

# **Project Plan for News Streams**

**Group May13-31**

**Group Members: Jamison Bradley, Lance Staley, Charles Litfin**

**Advisor: Srikanta Tirthupura**

**Client: IBM Rochester**

**Client Contact: Paul Bye**

## **Problem Statement**

In today's world there is a plethora of information that is available on the internet for users to consume. The problem is that this information is spread out over several different websites rather than being streamlined and presented to the user conveniently in one location on the web. This diffusion of data over many different sites forces users to access several different sites a day to pick out the news stories the user wants to read, which takes up valuable time that could have been spent doing other things, such as reading the stories rather than searching for the articles.

Another issue that arises from using multiple news sources is whenever there is a major news story, the story will be covered by almost all of the different websites the user uses, and the websites will all provide the user with more or less the same information about that event. This is not an ideal situation for most users, it would be better if the stories could be grouped together and presented to the user so that the user knows all the articles discuss the same story and can choose to read just one from the source that the user desires to read from.

## **Scope**

The scope of this project is to deliver a web-based application which will allow the user to access different news feeds from a range of news providers. The webpage should also all the user to view articles based on categories and to remove duplicate articles.

## **Goals**

During the first semester of the project our group would like to have all documentation and design done, as well as, have an understanding of both text analytic algorithms and have gained experience with web design. We would then hope to start implementation over the Christmas break.

At the end of the second semester we want to have a working web application which displays newsfeed, provides categorization of articles and removes duplicates articles.

Lastly we would want to have enough time to add additional functionality to our web application, such as customization of sources, time permitting.

## Concept Sketch



## System Description

Our system will gather streams of articles from several different sources on the internet to be displayed for the user of our site. The user can either accept the default sources that we choose to support, or the user can login to our website and customize the sources that are used, by unselecting sources, we support that the user does not want to receive articles from. Our team has an user table in the database which contains a username and a hash of a password as the means of login for our site. Since we are storing a hash of their passwords it will allow us to not store their sensitive data reducing security risks for us.

Our team has chosen the sites that are planned on being supported (Can be viewed below), baring technical difficulty in working with these sources. Our team will pull data from these sites that already categorize their data. This will allow our team to know, based on where the story comes from on a site, what category an article belongs to.

After an article is gather the article will be analyzed and compared to the stories in the database. If our team determines that the article is the same as another article from a different source, only one article will be displayed together on the website. This will save the user of the hassle of seeing the same story over and over again from different sources.

Below are the sources that our team are planning on supporting on our website initially. We decided to go with a wide range of sources and hit several different countries, so that there is a wide range of articles and opinions for our users.

## **Planned Sources**

#1 CNN (US)

#2 The Guardian (UK)

#3 BBC (UK)

## **Interface Description**

The user interface will be a website. There will be a main news box region (1) which shows the most recent articles, their titles and links to the articles. It will also show the time that the article was posted. There is a next (5) and previous (4) button which will show more articles or previously seen articles by chronological order.

There will be a menu on the left side (2) which contains categories. If a category is selected only articles from that category will be display in the main box(1). There is a search bar(3) in the upper right corner which will allow the user to search by article title and display relevant results in the main box(1).

## **System Block Diagram**

1) The website will use the database to populate the user interfaces with news stories for the user to read. The website will also store data about user preferences, so that we, the developer can tailor the experience to individual users.

2) The developers will have a database that will store the information about the article that will be presented on our website. This database will be populated by the backend and will be used by the web server to populate the website with stories.

3) The backend will be used to gather the stories from news websites. This will be done by using the java backend running on an Iowa State server. Our team will gather the data from the news websites and carry out the analysis of that information on this server and then populate the database with that information.

## **Operating Environment**

The operating system for the user frontend will be on any system that has an internet access and can use a web browser. Therefore our software will work on Windows, Linux, iOS, and mobile devices. Testing of our software will be done mostly on a Windows machine, the interface, which is the website, will be tested primarily on Firefox and Google Chrome web browsers.

The backend will work on any operating environment which supports java use.

# **Functional Requirements**

## **A). Java Backend**

### **1. Obtaining Articles**

The backend should be able to pull articles off of the internet and add them to the database. The backend will obtain various pieces of relevant information as laid out in the article table.

## **2. Scanning**

The backend should be continuously scanning the rss feeds for new articles to be added to the database.

## **3. Sorting**

The backend should be able to sort articles into stories based on date, location, and article content. These stories will be specific events on which multiple articles are reporting on.

## **B). Database**

### **1.) Article Information**

Each article should be stored with the information that is specified by the Article Table in the database.

### **2.) Categories**

Each category has multiple articles assigned to it, but each article can belong to no more than one category.

### **3.) Stories**

Articles belong to stories which are created when the first example of a given story's member articles is committed to the database. Stories groups of articles that all deal with the same event or subject matter, during the same time frame.

## **C.) User Interface**

### **1. User Access**

The user should be able to use a profile to login. Their preferences and excluded sources will be recorded along with their account information.

### **2. Article Access**

The user will be able to access articles through links listed under a number of fixed categories. These articles will be presented as the article title in text along with a link back to the original web source of the article.

### **3. Customization**

The user will be able to remove sources that are supported which the user does not want to see articles from. Also the user should be allowed to add any previously removed sources. If additional time is available at the end of development there will be plans to allow the user to add their own sources.

## **Non-Functional Requirements**

### **1). Security**

#### **A). Username**

A user account is associated with a single username and password, as well any modifications the user has made to their sources. Their passwords will be hashed in order to avoid compromising user security.

### **2). Performance Requirements**

#### **A). Speed**

The system should be able to obtain a news article quickly after the article is posted to an RSS feed. The current requirement will be adding an article within three hours of the article being added to an RSS feed.

**B.) Categorization and Aggregation**

The articles should be properly placed into stories at least 80% of the time. This is tested through individual user testing.

**Market and Literature survey**

There are a number of similar news aggregators currently on the market. Also the majority of other news aggregators do not allow the user to customize the sources from which the articles are gathered from. This provides our system with a major difference to other systems, which should give us a distinctive edge over current systems. Also our user interface should be clearer, simpler, and more concise than those currently available.

**Deliverable**

Our final deliverable will be a website that will aggregate news stories from several different sources. Supporting this website will be a database that contains the information needed to provide the articles to the user. Our team will also implement a backend using Java that will be used to populate the database with content for the website. This content will be categorized and checked for duplicate articles.

**Risks and Mitigations**

<i>Risk</i>	<i>Probability of Occurrence</i>	<i>Criticality (0-100)</i>	<i>Risk Factor (Prob. Of Occurrence x Criticality)</i>	<i>Mitigation strategy</i>
The IBM Streams software	.50	80	40	Hold a meeting with our

has not been used by any of our members and could be difficult to learn and implement				IBM contact about streams as well as prototyping small streams applications for experience
Lack of experience with web development	.40	60	24	Research web development as well as prototype some web development
Development of the article	.30	75	22.5	Research to see if there are

<p>comparison algorithm is more difficult than expected</p>				<p>any third party algorithms that could be used. Build a prototype of the algorithm to test early on. Add extra time to the schedule.</p>
<p>Usability in many different web browsers and operating systems</p>	<p>.40</p>	<p>50</p>	<p>20</p>	<p>Research and develop for a number of browsers</p>

We are unable to get streams to work with our java backend and web frontend.	.20	50	10	Research alternate design options as well as using streams with external code
--	-----	----	----	---