



Feed-me!: Automatic Identification and Crawling of Event Feeds from the Web and Social Media

Punpikorn Rattanawirojkul

School of Computing Science
Sir Alwyn Williams Building
University of Glasgow
G12 8RZ

A dissertation presented in part fulfillment of the requirements
of the Degree of Master of Science at the University of Glasgow

8th September 2018

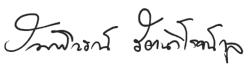
Abstract

Information are now available in the electronic format which makes records, like news articles, are now accessible online. News agencies are not the only sources of news, social networks also play an important role in the world of the Internet. There are many specific business areas that requires a lot of data. For example, the stock trading. On the other hand, news agencies themselves need data for their publishing and disseminating. However, gathering the information as much as possible of particular events from different sources together could be difficult and ineffective. Sources do not provide all data for free, offering only a limited quota. Information from different sources are not in the same format and the approach to access is also varying. On the other hand, the content of the event is changing overtime. For this reason, all aspects of that event may not be captured. Feed-me! is the event information aggregation web application that aims to aggregate information from different sources together as much as possible. To aggregate data from different sources, it lets users specify how to collect the data from a particular source by providing the mapping information from raw responses to a desire format of record, an authentication approach, including the period of crawling that allows the system to frequently get and update more records. It also features with query expansion, expanding from related data of a given query to another new query. Users could manually perform an expansion with the help of information provided by the system or let the system automatically perform an expansion periodically according to the given configuration. In addition, two graphs are available to let users gain insight from related records of a specific query. Bar graph shows the number of records from each source in timeline fashion while timeline graph explains the relationships between queries-records and queries-queries (expansion) in node-edge graph. Users can interact with the system via web pages that allow them to query for a specific event, see the result and visualisation of a query, track the existing queries, manage crawlers and manage query expansion. The project was evaluated using the software testing and the user evaluation. Software testing tests the correctness of core operations while the user evaluation evaluates overall functional requirements. For the user evaluation, evaluators will be asked to complete a set of tasks while the author records completion time and their suggestions. According to the evaluation result, query expansion is feasible and would be worked in practice. However, user interfaces improvement, especially indications and guides, is needed to help users understand the system without training sessions. Overall, the evaluation is successful by having more than 80% of passed evaluators and more than 80% of passed criteria in the questionnaire rating section.

Education Use Consent

I hereby give my permission for this project to be shown to other University of Glasgow students and to be distributed in an electronic form.

Name: Punpikorn Rattanawirojkul

Signature: 

Acknowledgements

I am particularly grateful for the assistance, support and solution to the problems given by Prof. Iadh Ounis and Dr. Richard McCreadie, my dissertation supervisors from the school of computing science, University of Glasgow. I wish to acknowledge the help provided by the evaluators of this project for their precious time, opinions and suggestions toward the project. Many thanks to Twitter and News API for such wonderful main data sources of this project. I wish to thank my friends and my parents for the support and encouragement throughout my study. Finally, I am honored to obtain an amazing opportunity for a study and a scholarship from the University of Glasgow.

Contents

Chapter 1 Introduction	1
1.1 Background and challenges	1
1.2 Solution.....	2
1.3 Structure.....	2
Chapter 2 Related Works and Tools.....	3
2.1 Related Works.....	3
2.1.1 Query Expansion.....	3
2.1.2 Sover!	3
2.2 Related Tools.....	4
2.2.1 Web Application Frameworks	4
2.2.2 JavaScript Framework	5
2.2.3 Frontend Frameworks.....	5
2.2.4 Visualisation frameworks	5
2.2.5 News APIs.....	6
2.2.6 Database.....	6
Chapter 3 Requirements	7
3.1 Requirements Gathering.....	7
3.2 User Definition.....	7
3.3 Requirements List.....	7
Chapter 4 Design	10
4.1 Overall System Design	10
4.2 Model-View-Controller (MVC)	10
4.2.1 Model	10
4.2.2 View	12
4.2.3 Controller.....	15
4.3 Query Expansion	16
Chapter 5 Implementation	17
5.1 Programming Languages	17
5.2 Model.....	17
5.3 View	18
5.3.1 Main views	18
5.3.2 Visualisation.....	20
5.3.2.1 Bar graph.....	20
5.3.2.2 Timeline graph.....	20
5.4 Controller	21
5.4.1 DbController.....	21
5.4.2 ResultController.....	22
5.5 Services.....	23
5.5.1 Crawler.....	23
5.5.2 Query Expansion.....	23
5.5.3 Scheduler for Crawler and Automatic Expansion.....	25

Chapter 6 Testing and Evaluation.....	26
6.1 Software Testing.....	26
6.2 User Evaluation	27
6.2.1 Evaluators.....	27
6.2.2 Preparation.....	27
6.2.2.1 Task-based Test Cases.....	27
6.2.2.2 Data Preparation.....	27
6.2.3 Materials.....	28
6.2.3.1 Project overview slides.....	28
6.2.3.2 Task sheets.....	28
6.2.3.3 Questionnaire.....	28
6.2.4 Evaluation Steps.....	29
6.2.5 Success Measurements.....	29
6.2.6 Result Summary	30
6.2.6.1 Average Time spent.....	30
6.2.6.2 Questionnaire: Section 1.....	31
6.2.6.3 Questionnaire: Section 2.....	32
6.2.7 Problems	33
6.2.8 Conclusion.....	33
Chapter 7 Conclusions.....	34
Chapter 8 References.....	36
Appendix A Initial UI Design	1
Appendix B Screenshots.....	2
Appendix C Detail of Software Test Cases	9
Appendix D User Evaluation Tasks.....	10
Appendix E Questionnaire section 1: Rating.....	12

Chapter 1 Introduction

1.1 Background and challenges

Nowadays, information is being published over the Internet in a digital form. The advantage of publishing information through Internet is a short distribution time. For example, we could see online news about the natural disaster published right after the incident happened. A notable example that follows this approach is news agencies. Previously, news was distributed through newspapers and television channels. However, agencies now currently have their own websites to persistently keep and distribute online news regarding to the increasing of information consumption through online medias. There are a lot of news sources, not only online news agencies, but also social networks where information is mostly published by ordinary people. Publishing information via social networks is faster compared with news agencies because of fewer steps made before the publishing, resulted in higher chance of being false information. News agencies have to look for events or information that could be a news first and verify the correctness then go through production steps further on. Information would spread from news sources through different channels when a single event is taking place. Developers or people who have an interest on data could get an information related to a specific event using web APIs (Application Programming Interface) provided by online news sources. News APIs usually come in 2 types. REST APIs provide access to past *records*, such as news article or tweet, which is collected in the provider's databases while streaming APIs provide access to a particular portion of real time records. Twitter is an example of a provider that has streaming APIs. However, they do not offer all of this data for free. News providers typically offer restricted rate limits, limited number of records in a response and few options to filter the records. Users have to pay if they need higher rates or better services and support. There are specific business areas that need lots of data, such as stock trading. Market data assets the decision of traders, especially day traders where real-time market data is impactful [1]. On the other hand, news providers themselves are organizations who need high amount of data. They need to track information emerging from different sources and publish as their content or disseminate to other needed clients. Bloomberg L.P is a company that provides financial news, media, software and data [2]. There are 5000 of stories from 120 countries published daily from the company's private news agency and about 325,000 clients subscript to the company professional service called Bloomberg Terminal [3]. Unimaginable amount of data has been ingested within the firm. Therefore, the solution to reduce the expenses on APIs, of a single provider, and to get lots of data as much as possible is to aggregate different information from different sources together which could be difficult due to variants of record formats and methods of access data.

There are many challenges to aggregate records together. Since they naturally came from different sources, approaches to get records from APIs could be different too. Clients need to authenticate themselves before using the API services. There are different authentication types available such as API Key and OAuth which also require different key information too. Responses from the web APIs are usually in JSON (JavaScript Object Notation) format. However, each source has its own object structure. Mapping from different structures of responses to the same structure is necessary to collect those records together in the same database collection.

Information from a particular event is continually published from the beginning of the event. So, records gathered at the beginning may not represent whole event. It is

important to frequently *crawl*, to get and update new records, more content over time in order to get all related data and better understand the nature of an event.

Content of an event changes from time to time. At the beginning, event like a natural disaster, articles portray destruction of housing and surrounding environment. After that, focus might move to the rescues. It is also possible to have sub-events such as actions of other countries or celebrities related to the event. As you can see, the event comprises of many interesting aspects. It may not be possible to cover every element of the event if *query*, terms that represent the event, is static. For instance, “cave thailand” is a query representing an event of Thai soccer team rescue. News related to other aspects like the help from Elon Musk and cheers from famous soccer players may not be included from the tracking of query “cave thailand” since the query gives overall information.

Humans understand information better with the help of visualisation. There could be insight within the relationships between records and queries. For example, it is interesting to see why a record is related to several queries. That record may give you the idea to generate a query to get more information related to that event.

Overall, identified problems are restrictions from free APIs, differences of APIs formats and accesses, the needs of frequent crawling, changes of content over time and lack of visualisation that improve human understanding. The project represents approaches to tackle these problems in the following section.

1.2 Solution

Feed-Me! is an event information aggregation web application powered by configurable *crawlers* (a unit that perform crawling from a specific source) with query expansion feature. A system main objective is to gather as much information as possible from given queries. It tackles the differences between news sources by providing an interface to define how to collect and map responses and authenticate API requests with key information to access API of a specific source. This enables the application to aggregate data from different sources together. The aforementioned configurations are called *source* information, where user can also enter the crawling period that let the system automatically schedule a *crawling*, getting and updating more records from source API, on tracked queries periodically. Moreover, the application provides several operations to view and manipulate collected data. User can browse through collected records and tracked queries, see visualisation of a query, manage crawlers, and manage query expansion. *Query expansion* is another main feature. “Expand” is to establish a link from one query to another query assuming that 2 queries are related in a particular aspect. This mechanism helps the system adapting to the changes of an event’s content and get more related data.

1.3 Structure

In the following section, we provide comparisons between similar systems and give brief overview of tools and libraries used in the development. Functional and non-functional requirements including how they are gathered are listed, prioritized by MoSCoW method, and indicated with implementation status in Section 3. Section 4 shows overall software architecture diagrams with the explanation of components and interactions also the UI storyboard. Section 5 explains system components in detail of how it was built, challenges and solution to overcome those. Detail of the software testing and user evaluation take place in Section 6. Finally, conclusions and future works are discussed in Section 7.

Chapter 2 Related Works and Tools

This section discussed about the related works, and tools that were used for the project. Section 2.1 describes the related works while Section 2.2 explained the related Tools.

2.1 Related Works

2.1.1 Query Expansion

In the area of Information Retrieval (IR), the general definition for a query expansion is the query reformation techniques that improves the performance of IR operations, while the expansion is specifically expanding a query using user's input (including terms in the query and other data) to get more related data in the search engines context [4].

The authors of the book stated that, in the query expansion, users could give an additional input by suggesting terms for a query, especially used in Web [5]. For example, the web suggests users to add terms like "Trump" or "Duck" when querying for "Donald". Additionally, the expansion could be *interactive or non-interactive*, the previous example is interactive since it lets users choose the term. *Automaticity* is also considered; the previous example is manual since users are responsible to choose the word. They remarked that the main problem for the query expansion is an approach to generate these suggested terms or an expanded query for users and suggested that there are commonly 2 approaches which are *global analysis* or *local analysis*. Global analysis corporates the use of a specific thesaurus, it expands the query by analysing each term and automatically adds the related term from the thesaurus. Local analysis expands a query by analysing related documents of the query instead.

For the local analysis, getting the related documents of a query requires the asset of document relevancy judgement. Users could give a relevancy directly to a result document, called *explicit relevance feedback*. However, collecting such relevancy is difficult because users need to mark all documents which is not practical. Therefore, automating the relevancy judgement of documents might be a more practical solution, called *pseudo relevance feedback* [6]. For each querying, the top "k" ranked documents in a set of result documents are assumed to be relevant. Then, the system analyses these documents and selects top ranked terms from the documents using the weighting techniques such as tf-idf (term frequency-inverse document frequency). These top ranked terms will be used as suggestion terms in the query expansion.

The query expansion approach of the dissertation expands a query, after the analysis, by *creating another query*. It does *not change or improve the first query*, however generating a new query to collect more relevant data. The query expansion of the dissertation performs a local analysis using the pseudo relevance feedback because of its practicality. More detail of the query expansion approach including the interactivity and automaticity will be explained in Chapter 4, in Section 4.3.

2.1.2 Sover!

Sover! is a system that monitors events real-time in Twitter using a technique called adaptive crawling and represents them in a feature-rich dashboard [7]. The adaptive crawler fetches data from the Twitter streaming API via 2 streams. First is the *filter stream*, which is the system main input, that filters tweets according to filter parameters. Another is the *sample stream* that simply provides 1 percent of the global real-time tweet stream. The system analyses tweets from the filter stream to get and

track new *query keywords*, keywords that are used as search queries, and eliminates *noise keywords*, keywords that come from other events, with the help of the sample stream. Dashboard provides content from different sources. It performs URL unshortening on tweets, groups them by sources and shows them in the dashboard. The dashboard also shows images and videos according to the source of a tweet.

The key differences between Sover! and Feed-Me! are listed in table 1.

Criteria	Sover!	Feed-me!
Goal	Real-time event monitoring with adaptive crawling	Event monitoring and information aggregation from different sources with query expansion
Data source	Mainly from the Twitter and enriched by the source of tweets	Any sources that provide API with supported authentications (OAuth and API Key)
Crawler	A crawler is specific to a seed keywords and parameter setting and keyword is updated and added from the adaptation feature.	A crawler is specific to a source, crawling new records for every tracked query.
Adaptation feature	Adaptive via adaptive crawling by analysing tweets from the filter stream to get new query keywords and eliminate noise keywords with the help of the sample stream	Adaptive via query expansion by analysing related records of tracked queries using any text analysis techniques (provide interface to let users implement their own technique) to get new expanded queries
Assistance from human	Manually adding keywords to a crawler	Manually perform the query expansion with the aid of statistical information generated by the system
Dashboard	Rich with content from different sources including videos and images	Query result page shows the related records, in text, of a specific query allowed user to filter and delete records
Visualisation	Histogram of popular hashtags and users in dashboard	2 visualisations: Bar graph shows the number of related records from each source of a specific query and timeline graph shows the relationship between queries-records and queries-queries (expansion)

Table 1: Key differences of Sover! And Feed-me!

2.2 Related Tools

2.2.1 Web Application Frameworks

Play: Play is a modern web application framework for Scala and Java [8]. Play adopts Model-View-Controller (MVC) architectural pattern. It has 3 components. Model (M) is responsible for the data management which accepts parsed inputs from the controller. View (V) is a user interface which is a representation of model. Controller (C) is a middleman that controls data flows between view and model [9]. Implementation using Play is a Must-have requirement from the beginning.

2.2.2 JavaScript Framework

The JavaScript frameworks offer powerful tools that enhances the frontend web development. They provide methods to efficiently manipulate DOM (Document Object Model), elements that build up a web document. According to Matt Burgess, the purpose of these frameworks is “*to map the application state to the DOM*” [10]. The following is the list of related tools that have been considered for the development.

jQuery: According to its official website, it is a simple, fast and light-weight JavaScript framework providing utility tools for DOM manipulation, animations, making asynchronous request and etc. [11].

Vue.js: Vue is a progressive framework for web view development which is easy to integrate with other libraries [12].

The author decided to use Vue as a main implementation of views for several reasons. First, it was a good chance to learn this new and popular JavaScript framework. It demonstrates a new style of web application development since the author still had an experience on the old era of the development using just jQuery. Another reason is that the author has bad experience on jQuery, having used it during the team project last semester. Developing using pure jQuery is not appropriate with complex web application. The codes are very long and hard to refactor, unlike Vue which highly adopts the concept of code reusability. However, jQuery is still required to be imported for some frameworks such as Bootstrap.

2.2.3 Frontend Frameworks

Frontend frameworks provide tools to enhance look and feel of the web application. The author decided to use **Bootstrap** because of prior experiences and its simplicity. Bootstrap is a HTML, JavaScript and CSS framework that provides lots of nice reusable components such as navigational bar, jumbotron and etc. [13]. It also offers various styles of existing components such as buttons and give a way to customize them easily. Since the framework is very popular, there is a customized version for Vue called **Bootstrap Vue**. Bootstrap Vue integrates Bootstrap components with Vue. It increases the ease of calling and controlling Bootstrap components within Vue codes. The author used both frameworks on this project but depending on the situation considered by the implementation simplicity.

2.2.4 Visualisation frameworks

The development of a complex component such as visualisation would be difficult without helps from frameworks. According to the design of visualisation in Chapter 4, the system will have two graphs which are bar graph and timeline graph. Bar graph shows number of related records from each source of a query and is developed using **Amchart**. Amchart offers many types of rich-features chart, including bar graph. Additional features are added to the graph to enhance data browsing such as a zoomable and scrollable chart [14]. Timeline graph is a node-edge graph that reveals the relationships between records-queries and queries-queries, showing the expansions from a query to others. It is not an ordinary graph which is needed to develop by the deep-to-fundamental framework like **D3.js**. D3.js allows developers to create visualisation from the core concept, by binding data with DOM [15]. It is flexible because the developers are not bound into any pre-implemented visualisation features, instead they need to develop their own visualisations by manipulating DOM. In timeline graph, nodes and links were built by D3.js.

2.2.5 News APIs

There are various of APIs for news. However, not all of them are free to use and even free one offers only a limited rate limit. **Twitter APIs** is very popular and offers free standard APIs for developers. There are many APIs available, but the system only uses the search tweets API. Rate limit for the API that authenticates with Twitter user account is limited to 180 requests per 15-minutes window which is enough for the system [16]. However, news from news agencies are necessary too. **News API** offers free API for getting information of a news from over 30000 news sources [17]. Information includes the headline, short description, authors, published date and etc. The rate limit for non-commercial project is 1000 requests per day which is plenty to satisfy the requests made for this project [18]. Twitter and News API are the main data sources for the system crawler. Especially, News API is used for getting articles from different sources such as BBC and CNN.

2.2.6 Database

The author and supervisors agreed to use NoSQL database because we could have a lot of records and queries in the future. With NoSQL, the system could be easily scalable. However, Play does not provide native APIs for NoSQL, unlike SQL. After some studies, the author found that **MongoDB** is the only NoSQL database listed in the Play official driver list and decided to use MongoDB through the driver called **ReactiveMongo** [10]. MongoDB keeps and represents data as a document, in JSON-like format [20]. ReactiveMongo is the well documented and has an active community, compared with other MongoDB drivers. It also supports Scala. The main features are asynchronicity and non-blocking operations [21]. The system uses Reactivemongo widely in every part that need to communicate with the database.

Chapter 3 Requirements

The chapter starts with how we gather the requirements in Section 3.1. Then, the definition of user is defined in Section 3.2. Section 3.3 shows the list of discovered requirements identified as functional or non-functional, categorized by MoSCoW method followed with the implementation status.

3.1 Requirements Gathering

Requirements were gathered from the weekly meeting in the form of ideas and discussions given to the deliverables of that week by supervisors. The deliverables could be ideas from the author or pieces of implemented requirement that have been developed during the week. The author normally gave the presentation on deliverable at the beginning of the meeting. The ideas and discussions on deliverables given by supervisors are mostly refinements on that deliverable, which could be extracted into requirements. After the meeting, the author transformed these ideas into requirements, identified as functional or non-functional and categorized by the MoSCoW method. MoSCoW categories helped ordering the features by labeled the priority on them, beginning with Must-have (M), Should-have (S), Could-have (C) and Won't-have-this-time (W) feature. To be more specific, M requirements are stated at the beginning in the introduction slides of 1st week while other categories were discovered in the weekly meeting of following weeks. The author then decided to choose a task to design and implement considering these criteria, for instance picking M requirements first.

3.2 User Definition

The supervisors of this project are *real users* for the system, they actually use this system in the future. However, the general definition of user is still needed to give the idea of target group. Users in our system are people who have an interest in the event data collection. However, users are recommended to have a certain level of programming skill because several main features require knowledge of web technologies and the understanding of how features work. For an instance, users need to specify technical information, such as JSON structure and API URL, when adding a crawler. One of the important attributes is a mapping information from a raw response to mapped records, which represented in JSON structure. In addition, the detail of an API like the authentication, an API URL and the access attribute are needed too. Another example is a query expansion feature that allows users to compose their own query expansion technique. But, they should understand how each component of the query expansion works and what will be the result of those. Overall, users without programming background could find that it is difficult to understand and use the system. Additionally, even evaluators where everyone is a programmer, found that the system was hard to use if there is no detailed guides and tutorials.

3.3 Requirements List

Functional requirements will be written in a format “*<who> can <what>*” while there is no writing format for non-functional requirement. MoSCoW method were applied in order to rank requirements in term of priority. Discovered requirements will be grouped into section by features, identified as functional or non-functional, categorized into MoSCoW categories, and indicated with the implementation status in table 2. Description of ranking criteria is in table 3 and summary of requirements grouped in MoSCoW categories are in table 4.

No.	Functional/ Non-functional	Description	MoSCoW Category	Imple- mented
A. Query and records				
1	Functional	User can query for records from sources	M	Yes
2	Functional	User can see information of tracked queries	M	Yes
3	Functional	User can delete a query	M	Yes
4	Functional	User can check if a query is tracked in the system	S	Yes
5	Functional	User can turn specific crawler(s) of a query on/off	S	Yes
6	Functional	User can see statistical information of a query	S	Yes
7	Functional	User can group queries together in term of expansion	C	No
8	Functional	User can export selective data of a query	W	No
9	Functional	User can filter out duplicate records	C	No
10	Functional	User can query for multiple queries using semantic keywords like and/or	W	No
11	Non-functional	System can link similar records together	C	No
12	Non-functional	System can keep related images from a record in the database	W	No
13	Functional	User can delete a record	S	Yes
B. Crawler				
1	Functional	User can add a crawler	M	Yes
2	Functional	User can delete a crawler	M	Yes
3	Functional	User can edit a crawler	M	Yes
4	Functional	User can see information of a crawler	M	Yes
5	Functional	User can force a crawler to perform crawling before a schedule	S	Yes
6	Non-functional	System can schedule a task for a crawler periodically	M	Yes
7	Non-functional	System can handle different API authentications	S	Yes
8	Non-functional	System can handle limited number of API rate limit efficiently	C	No
C. Visualisation				
1	Functional	User can see visualisations of a query statistic information	M	Yes
2	Functional	User can see relationships between records and queries	M	Yes

3	Functional	User can see relationships between queries and expanded queries	M	Yes
D. Query expansion				
1	Functional	User can perform query expansion on a query	M	Yes
2	Non-functional	System can perform query expansion on a query	M	Yes
3	Functional	User can control how system will perform query expansion	M	Yes
E. Other				
1	Non-functional	System must be implemented using Play framework	M	Yes
2	Non-functional	System must aggregate records from different sources together	M	Yes
3	Non-functional	System should keep raw responses in the database	S	Yes

Table 2: List of requirements grouped by features, identified as function or non-functional, categorised by MoSCoW categories and indicated with the implementation status

Category	Criterion
(M) Must have	The requirement corresponds with main goals and is obviously stated at initial meeting. The system will be incomplete if this is missing.
(S) Should have	The requirement is derived from suggestion during weekly meeting and significantly improves user experience of must-have requirements if implemented.
(C) Could have	The requirement is derived from suggestion during weekly meeting. It could improve user experience of must-have requirements if implemented but is estimated to spend effort excessively compared with should-have requirements.
(W) Won't have this time	The requirement is derived from suggestion during weekly meeting and is considered as optional because it is out of the project scope. Additionally, it could be relatively hard to implement.

Table 3: Criteria for categories of MoSCoW

Category	List of requirements	Total requirements	Implemented
(M) Must have	A1, A2, A3, B1, B2, B3, B4, B6, C1, C2, C3, D1, D2, D3, E1, E2	16	16
(S) Should have	A4, A5, A6, A13, B5, B7, E3	7	7
(C) Could have	A7, A9, A11, B8	4	0
(W) Won't have this time	A8, A10, A12	3	0

Table 4: Summary of requirements grouped in MoSCoW categories

Chapter 4 Design

Designs of the system will be discussed within this chapter. Section 4.1 shows the high-level system diagram derived from related tools in Section 2.2. System design in the aspect of MVC architectural pattern is explained in Section 4.2. Database schema will be shown in Section 4.2.1 under model while storyboard UI will be shown in Section 4.2.2 under view. Finally, the system diagram in term of MVC is displayed in Section 4.2.3 under controller.

4.1 Overall System Design

The diagram in Figure 1 explains high-level relationships between each related tool, from Section 2.2, in term of components that will be implemented, including the interactions of a user.

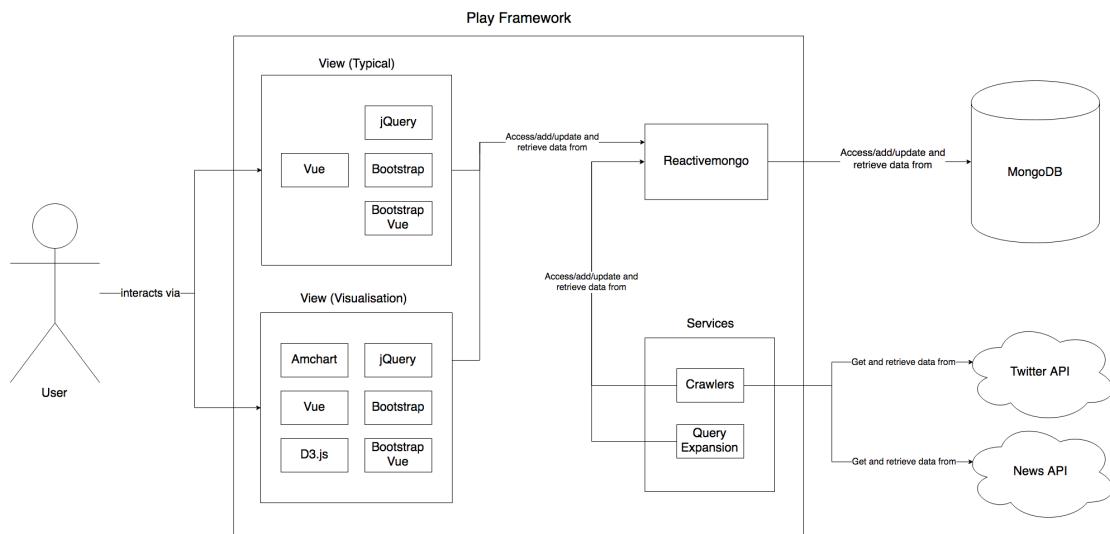


Figure 1: System diagram showing related tools, interacted with a user

Considering the requirement E1, the project must be implemented using Play framework. Play has a model-view-controller (MVC) architectural pattern as stated in Section 2.2.1. Section 4.2 describes the system design in 3 aspects according to MVC components.

4.2 Model-View-Controller (MVC)

4.2.1 Model

Model design is focused on how it will be kept and represented in the database. There is also a scenario when a model is not actually kept in database, source and query expansion unit are wrapper for system services. Figure 2 explains the detail of all attributes and interactions of models via the database schema. Concepts that could be extracted into models, highlighted attributes and ideas behind the model are described as follows.

SearchQuery: Query could be redundant in many concepts, specifically *query parameter* in the APIs. So, the term **SearchQuery** would be more explicit. Other than the text and typical model attributes, it includes previous queries (*prevQueries*) and expanded queries (*expQueries*) to conform the query expansion requirements (group D). The advantage is that it could be represented in graph, nodes and edges, later on in the visualisation part. The query is also designed to be independent of crawlers

through `EnabledSources` attribute. A query may be specific to particular crawlers, such as a hashtag query which is specific to the Twitter.

Record: Record is a model of the information unit in the system, such as news articles or tweets. There is no single static format of records because they come from different sources. However, they still have some attributes in common such as the text, created time, URL and etc. On the other hand, there are also attributes that are specific to the source. For example, hashtags, tweet ID and user ID are specific to Twitter while records from news agencies have authors and news source name. Additionally, the system adds other necessary attributes such as hash and database ID.

Source: A source holds information of a particular source for a crawler. It is special because it is not saved into database, it is saved in a configuration file instead. Users need to specify 2 types of collection names. `collectionName` is the name of collection where mapped records will be kept while `rawCollectionName` is the name of collection that collects raw records, satisfy the non-functional requirement E3. According to the requirement B7, the system supports 2 authentication types which are APIKey and OAuth. Information for such authentications, such as API key for the first type, is kept as JSON in a source model. The advantage is that it supports any authentication types, assuming that all information is correct and corresponding to the type. `APIFormat` specifies where the system needs to insert a query and authentication information in the API URL, depending on the authentication type. Response from API is in JSON and the system might need to access specific attribute name to get actual records, which is the responsibility for `recordAccess`. The most important attribute for a source is a `recordMapping`. It provides a way to map attributes from raw responses to mapped records which will be saved into the database collection (`collectionName`). There are few mandatory fields which are `text`, `title`, `createdTimestamp` and `source url`. Title could be “none”, in case of tweets, and user could give a string pattern for source url if it is not available directly in a raw response. Example of a mapping is shown in Figure 6 in Section 5.2.

QueryExpansion (Unit): It is a wrapper for query expansion components. Query expansion is designed to be flexible in order to satisfy the requirement D3. Basically, the expansion has 2 main components which are generator and selector and will be described in detail in Chapter 5. The unit is a reference of the generator and selector that users chose and it will be used by the query expansion service. The model is not saved into the database.

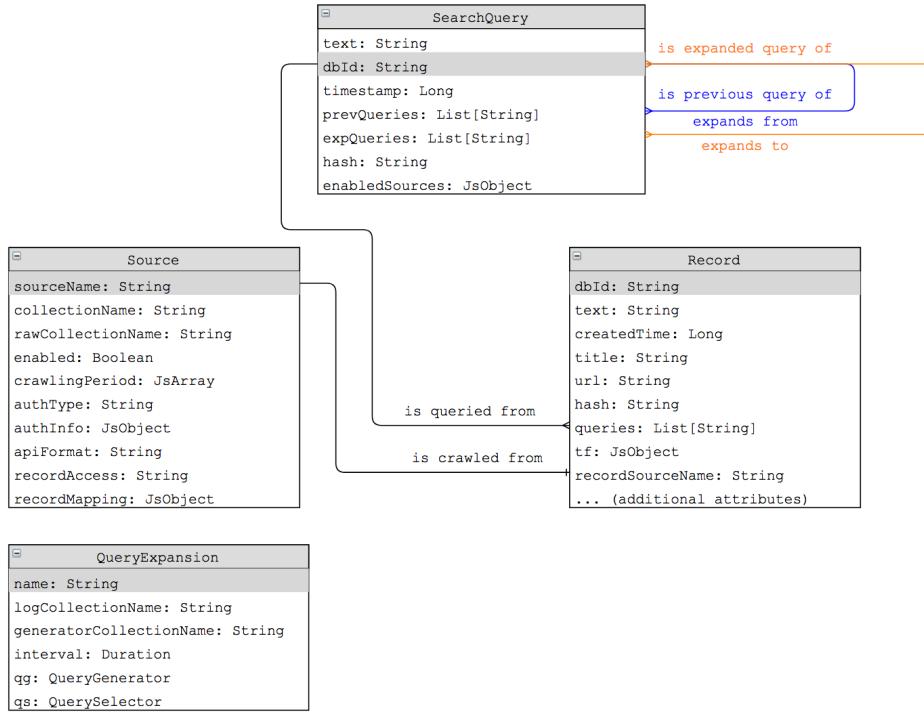


Figure 2: Database schema for models

4.2.2 View

View is simply a user interface representing the model, or the interface that receives data from user. A main goal of views is to let users interact with the system smoothly. The author adopted several core concepts when designing views for the system.

1. User interfaces should be easy to use, simple and to the point
2. Minimalism is a must, no excess decoration
3. Information should be displayed in the appropriate user interface components, such as cards and modals

The first views design was introduced on the 4th week, shown in an Appendix A. It comprises of index, result, track query, and manage crawler page. However, there were important changes listed as follows.

- **Card layout:** There are different cards in the design and each of them packs information of the model and binds operations with buttons in the card. The card layouts in initial design were changed due to the difficulties of an actual implementation. To be specific, most components on right side were moved to the bottom instead.
- **Query setting:** During the implementation of the requirement A5, the author found that users might not remember which queries were tracked in the system and especially their current query settings of which crawler is enabled or disabled.
- **Manage QE page:** Query Expansion (QE in short) is designed after the 5th week. At that time, we still didn't know how the feature is working and the appearance of UI was still unknown. It was later designed on the 8th week.

Overall, the ideas behind current designs for each view are listed as follows.

- **Index page:** It is the first view of the system. Basically, users can query and check for the existence of a query in the database if they cannot remember it. By checking the query, a previous query setting will be automatically applied

to the current query setting if the query is existed. In other views, a check button is located at top-right button group besides search bar (“C” button in the storyboard UI)

- **Result page:** The list of related records associated with the submitted query is listed here, displaying each in a card layout. This page allows users to browse the records, see more information, delete irrelevant records, filter and order the records. Submitting a query via top-right search bar will redirect the user to this page.
- **Track query page:** The view shows information of each tracked query in a card. User can delete queries and change query setting of a query via switch-like buttons. Users perform manual expansion in this page through the expand query modal. Information from each QE unit that showed in the modal also helps users to make decision on manual expansion. A generate now button is there to let users force the system to generate information for them immediately, in case of the system was just started and the automatic expansion has not been performed yet. Clicking at the visualisation button navigates users to the visualisation page of that query.
- **Visualisation page:** This page shows visualisations of a specific query. Users can switch between visualisations by clicking the tabs under the query text.
- **Manage crawler page:** Crawlers are displayed in this page as a form of cards. Users can edit the information of a crawler (source information), add new crawler via the add/edit crawler modal and turning a crawler on/off via a switch-like button. By clicking the crawl button on a crawler card, the system forces a crawler to start before its schedule time and crawling result, number of new records and total records, will be shown in the modal.
- **Manage QE page:** QE units are listed in cards. A card shows information of the QE unit including the expand and logs buttons. Users can force the system to perform an automatic expansion immediately before the schedule. After the expansion is completed, the expand summary will be shown in a modal. Clicking at logs button allows user to see the previous automatic expansions logs.

Figure 3 shows the latest storyboard UI of the system, comprises of pages design, interactions flow and main components design. Bold texts are the *main pages* (or *main views*), which are index, result, manage crawler, manage QE and track query pages. Bold arrow lines represent main page navigations where each page can navigate to any main pages from the navigation buttons on the left or top panel. The exception is for a visualisation page, which is not a main page, but navigation buttons are available. Dotted lines indicate a magnification of the component in that page, such as the layout of cards. Normal arrow lines are interactions after triggering a button, the annotations are explained in the following list:

- A1:** Submitting a query by clicking a search button redirects to the result page
A2: Check if query is existed in the database by clicking check button
A3: See more record information by clicking a more info button
A4: Delete a record, removing a card from the record cards list
B1: Show the add/edit crawler modal after clicking an add crawler button
B2: Show the add/edit crawler modal after clicking an edit crawler button
B3: Force a crawler to start before a schedule by clicking a crawl button then a crawl result modal will be shown
B4: Delete a crawler, removing a card from the crawler cards list
C1: Force a QE unit to start before a schedule by clicking expand button then an expand result modal will be shown
C2: Expand logs modal will be shown after clicking a logs button

D1: Clicking an expand button lets user perform a manual query expansion via the expand query modal

D2: Clicking a visualisation button in a query card navigates user to the visualisation page of that query

D3: Delete a query, removing a card from the query cards list

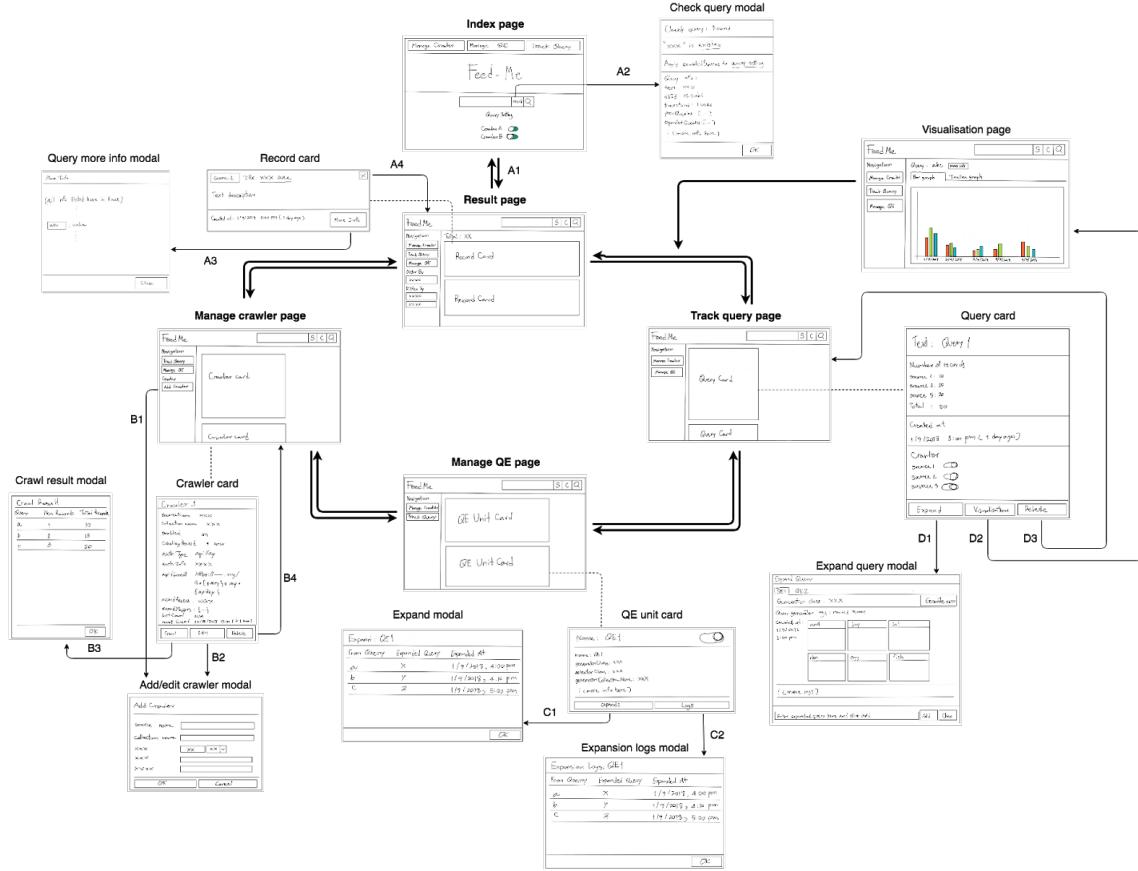


Figure 3: Storyboard UI

Another focus of the system is the visualisation. The aim of visualisation is to provide the insight into a query that cannot be achieved by simply looking at related records of a query. There are 2 graphs, which are bar graph and timeline graph.

- **Bar graph:** A query could be related with thousands of records from different point of time. The idea behind bar graph is to summarize the number of records from each source in timeline fashion.
- **Timeline graph:** Relationship between queries and records is interesting. We need to represent those relationships somehow in term of the query expansion. In other words, showing how queries are expanded. Also, there could be some particular records that have been crawled many times from different queries, these records link with many queries. This kind of scenario is important for records as it could be considered like gaining more reputation from queries. The goal of this graph is to show relationships between queries-records and queries-queries (query expansion). Initial design of the graph is shown in Figure 4.

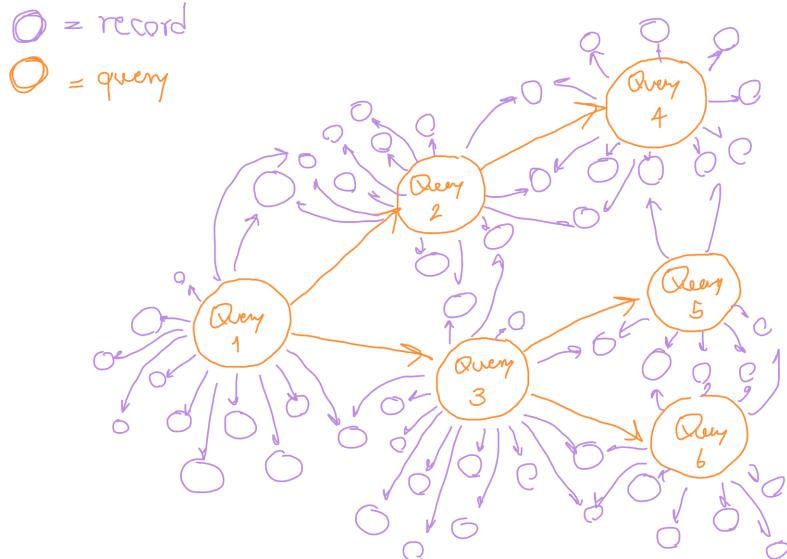


Figure 4: Initial timeline graph design

4.2.3 Controller

A controller manages the data flow from views to models and also interacts with other components in order to perform operations on the input data. The system basically has a single controller for a main view. It mainly handles interactions, or requests, from users in that particular view. For example, the index page controller has to handle a query request by parsing the input query text to a search query (model) and save it to the database. Visualisation page also has a controller since it needs to process and transform lots of data for graphs before returning back to users, even it is not a main view. There is a special controller called `DbController`. During the implementation, the author found that code chunks that are responsible for the database connection were very long and redundant. So, those chunks were refactored to a controller in order to maintain the cleanliness. The summary diagram of system interactions including models, views, controllers and database is in a Figure 5. The relationships between each controller and model(s) are explained as follows. Again, please note that DB controller handles all connections with the database.

Result controller

- `SearchQuery`: composes a query from the text input and save it to the database
- `Record`: retrieves records from the database and displays in the result page

Manage crawler controller

- `Source`: displays information of the crawler, edits, adds and deletes a crawler by manipulating source model

Manage QE controller

- `QE Unit`: shows QE units information, including turning on/off

Track query controller

- `SearchQuery`: lists all queries from the database
- `QE Unit`: shows the information of each QE Unit in the manual expand modal

Visualisation controller

- `SearchQuery`: gets the information of a specific query from the database
- `Records`: finds related records of a query and process those records into graphs

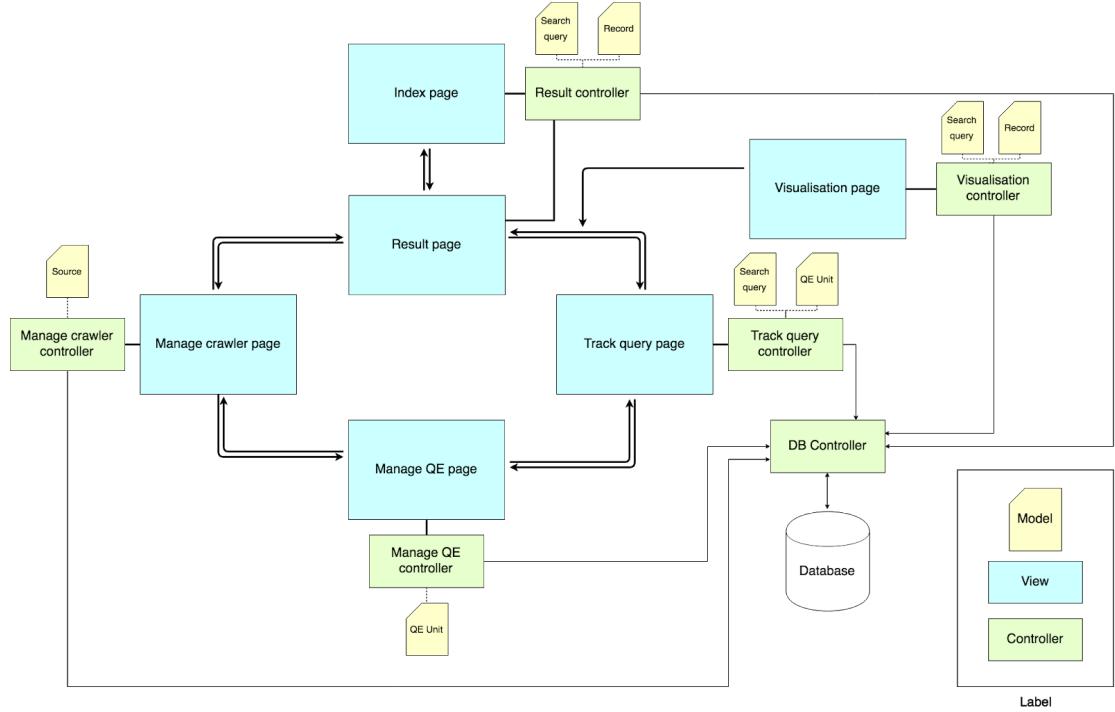


Figure 5: Summary diagram of system interactions

4.3 Query Expansion

We believe that it is possible to get more records by generating more queries from the existing data. Basically, the input for an expansion is the data of related records, from a given query and outputs an expanded query. The requirement D1 and D2 give explicit idea on designing query expansion approaches. There will be 2 approaches for query expansion which are *automatic expansion* and *manual expansion*. Automatic expansion satisfies the requirement D2, the system automatically performs a query expansion on the tracked queries, this is *non-interactive*. In each round of an automatic expansion, the system looks for related records of a query, analyses and processes them into statistical data, and composes a new expanded query. Manual expansion is *interactive*. It satisfies the requirement D1, allowing users to perform expansion on queries by their own judgements. The author also believes that human could identify potential insight from data better than the machine and would be good if there is an aid from the system. A help from the system is in a form of information. Actually, the system has that information in hands. According to the design of an automatic expansion above, the processed information from related records could be helpful to users, showing that information could assets the decision on manual query expansion.

The idea is a good new query comes from a good data analysis, using various techniques from the text analysis. Therefore, the author designed the software architecture of the query expansion that supports any text analysis on the collection of records. We divided query expansion into 2 components, *query generator* and *query selector* (generator and selector in short). Generator generates ranked terms and helpful information to be used in the selector. Selector selects the best term and composes an expanded query from the given ranked terms. More detail including the class diagram of the query expansion will be discussed in Chapter 5, in Section 5.5.2.

Chapter 5 Implementation

The detail of implementation starts from the discussion of programming language choices in Section 5.1. As mentioned in previous chapter, the MVC architectural pattern were used to describe the design. Therefore, the implementation will also be described in three aspects of the MVC components which are Model, View and Controller in the Section 5.2, 5.3 and 5.4 respectively. Finally, Section 5.5 explains how the system services works and implemented.

5.1 Programming Languages

We decided to develop the project with Scala for several reasons. The author already has certain experience with Java and it would be a great opportunity to learn how Play works in Scala. In addition, most developers in a community recommended using Scala if the project is not necessary to build upon the infrastructure that only supports Java [22]. The current version of Play itself is also written in Scala [8] which means Scala is native for the framework and is a main reason for me to go with Scala.

5.2 Model

Most models, except record, are easy to implement using Scala case class. A case class is similar to a normal class, but it is good for modeling an immutable data [23]. In order to save the object to the database or read it from the retrieved JSON document to a specific type, the object type needs to be written in a Scala case class. Writing a case class is very easy, it is simply a class with only attributes however methods are not allowed in a class. There are interesting aspects of specific models listed as follows.

- **Source:** Source is a model that act as a medium for source related operations such as adding and editing a crawler. Several components are related with the source model. A `sources.json` is a configuration file located in a server, while source model is a format of data kept in JSON configuration file. The file basically keeps all sources information. The idea for a file is that users could monitor and manipulate the files directly if they do not want to interact with the frontend. There is also a class called `Sources` which is responsible for holding live sources reading from the configuration file. It also handles incoming operations on `sources.json` which are reading, writing and updating the crawler.
- **Record:** There is no case class for records because the record format is not static. A `RecordMapper` class is responsible for creating records. It maps raw responses to mapped records according to a mapping setting of a particular source in the configuration file. Figure 6 shows the example of a mapping. As stated in Chapter 3, there are several mandatory attributes which are `text`, `title`, `createdTimestamp` and `url`. All 4 attributes need to be stated in the mapping setting. Special cases are for title and URL. Title's value could be "None", in case of no title is available like tweets. URL could be specified as a pattern, specified with `{ ... }` where the value in a bracket is a raw response attribute name. The example of the URL pattern is also shown in Figure 6. Additional attributes that are specific to a source are also supported such as tweet ID, hashtags and etc.

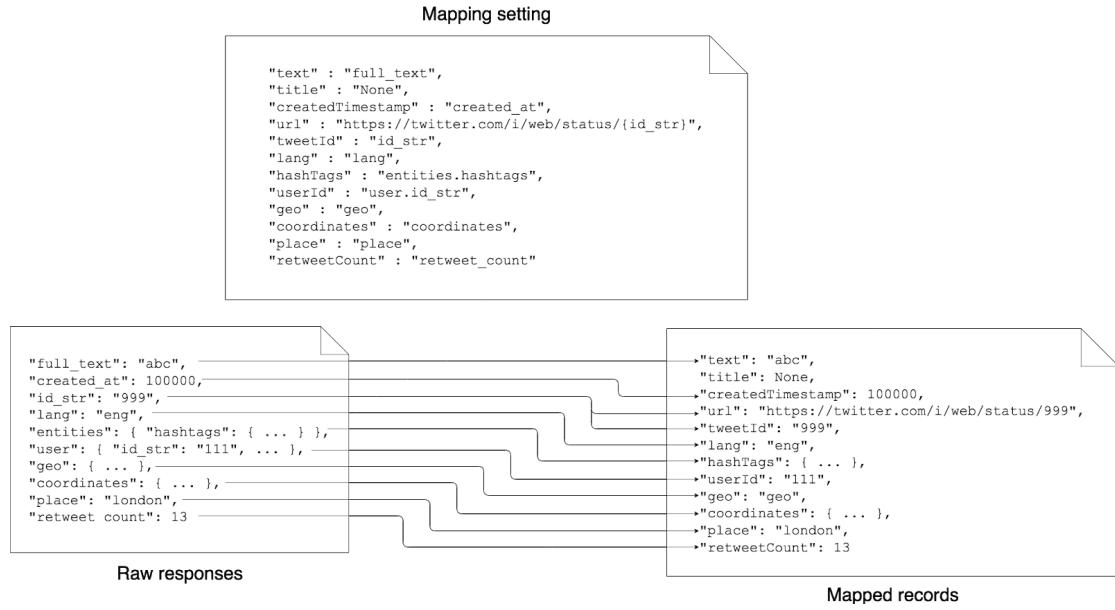


Figure 6: Example of a record mapping according to the mapping setting above

5.3 View

5.3.1 Main views

Tools: Native Play twirl template, Vue.js, Bootstrap and Bootstrap Vue

Overall: A main view is consisted of a template file, JavaScript files and CSS files. The template file is a Twirl template and it mainly has fundamental html components and imports of files and packages. JavaScript files are the main actors. There will be a single Vue application file and the rest would be Vue components reusing throughout the view. The view mechanism is controlled by the Vue application. Figure 7 explains the structure of a main view in term of files.

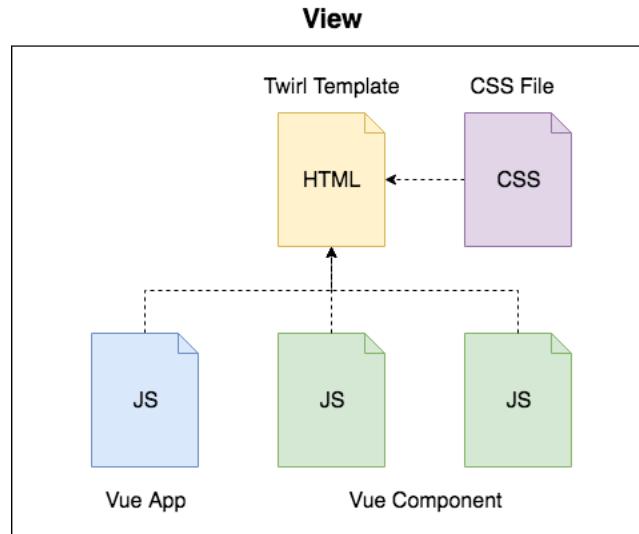


Figure 7: Structure of a main view explained in term of files

Figure 8 shows an example of implemented main views, a result page. All screenshots of implemented main views are in Appendix B and it is highly recommended for the readers to look at the appendix to understand the system clearer.

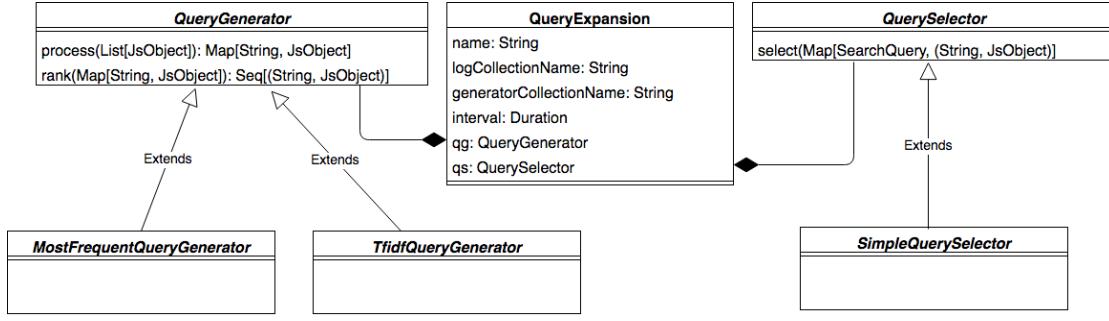


Figure 15: Class diagram of query expansion and pre-implemented class

5.5.3 Scheduler for Crawler and Automatic Expansion

Overall: The crawler and automatic query expansion services are registered with Play module. A module starts right after Play application started, including both registered services. The implementation approach of both services is similar. The system schedules a crawling and expansion task using Akka scheduler. It is similar to scheduling a runnable task in Java but based on Akka, which is a Play official supported concurrent platform [31].

Challenges and solutions: The challenge is a performance issue. The system needs to control number of threads for a task properly to prevent memory leak, which resulted in an internal crash. Burst number of threads could happen because there could be many asynchronous operations within a single task. The solution is a thread blocking. Scheduled tasks that performed within a thread and background operations do not have to be fast. Therefore, blocking a task to perform within a single thread is a good choice to prevent the crash.

6.2 User Evaluation

User evaluation is the evaluation that performs on users, called evaluators, that aims to test main functional requirements of the system. It is an important part of the project since it allows the developers to observe interactions between users and the system. Developers could also ask for feedbacks from the evaluators, qualitative results including likes, dislikes and suggestions.

Overall, the author first listed all implemented M and S requirements and then generated the evaluation tasks from them, grouped by main views. After that, database and materials used during the evaluation were prepared. Next, the author looked for suitable evaluators, and asked them for the meeting date, time and location. The first evaluator was a test pilot for this evaluation. The author used the observations from the first evaluation and suggestions given by the evaluator to improve and change materials that could enhance the experience of other evaluators.

Section 6.2.1 explains the overview of evaluators and their background information. Preparations needed before the evaluation are discussed in Section 6.2.2. All materials that were used during the evaluation are described in Section 6.2.3. After that, Section 6.2.4 describes the steps for the user evaluation. The success measurements are determined in the Section 6.2.5. The result and discussion are shown in Section 6.2.6. Lastly, discovered problems were explained in Section 6.2.7.

6.2.1 Evaluators

There are 6 evaluators. All of them are programmers, they have a certain level of programming experience. Five of them studies here at University of Glasgow while three studies the Data Science, one studies the Information Technology and another studies the Computer Science. Another evaluator studies the Information Technology at University of Reading. Their ages range from 23 to 35. An evaluator from University of Reading do not have any working experience while the rest have worked in IT sectors for at least 2 years.

6.2.2 Preparation

6.2.2.1 Task-based Test Cases

Test cases focus on the evaluation of must-have and should-have requirements and are grouped by main views. They are task-based which means a test case is a single task and evaluators need to complete the given task in order to pass a test case. Completion of a test case is specified by its success measurement. For example, the success of a query deletion task is that the query is removed from the view. The task could be anything that could measure a success. There could be a variable in a task. For example, task A1 asked evaluators to pick a query, named query “A”. They are free to choose any interesting event and represented as a query used in the task. The detail of all test cases is very long, please having a look in Appendix D for all test cases.

6.2.2.2 Data Preparation

The database was mocked in order to provide the same environment for all evaluators. Models were mocked to conform the tasks from Section 6.2.2.1. Models that were mocked are queries, records, sources (crawlers) and query expansion units. List and description of prepared models are as follows.

- Queries: 7 in total
 - Nvidia: represented an announcement of new 2000 graphic card series and used for evaluating querying requirements
 - 2080: expanded from Nvidia
 - Morandi bridge: represented a collapse of highway bridge in Italy and used for evaluating visualisation requirements
 - Morandi victims: expanded from Morandi bridge
 - Collapse: expanded from Morandi bridge
 - Building: expanded from collapse
 - Cave Thailand: represented old query for deletion requirement
- Crawlers: 3 in total which are from Twitter, BBC and CNN sources
- Records: 346 records from 3 sources
- Query expansion unit: 2 in total
 - QE1: using TfIdfQueryGenerator + SimpleQuerySelector
 - QE2: using MostFrequentQueryGenerator + SimpleQuerySelector

6.2.3 Materials

There are numbers of materials for the evaluation listed as follows.

6.2.3.1 Project overview slides

It aims to provide background information of the project. The slides give ideas on introduction, problems, system main features, important terms and the main pages overview.

6.2.3.2 Task sheets

Task sheets give background information, especially intuition to help evaluators understand what they need to do in order to complete tasks. It also improves the continuity of the evaluation by linking corresponding tasks together in the same sheet. A sheet could link 1-2 tasks together, the example is shown in Figure 16 below. The author showed task sheets in a form of presentation that displayed on iPad besides a laptop.

Background: You wonder how query "nvidia" is being tracked in the system and would like to see the number of related records . Since event "nvidia" is popular on Twitter, you find that there are few records from "bbc" and "cnn" sources. So, you are going to turn off "bbc" and "cnn" crawler on query "nvidia" .	
1	Find the number of related records of query "nvidia"
2	Turn "bbc" and "cnn" crawlers off on query "nvidia"

Figure 16: Task sheet example

6.2.3.3 Questionnaire

It aims to quantitatively evaluate the interesting aspects of the system and qualitatively asked for evaluators opinions. It is in a Google Form, which is easy to aggregate results. The questionnaire is divided into 2 sections.

Section 1: Rating

The score ranges from 1 to 4 where one is poor, two is passable/acceptable, three is good and four is excellent. These are the list of rating criteria followed by the interesting aspects of the system. Please have a look at Appendix E to find more description of each aspect.

1. *System overall:* Ease of use, feasibility, directions/guides within the system, and responsiveness
2. *Visualisation [Bar graph]:* Insight gaining from the graph and ease of understand
3. *Visualisation [Timeline graph]:* Insight gaining from the graph and ease of understand
4. *Expansion:* Feasibility and helpfulness of generated information for manual expansion
5. *Evaluation:* Directions and enjoyment

Section 2: Qualitative section

This section asks evaluators about the *likes*, *dislikes* and *suggestions* they have toward the system.

6.2.4 Evaluation Steps

1. The author gives a presentation on project overview
2. Starts the evaluation by
 - a. The author shows a task sheet that comprises of the task background information and list of tasks that an evaluator needs to complete
 - b. The author asks the evaluator to read and understand the task thoroughly and he/she could ask questions about the task if it was not clear
 - c. Let the evaluator starts working on a task until completed, while the author records time spent for the task
 - d. Starts evaluating next task by repeating step a to c until all tasks are completed
 - e. The author only gives a suggestion only when the evaluator is misled
 - f. Task is incomplete if the evaluator gave up

6.2.5 Success Measurements

As mentioned in Section 6.2.2.1, the task would be marked as completed if the evaluator successfully completes the task. Task will be failed if the evaluator cannot find the way to meet the success measurement of that task and gave up. An evaluator would successfully pass the evaluation if he/she has at least 80% pass rate, calculating from all tasks.

A criterion in the first section of questionnaire would be considered as success if the average score of that criterion from all evaluators is higher than 2.5 (out of four).

The evaluation is successful if there are at least 80% of passed evaluators and 80% of passed criteria in rating section.

Expansion: They believes that the expansion feature would be practically worked, which is a promising sign.

6.2.6.3 Questionnaire: Section 2

The summary of likes, dislikes and suggestions are listed as follows. However, the suggestions from the questionnaire is too short and not descriptive. So, the author also adds the recorded suggestions given by the evaluators during the evaluation into the list and grouped them by tasks.

Likes

- The system could get reliable news from many sources
- Query expansion can be well explained by node graph, which reveals potential relations between events
- The system would be good for many developers
- Query expansion feature
- Timeline graph
- Idea of expanding query from time to time

Dislikes

- The system is somehow lag at the beginning of the evaluation
- There is only connected nodes and edges. There is no text label or each relation display on the graph as it requires user to click see in detail each by each
- There are some unclear UI, especially icons that are hard to understand
- Some evaluation tasks are complicated and have complex interactions
- The location of query setting
- There are lots of text and sometimes it is hard to find what is needed to do

Suggestions

- A4: Remove irrelevant record(s)
 - It is good to have duplicate filter or retweet filter
- A5: Disable a “twitter” crawler and search for the event “A” again
 - Icon is hard to find, and gear icon is not appropriate
 - Icon is too dark, gear is not appropriate, and some keywords are needed
 - Query setting could be moved to the left (on filter panel) and gear icon looks dangerous
 - Query setting could be moved to the empty right panel
- B1: Find the number of related records of query “nvidia”
 - It would be great to show number of related records in check query too
 - Total records could be moved to the first line instead
- B2: Turn “bbc” and “cnn” crawler off on query “nvidia”
 - Query card could be collapsible to save more spaces and see more queries
- C1: [Bar graph] Find the number of records from each source on 18/8/2018
 - The bar graph is not easy to use
- C2: [Timeline graph] Identify the records that have multiple links and interpret their relationships between queries
 - Node type in info box is too big, legends are needed for nodes, and also indication of colors (source)
 - Query nodes can be represented with pictures/symbol/text
 - Legends for node types in timeline graph are needed, also for arrows (expansion)
- D1: Add new crawler to the system
 - Newly added crawler could be added at top of the list, instead of bottom

- D2: Figure out when is the next crawling time of a specific crawler
 - Highlighting next crawl and last crawl attributes could make it clearer

6.2.7 Problems

This section highlights interesting problems occurred during the evaluation.

- The evaluation on a first evaluator was the most troublesome. It was like a test pilot. The materials, especially the task sheets, were ambiguous. The direction of a task was not descriptive enough. In addition, the overview slides did not provide enough information for the first evaluator to help him understand the system. A good example would be the selection of a query “A”, mentioned in Section 6.2.6.1. He picked “Ed Sheeran”, which is not actually represented any specific event at that time. Therefore, before evaluating the rest, the author refined all materials again by adding more information to the overview slides such as important terms and definition of a query, highlighted the task with colors, changed some ambiguous words and phrases and rehearse more for the overview slides.
- There might be too few records for the evaluation. The system is designed to collect a lot of data so the visualisations could reveal potential insight from them. An evaluation of the bar graph was affected by this problem. Bar graph had a weird behavior during the evaluation. The time axis (x) was not separated by day (1/9/2018 following by 2/9/2018 and so on). But, it was separated by day-noon time-day order (1/9/2018, 12:00 of 1/9/2018, 2/9/2018 and so on). As a result of this, task C1, which required evaluators to count number of all records on a specific date, was harder than expected.

6.2.8 Conclusion

In summary, the evaluation was *successful* by meeting both 2 statements, *having at least 80% of passed evaluators and 80% of passed criteria*. The evaluation process was smooth as a result of efforts from all evaluators. Identified problems in the above section could also be considerations for other evaluations of the similar projects in the future.

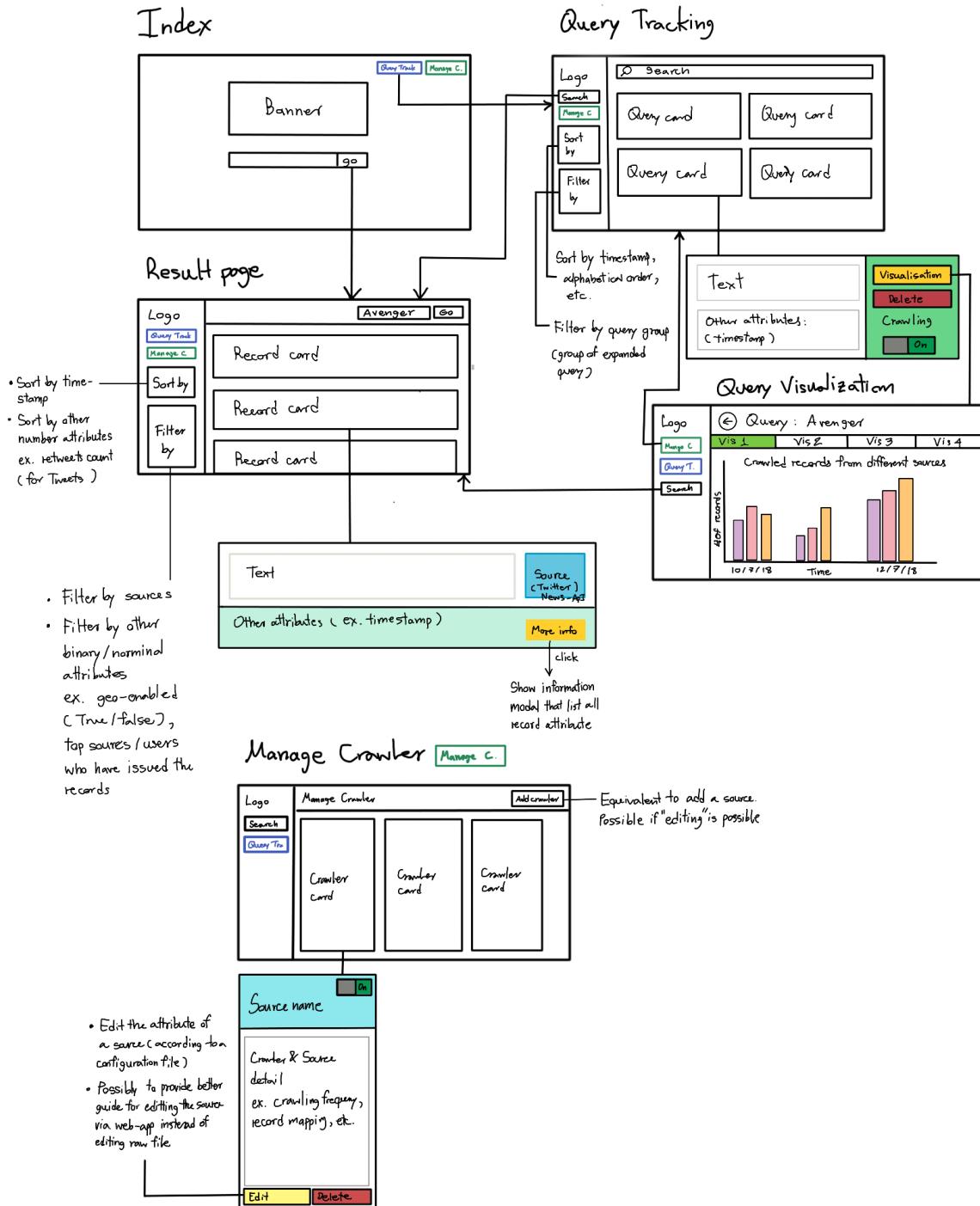
Chapter 7 Conclusions

Feed-me! is the event data aggregation web application that collects data from different sources together. It features customisable crawlers, query expansion and visualisation, that allows users to draw more insight from the data. The system met all must-have and should-have requirements, which is about 76% of all requirements. Although no requirement from could-have and won't-have were met, the system still works with its full potential to achieve the main goal. System overview could be identified very early as a result of initial view designs and all designed main views are successfully implemented. Starting the system design from views provided a high-level understanding and it could be used as a great medium for a discussion with stakeholders. There are a lot of detail on the implementation. Especially when building a web application, many tools were used in different components. Integration of Vue.js, Bootstrap and Bootstrap Vue significantly helped and boosted the development of the frontend. Visualisations were completely developed with the aid of visualisation frameworks. Even Play does not officially provide APIs for NoSQL database connection, the database driver like ReactiveMongo handles this task perfectly. However, the knowledge of asynchronous programming style is heavily used to achieve the fast and non-blocking access to the database. Learning the fundamental knowledge of Scala asynchronicity is highly recommended before working with the driver to prevent the confusion and further unknown errors. Query expansion was developed with the idea that allows users to develop their own expansion technique by extending the generator and selector abstract classes. The current pre-implemented classes for generator and selector might not have sophisticated and highly effective logics. But, this architecture gives the promising future that could support more advanced techniques. The user evaluation result is promising. The feasibility of query expansion was questionable at the beginning, however all users believed that it could be worked in practice according to the user evaluation. On the other hand, guides and directions within the system should be improved hence allowing users to learn how to use the system without the need of any user training beforehand. In summary, the evaluation was successful by having all evaluators passed the evaluation and 92% of passed criterion.

There are many interesting could-have and should-have requirements that could enhance existing features. Grouping queries together in term of query expansion (requirement A7) is a good idea since this could improve the query classification in timeline graph, also in query tracking page. The idea of keeping images (requirement A12) of raw responses should help enriching the content displayed together with the records in result page. In addition, URLs within the records would be beneficial to the system since they could be used to extract more information from the source of URL, or even extracting new queries out of it. Crawler is still too simple in term of the resource utilisation. It should manage limited rate limits more efficiently (requirement B8) by computing crawling period from the available rate limit. An important feature of the query expansion that is still missing is adding query expansion unit from frontend. Users could conveniently compose their query expansion unit from manage QE page if there are many implemented generators and selectors available. In addition, it could be good to have more pre-implemented generators and selectors. The frontend needs minor changes on user interfaces and components position, many suggestions from the evaluation were given. By doing this, it would significantly improve overall user experiences. The missing feature for visualisation was the indications, like legends for timeline graph, however the author had already implemented this in the final version as a last development task.

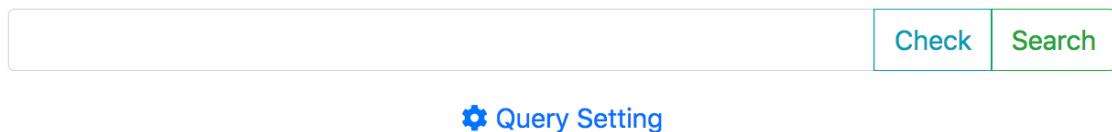
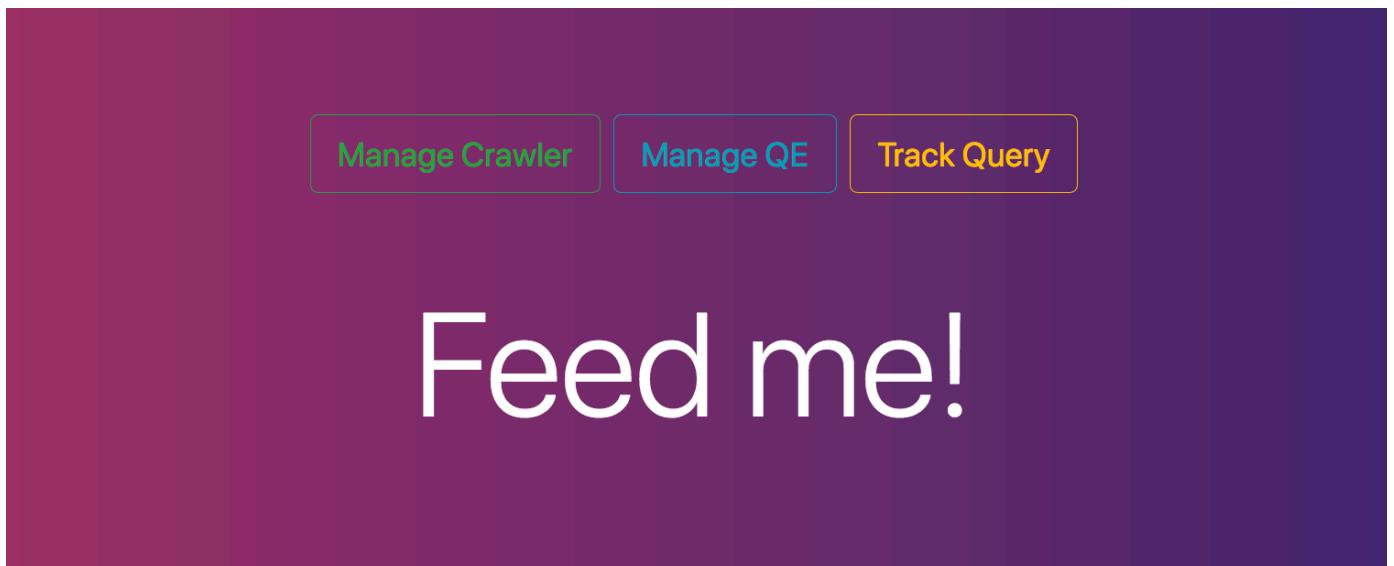
Timeline graph should handle more records and queries in the future, and at the same time approach to improve the performance of rendering is also necessary. Lastly, the export feature could benefit the uses of gathered data. Users might be able to export selective portion of data in several formats such as JSON.

Appendix A Initial UI Design



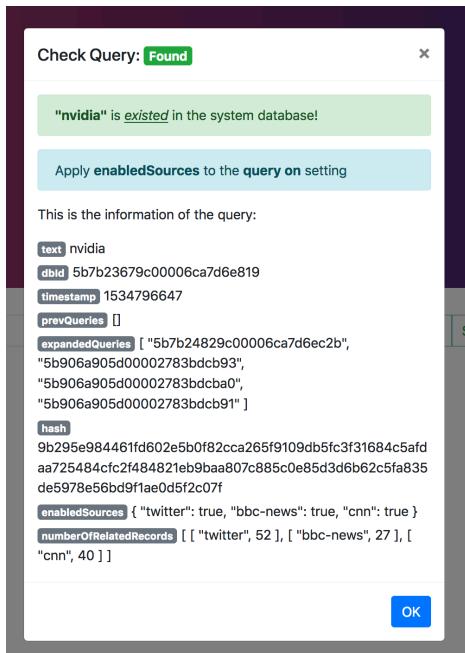
Appendix B Screenshots

Index page



Index page

Users can navigate to other main views using the top navigation, checking if the query is existed in the database using check button, change query setting from the blue hyperlink and querying by entering text in search box then click search.



After clicking check query button, the modal appears. It shows the status of that query and some of related information.

Check modal

Result page

The screenshot shows the 'Feed-Me!' search interface. On the left, there's a sidebar with navigation buttons (Track Query, Manage Crawler, Manage QE), visualization options (nvidia, Visualisation), filtering by source (All, cnn (40), bbc-news (27), twitter (66)), and ordering (Title, Text, Created time). The main area displays search results for 'nvidia' with a total of 133 records (1-10 shown). Each record is presented in a card format with a delete icon. The first card contains a tweet from 'twitter' with the title 'No title'. The text of the tweet discusses Parabricks, an NVIDIA Inception startup, and how GPUs can shrink DNA analysis time. The second card is similar, mentioning GTC18 robotics innovations. The third card is about a NVIDIA GeForce GTX 1070Ti video card. Each card has a 'More info' button at the bottom right.

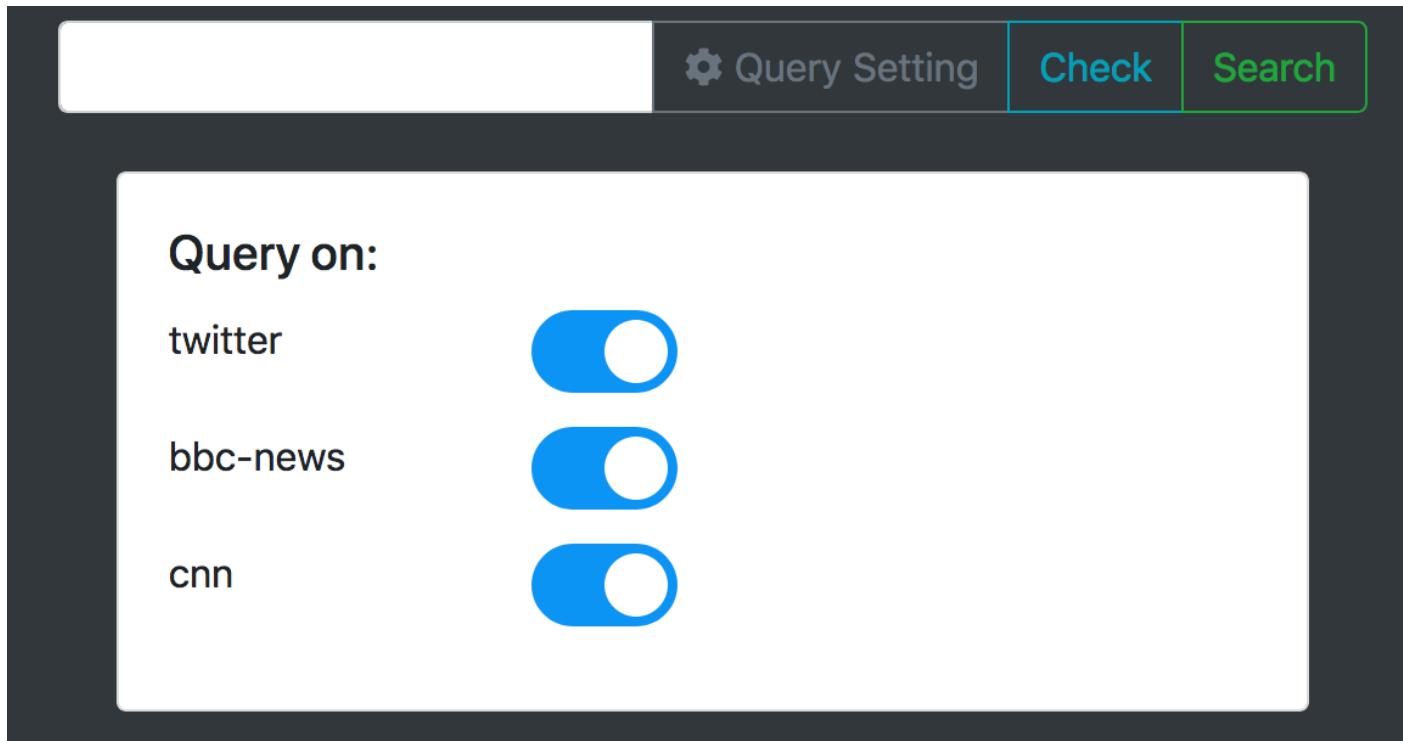
Result page

After querying, the result page is shown. There is a side bar on the left, consisting of navigation buttons, button linked to the visualisation page of a query, filtering and ordering menu. It shows the related records, each in a card, of the query. Records are queried from a new search and existing records in the database. Users can delete a record if it is duplicated or irrelevant and can click at more info button to see more information of that record. More info modal is shown after clicking.

This screenshot shows the same 'Feed-Me!' interface as above, but with a 'More info' modal window open over the third search result card. The modal is titled 'More info' and contains a detailed JSON-like list of the tweet's metadata. The list includes fields like text, title, createdTimestamp, url, tweetId, lang, hashtags, userId, geo, coordinates, place, retweetCount, and concatValue. It also lists hash and query values. The modal has a close button in the top right corner.

More info modal

Search bar



Search bar

Search bar appears on the top of all pages, except index which is located in the center. Users could enter a query and click check or search from here. In addition, it shows the list of crawlers by clicking at query setting. This allows users to alter the statuses of crawlers of that query.

Track query page

Text: morandi bridge

Number of records:

- twitter: 18
- bbc-news: 10
- cnn: 9

Total records: 37

Created at:
20th August 2018, 7:05:49 pm (19 days ago)

Crawler

twitter	<input checked="" type="checkbox"/>
bbc-news	<input checked="" type="checkbox"/>
cnn	<input checked="" type="checkbox"/>

Text: morandi victims

Expand **Visualisation** **Delete**

Track query page

This page lets users see the list of tracked queries listed each in a card. Users can see number of records, creating date and time, and the query setting (annotated by crawler). They can also delete a query, see the visualisation page of that query and perform manual expansion on the query. By clicking expand button, the manual expansion modal will be shown. Users will see the available query expansion units listed in the tabs. They could select any of them to see helpful information generated by the unit at a specific time, showing ranked terms with their text statistics. Additionally, they can force the unit to generate information for them now. To perform a manual expansion, users enter the expanded query and click add.

Text: morandibridge

Number of records:

- twitter: 18
- bbc-news: 10
- cnn: 9

Total records: 37

Created at:
20th August 2018, 7:07:58 pm (19 days ago)

Crawler

twitter	<input checked="" type="checkbox"/>
bbc-news	<input checked="" type="checkbox"/>
cnn	<input checked="" type="checkbox"/>

Text: morandibridge

Number of records:

- twitter: 13
- bbc-news: 7
- cnn: 4

Expand query

QE1 QE2

Query Generator Class: TfIdfQueryGenerator **Generate now**

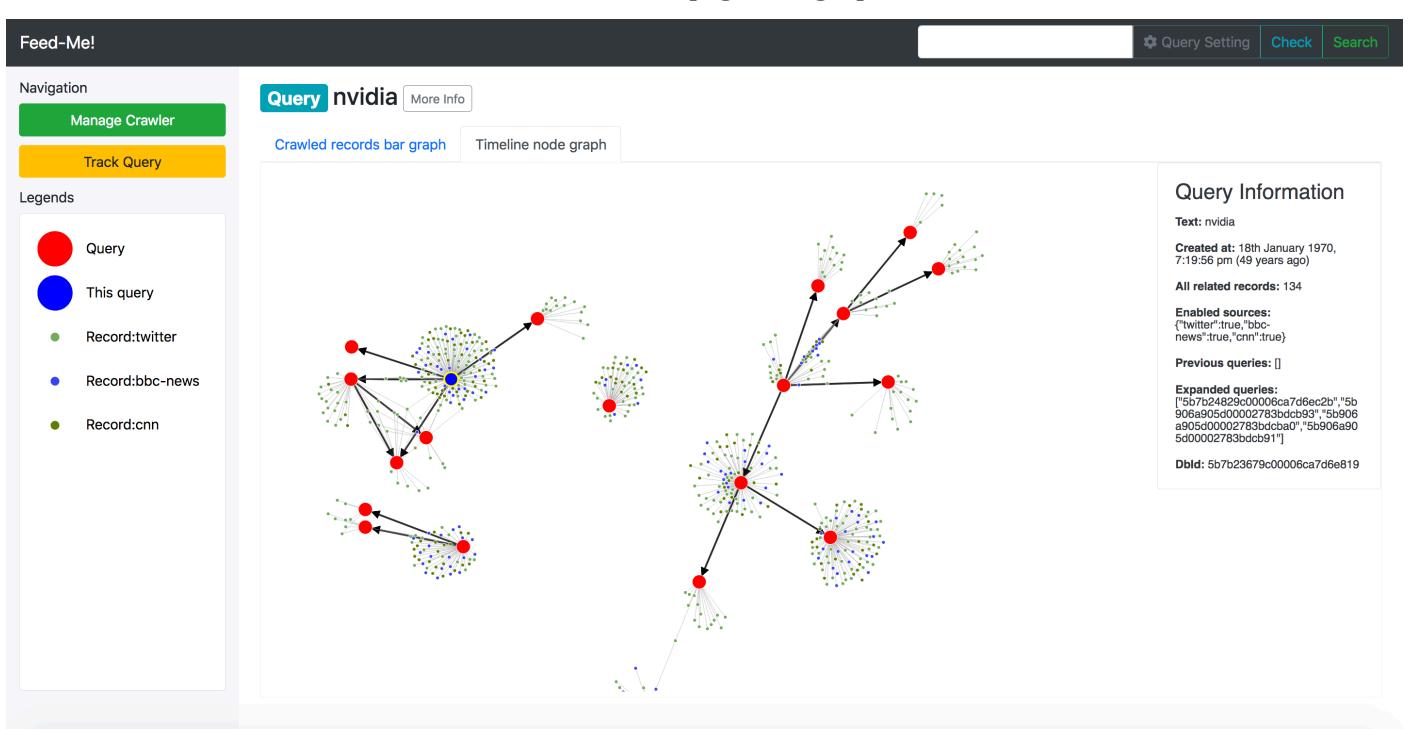
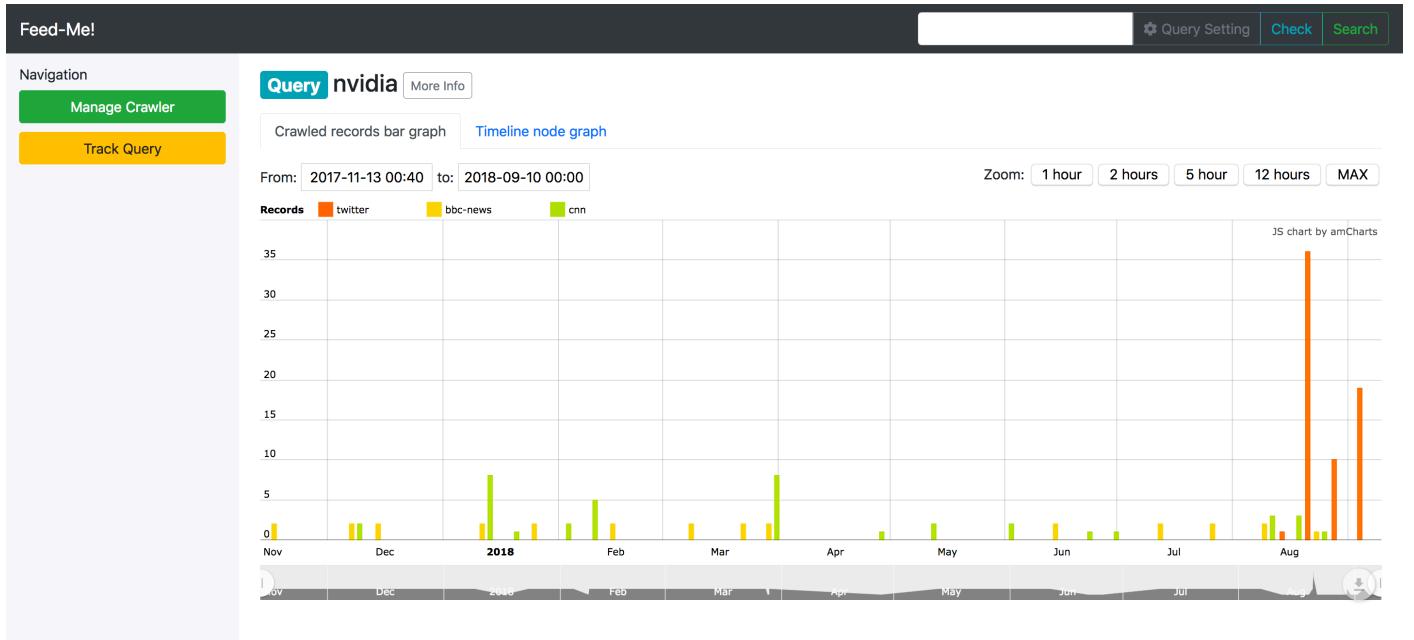
Query Generator Logs: Ranked Terms

CreatedAt:	state	victims	funeral
20th August 2018, 7:07:58 pm (19 days ago)	tf-idf: 12.648 count: 9 idf: 1.405	tf-idf: 11.266 count: 7 idf: 1.609	tf-idf: 11.075 count: 6 idf: 1.846
	families	italian	will
	tf-idf: 10.658 count: 5 idf: 2.132	tf-idf: 9.657 count: 6 idf: 1.609	tf-idf: 9.657 count: 6 idf: 1.609
	collapsed	collapse:	
	tf-idf: 8.432 count: 6 idf: 1.405	tf-idf: 8.047 count: 5 idf: 1.609	

Enter expanded query here and click add **Add** **Close**

Expand query modal

Visualisation page



After clicking at the visualisation button in a result page or track query page, the users will be navigated to visualisation page of that specific query. There are 2 graphs available. Users can change the graph by clicking the tabs above. For timeline graph, there is a legend description on the left. As stated in the main report sections, the graph shows relationships between queries and records. They can also drag and zoom the graph if they want. When clicking at a node (any query or record), information box will appear on the top right of the graph.

Manage crawler page

Appendix E Questionnaire section 1: Rating

1. System overall
 - a. *Ease of use*: Was the system easy to use?
 - b. *Feasibility*: Do you think that the system would be worked in practice?
 - c. *Directions/guides within the system*: Do you think that there are plenty of guides that help you easily understand the system?
 - d. *Responsiveness*: Was the system fast enough to response you as expected? Do you experience any lags?
2. Visualisation: Bar graph
 - a. *Insight gaining from the graph*: Do you feel that the graph provides useful information vs. raw data?
 - b. *Ease of understand*: Was it easy for you to understand the graph?
3. Visualisation: Timeline graph
 - a. *Insight gaining from the graph*: Do you feel that the graph provides useful information vs. raw data?
 - b. *Ease of understand*: Was it easy for you to understand the graph?
4. Expansion
 - a. *Feasibility*: Do you think that this feature would be worked in practice?
 - b. *Manual expansion*: Did generated information help you decide the query terms for manual expansion?
5. Evaluation
 - a. *Directions*: Do you think that the evaluation procedure is connected flawlessly?
 - b. *Enjoyment*: Did the evaluation entertain you in any aspect?