# Automatic Sychronisation of Subtitle Track With Live Audio

Joshua Fenech

MLDM
Université de Jean Monnet
Saint-Étienne, France

M1 Masters Thesis 2018
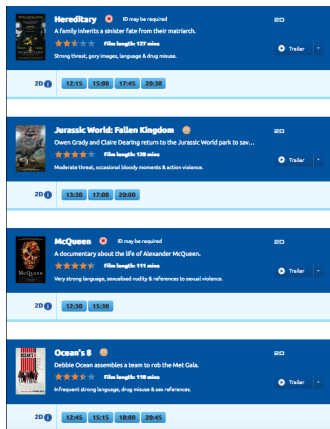
## Problem Motivation



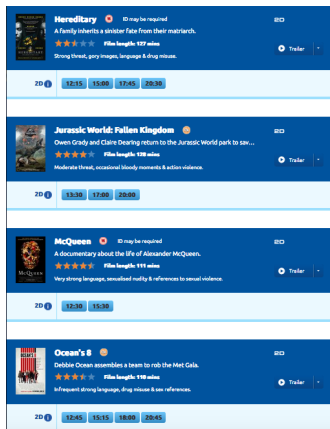FIGURE – Full FIlm Showings 1 Day
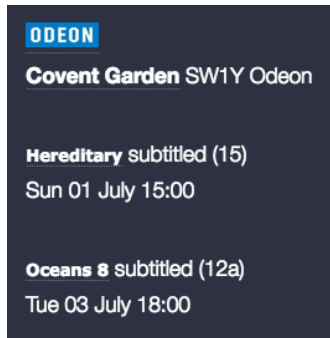
# Problem Motivation



FIGURE – Full Film Showings 1 Day



FIGURE – Subtitled FIlm Showings 1 Week

Figure – Full Film Showings 1 Day

Figure – Subtitled Film Showings 1 Week

Number of deaf people

Deaf people feeling excluded

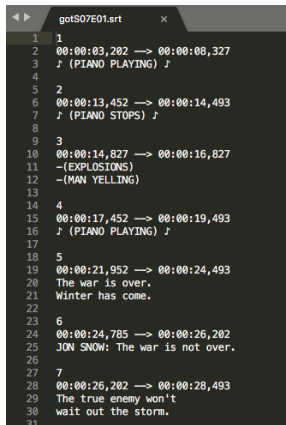Tourists

# Aim

- Develop a method to watch subtitles on a phone
- Problem : Synchronising the subtitles to the film
- Therefore, must identify the time in the film based on audio signals

# Prior Knowledge



FIGURE – SubRip .srt File

Figure – SubRip .srt File

Who here has pirated a film ?

Used subtitles ?

srt

Subrip files contain list of entries indicating start time, stop time and text to be displayed

## General Method

- Record audio, compressed using MP3
- Split signal into frames of duration 25ms - consider signal constant over this period
- Take frames every 10ms, so frames overlap
- Extract Mel Frequency Cepstral Coefficients (MFCC's) from each frame
- Use MFCC's as predictive feature of whether speech is present in a frame or not
- Match these predictions to the truth array, defined by a srt file

## Prior Knowledge

- How do you use subtitles on a laptop?
- SubRip Subtitle file (.srt)

Introduction
**Preprocessing**
Learner
Synchronisation

Audio Features
Probability array
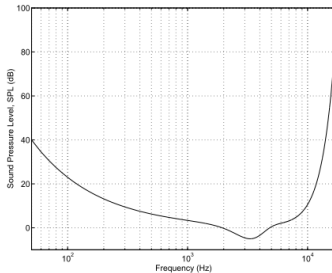
## MP3 Compression



FIGURE – Frequency response
of human hearing [3]. Curve
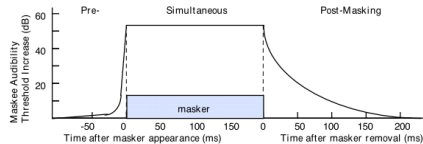indicates amplitude required to
detect tone at a given
frequency.



FIGURE – Frequencies masked
by more prevalent
frequencies[3].

Introduction
**Preprocessing**
Learner
Synchronisation

Audio Features
**Probability array**

- Create array of length appropriate to video with each entry corresponding to a frame

- Compare entries of srt to these start/stop times and ascribe a 1 if subtitles are present

Introduction
**Preprocessing**
Learner
Synchronisation

Audio Features
**Probability array**

**Algorithm 1** pb_array_fill

1: **procedure**
2:     $i \leftarrow 0$
3:     $j \leftarrow 0$
4:     $m \leftarrow pb\_array\_length$
5:     $n \leftarrow subs\_array\_length$
6:     **while** True **do**
7:         **if** $i > m$ **then**
8:         **if** $j > n$ **then**
9:         **if** *pb_array[i] start time* $\geq$ *subs[j] start time* **then**
10:             **if** *pb_array[i] end time* $<$ *subs[j] end time* **then**
11:                 $pb\_array[i] \leftarrow 1$
12:                 $i \leftarrow i + 1$
13:             **else**
14:                 $j \leftarrow j + 1$
15:
16:

Introduction
**Preprocessing**
Learner
Synchronisation

Audio Features
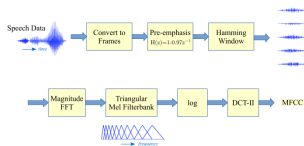Probability array

# MFCC Audio Features



FIGURE – Steps of MFCC[2]

- Process based on psychoacoustics to represent features most important to human hearing
- Split audio file into small sections, consider features constant over this period of time
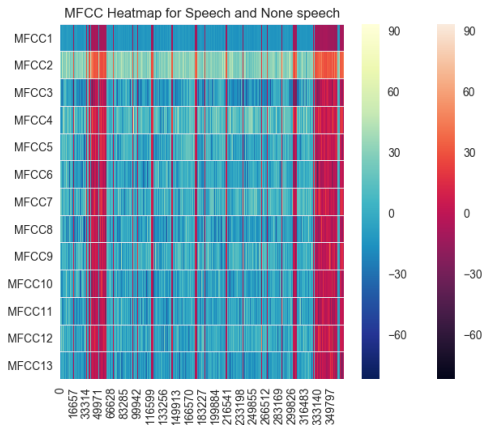- Apply a series of transformations
- Reduce stuff

Introduction
Preprocessing
Learner
Synchronisation

Audio Features
Probability array

# MFCC Audio Features



FIGURE – MFCC's Game of Thrones

## Learner Architecture



FIGURE – Model architecture[4]



FIGURE – 1d convolutions, no padding [1]

# Results



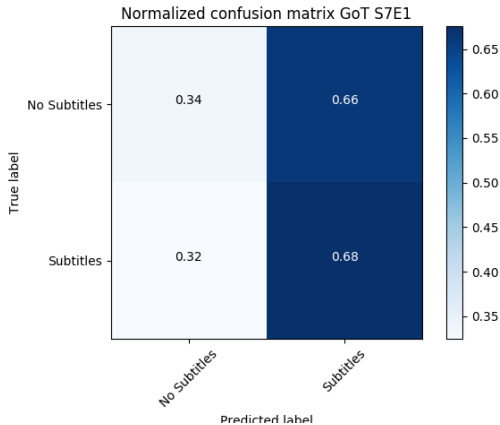Normalized confusion matrix GoT S7E1

FIGURE – Confusion Matrix Game of Thrones
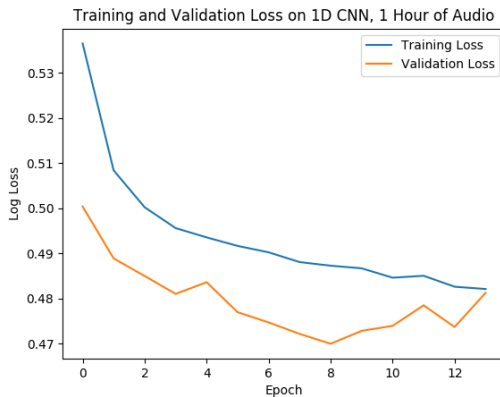
# Results



FIGURE – Training error
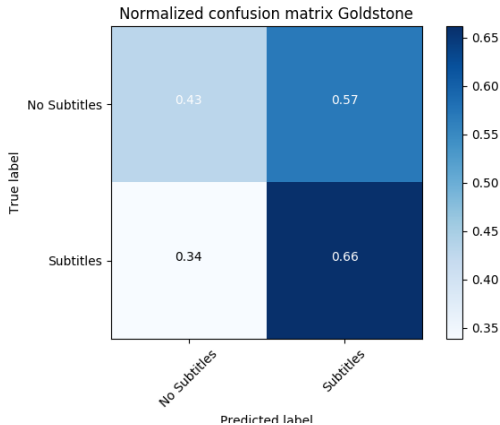
# Results



FIGURE – Test Time

# Results



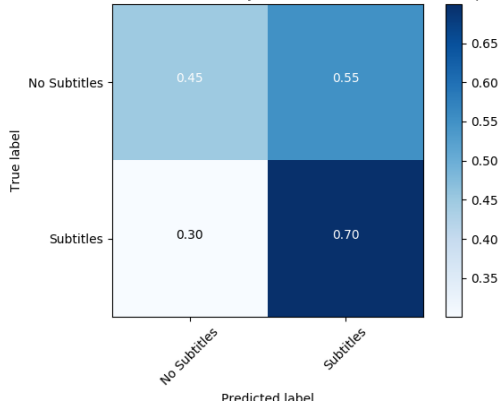Normalized confusion matrix noisy Game of Thrones Season 7 Episode 1

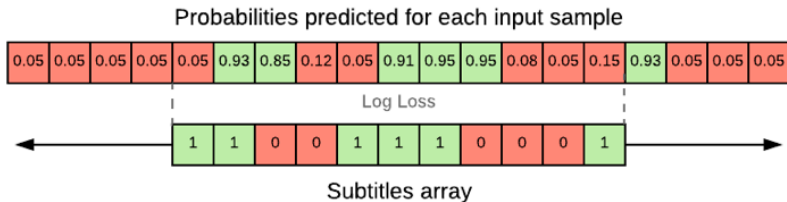FIGURE – Test Time

# Array Matching



FIGURE – Match predictions with truth array using log loss [4]

## Synchronisation

- Access to dataset granted incrementally as new audio is recorded
- Initially attempted to match a window of predicted probabilities with a similar array generated from srt
- Problem : Beginning of film often has no subtitle
- Solution : Continue recording data until speech is detected, and identify this as start of subtitle track

# Future Work

- Improve accuracy
- Remove nonspeech subtitles
- More efficient search algorithm
- Implement multithreading so that audio can be recorded and features extracted concurrently
- Alternative languages