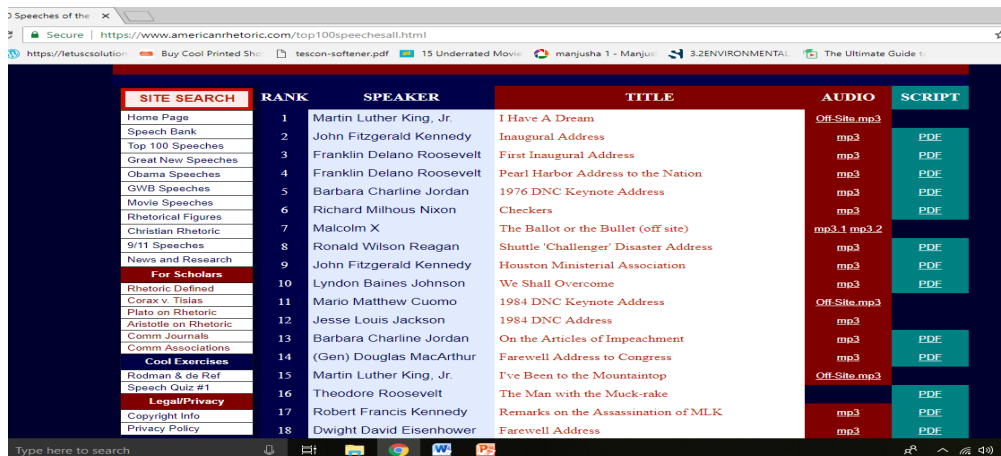USC Viterbi School of Engineering

# DeciQuest

Mansi Ganatra (4669963813)

Radhika Rao Annamraju (2800000280)

## 1 Project Idea

DeciQuest is a web application that allows the user to search over a collection of audio files such as all the iconic speeches from the past as well as the present times, podcasts, interviews etc. using the metadata associated with it. It also allows users to download the audio files and their PDF transcripts.

## 2 Data Source

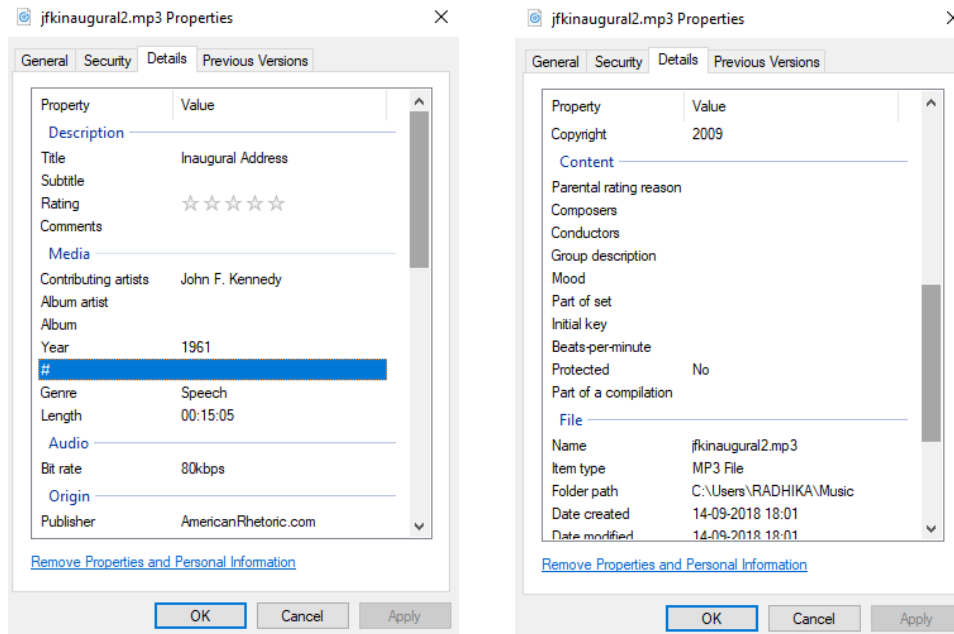We have obtained the audio files from the online audio repository:
https://www.americanrhetoric.com/top100speechesall.html

# 3  DETAILS OF METADATA

The primary metadata that we have extracted is the name of the artist, original date, file size, duration, genre, copyrights, title, type of the audio file etc. We also parsed the PDF transcripts of the speeches to obtain the actual contents.



# 4  PROJECT IMPLEMENTATION

## 4.1  TECHNOLOGIES

Python, JavaScript, Jquery, PyTika, HTML, CSS, Booststrap

## 4.2  METADATA EXTRACTION

We have used PyTika - Apache Tika library for python to extract the metadata from audio files. Apache Tika also handles extraction of content from PDF - hence it was an ideal choice.

## 4.3  KEYWORDS INDEX

We have also created a reverse index for keywords-documents to handle the search functionality of the application.
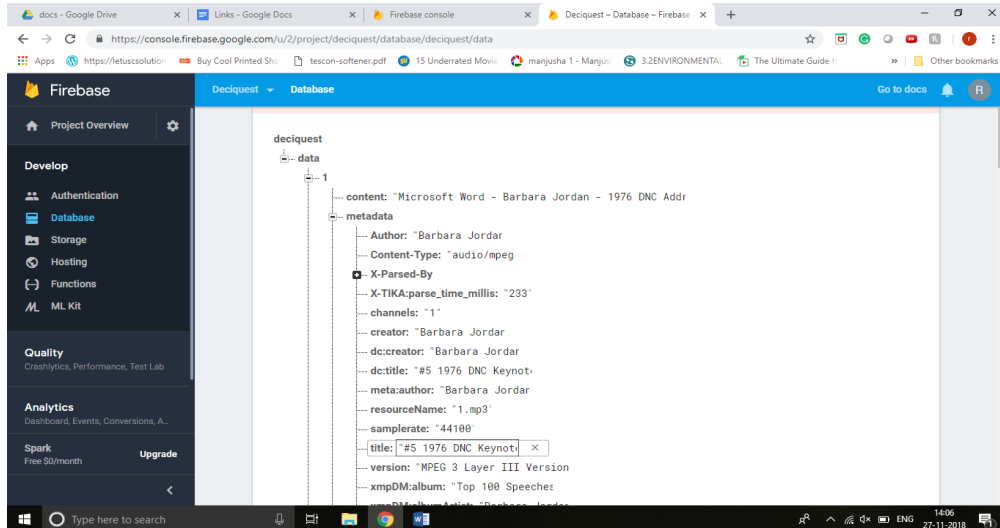
## 4.4  UPLOAD TO FIREBASE

We have written python scripts to upload the metadata and keywords to Firebase.

## 4.5 Codebase

Our code is present at: https://github.com/mansiganatra/DeciQuest

## 5 Metadata Storage

The extracted metadata and Keywords are stores in Firebase at: https://deciquest.firebaseio.com/



## 6 Data Storage

The audio fles and the PDF transcripts are stored in Amazon S3 at: http://deciquest.s3.us-east-2.amazonaws.com/

# 7 USER INTERFACE

The user interface for the application has been developed using Javascript, Jquery, HTML, CSS and Bootstrap. The UI allows the user to filter information using facets of Author, Year and Genre. The user can also sort the data using the same categories. The search bar enables the user to perform keyword search.

# 8 Individual Contribution

**Radhika Rao Annamraju**:

- Data collection from different sources

- Design and development of facets

- Faceted Search from UI

- Keyword Search from UI

**Mansi Ganatra**:

- Extract metadata and pdf content

- Generate keywords index

- Upload metadata to Firebase

- Upload data to S3