

# **Master's thesis:**

## Real-time signal synchronisation with acoustic fingerprinting

Ward Van Assche  
*Student information engineering technology*  
UGent

27 may 2016



Hi, my name is Ward Van Assche and I'm going explain to you what I've created in the research for my Master's thesis.

I'm going to start with a small demonstration. Firstly I will summarize what you are going to see during the demo.

## Demonstration: prelude

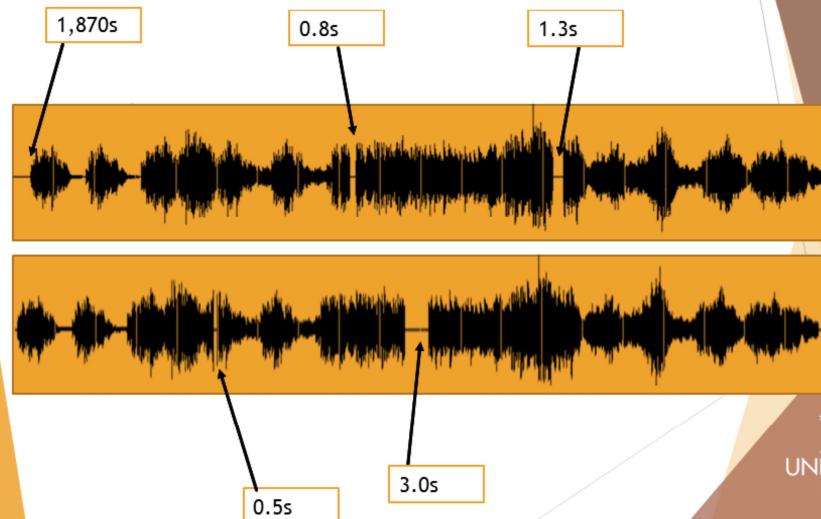
- ▶ 2 wave files
- ▶ To each file: added several moments of silence
- ▶ Audio is shifted to the right
- ▶ After each silence: latency changed
- ▶ Duration of waveforms (on next slide):  
90 seconds



2

For the demo I took two wavefiles and added some silences at random places in both files. After adding the silences the sound is shifted to the right, no samples are lost. Because of this the latency between the files will change after each silence.

## Demonstration: silences



3

On this slice you can see the added silences. You can use this slide to have an idea what is happening during the demo.

## Demonstration

- ▶ Max/MSP module:
  - ▶ Input: signals from wavefiles
  - ▶ Output: synchronized signals
- ▶ Synchronization:
  - ▶ No extra timing information
  - ▶ Only by using audio features
- ▶ First audiosignal: left speaker
- ▶ Second audiosignal: right speaker



4

The wavefiles will be used as input for the Max/MSP module I will show during the demo. The module is responsible for the synchronization and will send the synchronized signals as output. The module doesn't have access to extra timing information. The first synchronized signal will be sent to the left speaker. The second synchronized signal to the right speaker.

## Note

- ▶ Any latency change can be detected
  - ▶ (The algorithms do not just look for silence)
- ▶ The audiostreams do not have to be exactly the same
  - ▶ (This system also works with different recorded audiostreams)



5

Quite impressive isn't it? I just want to note that the used algorithms don't just look for silence. Any latency change caused by any reason can be detected. It's also important to know that the audiostreams don't have to be exactly the same. The system works with recorded streams from different kinds of microphones.

## Why is it useful?

- ▶ IPEM experiments
  - ▶ Audio
  - ▶ Video
  - ▶ Sensors
  - ▶ ...
- ▶ Problem: latency
  - ▶ Synchronisation: necessary



6

You are now probably wondering why this is useful. So, at IPEM a lot of experiments are done by using various sensors like accelerometers and pressure sensors. The problem is each sensor has its own unpredictable latency. Before analysis is possible the data has to be synchronized.

## The solution

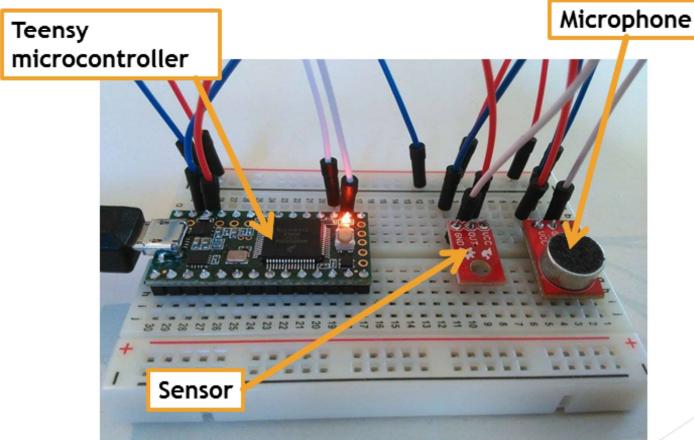
- ▶ Attach a microphone to each sensor
- ▶ Latency recorded sound = latency sensordata
- ▶ Microphone: records environment sound
- ▶ Latency recorded sound = latency sensordata
- ▶ Much easier to detect:
  - ▶ Recorded environment sound: almost the same for each recording



7

This can be done by attaching a microphone to each sensor so the latency of the sensor and microphone are the same. The difficulty of the initial problem can be reduced by recording the environment sound. This because the environment sound is the same for each microphone.

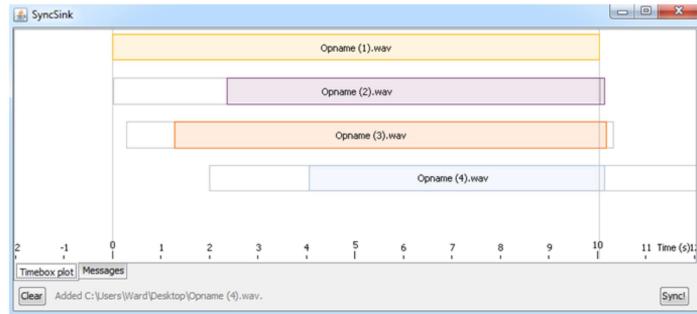
## Attaching sensor to microphone



8

This is an example of a sensor and microphone set-up. By attaching them to the same microcontroller the latency between them is negligible.  
(negligible)

## What about SyncSink? [2]



9

Some people probably already used SyncSink. It's a system using the same fundamentals I just explained. But there are some remarkable differences with what I've developed.

## What about SyncSink? [2]

- ▶ SyncSink:
  - ▶ Executed after the experiment
- ▶ Desirable:
  - ▶ Real-time synchronisation
  - ▶ A more user-friendly system (like a Max/MSP module)



10

SyncSink has to be manually executed after the experiment. This is not very user-friendly and can be time-consuming.

For my master's thesis Joren asked me to develop a system which can do this in real-time. After developing the system I wrapped the system in a Max/MSP module.

## How does it work?

- ▶ Acoustic fingerprinting [1]
  - ▶ Determining fingerprints of audio clip
  - ▶ Based on spectral peaks:
    - ▶ Generate spectrogram
    - ▶ Find peaks
    - ▶ Create fingerprints



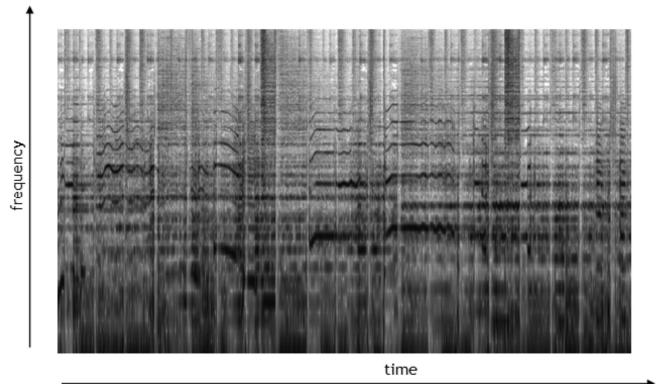
11

The latency between the audiostreams can be detected by using the acoustic fingerprinting algorithm.

In order to do this the fingerprints of each audio fragment have to be determined. This is done by connecting the spectral peaks of the spectrograms of the audio fragments.

An example will make this terrible explanation more clear.

## The spectrogram



time

frequency

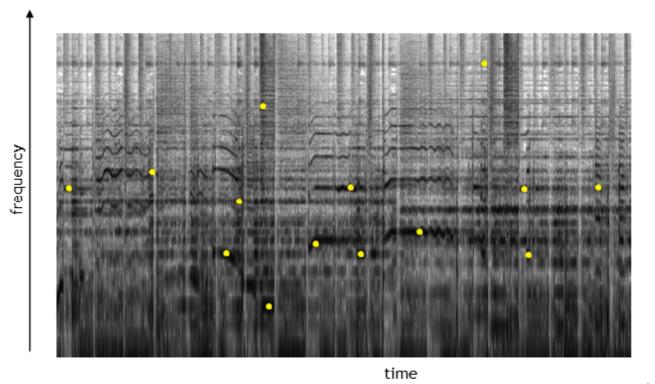


12

So, this is an example of a spectrogram. A spectrogram is a graph where the energy of each frequency-range is plotted against the time. On this slide the dark places have more energy than the less dark places.

The spectral peaks are places where the energy is higher than the energy of the neighbours. On the next slide I marked some spectral peaks.

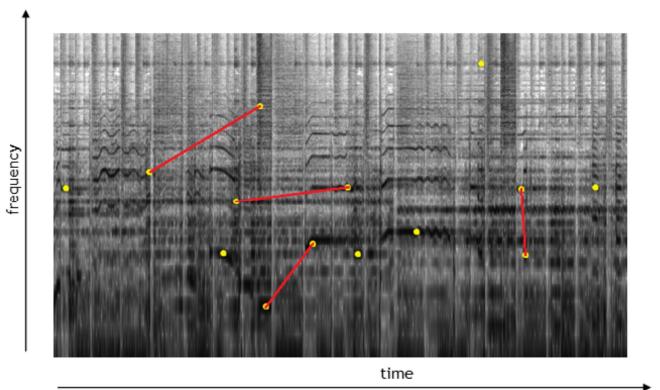
## Extraction of spectral peaks



13

Here you can see some spectral peaks. The fingerprints are nothing more than a connection between two spectral peaks. Which spectral peaks are connected depends of the parameters of the algorithm. On the next slide I have marked some possible fingerprints.

## Creation of the fingerprints



14

## Acoustic fingerprinting [1]

- ▶ Detecting the latency:
  - ▶ Each fingerprint contains timing information
  - ▶ Find matching fingerprints
  - ▶ Latency = difference between time of fingerprints
- ▶ Accuracy: default 32ms
  - ▶ Not sufficient



15

I forgot to say that each fingerprint contains information about the time where the spectral peaks occurred. When many fingerprints can be matched this means that the audio streams contain resembling elements. The latency between the audio fragments can be calculated from the timing information of the fingerprints.

The biggest problem of this system is its accuracy. Fortunately there is another algorithm which can solve this problem.

## Another technique: cross-covariance [2]

- ▶ Cross-covariance: degree of similarity between two signals
- ▶ Calculate for each shift
- ▶ Shift with highest crosscovariance value: latency
- ▶ Very accurate: < 1ms
- ▶ Pitfall: SLOW
  - ▶ First acoustic fingerprinting
  - ▶ Refine result with this technique

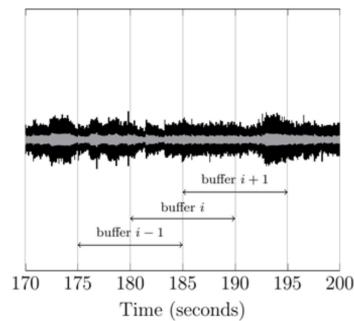


16

The cross-covariance is a calculation which measures the degree of similarity between two signals. By calculating this number for each possible shift between the signals it's easy to find out the latency. This algorithm is very accurate but also very slow. To solve we can first use the accoustic fingerprinting algorithm. Then the result can be refined by calculating the crosscovariance between the signals.

## Real-time: impossible

- ▶ Algorithms: need audio to analyze
- ▶ Real-time streams: buffered



17

It's practically impossible to synchronize the streams in real-time. The latency-detecting algorithms need some audio to analyse. That's why the streams are buffered.

## Buffering of streams

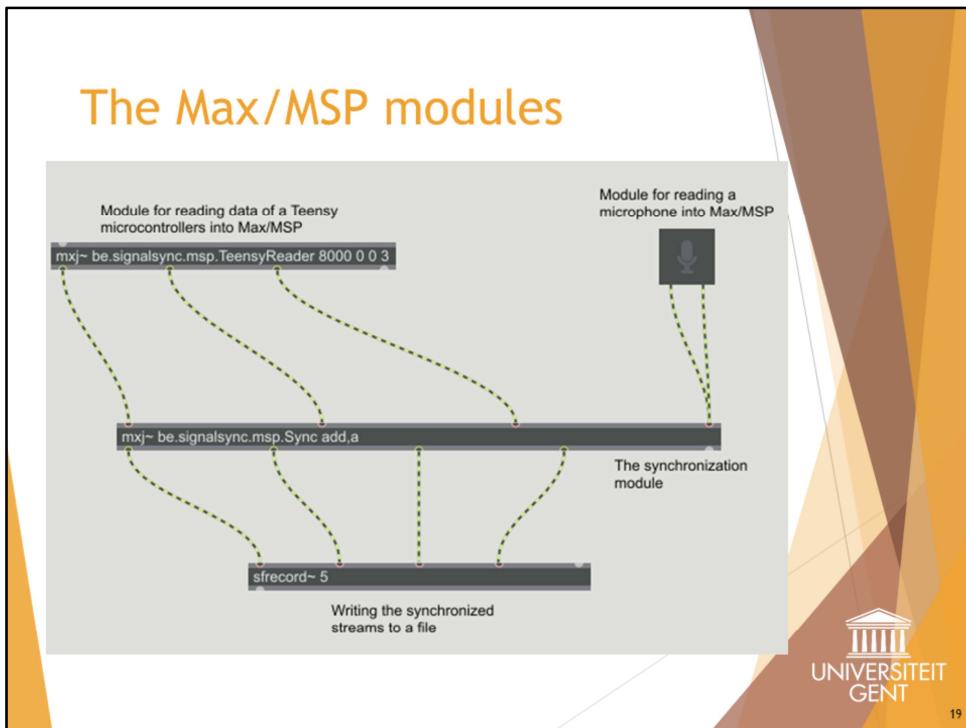
- ▶ Algorithms: executed on the consecutive buffers of each stream
- ▶ After determining latency:
  - ▶ Adding silence to audiostreams + attached datastreams
  - ▶ Streams: synchronized



18

The algorithms are executed on the consecutive buffers. The resulting latencies are converted to the amount of silence which has to be added to each stream. After adding the silence to the streams, the streams are synchronized.

## The Max/MSP modules



So this is an example of a Max/MSP patch which can be used during an experiment. The TeensyReader patch above reads 1 audiostream and 2 datastreams from a Teensy. The samples are sent to the Sync module. The data will be synchronized with the sound recorded with a microphone. The synchronized streams are written to a file by using the sfrecord module.

## Conclusion

- ▶ IPFM experiments:
  - ▶ Sensors attached to microphones
  - ▶ Audiostreams are almost the same
- ▶ Max/MSP patch:
  - ▶ Detects latency between audiostreams
  - ▶ Synchronizes datastreams
  - ▶ ...in real time



20

So never forget when you are doing an experiment to attach a microphone to your sensors. Because of the recorded environment audio is the same the latency can be detect easily. The latency is used to add whitespace to the attached streams resulting synchronized streams.

All this is possible in real-time.

## Future improvements

- ▶ Synchronization currently only possible with Max/MSP signals
- ▶ Impossible: synchronization of video or MIDI streams:
  - ▶ Should be converted to a Max/MSP signal.



21

The current synchronization system only supports Max/MSP signals. Everything can be synchronized but before it has to be converted to a Max/MSP signal. A possible improvement of the current system could be providing support for MIDI and video streams.

# Questions?



22

## References

- ▶ [1] Avery L. Wang. An Industrial-Strength Audio Search Algorithm.  
*In Proceedings of the 4th International Symposium on Music Information Retrieval (ISMIR 2003)*, pages 7-13, 2003.
- ▶ [2] Joren Six and Marc Leman. Synchronizing Multimodal Recordings Using Audio-To-Audio Alignment.  
*Journal of Multimodal User Interfaces*
- ▶ [3] Joren Six and Marc Leman. Panako - A Scalable Acoustic Fingerprinting System Handling Time-Scale and Pitch Modification.  
*In Proceedings of the 15th ISMIR Conference (ISMIR 2014)*, 2014.