# Handling Archived Data:  Big Data Team
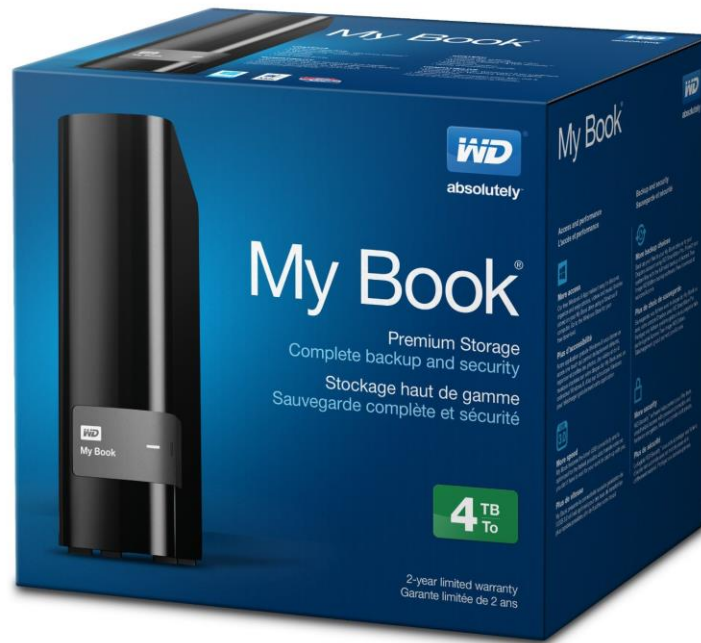




Advisor: Erik Rozolis                    Member: Jin Bai

# TABLE OF CONTENTS

## 1.0    Briefing into Big Data Team

This is Jin Bai's first semester on the Big Data Team. In the initial Sunday meetings, new members of Big Data were to first choose which group they wanted to work in for the semester. Jin choose to join the Web Interface and System/Cloud group.

The System/Cloud group ended up only meeting a handful of times and is not the focus of this report; from Appendix A, you can see what Jin has learned on the System/Cloud group such as installing and setting up Azure (Appendix A: figure 1,2), generating secure keys (Appendix A: figure 3), and creating Azure virtual machine instances as well as running virtual machine commands such starting and stopping a VM(Appendix A: figure 4).

In the earlier Thursday meetings of the semester from around 8/31 to 9/18, the Web Interface group mainly focused on teaching new members about how the website works, setting us up with access to the GitHub and Google Drive, helping us SSH into the server, and showed us how our database looked. These items were not delved into in major detail. Members took notes on these items.

Each member was then assigned a different task for the semester and beyond. Jin was assigned to handle archived data so that people could access and analyze the data we had in our archives.

## 2.0　Introduction to Handling Archived Data

We expect to have several hard disks in the future that store various types of archived data. Currently, we don't have a way to utilize that valuable data. The idea is to eventually access them easily and incorporate that into the schema of our current website.

The long term goal is to be able to access archived data and allow people to analyze that. There are currently several hard disks which are several terabytes, each, full of archived data from national parks. The advisors of the Web Interface team, Erik and Everett, would like to be able to pull images, and eventually videos, from these disks and display them on the website for people to analyze. This is actually part of a much larger project, but Jin's job is to focus on setting up the drives so that it is possible to retrieve any image from any one of them.

Even though Jin's job was to simply set up the hard drives so that we can access the images, he went above and beyond that and is hopeful that we can already incorporate those image files into the website as a new prototype by next semester.

## 3.0    Initial Ideas, Problems and Analysis of Various Use Cases

There were a lot of problems initially. It wasn't even possible to access our archived images in the WD storage device. My Book, the directory with the data could not be cd into, permission was denied and change mode had no effect on it since it was an external device. The only way that My Book was able to be accessed was as super user (Appendix B: figure 1). Another issue was that we could not view the images even as super user.

This problem was finally solved with Ahmed's help by creating a new user called "archive" on the archive machine. By connecting the drive when logged in as archive, the new user was granted permission to access My Book. It is no longer necessary to be super user to access images (no need for permission changes either). Archive will now be the primary account used by the team for accessing the archived images (Appendix B: figure 2).

The initial idea was to set up a two way socket form of communication between the archive drive machine that is connected to the WD media storage drive and the development machine. Professor Lu informed the team that a socket was a very ambitious task that takes a ton of time so should look else-where perhaps and HTTP server. After that, Professor Lu suggested that we find a way to index our data to account for various use cases such as for past crime incidents, vacation location history, historical weather analysis, and past events.

**Implementation and Issues**

Professor Lu suggested the use of mapping tables and after viewing the data with Michael, it was decided that mapping tables may not be necessary as our data is extremely organized (Appendix C: figure 1). I thought immediately of associating key to value through the use of an organized Python dictionary since our images are all named in a similar fashion (Appendix C: figure 2a, figure 2b). The format of the National Park Data (Appendix C: figure 3) is the basis of this idea.

Jin implemented this idea by first making a text file that mapped the park name abbreviation to the actual name of the park with the idea in mind that users can eventually select a park name which would map to our data where the name is simply an abbreviation (Appendix C: figure 4). Then Jin implemented this idea in Python by writing a script that parsed a user request for a park through the command line (Appendix C: figure 5). Then the script would cd into the correct folder with the text file, and create the dictionary. From there, the dictionary would be used to cd into the actual folder of the park that the user wanted.

Professor Lu thought this idea was too narrow and wouldn't cover all

the use cases so Jin had to further think of another idea to store all the possible

future data. He began to think about making another table in our database for

National Park information for easy access and integration with the site. He

created an archived_data folder (Appendix D: figure 1) in camcam directory of the

development machine and found a way to manually transfer single images to that

folder (Appendix D: figure 2) and set up a secure communication for file

transferring with ssh keys. Eventually, the idea would be to send these images

from the development machine to the site.

A new table called archived_data was also created in the cam2

database (Appendix D: figure 3). The fields park_name and snapshot_time are the

main parts of the table. Many values were manually input into the table

(Appendix D: figure 4). The idea is to eventually use this crucial information to find

images requested from the development machine to get them from the archived

data machine. The next step was to find a way to auto-populate these tables

instead of manually inputting values. Erik helped with the process by showing Jin

how to output all image filename data from a directory into a text file (Appendix

D: figure 5). Jin utilized this on each individual archive national park directory to

obtain text files of all the names of parks and verified the information was correct

(Appendix D: figure 6). The issue is to use this in a Python script; Jin learned how

to auto-populate the database from old scripts (Appendix D: figure 7, figure 8).

**Next Semester Goals**

This semester was a huge success in the members all completed the task and went beyond what was expected from them. The next step is to incorporate our back end communications into the website seamlessly. Jin has provided his ideas in two Sunday presentations and many great discussions stirred from that.

The idea Jin presented is to have a separate category specifically for archived data on our site which lets users browse our archives. The idea is to have users select that they want to see historical data, then let them select what kind of historical data to see to account for all use cases. For now we have data on national parks, but in the future when the Big Data team receives even more archive data such as from NASA or other organizations, there should be a way for users to select which one they want to view. From the URL in (Appendix E: figure 1) the Appendix E, "natl_parks" could be just one option; NASA for example could be there instead. After the user selects, further options appear (Appendix E: figure 2). From there, date and time options are displayed which allow users to further specify their request (Appendix E: figure 3) with a loading bar that pops out. Lastly, the correct images requested are displayed on the page (Appendix E: figure 4).

**Summary and Conclusions**

This semester Jin has definitely learned a lot through doing VIP with the Big Data team. Working with archived data has allowed Jin to work with both hardware and software. He has also learned a lot of valuable presentation skills from weekly Sunday presentations. He wants to work with the Big Data team again next semester.

The goal for archived data in the future has become much more clear due to the work done on it this semester. Once the back end is incorporated into the site, the team will have much more flexibility in determining how they specifically want this data to be visible to the user. Many ideas came from the Sunday meetings. Ahmed and Professor Lu suggested possibly encapsulating all the data until the data is requested on the map and then have the archived images appear on the map in a slightly different format as the cameras that show up to differentiate this previous historical data from our current video camera analysis selections.

Jin has learned how to work with an SQL database, practiced his Python coding proficiency, learned how hard drives work and researched them to provide accessibility for our team. He is looking forward to continuing his research on archive data.

# References

"Cyberduck Help". RetrievedNovember , 2014 Available:
https://trac.cyberduck.io/wiki/help/en

"Django Documentation". RetrievedSeptember , 2014 Available:
https://docs.djangoproject.com/en/1.7/intro/tutorial01/

"is not in the sudoers file. This incident will be reported.". Retrieved October , 2014 Available: http://www.linuxforums.org/forum/red-hat-fedora-linux/144428-solved-not-sudoers-file-incident-will-reported.html

"Permission Denied when I try to access my MyBook". Retrieved October , 2014 Available: http://ubuntuforums.org/showthread.php?t=887771&page=2

Rozolis, E , "Twith2/Script.py". RetrievedOctober , 2014 Available:
https://github.com/Erozolis/Twitch2/blob/master/Script.py

W , "Linux - Mount device with specific user rights". Retrieved October , 2014 Available: http://superuser.com/questions/320415/linux-mount-device-with-specific-user-rights

# Appendix A:

# System and Cloud

Figure 1. Installing azure on Ubuntu virtual box



Figure 2: Associating Azure account with Big Data virtual machines

Figure 3. Generating secure keys for Azure



Figure 4. Created an Azure instance, running various virtual machine commands

# Appendix B:

# Accessing Archived Data on My Book

Figure 1: Only able to access archive images as root user



Figure 2: Archive user was created  while connected to archive drive which granted data access

# Appendix C:

# Initial Ideas

Figure 1: Structure of My Book storage directories are very well organized



Figure 2a: Image files are all named in a similar fashion

Figure 2b: Image files are all named in a similar fashion



Figure 3: All names are in the format <Park Abbreviation><Year><Date><Time>

Figure 4: Dictionary key to value mapping idea, text file to store data



Figure 5: Implementation in Python

# Appendix D:

# Database Auto-population and Changes

Figure 1: Added archived_data directory to camcam folder



Figure 2: Able to transfer single images from the archived machine into the dev machine



Figure 3: Created archived_data table in DB



Figure 4: Manually input values into the table

Figure 5: Outputting all the image filename data into a text file



Figure 6: Performed same ls –R > outfile command on inner directory to check lines to confirm that we have all the files

Figure 7: Learn how to auto-populate the database through these files



Figure 8: State abbreviations load into DB file example

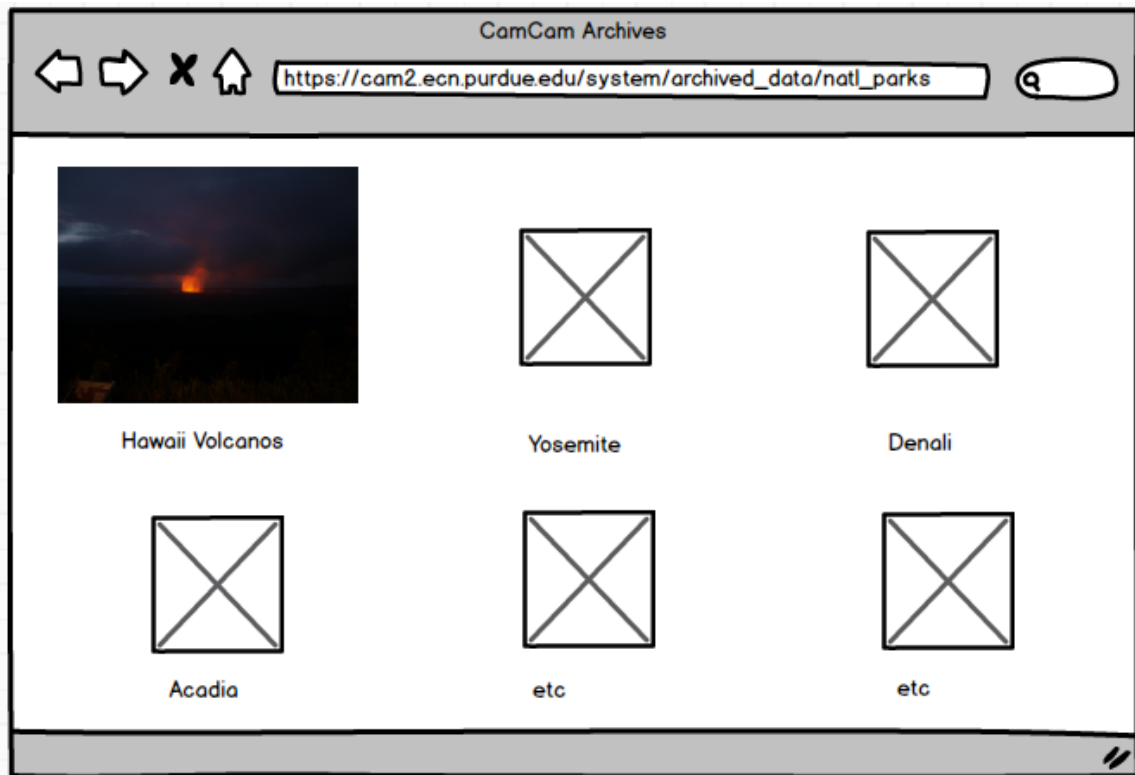# Appendix E:

# Front End Design Possibilities

Figure 1: Front end idea, select category archived_data then nat'l parks, display parks
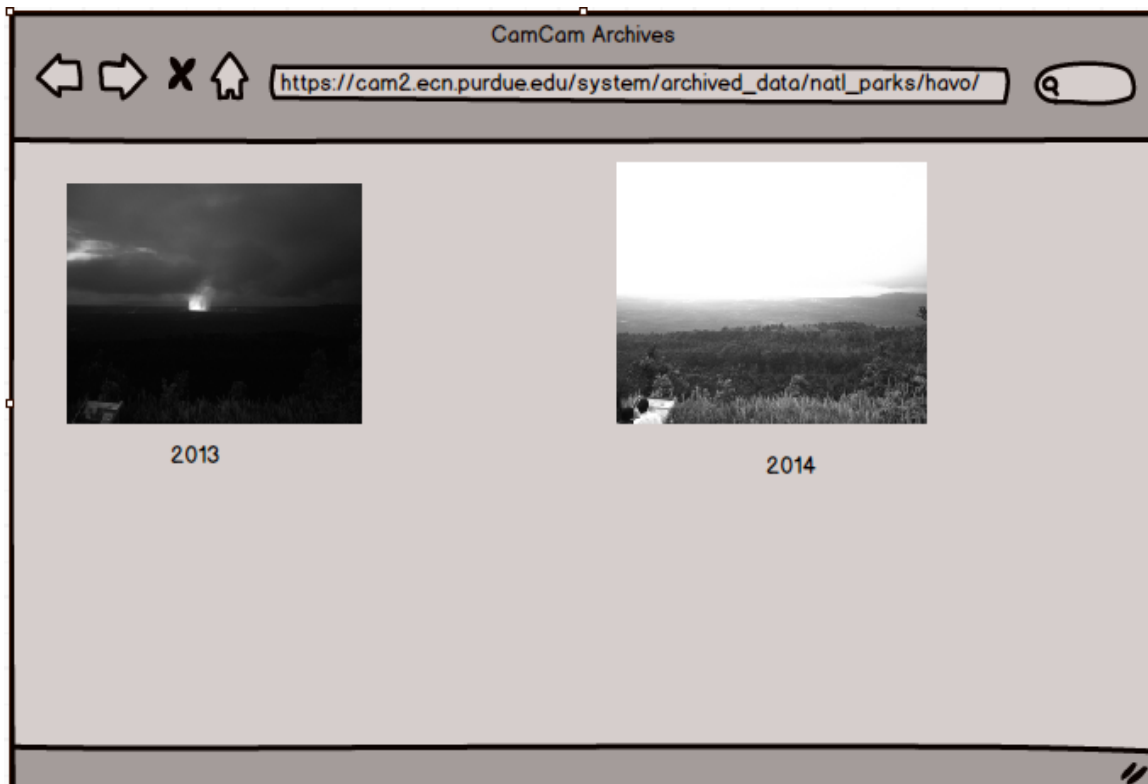


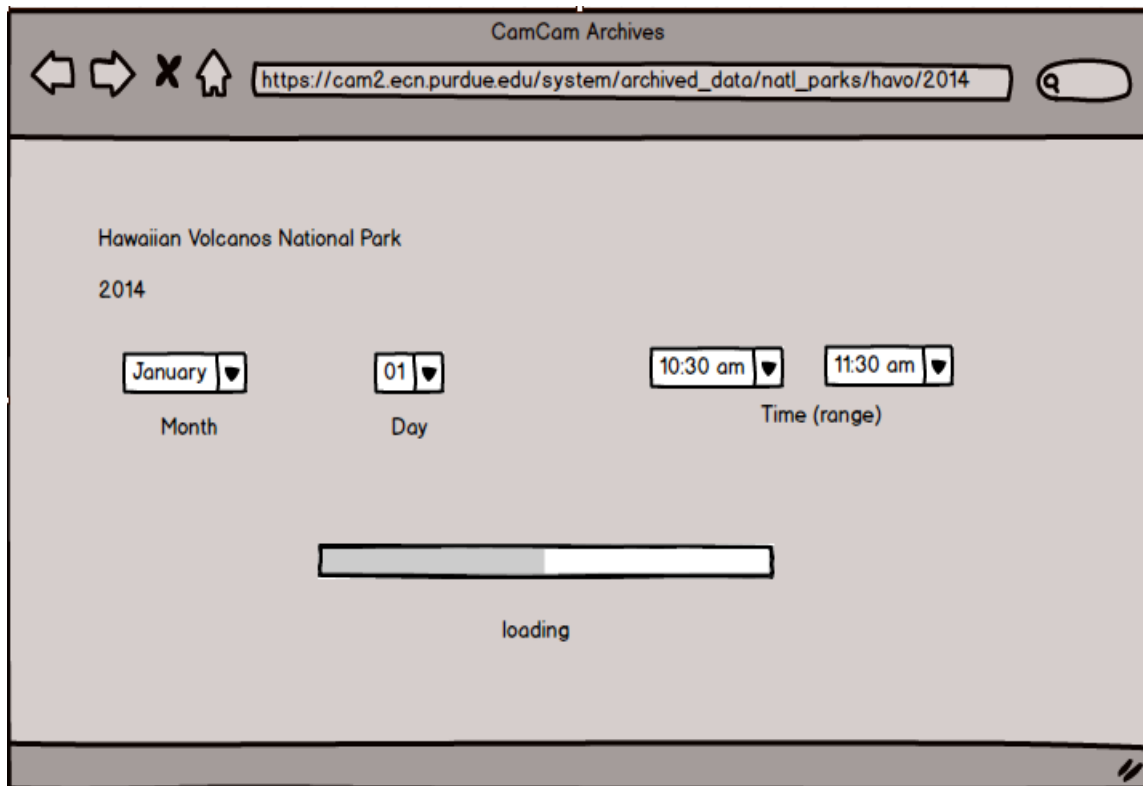Figure 2: After selection of park, show year options
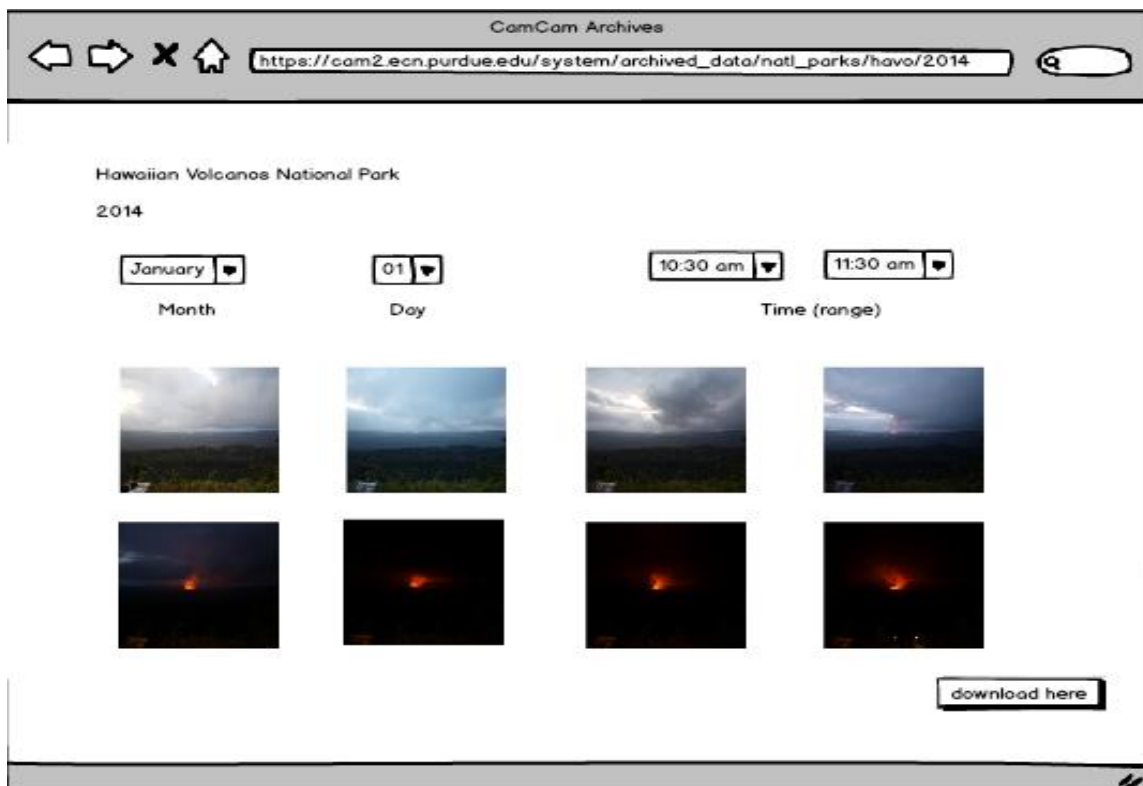
Figure 3: After selection of year, show month, day, and time options



Figure 4: After desired options selected, show images