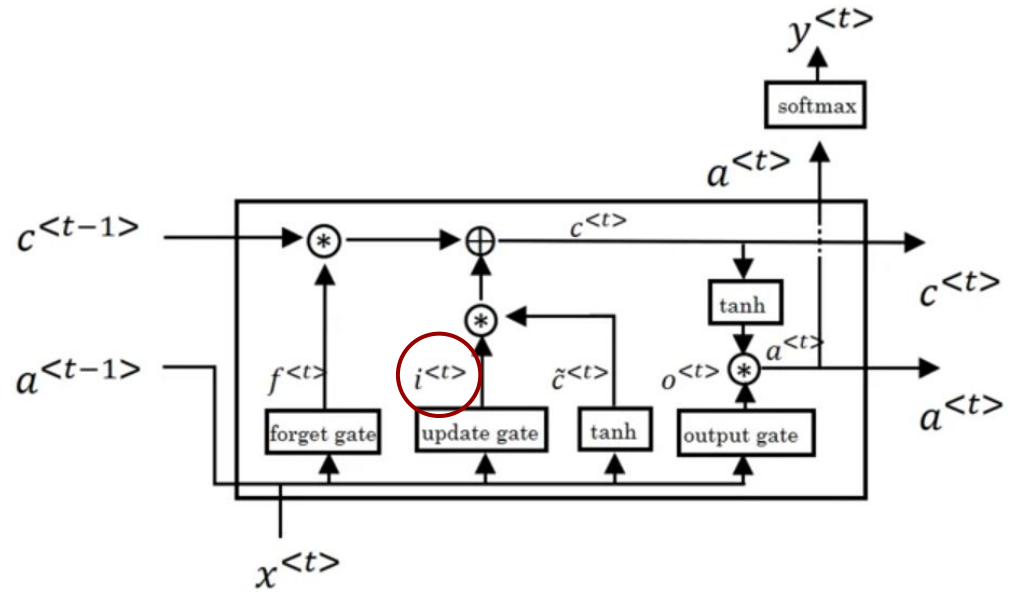
$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$



LSTM - Diagram

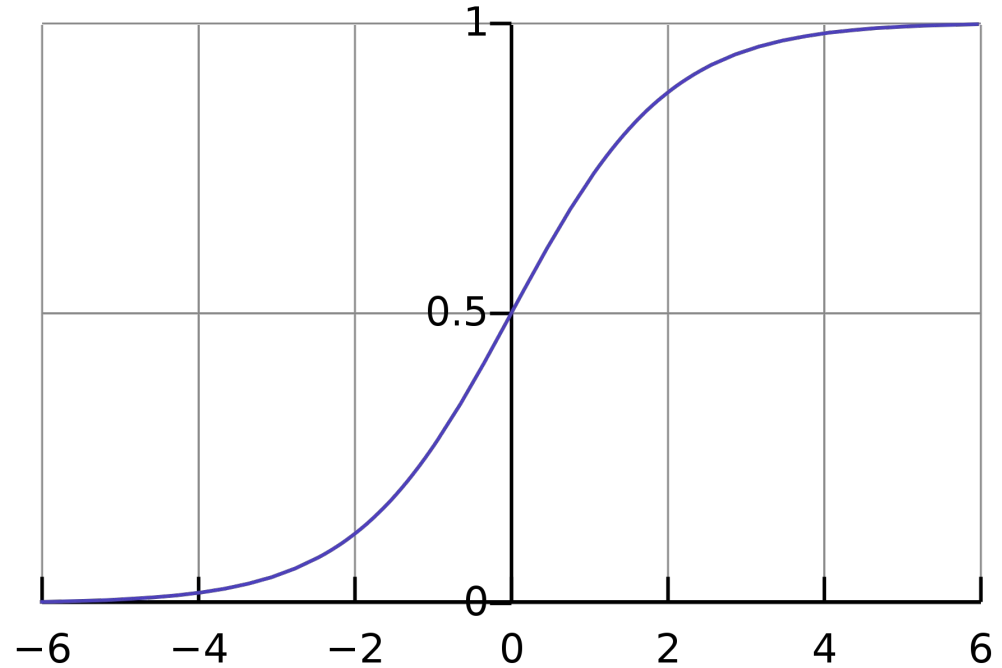
$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$



Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Source:

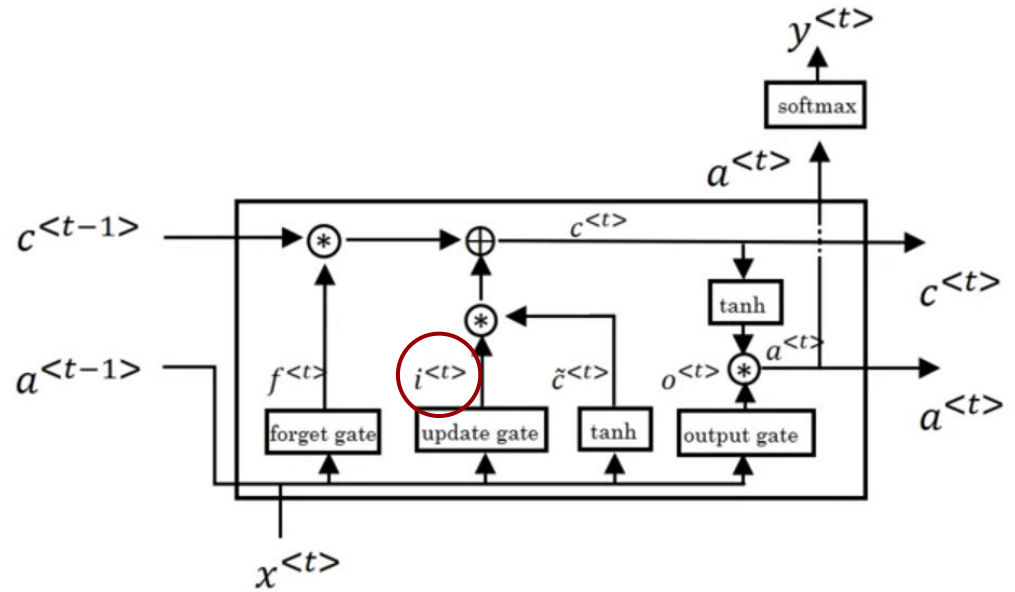
https://en.wikipedia.org/wiki/Sigmoid_function



LSTM - Diagram

$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

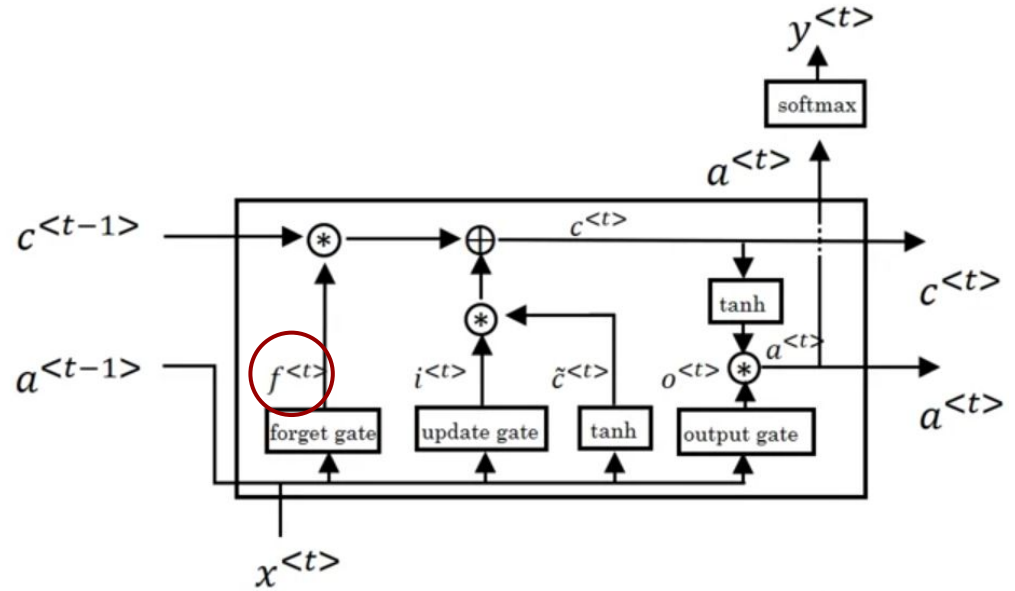


LSTM - Diagram

$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

$$\Gamma_f = \sigma(W_{fa}a^{<t-1>} + W_{fx}x^{<t>} + b_f)$$



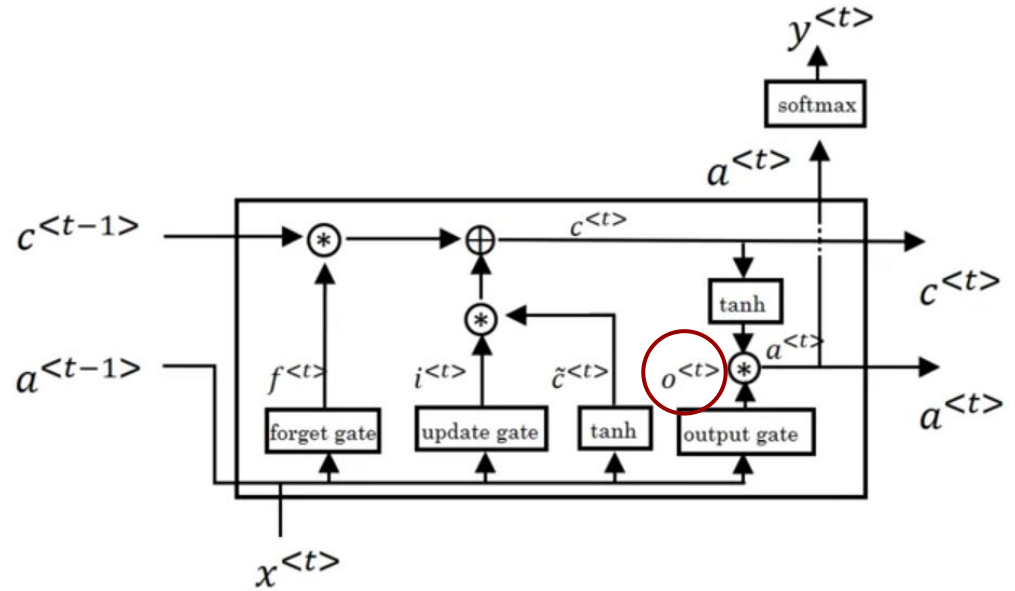
LSTM - Diagram

$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

$$\Gamma_f = \sigma(W_{fa}a^{<t-1>} + W_{fx}x^{<t>} + b_f)$$

$$\Gamma_o = \sigma(W_{oa}a^{<t-1>} + W_{ox}x^{<t>} + b_o)$$



LSTM - Diagram

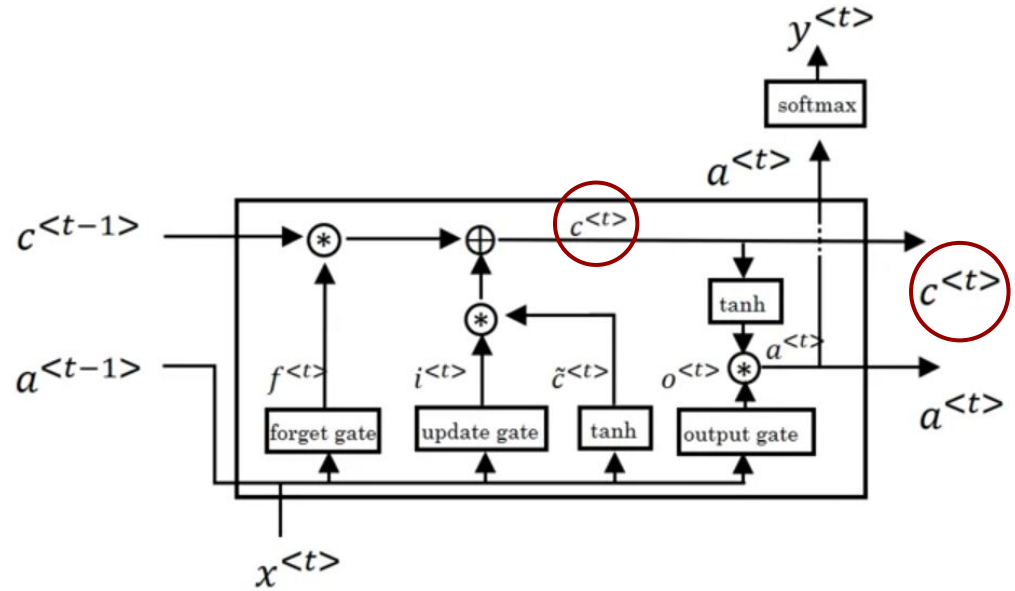
$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

$$\Gamma_f = \sigma(W_{fa}a^{<t-1>} + W_{fx}x^{<t>} + b_f)$$

$$\Gamma_o = \sigma(W_{oa}a^{<t-1>} + W_{ox}x^{<t>} + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$



LSTM - Diagram

$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

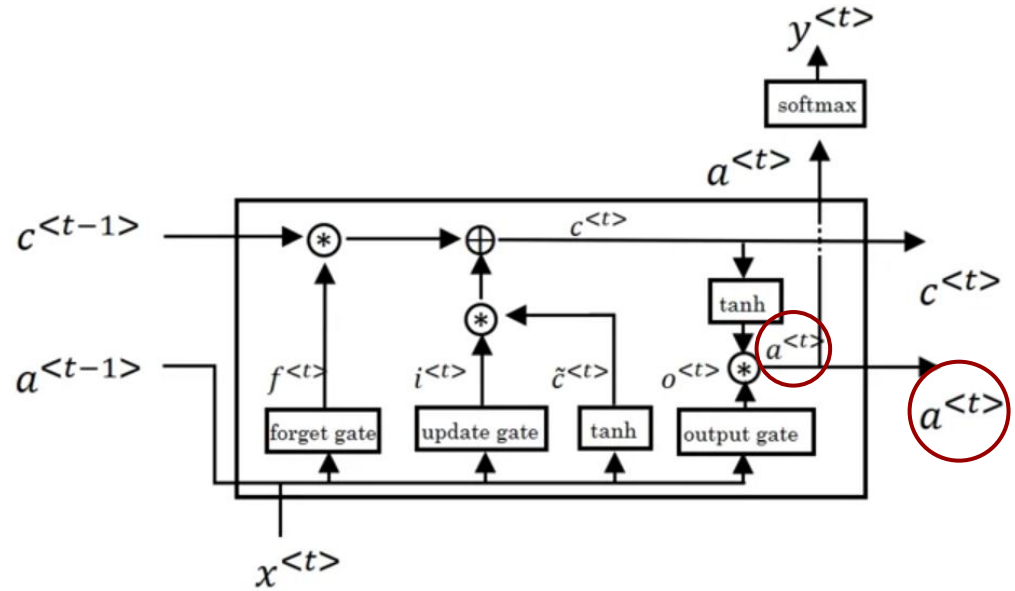
$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

$$\Gamma_f = \sigma(W_{fa}a^{<t-1>} + W_{fx}x^{<t>} + b_f)$$

$$\Gamma_o = \sigma(W_{oa}a^{<t-1>} + W_{ox}x^{<t>} + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$



LSTM - Diagram

$$\tilde{c}^{<t>} = \tanh(W_{ca}a^{<t-1>} + W_{cx}x^{<t>} + b_c)$$

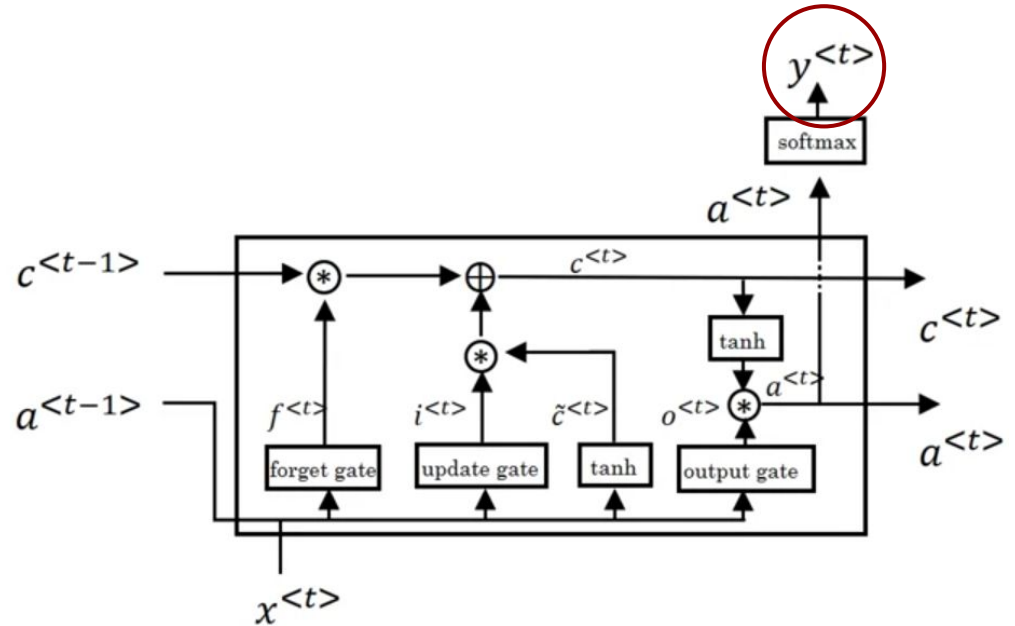
$$\Gamma_u = \sigma(W_{ua}a^{<t-1>} + W_{ux}x^{<t>} + b_u)$$

$$\Gamma_f = \sigma(W_{fa}a^{<t-1>} + W_{fx}x^{<t>} + b_f)$$

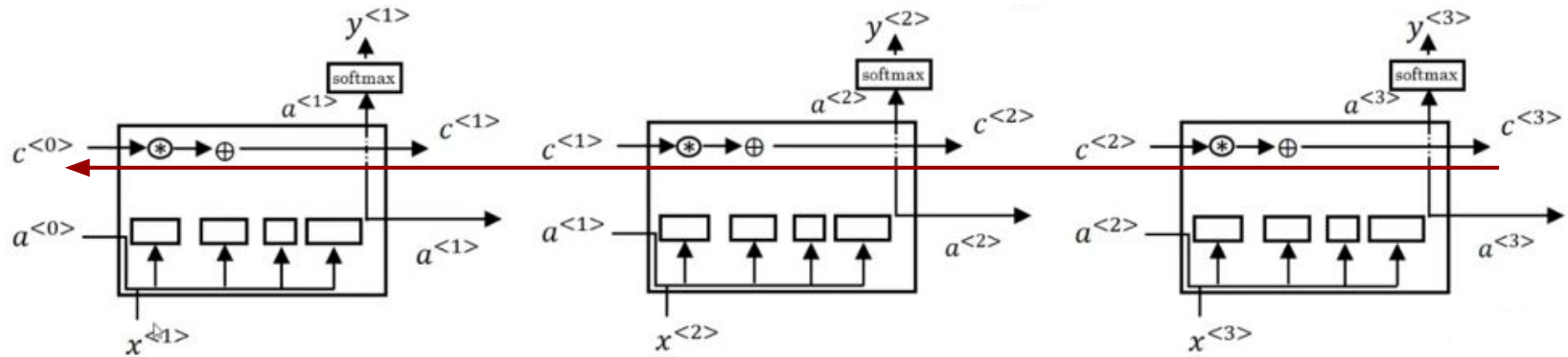
$$\Gamma_o = \sigma(W_{oa}a^{<t-1>} + W_{ox}x^{<t>} + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$



LSTM - Diagram



LSTM - Problems

- Difficult to train
 - Very long gradient paths
- Transfer learning doesn't really work
 - New dataset for every task
- Slow
 - Parallelization is impossible



Transformer

- Proposed in 2017 addressing machine translation
- Based on attention mechanism
- Parallelization is now possible
- **Transfer learning works**



<https://jalammar.github.io/illustrated-transformer/>



Source:

<https://transformers.fandom.com/pl/wiki/Bumblebee>



Kahoot time!

Kahoot!



Thank you & Q&A

