

University POLITEHNICA of Bucharest

Faculty of Automatic Control and Computers,  
Computer Science and Engineering Department



# BACHELOR THESIS

## Clustering Conversations in Social Networks

**Scientific Adviser:**

Costin Chiru

**Author:**

Victor Andrei Oprea

Bucharest, 2015

Maecenas elementum venenatis dui, sit amet  
vehicula ipsum molestie vitae. Sed porttitor  
urna vel ipsum tincidunt venenatis. Aenean  
adipiscing porttitor nibh a ultricies. Curabitur  
vehicula semper lacus a rutrum.

Quisque ac feugiat libero. Fusce dui tortor,  
luctus a convallis sed, lacinia sed ligula.  
Integer arcu metus, lacinia vitae posuere ut,  
tempor ut ante.

# Abstract

Here goes the abstract about Streamer. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

# Contents

<b>Acknowledgements</b>	<b>i</b>
<b>Abstract</b>	<b>ii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Project Description . . . . .	1
1.1.1 Project Scope . . . . .	1
1.1.2 Project Objectives . . . . .	3
1.1.3 Related Work . . . . .	4
1.1.4 Demo listings . . . . .	4
1.1.5 Tables . . . . .	5
<b>2 State of the Art</b>	<b>6</b>
2.1 Background . . . . .	6
2.2 Existing Solutions . . . . .	7
2.2.1 Analytics . . . . .	7
2.2.2 3rd Party Clients . . . . .	7
2.2.3 Trends . . . . .	7
2.3 Related Work . . . . .	8
<b>3 Design</b>	<b>9</b>
3.1 Implementation . . . . .	10
3.1.1 Clustering Algorithm . . . . .	12
3.1.2 Data Visualization . . . . .	15
3.1.3 Deployment . . . . .	16
3.2 Testing and Evaluation . . . . .	16
<b>4 Conclusion</b>	<b>20</b>
<b>5 Further work</b>	<b>22</b>
<b>A Project Build System Makefiles</b>	<b>24</b>
A.1 Makefile.test . . . . .	24

# List of Figures

1.1	Reporting Framework . . . . .	4
3.1	Design for the data processing pipeline . . . . .	9
3.2	Design for queue processing in the pipeline . . . . .	11
3.3	Example of cluster rendering with outlier . . . . .	14
3.4	Example of cluster rendering . . . . .	15
4.1	Automated messages sent by bots . . . . .	21

# List of Tables

1.1	Generated reports - associated Makefile targets and scripts . . . . .	5
-----	---	---

# Chapter 1

## Introduction

Social media platforms are no longer an emerging field, they have become well established and millions of messages are exchanged daily. Services try to keep up with this trend by promoting popular content, either by number of clicks, views, favorites or other metrics. In this paper, we present a study on the clustering of messages from the Twitter platform, also known as "tweets". Inspired by the website's "Trends" which presents the most popular subjects either worldwide or in a region, the aim of our paper is to explore more of the popular subjects and cluster the conversations on more than just a keyword. The purpose of the clustering is to offer an in depth view of the conversation on a particular topic where peoples opinions may differ greatly.

The process is separated into three parts: data collection via the Twitter public API, message annotation with part-of-speech tags and message clustering. The messages provided by the website's API already provide a filtering option, you can specify keywords that you want to be part of the messages you get back. This might provide some indication of the conversation topic but getting an overview of the different conversations on the same subject is not a trivial task because messages have no obvious order.

By clustering the conversations together and offering an interface in which to explore the information, it is possible to view multiple points of view on the same subject and get more information than a popular topic might provide.

### 1.1 Project Description

In this paper we first discuss related work in the field. The popularity of the medium and the large quantity of messages have made this subject a popular topic of research. We then continue by giving a detailed explanation of the clustering process, the system architecture and the parsing and message aggregation. The next section will go into details regarding testing and the results reached and the last part will offer suggestions for further improvements and research.

#### 1.1.1 Project Scope

We will first offer a broader explanation of the problem we wish to solve.

Twitter<sup>1</sup> is an online micro blogging platform and social networking website. Users communicate through short 140 character messages called "tweets". To ease communication people use

---

<sup>1</sup><https://support.twitter.com/articles/49309-using-hashtags-on-twitter>

mentions; the @ character followed by a person's name. This is a way to involve another account into the conversation. Another feature of the conversation is the hashtag; the # sign followed by a word this is used to highlight key parts in the message, either a feeling or subject. Popular such hashtags are included in the Trending Topics. These are popular subjects automatically generated from conversations taking place worldwide or in a certain region. This is why people often times include hashtags in order to associate their message with a popular topic. Users can also favorite a tweet, and "retweet" it, meaning they share it with the people that follow them while still attributing the message to the original author. These two metrics contribute to the overall popularity of a tweet.

Twitter is an interesting platform for research because of the large number of messages exchanged daily. There are 500 million tweets sent daily by its 302 million monthly active users<sup>1</sup> with a record of 143,199 tweets per second<sup>2</sup>. People usually turn to Twitter during major natural events, sporting events, award ceremonies and so on. With such a high amount of information coming in every second it is almost impossible to keep track of everything that is discussed.

The "Trends" category highlights all popular topics and hashtags but exploring any one topic reveals a very large number of messages with just as more coming in every second (for popular subjects). This makes it very hard to see different points of view, different opinions on the matter. The aim of Streamer is to offer a high level overview of the conversation, where different opinions on the same topic are grouped into different clusters. This would make it easy to identify what kind of messages you are likely to find in a cluster just by reading some of them and would allow quick understanding all aspects of a developing story.

Consider the following scenario: during a global sporting event such as the football world cup which takes place over the course of several weeks, you want to keep track of the public opinion for your country's national team. If they are playing a game it's most likely that they are trending and a hashtag has already become popular between the people exchanging messages about the game. But how do you track the information? What if the hashtag includes both team (for example #GERvsBRA used to track Germany - Brazil game) what if you want to see information about the first goal, or a certain player in the game. Clustering all available messages would easily reveal the ones referring to one of the teams or an individual player. This is where Streamer wants to step in and change the way information is being discovered.

Ce nu ar trebui sa lipseasca: - cum ai gandit rezolvarea problemei - care este arhitectura aplicatiei - ce ai facut tu ca implementare - ce metoda de testare ai folosit - care sunt rezultatele obtinute

Chestiile astea ar trebui sa fie in jur de 30 de pagini. In plus, inainte ar trebui sa ai o introducere in care sa pui: - descrierea problemei 1-2 pagini - state-of-the-art (ce s-a mai facut similar) -3-4 pagini - scurta descriere Twitter (jumătate de pagina cred ca ar fi suficient) + - ce instrumente externe ai folosit (maxim 3-4 pagini)

The scope of the project **Streamer** is to provide close to realtime clustering of conversations that take place in the online medium. My choice for a social network is Twitter. Twitter has around 302 million active users (May 2015)<sup>3</sup> who send 500 million tweets each day mostly from their mobile phones. A tweet is a 140 character long message and because of this conversations are hard to keep track of and provide little to no context on their subject. This makes them an excellent candidate for a clustering application like **Streamer** which aims to provide an overview for conversations spanning over all the topics the user of the application provided.

This thesis presents the **Streamer**.

This is an example of a footnote<sup>4</sup>. You can see here a reference to [Section 1.1.2](#).

<sup>1</sup><https://about.twitter.com/company>

<sup>2</sup><https://blog.twitter.com/2013/new-tweets-per-second-record-and-how>

<sup>3</sup><https://about.twitter.com/company>

<sup>4</sup>[www.google.com](http://www.google.com)



Here we have defined the CS abbreviation. and the UPB abbreviation.

The main scope of this project is to qualify xLuna for use in critical systems.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

### 1.1.2 Project Objectives

With **Streamer** we want to provide close to realtime clustering of conversations that take place on the Twitter social social platform.

To achieve this objective there are several components that go into the architecture:

- Data collection — This step involves getting the public messages matching a particular query from Twitter. This is achieved via Twitter's streaming API which serves a percentage amount of all messages available. The system requires that all requests be done by authenticated clients. After registering with the service it automatically starts sending realtime messages that match your search terms.
- Message parsing — The tweets that have been acquired need to be parsed and meta information is added to them before they can enter the clustering algorithm. In this step each word is augmented with its corresponding part of speech tags are added to each of them and weights using TF-IDF.
- Message clustering — The new messages obtained as a result of the previous step are clustered using K-Means algorithm. The algorithm requires a weight function, for this we have chosen the cosine similarity function.
- Cluster API — The backend exposes the clusters of messages via an HTTP endpoint. Any number of clients can connect and receive the information. The endpoint serves the messages in JSON format, a popular solution of web applications.
- Frontend — A web application that connects to the endpoint and presents the raw JSON in a useful format. It draws the clusters, displays the messages and allows to explore the information.

We have now included [Figure 1.1](#).

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

We can also have citations like [?].

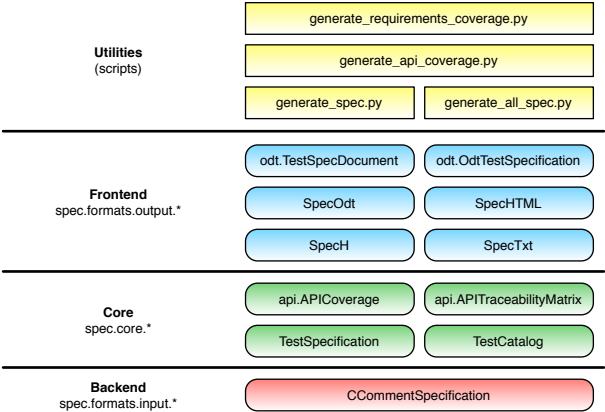


Figure 1.1: Reporting Framework

1.1.3 Related Work

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean aliquam lectus vel orci malesuada accumsan. Sed lacinia egestas tortor, eget tristique dolor congue sit amet. Curabitur ut nisl a nisi consequat mollis sit amet quis nisl. Vestibulum hendrerit velit at odio sodales pretium. Nam quis tortor sed ante varius sodales. Etiam lacus arcu, placerat sed laoreet a, facilisis sed nunc. Nam gravida fringilla ligula, eu congue lorem feugiat eu.

We are now discussing the **Ultimate answer to all knowledge**. This line is particularly important it also adds an index entry for *Ultimate answer to all knowledge*.

1.1.4 Demo listings

We can also include listings like the following:

```
CSRCS = app.c
SRC_DIR = ..
include $(SRC_DIR)/config/application.cfg
```

Listing 1.1: Application Makefile

Listings can also be referenced. References don't have to include chapter/table/figure numbers... so we can have hyperlinks [like this](#).

### 1.1.5 Tables

We can also have tables... like [Table 1.1](#).

Table 1.1: Generated reports - associated Makefile targets and scripts

Generated report	Makefile target	Script
Full Test Specification	full_spec	generate_all_spec.py
Test Report	test_report	generate_report.py
Requirements Coverage	requirements_coverage	generate_requirements_coverage.py
API Coverage	api_coverage	generate_api_coverage.py

## Chapter 2

# State of the Art

### 2.1 Background

The internet is offering a voice to millions of people that have access to it. Creating content which is accessible by everyone in the network is possible with little or no barriers. Anyone can report on events, share their thoughts and ideas and because there are no limits to this process the results are massive amounts of information.

Facebook<sup>1</sup> and Twitter<sup>2</sup> are two examples of Internet platforms that have made it increasingly easier for people to generate content in a multitude of different formats: from text and pictures to rich media such as audio and video.

People send on average 50 million tweets with a peak record of 6939 TPS record. Facebook has even bigger numbers, in an average of 20 minutes, 1 million links are shared, 10 million comments are posted and 1,6 million wall posts are made<sup>3</sup>.

**TODO:**

add number of users for each platform

. These are just two examples of popular social media websites and the rate at which content is being generated.

At the same time Facebook offers no way for a user to search through those comments and posts.

People communicate in short gists with the help of special annotations made possible on the platform. Mentions are a way of including another Twitter user into the conversation, they are formed by prepending the character "@" to the string that represents the user. Hashtags are a method of creating channels of communication, they are a way to distill the idea or felling of your tweet to a single word, and by doing so you ensure the inclusion of your message to a certain ongoing conversation. Hashtags are created by prepending the character "#" in front of words. Users can also reply to tweets, their message is grouped with the one they are replying to and context is preserved this way. Due to its short message format of 140 characters per message, Twitter has become a popular micro blogging platform for reporting on news and events as they occur. As a user you can always view the 10 most popular hashtags in different regions around the world or worldwide and participate in the conversation. You can also search for a certain query and messages that match gets returned either if it contains the string as a hashtag or in the message body.

The massive amount of information being generated especially on popular topics make it difficult to keep track of conversations as they happen. Unlike Facebook where the people you

---

<sup>1</sup>[www.facebook.com](http://www.facebook.com)

<sup>2</sup>[www.twitter.com](http://www.twitter.com)

<sup>3</sup><http://highscalability.com/blog/2010/12/31/facebook-in-20-minutes-27m-photos-102m-comments-46m-messages.html>

interact with are mostly people you know and that number can be within reasonable limits, on Twitter there are no barriers in communication and you have access to all the messages produced by every user of the platform. Even though hashtags help filter conversations they are too inclusive, there are no constraints over how to use them or how many to use so messages are included to any conversation as long as they have the hashtag. There are a number of services that use Twitter data and attempt to solve some of these problems. We will be presenting some of them in the following section.

## 2.2 Existing Solutions

### 2.2.1 Analytics

There are a number of analytics services that provide information regarding Twitter data. Most of them are businesses which offer information about the engagement of followers with the content created. Their goal is to help increase the visibility of tweets for businesses and therefor the metrics are related to the followers and focus less on exploring content.

This is also the solution offered by the Twitter analytics

**TODO:**  
add link

some of the information it provides is the number of user that views your messages, how many new users are now following your account.

One such example is SproutSocial

**TODO:**  
add link

which allow you to publish content from their application to Twitter, monitor your content for engagement and offer analytics on the users which interacted with your content.

Another example that tries to solve a similar problem to Streamer is TweetArchivist

**TODO:**  
add link

. You are able to query specific time frames and see top users and words related to certain search

terms, as well as the most shared URLs and the most influential

**TODO:**  
explain how they are influential

users that have send messages.

**TweetMotif**<sup>1</sup> retrieves tweets using the Twitter API based off of a user provided query. Then using n-grams extracts a certain number of topics and groups messages behind those topics, therefor giving an overview of what people are saying.

### 2.2.2 3rd Party Clients

There are a number of 3rd party clients. They allow for filtering of content based on a particular hashtag and popularity (this is rated by number of retweets and favorites). This is a good alternative for finding popular opinions, you can judge it by how popular that certain tweet is but it conveys either the voice of popular users which get lots of favorites and retweets or some tweets which happen to gain popularity by accident.

### 2.2.3 Trends

Twitter website offers access to world-wide trends and also custom trends. First off world trends represents a list of key words present in tweets in a certain region. This allows you to browse all the tweets with those key words in real time. You do not have any other type of control over the data. The data is not grouped by any other means so exploring it means going through each

<sup>1</sup><http://tweetmotif.com/about>

tweet and reading it and taking into account the volume of tweets some trends may produce (as presented in the introduction of this chapter) this task may be impossible.

## 2.3 Related Work

*Politics, Twitter, and information discovery* by **Moritz Sudhof**.<sup>1</sup>

The aim of the paper is to cluster Twitter users into groups based on the opinions they expressed regarding a political controversy. The corpus is fixed and contains tweets from the time the events occurred, they have been selected due to using the same hashtag specific to the event. Several different attempts are made at clustering the users using different methods. *Tf-idf* and *odds weighting* are used to extract relevant key words from messages. Multiple keywords shared between tweets are an indication of how similar they are and thus link the users together. *Mentions*, referencing one or more users in your tweet, are also used. Mentioning someone in your tweet means that they are relevant to your opinion or somehow involved. Finally *hashtags* are taken into consideration the idea behind it being that users who send out messages using the same hashtags share similar opinions, again the more hashtags users share the similar they must be.

*Topical Clustering of Tweets* by **Kevin Dela Rosa, Rushin Shah**.<sup>2</sup>

The scope of this paper is to classify Twitter messages into different categories. The authors consider hashtags an approximate indication of the message topic and use it to improve results. The topics categories in which the messages are sorted are predefined, and the corpus is fixed, composed of selected tweets that cover the predefined topics.

Before being able to cluster the messages they undergo an intermediate processing step. The processing step includes normalization in which several variations are experimented: tokenization, removal of rare terms, conversion to lowercase each in different combinations and results are tracked.

Both unsupervised and supervised methods are used. K-Means is used in combination with TF-IDF as a weight function for the unsupervised clustering. Rocchio classifier is used for the supervised clustering. Results are compared and it is noted that the supervised method has better results.

*TweetMotif: Exploratory Search and Topic Summarization for Twitter* by **Brendan O'Connor, Michel Krieger, David Ahn**.<sup>3</sup>

The paper presents the implementation details of TweetMotif. A platform that allows fetching tweets from Twitter Search API, generates 2-3 words topics from the newly formed corpus and associates messages to these topics.

The interface allows for a recursive drilldown into topics the goal being to offer a concise summary of the topics generated. Topic generation is achieved using n-grams with certain heuristics such as disregarding unigrams that are function words, or bigrams, trigrams that cross syntactic boundaries. Topics are merged and their sets of messages are combined. The user is presented with a limited number of topics to preserve cognitive load.

---

<sup>1</sup>[http://web.stanford.edu/group/journal/cgi-bin/wordpress/wp-content/uploads/2012/09/Sudhof\\_Eng\\_-2012.pdf](http://web.stanford.edu/group/journal/cgi-bin/wordpress/wp-content/uploads/2012/09/Sudhof_Eng_-2012.pdf)

<sup>2</sup><http://www.cs.cmu.edu/~kdelaros/sigir-swsm-2011.pdf>

<sup>3</sup>[http://brenocon.com/oconnor\\_krieger\\_ahn.icwsm2010.tweetmotif.pdf](http://brenocon.com/oconnor_krieger_ahn.icwsm2010.tweetmotif.pdf)

## Chapter 3

# Design

Streamer is composed of two parts:

1. A backend that communicates with the Twitter API, retrieves tweets that match the query provided by the user, performs parsing and clustering. The result is a list of messages annotated with the cluster id they belong to. This information is made available with the help of an HTTP server that exposes an API. We have chosen JSON as the format to export the data via the API.
2. A frontend that takes the JSON and renders clusters of tweets as well as provides an interface for the user to explore the conversations

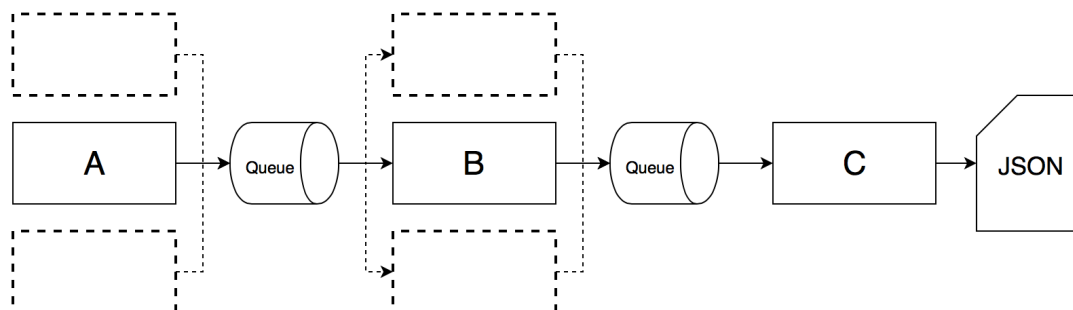


Figure 3.1: Design for the data processing pipeline

The data processing pipeline is composed of three main parts. Data acquisition (marked with A in the design) is responsible for fetching data that will later be processed, the component is agnostic of the data that passes through but in our case it is Twitter messages (tweets). It uses a library <sup>1</sup> to communicate with Twitter API, it retrieves the tweets and stores them in a queue. The queue library is provided by Apache Kafka <sup>2</sup>. As shown in the design, the reason for using queues between intermediary steps is that they allow for the producer and consumer to operate at different frequencies. The data acquisition segment could launch several clients regardless of what is happening further down the pipeline.

The second section (B), data processing, parses the raw tweets and converts them into shorter messages with key terms. Messages are read from the queue that is being filled by the data acquisition layer (A). The processed messages are written to a new queue, thus allowing this

<sup>1</sup><http://twitter4j.org/en/index.html>

<sup>2</sup><http://kafka.apache.org>

layer to scale just as the previous one. Tweets are parsed and using StanfordNLP library<sup>1</sup> each word is categorized with its own part-of-speech tag. Tags such as personal pronouns, possessive pronouns, prepositions, conjunctions are removed because they are too common. We could easily spin up several clients that consume messages because reading and writing is performed using two queues and thus the layer is decoupled from the other components of the pipeline. One issue related to parsing conversations especially ones from social media is the jargon used and possible spelling errors. This issue is exacerbated by the fact that Twitter conversations have such hard limits on the number of characters, on average a message does not have more than 10-12 words.

The last part that processes data is responsible for clustering the messages based on the keywords generated in the previous step. The clustering algorithm uses K-Means with cosine similarity as a distance function, which I will go into more detail in the following section.

The visualization (Streamer-Frontend) is rendered in the browser. This offers the advantage of being able to explore the profile of Twitter users and see the messages and their context. It works by polling the webserver that is serving a JSON file through its API. When new data becomes available it renders the clusters and the corresponding adjacent nodes. The polling process will continue in the background. From the user interface you are able to see the clusters and quickly identify large clusters. You are able to see all the clusters that belong to it either by hovering over the nodes or clicking the cluster and getting an expanded view with all messages.

## 3.1 Implementation

In the following paragraphs I will go into further details on how the system is implemented. The implementation is done in Scala<sup>2</sup>. The reason behind this choice is the fact that Scala enables us to use a functional programming paradigm and at the same time provides a type system that makes implementation easier. Many of the operations in the application include transformations of lists of messages, something that functional programming is very good at. At the same time the Scala source code is intended to be compiled into Java bytecode and run on the JVM. This allows us to include any Java library into the project as Scala was designed with Java interoperability in mind.

Another benefit of Scala is their implementation of parallel processing into the standard library. The aim of the language designers was achieving a high level abstraction that is easy to use thus achieving efficient parallel computations over collections in a transparent manner.

---

```
// Example of using parallel collections in Scala

list.map(_ + 42) // regular, sequential map over a collection
list.par.map(_ + 42) // parallel processing of the collection
```

---

Listing 3.1: Example of parallel collection usage in Scala

Using a similar approach we can speed up message parsing and also the clustering step. This decision has had beneficial results for the overall processing time of the Twitter messages. We will go into further details about the running time and speed benefits of parallel processing in the pipeline in the Testing and Evaluation section.

Accessing the Twitter API requires a developer account and an application created on their website dev.twitter.com<sup>3</sup>. The application gives you access to public, user or site streams. We will be using the public streams which returns data flowing through Twitter in real time given

---

<sup>1</sup><http://nlp.stanford.edu/software/index.shtml>

<sup>2</sup><http://www.scala-lang.org>

<sup>3</sup><https://dev.twitter.com/streaming/overview>



a certain array of keywords. This is the most useful endpoint for our data mining usecase. Using the provided API authentication tokens you can configure twitter4j library to retrieve tweets that match a specific keyword (or multiple keywords). The library comes with OAuth support and handles authentication with the endpoint. All messages received are passed to queue.

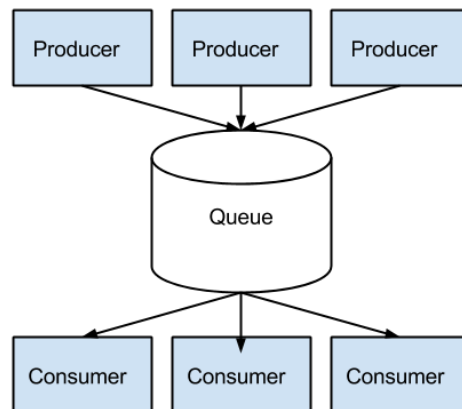


Figure 3.2: Design for queue processing in the pipeline

The messages retrieved are stored in a queue provided by Apache Kafka<sup>1</sup>. The queue, has a configurable storage period for the messages which allows us to use it for temporary storage. Communication between producers, consumers and the queue is done via a binary protocol over TCP so message delivery is always assured. The protocol does not require a handshake step in order to get or put messages in the queue, after a socket connection is established the client simply writes the messages it requires as request and then reads them. Kafka also guarantees message order and load balancing over a group of consumers but these features are not of interest for our project. Using a queue provides an advantage over the fact that messages arrive and are consumed at different frequencies. Depending on the popularity of the keywords specified tweets can come in at different rates. Based on a series of tests the rate of messages has been up to 150 per minute. At the same time the StanfordNLP has a parsing speed of one tweet per second. Using a queue also means that it is possible to start several consumers and producers at the same time. A consumer is concerned with getting the messages out of the Apache Kafka queue and placing them in a new queue where they will eventually be parsed. The Apache Kafka documentation lists "Stream Processing" as an ideal use case for the library, and the processing of the Twitter API feed is exactly this sort of use case, thus confirming our design decision.

Before applying the clustering algorithm the tweets are first parsed. Parsing the tweets means

<sup>1</sup><http://kafka.apache.org/documentation.html#introduction>

removing all non alphanumerical characters or punctuation: such as unicode characters. One of the reasons for removing non-alphanumerical characters is that the StanfordNLP library cannot parse emojis <sup>1</sup>.

The new messages are annotated using StanfordNLP part-of-speech tagger library. This library reads the text and assigns a part of speech or other token to each word. The set of part of speech tags follows the Penn Treebank Project <sup>2</sup>. The resulting output is a list of tuples containing of the word, its POS tag and its level in the dependency graph. The dependency tree is constructed by drawing an edge between a token and the all others that it determines. The tagger has an accuracy of 97.24

The next step is to filter out words based on the part-of-speech tagging. Tags such as personal pronouns, possessive pronouns, prepositions, conjunctions are removed because they are too common and might introduce false positives for the clustering algorithm.

These parsed messages are pushed to a new queue. The reason for this is to completely decouple the 3 different stages of the pipeline:

1. Retrieving messages in realtime from Twitter using its API.
2. Tagging the messages with their part-of-speech tag and using this to filter out common terms.
3. Running the clustering algorithm using the parsed messages as input.

The final stage of the pipeline takes all the parsed messages from the queue and applies the clustering algorithm. The algorithm used for this step is K-Means clustering and the distance function is the cosine similarity. Because the design of the K-Means clustering algorithm is to start the iterations with a random seed consisting of randomly selected points (in our case messages), we run the clustering using different starting seeds and choose the highest value out of all the runs as our best clustering option. This in turn will affect the overall running time of the application as we will see. Depending on our goals either speed or precision we can choose to just rely on the first random seed provided. The name K-Means, first introduced in **Some methods for classification and analysis of multivariate observations** by J. MacQueen. K-Means is an important clustering algorithm; its objective is to minimize the average distance between a point and its cluster center where a cluster center is the mean of all points in that cluster.

### 3.1.1 Clustering Algorithm

The first step of the algorithm is to compute the TF (term frequency) and IDF (inverse document frequency) for each of the tweets. TF-IDF is a numerical statistic intended to provide information about the importance of a word in a document. TF is a direct measure for the number of times a word appears in a document while IDF helps scale down that measure for words that tend to appear often and scale up for words that are rare.

$$\text{TF}(D, t) = \frac{\sum \text{occurrences of } t \text{ in } D}{\sum \text{number of terms in } D} \quad (3.1)$$

The TF of a term  $t$  in a document  $D$  is the total number of occurrences of that term divided by the total number of terms in the document. A document in this case would constitute a tweet.

<sup>1</sup>An emoji is a pictogram, similar to ASCII emoticons, which have been incorporated into Unicode meaning a wide adoption in online communication. Emoji symbols are two-byte sequences and support for them in browsers and mobile devices varies.

<sup>2</sup>[https://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn\\_treebank\\_pos.html](https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html)

And the terms we are computing TF for are all the terms that remain after the filtering in the previous stage.

$$\text{IDF}(D, t) = \log \left( \frac{\sum \text{total number of terms in } D}{\sum \text{total number of occurrences of } t} \right) \quad (3.2)$$

The IDF of a term  $t$  in a document  $D$  is the total number of terms in  $D$  divided by the total number of occurrences of  $t$  in  $D$ . The document in this case constitutes all of the available tweets that we are attempting to cluster. The number of occurrences is also computed taking into account all the available messages.

Using TF and IDF we transform each tweet into a vector with weights for each of the terms. This way we can use the cosine similarity as a distance function in the clustering algorithm. Cosine similarity is a way of computing the similarity of two vectors by measuring the cosine of the angle between them. The cosine function output is in the range  $[-1, 1]$  where -1 means that the vectors have opposite directions and 1 means that they have the same direction.

$$\text{cosine similarity}(A, B) = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}} \quad (3.3)$$

The  $A$  and  $B$  vectors in our equation constitute two vectors i.e. two tweets that have been turned into vectors after computing their TF and IDF. And the result of the *cosine similarity* function tells us how similar the two vectors are, in return how similar two tweets are.

The clustering algorithm uses the K-Means clustering algorithm. To define some terms that are used in the implementation of the algorithm, a cluster center is the mean of all points in that cluster, or more formally

$$\mu(w) = \frac{1}{|w|} \sum_{x \in w} x \quad (3.4)$$

Where  $w$  is the cluster. Ideally the clusters that are formed do not overlap. To get a sense of how well a centroid, the mean point of a cluster, represents all of its members, we define RSS *residual sum of squares*

$$\text{RSS} = \sum_{x \in w_k} |x - \mu(w)|^2 \quad (3.5)$$

The objective of the algorithm, and therefor our objective is minimizing this RSS function. By reducing this value we reduce the average squared distance which is a measure of how well points in a cluster are represented by their centroid.

The algorithm starts with an initial set of clusters chosen at random from all the tweets. It then uses the distance function (the cosine similarity described above) to assign the rest of the tweets to those clusters. For each cluster now formed it computes a weighted average. These steps are repeated only now we use the weighted average centroids instead of the random points until the centroids remain the same between 2 consecutive steps or until an upper bound of steps is reached. This upper bound is added in to ensure the algorithm halts if it fails to reach convergence in a fixed number of steps.

There are a number of termination conditions we can apply:

1. Terminate when RSS falls below a certain  $\epsilon$
2. Centroids do not change between iterations
3. Assignments of points to clusters do not change over iterations

4. A fixed number of iterations has completed. This will always ensure the termination of our algorithm but it might also affect the quality of our results.

Several issues we might run into while running the algorithm are data outliers. An outlier is a data point on a graph or in a set of results that is very much bigger or smaller than the next nearest data point. Therefore such points do not fit well in any cluster being too far away from the rest of the data. One such example would be having a cluster with only one message.



Figure 3.3: Example of cluster rendering with outlier

Another such example of clustering done sub optimally is having a cluster with no actual messages in it. This is due to the methods used to select the initial seeds for the clustering algorithm. Selecting a "bad" seed, by which we mean an outlier will always lead to sub optimal results. Heuristics used to improve results are using the RSS function to choose the seeds. By using this method we would go through different seed candidates and choose one. Another possible method is simply removing outliers from our sample data.

Different assignments yield different clusters which in turn have different weights. A higher weight means better similarity between the node and the cluster they belong to and therefore a better result. This similarity is computed with the help of the *cosine similarity* function after the K-Means algorithm has finished its assignment. For better results several iterations of the algorithm are used and the best score is chosen, at the cost of time spent clustering.

Another problem when using K-Means is choosing the correct  $K$  value, number of clusters. Most papers suggest having domain knowledge over the data that is being clustered but in our case it is not entirely possible due to the fact that data is coming in real time. Using RSS function and attempting to select a  $K$  value that minimizes it is a naive approach because RSS will reach its minimum for  $K = N$  meaning one cluster for every tweet in the dataset.

The result of the clustering algorithm as well as the tweet message and tweet author are combined and converted to a JSON data structure. This is in turn written to disk to be consumed by the endpoint accessing the server.. Using *Finagle*<sup>1</sup> an open source library from Twitter that provides an HTTP server implementation we can create the API endpoint.

Having the data accessible via HTTP it can easily be consumed by any application regardless of the programming language used, also the reason we have chosen JSON as a format is because it is lightweight and most languages have an implementation of a JSON parser. The connection to the API endpoint is not kept alive, all clients are expected to poll the endpoint and react to changes. As an optimization the server could reply with status code 304 Not Modified<sup>2</sup> this would not include a message body and would let the client know that no new data is currently available. This would reduce the payload that has to be transmitted. On the client side this can be improved by using an exponential backoff approach for long polling. The polling period

<sup>1</sup><https://twitter.github.io/finagle/>

<sup>2</sup><http://www.w3.org/Protocols/rfc2616/rfc2616-sec10.html>

would always double when a 304 status code is presented and it would reset to a lower default value when the endpoint has new data available.

### 3.1.2 Data Visualization

Visualizing the data is made possible in the browser using JavaScript. We chose the web as a platform because it made more sense for a number of reasons:

1. It is easily accessible from a number of different devices such as laptops, phones and tablets much like the content we are clustering.
2. There are a multitude of libraries that implement visualizations, graphs, plots for JavaScript. Transforming the JSON file exposed by the API in a graphical representation is faster this way.
3. Using URLs the tweets can easily be traced back, the user profile and be viewed and content can be explored without having to implement it.

The visualization for Streamer named **Streamer-Frontend**<sup>1</sup> was built using two libraries.

1. React<sup>2</sup> an Open Source JavaScript library from Facebook that handles updating the DOM (Document Object Model) every time new data comes in.
2. D3.js<sup>3</sup> which stands for Data Driven Documents, this is an Open Source library used for plotting, graphs and other types of data transformations.

200 tweets

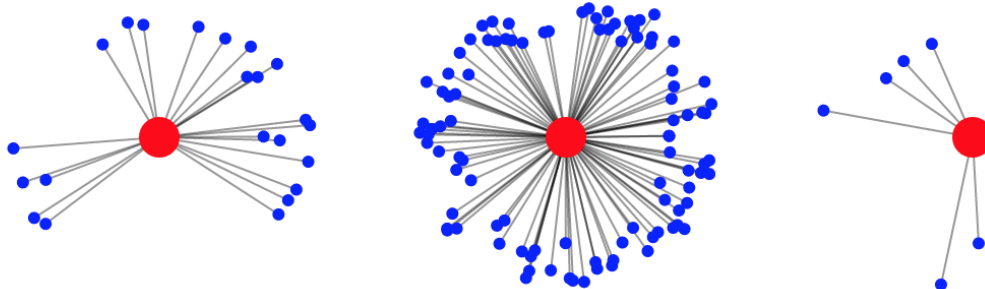


Figure 3.4: Example of cluster rendering

Streamer-Frontend polls the server API at a fixed time intervals and gets the latest version of the tweets as well as their respective cluster. Using an appropriate data structure we group the tweets together based on cluster id and draw the clusters to the webpage. Drawing is done using an SVG element with the help of the D3 library.

The interface presents the user a more meaningful representation of the data: tweets that belong to the same cluster are shown drawn together, also it is possible to explore the cluster and see all the messages that compose it.

The clusters will automatically update when new data comes in without the need to refresh the webpage. The application polls the API and when new messages become available it will redraw the clusters. A small optimization to the polling requests has been added: the polling

---

<sup>1</sup>add link to Github

<sup>2</sup><https://facebook.github.io/react/>

<sup>3</sup><http://d3js.org>

interval doubles every time the request does not return with new data. This prevents the server from receiving too many requests if multiple Streamer-Frontend clients are running and also works similarly to the data processing step which will always take longer as more tweets are fed into the pipeline.

An explanation of Figure 3.3: We have tweets (blue dots) and centroids (the red dots). The centroid is only used as a visual cue, making it obvious which tweets belong together. Hovering over any of the blue dots to view the content of the tweet in an overlay that will appear or clicking on the red dot reveals a side menu on the right side with all the tweets that make up the cluster as well as a total.

### 3.1.3 Deployment

For automated deployment of the project I used Docker<sup>1</sup>, it is an Open Source project that runs applications inside of software containers. A container is similar to a virtual machine but avoids the overhead of it by sharing the same kernel as the host. Dockerfiles are configuration files for Docker containers, they are scripts that describe the steps necessary to create the environment for the application. A Dockerfile for Streameris also available as a gist<sup>2</sup>. Some of the requirements of the project that the container automatically configures:

1. Scala 2.10
2. sbt 0.13.1 (Scala built tool) - Allows for task automation, running tasks and access to the Maven Central Repository for installing and updating packages
3. Kafka 2.9.1
4. Java 8
5. These are on top of Ubuntu 12.04

The script installs the required software with appropriate versions, updates all packages and it also created all the necessary configurations. Using this setup benchmarking was greatly simplified and we were able to try out different parameters for the clustering algorithm at the same time and observe how that affected the clustering performance and quality. Much of the testing involved in the development of the application was done on DigitalOcean. They are a PaaS<sup>3</sup> that offer virtual private servers, and one of the most useful feature they provide is the ability to start up several machines at the same time capable of running the project, each with custom performance capabilities. The servers used were equipped with 4 CPUs Intel(R) Xeon(R) CPUs and 8GB of RAM. Using virtual machines was an ideal setup for benchmarking because even though we used several different instances we were able to use the same snapshot across all machines thus ensuring a uniform testing environment. Once a container is started it will automatically begin fetching new data from Twitter and after an initial waiting period (in which the queue fills up with messages) the processing pipeline starts as well.

## 3.2 Testing and Evaluation

One of the advantages provided by the Scala programming language is its ease of manipulating collections using functions familiar to most functional programming languages such as map, reduce, filter, fold. This was an advantage that made it a perfect choice for manipulating

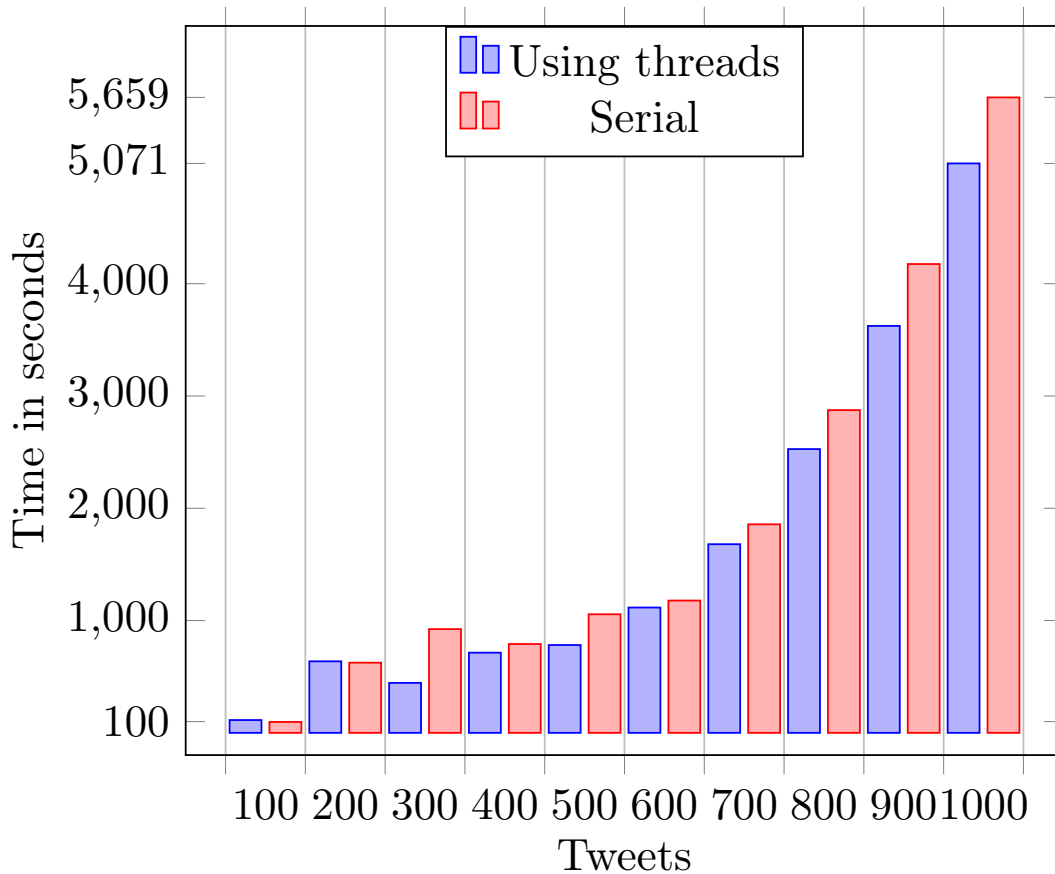
<sup>1</sup><https://www.docker.com>

<sup>2</sup><https://gist.github.com/piatra/53c0b6d185d10eeef8b>

<sup>3</sup>Platform as a Service

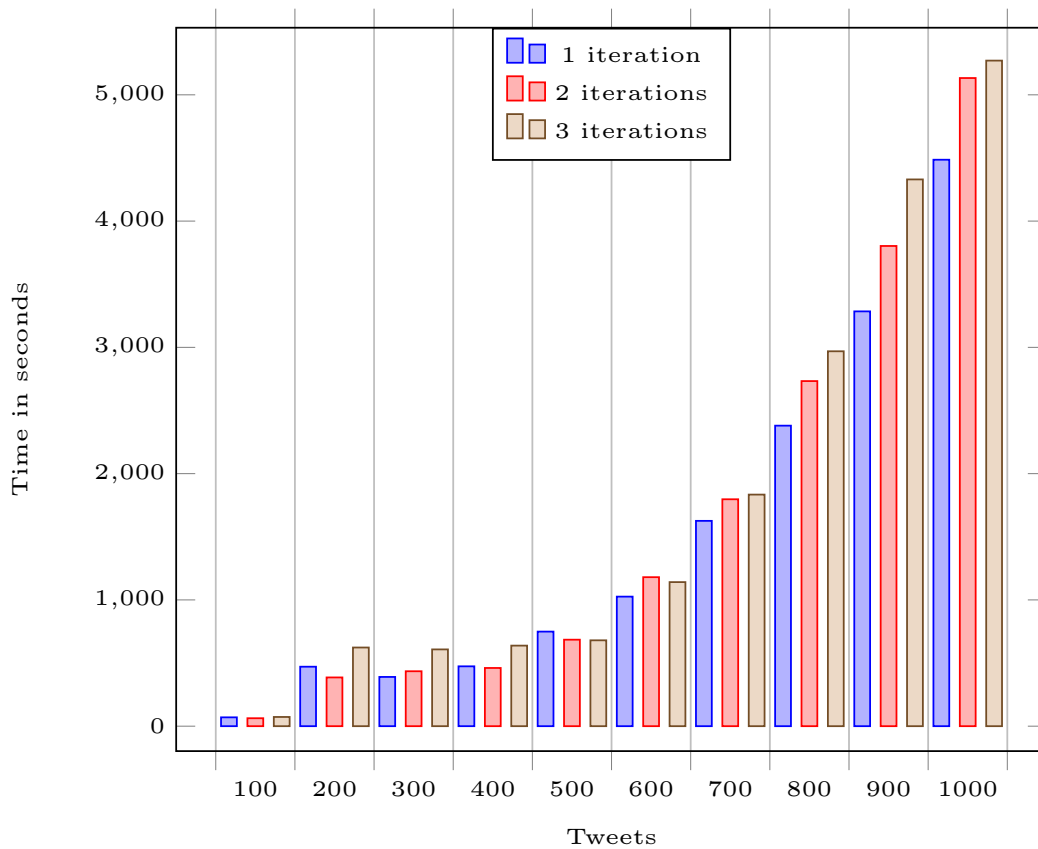
large collections of messages. At the same time Scala provides *parallel collections*<sup>1</sup>. These are special collections included in the standard library that offer a high level API that facilitates parallelization. Calling the `.par` method on a regular collection such as **Lists** or **Vectors** transforms it into a parallel collection and one can continue using that as a regular collection but now the transformations are done in parallel.

Below is a chart that outlines the performance improvements of using parallel collections. For the final set of data containing 1000 messages the performance improvement is of 588s (9,8 minutes).



As mentioned in the **Implementation** section, each cluster generated by the K-Means algorithm has a different weight and running the algorithm multiple times can yield better results, at the cost of running time.

<sup>1</sup><http://docs.scala-lang.org/overviews/parallel-collections/overview.html>

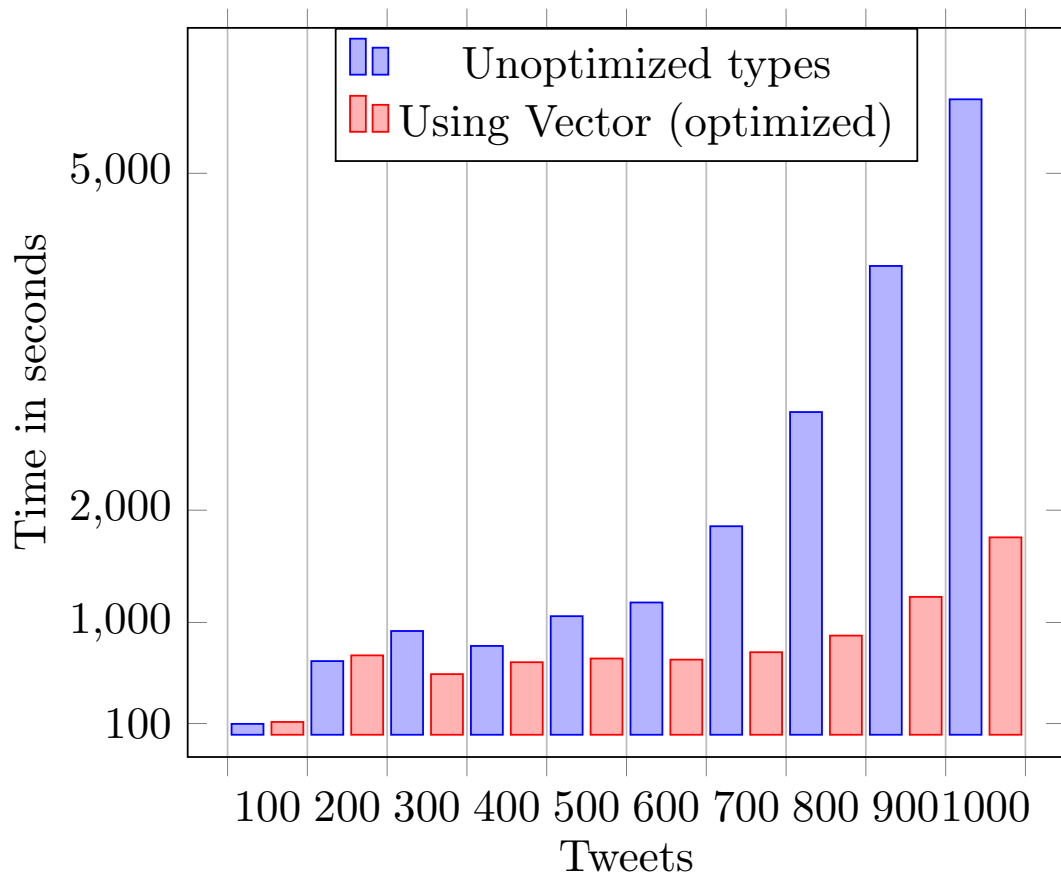


Above we can see how running 1, 2 or 3 iterations impacts the running time.

An advantage of using queues means that the consumer and producer processes that interact with the queue are completely decoupled. Therefore we are able to spin up several producers and consumers at the same time. Producers could either be Twitter API producers that provide the tweets, this would speed up the cold start of the application, currently it waits for 6 minutes for the queue to fill up with messages. Another producer in the pipeline is the parser that handles part-of-speech tagging and filtering of messages, this consumes raw tweets and produces tokenized messages that are stored in another queue.

Types in Scala are another optimization that yielded improved performance. Parallel collections offer improved running time by delegating the work to several threads but using incorrect types can yield unfavorable results. **Lists** in Scala are collections optimized for LIFO (last-in-first-out) type of access. Although random access is possible it comes with a performance loss. This disadvantage is also present when converting to a parallel collection *ParVector*, the overhead comes from the fact that all sequential collections such as lists have to be copied over to each thread. Using a different data structure improves performance especially when the quantity of messages increases.





From the chart above we can notice a 3x speed improvement when choosing the appropriate data structures.

## Chapter 4

# Conclusion

This is the conclusion. ce ați obținut, ce obiective ați avut, cum este relevant proiectul, ce rezultate ați obținut, cum ați evaluat.

The Internet allows anyone to create content, publish and spread information. Social platforms are lowering the barrier of entry and are making it especially easy for everyone to have a voice on the Internet and as a result millions of messages and content of all forms is generated daily. Making sense of everything, and keeping track is becoming difficult and therefore it is time for tools that help understand the content to evolve and adapt to these mediums and make content exploration as easy as it is to post a message.

Streamer attempts to solve this problem of content discovery and exploration. Streamer endeavors to understand the topics being published and presents it to the user in a way that is accessible and easy to use. It creates clusters of messages by interpreting content and presents it to the user in a web interface that allows for him to browse through a large number of messages efficiently.

The feature that makes Streamer relevant for the fast passed rate of tweets is its ability to parse the messages in real time. The data is not based on an archived corpus of documents but on streaming tweets as they happen. This way popular events, news and messages get reported in the interface and the user is able to keep in touch with what is happening right now.

- Getting real time data from Twitter using its API based on user queries.
- Parsing tweets as they arrive. Using part-of-speech tagging to make annotations that help filter messages and extract important information.
- Using a clustering algorithm that is able to group messages, that has the ability to configure precision and that can run in parallel for a choice between speed and precision.
- Building a decoupled system that can easily scale through the use of queues which allow for different rates of consumption and multiple consumers that can process the workload in parallel.
- Presenting the information through an accessible medium: the web browser, with an easy to use interface that allows the information to be explored.

Evaluating the results was done through the use of the frontend Streamer-Frontend. We expanded the clusters and started exploring the different topics.

**TODO:**  
more info

One unexpected answer we got by running the application was discovering topic influencers. What we refer to as influencers are people with large following in social media that also are really engaged with their followers through comments, retweets and favorites. The way it comes

up in Streamer is that a lot of people will end up retweeting a certain tweet written by one of these accounts, and it will show up as a cluster composed of the same tweet. There is also an exception, meaning that cluster can be composed of a single tweet of an account with few followers if their message was in turn retweeted by an influencer. In this case the cluster is formed entirely out of one popular tweet.

---

@OdgerFinnegan: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

---

@LoulsOuou: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

---

@charen\_lyno05: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

---

@GlavinArlenearl: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

---

@cipriani\_takoma: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

---

@FerncoteHallman: RT @Facebook\_Poker: Where do you buy the cheapest zynga #poker chips on the internet? <http://t.co/jBqRaaXx0R> #WSOP

Figure 4.1: Automated messages sent by bots

Due to the way the Twitter stream API works: it sends a percentage of all tweets currently being exchanged that match your query, we were able to make another interesting observation. A large number of Twitter topic clusters are formed from messages from bots. Twitter bots produce automatic messages usually with spam or promotional links. For the topics we tested on, mostly programming language topics, we noticed a large number of clusters related to job advertisements which turned out to be automated messages. In the figure we extracted an example of bot messages that are using a popular hashtag (World Series of Poker) to advertise a game. This also provides some indication that relying completely on hashtags in order to generate topics and cluster messages is not always very efficient. Bots might use trending topics to their advantage: most Twitter clients offer limited discovery features and usually just show tweets that match a certain hashtag, by using the same one, spam such as this can end up in your message list.

## Chapter 5

# Further work

The system is perfectly usable and but requires that each user deploys its own version of Streamer. Although the use of containers for automatic deployment via Docker reduces this task to running a script, the overhead involved with managing your own machine or environment where this can be done is a drawback. Wanting to improve on this issue but also increase the overall performance of the project, we have identified some areas where it can be improved:

- Having a centralized solution that offers immediate access to Streamer: it would offer Streamer as a service, allowing users to insert keywords or preferred sources of data and provide in return the clustered messages. This would remove the overhead of taking care of deployment and would make it easier to experiment with the project.
- Further improving the clustering algorithm both in terms of speed and precision. Currently a bottleneck of the project is the part-of-speech library which has an average parsing speed of 1 message per second. This could be replaced with other methods for determining the relevant keywords in a message. One example is using G-test log likelihood to remove statistically insignificant terms and improve clustering precision. This solution could run on multiple threads and improve the overall performance, but tests are required to ensure the project does not suffer a loss of precision.
- Using a distributed solution over multiple containers or even multiple machines. The container used in deployment contains both Streamer and the queues that hold the data. Extracting the queues would mean they can be shared between multiple instances and would be oblivious to any restarts to the server and would not lose data.
- Expanding the pipeline to add multiple consumers for the data, using multiple part-of-speech-taggers, and having more than one producer that retrieves data from the Twitter API. This would allow the system to retrieve more data as well as improve the time it takes to process it.
- Use websockets to notify Streamer-Frontend of new data that has been clustered. It would improve the user experience: right now when the API has new data the whole interface is redrawn including the clusters meaning that clusters may change position on the page making it harder to identify them. With a websocket implementation Streamer-Frontend can use event listeners and simply append new information to the page.
- Adding a storage layer or a caching layer in order to speed up the results. The storage could be shared between all running instances of Streamer and provided that the same query is used clustered results could be returned instantly. The storage layer could also hold raw or parsed tweets but basic check to ensure that the data is relatively new are

required. Caching would make more sense due to the fact that real time messages are expected.

## Appendix A

# Project Build System Makefiles

### A.1 Makefile.test

---

```
# Makefile containing targets specific to testing

TEST_CASE_SPEC_FILE=full_test_spec.odt
API_COVERAGE_FILE=api_coverage.csv
REQUIREMENTS_COVERAGE_FILE=requirements_coverage.csv
TEST_REPORT_FILE=test_report.odt

# Test Case Specification targets

.PHONY: full_spec
full_spec: $(TEST_CASE_SPEC_FILE)
    @echo
    @echo "Generated full Test Case Specification into \"${^}\""
    @echo "Please remove manually the generated file."

.PHONY: $(TEST_CASE_SPEC_FILE)
$(TEST_CASE_SPEC_FILE):
    $(TEST_ROOT)/common/tools/generate_all_spec.py --format=odt -o $@
    $(TEST_ROOT)/functional-tests $(TEST_ROOT)/performance-tests
    $(TEST_ROOT)/robustness-tests

#

# ...
```

---

Listing A.1: Testing Targets Makefile (Makefile.test)