## Overview:

This research project focuses on gaining comprehensive insights into the emergence, evolution, and diagnosis of the SARS-CoV-2 virus by leveraging advanced text mining and natural language processing techniques. The dataset utilized for this analysis is the CORD-NER dataset, a machine-readable collection of over 29,000 coronavirus-related articles, with more than 13,000 having full text.

## Data Preprocessing

**Data Collection**: Gathered an extensive dataset from the CDC-coordinated effort, providing a rich source of information on COVID-19.
Cleaning and Formatting: Conducted thorough data cleaning to ensure consistency and removed any irrelevant information.

**Document Ranking**: Implemented a ranking system to prioritize articles based on relevance to the SARS-CoV-2 genome.

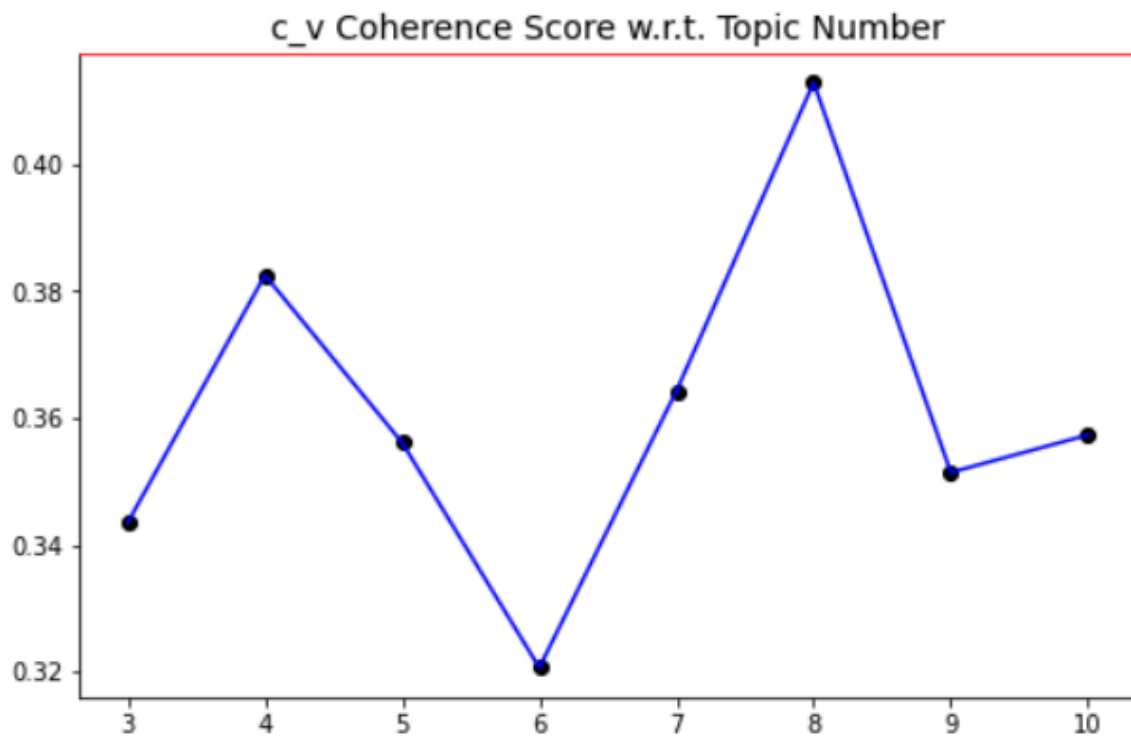## Text Mining and NLP

### Tokenization and Lemmatization

**Tokenization:** Utilized advanced tokenization techniques to break down documents into meaningful units.

**Lemmatization:** Applied lemmatization to reduce words to their base form, enhancing the quality of analysis.

### Coherence Scores

**Topic Modeling:** Employed topic modeling to extract meaningful topics from the dataset.

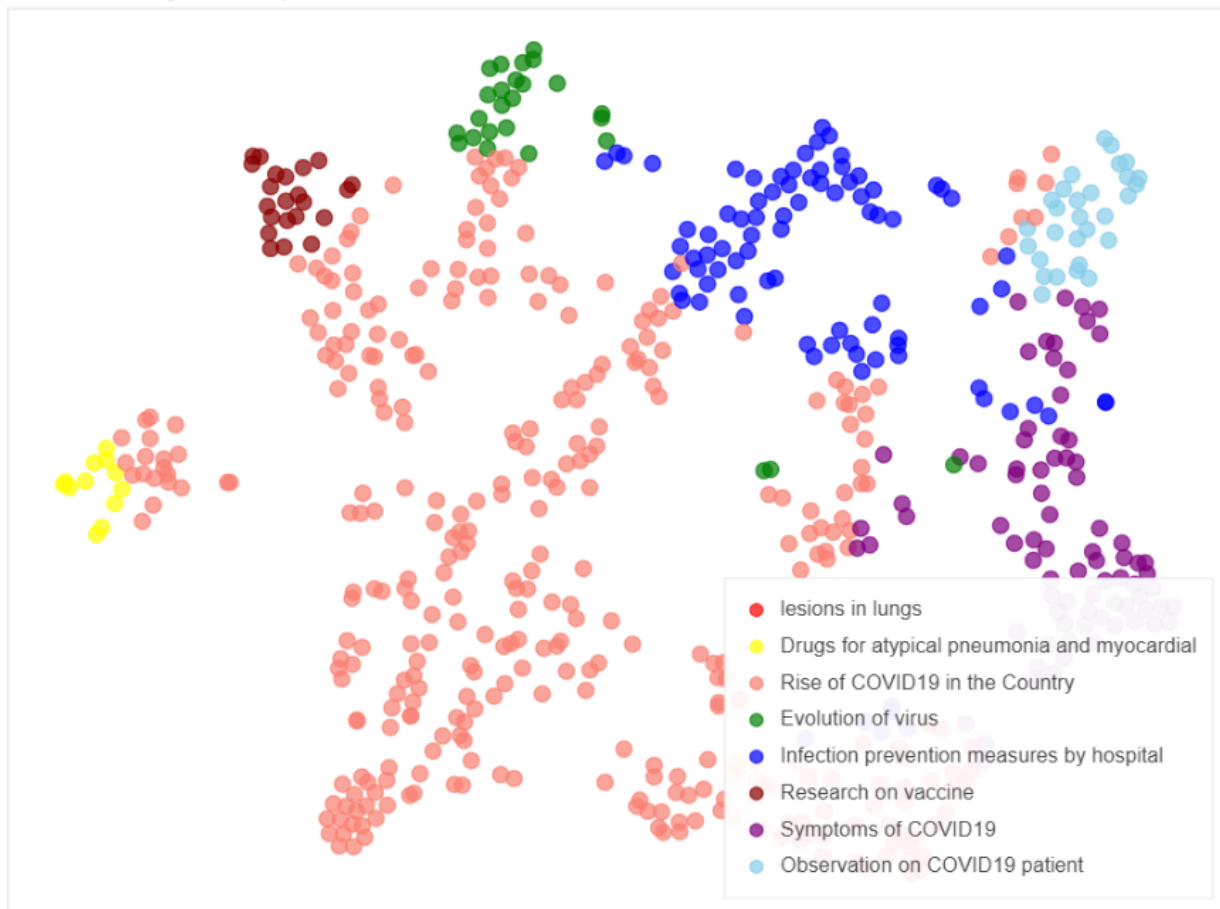**Coherence Scores**: Calculated coherence scores to evaluate the interpretability and relevance of identified topics

c_v Coherence Score w.r.t. Topic Number

## Dimensionality Reduction

**t-SNE**: Implemented t-Distributed Stochastic Neighbor Embedding (t-SNE) to visualize high-dimensional data, facilitating a better understanding of document relationships.

T-SNE Clustering of LDA Topics

- lesions in lungs
- Drugs for atypical pneumonia and myocardial
- Rise of COVID19 in the Country
- Evolution of virus
- Infection prevention measures by hospital
- Research on vaccine
- Symptoms of COVID19
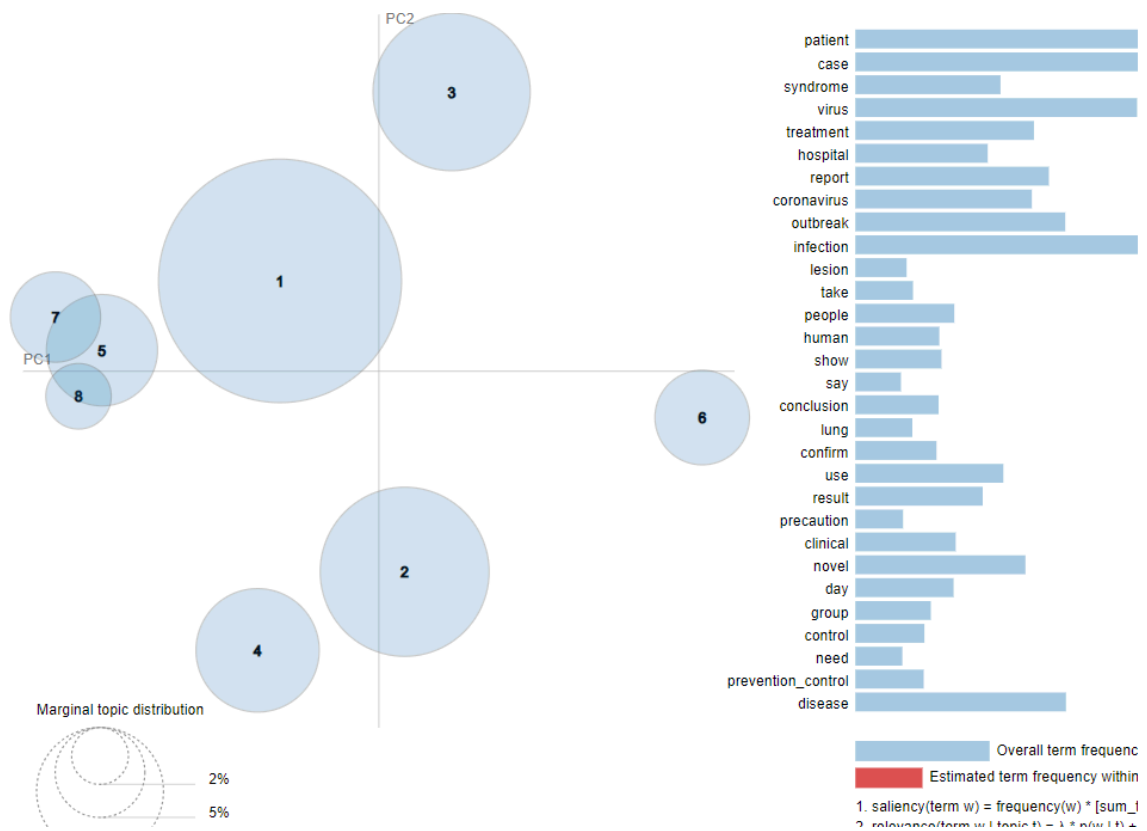- Observation on COVID19 patient

## Word Embeddings

**Word2Vec Embedding:** Utilized Word2Vec embeddings to represent words in vector space, capturing semantic relationships.

## Intertopic Graph

**Intertopic Graph Analysis:** Constructed a graph to visualize relationships and interactions between identified topics, providing a holistic view of the dataset.

PC2

| Term | |
|------|---|
| patient | |
| case | |
| syndrome | |
| virus | |
| treatment | |
| hospital | |
| report | |
| coronavirus | |
| outbreak | |
| infection | |
| lesion | |
| take | |
| people | |
| human | |
| show | |
| say | |
| conclusion | |
| lung | |
| confirm | |
| use | |
| result | |
| precaution | |
| clinical | |
| novel | |
| day | |
| group | |
| control | |
| need | |
| prevention_control | |
| disease | |

Marginal topic distribution

2%

5%

Overall term frequenc

Estimated term frequency within

1. saliency(term w) = frequency(w) * [sum_t
2. relevance(term w | topic t) = λ * p(w | t) +

## Spacy Parser

**Spacy Integration:** Incorporated Spacy for advanced natural language processing tasks, enhancing document understanding.