

Read the news you want to read

Authors: Tom Dooner, Mika Little, Brian Stack
Project: Project October
Date: April 20, 2013

Contents

| | |
|--|-----------|
| 1 Abstract | 4 |
| 2 Project October: Recommender System | 5 |
| 2.1 Background | 5 |
| 2.2 Intended Audience | 5 |
| 2.3 Progress Since Progress Report 2 | 5 |
| 2.4 This Report | 6 |
| 3 Software Design Requirements | 6 |
| 3.1 Homepage | 6 |
| 3.2 User Creation, Login, and Display | 6 |
| 3.3 User Profile Page | 6 |
| 3.4 User Profile Edit Page | 6 |
| 3.5 Article Posting | 7 |
| 3.6 Article Display Page | 7 |
| 3.7 RSS Subscription | 7 |
| 3.8 Frontend ↔ Backend API | 7 |
| 3.9 Backend | 7 |
| 4 Application Architecture | 7 |
| 4.1 Frontend | 7 |
| 4.1.1 RSS Subscription | 8 |
| 4.1.2 Homepage Article Display Sizes | 8 |
| 4.1.3 Homepage Article Layout Algorithm | 9 |
| 4.1.4 Backend Proxy | 10 |
| 4.1.5 Sending Email to Users | 10 |
| 4.1.6 User Tracking + Analysis | 10 |
| 4.2 Web Scraping | 10 |
| 4.3 Backend | 12 |
| 4.3.1 Vectorization of Posts | 12 |
| 4.3.2 Vectorization of Users | 12 |
| 4.3.3 Tf-Idf | 12 |
| 4.3.4 Signals | 12 |
| 4.3.5 Time Scaling | 12 |
| 4.4 Frontend/Backend API | 12 |
| 5 Project Management–Administrative Details | 12 |
| 5.1 Release v1.0 | 13 |
| 6 Testing & Evaluation | 13 |
| 7 Project Progress | 13 |
| 7.1 Discussions & Conclusions | 13 |

| | |
|---|-----------|
| 8 Lessons Learned | 15 |
| 8.1 Advantages & Disadvantages of a Service Oriented Architecture | 15 |
| 9 Appendix | 15 |
| A Thrift API Definition | 17 |
| B Database Design | 18 |
| B.1 Frontend Database | 18 |
| B.2 Backend Database | 18 |
| C User Manual | 18 |
| D Programmers Manual | 18 |

List of Figures

| | |
|--|----|
| 1 Webpage Sections | 11 |
| 2 Project October Gantt Chart | 12 |
| 3 Project October Homepage, pre-Alpha (Headlines courtesy The Onion) | 13 |
| 4 Project October Homepage, version 1.0 | 14 |
| 5 Frontend Database Relational Diagram | 19 |

1 Abstract

Longtime users of modern news aggregation services (such as Reddit, Slashdot, Digg, and Hacker News) report noticing a marked decrease in quality of discourse as the services gain mainstream attention. This gradual, irreversible decline was noted as early as the newsgroup era, when new college students would log in for the first time in September, causing an influx of new users and subsequently diluting discussion upon the sites which they joined. The situation worsened after AOL opened newsgroups to the masses, causing the community standards to continue to devolve, according to newsgroup veterans[1]. This gradual decline of content quality due to large numbers of new users came to be known as “Eternal September”, being “Eternal” as the influx of new users was no longer restricted to September, but persisted through all months of the year. This phenomenon indicates that balancing a large community with high quality content is a challenge.

October is a solution to this problem. By providing users with an automatically-customized experience, we believe we can provide a large community thoughtful discourse and interesting articles while avoiding the effects brought on by Eternal September. October will employ a technique to recommend news articles which, to our knowledge, is not currently applied on any other major news aggregators. With custom recommendations for each user, we believe that the October community can attain large size without a sacrifice in quality.

2 Project October: Recommender System

Due to the ever-growing amount of information available online, the need for a highly developed personalization and filtering system is growing significantly. Recommender systems constitute a specific type of information filtering that attempt to present items according to the interests expressed by a user. Most web recommenders are employed for e-commerce applications or customer adapted websites, which assist users in decision making by providing personalized information, but the same techniques that suggest related items on e-commerce websites can recommend news articles to users as well. We believe Project October is the first attempt to apply recommendation techniques to social news aggregation.

2.1 Background

It is our hypothesis that the recommendation of news sources will provide a scalable community experience that can be tailored to each person's interests. Providing an automatic, customized, recommendation of articles will prevent the community from being diluted by new users and thus prevent the Eternal September phenomenon. Each user will have their own viewing environment curated for them, with different mindsets and preferences forming their own communities in which like-minded individuals can partake in discussion, eliminating the cross-contamination of user bases.

2.2 Intended Audience

Project October will be open to the public. At launch, we will be inviting our friends from the CWRU community to try it out first, so we expect the user base to be somewhat technical in nature.

2.3 Progress Since Progress Report 2

Since Progress Report 2, we have implemented all our design requirements, released a beta version to the CWRU EECS department, and collected user metrics. More details about the v1.0 minimum viable product can be found in Section 5.1.

For the sake of releasing a minimum viable product within the semester, we have removed some complexity from our requirements in Progress Report 2. Specifically, we have made the following changes to the report, or the scope of our project:

- We have specified the plan for the backend recommendation engine in tangible terms.
 - We have removed the ability to upvote, downvote, and tag comments. Those abilities will be easy to add later but are not core to the application, especially until October gains critical mass and more people comment.
 - We have added an RSS Subscription feature, specified in Section 3.7 and described in Section 4.1.1.
 - We have made numerous copy tweaks to the report, such as
 - Reordering the report to list the requirements before the details of the project's implementation
 - Adding Sections 4.1.1-4.1.6 about the frontend and Section 7 about the overall final status.
 - Moving the User Interface section into the requirements.
 - Describing our testing process in Section 6
-

2.4 This Report

3 Software Design Requirements

3.1 Homepage

There are two cases for the homepage: logged in users and logged out users.

Logged out users arrive at the homepage and are presented with a splash page containing a description of October's features. There is a registration link so that interested users can register for their own account.

Logged in users are, upon landing on the homepage, presented with a display of news articles. Articles are displayed with pictures (when available) along with headlines. Each article is accompanied with links to positively or negatively influence the article's recommendation rank, and also to comment on the article. The source of the article is also provided, as is the time that it was posted.

The homepage is an especially crucial interface to get right. October should mimic a newspaper feel with articles staggered in a rough column grid. This will allow readers to be presented with many recommendations of different weights simultaneously in a manner in which they are accustomed to. Articles will have two links (upvote and downvote) into the article area, along with the link to that article's comments.

The search bar, along with a 'Post An Article' button, a link to the user's profile page, and a Logout link will be located at the top of the site.

3.2 User Creation, Login, and Display

New users are able to create October accounts with but three pieces of information – their desired username, their email, and a password. After creating an account, an email is sent to the user welcoming them to October.

When a user first logs into October, the backend recommender won't know anything about their interests. Thus, users should be shown various sample topics to browse through. They are also able to (via their account page) manually type in relevant keywords to aid the recommender in providing them with articles.

Every mention of a user's name (i.e. the credit for who posted an article or comment) is a link to a user profile page. More detail about the user profile page is given in Section 3.3.

3.3 User Profile Page

Each user has a profile page containing information related to that user. The profile page contains a list of comments that the user has made on posts and a list of articles the user has posted. All information on this page is public, as it is merely a summary of content available elsewhere on the public website.

3.4 User Profile Edit Page

In contrast to the profile page, which is a summary of a user's public activity, the profile edit page allows a logged in user to edit their own information. Settings such as their email address and password can be changed here.

The other significant part of the user profile edit page is a section which shows the user their own recommendation settings. This shows up in the format of a list of keywords that the backend has ascribed to that user. The user can remove keywords that they do not wish to be associated with.

Furthermore, there should be a field where a user can submit additional keywords they wish to see more news from. The user will be associated with the keyword when they submit the form.

3.5 Article Posting

Logged in users are able to click a link on the homepage to submit a story. The submit story page will ask for the link to the desired story, and automatically show the results of scraping that story's contents (see Section 4.2 for more about scraping). If the post looks good to the user, he or she can click a "Post!" button, which will persist the news story in all appropriate databases.

3.6 Article Display Page

The article display page is a page dedicated to a posted news article. Every news article has an associated article page. The intent of the article page is to facilitate sharing of opinions around the posted news story through comments.

A user arriving on an article's page arrives at a page which contains the title of the news article along with a threaded comment list. Users may leave comments as children of other comments, and thus may engage in a lively discussion.

3.7 RSS Subscription

Users should be able to subscribe to RSS feeds by entering the URL to a feed in their User Profile Edit page. This will allow users to customize the news sources present in October without painstakingly submitting things until the community grows larger.

Subscribing to an RSS feed will automatically post the articles from the feed to October in a timely manner. Users should be able to click the feed name on an article displayed and be taken to a list of posts from that RSS feed (much in the same way that users can see a list of another user's posts on the user profile page).

3.8 Frontend ↔ Backend API

Since the frontend and backend are divided with a Service Oriented Architecture (with Apache Thrift bridging the gap¹), there must be an API between them.

The API should be as simple as possible to maintain the functionality of October. This is not a public API, so it needs to be robust to volume but we do not need to implement API endpoints for the functionality mentioned elsewhere in the requirements.

The API is defined in a Thrift definitions file in the project-october-api repository on GitHub[8].

3.9 Backend

The backend is a recommender black-box which provides suggestions via an API to the frontend. The actual implementation of the recommendation service is arbitrary, only that it implements the services of the API effectively. For more about the implementation of the backend, consult Section 4.3.

4 Application Architecture

4.1 Frontend

The frontend is similar in function to a standard social news aggregator. Upon going to the homepage, new and logged-out users are presented with a splash page that describes the nature of Project October and are

¹See Section 4.4 for more about Thrift

given a registration link. Users can register for the application by selecting a username and password. A registration confirmation is emailed to newly registered users.

When logged in, users are presented with their personalized content. The design is modelled after the front page of a newspaper – articles that are evaluated to be more interesting to a user are placed in more prominent positions while articles that are deemed less relevant to the user's interest will fill the side columns, be located further down the page, and occupy less space. Users are able to click on the headline (or associated image, if existent) to be taken to the original news source.

Each news article features a link to view the associated comments as well as icons which allow the user to plus and minus the article. These are analogous to upvoting in Reddit, except that plusing or minusing an article on October will inform the recommendation engine of your individual preference, rather than directly impacting the weighting of an article by a predefined formula. Therefore, voting serves two purposes: to inform other users of the quality of the article and to inform the recommender system of the user's interests.

Users will be able to click a link on the homepage which takes them to an article submission form. The submission form requests a URL and a headline for that news story. Upon entering the URL, the news story is scraped in the background using a web scraping tool which attempts to extract properties about the page automatically – relying upon the `meta` tags and heuristics to determine the article's body text. The user is also presented with a list of tags that are garnered from the article, and they can delete or add their own tags to better represent the content of the article. A list of scraped images is also presented, and the user may select which one (if any) they want to head the article with.

4.1.1 RSS Subscription

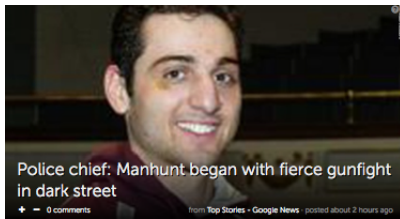
In order to implement the RSS Subscription feature (described in Section 3.7), we created a background task that scrapes each user-submitted RSS feed and posts each RSS entry as a new story. Since this means there are two types of entities that can post articles, and we wish to share logic between both, we use ActiveRecord's Single Table Inheritance[9] to make Users and Feeds share attributes in a single table **posters** (see the database schema in Figure 5). The most important reason for this is because Users and Feeds must share the keyspace for their primary key (their `id` columns) so the backend can treat them the same when dealing with following relationships (a user should be able to follow both a user and a feed, under the same logic).

This will also allow other types of posting entities in the future if we decide to add additional integrations.

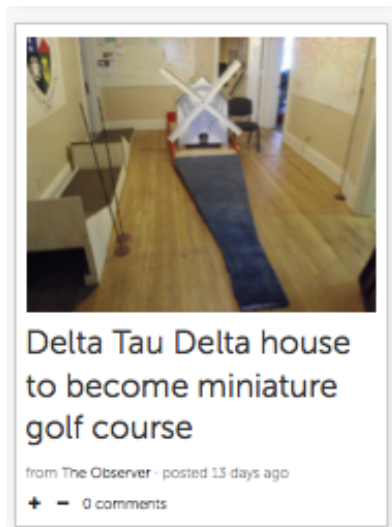
4.1.2 Homepage Article Display Sizes

According to the specification (in Section 3.1), the homepage should display articles in a visually appealing fashion. Towards this end, we split the homepage into three equal-width columns and display each article with one of the following types:

- featured type articles take up two columns and contain a large image from the post. The headline is large and stands out. We intend these featured articles to be the eye-catching cornerstone(s) of the homepage.
 - primary type articles are displayed in a white box to emphasize them from the background. The headline is large and flows with the picture to create the newspaper feel.
 - secondary type articles are displayed with a grey headline and grey link text, to blend in with the background and offer the reader additional options if they desire to look harder to find them.
-



featured article (scaled down to 50%)



primary article



secondary article

primary and secondary articles may take up any amount of vertical space, but are constrained to one column in width. If the homepage is not wide enough to fit three columns (on clients of less than 768px), the homepage goes into single-column mode where all articles are displayed linearly.

4.1.3 Homepage Article Layout Algorithm

The article layout algorithm takes place in Ruby (the backend of the frontend) and determines the size and approximate position of the articles to be displayed to the user. The algorithm is complex enough to make intelligent decisions, yet simple enough that it can be performed quickly and repeatedly. The algorithm is, given a list of n article IDs and weights for each article from the recommender, take the following steps:

1. Sort the article list by weight, so the highest-weighted article is first.
2. Take the top two highest-weighted articles with pictures. Assign them the type featured.
3. Take the remaining $\text{round}(\frac{n}{2}) - 2$ items and assign them the type primary.
4. Assign the remaining items the type secondary.
5. Introduce some randomness by naïvely grouping all posts into four categories by ID.
6. Insert the featured articles into the first and the $\frac{n}{4}$ th indexes of the list.
7. Return the list of $(item, type)$ tuple to the frontend.

The frontend will output HTML for each item and type in the order returned from the algorithm. The CSS for the homepage contains basic styling instructions, but not enough to create the desired homepage effect. For that, the user's browser will invoke the JavaScript masonry.js library[10]. Masonry.js takes a set of arbitrarily-sized DOM nodes and displays them in a compact way, thus giving us the desired newspaper front-page effect.

Using custom responsive CSS and JavaScript event handlers on window resize, we ensure that the newspaper feel is consistent and usable on any size screen and when the screen is resized after the page originally loads.

4.1.4 Backend Proxy

Since the connection to the backend is via a web service, the frontend cannot always assume a working connection to the backend. Towards this end, we have written a Ruby proxy around Thrift which keeps the frontend function even in case of problems with the backend. It will try again to establish the connection in case of momentary connection break, and will switch to a mock backend if the real backend still does not reply. The mock backend is API-compatible with the real backend, but is not capable of making informed decisions so we display a message to October users that October is operating at diminished capacity (“amnesia mode”).

4.1.5 Sending Email to Users

October uses Ruby on Rails’s packaged ActionMailer to send emails. We use the premailer Ruby gem to speed up development – it allows for the use of real stylesheets which are then compiled inline to the appropriate HTML elements of the emails before they are sent.

Currently, the only email that is delivered is a welcome email for the user offering basic onboarding instructions. However, in the future, we intend to engage users via email with customized recommendation emails (i.e. “We just found this article that we know you will like!”)

4.1.6 User Tracking + Analysis

In order to understand how effective our users find October, we have instrumented the site with a heavy dosage of analytics. To track basic page views and real time events, we have installed Google Analytics. But, Google Analytics does not provide the granularity that we would need to adequately understand user behavior. To this end, we have instrumented the site with Mixpanel[7], which specializes in tracking user-triggered *events* rather than simple pageviews.

Mixpanel helps us track user behavior by tracking an event every time a user performs one of the following actions: log in, sign up, receive recommendations, vote, click article headline, or post a new article.

We currently use Mixpanel in a couple ways beyond simply tracking user events, however. Primarily, we track how many milliseconds the user waits to receive recommendations from the backend. We also track the direction of users’ votes, so we can tell how satisfied users are, at a high level, with their article recommendations.

4.2 Web Scraping

We use a library called Pismo [4] to scrape webpages. While this may seem like a trivial operation that we should implement ourselves, there is much that must be done in order to determine the actual content of a webpage. We have selected to use the “cluster” approach Pismo provides to determine which parts of the webpage are actually content. This approach ranks each <div> element on a page according to the following heuristics:

1. element must contain text
 2. earlier elements are more likely to be in the body of text
 3. each element must be over some minimum length (80 characters by default)
 4. more text in an element and more “.” delimited sentences raise chance of being content
-



Figure 1: Sections of webpage colored by desirability.

Green Sections are best for determining content of article. Also includes potential images to use as a headline image for the article.

Yellow Sections may contain useful information but are not part of the body of the article.

Red Sections do not contain useful information and will harm results if included.

This works well for our purposes. It is possible that as time progresses, we can make use of HTML5 semantic tags such as `<article>` in order to more accurately extract text.

Once this text is extracted, Pismo relies upon the Phrasie [6] library to determine keyword tokens. The initial version of Pismo we used was based upon a naive list of (token, frequency) pairs, however Phrasie uses parts-of-speech to determine keywords and gives them a weight rather than a straight frequency. The strength of this approach is that it results in a much higher signal-to-noise ratio than the naive approach. A weakness is that we've now reduced our possible languages down to just english. As time goes on we can modify this approach to either support more languages or just work in a general sense on any language.

We also make use of Pismo's ability to extract images on the page that have a high likelihood of being content rather than UI elements.

Figure 1 gives a visual representation of the content on a page we want to scrape.

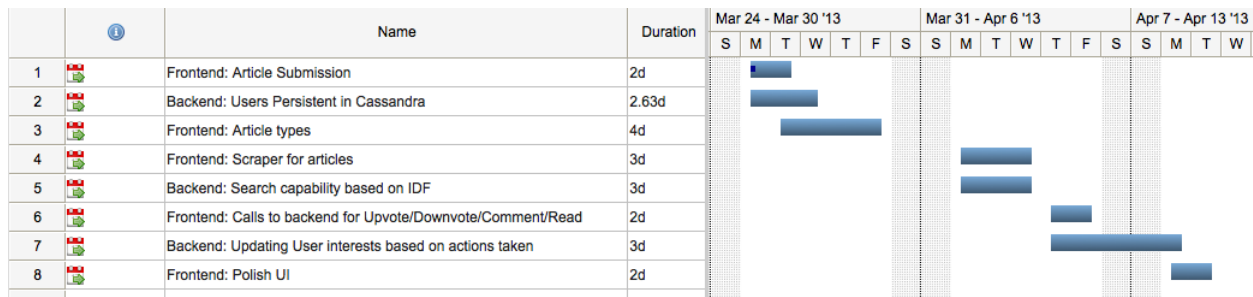


Figure 2: Project October Gantt Chart

4.3 Backend

4.3.1 Vectorization of Posts

The scraper from section 4.2

4.3.2 Vectorization of Users

4.3.3 Tf-Idf

4.3.4 Signals

4.3.5 Time Scaling

4.4 Frontend/Backend API

October will employ an API to promote separation between the frontend and the backend recommender system. This will provide a clear interface and facilitate easy simultaneous development of both parts of October.

To implement the API, we will use Apache Thrift, an Interface Description Language[2]. The essence of the API is simple, featuring primarily two types of calls:

Give n recommendations for user u This is the main output from the backend, returning recommended news stories or comments for a given user. Ancillary parameters will be added to this to facilitate the frontend placement of articles, e.g. the recommendation confidence and individual article weightings.

User u took action n This is the main input to the backend, allowing it to adjust recommendations according to user action. The parameters to this API call can be of many types. For example, "User commented on article $\#n$ ", "User plused or minused comment $\#n$ ", and "User visited link $\#n$ " are all valid parameters for this API call.

These two calls manifest themselves in Thrift with a few simple object definitions and an interface. The abbreviated version 0.1.0 is attached as Appendix A.

5 Project Management—Administrative Details

We have split the team into two parts, frontend and backend. With respect to the frontend, Tom and Mika have worked on the development and design respectively. The biggest difficulty here is to determine how to represent the news in such a way that the user will be inclined to continue using the service. As for the backend, Brian has been tasked with creating both the service architecture and the recommender engine. This includes

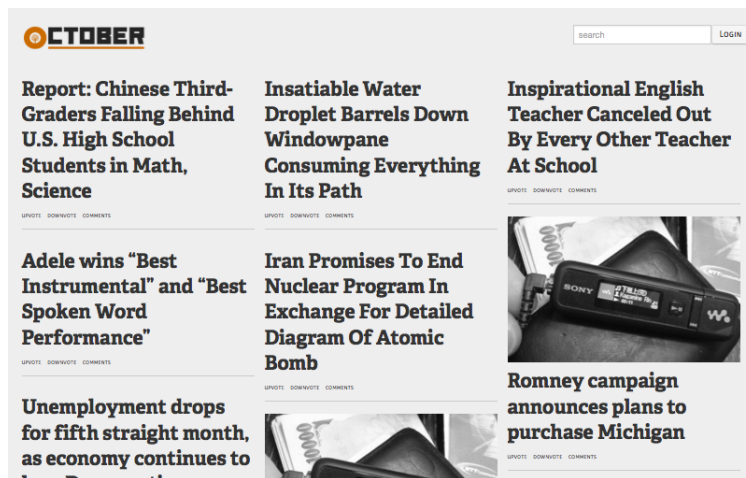


Figure 3: Project October Homepage, pre-Alpha (Headlines courtesy The Onion)

the implementation of our backend's TF-IDF algorithm and the Scala connector to the Thrift API and Mongo that computes recommendations.

5.1 Release v1.0

In Project Report 2, we promised a fully-functional beta by “early to mid-April”. On April 15th, we released the beta to the CWRU community, and it has been met with a warm reception.

The user interface for other pages will be in the same general content area, however there are no further screenshots to share at this time.

6 Testing & Evaluation

Testing and evaluation will be performed as we progress through the project. Currently, only very minimal testing is implemented.

7 Project Progress

The project is progressing satisfactorily. We are attempting to follow Agile Design Patterns, and we have been effective at meeting for short periods frequently. On the frontend, we have completed the user creation + login processes, mocked up the homepage, and connected the Thrift library. On the backend, we have implemented the API so it responds to requests from the frontend. The recommendation engine is still in-progress. We remain committed to maintaining team communication via scrum and continue working.

7.1 Discussions & Conclusions

Overall, the project is still on schedule. We still intend to complete a fully-featured product on-time.

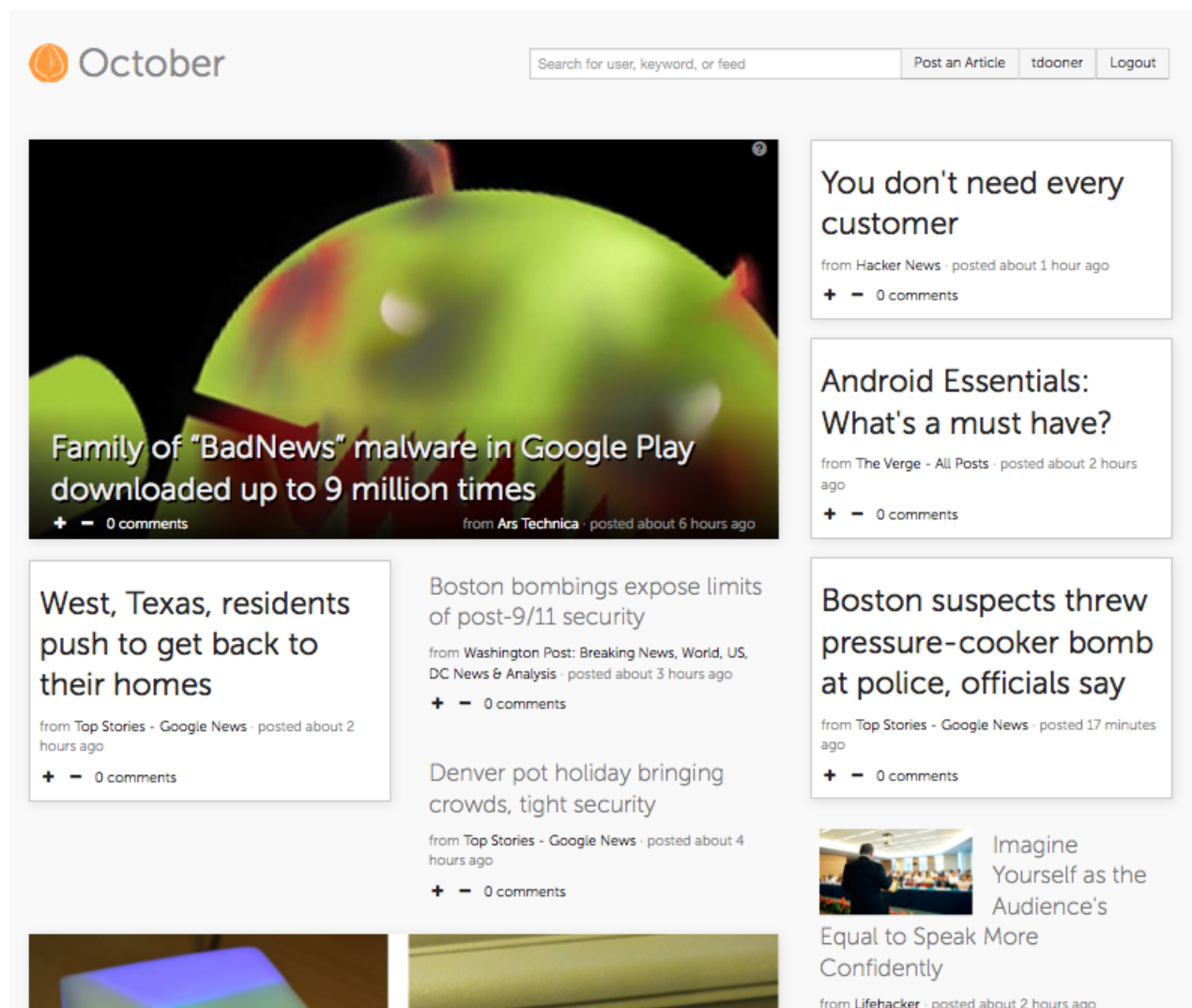


Figure 4: Project October Homepage, version 1.0

8 Lessons Learned

Throughout the development of October, we gained insight as how to properly design and create a product that will be used by many. Foremost among lessons learned, is that agile development ended up not being suitable for our work agenda. It required a certain rigidness in our schedules that we could not maintain. Instead, periodic scrum meetings sufficed for communicating progress and sharing ideas.

Alongside the ineffectiveness of agile development in a college setting, an important lesson was learned when October was first released to the public in early April. Based on user feedback, we learned that October was, in fact, a viable web service, and we have had moderate user retention. We have enjoyed working on this, and we hope to continue refining it after the class is over. We will be keeping in contact with users to build upon October, enhancing the service and hopefully increasing user retention.

8.1 Advantages & Disadvantages of a Service Oriented Architecture

9 Appendix

References

- [1] http://en.wikipedia.org/wiki/Eternal_September
 - [2] http://en.wikipedia.org/wiki/Apache_Thrift
 - [3] <http://en.wikipedia.org/wiki/Tf-idf>
 - [4] <https://github.com/peterc/pismo>
 - [5] <http://www.mongodb.org/>
 - [6] <https://github.com/ashleyw/phrasie>
 - [7] <http://mixpanel.com>
 - [8] <https://github.com/bis12/project-october-api>
 - [9] <http://api.rubyonrails.org/classes/ActiveRecord/Base.html#label-Single+table+inheritance>
 - [10] <http://masonry.desandro.com/>
-

A Thrift API Definition

```

const string VERSION = "1.0.0"
struct Post {
    1: required i64 post_id,
    2: optional double weight,
}
struct PostList {
    1: optional double confidence,
    2: required list<Post> posts,
}
struct User {
    1: required i64 user_id,
}
struct Token {
    1: required string t,
    2: required i32 f,
}
enum Action {
    READ,
    VOTE_UP,
    VOTE_DOWN,
    VOTE_UP_NEGATE,
    VOTE_DOWN_NEGATE,
    POST,
    COMMENT,
    REPORT,
    TAG,
    FOLLOW,
}
service Recommender {
    string ping() throws (1: TimeoutException te),
    PostList recPosts(1: required i64 user_id, 2: required i32 limit, 3: required i32 skip)
        throws (1: NotFoundException nfe, 2: EngineException ee, 3: TimeoutException te),
    bool addUser(1: required i64 user_id)
        throws (1: EngineException ee, 2: TimeoutException te),
    bool addPost(1: required i64 user_id, 2: required i64 post_id, 3: required list<Token> raw_fr
        throws (1: EngineException ee, 2: TimeoutException te, 3: NotFoundException nfe),
    bool userToPost(1: required i64 user_id, 2: required Action verb, 3: required i64 post_id)
        throws (1: NotFoundException nfe),
    bool userToComment(1: required i64 user_id, 2: required Action verb, 3: required i64 comment_
        throws (1: NotFoundException nfe),
    bool userToUser(1: required i64 actioner_id, 2: required Action verb, 3: required i64 actione
        throws (1: NotFoundException nfe),
    map<string, i64> userTopTerms(1: required i64 user_id, 2: required i32 limit)
        throws (1: NotFoundException nfe),

```

```
map<i64, double> textSearch(1: required list<string> tokens, 2: required i32 limit, 3: requir
  throws (1: EngineException ee),
bool addUserTerms(1: required i64 user_id, 2: required list<string> terms)
  throws (1: NotFoundException nfe),
bool removeUserTerms(1: required i64 user_id, 2: required list<string> terms)
  throws (1: NotFoundException nfe),
}
```

B Database Design

B.1 Frontend Database

Please see Figure 5.

B.2 Backend Database

MongoDB [5] is a schemaless database and as such there is no schema to show here. However, the backend uses “models” to access the data and as such there is a pseudo-schema which we can present here. The notation is simply cleaned up Scala without the actual logic attached.

```
class User(
  id: Long,
  tokens: Map[String, Long],
  friends: Seq[Long])

class Token(
  id: String,
  df: Long,
  posts: Seq[Long])

class Post(
  id: Long,
  posted: DateTime,
  tokens: Map[String, Long])
```

C User Manual

D Programmers Manual

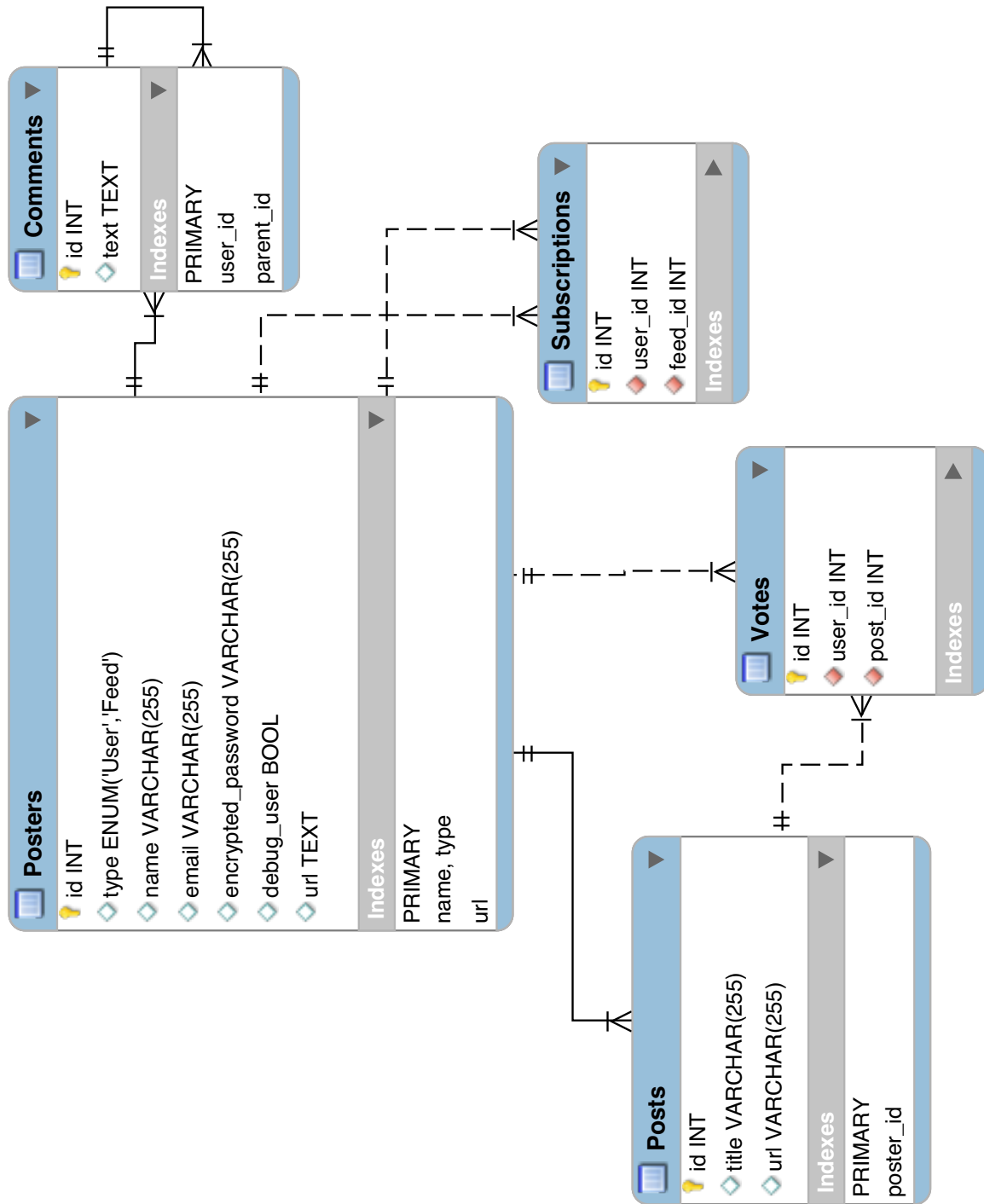


Figure 5: Frontend Database Relational Diagram