# NEWS AGGREGATOR

Alaa Mohamed– Marwan Ibrahim– Mayar Yasser– Mohamed Ayman

**SUPERVISED BY:**

**Dr. Walaa Hassan and Eng. Mennatullah Gamil**

# Introduction 1/2

- In the last few years, The world had an incredible and huge growth of the rate of news.

- 44 percent of U.S. consumers cited some sort of online publication as their <u>main source of news</u> in 2017[1].

- At February 2020, Yahoo News had 175 million unique monthly visitors and CNN 95 million unique visitors per month[2].

[1]https://www.statista.com/topics/1640/news/
[2]https://www.statista.com/statistics/381569/leading-news-and-media-sites-usa-by-share-of-visits/

# Introduction 2/2

- Our project aims to aggregate news from as much trusted sources as possible in one place.


- Main tasks for our project is Summarizing those aggregated articles and to Recommend articles according to user's keyword or key phrase.

# Related Work

| Paper Name | NewsOne- An Aggregation System For News Using Web Scrapping Method[3] | An Analysis of Recommender Algorithms for Online News[4] | A news-topic recommender system[5] |
|---|---|---|---|
| **Year** | • 2017 | • 2014 | • 2018 |
| **Aim** | • Collect the news from multiple sites and merge them all in a summarized website.<br>• Deals with the news URL and save them in the database. | • Implement a recommender system for online news articles. | • Make a news-topic recommender system based on keywords extraction. |
| **Tools** | • Web scraping (To get more accurate results).<br>• RSS fetcher (Summary reports from specific websites). | • Content-based filter.<br>• Elastic Search.<br>(To get data that is similar to the articles and events that the user likes or searches). | • Use RAKE(Rapid automatic keyword extraction) operates on individual documents to extract keywords.<br>• Use Keyword scoring to extract topic keywords for a specific time. |
| **Disadvantage** | • Provides only the source URL of the news.<br>• Provides only the title of the news. | • Use a static dataset from 25/5/2014 to 30/6/2014. | • Get news from specific(static) Sources as CNN ,Fox news and The New York Times. |

[3]K. Sundaramoorthy, R. Durga, and S. Nagadarshini, "Newsone—an aggregation systemfor news using web scraping method," in2017 International Conference on TechnicalAdvancements in Computers and Communications (ICTACC). IEEE, 2017, pp. 136–140.
[4]D. Doychev, A. Lawlor, R. Rafter, and B. Smyth, "An analysis of recommender algo-rithms for online news," inCLEF 2014 Conference and Labs of the Evaluation Forum:Information Access Evaluation Meets Multilinguality, Multimodality and Interaction,15-18 September 2014, Sheffield, United Kingdom, 2014, pp. 177–184.
[5]Z. Wang, K. Hahn, Y. Kim, S. Song, and J.-M. Seo, "A news-topic recommender systembased on keywords extraction,"Multimedia Tools and Applications, vol. 77, no. 4, pp.4339–4353, 2018.
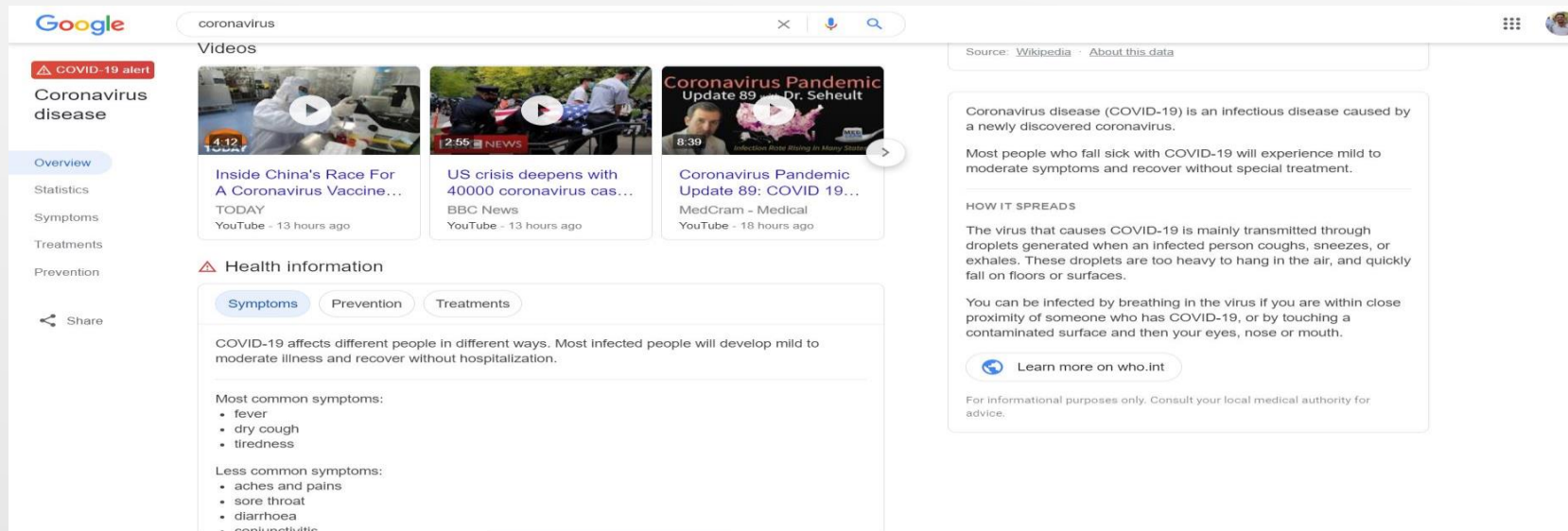
# Problem Statement

- **Aggregate/Recommend** highly relevant news to a certain input keyword or key phrases from heterogeneous sources and **Summarize** them to reduce user's read time.

# Market Need 1/2

- Especially this year, news rate are higher than expected.

- Coronavirus is live example and a main reason for this growth.

- People need to be kept updated with all of the news around the world.



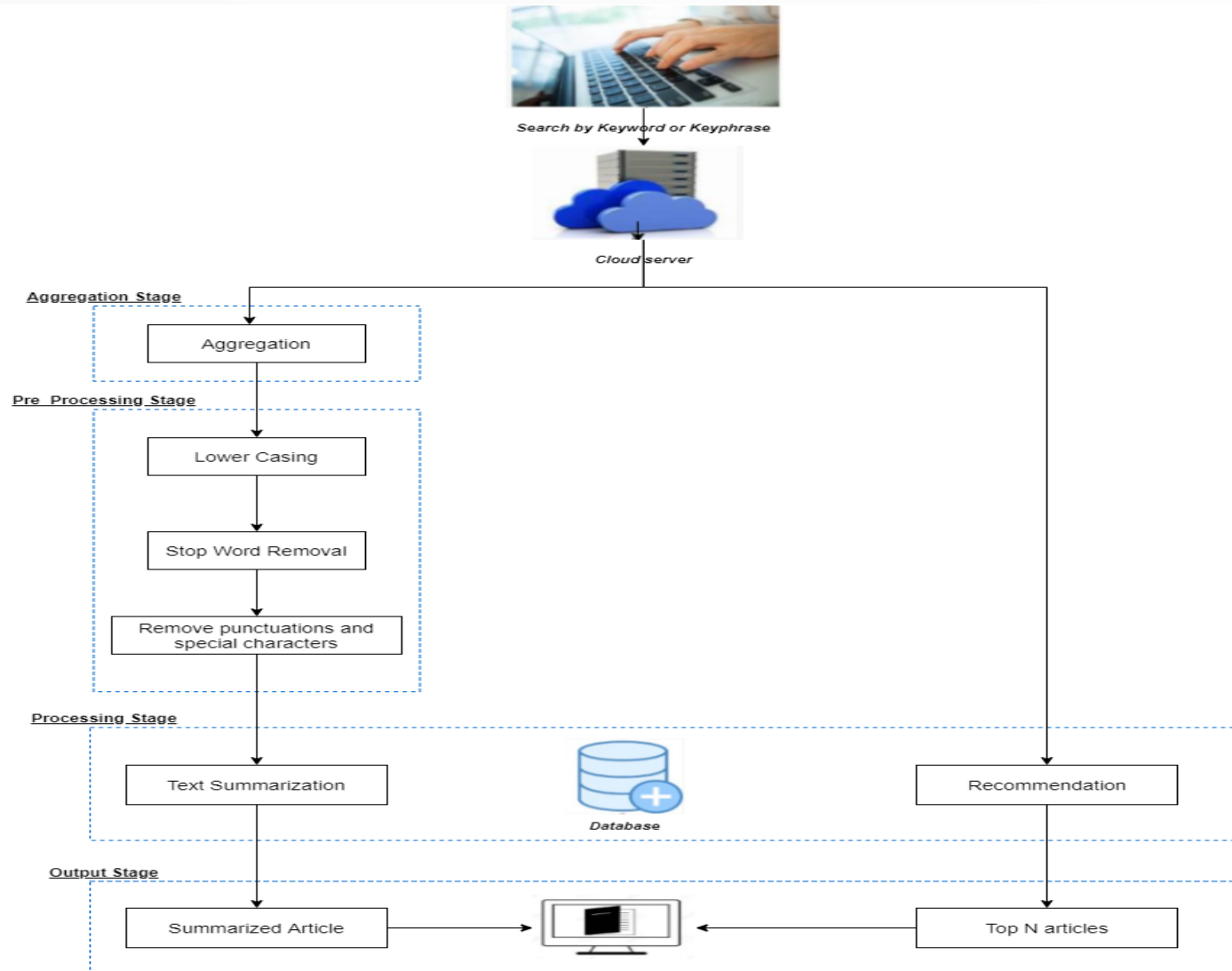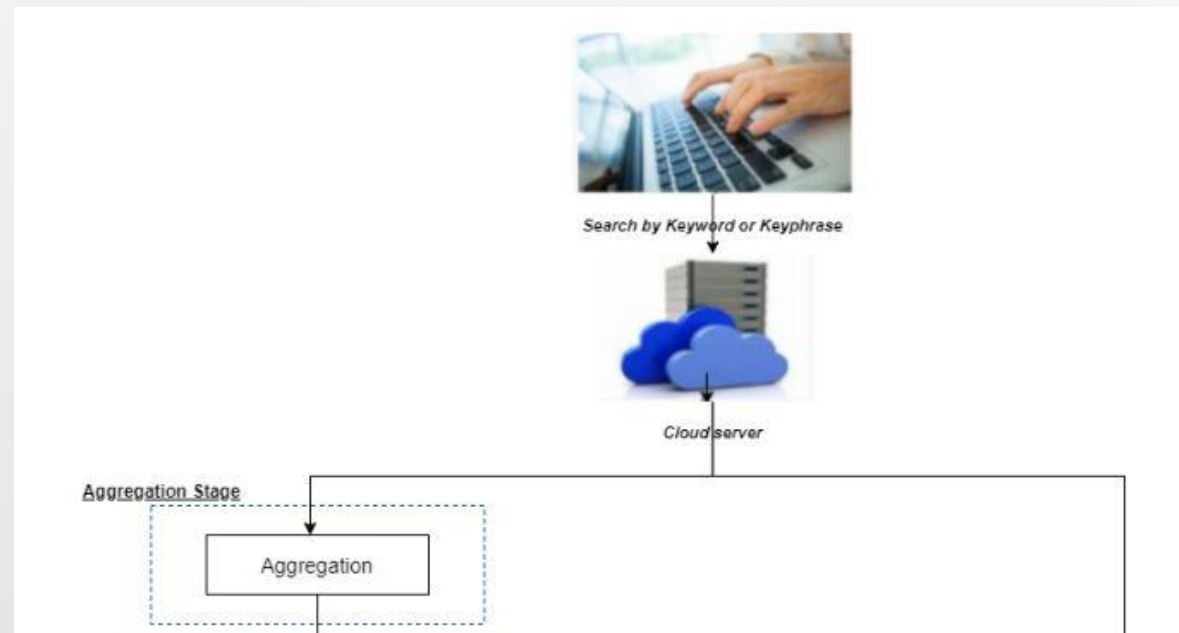So , what will the system provide as a solution ?

# Market Need 2/2



Understandable summary.



Recommended articles from trusted sources.

# System Overview

# Aggregation Stage

- Rich Site Summary (RSS) Feeds.

- Web Scraping with Beautiful Soup library.

# Pre-processing Stage

- Get Aggregated online content.

- Proposed pre-processing techniques:
  - ❑ Lowercasing
  - ❑ Stop Word removal
  - ❑ Remove punctuations and special characters.

- All articles are in the same format.



Pre Processing Stage

Lower Casing

Stop Word Removal

Remove punctuations and special characters

# Processing Stage

- Applying Summarization technique to the aggregated/related articles.

- Recommend articles for the search keyword/key-phrase.

# Output Stage

- Getting an understandable summarized article.

- Recommended articles for the search keyword/key-phrase.

# Main Algorithms

❑Two main algorithms in the system :

1. Summarization Algorithm.

    • Text-Rank.

2. Recommendation Algorithm.

    • Content-based filter.

# Summarization algorithm 1/2

## What is Text-Rank?

TextRank Algorithm is a graph-based ranking algorithm that is used for articles summarization.

# Summarization algorithm 2/2

- ## Vectors:
  Using global vector(GloVe)
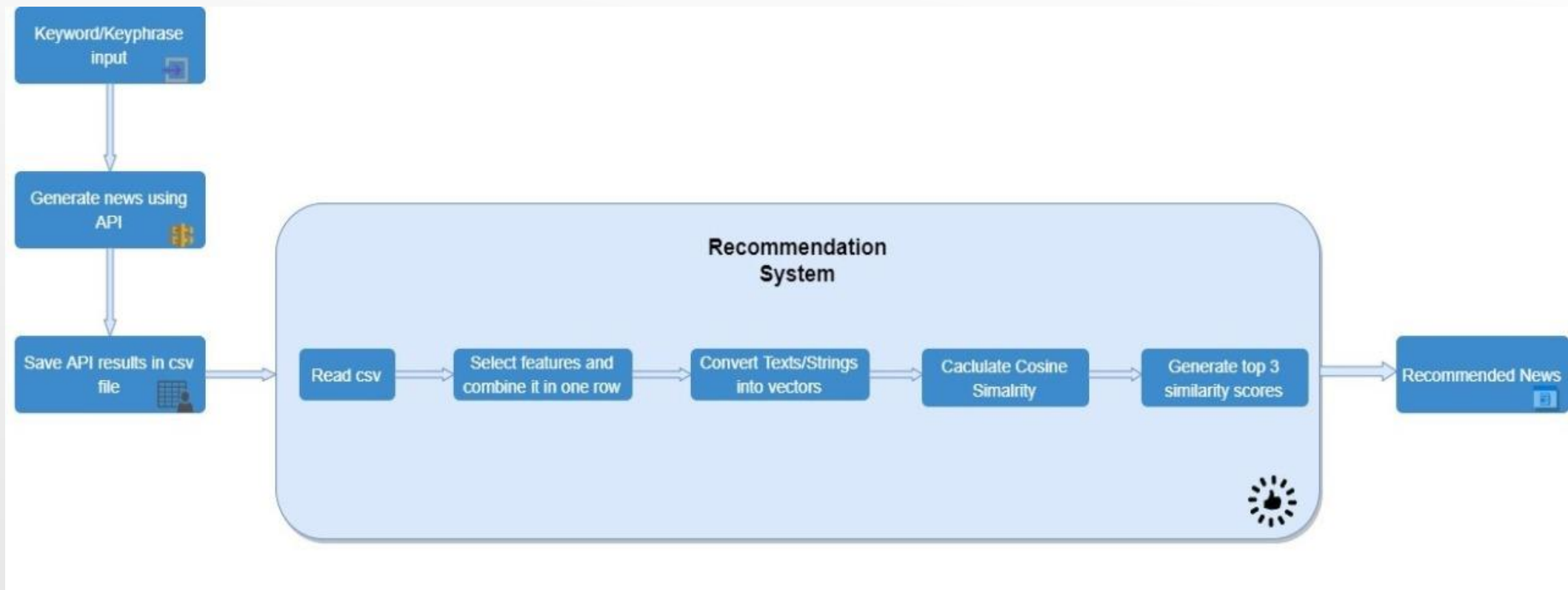
- ## Similarity Matrix:
  Using cosine similarity

[6]https://www.pdfs.semanticscholar.org/b20e/a99a35378f4ba09dcd33b2f897b5c8025b8e.pdf
[7]https://medium.com/the-artificial-impostor/use-textrank-to-extract-most-important-sentences-in-article-b8efc7e70b4

# Recommendation algorithm 1/2 [8]

<span style="color:red">**Content-based filter algorithm**</span>

Uses item features to recommend other items similar to user's search.

[8] https://developers.google.com/machine-learning/recommendation/content-based/basics

# Recommendation algorithm 2/2 [9]



[9]https://medium.com/code-heroku/building-a-movie-recommendation-engine-in-python-using-scikit-learn-c7489d7cb145

# Experiments & Results

❑Two experiments done on different approaches in the system :

1. Summarization experiment on two different algorithms:

   I.    Word-frequency algorithm.

   II.   Text-Rank algorithm.

2. Recommendation experiment.

# Experiment 1/2 : Summarization

❑Two approaches were tested to get the most <span style="color:red">efficient</span> summary :

| Input keyword : Football | |
|---|---|
| **Word-frequency Summary** | **Text-Rank Summary** |
| Football (or soccer as the game is called in some parts of the world) has a long history. Football in its current form arose in England in the middle of the 19th century. But alternative versions of the game existed much earlier and are a part of the football history. The first known examples of a team game involving a ball, which was made out of a rock, occurred in old Mesoamerican cultures for over 3,000 years ago. Cuju was played with a round ball on an area of a square. The ball was made by shreds of leather filled with hair (the first documents of balls filled with air are from the 7th century). In the Ancient Rome, games with balls were not included in the entertainment on the big arenas, but could occur in exercises in the military. It was the Roman culture that would bring football to the British island (Britannica). It is, however, uncertain in which degree the British people were influenced by this variety and in which degree they had developed their own variants. The sport we in the United States know as football is more properly called gridiron football, for the vertical yard lines that mark the field. Closely related to two English sports_rugby and soccer (or association football)_gridiron football originated at universities in North America, primarily the United States, in the late 19th century. On November 6, 1869, players from Princeton and Rutgers held the first intercollegiate football contest in New Brunswick, New Jersey, playing a soccer-style game with rules adapted from the London Football Association. In May 1874, after a match against McGill University of Montreal, the Harvard players decided they preferred McGills rugby-style rules to their own. Even more importantly, he was the guiding force on the rules board of the newly formed Intercollegiate Football Association (IFA). | Summary: In 1875, Harvard and Yale played their first intercollegiate match, and Yale players and spectators (including Princeton students) embraced the rugby style as well.The man most responsible for the transition from this rugby-like game to the sport of football we know today was Walter Camp, known as the Father of American Football.

Football (or soccer as the game is called in some parts of the world) has a long history.

But alternative versions of the game existed much earlier and are a part of the football history.

While a number of other elite Northeastern colleges took up the sport in the 1870s, Harvard University maintained its distance by sticking to a rugby-soccer hybrid called the Boston Game.

The sport we in the United States know as football is more properly called gridiron football, for the vertical yard lines that mark the field.

The first known examples of a team game involving a ball, which was made out of a rock, occurred in old Mesoamerican cultures for over 3,000 years ago.

Closely related to two English sports_rugby and soccer (or association football)_gridiron football originated at universities in North America, primarily the United States, in the late 19th century.

Thanks to Camp, the IFA made two key innovations to the fledgling game: It did away with the opening scrummage or scrum and introduced the requirement that a team give up the ball after failing to move down the field a specified yardage in a certain number of downs. Among the other innovations Camp introduced were the 11-man team, the quarterback position, the line of scrimmage, offensive signal-calling and the scoring scale used in football today.

On November 6, 1869, players from Princeton and Rutgers held the first intercollegiate football contest in New Brunswick, New Jersey, playing a soccer-style game with rules adapted from the London Football Association.

According to the sources, the ball would symbolize the sun and the captain of the losing team would be sacrificed to the gods.The first known ball game which also involved kicking took place In China in the 3rd and 2nd century BC under the name Cuju. |

# Experiment 1/2 : Results

❑ A survey was applied using our social network, asking to rate the two summaries out of 5.



**Word-Frequency chart**

**Text-Rank chart**

# Experiment 1/2 : Results

❑ To further emphasize that, the output summary from the Text-Rank algorithm was revised by an expert besides the normal readers rating.

❑ The feedback considers that the output summary fulfills all the requirements of a good summary such as:

- The length is about 10% of the original.

- Short paragraphs that contain the key ideas and to the point.

- Could be read and clearly understood without referring to the original

  article.

- Used the appropriate language just like that used in the original article.

❑ Text-Rank algorithm was the chosen approach to be applied in the summarization system over Word Frequency algorithm. The reason behind applying Text-Rank algorithm was simply that Text-Rank gives more efficient and understandable summary for the reader.

# Experiment 2/2 : Recommendation

- In this experiment, we have compared between **n-grams (Uni-gram and Bi-gram)** using count vectorizer to get the best results.

Definitions :

❑ **Uni-gram** :  The occurrence of each word is independent of its previous word.

❑ **Bi-gram** :    The occurrence of each word depends only on its previous word.

# Experiment 2/2 : Results

Covid-19

**Unigram**

- Imperial County Advised To Reinstitute Stay-At-Home Order As COVID-19 Cases Overwhelm System
  https://sanfrancisco.cbslocal.com/2020/06/26/imperial-county-stay-at-home-order-covid19-coronavirus-overwhelm-system/amp/

- UPDATE: 39444 COVID-19 cases, 577 deaths reported in Tennessee; - WRCBtv.com
  https://www.wrcbtv.com/story/41867206/update-39444-covid19-cases-577-deaths-reported-in-tennessee-hamilton-co-2356

- Live updates: COVID-19 clusters linked to 4 bars in Minneapolis, Mankato
  https://www.kare11.com/amp/article/news/health/coronavirus/covid-19-coronavirus-in-minnesota-minneapolis-st-paul-wisconsin-live-updates-june-26-2020/89-1085650a-869e-4bf1-9730-ac48c2ec43e8

**Bigram**

- Imperial County Advised To Reinstitute Stay-At-Home Order As COVID-19 Cases Overwhelm System
  https://sanfrancisco.cbslocal.com/2020/06/26/imperial-county-stay-at-home-order-covid19-coronavirus-overwhelm-system/amp/

- UPDATE: 39444 COVID-19 cases, 577 deaths reported in Tennessee; - WRCBtv.com
  https://www.wrcbtv.com/story/41867206/update-39444-covid19-cases-577-deaths-reported-in-tennessee-hamilton-co-2356

- Live updates: COVID-19 clusters linked to 4 bars in Minneapolis, Mankato
  https://www.kare11.com/amp/article/news/health/coronavirus/covid-19-coronavirus-in-minnesota-minneapolis-st-paul-wisconsin-live-updates-june-26-2020/89-1085650a-869e-4bf1-9730-ac48c2ec43e8

Playstation 5

**Unigram**

- Get PS5 ready with Official PlayStation Magazine's 180-page PlayStation 5 special issue
  https://www.gamesradar.com/get-ps5-ready-with-official-playstation-magazines-180-page-playstation-5-special-issue-on-sale-now

- One Of PlayStation 4's Best Games Is Now Only $5 For A Limited Time
  https://wegotthiscovered.com/gaming/gamestop-selling-playstation-copy-limited-time

- Why Fans Are Excited About PlayStation 5's 'Spider-Man ...
  https://studybreaks.com/tvfilm/spider-man-miles-morales

**Bigram**

- Get PS5 ready with Official PlayStation Magazine's 180-page PlayStation 5 special issue
  https://www.gamesradar.com/get-ps5-ready-with-official-playstation-magazines-180-page-playstation-5-special-issue-on-sale-now

- Pricing for both PS5 consoles and every single accessory might've just leaked
  https://bgr.com/2020/06/26/ps5-release-date-price-dualsense-controller-accessories

- PlayStation 5 could revive Xbox One's canned Snap feature
  https://www.videogameschronicle.com/news/playstation-5-could-revive-xbox-ones-canned-snap-feature

# Experiment 2/2 : Results



As for the following results, we used both **Uni-gram** and **Bi-gram** as an **hybrid** solution for this experiment.

# Experiment 2/2 : Results [10]

News. Global. 2019

Year in Search 2019

1. Copa America
2. Notre Dame
3. ICC Cricket World Cup
4. Hurricane Dorian
5. Rugby World Cup
6. Sri Lanka
7. Area 51
8. India election results
9. 台風19号
10. Fall of Berlin Wall

Google Trends

- According to Google News Search in 2019, most of the user's common search contains two words.

[10]https://trends.google.com/trends/yis/2019/GLOBAL/

# Achievements

Published paper entitled <span style="color:red">"News Aggregator And Efficient Summarization System"</span> in *International Journal of Advanced Computer Science and Applications (IJACSA), Vol. 11, No. 6,2020* .

# Demo 1/2

# Demo 2/2

| Input keyword: Coronavirus | |
|---|---|
| Summary | Recommended Links |
|  |  |

# Users' Feedback

# Any Questions ?

# APPENDIX

# Tools used

# Future directions:

▸ In the near future, we aim that our system can deal with all types of media like photos or videos in summarization and recommendation phases by image processing techniques.
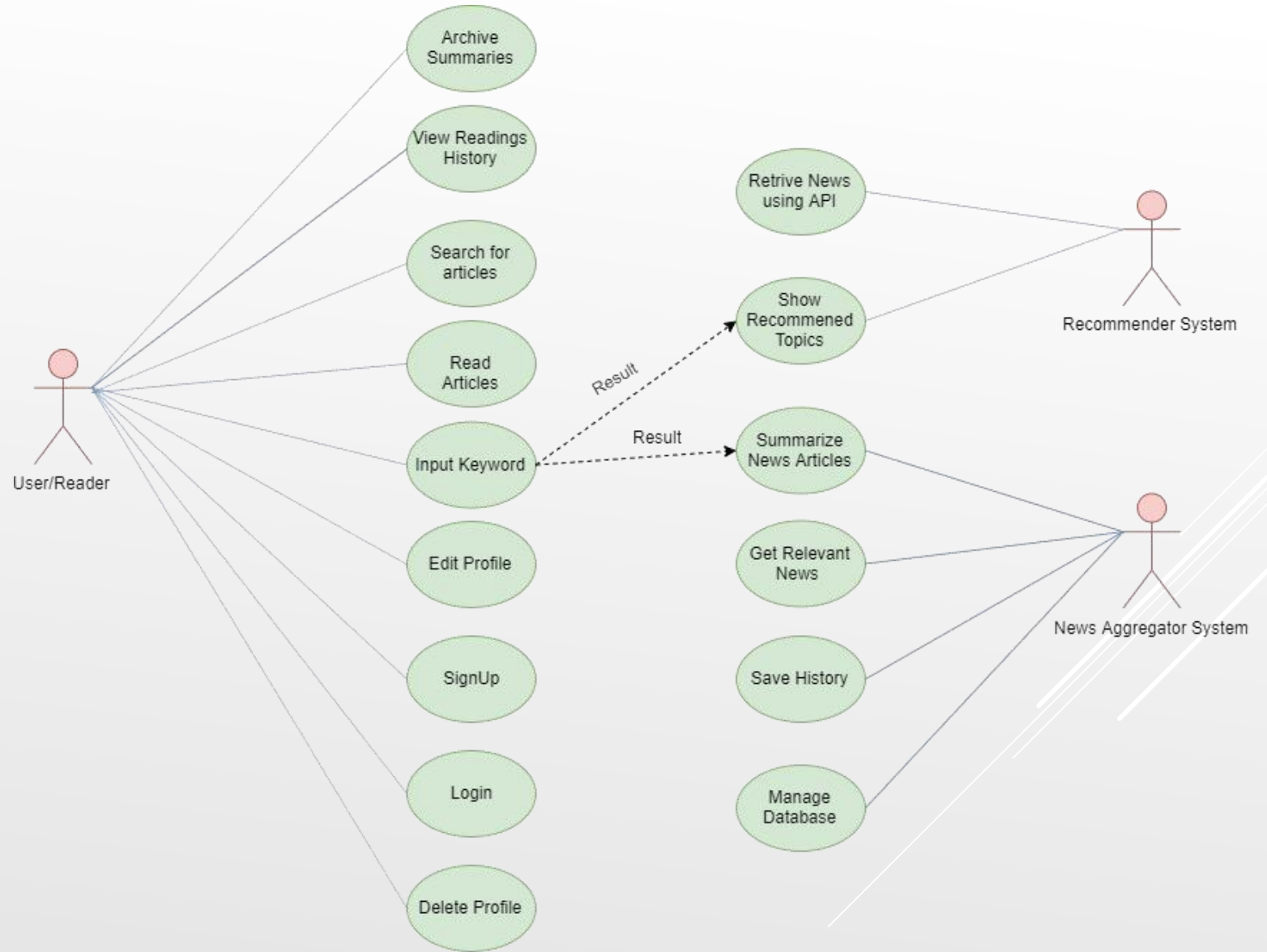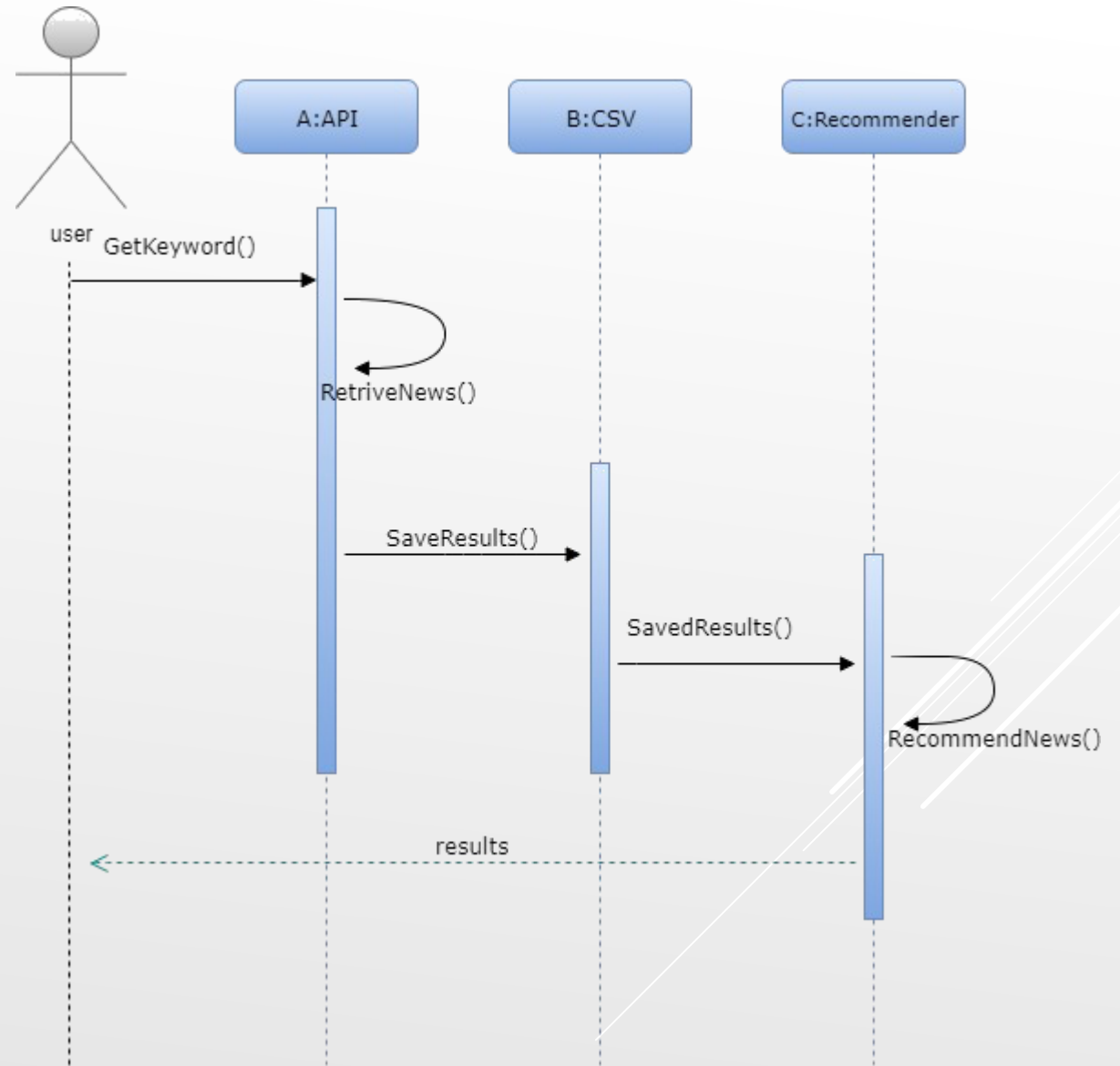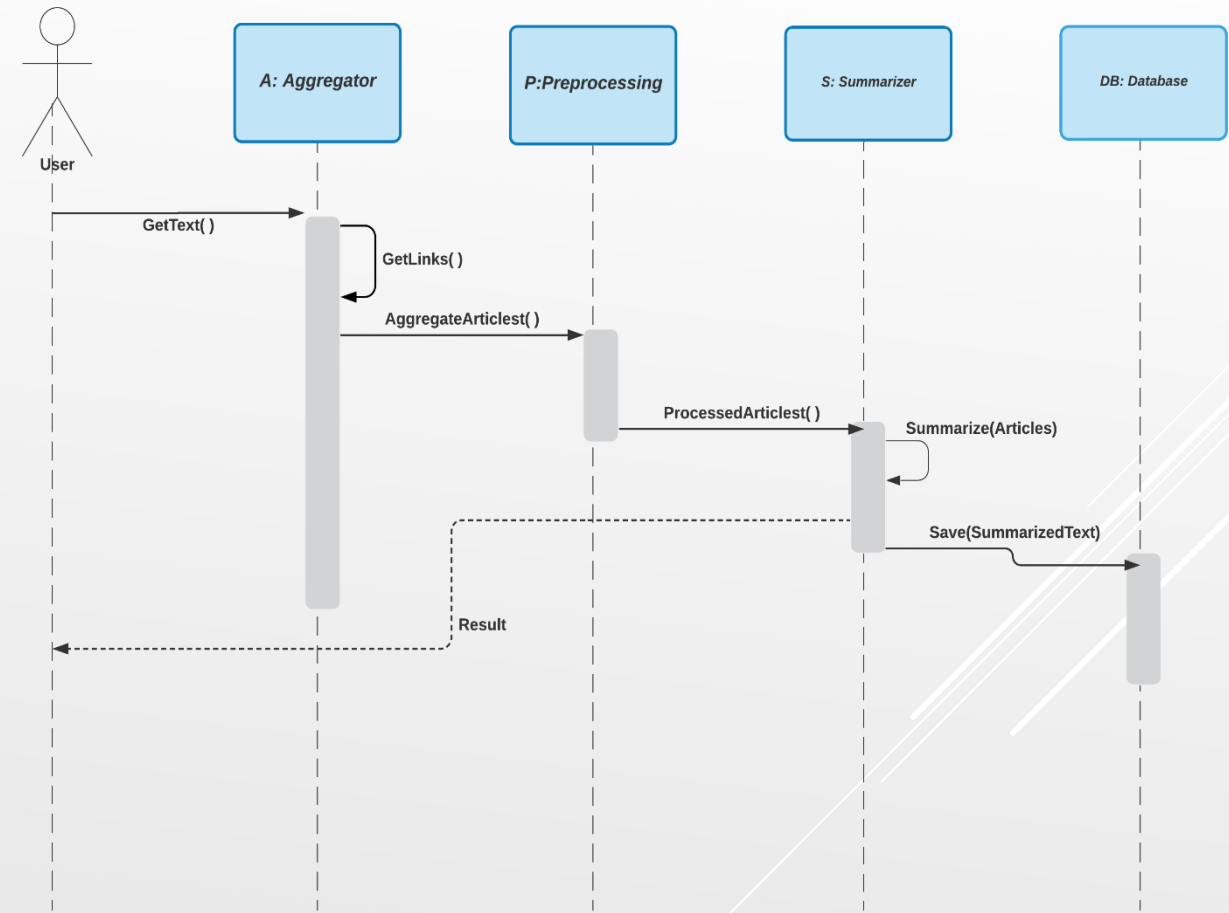
# Architecture Diagram

# Class Diagram
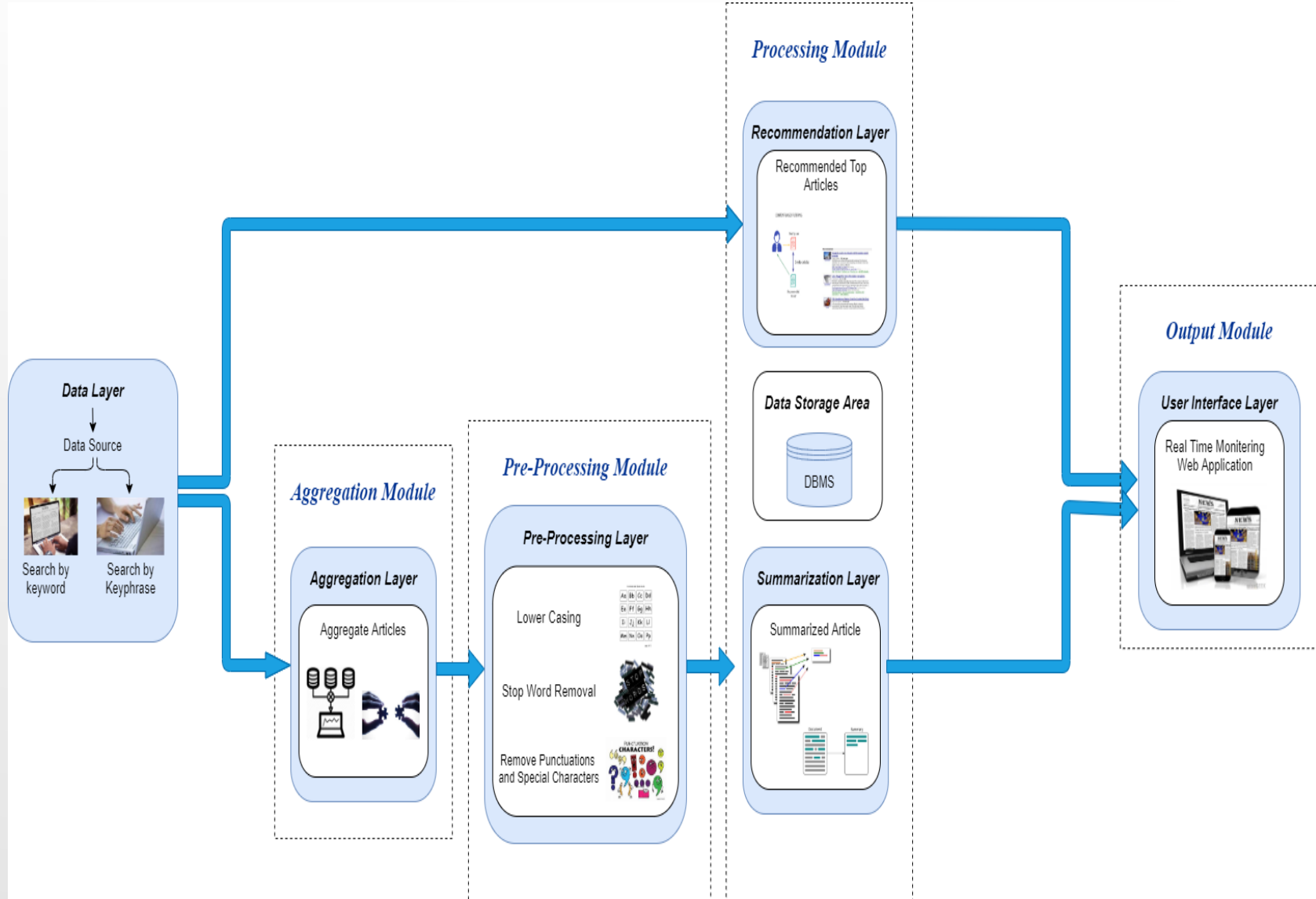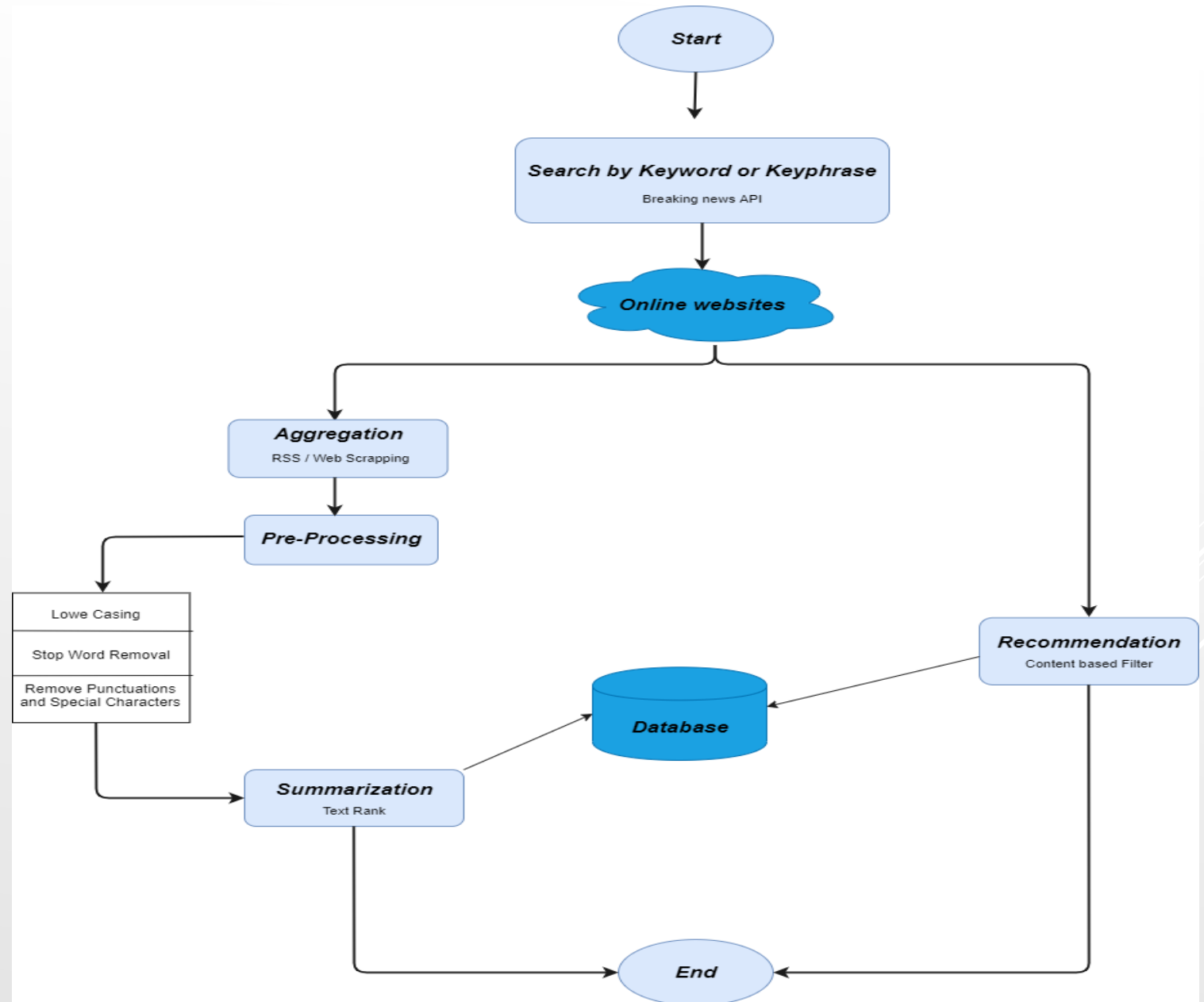
# Use-case Diagram

Sequence Diagram 1/2

# Sequence Diagram 2/2

# Block Diagram

# Activity Diagram

Context Diagram