# Artificial Intelligence for Text Analytics
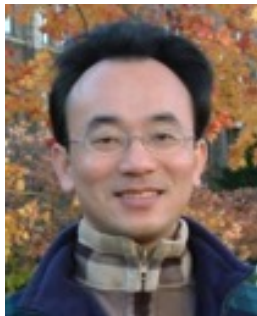
# Text Summarization and Topic Models

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

**Institute of Information Management**, **National Taipei University**

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

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

2022-04-26

# 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

# Text Summarization
# and
# Topic Models

# Outline

- **Text Summarization**
  - **Extractive Text Summarization**
  - **Abstractive Text Summarization**
    - **PEGASUS: Abstractive Summarization**
- **Topic Models**
  - **Topic Modeling**
  - **Latent Dirichlet Allocation (LDA)**
  - **BERTopic**

# Text Summarization

# Text Summarization

https://huggingface.co/tasks/summarization

# Text Summarization

Summarization                                    Examples ⌄

The tower is 324 metres (1,063 ft) tall, about the same height as an 81-storey building, and the tallest structure in Paris. Its base is square, measuring 125 metres (410 ft) on each side. During its construction, the Eiffel Tower surpassed the Washington Monument to become the tallest man-made structure in the world, a title it held for 41 years until the Chrysler Building in New York City was finished in 1930. It was the first structure to reach a height of 300 metres. Due to the addition of a broadcasting aerial at the top of the tower in 1957, it is now taller than the Chrysler Building by 5.2 metres (17 ft). Excluding transmitters, the Eiffel Tower is the second tallest free-standing structure in France after the Millau Viaduct.

Compute

Computation time on cpu: cached

The tower is 324 metres (1,063 ft) tall, about the same height as an 81-storey building . It was the first structure to reach a height of 300 metres . It is now taller than the Chrysler Building in New York City by 5.2 metres (17 ft) Excluding transmitters, the Eiffel Tower is the second tallest free-standing structure in France .
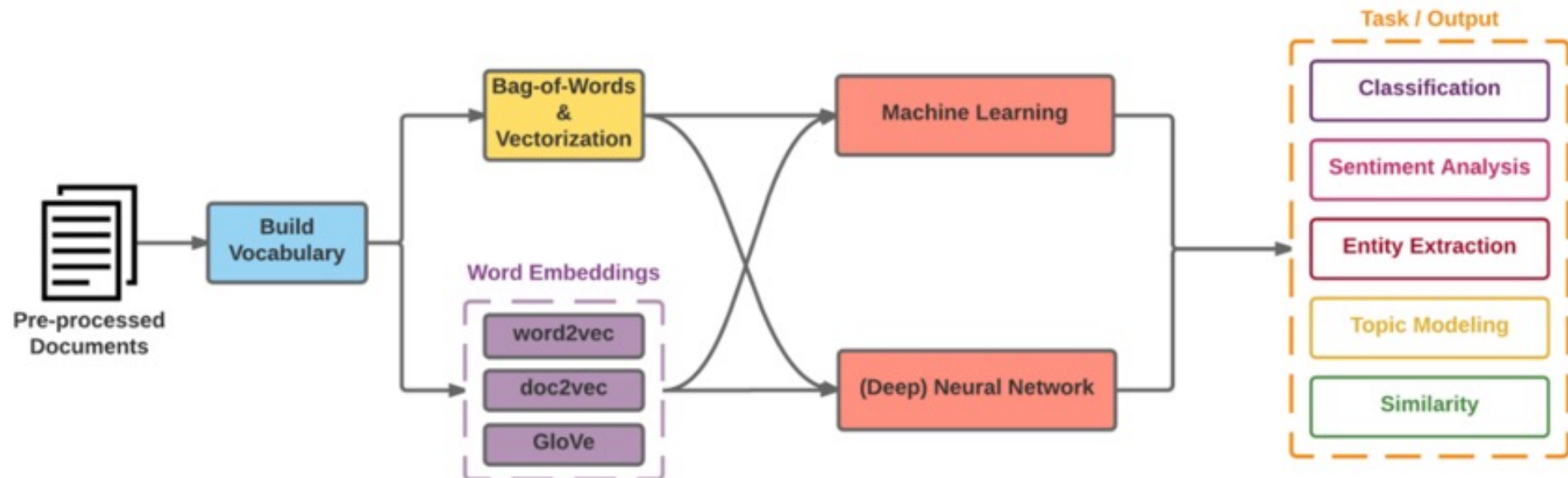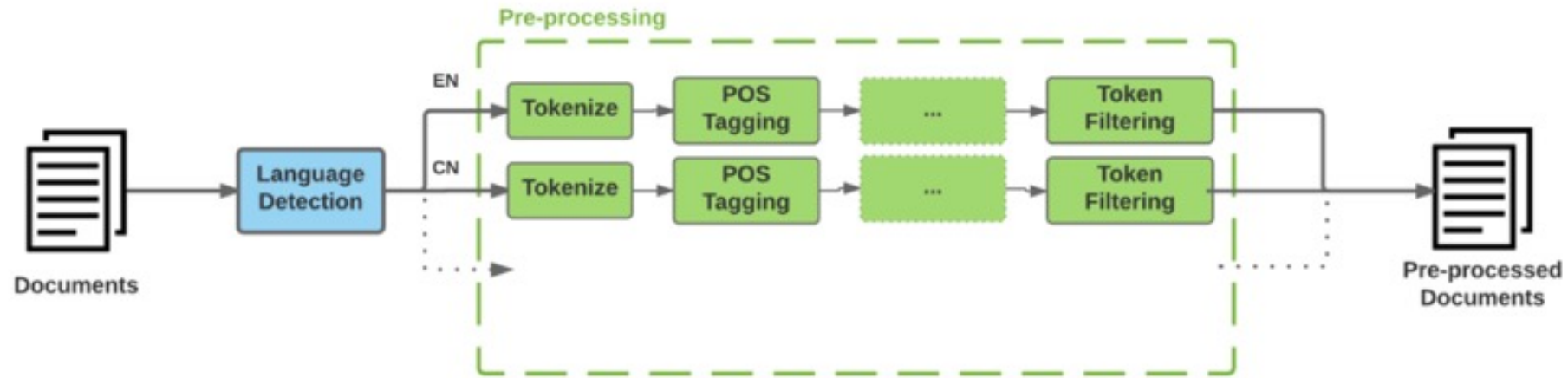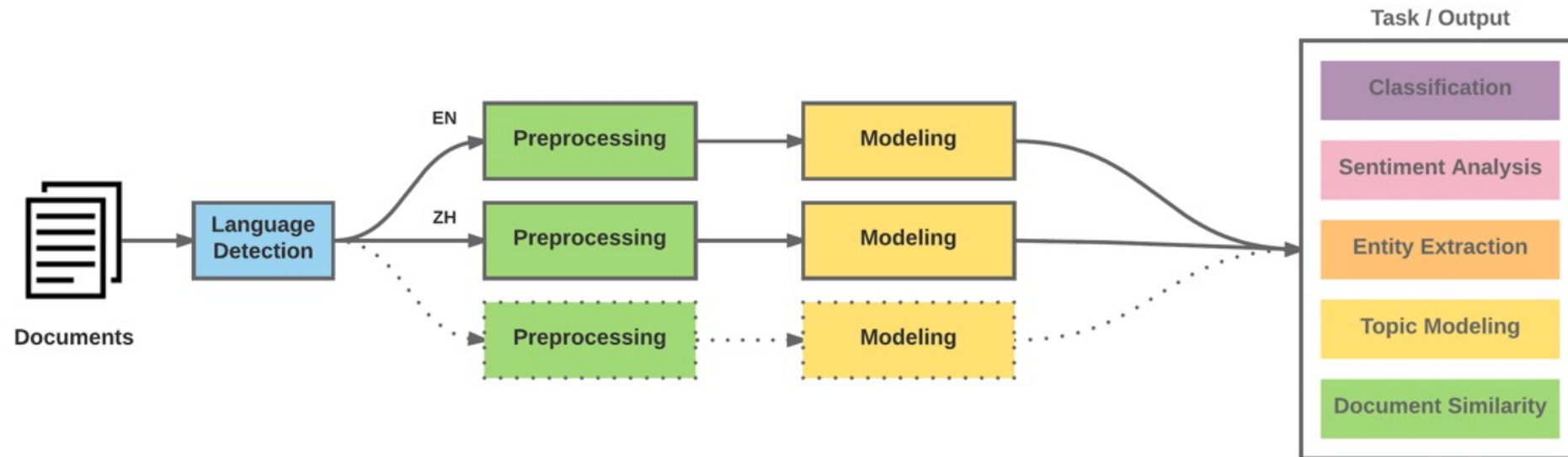
https://huggingface.co/tasks/summarization

# NLP

# Modern NLP Pipeline

Source: https://github.com/fortiema/talks/blob/master/opendata2016sh/pragmatic-nlp-opendata2016sh.pdf

# Modern NLP Pipeline

Source: http://mattfortier.me/2017/01/31/nlp-intro-pt-1-overview/

# Deep Learning NLP

# T5
# Text-to-Text Transfer Transformer



Source: JColin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu (2019). "Exploring the limits of transfer learning with a unified text-to-text transformer." arXiv preprint arXiv:1910.10683 (2019).

# Text Summarization and Information Extraction

- **Key-phrase extraction**
  - **extracting key influential phrases from the documents.**

- **Topic modeling**
  - **Extract various diverse concepts or topics present in the documents, retaining the major themes.**

- **Document summarization**
  - **Summarize entire text documents to provide a gist that retains the important parts of the whole corpus.**

# Natural Language Processing (NLP) and Text Mining

| Raw text |
|---|

| Sentence Segmentation |
|---|

| Tokenization |
|---|

| Part-of-Speech (POS) |
|---|

| Stop word removal |
|---|

| **Stemming** / **Lemmatization** |
|---|

| Dependency Parser |
|---|

| String Metrics & Matching |
|---|

word's stem
am → am
having → hav

word's lemma
am → be
having → have

# Text Summarization

Source: Vishal Gupta and Gurpreet S. Lehal (2009), "A survey of text mining techniques and applications,"
Journal of emerging technologies in web intelligence, vol. 1, no. 1, pp. 60-76.

# Topic Modeling

Source: Blei, David M. "Probabilistic topic models." *Communications of the ACM* 55, no. 4 (2012): 77-84.

# Automatic Text Summarization



(a) Single-document or (b) Multi-document, automatic text summarizer
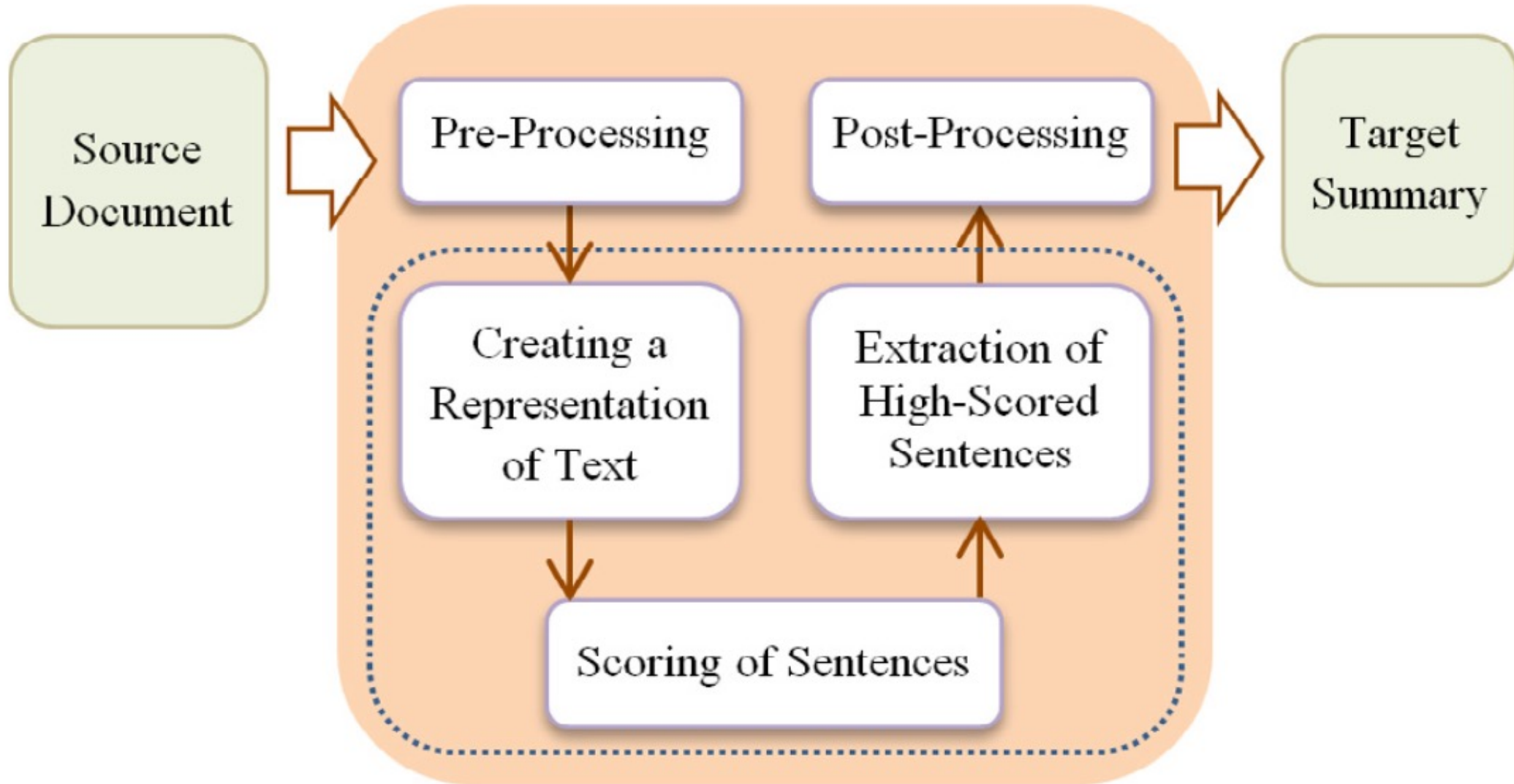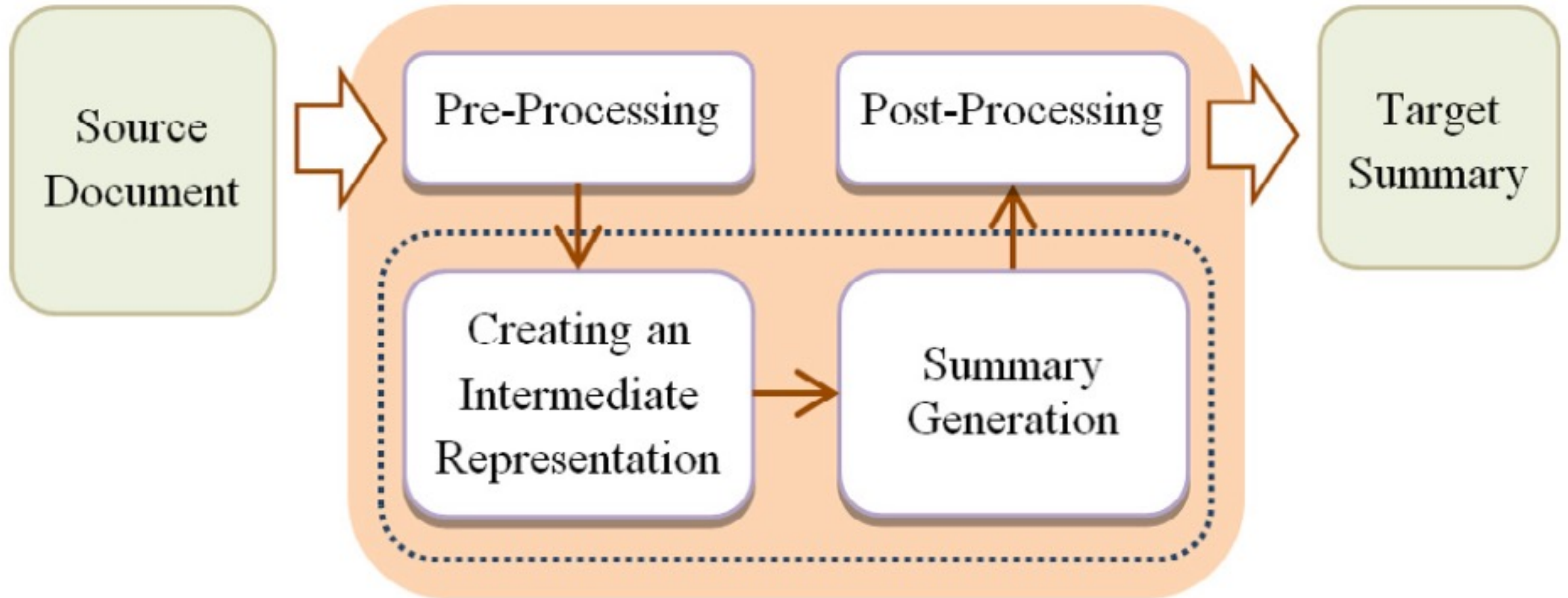
# Classification of Automatic Text Summarization Systems



Source: Wafaa S. El-Kassas, Cherif R. Salama, Ahmed A. Rafea, and Hoda K. Mohamed (2021). "Automatic text summarization: A comprehensive survey." Expert Systems with Applications 165 (2021): 113679.

# Automatic Text Summarization Approaches



Source: Wafaa S. El-Kassas, Cherif R. Salama, Ahmed A. Rafea, and Hoda K. Mohamed (2021). "Automatic text summarization: A comprehensive survey." Expert Systems with Applications 165 (2021): 113679.

# Extractive Text Summarization System

# Abstractive Text Summarization System

# Hybrid Text Summarization System

# Single-sentence and Multi-sentence Text Summarization Operations

# Automatic Text Summarization Building Blocks



Automatic Text Summarization Building Blocks

- Text Representation Models
  - Graph Model
    - Lexical Graph
    - Semantic Graph
  - Vector Model
    - Bag-of-Words
    - Vector Space Model
    - Word Vector
  - N-Gram Model
    - Bi-grams
    - Tri-grams
    - Quad-grams
  - Topic Model
    - LDA
    - PLSA
  - Meaning Representation
    - Lambda Calculus
    - AMR

- Linguistic Analysis and Processing Techniques
  - Pre-Processing Techniques
    - Noise Removal
    - Sentence Segmentation
    - Removal of Punctuation Marks
    - Word Tokenization
    - Named Entity Recognition
    - Removal of Stop-Words
    - Stemming
    - POS
    - Frequency Computation
  - Parsing Techniques
    - Syntactic Parsing
    - Text Chunking
    - Semantic Parsing
    - Shallow semantics
  - Semantic-Based Techniques
    - WSD
    - Anaphora Resolution
    - LSA
    - Textual Entailment
    - Lexical Chain
  - Discourse Analysis
    - Rhetorical Structure Theory
  - Sentence Similarity
    - Syntactic Similarity
    - Semantic-based
    - Hybrid
  - Natural Language Generation

- Soft Computing Techniques
  - Machine learning
    - Supervised Learning
      - Support Vector Machine
      - Naïve Bayes Classification
      - Mathematical Regression
      - Decision Trees
      - Neural Networks
    - Unsupervised Learnig
      - Clustering
      - Hidden Markov Model
    - Semi-Supervised Learning
  - Optimization Algorithms
    - Genetic Algorithm
    - Particle Swarm Optimization
  - Fuzzy Logic
    - Fuzzy Logic System

# PEGASUS:
## Pre-training with Extracted Gap-sentences for Abstractive Summarization

# Topic Modeling

# Topic Model in Bioinformatics

# Topic Modeling

# Topic Modeling (Unsupervised Learning) vs. Text Classification (Supervised Learning)

# Topic Modeling
# Term Document Matrix to
# Topic Distribution



Source: Avinash Navlani (2018), Latent Semantic Analysis using Python,
https://www.datacamp.com/community/tutorials/discovering-hidden-topics-python

# Topic Modeling
# Latent Dirichlet Allocation
# (LDA)



$$\alpha \quad \theta_d \quad z_{d,n} \quad w_{d,n} \quad N \quad D \quad \beta_k \quad K \quad \eta$$

***D* documents**
***N* words**
***K* topics**

Source: Hoffman, Matthew D., David M. Blei, Chong Wang, and John Paisley.
"Stochastic variational inference." The Journal of Machine Learning Research 14, no. 1 (2013): 1303-1347.

# Latent Dirichlet Allocation (Blei et al., 2003)

# Latent Dirichlet Allocation

**David M. Blei**    BLEI@CS.BERKELEY.EDU
*Computer Science Division*
*University of California*
*Berkeley, CA 94720, USA*

**Andrew Y. Ng**    ANG@CS.STANFORD.EDU
*Computer Science Department*
*Stanford University*
*Stanford, CA 94305, USA*

**Michael I. Jordan**    JORDAN@CS.BERKELEY.EDU
*Computer Science Division and Department of Statistics*
*University of California*
*Berkeley, CA 94720, USA*

**Editor:** John Lafferty

## Abstract

We describe *latent Dirichlet allocation* (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying set of topics. Each topic is, in turn, modeled as an infinite mixture over an underlying set of topic probabilities. In the context of text modeling, the topic probabilities provide an explicit representation of a document. We present efficient approximate inference techniques based on variational methods and an EM algorithm for empirical Bayes parameter estimation. We report results in document modeling, text classification, and collaborative filtering, comparing to a mixture of unigrams model and the probabilistic LSI model.

Source: Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." Journal of machine Learning research 3, no. Jan (2003): 993-1022.

# Topic Modeling Using Latent Dirichlet allocation (LDA)



Information Retrieval

Text Summarization

Topic Modeling

LDA

Source: Uttam Chauhan, and Apurva Shah (2021). "Topic modeling using latent Dirichlet allocation: A survey." ACM Computing Surveys (CSUR) 54, no. 7 (2021): 1-35.

# Topic Modeling Technique

# The Generative Process of Latent Dirichlet Allocation (LDA)



Source: Uttam Chauhan, and Apurva Shah (2021). "Topic modeling using latent Dirichlet allocation: A survey." ACM Computing Surveys (CSUR) 54, no. 7 (2021): 1-35.

37

# Topic Visualization as Word Clouds

# LDAvis: Gensim Topic Model Visualization



Source: Uttam Chauhan, and Apurva Shah (2021). "Topic modeling using latent Dirichlet allocation: A survey." ACM Computing Surveys (CSUR) 54, no. 7 (2021): 1-35.

# BERTopic

## Neural topic modeling with a class-based TF-IDF procedure



Maarten Grootendorst (2022). "BERTopic: Neural topic modeling with a class-based TF-IDF procedure."
arXiv preprint arXiv:2203.05794 (2022).

https://github.com/MaartenGr/BERTopic

# gensim

# spaCy

# Python in Google Colab (Python101)

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

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

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

https://tinyurl.com/aintpupython101

# NLP Benchmark Datasets

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

# gensim

# spaCy

# Hugging Face Tasks
# Natural Language Processing

**Text Classification**

3345 models

**Token Classification**

1492 models

**Question Answering**

1140 models

**Translation**

1467 models

**Summarization**

323 models

**Text Generation**

3959 models

**Fill-Mask**

2453 models

**Sentence Similarity**

352 models

https://huggingface.co/tasks

# NLP with Transformers Github



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

# Python in Google Colab (Python101)

# Python in Google Colab (Python101)

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



https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

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

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

python101.ipynb

File   Edit   View   Insert   Runtime   Tools   Help   Saving...

+ Code   + Text

## Text Summarization

- Source: Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022), Natural Language Processing with Transformers: Building Language Applications with Hugging Face, O'Reilly Media.
- Github: https://github.com/nlp-with-transformers/notebooks

```python
#Source: https://huggingface.co/tasks/summarization
!pip install transformers
from transformers import pipeline
classifier = pipeline("summarization")
text = "Paris is the capital and most populous city of France, with an estimated population of 2,175,601 residents as of 2018, in an area of more than
classifier(text, max_length=30)
```

```
No model was supplied, defaulted to sshleifer/distilbart-cnn-12-6 (https://huggingface.co/sshleifer/distilbart-cnn-12-6)
Your min_length=56 must be inferior than your max_length=30.
[{'summary_text': ' Paris is the capital and most populous city of France, with an estimated population of 2,175,601 residents . The City of Paris'}]
```

```python
#!pip install transformers
text = """Dear Amazon, last week I ordered an Optimus Prime action figure \
from your online store in Germany. Unfortunately, when I opened the package, \
I discovered to my horror that I had been sent an action figure of Megatron \
instead! As a lifelong enemy of the Decepticons, I hope you can understand my \
dilemma. To resolve the issue, I demand an exchange of Megatron for the \
Optimus Prime figure I ordered. Enclosed are copies of my records concerning \
this purchase. I expect to hear from you soon. Sincerely, Bumblebee."""
from transformers import pipeline
summarizer = pipeline("summarization")
outputs = summarizer(text, max_length=45, clean_up_tokenization_spaces=True)
print(outputs[0]['summary_text'])
```

https://tinyurl.com/aintpupython101

# Text Summarization

```
text = """Dear Amazon, last week I ordered an Optimus Prime action figure \
from your online store in Germany. Unfortunately, when I opened the package, \
I discovered to my horror that I had been sent an action figure of Megatron \
instead! As a lifelong enemy of the Decepticons, I hope you can understand my \
dilemma. To resolve the issue, I demand an exchange of Megatron for the \
Optimus Prime figure I ordered. Enclosed are copies of my records concerning \
this purchase. I expect to hear from you soon. Sincerely, Bumblebee."""
```

```python
from transformers import pipeline
summarizer = pipeline("summarization")
outputs = summarizer(text, max_length=45, clean_up_tokenization_spaces=True)
print(outputs[0]['summary_text'])
```

Bumblebee ordered an Optimus Prime action figure from your online store in Germany. Unfortunately, when I opened the package, I discovered to my horror that I had been sent an action figure of Megatron instead.

# Summary

- **Text Summarization**
  - **Extractive Text Summarization**
  - **Abstractive Text Summarization**
    - **PEGASUS: Abstractive Summarization**
- **Topic Models**
  - **Topic Modeling**
  - **Latent Dirichlet Allocation (LDA)**
  - **BERTopic**

# 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.
- Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress.
- Wafaa S. El-Kassas, Cherif R. Salama, Ahmed A. Rafea, and Hoda K. Mohamed (2021). "Automatic text summarization: A comprehensive survey." Expert Systems with Applications 165 (2021): 113679.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu (2019). "Exploring the limits of transfer learning with a unified text-to-text transformer." arXiv preprint arXiv:1910.10683 (2019).
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu (2020). "Pegasus: Pre-training with extracted gap-sentences for abstractive summarization." In International Conference on Machine Learning, pp. 11328-11339. PMLR, 2020.
- Uttam Chauhan, and Apurva Shah (2021). "Topic modeling using latent Dirichlet allocation: A survey." ACM Computing Surveys (CSUR) 54, no. 7 (2021): 1-35.
- Maarten Grootendorst (2022). "BERTopic: Neural topic modeling with a class-based TF-IDF procedure." arXiv preprint arXiv:2203.05794 (2022).
- Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python: Enabling Language-Aware Data Products with Machine Learning, O'Reilly.
- Selva Prabhakaran (2020), Topic modeling visualization – How to present the results of LDA models?, https://www.machinelearningplus.com/nlp/topic-modeling-visualization-how-to-present-results-lda-models/
- 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.
- The Super Duper NLP Repo, https://notebooks.quantumstat.com/
- Jay Alammar (2018), The Illustrated Transformer, http://jalammar.github.io/illustrated-transformer/
- Jay Alammar (2019), A Visual Guide to Using BERT for the First Time, http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/
- NLP with Transformer, https://github.com/nlp-with-transformers/notebooks
- Min-Yuh Day (2022), Python 101, https://tinyurl.com/aintpupython101