UKRAINIAN CATHOLIC UNIVERSITY

BACHELOR THESIS

# Aggregation system of news search with configurable notifications by keywords

*Author:*
Arsen Ilchyniak

*Supervisor:*
Dmytro Pryimak

*A thesis submitted in fulfillment of the requirements*
*for the degree of Bachelor of Science*

*in the*

Department of Computer Sciences
Faculty of Applied Sciences

Lviv 2022

# Declaration of Authorship

I, Arsen Ilchyniak, declare that this thesis titled, "Aggregation system of news search with configurable notifications by keywords" and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.

- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.

- Where I have consulted the published work of others, this is always clearly attributed.

- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.

- I have acknowledged all main sources of help.

- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

_____

Date:

_____

*"quote"*

Author

UKRAINIAN CATHOLIC UNIVERSITY

Faculty of Applied Sciences

Bachelor of Science

**Aggregation system of news search with configurable notifications by keywords**

by Arsen Ilchyniak

# *Abstract*

This thesis focuses on analysing current news resources and optimization of accessibility. Other focus is performance of program, algorithms and tools connected to it. One of main things is implementation of solution with usage of Java programming language and spring framework. Data integrity is implemented with Apache Kafka. Two types of algorithm are proposed to implement main logic. As result, one of them is chosen and effective telegram bot with its implementation is created.

# *Acknowledgements*

# Contents

# List of Figures

# Chapter 1

# Introduction

## 1.1 Motivation

Last month I spent hours reading the news and becoming more dependent on them. I decided to watch through information connected to news resources. As a result, short research helped me discover that many users have the same problem and became a new idea for my work.

## 1.2 Reasons of resource improvement

News always was a popular theme. However, in our country, the number of people interested in the news has significantly increased over the last time. This has led to creating an uncountable number of news resources. However, in modern reality, not everyone can spend a significant amount of time reviewing each resource separately and looking through information that interests him. Web users have become much pickier than they used to be. That is why many people began to switch from the usual sources of information such as websites to alternatives, the central place among which took Telegram, including channels, chats and bots. However, users still spend much time looking through different sources to find the information that interests them.

## 1.3 Goals

- Analyze and collect information about the available media resources.

- Plan and architect the solution based on advantages gotten.

- Implement an alternative solution.

# Chapter 2

# Analysis of existing sources

## 2.1 Introduction

In this part different analytical tools are used to show trends and demand of different news resources.

## 2.2 Google trends

Google trends - https://trends.google.com/

A resource allows viewing and analyzing specific queries that users have typed in Google and other related systems. This tool will help see a reasonably accurate picture of changes in the popularity of specific topics and search trends for each day. The popularity of news in Ukraine has grown significantly since the war started, as indicated by statistics, see Fig. 2.1. Firstly, it was a panic growth, but now we can talk about a steady audience growth, which is currently three times larger than before the war.

Many people are surfing the internet searching for information during this period. They want to be confident in the information and choose a source for themselves or find a place where more compelling or exclusive news is. Another important factor is which queries are still the most popular. People are more likely to query specific news than watch news in general, as shown in Fig. 2.2. The general conclusion is that individual moments and details are much more interesting to people than giving them a broad picture. These results suggest that the information distributed in individual parts is more popular than when all the news follows one another. This information will be taken and compared with the results of further analysis.
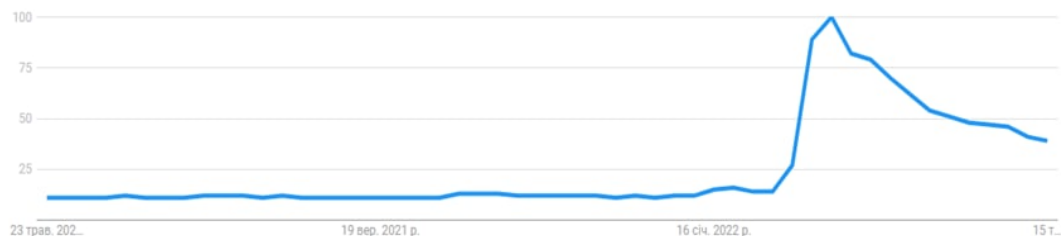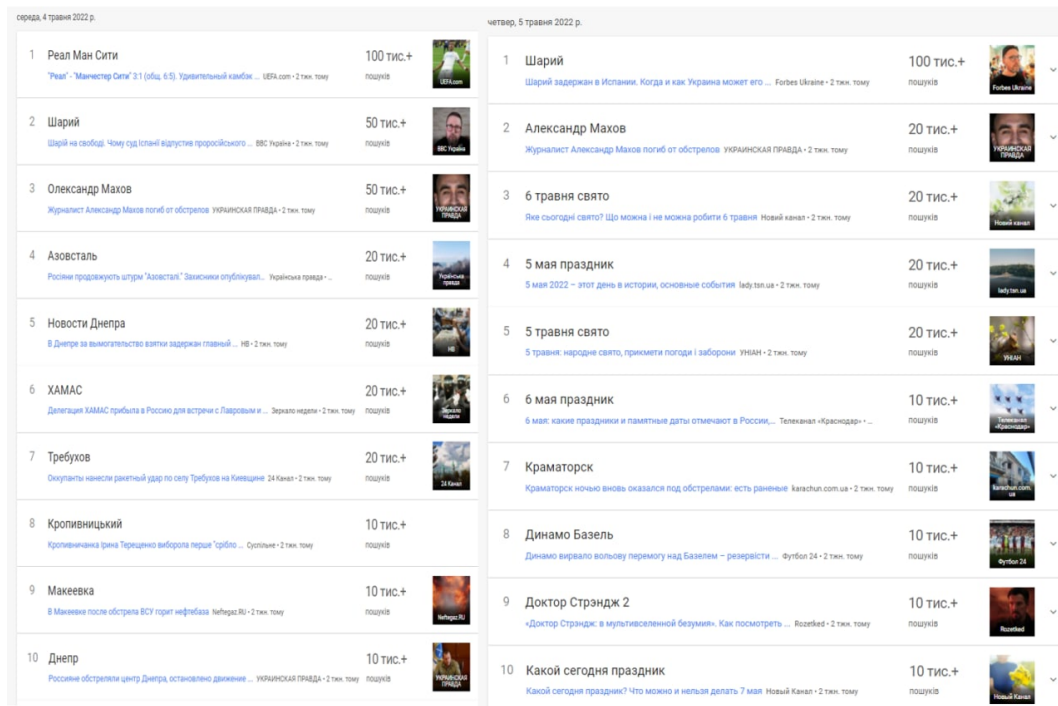


FIGURE 2.1: My architecture

FIGURE 2.2: Google analytic about most popular search topics

## 2.3 Similarweb

Similarweb - https://www.similarweb.com/

This part uses several parameters for evaluation, namely the number of visitors, average visit time, page attendance, bounce rate, features. Censor.net is used to be the most rated system among others. With only ninety million visits per month and most rated among other websites, it has several features that make it so favoured. First of all, it is an independent resource where the editor-in-chief is Yuriy Butusov, who is known for being on the front lines of war. Secondly, it can find texts by predefined keywords. Compared to others, it has better access to the previously mentioned feature. Lastly, it has pretty unique materials about the situation on the front line and is different from the official resources view. On the other side, it has a bigger than average bounce rate caused by people checking the website for new information and closing when they do not see anything new. The same thing can be said about all information resources.

TSN is the most visited website. It has two hundred millions visits per month. A pretty similar one has Pravda. Both resources have a small average page visit rate. At the same time, the main difference is the interface. Pravda has a simpler one, which causes a lower bounce rate, while TSN has the biggest one caused by site complexity. TSN got its audience from TV, partly because of its connection to the president in the past. One interesting feature that both of them have is the up bar of losses of the occupiers' army.

UKR.NET has the biggest average visit time and the number of pages visited, but it can be due to its usage as email. However, Similarweb still rates it as second among news resources. So, it is not only about email usage but about the news too. Furthermore, the primary thing that makes it popular is topic-based information. Usage of diverse information resources is a reason too.
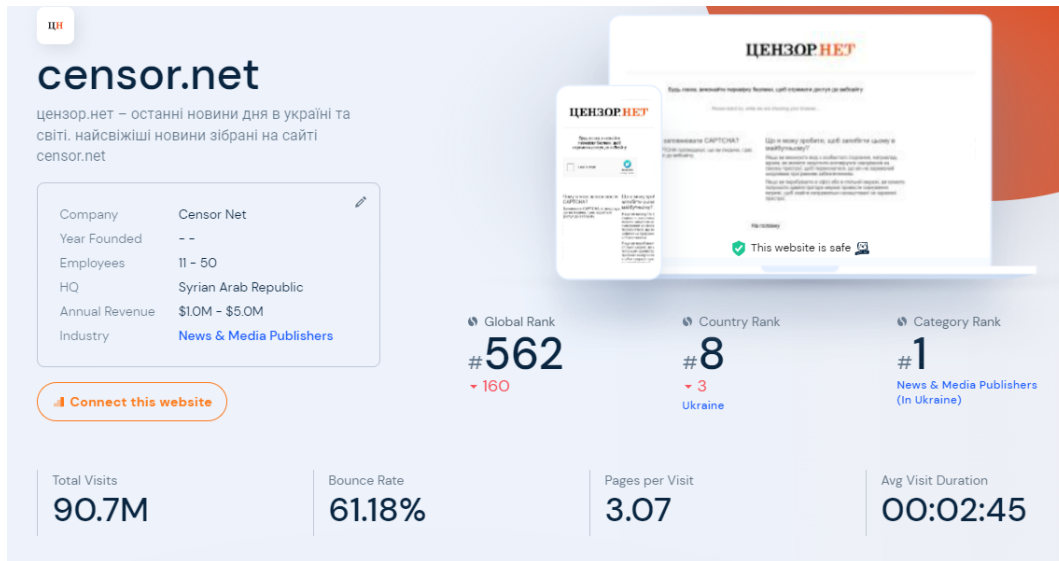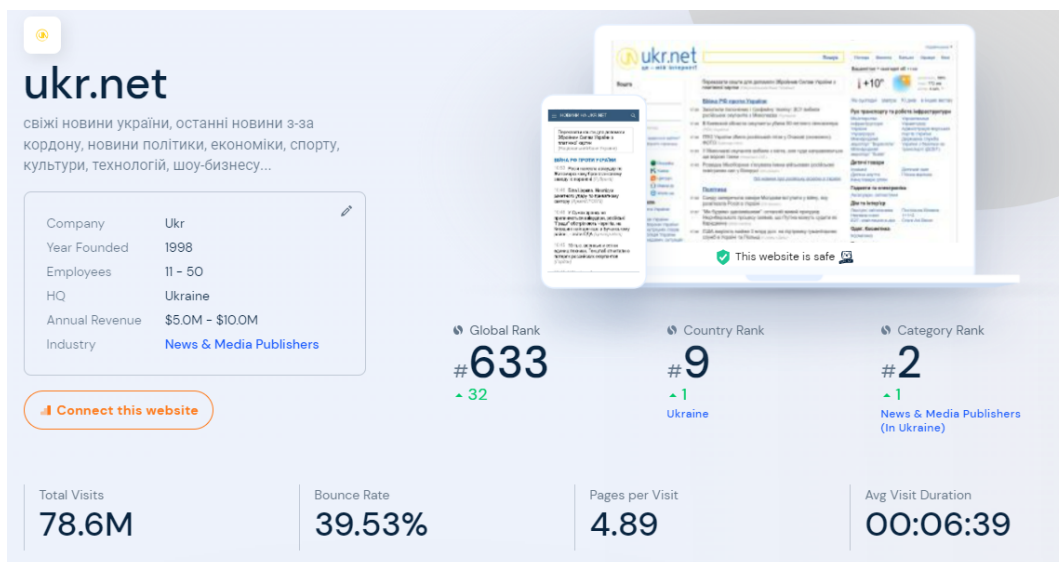
FIGURE 2.3: Censor.net statistics



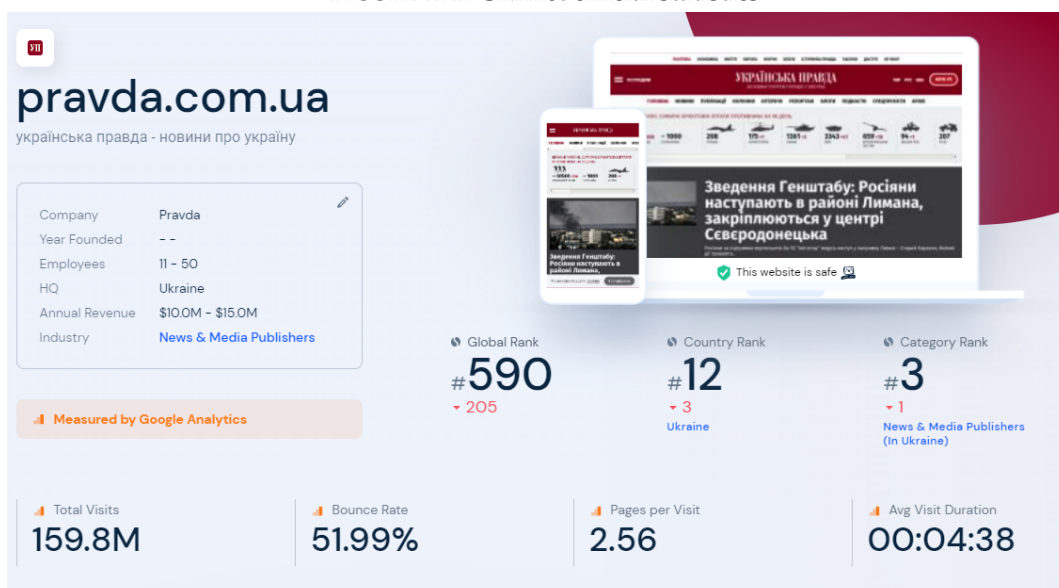FIGURE 2.4: Ukr.net official statistics
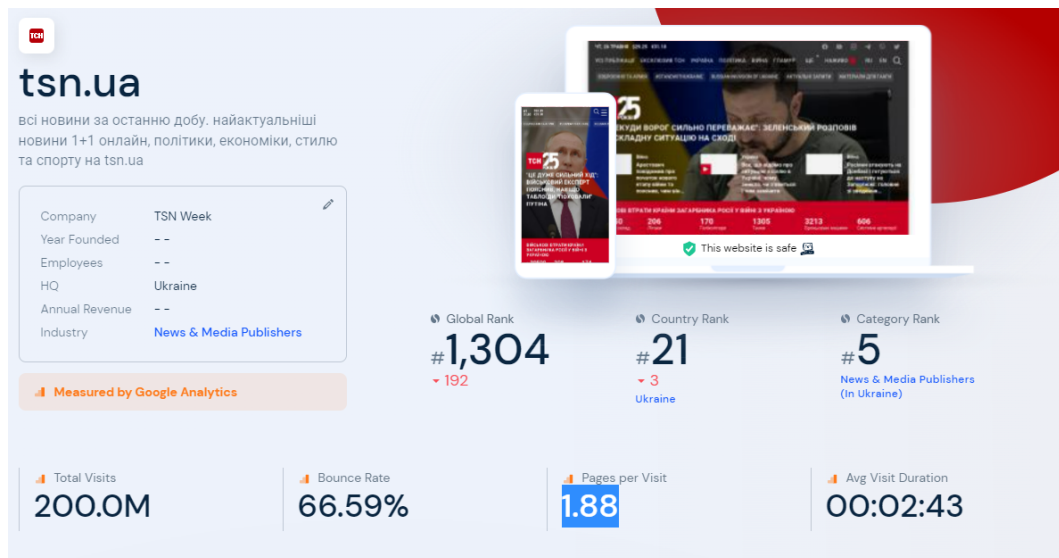


FIGURE 2.5: Pravda statistics

FIGURE 2.6: TSN statistics

Based on previous results, several conclusions about the efficiency and disadvantages of different resources were made. First of all, this is the number of visited pages. Web sites with separated topics have a lower bounce rate. Due to that, UKR.NET and Censor.net have higher results by this parameter. Moreover, Pravda has a bigger audience due to its straightforward interface. Despite the fact that TSN is the most visited one, its bounce rate and small average page visit rate tell us about its weaknesses and show that a significant part of the audience will leave soon.

## 2.4   tgstat

Tgstat - https://uk.tgstat.com/ratings/channels/news?sort=members

For subscribers' analysis and news popularity in Telegram, the work compares the two biggest channels in Ukraine with several official ones. U_Now provides information in quite an ordinary manner. All information is verified and comes almost already after official publication on the internet. On another side, Truexanewsua quite often publishes not verified information. However, this allows him to publish many exclusive materials at the exact moment after some event happens, making it popular among users.

This work takes @V_Zelenskiy_official as the biggest official information resource. Generally, the subscribers' number of official information resources is between 0.6 and 1.3 million when biggest unofficial channels have close to 2 million subscribers. Mostly it is connected to the ability of unofficial channels to public information from as many sources as they can process. However, the total number of views of official posts is often higher than others. It is connected to resend of official sources. For example, views of one post in @V_Zelenskiy_official is 1.15 higher than subscribers. At the exact moment, @Truexanewsua have only 40 per cent views compared to its subscribers. Despite this, the number of subscribers of @Truexanewsua continues to rise when @V_Zelenskiy_official decreases.

The other resource that is popular in telegram is @UASmartNewsBot. Opened information about these bot users is missing. However, due to its spread among different media when KitSoft launched it and its new way of accessing data, it is assumed that the audience is quite significant. The main advantages are different
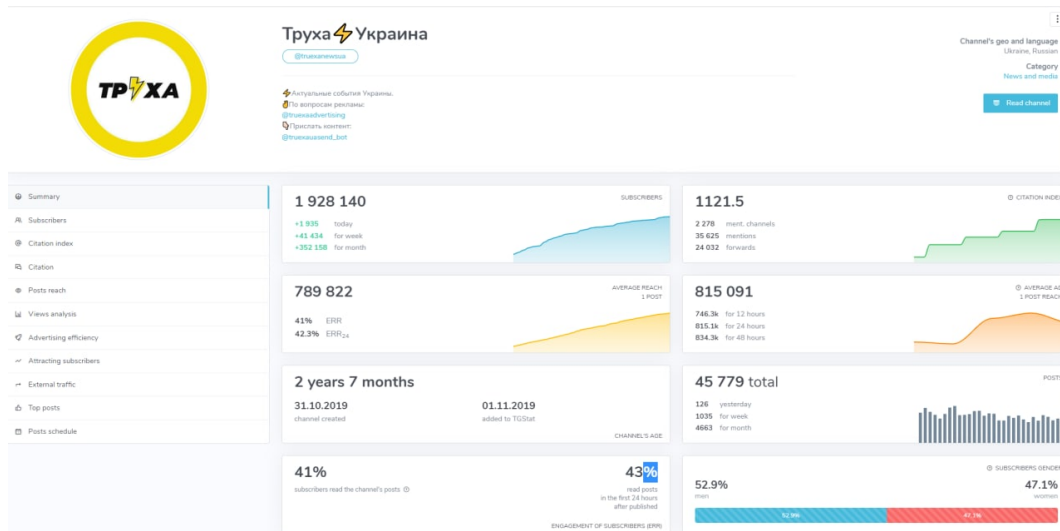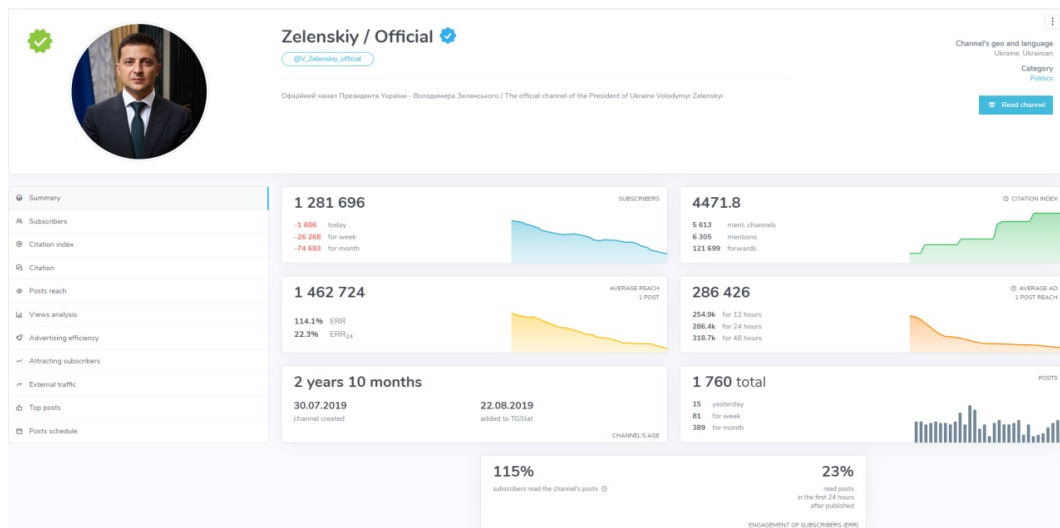
FIGURE 2.7: Truexanewsua statistics



FIGURE 2.8: Zelenskiy official statistics

sources of information collected in one place and that customers can write and get news by keywords that interest them. Because of the decision to take information only from official resources and not add new resources after launch, the channel became worthless resource for a significant part of the audience. Another big drawback is its base of 86 official channels that post similar information, and users get notifications about the same action but from many similar sources.

## 2.5 Conclusions

- People require structured information.

- People want to see information from different resources.

- Alternative resources have become more popular and continue to improve.

- Telegram is one of the significant resources.

- Information delivery has to be fast to interest subscribers.

- A variety of sources make information sound more confident.

- Exclusive materials make the resource more attractive.

# Chapter 3

# Architecture and implementation

## 3.1 Architecture planning

First of all, I have thought about the type of information delivery. Channels this time have too much information and are too generalized. I have chosen the "Telegram" bot, which has potential. However, as mentioned in the last part, "@UASmartNews-Bot" has not used this potential. Commonly, users have two primary needs. Firstly, subscribers want to get notifications when information that interests them is published. Secondly, they want to get information about some topic in the past.

As a result, this work needs to have a database to save information about two different types of information. The first one is subscribers' keywords needed to make information delivered user-specific. Another one is articles that have to be delivered to users interested in past events.

It is possible to make a new resource that will interest people, but this work is about improving news accessibility. Therefore, I have decided to scrap and process different popular resources.

Another mandatory part is the effectiveness of this work. The number of subscribers and keywords will increase with time, and naively making this work will lead to failure. According to this, I have tested several algorithms for keyword checking in texts. Last but not least, according to data consistency and modules that process the information, I have to add a message broker for the transmission user requests reliability.

There are a set of commands that users will use:

- Add keyword. To do this, send a keyword to the bot. Then the user is subscribed to the news that contains this word. Later, he has to wait for when new news containing his keywords appear

- Delete keyword. Works with /delete keyword.

- Get posts by keyword and dates. A user enters command /get_posts_by_keyword_and_dates. After that, he has to enter a keyword and two dates in the next message.

## 3.2 Used tools and frameworks

Java was chosen as the programing language.

Java is a programming language and computing platform first released by Sun Microsystems in 1995. It has evolved from humble beginnings to power a large share of today's digital world, by providing the reliable platform upon which many services and applications are built. New, innovative products and digital services designed for the future continue to rely on Java, as well.
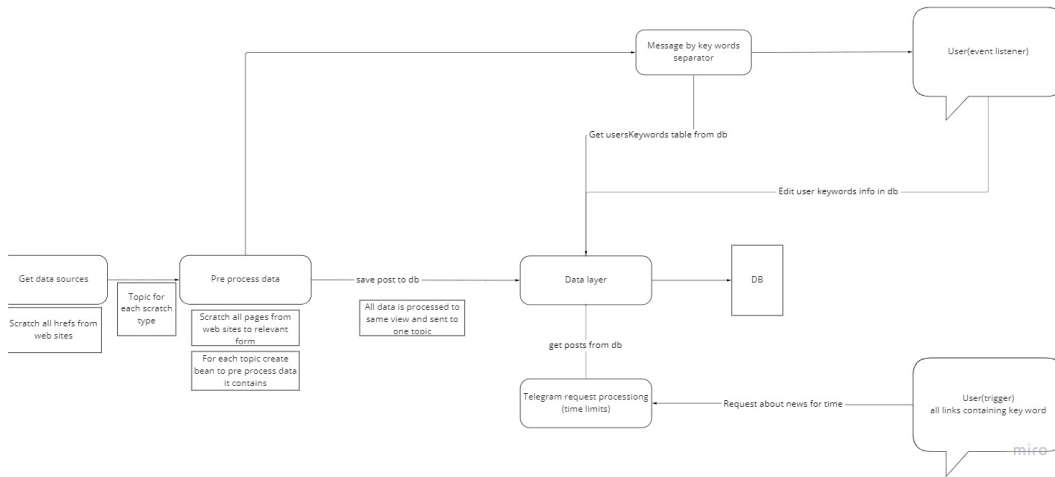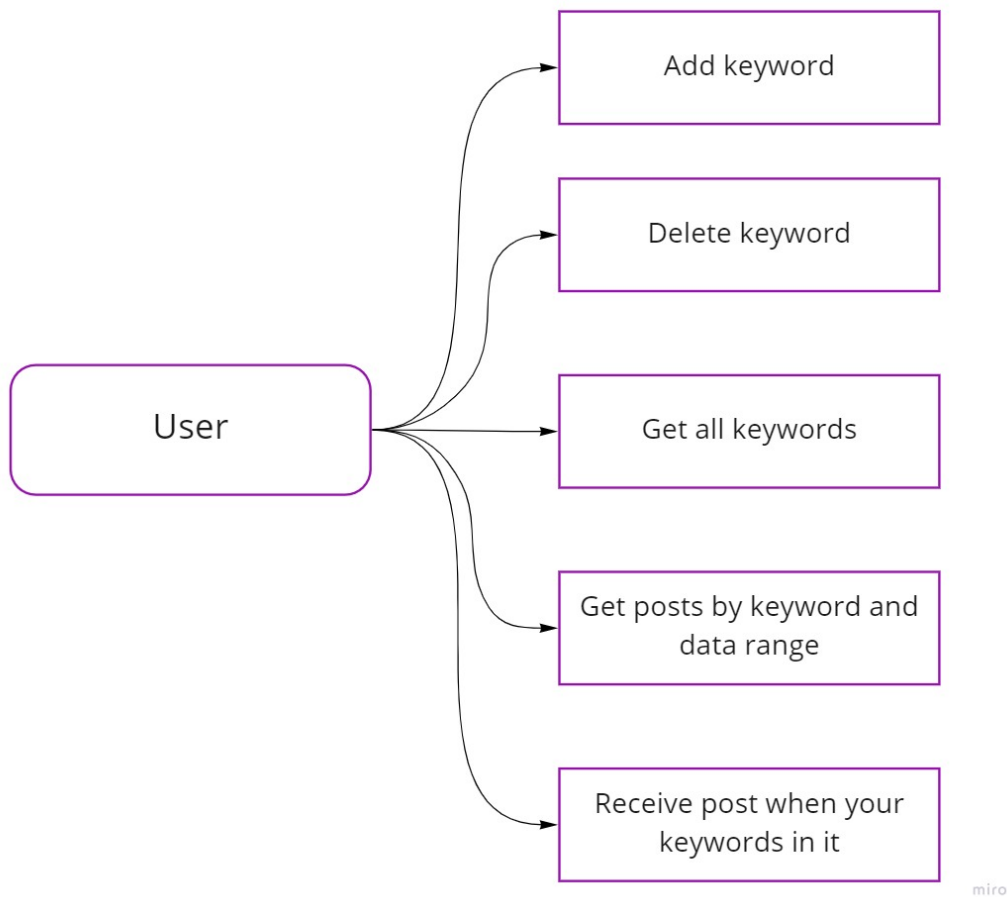
FIGURE 3.1: Architecture



FIGURE 3.2: User interactions

Its libraries are pretty helpful and easy to use.

It has many tools to simplify work and use different annotations to replace tones of code.

Kafka was chosen as a message broker. Combining it and Spring for Apache Kafka makes it the most suitable option.

Apache Kafka — Robust Queue Broker

" Kafka combines three key capabilities so you can implement your use cases for event streaming end-to-end with a single battle-tested solution:

- To publish (write) and subscribe to (read) streams of events, including continuous import/export of your data from other systems.

- To store streams of events durably and reliably for as long as you want.

- To process streams of events as they occur or retrospectively.

And all this functionality is provided in a distributed, highly scalable, elastic, fault-tolerant, and secure manner. Kafka can be deployed on bare-metal hardware, virtual machines, and containers, and on-premises as well as in the cloud. You can choose between self-managing your Kafka environments and using fully managed services offered by a variety of vendors. " Description of kafka at official page.

" The Spring for Apache Kafka (spring-kafka) project applies core Spring concepts to the development of Kafka-based messaging solutions. It provides a "template" as a high-level abstraction for sending messages. It also provides support for Message-driven POJOs with @KafkaListener annotations and a "listener container". These libraries promote the use of dependency injection and declarative. " Description of spring-kafka at official page.

I have chosen "PostgreSQL" as a database because it is standardized and easy to work with.

"PostgreSQL is a powerful, open source object-relational database system with over 30 years of active development that has earned it a strong reputation for reliability, feature robustness, and performance. ". Description of PostgreSQL at official page.

## 3.3   Algorithms testing

There were two main algorithms that this work tests.

The first one is - "Aho-Corasick algorithm".

" Let there be a set of strings with the total length M (sum of all lengths). The Aho-Corasick algorithm constructs a data structure similar to a trie with some additional links, and then constructs a finite state machine (automaton) in O(MK) time, where K is the size of the used alphabet.

It was considered that code would need to process big text files. However, practice shows that the average amount of chars in one file is less than 1000. Moreover, there is a variety of letters in the alphabet. As a result, this algorithm was the worst among others. ". Description of suffix-automation from cp-algorihms.

The second one is - "Suffix Automation".

" A suffix automaton is a powerful data structure that allows solving many string-related problems.
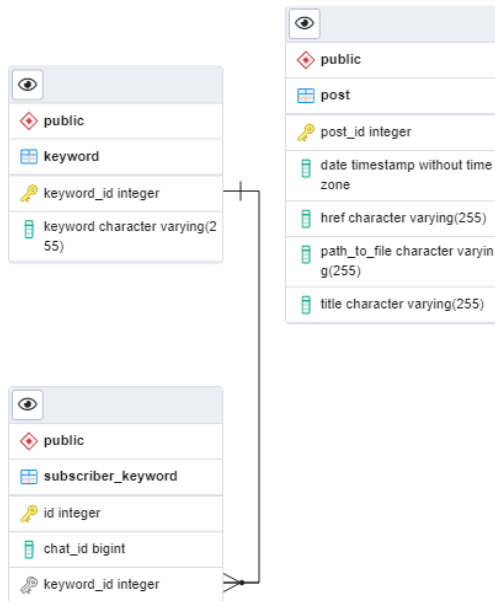
FIGURE 3.3: Database model

For example, you can search for all occurrences of one string in another, or count the amount of different substrings of a given string. Both tasks can be solved in linear time with the help of a suffix automaton.

Intuitively a suffix automaton can be understood as a compressed form of all substrings of a given string. An impressive fact is, that the suffix automaton contains all this information in a highly compressed form. For a string of length n it only requires O(n) memory. Moreover, it can also be built in O(n) time (if we consider the size k of the alphabet as a constant), otherwise both the memory and the time complexity will be O(n*log(k)).

The linearity of the size of the suffix automaton was first discovered in 1983 by Blumer et al., and in 1985 the first linear algorithms for the construction was presented by Crochemore and Blumer.

". Description of aho-corasick from cp-algorihms.

Naive method. String method - ".contains()" is used to check whether a keyword is in text. Its complexity is O(n*k), where n - text length, k - the number of keywords.

Comparing these two left algorithms, we can easily see that Suffix automation is faster with a larger number of keywords. However, with a small number of keywords, it deals better. This is because a naive algorithm does not have to prepare some structure to give results about keywords much faster. According to this I can say that naive algorithm will spend too much time and resources to execute the same piece of work. If to watch from perspective of small program with limited number of users it is quite satisfying. General results can be even better then Suffix Automation, but users would not see them. And when the number of subscribers will increase and number of keywords too, the program will need much more time to execute all queries that come. As result information will not be delivered fast and exclusive as it was planed for.

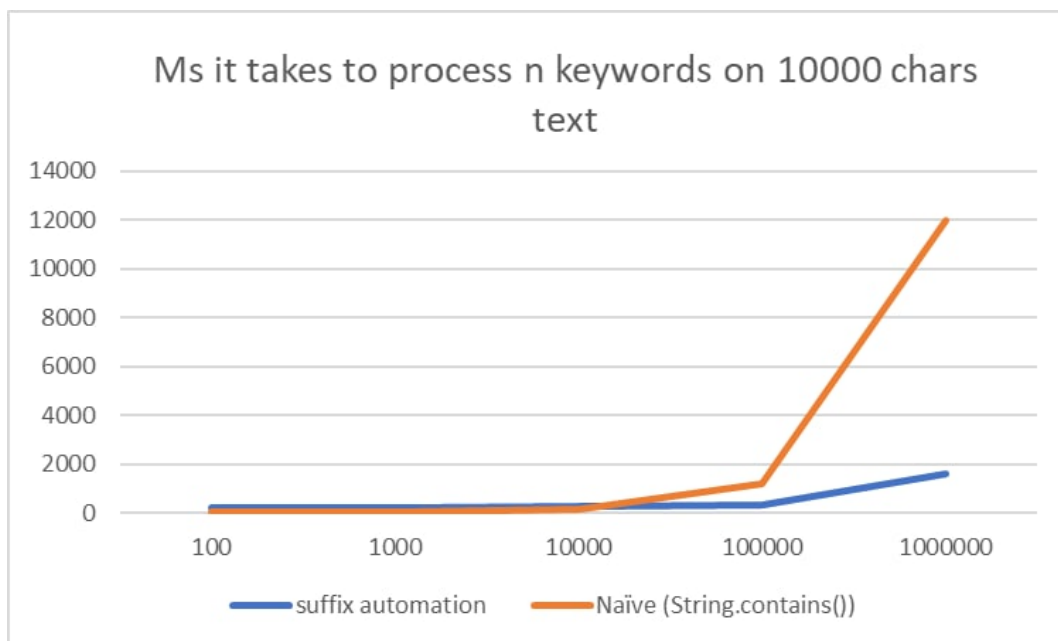As the result, Suffix automation were taken, because of its overall result and good complexity.

FIGURE 3.4: Milliseconds it takes to process n words on 10000 characters.

# Chapter 4

# Conclusions

As a result, I got a ready program that sends new news to subscribers by keywords they entered. It has a simple interface and can run on any machine globally with Telegram API. This work analyses algorithms and detects their weak and strong sides. Aho-Corasick, which I have used much before, turned out to be useless here because of its specification and data amount that incomes from news resources.

I have already partly launched it and gave bots link to several people to test it. As a result, several of my acquaintances started to browse news sources less if to believe their words. The most effect it had was on my acquaintance, who is sensitive to everything, and this bot helps her see only good news and avoid bad by entering specific keywords. Hope after deploy it will attract more people and I will show statistics of its usage during presentation.

# Chapter 5

# References

## 5.1   References

- Telegram api documentation - https://core.telegram.org/bots/api

- Best message brokers - https://blog.containerize.com/2021/07/09/top-5-open-source-message-queue-software-in-2021

- Aho-corasick algorithm - https://cp-algorithms.com/string/aho_corasick.htmlreferences

- Suffix-automation - https://cp-algorithms.com/string/suffix-automaton.html

- General Suffix Automaton Construction - https://static.googleusercontent.com/media/research.goo

- Kafka documentation - https://kafka.apache.org/documentation/

- Spring framework for Kafka documentation - https://docs.spring.io/spring-kafka/docs/current/reference/html/

- Spring data jpa documentation - https://docs.spring.io/spring-data/jpa/docs/current/reference/ht

- Lombok documentation - https://projectlombok.org/features/all

- Jsoup documentation - https://jsoup.org/

- Date formats - https://jenkov.com/tutorials/java-date-time/parsing-formatting-dates.html

## 5.2   Analytical websites

- Telegram statistics - https://uk.tgstat.com/ratings/channels/news?sort=members

- Similarweb - https://www.similarweb.com/

- Google trends - https://trends.google.com/

- Pravda advertise - https://www.pravda.com.ua/cdn/cd1/advertising/