Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

# Real-time Onset Detection in Musical Signals
## EE779 : Course Project

Yash Bhalgat (13D070014)
Kalpesh Patil (130040019)

October 27, 2016

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

# Overview

1. Problem Statement
2. Motivation
3. Background
4. Multi-step approach
   - Onset Detection Function
   - Peak Detection
   - Thresholding
5. Methods
   - Energy Method
   - Spectral Difference Method
   - Complex Domain Method
   - Linear Prediction Method
   - Sinusoidal Modelling
6. Observations
7. Results

## Problem Statement

The aim of the project is to:

- Study and present various state-of-art methods used in **real-time onset detection** in audio (music) signals.
- Formulation of ODFs (Onset-detection functions) by different techniques and using their peaks for the localization of onsets in the signals.

## Motivation

- Onset detection is very useful in temporal segmentation of audio signals, which is an important step in **beat-synchronous analysis**

- Onsets also play a crucial role in **score following**. It is the synchronisation of a computer with a performer playing a known musical score.
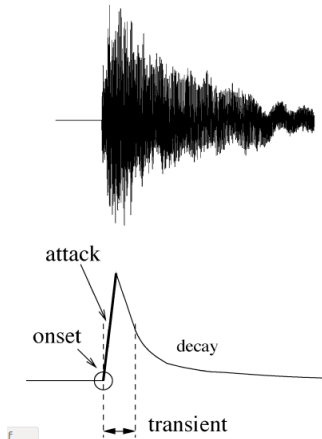
Score following is still an active area of research which stands at the heart of *AI, pattern recognition* and *signal processing*!

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

## Background

**Onsets**: instant marking the beginning of a transient or a note.
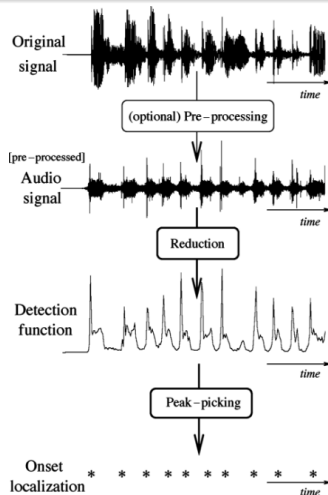
**Transient**: short interval during which the signal behaves in a relatively unpredictable way.

**Attack**: The time interval during which the amplitude envelope increases

Problem Statement
Motivation
Background
**Multi-step approach**
Methods
Observations
Results

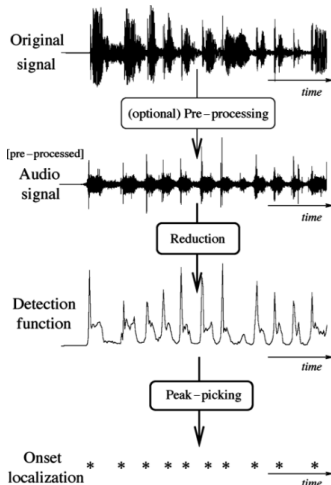Onset Detection Function
Peak Detection
Thresholding

# Onset Detection Function

- **Onset Detection Functions**
  - Highly subsampled detection sequence
  . showing the occurrence of transients

  - A measure of the unpredictability of
  . that frame indicating onsets

  - Vector passed to peak detection
  . algorithm for onset detection

- **Peak Detection**

- **Threholding**

Problem Statement
Motivation
Background
**Multi-step approach**
Methods
Observations
Results

Onset Detection Function
Peak Detection
Thresholding

## Peak Detection

- **Onset Detection Functions**

- **Peak Detection**
  - Appropriate ODFs give identifiable peaks at the onsets or abrupt events
  - Identify local maxima in ODF
  - Comparing current ODF value with neighbouring sample values (this results in latency)
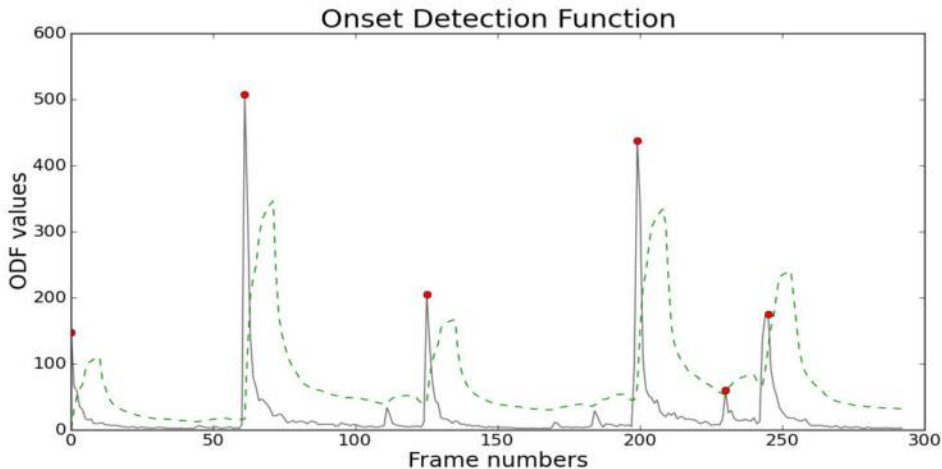
- **Threholding**

Problem Statement
Motivation
Background
**Multi-step approach**
Methods
Observations
Results

Onset Detection Function
Peak Detection
Thresholding

## Thresholding

- **Onset Detection Functions**

- **Peak Detection**

- **Threholding**

    $$thresh_n = \lambda \times median(O[n_m]) + \alpha \times mean(O[n_m]) + N$$

    $O[n_m]$ *is the previous m values of the ODF at frame n*

Every ODF peak that is above this threshold is taken to be a
note onset. See figure below for demo illustration.

Problem Statement
Motivation
Background
**Multi-step approach**
Methods
Observations
Results

Onset Detection Function
Peak Detection
**Thresholding**

## Thresholding



Onset Detection Function

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Energy Method

- Most simple and computationally efficient

$$E[n] = \sum_{m=0}^{N} x(m)^2$$

- Onsets often have more energy than the steady-state component of the note, because that is when the excitation is applied

- Larger changes in the amplitude envelope of the signal are expected to coincide with onset locations

$$ODF_E(n) = |E(n) - E(n-1)|$$

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Spectral Difference Method

- Successful in detecting onsets in polyphonic signals (multiple notes simultaneously) and 'soft' onsets created by instruments which do not have a percussive attack

$$X(k, n) = \sum_{m=0}^{N-1} x(m)w(m)e^{\frac{-2j\pi mk}{N}}$$

- Spectral difference ODF (ODFSD) is calculated by examining frame-to-frame changes in the Short-Time Fourier Transform

$$ODF_{SD}(n) = \sum_{k=0}^{N/2} ||X(k, n)| - |X(k, n-1)||$$

Problem Statement
Motivation
Background
Multi-step approach
**Methods**
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Complex Domain Method

- Combined magnitude and phase information
- Magnitude prediction remains same. Assuming a constant rate of phase change between frames for phase prediction:

$$\hat{R}(k, n) = |X(k, n-1)|$$

$$\hat{\phi}(k, n) = princarg[2\phi(k, n-1) - \phi(k, n-2)]$$

- Euclidian distances in predicted and actual phasors gives ODF:

$$\Gamma(k, n) = \sqrt{\hat{R}(k, n)^2 + R(k, n)^2 - 2R(k, n)\hat{R}(k, n)\cos(\phi(k, n) - \hat{\phi}(k, n))}$$

$$ODF_{CD}(n) = \sum_{k=0}^{N/2} \Gamma(n, k)$$

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Linear Prediction Method

- LP "error"' is useful in creating ODFs.

- ODF is then the absolute value of the differences between the actual frame measurements and the LP predictions

- The ODF values are low when the LP prediction is accurate, but larger in regions of the signal that are more **unpredictable**, which should correspond with note **onset locations**.

- All the previous methods can be combined with LP prediction to give better results.

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Linear Prediction Method

$$\hat{x}(n) = \sum_{k=1}^{p} a_k x(n-k)$$

Hence, if we apply LP to each of the methods, we get:

$$ODF_{ELP}(n) = |E(n) - P_E(n)|$$

$$ODF_{SDLP}(n) = \sum_{k=0}^{N/2} ||X(k,n)| - |P_{SD}(k,n)||$$

$$ODF_{CDLP}(n) = \sum_{k=0}^{N/2} ||\Gamma(k,n)| - |P_{CD}(k,n)||$$

Problem Statement
Motivation
Background
Multi-step approach
**Methods**
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
**Sinusoidal Modelling**

## Sinusoidal Modelling

**Fourier's Theorem**: Any periodic waveform can be modelled as the sum of sinusoids at various amplitudes and harmonic frequencies.

**Additive synthesis**: adding together many sinusoidal components modulated by relatively slowly varying amplitude and frequency envelopes

$$y(t) = \sum_{i=1}^{N} A_i(t) sin\left[\theta_i(t)\right]$$
$$where \quad \theta_i(t) = \int_0^t \omega_i(t)\, dt + \theta_i(0)$$

**Partial tracking**: Calculating these parameters for each frame is referred to as peak detection, while the process of connecting these peaks between frames is called partial tracking.

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Offline Processing Technique

- Transient signals in the time domain can be mapped onto sinusoidal signals in a frequency domain using DCT.
- Not suitable for real time since DCT frame length required is very large
- A multi-resolution sinusoidal model is then applied to the signal to isolate the harmonic component of the sound
- Onset Location is determined by abrupt increase in energy of the frame.

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

Energy Method
Spectral Difference Method
Complex Domain Method
Linear Prediction Method
Sinusoidal Modelling

## Online Processing Technique

*Proposed in the paper*

- During the steady state of a musical note, the harmonic signal component can be well modelled as a sum of sinusoids, whose amplitude and frequency are slowly evovling in time.

- Absolute values of the frame-to-frame differences in the sinusoidal peak amplitudes and frequencies should be quite low for steady state

- Amplitudes of detected sinusoidal partials increases during attack region

- Large amplitude and/or frequency deviations in the partials implies existence of Onsets
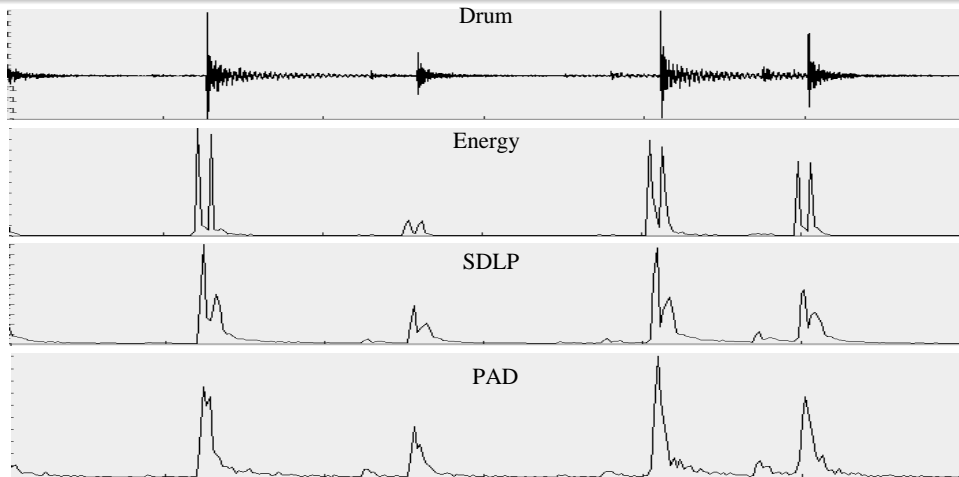
Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

## Pre-processing

- The signal is converted to overlapping, fixed-sized frames of audio, each having 4 buffers - 512 x 4 in duration.
- We used a 50% overlap for the frames and a samples.
- It consists of the most recent audio buffer which is passed directly to the algorithm, combined with the previous three buffers which are saved in memory.
- The time taken by the algorithm to process one frame of audio must be less than the duration of audio that is held in each buffer. Sampling rate $= 44$kHz and buffer size $= 512$ samples. So, the algorithm must be able to process a frame in 11.6 ms or less when operating.
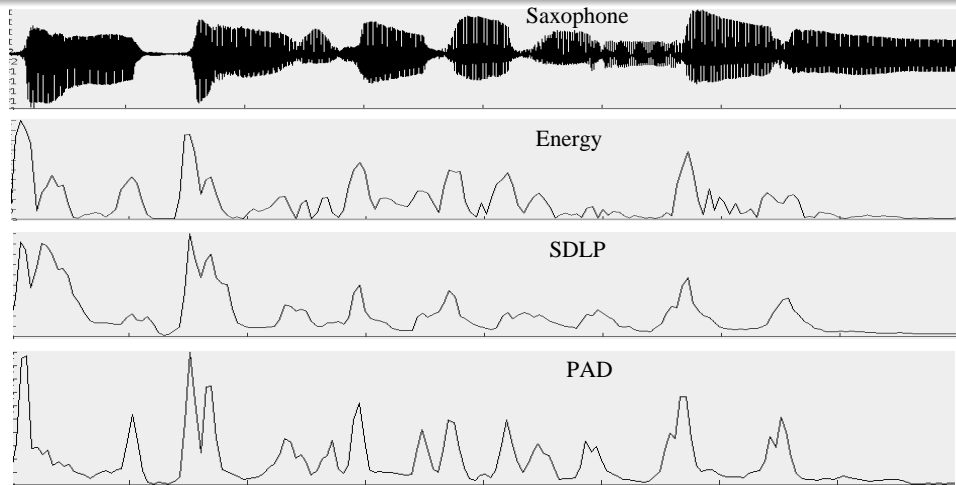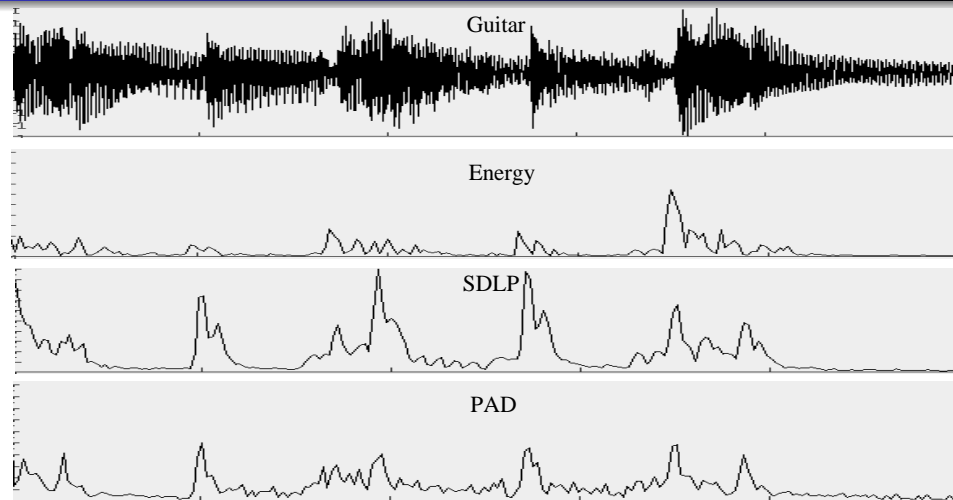
Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

# Pre-processing

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

## Observations

Problem Statement
Motivation
Background
Multi-step approach
Methods
**Observations**
Results

## Observations



Saxophone

Energy

SDLP

PAD

Problem Statement
Motivation
Background
Multi-step approach
Methods
**Observations**
Results

## Observations

Problem Statement
Motivation
Background
Multi-step approach
Methods
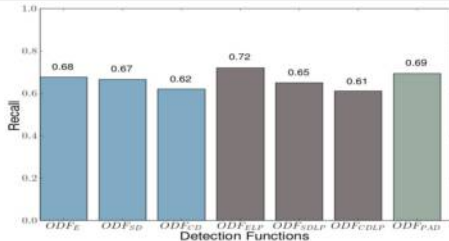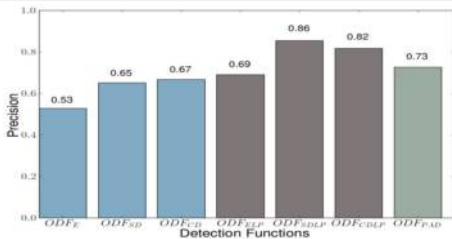Observations
Results

## Detection Results

- The detection accuracy of the ODFs was measured by comparing the onsets detected using each method with the reference samples in the Modal database.
- To be marked as 'correctly detected', the onset must be located within 50 ms of a reference onset.

$$P = \frac{C}{C + f_p}$$

$$R = \frac{C}{C + f_n}$$

$$F = \frac{2PR}{P + R}$$

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
**Results**

# Results

Problem Statement
Motivation
Background
Multi-step approach
Methods
Observations
Results

## References

[1] *John Glover, Victor Lazzarini, Joseph Timoney* - Real-time detection of musical onsets with linear prediction and sinusoidal modeling, EURASIP '11

[2] *Robert McAulay, Thomas Quatieri* - Speech Analysis/Synthesis based on sinusoidal representation, IEEE Transactions '86

[3] Sinusoidal Modelling - http://www.music.mcgill.ca/ ich/classes/dafx_book.pdf

[4] *Juan Pablo Bello, Laurent Daudet et al* - A Tutorial on Onset Detection in Music Signals, IEEE Transactions '05

# Thank You!