

sequence data :-

Example :-

- Speech recognition
- Music generation
- Sentiment classification
- DNA sequence analysis
- Machine translation
- Video activity recognition
- Name entity recognition.

Notation :-

Let say we want to know the position of person name in a sentence.

$X$  :- Harry porter and Harmony Granger invented a new spell.

$x^{(1)}$     $x^{(2)}$     $x^{(3)}$    . . . . .  $x^{(9)}$

$Y$  :-   1   1   0   1   1   0   0   0   0

$y^{(1)}$     $y^{(2)}$     $y^{(3)}$    . . . . .  $y^{(9)}$

Each of the word is represented by  $x^{(i)}$  and their output  $y^{(i)}$

- Number of word  $\rightarrow T_x = 9$

$T_y = 9$

$x^{(i)(t)}$   $\rightarrow$   $i$ th training example  $t$ -th word.

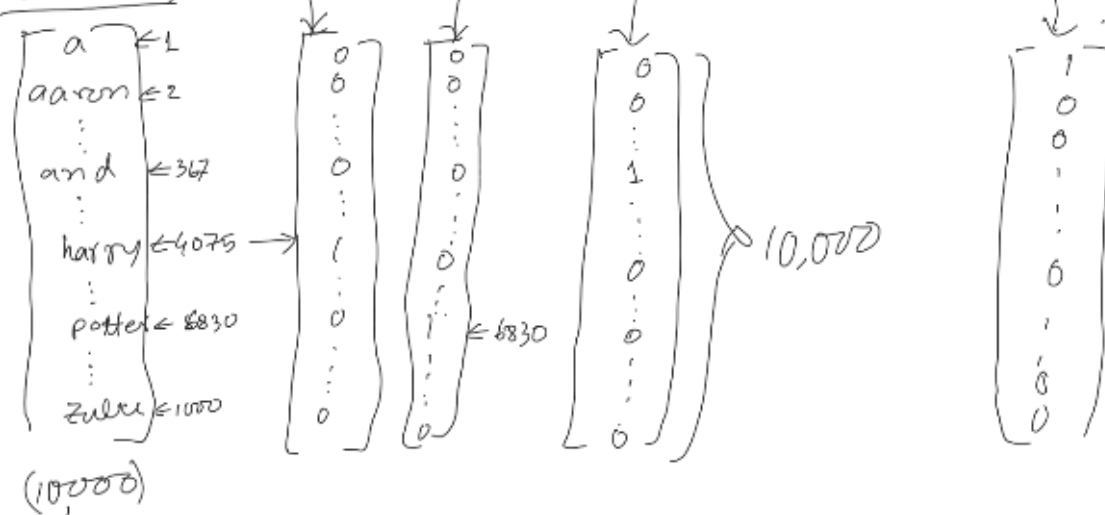
$T_x^{(i)}$   $\rightarrow$  length of  $i$ th training example.

## Representing words:-

We maintain a vocabulary of words & use a one-hot vector to represent each word of a sentence. One-hot means there will be one "1" in the vector & rest are zero. And that "1" will be in the same position of the word in the vocabulary.

X:- Harry Potter and Harmony Granger invented a new spell.

Vocabulary:-



This single line will be a  $(10000, 9)$  matrix.

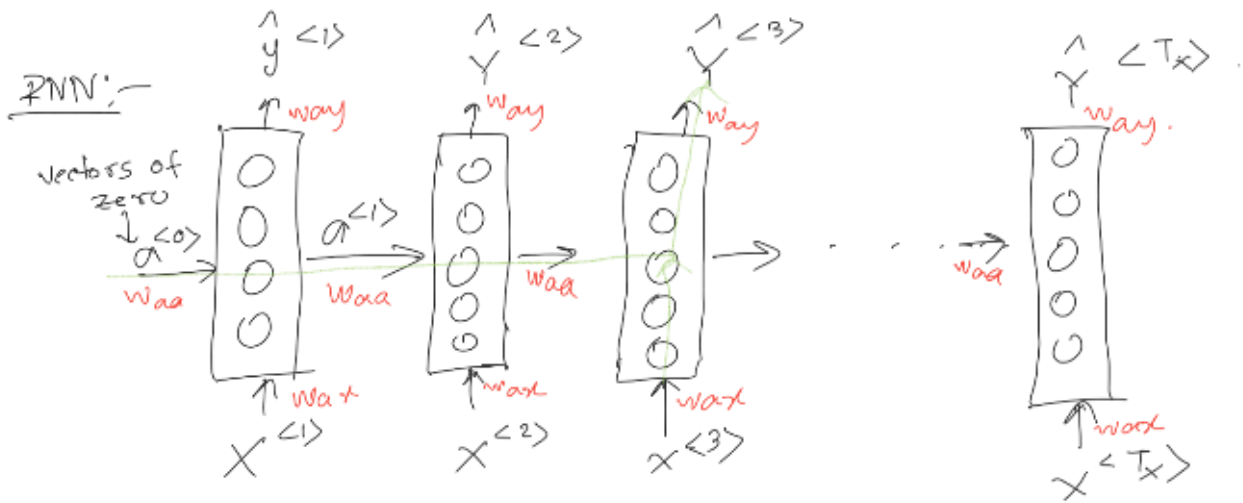
If a word doesn't present in the vocabulary then all those word will contain one position, which will refer to all unknown words.

## RNN (Recurrent neural network):-

Why not standard network:-

- Input, outputs can be different lengths in different example.
- Doesn't share features learned across different position of text.

- Inputs are pretty big.



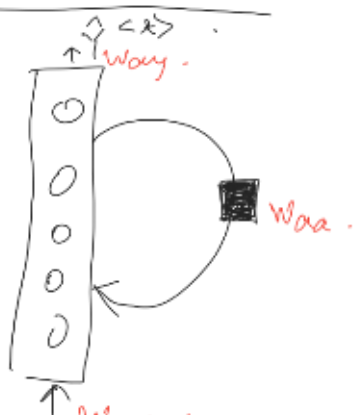
if we want to predict  $\hat{y}^{<3>}$  we use the previous data as well ( $x^{<1>}$ ,  $x^{<2>}$ ) those informations are passed through recurrent network. But we cannot use latter information ( $x^{<4>}$ ,  $x^{<5>}$ ...). And it's a biggest disadvantages of RNN

For example:- we want to extract the name from following two sentences.

- He said, "Teddy roosevelt was a great president"
- He said "Teddy bears are on sale".

RNN have only idea about previous word (He said). Both the teddy in the two sentences will give same result. But those are different. If RNN have the idea about later information it could have guess the word correctly.

Other interpretation:-

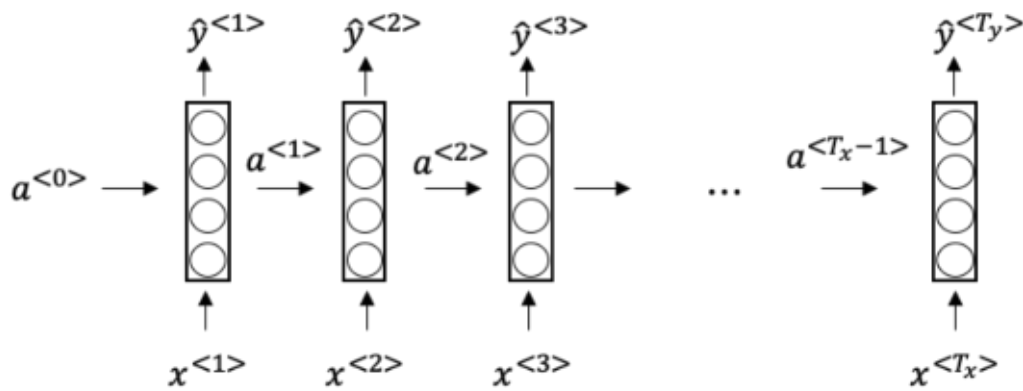


\*Loop denote their recurrent connection

\*shaded box denote the time delay of one steps.

$$x^{<t>} \quad \text{wax}$$

Forward propagation:-



$$a^{<0>} = \vec{0} \quad \left\{ \begin{array}{l} a^{<1>} = g(w_{aa}a^{<0>} + w_{ax}x^{<1>} + b_a) \leftarrow \text{tanh/Relu} \quad \leftarrow \text{more common} \\ \hat{y}^{<1>} = g(w_{ya}a^{<1>} + b_y) \leftarrow \text{sigmoid/softmax} \\ \vdots \\ a^{<t>} = g(w_{aa}a^{<t-1>} + w_{ax}x^{<t>} + b_a) \\ \hat{y}^{<t>} = g(w_{ya}a^{<t>} + b_y) \end{array} \right.$$

Notation:-

$w_{ax} \rightarrow$  it will be multiplied with  $x$  like quantities to compute  $a$  like quantities.

simplified RNN notation:-

$$a^{<t>} = g(\underbrace{w_{aa}a^{<t-1>}}_{(100,100) \times 100} + \underbrace{w_{ax}x^{<t>}}_{(100,10000) \times 10,000} + b_a)$$

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

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

$$= g(w_y a^{<t>} + b_y)$$

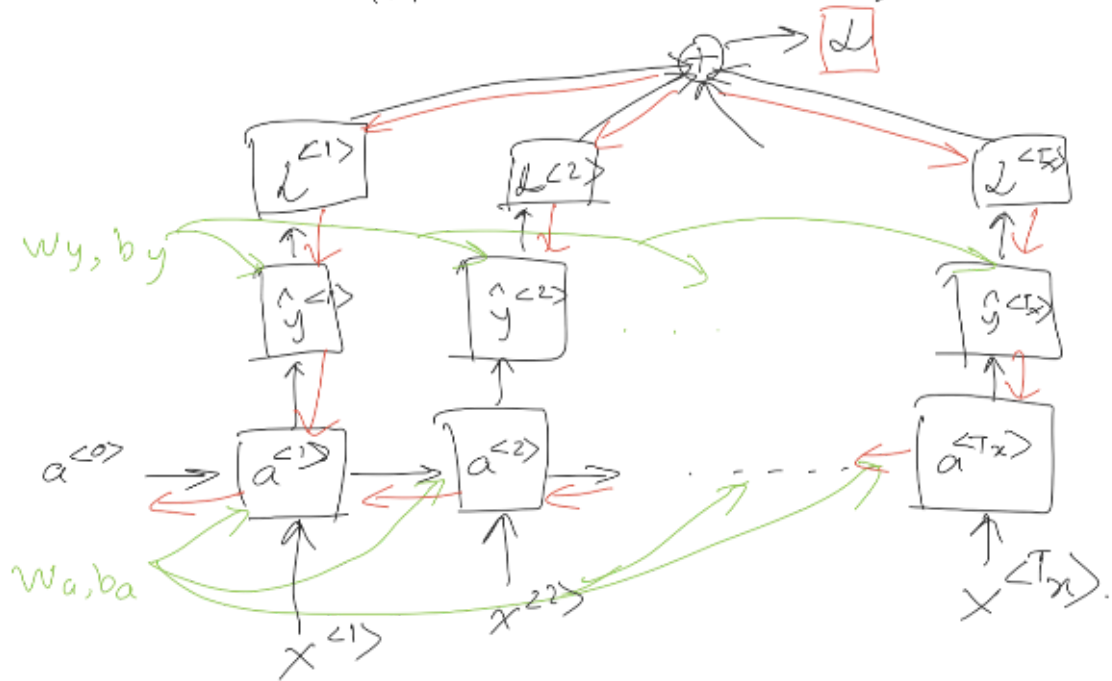
$$\left\{ \begin{array}{l} w_a = \begin{bmatrix} w_{aa} & \vdots & w_{ax} \end{bmatrix} \begin{matrix} \uparrow 100 \\ \leftarrow 100 \quad \leftarrow 10,000 \end{matrix} \\ \begin{bmatrix} a^{<t-1>} & x^{<t>} \end{bmatrix} = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \begin{matrix} \uparrow 100 \\ \leftarrow 10,000 \end{matrix} \\ \begin{bmatrix} w_{aa} & \vdots & w_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = \end{array} \right.$$

$$| w_{aa} a^{<t>} + w_{ax} x^{<t>}$$

## Backpropagation:-

$$\mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log(1-\hat{y}^{<t>})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$



This called backpropagation through time

## Different types of RNN:-

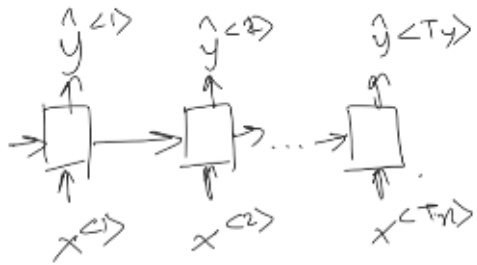
The examples we have seen so far is :-  $T_x = T_y$ .

But there are many cases where  $T_x \neq T_y$ . For example machine translation:- translated word might not have the same number of word as original sentence.

sentiment analysis :- Input can have multiple word but output might have integer value.

| sentiment classification |

$$T_x = T_y$$

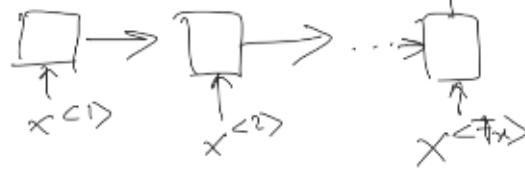


Many-to-many

sequence classification

$x = \text{text}$

$y = 0/1 / 1 \dots 5$



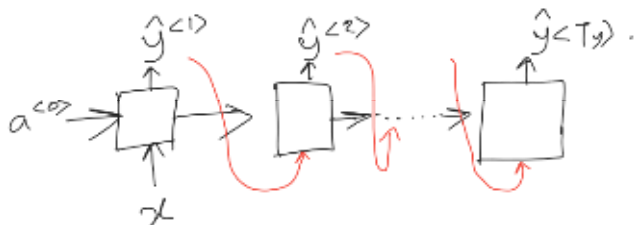
many-to-one



one-to-one

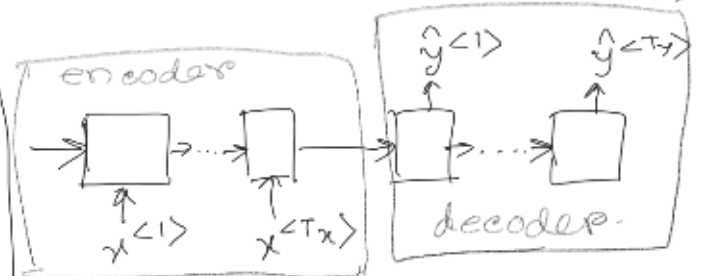
Music generation.

$\rightarrow$  can be music generator  
 $x \rightarrow y^{<1>} y^{<2>} \dots y^{<T_x>}$



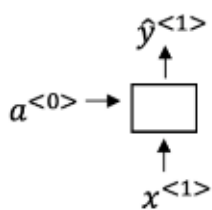
one-to-many

$T_x \neq T_y$  (e.g. Machine translation)

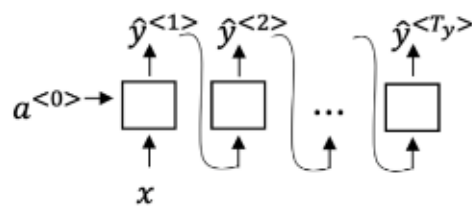


Many-to-many

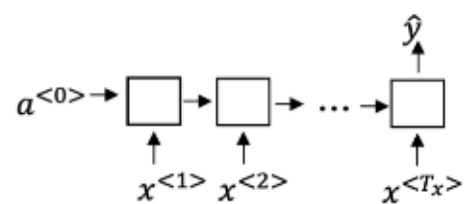
Summary :-



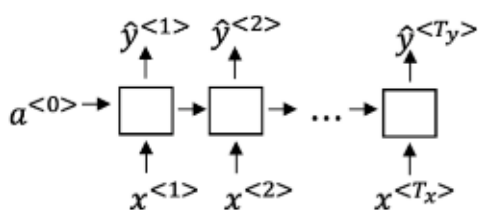
One to one



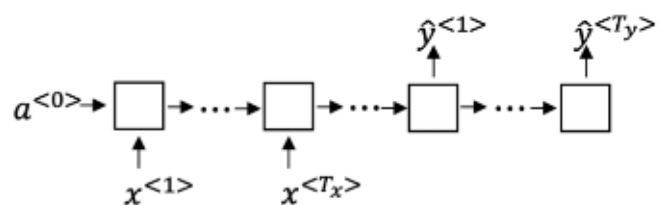
One to many



Many to one



Many to many



Many to many

## Language modeling:-

While we speak some of the word might have the same pronunciation. Language modeling provides the probability of each of the sentence that have similar pronouncing word.

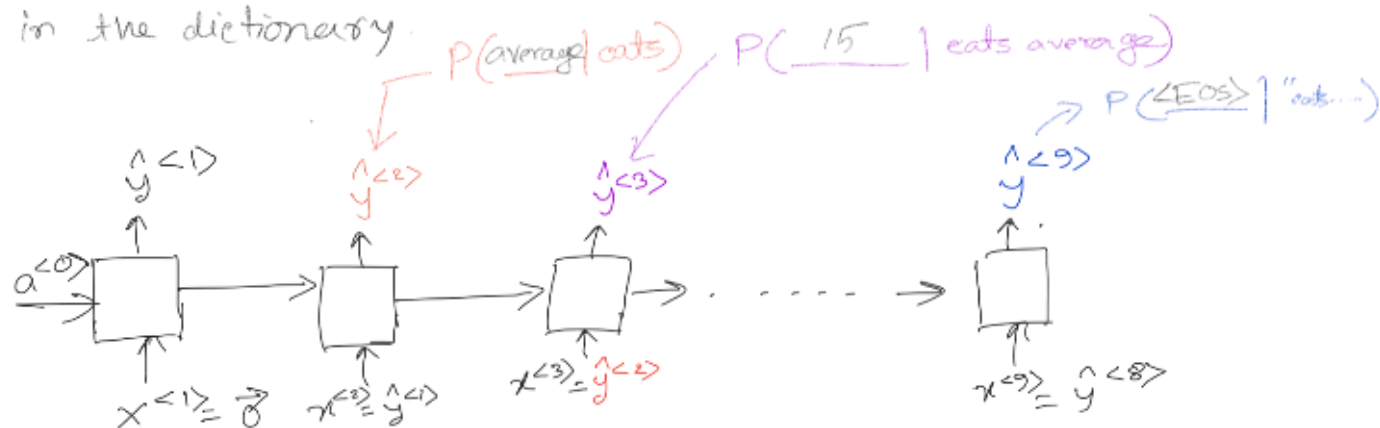
For example:-

$$P(\text{Apple of pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{Apple of pear salad}) = 5.7 \times 10^{-10}$$

## RNN model:-

Cats average 15 hours of sleep a day  $\langle \text{EOS} \rangle$   
softmax, that will predict the probability of each word happening in the dictionary.



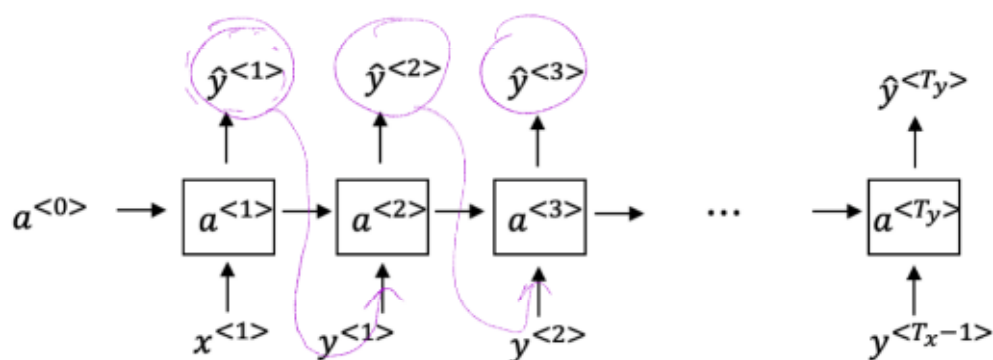
Cost function:-

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Sampling a sequence from a trained RNN:-

this will provide a softmax output. Instead of using the single word from softmax. We will use a random word based on the probability that softmax provided.



### Character level language model:-

Instead of using word in the vocabulary, we will use characters in the vocabulary. So each RNN block will compute a single character instead of the a word.

### Vanishing gradient with RNN:-

The RNN that we have seen so far are not super good to capture long term dependency in a sentence. For example.

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

The **cats**, which already ate . . . . ., **were** full.

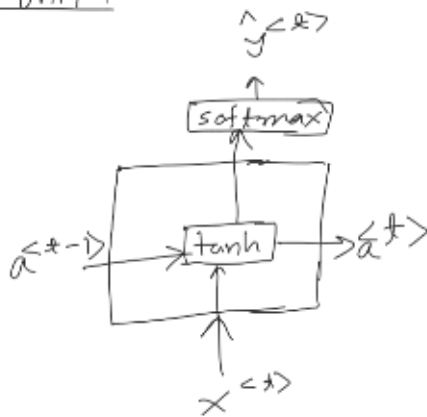
Exploding gradient is easier to solve. If the gradient surpasses a threshold value then we can normalize the gradients this is also called gradient clipping



## Gate Recurrent Unit (GRU) :-

- Capture long range connection.
- Mitigate the vanishing gradient problem.

### RNN unit:-



$$a^{<t>} = \tanh(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

### GRU simplified

$c$  = memory cell.

$$\Rightarrow c^{<t>} = a^{<t>}$$

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

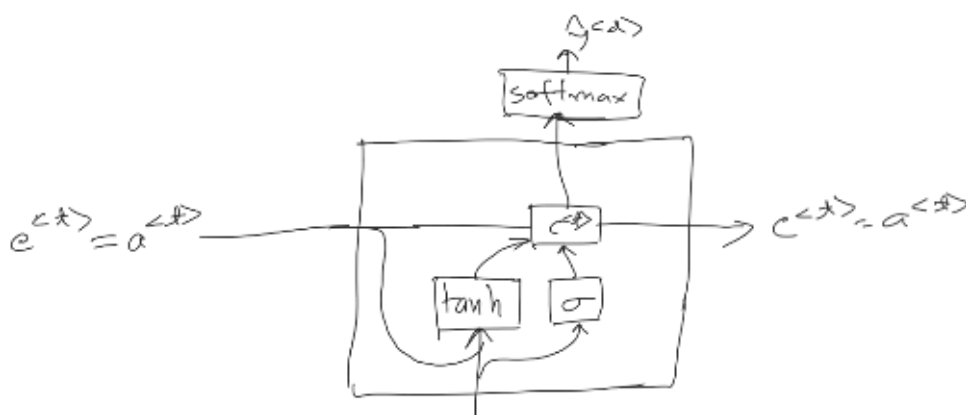
$$\Rightarrow \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

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

update the  $c^{<t>}$  if  $\Gamma_u$  becomes 1. Otherwise stick with the previous one.

$\Gamma_u = 1$     $\Gamma_u = 0$     $\Gamma_u = 0$     $\Gamma_u = 0$    ...    $c^{<t>}$  remembers if the cat was singular or not.  
 $c^{<t>} = 1$     $c^{<t>} = 1$     $c^{<t>} = 1$    ...    $= 1$  if  $c^{<t>} = 1$  it's singular.

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



$x^{<t>}$

Full GRU:-

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

How relevant  $c^{<t-1>}$   
to compute next  
candidate  $c^{<t>}$ .

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

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

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

Long-short term memory:-

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

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

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

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

$$a^{<t>} = c^{<t>}$$

LSTM

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

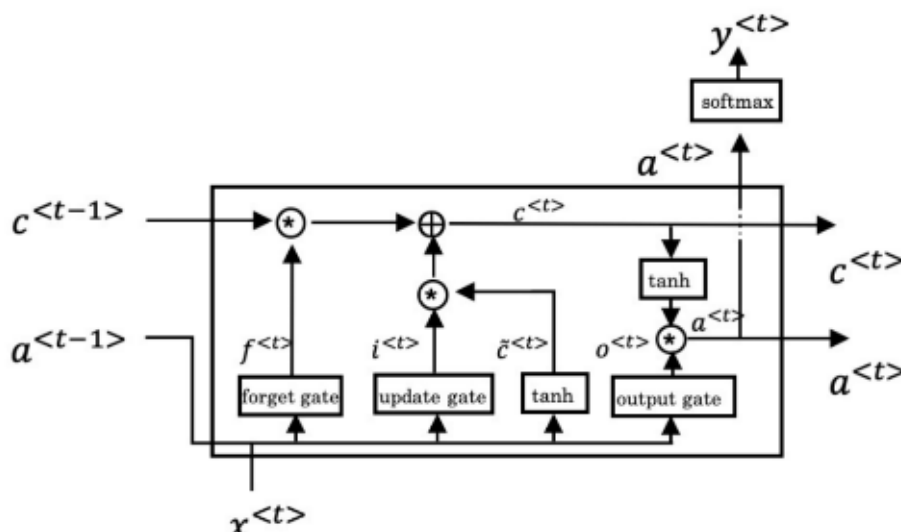
$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

forget  $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

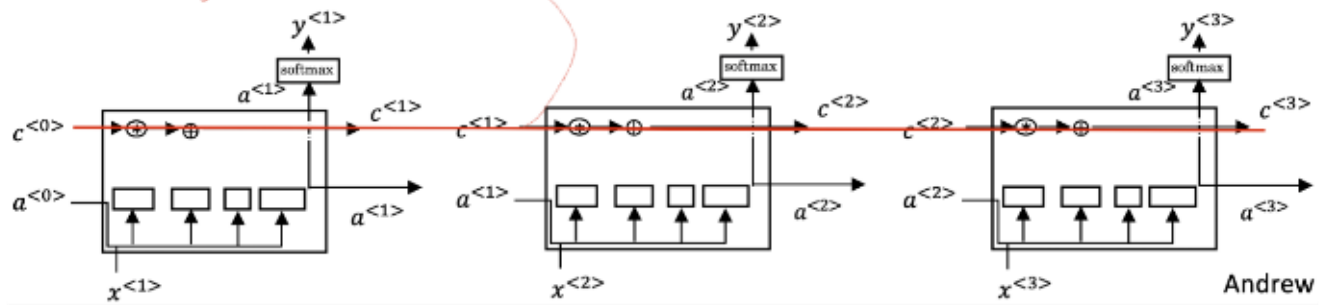
output  $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$

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

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



This line enables LSTM to remember certain value even for many many time steps.



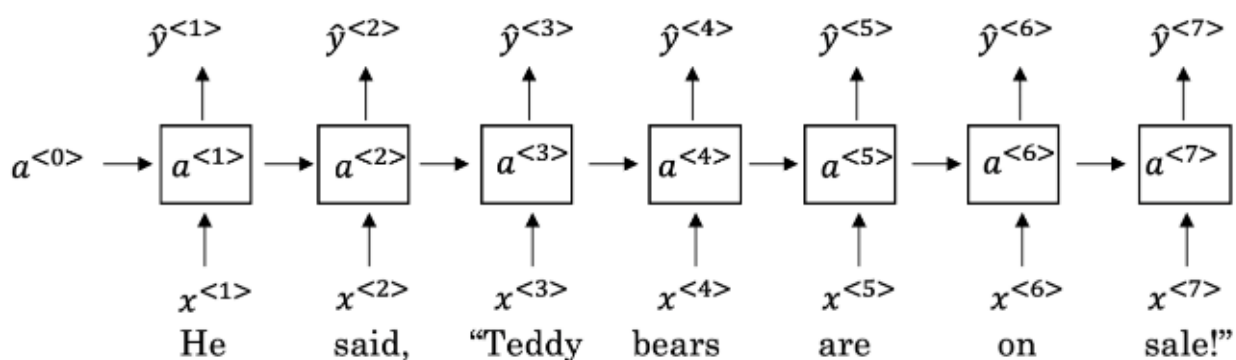
## Bidirectional RNN (BRNN):-

Problem with unidirectional:

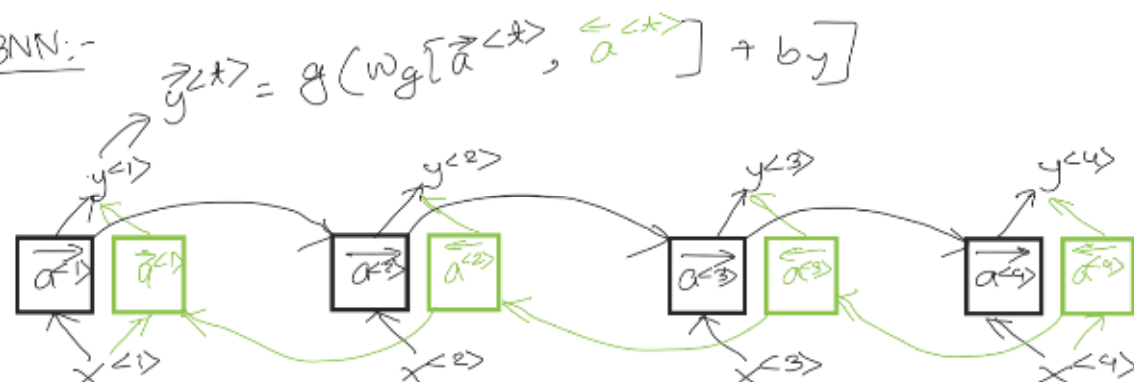
This is a name entity problem, that extract name from sentence. Unidirectional RNN doesn't have enough information to classify teddy as name.

He said, "Teddy bears are on sale!"

He said, "Teddy Roosevelt was a great President!"



BNN:-

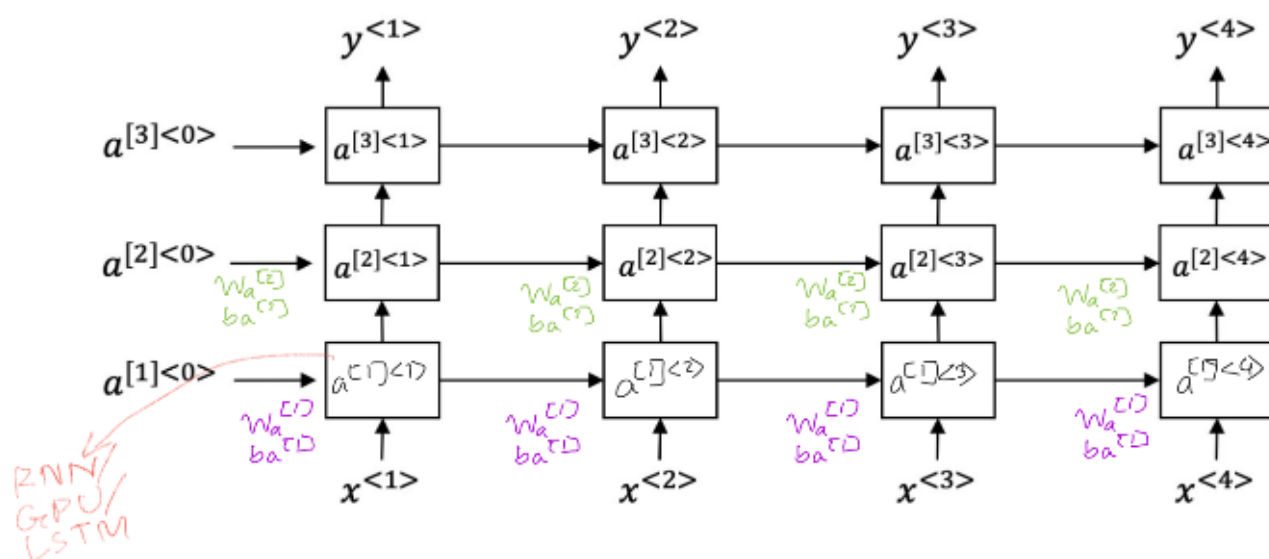


→ If we build a speech recognition system using BRNN then we have to wait until a person stops talking.

## Deep RNN:-

$$a^{[l]} < t >$$

$l = \text{layers}$   
 $t = \text{timestamp}$



$$a^{[2]} < t > = g(w_a^{[2]} [a^{[1]} < t-1 >, a^{[1]} < t >] + b_a^{[2]})$$

every layer will have their own weights. which will be used in each of the timestamp.