



# Deep Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



[https://github.com/safayani/deep\\_learning\\_course](https://github.com/safayani/deep_learning_course)

# Examples of sequence data

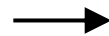
Speech recognition



“The quick brown fox jumped  
over the lazy dog.”

Music generation

∅



Sentiment classification

“There is nothing to like  
in this movie.”



DNA sequence analysis

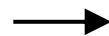
AGCCCCTGTGAGGAAGTAG



AG**CCCCTGTGAGGAAGTAG**

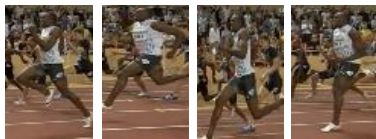
Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Video activity recognition



Running

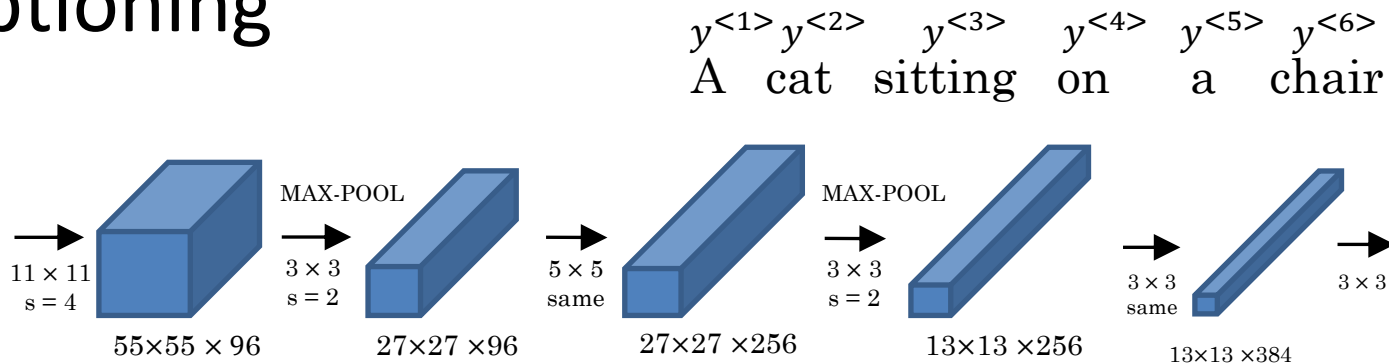
Name entity recognition

Yesterday, Harry Potter  
met Hermione Granger.

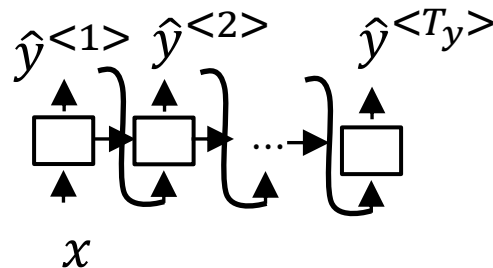
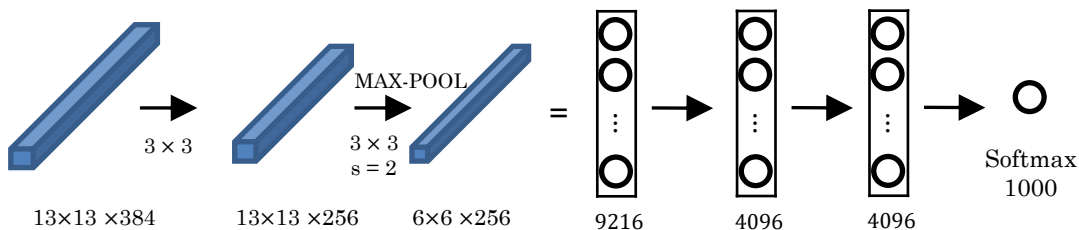


Yesterday, **Harry Potter**  
met **Hermione Granger**.

# Image captioning



$y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>} y^{<6>}$   
A cat sitting on a chair



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

# Recurrent Neural Networks

- Now, given this input X let's say that you want a model to operate Y that has one outputs per input word and the target output the design Y tells you for each of the input words is that part of a person's name.

- X: (Ali Ahmadi ) and (Hassan Hamidi ) invented a new drug.

$$X^{(1)} \quad X^{(2)} \quad X^{(3)} \quad \dots \quad \dots \quad X^{(t)} \quad \dots \quad \dots \quad X^{(9)} \quad T_x=9$$

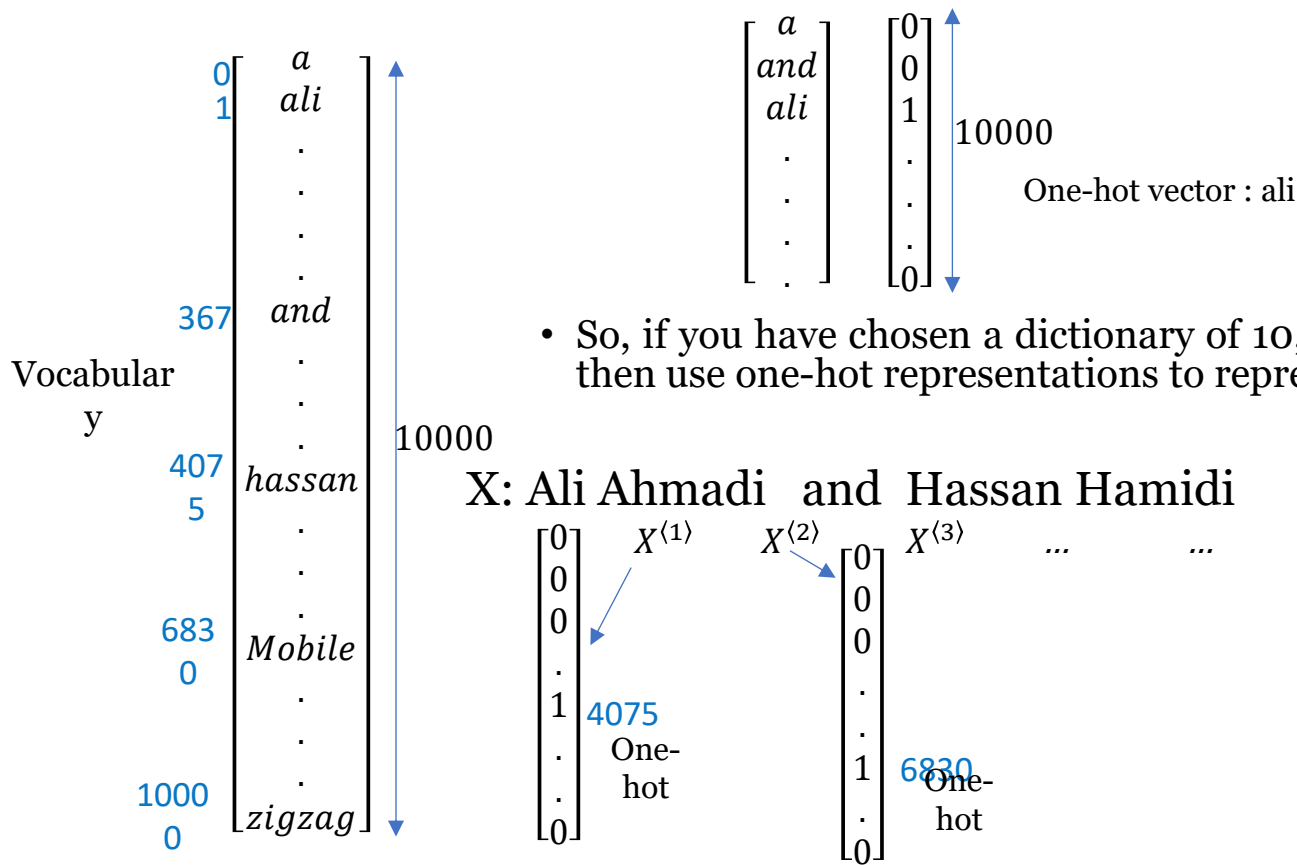
- y:  $\begin{matrix} 1 & 1 & 0 & 1 & & 1 & & 0 & 0 & 0 & 0 \\ y^{(1)} & y^{(2)} & y^{(3)} & & & & & & & & y^{(9)} \end{matrix} \quad T_y=9$

- $X^{(i)(t)}$   
  
 داده آموزشی i ام
- $T_x^{(i)} = 9$   
 $T_y^{(i)}$

- This is our first serious foray into NLP or Natural Language Processing.

# Representing words

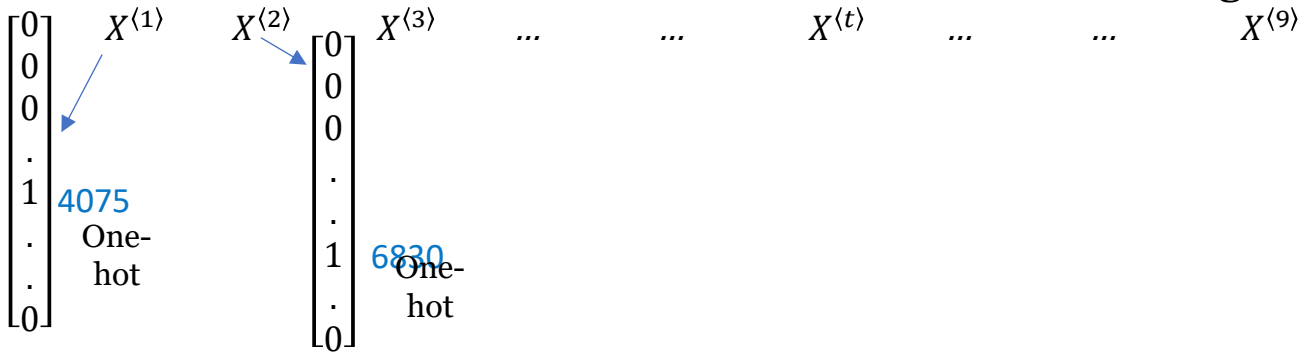
- Dictionary: 30000, 50000



- So, if you have chosen a dictionary of 10,000 words, what you can do is then use one-hot representations to represent each of these words.

X: Ali Ahmadi and Hassan Hamidi

invented a new drug.

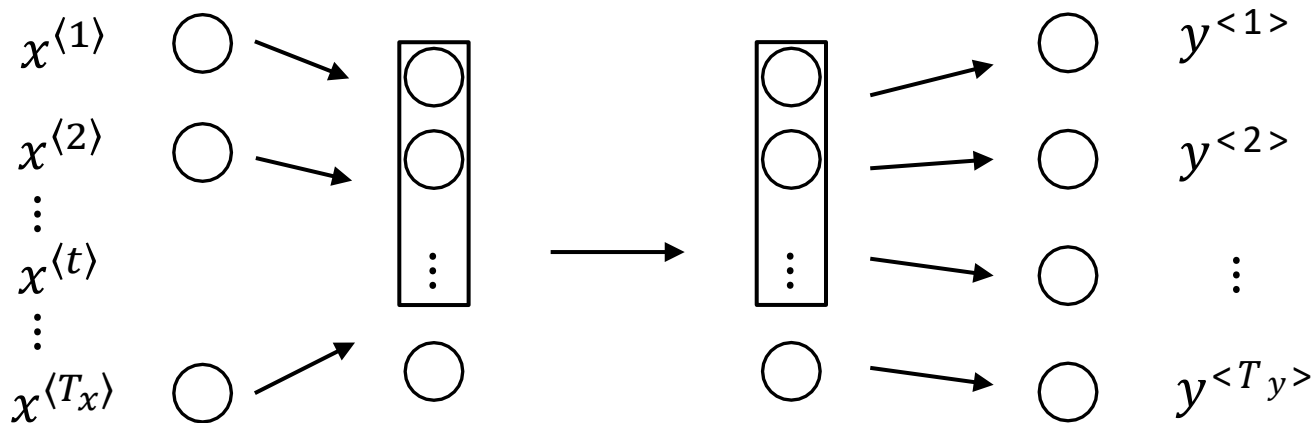


# Representing words

- Utilizing a sequence model for supervised learning to map input  $X$  to output  $Y$ .
- Introduction of an "Unknown Word" token for handling out-of-vocabulary words.
- Describing a notation for training sets in sequence data.

# Recurrent Neural Network Model

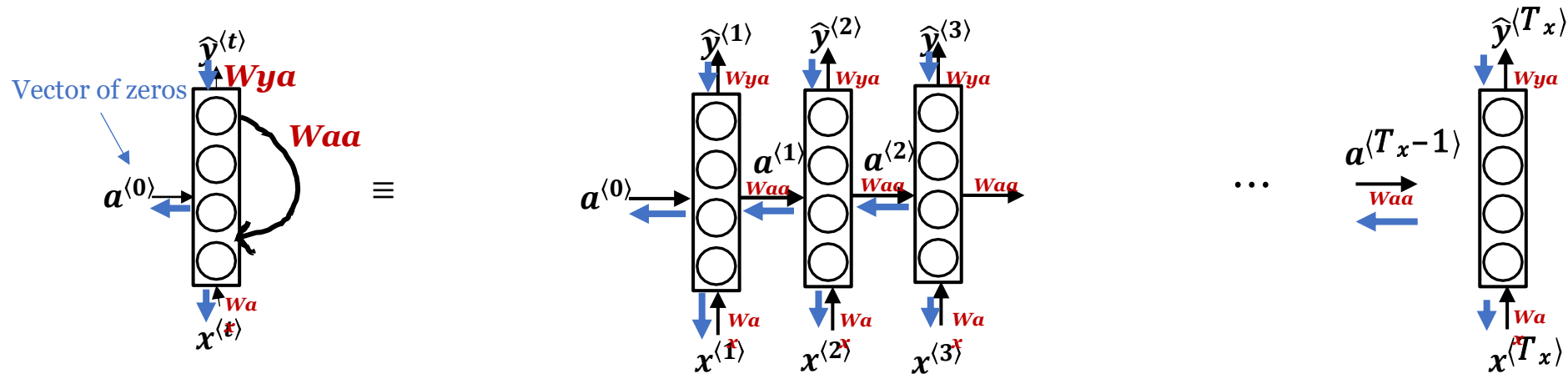
Why not a standard network?



Problems:

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

# Recurrent Neural Networks

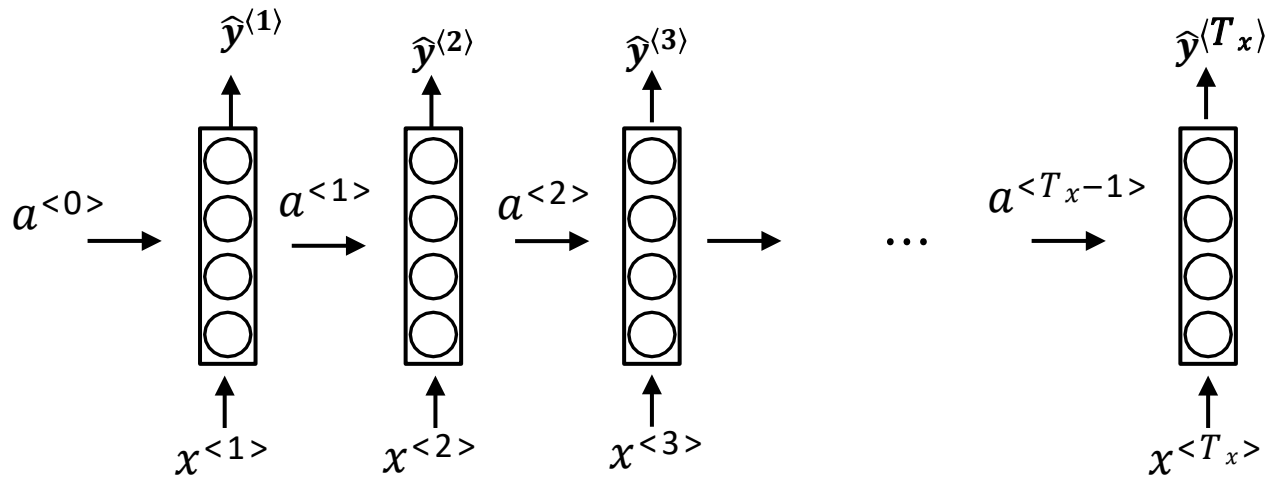


- Activation at time zero in neural networks is often initialized as a vector of zeros.
- Some researchers prefer initializing it as  $a^{(0)}$  randomly.
- There are alternative methods to initialize the activation at time zero

Backward propagation through time



# Forward Propagation



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g_1(w_{aa}a^{<0>} + w_{ax}x^{<1>} + b_a) \quad \leftarrow \text{tanh|Relu}$$

$$\hat{y}^{<1>} = g_2(w_{ya}a^{<1>} + b_y) \quad \leftarrow \text{sigmoid}$$

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

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

# Simplified RNN notation

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

Diagram annotations for the first equation:   
 -  $w_{aa}$  is annotated with  $(100, 100)$  and a double-headed arrow labeled 100.   
 -  $w_{ax}$  is annotated with  $(100, 10000)$  and a double-headed arrow labeled 10000.   
 -  $b_a$  is annotated with 100.

$$a^{(t)} = g(w_a[a^{(t-1)}, x^{(t)}] + b_a)$$

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

$$\bullet a^{(t)} = g(w_a[a^{(t-1)'}', u^{(t)'}'] + b_a)$$

$$\begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

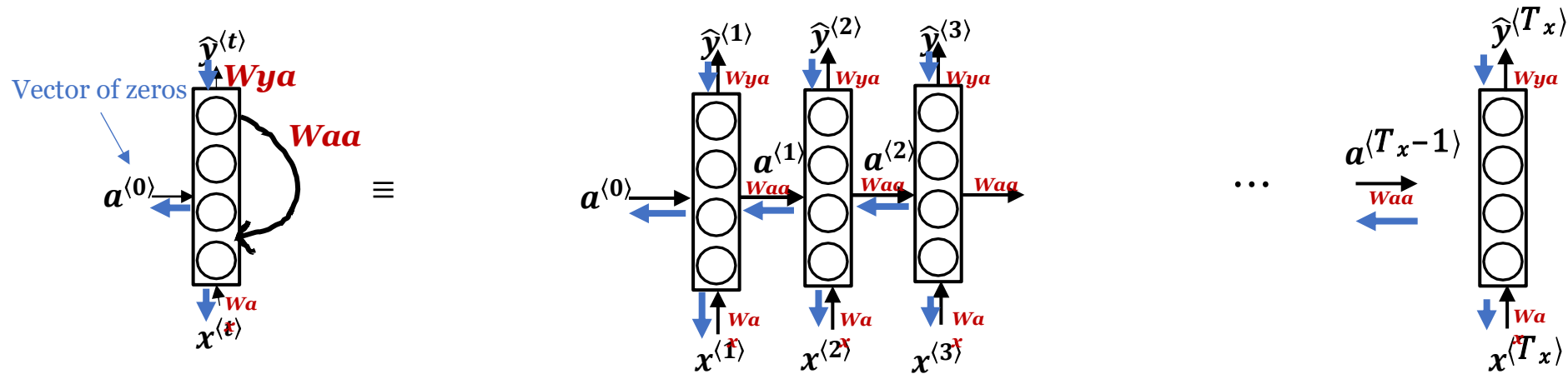
Diagram annotations for the vector:   
 -  $a^{(t-1)}$  has a vertical double-headed arrow labeled 100.   
 -  $x^{(t)}$  has a vertical double-headed arrow labeled 10000.   
 - The entire vector has a vertical double-headed arrow labeled 10100.

$$w_a = [w_{aa} : w_{ax}]$$

Diagram annotations for the weight vector:   
 -  $w_{aa}$  has a vertical double-headed arrow labeled 100.   
 -  $w_{ax}$  has a vertical double-headed arrow labeled 100.   
 - The entire vector has a vertical double-headed arrow labeled 10000.   
 - The vector is annotated with  $(100, 10100)$  and horizontal double-headed arrows above it.

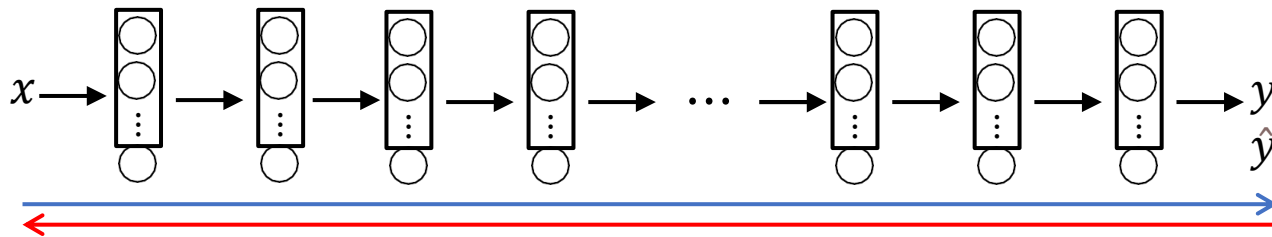
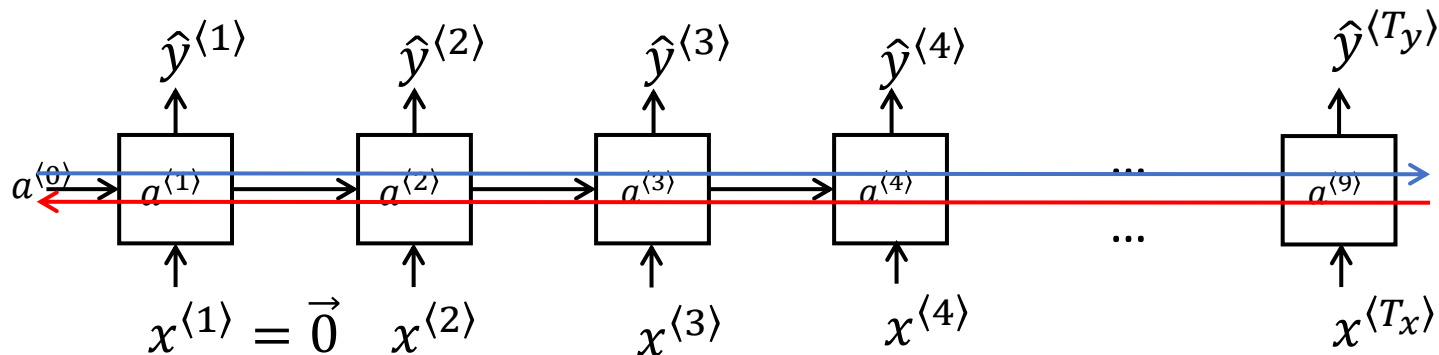
$$[w_{aa} : w_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa}a^{(t-1)} + w_{ax}x^{(t)}$$

# Recurrent Neural Networks



# Vanishing gradients with RNNs

- The cat which already ate bunch of food was full
- The cats ... were full



Exploding gradients.

NaN

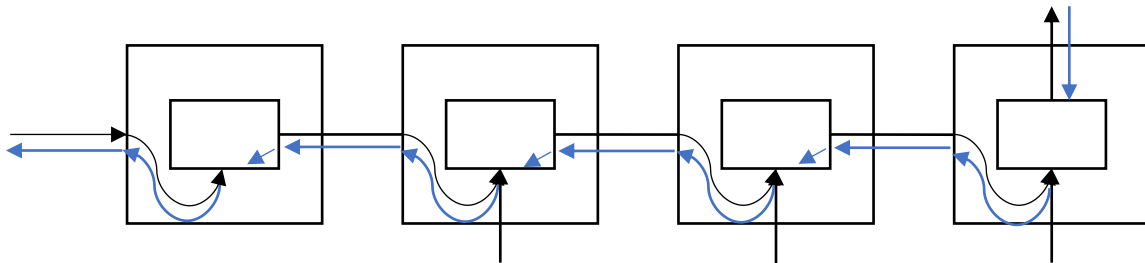
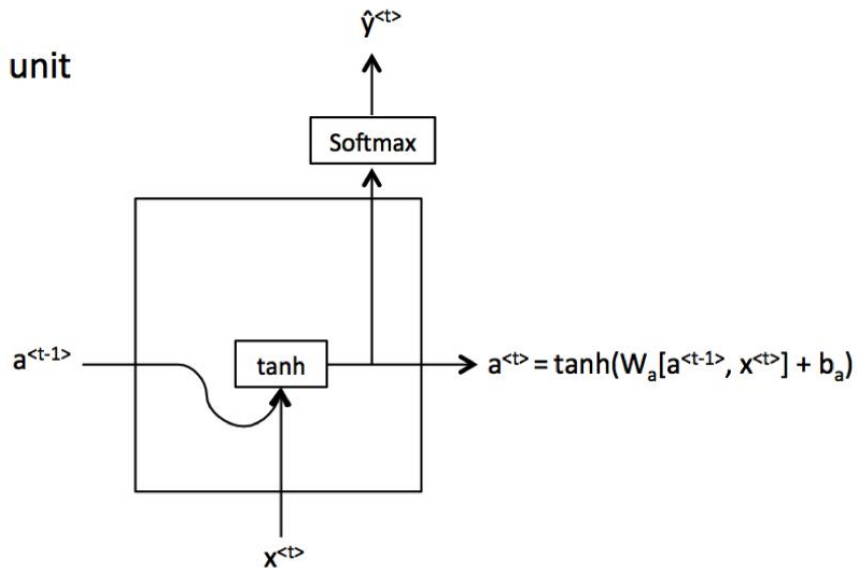
gradient clipping

# RNN Unit

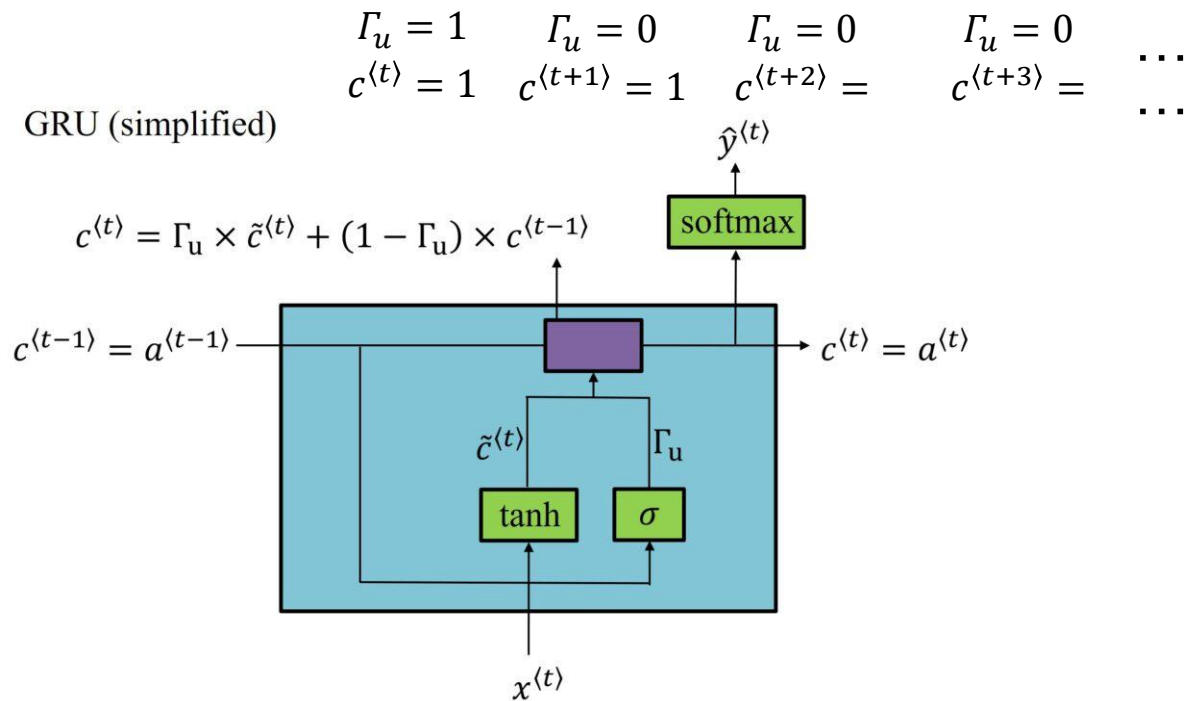
- $a^{(t)} = g(W_a[a^{(t-1)}, x^{(t)}] + b_a)$

↑  
tanh

RNN unit

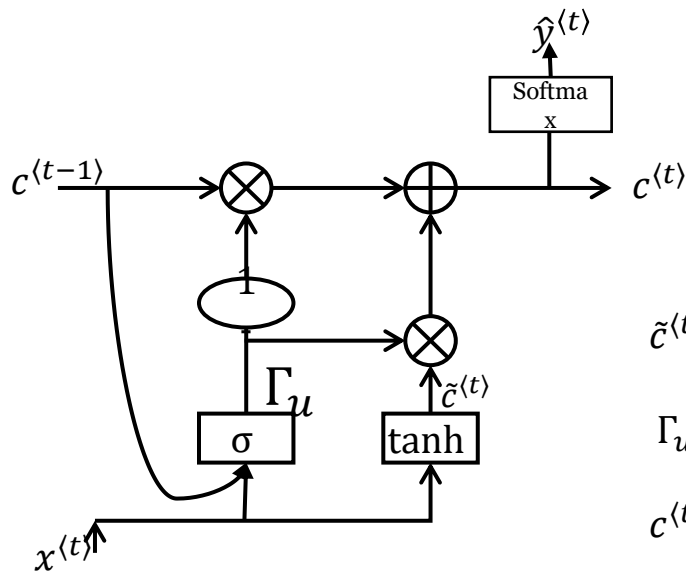


# GRU (simplified)



The cat, which already ate ..., was full.

# GRU (Simplified)



C : memory cell

$$\tilde{c}^{(t)} = \tanh(\omega_c[c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(w_u[c^{(t-1)}, x^{(t)}] + b_u)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

# Full GRU

- $\hat{c}^{(t)} = \tanh(w_c[\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c)$
- $\Gamma_u = \sigma(w_u[c^{(t-1)}, x^{(t)}] + b_u)$
- $\Gamma_r = \sigma(w_r[c^{(t-1)}, x^{(t)}] + b_r)$
- $c^{(t)} = \Gamma_u * \hat{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$

