

# Introduction to Artificial Intelligence for Text Analytics

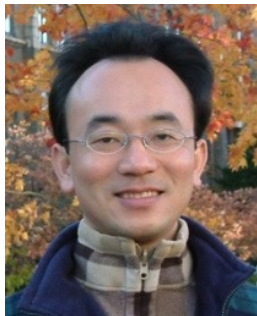
1102AITA01

MBA, IM, NTPU (M5026) (Spring 2022)

Tue 2, 3, 4 (9:10-12:00) (B8F40)



<https://meet.google.com/paj-zhji-mya>



Min-Yuh Day, Ph.D,  
Associate Professor

[Institute of Information Management, National Taipei University](https://web.ntpu.edu.tw/~myday)

<https://web.ntpu.edu.tw/~myday>





# Min-Yuh Day, Ph.D.



2020 Cohort

**Associate Professor, Information Management, NTPU**

**Visiting Scholar, IIS, Academia Sinica**

**Ph.D., Information Management, NTU**

**Director, Intelligent Financial Innovation Technology, IFIT Lab, IM, NTPU**

**Artificial Intelligence, Financial Technology, Big Data Analytics,  
Data Mining and Text Mining, Electronic Commerce**



2020 Cohort



Accredited  
Educator



Solutions  
Architect  
Associate



Cloud  
Practitioner



# Course Syllabus

## National Taipei University

### Academic Year 110, 2<sup>nd</sup> Semester (Spring 2022)

- **Course Title: Artificial Intelligence for Text Analytics**
- **Instructor: Min-Yuh Day**
- **Course Class: MBA, IM, NTPU (3 Credits, Elective)**
- **Details**
  - **In-Class and Distance Learning EMI Course (3 Credits, Elective, One Semester) (M5026)**
- **Time & Place: Tue, 2, 3, 4, (9:10-12:00) (B8F40)**
- **Google Meet: <https://meet.google.com/paj-zhhj-mya>**



<https://meet.google.com/paj-zhhj-mya>



# Course Objectives

1. Understand the **fundamental concepts and research issues of Artificial Intelligence for Text Analytics**.
2. Equip with Hands-on practices of **Artificial Intelligence for Text Analytics**.
3. Conduct **information systems research in the context of Artificial Intelligence for Text Analytics**.

# Course Outline

- This course introduces the **fundamental concepts, research issues, and hands-on practices of Artificial Intelligence for Text Analytics.**
- Topics include:
  1. Introduction to Introduction to Artificial Intelligence for Text Analytics
  2. Foundations of Text Analytics: Natural Language Processing (NLP)
  3. Python for Natural Language Processing
  4. Natural Language Processing with Transformers
  5. Text Classification and Sentiment Analysis
  6. Multilingual Named Entity Recognition (NER), Text Similarity and Clustering
  7. Text Summarization and Topic Models
  8. Text Generation
  9. Question Answering and Dialogue Systems
  10. Deep Learning, Transfer Learning, Zero-Shot, and Few-Shot Learning for Text Analytics
  11. Case Study on Artificial Intelligence for Text Analytics

# Core Competence

- **Exploring new knowledge in information technology, system development and application 80 %**
- **Internet marketing planning ability 10 %**
- **Thesis writing and independent research skills 10 %**

# Four Fundamental Qualities

- **Professionalism**
  - **Creative thinking and Problem-solving 40 %**
  - **Comprehensive Integration 40 %**
- **Interpersonal Relationship**
  - **Communication and Coordination 10 %**
  - **Teamwork 5 %**
- **Ethics**
  - **Honesty and Integrity 0 %**
  - **Self-Esteem and Self-reflection 0 %**
- **International Vision**
  - **Caring for Diversity 0 %**
  - **Interdisciplinary Vision 5 %**

# College Learning Goals

- **Ethics/Corporate Social Responsibility**
- **Global Knowledge/Awareness**
- **Communication**
- **Analytical and Critical Thinking**



# Department Learning Goals

- **Information Technologies and System Development Capabilities**
- **Internet Marketing Management Capabilities**
- **Research capabilities**

# Syllabus

**Week Date Subject/Topics**

- 1 2022/02/22 Introduction to Artificial Intelligence for Text Analytics**
- 2 2022/03/01 Foundations of Text Analytics:  
Natural Language Processing (NLP)**
- 3 2022/03/08 Python for Natural Language Processing**
- 4 2022/03/15 Natural Language Processing with Transformers**
- 5 2022/03/22 Case Study on Artificial Intelligence for Text Analytics I**
- 6 2022/03/29 Text Classification and Sentiment Analysis**

# Syllabus

Week	Date	Subject/Topics
7	2022/04/05	Tomb-Sweeping Day (Holiday, No Classes)
8	2022/04/12	Midterm Project Report
9	2022/04/19	Multilingual Named Entity Recognition (NER), Text Similarity and Clustering
10	2022/04/26	Text Summarization and Topic Models
11	2022/05/03	Text Generation
12	2022/05/10	Case Study on Artificial Intelligence for Text Analytics II

# Syllabus

**Week Date Subject/Topics**

**13 2022/05/17 Question Answering and Dialogue Systems**

**14 2022/05/24 Deep Learning, Transfer Learning,  
Zero-Shot, and Few-Shot Learning for Text Analytics**

**15 2022/05/31 Final Project Report I**

**16 2022/06/07 Final Project Report II**

**17 2022/06/14 Self-learning**

**18 2022/06/21 Self-learning**

# Teaching Methods and Activities

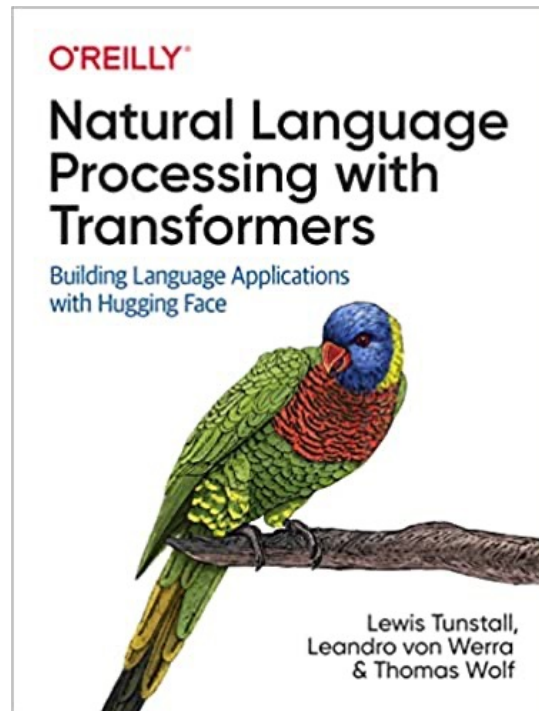
- **Lecture**
- **Discussion**
- **Practicum**

# Evaluation Methods

- **Individual Presentation 60 %**
- **Group Presentation 10 %**
- **Case Report 10 %**
- **Class Participation 10 %**
- **Assignment 10 %**

# Required Texts

- **Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022),  
Natural Language Processing with Transformers:  
Building Language Applications with Hugging Face,  
O'Reilly Media.**



# Reference Books

- **Denis Rothman (2021), Transformers for Natural Language Processing: Build innovative deep neural network architectures for NLP with Python, PyTorch, TensorFlow, BERT, RoBERTa, and more, Packt Publishing.**
- **Savaş Yıldırım and Meysam Asgari-Chenaghlu (2021), Mastering Transformers: Build state-of-the-art models from scratch with advanced natural language processing techniques, Packt Publishing.**
- **Sudharsan Ravichandiran (2021), Getting Started with Google BERT: Build and train state-of-the-art natural language processing models using BERT, Packt Publishing.**
- **Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta (2020), Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems, O'Reilly Media.**

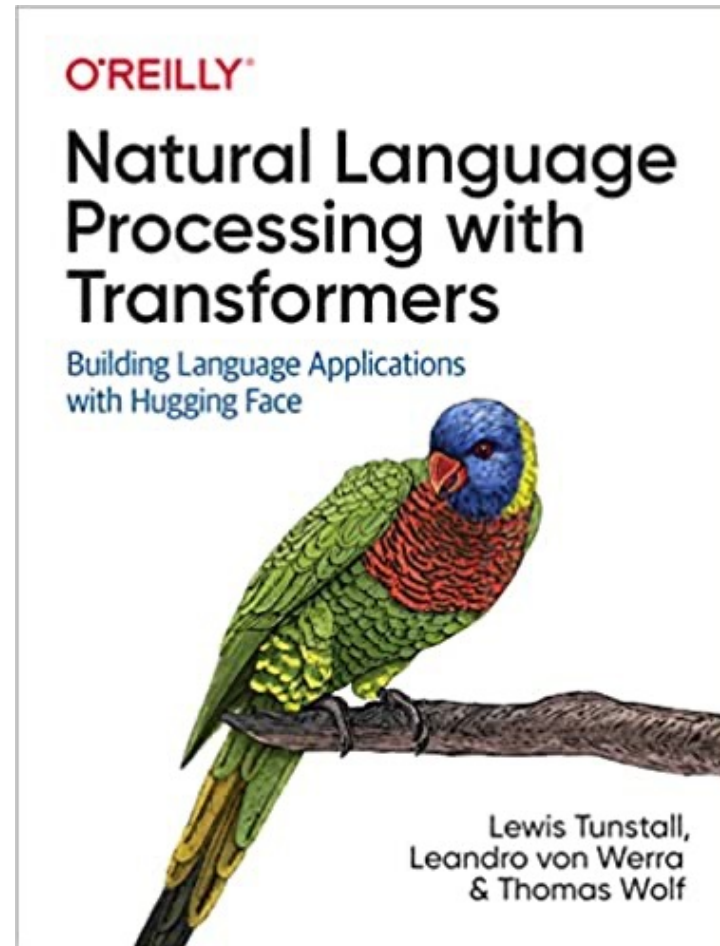


# Other References

- **Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress.**
- **Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python: Enabling Language-Aware Data Products with Machine Learning, O'Reilly.**
- **Charu C. Aggarwal (2018), Machine Learning for Text, Springer.**
- **Gabe Ignatow and Rada F. Mihalcea (2017), An Introduction to Text Mining: Research Design, Data Collection, and Analysis, SAGE Publications.**
- **Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media.**
- **Frederick Kaefer and Paul Kaefer (2020), Introduction to Python Programming for Business and Social Science Applications, SAGE Publications**
- **Vic Anand, Khrystyna Bochkay, and Roman Chychyla (2020), Using Python for Text Analysis in Accounting Research, Now Publishers.**

Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022),  
**Natural Language Processing with Transformers:**

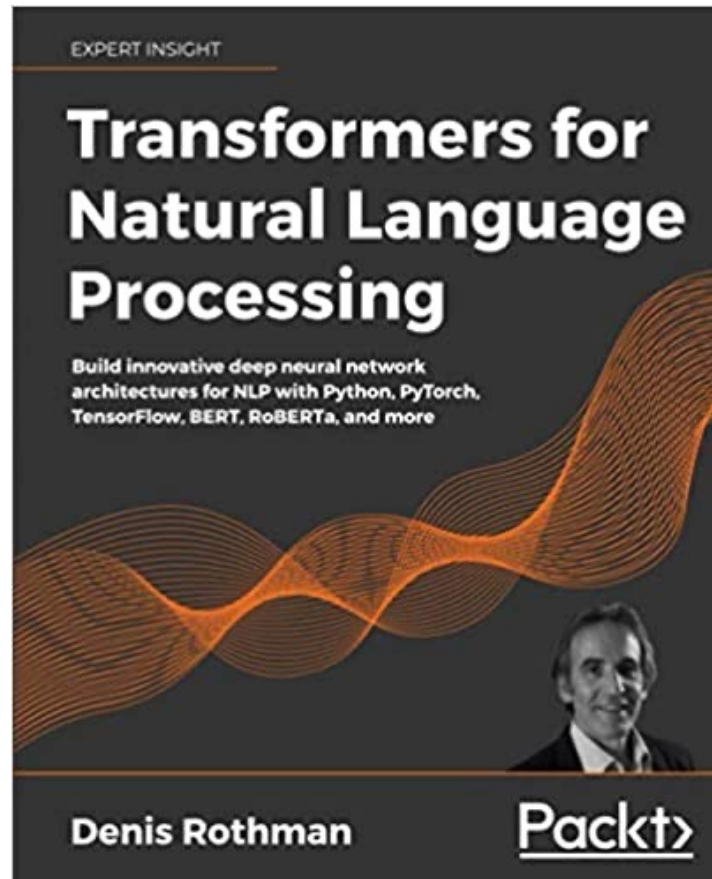
Building Language Applications with Hugging Face,  
O'Reilly Media.



Denis Rothman (2021),

# Transformers for Natural Language Processing:

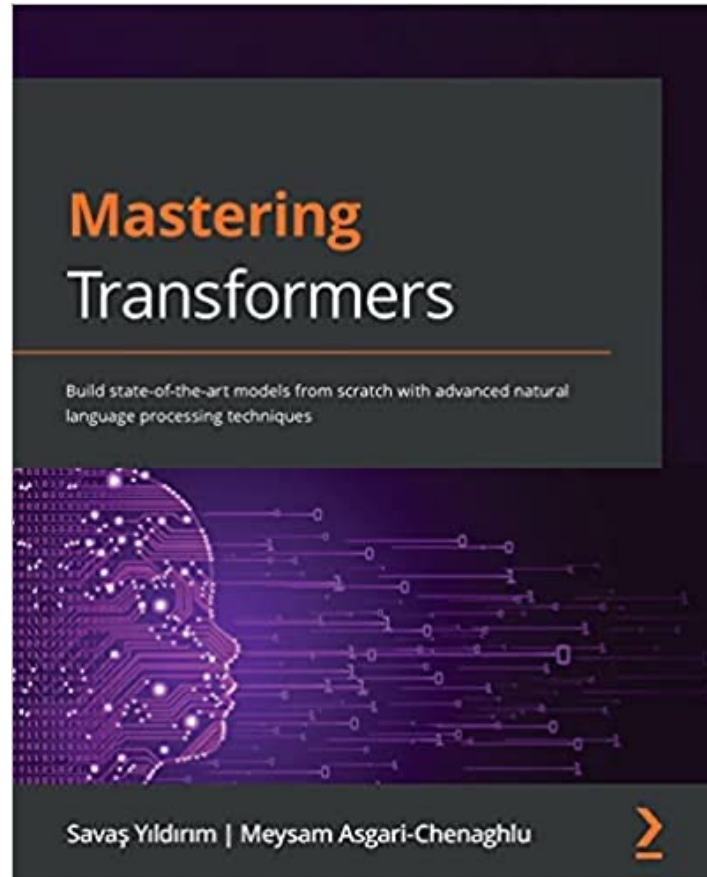
Build innovative deep neural network architectures for NLP with Python, PyTorch, TensorFlow, BERT, RoBERTa, and more,  
Packt Publishing.



Savaş Yıldırım and Meysam Asgari-Chenaghlu (2021),

## Mastering Transformers:

Build state-of-the-art models from scratch with advanced natural language processing techniques,  
Packt Publishing.

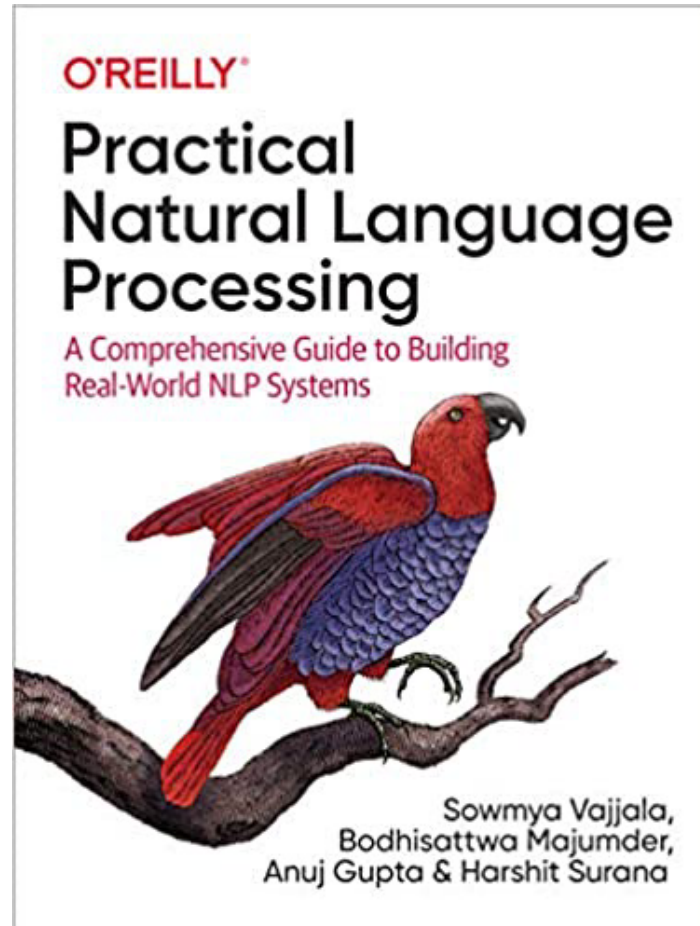


Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta (2020),

# Practical Natural Language Processing:

A Comprehensive Guide to Building Real-World NLP Systems,

O'Reilly Media.



O'REILLY®

# Practical Natural Language Processing

A Comprehensive Guide to Building Real-World NLP Systems



Sowmya Vajjala,  
Bodhisattwa Majumder,  
Anuj Gupta & Harshit Surana

## FOUNDATIONS

Covered in  
Chapters 1 to 3



ML for NLP



NLP Pipelines



Data  
Gathering



Multilingual  
NLP



Text  
Representation

## CORE TASKS

Covered in  
Chapters 3 to 7



Text  
Classification



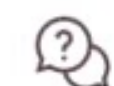
Information  
Extraction



Conversational  
Agents



Information  
Retrieval



Question  
Answering

## GENERAL APPLICATIONS

Covered in  
Chapters 4 to 7



Spam  
Classification



Calendar Event  
Extractor



Personal  
Assistants



Search  
Engines

**JEOPARDY!**

Jeopardy!

## INDUSTRY SPECIFIC

Covered in  
Chapters 8 to 10



Social Media  
Analysis



Retail Data  
Extraction



Health Records  
Analysis



Financial  
Analysis



Legal Entity  
Extraction

## AI PROJECT PLAYBOOK

Covered in  
Chapters 2 & 11



Project  
Processes



Best  
Practices



Model  
Iterations



MLOps



AI Teams  
& Hiring

# Artificial Intelligence (AI)

# **Text Analytics**

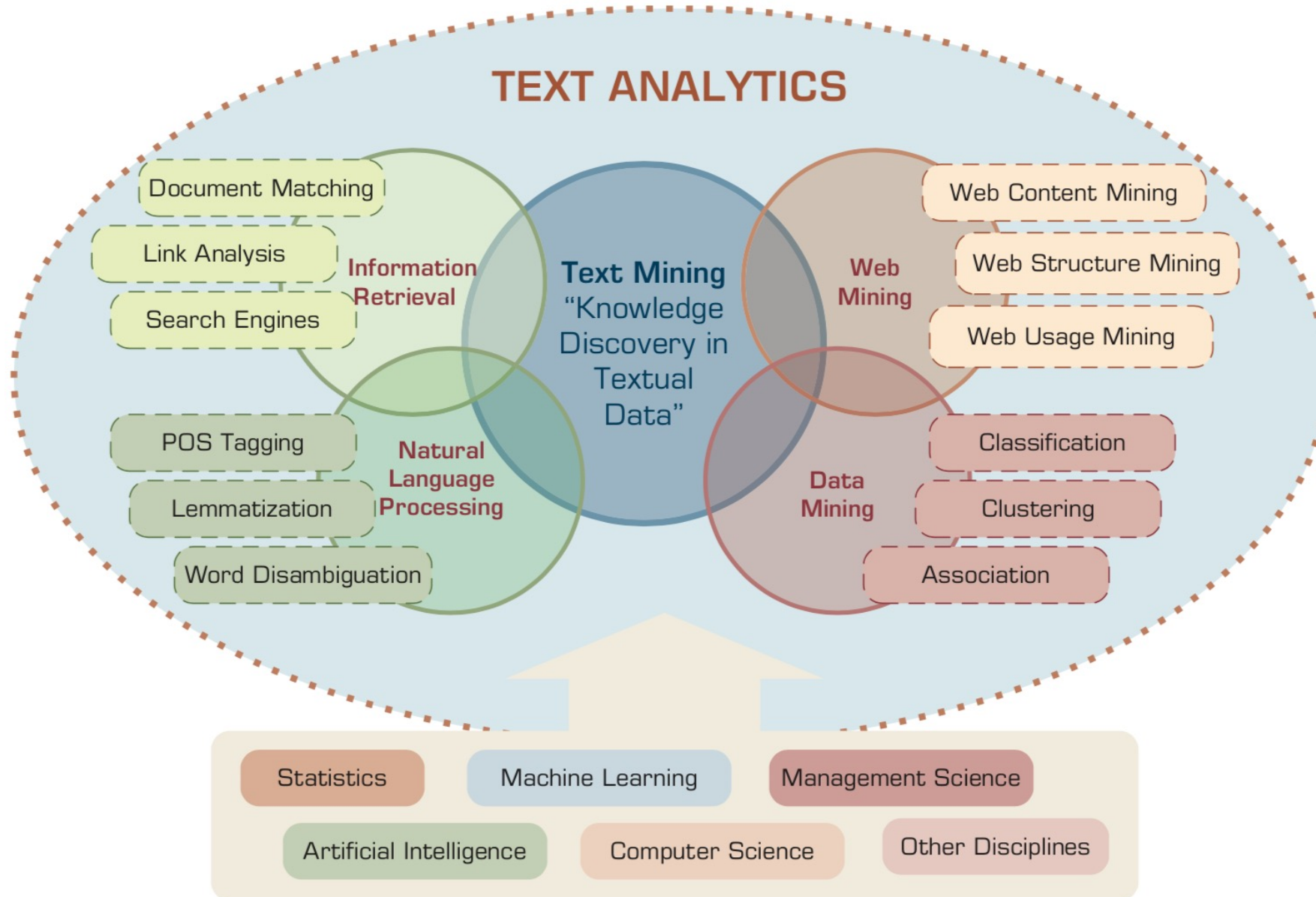
## **(TA)**



# Text Mining (TM)

# Natural Language Processing (NLP)

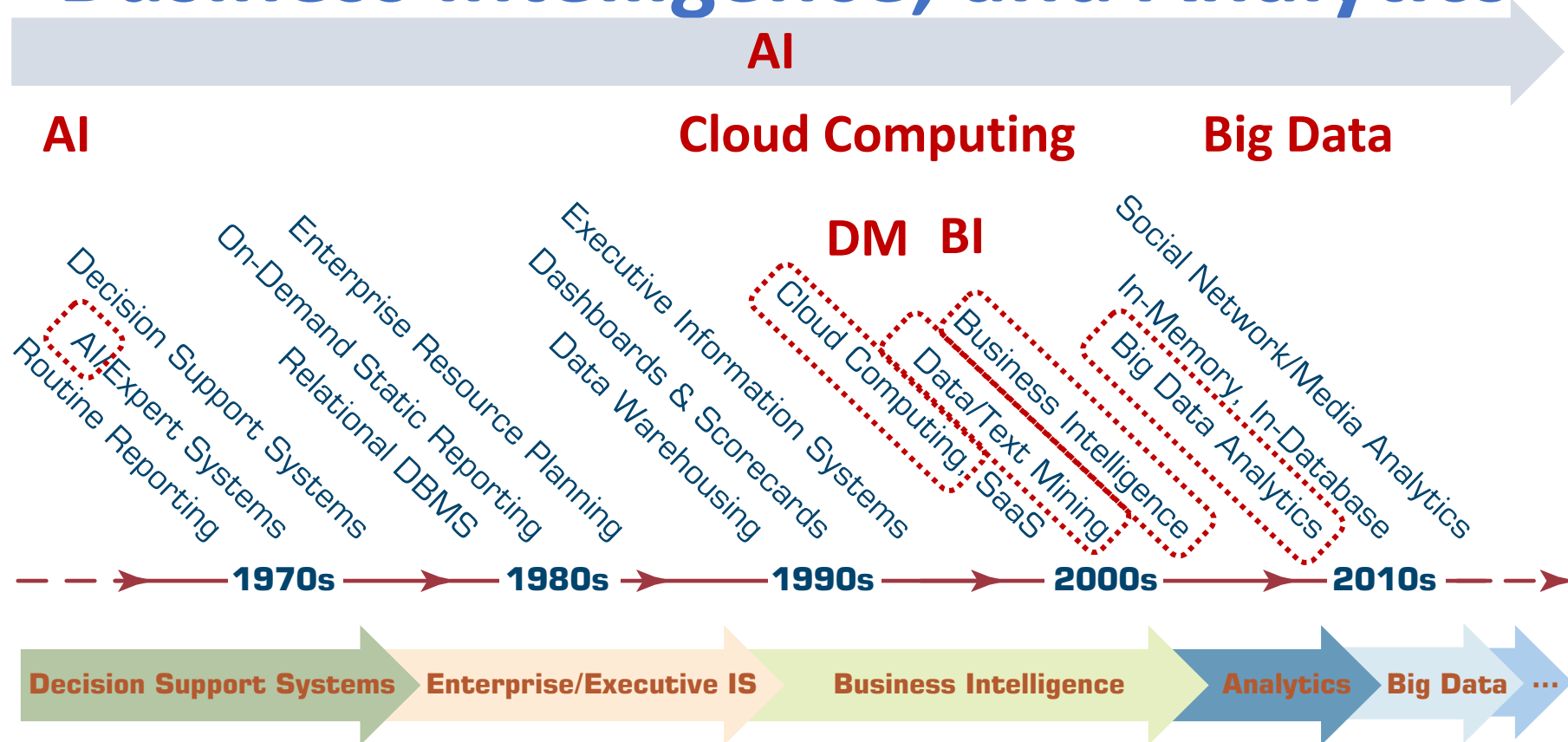
# Text Analytics and Text Mining



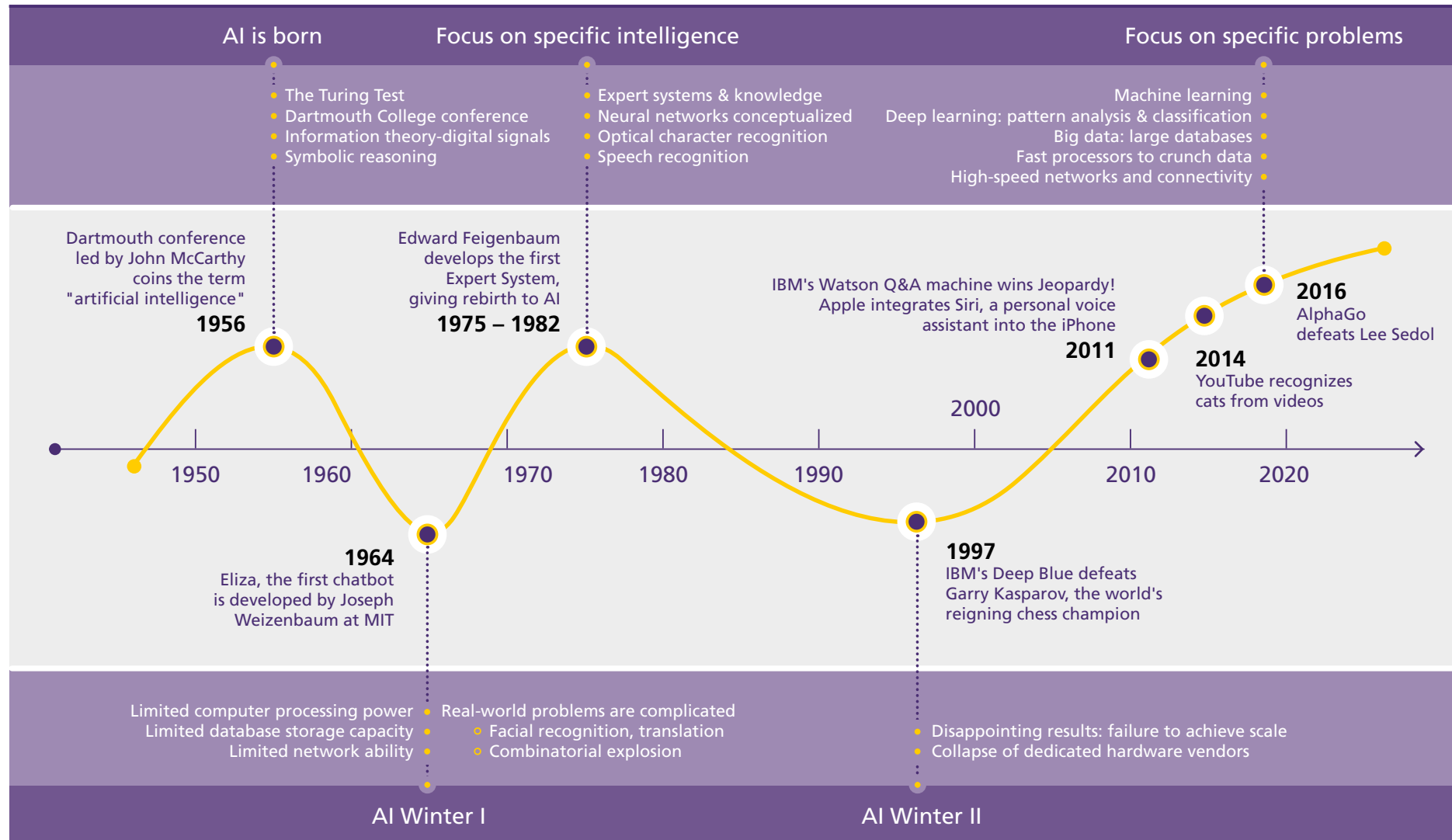
# Artificial Intelligence (AI)

# AI, Big Data, Cloud Computing

## Evolution of Decision Support, Business Intelligence, and Analytics



# The Rise of AI



# **Definition of Artificial Intelligence (A.I.)**

# Artificial Intelligence

**“... the science and  
engineering  
of  
making  
intelligent machines”**

**(John McCarthy, 1955)**



# Artificial Intelligence

**“... technology that  
thinks and acts  
like humans”**

# Artificial Intelligence

**“... intelligence  
exhibited by machines  
or software”**

# 4 Approaches of AI

<b>Thinking Humanly</b>	<b>Thinking Rationally</b>
<b>Acting Humanly</b>	<b>Acting Rationally</b>

# 4 Approaches of AI

<p><b>2.</b> <b>Thinking Humanly: The Cognitive Modeling Approach</b></p>	<p><b>3.</b> <b>Thinking Rationally: The “Laws of Thought” Approach</b></p>
<p><b>1.</b> <b>Acting Humanly: The Turing Test Approach</b> (1950)</p>	<p><b>4.</b> <b>Acting Rationally: The Rational Agent Approach</b></p>

# AI Acting Humanly: The Turing Test Approach (Alan Turing, 1950)

- Knowledge Representation
- Automated Reasoning
- Machine Learning (ML)
  - Deep Learning (DL)
- Computer Vision (Image, Video)
- Natural Language Processing (NLP)
- Robotics

# **Text Analytics**

## **(TA)**

# Text Analytics

- **Text Analytics =**  
**Information Retrieval +**  
**Information Extraction +**  
**Data Mining +**  
**Web Mining**
- **Text Analytics =**  
**Information Retrieval +**  
**Text Mining**

# Text Mining

- **Text Data Mining**
- **Knowledge Discovery in Textual Databases**



# Application Areas of Text Mining

- **Information extraction**
- **Topic tracking**
- **Summarization**
- **Categorization**
- **Clustering**
- **Concept linking**
- **Question answering**

# Emotions



Love

Anger

Joy

Sadness

Surprise

Fear



## Example of Opinion: review segment on iPhone



**“I bought an iPhone a few days ago.**

**It was such a nice phone.**

**The touch screen was really cool.**

**The voice quality was clear too.**

**However, my mother was mad with me as I did not tell her before I bought it.**

**She also thought the phone was too expensive, and wanted me to return it to the shop. ... ”**

# Example of Opinion: review segment on iPhone

“(1) I bought an iPhone a few days ago.

(2) **It was such a nice phone.**

(3) **The touch screen was really cool.**

(4) **The voice quality was clear too.**

(5) **However, my mother was mad with me as I did not tell her before I bought it.**

(6) **She also thought the phone was too expensive, and wanted me to return it to the shop. ...”**

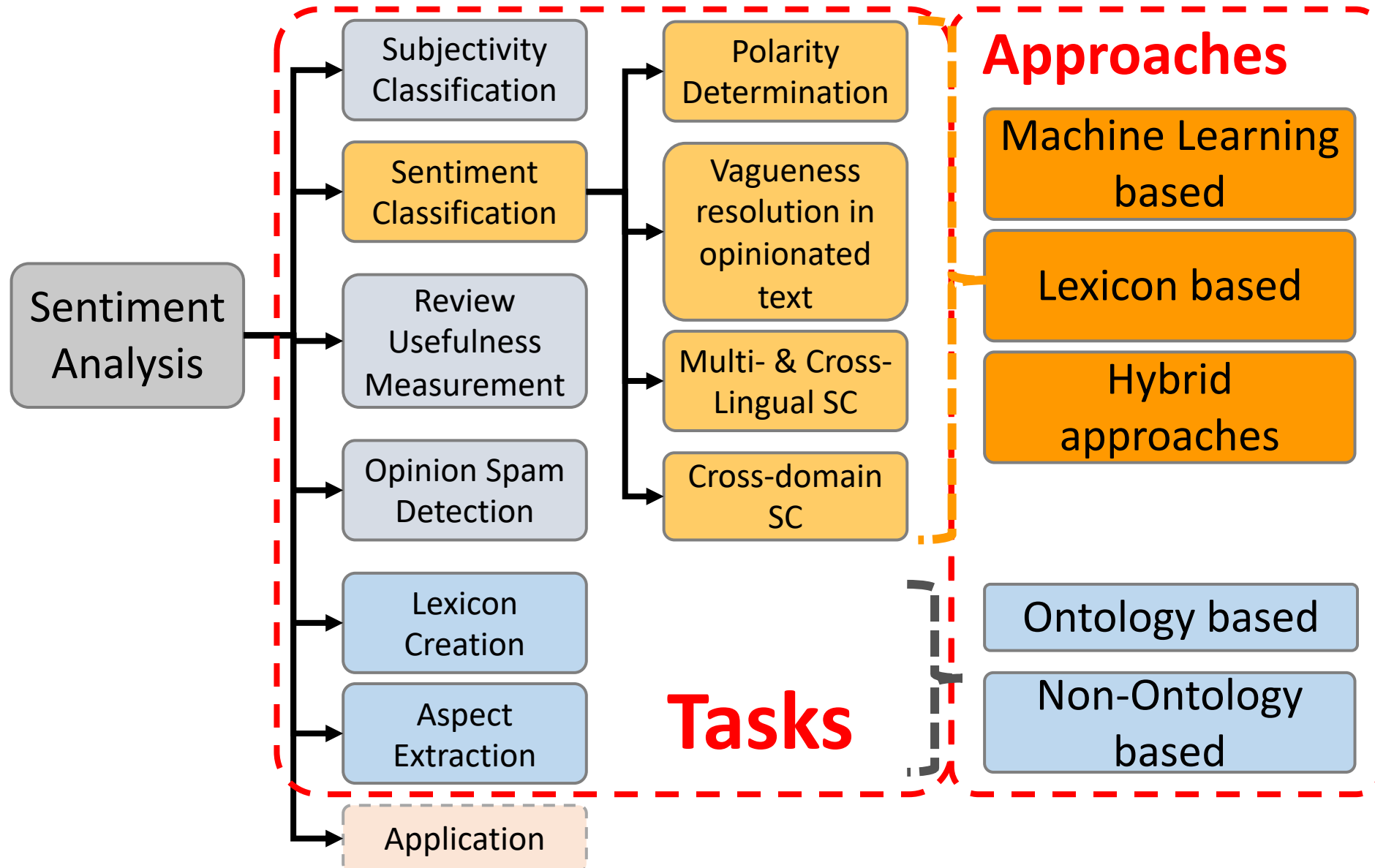


**+Positive  
Opinion**

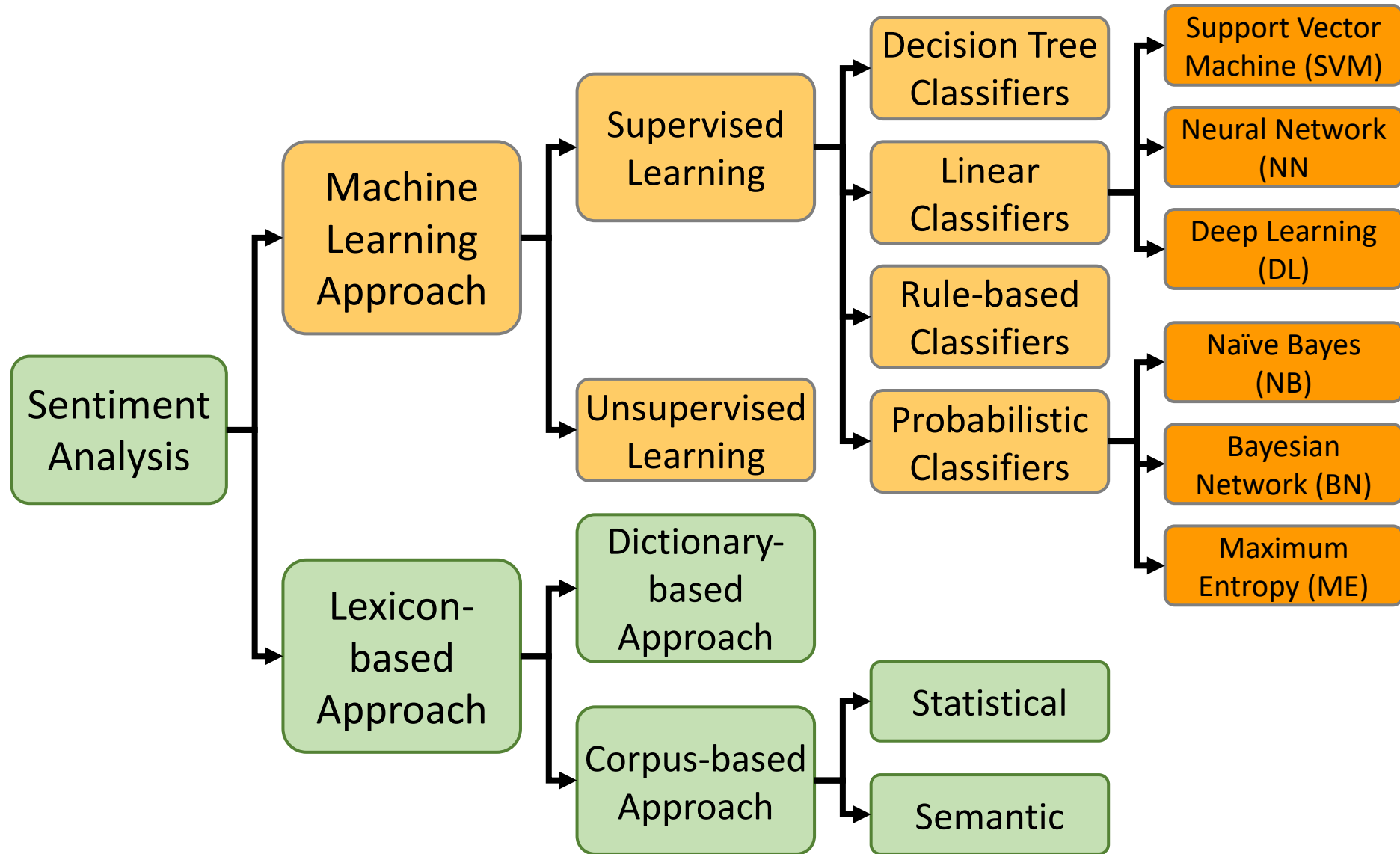


**-Negative  
Opinion**

# Sentiment Analysis



# Sentiment Classification Techniques



# Text Mining Technologies

# **Text Mining (TM)**

## **Natural Language Processing (NLP)**



Text mining

Text Data Mining

Intelligent Text Analysis

Knowledge-Discovery in Text (KDT)

# **Text Mining**

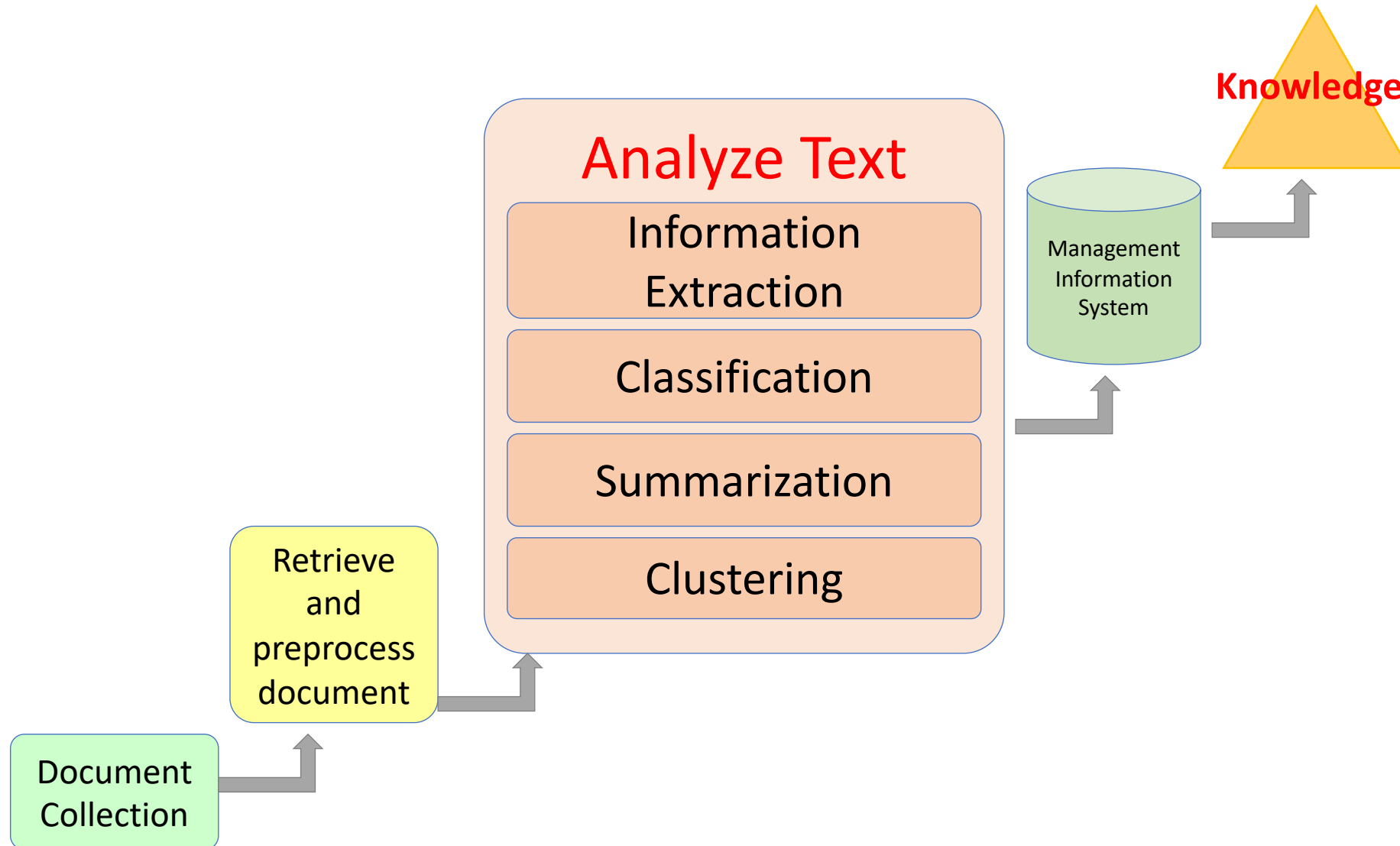
## **(text data mining)**

**the process of  
deriving  
high-quality information  
from text**

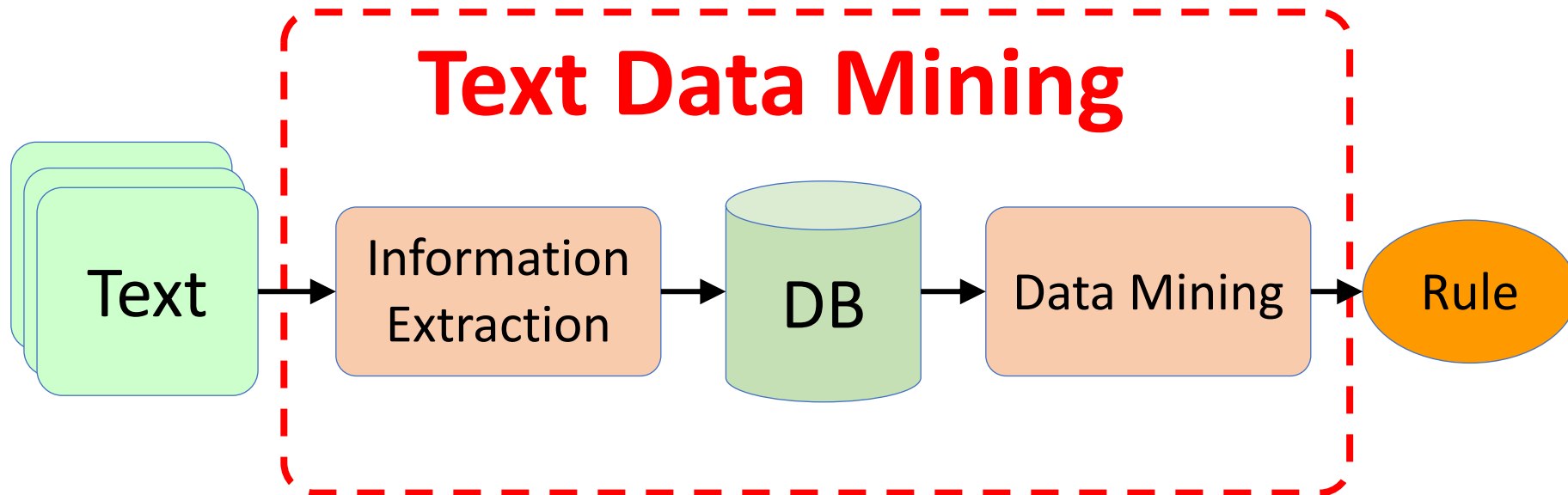
**Text Mining:**  
the process of extracting  
interesting and non-trivial  
information and knowledge  
from unstructured text.

**Text Mining:**  
discovery by computer of  
**new, previously**  
**unknown information,**  
**by automatically**  
**extracting information**  
**from different written resources.**

# An example of Text Mining



# Overview of Information Extraction based Text Mining Framework



# Natural Language Processing (NLP)

- **Natural language processing (NLP) is an important component of text mining and is a subfield of artificial intelligence and computational linguistics.**

# Natural Language Processing (NLP)

- **Part-of-speech tagging**
- **Text segmentation**
- **Word sense disambiguation**
- **Syntactic ambiguity**
- **Imperfect or irregular input**
- **Speech acts**

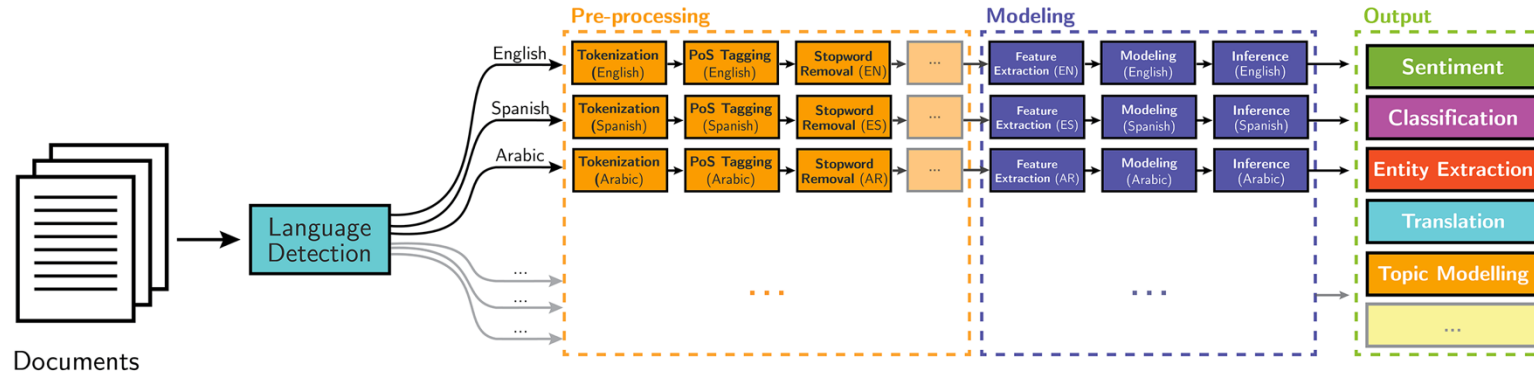


# NLP Tasks

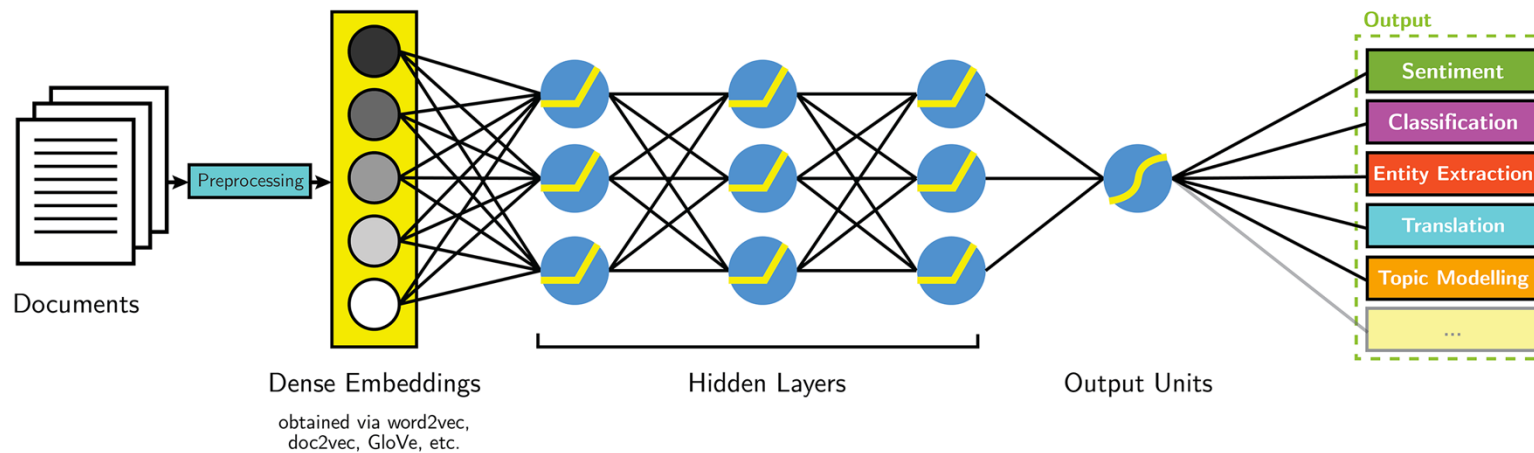
- **Question answering**
- **Automatic summarization**
- **Natural language generation**
- **Natural language understanding**
- **Machine translation**
- **Foreign language reading**
- **Foreign language writing.**
- **Speech recognition**
- **Text-to-speech**
- **Text proofing**
- **Optical character recognition**

# NLP

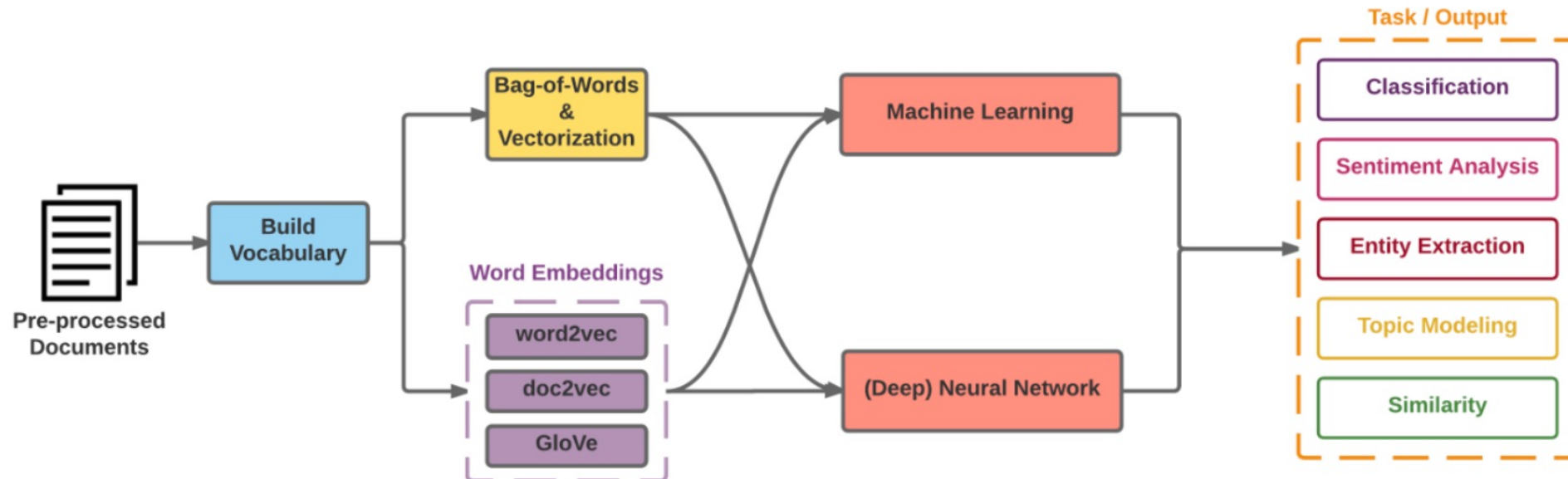
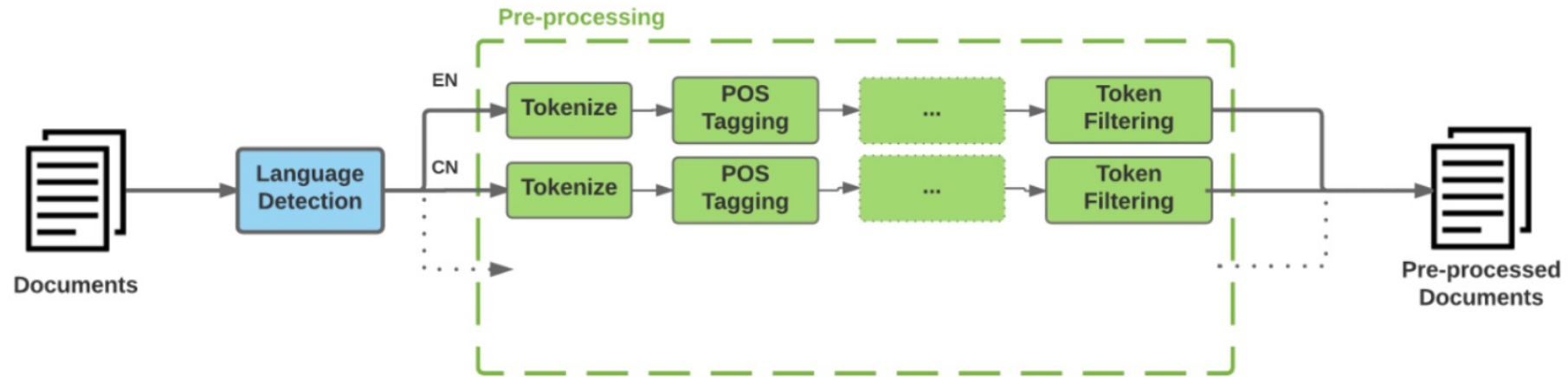
## Classical NLP



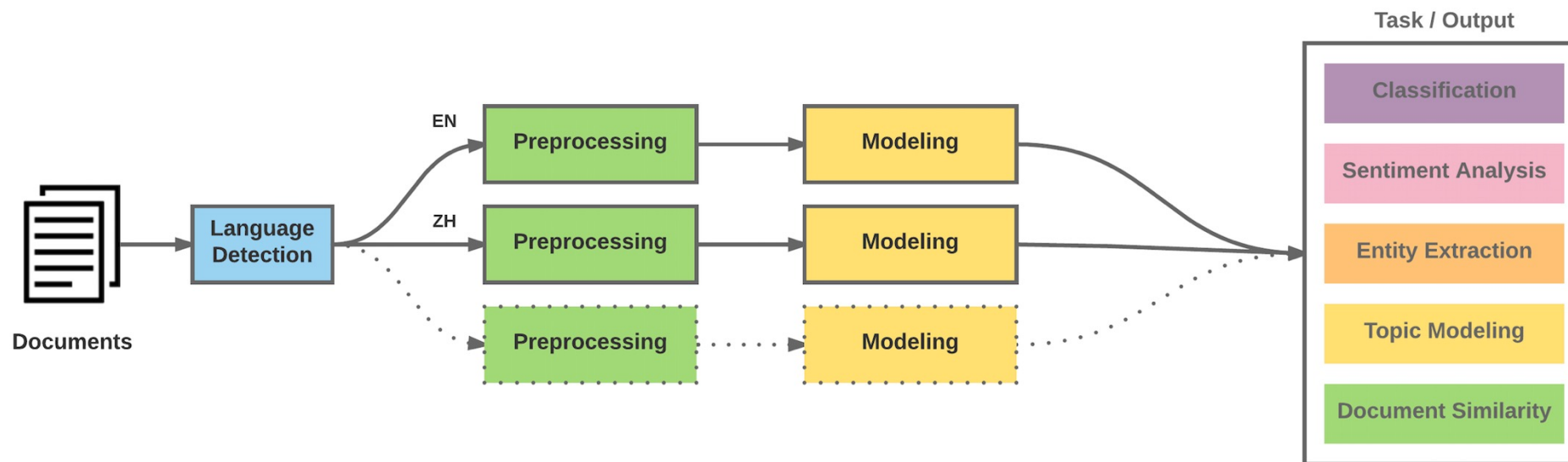
## Deep Learning-based NLP



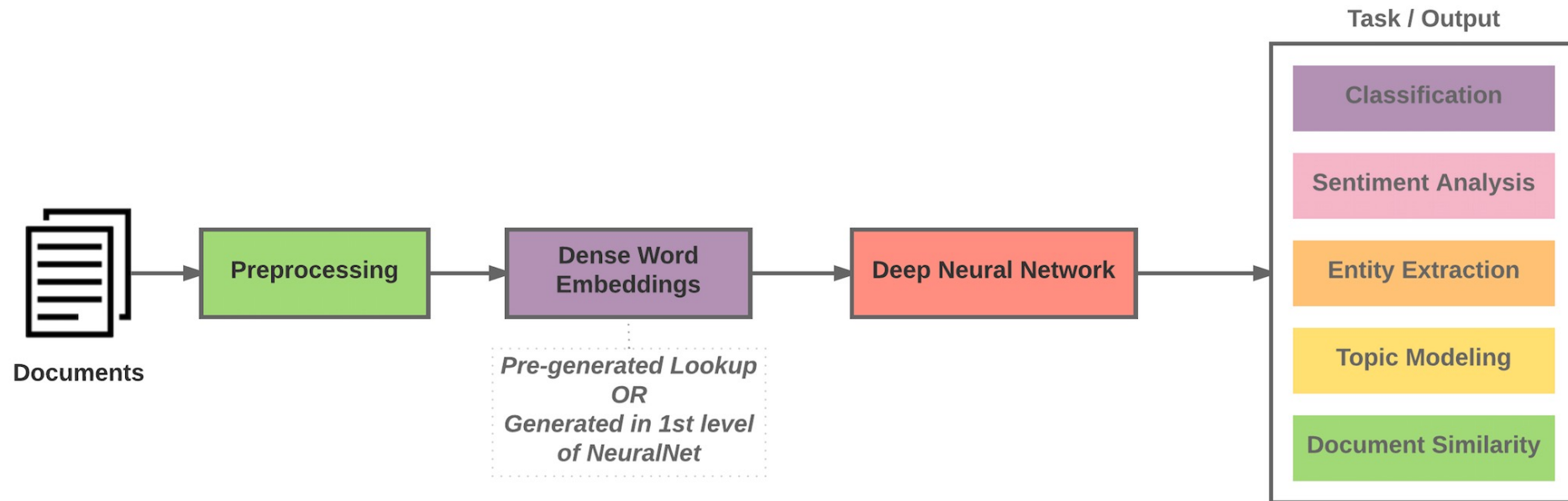
# Modern NLP Pipeline



# Modern NLP Pipeline

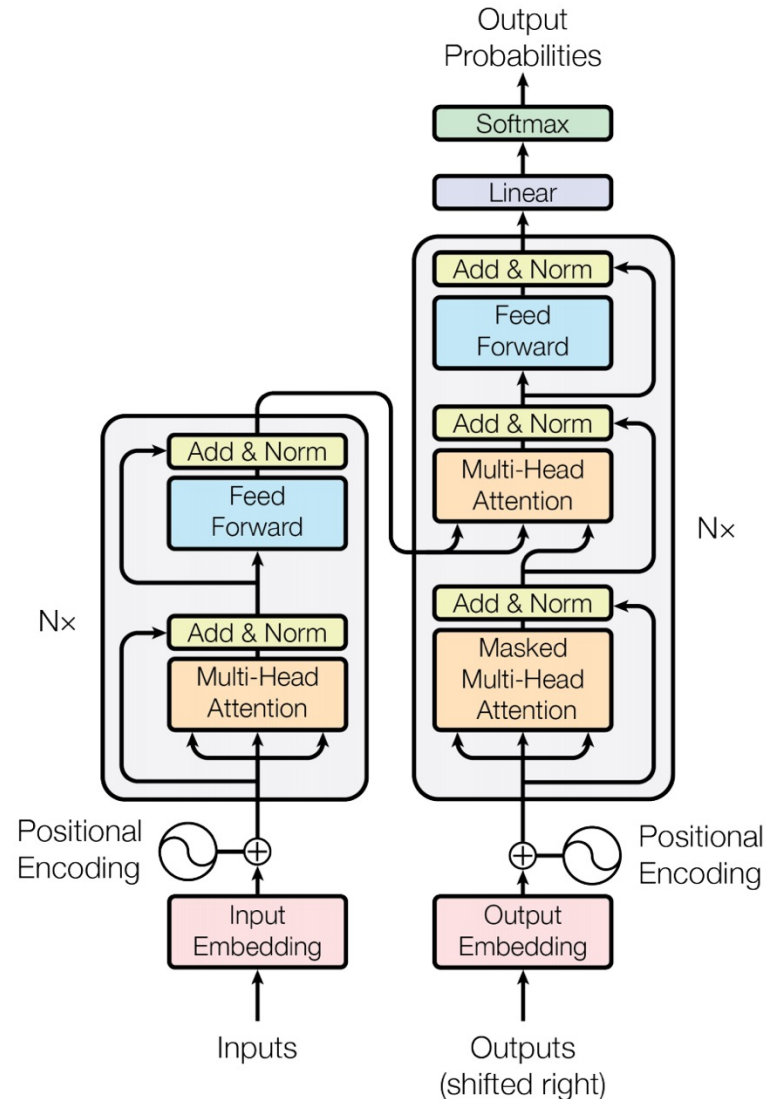


# Deep Learning NLP



# Transformer (Attention is All You Need)

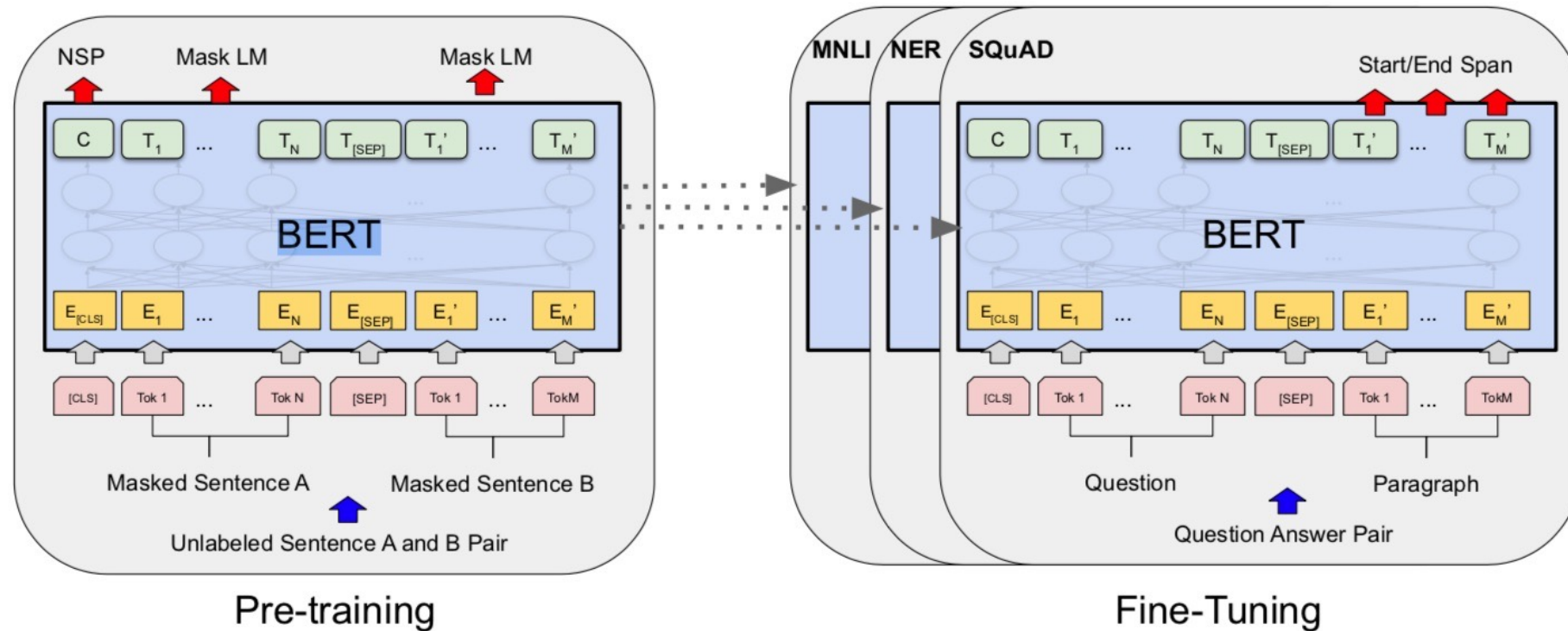
(Vaswani et al., 2017)



# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT (Bidirectional Encoder Representations from Transformers)

Overall pre-training and fine-tuning procedures for BERT



# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

**BERT: Pre-training of Deep Bidirectional Transformers for  
Language Understanding**

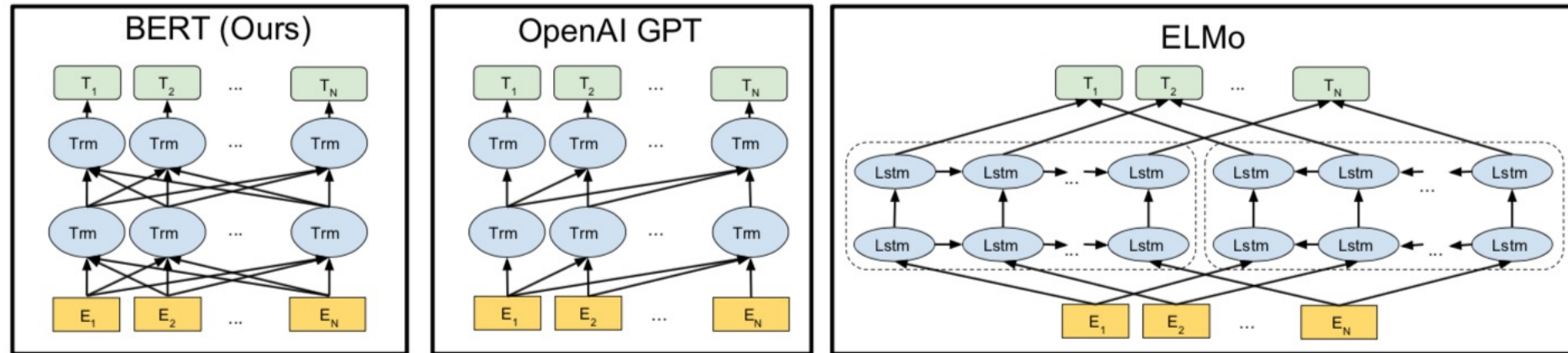
**Jacob Devlin   Ming-Wei Chang   Kenton Lee   Kristina Toutanova**  
Google AI Language

`{jacobdevlin, mingweichang, kentonl, kristout}@google.com`



# BERT

## Bidirectional Encoder Representations from Transformers



### Pre-training model architectures

**BERT** uses a bidirectional Transformer.

**OpenAI GPT** uses a left-to-right Transformer.

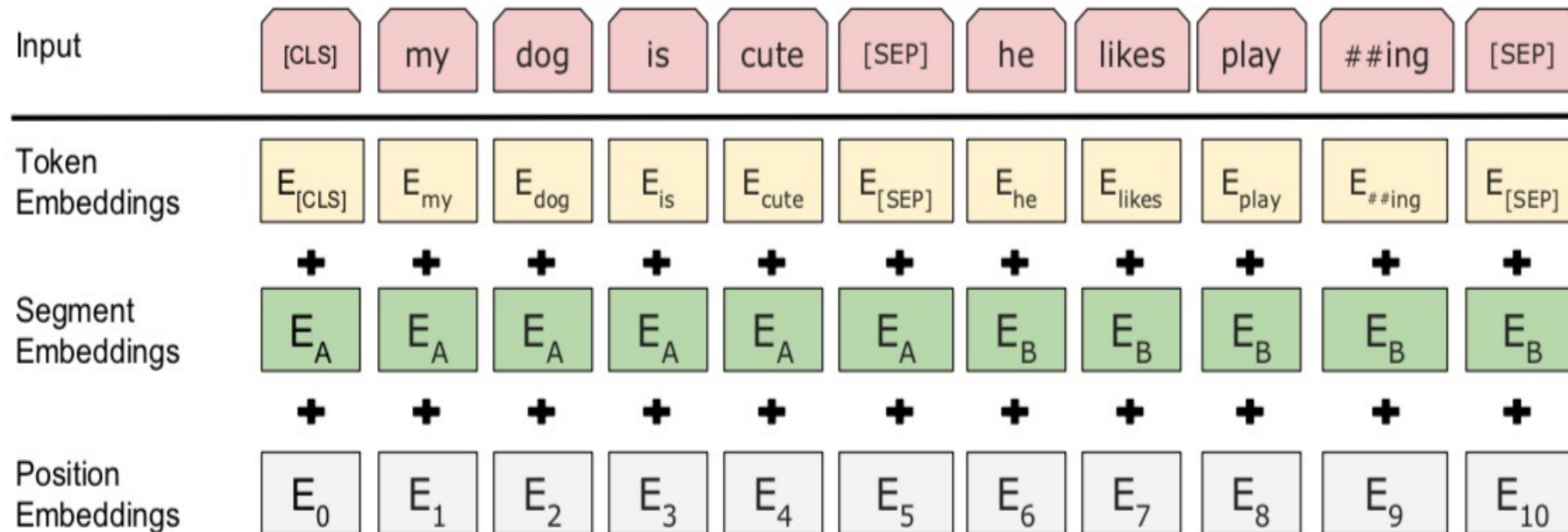
**ELMo** uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks.

Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

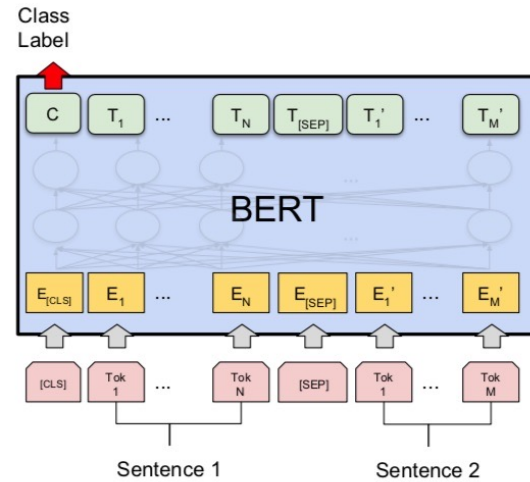
BERT (Bidirectional Encoder Representations from Transformers)

## BERT input representation

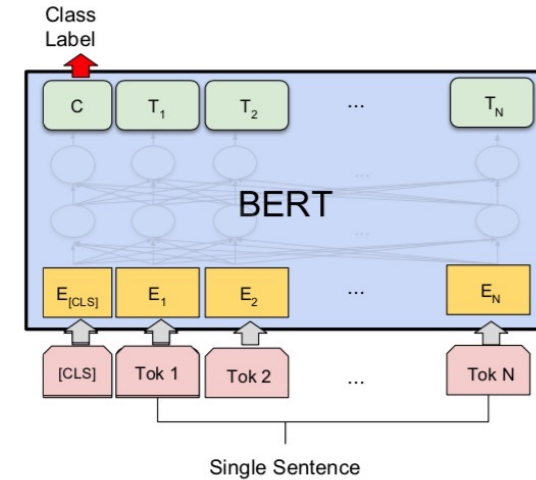


The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

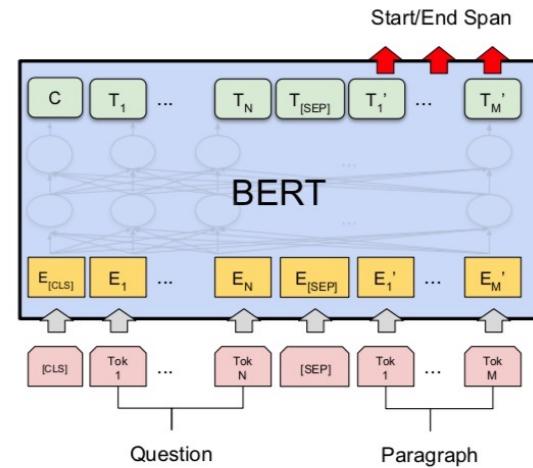
# Fine-tuning BERT on NLP Tasks



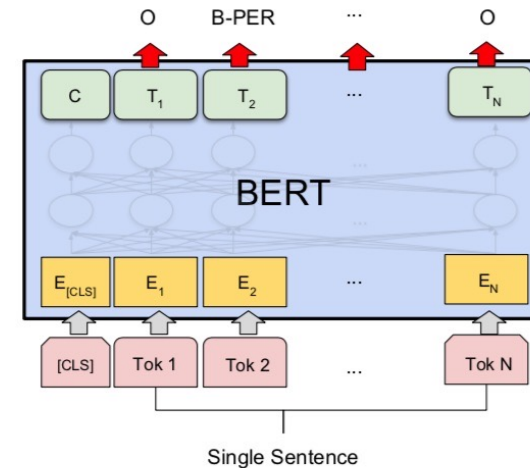
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1

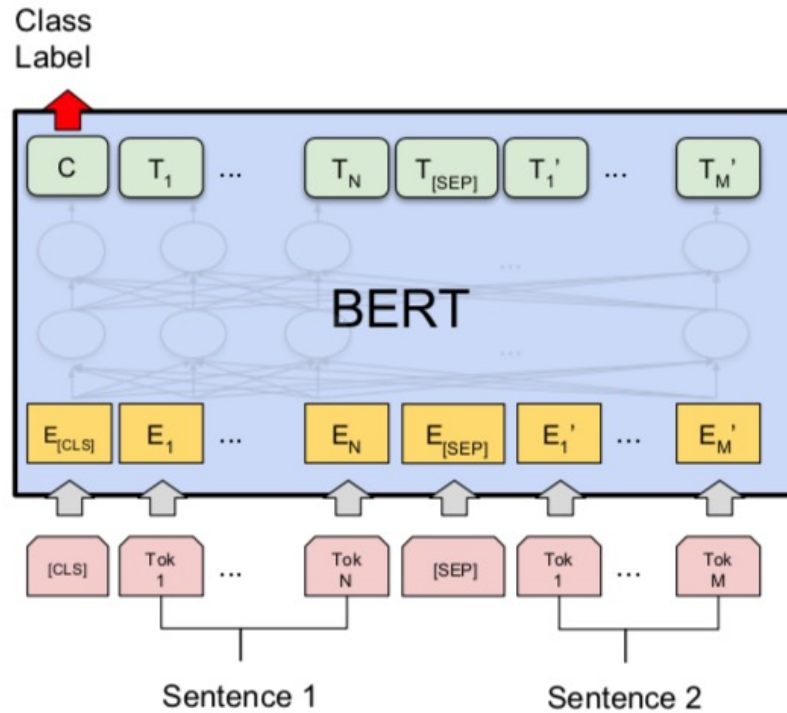


(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

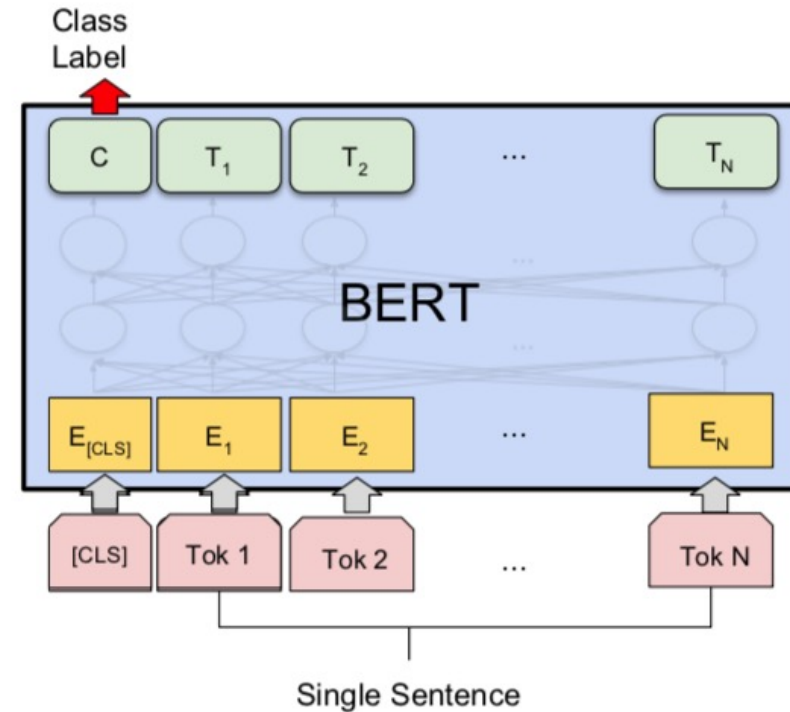
Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805

# BERT Sequence-level tasks

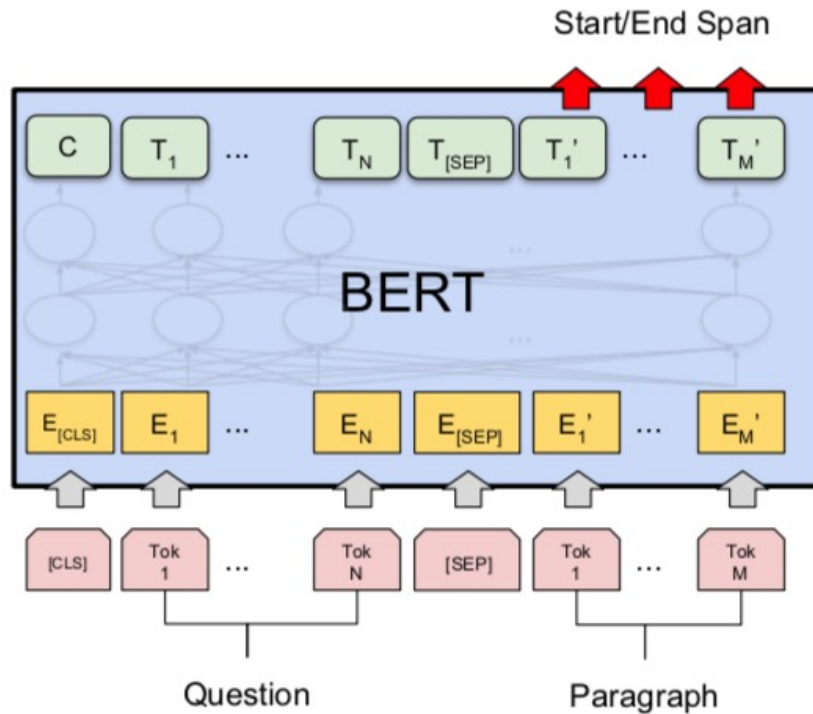


(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG

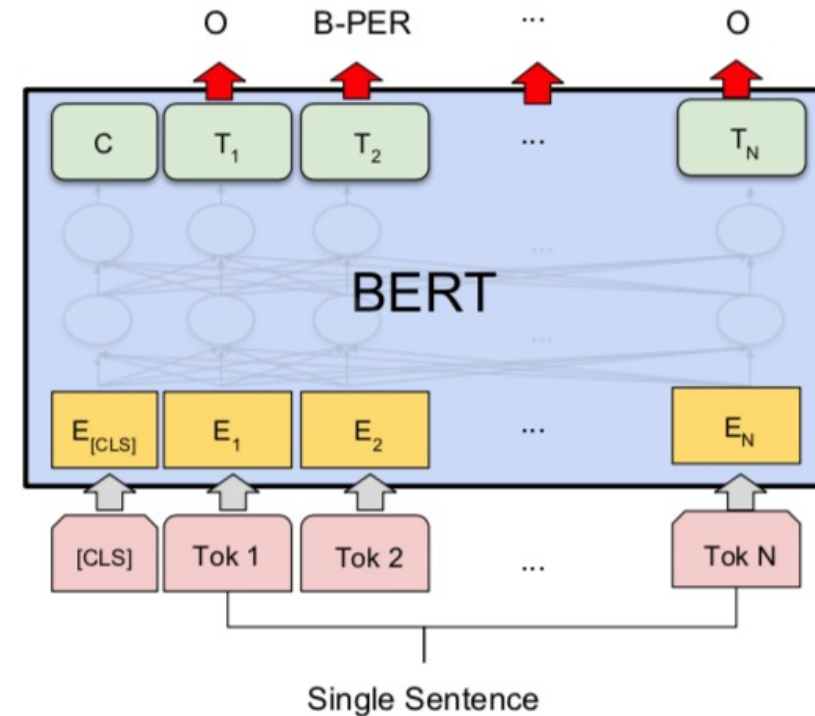


(b) Single Sentence Classification Tasks:  
SST-2, CoLA

# BERT Token-level tasks



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# General Language Understanding Evaluation (GLUE) benchmark

## GLUE Test results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

**MNLI:** Multi-Genre Natural Language Inference

**QQP:** Quora Question Pairs

**QNLI:** Question Natural Language Inference

**SST-2:** The Stanford Sentiment Treebank

**CoLA:** The Corpus of Linguistic Acceptability

**STS-B:** The Semantic Textual Similarity Benchmark

**MRPC:** Microsoft Research Paraphrase Corpus

**RTE:** Recognizing Textual Entailment

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805



# Transformers Transformers

## State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch

- **Transformers**
  - **pytorch-transformers**
  - **pytorch-pretrained-bert**
- **provides state-of-the-art general-purpose architectures**
  - **(BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, CTRL...)**
  - **for Natural Language Understanding (NLU) and Natural Language Generation (NLG) with over 32+ pretrained models in 100+ languages and deep interoperability between TensorFlow 2.0 and PyTorch.**

# **Text Analytics with Python**





# spaCy:

## Natural Language Processing

spaCy

USAGE

MODELS

API

UNIVERSE



Search docs

# Industrial-Strength Natural Language Processing

IN PYTHON

### Get things done

spaCy is designed to help you do real work — to build real products, or gather real insights. The library respects your time, and tries to avoid wasting it. It's easy to install, and its API is simple and productive. We like to think of spaCy as the Ruby on Rails of Natural Language Processing.

### Blazing fast

spaCy excels at large-scale information extraction tasks. It's written from the ground up in carefully memory-managed Cython. Independent research in 2015 found spaCy to be the fastest in the world. If your application needs to process entire web dumps, spaCy is the library you want to be using.

### Deep learning

spaCy is the best way to prepare text for deep learning. It interoperates seamlessly with TensorFlow, PyTorch, scikit-learn, Gensim and the rest of Python's awesome AI ecosystem. With spaCy, you can easily construct linguistically sophisticated statistical models for a variety of NLP problems.

<https://spacy.io/>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

The screenshot shows a Google Colab notebook titled "python101.ipynb". The table of contents on the left lists various NLP topics. The main code cell contains the following Python code:

```
1 text = "Steve Jobs and Steve Wozniak incorporated Apple Computer on January 3, 1977, in Cupertino, California."
2 doc = nlp(text)
3 displacy.render(doc, style="ent", jupyter=True)
```

The visual output of the code shows the sentence with entities highlighted: "Steve Jobs" (PERSON), "Steve Wozniak" (PERSON), "Apple Computer" (ORG), "January 3, 1977" (DATE), "Cupertino" (GPE), and "California" (GPE).

```
[ ] 1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Stanford University is located in California. It is a great university.")
4 import pandas as pd
5 cols = ("text", "lemma", "pos", "tag", "pos_explain", "stopword")
6 rows = []
7 for t in doc:
8     row = [t.text, t.lemma_, t.pos_, t.tag_, spacy.explain(t.pos_), t.is_stop]
9     rows.append(row)
10 df = pd.DataFrame(rows, columns=cols)
11 df
```

The output of the DataFrame is as follows:

	text	lemma	pos	tag	pos_explain	stopword
0	Stanford	Stanford	PROPN	NNP	proper noun	False
1	University	University	PROPN	NNP	proper noun	False
2	is	be	VERB	VBZ	verb	True
3	located	locate	VERB	VBN	verb	False
4	in	in	ADP	IN	adposition	True
5	California	California	PROPN	NNP	proper noun	False
6	.	.	PUNCT	.	punctuation	False
7	It	-PRON-	PRON	PRP	pronoun	True

<https://tinyurl.com/aintpupython101>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

The screenshot shows a Google Colab notebook titled "python101.ipynb". The interface includes a top menu bar with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help", along with a status indicator "All changes saved". On the right, there are icons for "Comment", "Share", "Settings", and a user profile "A".

The left sidebar displays a "Table of contents" for the notebook, listing various topics under "Text Analytics and Natural Language Processing (NLP)", such as "Python for Natural Language Processing", "spaCy Chinese Model", "Open Chinese Convert", "Jieba", "NLTK", "Stanza", "Text Processing and Understanding", and "NLP Zero to Hero".

The main content area shows a code cell with the following text:

- Text Analytics and Natural Language Processing (NLP)
- Python for Natural Language Processing
  - spaCy

Below the text, there are two code execution blocks:

```
[1] 1 !python -m spacy download en_core_web_sm
```

```
[3] 1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Apple is looking at buying U.K. startup for $1 billion")
4 for token in doc:
5     print(token.text, token.pos_, token.dep_)
```

The output of the second code block is a list of tokens with their part-of-speech tags and dependency labels:

```
Apple PROPN nsubj
is AUX aux
looking VERB ROOT
at ADP prep
buying VERB pcomp
U.K. PROPN compound
startup NOUN dobj
for ADP prep
$ SYM quantmod
1 NUM compound
billion NUM pobj
```

<https://tinyurl.com/aintpupython101>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>



python101.ipynb ☆

File Edit View Insert Runtime Tools Help [All changes saved](#)

Comment Share Settings Profile



+ Code + Text

RAM Disk Editing

```
[ ] 1 import spacy
    2 nlp = spacy.load("en_core_web_sm")
    3 doc = nlp("Apple is looking at buying U.K. startup for $1 billion")
    4 import pandas as pd
    5 cols = ("text", "lemma", "POS", "explain", "stopword")
    6 rows = []
    7 for t in doc:
    8     row = [t.text, t.lemma_, t.pos_, spacy.explain(t.pos_), t.is_stop]
    9     rows.append(row)
   10 df = pd.DataFrame(rows, columns=cols)
   11 df
```

	text	lemma	POS	explain	stopword
0	Apple	Apple	PROPN	proper noun	False
1	is	be	VERB	verb	True
2	looking	look	VERB	verb	False
3	at	at	ADP	adposition	True
4	buying	buy	VERB	verb	False
5	U.K.	U.K.	PROPN	proper noun	False
6	startup	startup	NOUN	noun	False
7	for	for	ADP	adposition	True
8	\$	\$	SYM	symbol	False
9	1	1	NUM	numeral	False
10	billion	billion	NUM	numeral	False

<https://tinyurl.com/aintpupython101>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

python101.ipynb ☆

File Edit View Insert Runtime Tools Help [All changes saved](#)

+ Code + Text

```
[ ] 1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Stanford University is located in California. It is a great university.")
4 import pandas as pd
5 cols = ("text", "lemma", "POS", "explain", "stopword")
6 rows = []
7 for t in doc:
8     row = [t.text, t.lemma_, t.pos_, spacy.explain(t.pos_), t.is_stop]
9     rows.append(row)
10 df = pd.DataFrame(rows, columns=cols)
11 df
```

	text	lemma	POS	explain	stopword
0	Stanford	Stanford	PROPN	proper noun	False
1	University	University	PROPN	proper noun	False
2	is	be	VERB	verb	True
3	located	locate	VERB	verb	False
4	in	in	ADP	adposition	True
5	California	California	PROPN	proper noun	False
6	.	.	PUNCT	punctuation	False
7	It	-PRON-	PRON	pronoun	True
8	is	be	VERB	verb	True
9	a	a	DET	determiner	True
10	great	great	ADJ	adjective	False
11	university	university	NOUN	noun	False
12	.	.	PUNCT	punctuation	False

<https://tinyurl.com/aintpupython101>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>



python101.ipynb ☆

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text



```
[ ] 1 import spacy
     2 nlp = spacy.load("en_core_web_sm")
     3 text = "Stanford University is located in California. It is a great university."
     4 doc = nlp(text)
     5 for ent in doc.ents:
     6     print(ent.text, ent.label_)
```

↳ Stanford University ORG  
California GPE

```
[ ] 1 from spacy import displacy
     2 text = "Stanford University is located in California. It is a great university."
     3 doc = nlp(text)
     4 displacy.render(doc, style="ent", jupyter=True)
```

↳ Stanford University ORG is located in California GPE . It is a great university.

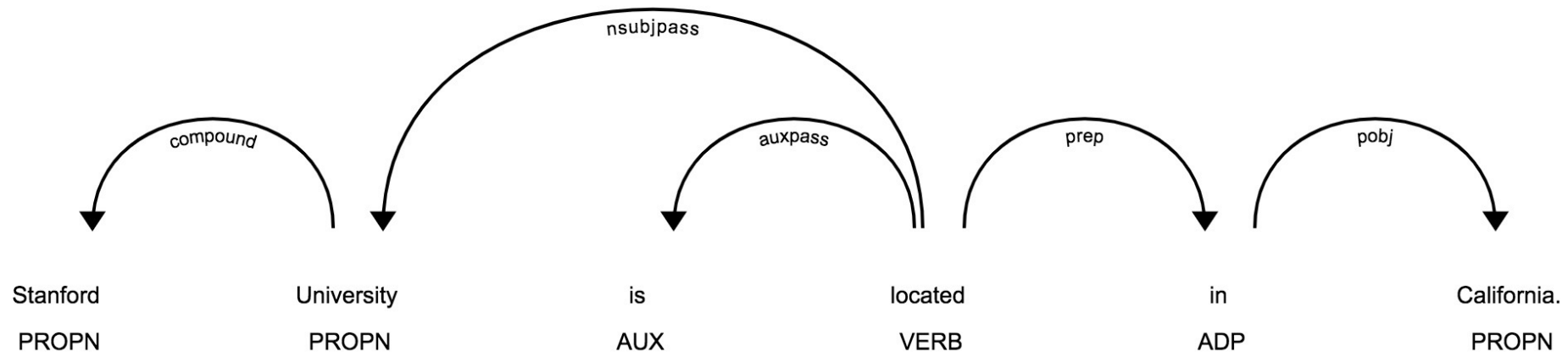
<https://tinyurl.com/aintpupython101>

# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

```
1 from spacy import displacy
2 text = "Stanford University is located in California. It is a great university."
3 doc = nlp(text)
4 displacy.render(doc, style="ent", jupyter=True)
5 displacy.render(doc, style="dep", jupyter=True)
```

Stanford University **ORG** is located in **California GPE** . It is a great university.



<https://tinyurl.com/aintpupython101>



# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

The screenshot shows a Google Colab notebook titled "python101.ipynb". The interface includes a top menu bar with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". A "Table of contents" sidebar on the left lists various NLP topics. The main code cell contains two Python snippets. The first snippet uses spaCy's displacy.render to visualize named entities in a sentence. The second snippet uses spaCy's nlp to process a sentence and pandas to create a DataFrame of token features. Below the code, a table displays the output of the second snippet, showing token index, text, lemma, part of speech (pos), tag, explanation, and whether it is a stopword.

```
1 text = "Steve Jobs and Steve Wozniak incorporated Apple Computer on January 3, 1977, in Cupertino, California."
2 doc = nlp(text)
3 displacy.render(doc, style="ent", jupyter=True)
```

Steve Jobs PERSON and Steve Wozniak PERSON incorporated Apple Computer ORG on January 3, 1977 DATE , in Cupertino GPE , California GPE .

```
[ ] 1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Stanford University is located in California. It is a great university.")
4 import pandas as pd
5 cols = ("text", "lemma", "pos", "tag", "pos_explain", "stopword")
6 rows = []
7 for t in doc:
8     row = [t.text, t.lemma_, t.pos_, t.tag_, spacy.explain(t.pos_), t.is_stop]
9     rows.append(row)
10 df = pd.DataFrame(rows, columns=cols)
11 df
```

	text	lemma	pos	tag	pos_explain	stopword
0	Stanford	Stanford	PROPN	NNP	proper noun	False
1	University	University	PROPN	NNP	proper noun	False
2	is	be	VERB	VBZ	verb	True
3	located	locate	VERB	VRB	verb	False
4	in	in	ADP	IN	adposition	True
5	California	California	PROPN	NNP	proper noun	False
6	.	.	PUNCT	.	punctuation	False
7	It	-PRON-	PRON	PRP	pronoun	True

<https://tinyurl.com/aintpupython101>

# Teaching



- **Artificial Intelligence for Text Analytics**
  - Spring 2022
- **Software Engineering**
  - Fall 2020, Fall, 2021, Spring 2022
- **Artificial Intelligence in Finance and Quantitative**
  - Fall 2021
- **Artificial Intelligence**
  - Spring 2021
- **Data Mining**
  - Spring 2021
- **Big Data Analytics**
  - Fall 2020
- **Foundation of Business Cloud Computing**
  - Spring 2021, Spring 2022

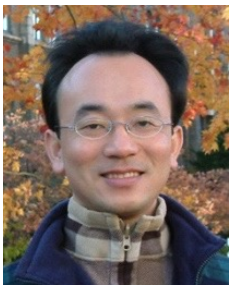
# Research Project



- **Applying AI technology to construct knowledge graphs of cryptocurrency anti-money laundering: a few-shot learning model**
  - MOST, 110-2410-H-305-013-MY2, 2021/08/01~2023/07/31
- **AI for Corporate Sustainability Assessment and Cross Language Corporate Sustainability Reports Generative Mode**
  - NTPU, 111-NTPU\_ORDA-F-001 , 2022/01/01~2022/12/31
- **Artificial Intelligence for FinTech Knowledge Graph from Patent Textual Analytics**
  - NTPU, 111-NTPU\_ORDA-F-003, 2022/01/01~2022/12/31

# Summary

- This course introduces the **fundamental concepts, research issues, and hands-on practices of Artificial Intelligence for Text Analytics.**
- Topics include:
  1. Introduction to Introduction to Artificial Intelligence for Text Analytics
  2. Foundations of Text Analytics: Natural Language Processing (NLP)
  3. Python for Natural Language Processing
  4. Natural Language Processing with Transformers
  5. Text Classification and Sentiment Analysis
  6. Multilingual Named Entity Recognition (NER), Text Similarity and Clustering
  7. Text Summarization and Topic Models
  8. Text Generation
  9. Question Answering and Dialogue Systems
  10. Deep Learning, Transfer Learning, Zero-Shot, and Few-Shot Learning for Text Analytics
  11. Case Study on Artificial Intelligence for Text Analytics



# Artificial Intelligence for Text Analytics



2020 Cohort



Accredited  
Educator



Solutions  
Architect  
Associate



Cloud  
Practitioner

## Contact Information

**Min-Yuh Day, Ph.D.**

**Associate Professor**

[Institute of Information Management, National Taipei University](#)

Tel: 02-86741111 ext. 66873

Office: B8F12

Address: 151, University Rd., San Shia District, New Taipei City, 23741 Taiwan

Email: [myday@gm.ntpu.edu.tw](mailto:myday@gm.ntpu.edu.tw)

Web: <http://web.ntpu.edu.tw/~myday/>

