# SAIVT WebDemo

# Project Report

Francesco Ferraioli                                    n8323143

## TABLE OF CONTENTS

## ACHIEVEMENTS

The project proposal listed various tasks for the entity extraction webdemo project from SAIVT. This section will outline my progress and achievements on those tasks.

The tasks the project proposal mentioned were categorised as three mandatory and one desirable. The tasks are listened below, the first three are mandatory and the last was the desirable one:

1. User Interface Enhancements
2. Add Video Functionality
3. Entity Display Enhancements
4. Implementation of Authentication and Authorization

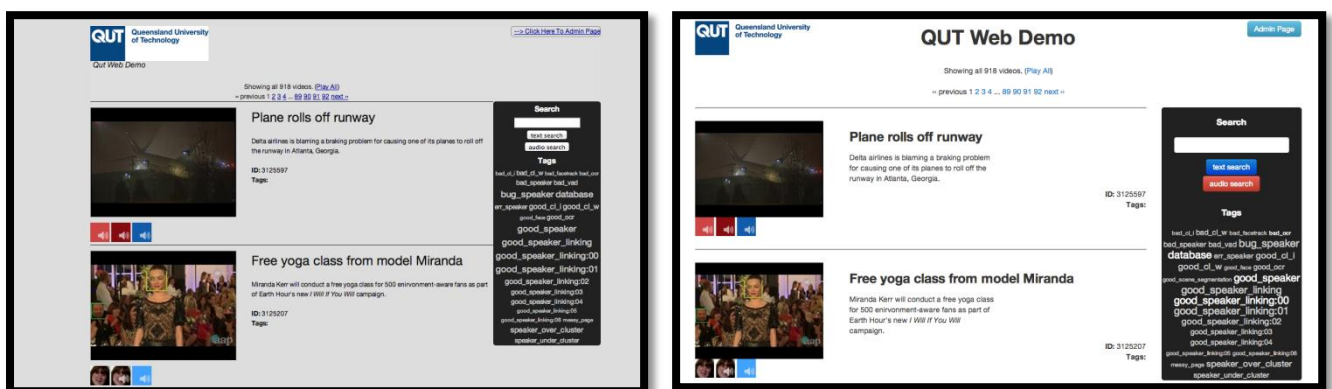I will now discuss my progress on these tasks individually.

## USER INTERFACE ENHANCEMENTS

The two main purpose's of this task were firstly to ensure that the code for the front end is written well, easy to understand and renders as intended and secondly to make the UI more user friendly and pleasing to look at.
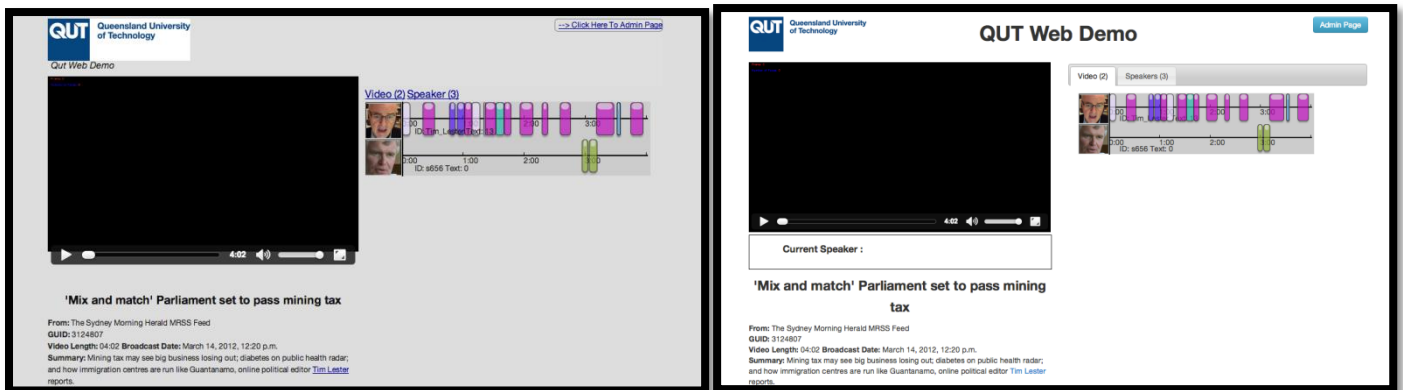
I started off by taking a lot of the inline styling out of elements into classes and put those classes into a separate css file which I included in the file header. This made the html much easier to read and manage and furthermore implementing best practices. Moreover, the base html file was written very poorly, having various blocks for each page to implement, whereas it should simply contain two blocks, a head block and a content block and the html which should remain constant for all the pages that inherit from it. This allows for all those page to only have to specify two blocks making it very clear to the developer how the page is going to render in a browser.

The next step was to introduce the Bootstrap Framework for further enhancements to the UI. Bootstrap has a very nice layout grid template using divs and I used this to style the elements where I wanted them to be, taking away a lot of the previously unmanageable styling which was used to place elements on the page. Bootstrap also has a very nice list of classes to style buttons and I used these to my advantage to make the page much more pleasant to look at from a user perspective.

The images below shows the old look compared with the new look of the index page which is the page which lists out videos.
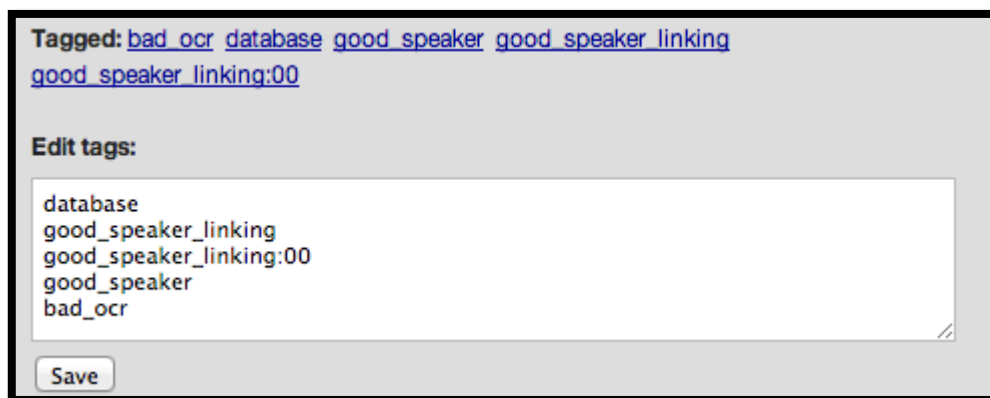


When I finished styling the index page I moved on to the video play page which is the page that renders when a particular video has been chosen by the user to play and to see what entities belong to this particular video.

The above images compare the old and the new look of the video play page. The main styling differences are the added tabs on the right hand side showing the Faces and the Speakers. Previously to switch between Faces and Speakers the user needed to click the links above and that would hide and show the respective entities. Now the Faces and Speakers are separated in tabs and thus it's a simple as clicking the tabs to switch between the entities. The tabs are much nicer to look at than the individual links that were previously available.
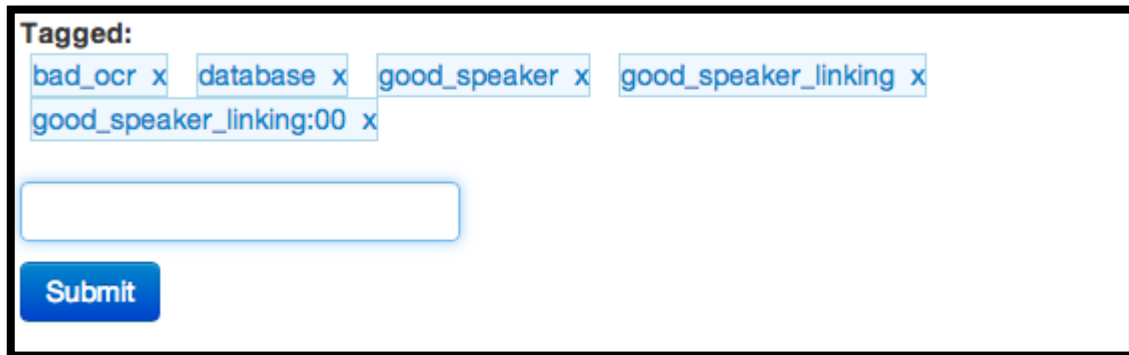
The last major enhancement to the UI and user friendliness of the web pages is the way we are currently allowing the user to edit the tags of a video. Previously this was done through the use of a textarea where all the tags were put on the textarea, each tag on one line. If a user wished to edit them they would need to edit them straight on the textarea and ensure that they are kept each in one line and there aren't any spaces as there is no validation. The image below shows how the editing of tags was previously done.



The new method of editing tags that I implemented is very different. All the tags are listed out with a cross next to them that when clicked will delete that tag. Below the list is a "Add a tag" button, this is all shown below.

Clicking the button will reveal a textbox to enter the new tag and a button to submit.

Unlike the previous method, this method has validation to ensure a good tag is added to the list. Validation includes:

1. Presence

You must enter a value

2. No spaces

test space — You cannot have white spaces

3. No duplicates

bad_ocr — You cannot have the same tag twice

Once an input from the user is given that passes all these validation tests, it is added to the list.

This is all achieved via ajax, so no reload to the page is necessary. This allows for a very enjoyable and pleasant user experience.

## ADD VIDEO FUNCTIONALITY

The next task that I tackled was that off implementing the "Add Video" functionality. I started off by creating a page that contains the form to add a video, the page is displayed below.



In this page the user can enter the title and the summary they wish to have for the video as well as choosing the video itself. Once the user has finished filling out the form they click the Add Video button to submit the form. This form has different validation requirements that must be met for the video to be uploaded.

1. Presence



2. File Format



The user is able to choose any file they wish but the form will not submit unless the file format matches one of the formats listed there. The list of accepted extensions are kept in one centralized place so that changing the list will change the html and the validation at the same time, making it very easy to add and remove accepted extensions.

Once a user successfully passes the validation of the form the form is submitted and the video starts uploading. Depending on the size of the file it may take longer to upload. I have implemented a progress bar that shows once the form is submitted. The progress bar shows how much of the file is uploaded to the user. This ensure that the user is kept notified of the progress and also enhances user experience. Below are screenshots taken of the progress bar at different stages.

Once the upload is complete, the user is notified and a button is presented to them to take them to the page which will show them the newly uploaded video.



The video however has just been uploaded and has not gone through any processing for entity extraction. A video that has not been processed yet will have the processed flag not set and the user will be shown if the video has been processed or not. The image below shows what the user is shown when they click the button. The page they are shown is the video play page with the video that was just uploaded and the text to show that it has not yet been processed. The video is also present on the index page. Again, on this shows that the video is still to be processed.





Now that the uploaded video is saved into the database and shown to the user the next step was to send the video to the High Performance Computer (HPC) to undertake the processing on that video. However there were technical difficulties in setting up my account on the HPC and thus I moved onto the next task.

## ENTITY DISPLAY ENHANCEMENTS

The next step was to implement the entity display enhancements. This task involves extracting the information regarding the current speaker and current faces and displaying it on the page. The system stores information that relates an entity to a time in seconds in a video. This information was needed to calculate who is the current speaker and who are the current faces at a particular point in time for a particular video that the user has decided to view. For the faces more information was necessary including the position of the face in the video. This information was stored and various calculations were done to determine where the face was relative to the particular web page.

The current speaker's image and id is displayed underneath the video. As the video plays, the current speaker is calculated from stored values and shown to the user. When a user hovers over the image, more information is displayed about the entity.



On the other hand the current faces are displayed straight onto the video. A hollow rectangle is place on top of the video to show where on the video the face is shown. Calculations are made to determinate the position of the face on the video as well as the height and the width. On the bottom the box is the image of the current face and again, like with the speaker, once a user hovers over that image, more information is displayed about that entity.

This section will discuss the projections for the next period in terms of tasks, which is semester 2.

Having successfully completely the styling enhancements and entity display enhancements and partly completed the add video functionality, I have two options.

I could continue working on the add video functionality or I could start working on the desirable task of implementing authentication and authorization. Having had a large success in this semester with completing tasks, I believe that completing even the desirable task is achievable.

However, the first task I believe I should start working on is that of continuing the implementation of the add video functionality.  The current state of the task is that a user can upload a video using the add video form discussed in the previous section and this video, along with the title and the summary are saved into the database and the video stored in the file system.

The next step is to implement the file transfer from the webdemo server to the High Performance Computer (HPC).  The HPC has a protocol that it follows to undertake entity extraction. Entity extraction consists of various steps and the HPC has all these steps stored and ready to be run. Nevertheless, especially the face and speaker clustering steps are implemented to have all the video's as inputs to be able to successfully cluster the entities. This works fine when the entity extraction steps are done against all the video's in the database and before the implementation of the add video functionality was in place was the only way the entity extraction was undertaken. However, running the process against all the video's in the database every time a user adds one single video is very expensive in terms of time and resources. A major sub task of the add video functionality task is to change the implementation of the entity extraction process to allow for only one video to be the input.

Once the implementation has been successfully changed and tested, the next step is to take the video, along with all the data pertaining to the entities in the video that have been extracted, are all sent back to the webdemo server to store both in the database and in the file system. This would complete the add video functionality.

Moving on from that, I would now have time to commence work on the desirable task of implementing authentication and authorization into the system. DJango comes with an inbuilt admin page and this task would involve delving into how django does this and how to tap into this functionality. Once I gain a better understanding I can begin to allow and disallow certain actions depending on whether or not a user is signed in and furthermore, whether or not the user has the permission to perform the task. Initially this could be in the form of only allowing a signed in user to edit tags on a video and a guest user (i.e. someone who has not identified themselves via login) only to view the tags of a video but not have the functionality and ability to edit them in any way.

The next page shows two gantt charts, one for semester 1 and one for semester 2. These are very similar to the ones in the project proposal but have been updated according to what has happened in semester 1. This means that the semester 1 chart shows exactly what tasks I was working on and when, whereas the semester 2 chart shows the updated version of my projections for semester 2 which have certainly changed due to the results that have risen from semester 1.

## SEMESTER 1 - BEB801 (COMPLETED)

| Task | Sub-Task | Project Week 1 24/02/14 Week 1 | Project Week 2 03/03/14 Week 2 | Project Week 3 10/03/14 Week 3 | Project Week 4 17/03/14 Week 4 | Project Week 5 24/03/14 Week 5 | Project Week 6 31/03/14 Week 6 | Project Week 7 07/04/14 Week 7 | Project Week 8 14/04/14 Week 8 | Project Week 9 21/04/14 Mid-Sem Break | Project Week 10 28/04/14 Week 9 | Project Week 11 05/05/14 Week 10 | Project Week 12 12/05/14 Week 11 | Project Week 13 19/05/14 Week 12 | Project Week 14 26/05/14 Week 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **UI Enhancements** | Clean up | █ | | | | | | | | | | | | | |
| | Bootstrap Integration | | █ | █ | | | | | | | | | | | |
| | Tag CRUD | | | | █ | | | | | | | | | | |
| **Add Video** | Page setup | | | | | █ | █ | | | | | | | | |
| | Back-end setup | | | | | | | █ | █ | | | | | | |
| | File transfer to Lyra | | | | | | | | | | | | | | |
| | Processing | | | | | | | | | | | | | | |
| | File transfer to server | | | | | | | | | | | | | | |
| | Back-end setup | | | | | | | | | | | | | | |
| **Entity display enhancements** | Current speaker label | | | | | | | | | █ | █ | █ | | | |
| | Face tracking html element | | | | | | | | | | | | █ | █ | █ |
| **Implementing Authentication and Authorization** | Authentication | | | | | | | | | | | | | | |
| | Authorization | | | | | | | | | | | | | | |

## SEMESTER 2 - BEB802 (TODO)

| Task | Sub-Task | Project Week 15 21/07/14 Week 1 | Project Week 16 28/07/14 Week 2 | Project Week 17 04/08/14 Week 3 | Project Week 18 11/08/14 Week 4 | Project Week 19 18/08/14 Week 5 | Project Week 20 25/08/14 Week 6 | Project Week 21 01/09/14 Week 7 | Project Week 22 08/09/14 Week 8 | Project Week 23 15/09/14 Week 9 | Project Week 24 22/09/14 Week 10 | Project Week 25 29/09/14 Mid-Sem Break | Project Week 26 06/10/14 Week 11 | Project Week 27 13/10/14 Week 12 | Project Week 28 20/10/14 Week 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **UI Enhancements** | Clean up | | | | | | | | | | | | | | |
| | Bootstrap Integration | | | | | | | | | | | | | | |
| | Tag CRUD | | | | | | | | | | | | | | |
| **Add Video** | Page setup | | | | | | | | | | | | | | |
| | Back-end setup | | | | | | | | | | | | | | |
| | File transfer to Lyra | █ | █ | | | | | | | | | | | | |
| | Processing | | | █ | █ | █ | █ | | | | | | | | |
| | File transfer to server | | | | | | | █ | █ | | | | | | |
| | Back-end setup | | | | | | | | | █ | █ | | | | |
| **Entity display enhancements** | Current speaker label | | | | | | | | | | | | | | |
| | Face tracking html element | | | | | | | | | | | | | | |
| **Implementing Authentication and Authorization** | Authentication | | | | | | | | | | | █ | █ | | |
| | Authorization | | | | | | | | | | | | | █ | █ |

## LITERATURE REVIEW

As discussed, the system is a website which showcases a lot of SAIVT's work. SAIVT has taken numerous news related video's and has run them through many different script that they have engineered to extract entities from videos. These script include a combination of various speaker extraction, face extraction and optical character recognition scripts as well as speaker and face clustering to link the video's together as best as possible. I will now discuss some of my findings on both speaker and face recognition.

### SPEAKER RECOGNITION:

Speaker recognition is defined as not only detecting when an individual is speaking but moreover the identification of the person, the individual, who is speaking.  Over the years there has been much research in the area of the use of machines to automate speaker recognition (Campbell, J., 1997). There have been a vast quantity of algorithms and techniques implemented for automated speaker recognition, even over coming major difficulties like background noise (i.e. recognizing the speaker in a noisy environment), some more robust than others (Campbell, J., 1997).

A method used for speaker recognition, and one that is highly researched and used at SAIVT is that of Joint Factor Analysis (JFA) Modeling. This technique involves combining the appropriate estimation of speaker variability and session variability subspaces. (Vogt, R., et al, 2008). Each subspace are then "trained on a database containing a large number of speakers each with several independently recorded sessions." (Vogt, R., et al, 2008).

A similar method is called "i-vector speaker modeling" (Ghaemmaghami, H., et al, 2013) and also this is heavily researched and used at SAIVT. I-vector speaker modeling differs from JFA modeling in that it "represent the GMM super-vector by a single total-variability space" (Kanagasundaram, A., et al, 2011). This total variability space contains both the speaker and the channel variability but requires "additional intersession compensation approaches" (Kanagasundaram, A., et al, 2011). There have been some existing approaches that have shown effective for this task including WCCN, LDA and NAP (Kanagasundaram, A., et al, 2011).

### FACE RECOGNITION

Face recognition has been a hot topic, especially in the field of Human Computer interactions (Castrillon-Santana, M., et al, 2008). Face recognition was defined as "visual pattern recognition problem where the face, represented as a three-dimensional object that is subject to varying illumination, pose, expression and other factors, needs to be identified based on acquired images" (Jain, A., 2005).

The Viola-Jones object  detection framework has drastically increased the number of facial processing applications (Castrillon-Santana, M., et al, 2008). All these applications have provided very reliable approaches to the problem of face recognition.

The basic idea behind the Viola-Jones framework, more specifically the Haar-Cascade technique, is based on the idea of "boosted cascade of weak classifiers" (Castrillon-Santana, M., et al, 2008). It is the process of sequentially testing the image against numerous classifiers and rejecting those that fail. These classifiers include multiple "Haar-like features" which is an encapsulation of rectangular patterns (Kasinski, A, 2010). There are three types of patterns, edge features, line features and centre-surround features (Kasinski, A, 2010).

## SPEECH CLUSTERING

Speech Clustering is the process of linking the entities found in an audio file based on identity (Ghaemmaghami, H., et al, 2013). This is achieved "using speaker diarization and speaker linking" (Ghaemmaghami, H., et al, 2013).

The most extensively utilized method for speaker clustering in speaker diarization is agglomerative clustering with retraining (ACR) (Ghaemmaghami, H., 2013). Even though ACR is the most widely used technique, it was found to be very inefficient. For this reason, when conducting speaker linking in large dataset, another method is preferably used (Ghaemmaghami, H., 2013).

This more efficient method is called complete-linkage clustering (CLC). This is achieved through merging based on a highest similarity, or lowest distance, score. (Ghaemmaghami, H., et al 2013). Moreover, another technique for speech clustering is named cross likelihood ratio (CLR) (Ghaemmaghami, H., 2013). This method has shown to be robust and effective especially when "incorporated into the JFA modeling framework" (Ghaemmaghami, H., 2013).

## OPTICAL CHARACTER RECOGNITION (OCR)

Optical Character Recognition (OCR) is the process of converting typewritten or printed images to computer readable text (Mori S., 1997). This technique is mostly used in data entry, where it is necessary to import vast information stored in scanned images that needs to be computer readable. However, this process can be also very useful in entity extraction. (Mori S., 1997).

The most widely used engine for OCR is Tesseract, not only because it is open source, but also because it is very efficient (Smith R., 2007). Tesseract performs OCR by running various processes including: Line and Word Finding and Word Recognition (Smith R., 2007). Line and Word Finding are algorithms which essentially tell the computer where the lines of words are in the image. The Word Recognition process ultimately tries to understand exactly what computer readable word is a particular imaged word. (Smith R., 2007)

## REFERENCES

Campbell Jr, J. P. (1997). Speaker recognition: A tutorial. Proceedings of the IEEE, 85(9), 1437-1462.

Castrillón-Santana, M., Déniz-Suárez, O., Antón-Canalís, L., & Lorenzo-Navarro, J. (2008). Face and facial feature detection evaluation performance evaluation of public domain haar detectors for face and facial feature detection.

Ghaemmaghami, H. (2013). Robust automatic speaker linking and attribution.

Ghaemmaghami, H., Dean, D., & Sridharan, S. (2013, August). Speaker attribution of Australian broadcast news data. In Proceedings of the First Workshop on Speech, Language and Audio in Multimedia (SLAM): CEUR Workshop Proceedings, Volume 1012 (pp. 72-77). Sun SITE Central Europe.

Jain A., Li, S. (2005). Handbook of Face Recognition.

Kanagasundaram, A., Vogt, R., Dean, D. B., Sridharan, S., & Mason, M. W. (2011, August). I-vector based speaker recognition on short utterances. In Proceedings of the 12th Annual Conference of the International Speech Communication Association (pp. 2341-2344). International Speech Communication Association (ISCA).

Kasinski, A., & Schmidt, A. (2010). The architecture and performance of the face and eyes detection system based on the Haar cascade classifiers. Pattern Analysis and Applications, 13(2), 197-211.

Mori, S., Nishida, H., & Yamada, H. (1999). Optical character recognition. John Wiley & Sons, Inc..

Smith, R. (2007, September). An Overview of the Tesseract OCR Engine. In ICDAR (Vol. 7, pp. 629-633).

Vogt, R. J., Baker, B. J., & Sridharan, S. (2008). Factor analysis subspace estimation for speaker verification with short utterances.

Waheed, K., Weaver, K., & Salam, F. M. (2002, August). A robust algorithm for detecting speech segments using an entropic contrast. In Circuits and Systems, 2002. MWSCAS-2002. The 2002 45th Midwest Symposium on (Vol. 3, pp. III-328). IEEE.