

# Masterproef

## Realtime signaalsynchronisatie met accoustic fingerprinting

Ward Van Assche  
28 april 2016

Onderzoeksgroep: IPEM





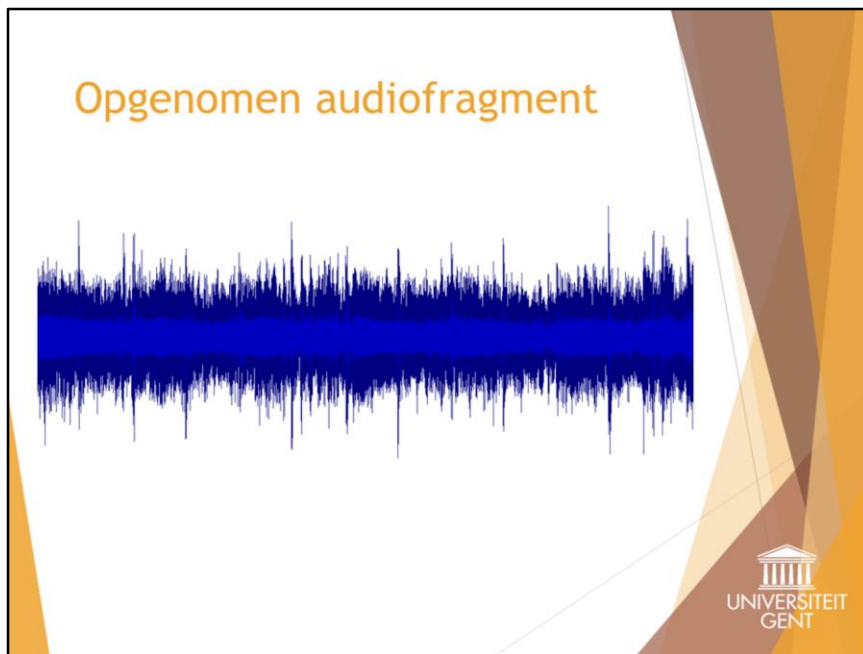
Wie kent er deze app? Iedereen waarschijnlijk?

En wie weet er hoe deze app er in slaagt om zo snel deze liedjes te detecteren?

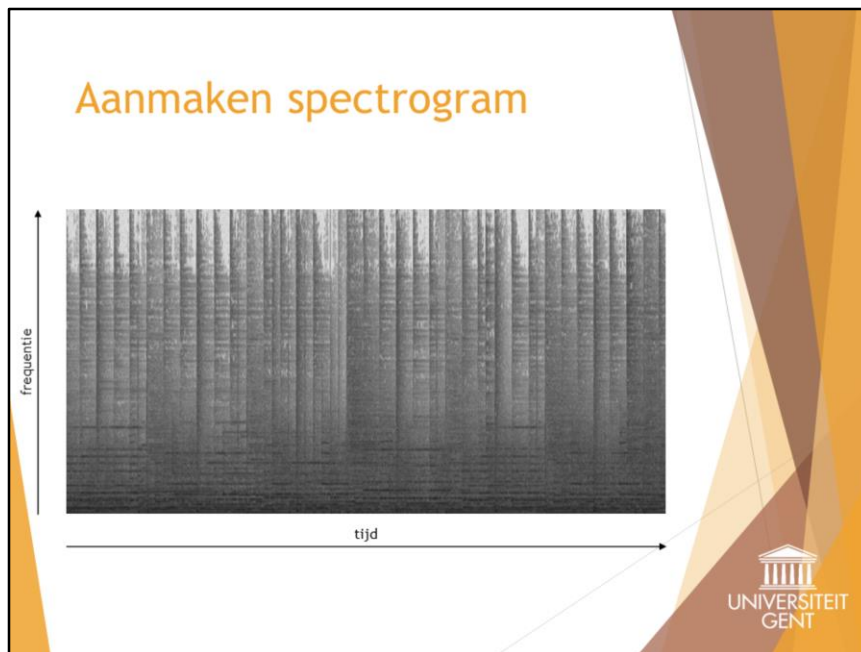


Dit gebeurt met accoustic fingerprinting.

Bij accoustic fingerprinting wordt er op basis van een audiofragment (de opname) een verzameling fingerprints berekend. Die verzameling fingerprints kan gezien worden als de DNA van het liedje. Het zoeken naar fingerprints gebeurt op basis van de spectrale pieken. Dit zijn pieken het spectrogram van de audio. Dit klinkt vrij complex maar dit is het eigenlijk niet. Dit zal duidelijker worden na een voorbeeldje.



Stel je zit op café en hoort een liedje waarvan je kost wat het kost de titel en band van wil te weten komen. Je pakt je smartphone en laat Shazam er op los. Er wordt een opname gemaakt van pakweg 10 seconden. Dit fragmentje wordt doorgestuurd naar de Shazam servers en daar begint het echte werk: het acoustic fingerprinting.

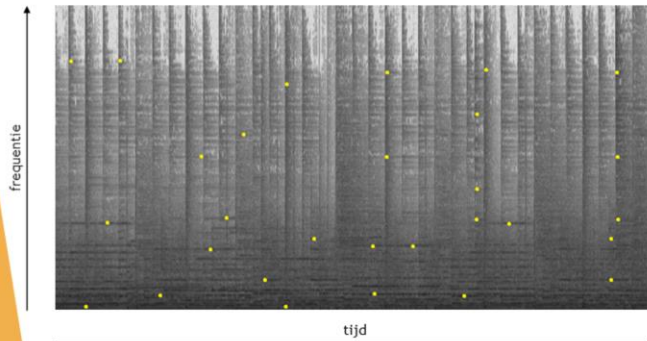


Eerst en vooral wordt er van het audiofragment een spectrogram aangemaakt. Dit is een soort van grafiek waarbij de frequentie is uitgezet ten opzichte van de tijd. Hoe donker een bepaald vlekje is, met hoe meer energie deze frequentie op dat moment voorkomt in de muziek.

Veel eigenschappen van de muziek zijn zichtbaar in het spectrogram. Uit de verticale lijnen kan je het ritme afleiden, de horizontale lijnen worden bepaald door de melodie.

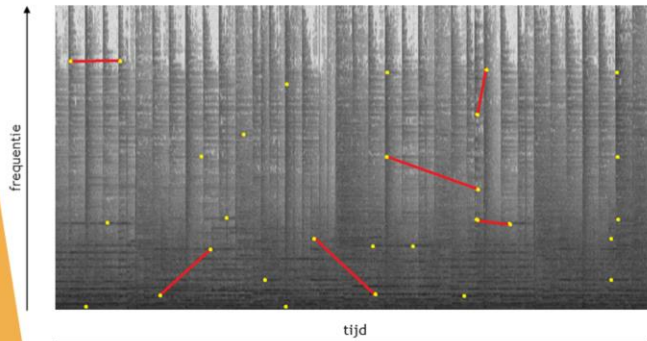
Daarstraks hebben ik gezegd dat fingerprints worden opgebouwd met behulp van spectrale pieken. Dit zijn plaatsen in het spectrogram waar de energie hoger is dan alle nabijgelegen plaatsen. In het spectrogram van onze opname heb ikzelf ook handmatig enkele spectrale pieken aangeduid.

## Extractie spectrale pieken



Nadat de spectrale pieken bepaald zijn kunnen de fingerprints worden aangemaakt. Een fingerprint is een verbinding tussen twee spectrale pieken. Hoe er precies beslist wordt welke spectrale pieken verbonden worden hangt af van heel wat parameters, daar ga ik niet dieper op ingaan.

## Aanmaken fingerprints



Nadat de fingerprints geconstrueerd zijn kan er gezocht worden naar welk liedje we nu eigenlijk hebben opgenomen. Hiervoor heeft shazam een enorme database met daarin de fingerprints van miljoenen liedjes. De fingerprints van het opgenomen geluidsfragment worden vergeleken met de fingerprints uit deze database, en zo probeert shazam te achterhalen welk liedje we hebben opgenomen. Elke fingerprint kan héél compact, als een getal, worden opgeslaan. Dit getal wordt een hash genoemd. Aangezien het vergelijken van getallen door computers heel snel kan gebeuren kan shazam zo vlug resultaten vinden.

## Signaalsynchronisatie?

- ▶ IPEM: onderzoek naar **musicologie** / **psychoakoestiek**
- ▶ Experimenten: invloed muziek op gedrag/beweging
  - ▶ Sensoren
  - ▶ Video-opnames
  - ▶ Probleem: latency
  - ▶ Synchronisatie is nodig



Allemaal heel interessant, maar wat heeft dit nu met mijn masterproef te maken? Wel, ik doe mijn masterproef bij het IPEM, dit is een onderzoeksinstelling aan universiteit gent, die onderzoek doet naar alles wat met muziek en audio te maken heeft. Er worden veel experimenten uitgevoerd waarbij er gekeken wordt naar wat de invloed is van muziek op gedrag of beweging van de mens. Om het gedrag of de beweging van personen in zo'n experiment te analyseren wordt er gebruik gemaakt van allerlei sensoren en video-opnames. Na zo'n experiment beschikken de onderzoekers over grote hoeveelheden data. Het grote probleem is dat elke sensor of opname een bepaalde vertraging of latency heeft. Voordat de gegevens door de onderzoekers geanalyseerd kunnen worden moet er voor worden gezorgd dat ze synchroon zijn.



## Synchronisatie: met audio

- ▶ Elke sensor of video-opname: koppelen met synchrone audiostream.
- ▶ Audiostream: opname van omgevingsgeluid
- ▶ Synchronisatie datastreams = synchronisatie audiostreams.



Bij het IPEM heeft men een synchronisatie systeem ontwikkeld dat gebruik maakt van audio. Daarom wordt er aan elke sensor een microfoon gekoppeld die het omgevingsgeluid opneemt. Ook wordt elke video-opname voorzien van een audiospoor. Door er voor te zorgen dat de opname zeer dicht bij de sensor gebeurt kan men er vanuitgaan dat de latency van de geluidsopname hetzelfde is als de latency van de datastream van de sensor.

Het probleem van het synchroniseren van de datastreams hebben we eigenlijk kunnen omzetten naar het synchroniseren van opnames van het omgevingsgeluid. Aangezien het omgevingsgeluid van elke sensor hetzelfde is: de afgespeelde muziek is het probleem veel eenvoudiger op te lossen.

## Synchronisatie: met audio

- ▶ Accoustic fingerprinting!
  - ▶ Geen match zoeken in database
  - ▶ Opnames van omgevingsgeluid met elkaar matchen
  - ▶ Nauwkeurig tot op enkele tientallen milliseconden

Heel toevallig bestaat er een methode waarmee audiofragmenten met elkaar kunnen worden vergeleken: accoustic fingerprinting! In tegenstelling tot de manier waarop Shazam het algoritme gebruikt gaan we de opnames niet matchen met een gigantische database maar worden de opnames van de verschillende sensoren met elkaar gematcht. Dit heb ik daarstraks niet vermeld maar een fingerprint bevat onder meer informatie over de tijd wanneer het voorkomt in het audiofragment. Wanneer verschillende fingerprints overeenkomen kunnen we deze informatie gebruiken om de latency tussen de opnames te bepalen.

## Synchronisatie: met audio

- ▶ Nog nauwkeuriger: kruiscovariantie
- ▶ Berekent mate van gelijkheid tussen twee signalen
- ▶ Zéér nauwkeurig:  $<1\text{ms}$
- ▶ Traag
  - ▶ Eerst acoustic fingerprinting
  - ▶ Verfijnen door berekenen van kruiscovariantie

Eigenlijk is deze methode nog niet echt nauwkeurig genoeg. Een resultaat tot op één milliseconde nauwkeurig is zeer gewenst. Daarom wordt het resultaat dat bepaalt is met het acoustic fingerprinting nog verfijnt met een tweede methode namelijk door het berekenen van de kruiscovariantie. Dit is een methode waarbij de gelijkheid tussen twee signalen berekend wordt. Aan de mate van gelijkheid wordt een bepaalde score toegekend. Door twee opnames over elkaar te verschuiven en telkens de kruiscovariantie te berekenen kunnen we zeer precies de latency tussen beide audiofragmenten bepalen. Deze berekening is zeer traag, daarom gebruiken we eerst acoustic fingerprinting om tot een ruw resultaat te komen. Vervolgens kan de latency verfijnt worden door de kruiscovariantie te berekenen.

**Mijn opdracht**

- ▶ Nu:
  - ▶ Alles handmatig, naverwerking
  - ▶ Niet gebruiksvriendelijk, tijdrovend
- ▶ Doel:
  - ▶ Realtime synchronisatie
  - ▶ Optimaliseren algoritmes
  - ▶ Ontwikkelen grafische interface

UNIVERSITEIT GENT

Op dit moment moeten deze methode's na een experiment handmatig op de geluidsopnames worden uitgevoerd. Dit is niet alleen erg ongebruiksvriendelijk maar ook zeer tijdrovend. Daarom heb ik de opdracht gekregen om met behulp van deze methodes een real-time systeem te bouwen. Dit systeem moet in staat zijn om tijdens het experiment de latency onmiddellijk te bepalen en de synchronisatie direct uit te voeren. Hiervoor heb ik verschillende wijzingen moeten aanbrengen aan de synchronisatiealgoritmes. Ook is het de bedoeling dat ik een plugin schrijf voor een audio verwerking software pakket dat het synchroniseren volledig automatiseert voor de onderzoekers van het IPEM.