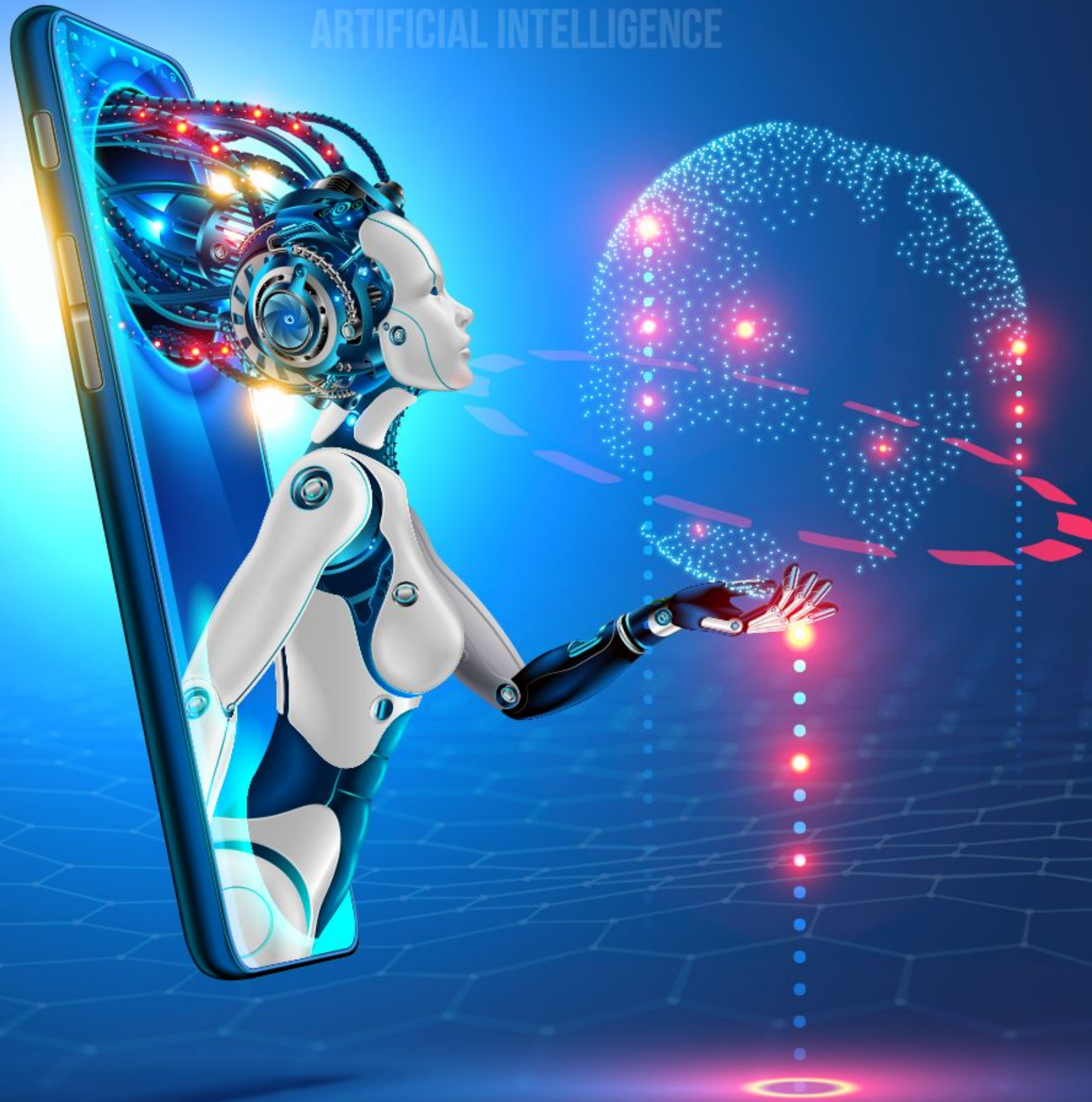


# DATA AND ARTIFICIAL INTELLIGENCE



simplilearn

PURDUE  
UNIVERSITY

## Natural Language Processing



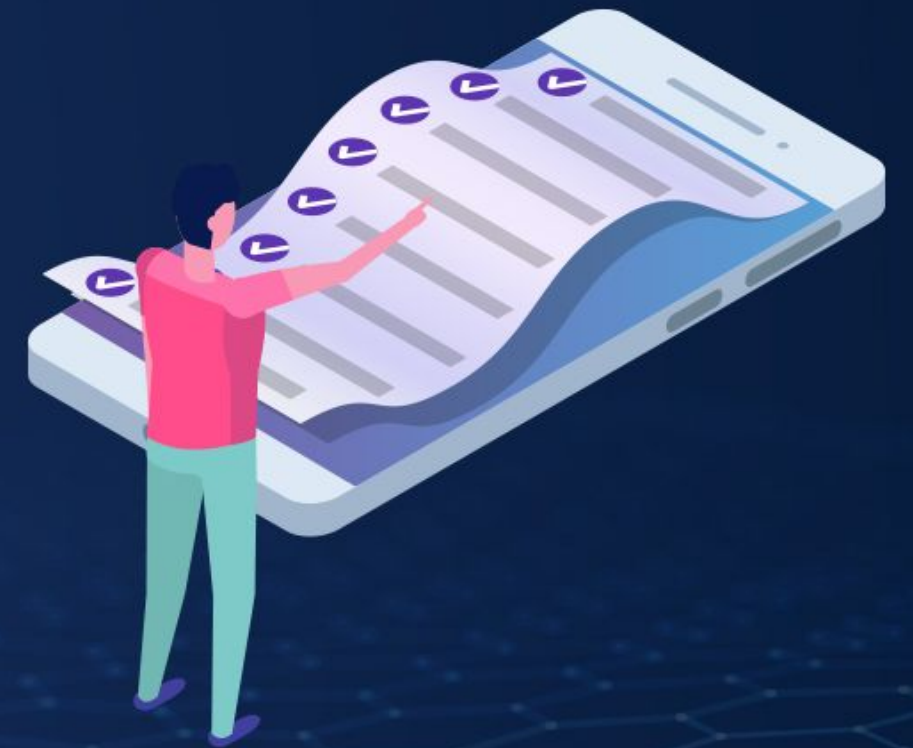
## Introduction to Natural Language Processing



# Learning Objectives

By the end of this lesson, you will be able to:

- 👁 Describe natural language processing and its components
- 👁 Explain the different applications of NLP
- 👁 Define and demonstrate text processing



## Introduction to NLP

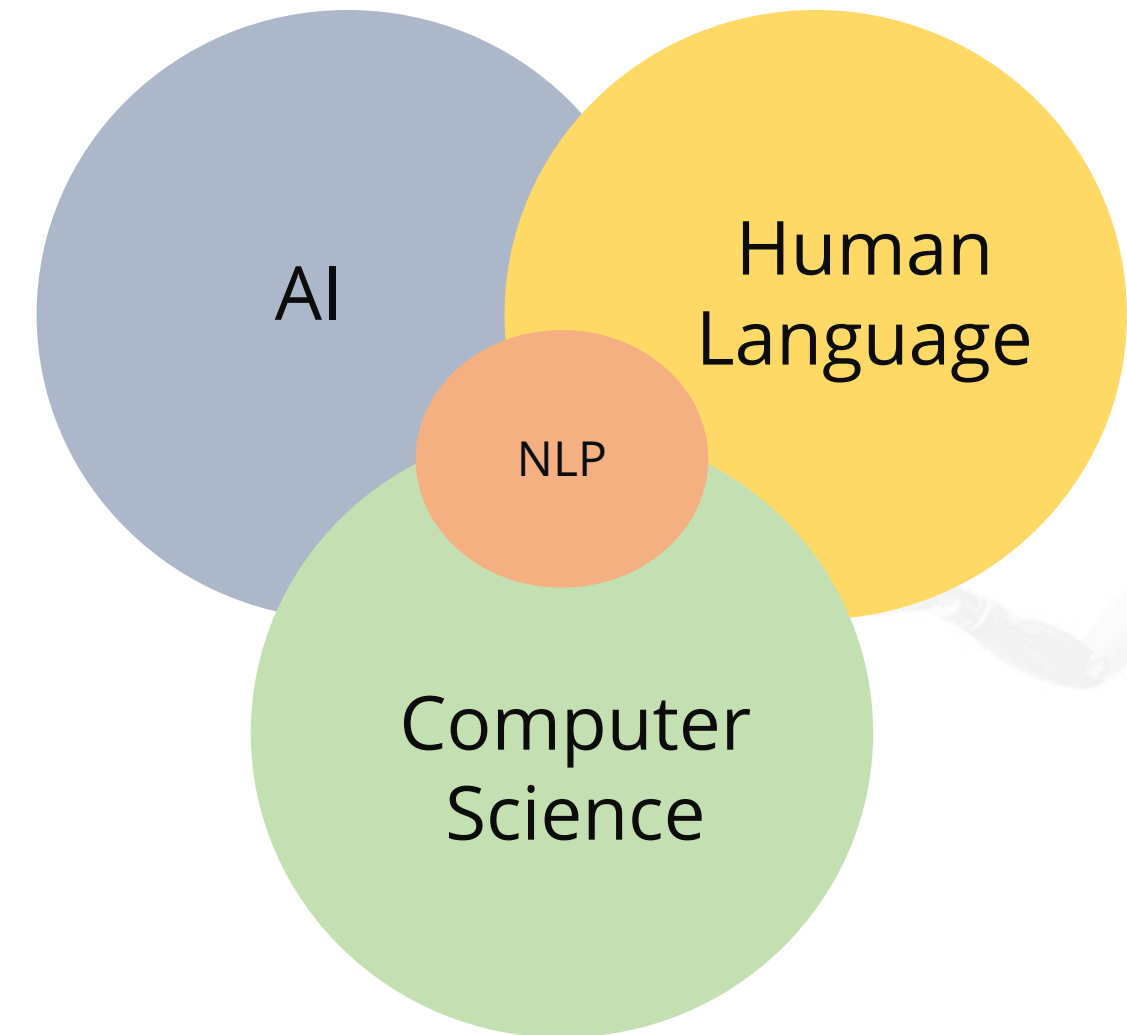
# What Is NLP?

Natural Language Processing (NLP) is a branch of AI.

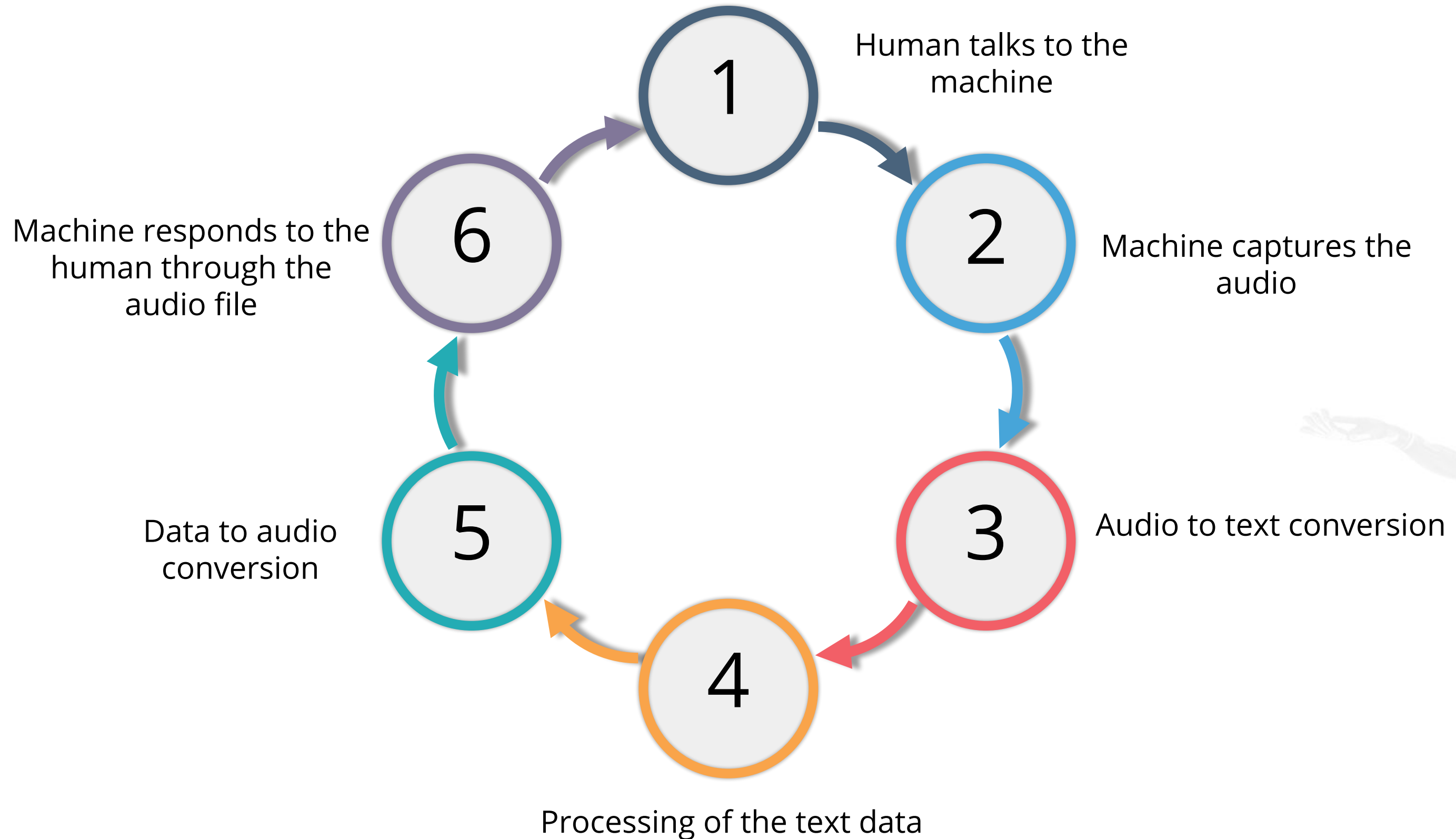
It helps machine to deal with human languages.

It helps machine to understand, interpret, and manipulate human languages.

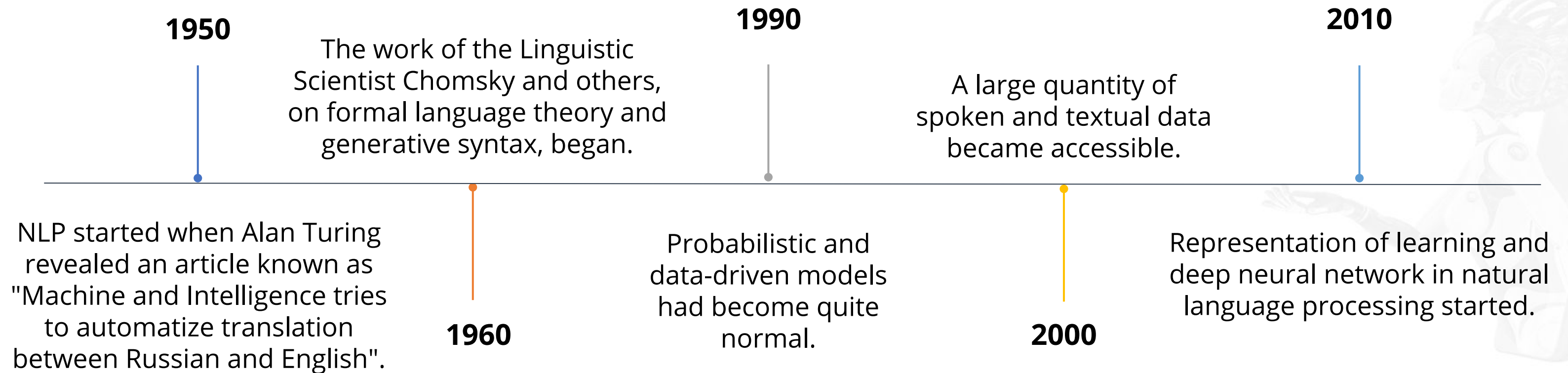
Most of the Natural Language Processing techniques depend on machine learning to derive meaning from human languages.



# Interaction between Humans and Machines Using NLP

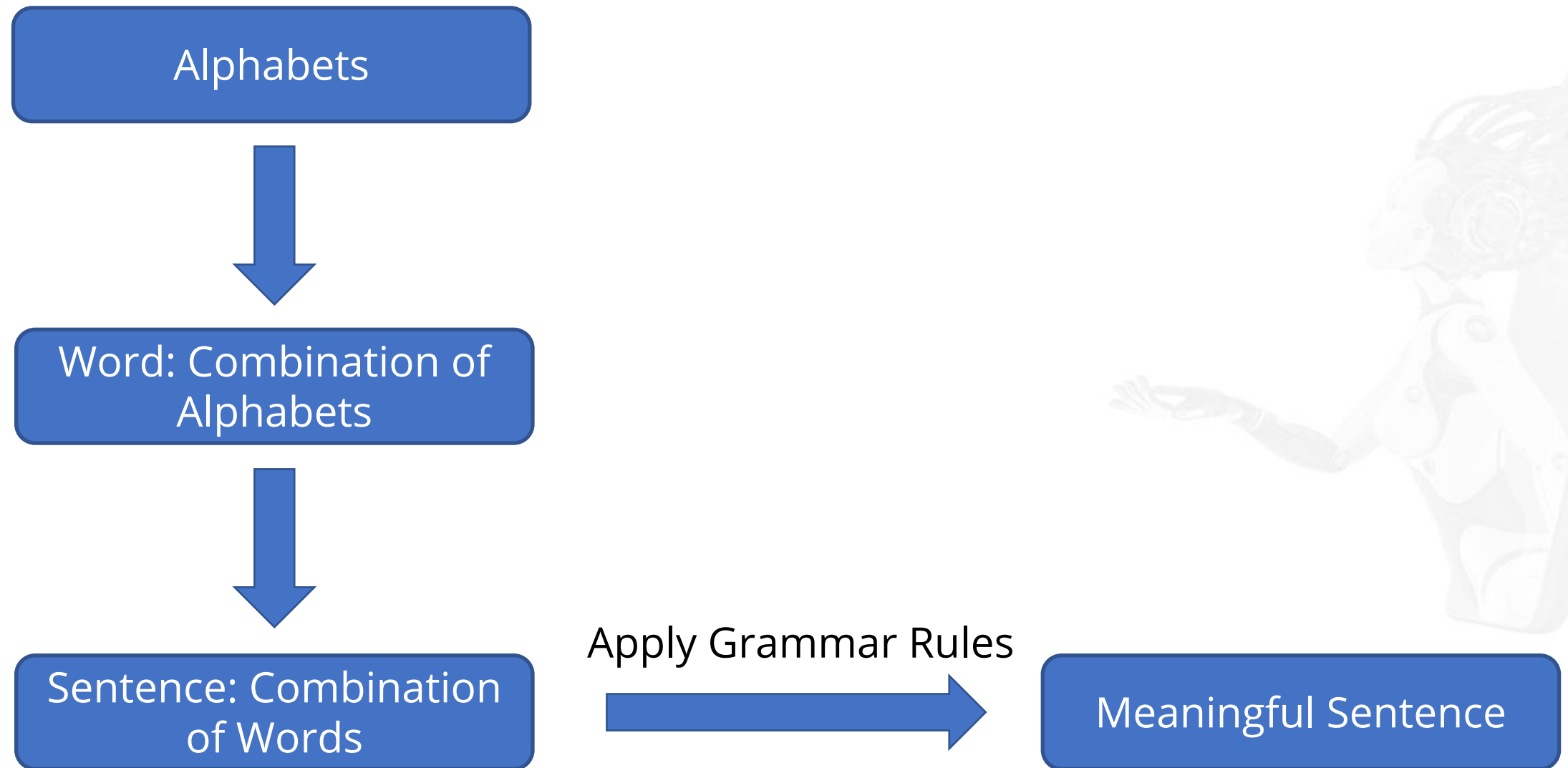


# History of NLP



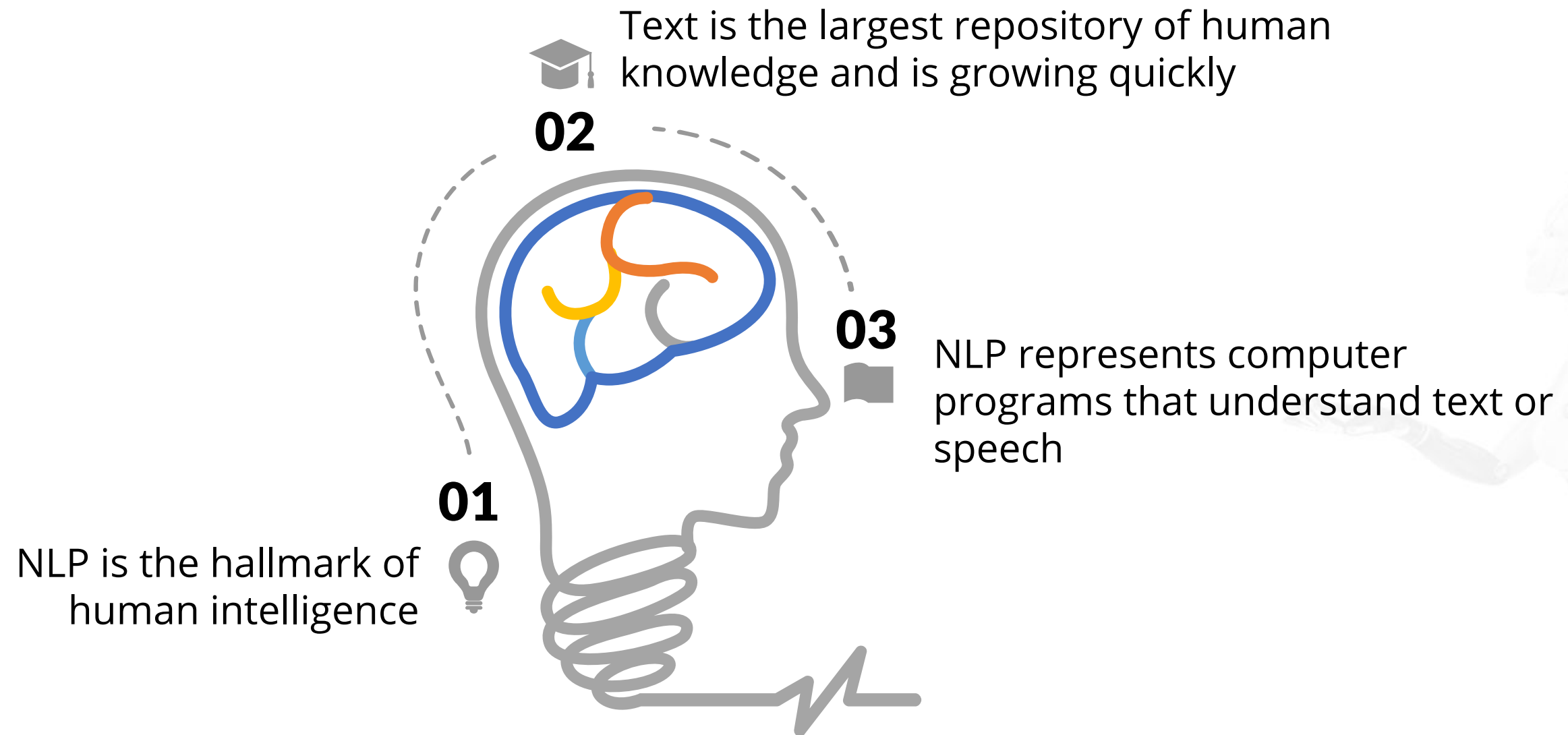
# Human Language

To understand NLP,  
let us first  
understand the  
human language.





# Why NLP



# Need for NLP

Collections  
Organization Words  
Text **Unstructured** Packaged  
Facilitate **Data** Big Data  
Patterns  
Originates  
none



Data  
Sources



How to analyze this unstructured data?  
Use Text Mining

# Understanding Text Mining

It is also called text analysis.

Process of deriving insights from natural language text



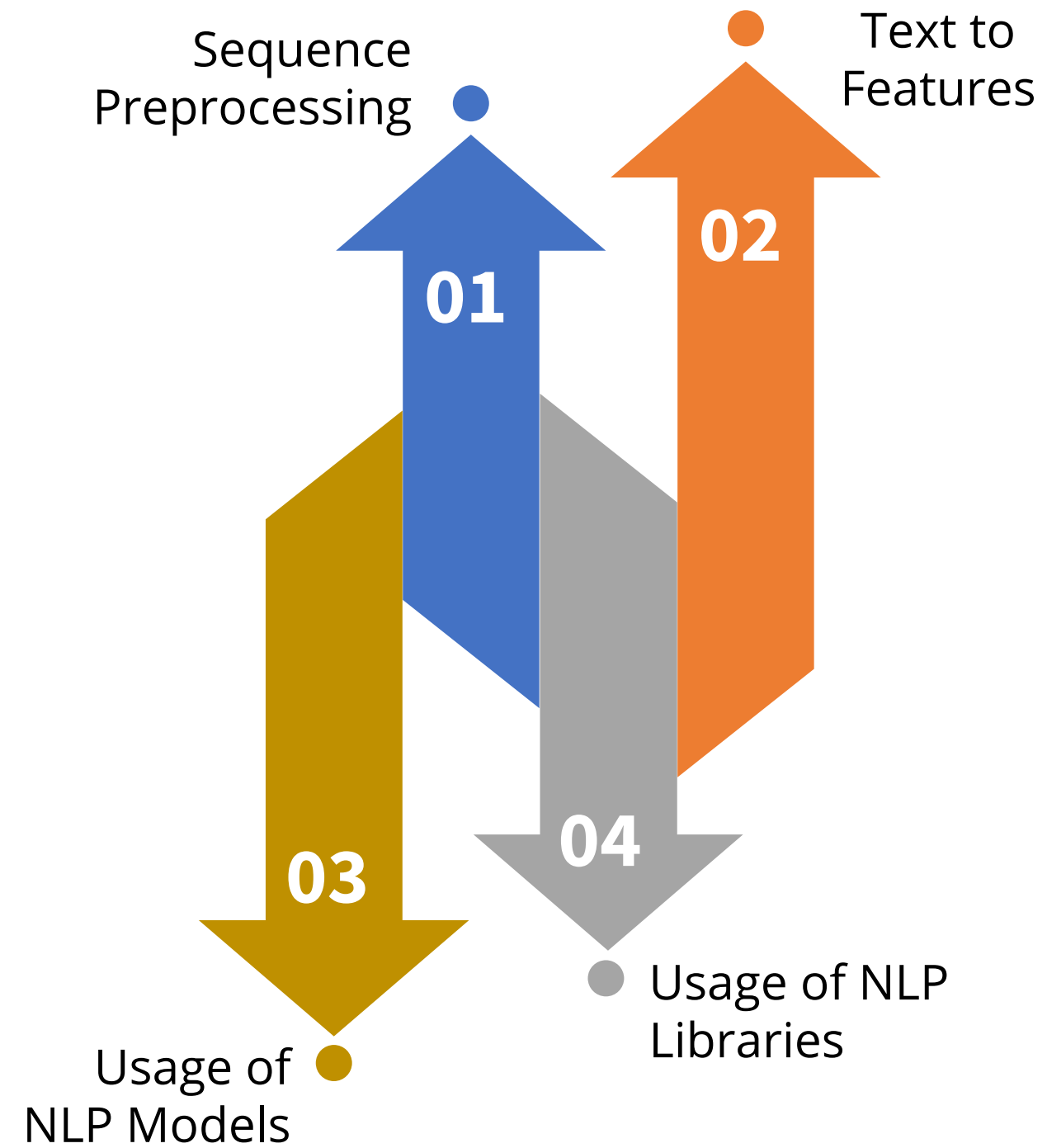
This is  
where NLP  
helps

Structure the input  
text

Derive pattern

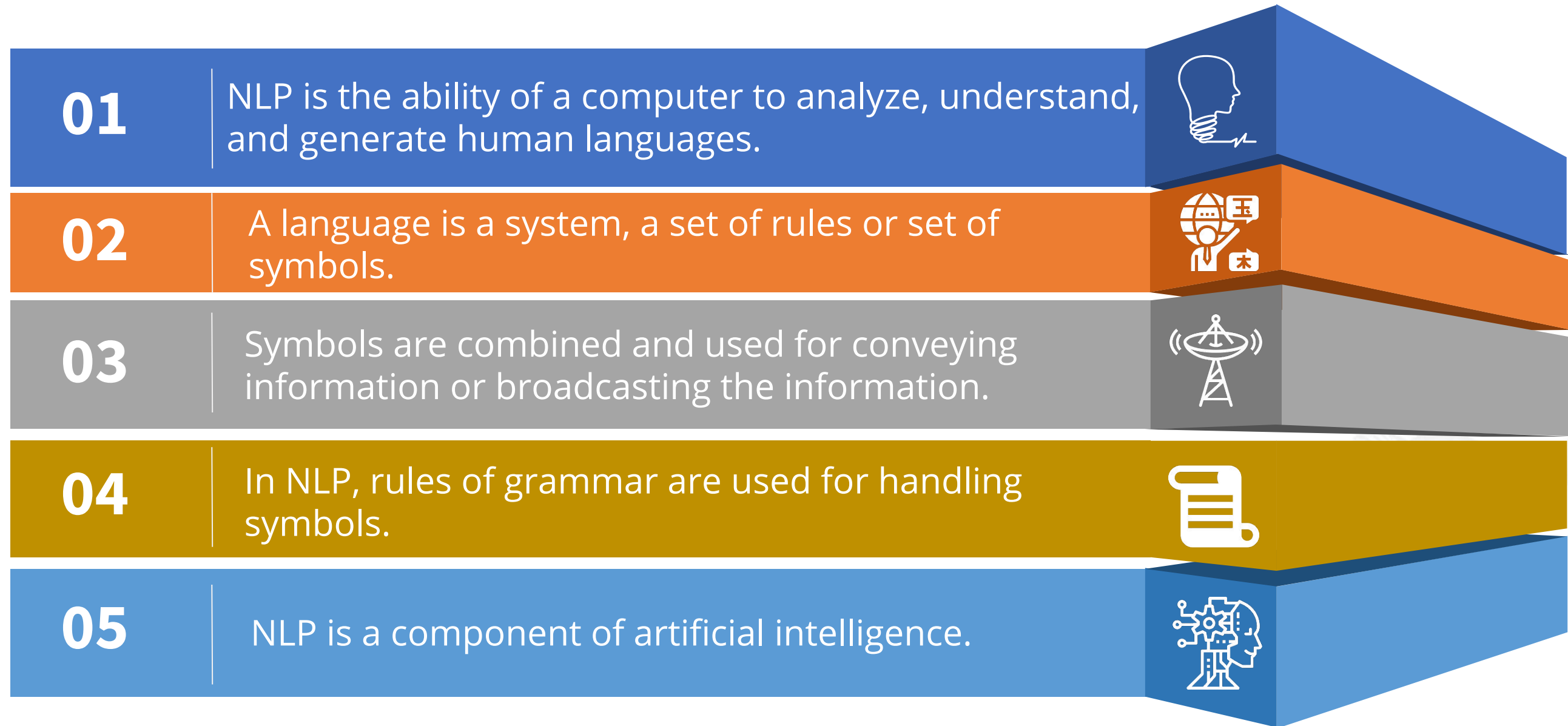
Evaluate output

# How NLP Works





# Different Aspects of NLP



# Categories of NLP

## Rule-Based NLP

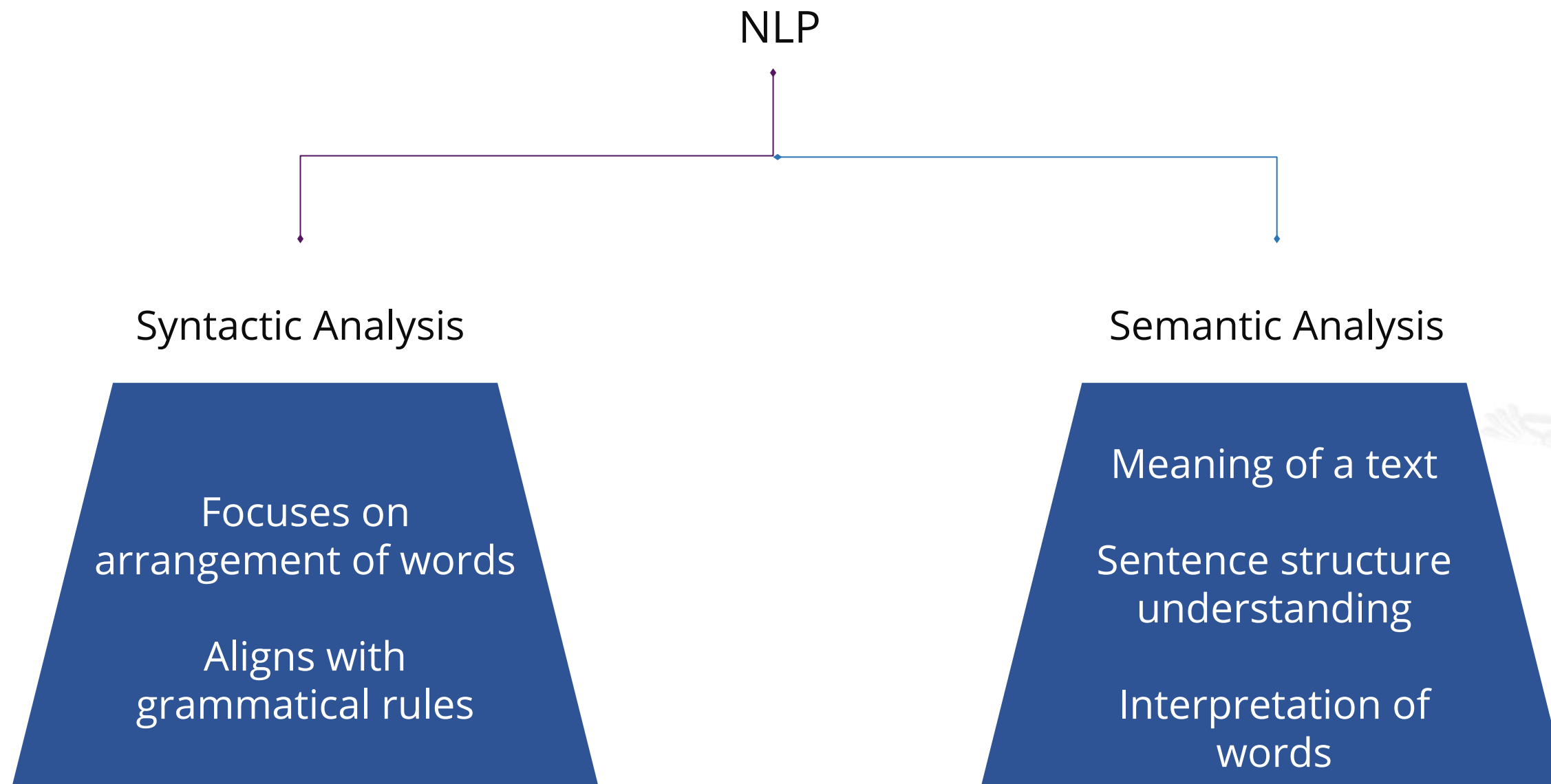
- Designed by creating a set of rules
- Developed by heuristic rules

Statistical Revolution

## Statistical NLP

- Relies heavily on machine learning
- Applies automatic learning procedure

# Techniques Used in NLP



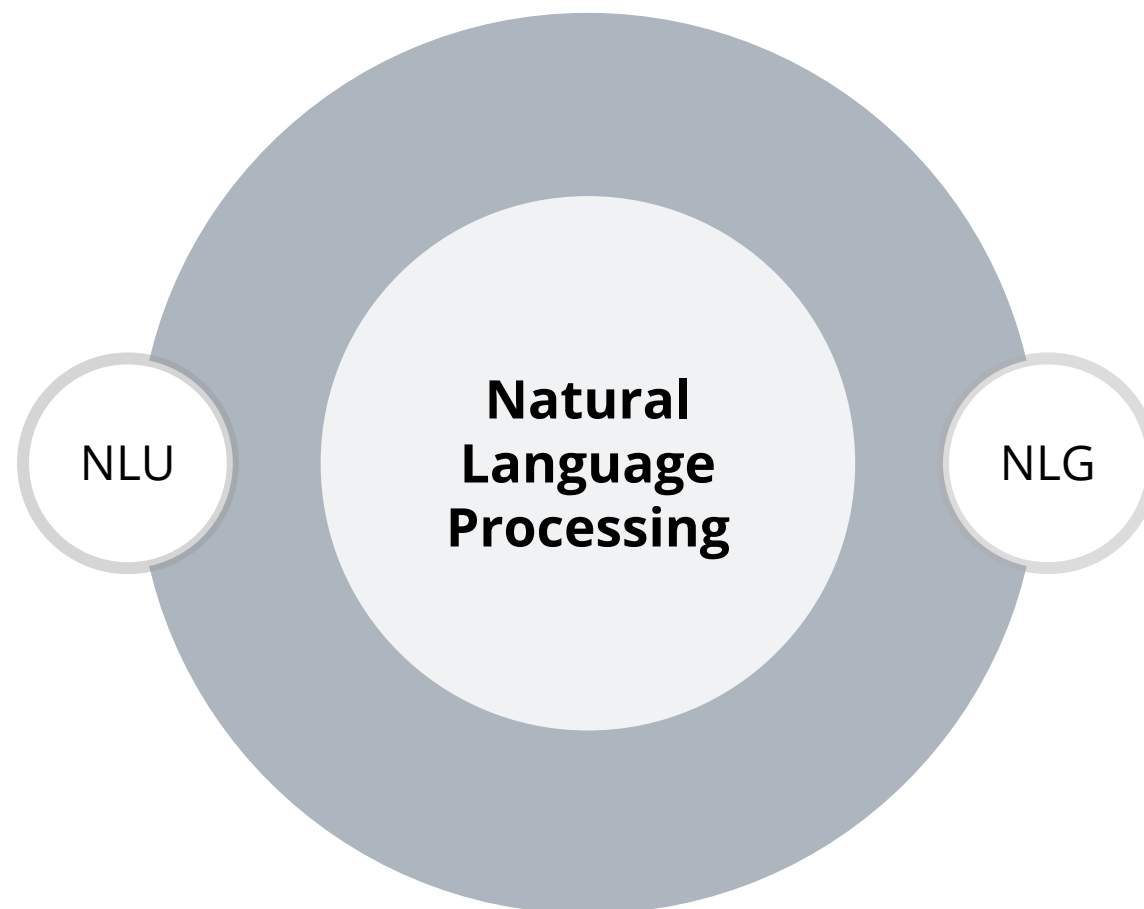
## **Components of Natural Language Processing**



# Components of Natural Language Processing

1

**Natural Language Understanding (NLU)**



**Natural Language Generation (NLG)**



2

# Components: Natural Language Understanding (NLU)



Taking some sentences and  
finding out what they mean

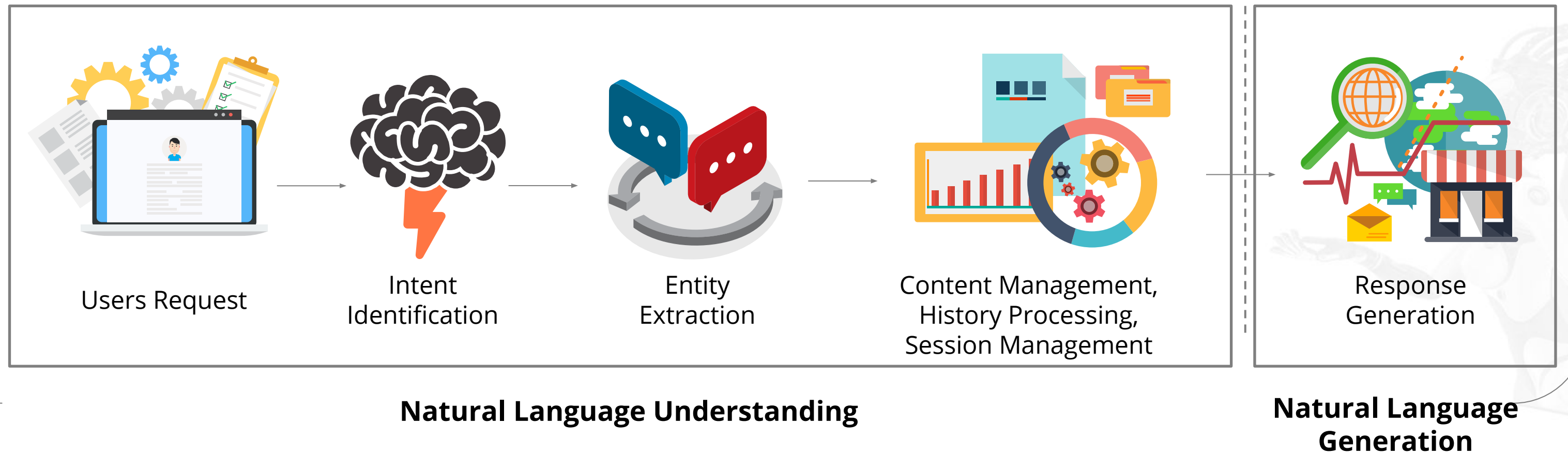


## Components: Natural Language Generation (NLG)

- 1 Taking some formal representation of what you want to say and working out a way to express it in a natural language
- 2 Mapping the given input in the natural language with a useful representation
- 3 Producing output in the natural language from some internal representation
- 4 Different level of analysis: morphological analysis, syntactic analysis, semantic analysis, and discourse analysis

# Uses of NLP

Use of NLP in conversational bot in each step:



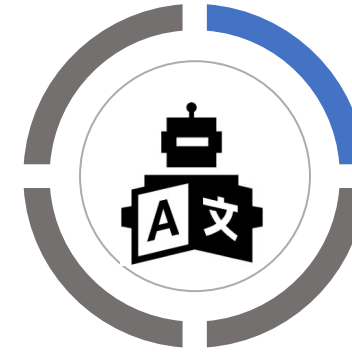


## **Applications of Natural Language Processing**

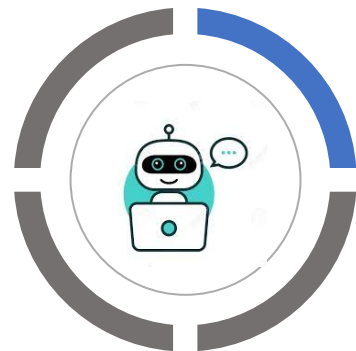
# NLP in Real-Life



**Speech  
Recognition**



**Machine  
Translation**



**Chatbot**



**Information  
Retrieval**



# NLP in Real Life



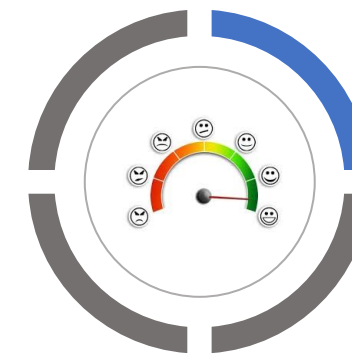
**Information  
Extraction**



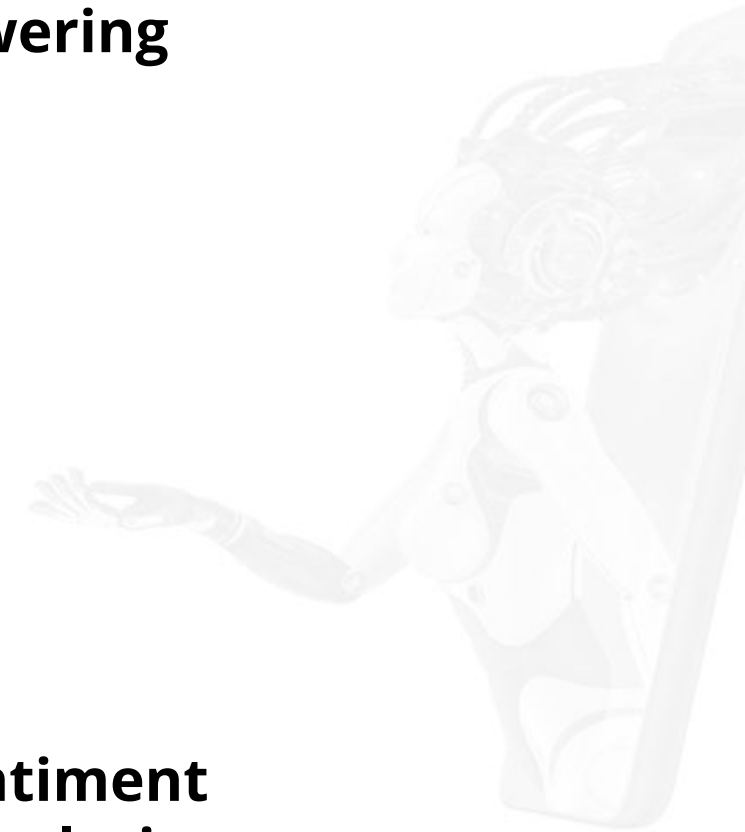
**Question  
Answering**



**Spell Check**



**Sentiment  
Analysis**



# NLP in Real-Life: Business Usage

## Improve user experience

- Spellcheck
- Autocomplete
- Autocorrect

## Automate support

- Chatbot
- Product ordering

## Monitor and analyze feedback

- Generate actionable insight from huge amount of review or feedback

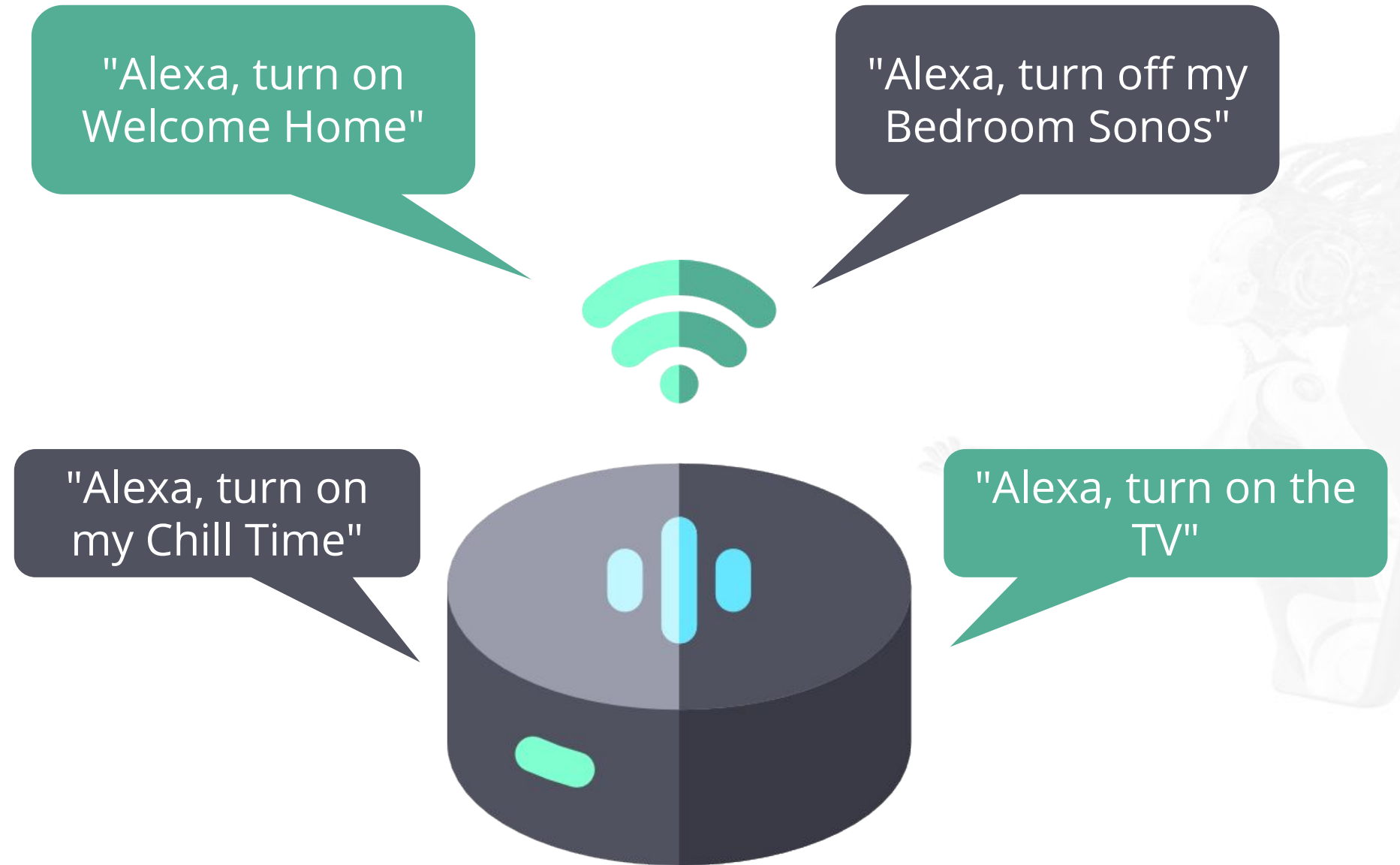




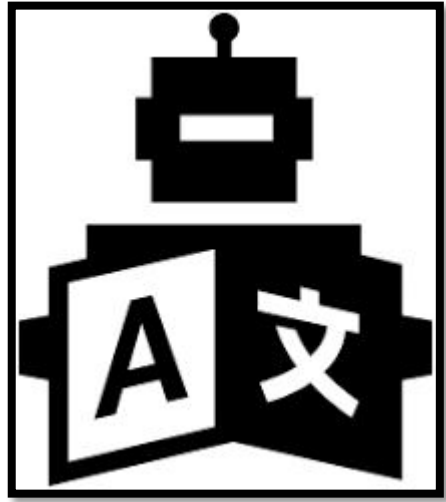
# NLP in Real-Life: Speech Recognition



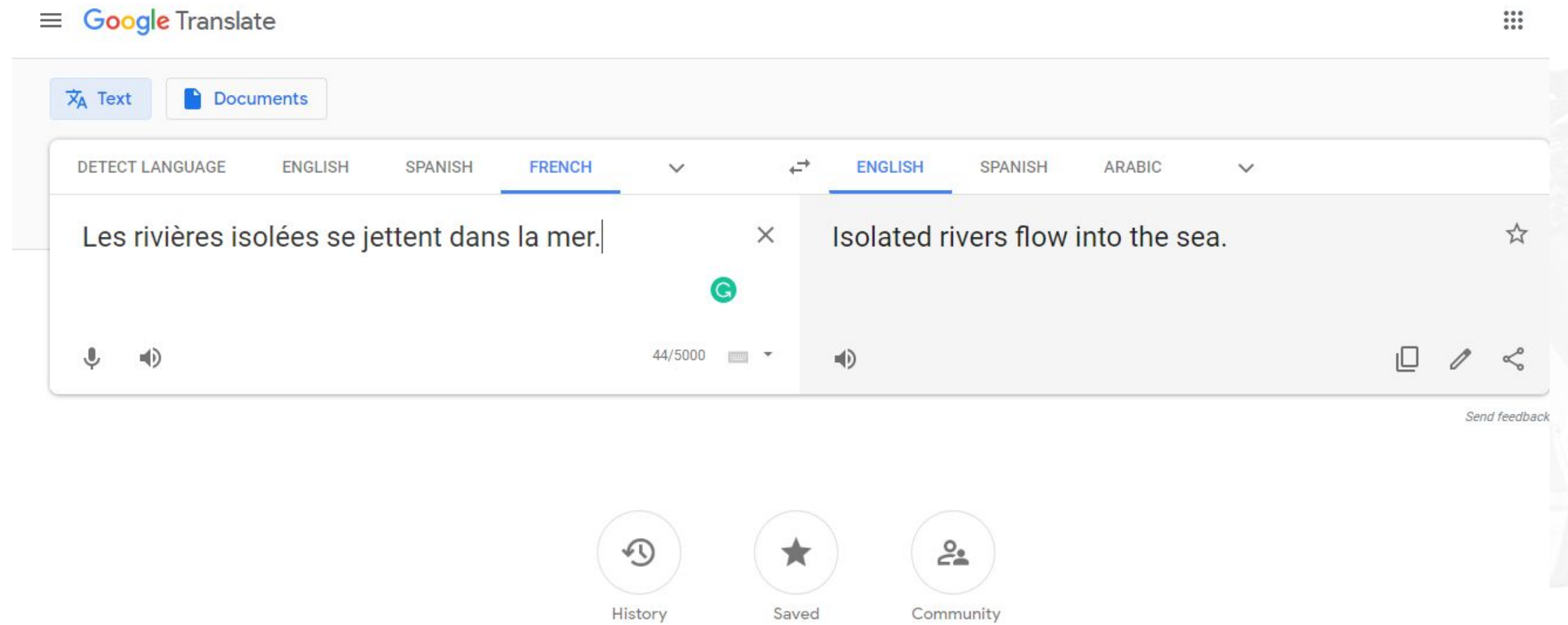
- Google Assistant
- Siri
- Alexa
- Cortana



# NLP in Real-Life: Machine Translation

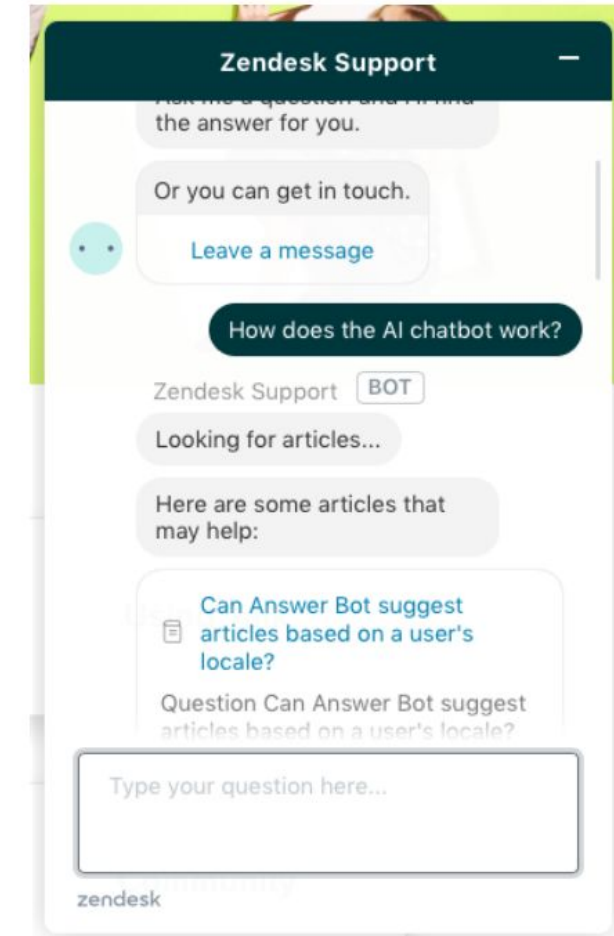
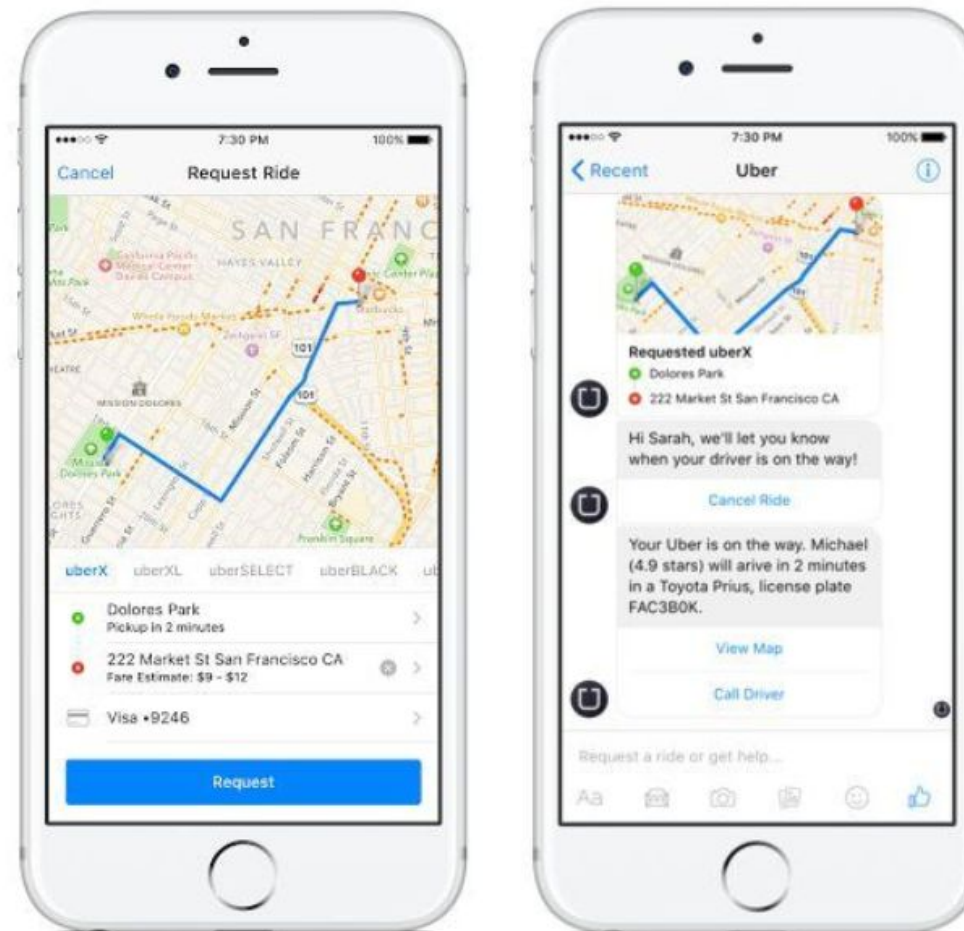
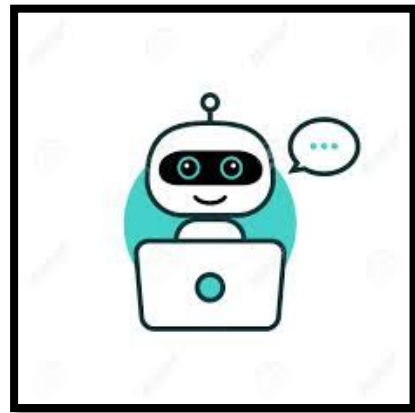


Google translator



# NLP in Real-Life: Chatbots

Uber, Facebook Messenger, and Zendesk are some of the companies who have implemented chatbots using NLP.



Source

# NLP in Real-Life: Information Retrieval

Find information according to the given query



Collections Audio Video  
Organization Words Packaged Involve  
Sentences Text **Data** Big Data  
Facilitate Unstructured Patterns  
none Originates

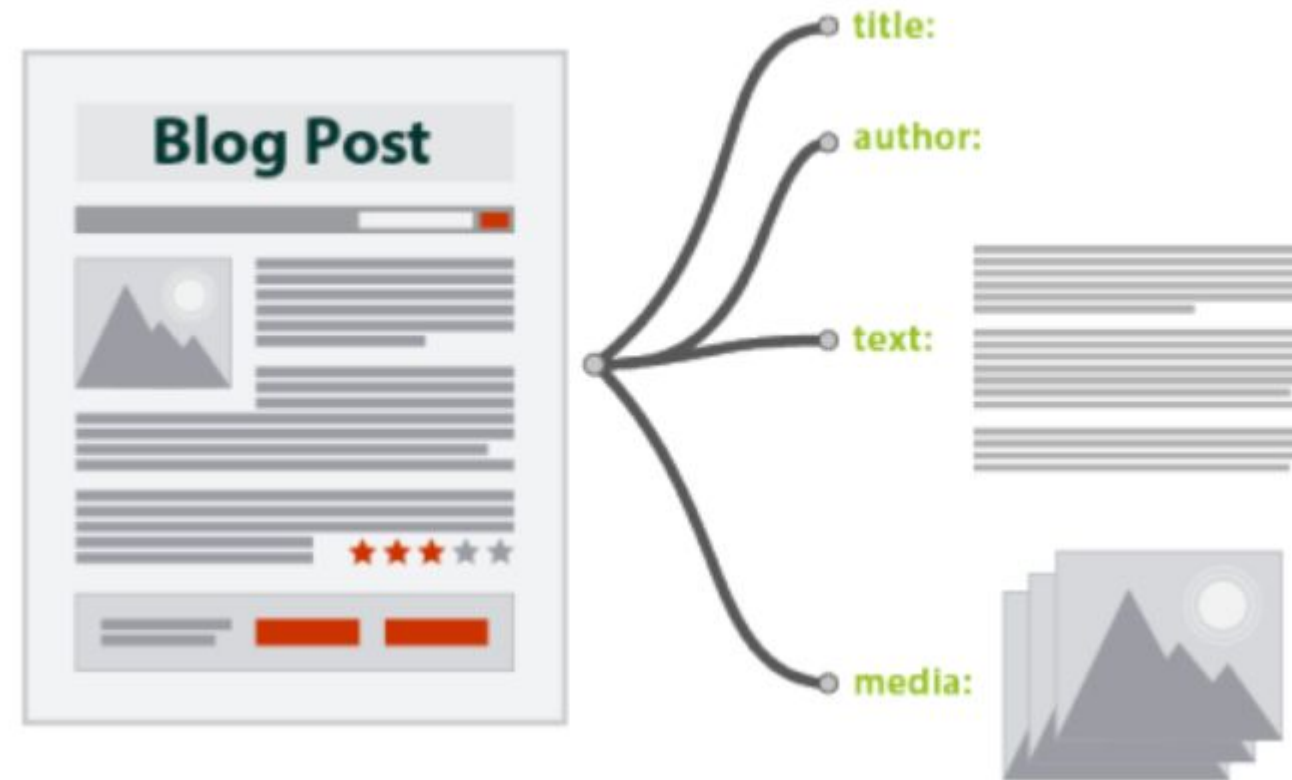
NLP techniques used in IR are:

- Stemming
- Part-of-Speech Tagging
- Compound Recognition
- Decompounding
- Chunking
- Word-Sense Disambiguation

Google finds relevant and similar results using Information Retrieval.

# NLP in Real-Life: Information Extraction

Automatic extraction of structured information from unstructured or semi-structured machine-readable documents

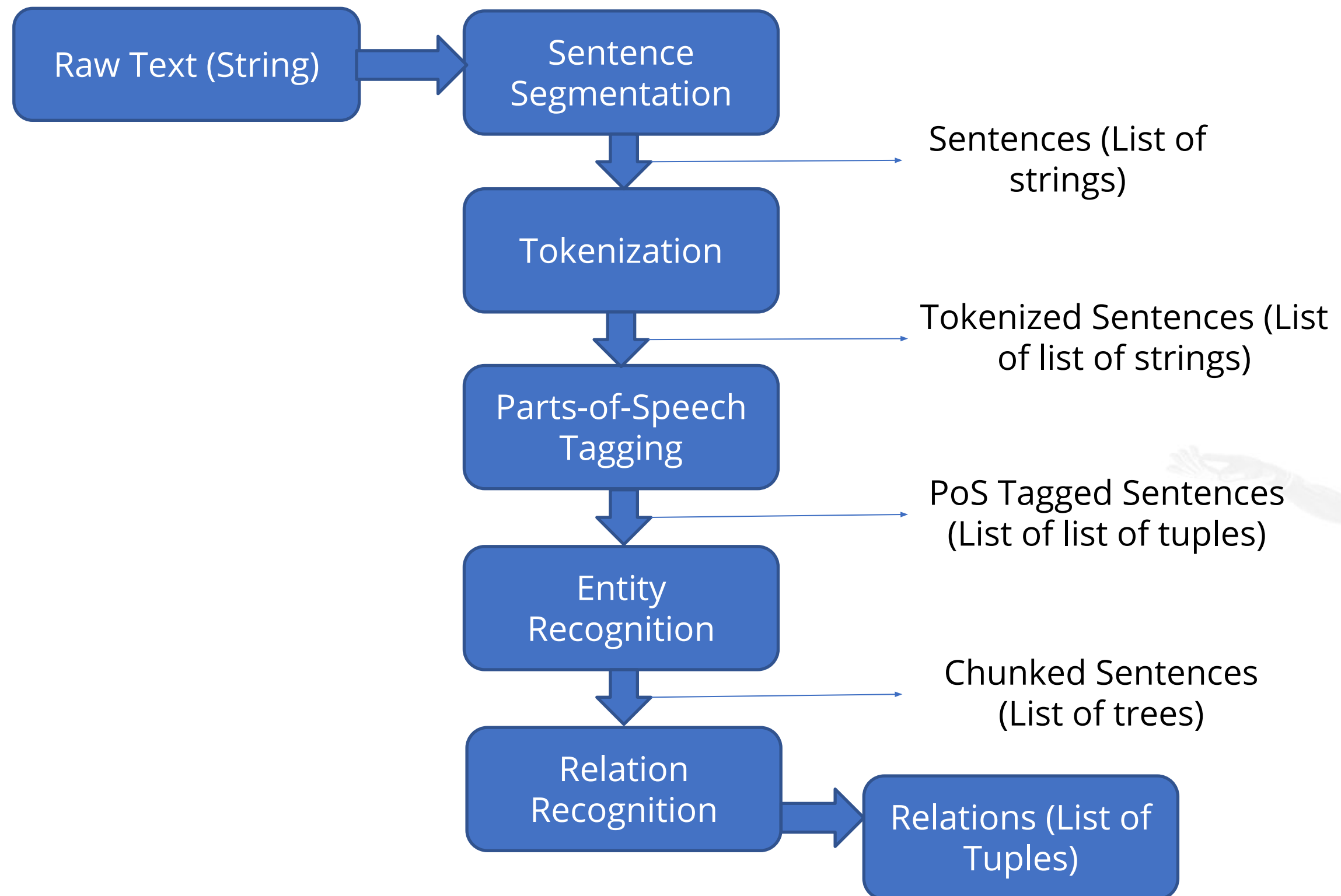


Gmail structures events from e-mails

Source

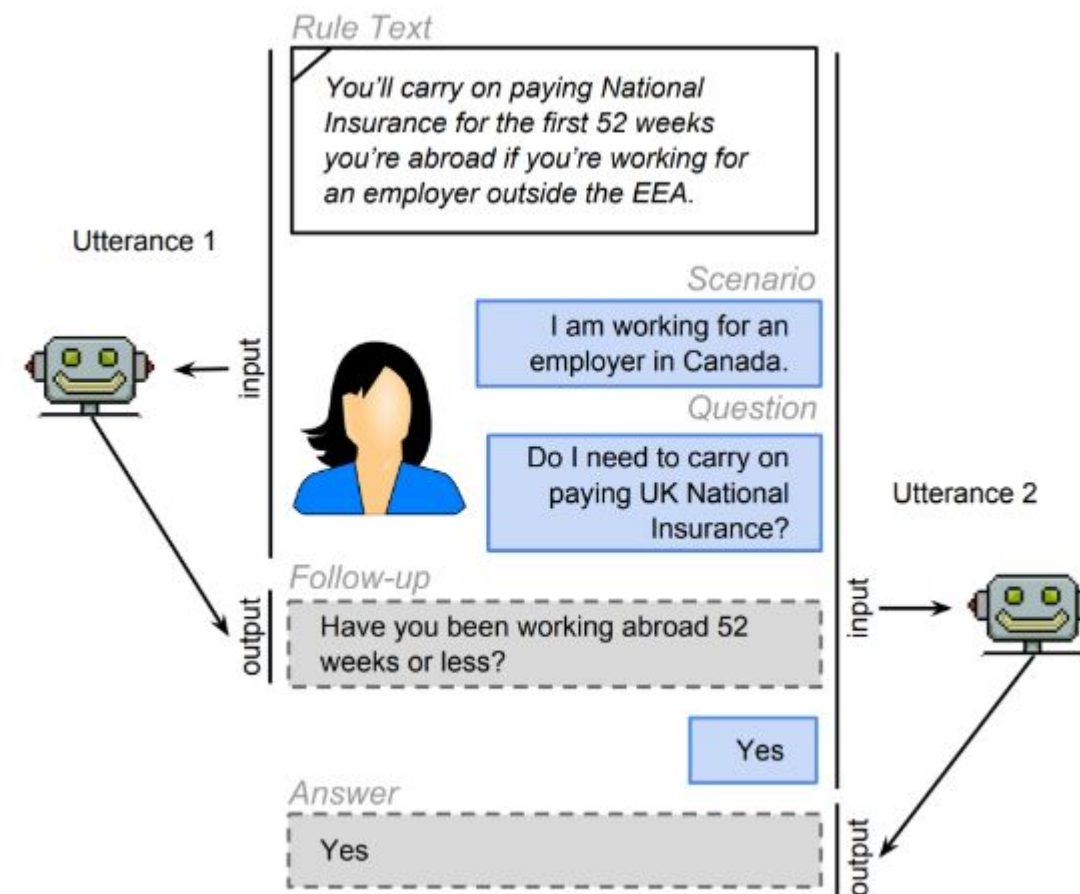


# NLP in Real-Life: Information Extraction



# NLP in Real-Life: Question Answering

System that automatically answers questions



Source

# NLP in Real-Life: Spell Check

Salesforce implemented spell check in the contact forms using NLP.



**salesforce** CONNECT TO YOUR CUSTOMERS IN A WHOLE NEW WAY

**Support** Get help from our technical team >

**Community** Get answers from other customers >

Hellow  
Hello x

994

Easy ☐ ☐ ☐ ☐ ☐

Effective ☐ ☐ ☐ ☐ ☐

Enjoyable ☐ ☐ ☐ ☐ ☐

Overall Page Rating ☐ ☐ ☐ ☐ ☐

Email address for follow-up (optional):

What was the purpose of your visit today? Please choose one...

Were you able to find what you were looking for today? Please choose one...

Would you recommend Salesforce to a friend or colleague? (0 = Very unlikely; 10 = Very likely)

0 1 2 3 4 5 6 7 8 9 10

**opinionlab** A VERINT Company Privacy Policy About OpinionLab  
© OpinionLab, Inc. All rights reserved; Patented

Submit



Grammarly: A grammar checking SW built on NLP

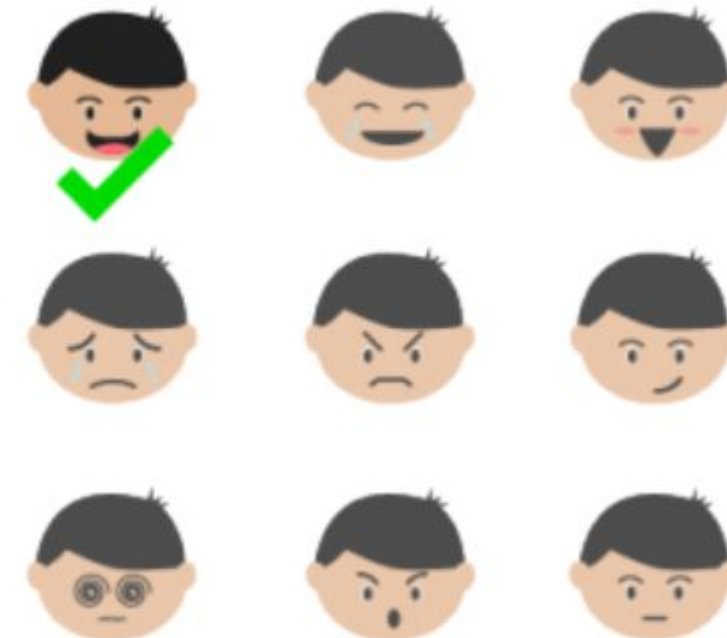
Source

# NLP in Real-Life: Sentiment Analysis

To extract subjective information from a piece of text  
Example: Whether an author is being subjective or objective or even positive or negative



NLP is used here



Source

## Challenges and Scope

# Why NLP Is Difficult

Nature of the human language



Rules that dictate the passing of information using natural languages are not easy for computers to understand.

Human language is unstructured data



Only 21% of the data is structured data, and a lot of information in the world is unstructured.

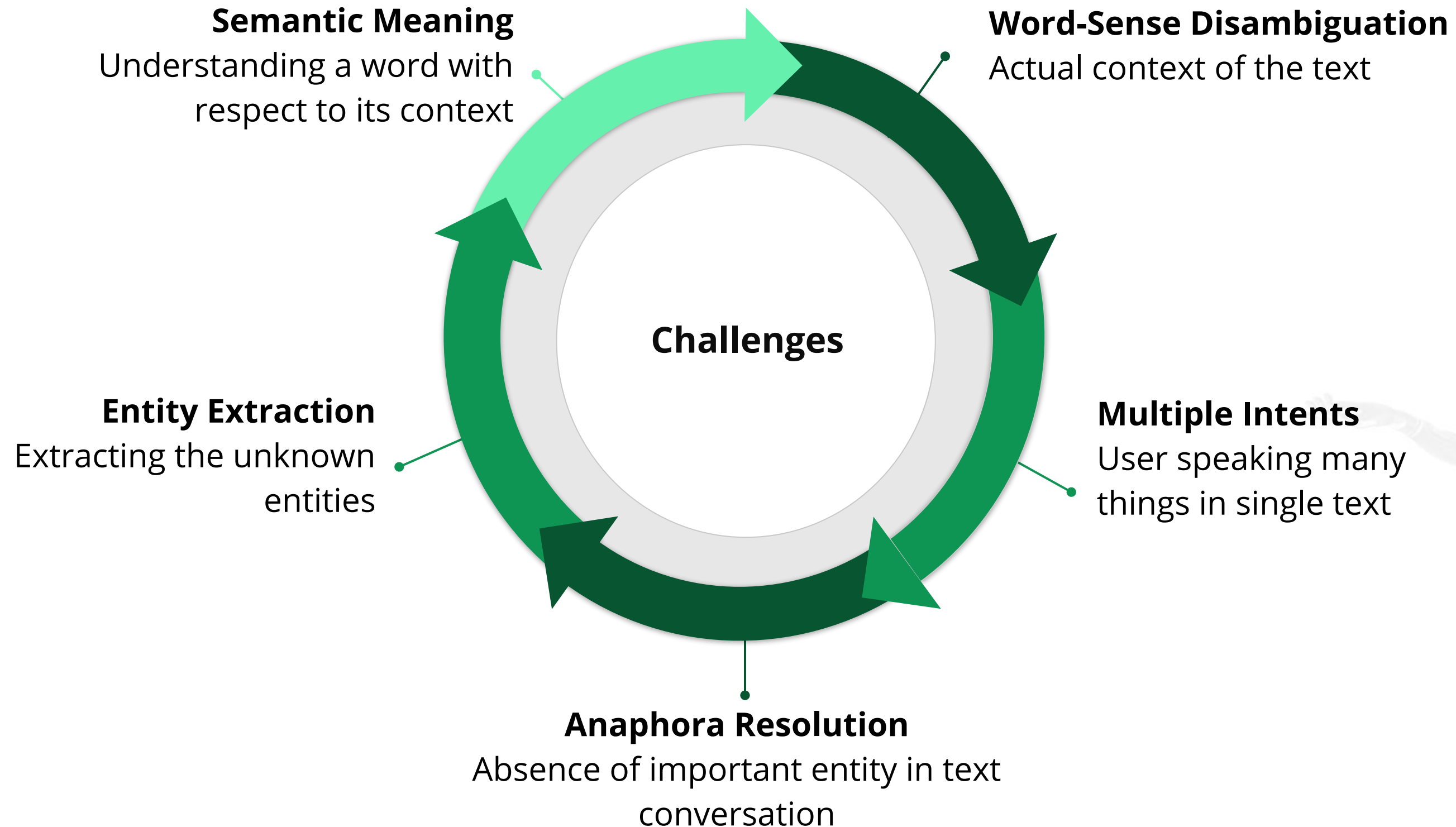
Tough to extract meaning from text



Process of reading and understanding English is very complex.



# Challenges and Scope



## Challenges and Scope: Semantic Meaning

There are many good properties available on HDFC Red portal.

I want to purchase a red carpet from a store.

Word RED has different meanings in these contexts.

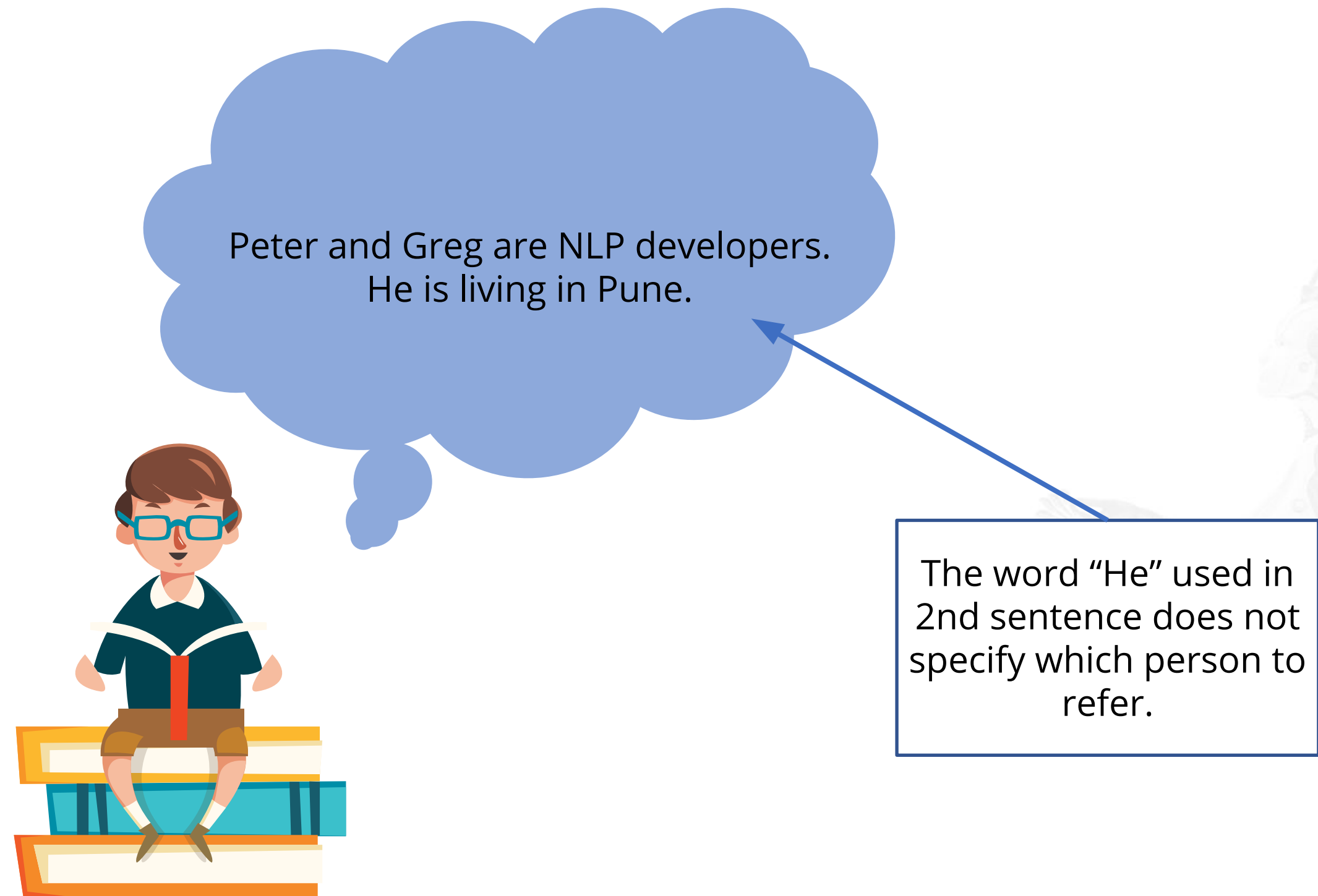
# Challenges and Scope: Understanding Entities

A 2M solution of  $\text{CaCl}_2$  consists of 221.82g of  $\text{CaCl}_2$  dissolved in enough water to make one liter of solution.

Understanding and extraction of  $\text{CaCl}_2$  as entity in this context is complex.



# Challenges and Scope: Anaphora Resolution



## Challenges and Scope: Multiple Intents

My bank account is functional. Please provide me resolution process and I want to buy Laptop from Flipkart.

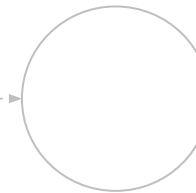
The word "He" used in 2nd sentence does not specify which person to refer.



# NLU Challenges: Ambiguity



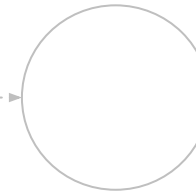
Lexical  
Ambiguity



**More than 1 meaning of a word in a sentence**

Example: The fisherman went to the **bank**.

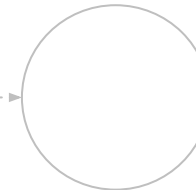
Syntactic or  
Grammatical  
Ambiguity



**More than 1 meaning of a sentence**

Example : **Visiting relatives** can be boring.

Referential  
Ambiguity



**Reference of a pronoun**

The boy told the father about the theft. **He** was very upset.

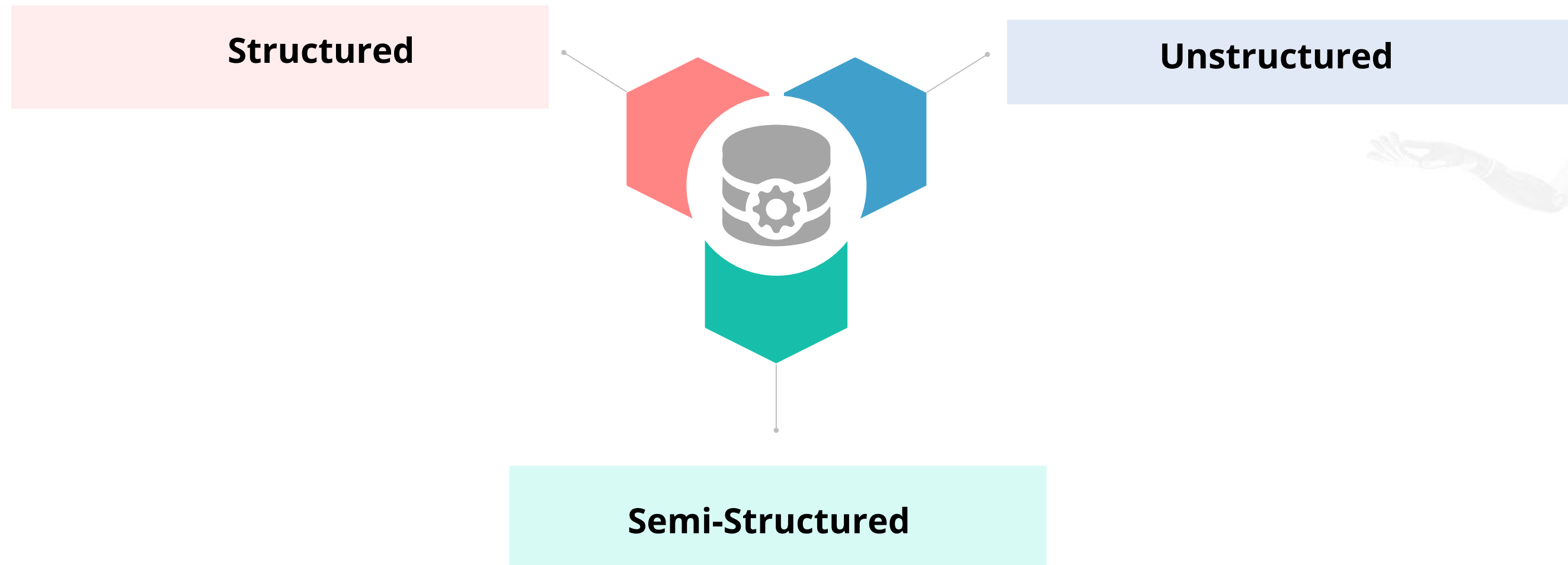


## Data Formats

# Data Formats

To apply NLP on data, we need to have the data which is available on different kinds of sources in different formats.

Below are the types of data formats:



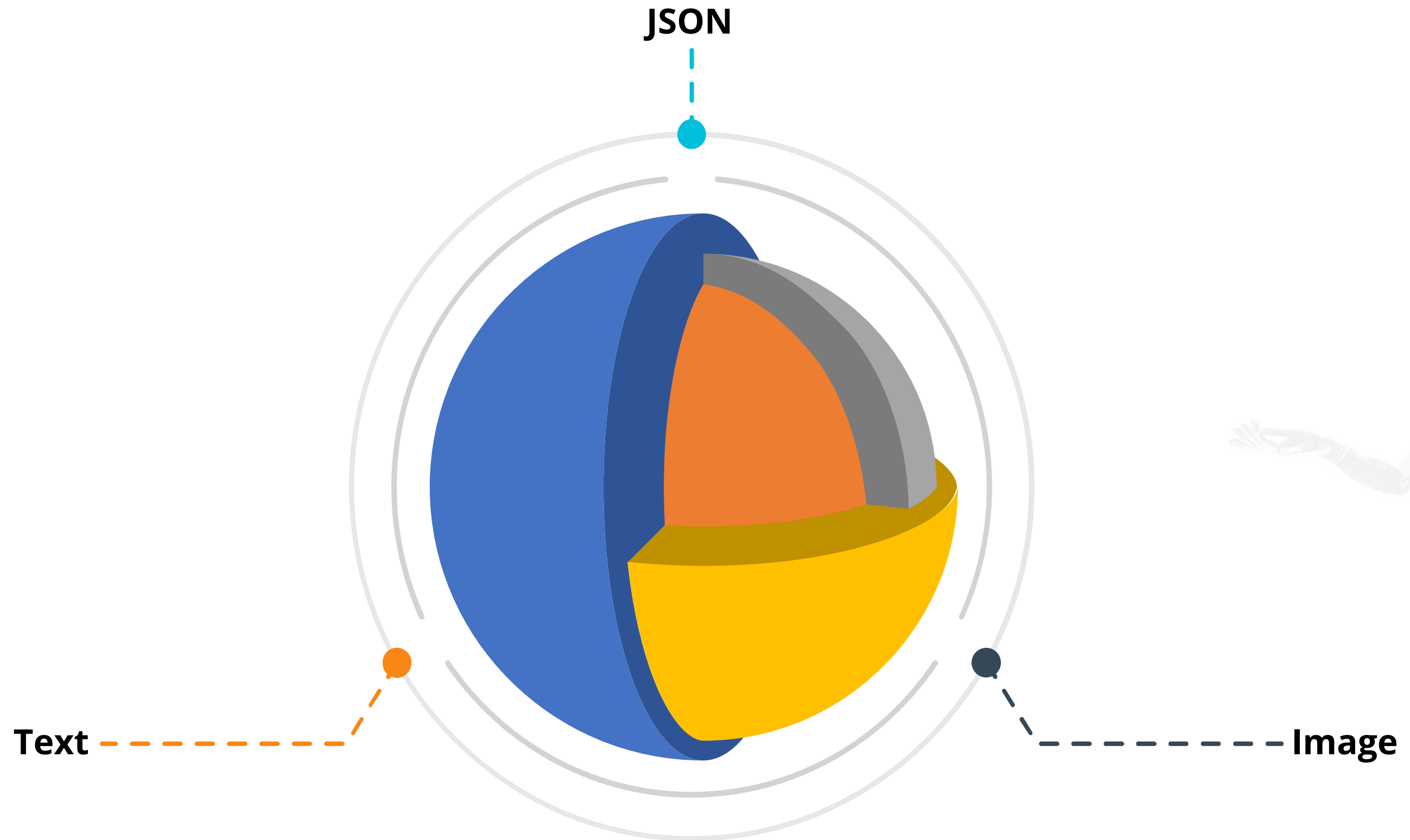
# Data Formats: Structured

	Excel, CSV

	SQL Data	

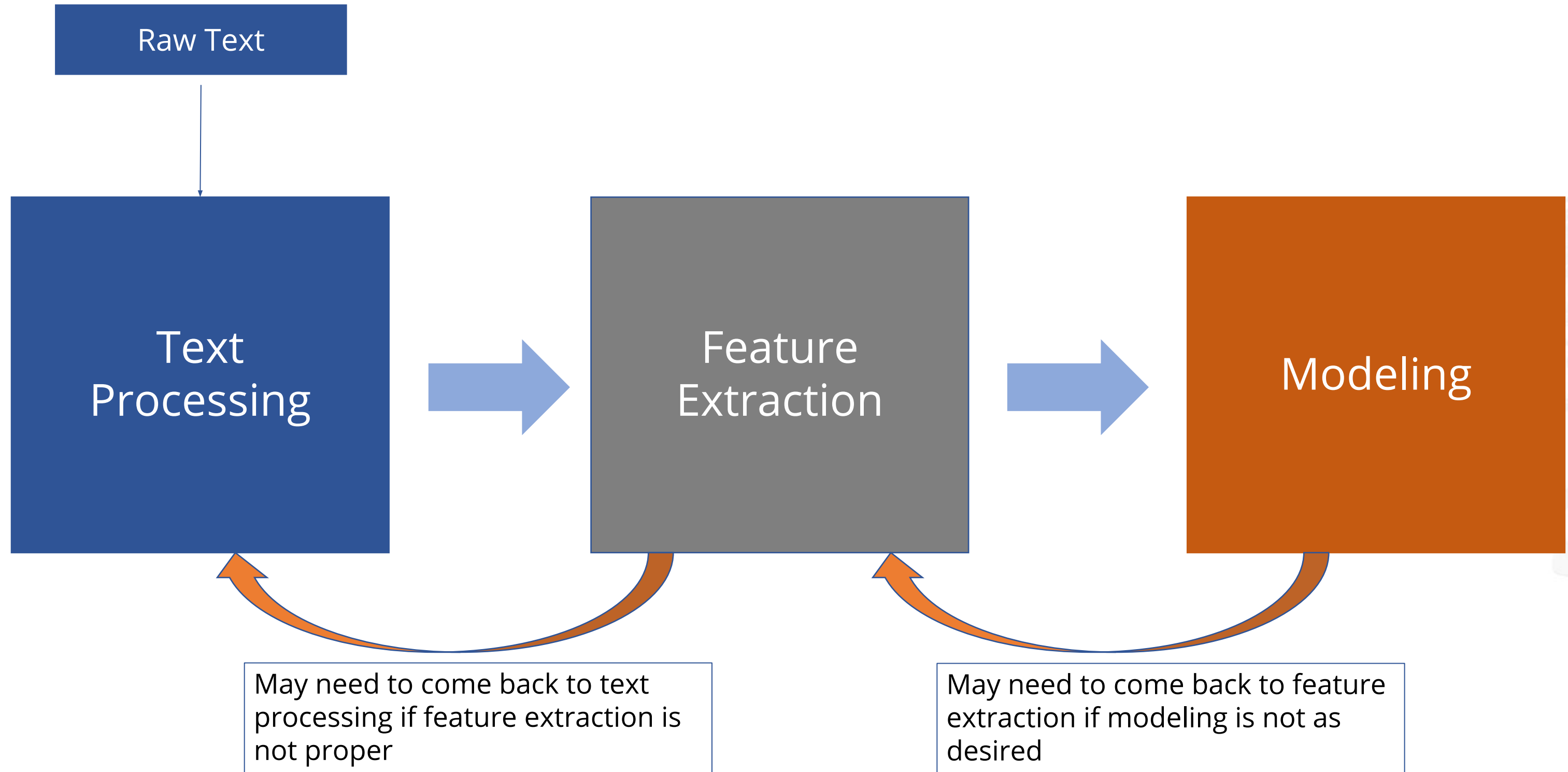


# Data Formats: Unstructured and Semi-Structured



## NLP Pipeline

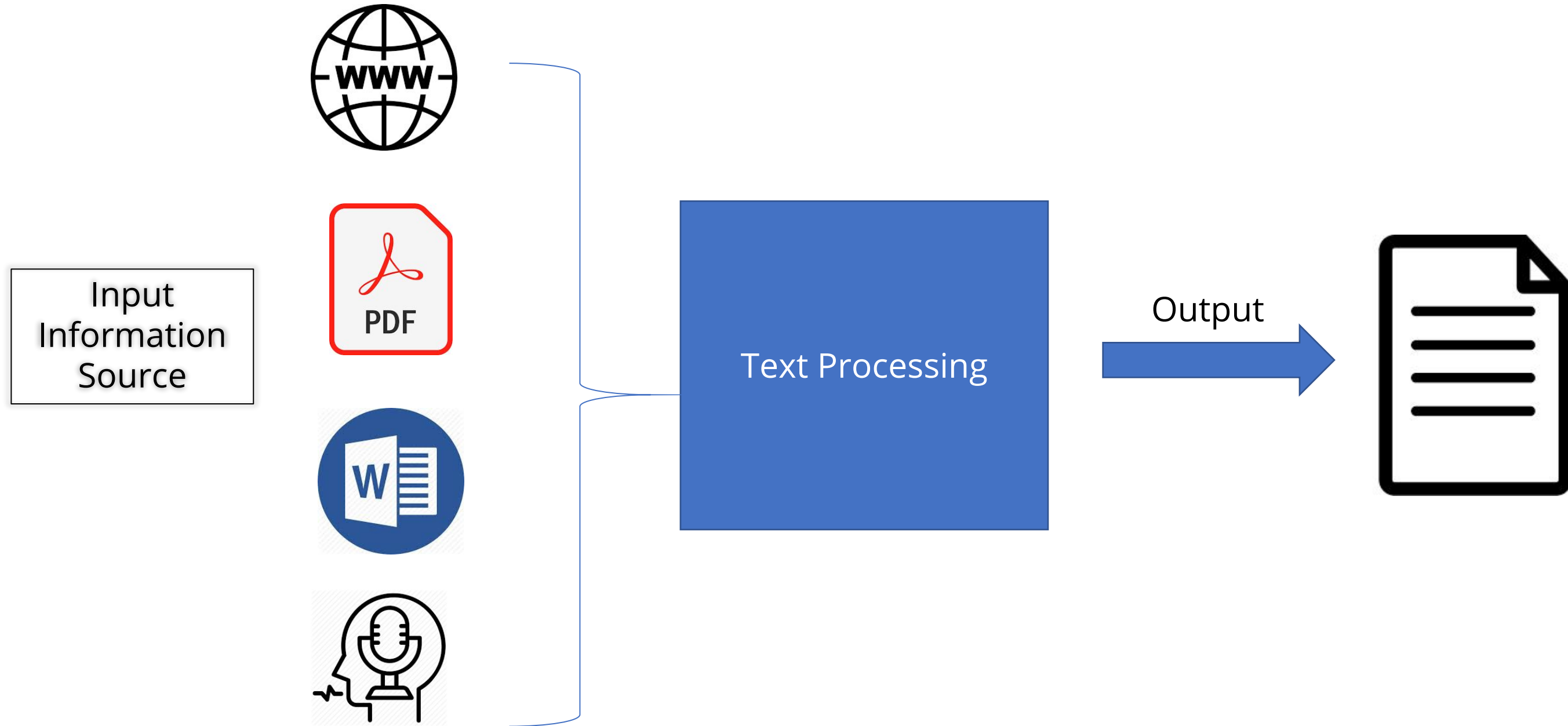
# NLP Pipeline



## Text Processing

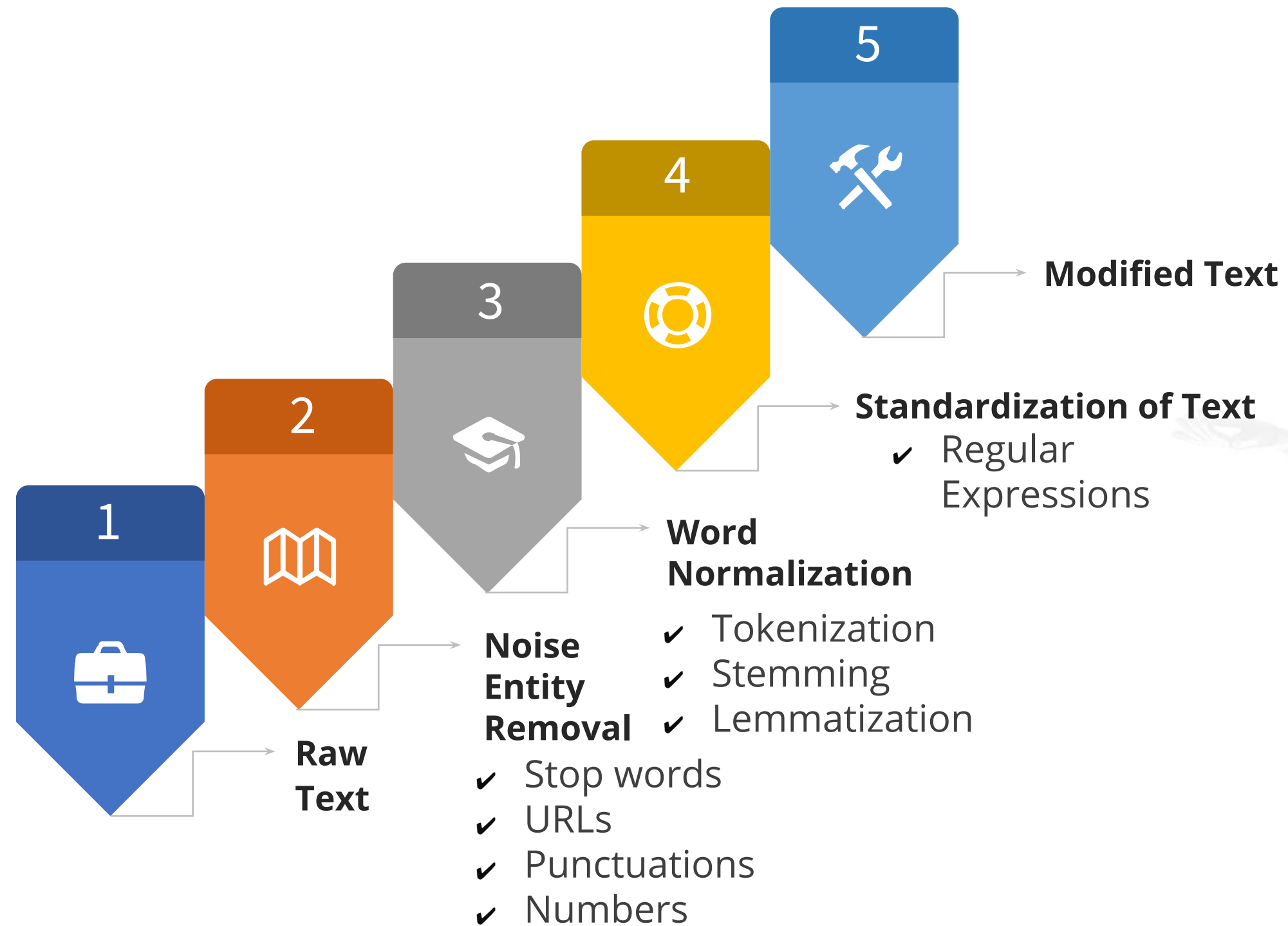


# Text Processing

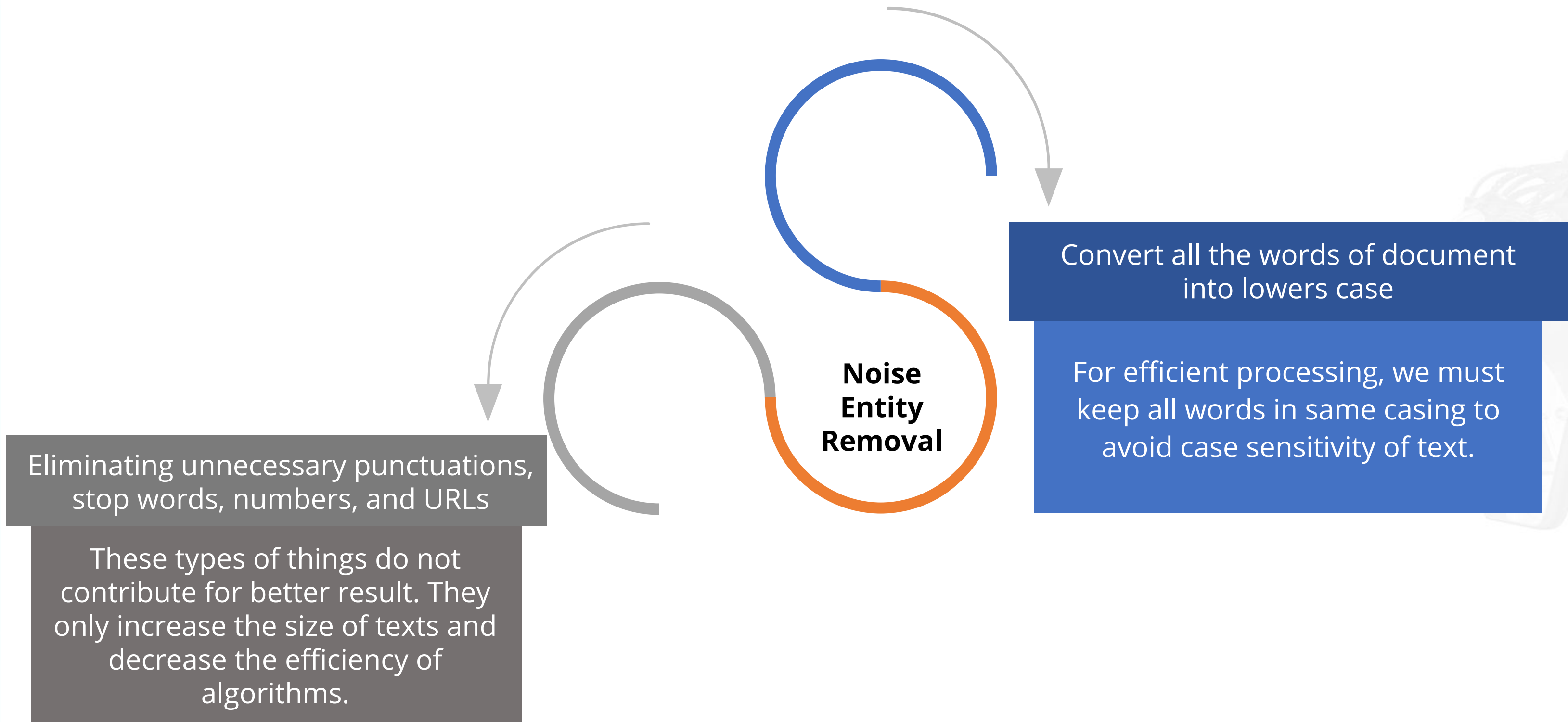


# Sequence or Text Processing

Sequence or Text Processing has the following steps:-



# Sequence or Text Processing: Noise Entity Removal



# Sequence or Text Processing: Tokenization

## Tokenization

- Break the sentence into separate words.
- These words are called tokens.
- Split words whenever there is a space between them.
- Treat punctuation marks as separate tokens since punctuation also has meaning.

### Example:

Sentence	Word
London is the capital and the most populous city of England and the United Kingdom	"London", "is", " the", "capital", "and", " the", "most", "populous", "city", "of", "England", "and", "the", "United", "Kingdom"



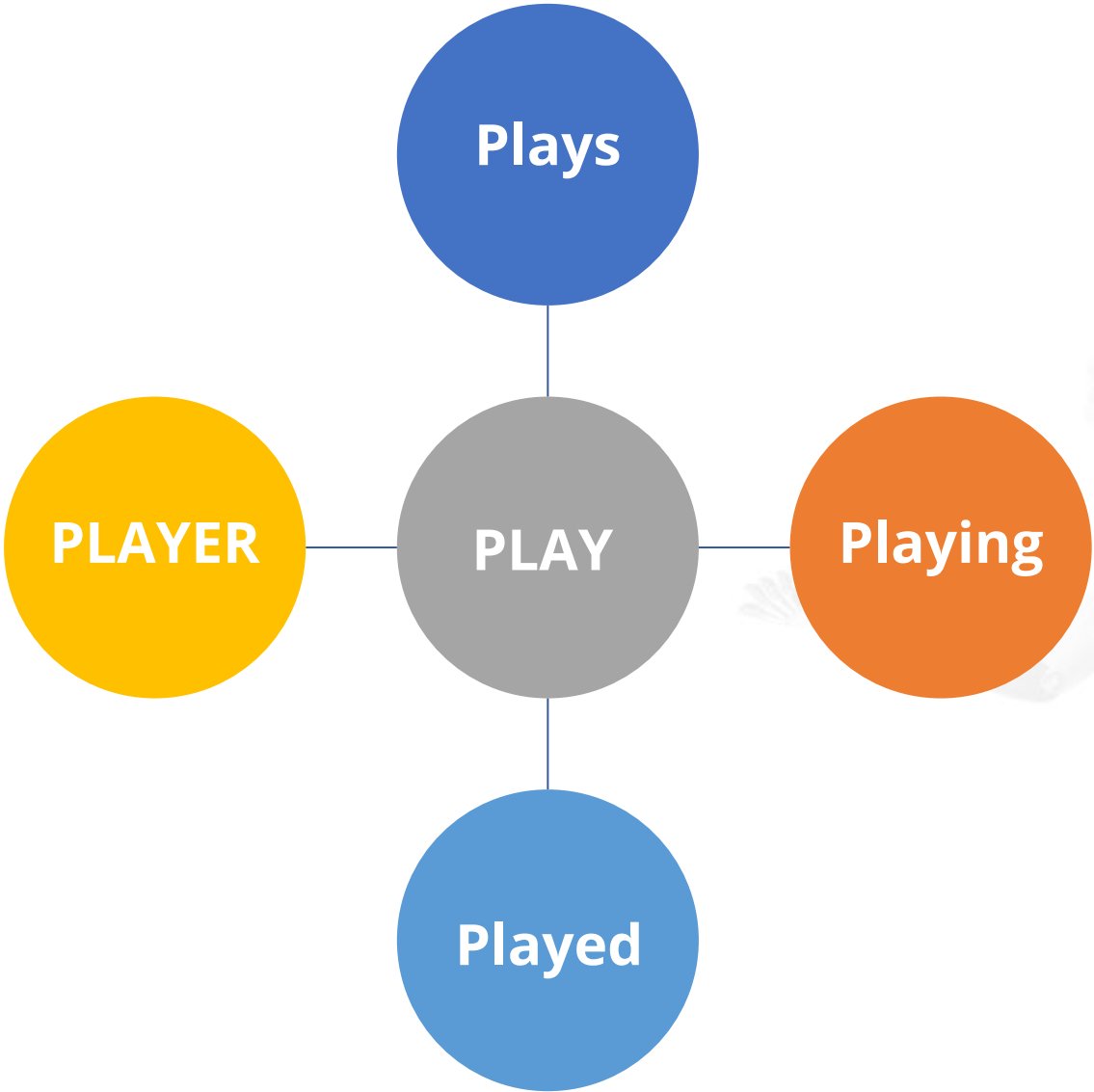
# Sequence or Text Processing: Stemming

## Stemming:

- It takes the root of the words.
- It removes the last few words or suffix of a word where it misspelt or incorrect words.

### Example:

Word	Suffix	Stem
studies	-es	studi
ninez	-ez	nin



# Sequence or Text Processing: Lemmatization

## Lemmatization:

It converts the text to meaningful base form by considering its context.

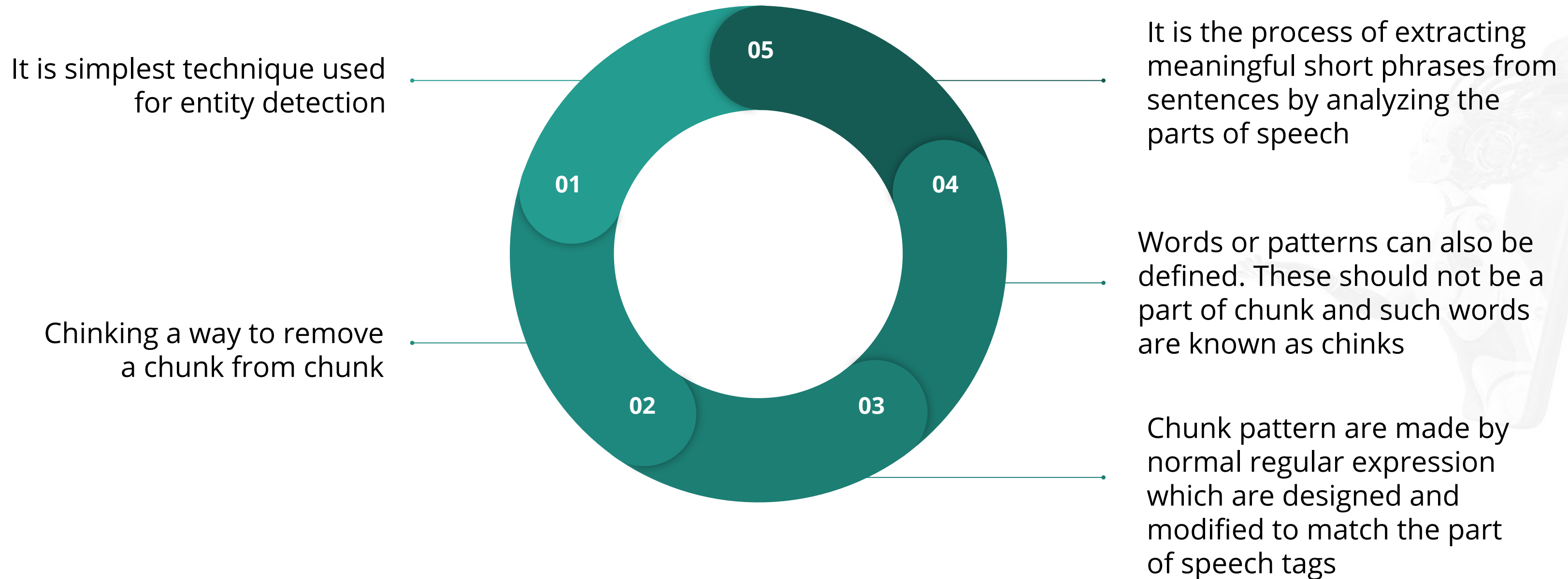
### Example:

Word	Morphological Information	Lemma
Studying	Gerund of the word study	Study
Ninez	Singular number of nine	Ninez

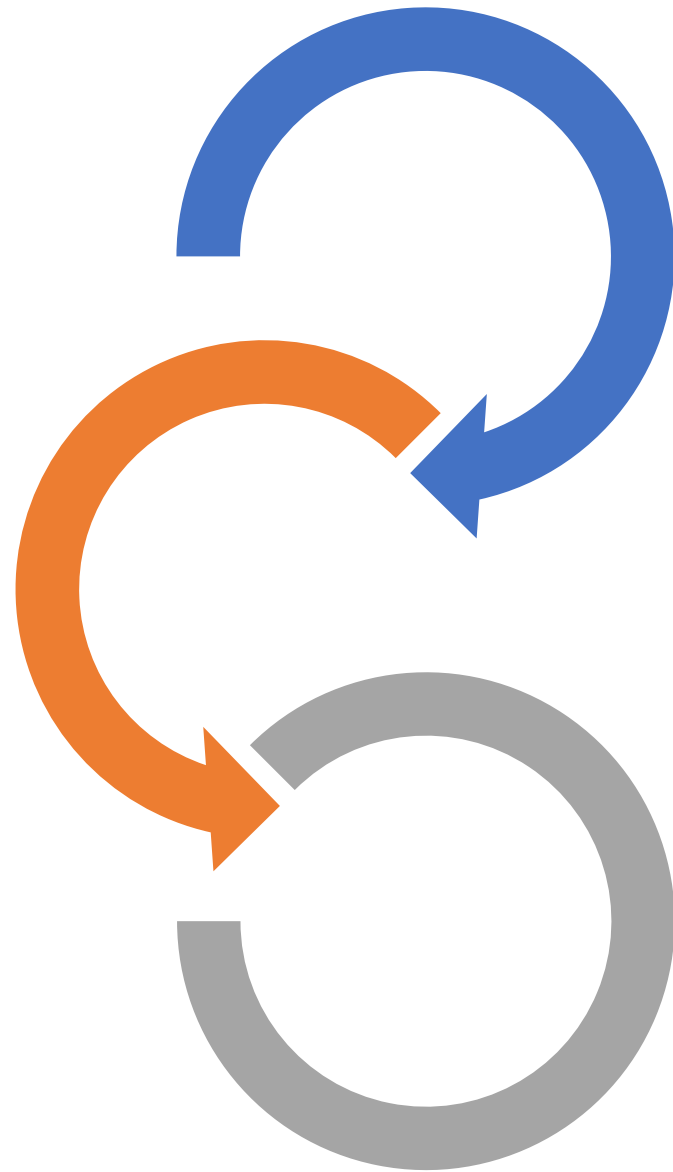




# Sequence or Text Processing: Chunking and Chinking



# Sequence or Text Processing: Regular Expression



## Object Standardization:

- Some words or symbols which are not present in standard dictionary are also not recognized by any search processes.
- Examples: hashtags, acronyms, and colloquial slangs

**Note:** With the help of regular expression, we can remove these things.

# Sequence or Text Processing: Regular Expression

## Regular Expression (Regex):

It is a sequence of characters that define pattern-matching, search-and-replace, and elimination functions. All type of noises can be removed with the help of regular expressions.

### Regex Examples:

Expression	Description
[abc]	Find any character between the brackets
[^abc]	Find any character that is not between the brackets
[0-9]	Find any character between the brackets (digit)

# DATA AND ARTIFICIAL INTELLIGENCE

## NLTK

# NLTK: Introduction

1

This tool is used for manipulation or understanding text or speech by any software or machine.

2

This is one of the most usable and mother of all NLP libraries.

3

It is a platform used for building Python programs that work with human language data for application in statistical Natural Language Processing (NLP).

# NLTK: Introduction

Following are text processing libraries:

Tokenization

Lemmatization

Parsing

Classification

Stemming

Tagging

Semantic Reasoning

# NLTK: Syntax and library

## System Requirement:

Operating System:

macOS / OS X · Linux · Windows (Cygwin, MinGW, Visual Studio)

Python Version:

Python 2.7, 3.5+ (only 64 bit)

```
>> import nltk
```





# NLTK: Lemmatization

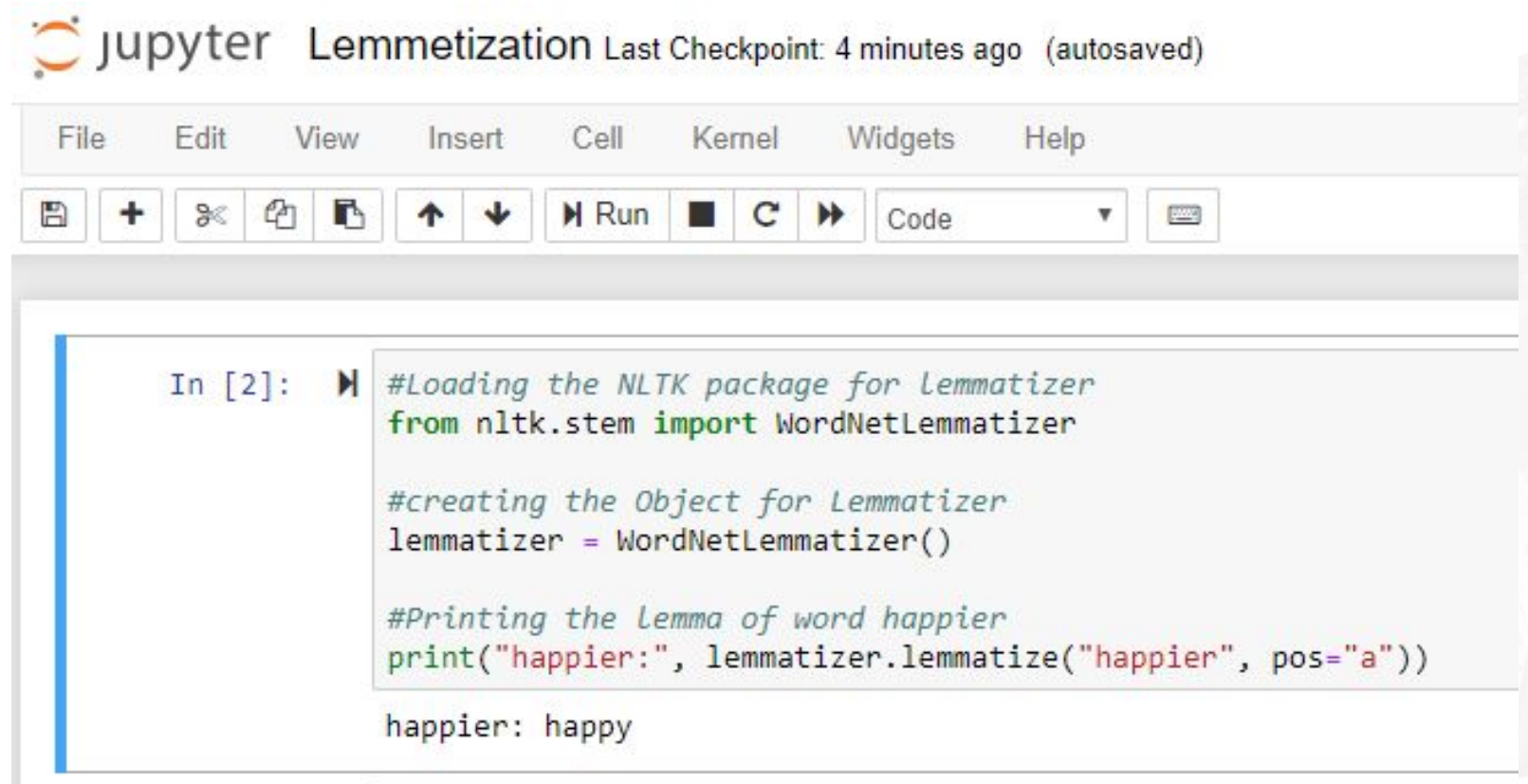
For grammatical purpose, documents are going to use different forms of a word, for example:

```
#Loading the NLTK package for
lemmatizer
from nltk.stem import WordNetLemmatizer

#creating the Object for Lemmatizer
lemmatizer = WordNetLemmatizer()

#Printing the lemma of word happier
print("happier:",
      lemmatizer.lemmatize("happier",
                           pos="a"))
```

**Output:** happier: happy

A screenshot of a Jupyter Notebook interface. The title bar says "jupyter Lemmetization" (note the typo) and "Last Checkpoint: 4 minutes ago (autosaved)". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below the menu is a toolbar with icons for saving, adding cells, undo, redo, and running code. The main area shows a code cell labeled "In [2]:". The code in the cell is the same as the one in the previous block, but with some lines commented out or in a different color. The output of the cell is "happier: happy".

```
jupyter Lemmetization Last Checkpoint: 4 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

[Icons] [Run] [Code]

In [2]: #Loading the NLTK package for lemmatizer
        from nltk.stem import WordNetLemmatizer

        #creating the Object for Lemmatizer
        lemmatizer = WordNetLemmatizer()

        #Printing the lemma of word happier
        print("happier:", lemmatizer.lemmatize("happier", pos="a"))

        happier: happy
```

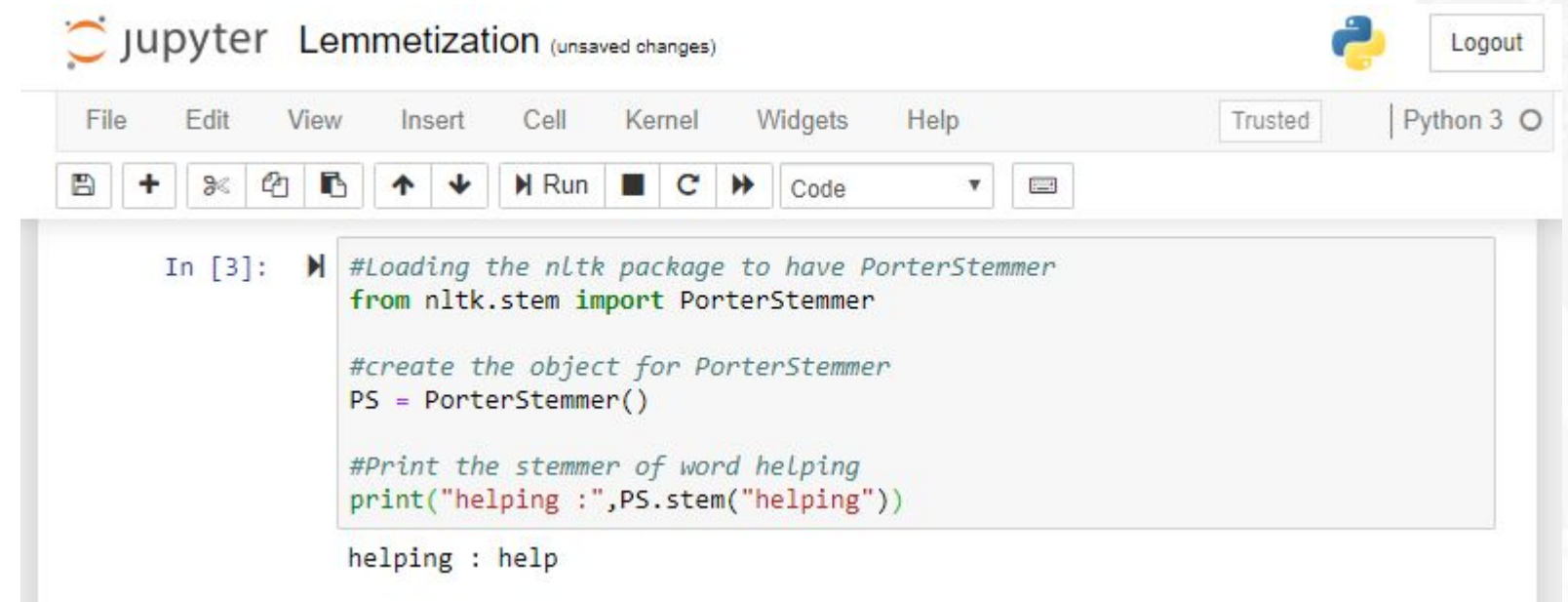
# NLTK: Stemming

```
#Loading the nltk package to have
PorterStemmer
from nltk.stem import PorterStemmer

#create the object for PorterStemmer
PS = PorterStemmer()

#Print the stemmer of word helping
print("helping :",PS.stem("helping"))
```

**Output:** helping: help

A screenshot of a Jupyter Notebook interface. The title bar shows 'jupyter Lemmetization (unsaved changes)' with a Python logo and a 'Logout' button. The menu bar includes 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. Below the menu is a toolbar with icons for saving, adding, deleting, and running code. The main area shows a code cell with the following Python code: 

```
In [3]: #Loading the nltk package to have PorterStemmer
from nltk.stem import PorterStemmer

#create the object for PorterStemmer
PS = PorterStemmer()

#Print the stemmer of word helping
print("helping :",PS.stem("helping"))
```

 The output of the code is displayed below the cell: 

```
helping : help
```

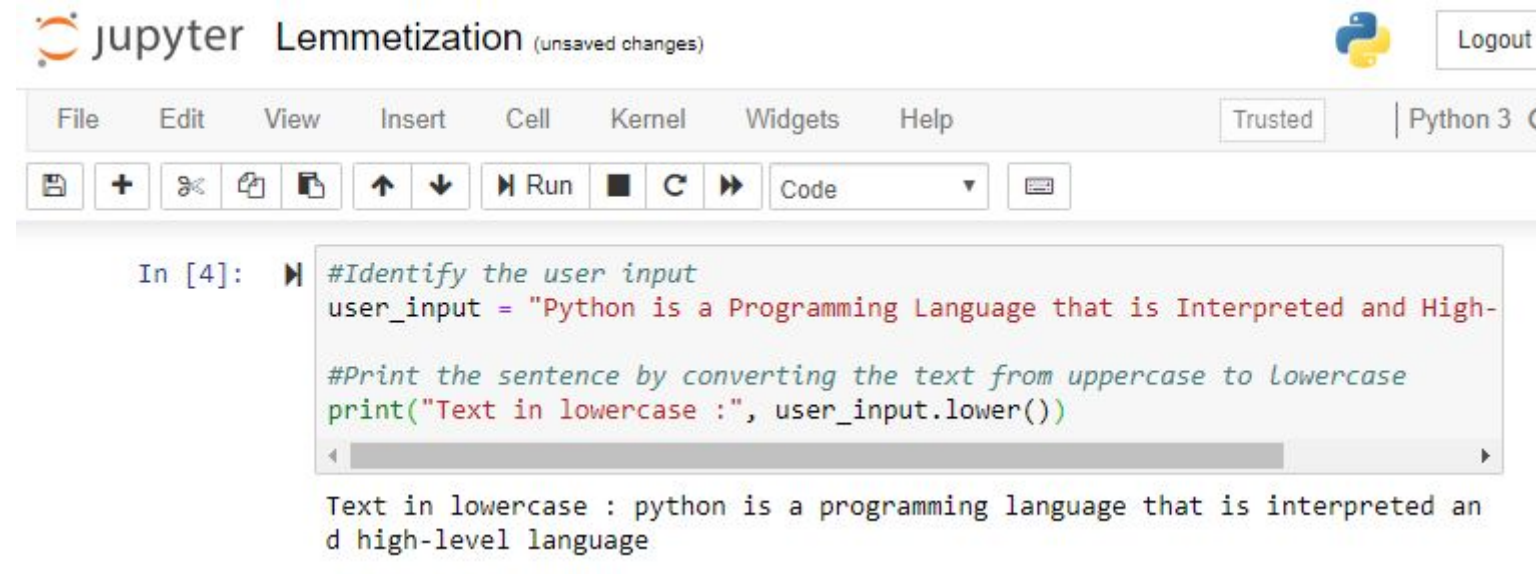
# NLTK: Processing Raw Text

Text processing includes:  
Converting all letters to lower or upper case

```
#Identify the user input
user_input = "Python is a Programming Language that is Interpreted and
High-Level language"

#Print the sentence by converting the text from uppercase to lowercase
print("Text in lowercase :", user_input.lower())
```

**Output:** Text in lowercase : python is a programming language that is interpreted and high-level language.

A screenshot of a Jupyter Notebook interface. The title bar shows 'jupyter Lemmetization (unsaved changes)' and a 'Logout' button. The menu bar includes 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. Below the menu is a toolbar with icons for saving, adding cells, and running code. The main area shows a code cell with the same Python code as in the previous block. Below the code cell, the output is displayed: 'Text in lowercase : python is a programming language that is interpreted and high-level language'.

```
jupyter Lemmetization (unsaved changes) Logout
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3
In [4]: #Identify the user input
user_input = "Python is a Programming Language that is Interpreted and High-
#Print the sentence by converting the text from uppercase to lowercase
print("Text in lowercase :", user_input.lower())
Text in lowercase : python is a programming language that is interpreted and
high-level language
```

# NLTK: Processing Raw Text

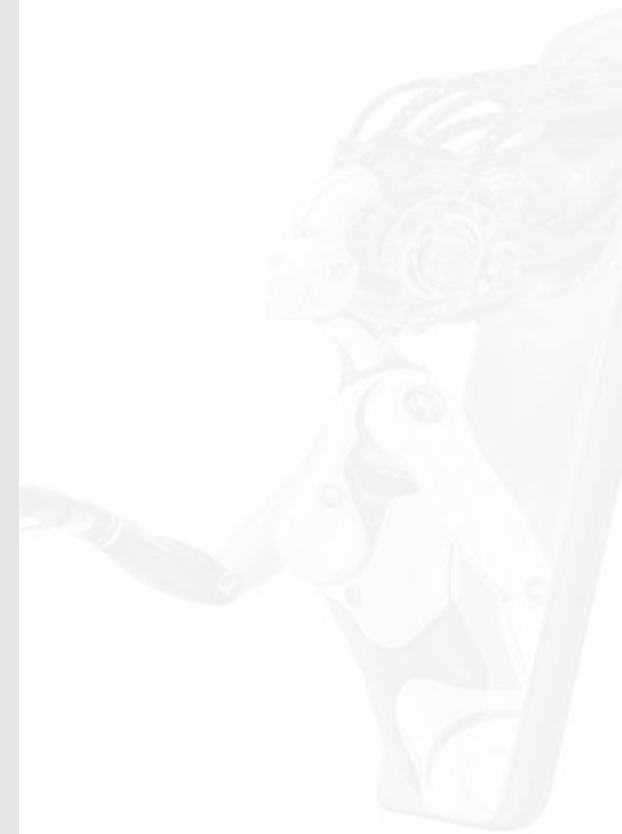
Converting numbers into words or removing numbers

```
#Loading the regex package to find number
import re

#identify user input
input_str = "Team A has 6 batsman and 5 bowlers, while team b has
5 batsman and 6 bowlers"

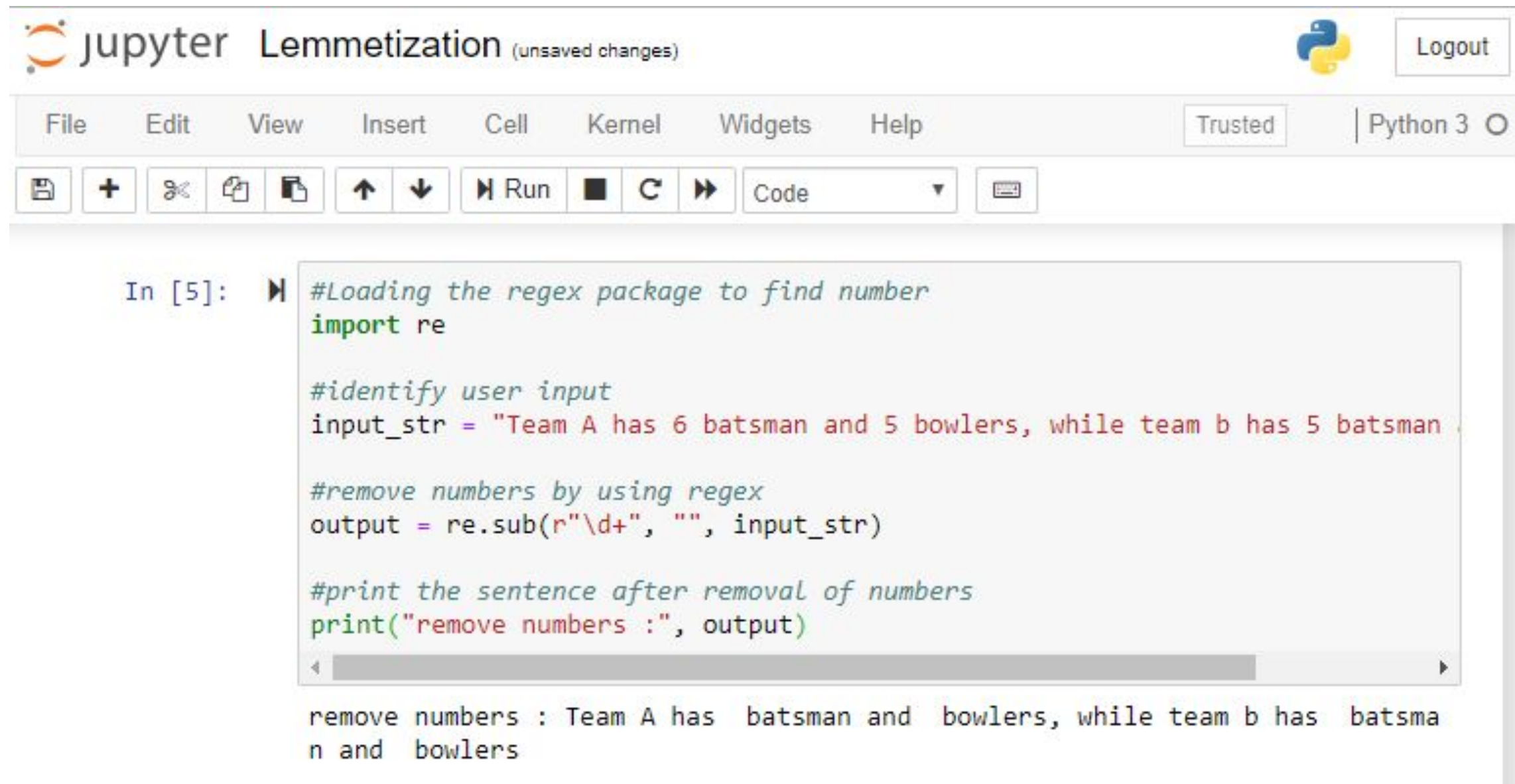
#remove numbers by using regex
output = re.sub(r"\d+", "", input_str)

#print the sentence after removal of numbers
print("remove numbers :", output)
```



# NLTK: Processing Raw Text

**Output:** remove numbers : Team A has batsman and bowlers, while team b has batsman and bowlers



The image shows a Jupyter Notebook interface with the title "Lemmetization (unsaved changes)". The notebook has a menu bar with "File", "Edit", "View", "Insert", "Cell", "Kernel", "Widgets", and "Help". Below the menu bar is a toolbar with icons for saving, adding cells, undo, redo, and running code. The code cell is labeled "In [5]:" and contains the following Python code:

```
#Loading the regex package to find number
import re

#identify user input
input_str = "Team A has 6 batsman and 5 bowlers, while team b has 5 batsman"

#remove numbers by using regex
output = re.sub(r"\d+", "", input_str)

#print the sentence after removal of numbers
print("remove numbers :", output)
```

The output of the code is displayed below the cell:

```
remove numbers : Team A has  batsman and  bowlers, while team b has  batsma
n and  bowlers
```



# NLTK: Processing Raw Text

Removing accent, punctuations marks, and other diacritics

```
#Load the regex package and string package
import re, string

#define user input
input_str = "Sentence. having. string with. Punctuation?"

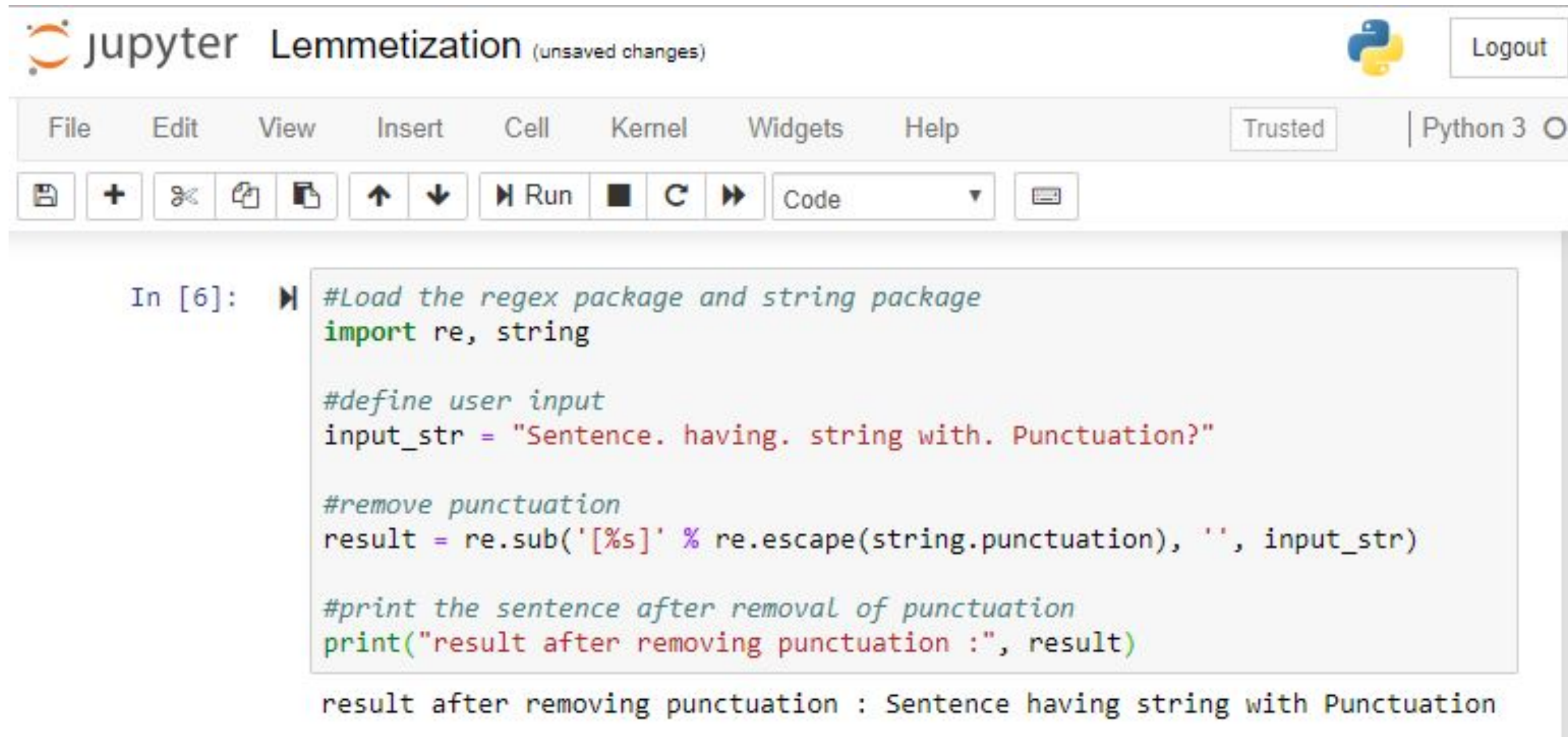
#remove punctuation
result = re.sub('[%s]' % re.escape(string.punctuation), '', input_str)

#print the sentence after removal of punctuation
print("result after removing punctuation :", result)
```



# NLTK: Processing Raw Text

**Output:** result after removing punctuation : Sentence having string with Punctuation



The image shows a Jupyter Notebook interface with the title 'Lemmetization (unsaved changes)'. The interface includes a top bar with the Jupyter logo, a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), and a toolbar with icons for saving, adding, deleting, and running code. The code cell contains the following Python code:

```
In [6]: #Load the regex package and string package
import re, string

#define user input
input_str = "Sentence. having. string with. Punctuation?"

#remove punctuation
result = re.sub('[%s]' % re.escape(string.punctuation), '', input_str)

#print the sentence after removal of punctuation
print("result after removing punctuation :", result)
```

The output of the code is displayed below the cell:

```
result after removing punctuation : Sentence having string with Punctuation
```



# NLTK: Processing Raw Text

## Removing white spaces:

```
#Load the regex and string package
import re

#define input from user
input_str = 'pythonis programming    language \t\n\r\tHello \t'

#Print the sentence after removing the spaces
print('Remove spaces using regex :', re.sub(r"\s+", "", input_str), "\n", sep='')

#Print the sentence after removing the landing spaces
print('Remove landing spaces using regex :', re.sub(r"^\s+", "",
input_str), "\n", sep='')

#Print the sentence after removing the trailing spaces
print('Remove trailing spaces using regex :', re.sub(r"\s+$", "",
input_str), "\n", sep='')

#Print the sentence after removing the leading and trailing spaces
print('Remove landing spaces using regex :', re.sub(r"^\s+|\s+$", "",
input_str), "\n", sep='')
```

# NLTK: Processing Raw Text

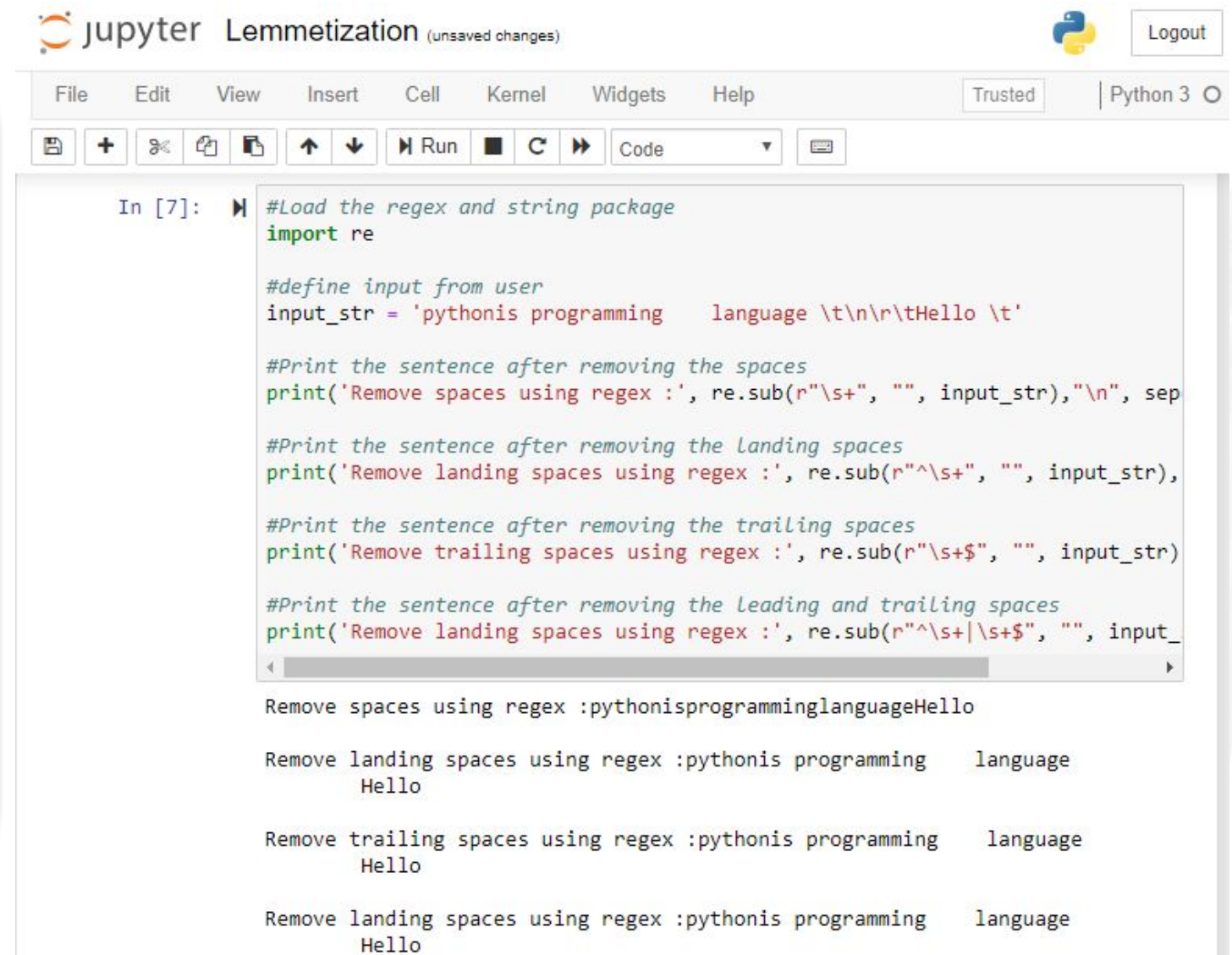
## Output:

Remove spaces using regex  
:pythonisprogramminglanguageHello

Remove landing spaces using regex :pythonis programming  
language  
Hello

Remove trailing spaces using regex :pythonis programming  
language  
Hello

Remove landing spaces using regex :pythonis programming  
language  
Hello



The image shows a Jupyter Notebook interface with the title 'Lemmetization (unsaved changes)'. The notebook has a menu bar with 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. Below the menu bar is a toolbar with icons for saving, adding cells, undo, redo, and running code. The code cell contains the following Python code:

```
In [7]: #Load the regex and string package
import re

#define input from user
input_str = 'pythonis programming    language \t\n\r\tHello \t'

#Print the sentence after removing the spaces
print('Remove spaces using regex :', re.sub(r"\s+", "", input_str), "\n", sep='')

#Print the sentence after removing the landing spaces
print('Remove landing spaces using regex :', re.sub(r"^\s+", "", input_str), "\n", sep='')

#Print the sentence after removing the trailing spaces
print('Remove trailing spaces using regex :', re.sub(r"\s+$", "", input_str), "\n", sep='')

#Print the sentence after removing the leading and trailing spaces
print('Remove landing spaces using regex :', re.sub(r"^\s+|\s+$", "", input_str), "\n", sep='')
```

The output of the code cell is displayed below the code:

```
Remove spaces using regex :pythonisprogramminglanguageHello

Remove landing spaces using regex :pythonis programming    language
Hello

Remove trailing spaces using regex :pythonis programming    language
Hello

Remove landing spaces using regex :pythonis programming    language
Hello
```

# NLTK: Stopwords

```
#Load the stopwords package
from nltk.corpus import stopwords

#Load the word tokenizer package
from nltk.tokenize import word_tokenize

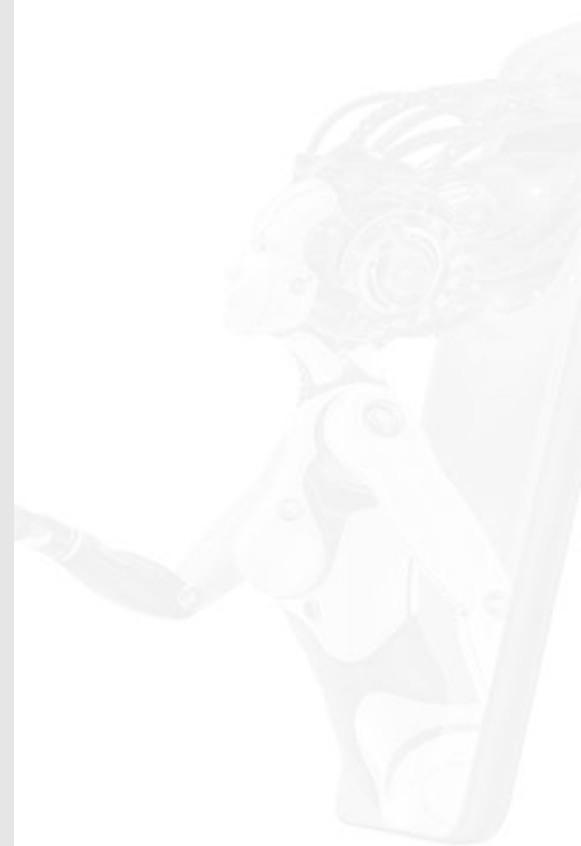
#define the user input
input_str = "Stop words are the words that are filtered before
and after processing of text."

#create object for stopwords
stop_word = set(stopwords.words("english"))

#convert word into tokens
token = word_tokenize(input_str)

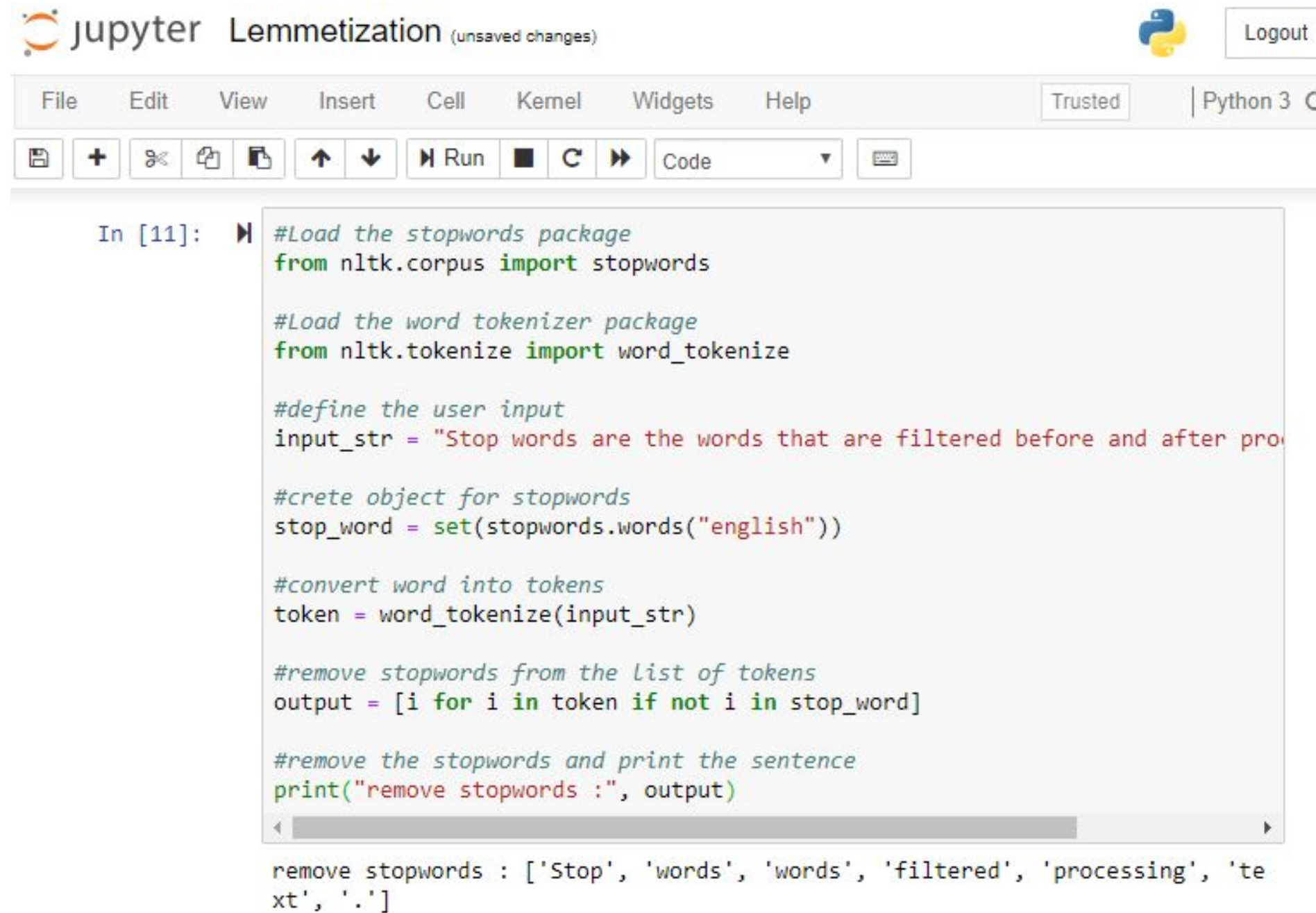
#remove stopwords from the list of tokens
output = [i for i in token if not i in stop_word]

#remove the stopwords and print the sentence
print("remove stopwords :", output)
```



# NLTK: Stopwords

**Output:** remove stopwords : ['Stop', 'words', 'words', 'filtered', 'processing', 'text', '.']



The image shows a Jupyter Notebook interface with the title 'Lemmetization (unsaved changes)'. The notebook contains a single code cell with the following Python code:

```
In [11]: #Load the stopwords package
from nltk.corpus import stopwords

#Load the word tokenizer package
from nltk.tokenize import word_tokenize

#define the user input
input_str = "Stop words are the words that are filtered before and after pro

#create object for stopwords
stop_word = set(stopwords.words("english"))

#convert word into tokens
token = word_tokenize(input_str)

#remove stopwords from the list of tokens
output = [i for i in token if not i in stop_word]

#remove the stopwords and print the sentence
print("remove stopwords :", output)
```

The output of the code is displayed below the cell:

```
remove stopwords : ['Stop', 'words', 'words', 'filtered', 'processing', 'te
xt', '.']
```

# NLTK: Tokenizers

```
#Load the package for tokenizer
from nltk.tokenize import sent_tokenize

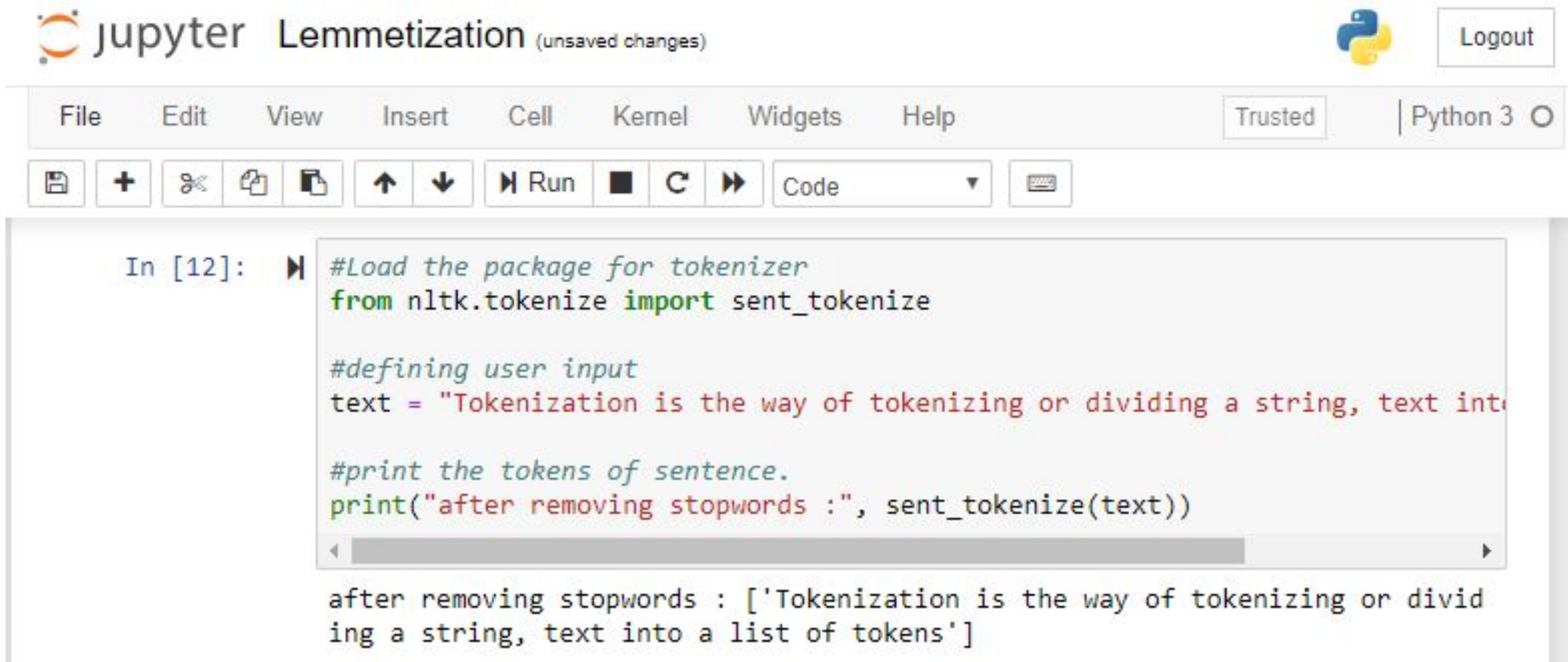
#defining user input
text = "Tokenization is the way of tokenizing or dividing
a string, text into a list of tokens"

#print the tokens of sentence.
print("after removing stopwords :", sent_tokenize(text))
```



# NLTK: Tokenizers

**Output:** after removing stopwords : ['Tokenization is the way of tokenizing or dividing a string, text into a list of tokens']



The image shows a Jupyter Notebook interface with the title 'Lemmetization (unsaved changes)'. The top bar includes a 'Logout' button and a 'Python 3' selector. The menu bar contains 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. Below the menu is a toolbar with icons for saving, adding cells, undo, redo, and running code. The code cell, labeled 'In [12]:', contains the following Python code:

```
#Load the package for tokenizer
from nltk.tokenize import sent_tokenize

#defining user input
text = "Tokenization is the way of tokenizing or dividing a string, text into a list of tokens"

#print the tokens of sentence.
print("after removing stopwords :", sent_tokenize(text))
```

The output of the code is displayed below the cell:

```
after removing stopwords : ['Tokenization is the way of tokenizing or dividing a string, text into a list of tokens']
```

# NLTK: Ngram

```
#Load the package for ngrams
from nltk import ngrams

#define user input
usr_input = 'i want to ngramize the foo bar
sentences'

#define number of gram
n = 3

#split the sentence to make grams
sixgrams = ngrams(usr_input.split(), n)

for grams in sixgrams:
    #Print the 3 grams of user-input
    print(grams)
```

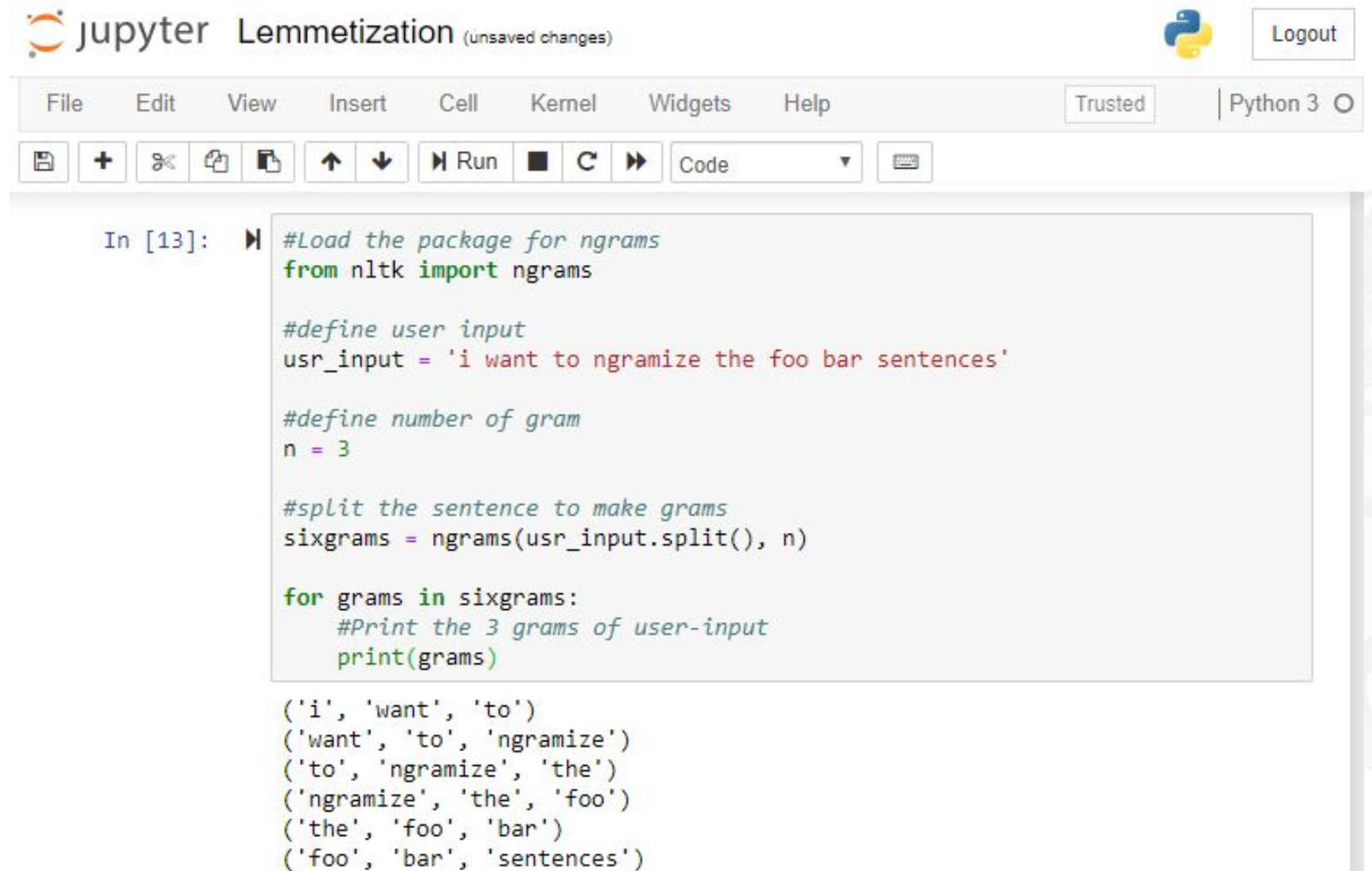




# NLTK: Ngram

## Output:

('i', 'want', 'to')  
('want', 'to', 'ngramize')  
('to', 'ngramize', 'the')  
('ngramize', 'the', 'foo')  
('the', 'foo', 'bar')  
('foo', 'bar', 'sentences')



The image shows a Jupyter Notebook interface with the title 'Lemmetization (unsaved changes)'. The top bar includes a menu (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a 'Trusted' status indicator, and a 'Python 3' version selector. Below the menu is a toolbar with icons for file operations and execution. The main area contains a code cell with the following Python code:

```
In [13]: #Load the package for ngrams
          from nltk import ngrams

          #define user input
          usr_input = 'i want to ngramize the foo bar sentences'

          #define number of gram
          n = 3

          #split the sentence to make grams
          sixgrams = ngrams(usr_input.split(), n)

          for grams in sixgrams:
              #Print the 3 grams of user-input
              print(grams)
```

The output of the code cell is displayed below the code:

```
('i', 'want', 'to')
('want', 'to', 'ngramize')
('to', 'ngramize', 'the')
('ngramize', 'the', 'foo')
('the', 'foo', 'bar')
('foo', 'bar', 'sentences')
```

# NLTK: Limitations

1

Does not support word vectors

2

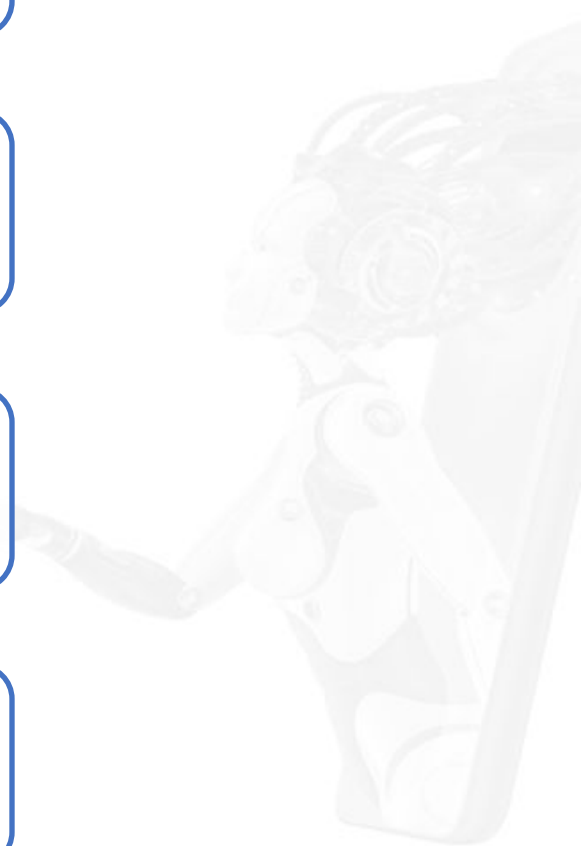
Is slow

3

Not for production purpose

4

Good only for English and difficult for other languages



# DATA AND ARTIFICIAL INTELLIGENCE

**Re**

## Re: Introduction

- Re is an inbuilt library which comes with python.
- It uses a set of symbols to identify the patterns from the text.  
Example: email address `^([a-zA-Z0-9_\-\.]+)@([a-zA-Z0-9_\-\.]+)\.([a-zA-Z]{2,5})$`
- It is used in information retrieval: **import nltk**

```
import re
```

# Text Processing Using Stemming and Regular Expression



**Problem Statement:** Demonstrate text processing using stemming and regular expression.

**Access:** Click on the **Practice Labs** tab on the left side panel of the LMS. Copy or note the username and password that is generated. Click on the **Launch Lab** button. On the page that appears, enter the username and password in the respective fields, and click **Login**.

ASSISTED PRACTICE



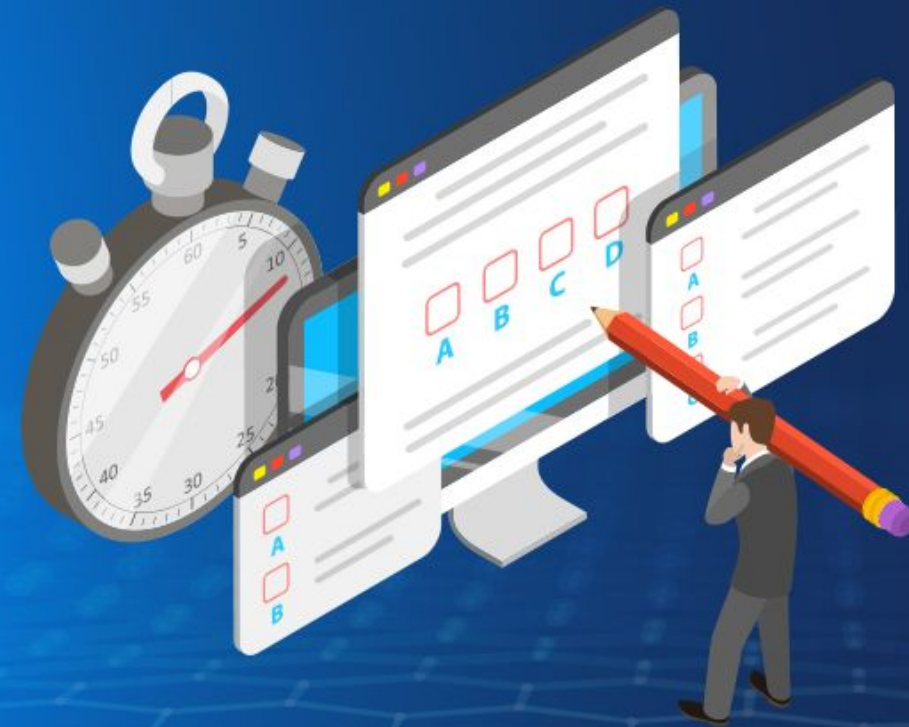
# Tweets Cleanup and Analysis Using Regular Expressions



**Objective:** Use regular expressions to work with messy tweets data: clean up the data, extract hashtags, analyze the most popular hashtags that occur along with a target hashtag (#economy).

**Problem Statement:** Social media is a gold mine of information. Brands, governments, or anyone can leverage their business with the help of the information contained. It can be information on the sentiments for a brand, or the themes being spoken about, or the associated trends for a particular hashtag. In this project, we will work on the tweets on Twitter. We will find other hashtags that occur frequently with our target hashtag. This will give us an understanding of which other topics people are associating this hashtag with.

# DATA AND ARTIFICIAL INTELLIGENCE



## Knowledge Check

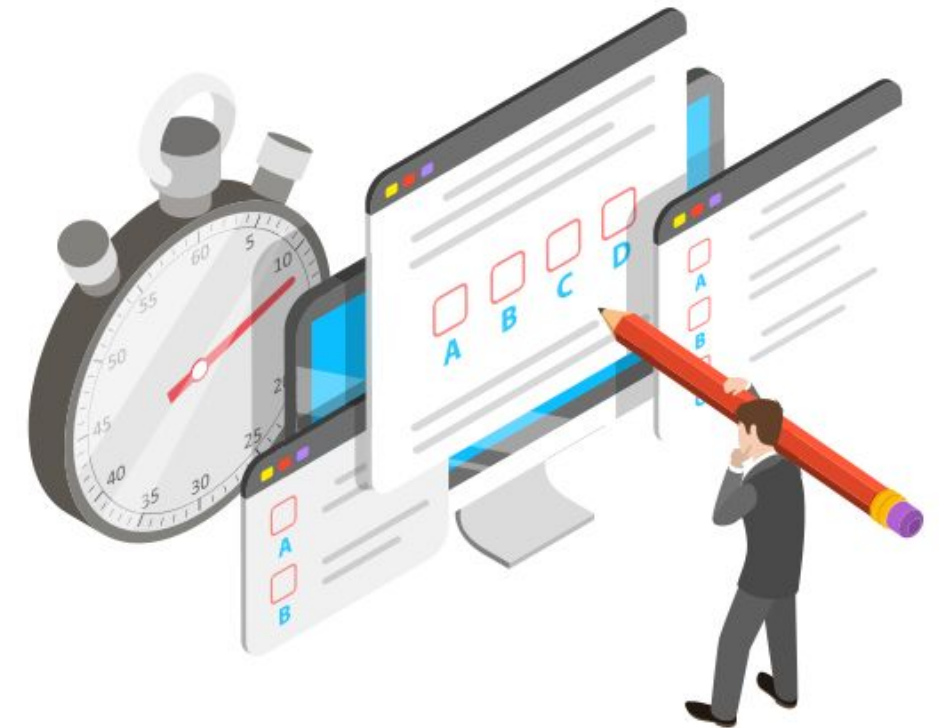


## Knowledge Check

1

One of the main challenges of NLP is \_\_\_\_\_.

- a. Handling ambiguity of sentences
- b. Handling tokenization
- c. Both a and b
- d. None of the above

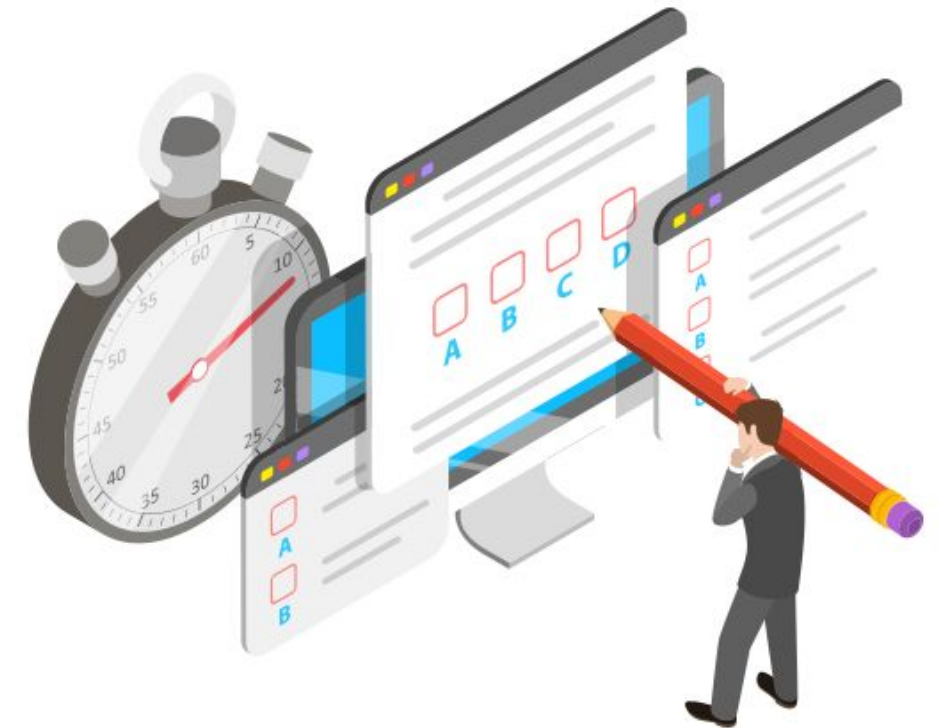


## Knowledge Check

1

One of the main challenge of NLP is \_\_\_\_\_.

- a. Handling ambiguity of sentences
- b. Handling tokenization
- c. Both a and b
- d. None of the above



The correct answer is **a.**

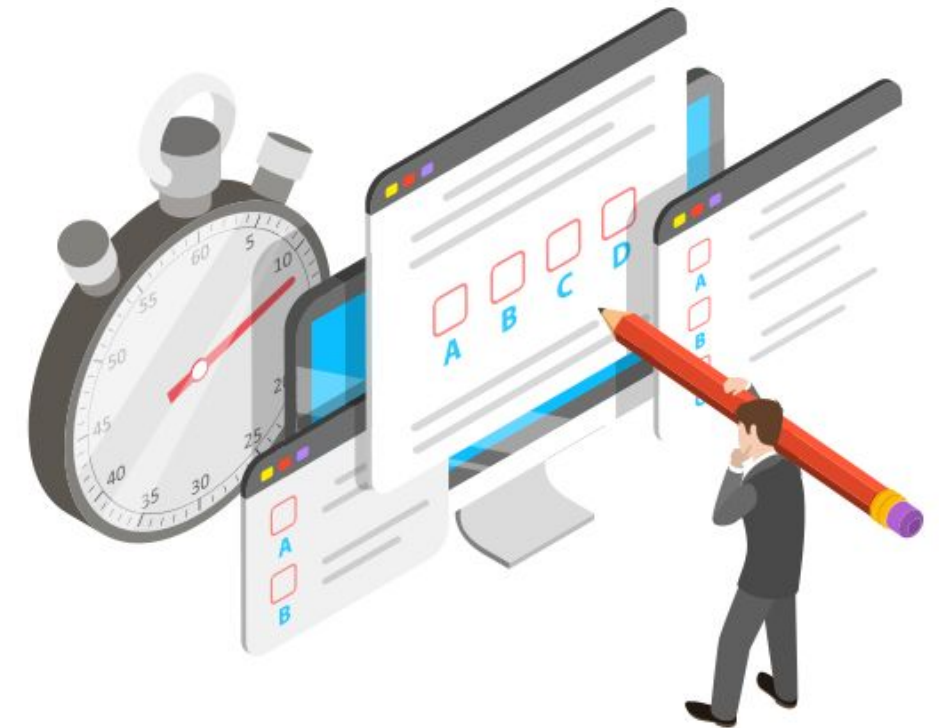
**One of the main challenges of NLP is handling ambiguity of sentences.**

## Knowledge Check

2

Regular expression is used for\_\_\_\_\_.

- a. Information retrieval
- b. Finding the pattern
- c. Database management
- d. Both a and b

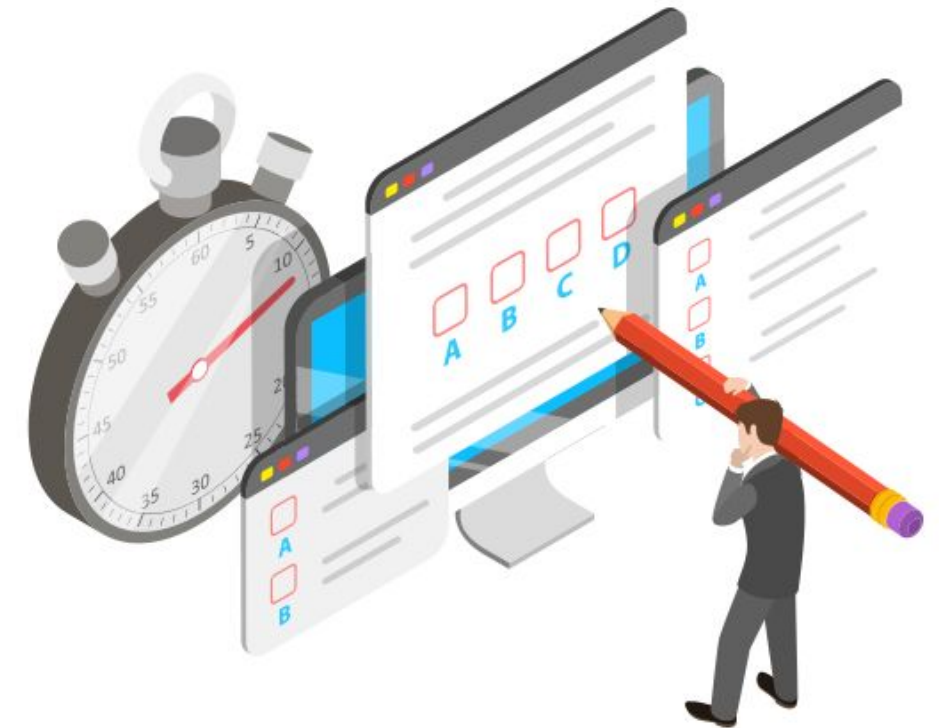


## Knowledge Check

2

Regular expression is used for\_\_\_\_\_.

- a. Information retrieval
- b. Finding the pattern
- c. Database management
- d. Both a and b



The correct answer is **d.**

**Regular expression is used for information retrieval and finding the pattern.**

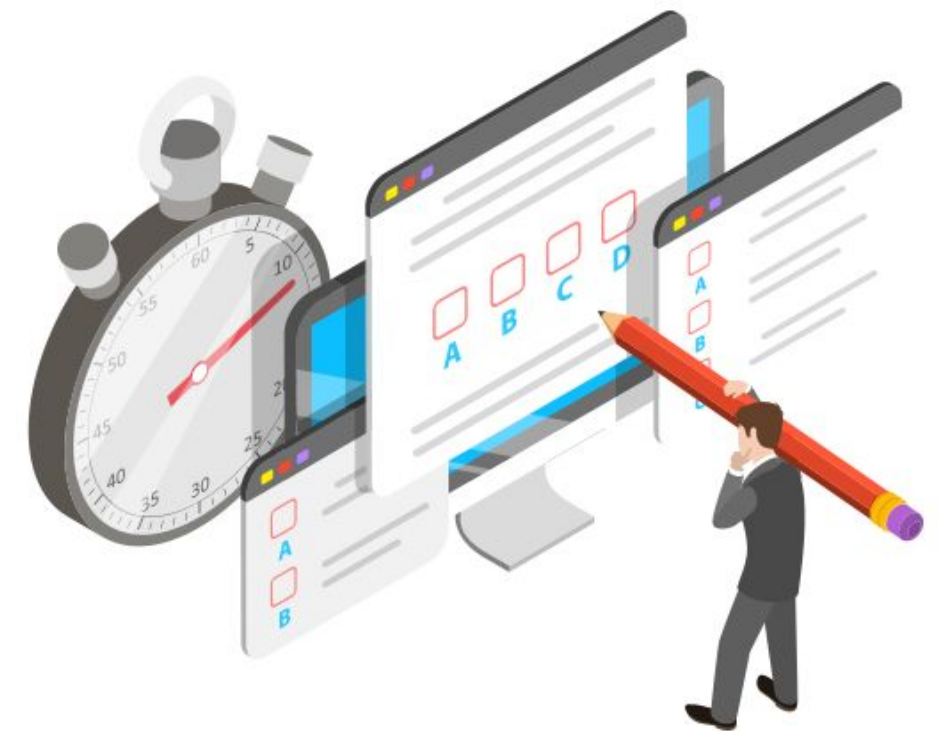
## Knowledge Check

3

NLP is the technique of interpretation of all types of languages which includes

\_\_\_\_\_.

- a. Human Language
- b. Assembly Language
- c. Machine Language
- d. Binary Data

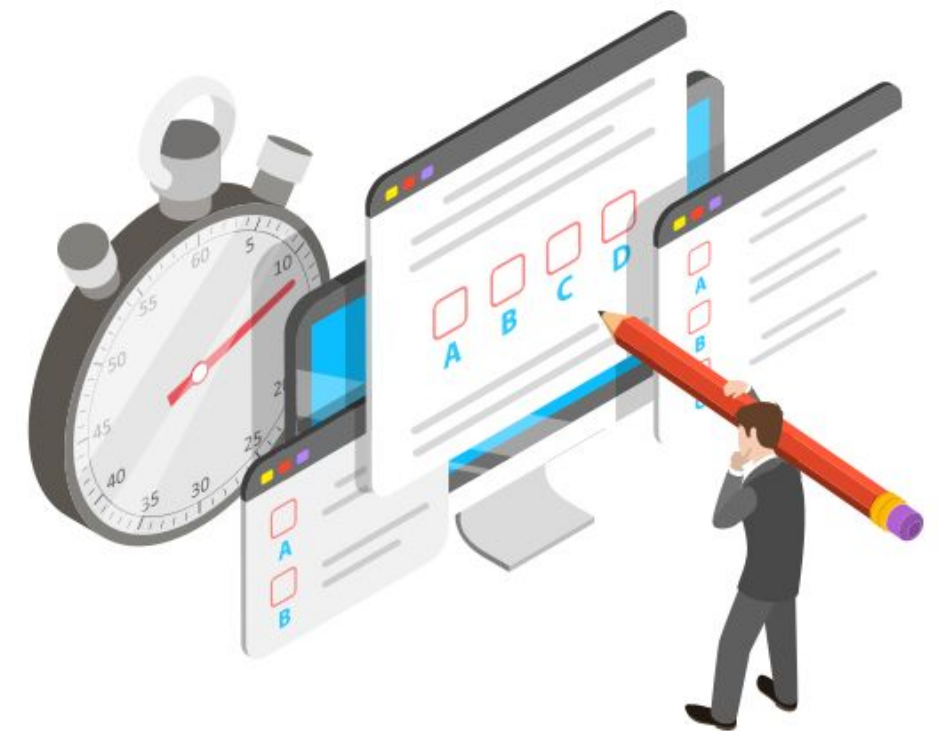


## Knowledge Check

3

NLP is technique for interpretation of all type of languages which includes \_\_\_\_\_.

- a. Human Language
- b. Assembly Language
- c. Machine Language
- d. Binary Data



The correct answer is **a.**

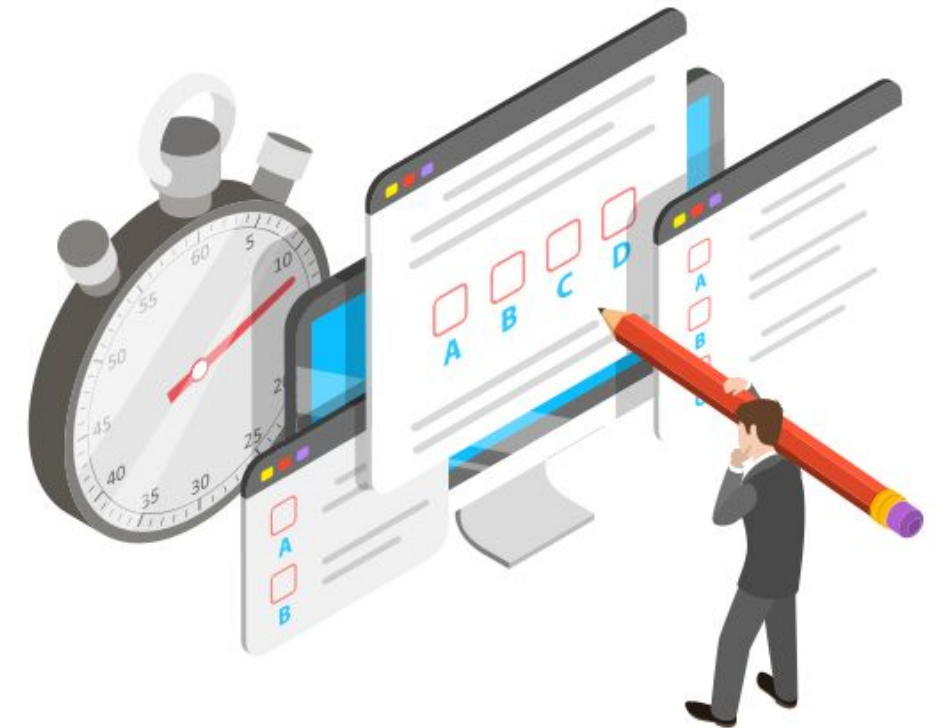
**NLP has its focus on understanding the human spoken or written language and converting that interpretation into machine understandable language.**

## Knowledge Check

4

Natural Language Processing (NLP) is a field of \_\_\_\_\_.

- a. Computer Science
- b. Artificial Intelligence
- c. Linguistics
- d. All of the above



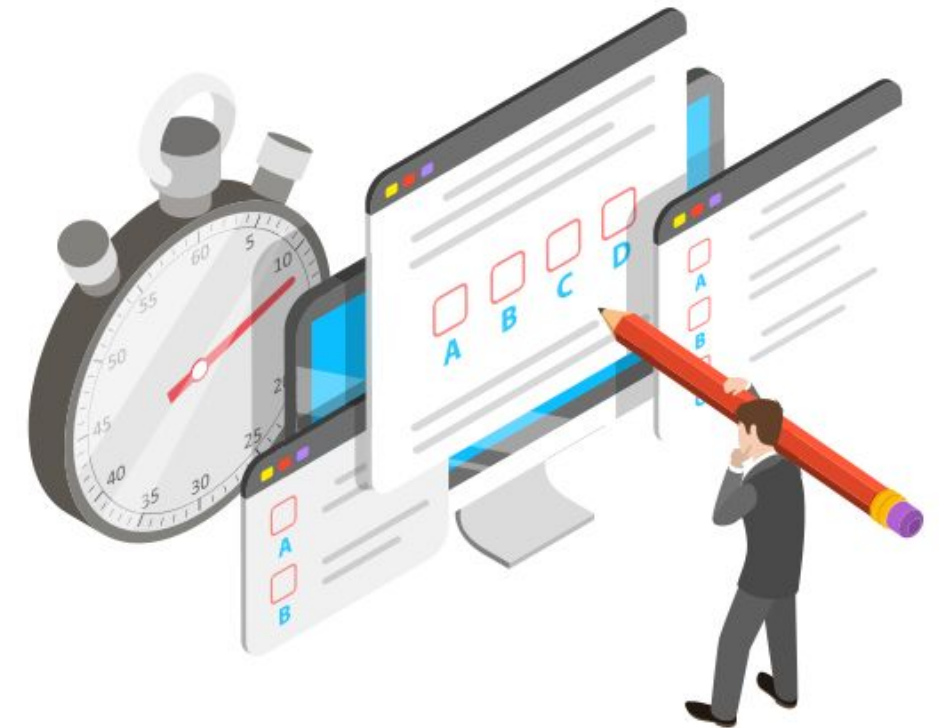


**Knowledge  
Check**

**4**

Natural Language Processing (NLP) is a field of \_\_\_\_\_.

- a. Computer Science
- b. Artificial Intelligence
- c. Linguistics
- d. All of the above



The correct answer is **d.**

**Natural Language Processing is a field of computer science, artificial intelligence, and linguistics.**

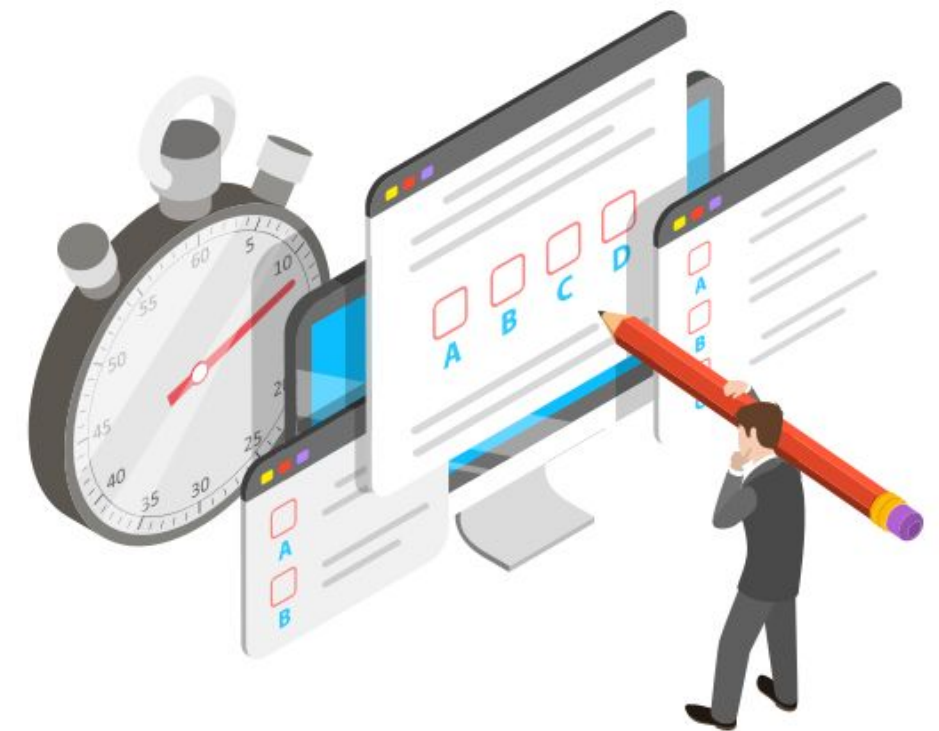
## Knowledge Check

5

Which of the following techniques can be used for the purpose of keyword normalization?

1- Lemmatization 2- Levenshtein 3- Stemming 4- POS

- a. 1 and 2
- b. 2 and 4
- c. 1 and 3
- d. 1,2, and 3



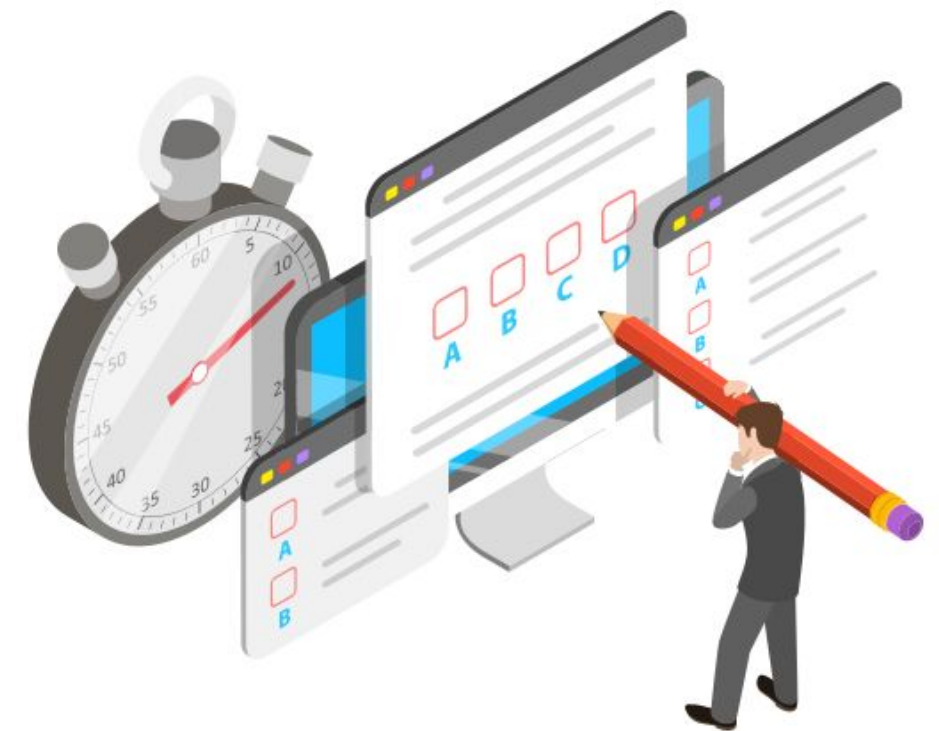
**Knowledge  
Check**

**5**

Which of the following techniques can be used for the purpose of keyword normalization?

1- Lemmatization 2- Levenshtein 3- Stemming 4- POS

- a. 1 and 2
- b. 2 and 4
- c. 1 and 3
- d. 1,2, and 3



The correct answer is **c.**

**Lemmatization and stemming are the techniques of keyword normalization.**

# Key Takeaways

You are now able to:

- Describe natural language processing and its components
- Explain the different applications of NLP
- Define and demonstrate text processing

