

Artificial Intelligence for Text Analytics

Foundations of Text Analytics: Natural Language Processing (NLP)

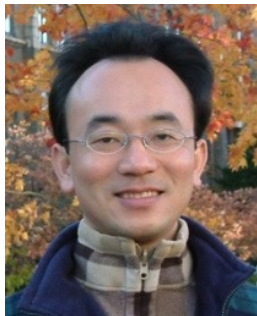
1102AITA02

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>



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

Foundations of Text Analytics: Natural Language Processing (NLP)

Outline

- **Text Analytics and Text Mining**
- **Natural Language Processing (NLP)**
- **Text Analytics with Python**

Artificial Intelligence (AI)

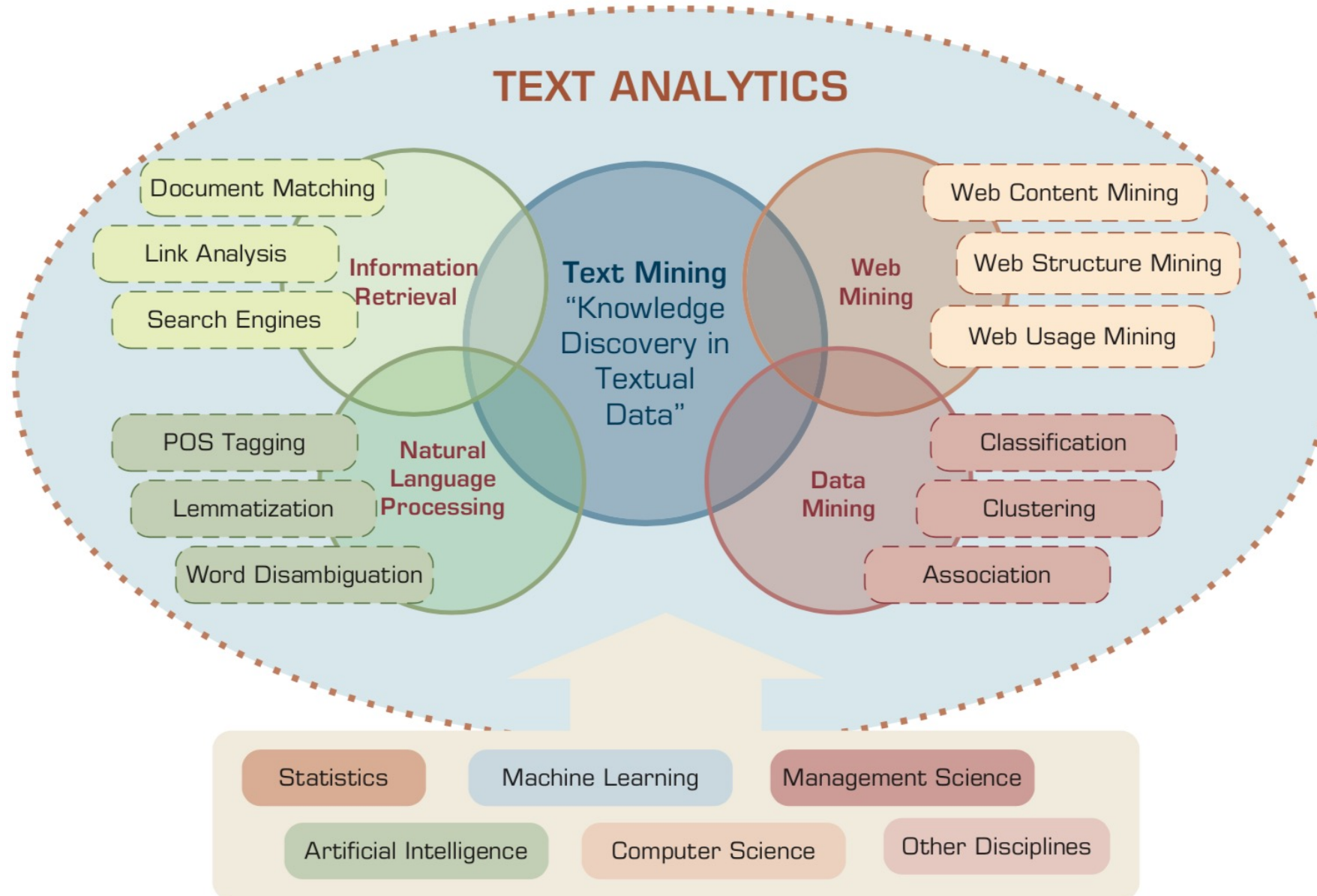
Text Analytics

(TA)

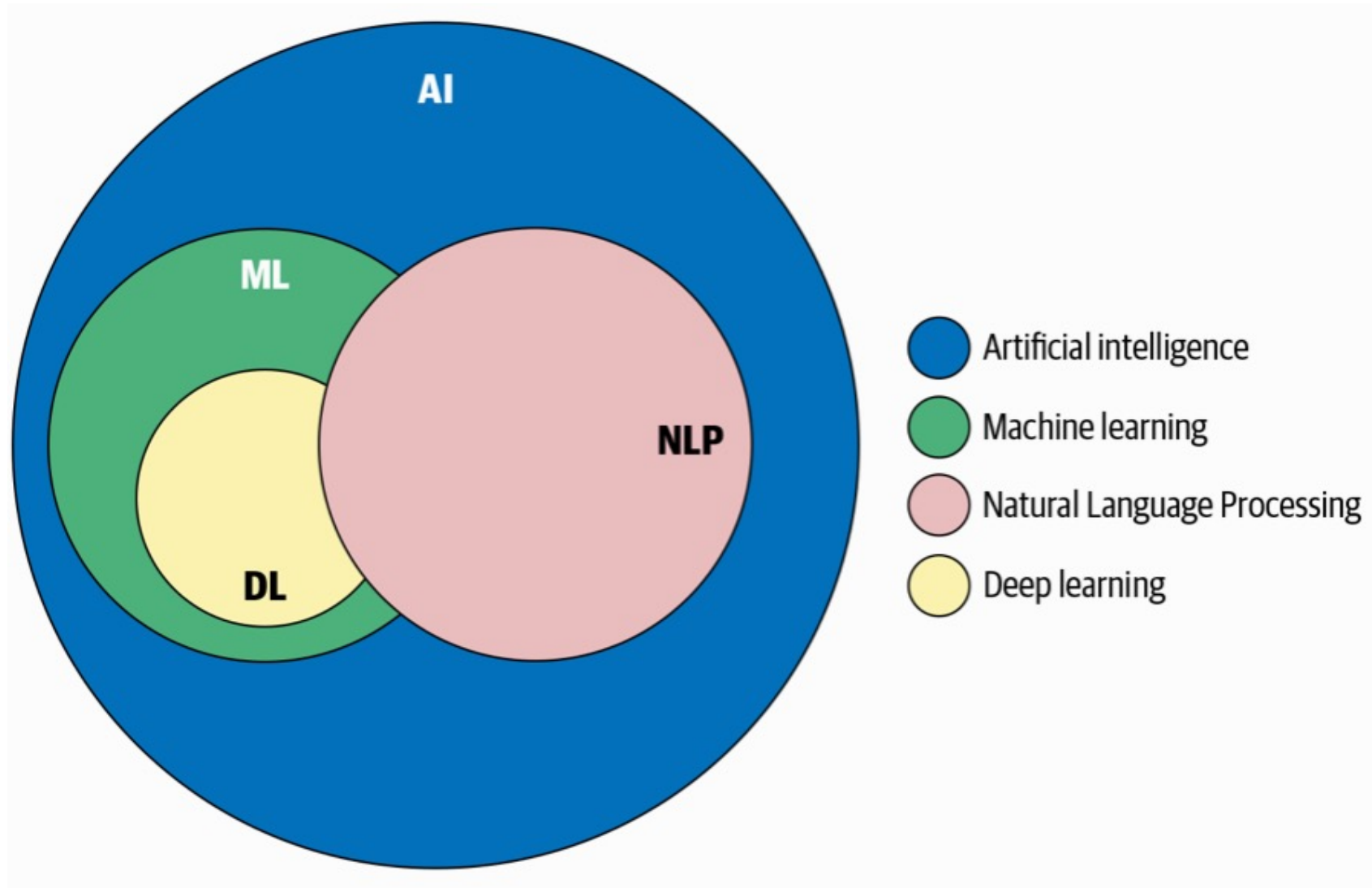
Text Mining (TM)

Natural Language Processing (NLP)

Text Analytics and Text Mining



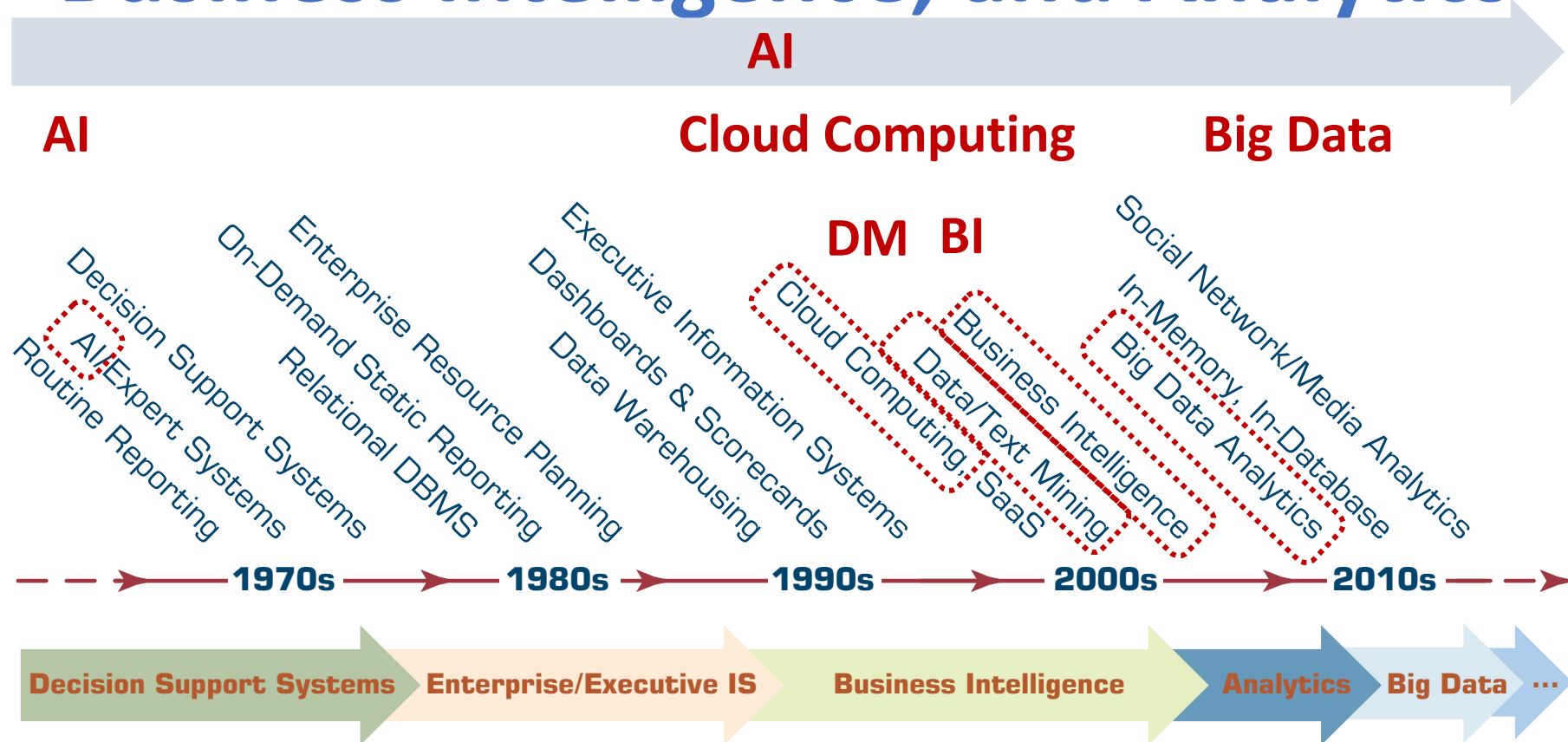
AI, NLP, ML, DL



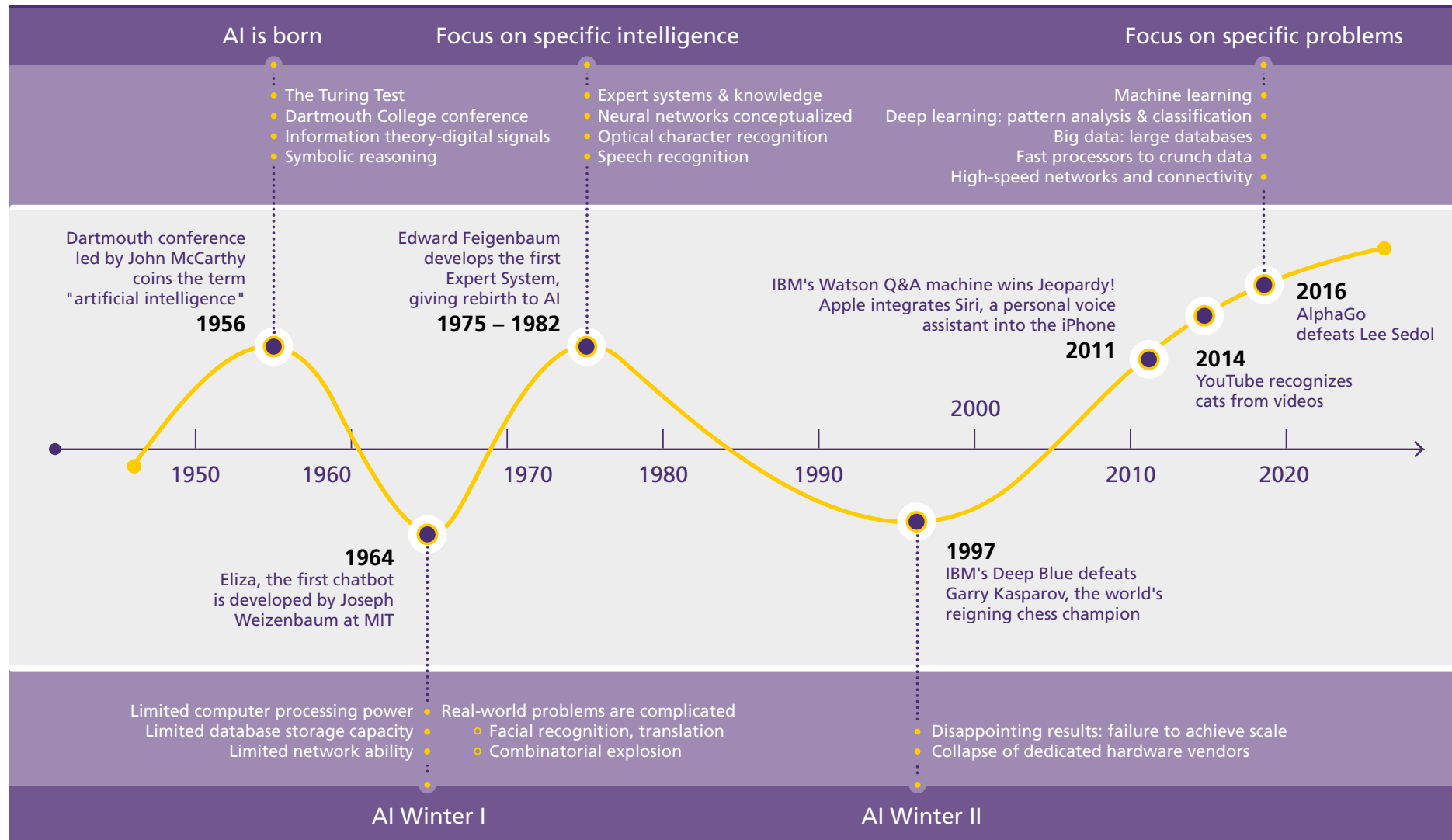
Artificial Intelligence (AI)

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

Thinking Humanly	Thinking Rationally
Acting Humanly	Acting Rationally

4 Approaches of AI

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

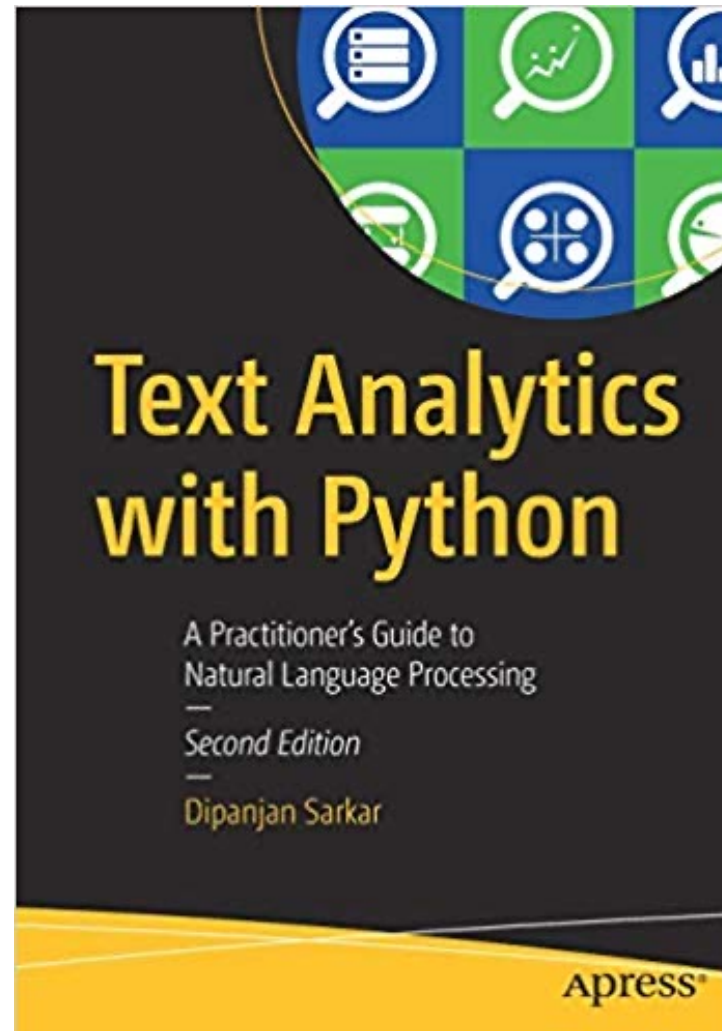
AI Acting Humanly: The Turing Test Approach (Alan Turing, 1950)

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

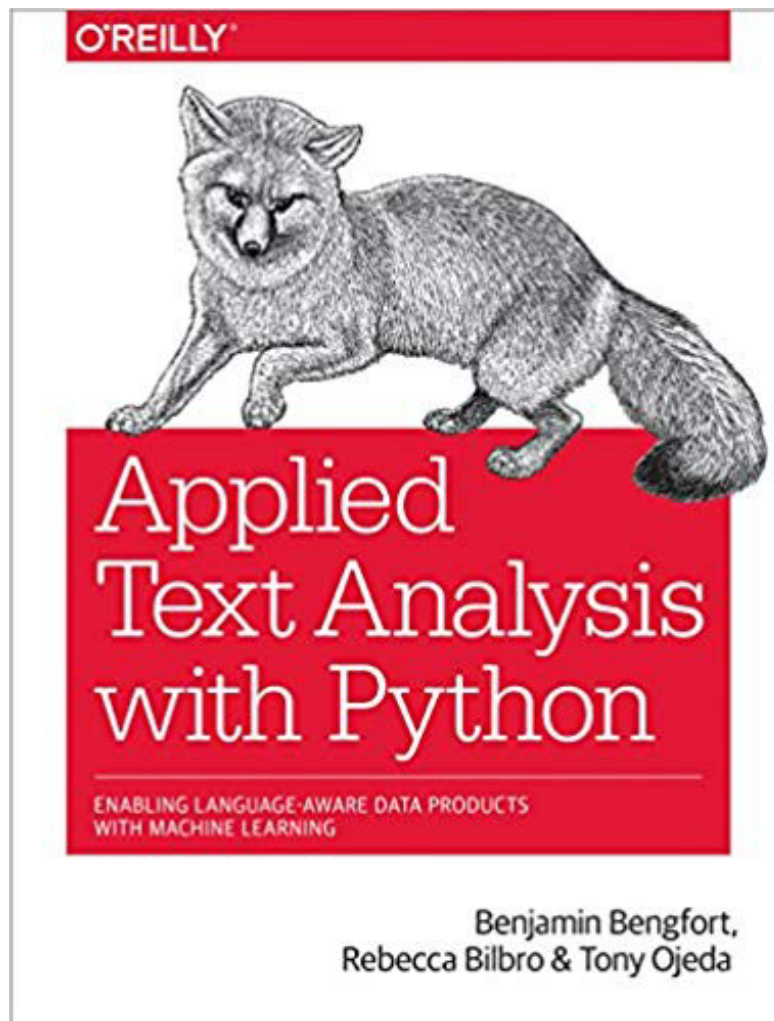
Text Analytics and Text Mining

Dipanjan Sarkar (2019),

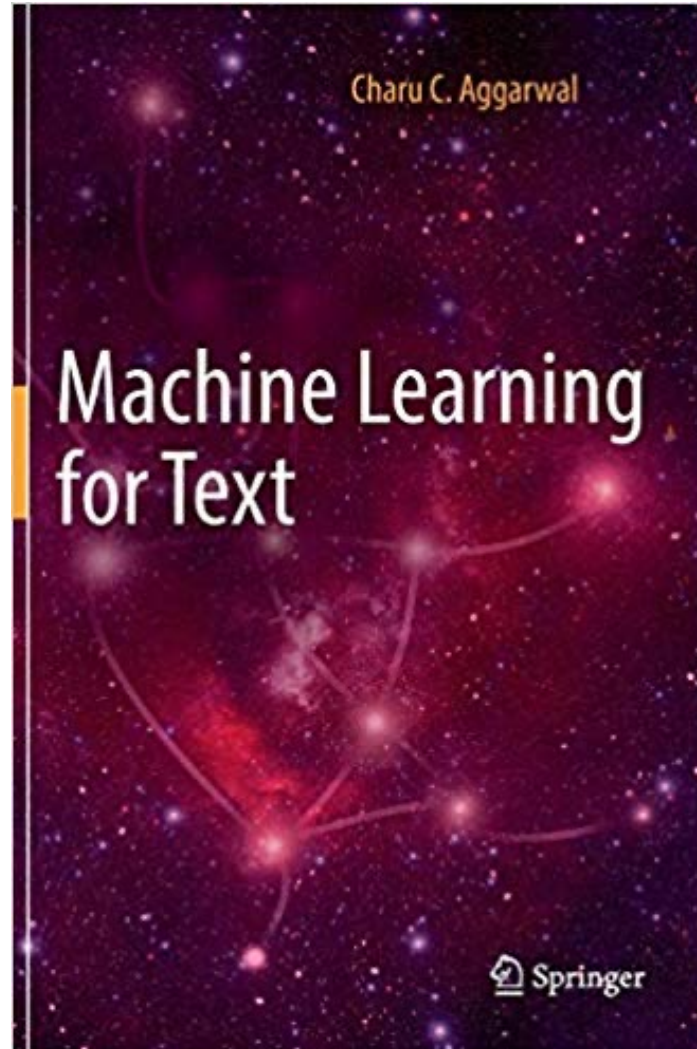
Text Analytics with Python:
A Practitioner's Guide to Natural Language Processing,
Second Edition. APress.



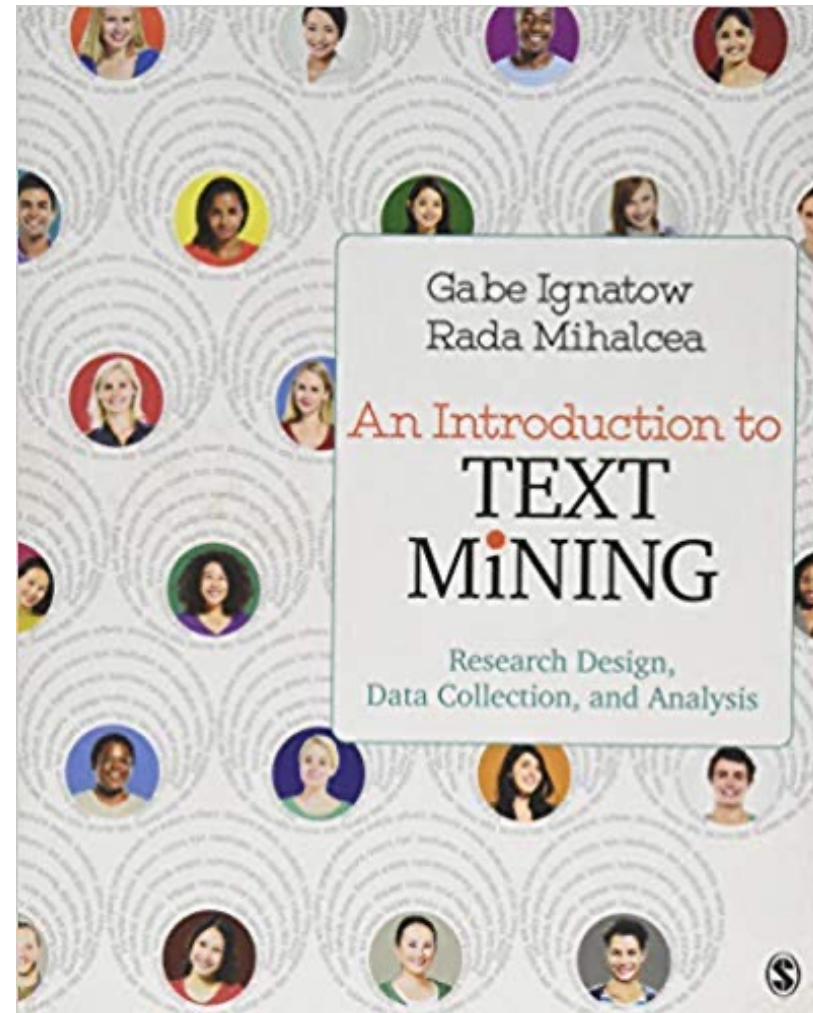
Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018),
Applied Text Analysis with Python:
Enabling Language-Aware Data Products with Machine Learning,
O'Reilly.



Charu C. Aggarwal (2018),
Machine Learning for Text,
Springer

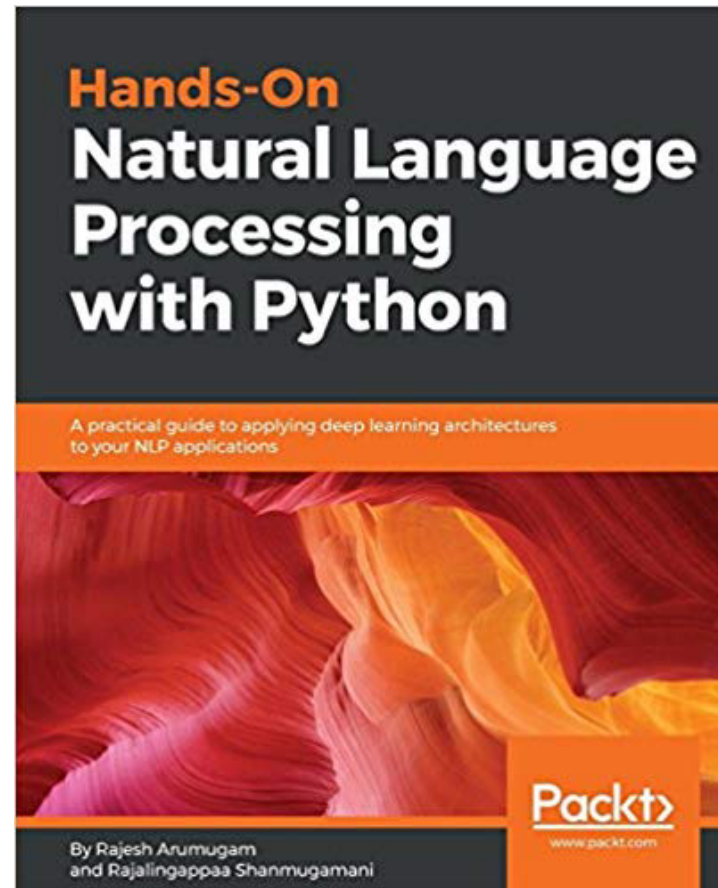


Gabe Ignatow and Rada F. Mihalcea (2017),
An Introduction to Text Mining:
Research Design, Data Collection, and Analysis,
SAGE Publications.



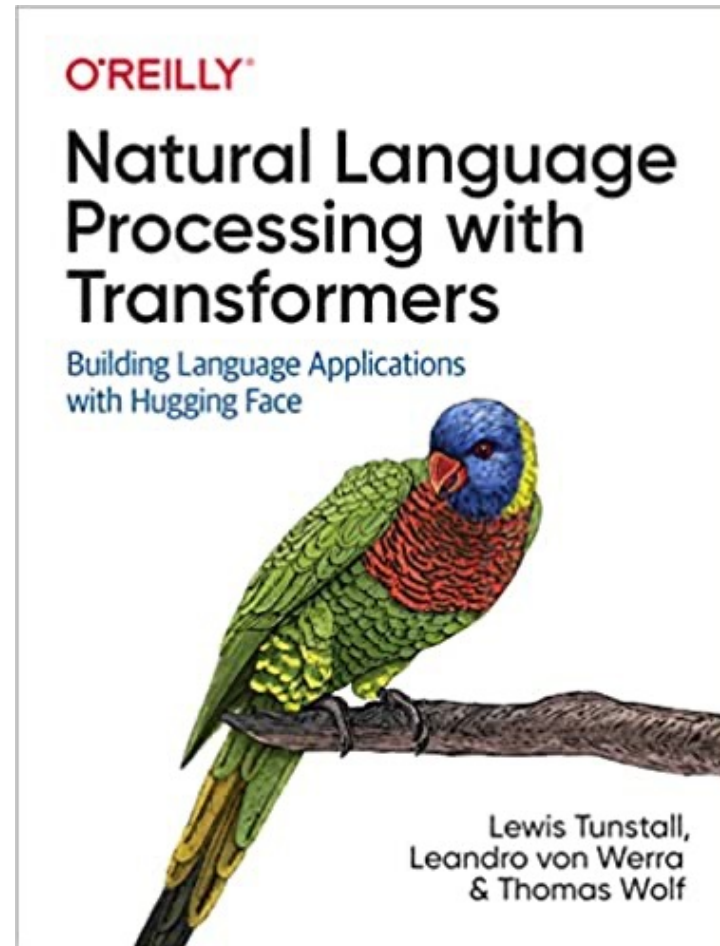
Rajesh Arumugam (2018),
**Hands-On Natural Language Processing
with Python:**

A practical guide to applying deep learning architectures to your NLP applications,
Packt



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

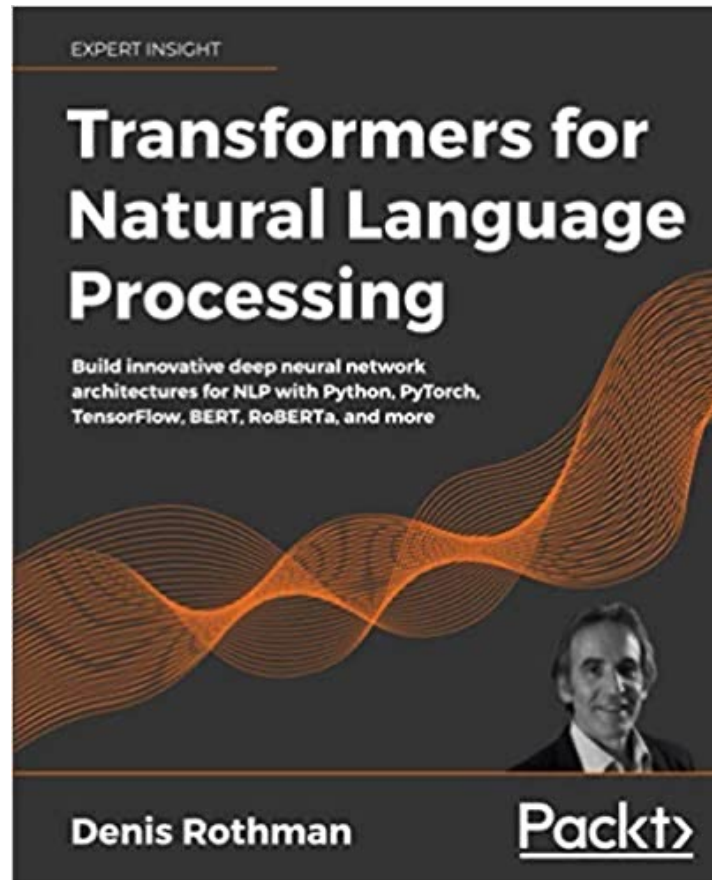
Building Language Applications with Hugging Face,
O'Reilly Media.



Denis Rothman (2021),

Transformers for Natural Language Processing:

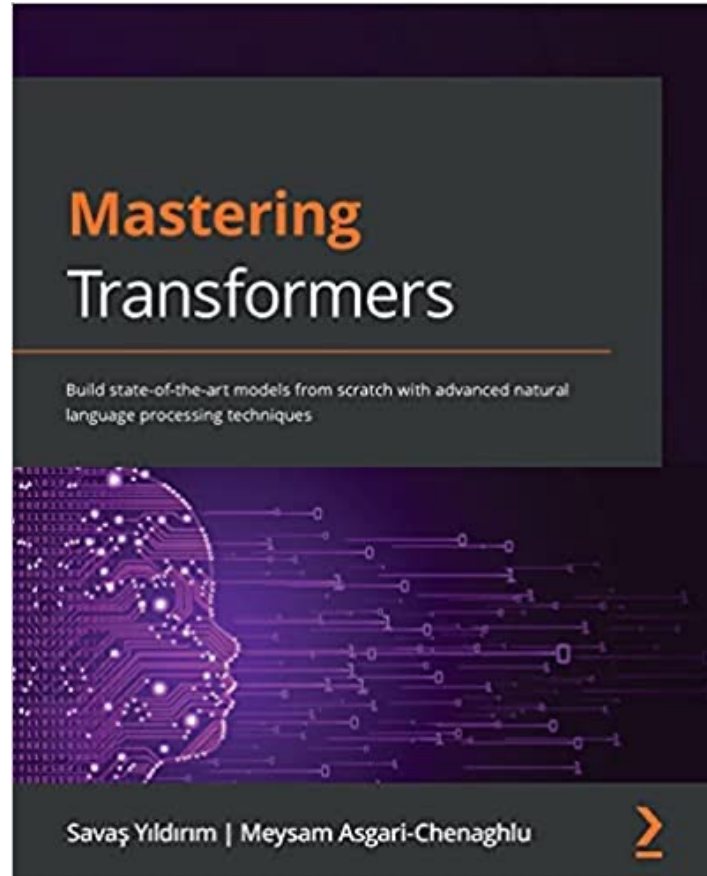
Build innovative deep neural network architectures for NLP with Python, PyTorch, TensorFlow, BERT, RoBERTa, and more,
Packt Publishing.



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.

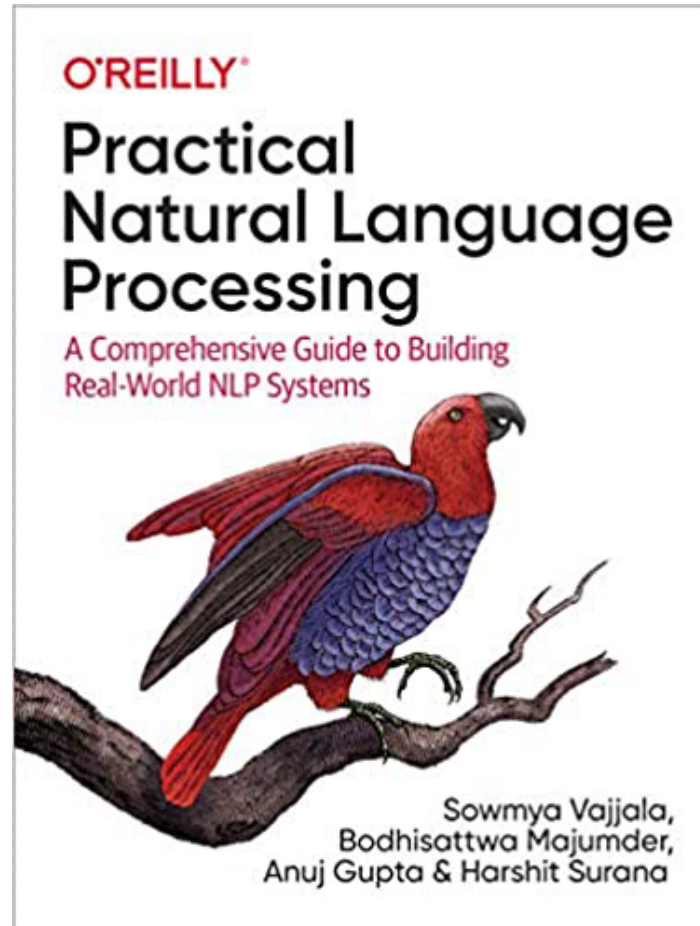


Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta (2020),

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



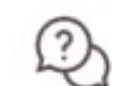
Information
Extraction



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

Text Analytics

(TA)

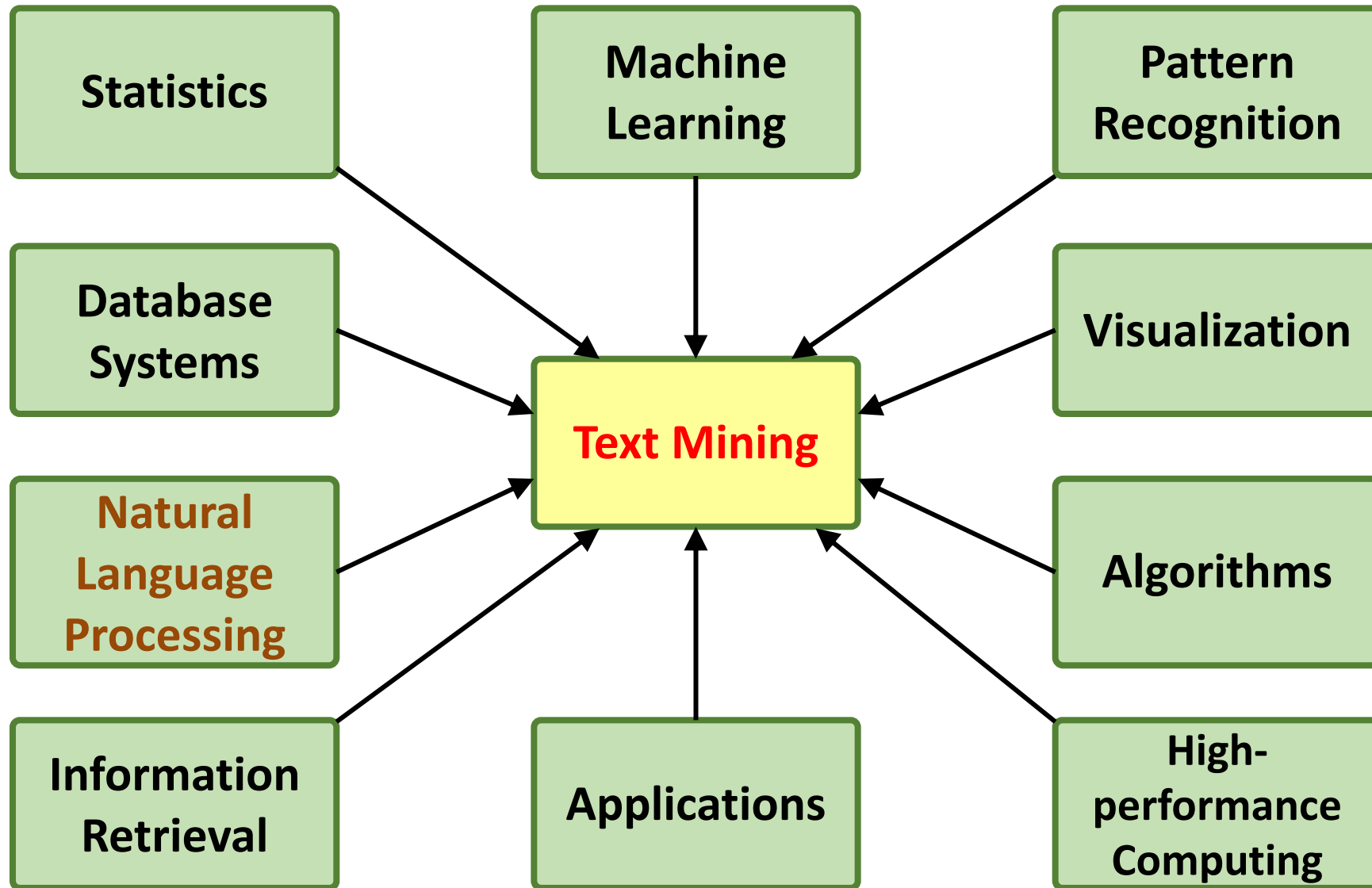
Text Analytics

- **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**

Text Mining Technologies



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. ...”**

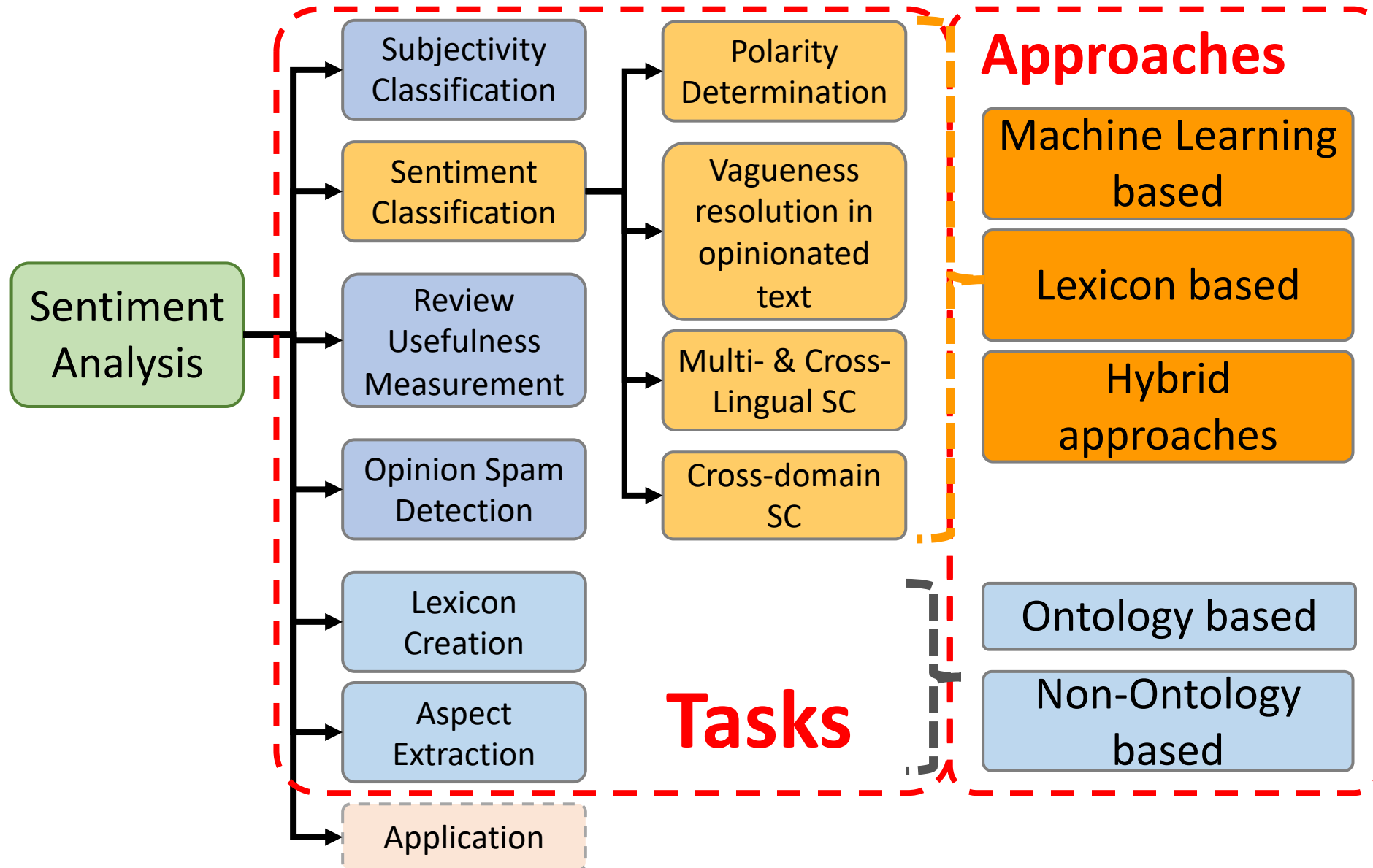


**+Positive
Opinion**

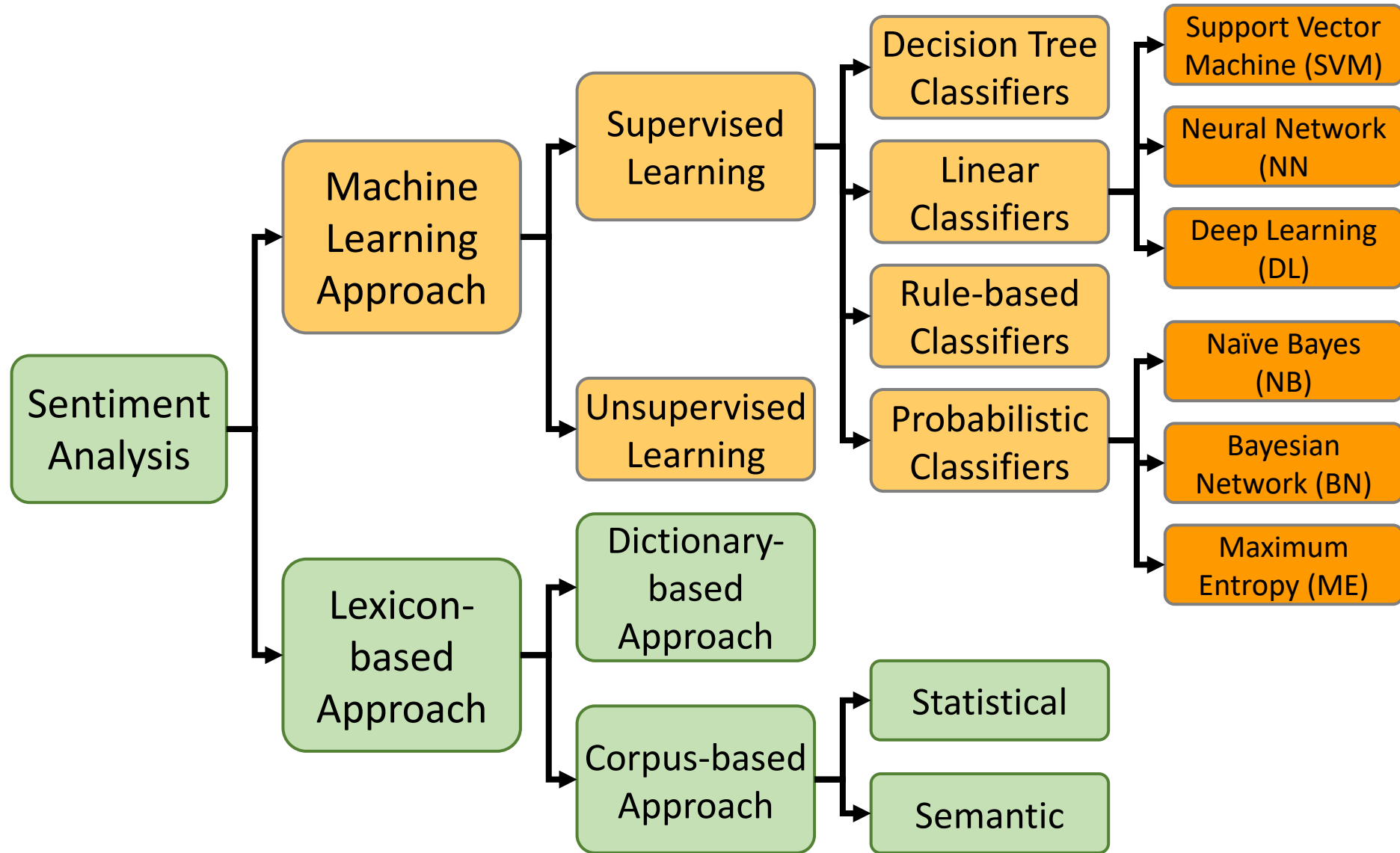


**-Negative
Opinion**

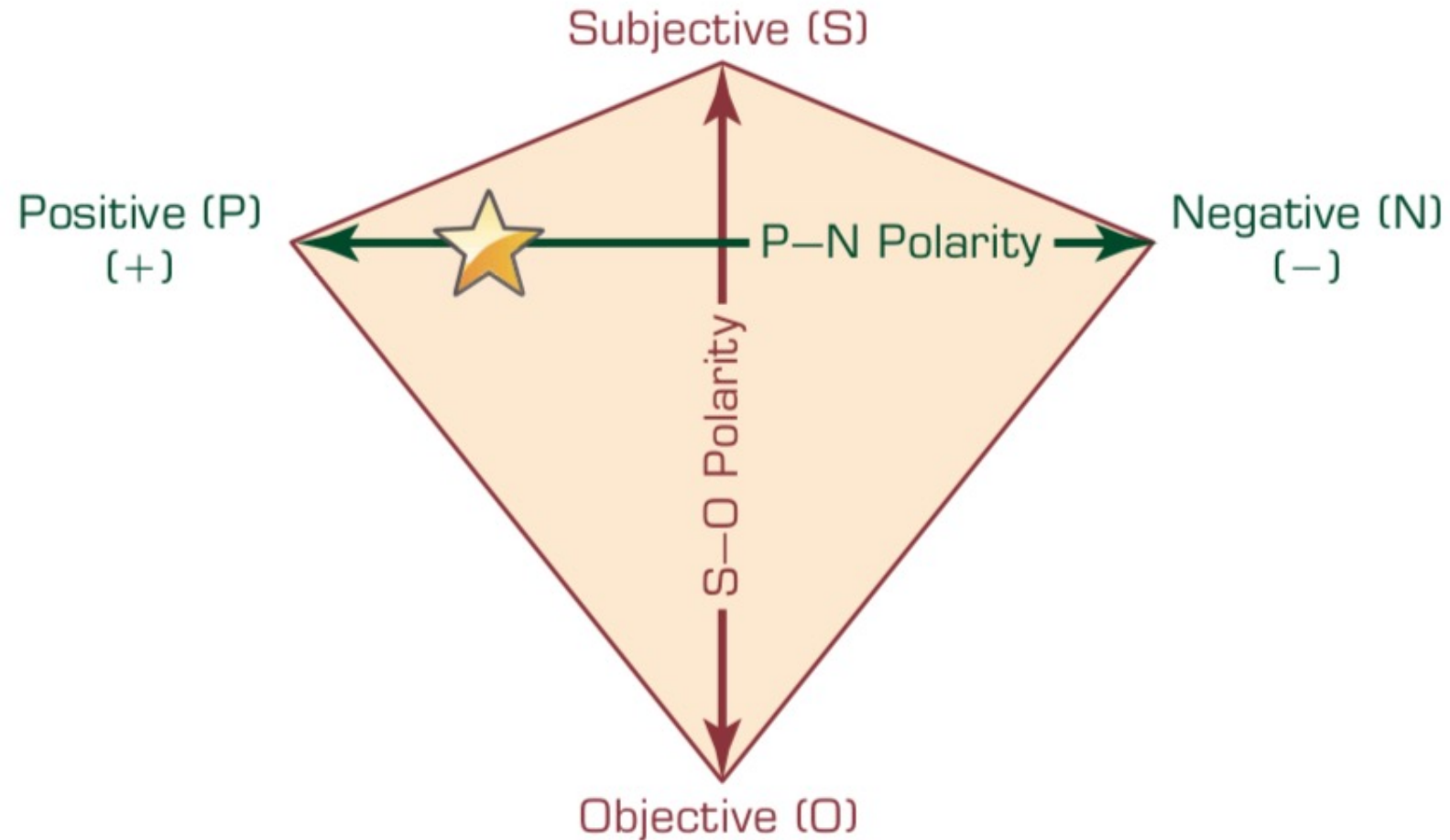
Sentiment Analysis



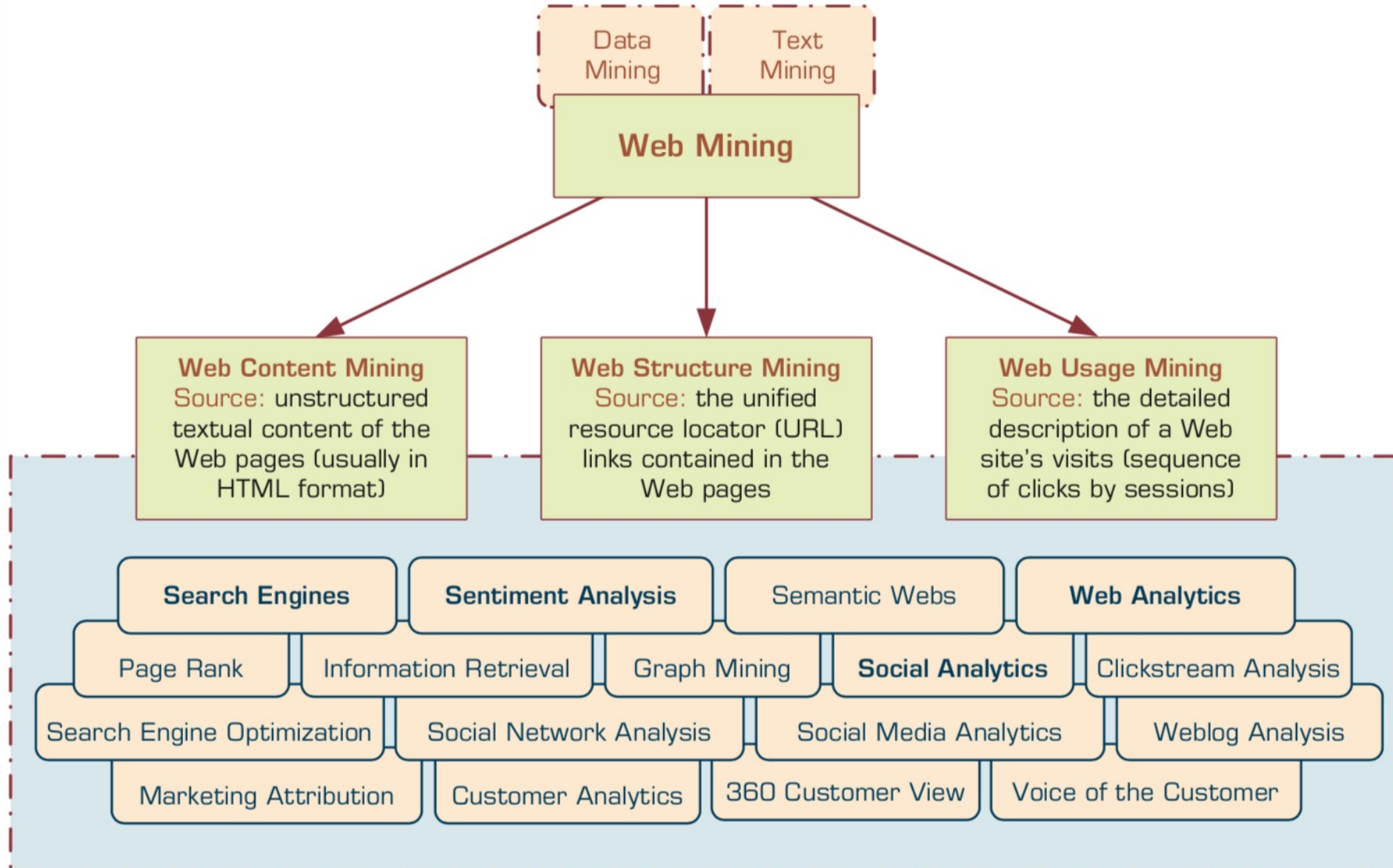
Sentiment Classification Techniques



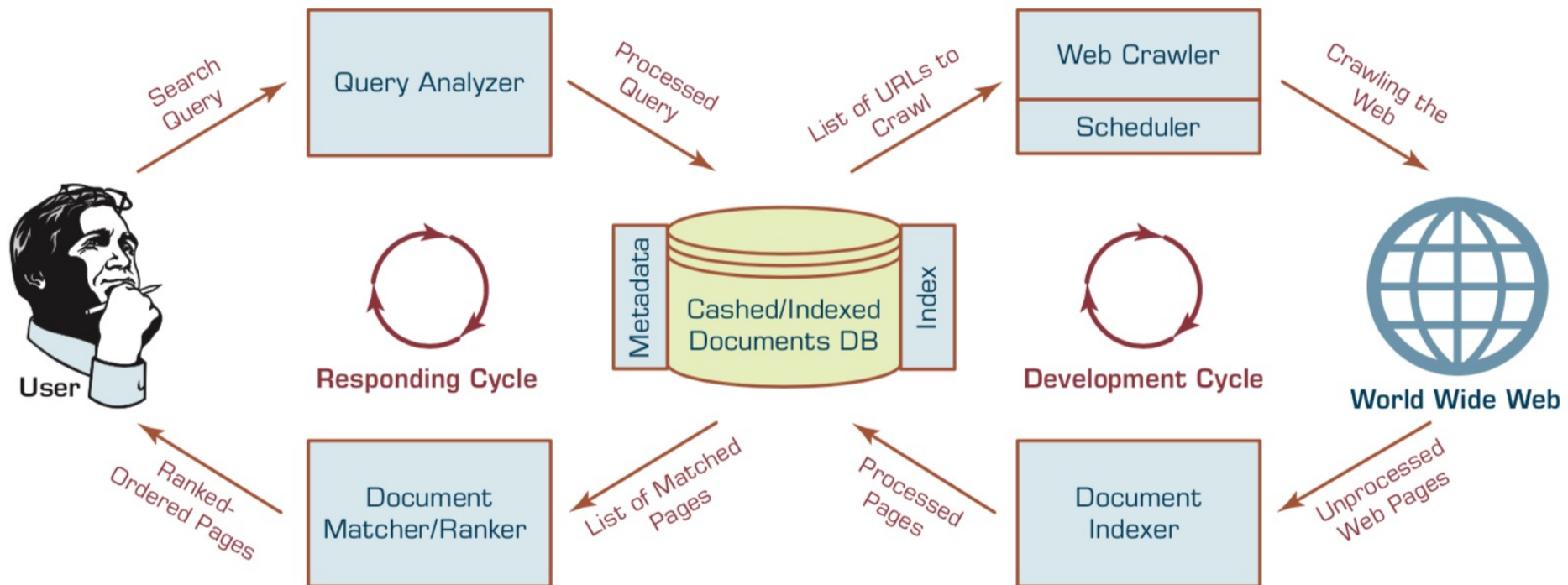
P–N Polarity and S–O Polarity Relationship



Taxonomy of Web Mining



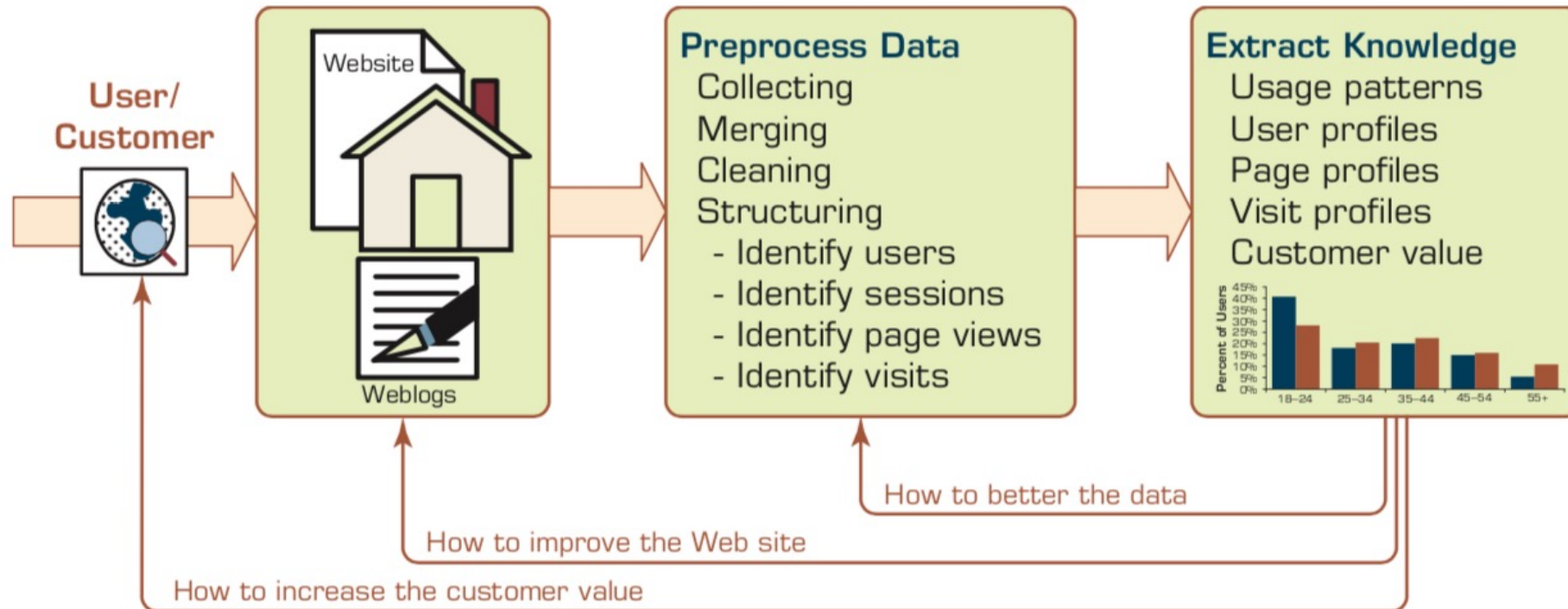
Structure of a Typical Internet Search Engine



Web Usage Mining (Web Analytics)

- **Web usage mining (Web analytics) is the extraction of useful information from data generated through Web page visits and transactions.**
- **Clickstream Analysis**

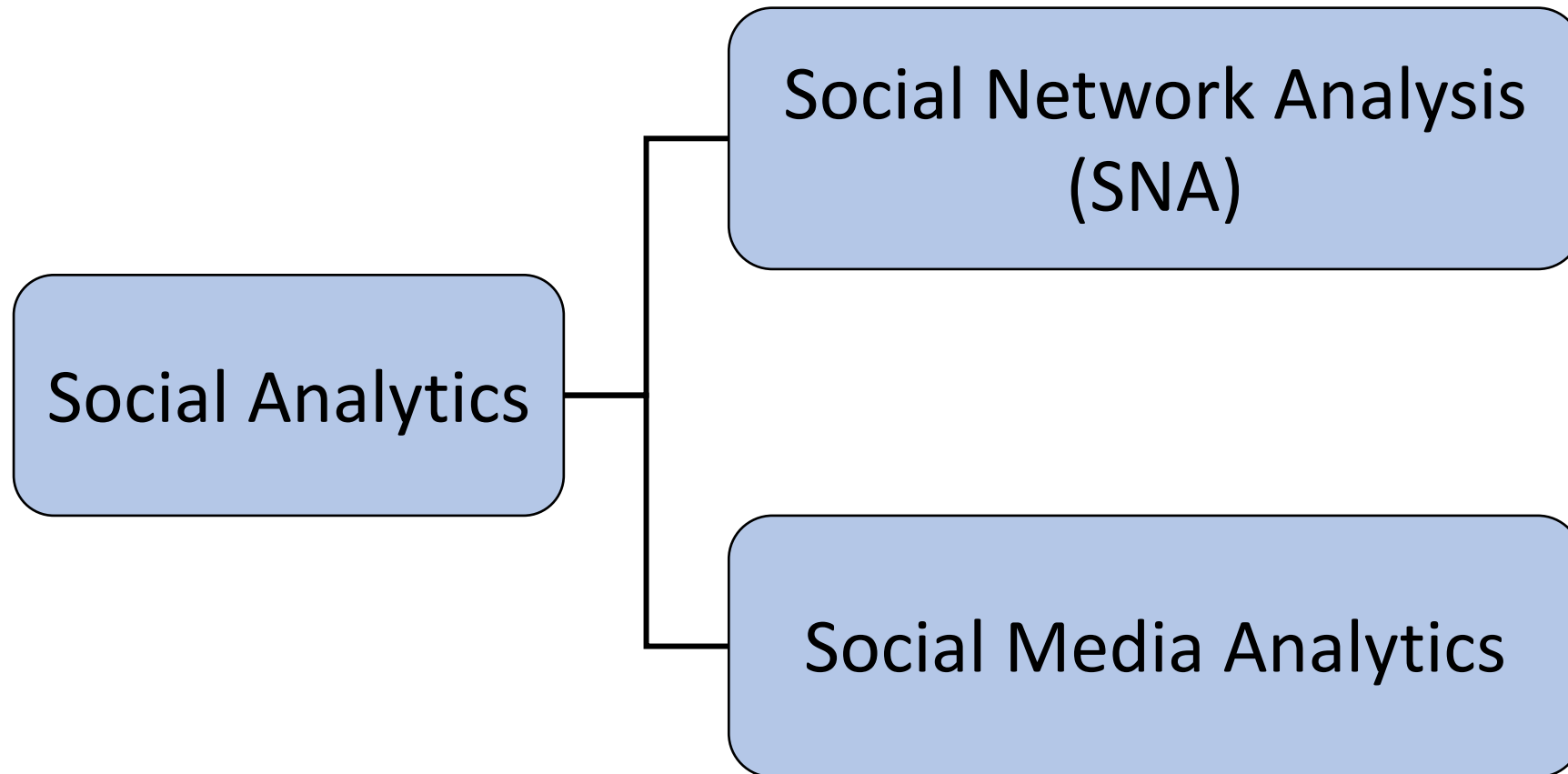
Extraction of Knowledge from Web Usage Data



Social Analytics

- **Social analytics is defined as monitoring, analyzing, measuring and interpreting digital interactions and relationships of people, topics, ideas and content.**

Branches of Social Analytics



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

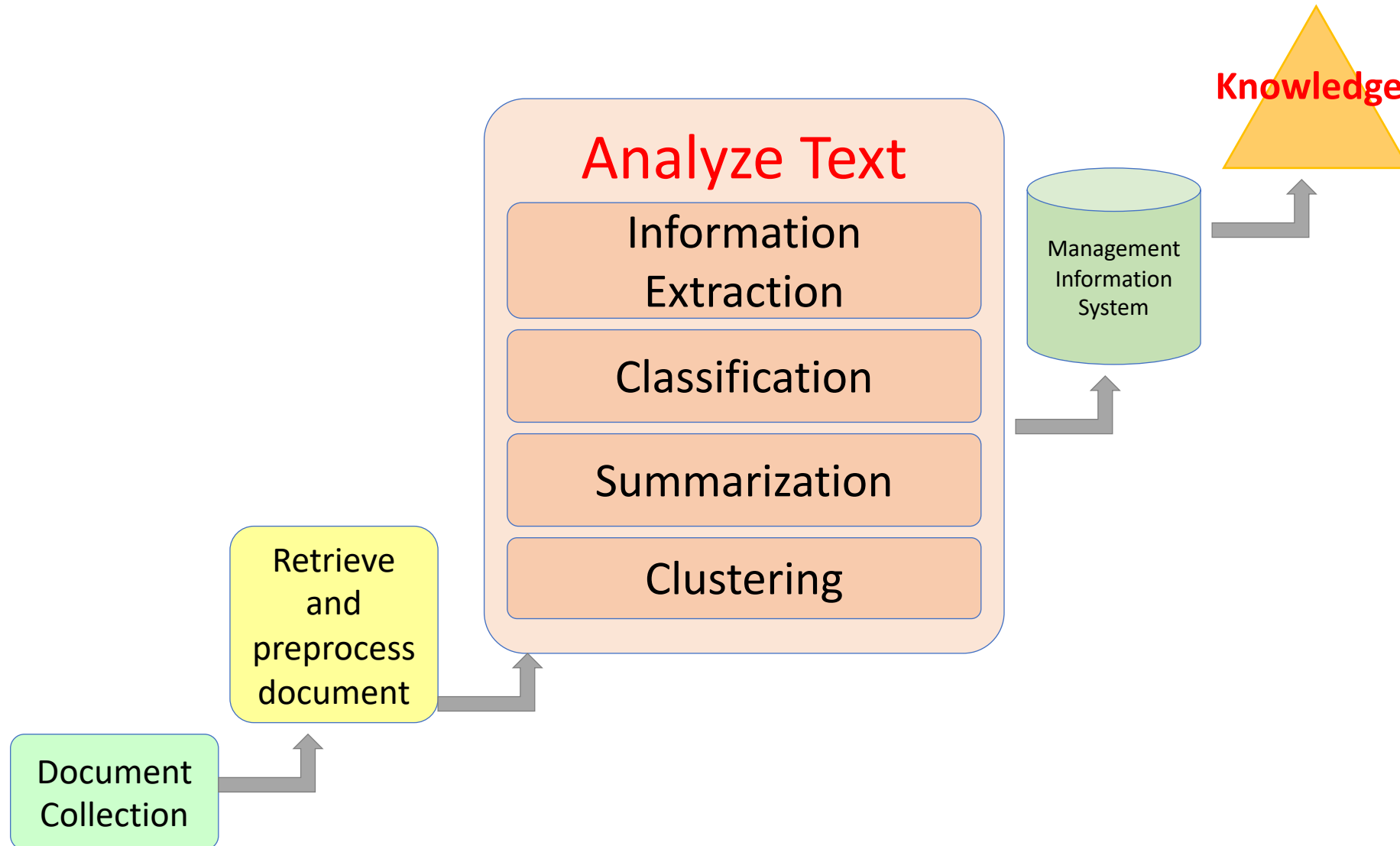
(text data mining)

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

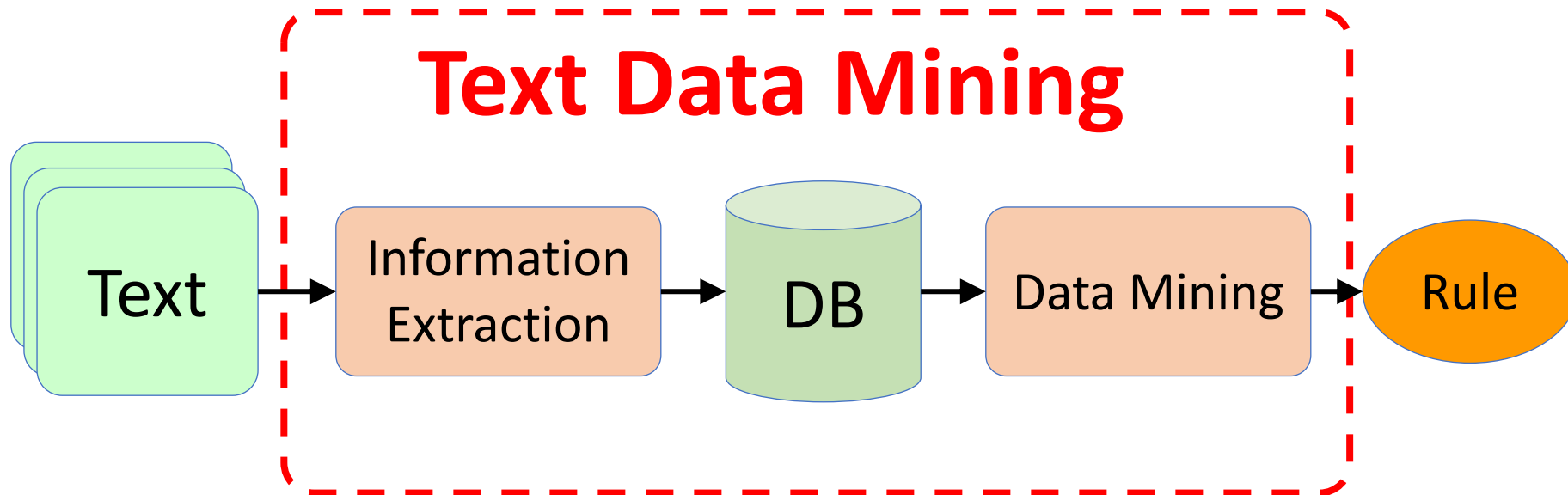
Text Mining:
the process of extracting
interesting and non-trivial
information and knowledge
from unstructured text.

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

An example of Text Mining



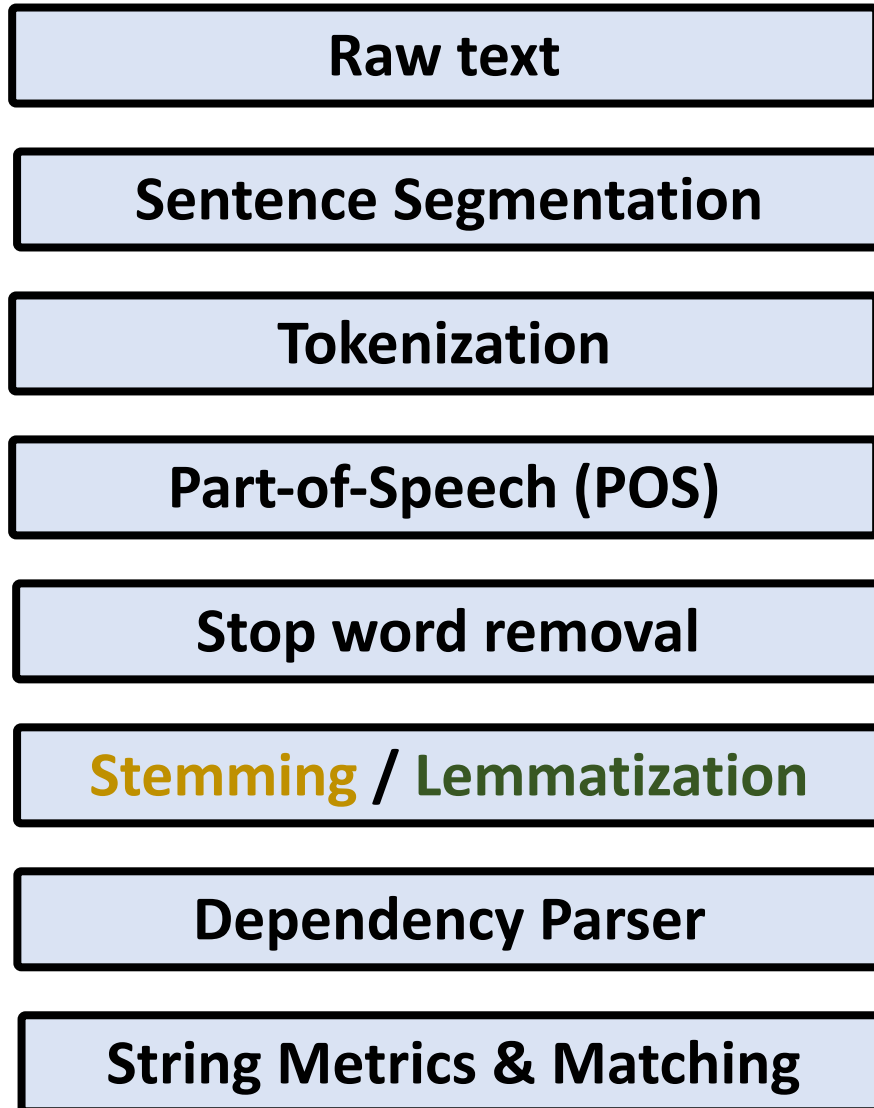
Overview of Information Extraction based Text Mining Framework



Natural Language Processing (NLP)

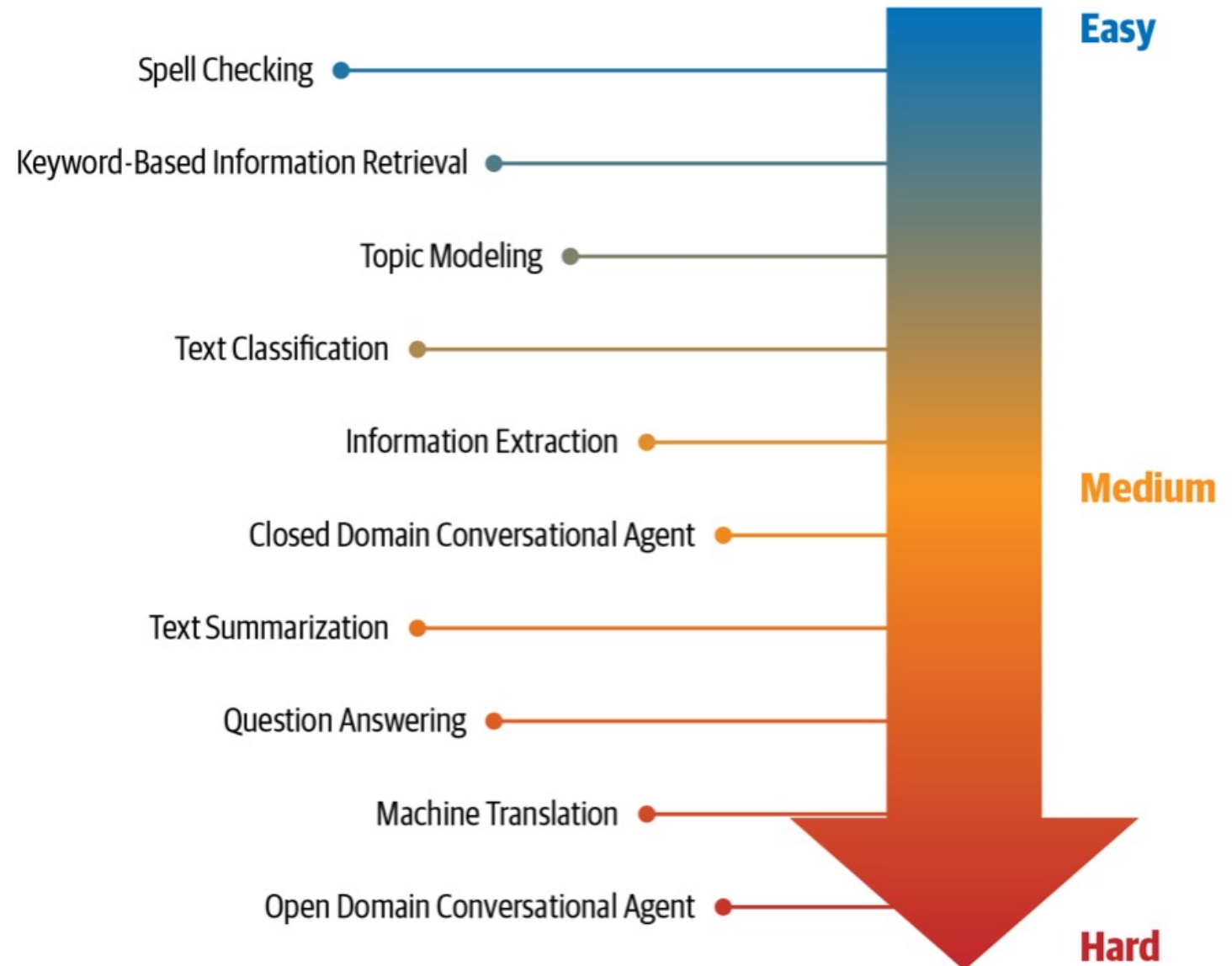
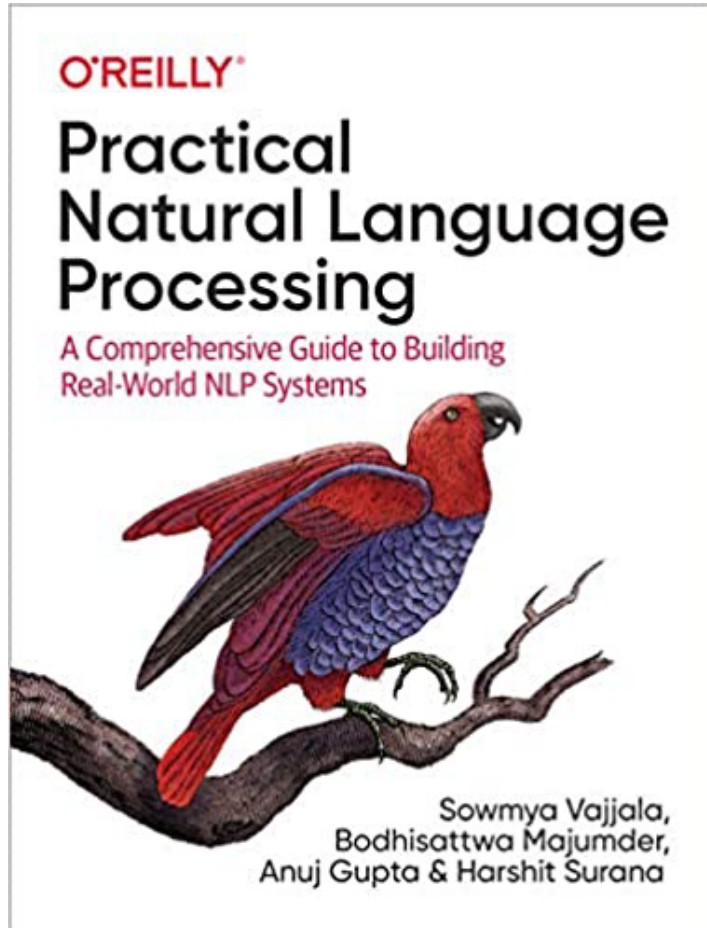
- **Natural language processing (NLP) is an important component of text mining and is a subfield of artificial intelligence and computational linguistics.**

Natural Language Processing (NLP) and Text Mining

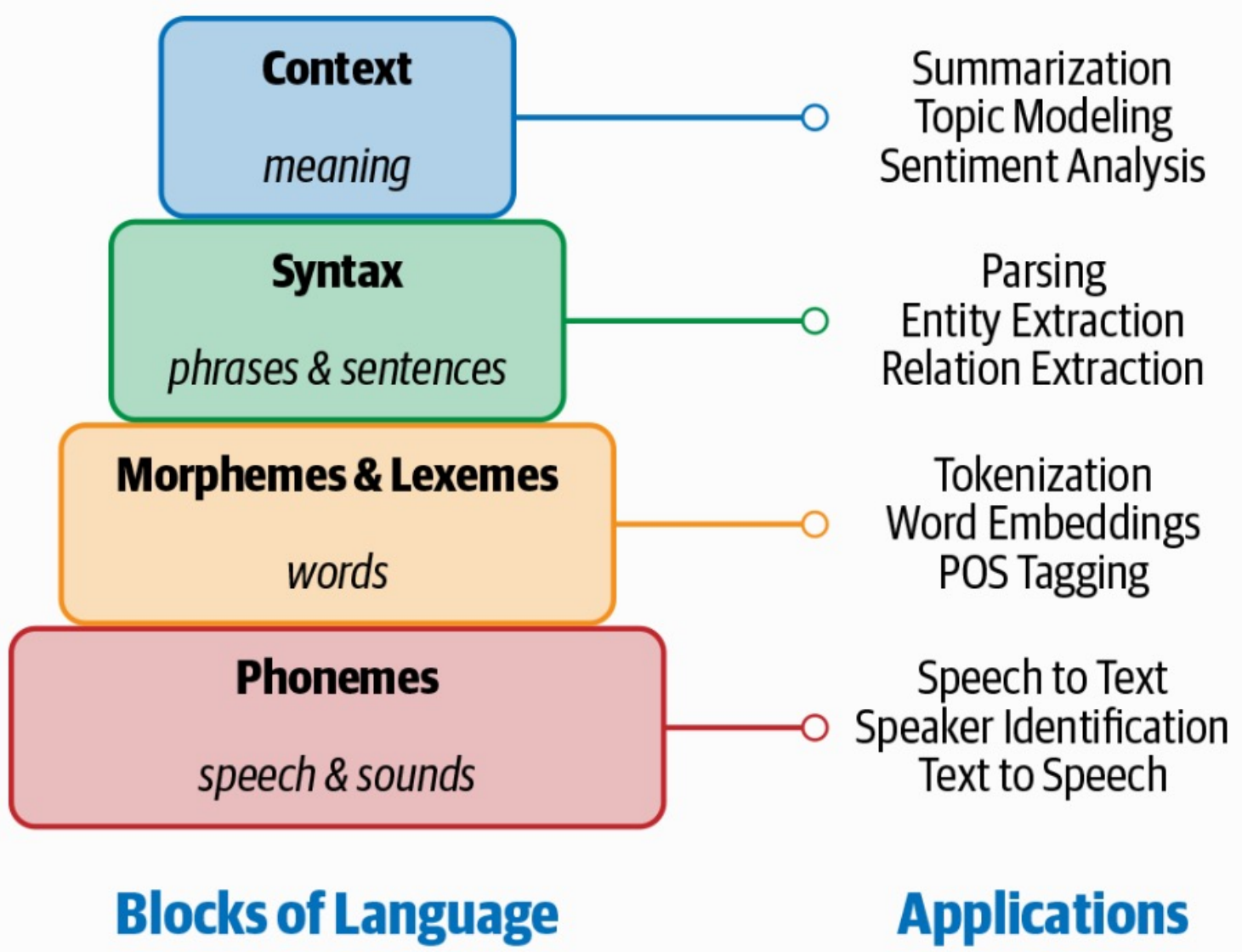


word's stem	word's lemma
am → am	am → be
having → hav	having → have

NLP Tasks



Building Blocks of Language and Applications



Source: Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta (2020), Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems, O'Reilly Media.

Morpheme Examples

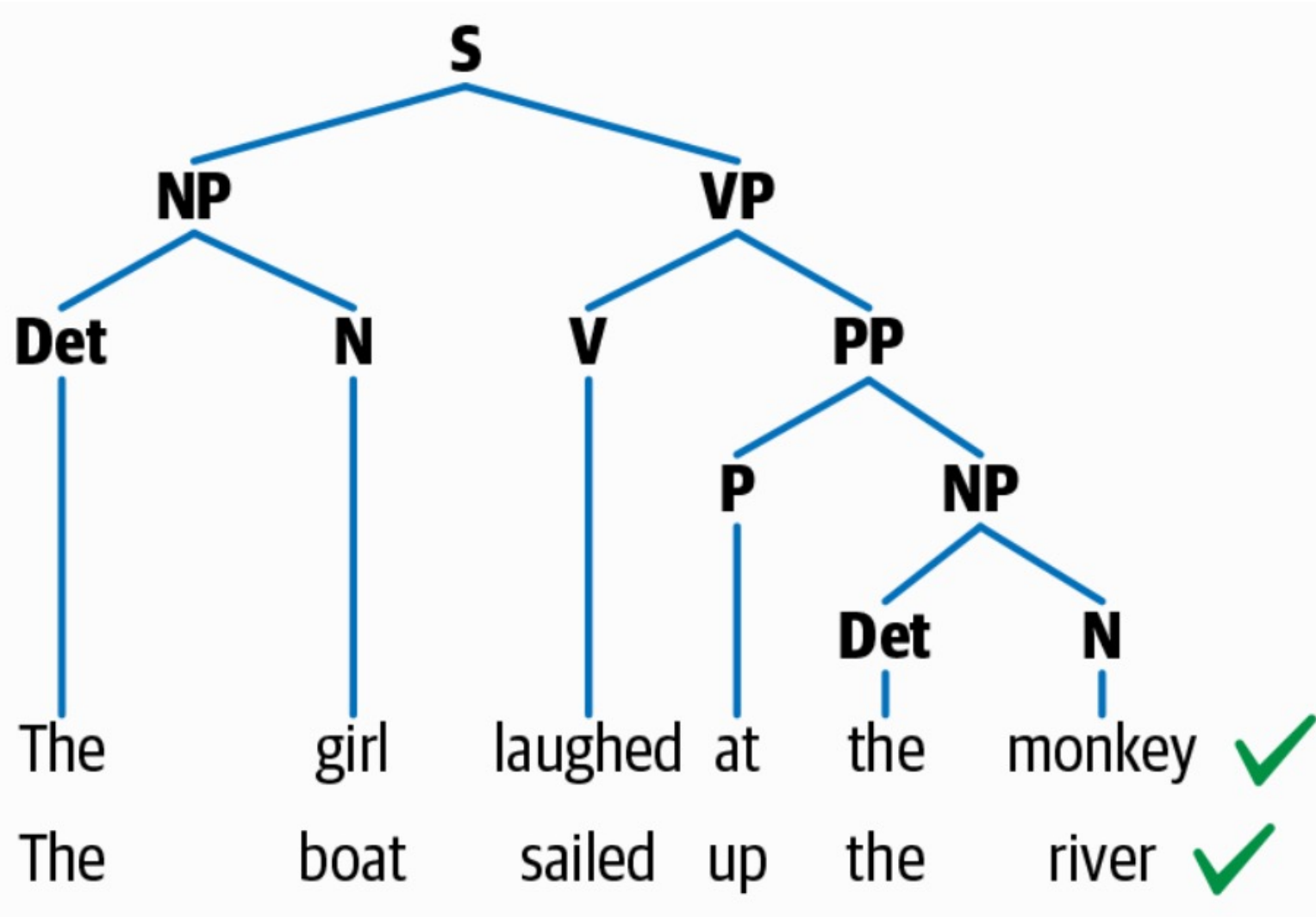
unbreakable
un + break + able

cats
cat + s

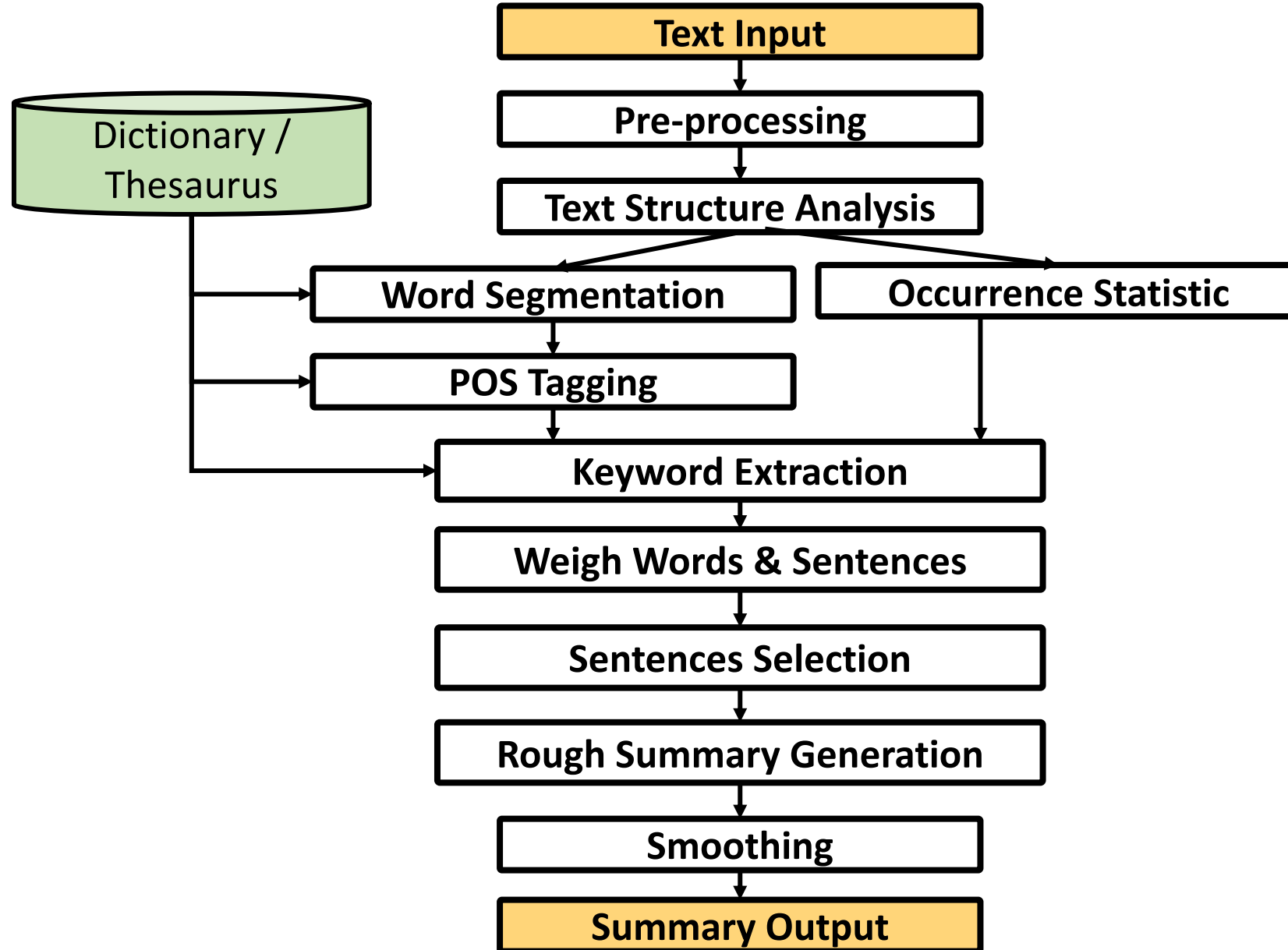
tumbling
tumble + ing

unreliability
un + rely + able + ity

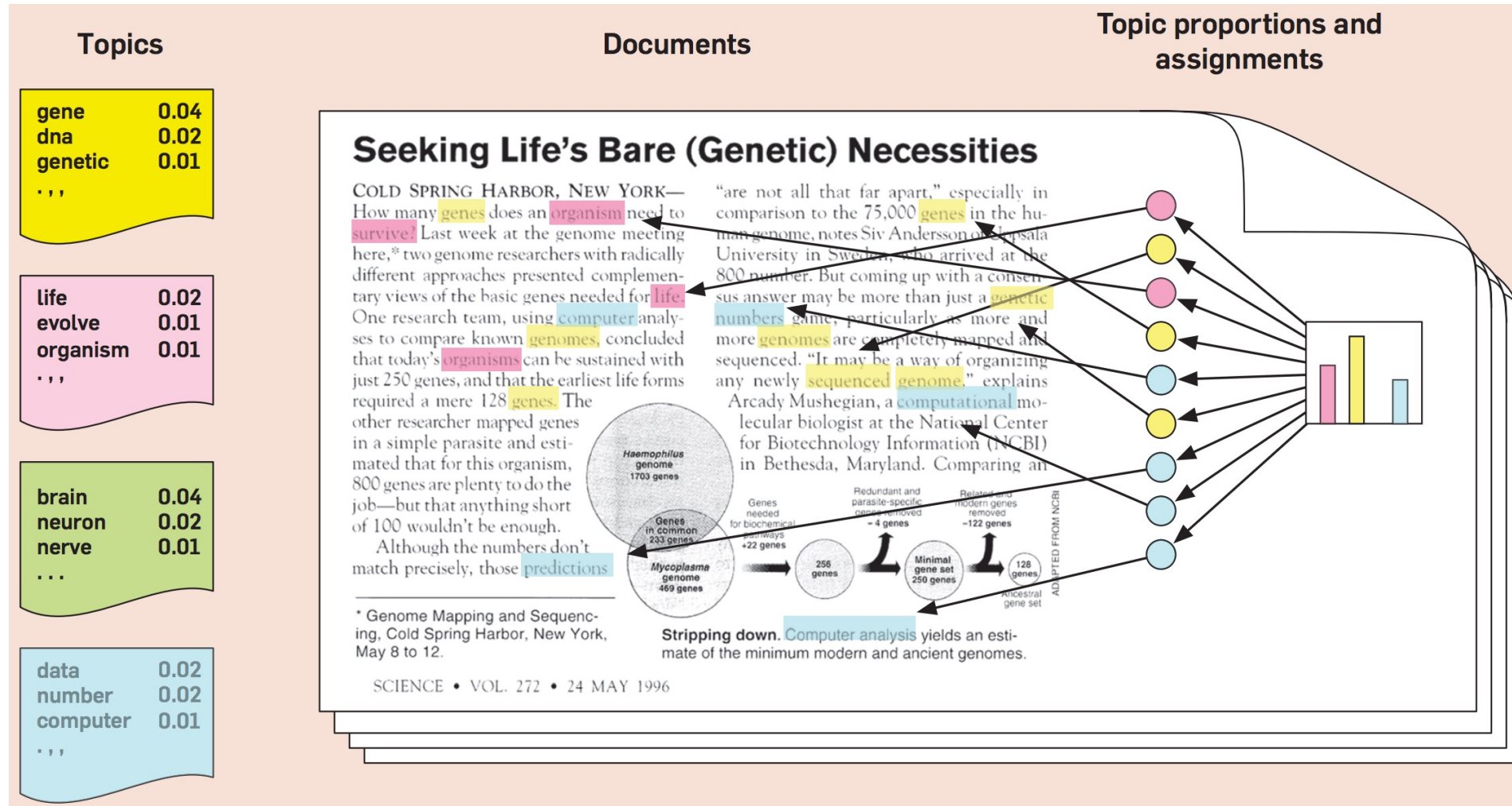
Syntactic Structure



Text Summarization



Topic Modeling



Natural Language Processing (NLP)

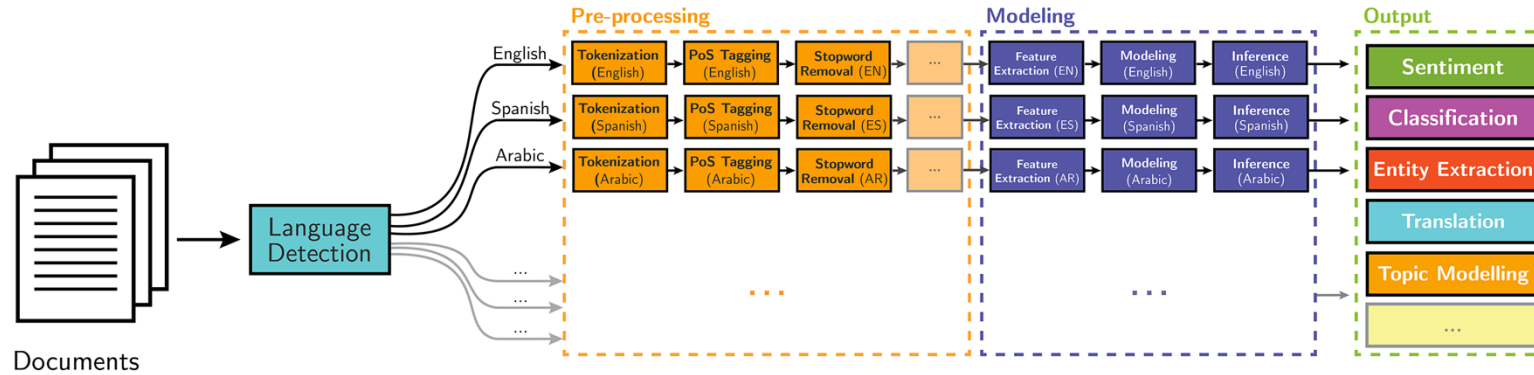
- **Part-of-speech tagging**
- **Text segmentation**
- **Word sense disambiguation**
- **Syntactic ambiguity**
- **Imperfect or irregular input**
- **Speech acts**

NLP Tasks

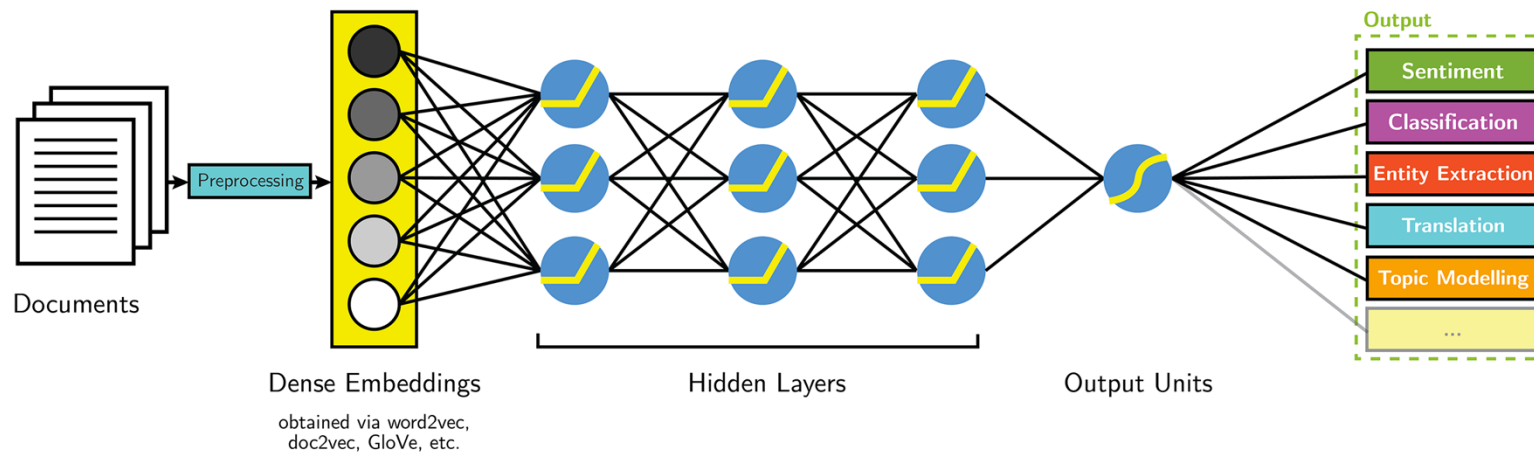
- **Question answering**
- **Automatic summarization**
- **Natural language generation**
- **Natural language understanding**
- **Machine translation**
- **Foreign language reading**
- **Foreign language writing.**
- **Speech recognition**
- **Text-to-speech**
- **Text proofing**
- **Optical character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

NLP

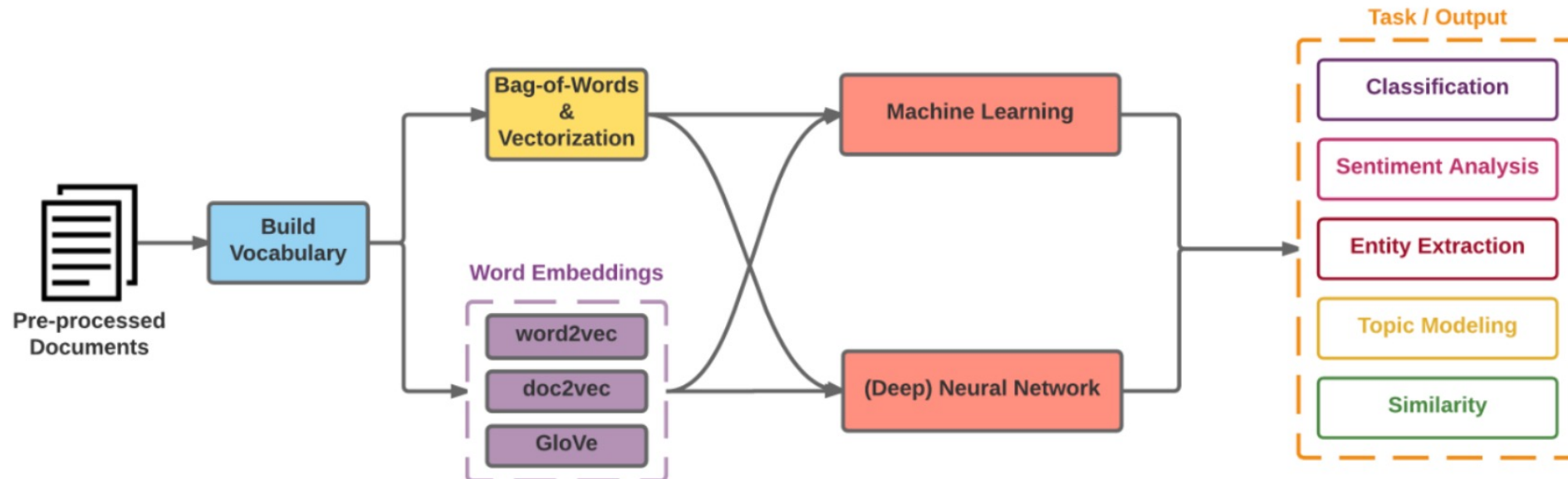
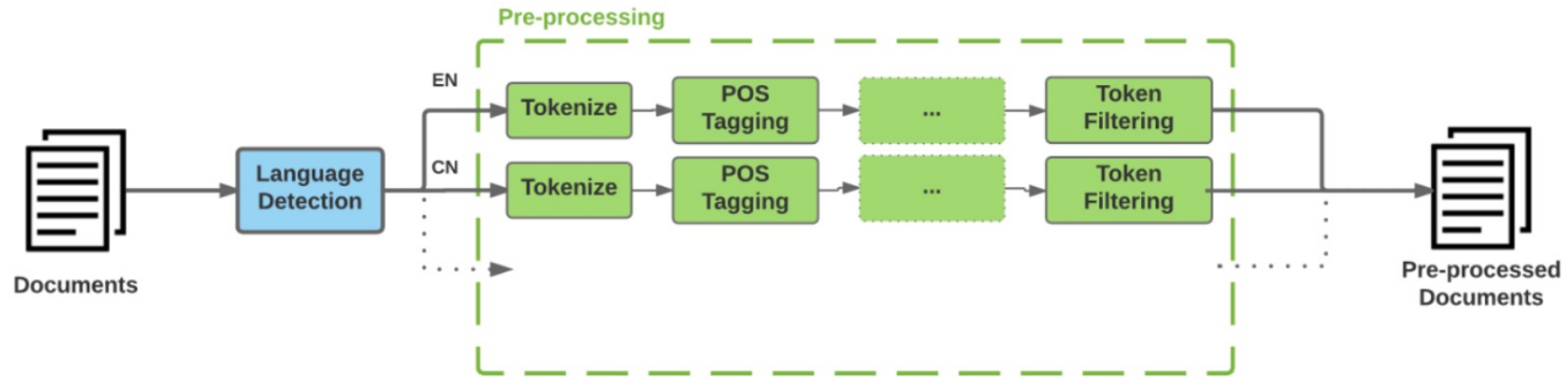
Classical NLP



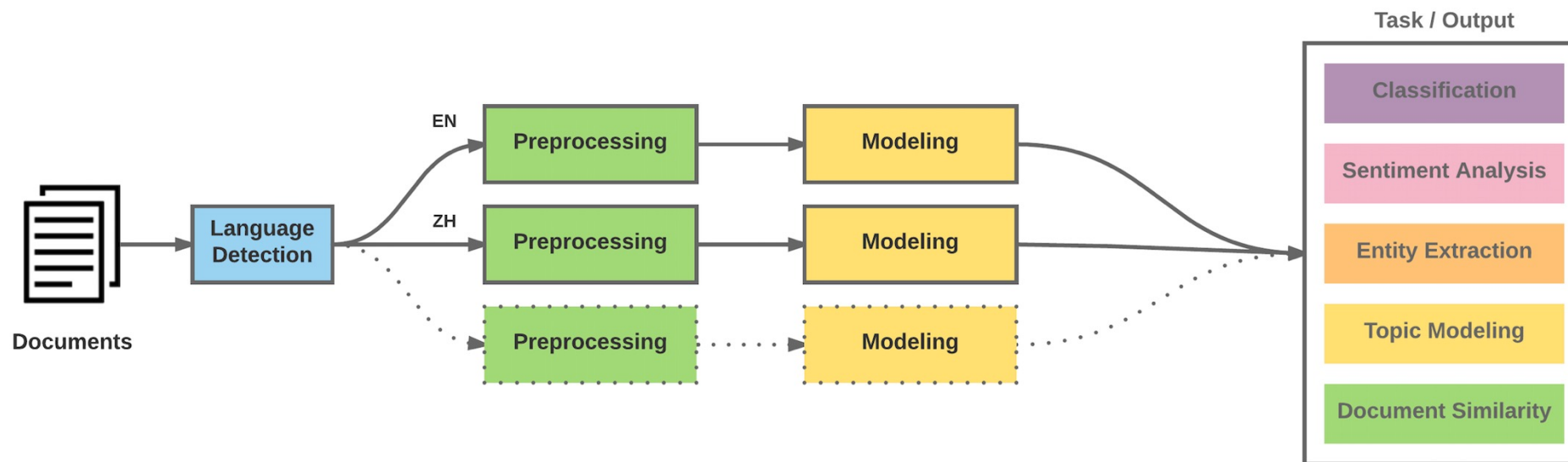
Deep Learning-based NLP



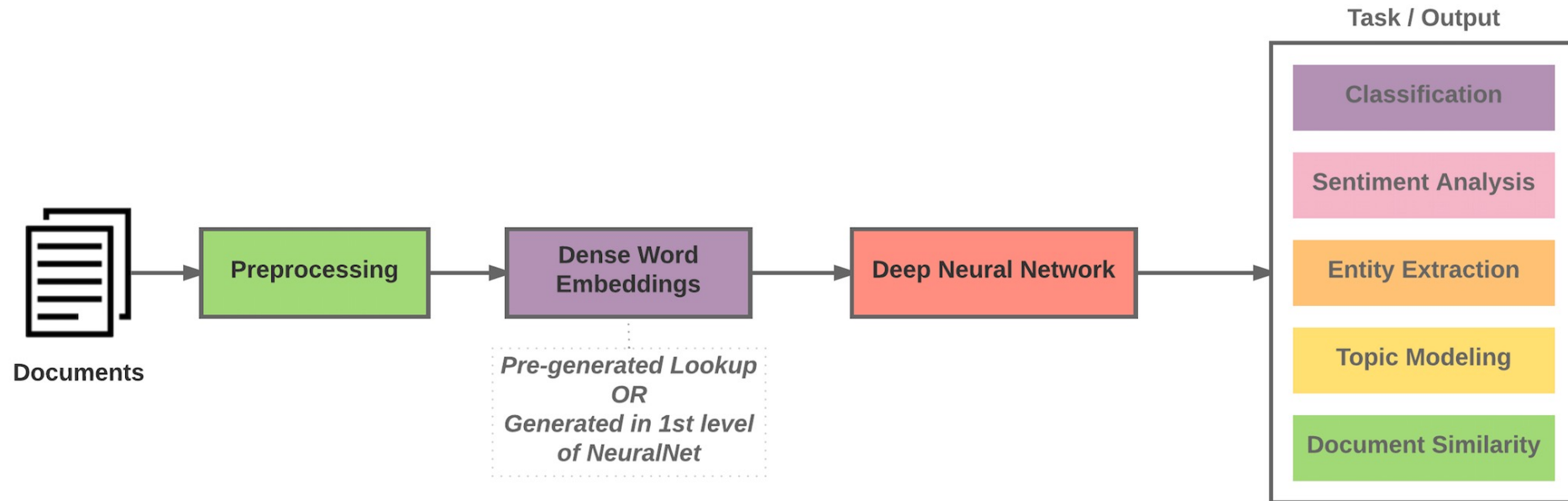
Modern NLP Pipeline



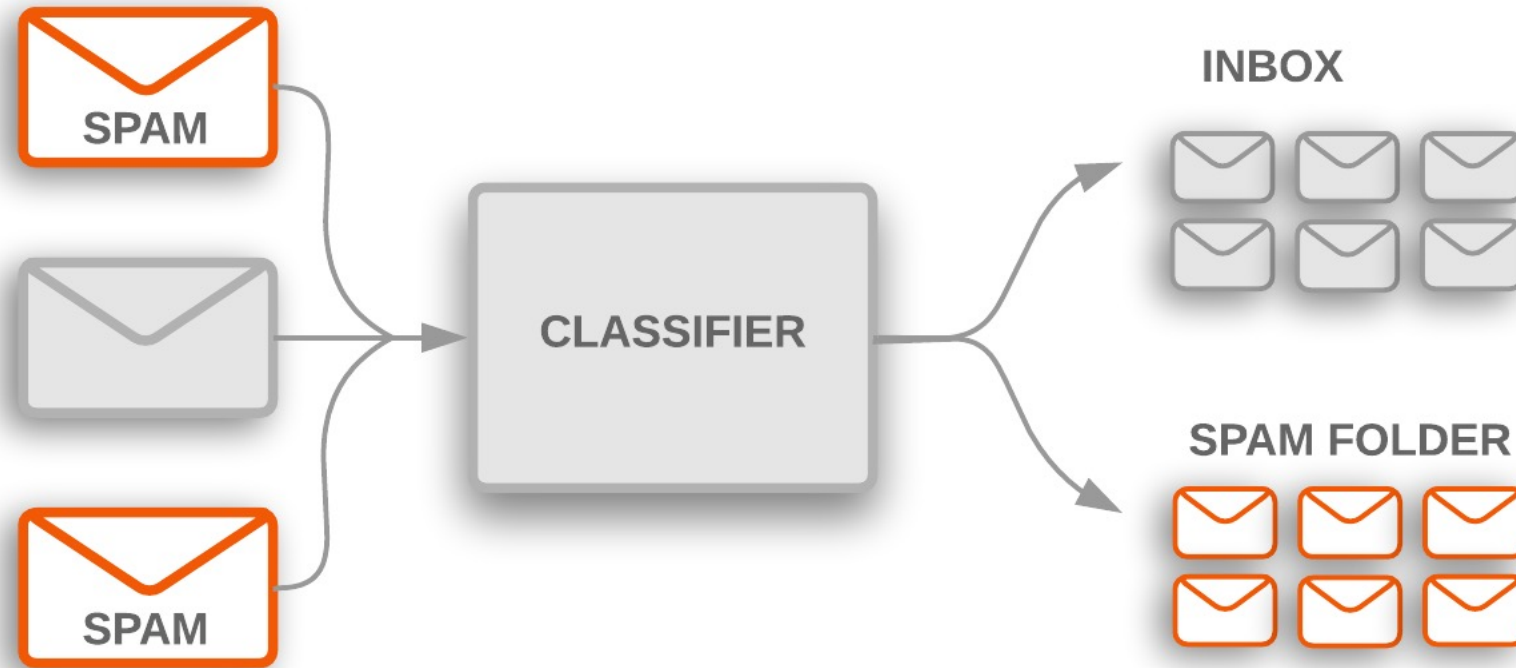
Modern NLP Pipeline



Deep Learning NLP



Text Classification

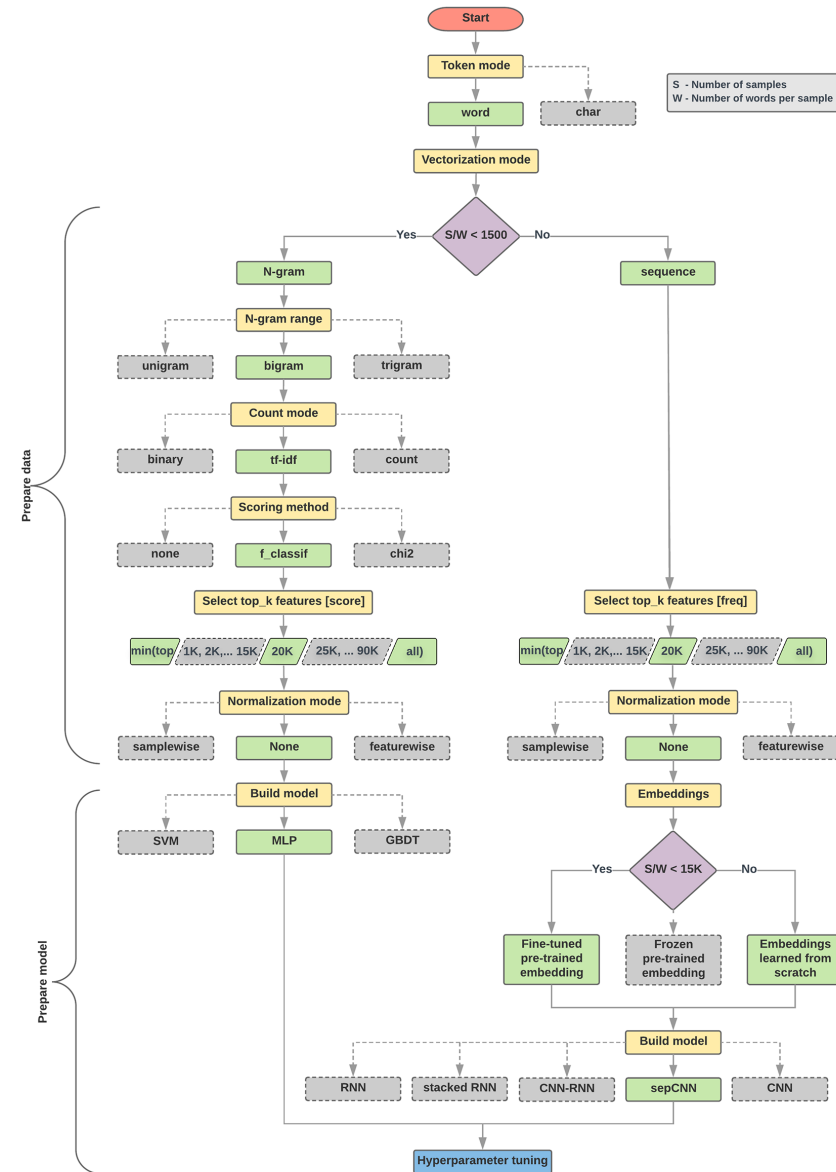


Text Classification Workflow

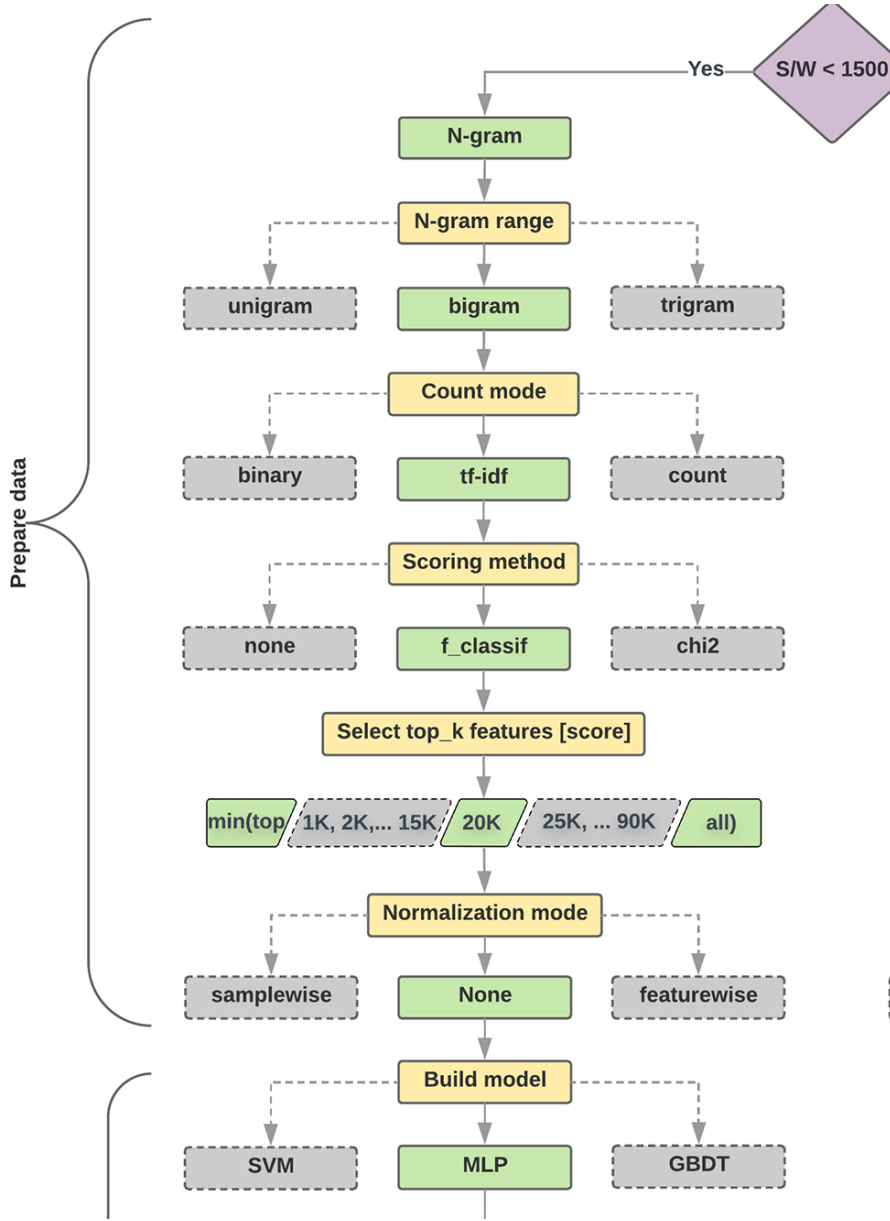
- **Step 1: Gather Data**
- **Step 2: Explore Your Data**
- **Step 2.5: Choose a Model***
- **Step 3: Prepare Your Data**
- **Step 4: Build, Train, and Evaluate Your Model**
- **Step 5: Tune Hyperparameters**
- **Step 6: Deploy Your Model**



Text Classification Flowchart

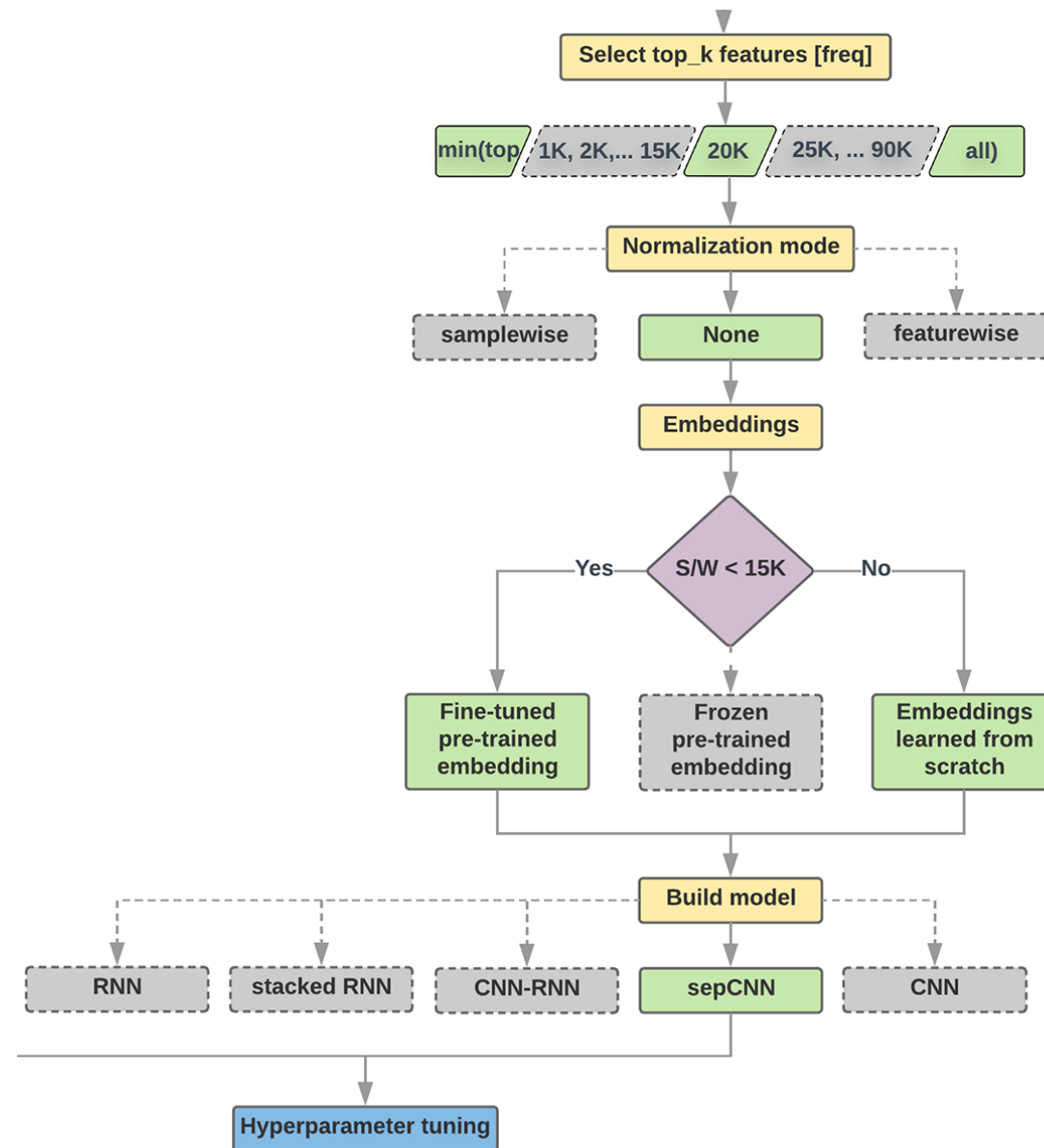


Text Classification S/W < 1500: N-gram



Source: <https://developers.google.com/machine-learning/guides/text-classification/step-2-5>

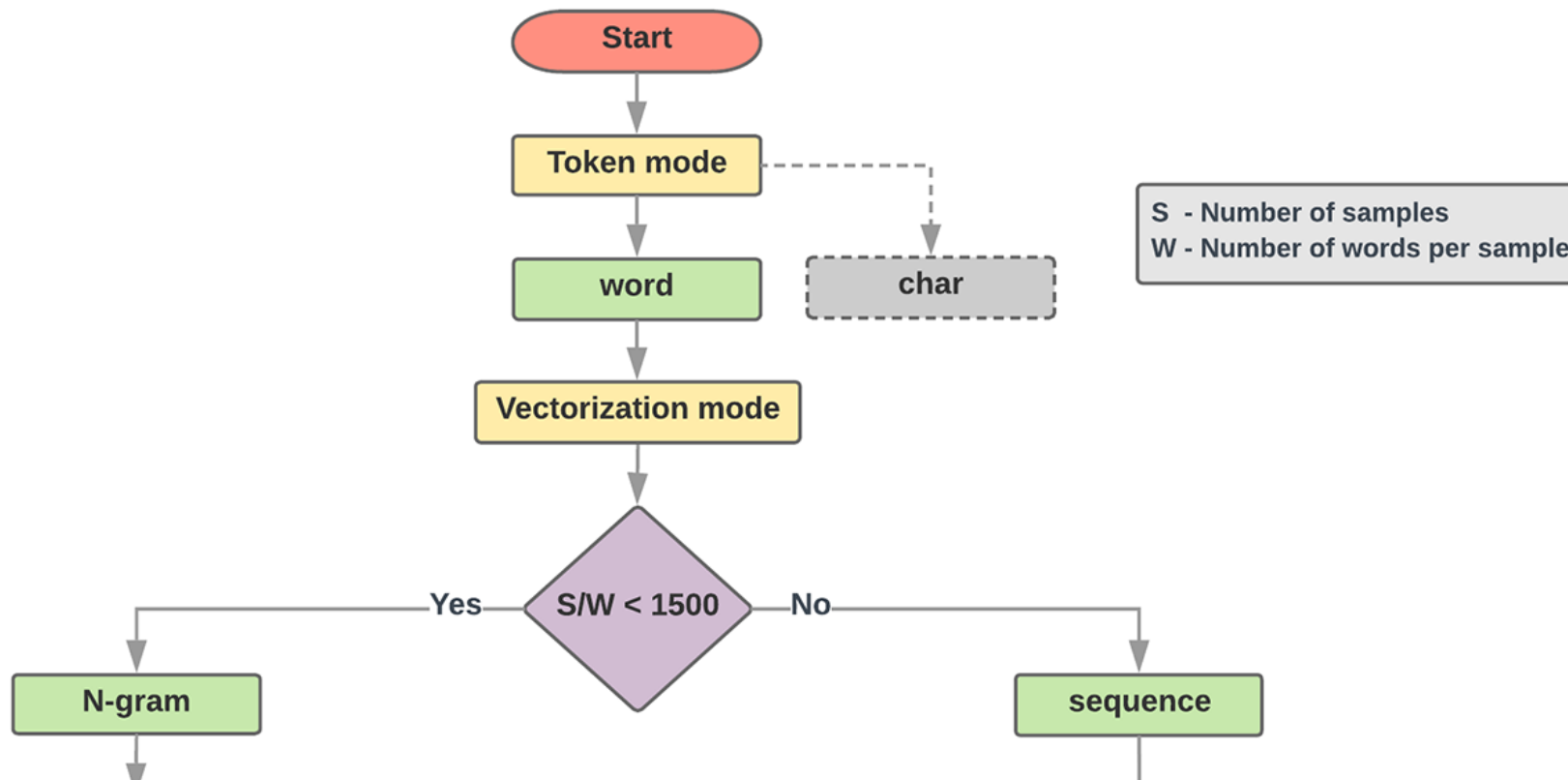
Text Classification $S/W \geq 1500$: Sequence



Step 2.5: Choose a Model

Samples/Words < 1500

150,000/100 = 1500

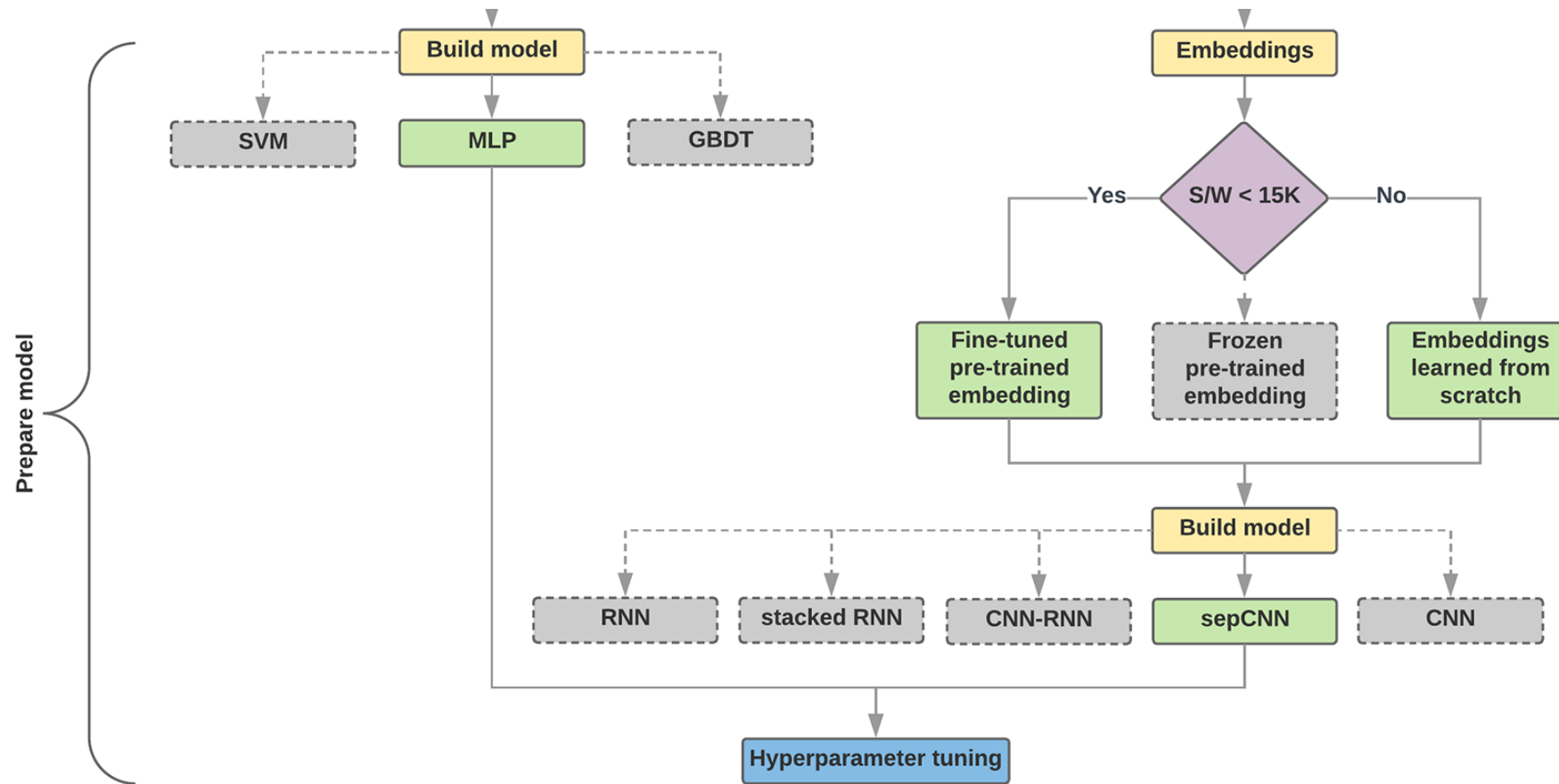


IMDb review dataset,
the samples/words-per-sample ratio is ~ 144

Step 2.5: Choose a Model

Samples/Words < 15,000

1,500,000/100 = 15,000



Step 3: Prepare Your Data

Texts:

T1: 'The mouse ran up the clock'

T2: 'The mouse ran down'

Token Index:

```
{'the': 1, 'mouse': 2, 'ran': 3, 'up': 4, 'clock': 5, 'down': 6,}
```

NOTE: 'the' occurs most frequently,
so the index value of 1 is assigned to it.
Some libraries reserve index 0 for unknown tokens,
as is the case here.

Sequence of token indexes:

T1: 'The mouse ran up the clock' =

```
[1, 2, 3, 4, 1, 5]
```

T2: 'The mouse ran down' =

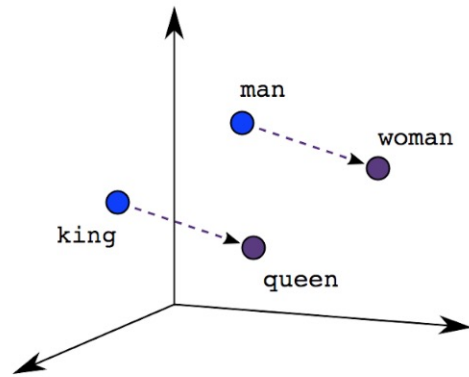
```
[1, 2, 3, 6]
```


One-hot encoding

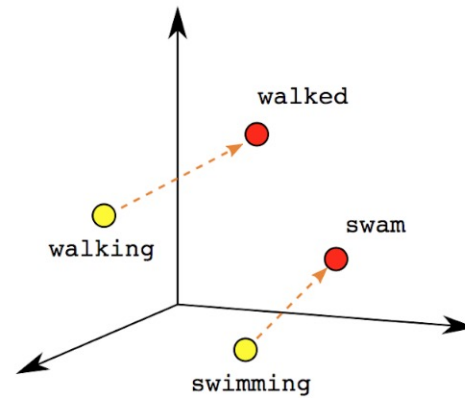
'The mouse ran up the clock' =

The	1	[[0, 1, 0, 0, 0, 0, 0],
mouse	2		[0, 0, 1, 0, 0, 0, 0],
ran	3		[0, 0, 0, 1, 0, 0, 0],
up	4		[0, 0, 0, 0, 1, 0, 0],
the	1		[0, 1, 0, 0, 0, 0, 0],
clock	5		[0, 0, 0, 0, 0, 1, 0]]
			[0, 1, 2, 3, 4, 5, 6]

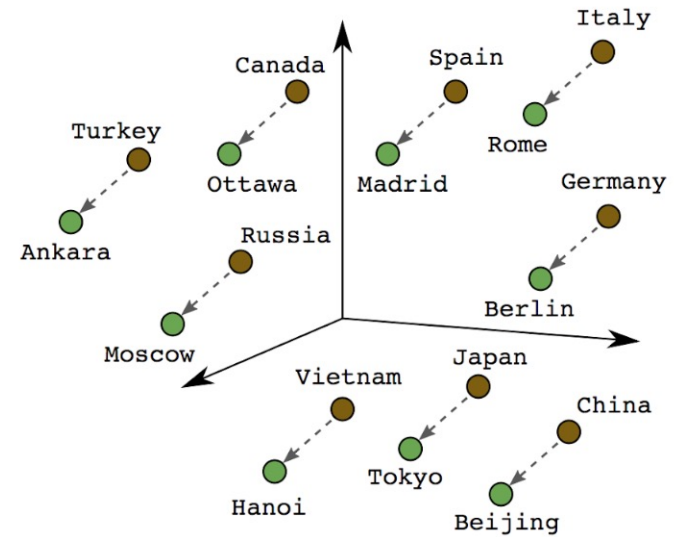
Word embeddings



Male-Female

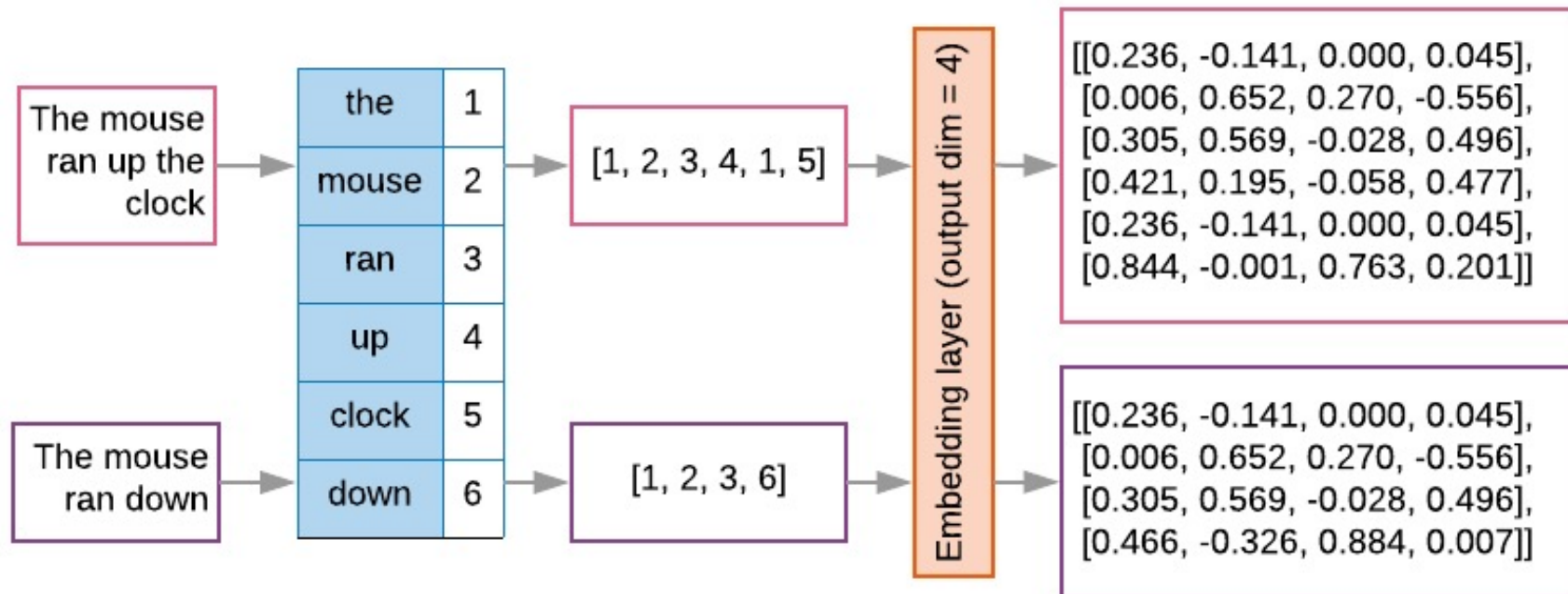


Verb Tense



Country-Capital

Word embeddings



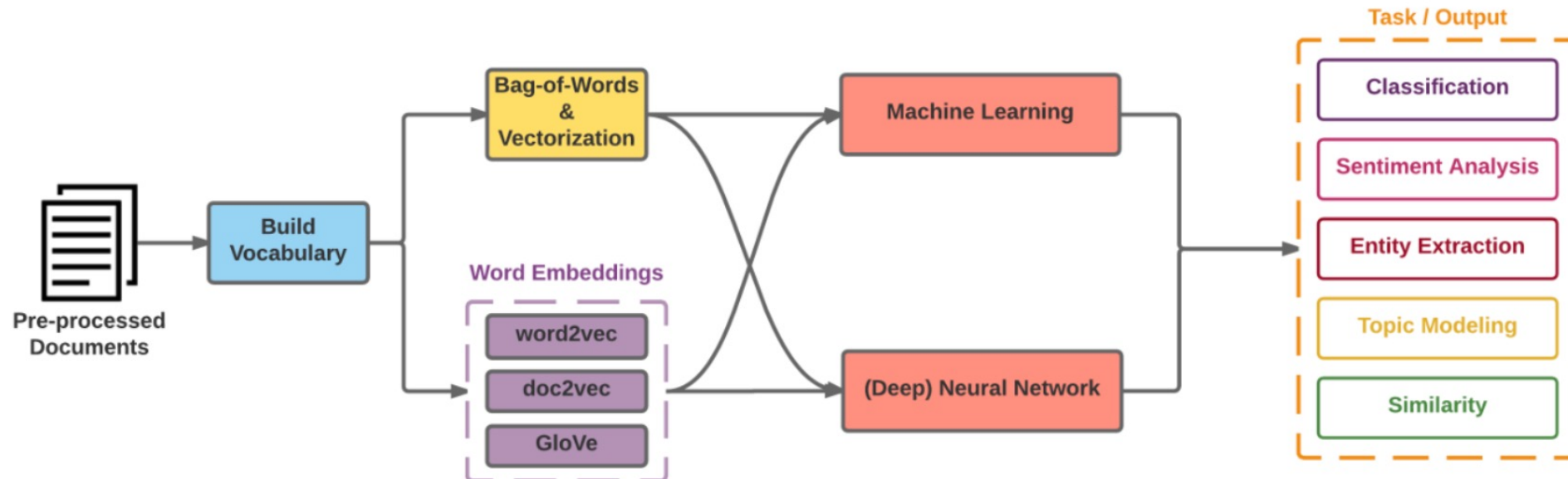
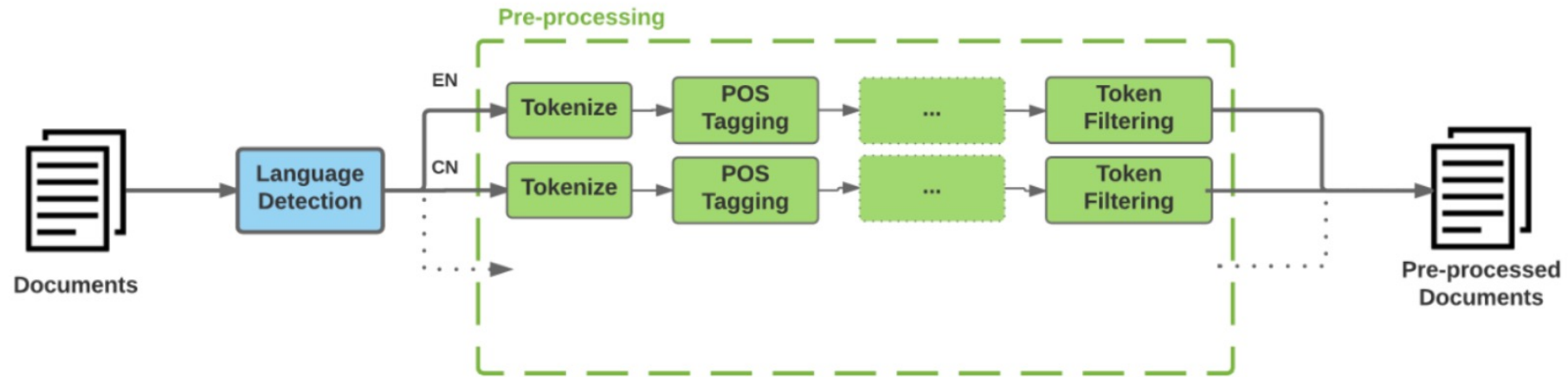
Vector Representations of Words

Word Embeddings

Word2Vec

GloVe

Modern NLP Pipeline



Facebook Research FastText

Pre-trained word vectors

Word2Vec

wiki.zh.vec (861MB)

332647 word

300 vec

Pre-trained word vectors for 90 languages,
trained on Wikipedia using fastText.

These vectors in dimension 300 were obtained using the
skip-gram model with default parameters.

<https://github.com/facebookresearch/fastText/blob/master/pretrained-vectors.md>

Source: Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov. "Enriching word vectors with subword information." *arXiv preprint arXiv:1607.04606* (2016).

Facebook Research FastText Word2Vec:

wiki.zh.vec

(861MB) (332647 word 300 vec)

wiki.zh.vec

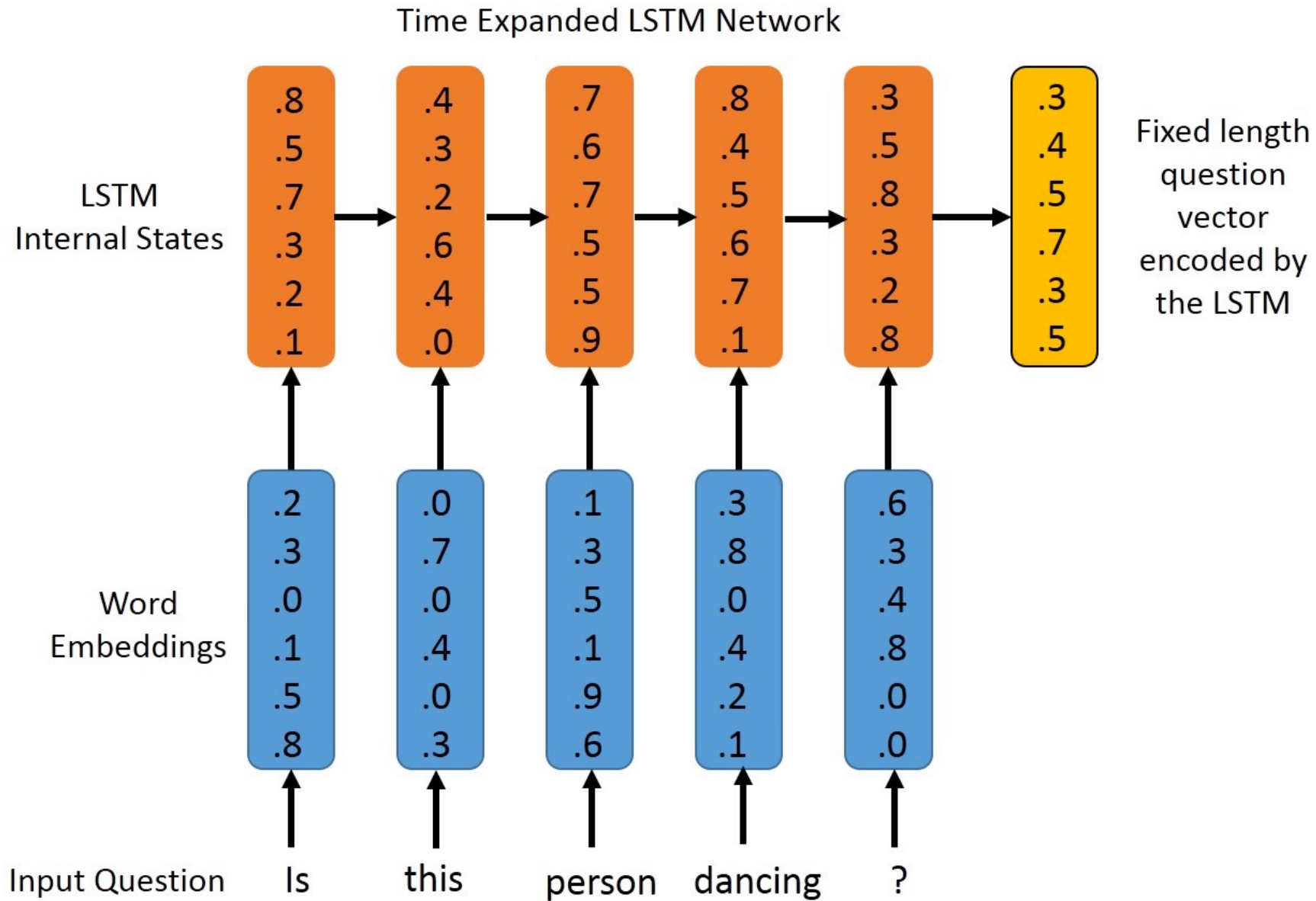
31845	yg	-0.3978	0.49084	-0.54621	0.078991	0.8584	-0.26163	-0.45787	0.060828	0.36513	-0.03771	0.80791	0.16613	1.4828	-0.89862	0.085965
31846	迴圈	-0.034834	0.71651	-0.4377	0.48344	0.31117	-0.51783	-0.40156	-0.057097	0.31535	-0.088301	0.23436	0.30884	1.2932	-0.6704	0.215
31847	ぶっ	-0.23267	0.39349	-0.90806	-0.53805	0.59308	-0.31819	-0.64229	0.16871	0.10086	0.09342	1.0914	-0.16019	1.6954	-0.70604	-0.2186
31848	三公	0.54129	0.55641	-0.4348	0.25094	0.1631	-0.10326	-0.54099	0.064742	0.13175	0.10217	0.84938	-0.10287	1.312	-0.74969	0.24025
31849	水貨	-0.14451	0.80455	-0.6145	0.55905	0.58307	-0.02559	-0.41088	-0.19056	-0.09178	0.33935	1.1927				
31850	刚才	0.19347	0.553	-0.64736	0.26358	0.83816	-0.24098	-0.83997	-0.16232	-0.024786	-0.2483	0.69732				
31851	無知	-0.0089777	0.90866	-0.25306	0.72983	0.67791	-0.3285	-0.63835	0.075295	0.4774	-0.04134	0.7210				
31852	好轉	-0.026068	0.92676	-0.47469	0.50129	0.67343	-0.32509	-0.32917	0.066499	0.3875	0.0011722	0.66				
31853	紀事	0.40541	0.67654	-0.5351	0.30329	0.43042	-0.24675	-0.19287	0.34207	0.35516	-0.076331	0.85916				
31854	變回	-0.089933	0.88136	-0.43524	0.59963	0.6403	-0.70981	-0.56788	-0.074018	0.16905	-0.086594	0.6				
31855	牟尼	-0.26578	0.6434	0.028982	-0.044001	0.88297	-0.17646	-0.64672	0.040483	0.43653	0.084908	0.74				
31856	埋藏	-0.0985	0.85082	-0.33363	0.24784	0.71518	-0.59054	-0.73731	0.050949	0.36726	-0.076886	0.817				
31857	正大	0.21069	0.27605	-0.83862	-0.099698	0.47894	-0.32196	-0.38288	-0.01892	0.40548	-0.029619	0.7				
31858	kis	-0.30595	0.18482	-0.71287	-0.314	0.44776	-0.44245	-0.36447	-0.23723	0.00098801	-0.2528	0.606				
31859	合奏	0.1841	0.60874	-0.51376	-0.48002	0.21506	-0.55515	-0.71746	0.030735	0.39508	-0.40856	0.6226				
31860	精兵	0.25619	0.77186	-0.48847	0.23118	0.27254	0.21305	-0.3517	0.47305	0.24882	-0.34756	1.025	0.11			
31861	疲勞	-0.072521	1.0381	-0.51933	0.19421	0.67573	-0.45204	-0.20126	0.22704	0.44196	0.018401	0.3473				
31862	襪	-0.11771	1.4272	-1.0849	0.77532	0.87026	-0.6892	-0.3521	0.036517	0.42727	-0.1871	0.82789	-0.0			
31863	小貓	-0.21554	0.73988	-0.39628	0.044656	1.0602	-0.67047	-0.54102	0.11888	0.1693	0.19343	1.0841	0.			
31864	lai	-0.25451	0.31596	-0.29228	-0.19144	0.99059	-0.24459	-0.66342	0.063093	-0.061142	-0.22749	0.6				
31865	偏東	-0.50835	1.0943	0.043918	0.29173	1.0161	-0.32493	-0.27305	0.026946	0.46811	-0.3874	1.4049	0.			
31866	大约是	-0.35726	-0.03476	-0.28672	0.075447	0.18175	-0.39421	-0.32088	0.025225	0.34808	0.074744	0.				
31867	franch	-0.6046	-0.3235	0.024041	-0.2756	0.74761	-0.14654	0.0082566	-0.10071	0.53593	-0.17374	0.2				
31868	brazilian	-0.54029	-0.63905	-0.094006	-0.68768	0.33263	-0.1583	-0.060424	0.20644	0.46234	-0.0764					
31869	夹竹桃	-0.4361	0.011429	-0.078896	-0.078186	0.37747	-0.052101	-0.096683	0.10769	0.62661	-0.37252					
31870	continent	-0.37761	-0.72151	-0.42248	-0.81768	0.5016	-0.48569	0.13464	0.12644	0.32292	0.18099	0.				
31871	我还是	0.097443	0.28929	-0.14202	0.034027	0.50621	-0.1647	-0.45849	-0.16198	0.13965	-0.33451	0.61				
31872	vienna	-0.25827	-0.050966	0.050502	-0.63466	0.4949	-0.17448	-0.59978	0.20269	0.37532	0.059419	0.				
31873	固态	-0.12678	0.4556	-0.27108	0.12506	0.52106	-0.058477	-0.69296	0.12162	0.26508	-0.089028	0.752				
31874	吉普	-0.33693	0.48335	-0.58455	0.13722	0.74856	-0.24529	-0.41125	-0.13832	0.33871	-0.12051	0.864				
31875	實物	0.030096	0.65756	-0.67982	0.22203	0.38492	-0.19001	-0.53136	-0.10322	0.24523	0.15287	0.92591				
31876	教职	0.11559	0.67087	-0.5111	0.14955	0.61417	-0.51571	-0.47901	0.29445	0.37629	-0.24232	0.4608	-0.			
31877	惕	0.50469	1.5357	-0.64393	0.48668	0.69479	-0.23443	-0.47863	0.16288	0.3347	-0.51673	0.86777	0.0			
31878	岸上	0.088323	0.85815	-0.485	0.30383	0.75965	-0.25031	-0.76678	0.12805	0.37641	-0.088752	0.65012				
31879	议和	0.26835	0.94854	-0.27972	0.097623	0.43305	-0.031361	-0.57406	0.21608	0.3324	-0.36823	0.6987				
31880	aka	-0.21332	0.11216	-0.48872	-0.18531	0.79093	-0.34221	-0.51122	0.10067	0.29963	-0.075253	0.642				
31881	滑鐵盧	-0.28726	0.88014	-0.39751	-0.056992	0.37408	-0.16967	-0.20673	-0.048533	-0.1978	-0.13107	0.				

Models

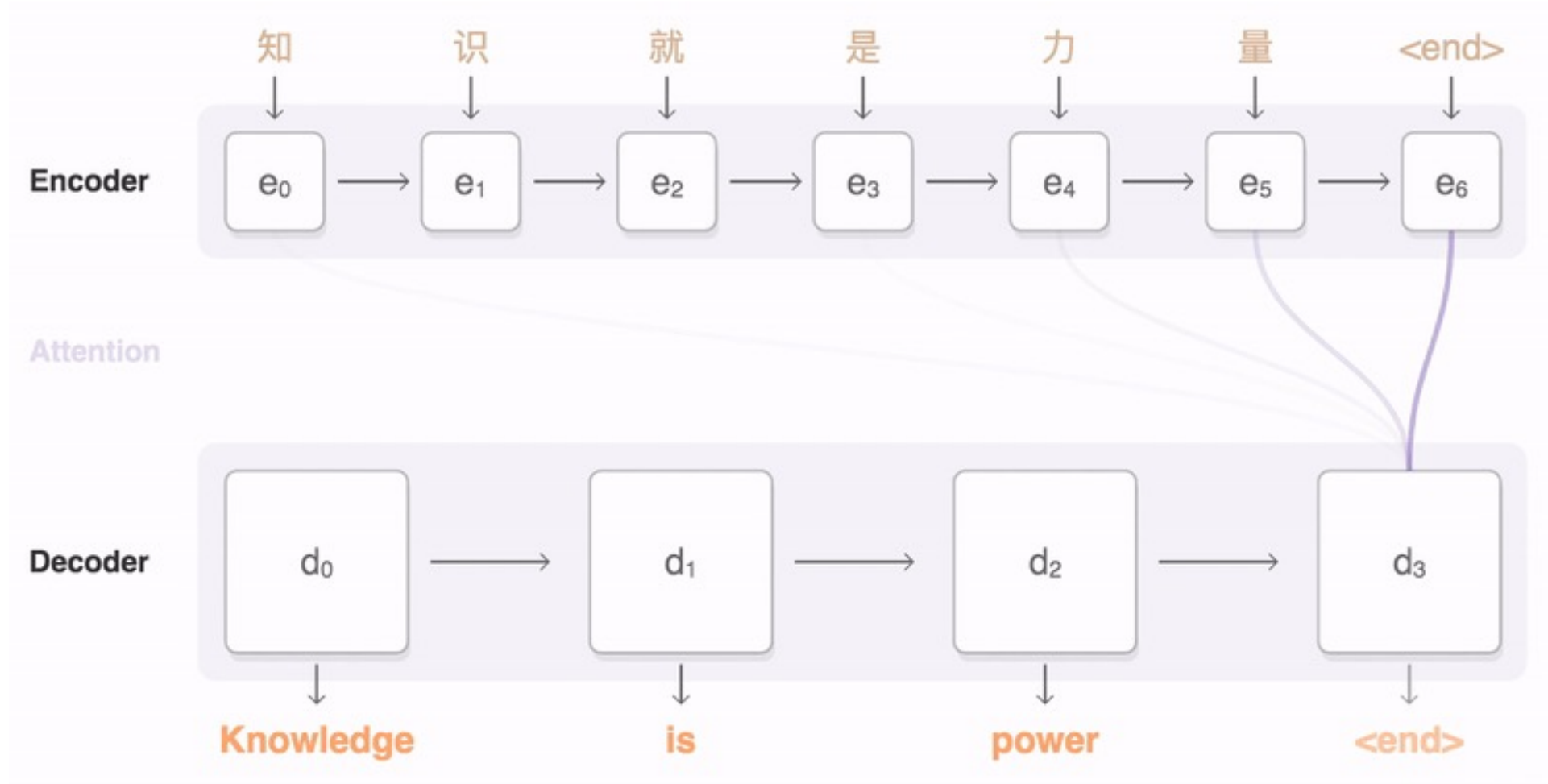
The models can be downloaded from:

- Afrikaans: [bin+text, text](#)
- Albanian: [bin+text, text](#)
- Arabic: [bin+text, text](#)
- Armenian: [bin+text, text](#)
- Asturian: [bin+text, text](#)
- Azerbaijani: [bin+text, text](#)
- Bashkir: [bin+text, text](#)
- Basque: [bin+text, text](#)
- Belarusian: [bin+text, text](#)
- Bengali: [bin+text, text](#)
- Bosnian: [bin+text, text](#)
- Breton: [bin+text, text](#)
- Bulgarian: [bin+text, text](#)
- Burmese: [bin+text, text](#)
- Catalan: [bin+text, text](#)
- Cebuano: [bin+text, text](#)
- Chechen: [bin+text, text](#)
- Chinese: [bin+text, text](#)
- Chuvash: [bin+text, text](#)
- Croatian: [bin+text, text](#)
- Czech: [bin+text, text](#)

Word Embeddings in LSTM RNN

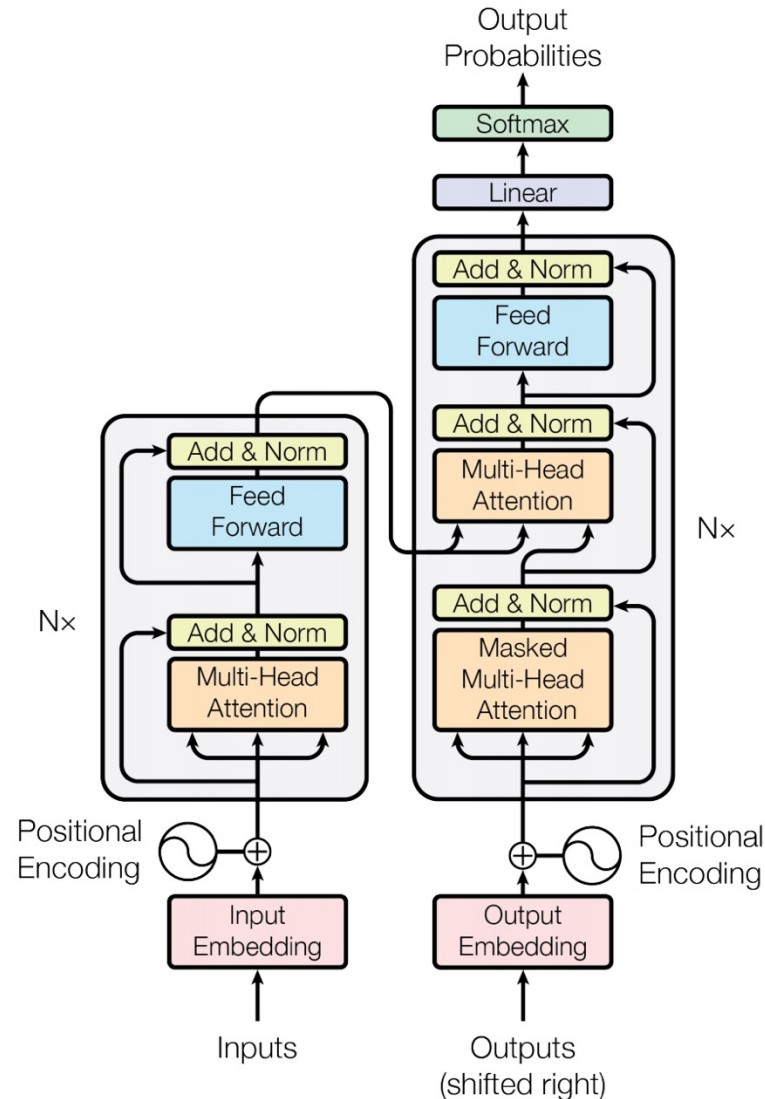


Sequence to Sequence (Seq2Seq)



Transformer (Attention is All You Need)

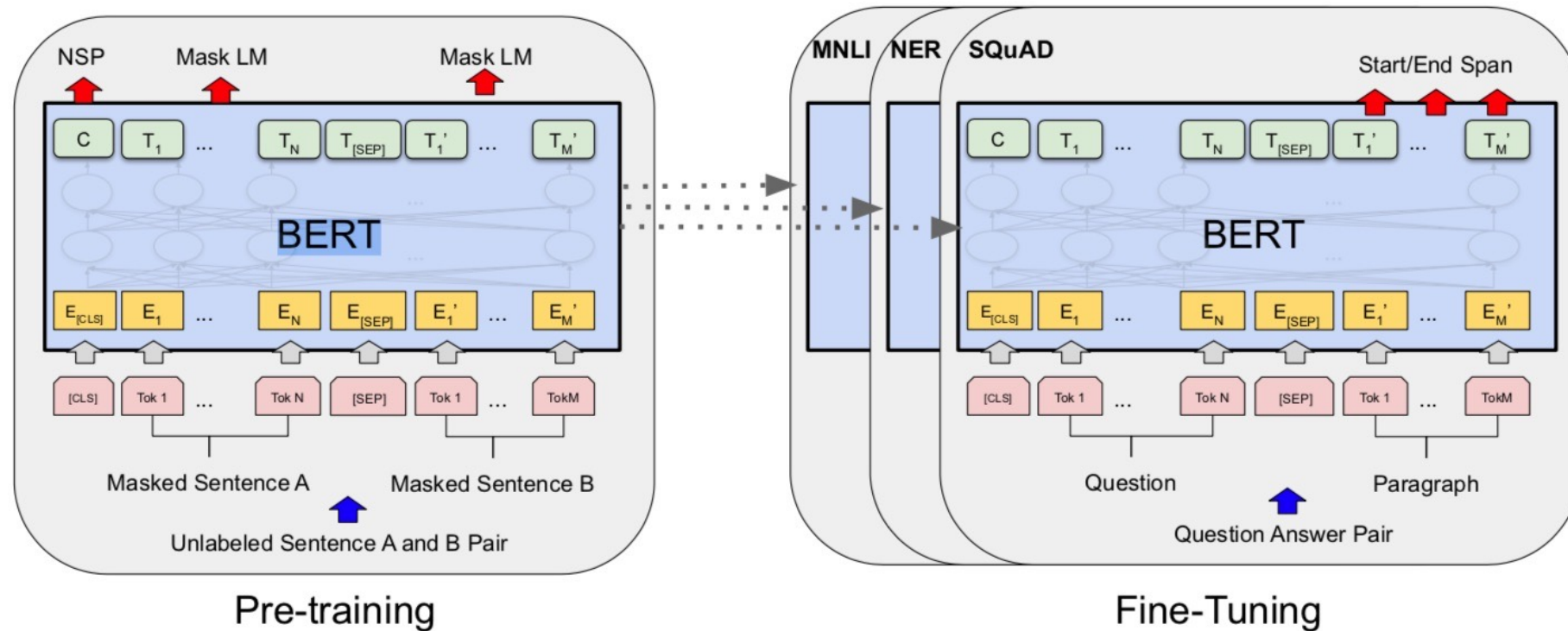
(Vaswani et al., 2017)



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



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:

Pre-training of Deep Bidirectional Transformers for Language Understanding

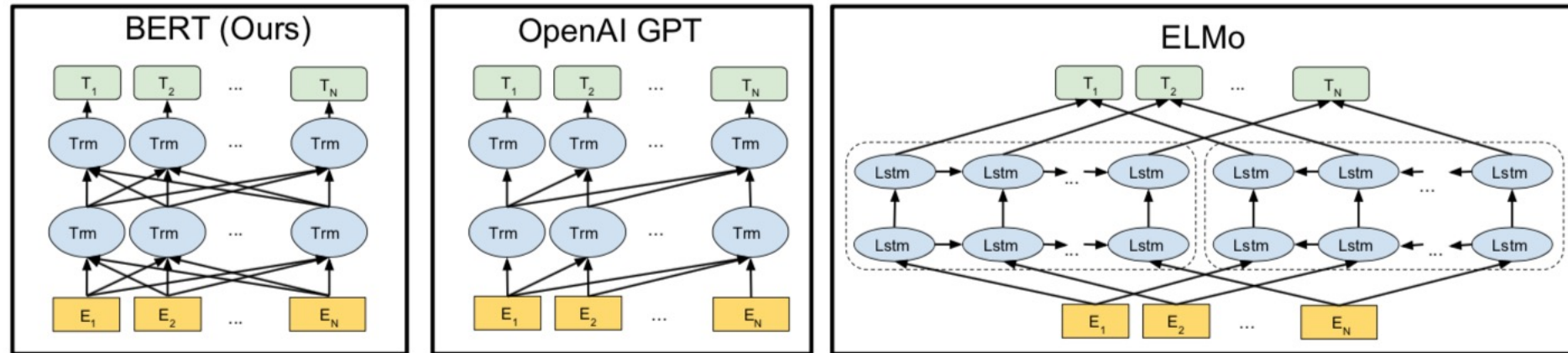
BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

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.

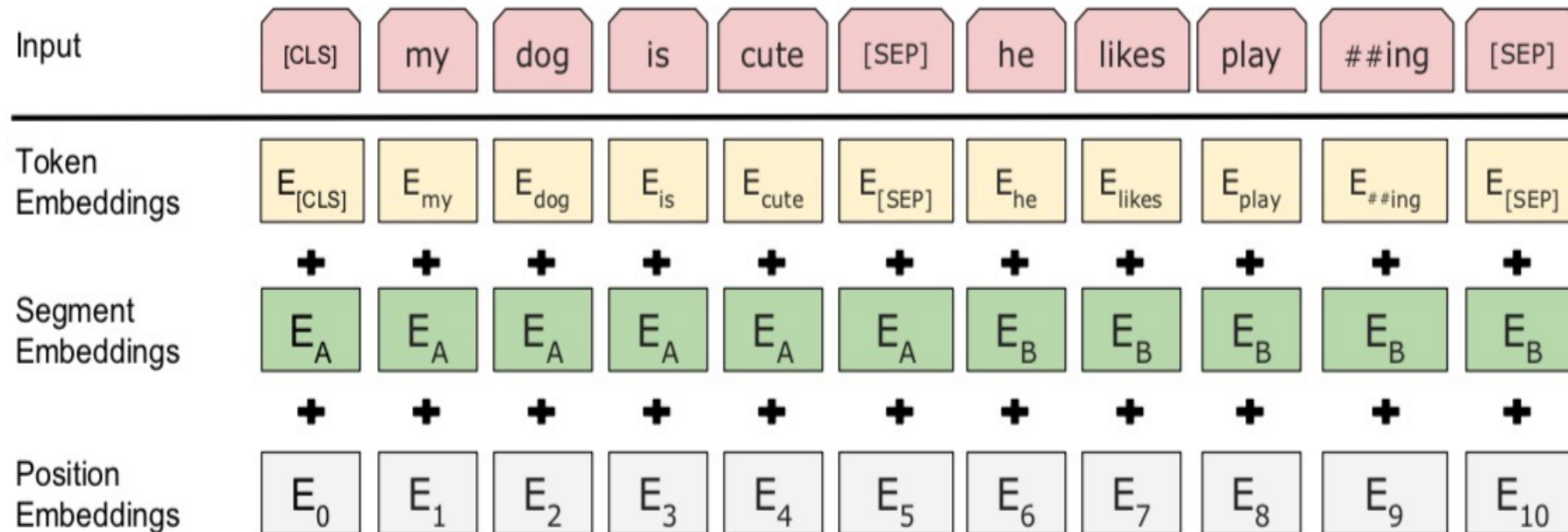
ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks.

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

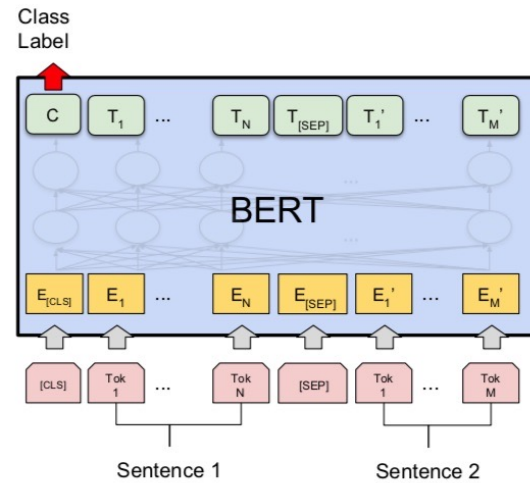
BERT (Bidirectional Encoder Representations from Transformers)

BERT input representation

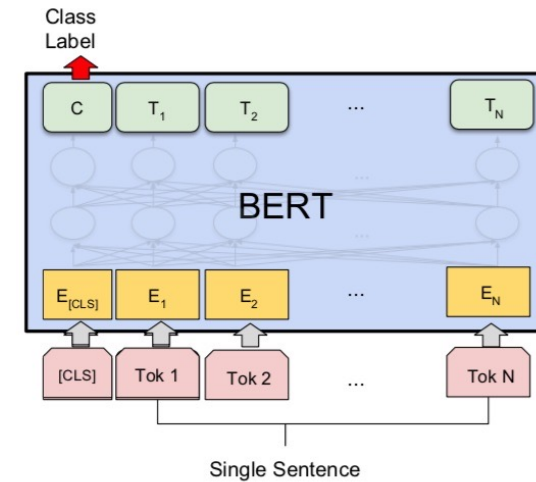


The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

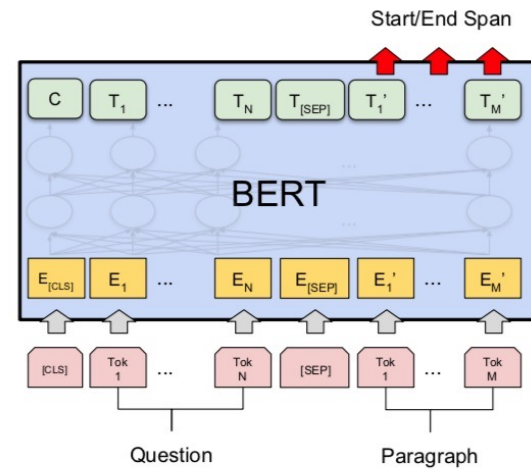
Fine-tuning BERT on NLP Tasks



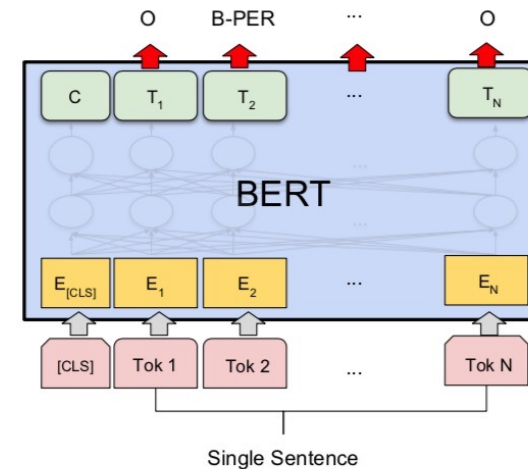
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



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



(c) Question Answering Tasks:
SQuAD v1.1

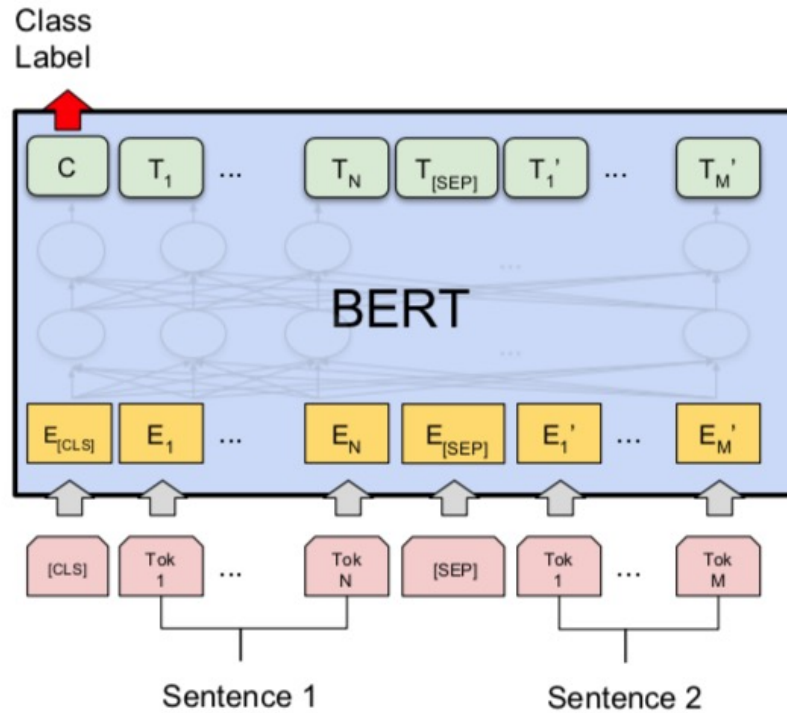


(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

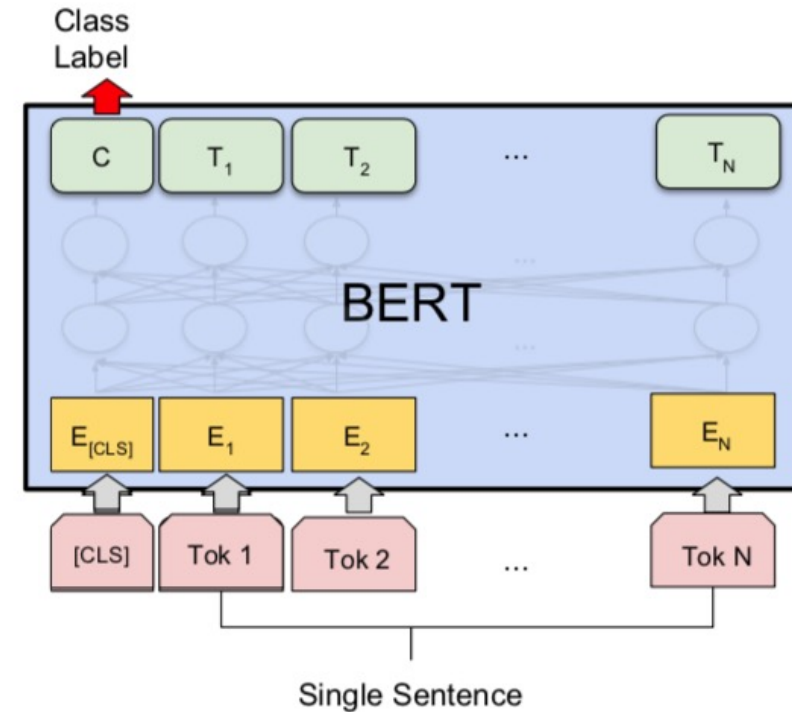
Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

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

BERT Sequence-level tasks

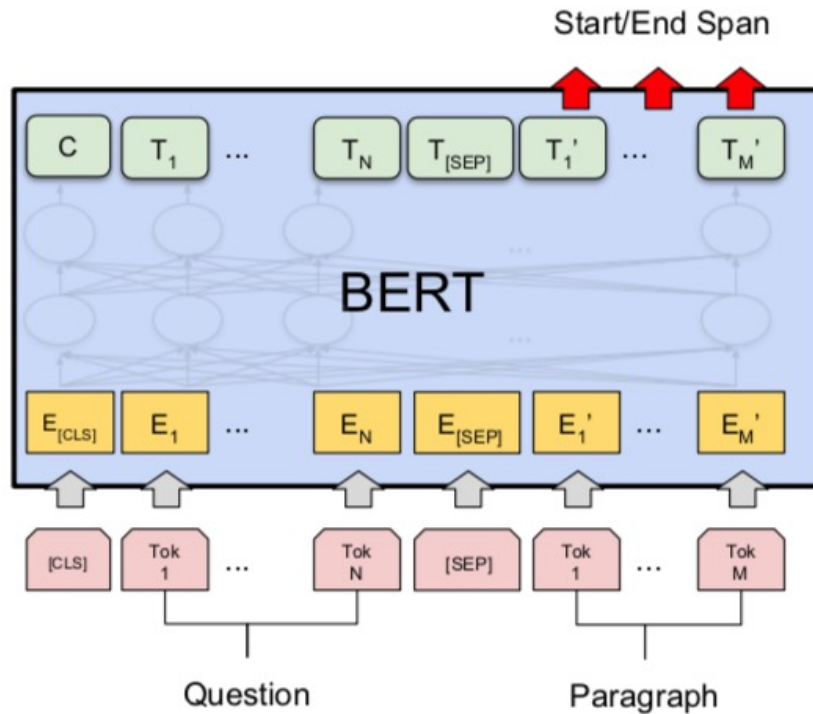


(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

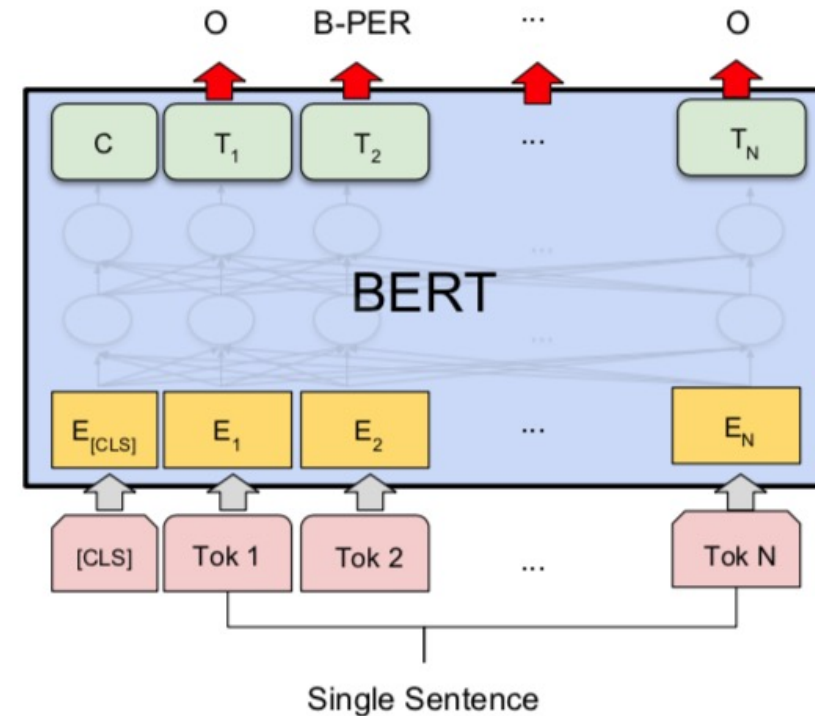


(b) Single Sentence Classification Tasks:
SST-2, CoLA

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 _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

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

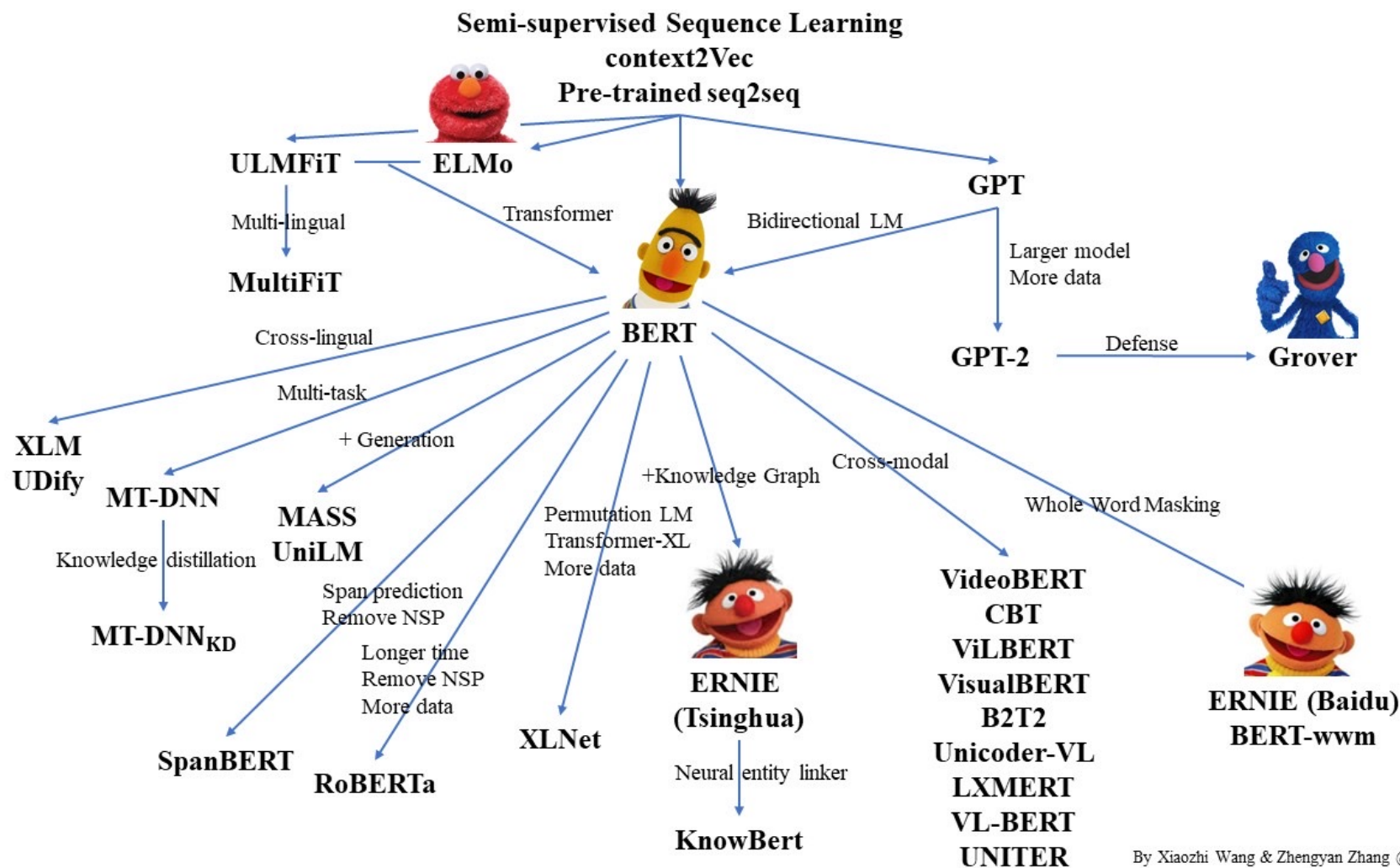
MRPC: Microsoft Research Paraphrase Corpus

RTE: Recognizing Textual Entailment

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

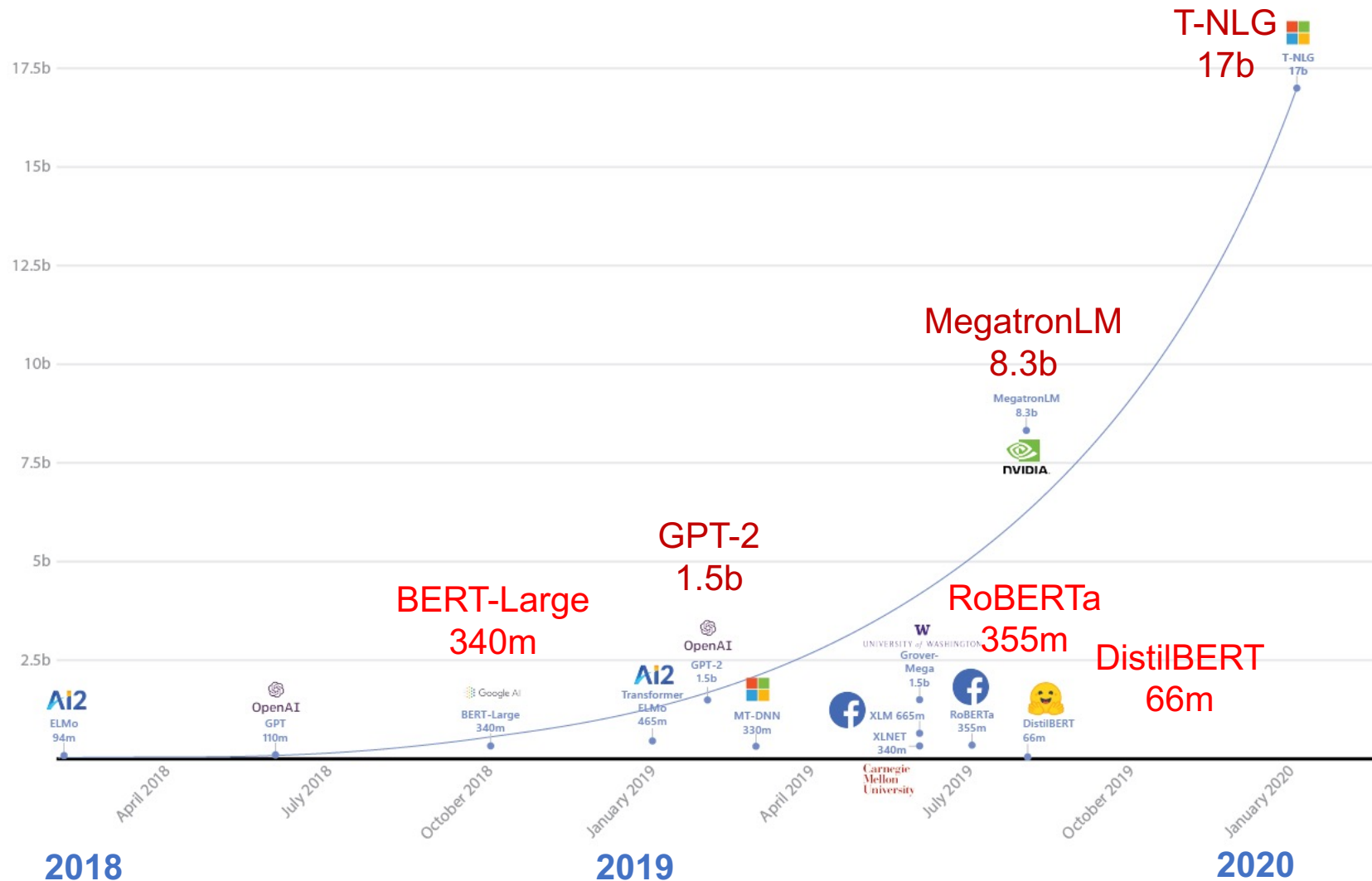
"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805

Pre-trained Language Model (PLM)



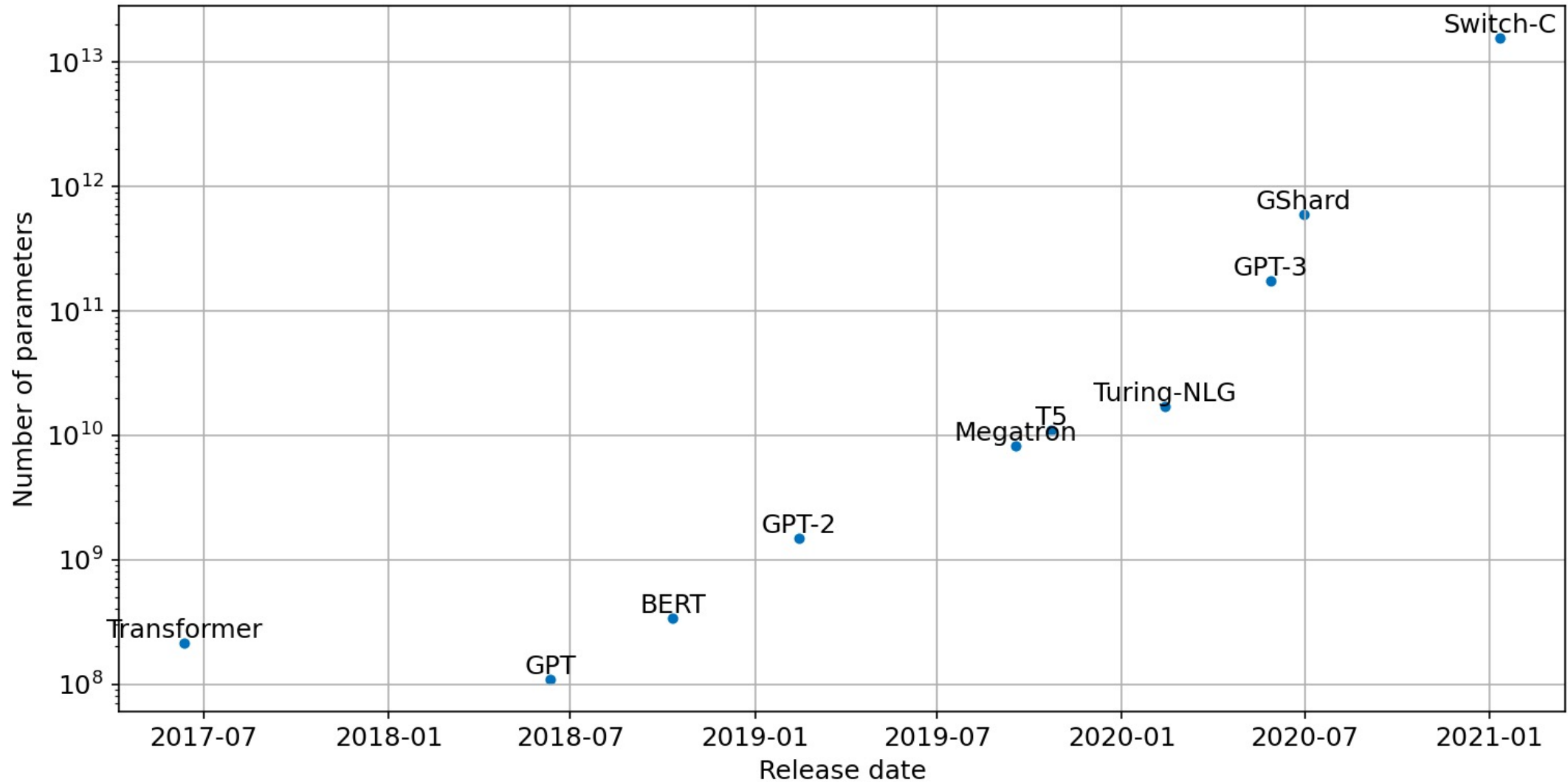
By Xiaozhi Wang & Zhengyan Zhang @THUNLP

Transformers Pre-trained Language Model

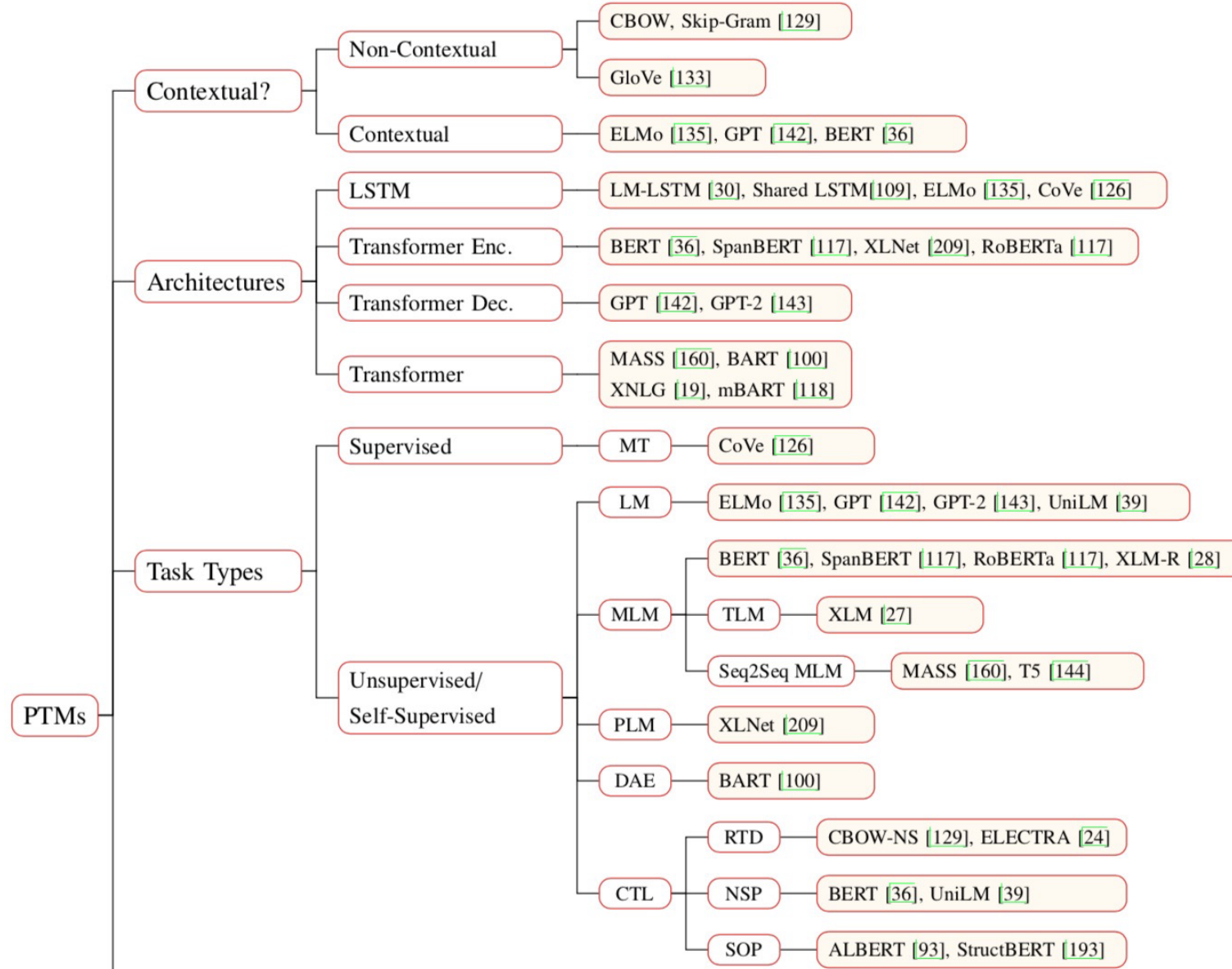


90+ Models

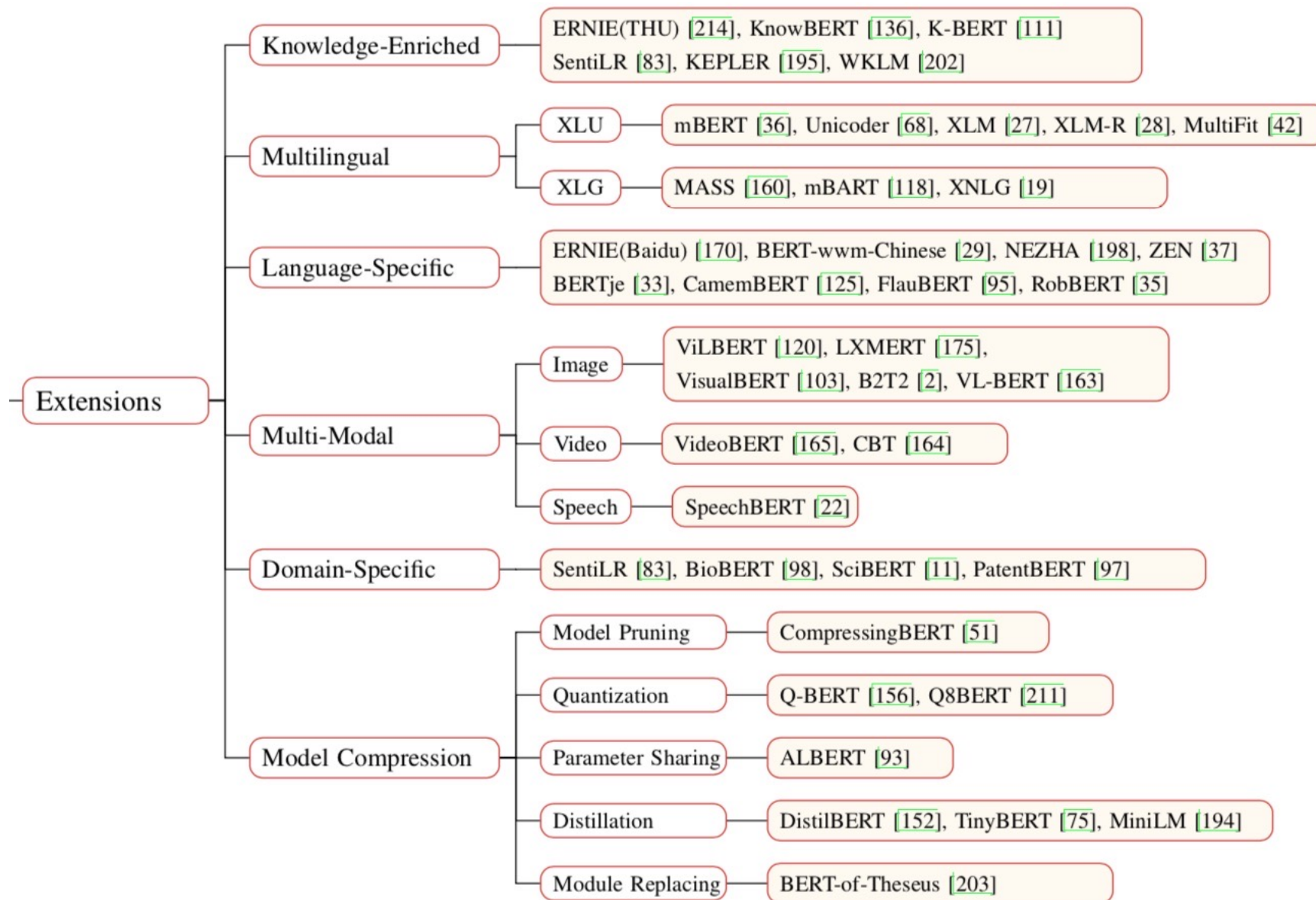
Scaling Transformers



Pre-trained Models (PTM)



Pre-trained Models (PTM)





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.**

Hugging Face



Hugging Face

Search models, datasets

Models

Datasets

Spaces

Docs

Solutions

Pricing



Log In

Sign Up



The AI community building the future.

Build, train and deploy state of the art models powered by
the reference open source in machine learning.



Star

58,696

<https://huggingface.co/>

Hugging Face Transformers

[Models](#)[Datasets](#)[Spaces](#)[Docs](#)[Solutions](#)[Pricing](#)[Log In](#)[Sign Up](#)

Transformers

[V4.16.2](#)[EN](#)

58,697

GET STARTED

[Transformers](#)[Quick tour](#)[Installation](#)[Philosophy](#)[Glossary](#)

USING TRANSFORMERS

[Summary of the tasks](#)[Summary of the models](#)[Preprocessing data](#)[Fine-tuning a pretrained model](#)[Distributed training with 🤗](#)[Accelerate](#)

🤗 Transformers

State-of-the-art Machine Learning for Jax, Pytorch and TensorFlow

🤗 Transformers (formerly known as *pytorch-transformers* and *pytorch-pretrained-bert*) provides thousands of pretrained models to perform tasks on different modalities such as text, vision, and audio.

These models can applied on:

- 📄 Text, for tasks like text classification, information extraction, question answering, summarization, translation, text generation, in over 100 languages.
- 🖼️ Images, for tasks like image classification, object detection, and segmentation.
- 🗣️ Audio, for tasks like speech recognition and audio classification.

Transformer models can also perform tasks on **several modalities combined**, such as table question answering, optical character recognition, information extraction from scanned documents. video classification. and visual question answering.

<https://huggingface.co/docs/transformers/index>

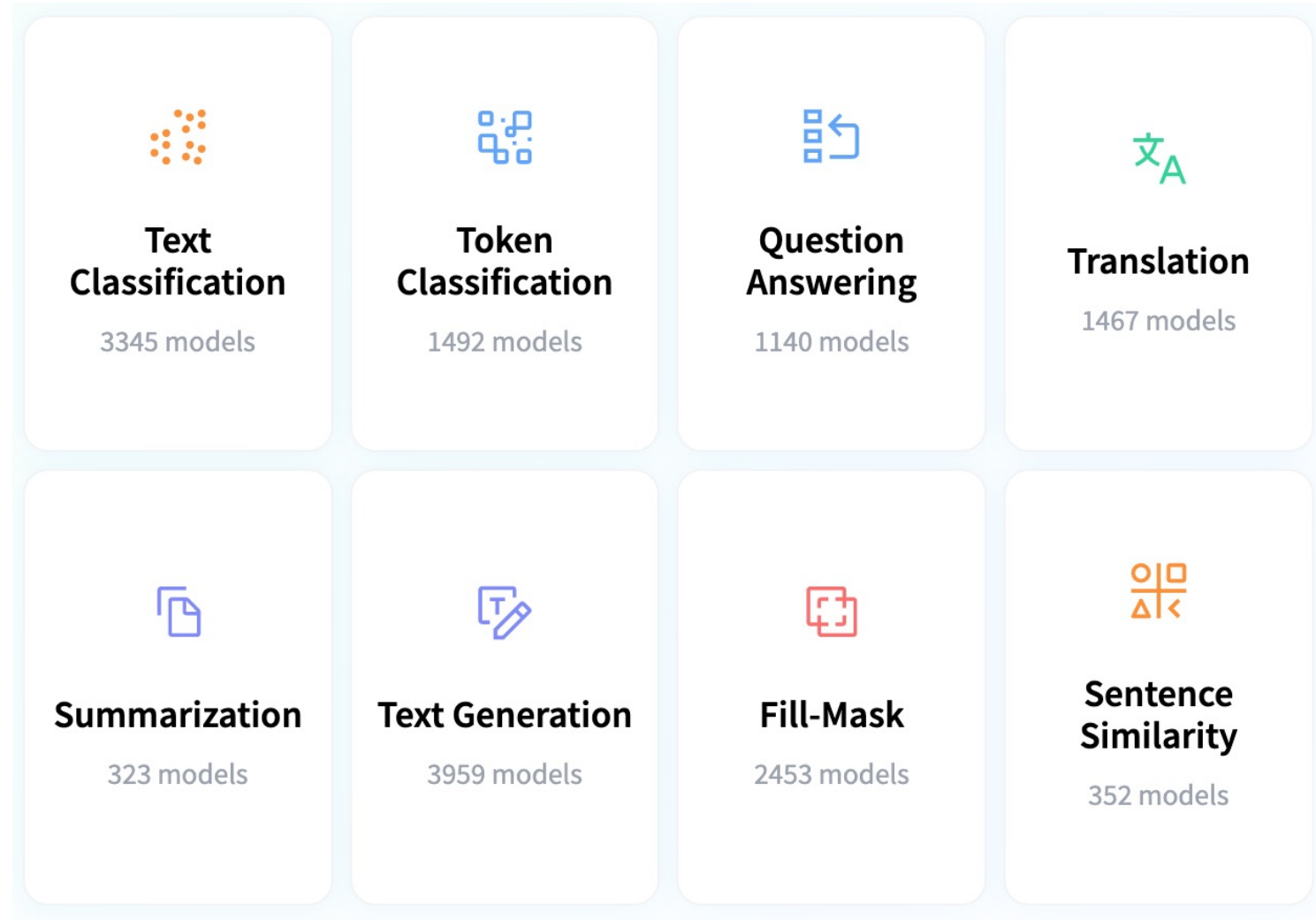
🤗 Transformers

If you are looking for custom support from the Hugging Face team

[Features](#)[Contents](#)[Supported models](#)[Supported frameworks](#)

Hugging Face Tasks

Natural Language Processing



<https://huggingface.co/tasks>

Text Analytics with Python

NLP Libraries and Tools



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/>

Natural Language Processing with Python

– Analyzing Text with the Natural Language Toolkit

← → ↻ ⓘ www.nltk.org/book/

Natural Language Processing with Python

– Analyzing Text with the Natural Language Toolkit

NLTK

Steven Bird, Ewan Klein, and Edward Loper

This version of the NLTK book is updated for Python 3 and NLTK 3. The first edition of the book, published by O'Reilly, is available at http://nltk.org/book_1ed/. (There are currently no plans for a second edition of the book.)

0. [Preface](#)
1. [Language Processing and Python](#)
2. [Accessing Text Corpora and Lexical Resources](#)
3. [Processing Raw Text](#)
4. [Writing Structured Programs](#)
5. [Categorizing and Tagging Words](#) (minor fixes still required)
6. [Learning to Classify Text](#)
7. [Extracting Information from Text](#)
8. [Analyzing Sentence Structure](#)
9. [Building Feature Based Grammars](#)
10. [Analyzing the Meaning of Sentences](#) (minor fixes still required)
11. [Managing Linguistic Data](#) (minor fixes still required)
12. [Afterword: Facing the Language Challenge](#)

[Bibliography](#)

[Term Index](#)

This book is made available under the terms of the [Creative Commons Attribution Noncommercial No-Derivative-Works 3.0 US License](#). Please post any questions about the materials to the [nltk-users](#) mailing list. Please report any errors on the [issue tracker](#).

<http://www.nltk.org/book/>

gensim

Fork me on GitHub



gensim

topic modelling for humans

[Download](#)
latest version from the Python Package Index

[Direct install with:
easy_install -U gensim](#)

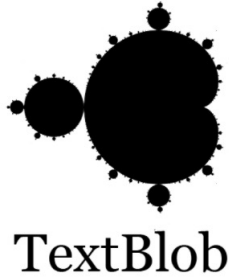
[Home](#) [Tutorials](#) [Install](#) [Support](#) [API](#) [About](#)

```
>>> from gensim import corpora, models, similarities
>>>
>>> # Load corpus iterator from a Matrix Market file on disk.
>>> corpus = corpora.MmCorpus('/path/to/corpus.mm')
>>>
>>> # Initialize Latent Semantic Indexing with 200 dimensions.
>>> lsi = models.LsiModel(corpus, num_topics=200)
>>>
>>> # Convert another corpus to the latent space and index it.
>>> index = similarities.MatrixSimilarity(lsi[another_corpus])
>>>
>>> # Compute similarity of a query vs. indexed documents
>>> sims = index[query]
```

Gensim is a FREE Python library

- ✓ Scalable statistical semantics
- ✓ Analyze plain-text documents for semantic structure
- ✓ Retrieve semantically similar documents

TextBlob



Star 3,777

TextBlob is a Python (2 and 3) library for processing textual data. It provides a consistent API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, and more.

Useful Links

[TextBlob @ PyPI](#)
[TextBlob @ GitHub](#)
[Issue Tracker](#)

Stay Informed

Follow @sloria

Donate

If you find TextBlob useful,

TextBlob: Simplified Text Processing

Release v0.12.0. ([Changelog](#))

TextBlob is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

```
from textblob import TextBlob

text = '''
The titular threat of The Blob has always struck me as the ultimate movie
monster: an insatiably hungry, amoeba-like mass able to penetrate
virtually any safeguard, capable of--as a doomed doctor chillingly
describes it--"assimilating flesh on contact.
Snide comparisons to gelatin be damned, it's a concept with the most
devastating of potential consequences, not unlike the grey goo scenario
proposed by technological theorists fearful of
artificial intelligence run rampant.
'''

blob = TextBlob(text)
blob.tags          # [('The', 'DT'), ('titular', 'JJ'),
                   # ('threat', 'NN'), ('of', 'IN'), ...]

blob.noun_phrases # WordList(['titular threat', 'blob',
                              # 'ultimate movie monster',
                              # 'amoeba-like mass', ...])

for sentence in blob.sentences:
    print(sentence.sentiment.polarity)
# 0.060
```

<https://textblob.readthedocs.io>

Polyglot

🏠 polyglot
latest

Search docs

- Installation
- Language Detection
- Tokenization
- Command Line Interface
- Downloading Models
- Word Embeddings
- Part of Speech Tagging
- Named Entity Extraction
- Morphological Analysis
- Transliteration
- Sentiment
- polyglot

Docs » Welcome to polyglot's documentation!

[Edit on GitHub](#)

Welcome to polyglot's documentation!

polyglot

downloads 17k/month pypi package 16.7.4 build passing docs passing

Polyglot is a natural language pipeline that supports massive multilingual applications.

- Free software: GPLv3 license
- Documentation: <http://polyglot.readthedocs.org>.

Features

- Tokenization (165 Languages)
- Language detection (196 Languages)
- Named Entity Recognition (40 Languages)
- Part of Speech Tagging (16 Languages)
- Sentiment Analysis (136 Languages)
- Word Embeddings (137 Languages)
- Morphological analysis (135 Languages)
- Transliteration (69 Languages)

scikit-learn



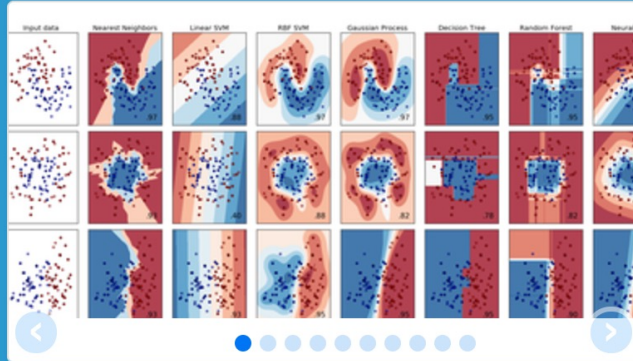
powered by Google

Home Installation Documentation Examples

Google Custom Search

Search

Fork me on GitHub



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ...

— Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ...

— Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ...

— Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction.

<http://scikit-learn.org/>

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 navigation bar with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help" menus. A "Table of contents" sidebar on the left lists various NLP topics. The main workspace contains two code cells. The first cell runs a Spacy pipeline to parse a sentence, and the second cell uses pandas to create a DataFrame from the Spacy output. The output of the first cell shows a visual representation of the sentence with entities highlighted in colored boxes. The output of the second cell is a DataFrame with columns for text, lemma, pos, tag, pos_explain, and stopwords.

Table of contents:

- Text Analytics and Natural Language Processing (NLP)
- Python for Natural Language Processing**
 - spaCy Chinese Model
 - Open Chinese Convert (OpenCC, 開放中文轉換)
 - Jieba 結巴中文分詞
 - Natural Language Toolkit (NLTK)
 - Stanza: A Python NLP Library for Many Human Languages
- Text Processing and Understanding
 - NLTK (Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit)
 - NLP Zero to Hero
 - Natural Language Processing - Tokenization (NLP Zero to Hero, part 1)
 - Natural Language Processing - Sequencing - Turning sentence into data (NLP Zero to Hero, part 2)
 - Natural Language Processing - Training a model to recognize sentiment in text (NLP Zero to Hero, part 3)
 - Keras preprocessing text
 - JSON File

```
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	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 contains a "Table of contents" panel with the following items:

- Text Analytics and Natural Language Processing (NLP)
 - Python for Natural Language Processing
 - spaCy Chinese Model
 - Open Chinese Convert (OpenCC, 開放中文轉換)
 - Jieba 結巴中文分詞
 - Natural Language Toolkit (NLTK)
 - Stanza: A Python NLP Library for Many Human Languages
- Text Processing and Understanding
 - NLTK (Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit)
 - NLP Zero to Hero
 - Natural Language Processing - Tokenization (NLP Zero to Hero, part 1)
 - Natural Language Processing - Sequencing - Turning sentence into data (NLP Zero to Hero, part 2)
 - Natural Language Processing - Training a model to recognize sentiment in text (NLP Zero to Hero, part 3)

Text Analytics and Natural Language Processing (NLP)

Python for Natural Language Processing

spaCy

- spaCy: Industrial-Strength Natural Language Processing in Python
- Source: <https://spacy.io/usage/spacy-101>

```
[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_)
```

```
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
[ ] 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

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.

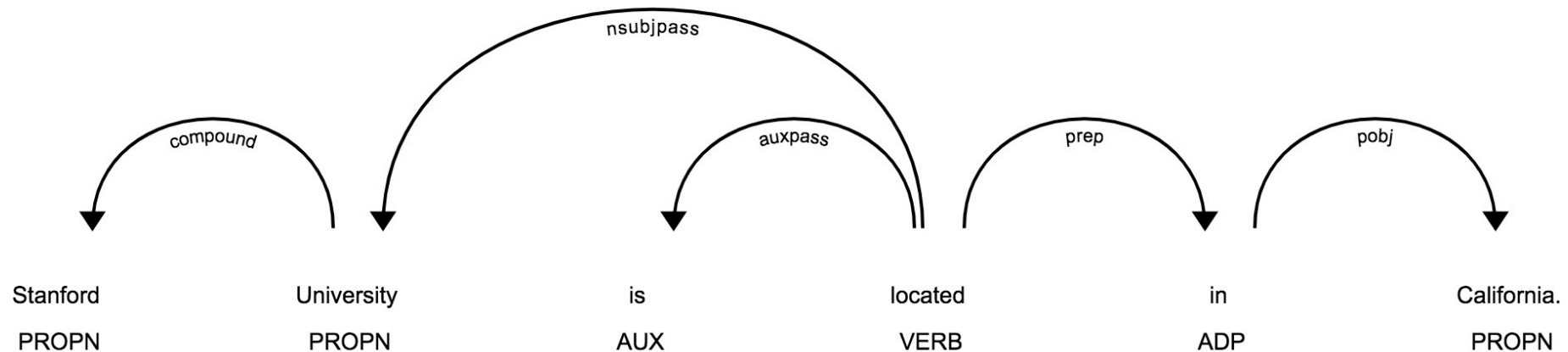
<https://tinyurl.com/aintpupython101>

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 Python code for NLP processing. The first cell shows a text snippet with entities highlighted in colored boxes. The second cell shows a pandas DataFrame with columns for text, lemma, pos, tag, pos_explain, and 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>

NLP with Transformers Github

Why GitHub? Team Enterprise Explore Marketplace Pricing Search / Sign in Sign up

nlp-with-transformers / notebooks Public

Notifications Fork 170 Star 1.1k

Code Issues Pull requests Actions Projects Wiki Security Insights

main 1 branch 0 tags Go to file Code

lewtun Merge pull request #21 from JingchaoZhang/patch-3	ae5b7c1 15 days ago	71 commits
.github/ISSUE_TEMPLATE	Update issue templates	25 days ago
data	Move dataset to data directory	4 months ago
images	Add README	last month
scripts	Update issue templates	25 days ago
.gitignore	Initial commit	4 months ago
01_introduction.ipynb	Remove Colab badges & fastdoc refs	27 days ago
02_classification.ipynb	Merge pull request #8 from nlp-with-transformers/remove-display-df	26 days ago
03_transformer-anatomy.ipynb	[Transformers Anatomy] Remove cells with figure references	22 days ago
04_multilingual-ner.ipynb	Merge pull request #8 from nlp-with-transformers/remove-display-df	26 days ago
05_text-generation.ipynb	Merge pull request #8 from nlp-with-transformers/remove-display-df	26 days ago

About

Jupyter notebooks for the Natural Language Processing with Transformers book

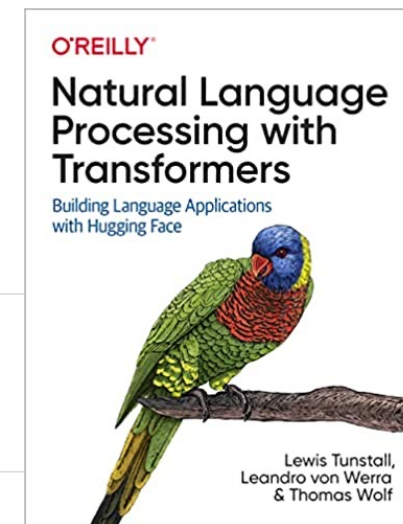
transformersbook.com/

- Readme
- Apache-2.0 License
- 1.1k stars
- 33 watching
- 170 forks

Releases

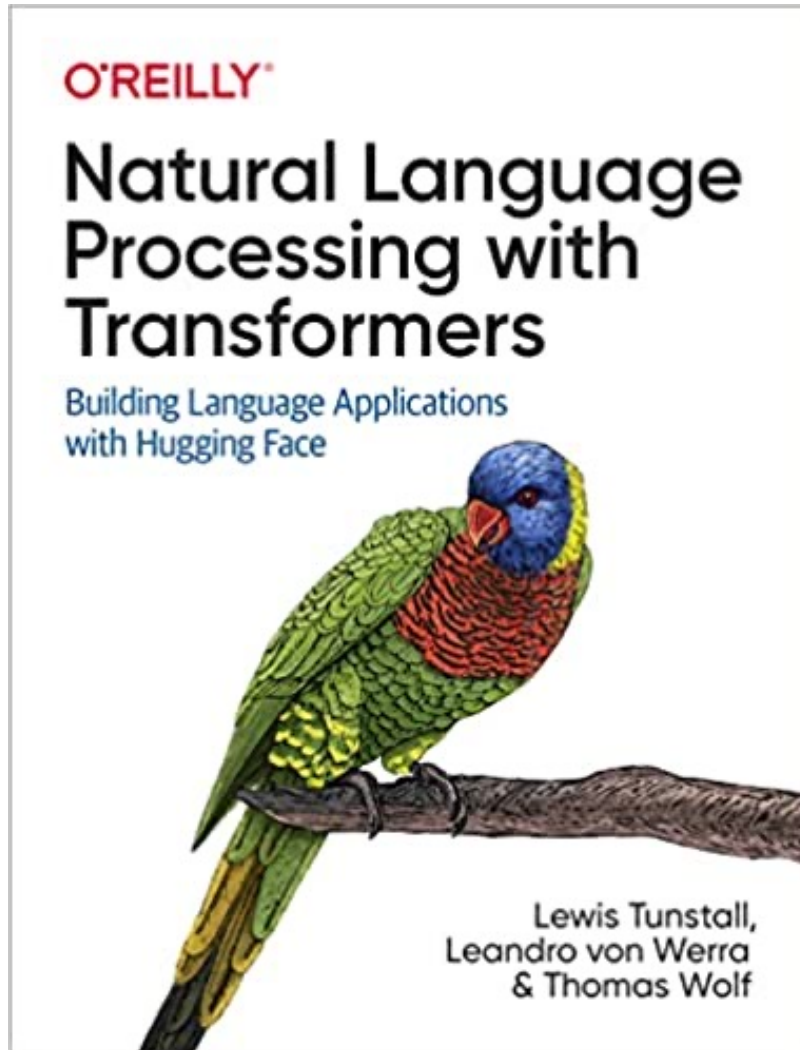
No releases published

Packages





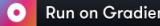


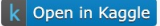
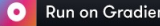

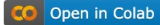

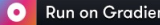
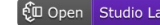
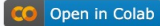
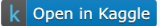
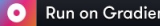
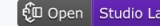
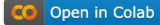
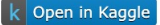
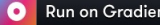
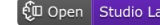


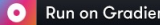



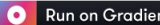
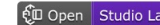


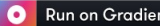



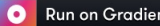



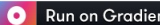



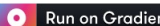

<https://github.com/nlp-with-transformers/notebooks>

NLP with Transformers Github Notebooks



Running on a cloud platform

To run these notebooks on a cloud platform, just click on one of the badges in the table below:

Chapter	Colab	Kaggle	Gradient	Studio Lab
Introduction	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Text Classification	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Transformer Anatomy	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Multilingual Named Entity Recognition	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Text Generation	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Summarization	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Question Answering	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Making Transformers Efficient in Production	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Dealing with Few to No Labels	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Training Transformers from Scratch	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab
Future Directions	 Open in Colab	 Open in Kaggle	 Run on Gradient	 Open Studio Lab

Nowadays, the GPUs on Colab tend to be K80s (which have limited memory), so we recommend using [Kaggle](#), [Gradient](#), or [SageMaker Studio Lab](#). These platforms tend to provide more performant GPUs like P100s, all for free!

<https://github.com/nlp-with-transformers/notebooks>

NLP Benchmark Datasets

Task	Dataset	Link
Machine Translation	WMT 2014 EN-DE WMT 2014 EN-FR	http://www-lium.univ-lemans.fr/~schwenk/cs1m_joint_paper/
Text Summarization	CNN/DM Newsroom DUC Gigaword	https://cs.nyu.edu/~kcho/DMQA/ https://summari.es/ https://www-nlpir.nist.gov/projects/duc/data.html https://catalog.ldc.upenn.edu/LDC2012T21
Reading Comprehension Question Answering Question Generation	ARC CliCR CNN/DM NewsQA RACE SQuAD Story Cloze Test NarrativeQA Quasar SearchQA	http://data.allenai.org/arc/ http://aclweb.org/anthology/N18-1140 https://cs.nyu.edu/~kcho/DMQA/ https://datasets.maluuba.com/NewsQA http://www.qizhexie.com/data/RACE_leaderboard https://rajpurkar.github.io/SQuAD-explorer/ http://aclweb.org/anthology/W17-0906.pdf https://github.com/deepmind/narrativeqa https://github.com/bdhingra/quasar https://github.com/nyu-dl/SearchQA
Semantic Parsing	AMR parsing ATIS (SQL Parsing) WikiSQL (SQL Parsing)	https://amr.isi.edu/index.html https://github.com/jkkummerfeld/text2sql-data/tree/master/data https://github.com/salesforce/WikiSQL
Sentiment Analysis	IMDB Reviews SST Yelp Reviews Subjectivity Dataset	http://ai.stanford.edu/~amaas/data/sentiment/ https://nlp.stanford.edu/sentiment/index.html https://www.yelp.com/dataset/challenge http://www.cs.cornell.edu/people/pabo/movie-review-data/
Text Classification	AG News DBpedia TREC 20 NewsGroup	http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html https://wiki.dbpedia.org/Datasets https://trec.nist.gov/data.html http://qwone.com/~jason/20Newsgroups/
Natural Language Inference	SNLI Corpus MultiNLI SciTail	https://nlp.stanford.edu/projects/snli/ https://www.nyu.edu/projects/bowman/multinli/ http://data.allenai.org/scitail/
Semantic Role Labeling	Proposition Bank OneNotes	http://propbank.github.io/ https://catalog.ldc.upenn.edu/LDC2013T19

Summary

- **Text Analytics and Text Mining**
- **Natural Language Processing (NLP)**
- **Text Analytics with Python**

References

- 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.
- Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson.
- 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.
- Rajesh Arumugam (2018), Hands-On Natural Language Processing with Python: A practical guide to applying deep learning architectures to your NLP applications, Packt.
- Jake VanderPlas (2016), Python Data Science Handbook: Essential Tools for Working with Data, O'Reilly Media.
- 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.
- Christopher D. Manning and Hinrich Schütze (1999), Foundations of Statistical Natural Language Processing, The MIT Press.
- Bruce Croft, Donald Metzler, and Trevor Strohman (2008), Search Engines: Information Retrieval in Practice, Addison Wesley, <http://www.search-engines-book.com/>
- Steven Bird, Ewan Klein and Edward Loper (2009), Natural Language Processing with Python, O'Reilly Media, <http://www.nltk.org/book/> , http://www.nltk.org/book_1ed/
- Bing Liu (2009), Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data, Springer.
- The Super Duper NLP Repo, <https://notebooks.quantumstat.com/>
- Min-Yuh Day (2022), Python 101, <https://tinyurl.com/aintpupython101>