

Ideal Climate Change

Michael Steele, Somn Kafley, Samyak Singh
CS 4624, Virginia Tech, Blacksburg, VA.
05/27/2016

Instructor: Edward A. Fox
Client: Mohammed M. Gharib Farag

Table Of Contents

Table Of Tables	4
Table Of Figures	5
Executive Summary	6
Chapter 1: User's Manual	7
I. Overview	7
A. Objectives	7
B. Communication Method	7
C. User Roles	7
D. Interactions With Other Systems	7
E. Production Rollout Considerations	7
F. Glossary	8
G. Interface	9
Chapter 2: Design & Requirements	14
I. Architectural Design	14
A. Decomposition Description	14
B. Design Rationale	14
II. Data Description	14
III. Screen Objects & Actions	14
IV. Functional Requirements	15
A. Functionality Statement	15
B. Scope	15
C. Performance	15
D. Usability	16
Chapter 3: Developer's Manual	17
I. Management Overview	17
A. Implementation Description	17
B. Major Tasks	17
C. Timeline	17
D. Security & Privacy	18
II. Implementation Support	18
A. Hardware, Software, Facilities, & Materials	18
B. Staffing Requirements	18

C. Implementation Impact	19
D. Performance Monitoring	19
III. Implementation Details	19
A. System Requirements	19
B. System Implementation	19
C. Acceptance Criteria	20
IV. Future Works	20

Chapter 4: Prototype & Refinement	22
I. Components	22
II. Data Flow & Functionality	22
A. Data Extraction Prototype	22
B. Data Indexing Prototype	22
C. User Interaction / Search Query Prototype	23

Chapter 5: Testing	25
I. Summary	25
II. Testing Scope	25
III. Types Of Testing Performed	25
A. Black-Box Testing	25
B. Functional Testing: Search	25
C. Functional Testing: Page Ranking	25
IV. Lessons Learned	26
V. Exit Criteria	26
VI. Conclusion	26

Bibliography	27
---------------------	-----------

Acknowledgements	28
-------------------------	-----------

Table Of Tables

Table No.	Title	Page No.
1	Glossary	5
2	Point of contacts	10
3	Timeline	11
4	Potential Refinements	18

Table Of Figures

Figure No.	Title	Page No.
1	An example of search display	9
2	Search display with facets	10
3	Display detail of each webpage	11
4	Display search history	12
5	Display Bookmarks	13
6	Data Flow 1: Data Extraction	15
7	Data Flow 2: Data Indexing	16
8	Data Flow 3: User Interface & Data Query	17
9	Sample Blacklight Interface 1: Landing Site	17
10	Sample Blacklight Interface 2: Search Result Display	18

Executive Summary

The Ideal Climate Change is a digital library and search engine project. It aims to provide an efficient way for graduate students and researchers to search and access archived tweets and web pages related to climate change. It will allow them to utilize the enormous collection of climate change data. Instead of being a complete (end-user ready) application, the project aims to be a starting point for developers who would like to expand on the current progress. The application consists of data containing tweets and webpages that have been extracted and indexed to SOLR. The results of user search is organized and displayed on the interface provided by Blacklight.

This report aims to highlight the design and software requirements of the application to demonstrate the scope, functional requirements, data decomposition, design rationale, and usability; a developer's manual to provide implementation descriptions, major tasks, timeline, hardware and software requirements, staff requirements, and acceptance criteria to future developers who want to expand on the current progress; a user's manual to inform stakeholders about user roles, communication methods, rollout considerations, and application glossary; prototype and refinement to explain various components and prototypes of the application; and finally a summary of tests performed and lessons learned.

It can be concluded that the user has been provided with an efficient tool, which can currently help them search a bulk of archived data. There are a lot of prospective features that can be implemented to enhance the application. For example, a personalized user experience with search recommendations based on search-history. Enhancements like this will be a great project for future Hypertext and Multimedia students to learn about searching, indexing, Artificial Intelligence, and Machine Learning.

Chapter 1: User's Manual

I. Overview

A. Objectives

The primary goal of the IDEAL climate change project is to provide an efficient way for technical and non-technical users to search and access archived tweets and websites related to climate change. The project serves as a tool for researchers who want to utilize the large pool of data behind these websites.

B. Communication Method

The mode of communication between the engineers and the client is currently via email, due to ease-of-access. Any functional or design requirements are sent to the engineers' Virginia Tech email. So far, this method of communication has been effective due to the high correspondence-rate between the parties.

C. User Roles

Our primary clients/users are researchers/scientists/engineers who want to utilize the large pool of data access through these websites. The main role of the user is to search tweets and websites by typing in a keyword on the given user-interface. The user is then responsible for using data-extracted from the underlying SOLR search engine for various research related to climate-control. Users can also report any search error or usability issues to the design team via email.

D. Interactions With Other Systems

Currently, the IDEAL climate change system interacts with an interface built on an Apache search platform called SOLR. SOLR is an excellent system because of its reliability, scalability, tolerance to faulty queries, distributed indexing, load-balanced querying, automated failover and recovery, and centralized configuration (Apache). SOLR controls the search and navigation, but it sits below a user-friendly framework called Blacklight.

E. Production Rollout Considerations

The final-product will be a web-application that allows users to successfully search climate-change related data, bookmark search results, and keep track of search history. It will be hosted on the Virginia Tech server at <http://nick.dlib.vt.edu:3000/> for use by end users.

F. Glossary

Apache: Popular web server software that owns SOLR.	Ruby: A dynamic, object-oriented, general-purpose programming language.
Blacklight: A Ruby on Rails engine plugin,	Ruby On Rails: A web application

which provides an application that runs on top of SOLR. A user-friendly framework that helps interaction between SOLR and the user.	framework written in Ruby, which provides data structures for a database, a web service, and web pages.
Climate Change: The change in global and regional climate patterns due to the increased levels of atmospheric carbon dioxide produced by fossil fuels.	SOLR: An open source enterprise search platform written in Java from Apache.
Java: A general-purpose programming language well known for class-based and object-oriented approach to problem solving.	Tweet: A status update from a user limited to 140 characters on Twitter.
Python: A general-purpose, high-level programming language focused on readability and light syntax.	Twitter: A free social networking platform, which allows users to follow, tweet, and reach out to other users.
URL: Acronym for “Uniform Resource Locator”, which is a reference to a resource on the internet.	User Interface: Part of software with which a human being interacts.

Table 1: Glossary

G. Interface

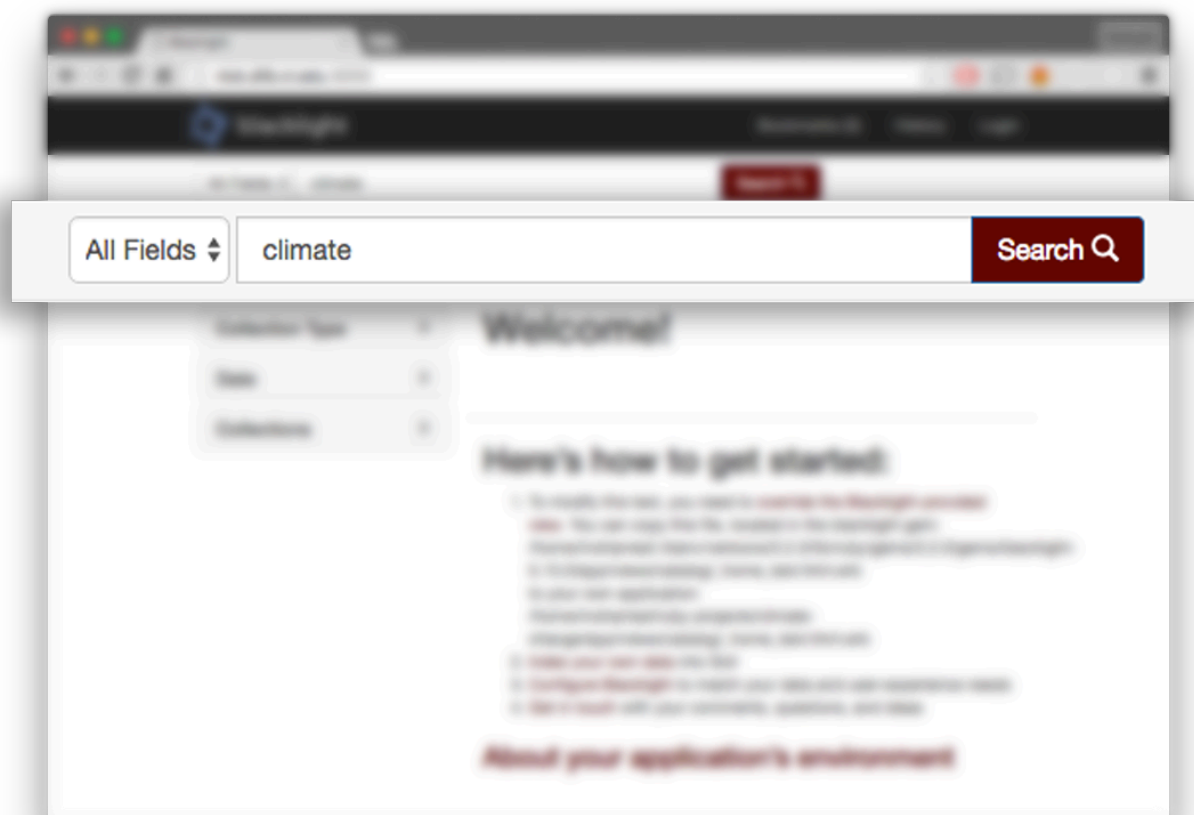


Figure 1: An example of search result display

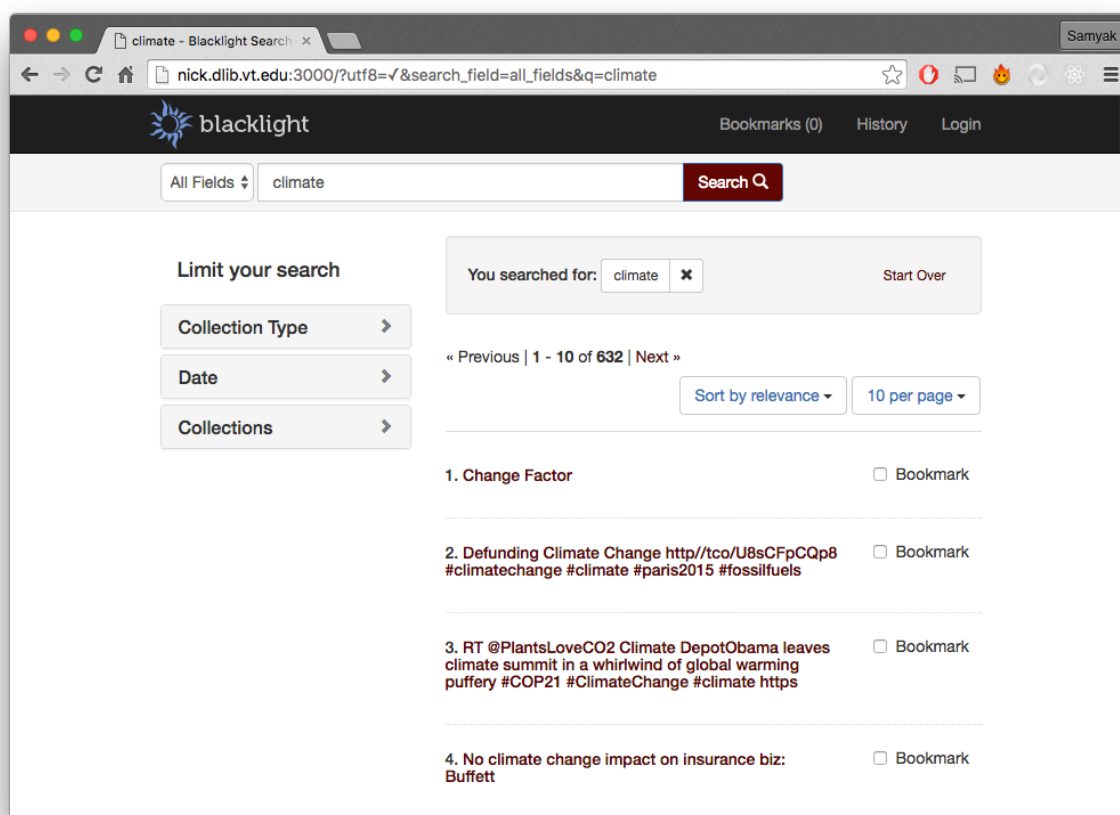


Figure 2: Search display with facets

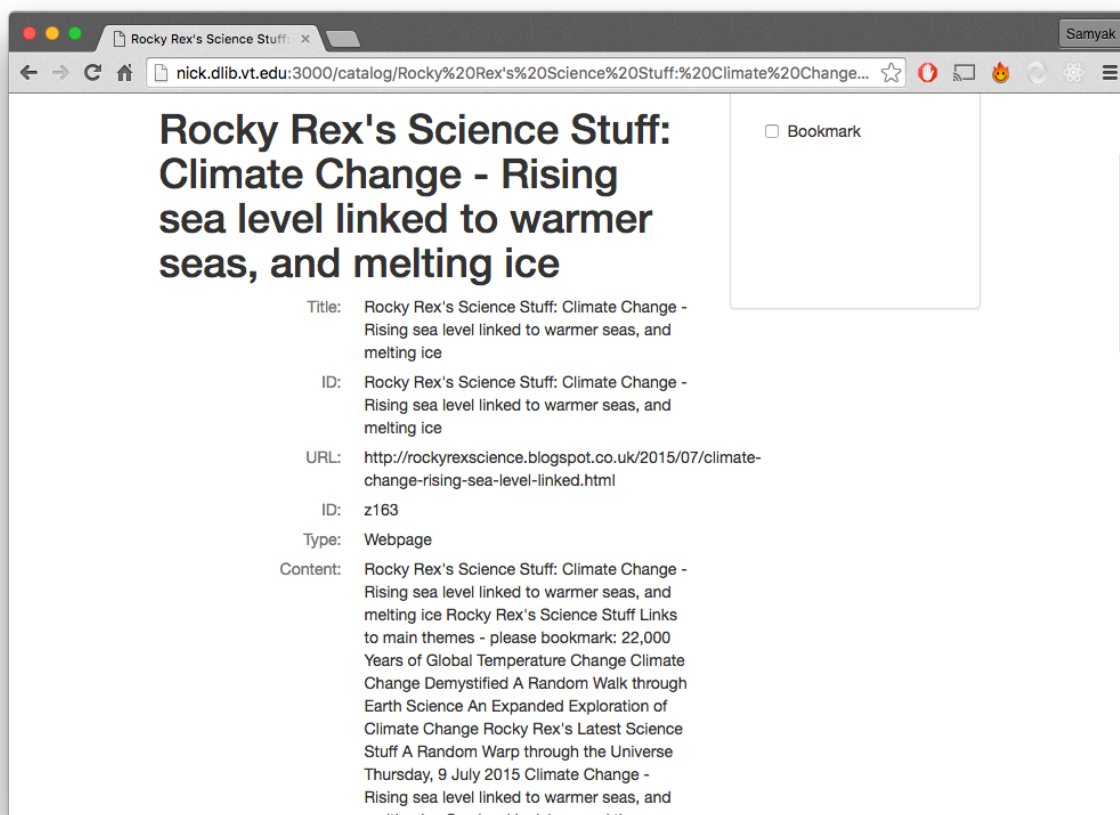


Figure 3: Display detail of each webpage

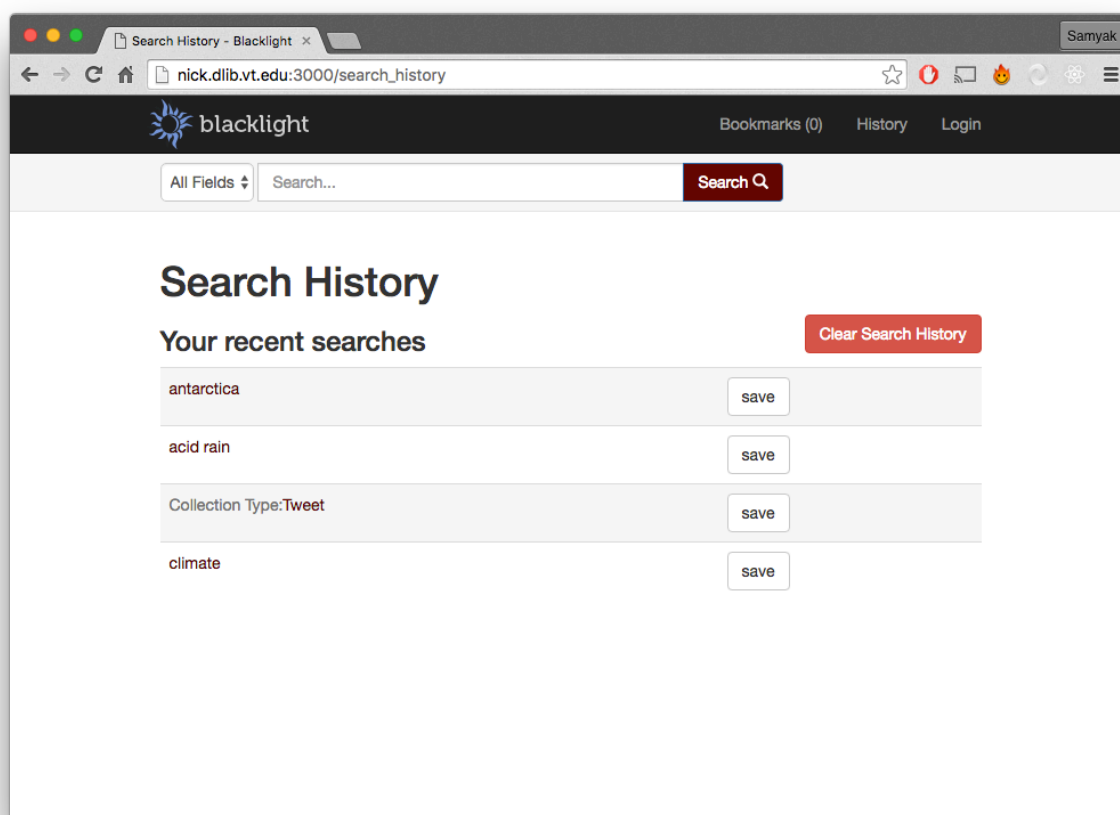


Figure 4: Display search history

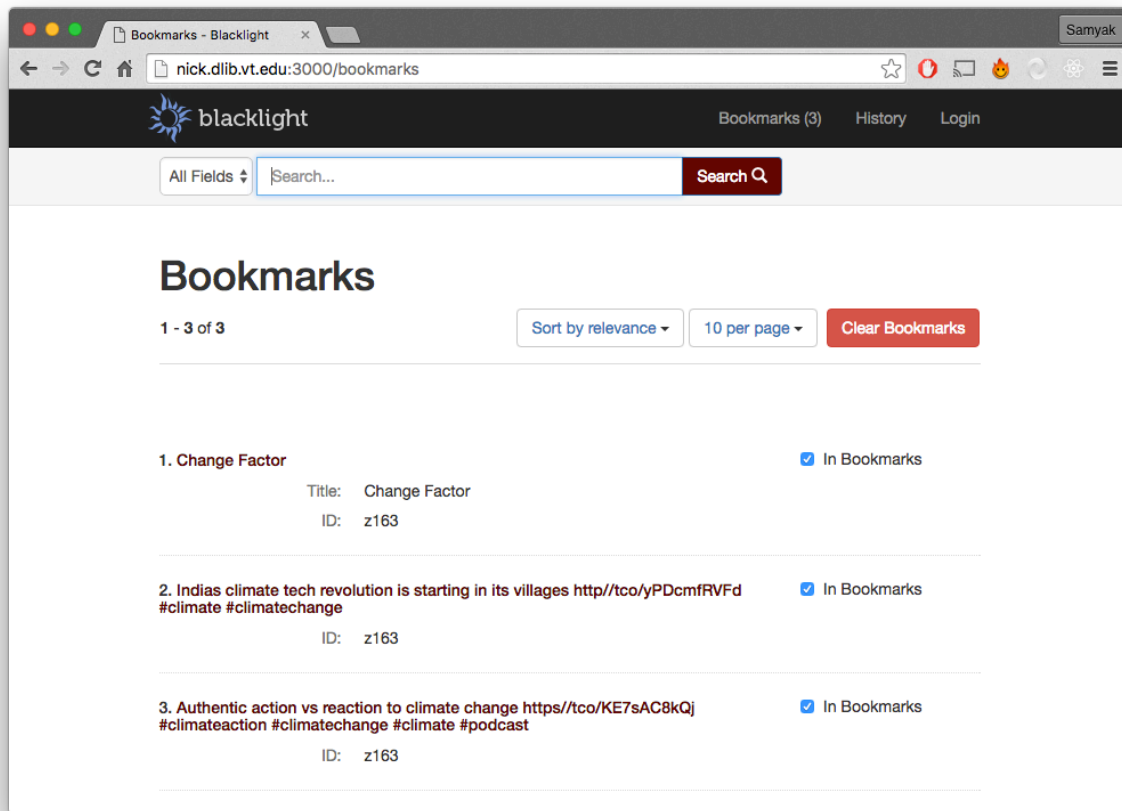


Figure 5: Display Bookmarks

The interface will display web pages from archived tweets. There are two phases of interface display. First one without displaying archived tweets and another one with detail such as place the tweets or web pages originate, time tweets first originate, a topic of the tweets, and frequency of the tweets. Displaying occurs only after adding functionality in Blacklight the search web pages and tweets will display on the interface.

Chapter 2: Design & Requirements

I. Architectural Design

A. Decomposition Description

Documents to index are converted to XML format. Once they have been sent to the index server, index tables are created that list where keywords are stored in each document. The SOLR server relays the search engine requests and returns the search results.

The Blacklight engines is an open source web application that provides a basic discovery interface for searching an Apache SOLR index providing searching box. Blacklight supports a collection of simple formats for sharing of search results and provides customizable interface. Our Blacklight engines will display origin of tweets, location and tweet time. For a demo and prototype, we use out-of-the-box vanilla demo of core Blacklight features.

B. Design Rationale

The architecture above was selected because of SOLR's powerful indexing, storing, and accessing efficiency. As opposed to a relational database such as SQL, in SOLR all searchable words are stored in an inverse index, which results in an incredibly faster search. Also, SOLR allows access to internal data structures that a relational database doesn't.

The biggest reason SOLR architecture was chosen was because of the predictability and speed of query speed; a relational database is very dependent on design and use case. Additionally, our client preferred to work with Blacklight and SOLR. Pleasing the client is a top priority, so we prefer to use SOLR.

II. Data Description

Majority of data this IDEAL Climate Change projects handles will be tweets, URLs, and websites extracted from the URLs. The websites will consist of relevant climate change content based on the user's search keyword.

The initial source of the data will mainly be organizations (e.g. Oxfam, WWF, Friends of the Earth, Greenpeace, Global Action Plan, etc.), politicians and government (e.g. Al Gore, Bernie Sanders, Ed Miliband, Department of Energy and Climate Change, etc.), news agencies (e.g. The Ecologist, Digg Environment, James Murray, The Climate Desk, etc.), bloggers (e.g. Grist, TreeHugger, Kate Sheppard, Julian Wong, etc.), and campaign groups (e.g. Stop Climate Chaos, 350.org, Plane Stupid, One Climate, Climate Camp, etc.).

III. Screen Objects & Actions

A user can click the search button after typing keywords after which screen objects related to climate change will display in the interface. A user with admin privilege can add content to the

website according to their needs. In addition, a user changes their search entry anytime and look for different results. By looking into the result a user can correlate trends in climate change webpages & tweets. It will make easier for user to see the change in the frequency of climate change. Thus, through web application users will benefit by being able to search by keyword and by tweet frequency, and will use this information for research.

IV. Functional Requirements

A. Functionality Statement

The software's top priority is to address every query with pinpoint accuracy. As a search based application, it is important that the user is presented with accurate and relevant information. The targeted users are individuals involved in research who may not have a lot of time to dwell on irrelevant search results. The software requires a user-interface that is efficient and fast; a clean and organized way to present the search results, and an accessible platform to be hosted on.

B. Scope

The scope of this project consists of several phases with various tasks due at the end of every phase, which will also be tested incrementally.

Phase 1: Environment setup.

This is when the SOLR server is set up, and Blacklight framework is installed. It also involves downloading dependencies like Ruby, Rails, and Java. After this phase, the environment will be tested so that data extraction and indexing can begin.

Phase 2: Data extraction from collection.

The data-extraction phase involves running a Python script to parse, scan, and extract website URLs from relevant tweets related to climate change. The collection of tweets is given to us by the client and resides in a database in the research lab, in Torgersen.

Phase 3: Data indexing to SOLR.

After the data extraction phase has been completed and properly tested, a second Python script is executed to index all the URLs and tweets to the SOLR search platform, from where the user can search by keyword. This is the last "back-end" step before the team starts working on the user-interface.

Phase 4: User interface design and development.

The user-interface phase is the last phase whose goal is to create an "easy-to-use" interface for the user. This is will be done only when Phases 1-3 because the interface provided by Blacklight is usable, but it is possible to have a better one.

C. Performance

The software needs to be consistently running so that busy researchers are not helplessly wasting precious time. The software doesn't need to render large images or files, which allows it to be swift and secure. Performance should be assessed by how fast the software responds to queries

and how fast it can display a large number of results. So, the main factors of performance are time and correct results.

D. Usability

The user-interface provided by Blacklight is good. However, the interface can be simplified by removing redundant and unnecessary features. The design goal is to provide a simple experience to achieve good performance. The elements required and interface-features required will be updated and refined as the team approaches Phase 4: User-Interface Design & Development.

Chapter 3: Developer's Manual

I. Management Overview

A. Implementation Description

Our goal is to design and implement a user interface where a non-technical user can search tweets and websites by typing a keyword related to climate change. We have two major phase of implementation: tweet extraction and archiving and indexing.

Instructor and Client

Role	Role Description	Name	Email
Instructor	Project assignment and supervision	Dr. Edward Fox	fox@vt.edu
Client	Research lead and supervision	Mohammed M. Gharib Farag	mmagdy@vt.edu

Developer and Research Team:

Developer	File extraction, organization design	Samyak M. Singh	Ssingh94@vt.edu
Developer	File indexing, organization, code analysis.	Michael D. Steele	derek22@vt.edu
Developer	File extraction, organization, and documentation.	Somn kafley	Somkk13@vt.edu

Table 2: Point of contacts

B. Major Tasks

Main tasks in the implementation of our projects involves receiving Python scripts and archived tweet files from research lead; running Python scripts on archived JSON file to extract targeted information from tweets related to climate change; using Python script to index “.txt” files containing extracted website to SOLR and Blacklight; changing SOLR web interface to properly indexing tweets into SOLR; demonstrating newly built climate change search engine to clients for feedback and testing.

C. Timeline

Date	Description
February 25, 2016	Meet with research-lead to gather archived tweet files.
March 1, 2016	Extract all data from archived file of tweets into text files.

March 4, 2016	Evaluate inconsistencies in Python scripts and finalize text files
March 31, 2016	Index files to SOLR, make it searchable and test
April 18, 2016	Set up Blacklight and have a search engine up and running.
April 25, 2016	Make changes in the interface, deliver to client

Table 3: Timeline

D. Security & Privacy

The server where our extracted data is hosted runs on the Virginia Tech server network, which is secured by the university. There is no risk of data invasion, or breach. Users will be able to safely access data once granted access to the system.

II. Implementation Support

This section describes the support systems for our project. Implementation support includes the support hardware, software used, facilities, and materials required for the project. This section also includes documentation, people involved, outstanding issues and impact, and monitoring of performance.

A. Hardware, Software, Facilities, & Materials

Hardware

The research-lead and developers all work on laptops for implementing the system. The SOLR server is installed on a desktop in the IDEAL research lab in Torgersen 2030. The rationale behind that was to maintain security of the server.

Software

On the MacBook and PC, terminal is used to run scripts and organize files. Sublime Text and Notepad++ are used to edit and comment Python scripts. On the web, Blacklight framework is used over SOLR. Additionally, Microsoft Word, Google Chrome, and Mozilla Firefox are used for research, documentation, and reports.

Facilities

Virginia Tech's research facilities in Torgersen have been very useful for the research, developer meetings, and client presentations. However, other McBryde 106 is where most of the implementation takes place.

B. Staffing Requirements

The staff required to implement this project must have the following skills:

- Code (Python, Java, etc.)
- Analyze written scripts
- Good sense of system design

- Experience with servers, database, or search-engines
- Excellent communication abilities

D. Implementation Impact

The implementation of this IDEAL climate-change search engine will have a positive impact on the research community. Due to the lack of a current system which allows a keyword search functionality, anybody interested in climate-change will now have access to the plethora of information indexed in this system. There will be a considerable amount of traffic once the system is up and running, however, it will not be big enough to affect the network infrastructure of Virginia Tech. The data made available for users will be regularly backed-up in case of network crash, which has a very minute possibility.

E. Performance Monitoring

The performance of the system will be measured and monitored by the research-lead and developers. A system/implementation/test plan has been created. Once the data extracted is indexed to the server, the team will start searching and accessing data via keywords.

The speed of the search, the accuracy of the results, and the functionality of the interface will measure performance. The monitoring will be simultaneous to the lifetime of the system with issue reporting, debugging, and enhancements done as the system lives on.

III. Implementation Details

This section contains information regarding any requirements of the system implementation. This involves software and hardware requirements, risks, and acceptance criteria.

A. System Requirements

Software

- Ruby Version 1.9 or higher
- Rails 3.2 or 4.x
- Java Version 1.7 or higher
- Apache SOLR
- Blacklight
- Terminal or Command Prompt

Hardware

Dedicated server with the following minimum specifications:

- Intel Core i5
- 8GB RAM
- 500GB Hard Drive
- Internet Access

B. System Implementation

Procedures

- Install Blacklight using [Quick Start](#)
- Write/amend Python scripts to extract websites and index information to SOLR. Script descriptions provided below.
- Test and confirm that SOLR is searchable with default configurations.
- Update scripts if test is not successful.
- Index tweets to SOLR by adding an additional core.
- Test and confirm that SOLR can search tweets and webpages.
- Enhance appearance and user interface for ease of use.
- Demo SOLR to social scientists, and gather their feedback.
- Re-configure SOLR based on the feedback from the clients.

Script Descriptions

downloadWebpagesFromTweets_Txts_Titles.py: This script extracts information from tweets specifically websites from links provided in the raw tweet data. The input of the script includes a JSON file that includes thousands of tweets related to climate change. It outputs text files containing website data.

indexToSOLR.py: This script takes all the text files generated by the first script and archives them to the SOLR database. The first line of each webpage contains its URL. Specifically, the script connects to the SOLR's Blacklight URL that accepts indexed files. It then sends a dictionary object containing each tag as the key and the value of each tag as the value.

All the scripts can be found at: <https://github.com/kafley10/ClimateChangeWebApp>

Implementation Verification & Validation

To ensure that the implementation was executed correctly, the research lead along with the developers will conduct an end-user test. This will happen after all the steps described in system implementation have been completed.

C. Acceptance Criteria

The system will label as “approved” or “accepted” for production after it has been thoroughly tested, verified, and validated described in system implementation and performance monitoring section. The primary exit criteria of the system are a qualitative measure of the user-friendliness of the system and a quantitative monitoring of the system performance. It is based on the display result of the search results. Future Work

IV. Future Work

Unlike many other projects, the climate change project is ongoing process. Someone will take charge from the point we left. So, it is important for us to show some of the possible future works. Our system currently addresses the scope of the project. However, like any other system, there are various aspects, which can be refined and updated to increase overall efficiency and

productivity of the system. Some refinements can be implemented but others require more time than what is available to us. Cumulatively, these refinements are a future upgrade to the system.

Present Scenario	Future work
Manually running scripts to download websites from tweets and then indexing them to SOLR is not only time consuming, but also tedious.	Automate this process so that these scripts can be run without the need of any developers. Maybe schedule the script execution so that it happens daily.
The scripts written only support Python 2.7.x, which is a problem because it is outdated. The latest version of Python is 3.5.x and offers new, updated features.	Discuss with the client and decide to upgrade the version of Python used. This will require updating both scripts entirely. The updates will enhance in search results and performance.

Table 4: Future work

Chapter 4: Prototype & Refinement

The purpose of our prototype was to have a system ready to be used such that users could start searching Climate Change subjects via keywords. To do this, we created a high-level navigational flow of the system to map the flow of data. Since running the extraction and indexing scripts didn't require any prototype, our primary goal was to have well-organized search functionality.

IV. Components

Apart from the server there are no hardware components in the prototype. Thus the prototype consists only software components. The software controls the interface, queries, data-retrieval, and display. The primary software components of the prototype involved a raw tweet data (JSON), two Python scripts, text files extracted from the JSON file containing website data, SOLR search platform, and Blacklight.

V. Data Flow & Prototype Functionality

A. Data Extraction Prototype

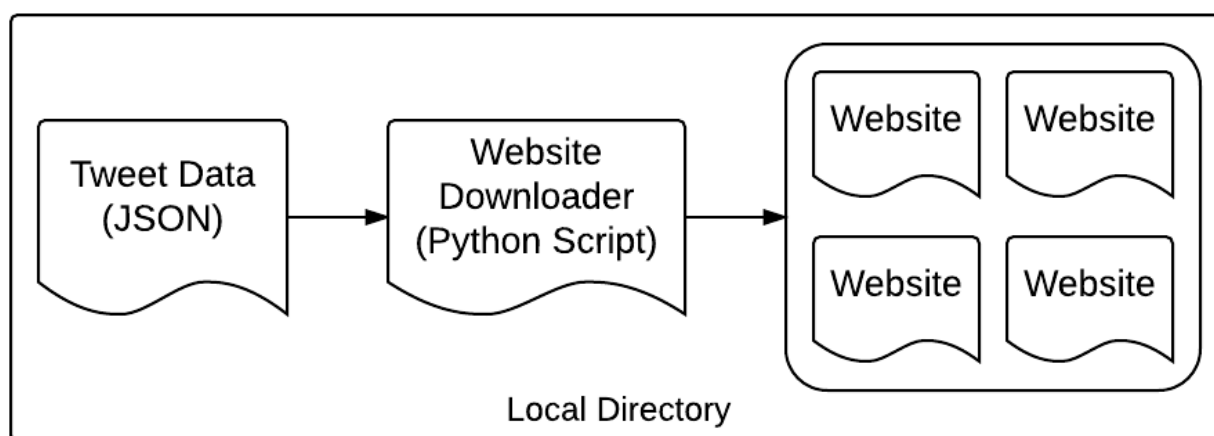


Figure 2: Data Flow 1: Data Extraction

The data flow for the first process involved an initial JSON file containing raw tweet data. This file was then taken by the Python script (1) to extract websites from individual tweet object, and then store it in a local directory. The websites were of text format ready to be archived to SOLR. More information about the specifics about the Python script can be found in Chapter 3, Section B.

B. Data Indexing Prototype

The indexed data is then archived inside SOLR with the help of the Python script (2). SOLR is configured to use ID tags and content tags. The content tags contain full text of each search result. This is used for 2 reasons: first to populate the index files with keywords and secondly to save and display the content of the search result to the user.

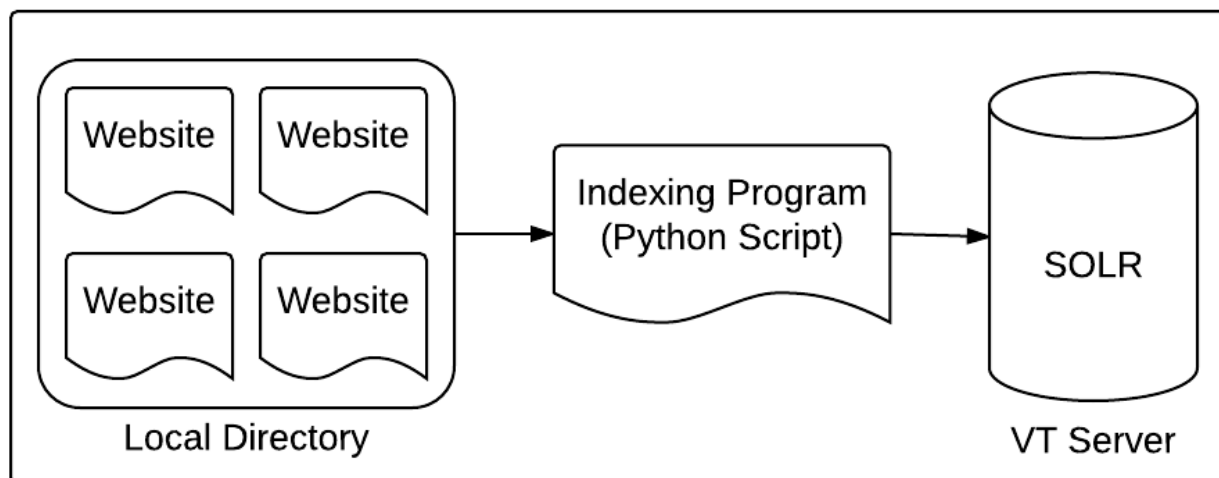


Figure 3: Data Flow 2: Data Indexing

We create a unique ID tag for each web page by hashing the URL of that webpage. As a result of running this script, the text of all archived web pages will be indexed and saved on the SOLR server and will be searchable. This includes all working URLs mentioned in all archived tweets.

C. User Interaction / Search Query Prototype

User interaction with SOLR directly is incredibly inefficient so we added a mid-layer component, Blacklight. Now, the user has a friendly interface where they can execute search queries to the platform. The prototype for this involves a two-way data-flow unlike the prototype for extraction and indexing due to the search query and results display.

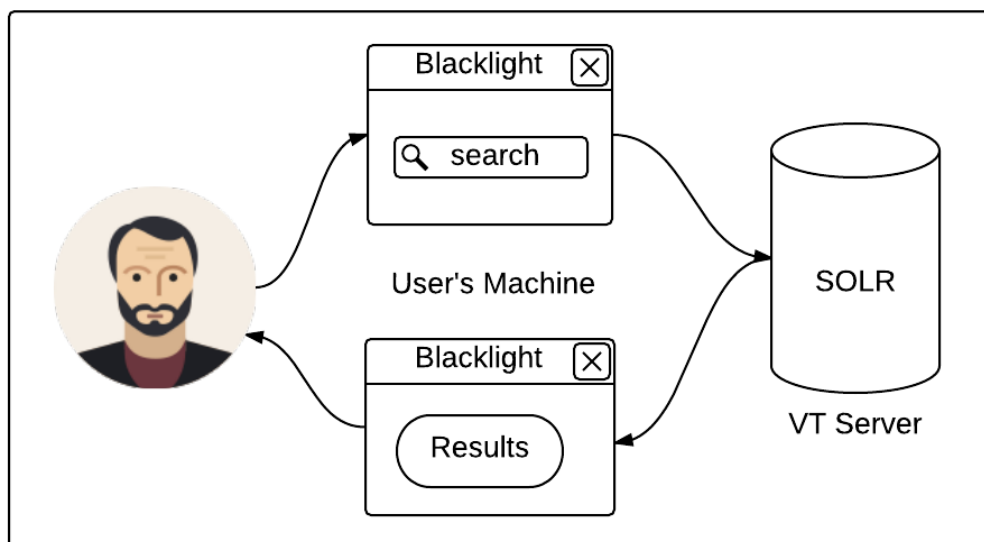


Figure 4: Data Flow 3: User Interaction & Data Query

Blacklight Interfaces

The following screenshots show the interface that the users will use to make search queries and get the results to that. It is an easier layer between the user and SOLR platform and makes the user experience more efficient.

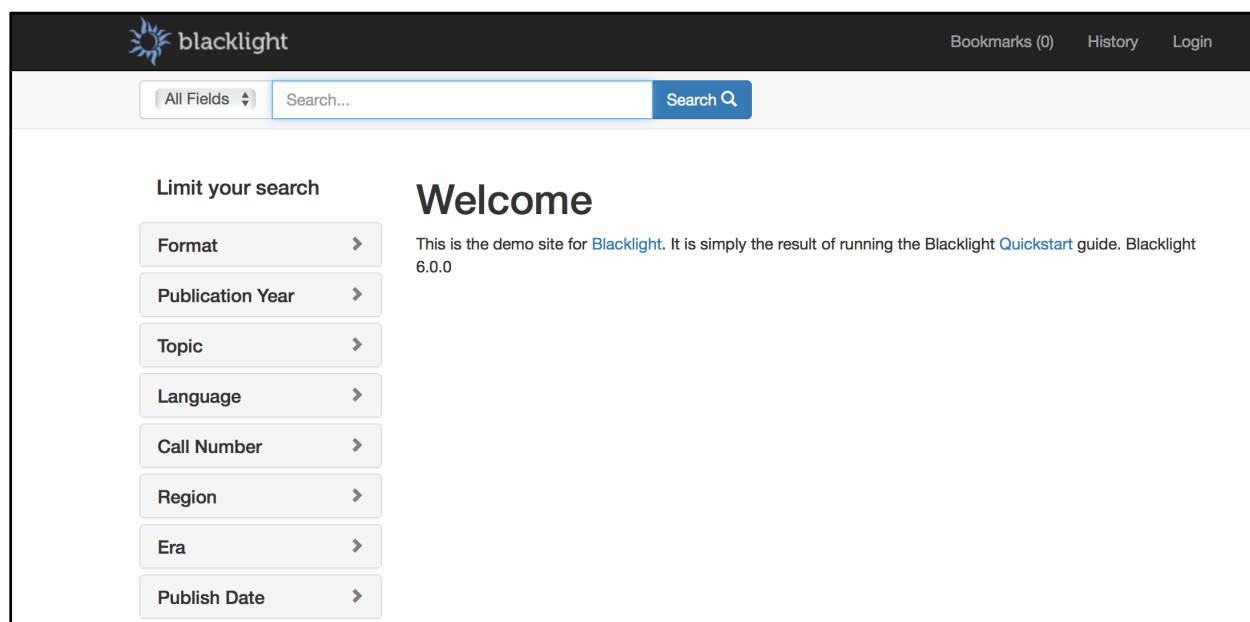


Figure 5: Sample Blacklight Interface 1: Landing Site.

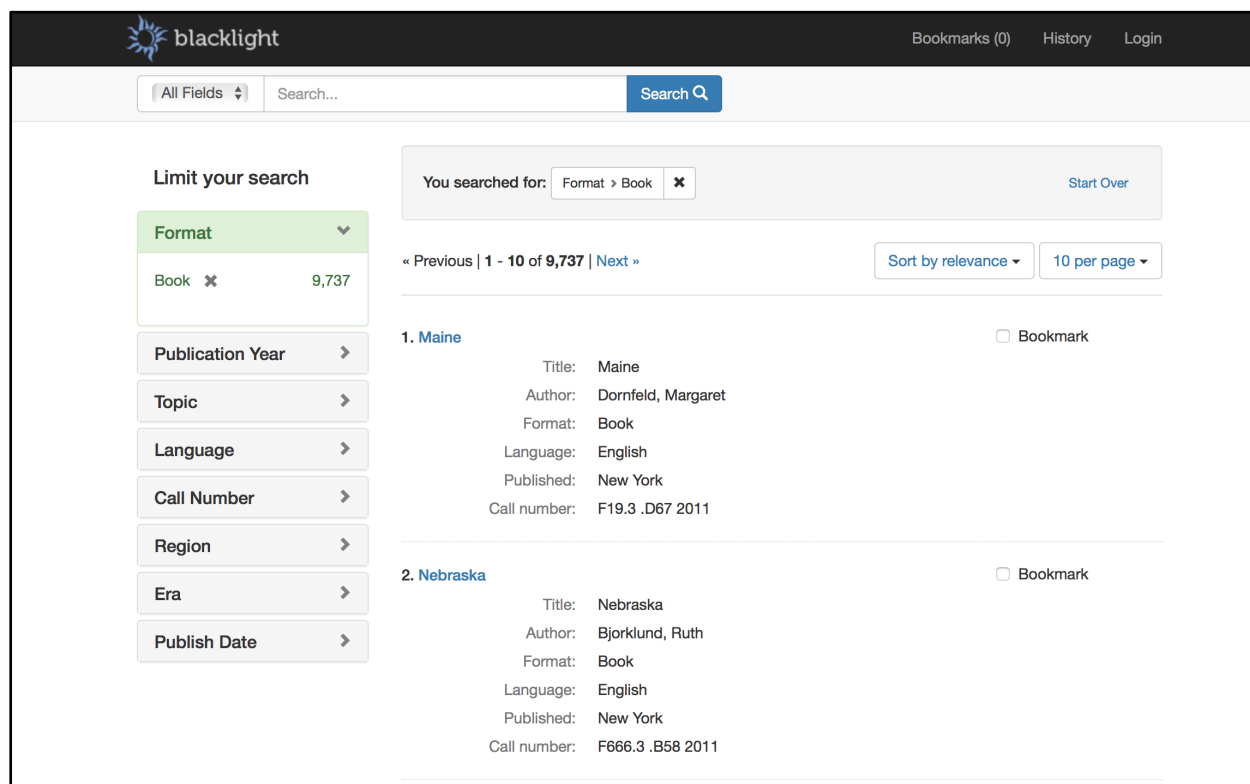


Figure 6: Sample Blacklight Interface 2: Search Result Display

Chapter 5: Testing

I. Summary

The climate change test report provides a summary of the results of test performed as outlined within this document. The main purpose of the climate change testing report is to obtained an evaluation of the project, find out the operating conditions of the climate change project, and find any irregularities that exist in the project. Another aspect of the testing report is to get feedback from the end user of climate change interface. We did three kind of testing: Black Box Testing, and two kind of functional testing. The detail of our testing procedure is outlined below.

II. Testing Scope

The testing scope of the IDEAL climate project was to efficiently test the user interface and the search functionality of the application. The user interface testing involved how well an end-user could execute a search and read the results. Navigating the application is fairly simple but testing was required to identify any possible errors. The functional testing of the application involved confirming that the search queries returned correct results.

III. Types Of Testing Performed

A. Black-Box Testing

Various users (undergraduate students) who have no knowledge of the internal workings of the project tested the user-interface. They were told about the overall goal of the project and were given the link to the application (<http://nick.dlib.vt.edu:3000/>) and asked to use the GUI to perform several search queries. They were also asked to browse the returned results of the query. These tasks would test the usability of the application by checking if users can perform tasks efficiently.

B. Functional Testing: Search

The developers tested the search functionality in the following way:

1. Go to the application at <http://nick.dlib.vt.edu:3000/>
2. Open one of the webpage text files from a local directory.
3. Choose a unique string that appears in only one webpage, and search it.
4. Confirm that there is only one result for the search query.

C. Functional Testing: Page Ranking

The developers tested the page ranking functionality in the following way:

1. Created a text file with a certain keyword present more than 10 times.
2. Created another text file with the same keyword present 5 to 10 times.
3. Created another text file with the same keyword present less than 5 times.
4. Then, those files were indexed to SOLR.

5. When that keyword was searched, the files with the highest number of that keyword present will appear on top.

IV. Lessons Learned

During the Usability Testing (via Black Box Testing), we realized that the interface design was not personalized to our project. It has the 'Blacklight' logo and title that confused the users about the interface. Removing the Blacklight logo and changing the color schemes can mend that. Users were curious if the search results could be sorted by date, title, etc. and if there were available filters that could help filter the results. This helped us learn that simply displaying results is not enough and that allowing users to interact with the organization is crucial.

V. Exit Criteria

The exit criteria include:

- 100% positive feedback from the black-box testers so that they can continue to efficiently perform search tasks and feel comfortable with the user-interface.
- Continue to pass all functional testing done on the search and page ranking, as the application is refined to match the black-box feedbacks.

VI. Conclusion

The primary types of testing done on the IDEAL Climate Change project were the black box testing for usability and functional testing for search results and page ranking. These tests exposed some user interface design problems that are currently being fixed. From these errors, we learned that although our goal was to provide a simple and efficient interface, it is important to plan the interface design from the beginning instead of attempting to build on top of the Blacklight interface.

Bibliography

1. Beer, Chris. "ProjectBlacklight/Blacklight." *GitHub*. Feb. 2016. Web. Mar. 2016.
<<https://github.com/projectBlacklight/Blacklight/wiki/Quickstart>>.
2. "SOLR Quick Start¶." Apache SOLR. Apache SOLR Foundation. Web. 14 Mar. 2016.
<<http://lucene.apache.org/SOLR/quickstart>>.
3. "SOLR Admin." *SOLR Admin*. Ed. Mohammed M. Garib Farag. Open Source. Web. Mar. 2016. <<http://nick.dlib.vt.edu:8983/SOLR/#/~cores/Blacklight-core>>.
4. Beer, Chris. "Blacklight." [Http://projectBlacklight.org](http://projectBlacklight.org). GitHub. Web. 01 Apr. 2016.<<http://projectBlacklight.org/#examples>>.

Acknowledgements

Thanks go to National Science Foundation grant IIS - 1319578, III: Small: Integrated Digital Event Archiving and Library (IDEAL). We would like to express our sincere appreciation to Dr. Edward Fox for providing us the climate changes project and guiding us with valuable and constructive suggestions. His willingness to give his time and walk us through the project is much more appreciated. In addition, we would like to thanks Mr. Mohammed M. Gharib Farag for providing us valuable information we needed.