

# Generating Text Using a Character RNN and LSTM

Lecture 11 (eng)

# Contents

---

- Understand of LSTM
- Text Generation Using a character RNN
- Generating Shakespearean Text Using a LSTM, GRU



*PANDARUS:*

*Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.*

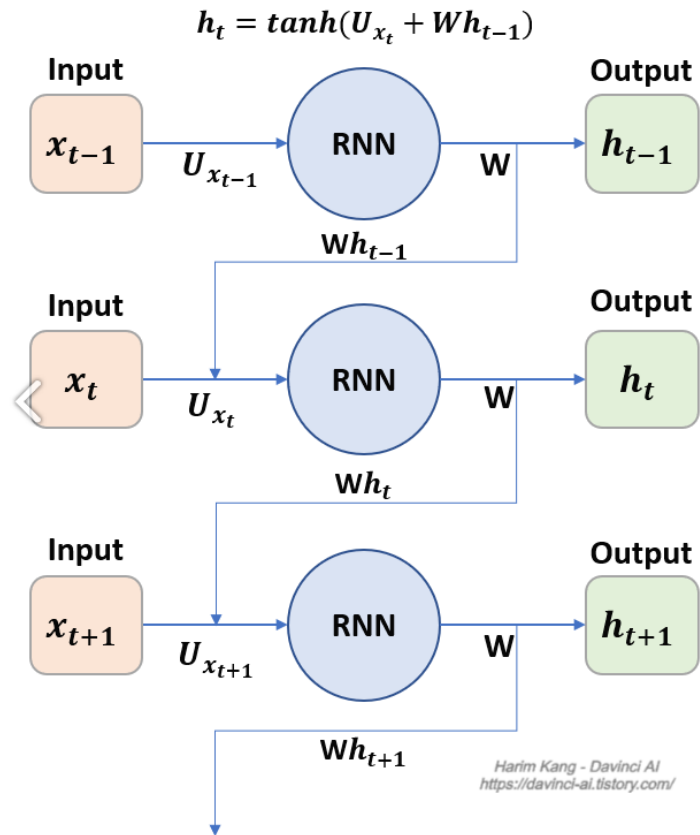
# RNN-NLP for UST AI

---

- download ppt and codes
- <https://github.com/hongsukyi/rnn-nlp-lectures>

# Review of Recurrent Neural Networks

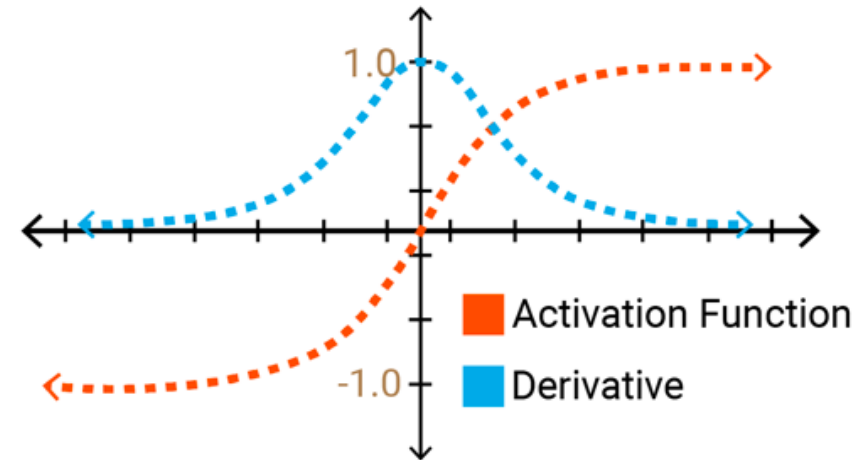
- In standard RNNs, this repeating module will have a single tanh layer.



$$h_t = \tanh(U_{x_t} + Wh_{t-1})$$

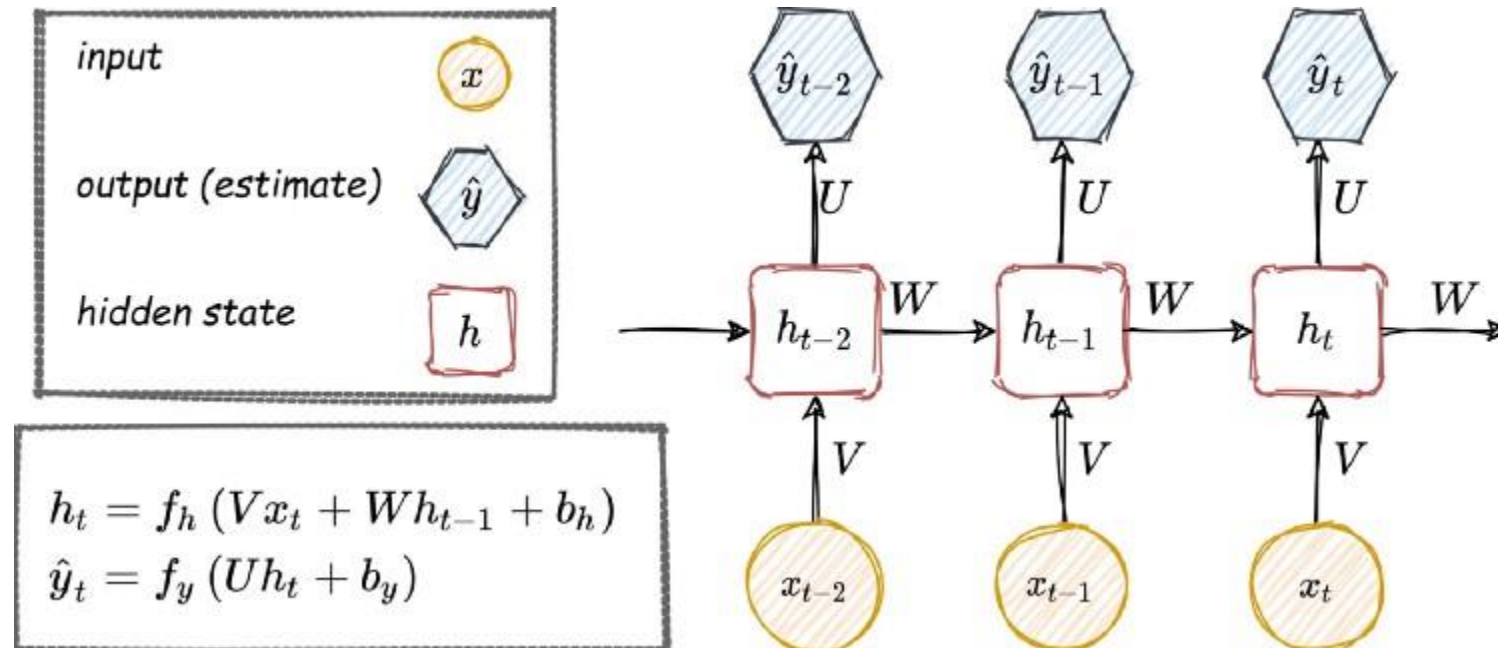
$$f(x) = \frac{2}{1 + e^{-2x}} - 1 = \tanh(x)$$

$$f'(x) = 1 - \tanh^2(x)$$

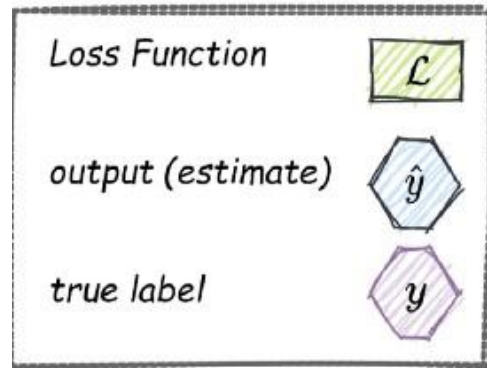


# Review of Recurrent Neural Networks

- A simple RNN architecture is shown below, where  $V$ ,  $W$ , and  $U$  are the weights matrices, and  $b$  is the bias vector.



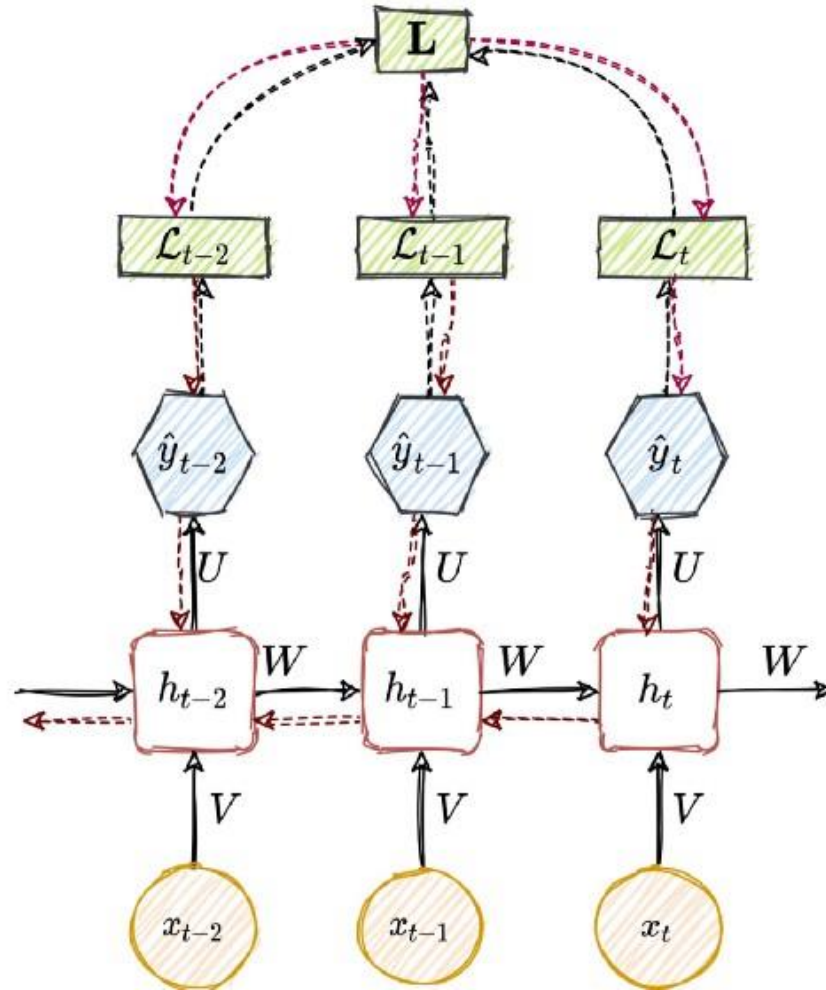
# Backpropagation Through Time (BPTT)



$$\mathbf{L} = \sum_i \mathcal{L}_i(\hat{y}_t, y_t)$$

Forward Pass:  
 $h_t, \hat{y}_t, \mathcal{L}_t, \mathbf{L}$

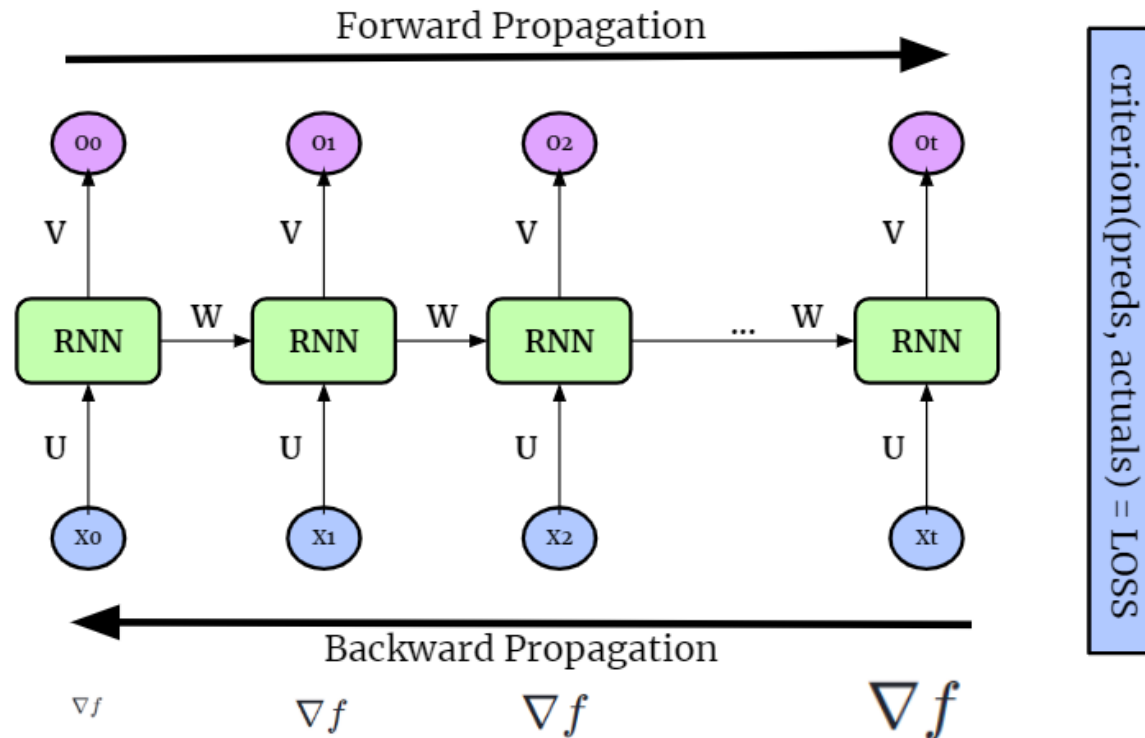
Backward Pass:  
 $\frac{\partial \mathbf{L}}{\partial U}, \frac{\partial \mathbf{L}}{\partial V}, \frac{\partial \mathbf{L}}{\partial W}, \frac{\partial \mathbf{L}}{\partial b_h}, \frac{\partial \mathbf{L}}{\partial b_y}$



$$\frac{\partial \mathbf{L}}{\partial W} \propto \sum_{i=0}^T \left( \prod_{i=k+1}^y \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W}$$

# Vanishing Gradient Problem

- it is due to the nature of backpropagation (during the optimization process)
  - ✓ if the adjustment in the previous layer is small, then that in the current layer will be smaller



# RNNs suffer from the problem of VGP and EGP

---

## ➤ VGP(Vanishing Gradient Problem)

- ✓ when the gradient becomes too small, the parameter updates become insignificant. This makes the learning of long data sequences difficult.

## ➤ EGP(Exploding Gradient Problem)

- ✓ If the slope tends to grow exponentially instead of decaying, when large error gradients accumulate during the training process.



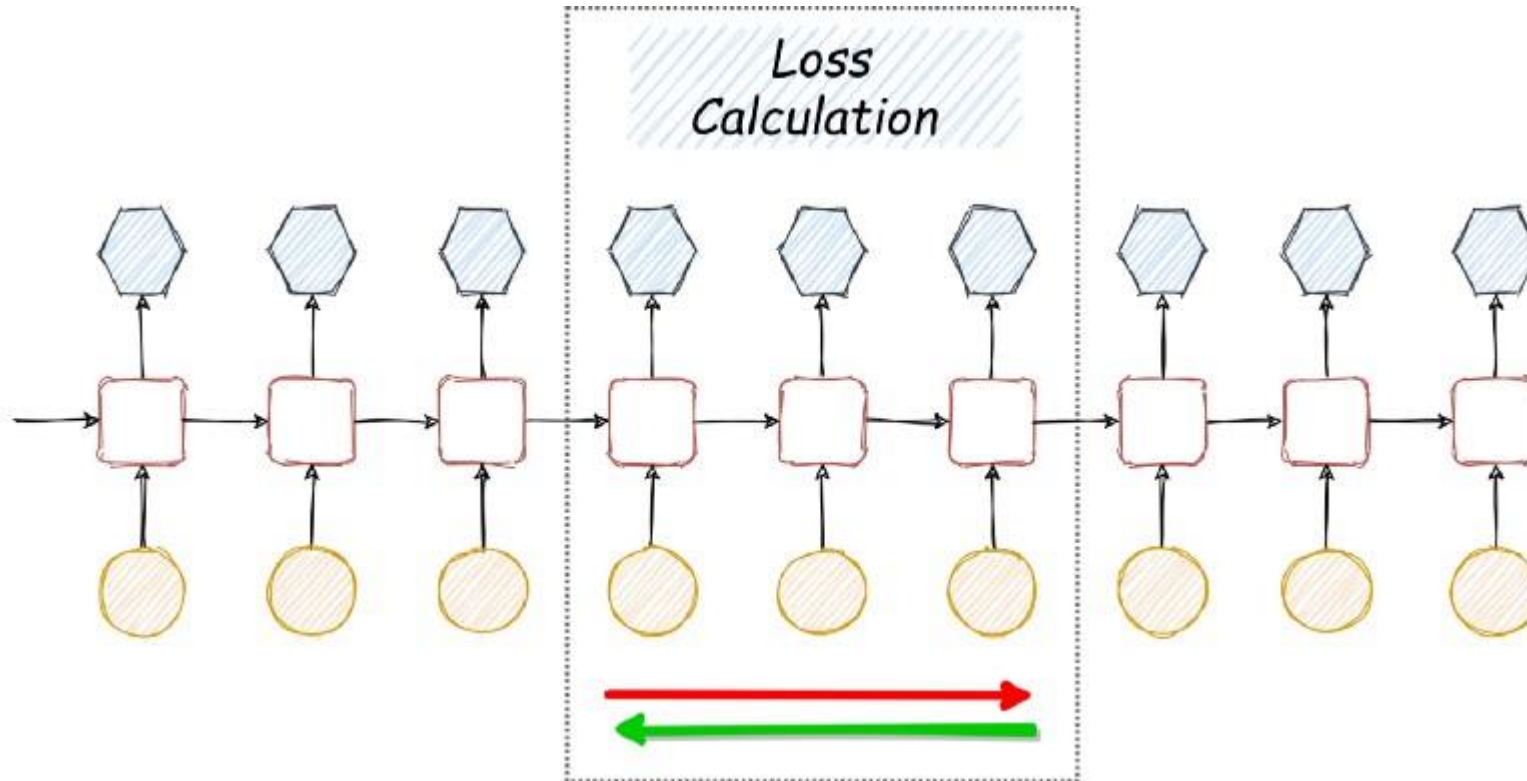
1. *Vanishing gradient*  $\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 < 1$

2. *Exploding gradient*  $\left\| \frac{\partial h_i}{\partial h_{i-1}} \right\|_2 > 1$



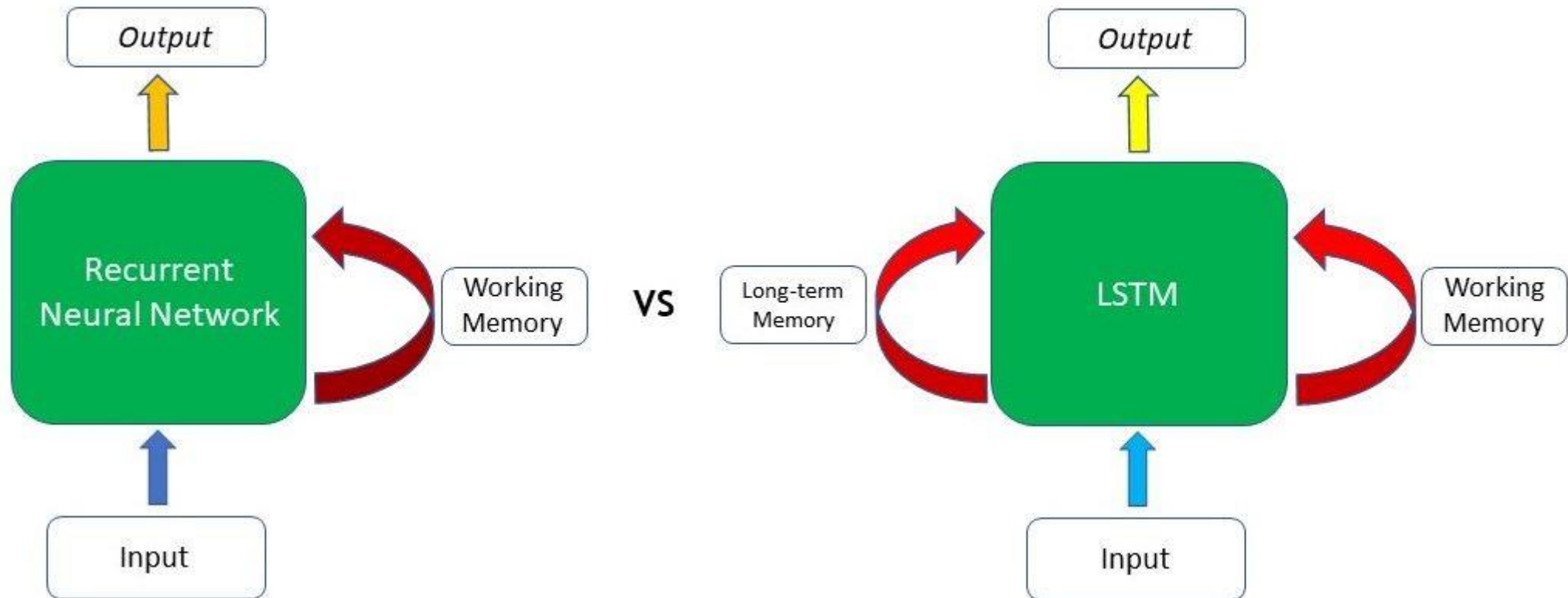
# Truncated Backpropagation Through Time (Truncated BPTT)

- Truncated BPTT trick tries to overcome the VGP
  - ✓ by considering a moving window through the training process



# Long Short-Term Memory networks (LSTMs)

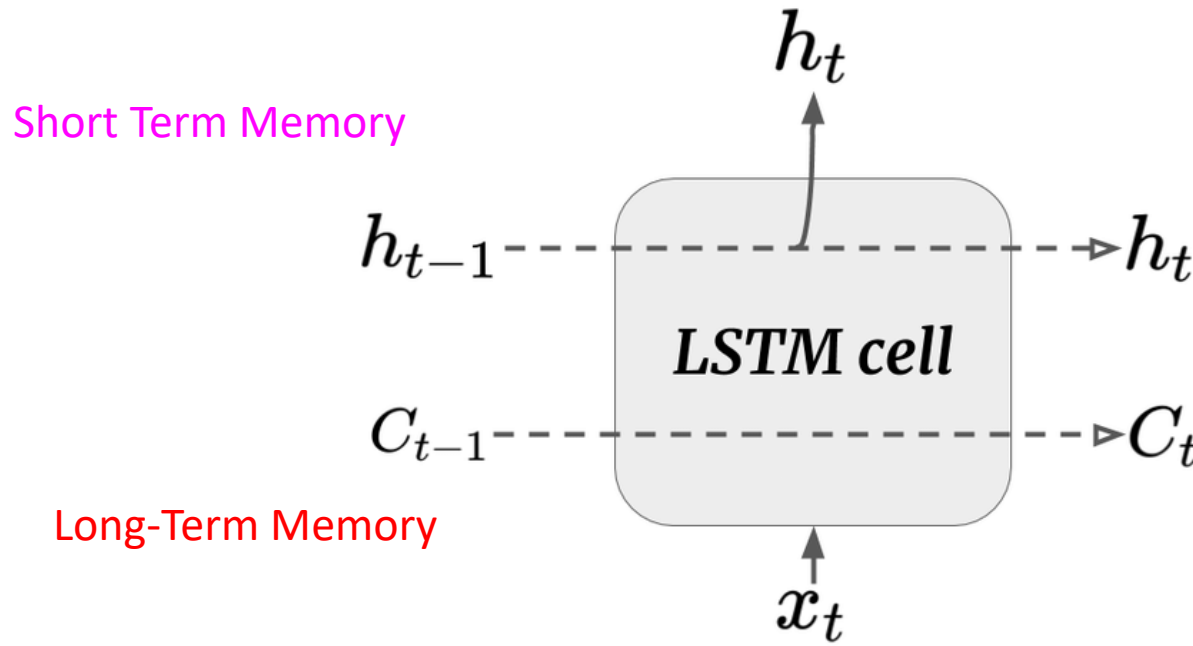
- Allows learning of long-term dependencies
  - ✓ A RNN addresses the vanishing/exploding gradient problem
  - ✓ Capable of learning long-term dependencies by remembering information



# Long Short Term Memory Network

## ➤ LSTM (Long Short Term Memory Network)

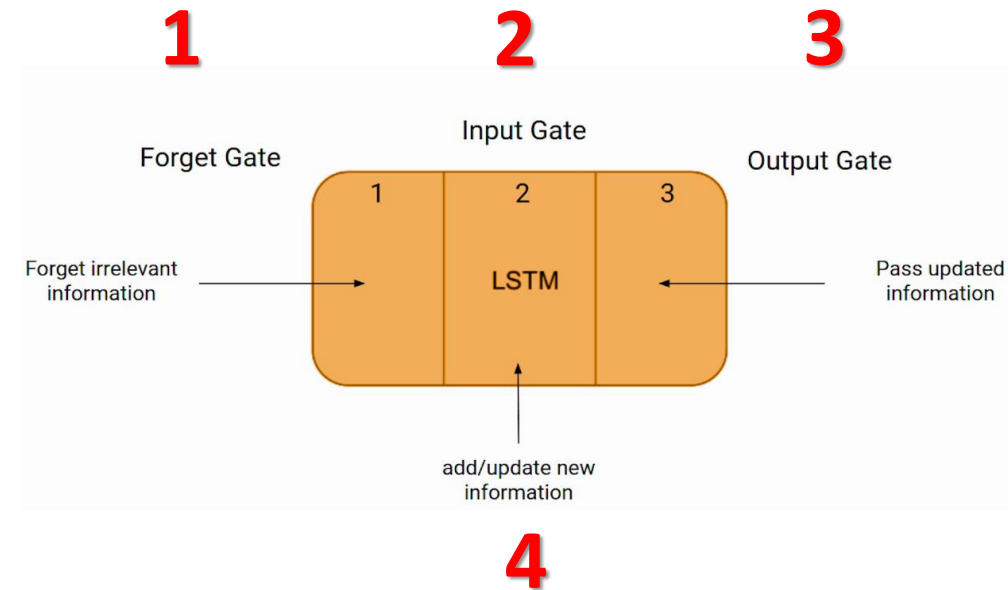
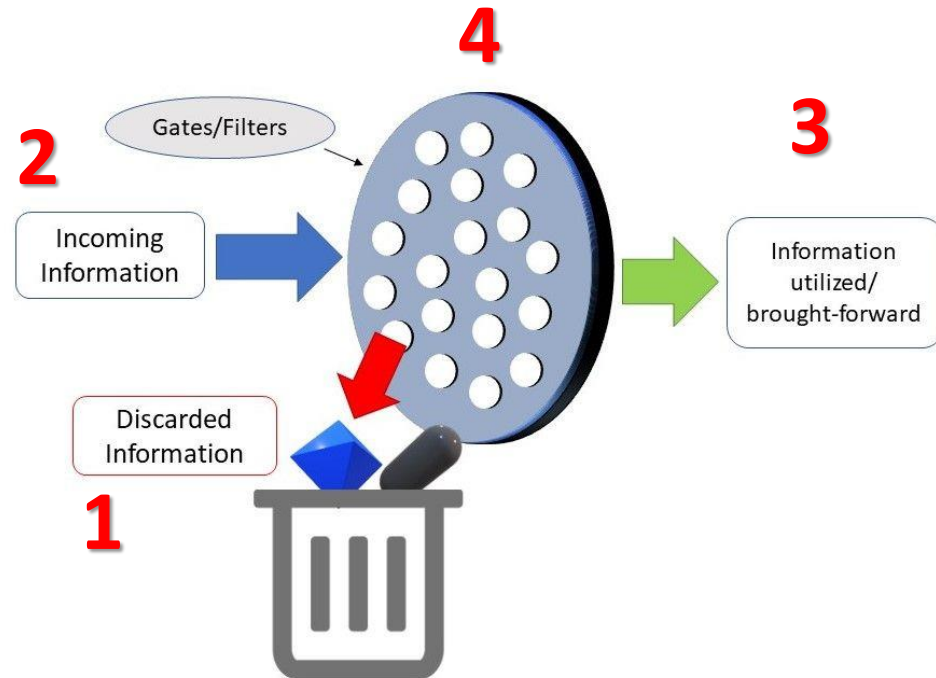
- ✓ LSTMs are explicitly designed to avoid long-term dependency problems
- ✓ The observed state  $x_t$  is combined with previous memory and hidden states to output a hidden state  $h_t$ .



# Gates

## ➤ Gates control the flow of information to/from the memory

- ✓ Gates are controlled by a concatenation of the output from the previous time step and the current input and optionally the cell state vector.



# Understanding the roles played by gates in LSTM architecture

---

## ➤ Forget Gate

- ✓ whether we should keep the information from the previous timestamp or forget it

## ➤ Input Gate

- ✓ Decide how much this unit adds to the current state

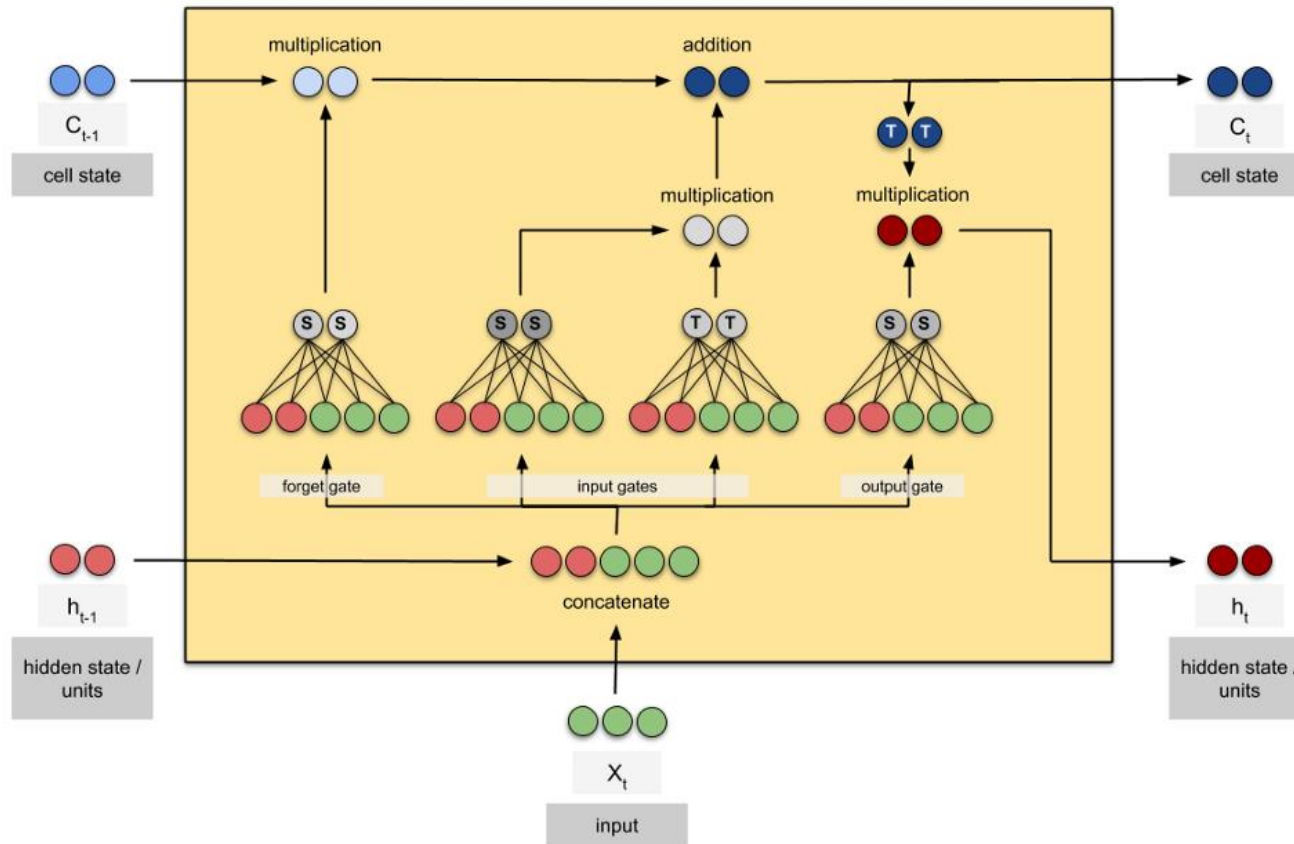
## ➤ New information: Memory Upgate

- ✓ The cell state vector aggregates the two components (old memory via the forget gate and new memory via the input gate)

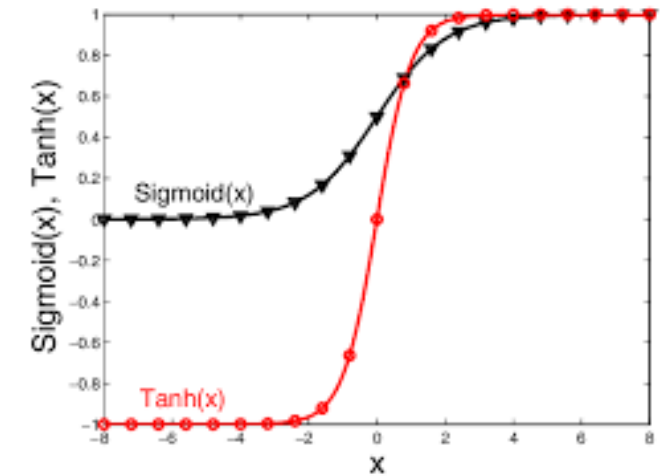
## ➤ Output Gate

- ✓ Decide what part of the current cell state makes it to the output

# LSTM Process



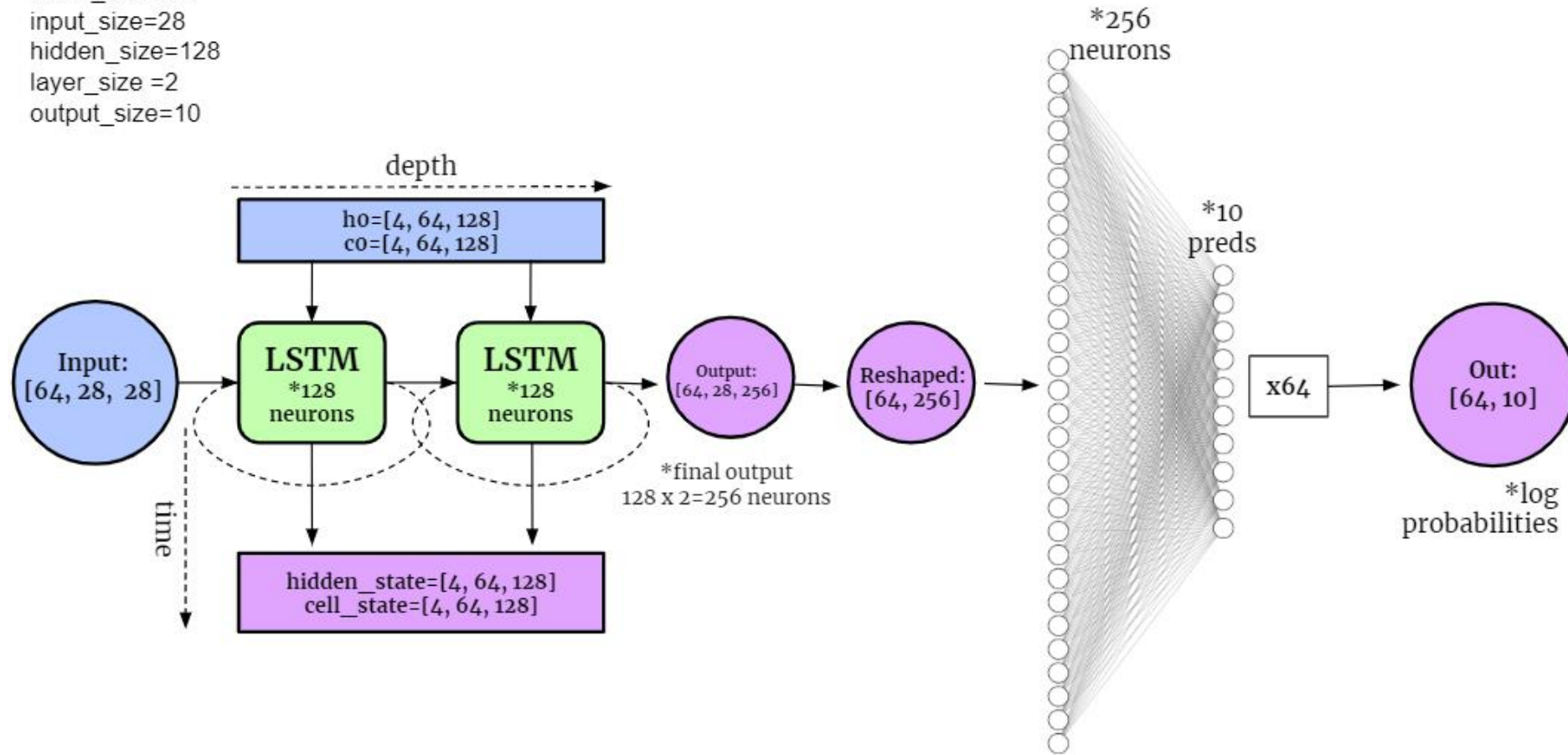
## Activation functions



# LSTM Example

## LSTM Example

batch\_size=64  
input\_size=28  
hidden\_size=128  
layer\_size =2  
output\_size=10



# Lab 11–1

# Text Tokenizer

Text Generation



# Understanding the principles

---

- Token : Language elements that we can't share anymore
- Tokenizer
  - ✓ work to input text data into the neural network.
  - ✓ The preprocessing process that converts it into an appropriate form through encoding
- One-hot encoding
  - ✓ In the case of text data, an embedding layer is basically used.



# Tokenizing words

---

- Word tokenization divides sentences based on spacing as follows.

“This book is for deep learning learners”

Tokenizing



This	book	is	for	deep	learning	learners
------	------	----	-----	------	----------	----------

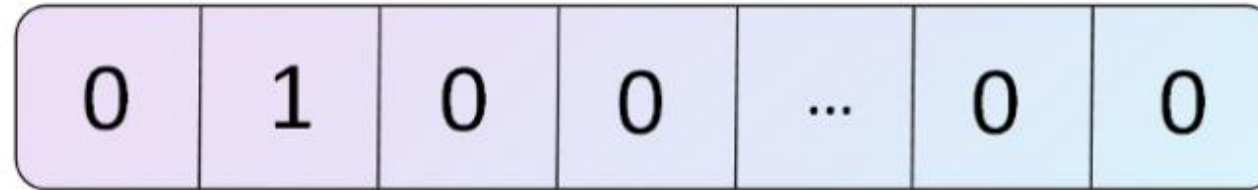
Copyright © Gilbut, Inc. All rights reserved.

# One-Hot Encoding

---






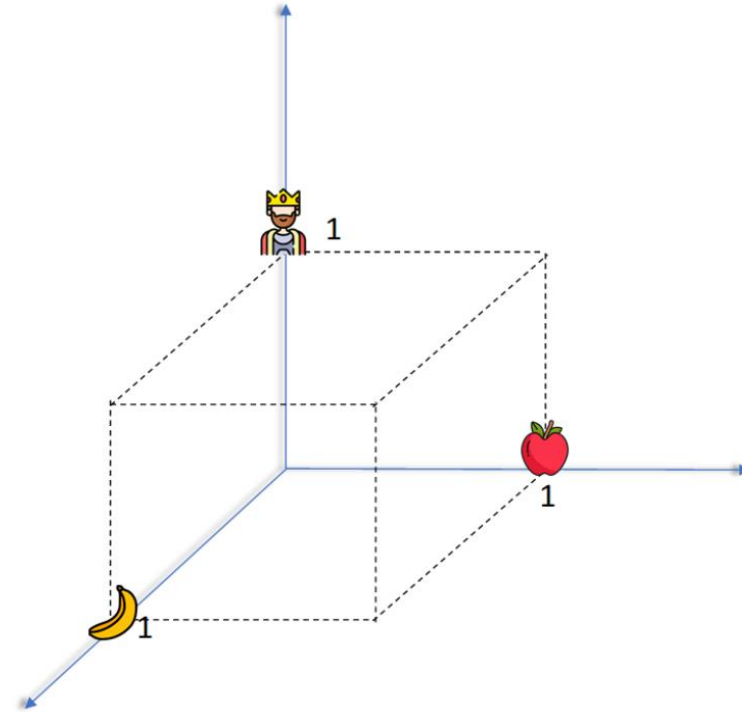
Index:    0        1        2        3        ...    99998   99999



Index:    0        1        2        3        ...    99998   99999

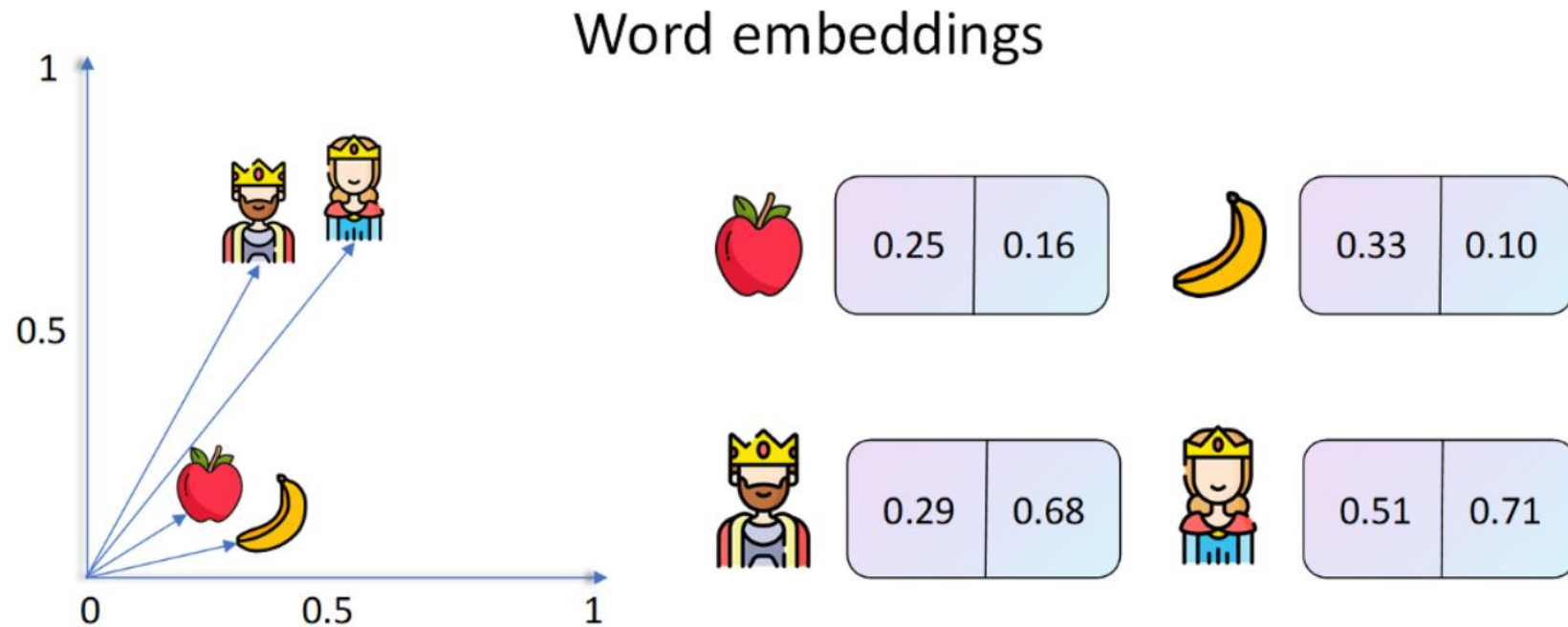
# Word Embedding

	<table><tr><td>1</td><td>0</td><td>0</td></tr></table>	1	0	0
1	0	0		
Index:	<table><tr><td>0</td><td>1</td><td>2</td></tr></table>	0	1	2
0	1	2		
	<table><tr><td>0</td><td>1</td><td>0</td></tr></table>	0	1	0
0	1	0		
Index:	<table><tr><td>0</td><td>1</td><td>2</td></tr></table>	0	1	2
0	1	2		
	<table><tr><td>0</td><td>0</td><td>1</td></tr></table>	0	0	1
0	0	1		
Index:	<table><tr><td>0</td><td>1</td><td>2</td></tr></table>	0	1	2
0	1	2		



# Word Embedding

- 2 dimensional word embedding representation of our example words



# 1. Word-based encoding

---

- The example below shows how to encode two sentences
  - ✓ 'You are the Best', and 'You are the Nice' based on words using TensorFlow.

## 1. Word-based encoding

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.utils import to_categorical
```

```
sentences = [
    'You are the Best',
    'You are the Nice'
]
```

# 1. Word-based encoding

- `fit_on_texts()` 메서드는 문자 데이터를 입력받아서 리스트의 형태로 변환합니다.

```
tokenizer = Tokenizer(num_words = 100, oov_token="<OOV>")
tokenizer.fit_on_texts(sentences)
word_index = tokenizer.word_index
```

Words that have not been indexed in advance are indexed as "OOV"

```
print(word_index)
print('-----')
total_words = len(tokenizer.word_index) + 1
print('total_words=', total_words)
```

```
{'<OOV>': 1, 'you': 2, 'are': 3, 'the': 4, 'best': 5, 'nice': 6}
-----
total_words= 7
```

The `word_index` attribute of the tokenizer returns a dictionary containing a pair of key-values of words and numbers.

The output results show that the upper case 'I' has been converted to the lower case 'i'.

## 2) Converting text into a sequence

---

텍스트를 시퀀스로 변환하기

```
sequences = tokenizer.texts_to_sequences(sentences)
```

```
print(word_index)
```

```
print(sequences)
```

```
{ '<OOV>': 1, 'you': 2, 'are': 3, 'the': 4, 'best': 5, 'nice': 6 }
```

```
[[2, 3, 4, 5], [2, 3, 4, 6]]
```



### 3. Setting up padding

---

➤ You have to padding to make the sentence the same length.

- ✓ Padding uses the `pad_sequences` function.

```
from tensorflow.keras.preprocessing.sequence import pad_sequences
```

Sequences are text sentences converted into sequences of integers

- Since the longest sequence is 7, it has all been converted into sequences of the same length

```
padded = pad_sequences(sequences)
```

```
print(word_index)
```

```
print(sequences)
```

```
print(padded)
```

```
{ '<OOV>': 1, 'you': 2, 'are': 3, 'the': 4, 'best': 5, 'nice': 6 }
```

```
[[2, 3, 4, 5], [2, 3, 4, 6]]
```

```
[[2 3 4 5]
```

```
 [2 3 4 6]]
```

### 3. Setting up padding

---

➤ padding parameter : 'pre', 'post'

- ✓ If the padding parameter is specified as 'post', padding is filled after the sequence. The default is "pre"

```
padded = pad_sequences(sequences, padding='post')  
print(padded)
```

```
[[2 3 4 5]  
 [2 3 4 6]]
```

# 이진 형태로 인코딩합니다.

```
binary_results = tokenizer.sequences_to_matrix(sequences, mode = 'binary')
print(f'binary_vectors:\n {binary_results}\n')
```

binary\_vectors:

[illegible][illegible]

## 4) Encoding in binary form

---

```
print(f'One-Hot Encoding:',to_categorical(sequences))
```

```
One-Hot Encoding: [[[0. 0. 1. 0. 0. 0. 0.]  
[0. 0. 0. 1. 0. 0. 0.]  
[0. 0. 0. 0. 1. 0. 0.]  
[0. 0. 0. 0. 0. 1. 0.]]]
```

```
[[[0. 0. 1. 0. 0. 0. 0.]  
[0. 0. 0. 1. 0. 0. 0.]  
[0. 0. 0. 0. 1. 0. 0.]  
[0. 0. 0. 0. 0. 0. 1.]]]
```

```
test_text = ['You are the One']  
test_seq = tokenizer.texts_to_sequences(test_text)
```

```
print(f'test sequences: {test_seq}')
```

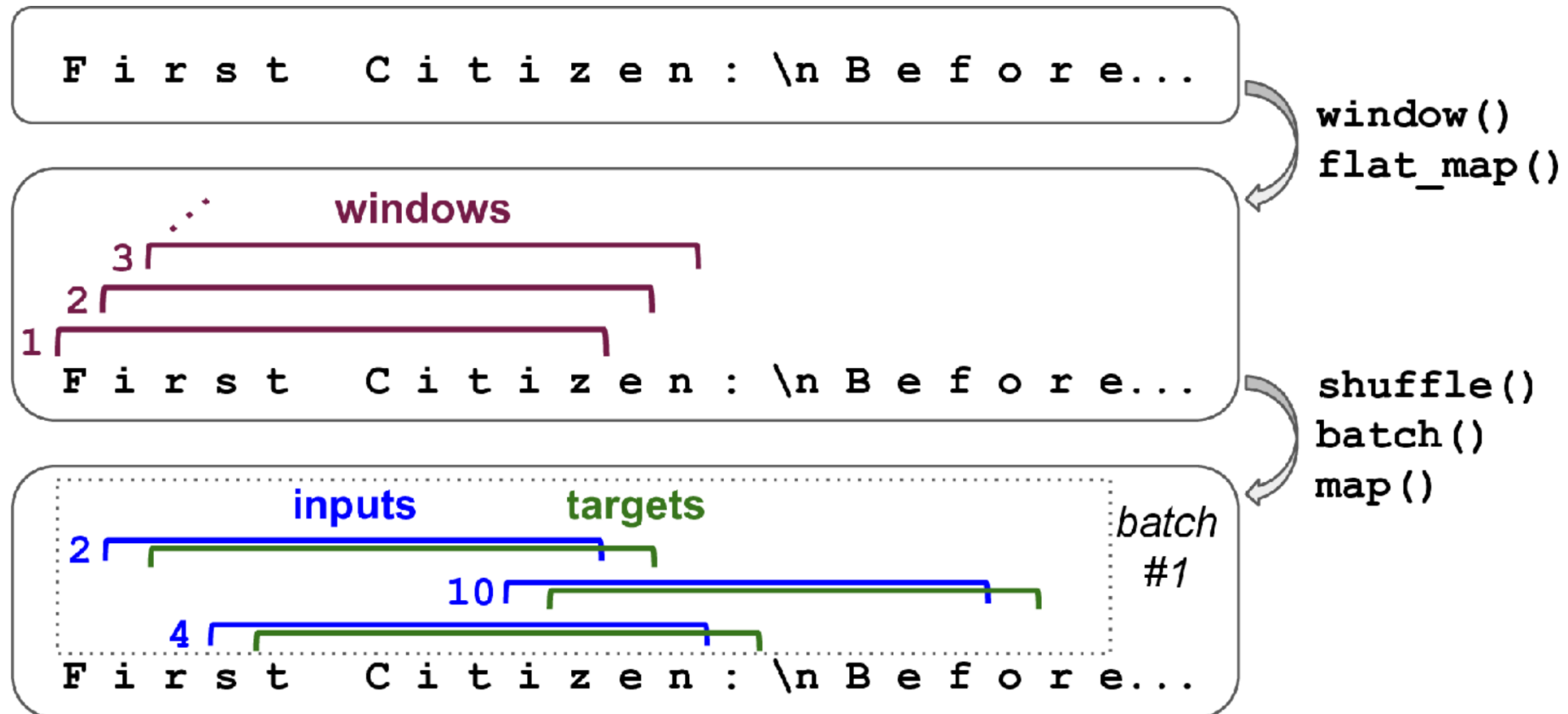
```
test sequences: [[2, 3, 4, 1]]
```

# Lab 11–2

## Generating Text Using a Character RNN

Text Generation

# Preparing a dataset of shuffled windows

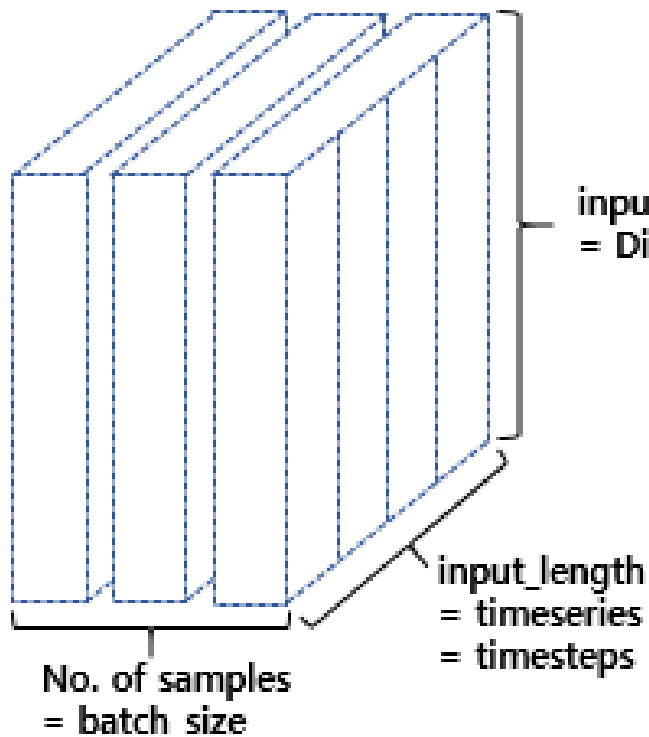


# Text Generation through LSTM model

Lab11-2

## 1) 데이터에 대한 이해와 전처리

원-핫 벡터의 차원은 글자 집합의 크기인 56이어야 하므로 원-핫 인코딩이 수행되었다.



```
print('train_X의 크기(shape) : {}'.format(train_X.shape))  
print('train_y의 크기(shape) : {}'.format(train_y.shape))
```

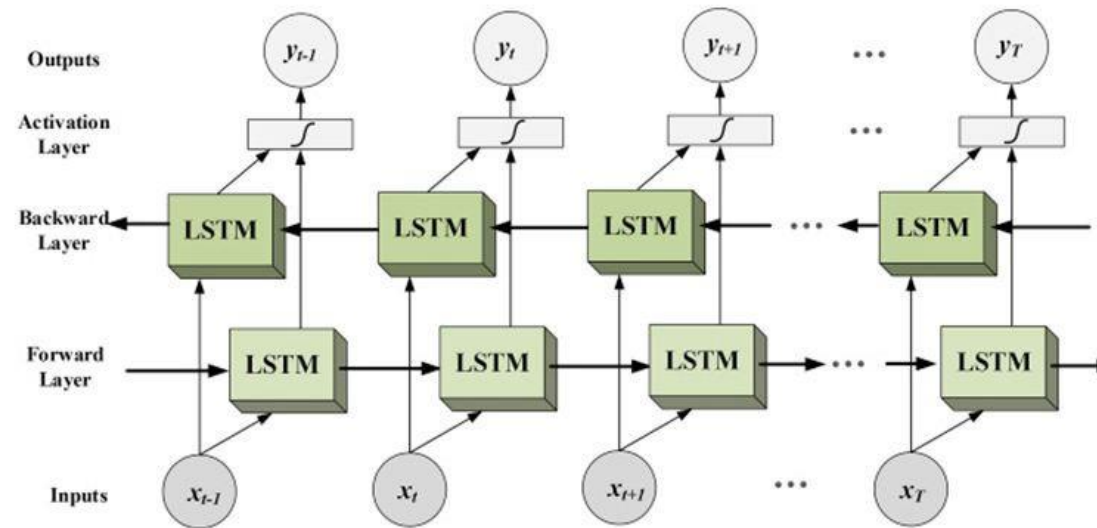
train\_X의 크기(shape) : (2658, 60, 56)

train\_y의 크기(shape) : (2658, 60, 56)



# 1) Bidirectional LSTMs

- Bidirectional LSTMs are an extension of traditional LSTMs that can improve model performance on sequence classification problems.
  - ✓ They train the model forward and backward on the same input (so for 1 layer LSTM we get 2 hidden and cell states)
  - ✓ First from left to right on the input sequence and the second in reversed order of the input sequence.



# 1) Bidirectional LSTMs

---

- implement this model in text generation

```
import numpy as np
import tensorflow as tf

from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.layers import Embedding, LSTM, Dense, GRU, Bidirectional
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.models import Sequential
from tensorflow.keras.optimizers import Adam
from tensorflow.keras import regularizers
from tensorflow.keras.utils import to_categorical, plot_model
```



## 2) Text pre-processing

---

```
print(corpus[:15])
```

```
['first citizen:', 'before we proceed any further, hear me speak.', '', 'all:', 'speak, speak.', '', 'first citizen:', 'you are a  
ll resolved rather to die than to famish?', '', 'all:', 'resolved. resolved.', '', 'first citizen:', 'first, you know caius marci  
us is chief enemy to the people.', '']
```

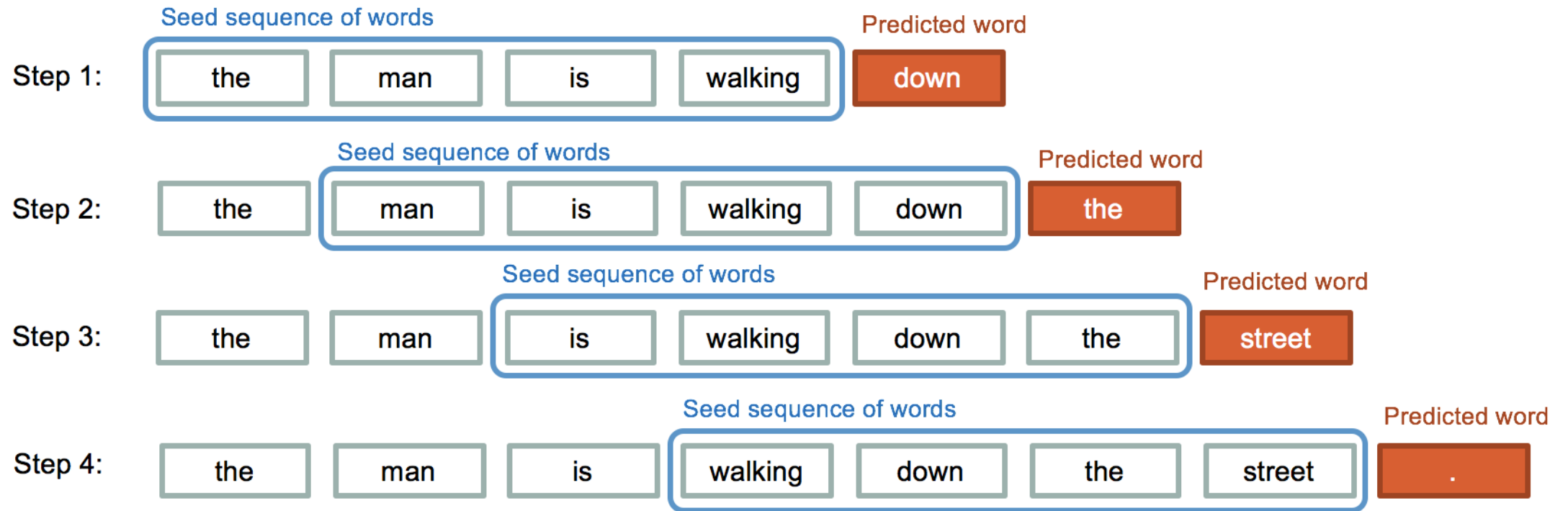
```
total_words = len(tokenizer.word_index) + 1  
print('total_words=', total_words)
```

```
total_words= 672
```

### 3) Creating Sequences

---

- For each word, an n-gram sequence is made and input sequences are updated.
  - ✓ It happens in the iteration for the next word and so on



### 3) Creating Sequences

#### ➤ For example

- ✓ in the sentence below first 'He' was extracted out then, 'He was ' was extracted, and then 'He was walking ' was extracted, and so on.

```
# create input sequences using list of tokens
input_sequences = []
for line in corpus:
    token_list = tokenizer.texts_to_sequences([line])[0]
    for i in range(1, len(token_list)):
        n_gram_sequence = token_list[:i+1]
        input_sequences.append(n_gram_sequence)
```

Can you please come **here** ?

The diagram shows the sentence "Can you please come here?". A bracket under "Can you please come" is labeled "History" with an upward arrow. An upward arrow points to the word "here", which is labeled "Word being predicted".

The diagram shows the sentence "He was walking in the garden and..." with blue boxes representing the sequence of n-grams extracted: "He", "He was", "He was walking", "He was walking in", "He was walking in the", "He was walking in the garden", and "He was walking in the garden and".

# Padding sequences

---

- The maximum length of the sentence is extracted and then the rest of the sentences are pre-padded as per the longest sentence.

```
# pad sequences
max_sequence_len = max([len(x) for x in input_sequences])
input_sequences = np.array(pad_sequences(input_sequences,
                                         maxlen=max_sequence_len,
                                         padding='pre'))
```

```
print('max_len', max_sequence_len)
print('total-words', total_words)
```

```
max_len 12
total-words 672
```

# create features and Labels

- Extract the last word of sequence and convert it to categorical from numerical

```
# create predictors and label
predictors, label = input_sequences[:, :-1], input_sequences[:, -1]

label = to_categorical(label, num_classes=total_words)
```

```
print(input_sequences)
```

```
[[ 0  0  0 ...  0 265 1628]
 [ 0  0  0 ... 265 1628 101]
 [ 0  0  0 ... 1628 101 235]
 ...
 [ 0  0  0 ... 2018 36 1]
 [ 0  0  0 ... 36 1 6745]
 [ 0  0  0 ... 1 6745 66]]
```

X

Y



## 4) Bidirectional LSTM Model

---

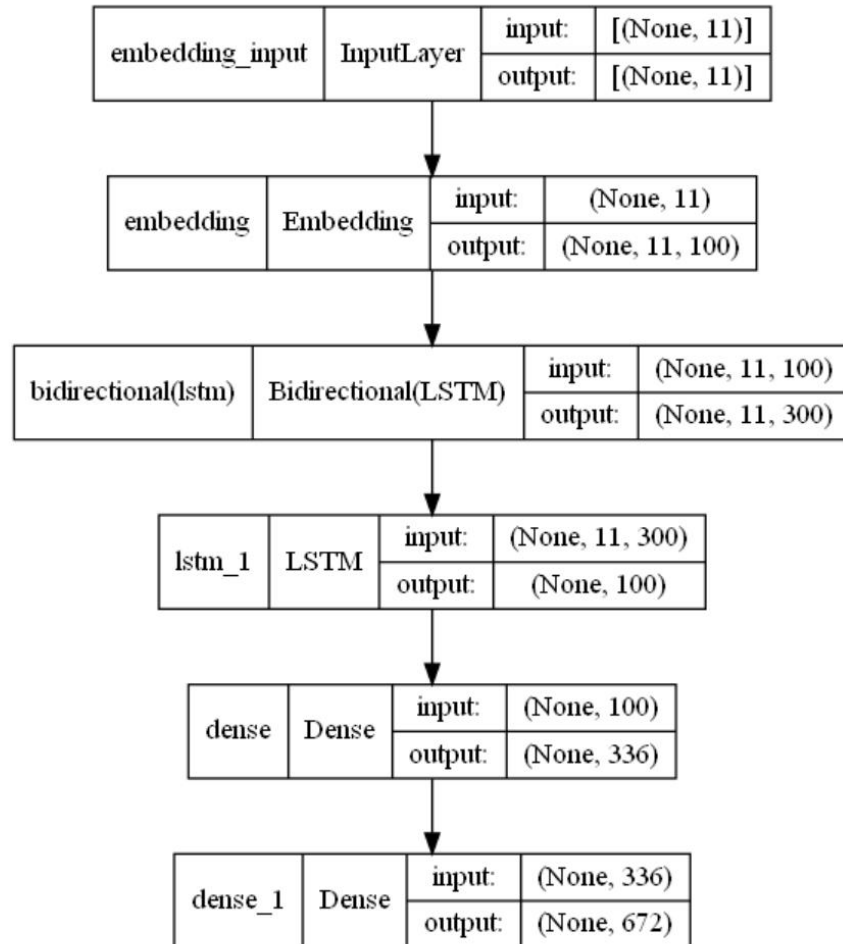
- Let's make a sequential model now with the first layer as the word embedding layer.
- 'return\_sequence' is marked as "True" so that the word generation keeps in consideration, previous and even the words coming ahead in the sequence.
  - ✓ The output layer has softmax to get the probability of the word to be predicted next.

```
model = Sequential()
model.add(Embedding(total_words, 100, input_length=max_sequence_len-1))
model.add(Bidirectional(LSTM(150, return_sequences = True)))
model.add(LSTM(100))
model.add(Dense(total_words/2, activation='relu' ))
model.add(Dense(total_words, activation='softmax'))
```

## 4) Bidirectional LSTM Model

```
plot_model(model, './fig_lab10_eng2.png', show_shapes=True)
```

```
]:
```



```
print(model.summary())
```

Model: "sequential\_3"

Layer (type)	Output Shape	Param #
=====		
embedding_3 (Embedding)	(None, 11, 100)	67200
bidirectional_1 (Bidirectional)	(None, 11, 300)	226800
gru_3 (GRU)	(None, 100)	120600
dense_4 (Dense)	(None, 336)	33936
dense_5 (Dense)	(None, 672)	226464
=====		
Total params: 675,000		
Trainable params: 675,000		
Non-trainable params: 0		
None		

## Model compile

---

```
model.compile(loss='categorical_crossentropy',  
              optimizer='adam',  
              metrics=['accuracy'])
```

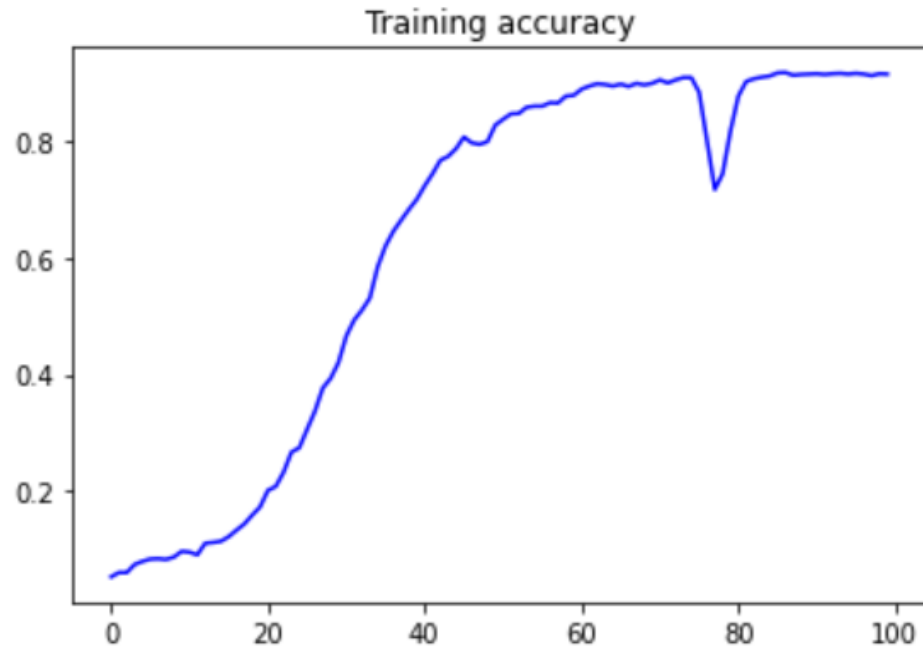
# Model.fit and Elapsed Times

---

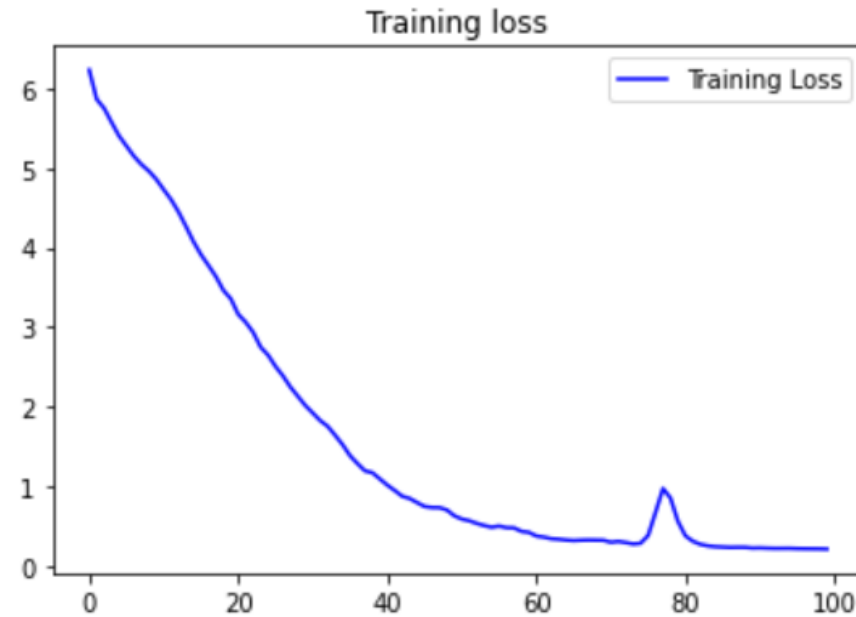
```
import time
start = time.perf_counter()
history = model.fit(predictors, label, epochs=100, verbose=1)
elapsed = time.perf_counter() - start
print('===== \n Elapsed %.3f seconds.' % elapsed)
```

```
Epoch 99/100
47/47 [=====] - 2s 48ms/step - loss: 0.2143 - accuracy: 0.9176
Epoch 100/100
47/47 [=====] - 2s 48ms/step - loss: 0.2117 - accuracy: 0.9169
=====
Elapsed 230.794 seconds.
```

# Training Accuracy and Loss



```
import matplotlib.pyplot as plt
acc = history.history['accuracy']
loss = history.history['loss']
epochs = range(len(acc))
plt.plot(epochs, acc, 'b', label='Training accuracy')
```



```
plt.figure()
plt.plot(epochs, loss, 'b', label='Training Loss')
plt.title('Training loss')
plt.legend()
plt.show()
```

## 5) Test for the next generation

---

➤ 100 next words are generated this way

- ✓ The seed will be taken at first and tokenized and padded on the token list.
- ✓ Model is then used to predict with the token list as input.
- ✓ Then most probable word is added to seed text and this happens for the next 100 words.

```
seed_text = "First Servingman: A strange one as ever I looked on: I cannot get him!"  
next_words = 100
```

## 5) Test for the next generation

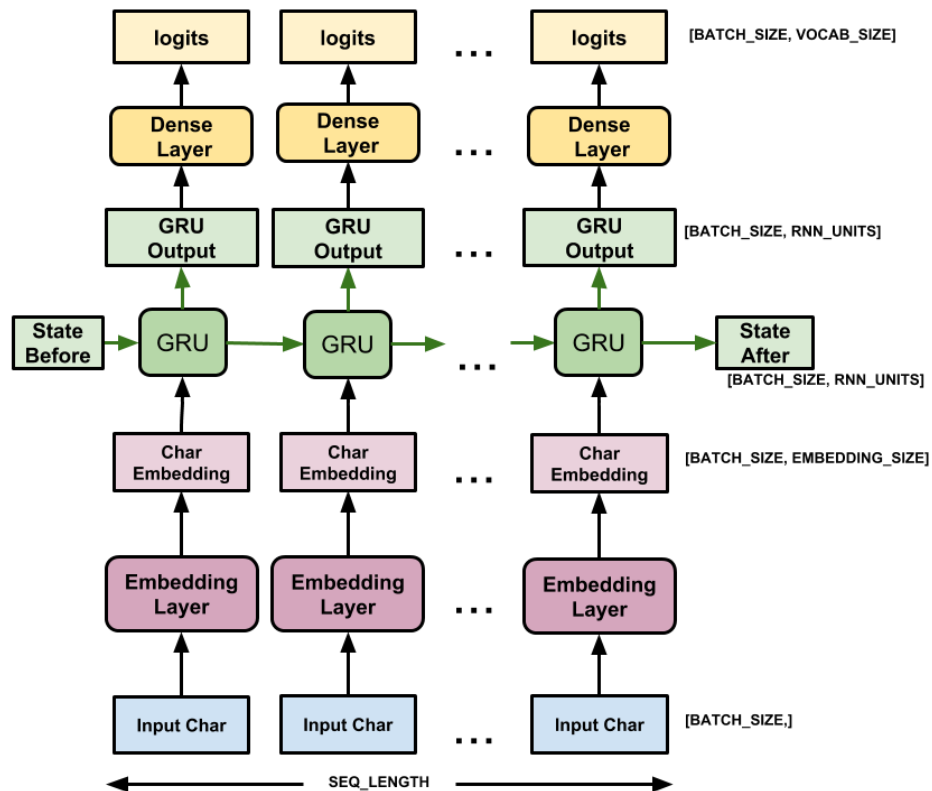
```
for _ in range(next_words):
    token_list = tokenizer.texts_to_sequences([seed_text])[0]
    token_list = pad_sequences([token_list], maxlen=max_sequence_len-1, padding='pre')
    p_x = model.predict(token_list, verbose=0)
    predicted=np.argmax(p_x,axis=1)

    output_word = ""
    for word, index in tokenizer.word_index.items():
        if index == predicted:
            output_word = word
            break
    seed_text += " " + output_word
print(seed_text)
```

First Servingman: A strange one as ever I looked on: I cannot get him! in pray  
you account remember you gods you you barren of all the rest were so mark  
me say he fathers fathers fathers fathers fathers fathers did venture pray you s  
ay there's which you to't own belly of man man which he remember remembe  
r please please please can which our own price please ye caps please ye price  
ye caps please ye caps please please ye trust ye 't please please ye 't please p  
lease ye 't please please ye 't please please ye 't please please ye 't please ple  
ase ye 't please please ye 't please please ye

# HW : Optimized Model Architecture

- The output is not perfect as for training we took only a few lines of text.
  - ✓ Hence we can very well fine-tune it.





# New GRU Model

---

```
model = Sequential()
model.add(Embedding(total_words, 100, input_length=max_sequence_len-1))
model.add(Bidirectional(GRU(150, return_sequences = True)))
model.add(GRU(100))
model.add(Dense(total_words/2, activation='relu' ))
model.add(Dense(total_words, activation='softmax'))
```

```
model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
```

```
import time
start = time.perf_counter()
history = model.fit(predictors, label, epochs=100, verbose=1)
elapsed = time.perf_counter() - start
print('=====\n Elapsed %.3f seconds.' % elapsed)
```

# Code in Colab

Thanks