ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
ARTIFICIAL INTELLIGENCE

# Graph-Based Keyword Extraction from Scientific Paper Abstracts using Word Embeddings

Author:
Dinno Koluh

Supervisor:
Prof. Paolo Torroni

Co-Supervisor
Dr. Federico Ruggeri

Bologna
16. December 2023.

## Problem statement

- Problem of keyword extraction
  - Important words in text
  - Inherently a ranking problem
- Application to scientific paper abstracts
- Why model the problem as a graph?
  - Well-established model
  - Model text as a graph
  - Ranking algorithms
- The role of word embeddings

# NLP Pipeline
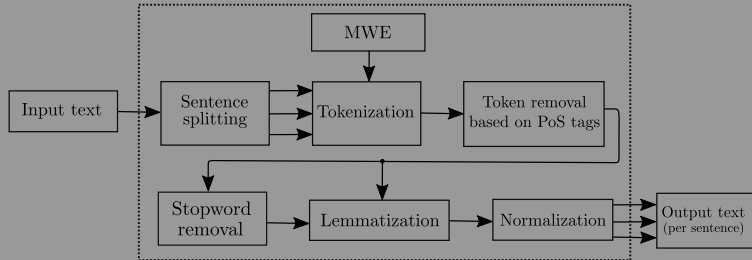
- Text preprocessing
- Broken down into blocks



Figure 1: The NLP pipeline

# Graph construction

- Distributional hypothesis



Figure 2: Co-occurrence matrix construction
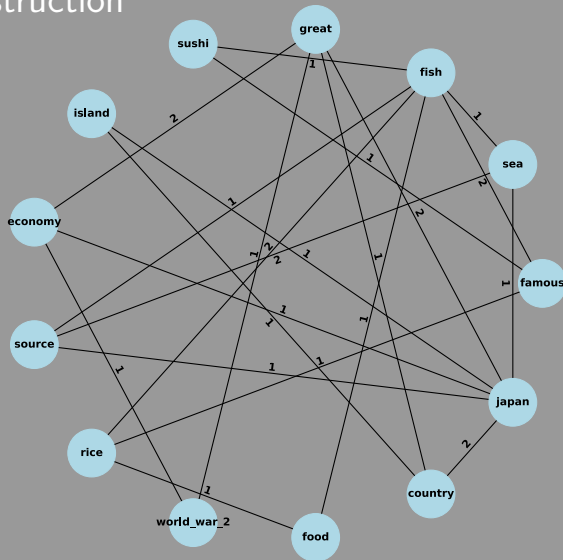
# Graph construction



Figure 3: Graph representation of the co-occurrence matrix

## Ranking algorithms

- Degree centrality (based on node degree)
- Closeness centrality (based on distance to other nodes)
- Betweenness centrality (based on number of shortest paths that pass through the node i.e. information flow)
- Eigenvector centrality (based on direct and neighbour connections)
- PageRank algorithm (Google's web page ranking algorithm applied to text)

# Ranking algorithms

- Application on examples sentences:
- TABLE

# Ranking algorithms
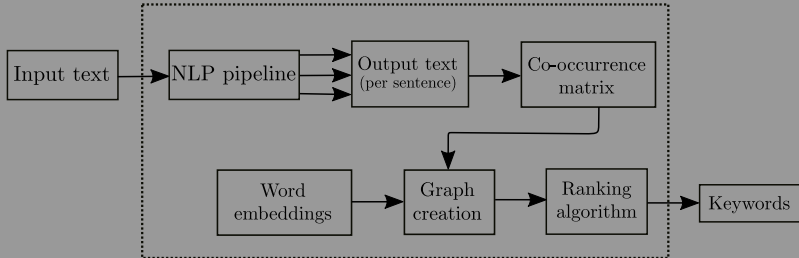
- Full pipeline for keyword extraction



Figure 4: Full pipeline of keyword extraction

# Experiments and Results

- Dataset: 5000 computer science paper abstracts with human assigned keywords, scraped from IEEE Xplore digital library
- Standard evaluation metrics (precision, recall, F-score)
- Important to address
    - Human assigned keywords are subjective and can be generated
    - Expected precision, recall and F-score from available literature in the range of $10 - 40\%$
    - Selecting the top $n$ keywords calculated by the model, where $n$ is the number of true keywords from the dataset

# Experiments and Results

- TABLE

Q&A