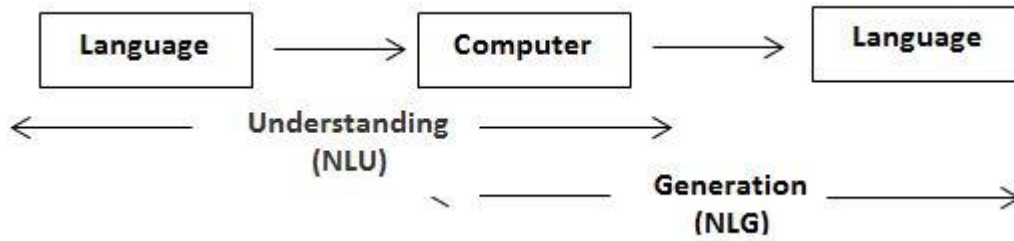
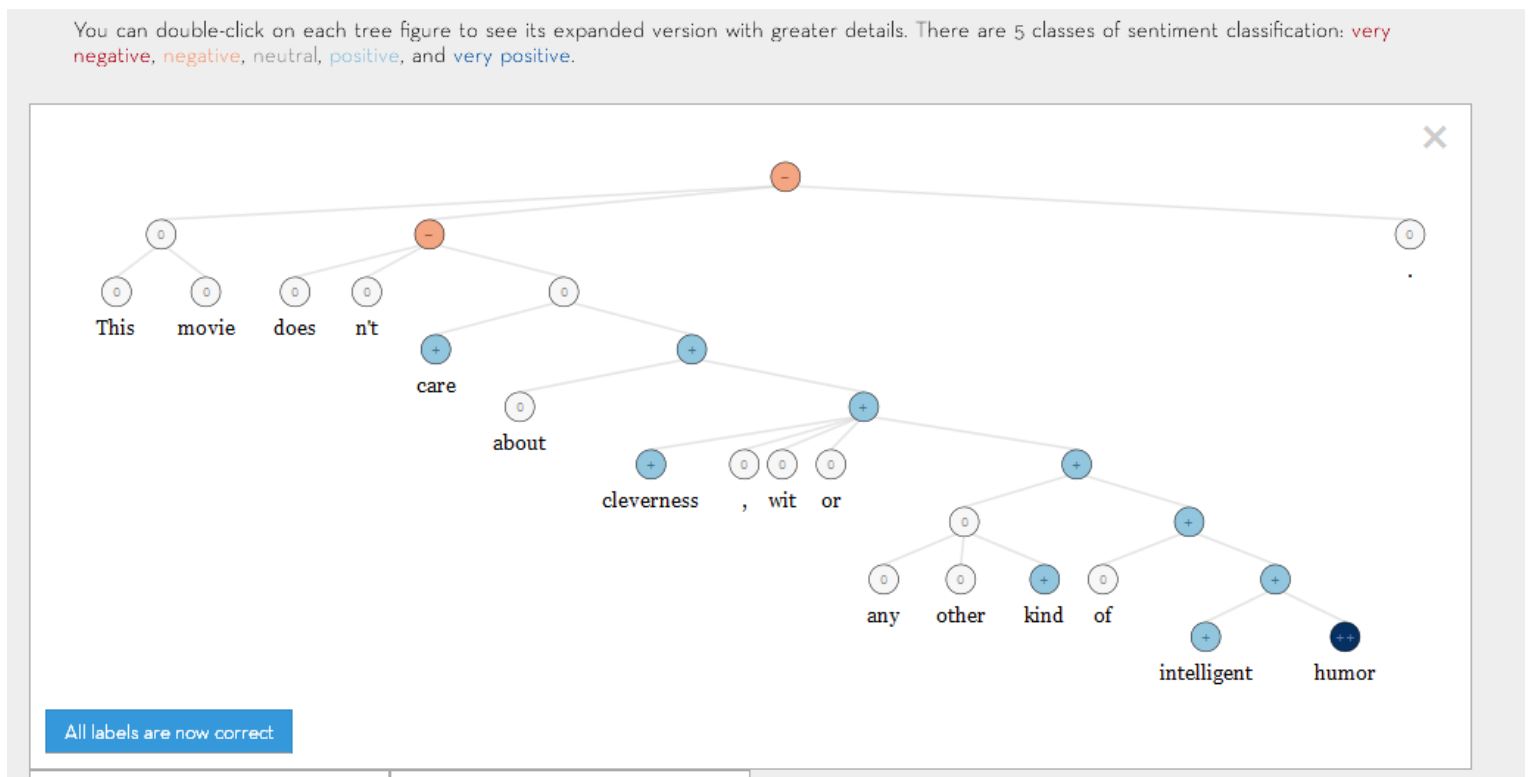


## Overview



- NLP usually includes 2 parts: NLU (Understanding) and NLG (Generation)
- Application
  - Spell Checking
  - Spam Detection
  - Part-of-speech tagging: Colorless(Adj) green(Adj) ideas(Noun) sleep(Verb) furiously(Adv)
  - Name Entity Recognition: Einstein(Person) met with UN(Organization) officials in Princeton(Location)
  - Sentiment Analysis
  - Coreference Resolution: Carter told Peter he shouldn't run again → he here is Carter or Peter
  - Word Sense Disambiguation: I need new batteries for my mouse
  - Parsing
  - Machine Translation: Need to process both NLU and NLG
  - Information Extraction: You're invited to our dinner party, Friday May 27 at 5:30 → Party May 27 add
  - Paraphrase, Summarization
  - Question and Answering
- Sentiment analysis



- Use RNN to encode input sentence and decode vector to output sentence?

## Frequency based Embedding

```
array([[0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
       [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]], dtype=int64)
```

- We have N unique tokens in D documents (D1, D2, ..., Dn). We have D x N matrix, each row in the matrix contains the frequency of tokens in document Di. For example,

D1: He is a lazy boy. She is also lazy

D2: Dick is a lazy person

The 2 x 6 matrix:

	He	She	lazy	boy	Dick	person
D1	1	1	2	1	0	0
D2	0	0	1	0	1	1

- In practice, the matrix can be very sparse and inefficient for computation

→ Vector for 'he': (1, 0)

### TFIDF for 3 documents

- Common words like 'is', 'the', 'a' etc. tend to appear quite frequently in comparison to the words which are important to a document
- For example, the document A about Messi must have lots of word Messi, but the common words like 'the' is also going to present in higher frequency
- TFIDF will penalize the common words by assigning them lower weights while giving more importance to words like Messi in a particular document

Document 1: The game of life is a game of everlasting learning

Document 2: The unexamined life is not worth living

Document 3: Never stop the learning

- Step 1: Term Frequency: measures the number of times word occurs in a document

#### TF for Document 1

Document1	the	game	of	life	is	a	everlasting	learning
Term Frequency	1	2	2	1	1	1	1	1

#### TF for Document 2

Document2	the	unexamined	life	is	not	worth	living
Term Frequency	1	1	1	1	1	1	1

#### TF for Document 3

Document3	never	stop	learning	the
Term Frequency	1	1	1	1

○ **Normalized TF for Document 1**

Document1	the	game	of	life	is	a	everlasting	learning
Normalized TF	$\frac{1}{10}$	$\frac{2}{10}$	0.2	0.1	0.1	0.1	0.1	0.1

○ **Normalized TF for Document 2**

Document2	the	unexamined	life	is	not	worth	living
Normalized TF	$\frac{1}{7}$	0.142857	0.142857	0.142857	0.142857	0.142857	0.142857

○ **Normalized TF for Document 3**

Document3	never	stop	learning	the
Normalized TF	$\frac{1}{4}$	0.25	0.25	0.25

- Step 2: Inverse Document Frequency:
  - **IDF(the) =  $\log_e(\text{Total Number Of Documents} / \text{Number Of Documents with term game in it}) = \log_e(3 / 3) = 0$**
  - If a word has appeared in all documents, then probably that word is not relevant to a particular document, but if it has appeared in a subset of documents then probably the word may be relevant to the documents it is present in
- Step 3: TF \* IDF

→ vector for 'the' = (0, 0, 0)

## Co-occurrence Matrix

- Idea: Similar words tend to occur together and will have similar context. E.g. Apple is a fruit. Mango is a fruit → Apple and mango tend to have a similar context
- Here we use Context Windows of 2 and the corpus: He is not lazy. He is intelligent. He is smart

	He	is	not	lazy	intelligent	smart
He	0	4	2	1	2	1
is	4	0	1	2	2	1
not	2	1	0	1	0	0
lazy	1	2	1	0	0	0
intelligent	2	2	0	0	0	0
smart	1	1	0	0	0	0

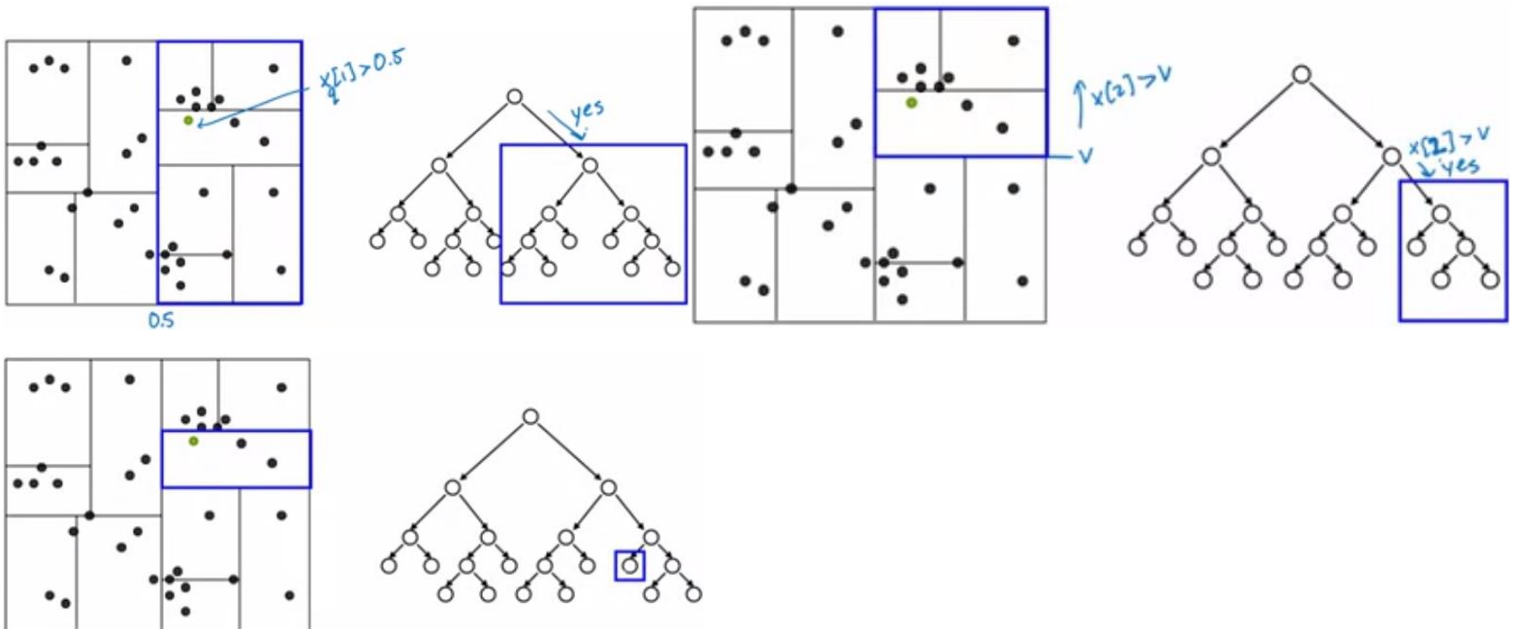
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart

- While the word 'lazy' has never appeared with 'intelligent' in the context window, therefore assigned 0 in the blue box

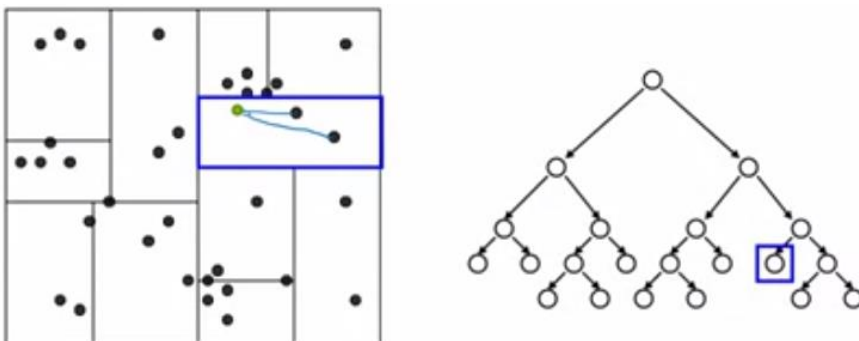
- If we have  $N$  unique tokens, vocabulary size =  $N$ , the co-occurrence matrix is  $N \times N$ , so it is very large. To partially solve it, we can construct the  $M \times N$  co-occurrence matrix with  $M$  is the subset of  $N$  by removing irrelevant words like stop words, but it is still very large and difficult to compute
- However, **this co-occurrence matrix is not the word vector representation**
- But, if it performs PCA on  $N \times N$  matrix, you can choose  $k$  components out of these  $N$  components so the matrix will be  $N \times k$  matrix
- PCA will decompose co-occurrence matrix into 3 matrices =  $U \cdot S \cdot V^T$ ,  $U \cdot S$  is the word vector representation
- Advantage
  - It preserves the semantic relationship between words. i.e man and woman tend to be closer than man and apple
  - It uses SVD at its core, which produces more accurate word vector representations than existing methods.

### KD Tree for Nearest Neighbor Search

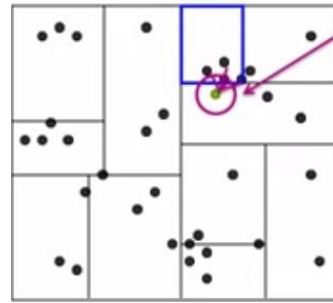
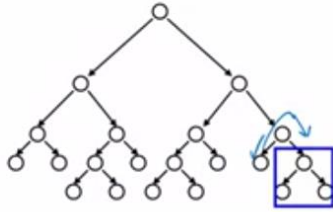
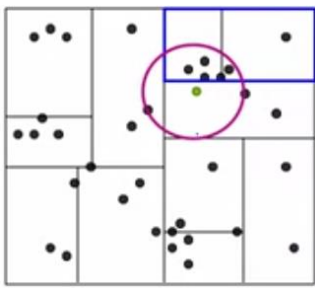
- For the dataset, you partition based on the  $x$  and  $y$  axis. Based on your partitioning, you draw the binary tree to keep track
- Step 1: Start by exploring leaf node containing query point



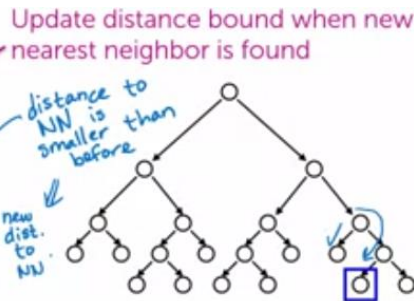
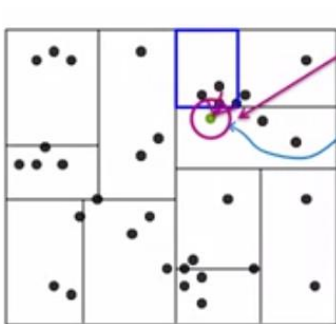
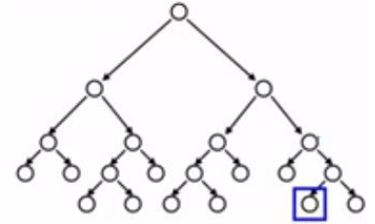
- Step 2; Compute distance to each other point at leaf node



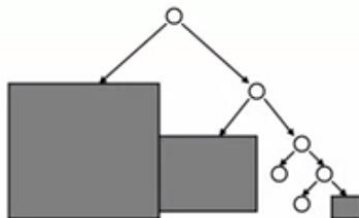
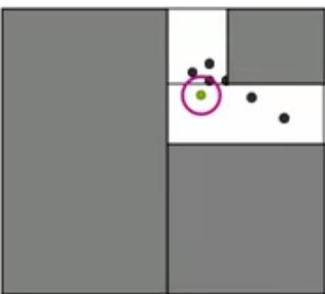
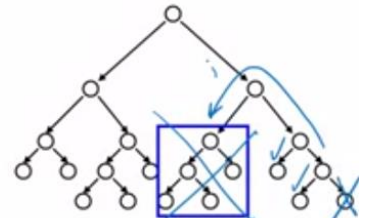
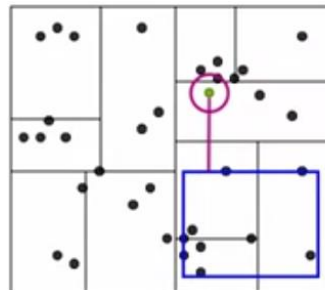
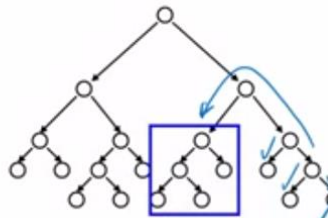
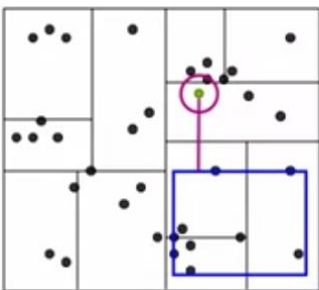
- But in this case, update threshold the distance between query point and points in your considering leaf nodes. Find another leaf node
- Step 3: Backtrack and try other branch at each node visited



Update distance bound when new nearest neighbor is found



- Step 4: Use your updated threshold to prune the tree



## Word2Vec/ Prediction based Vector

- **Transform single word to vector** is very important in NLP
- Word2vec helps you to do it, it includes 2 models:
  - Skip-gram
  - Continuous Bag-of-words
- Both of these are shallow neural network which map words to the input. They learn weights which act as **word vector representation**

### Continuous Bag of words (CBOW)

- The way CBOW work is that it tends to predict the probability of a word given a context (n words before and after). Consider example with context size of 1, we have corpus C = 'Hey, this is sample corpus using only one context word'
- The corpus may be converted into



- For example, Apple can be both a fruit and a company but CBOW takes an average of both the contexts and places it in between a cluster for fruits and companies

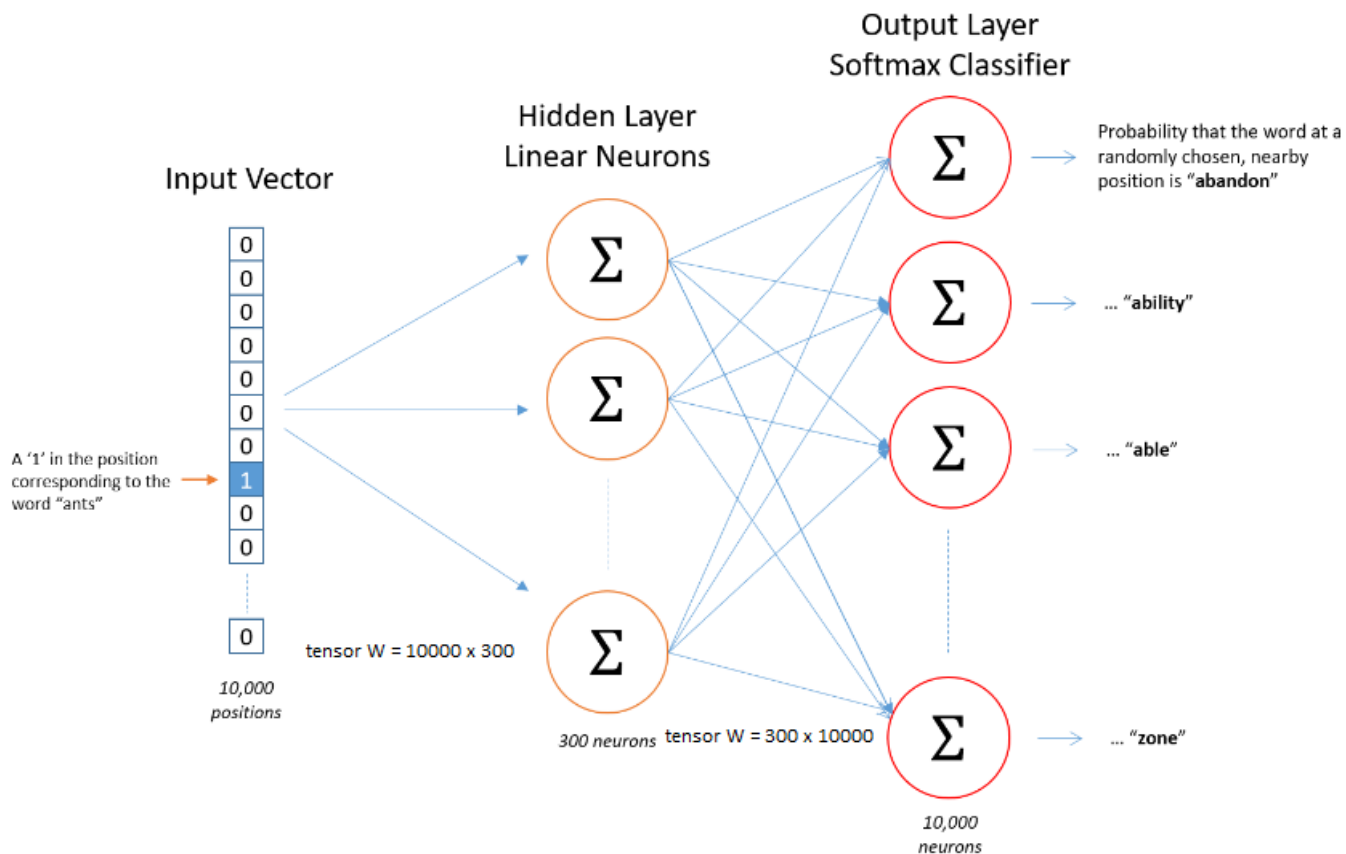
### Skip-gram

- Given the word, predict the n words before and n words after based on the probability
  - For example, if you gave the trained network the input word “Soviet”, the output probabilities are going to be much higher for words like “Union” and “Russia” than for unrelated words like “watermelon” and “kangaroo”.
  - The weight between the input and hidden layer is **taken as the word vector representation of the word**
- First of all, you cannot feed the raw text string to neural network. For example, you have corpus (list of sentence):
  - The ants eat my cake
  - You are my boy friend

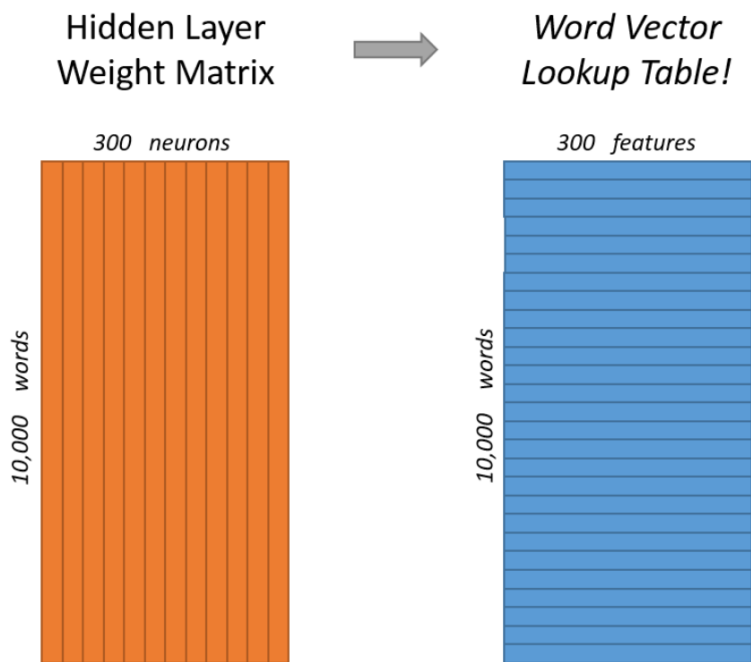
	The	ants	eat	my	cake	You	are	boy	friend
The	1	0	0	0	0	0	0	0	0
ants	0	1	0	0	0	0	0	0	0

- Create one-hot for ants:  $[0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$

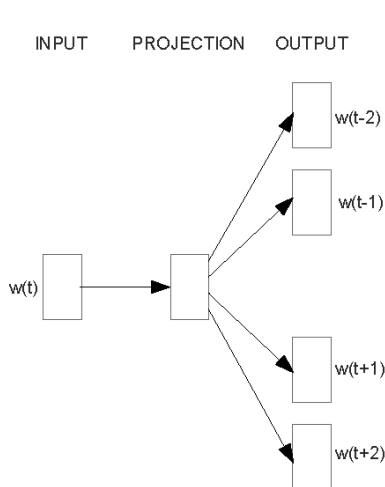
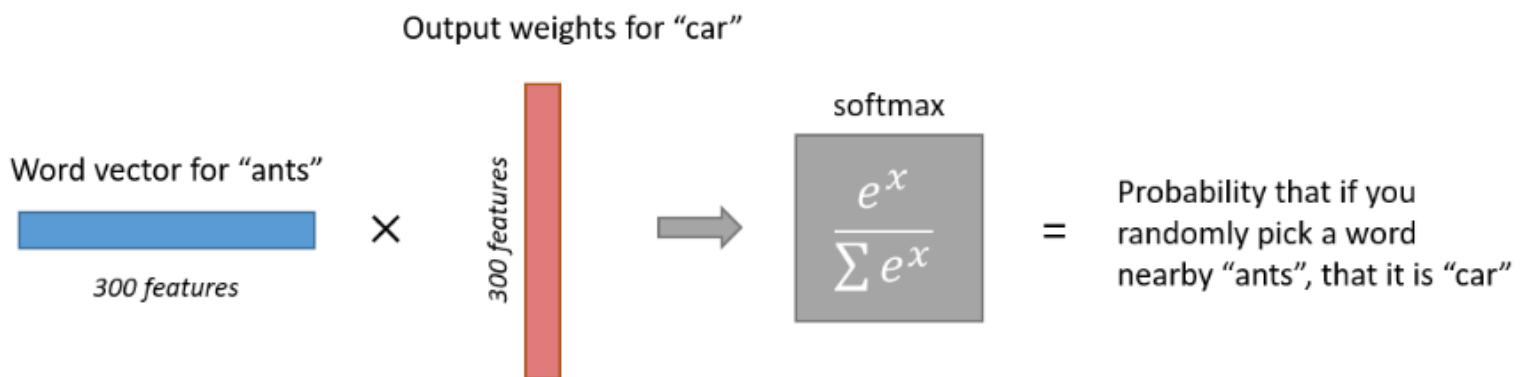
- In our example, we have 1 rows x 10000 components/ columns (one for every word in our vocabulary) for ants in input of NN, and output is single vector with 10000 components, the probability that a randomly selected nearby word is that vocabulary word



- There is **no activation function** on the hidden layer neurons, but the output neurons use softmax
- Every word is represented for class and 300 neurons means 300 features that a word include, 300 features is just hyperparameter, we need to tune model to find out



- So the end goal of all of this is really just to **learn this hidden layer weight matrix** – the output layer we’ll just toss when we’re done
- Because our input is one-hot vector which is almost all zeros, if you multiply 1 x 10000 one hot to hidden layer matrix 10000 x 300, the output is just the word vector for the input word
- For the output layer, here’s an illustration of calculating the output of the output neuron for the word “car”



**Skip-gram**

- So, we use 1 hidden layer and 1 softmax as output layer:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(\omega_{t+j} | \omega_t)$$

- T: number of subset 2c+1 in the corpus

- Another example, continue example in CBOW, take ‘this’ as input and ‘Hey’, ‘is’ as target



Context											Input-Hidden Weight				Hidden Activation			
C1	this	0	1	0	0	0	0	0	0	0	1	2	3	4	5	6	7	8
											1	2	3	4				
											5	6	7	8				
											9	10	11	12				
											13	14	15	16				
											17	18	19	20				
											21	22	23	24				
											25	26	27	28				
											29	30	31	32				
											33	34	35	36				
											37	38	39	40				

Hidden Activation											hidden-output weight matrix				output			
Thir	5	6	7	8							0.12	0.13	0.14	0.15	7.52	7.73	8.94	9.3
											0.22	0.23	0.24	0.25	7.52	7.73	8.94	9.3
											0.32	0.33	0.34	0.35				
											0.42	0.43	0.44	0.45				
											0.17	0.18	0.19	0.2				
											0.27	0.28	0.29	0.3				
											0.37	0.38	0.39	0.4				
											0.47	0.48	0.49	0.5				
											0.57	0.58	0.59	0.6				

refined probabilities											Target				Error				
											0.024	0.03019744	0.04097102	0.05197034	-0.98	-0.98	-0.98	-0.98	
											0.024	0.03019744	0.04097102	0.05197034	0.024	0.03019744	0.04097102	0.05197034	
											0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	
											0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	
											0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	
											0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	
											0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	
											0.024	0.03019744	0.04097102	0.05197034	-0.98	-0.98	-0.98	-0.98	
											0.024	0.03019744	0.04097102	0.05197034	0.024	0.03019744	0.04097102	0.05197034	
											0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	0.0674019	
											0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	0.11337104	
											0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	0.14703528	
											0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	0.19069461	
											0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	0.24731759	
											Sum of error								
											-0.98	0.06179417	-0.9191854	0.10394049	0.1348031	0.17403111	0.22574372	0.29407076	0.35131922

- The red row in second picture is hidden layer. You multiply it to hidden-output weight matrix, you will get 1 row, but because you have 2 targets, so you must duplicate this row (blue matrix)
- Error is calculated by subtracting the grey matrix(target) to green matrix(output) element-wise
- Advantage:
  - Skip-gram model can capture two semantics for a single word. i.e it will have two vector representations of Apple. One for the company and other for the fruit

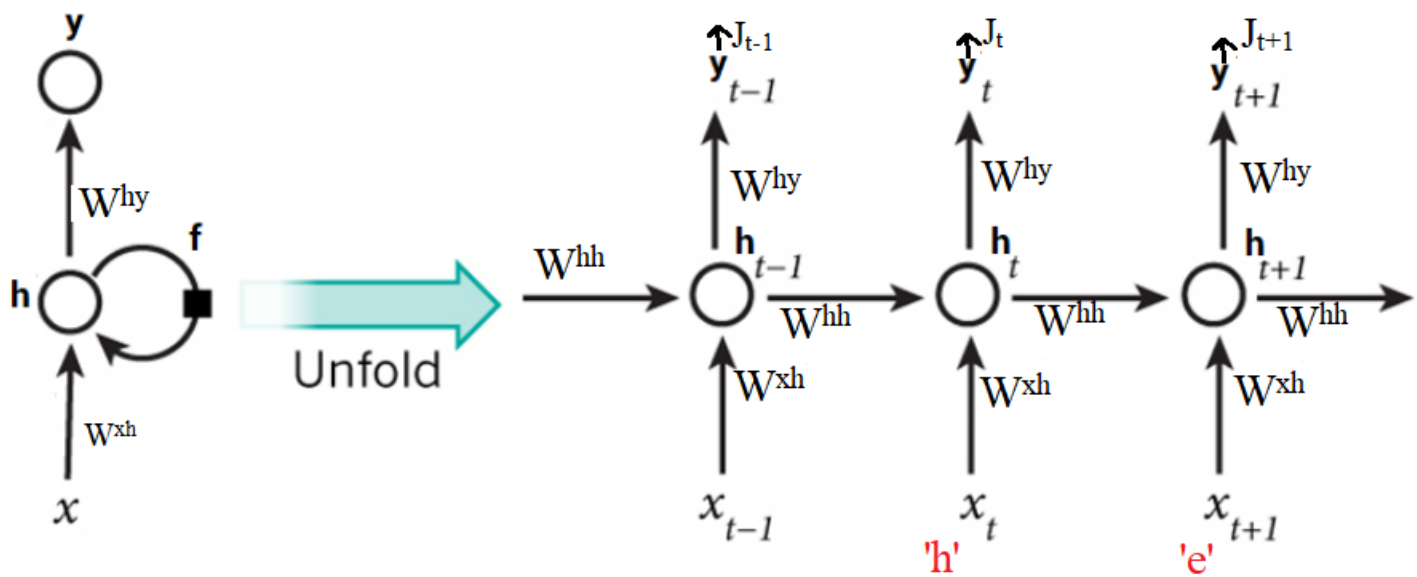
## Phrase Representation

$$score(\omega_i, \omega_j) = \frac{count(\omega_i \omega_j) - \delta}{count(\omega_i) count(\omega_j)}$$

## Recurrent Neural Network

### Language Model

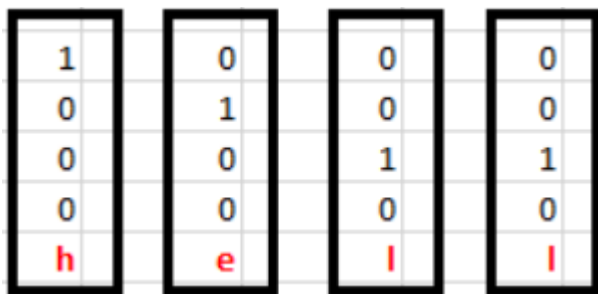
- Language Model calculates the probability of the string in the corpus:  $p(\omega_1, \dots, \omega_T)$
- Given i-1 words, predict the word #i, conditional probabilities:  $p(\omega_i | \omega_1, \dots, \omega_{i-1}) = \prod_{i=1}^m p(\omega_i | \omega_1, \dots, \omega_{i-1})$
- Markov Assumption helps to calculate  $p(\omega_2 | \omega_1) = \frac{count(\omega_1, \omega_2)}{count(\omega_1)}$  or  $p(\omega_3 | \omega_1, \omega_2) = \frac{count(\omega_1, \omega_2, \omega_3)}{count(\omega_1, \omega_2)}$
- Because Language Model need to calculate  $count(\omega_1, \omega_2), count(\omega_1, \omega_2, \omega_3)$  so the complex will be explode exponentially
- Markov Order: the number of word we consider to predict the target word
  - Markov-order = 1,  $p(\omega_2 | \omega_1)$
  - Markov-order = 2,  $p(\omega_3 | \omega_1, \omega_2)$



- Given a list of word vectors:  $x_1, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_T$
- At a single step

- $h_t = \sigma(W^{hh}h_{t-1} + W^{xh}x_t)$
- $\hat{y}_t = \text{soft max}(W^{hy}h_t)$
- $\hat{P}(x_{t+1} = v_j | x_t, \dots, x_1) = \hat{y}_{t,j}$

- Feed Forward in RNN: Instead of words, we just consider characters, given {h, e, l, l}, predict the next character and the target is {o}
- ddd  $h'_{t+1}$



- First of all, we need to determine the dimension of  $h_{t-1}, h_t, h_{t+1}$ , here we choose  $\mathbb{R}^{3 \times 1}$ , so  $h_{t-1} = [0 \ 0 \ 0]^T$ . From the input, we can see the dimension of  $x_t \in \mathbb{R}^{4 \times 1}$
- To fulfil the first formula  $h_t = \sigma(W^{hh}h_{t-1} + W^{xh}x_t) = \sigma(W[h_{t-1}, x_t] + b)$ , ( $\sigma$ : tanh function)  $W^{hh}$  could be followed 2 approaches:

- $W^{hh}$  includes 3x3 matrix:  $\begin{bmatrix} 0.427 & 0 & 0 \\ 0 & 0.427 & 0 \\ 0 & 0 & 0.427 \end{bmatrix}$  and 3x1 matrix bias:  $\begin{bmatrix} 0.567 \\ 0.567 \\ 0.567 \end{bmatrix}$ , so

$$W_{hh}h_{t-1} = \begin{bmatrix} 0.427 & 0 & 0 \\ 0 & 0.427 & 0 \\ 0 & 0 & 0.427 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.567 \\ 0.567 \\ 0.567 \end{bmatrix}$$

- $W^{hh} \in \mathbb{R}^{3 \times 4} = \begin{bmatrix} 0.427 & 0 & 0 & 0.567 \\ 0 & 0.427 & 0 & 0.567 \\ 0 & 0 & 0.427 & 0.567 \end{bmatrix}$  and add one to  $h_{t-1} = [0 \ 0 \ 0 \ 1]^T$

- In this case, to make it simpler, we will choose first approach

- Randomize the initial  $W^{xh}$  and multiply to  $x_t$

wxh				X	=	
0.287027	0.84606	0.572392	0.486813			
0.902874	0.871522	0.691079	0.18998			
0.537524	0.09224	0.558159	0.491528			

- $W \in \mathbb{R}^{3 \times 7} = \begin{bmatrix} 0.427 & 0 & 0 & 0.287 & 0.84 & 0.572 & 0.486 \\ 0 & 0.427 & 0 & 0.902 & 0.871 & 0.691 & 0.189 \\ 0 & 0 & 0.427 & 0.537 & 0.092 & 0.558 & 0.491 \end{bmatrix}$

- $W^{hh}h_{t-1} + W^{xh}x_t = \left( \begin{bmatrix} 0.427 & 0 & 0 \\ 0 & 0.427 & 0 \\ 0 & 0 & 0.427 \end{bmatrix} \mathbf{x} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.567 \\ 0.567 \\ 0.567 \end{bmatrix} \right) + \begin{bmatrix} 0.287 \\ 0.902 \\ 0.537 \end{bmatrix} = \begin{bmatrix} 0.85 \\ 1.46 \\ 1.10 \end{bmatrix}$

- $h_t = \sigma(W^{hh}h_{t-1} + W^{xh}x_t) = [0.7 \quad 0.81 \quad 0.75]^T$

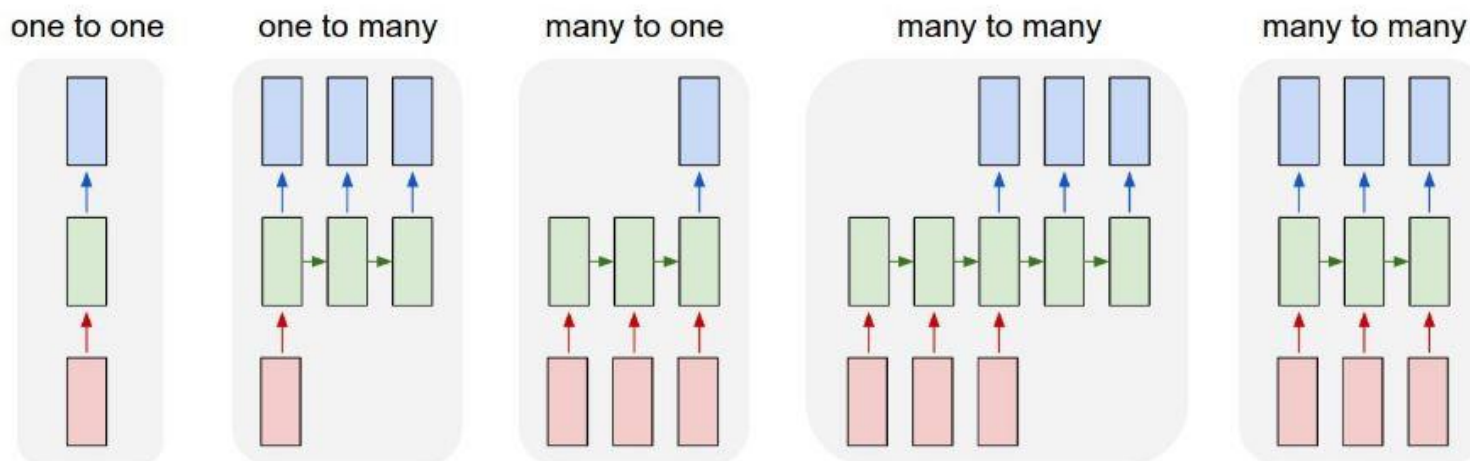
- Randomize initial weight  $W^{hy} = \begin{bmatrix} 0.37 & 0.97 & 0.83 \\ 0.39 & 0.28 & 0.65 \\ 0.64 & 0.09 & 0.33 \\ 0.91 & 0.32 & 0.14 \end{bmatrix}$ ,

$$\hat{y}_t = \text{soft max}(W^{hy} \mathbf{x} h_t) = \text{soft max} \left( \begin{bmatrix} 0.37 & 0.97 & 0.83 \\ 0.39 & 0.28 & 0.65 \\ 0.64 & 0.09 & 0.33 \\ 0.91 & 0.32 & 0.14 \end{bmatrix} \mathbf{x} \begin{bmatrix} 0.7 \\ 0.81 \\ 0.75 \end{bmatrix} \right) = \text{soft max} \left( \begin{bmatrix} 1.69 \\ 0.98 \\ 0.76 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0.41 \\ 0.2 \\ 0.16 \\ 0.21 \end{bmatrix}$$

- At  $t$ ,  $\begin{cases} y_t = [1 \quad 0 \quad 0 \quad 0]^T \rightarrow y_{t,0} = y_{t,h} = 1 \\ \hat{y}_t = [0.41 \quad 0.2 \quad 0.16 \quad 0.21] \rightarrow \hat{y}_{t,0} = \hat{y}_{t,h} = 0.41 \end{cases}$

$$J^{(t)} = -\sum_{j=1}^{|V|} y_{t,j} \log(\hat{y}_{t,j}) \rightarrow J = -\sum_{t=1}^T \sum_{j=1}^{|V|} y_{t,j} \log(\hat{y}_{t,j})$$

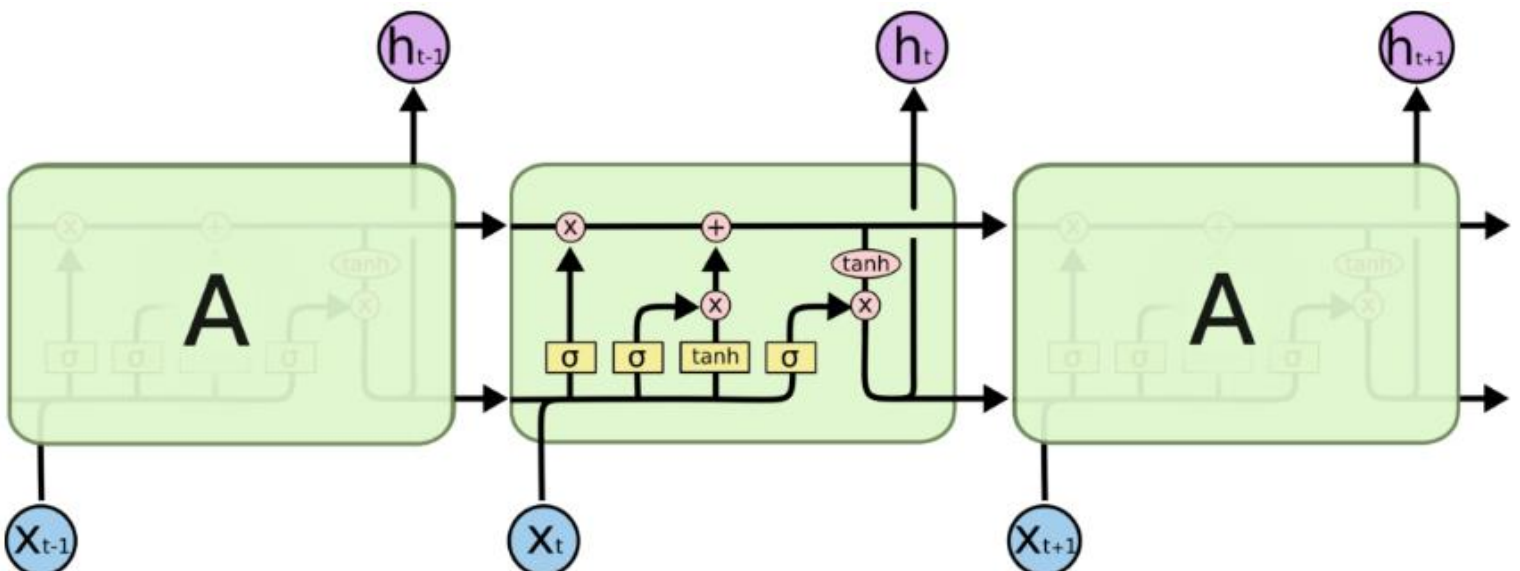
#### • Architecture



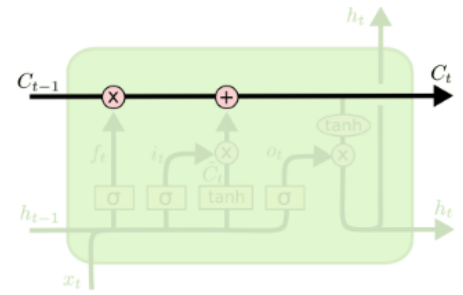
- The example above is many to many

### Long Term dependencies and Short Term attention

- Each lines carries an entire vector
- This figure is for standard RNN, below is LSTM

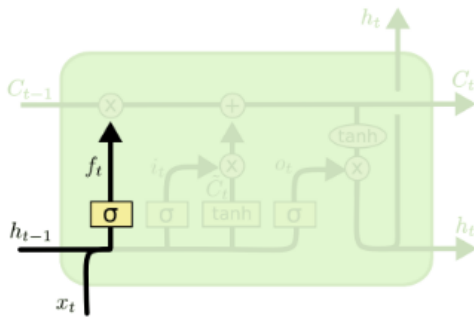


- LSTM also have the chain structure like standard RNN, but instead of single neural layer, they have four, interacting each other
- Another main difference is the **horizontal line through the top of the diagram (HL)** or the **cell state** in original diagram, so the LSTM can remove or add information from other parts of corpus
- Sigmoid layer describes how much the words should be added or removed
  - 0 means “completely get rid of this.”
  - 1 means “completely keep this”



### Step-by-Step LSTM

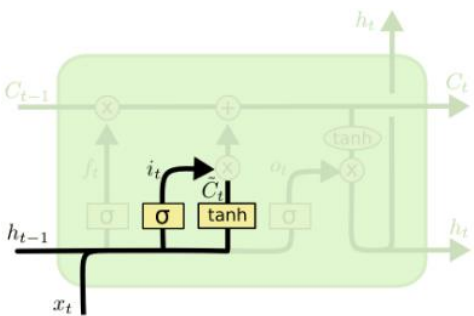
#### Forget gate



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

- In the first step, determine how much information in  $C_{t-1}$  should be forgot
- Example of a language model trying to predict the next word based on all the previous ones
  - When we see the new subject, we want to remove the gender of the old subject

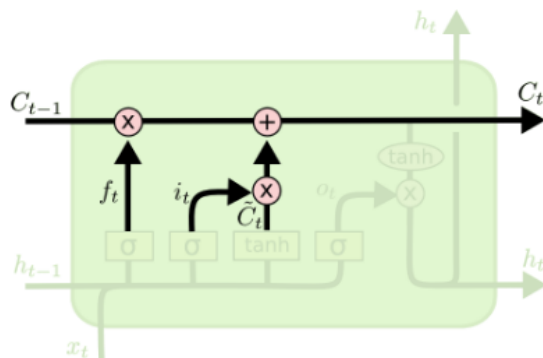
#### Input gate



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

- Next step, determine the new information should be saved in cell state, E.g. add the gender of the new subject to the hidden state, to replace the old one we're forgetting

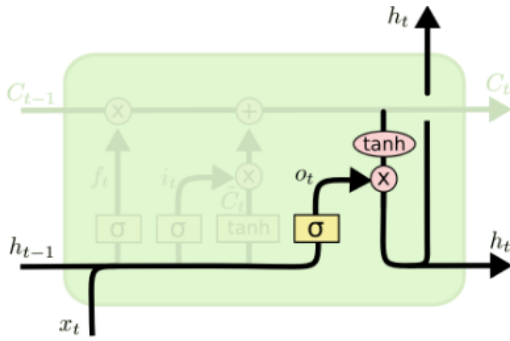


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- Update the old hidden state  $C_{t-1}$ , into the new cell state  $C_t$

# Output gate

- Determine  $h_t$  based on cell state  $C$



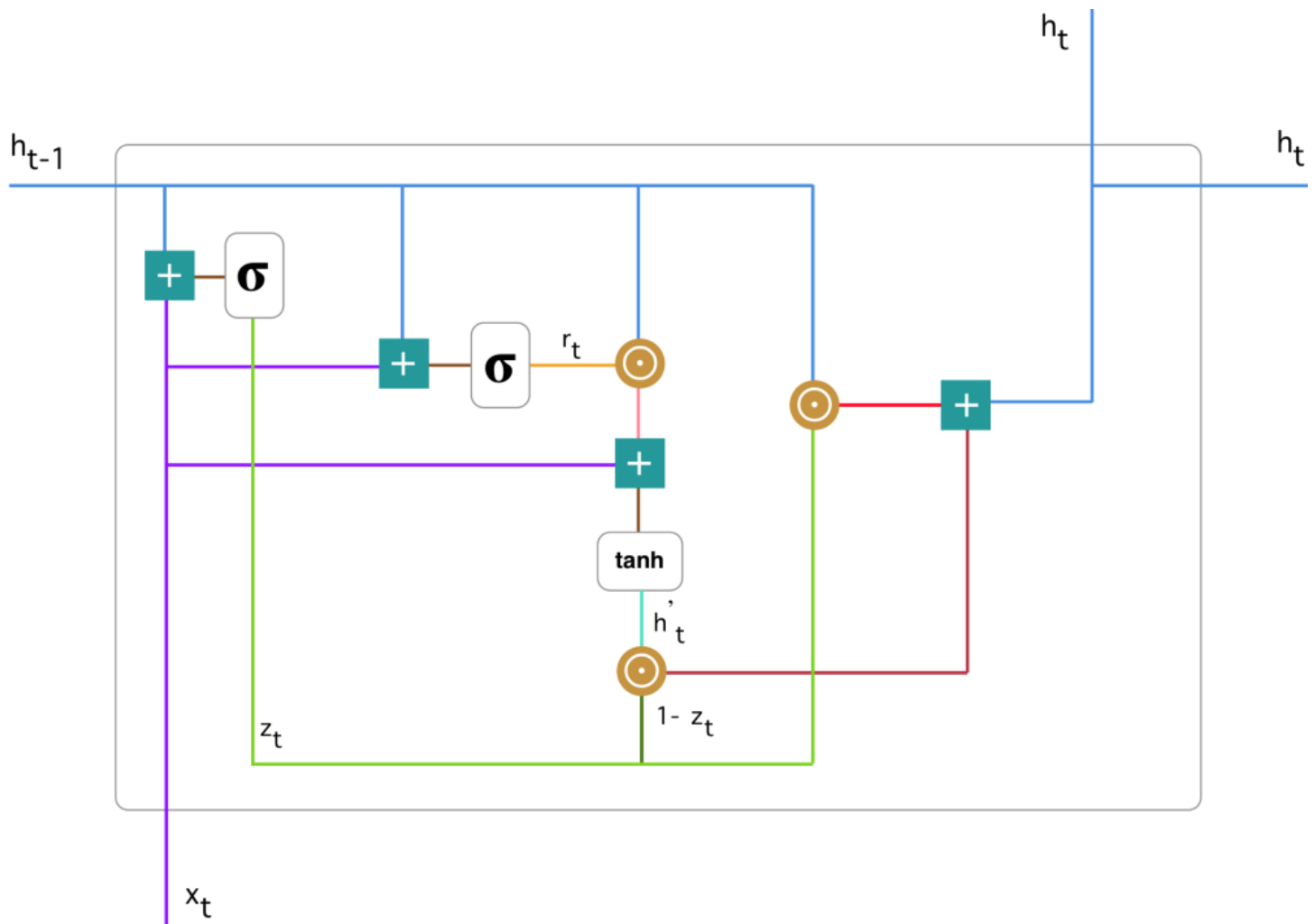
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

- Run a sigmoid layer which decides what parts of the cell state we're going to output
- Then, we put the cell state through  $\tanh$  (to push the values to be between  $-1$  and  $1$ ) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.
- For the language model example, since it just saw a subject, it might want to output information relevant to a verb, in case that's what is coming next
- For example, it might output whether the subject is singular or plural, so that we know what form a verb should be conjugated into if that's what follows next.

## GRU

- Instead of using Forget, Input and output gate, GRU uses Update and Reset gate. Basically, there are 2 gates which decide what information should be passed to the output
- They can be trained to keep relevant information from long ago without vanishing or remove though time





“plus” operation



“sigmoid” function



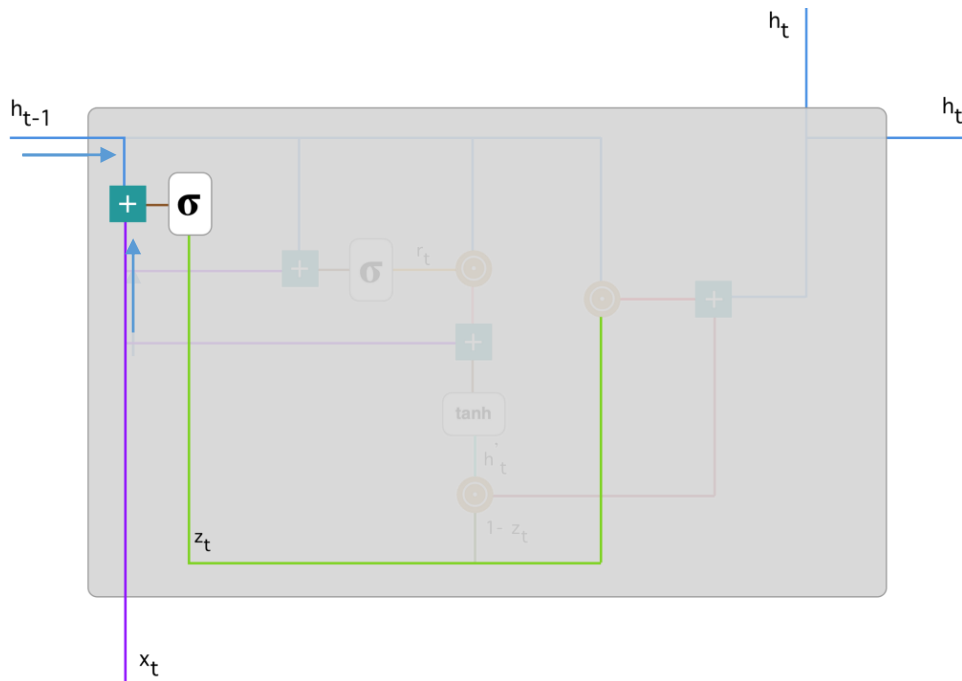
“Hadamard product” operation



“tanh” function

How the GRU works

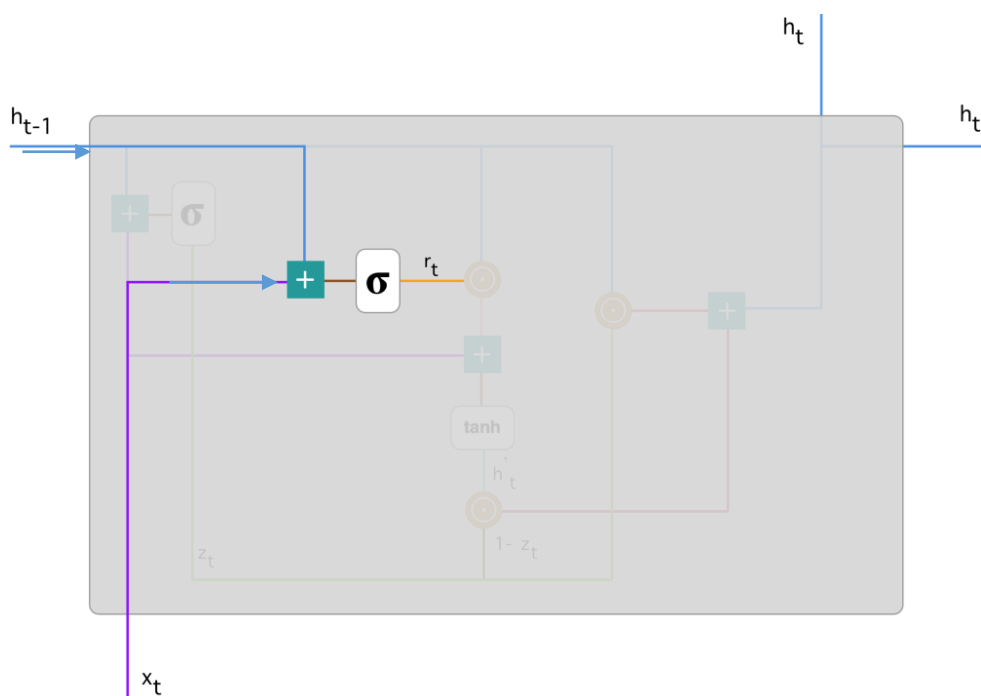
- Update gate



$$z_t = \sigma(W^z x_t + U^z h_{t-1})$$

- Update gate helps the model to determine how much of the past information needs to be passed along the future, so the model can decide to copy all information and reduce the vanishing problem

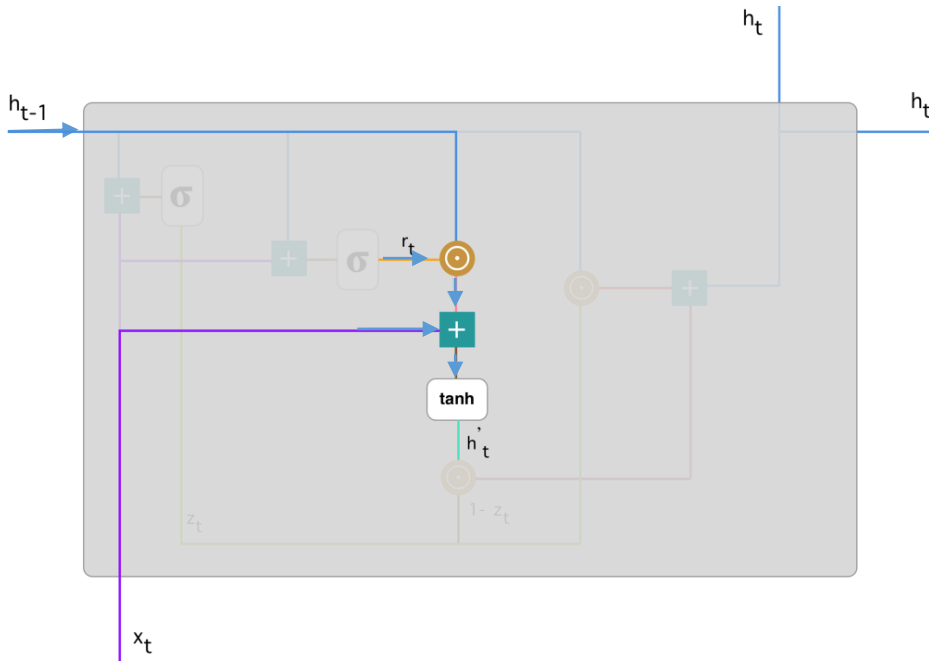
- Reset gate



$$r_t = \sigma(W^r x_t + U^r h_{t-1})$$

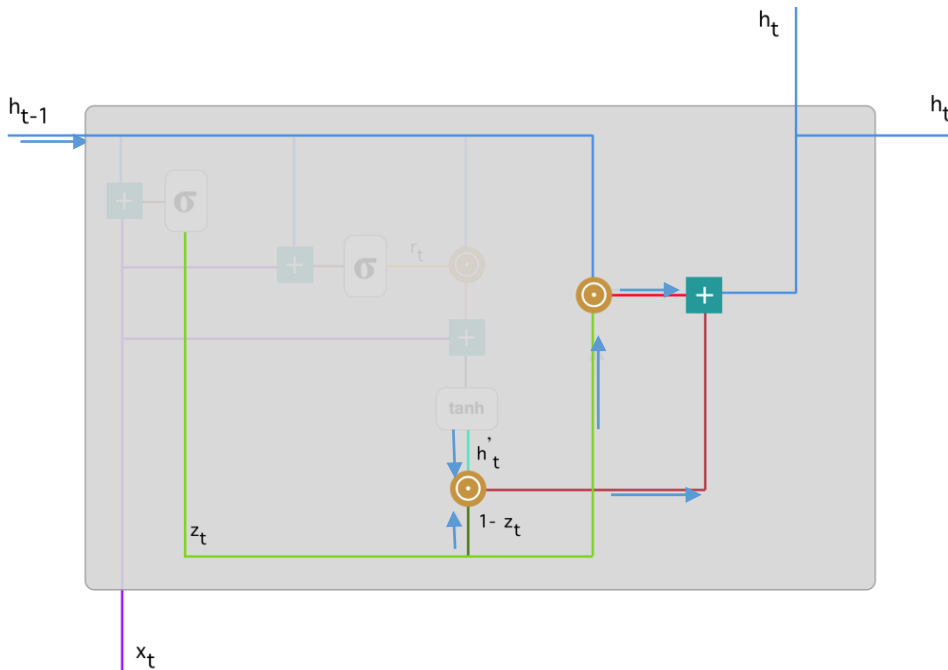
- This gate is used to decide how much of the past information to forget
- The formula is almost same as the update gate, but the weights

- Current Memory Content



$$h'_t = \tanh(Wx_t + r_t \odot Uh_{t-1})$$

- $r_t \odot Uh_{t-1}$  helps to remove the should-forget information, which determined in Reset gate. E.g. The text starts with “This is a fantasy book which illustrates...” and after a couple paragraphs ends with “I didn’t quite enjoy the book because I think it captures details which seem irrelevant to me.”. To determine the level of satisfaction, we only need last part of this review, so it will learn to  $r_t \approx 0$  to wash out the past and focus on the last sentence
  - If  $h'_t = 0$  means we don’t need to care about all words before, just care about the current word
- Final Memory



$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot h'_t$$

- GRU doesn’t use Cell State, but it uses  $h'_t$

### Vanishing Gradient

- Derivative of loss function  $J$  of variable  $W^{hh}$ :  $\frac{\partial J}{\partial W^{hh}} = \sum_t \frac{\partial J^{(t)}}{\partial W^{hh}}$
- Derivative of loss function  $J^{(t)}$  of variable  $W^{hh}$ :  $\frac{\partial J^{(t)}}{\partial W^{hh}} = \frac{\partial J^{(t+1)}}{\partial y_{t+1}} \frac{\partial y_{t+1}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial W^{hh}}$



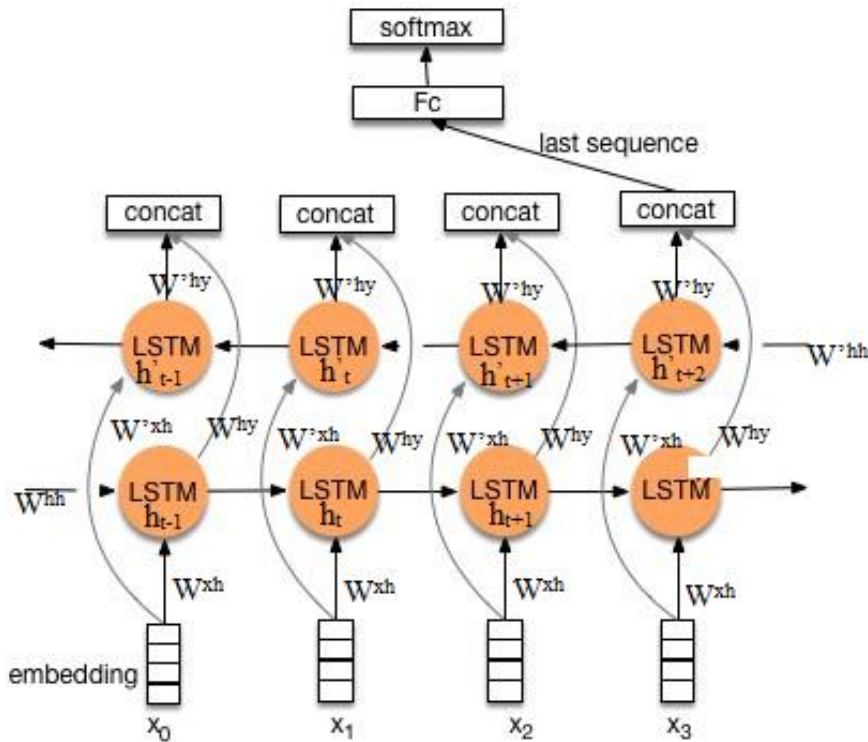
- However,  $\frac{\partial h_{t+1}}{\partial W^{hh}} = \frac{\partial h_{t+1}}{\partial h_t} \frac{\partial h_t}{\partial W^{hh}} + \frac{\partial h_{t+1}}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W^{hh}}$

$$\rightarrow \frac{\partial J_n}{\partial W^{hh}} = \sum_{k=0}^n \frac{\partial J_n}{\partial y_{t+1}} \frac{\partial y_{t+1}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_k} \frac{\partial h_k}{\partial W^{hh}}, \text{ at } k=0, \frac{\partial h_{t+1}}{\partial h_0} = \frac{\partial h_{t+1}}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \dots \frac{\partial h_1}{\partial h_0}$$

### Fasttext

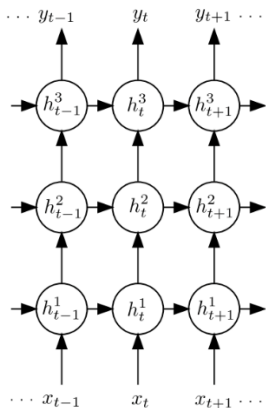
- An extension to Word2Vec proposed by Facebook
- Instead of feeding individual words into RNN, we break words into n-grams
- For instance, the tri-grams for the word apple is app, ppl, and ple

### sBidirectional RNN



- While doing the first row  $h_{t-1}, h_t, \dots$ , you also do the second row like the first row but in the reverse way

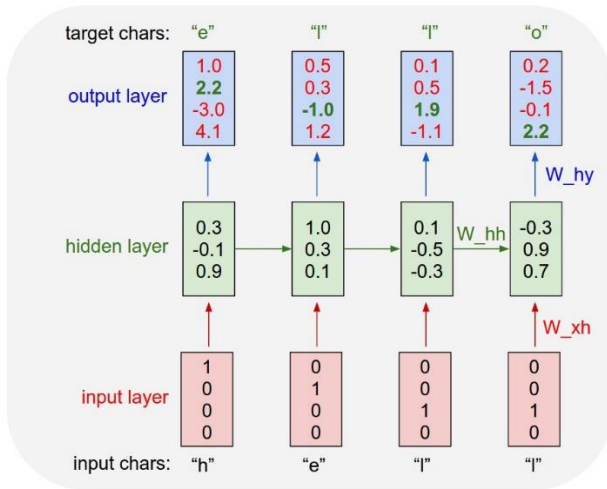
### Deep/ Stacked RNN



**Fig. 3.** Deep Recurrent Neural Network

### Character-level RNN

- Some words are Name, weird words, it's called OOV(Out-of-vocabulary). Character-level RNN are applied



- Character-level has covered all words but it cannot apply the pre-trained weights
- Moreover, there are some ways of word embedding, train weights
  - from scratch
  - apply pre-trained weight with
    - trainable = T: continue to train the weight with pre-trained
    - trainable = F: apply weights without any train

## Seq2seq

- I/O: Input is a sentence(corpus) and output is also a sentence(corpus)
- The problem Seq2seq can apply:
  - Machine translation
  - Image Captioning: input is image instead of corpus and output is still output, so in the encode session, we replace RNN by CNN
  - Summarization
  - Headline Generate
  - Add accent in Vietnamese
  - Punctuation

## Conditional Language Model

- Suppose you need to translate 'I am a student' to the French, you need to calculate:

$$P(y_1 y_2 \dots y_N | x) = \prod_{i=1}^N P(y_i | y_1 \dots y_{i-1}, x)$$

in which N: number of word in translated sentence,  $y_1, y_2, \dots, y_N$  is words in translated sentence

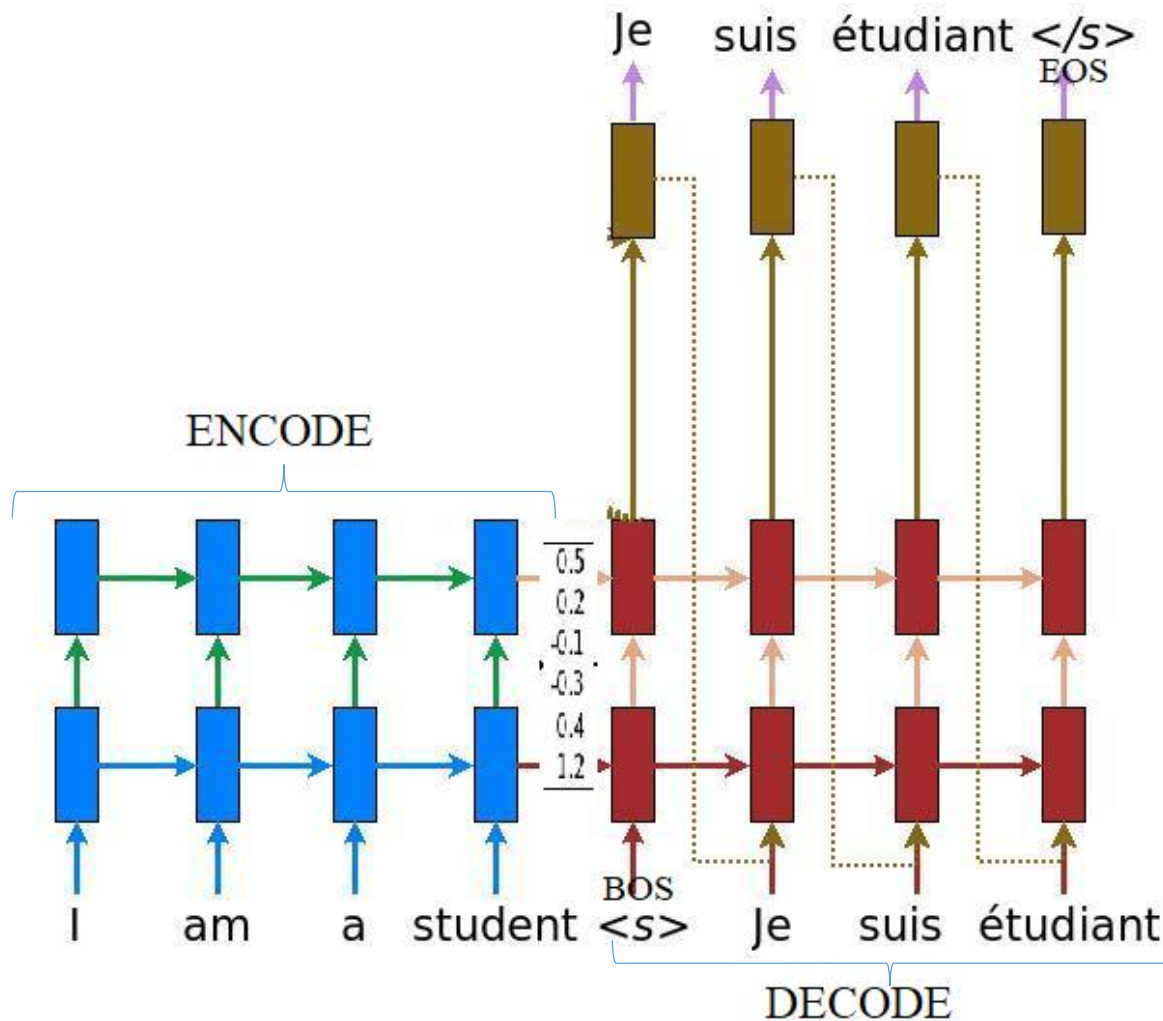
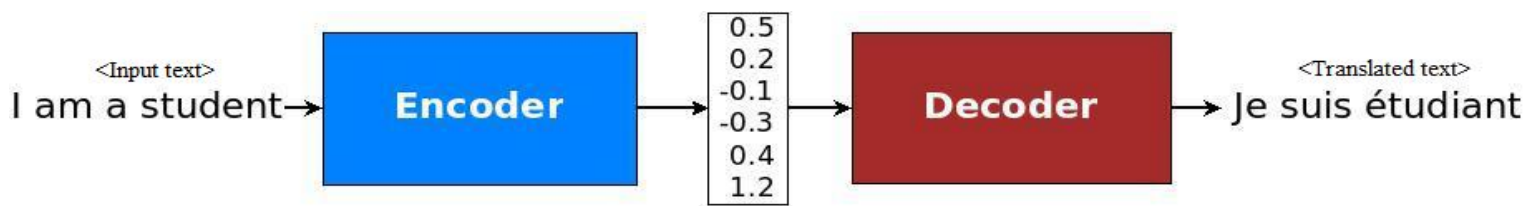
- Cost function:  $J(\theta) = -\sum_{j=1}^m \sum_{i=1}^N \log P(y_i^{(j)} | y_1^{(j)} \dots y_{i-1}^{(j)}, x^{(j)})$  in which m is number of sentence

- Inference/ Feed Forward Process:  $\arg \max_y P(y_1 y_2 \dots y_N | x) = \arg \max_y \prod_{i=1}^N P(y_i | y_1 \dots y_{i-1}, x)$

- However,  $\prod_{i=1}^N P(y_i | y_1 \dots y_{i-1}, x)$  easily  $\rightarrow 0$ , We just log them,

$$\arg \max_y \log P(y_1 y_2 \dots y_N | x) = \arg \max_y \sum_{i=1}^N \log P(y_i | y_1 \dots y_{i-1}, x)$$

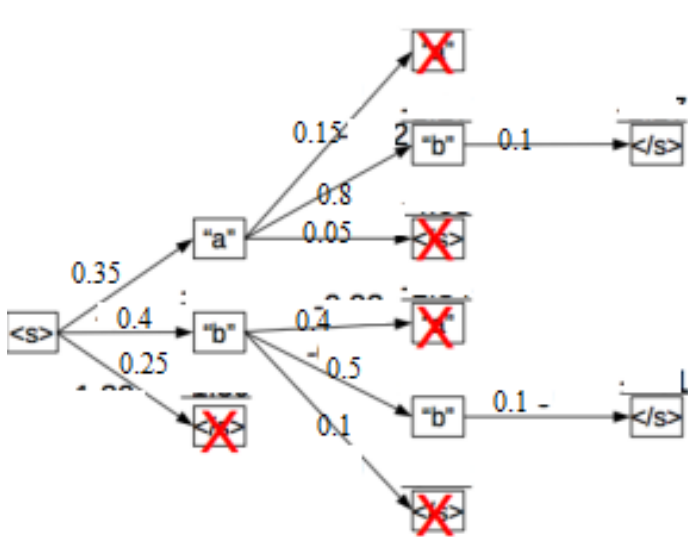
- Architecture



- The input is word vector of 'I', 'am', 'a' and 'student', the word vector is built on English words only
- The output for 'Je' is the list of all probability of all French words which have probability of 'Je' is the largest
- This architecture has modelled the Inference/ Feed forward process. E.g.  
 $P(Je | BOS, I, am, a, student) \times P(suis | Je, BOS, I, am, a, student) \times \dots \times P(EOS | étudiant, suis, Je, BOS, I, am, a, student)$
- The problem occurs when you just choose the largest probability to be your output because the model mention

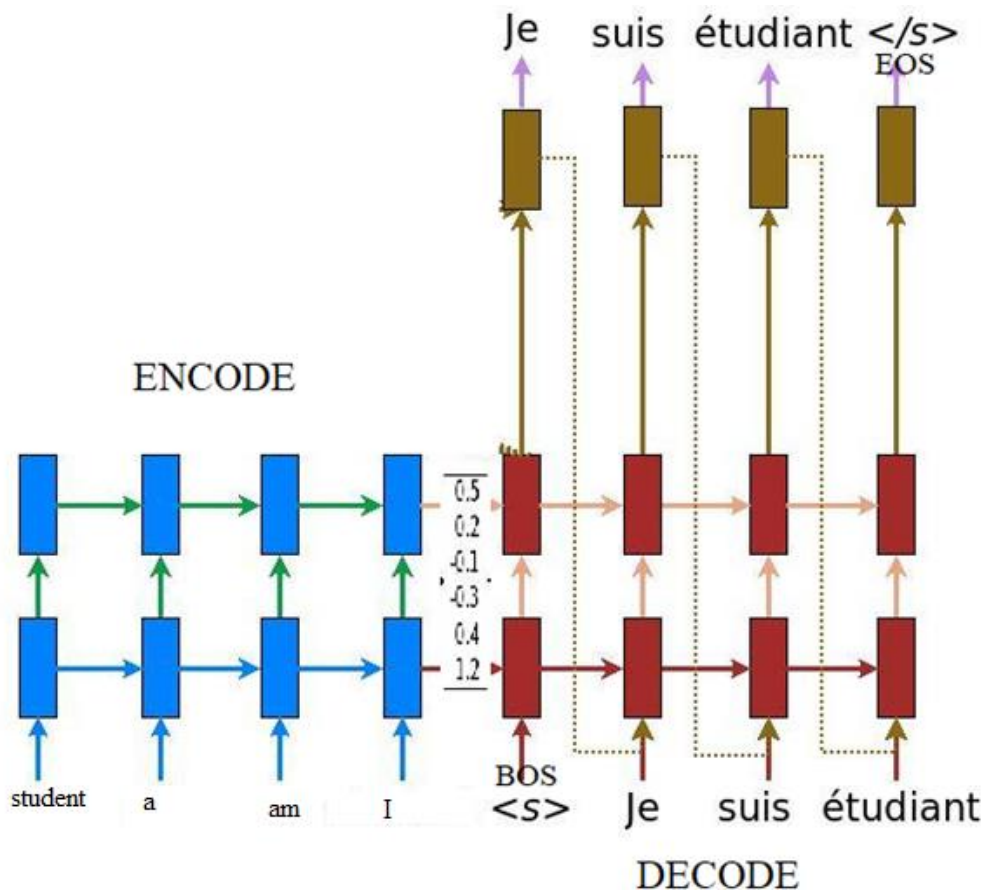
$\arg \max_y \prod_{i=1}^N P(y_i | y_1 \dots y_{i-1}, x)$ , we can solve by some options

- Random Sampling (Ancestral Sampling): in the step  $i$  in  $N$ , you choose the word based on the distribution which follows distribution of decode output at step  $i$ . Application in ChatBot
- Greedy Search: mentioned above which choose the largest probability in the output at step  $i$
- Beam Search: at step  $i$ , we just keep  $k$ / Beam Size/ Beam Width string which has length of  $i$ . E.g. Beam Size = 2



○ The string here:  $\langle s \rangle$  a b  $\langle /s \rangle$

- Another problem: for instance, you translate “I am a student” to French, so “I” should be translated to “Je”, but in the architecture we consider, “Je” has to go through “student”, “a”, “am”, so the weight of “I” impacts on prediction of “Je” is now a little. To solve it, we reverse the sentence

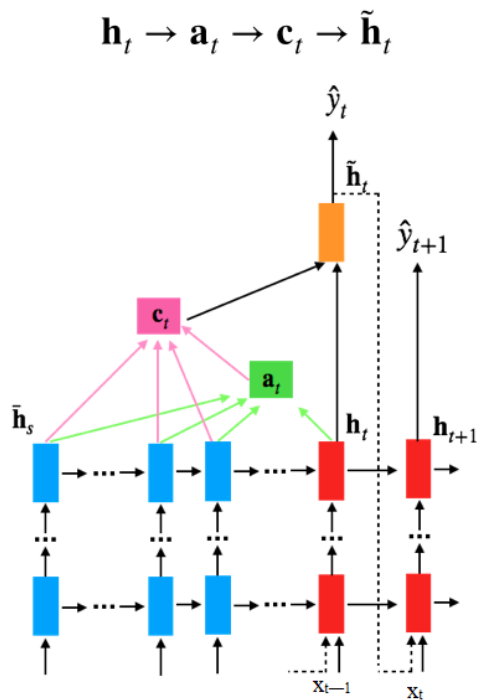


- However, we also have another method to solve completely this problem, called Attention, or Bahdanau Attention

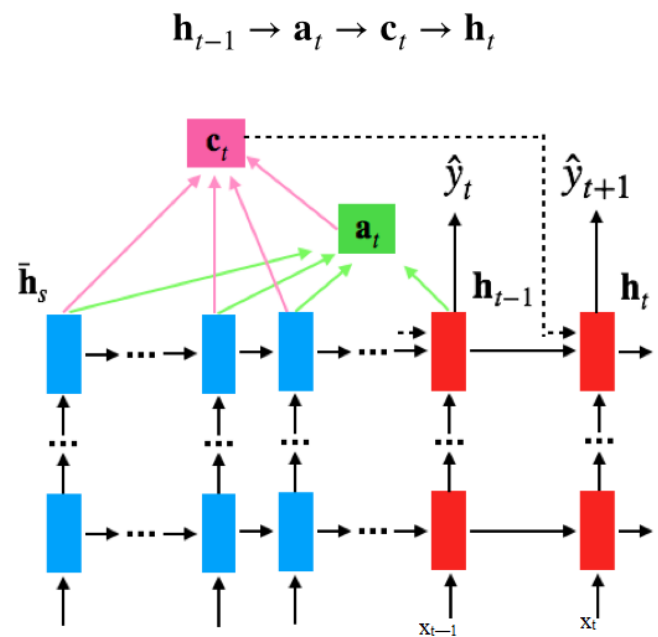
#### Attention

- Calculate the score to measure the alignment of target to the all words of input sentence
- Now, we want to find out  $h_t$ , which is  $x_t$  is ‘Je’

## Luong Attention Mechanism



## Bahdanau Attention Mechanism



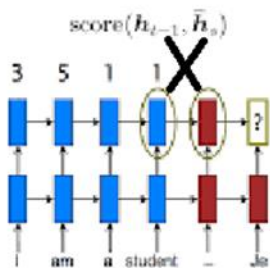
$$\begin{aligned} \mathbf{a}_t(s) &= \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \\ \mathbf{c}_t &= \sum \mathbf{a}_t \mathbf{h}_s \\ \tilde{\mathbf{h}}_t &= \tanh(W_c[\mathbf{c}_t; \mathbf{h}_t]) \end{aligned}$$

$$\begin{aligned} \mathbf{a}_t(s) &= \text{align}(\mathbf{h}_{t-1}, \bar{\mathbf{h}}_s) \\ \mathbf{c}_t &= \sum \mathbf{a}_t \mathbf{h}_s \\ \mathbf{h}_t &= \text{RNN}(\mathbf{h}_{t-1}^{l-1}, [\mathbf{c}_t; \mathbf{h}_{t-1}]) \end{aligned}$$

### Bahdanau Attention

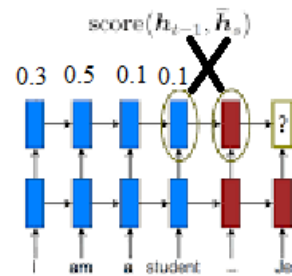
- $\text{align}(h_{t-1}, \bar{h}_s) = \text{score}(h_{t-1}, \bar{h}_s) = v_\alpha^T \tanh(W_1 h_{t-1} + W_2 \bar{h}_s)$
- $v_\alpha, W_1, W_2$  is variable like  $w_{hh}, w_{xh}$  and initialize these variable in the beginning

### Attention Mechanism - Scoring



- Compare target and source hidden states.

### Attention Mechanism - Scoring



- Compare target and source hidden states.

- Then, put the score vector of  $h_{t-1}$  to the softmax:
 
$$\alpha_{ts} = \frac{\exp(\text{score}(h_{t-1}, \bar{h}_s))}{\sum_{s'} \exp(\text{score}(h_{t-1}, \bar{h}_{s'}))}$$
- $$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s$$

- $h_t = f(h_{t-1}, x_t, c_t) = \sigma(W^{hh}h_{t-1} + W^{xh}x_t + Cc_t)$
- $f(h_{t-1}, x_t, c_t)$ , we can use LSTM, GRU or Vanilla as long as we add  $Cc_t$  when calculating  $h_t$

- Luong Attention

- Assume we have  $\tilde{h}_{t-1}$  and  $x_{t-1}$  and we use LSSTM/GRU or Vanilla to calculate  $h_t$

- $$align(h_t, \bar{h}_s) = score(h_t, \bar{h}_s) = \begin{bmatrix} h_t^T \bar{h}_s \\ h_t^T W \bar{h}_s \\ v_\alpha^T \tanh(W_1 h_t + W_2 \bar{h}_s) \end{bmatrix}$$

- Then, put the score vector of  $h_t$  to the softmax: 
$$\alpha_{ts} = \frac{\exp(score(h_t, \bar{h}_s))}{\sum_{s'} \exp(score(h_t, \bar{h}_{s'}))}$$

- $$c_t = \sum_s \alpha_{ts} \bar{h}_s$$

- $$\tilde{h}_t = \tanh(W_C [c_t, h_t])$$

- Now, we have  $\tilde{h}_t$  and  $x_t$  to calculate  $h_{t+1}$