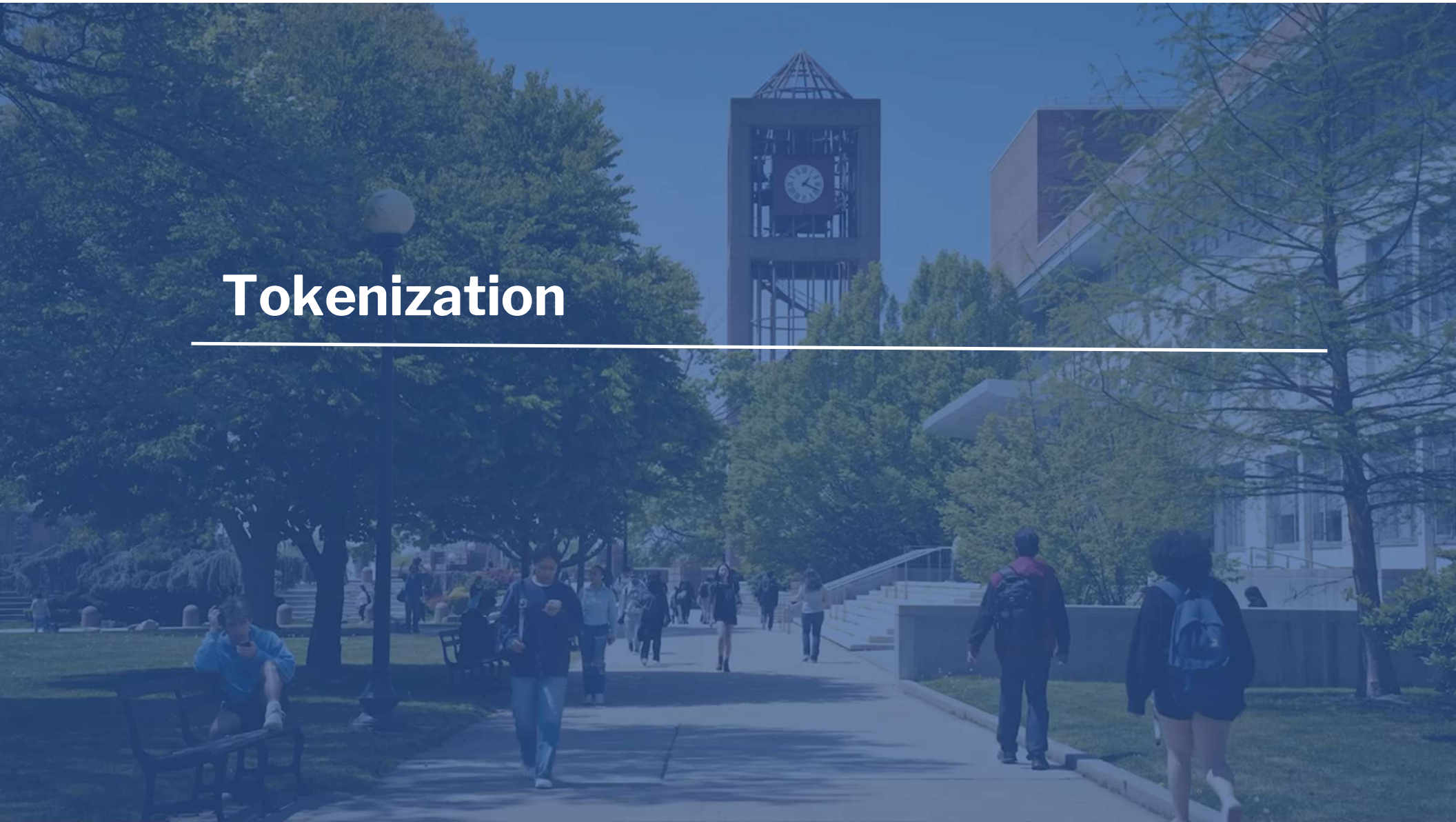


An aerial photograph of the New York City skyline, featuring numerous skyscrapers and the Hudson River in the foreground. A large blue rectangular box is centered over the image, containing white text. The sky is blue with scattered white clouds.

Introduction to Generative AI (GAI 602)

WEEK 2

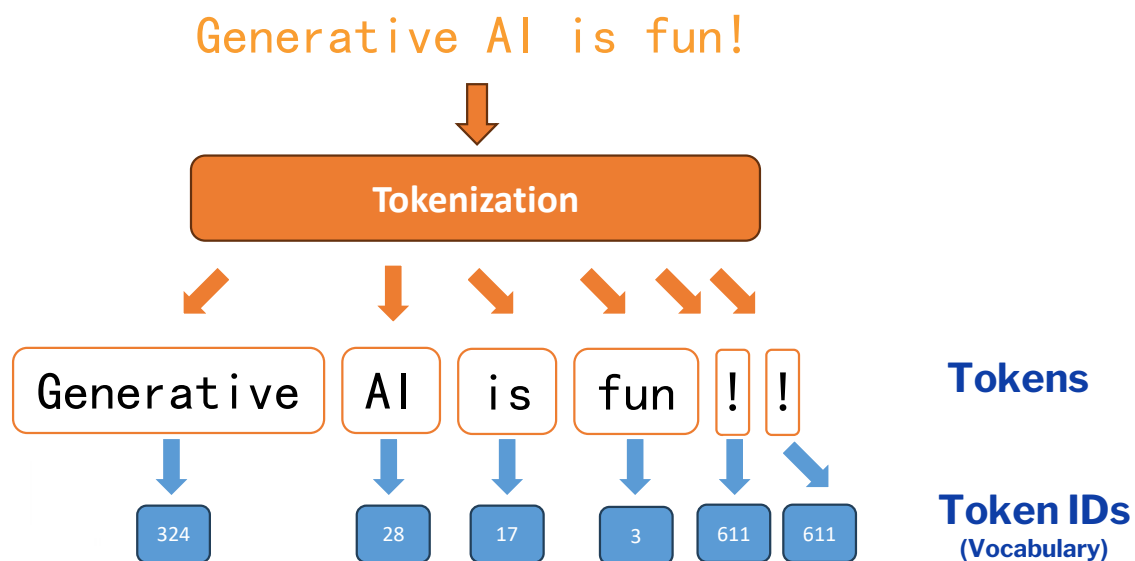
Tokenization



Tokenization

Tokenization is required for:

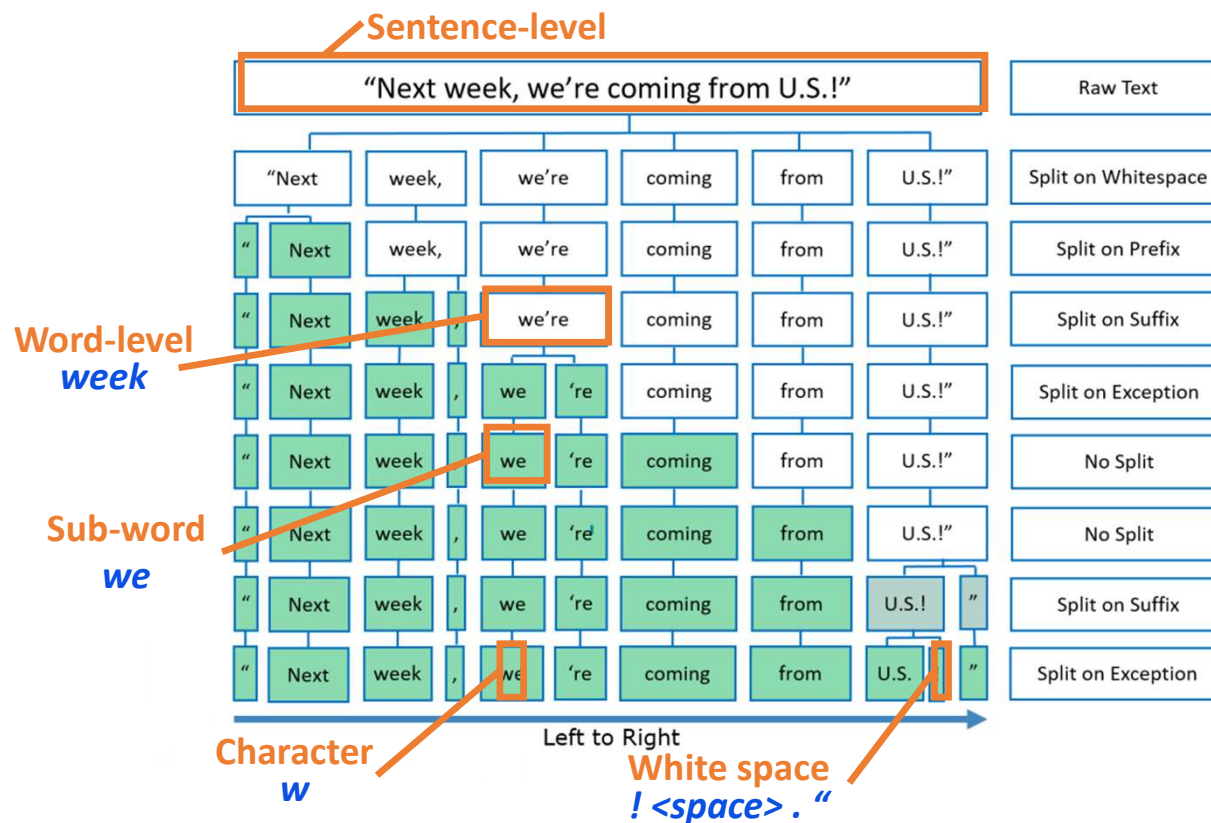
- Numerical representation (computers need numbers)
- Structure unstructured data (structure text)
- Surface meaning (text is messy)
- Efficiency



Tokenization Options

1. Sentence-level
2. Whitespace
3. Word
4. Sub-word (WordPiece)
5. Character-level
6. Rule-based

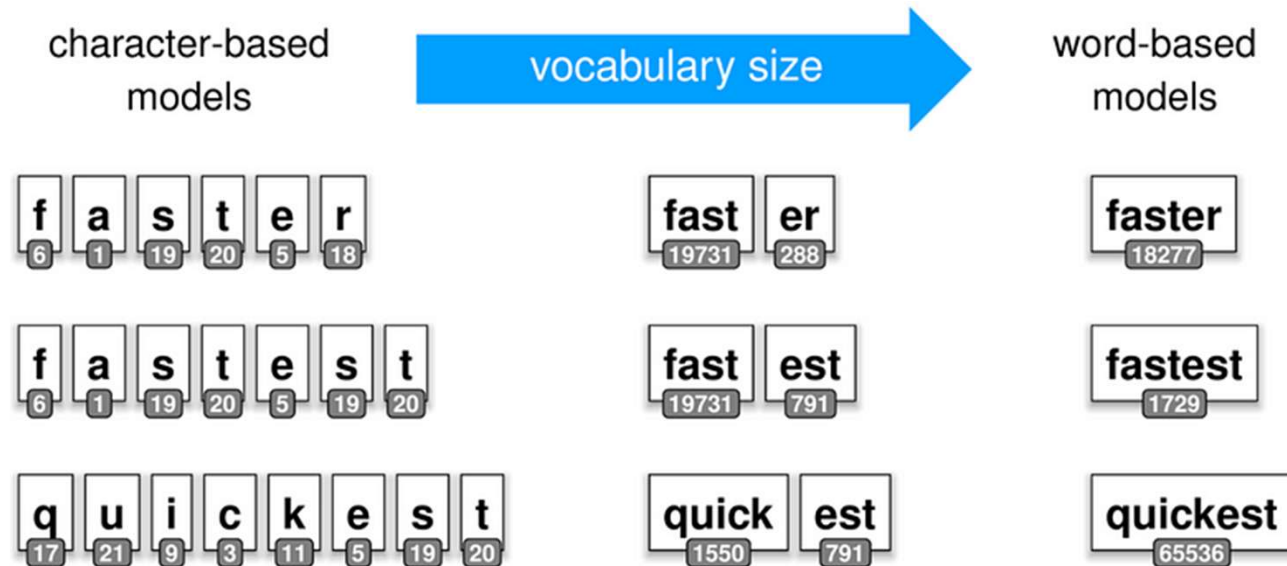
... and many more



Vocabulary Size

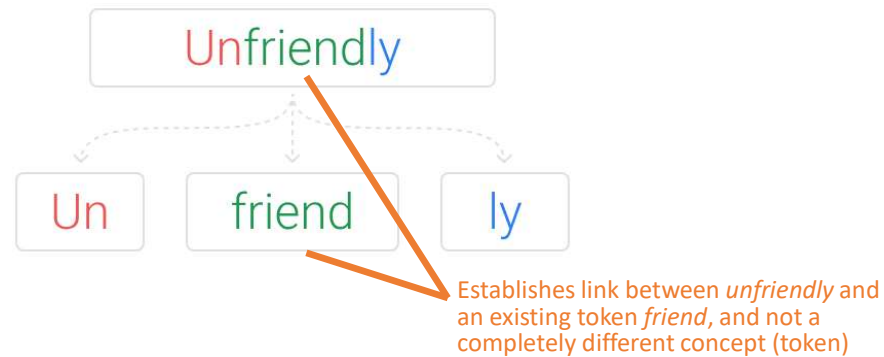
A vocabulary is the set of all tokens (characters, words, sentences, etc.) that are used to represent and analyze text.

There is a trade-off between token granularity (retention of meaning) and vocabulary size.

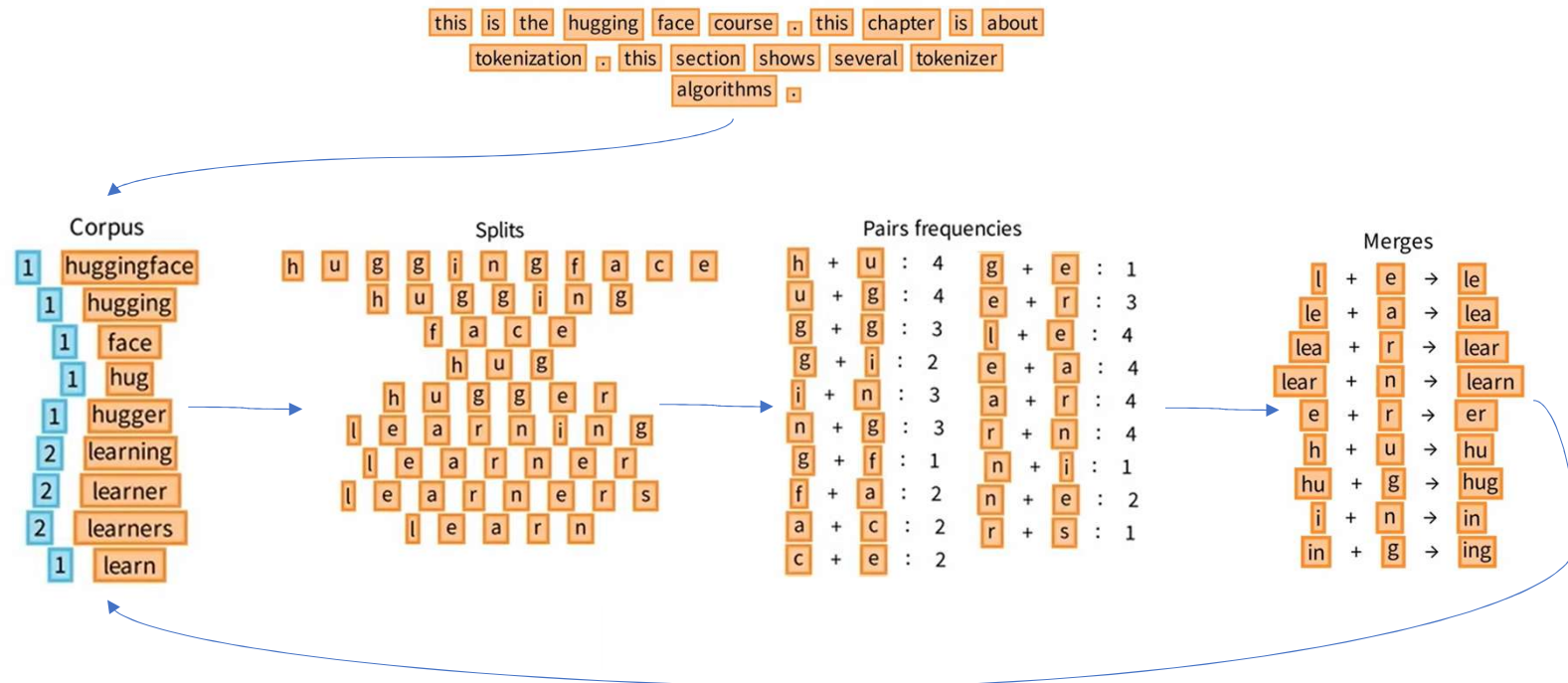


Sub-word

Sub-word is most popular, best balance of vocabulary size and retention of meaning

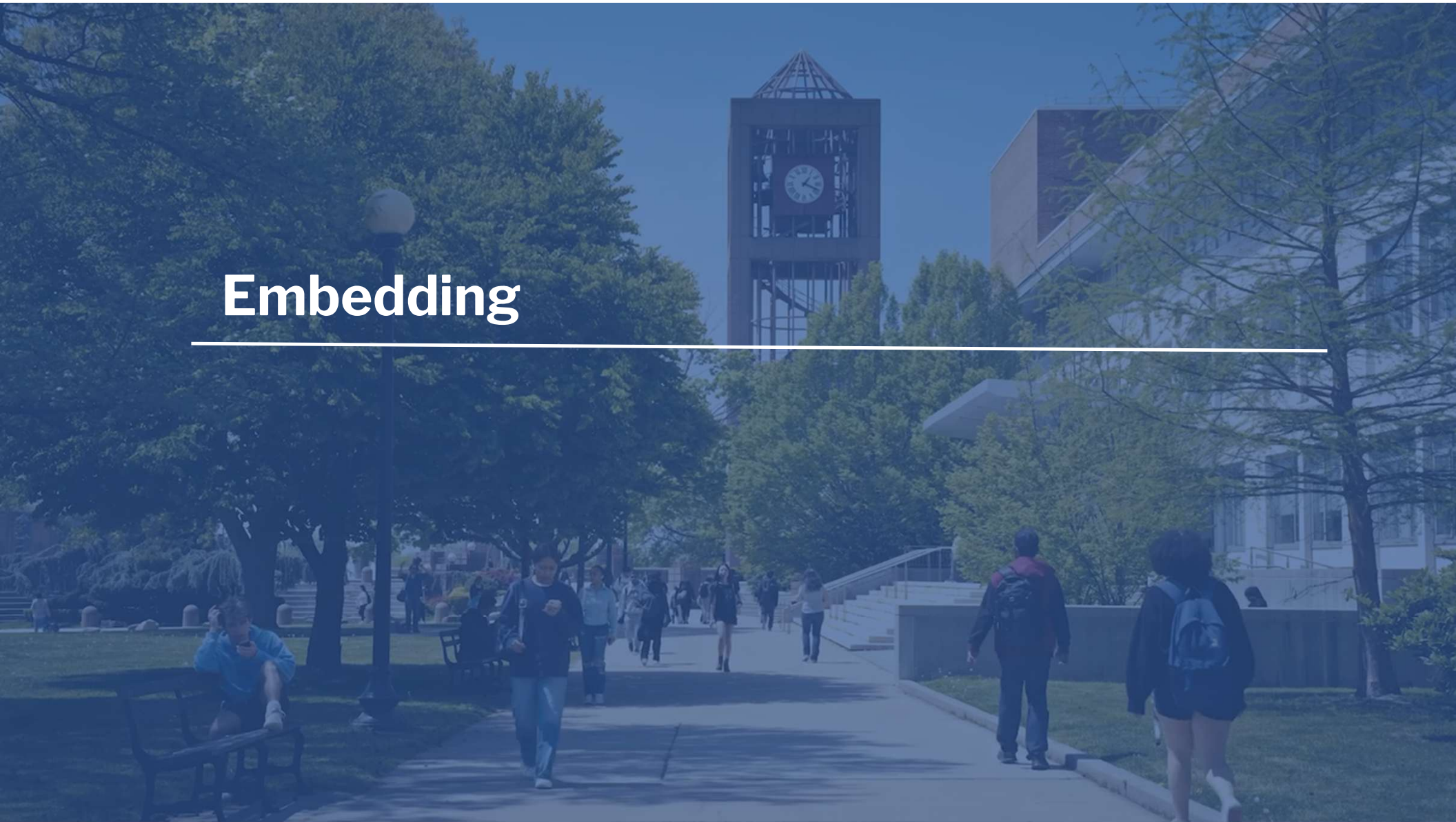


Byte-pair Encoding (BPE)



Source: Huggingface Byte Pair Encoding Tokenization - <https://www.youtube.com/watch?v=HEikzVL-IZU>

Embedding



Word Meaning

A word only has meaning in the context of other words

Words may have many meanings (polysemy).
The meaning of a word depends on its context.

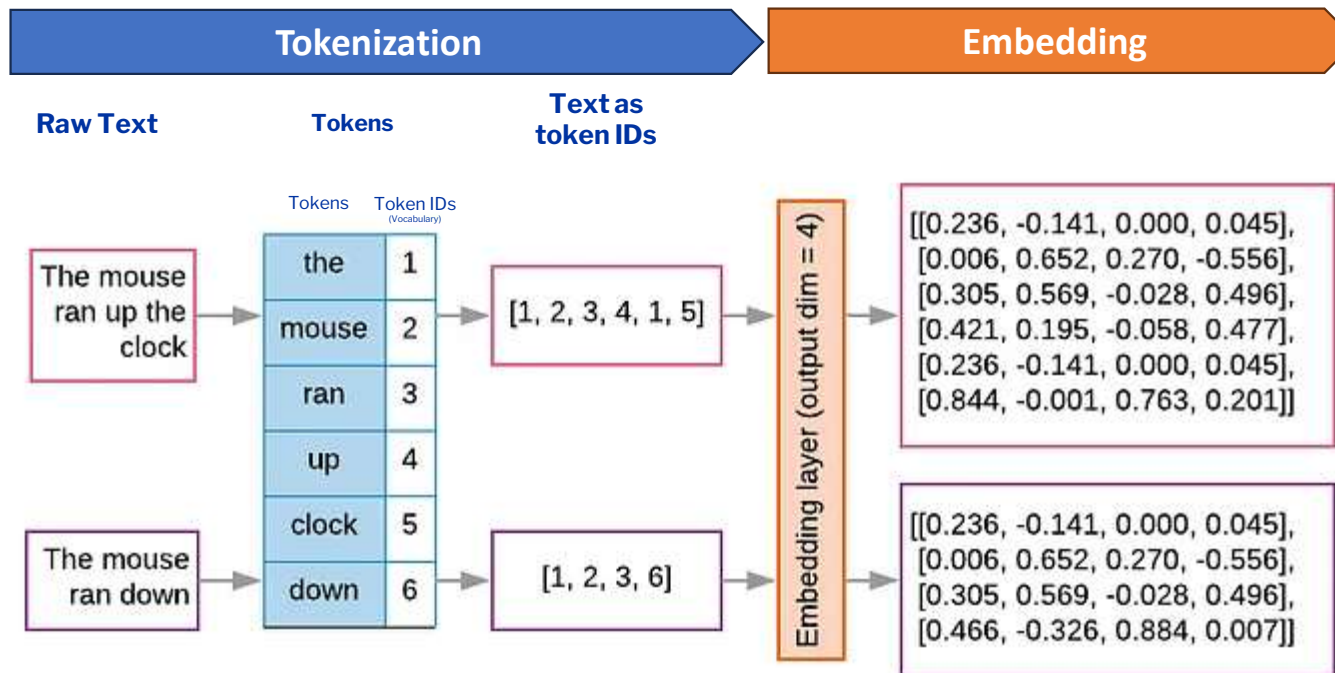
An Example: **bow**

- the front of a ship
- to bend forward in respect
- a weapon that shoots arrows
- a decorative knot
- used to play a violin

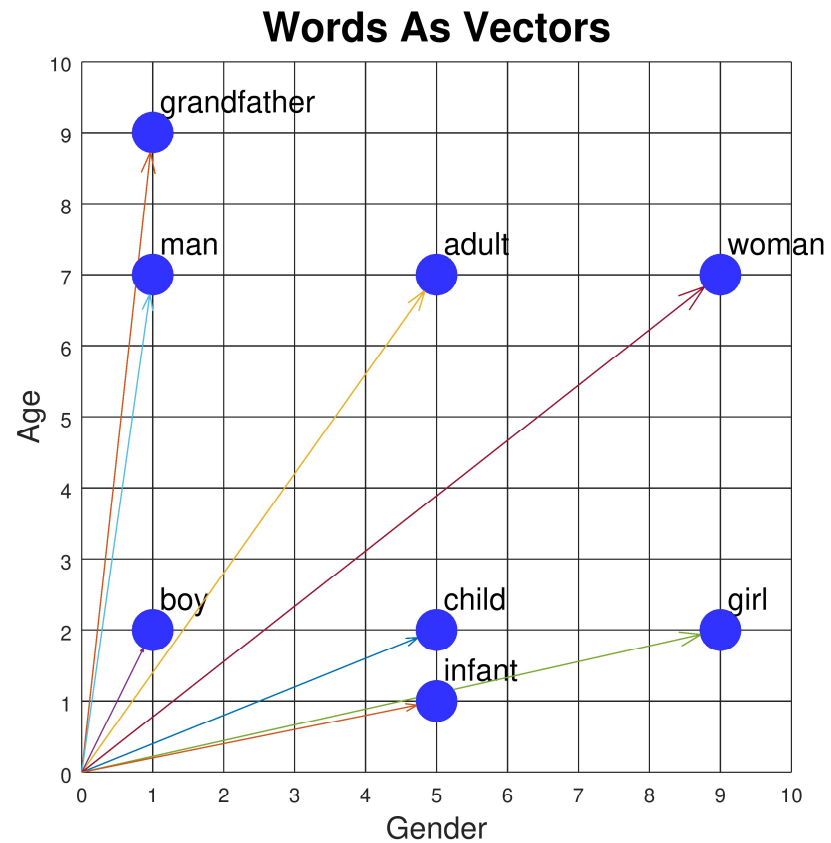


All of these have the **same token** (and token ID)
but they have **different meanings**

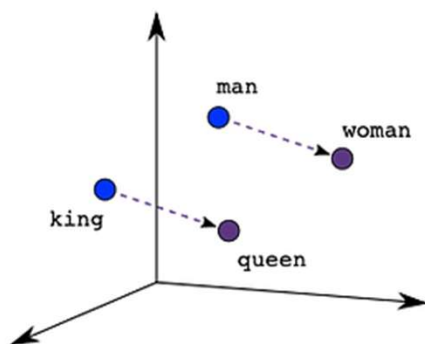
Embedding



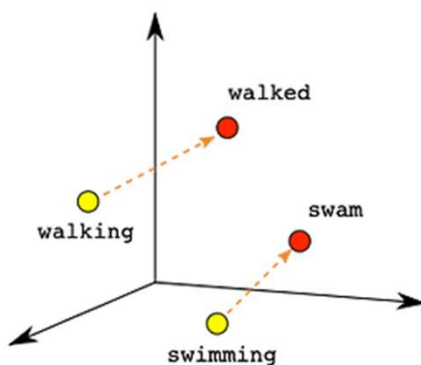
Word2Vec



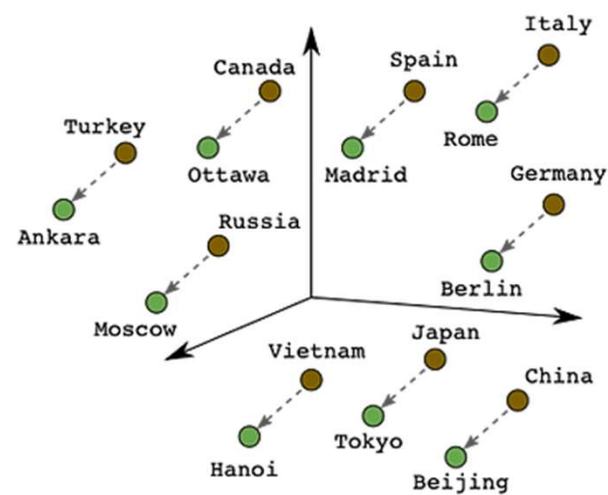
Embedding



Male-Female

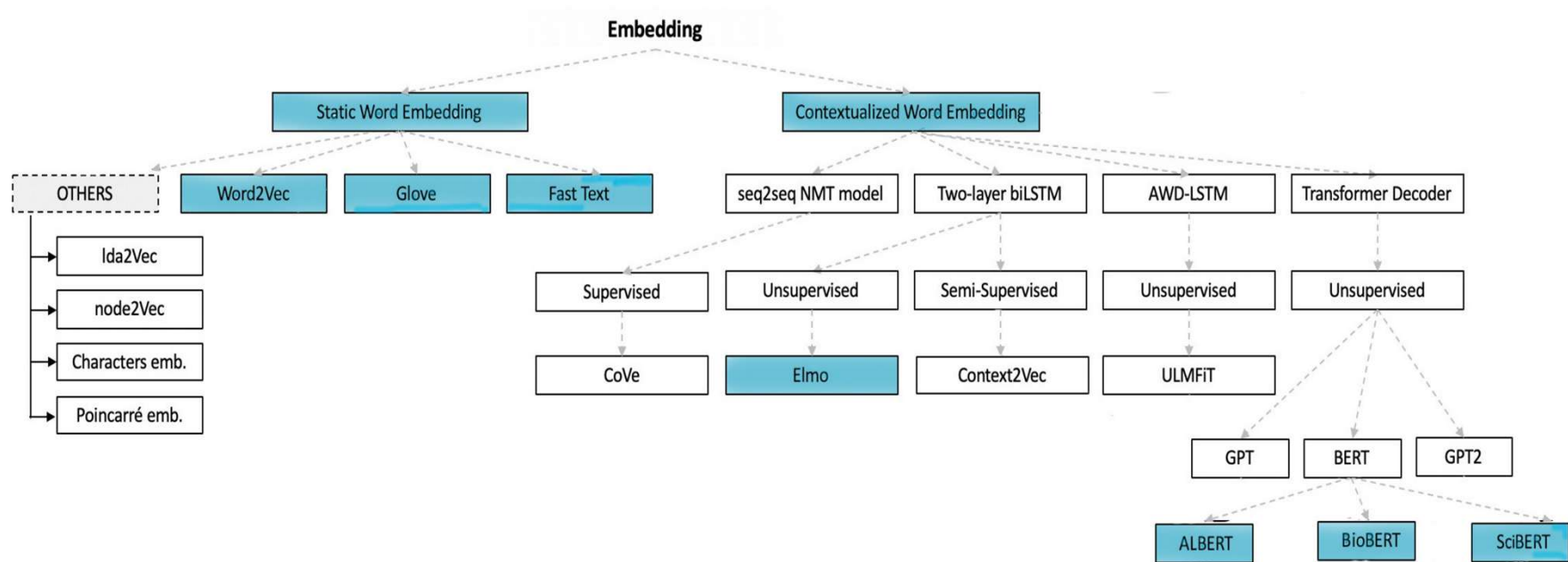


Verb Tense



Country-Capital

Embedding Techniques



Source: Mohiuddin Md Abdul Qudar and Vijay Mago. 2020. A Survey on Language Models.

A blue-tinted photograph of a university campus. In the background, a tall clock tower with a glass-enclosed upper section and a clock face is visible. The foreground shows a wide, paved pedestrian walkway lined with lush green trees. Several students are walking along the path; some are carrying backpacks. On the left, a person is sitting on a wooden bench. The overall scene is bright and clear, suggesting a sunny day.

Natural Language Processing (NLP)

Natural Language Processing Tasks

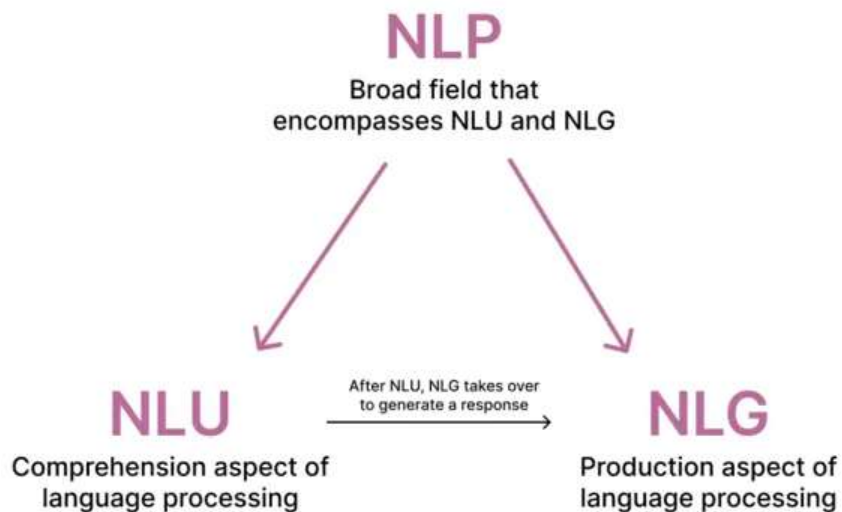
NLP is the process through which **AI is taught to understand the rules and syntax of language**, programmed to develop complex algorithms to represent those rules, and then made to **use those algorithms to carry out specific tasks** like these.

WORD TAGGING	SENTENCE PARSING	TEXT CLASSIFICATION	TEXT GENERATION	TEXT PAIR MATCHING
Tokenization	Constituency parsing	Sentiment analysis	Generative text modeling	Semantic textual similarity
Named entity recognition (NER)	Semantic labeling	Intent detection	Machine translation	Natural language inference (NLI)
Part-of-speech tagging	Dependency parsing	Topic classification	Summarization	Relation extraction
Lemmatization/ Stemming	Coreference parsing	Fake news detection	Personalized dialogue systems	
Word sense disambiguation	Clause boundary detection	Email classification	Report generation	
Keyword extraction		Customer feedback analysis	Question answering (QA)	

Source: Mobidev

Natural Language Understanding & Generation

The difference between **Natural Language Processing (NLP)**, **Natural Language Understanding (NLU)**, and **Natural Language Generation (NLG)**

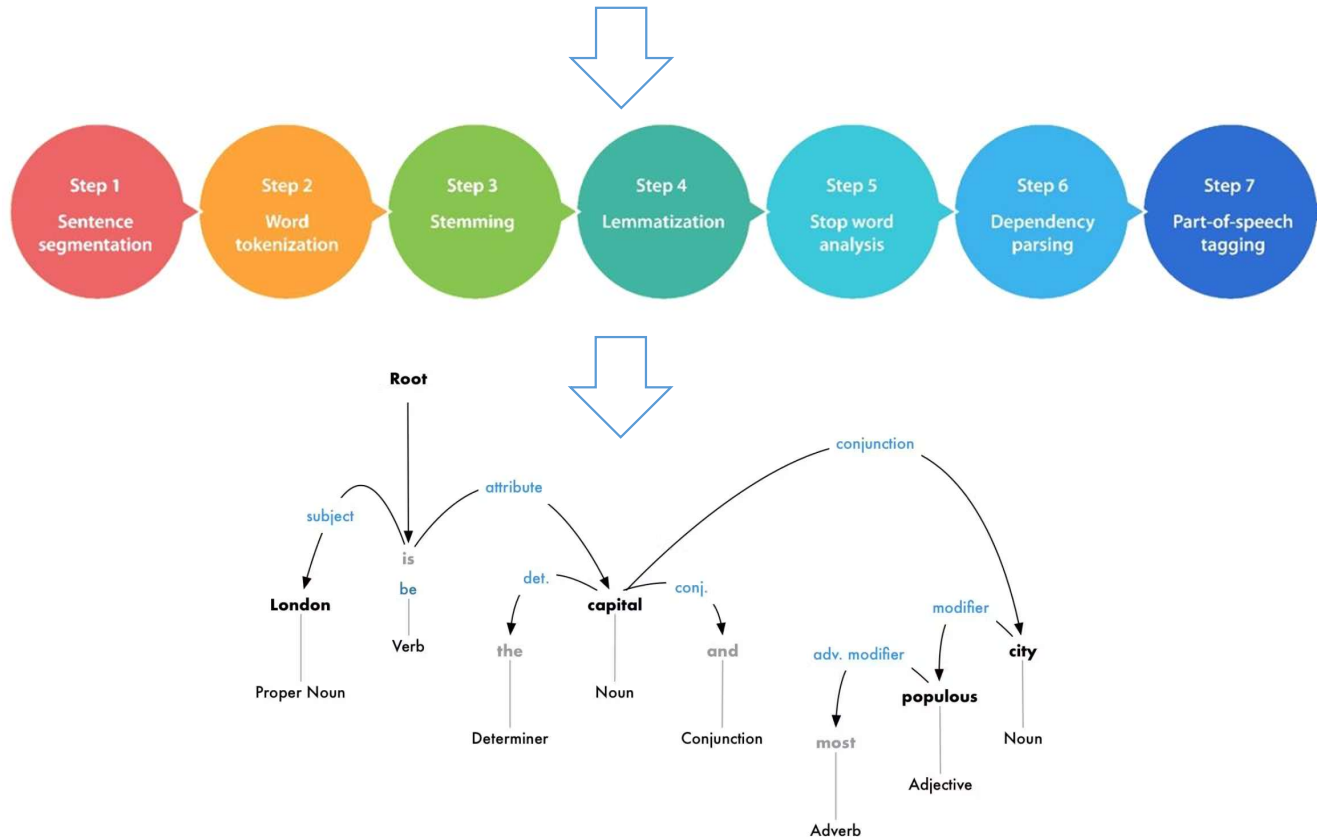


NLU	NLP	NLG
It is a narrow concept.	It is a broader concept.	It is a limited concept.
If we only talk about an understanding text, then it is enough.	But if we want more than understanding, such as decision-making, then it comes into play.	It generates a human-like manner text based on the structured data.
It is a subset of NLP.	It is a combination of it and NLG for conversational Artificial Intelligence problems.	It is a subset of NLP.
It is not necessarily that what is written or said is meant to be the same. There can be flaws and mistakes. It ensures that it will infer correct intent and meaning even if data is spoken and written with some errors. It is the ability to understand the text.	But if we talk about NLP, it is about how the machine processes the given data, such as making decisions, taking action, and responding to the system. It contains the whole End-to-end process. It doesn't need to have it every time.	It generates structured data, but the generated text is not necessarily easy for humans to understand. Thus, NLG ensures that it will be human-understandable.
It reads data and converts it to structured data.	It converts unstructured data to structured data.	NLG writes structured data.

Source: Botpress & Xenonstack

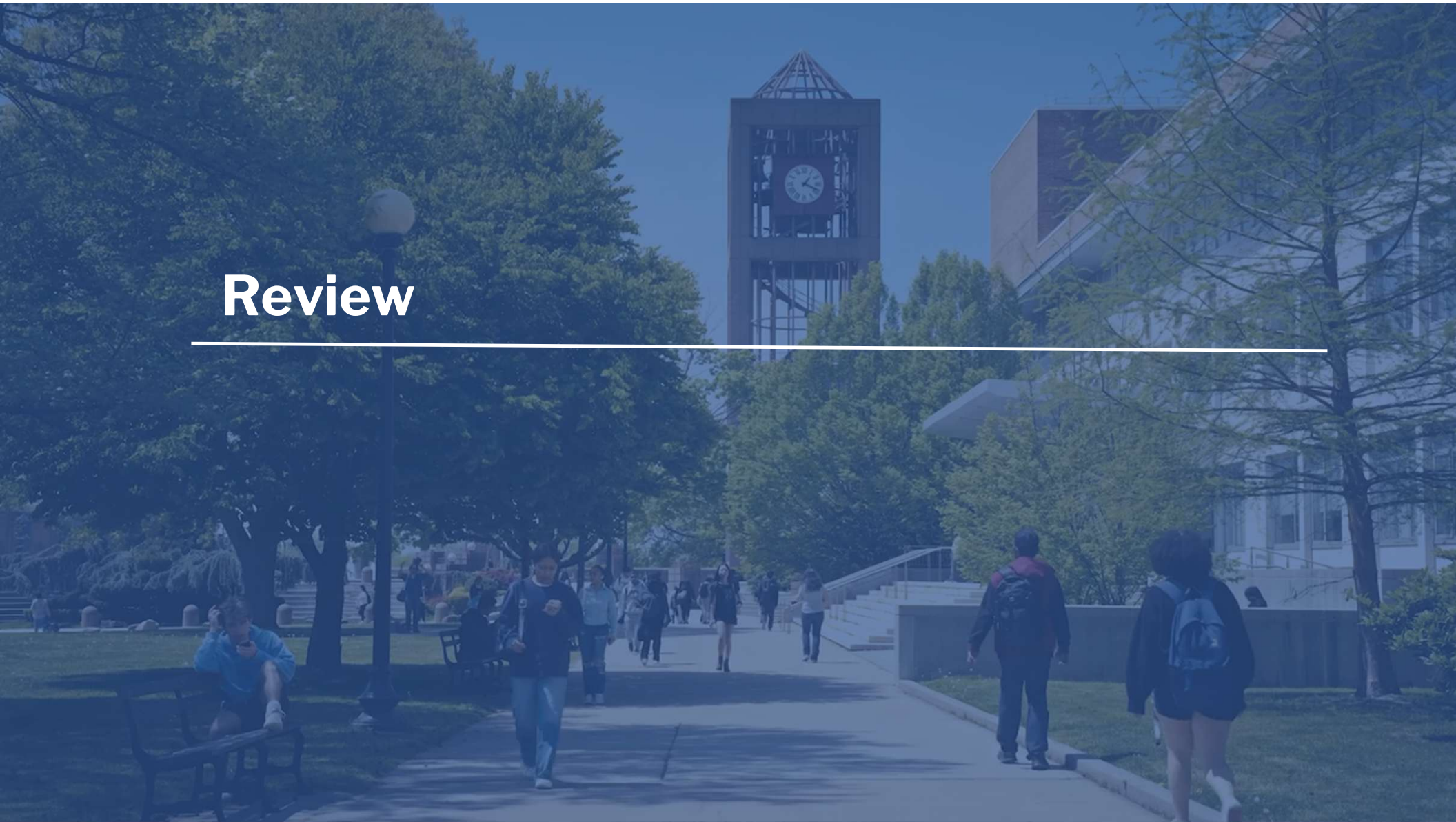
NLP Pipeline

“London is the capital and most populous city of England and the United Kingdom.”



Source: Turing

Review



This week we covered

Course Objectives/Topics (Objectives)

1. What is tokenization and why it matters
2. Bag-of-words vs. modern embedding techniques
3. Sub-word tokenization methods (e.g., byte pair encoding, Word Piece)
4. Language representation and vocabulary construction

Concepts

Tokenization

Bag-of-words

An aerial photograph of the Manhattan skyline, featuring the Manhattan Bridge in the foreground and the dense cityscape in the background. The bridge's stone towers and suspension cables are prominent. The water of the harbor is visible in the lower foreground. The sky is clear and blue.

CU | School of **NY** | Professional Studies