# IMPERIAL

# Word Representation in Biomedical Domain

## Natural Language Processing Team Project
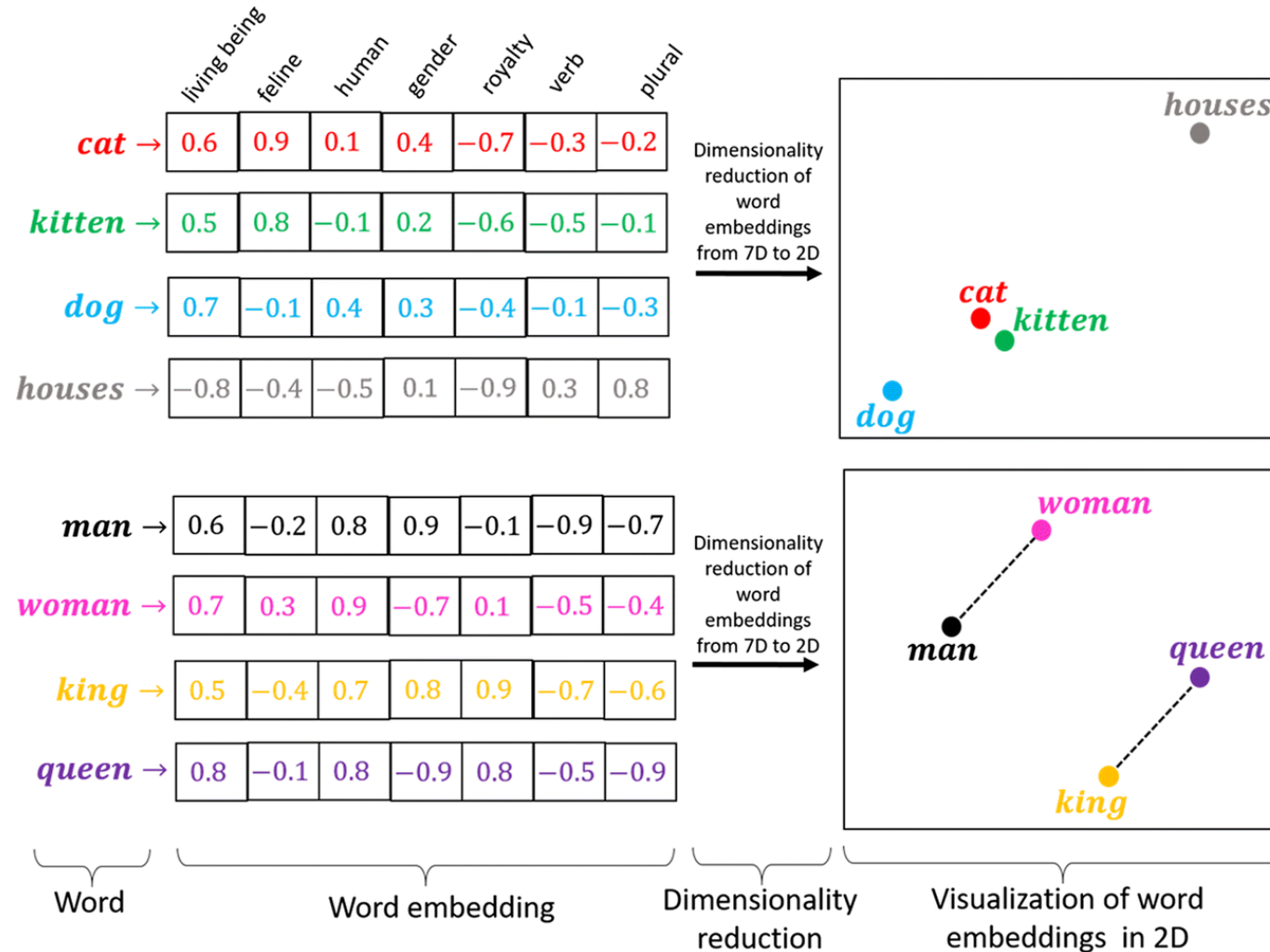
Weihang Zhang, Shuojie Fu

03/02/2025

# Natural Language Processing
## Applications of NLP

- Question Answering
- Text summarization
- Sentiment analysis
- Speech recognition
- Neural Machine Translation

- Search engine
- Smart assistant
- Language translation
- Recommendation system used by major companies
- Healthcare, finance....

# Natural Language Processing
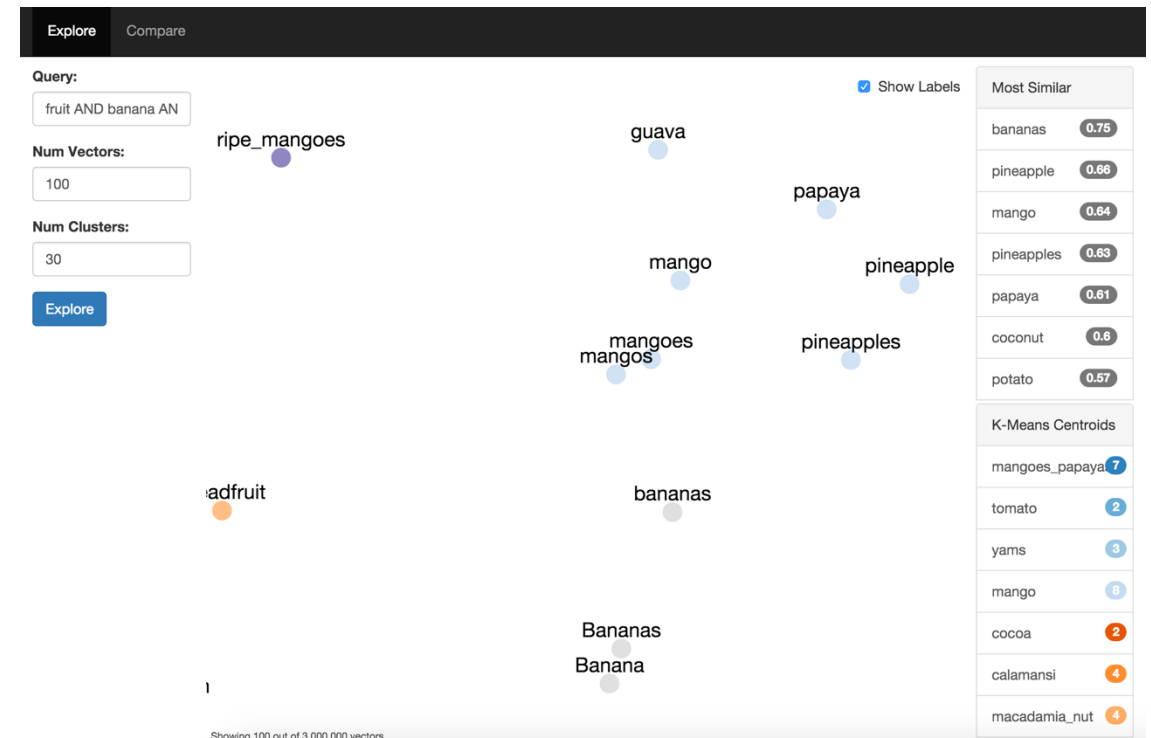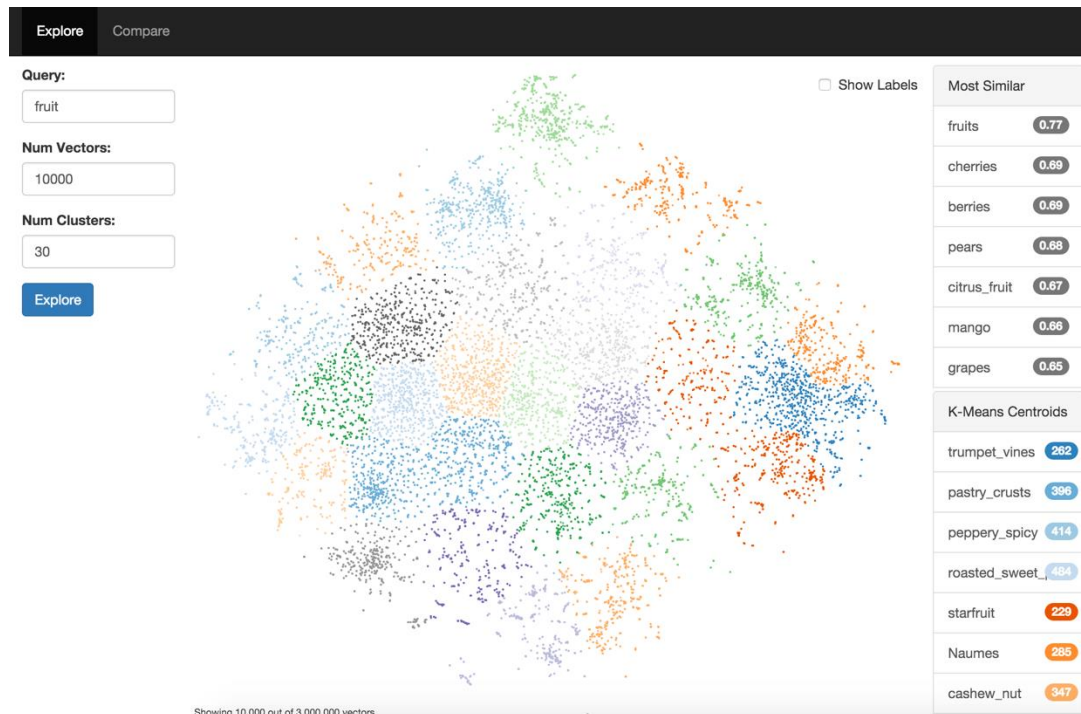## Word Representation



- How can machine learn to represent each word in the human language?
- How can we know if the machine has learnt the semantic meaning of the words?
- What can we do with the learnt word representations?

# Natural Language Processing
## Word Representation

- Mapping each word token to a vector that represents a point in a 'space'.
- Using limited dimensions to represent semantics of words.
- Considering semantic similarity and difference between words.

# Natural Language Processing
## Dataset – CORD-19

- COVID-19 Open Research Dataset (CORD-19)
  - Over 500,000 scholarly articles
  - including over 200,000 with full text
  - about COVID-19, SARS-CoV-2, and related coronaviruses

# Natural Language Processing
## Project Overview

- In this project, you will develop algorithms to learn word representations.
- Part 1: Parse the data
- Part 2: Tokenization
- Part 3: Build word representations
- Part 4: Explore word representations
- Part 5 (Bonus): Mine biomedical knowledge

# Natural Language Processing
## Jupyter Notebook

## Part 1 (20%): Parse the Data

The JSON files are located in two sub-folders in `document_parses` . You will need to scan all JSON files and extract text (i.e. `string` ) from relevant fields (e.g. body text, abstract, titles).

You are encouraged to extract full article text from body text if possible. If the hardware resource is limited, you can extract from abstract or titles as alternatives.

Note: The number of JSON files is around 425k so it may take more than 10 minutes to parse all documents.

For more information about the dataset: https://www.semanticscholar.org/cord19/download

Recommended output:

- A list of text ( `string` ) extracted from JSON files.

```
In [1]:   ##################
          # TODO: add your solution

          ##################
```

# Natural Language Processing
## Report

- For each part
  - Explain the method used
  - Show the results
  - Analyze and discuss the result

# Natural Language Processing
## Marking Scheme

- NLP Project – 40%
  - Completeness – 30%
  - Final report – 10%
  - It is fine to learn from online sources, tutorials, examples, etc. <span style="color:red">However, ensure that all relevant sources are properly cited</span>.
- Breakdown scores for each part (1-5) are available in Jupyter Notebook
- Group work, collaboration is important
- NLP project code and final report
- Cite other's work wherever appropriate

# IMPERIAL

Thank you