



Faculteit Ingenieurswetenschappen en Architectuur

Vakgroep Informatietechnologie

Voorzitter: Prof. Dr. Ir. Daniël De Zutter

Realtime signaal synchronisatie met acoustic fingerprinting

door

Ward Van Assche

Promotoren: Dr. Marleen Denert, Joren Six

Scriptiebegeleider: Prof. Helga Naessens

Masterproef ingediend tot het behalen van de academische graad van
Master of Science in de industriële wetenschappen: informatica

Academiejaar 2015–2016

Voorwoord

Todo: voorwoord schrijven

Ward Van Assche, juni 2016

Toelating tot bruikleen

“De auteur geeft de toelating deze scriptie voor consultatie beschikbaar te stellen en delen van de scriptie te kopiëren voor persoonlijk gebruik.

Elk ander gebruik valt onder de beperkingen van het auteursrecht, in het bijzonder met betrekking tot de verplichting de bron uitdrukkelijk te vermelden bij het aanhalen van resultaten uit deze scriptie.”

Ward Van Assche, juni 2016

Realtime signaal synchronisatie met acoustic fingerprinting

door

Ward Van Assche

Masterproef ingediend tot het behalen van de academische graad van
Master of Science in de industriële wetenschappen: informatica

Academiejaar 2015–2016

Promotoren: Dr. Marleen Denert, Joren Six

Scriptiebegeleider: Prof. Helga Naessens

Faculteit Ingenieurswetenschappen en Architectuur

Universiteit Gent

Vakgroep Informatietechnologie

Voorzitter: Prof. Dr. Ir. Daniël De Zutter

Samenvatting

Todo: samenvatting schrijven

Trefwoorden

synchronisatie, realtime, datastromen, musicologie, acoustic fingerprinting

Realtime signal synchronization with acoustic fingerprinting

Ward Van Assche

Supervisor(s): Joren Six, Marleen Denert

Abstract—Sed nec tortor in libero rutrum pellentesque et gravida turpis. Phasellus gravida neque vitae elit fringilla, a efficitur purus sollicitudin. Proin lacus est, suscipit sed nibh ac, hendrerit eleifend leo. Suspendisse quis semper leo. Duis non elit commodo, sodales ex non, venenatis diam. Sed libero tortor, hendrerit et sollicitudin ut, facilisis vitae odio. Fusce vitae mi odio.

Keywords—kernwoord1, kernwoord2, kernwoord 3, kernwoord 4

I. INTRODUCTION

SED Sed nec tortor in libero rutrum pellentesque et gravida turpis. Phasellus gravida neque vitae elit fringilla, a efficitur purus sollicitudin. Proin lacus est, suscipit sed nibh ac, hendrerit eleifend leo. Suspendisse quis semper leo. Duis non elit commodo, sodales ex non, venenatis diam. Sed libero tortor, hendrerit et sollicitudin ut, facilisis vitae odio. Fusce vitae mi odio. Cras vitae quam bibendum, elementum velit ut, varius enim. Donec sagittis elit ligula, laoreet viverra felis rhoncus nec. Donec mattis metus pretium, pulvinar enim a, luctus nunc. Pellentesque quis suscipit leo.

II. SECTIE

A. Subsectie

Sed nec tortor in libero rutrum pellentesque[1] et gravida turpis. Phasellus gravida neque vitae elit fringilla, a efficitur purus sollicitudin. Proin lacus est, suscipit sed nibh ac, hendrerit eleifend leo. Suspendisse quis semper leo. Duis non elit commodo, sodales ex non, venenatis diam. Sed libero tortor, hendrerit et sollicitudin ut, facilisis vitae odio. Fusce vitae mi odio. Cras vitae quam bibendum, elementum velit ut, varius enim. Donec sagittis elit ligula, laoreet viverra felis rhoncus nec. Donec mattis metus pretium, pulvinar enim a, luctus nunc. Pellentesque quis suscipit leo.

B. Andere subsectie

Sed nec tortor in libero rutrum pellentesque et gravida turpis. Phasellus gravida neque vitae elit fringilla, a efficitur purus sollicitudin. Proin lacus est, suscipit sed nibh ac, hendrerit eleifend leo. Suspendisse quis semper leo. Duis non elit commodo, sodales ex non, venenatis diam. Sed libero tortor, hendrerit et sollicitudin ut, facilisis vitae odio. Fusce vitae mi odio. Cras vitae quam bibendum, elementum velit ut, varius enim. Donec sagittis elit ligula, laoreet viverra felis rhoncus nec. Donec mattis metus pretium, pulvinar enim a, luctus nunc. Pellentesque quis suscipit leo.

Sed nec tortor in libero rutrum pellentesque et gravida turpis. Phasellus gravida neque vitae elit fringilla, a efficitur purus sollicitudin. Proin lacus est, suscipit sed nibh ac, hendrerit eleifend

leo. Suspendisse quis semper leo. Duis non elit commodo, sodales ex non, venenatis diam. Sed libero tortor, hendrerit et sollicitudin ut, facilisis vitae odio. Fusce vitae mi odio. Cras vitae quam bibendum, elementum velit ut, varius enim. Donec sagittis elit ligula, laoreet viverra felis rhoncus nec. Donec mattis metus pretium, *pulvinar* enim a, luctus nunc. Pellentesque quis suscipit leo.

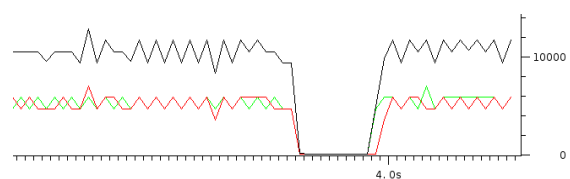


Fig. 1. Detailed capture of the stream at the moment of a handover between two simulated AP's with a strong signal. Notice the gap of 100 ms.

III. SECTIE

Aenean auctor congue nisi, volutpat porta urna lobortis in. Donec accumsan fermentum lectus, sed aliquet lectus gravida eget. Ut turpis quam, fermentum eget sem sed, euismod facilisis sem. Etiam a sollicitudin purus. Vestibulum quis nisl et nibh condimentum suscipit tristique eget quam. Aenean eget varius lectus. Maecenas sit amet mi augue. Nullam semper ex et facilisis ullamcorper. Cras volutpat ornare arcu, ac suscipit nisi pellentesque iaculis. Vivamus sit amet ipsum consequat, egestas massa varius, aliquam lorem. Aenean ornare iaculis dolor, eget efficitur elit. Maecenas ut massa ac tortor hendrerit pulvinar a in ante.

IV. CONCLUSION

The simulation results show a nice advantage for the moving cell[3] concept. The traditional handover problem can be avoided so a workable model for broadband access on trains seems realistic. But there needs to be done a lot of research to make RAU's and RoF as reliable and cheap as possible.

REFERENCES

- [1] Joren Six, Olmo Cornelis, and Marc Leman, "TarsosDSP, a Real-Time Audio Processing Framework in Java," in *Proceedings of the 53rd AES Conference (AES 53rd)*. 2014, The Audio Engineering Society.
- [2] 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.
- [3] Joren Six and Marc Leman, "Synchronizing Multimodal Recordings Using Audio-To-Audio Alignment," *Journal of Multimodal User Interfaces*, vol. 9, no. 3, pp. 223–229, 2015.
- [4] A. L. Wang, "An industrial-strength audio search algorithm," in *ISMIR 2003, 4th Symposium Conference on Music Information Retrieval*, 2003, pp. 7–13.

Inhoudsopgave

Extended abstract	4
Gebruikte afkortingen	iv
1 Introductie	1
1.1 Probleemschets	1
1.2 Evaluatiecriteria	4
1.3 Bestaande methoden	5
1.3.1 Event-gebaseerde synchronisatie	5
1.3.2 Synchronisatie met een kloksignaal	6
1.3.3 Dynamic timewarping	7
1.3.4 Accoustic fingerprinting	8
1.3.5 Kruiscovariantie	11
1.4 Doel van deze masterproef	12
2 Methode	14
2.1 Gebruikte technologieën en software	14
2.1.1 TarsosDSP	14
2.1.2 Panako	15
2.1.3 FFmpeg	16
2.1.4 SoX	16
2.1.5 Sonic Visualiser	16

2.1.6	Audacity	17
2.1.7	Max/MSP	17
2.2	Gebruikte algoritmen	18
2.2.1	Accoustic fingerprinting	18
2.2.2	Kruiscovariantie	18
2.3	Implementatie van een softwarebibliotheek	18
2.3.1	Bufferen van de audio	18
2.3.2	Synchronisatie	18
2.4	Implementatie van een Max/MSP module	18
3	Evaluatie	19
3.1	Unit testen	20
3.2	Stresstesten	20
3.3	Test in de praktijk	20
3.4	Usability testen	20
3.5	Analyse van de complexiteitsgraad	20
3.5.1	Accoustic fingerprinting	20
3.5.2	Kruiscovariantie	20
3.6	Praktische bruikbaarheid van het systeem	20
4	Conclusie	21
	Appendices	22
	A Resultaten DTW experiment	23
	Referentielijst	25
	Lijst van figuren	28
	Lijst van figuren	28

Lijst van tabellen	29
Lijst van tabellen	29
Lijst van codefragmenten	30
Lijst van codefragmenten	30

Gebruikte afkortingen

IPEM	Instituut voor Psychoakoestiek en Elektronische Muziek
DSP	Digital Signal Processing
FFT	Fast Fourier transform
ECG	Elektrocardiogram
DTW	Dynamic timewarping

Hoofdstuk 1

Introductie

1.1 Probleemschets

Het probleem dat in deze masterproef zal worden onderzocht doet zich heel specifiek voor bij verschillende experimenten die aan het IPeM worden uitgevoerd. Dit is de onderzoeksinstelling van het departement musicologie aan Universiteit Gent. De focus van het IPeM ligt vooral op onderzoek naar de interactie van muziek op fysieke aspecten van de mens zoals dansen, sporten en fysieke revalidatie.[1]

Om de relatie tussen muziek en beweging te onderzoeken worden er tal van experimenten uitgevoerd. Deze experimenten maken gebruik van allerlei sensoren om bepaalde gebeurtenissen omzetten in analyseerbare data.

Bij een klassieke experiment wordt onderzocht wat de invloed is van muziek op de lichamelijke activiteit van een persoon. Alle bewegingen worden geregistreerd met een videocamera en verschillende accelerometers.

Hierbij moeten minstens drie datastromen worden geanalyseerd: de videobeelden, de data van de accelerometer(s) en de afgespeelde audio. Een uitdaging hierbij is de synchronisatie van deze verschillende datastreams. Om een goede analyse mogelijk te maken is het

zeer gewenst dat men exact weet (tot op de milliseconde nauwkeurig) wanneer een bepaalde gebeurtenis in een datastroom zich heeft voorgedaan, zodat men deze gebeurtenis kan vergelijken met de gebeurtenissen in de andere datastromen. Door de verschillen in samplefrequentie en door de verschillende vertragingen van elke opname is dit zeker geen sinecure. [17]

Bij het IPEM maakt men gebruik van een systeem waarbij audio opnames het synchronisatieproces vereenvoudigen. Het principe werkt als volgt: men zorgt ervoor dat elke datastroom vergezeld van een perfect gesynchroniseerde audiostroom, afkomstig van een opname van het omgevingsgeluid. In het voorgaande experiment is dit eenvoudig te verwezenlijken. Bij de videobeelden kan automatisch een audiospoor mee worden opgenomen. De accelerometer kan geplaatst worden op een microcontroller vergezeld van een kleine microfoon.. Aangezien beide componenten zo dicht op de hardware geplaatst zijn is de latency tussen beide datastromen te verwaarlozen.¹ De afgespeelde audio kan gebruikt worden als referentie, aangezien dit uiteraard al een perfecte weergave is van het omgevingsgeluid.

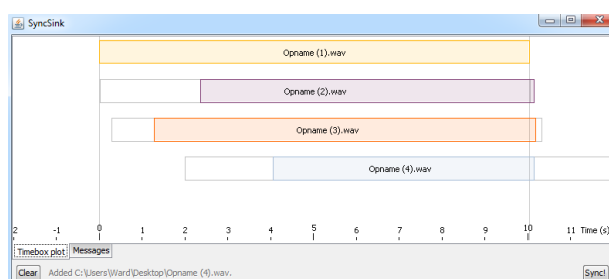
Na het uitvoeren van het experiment beschikt men dus over drie datastromen, waarbij we voor elke datastroom ook beschikken over een quasi perfect synchrone opname van het omgevingsgeluid. Aangezien het experiment in één ruimte is uitgevoerd zijn de verschillende opnames van het omgevingsgeluid zeer gelijkend. Het probleem van de synchronisatie van de verschillende datastromen kan bijgevolg gereduceerd worden tot het synchroniseren van de verschillende audiostromen.

Door de typisch eigenschappen van geluid is het helemaal niet zo moeilijk om verschillende audiostromen te synchroniseren. Bij het IPEM heeft men een systeem ontwikkeld dat in staat is om verschillende audiostreams te synchroniseren. Dit systeem maakt gebruik van *acoustic fingerprinting* en synchronisatie met *kruiscovariantie*.

Dit systeem heeft in de praktijk echter heel wat beperkingen. De grootste beperking is dat

¹De latency van de audioverwerking op een *Axoloti* microcontroller is vastgesteld op 0.333 ms. Meer informatie: <http://www.axoloti.com/more-info/latency/>

het synchronisatieproces pas kan worden uitgevoerd wanneer het experiment is afgelopen, en dit volledig handmatig. De opgenomen audiobestanden moet worden verzameld op een computer, vervolgens kan met behulp van de audiobestanden de latency van elke datastroom worden berekend. Vervolgens kunnen de datastromen worden gesynchroniseerd. Voor de musicologen die deze experimenten uitvoeren is deze werkwijze veel te omslachtig. Daarom is een eenvoudiger realtime systeem om de synchronisatie uit te voeren zeer gewenst.



Figuur 1.1: Huidige werkwijze om streams te synchroniseren: Een drag and drop interface waarin de opgenomen fragmenten gesleekt kunnen worden na afloop van het experiment. Vervolgens wordt de latency berekend.

Een ander probleem is iets vager en minder duidelijk te omschrijven. De resultaten van het kruiscovariantie algoritme bevatten soms afwijkingen die moeilijk te verklaren zijn. De precieze oorzaak hiervan, en hoe dit kan worden opgelost zal ook worden onderzocht. Ook is het kruiscovariantie algoritme in vergelijking met het acoustic fingerprinting algoritme véél gevoeliger voor storingen en ruis, veroorzaakt door slechte opnames. Aangezien de opnameapparatuur (zeker op microcontrollers) bij de uit te voeren experimenten vaak van slechte kwaliteit is, is het belangrijk om de algoritmes voldoende robuust te maken zodat ze hier niet over struikelen.

1.2 Evaluatiecriteria

Uit de probleemcontext wordt duidelijk dat het te ontwikkelen systeem moet voldoen aan heel wat vereisten. In deze sectie zullen de vereisten eenduidig geformuleerd en besproken worden zodat het mogelijk is om er in verdere analyses en testen naar te verwijzen.

Realtime synchronisatie

Een cruciale vereiste is dat de toepassing in *realtime* moet kunnen werken. Concreet wil dit zeggen dat het tijdens en van het experiment mogelijk moet zijn om alle latencies van de streams op te vragen.

Het is moeilijk om het opvragen van deze gegevens echt in realtime mogelijk te maken. Veel algoritmes vereisen een bepaalde hoeveelheid data om berekening op uit te voeren voordat ze tot een resultaat komen. De streams moeten dus gebufferd worden. Dit zorgt er in feite voor dat het systeem niet meer realtime is in de enge zin van het woord. Een buffer met als maximumgrootte de hoeveelheid data verzameld in tien seconden lijkt ons aanvaardbaar.

Detecteren van gedropte data

De beperkte resources van een microcontroller kan voor problemen zorgen bij het verwerken van streams. Zo kan het gebeuren dat er gegevens van streams verloren gaan. Bij de synchronisatie van streams met gedropte data leidt dit probleem voor een plotse verhoging van de latency. Hoewel het onmogelijk is om de gedropte data te reconstrueren is het wel gewenst dat de wijziging in latency gedetecteerd wordt en dat hier mee wordt rekening gehouden bij de verdere verwerking.

Detecteren van drift

Elke stream heeft een bepaalde samplefrequentie. Het is belangrijk dat de samplefrequentie gekend is om de gegevens correct en precies te kunnen verwerken. Het kan echter voorvallen dat de samplefrequentie bij de verwerking op microcontrollers minder nauwkeurig gekend is. Een stream waarbij de samplefrequentie 1Hz afwijkt van de theoretische waarde zal na 60 seconden een latency hebben opgebouwd van 60 samples. Bij een samplefrequentie van 8000 Hz komt dit overeen met 7,5 ms². Dit probleem mag in ons systeem zeker niet worden verwaarloosd.

1.3 Bestaande methoden

Er bestaan verschillende methoden om datastreams te synchroniseren. Welke methode te verkiezen is hangt volledig af van de toepassing.

In deze sectie komen de belangrijkste methoden aan bod en zullen ze worden getoetst aan de eerder beschreven evaluatiecriteria.

1.3.1 Event-gebaseerde synchronisatie

Deze methode wordt beschreven in [6, 17] en is een eenvoudige, intuïtieve methode om synchronisatie van verschillende datastreams uit te voeren. De synchronisatie gebeurt aan de hand van markeringen die in de verschillende streams worden aangebracht. In audiostreams kan een kort en krachtig geluid een markering plaatsen. Een lichtflits kan dit realiseren in videostreams. De latency wordt bepaald door het verschil te berekenen tussen de tijdspositie van de markeringen in de streams. De synchronisatie kan vervolgens zowel manueel als softwarematig worden uitgevoerd.

²Berekening: $60/8000\text{Hz} = 0.0075\text{s} = 7.5\text{ms}$

Deze methode kent heel wat beperkingen. Zo vormt bij de synchronisatie van een groot aantal streams de schaalbaarheid een probleem. Ook wanneer er in een stream samples gedropt worden of er drift ontstaat, leidt dit tot foutieve synchronisatie. De methode kan deze twee problemen niet detecteren tot er opnieuw markeringen worden aangebracht en de streams gesynchroniseerd worden. Verder laten ook niet alle sensoren toe om markeringen aan te brengen: zo is de synchronisatie van een ECG onmogelijk met deze methode.

Het manueel aanwenden van deze methode blijkt derhalve in een realtime situatie niet mogelijk. Wanneer de synchronisatie echter door software wordt uitgevoerd is deze methode wel in realtime bruikbaar. In dat geval moet er per tijdsinterval een markering worden aangebracht om de problemen veroorzaakt door drift en gedropte samples te overbruggen.

1.3.2 Synchronisatie met een kloksignaal

Artikel [11] beschrijft een methode waarbij door een kloksignaal realtime streams van verschillende soorten toestellen worden gesynchroniseerd. Hiervoor gebruikt men standaard audio en video synchronisatieprotocollen. Elk toestel kan gebruik maken van verschillende samplefrequenties en communicatieprotocollen.

De methode maakt gebruik van een *master time code* signaal dat verstuurd wordt naar elk toestel. Dit laat het realtime analyseren van elke stream toe. Bij deze analyse kan vervolgens meteen de samplefrequentie en latency bepaald worden.

Een groot nadeel van dit systeem is dat elk toestel een kloksignaal als input moet kunnen toelaten en verwerken. In het geval van de verwerking van videobeelden kan deze methode enkel gebruikt worden met zeer dure videocamera's. Bij goedkopere camera's (zoals webcams) zal men op zoek moeten gaan naar een alternatief. [17]

1.3.3 Dynamic timewarping

Dynamic timewarping (DTW) is een techniek die gebruikt wordt voor het detecteren van gelijkenissen tussen twee tijdreeksen³. Om een optimaal resultaat te bekomen wordt gebruikt gemaakt van een padkost van de ene tijdreeks naar de andere. Hierbij kunnen de tijdreeksen worden kromgetrokken door ze niet-lineair uit te rekken of in te krimpen ten opzichte van de tijdas [15]. De minimale kost kan in kwadratische tijd berekend worden door gebruik te maken van dynamisch programmeren [9]. DTW is een veelgebruikte techniek in domeinen zoals spraakherkenning, bio-informatica, data-mining, etc [14].

Aangezien DTW het toelaat om tijdreeksen krom te trekken is het gewenst dat zowel het verleden als de toekomst van de streams voor het algoritme toegankelijk is. Een uitbreiding op dit algoritme beschreven door Dixon laat toe één tijdreeks in realtime te streamen mits de andere stream op voorhand is gekend [9]. Toch houdt deze uitbreiding geen oplossing in voor het gestelde probleem. Alle streams komen immers in realtime toe en we willen zo snel mogelijk de latency tussen de streams achterhalen.

Het bufferen van de binnenkomende streams en vervolgens het DTW algoritme uit te voeren op de buffers leek een mogelijke manier om dit probleem te omzeilen.

Of het algoritme na deze aanpassing voldoet aan onze vereisten diende een klein experiment uit te wijzen. De resultaten hiervan zijn te vinden in appendix A: Resultaten DTW experiment.

Het experiment toonde evenwel aan dat DTW niet bruikbaar is voor de realtime stream synchronisatie. De resultaten bleken niet nauwkeurig genoeg, zeker niet wanneer we ook de performantie van het algoritme in beschouwing namen.

³Een tijdreeks is een sequentie van opeenvolgende datapunten over een continu tijdsinterval, waarbij de datapunten elkaar na telkens hetzelfde interval opvolgen.

1.3.4 Accoustic fingerprinting

De techniek van accoustic fingerprinting extraheert en vergelijkt fingerprints van audiofragmenten. Een accoustic fingerprint bevat gecondenseerde informatie gebaseerd op typische eigenschappen van het audiofragment. De kracht van dit algoritme schuilt in haar snelheid en robuustheid. Het is immers uitzonderlijk bestand tegen achtergrondgeluiden en ruis. Door deze eigenschappen is het algoritme in staat om in enkele seconden een database met miljoenen fingerprints van audiofragmenten te doorzoeken. De bekendste toepassing van accoustic fingerprinting is de identificatie van liedjes op basis van een korte opname⁴.

Het is onder meer deze techniek die het IPEM gebruikt om de opgenomen audiostreams van experimenten te synchroniseren. In tegenstelling tot *Shazam* gaat men uiteraard niet op zoek naar matches in een database maar zoekt men ze rechtstreeks tussen de audiofragmenten. Het uitgangspunt is immers dat er tussen de opnames gelijkenissen moeten gevonden kunnen worden.

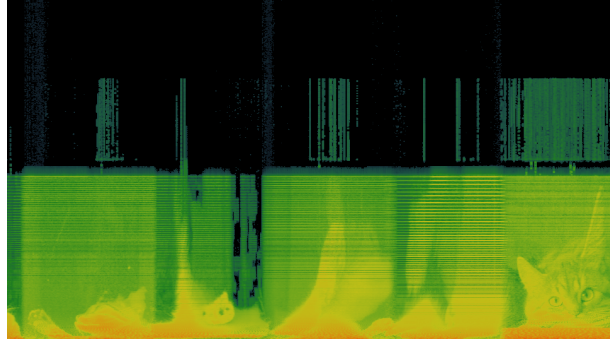
Werking

De cruciale stap bij de ontwikkeling van een accoustic fingerprinting systeem is het bepalen van de meest betrouwbare *feature* om de fingerprints op te baseren. Mogelijke features zijn frequentie, toonhoogte, tempo, ritme, dynamiek, etc. Veel features zijn echter moeilijk (softwarematig) te bepalen wat hen niet bruikbaar maakt in een robuust fingerprinting systeem. Een feature die wel geschikt is voor het bepalen van fingerprints zijn de pieken in het frequentiespectrum.

Een fingerprinter gebaseerd op de extractie van spectrale pieken gaat in verschillende stappen te werk: Eerst wordt er van elk audiofragment een spectrogram⁵ gegenereerd. Dit

⁴Het grootste voorbeeld hiervan is de smartphone app Shazam. Deze app is de eerste toepassing dat gebruik maakte van dit algoritme.

⁵Een spectrogram is een grafische voorstelling van de frequentie en intensiteit van geluid ten opzicht van de tijd [5]

Figuur 1.2: Spectrogram van Venetian Snares - Look

kan snel gebeuren met het Fast Fourier Transformation algoritme (FFT). In artikel [13] wordt deze methode uitgebreid besproken. Vervolgens worden de fingerprints bepaald door telkens twee pieken in het spectrogram te verbinden. Een tijd-frequentie punt in het spectrogram is een kandidaat-piek als het punt een hogere energetische waarde heeft dan al zijn burens [19]. Welke pieken precies met elkaar worden verbonden hangt af van verschillende parameters.

Na het bepalen van de fingerprints worden ze opgeslagen in een datastructuur waarin er snel naar matches kan worden gezocht. Van elke fingerprint worden volgende parameters bepaald:

- $f1$: de frequentie van de eerste spectrale piek van de fingerprint.
- $t1$: de tijd van de eerste spectrale piek van de fingerprint.
- Δf : het verschil van de frequenties van beide spectrale pieken van de fingerprint.
- Δt : het verschil in tijd van beide spectrale pieken van de fingerprint.

De fingerprints worden in volgende structuur bijgehouden: $(id; t1; hash(f1; \Delta f; \Delta t))$.

De hash wordt gebruikt om verschillende fingerprints te kunnen matchen. Deze bevat $f1$ en Δf omdat bij een match de beginfrequentie en het verschil in frequentie van beide fingerprints gelijk moet zijn. Enkel Δt wordt bijgehouden aangezien de begintijd van beide fingerprints waarschijnlijk niet zal overeenkomen. Het verschil in tijd tussen de fingerprints moet uiteraard wel overeenkomen.

Na het extraheren en opslaan van de fingerprints kunnen er matches gezocht worden door de hashwaarden van de fingerprints van beide audiofragmenten met elkaar te vergelijken. Van elke match wordt de offset berekend, wat het verschil is tussen t_1 van beide fingerprints. Wanneer beide fragmenten overeenkomen zal dit resulteren in een groot aantal matches met dezelfde offset.

Accoustic fingerprinting kunnen we toepassen op streams door ze te bufferen. Na het opbouwen van de buffers dient het algoritme hierop te worden uitgevoerd. De offset die gebruikt werd bij het matchen stelt de latency voor tussen beide streams.

Een uitgebreidere beschrijving is te vinden in het artikel van Wang. Deze methode is echter beperkt tot het vergelijken van audiofragmenten die in tijd noch toonhoogte gewijzigd zijn. Aan het IPEM is een aangepaste methode ontwikkeld die dit wel toelaat [16].

Toepassing in realtime

Accoustic fingerprinting lijkt in tegenstelling tot de eerder besproken methodes wel bruikbaar in een realtime toepassing. Afhankelijk van de gebruikte parameters en de kwaliteit van de audio is een buffergrootte van enkele seconden voldoende om de latency tussen de streams te bepalen. De toepasbaarheid van dit algoritme zal verder worden onderzocht.

De nauwkeurigheid van dit algoritme hangt af van de grootte van de *FFT bins*. De waarde hiervan is afhankelijk van de parameters van het FFT algoritme. Een nauwkeurigheid van 16 ms of 32 ms is standaard.

Gedropte samples en drift

De snelheid waarmee gedropte samples of drift gedetecteerd kunnen worden is afhankelijk van de implementatie van het algoritme. Een mogelijke strategie verwerkt enkel de offset met de meeste matches. In dat geval duurt het maximaal één buffertijd⁶ voordat

⁶In dit onderzoek hebben we besloten om de maximale buffertijd te beperken: zie Evaluatiecriteria

het algoritme zich aanpast aan drift of gedropte samples. Indien er bij het zoeken van matches meerdere offsets als resultaat worden toegelaten is een snellere detectie mogelijk. De snelheid hangt in dat geval af van de parameters van het algoritme.

1.3.5 Kruiscovariantie

Kruiscovariantie (ook wel kruiscorrelatie genoemd) is een berekening uit de signaalverwerking waarbij de gelijkenis tussen twee signalen wordt bepaald. De latency tussen twee audiofragmenten kan bepaald worden door deze berekening uit te voeren voor elke mogelijke verschuiving. De verschuiving waarbij de kruiscovariantie het hoogst is bepaalt de latency.

Stel twee audioblokken a en b bestaande uit s samples en verschuiving i . Voor elke i gaande van 0 tot s wordt de kruiscovariantie berekend met volgende formule:

$$\sum_{j=0}^s a_j * b_{(i+j) \bmod s} \quad (1.1)$$

De waarde van i waarbij de kruiscovariantie het hoogst is stelt de latency voor tussen beide audioblokken in aantal samples. De latency in seconden bepaalt men door dit resultaat te delen door de samplefrequentie.

De methode kan de latency tot op één sample nauwkeurig bepalen. De maximaal bereikbare nauwkeurigheid hangt dus af van de samplefrequentie van de audioblokken. Bij een samplefrequentie van $8000Hz$ is dit $1/8000Hz = 0.125ms$. Dit is zeker voldoende voor onze toepassing.

Een nadeel aan deze methode is de performantie. Het berekenen van de beste kruiscovariantie van twee audioblokken bestaande uit s samples kan gebeuren in $O(s^2)$. Het is dus belangrijk om bij deze berekening de grootte van de audioblokken te beperken.

In artikel 17 wordt deze techniek meer in detail besproken.

1.4 Doel van deze masterproef

Dit onderzoek wil drie zaken bereiken:

Selectie en optimalisatie van algoritmes

Dit houdt in: bepalen welke algoritmes we zullen gebruiken om het probleem op te lossen en het zoeken naar optimalisaties en de juiste parameters voor maximale efficiëntie. Als ons probleem dit toelaat zal er gebruik gemaakt worden van algoritmes waarvan al een IPER implementatie beschikbaar is. Deze moeten dan wel nog geoptimaliseerd worden voor onze toepassing.

Het beoogde doel is dat de algoritmes in staat zijn om audio opgenomen met een basic microfoon op een microcontroller te synchroniseren met een nauwkeurigheid van minstens één milliseconde.

Ontwerp en implementatie van een softwarebibliotheek

Het tweede doel van het onderzoek betreft het schrijven van een softwarebibliotheek. Deze bibliotheek zal gebruik maken van de geoptimaliseerde algoritmes om de audiostromen te synchroniseren. Deze bibliotheek moet vanuit andere software kunnen worden opgeroepen en gedetailleerde informatie teruggeven over de synchronisatie van de verschillende audiostreams.

Ontwerp en implementatie van een gebruiksvriendelijke interface

Uiteindelijk is het de bedoeling dat dit onderzoek resulteert in een gebruiksvriendelijke applicatie die toegankelijk is voor onderzoekers/musicologen zonder uitgebreide informatica kennis. De software moet in staat zijn om van verschillende binnenkomende datastromen

(vergezeld met audiostream) te synchroniseren en op één of andere manier weg te schrijven naar een persistent medium.

Hoofdstuk 2

Methode

2.1 Gebruikte technologieën en software

2.1.1 TarsosDSP

TarsosDSP is een Java bibliotheek voor realtime audio analyse en verwerking ontwikkeld aan het IPeM. De bibliotheek bevat een groot aantal algoritmes voor audioverwerking en kan nog verder worden uitgebreid. Deze bibliotheek wordt beschreven in artikel [18].

Een processing pipeline wordt voorgesteld als instantie van de klasse `AudioDispatcher`. Het aanmaken gebeurt met behulp van de klasse `AudioDispatcherFactory`. Deze bevat statische methodes om een `AudioDispatcher` aan te maken van een audiobestand, een `float` array of een microfoon. Aan de pipeline kunnen verschillende `AudioProcessors` worden toegevoegd. Een `AudioProcessor` is een interface met de methodes `process` en `processingFinished`. De `process` methode heeft als enige parameter een `AudioEvent`. Dit object bevat een audio blok, voorgesteld als `float` array met waarden tussen -1.0 en 1.0. De grootte van dit blokje audio, en de mate van overlapping tussen de opeenvolgende blokjes audio is instelbaar. Verder bevat dit object nog andere metadata zoals onder meer een *timestamp*.

Afhankelijk van de implementatie van de `process` methode kan de audiostroom op een bepaalde manier verwerkt, geanalyseerd of gewijzigd worden.

2.1.2 Panako

Panako is net zoals TarsosDSP een Java bibliotheek, door de zelfde auteurs ontwikkeld aan het IPEM. Panako bevat buiten implementaties van algoritmen ook enkele applicaties die hiervan gebruik maken. Deze bibliotheek wordt beschreven in artikel [16]. Panako heeft TarsosDSP als enige *dependency*.

Panako bevat een open-source implementatie van het accoustic fingerprinting algoritme beschreven in de paper van Avery Li-Chun Wang[19]. Dit algoritme is verder uitgebreid zodat audio waarbij de toonhoogte verhoogd of verlaagd is, of audio die sneller of trager is afgespeeld toch gedetecteerd kan worden.

De bibliotheek bevat verschillende applicaties die gebruik maken van dit algoritme. Zo is het mogelijk om de fingerprints van een geluidsfragment te bekijken, matches tussen verschillende geluidsfragmenten te visualiseren, en grafisch te experimenteren met de verschillende parameters.

Er is ook een applicatie beschikbaar om verschillende geluidsfragmenten te synchroniseren. Deze applicatie maakt behalve van het accoustic fingerprinting algoritme ook nog gebruik van het kruiscovariantie algoritme.

Wanneer de latency tussen de verschillende audiofragmenten bepaald is, dan kan de applicatie een shell script genereren dat met behulp van *FFmpeg* stukjes van de geluidsbestanden wegnipt of er stilte aan toevoegt. Het resultaat is dat na het uitvoeren van het script de geluidsbestanden gesynchroniseerd zijn.

2.1.3 FFmpeg

FFmpeg is een command-line multimedia framework dat gebruikt wordt voor encoderen, decoderen, multiplexen, demultiplexen, streamen en afspelen van audio en video. [12]

In dit onderzoek wordt FFmpeg voornamelijk gebruikt in scripts bij het geautomatiseerd genereren van testdata.

2.1.4 SoX

SoX is net zoals FFmpeg een command-line tool voor audioverwerking. Buiten de mogelijkheid om audiobestanden te converteren laat SoX ook minder triviale operaties toe. Zo is het onder meer mogelijk om het volume aan te passen, effecten toe te voegen, de bestanden bij te knippen of gegenereerde geluiden in een audiobestand te mixen. [7]

In dit onderzoek wordt SoX ook gebruikt in scripts bij het manipuleren van de testdata.

2.1.5 Sonic Visualiser

Sonic Visualiser is een gebruiksvriendelijke desktopapplicatie voor de analyse, visualisatie van audiobestanden. Sonic Visualiser laat toe om audiobestanden vanuit verschillende perspectieven te analyseren, zo kan zowel de waveform als het spectrogram van een audiobestand gevisualiseerd worden. Sonic Visualiser is uitbreidbaar met plug-ins in het Vamp formaat. [8]

Sonic visualiser is in dit onderzoek gebruikt om handmatig de latency tussen verschillende audiofragmenten te bepalen. De applicatie is ook gebruikt geweest om de principes achter de algoritmes te visualiseren.

2.1.6 Audacity

Audacity is een open-source desktopapplicatie voor het bewerken, opnemen en converteren van audio. Ook is het met Audacity mogelijk om verscheidene effecten aan audio toe te voegen.[3]

Alle opnames en handmatige bewerkingen op audiobestanden in dit onderzoek zijn uitgevoerd met Audacity.

2.1.7 Max/MSP

Max/MSP is een visuele programmeertaal voor muziek en multimedia. Het is een modulair systeem waarbij modules met elkaar verbonden kunnen worden om zo complexe systemen op te bouwen. Max/MSP beschikt ook over een API waarmee nieuwe modules mee ontwikkeld kunnen worden. [4]

Max/MSP kan realtime audio verwerken, daarom zullen we deze toepassing gebruiken voor het ontwikkelen van onze gebruikersinterface.

2.2 Gebruikte algoritmen

2.2.1 Accoustic fingerprinting

Werking

Optimalisaties

Parameters en hun invloed op het algoritme

2.2.2 Kruiscovariantie

Werking

Optimalisaties

Parameters en hun invloed op het algoritme

2.3 Implementatie van een softwarebibliotheek

2.3.1 Bufferen van de audio

2.3.2 Synchronisatie

2.4 Implementatie van een Max/MSP module

Hoofdstuk 3

Evaluatie

3.1 Unit testen

3.2 Stresstesten

3.3 Test in de praktijk

3.4 Usability testen

3.5 Analyse van de complexiteitsgraad

3.5.1 Accoustic fingerprinting

Tijdscomplexiteit

Ruimtecomplexiteit

3.5.2 Kruiscovariantie

Tijdscomplexiteit

Ruimtecomplexiteit

Hoofdstuk 4

Conclusie

Appendices

Bijlage A

Resultaten DTW experiment

In dit experiment proberen we de nauwkeurigheid van het DTW algoritme te bepalen wanneer streams gebufferd worden. Hiertoe bepaalden we eerst de latency tussen twee audiofragmenten. Vervolgens verkleinden we iteratief de duur van het fragment met 10 seconden waarop we het algoritme opnieuw uitvoerden. Tenslotte vergeleken we de buffergrootte en nauwkeurigheid van de resultaten.

We hebben gebruik gemaakt van twee audiofragmenten waarbij het ene fragment 2.390 seconden vertraging heeft ten opzichte van het andere fragment. Beide fragmenten hebben samplefrequentie van 8000 Hz. Eén van de twee fragmenten is een opname van het origineel en bijgevolg van matige kwaliteit.

Het experiment is uitgevoerd in *Sonic Visualiser* met behulp van de *Match Performance Aligner* plug-in. Deze plug-in laat synchronisatie toe met behulp van het DTW algoritme. De implementatie wordt uitgebreider besproken in artikel [10]. Voor dit experiment hebben we de default instellingen gebruikt. De plug-in bepaalt elke twintig milliseconden de latency tussen beide fragmenten.

De volgende tabel geeft de resultaten van het experiment weer. De eerste kolom bevat de lengte van de vergeleken fragmenten in seconden. Deze lengte stelt de buffergrootte voor

van een audiostream. De tweede kolom geeft aan hoeveel seconden van de stream moet worden verwerkt tot er een stabiel resultaat wordt bekomen. De derde kolom geeft het gemiddelde weer van de gevonden latencies. Deze waarde wordt berekend vanaf dat het algoritme een stabiel resultaat heeft gevonden. De vierde kolom bevat de standaardafwijking van dit resultaat.

Lengte	Tijd tot stabiel	Gemiddelde latency	Standaardafwijking
60s	2.540s	2,393s	0.048s
50s	2.540s	2,390s	0.095s
40s	2.540s	2,394s	0.020s
30s	2.540s	2,384s	0.145s
20s	2.540s	2,390s	0.108s
10s	2.540s	2,395s	0.025s

Uit bovenstaande resultaten kunnen we verschillende zaken concluderen. Ten eerste zien we aan de standaardafwijking dat de individuele resultaten (die iedere 20ms gegenereerd worden) niet nauwkeurig genoeg zijn om te gebruiken in onze toepassing. De gemiddelde waarde komt wel in de buurt van de werkelijke latency maar is nog steeds niet zo nauwkeurig. Ook moeten we bij de berekening van het gemiddelde rekening houden met het feit dat het algoritme pas na een bepaalde tijd een stabiel resultaat vindt, in dit geval 2.540s.

We hebben dit algoritme ook uitgetest op een fragment waaruit 500 ms hebben weggeknipt om het probleem met gedropte samples te simuleren. Het algoritme reageerde hier zeer snel op: de nieuwe latency werd na 240 ms gevonden. Het probleem is dat we zojuist hebben getracht de nauwkeurigheid te verbeteren door het gemiddelde te nemen van de resultaten. Dit heeft als gevolg dat wanneer er samples gedropt zijn het eindresultaat zich bevindt tussen de initiële en nieuwe latency.

Referentielijst

- [1] Ipem - systematic musicology. <https://www.ugent.be/lw/kunstwetenschappen/en/research-groups/musicology/ipem>. [Online; geraadpleegd 05-maart-2016].
- [2] Cycling ‘74 max. <https://cycling74.com/>. [Online; geraadpleegd 05-maart-2016].
- [3] Audacity. <http://audacity.sourceforge.net/>, 2015. [Online; accessed 12-March-2016].
- [4] CYCLING ‘74 MAX. <https://cycling74.com/>, 2016. [Online; accessed 12-March-2016].
- [5] Dictionary.com unabridged. Mar 2016. URL <http://www.dictionary.com/browse/sound-spectrogram>.
- [6] David Bannach, Oliver Amft, and Paul Lukowicz. Automatic event-based synchronization of multimodal data streams from wearable and ambient sensors. In *Smart sensing and context*, pages 135–148. Springer, 2009.
- [7] Benjamin Barras. Sox: Sound exchange. Technical report, 2012.
- [8] Chris Cannam, Christian Landone, and Mark Sandler. Sonic visualiser: An open source application for viewing, analysing, and annotating music audio files. In *Proceedings of the 18th ACM international conference on Multimedia*, pages 1467–1468. ACM, 2010.

-
- [9] Simon Dixon. Live tracking of musical performances using on-line time warping. In *Proceedings of the 8th International Conference on Digital Audio Effects*, pages 92–97. Citeseer, 2005.
 - [10] Simon Dixon and Gerhard Widmer. Match: A music alignment tool chest. In *ISMIR*, pages 492–497, 2005.
 - [11] Javier Jaimovich and Benjamin Knapp. Synchronization of multimodal recordings for musical performance research. In *NIME*, pages 372–374, 2010.
 - [12] Roman Kollár. Configuration of ffmpeg for high stability during encoding.
 - [13] Alan V Oppenheim. Speech spectrograms using the fast fourier transform. *IEEE spectrum*, 8(7):57–62, 1970.
 - [14] Chotirat Ann Ratanamahatana and Eamonn Keogh. Everything you know about dynamic time warping is wrong. In *Third Workshop on Mining Temporal and Sequential Data*. Citeseer, 2004.
 - [15] Stan Salvador and Philip Chan. Toward accurate dynamic time warping in linear time and space. *Intelligent Data Analysis*, 11(5):561–580, 2007.
 - [16] 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.
 - [17] Joren Six and Marc Leman. Synchronizing Multimodal Recordings Using Audio-To-Audio Alignment. *Journal of Multimodal User Interfaces*, 9(3):223–229, 2015. ISSN 1783-7677. doi: 10.1007/s12193-015-0196-1.
 - [18] Joren Six, Olmo Cornelis, and Marc Leman. TarsosDSP, a Real-Time Audio Processing Framework in Java. In *Proceedings of the 53rd AES Conference (AES 53rd)*. The Audio Engineering Society, 2014.

-
- [19] Avery Li-Chun Wang. An industrial-strength audio search algorithm. In *ISMIR 2003, 4th Symposium Conference on Music Information Retrieval*, pages 7–13, 2003.

Lijst van figuren

1.1	Huidige werkwijze om streams te synchroniseren: Een drag and drop interface waarin de opgenomen fragmenten gesleept kunnen worden na afloop van het experiment. Vervolgens wordt de latency berekend.	3
1.2	Spectrogram van Venetian Snares - Look	9

Lijst van tabellen

Lijst van codefragmenten