



**MATEMATICKO-FYZIKÁLNÍ
FAKULTA**
Univerzita Karlova

BAKALÁŘSKÁ PRÁCE

Jiří Balhar

Extrakce melodie pomocí hlubokého učení

Ústav formální a aplikované lingvistiky

Vedoucí bakalářské práce: Mgr. Jan Hajič **Ph.D.??**

Studijní program: Informatika

Studijní obor: Programování a softwarové systémy

Praha 2019

Prohlašuji, že jsem tuto bakalářskou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Poděkování.

.

Název práce: Extrakce melodie pomocí hlubokého učení

Autor: Jiří Balhar

Ústav: Ústav formální a aplikované lingvistiky

Vedoucí bakalářské práce: Mgr. Jan Hajič **Ph.D.??**, ústav

Abstrakt: Extrakce melodie patří mezi nejdůležitější a nejtěžší úlohy oboru Music Information Retrieval, jejíž šíře uplatnění zahrnuje zlepšení metod pro vyhledávání v hudebních datech (například pomocí zabroukání části skladby), automatický přepis improvizovaných sólových pasáží nebo zpracování zvukových nahrávek pro muzikologické studie. Právě melodie je totiž tím hlavním, co si člověk po poslechu skladby odnáší a zpětně vybavuje a z podstaty se tedy často jedná o její nejvýraznější rys. Přítomnost hudebního doprovodu, který melodii podbarvuje a dává možnost jí vyniknout, však pro dosavadní rigidní algoritmické metody znemožňuje její průběh spolehlivě zachytit. V posledních letech se proto obor posouvá směrem k využívání metod hlubokého učení, díky kterému je možné napodobit intuitivní schopnost člověka rozeznat melodii v nahrávce. Na tyto práce navazujeme a představujeme nové metody, zejména pak architekturu *Harmonic Convolutional Neural Network*, založenou na úpravě vnitřního uspořádání obvyklých konvolučních sítí, pro snažší zachycení harmonické povahy jednotlivých tónů, která překonává poslední state-of-the-art metody na většině veřejně dostupných datasetech.

Title: Melody Extraction using Deep Learning

Klíčová slova: Music Information Retrieval Extrakce Melodie Odhad F0 Hluboké učení Harmonická konvoluční neuronová síť

Author: Jiří Balhar

Institute: Institute of Formal and Applied Linguistics

Supervisor: Mgr. Jan Hajič **Ph.D.??**, institute

Abstract:

koukám, že šablona mi tu pasivně agresivně naznačuje, že abstrakt mám moc dlouhý, zítra to ještě seškrtnám.

Keywords: Music Information Retrieval Melody Extraction F0 estimation Deep Learning Harmonic Convolutional Neural Network

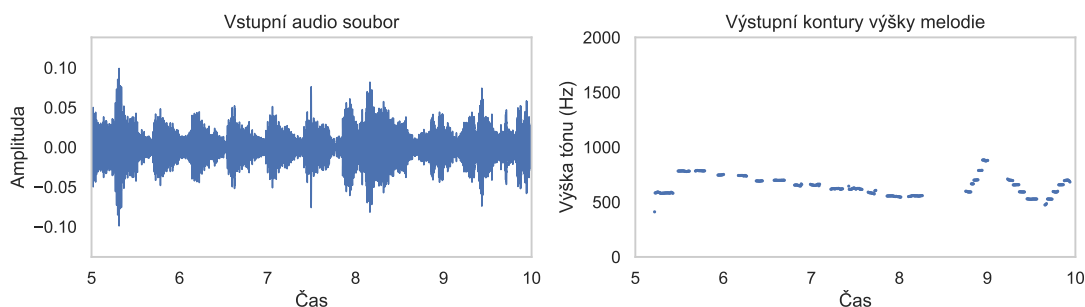
Obsah

1 Úvod	3
Úvod	3
1.1 Analýza hudebního signálu	4
1.2 Definice melodie	7
1.3 Metody extrakce melodie	9
1.4 Hluboké učení	11
1.5 Přínosy práce	12
1.6 Struktura práce	12
2 Související práce	13
2.1 Průzkum existujících metod	13
2.1.1 Spektrální analýza	14
2.1.2 Funkce salience	18
2.1.3 Hledání melodie	21
2.2 Srovnání existujících metod	22
2.2.1 Replikace výsledků	23
3 Datasets	25
3.1 Struktura dostupných dat a jejich přehled	26
3.2 MedleyDB	27
3.3 Orchset	28
3.4 MIREX datasety	29
3.5 Weimar Jazz Database	29
3.6 MDB-synth	29
3.7 Dataset RWC	30
4 Metody evaluace	31
4.1 MIREX	31
4.1.1 Formát výstupu v soutěži MIREX	31
4.2 Trénovací, validační a testovací množina	31
4.3 Metriky	33
4.3.1 Definice metrik	33
4.3.2 Další metriky	36
5 Experimenty	38
5.1 Architektura CREPE	38
5.1.1 Replikace výsledků CREPE	40
5.1.2 CREPE pro extrakci melodie	41
5.1.3 Vliv rozlišení diskretizace výšky noty	42
5.1.4 Vliv rozptylu cílové pravděpodobnostní distribuce výšky noty	43
5.1.5 Vliv šířky vstupního okna	44
5.1.6 Vliv násobného rozlišení první konvoluční vrstvy	45
5.2 Architektura WaveNet	46
5.2.1 Baseline na základě Martak a kol. (2018)	47
5.2.2 Vliv počtu filtrů dilatačních vrstev a skip propojení	49

5.2.3	Systematické prohledávání počtu dilatačních vrstev a bloků	49
5.2.4	Vliv velikosti šířky kernelu dilatací	51
5.2.5	Vliv výstupní transformace skip propojení	52
5.2.6	Vliv velikosti první konvoluce	52
5.3	Architektura HCNN	54
5.3.1	Harmonické transformace	57
5.3.2	Parametr <code>hop_size</code>	58
5.3.3	Vícekanálový vstup CQT	58
5.3.4	Kontext	59
5.3.5	Augmentace vstupních spektrogramů	62
6	Výsledky	64
6.1	Výběr testovaných modelů	64
6.1.1	Architektura CREPE	64
6.1.2	Architektura WaveNet	64
6.1.3	Architektura HCNN	64
6.1.4	Detekce melodie	66
6.2	Kvantitativní srovnání	66
6.2.1	Popis výsledků	66
6.3	Kvalitativní srovnání	68
6.4	Interpretace výsledků	74
	Závěr	75
	Seznam použité literatury	77
	Seznam obrázků	84
	Seznam tabulek	87
	Seznam použitých zkratk	89
	Přílohy	90
6.5	Architektura HCNN, Vliv úpravy architektury ovlivňující recep- tivní pole modelu	90
6.6	Výsledky detekce melodie testovaných metod	90

1. Úvod

Spolu s harmonií a rytmem představuje melodie základní stavební kámen většiny existující hudby. V průběhu vývoje od folklórních zpěvů přes orchestrální skladby po soudobou elektroniku si melodie téměř vždy zachovávala své dominantní postavení nositele esence jednotlivých písní. Melodie je to hlavní, co si člověk po poslechu skladby odnáší a nejsnadněji vybaví, a její důležitost je zejména v našem kulturním kontextu natolik jednoznačná, že se hudba, která ji postrádá, zřídka dostává do širokého povědomí.



Obrázek 1.1: Příklad vstupu a výstupu metody pro extrakci melodie. [přidat zvukový příklad do přílohy](#)

Tato práce se zabývá metodami odhadu fundamentální frekvence melodie ze zvukové nahrávky. Jinými slovy je naším cílem získat v každém bodě vstupní skladby informaci o tom, zda melodie v daném okamžiku zní a její případnou výšku. Jde o jednu z nejdůležitějších a zároveň nejtěžších úloh z oboru *Music Information Retrieval* (MIR), jejíž rozsah využití v této doméně pokrývá významnou část aktivně řešených otevřených problémů.

Spolehlivý přepis melodie by usnadnil vyhledávání v hudebních datech, ať už na základě notového zápisu (*Symbolic Melodic Similarity*), pomocí nekvalitní nahrávky z rádia (*Audio Fingerprinting*), pomocí broukání (*Query by Singing/Humming*) nebo dokonce pomocí coveru hledané písně (*Audio Cover Song Identification*). Mimo vyhledávání by byl algoritmus užitečný pro další zpracování zvukového signálu, ať už pro manipulaci a úpravu melodického hlasu (jako se pokouší například software Melodyne), nebo naopak jeho odstranění a vytvoření karaoke doprovodu (*Informed Source Separation*). V neposlední řadě by extrakce melodie pomohla při kategorizaci hudebních dat, například podle žánru (*Genre Classification*) nebo podle zpěváka (*Singer Characterization*). A konečně široké spektrum využití by našla i v muzikologii (případně etnomuzikologii) pro kvantitativní i kvalitativní studii hudebních motivů a postupů (v jazzu například Pfeiderer a kol.).

Extrakce melodie však nemusí sloužit pouze jako mezikrok pro řešení jiné úlohy. Užitečný je i samotný výstup algoritmu, znázorněný na obrázku 1.1. Motivacním příkladem použití může být pomoc při transkripci: Představíme-li si začínajícího hráče na saxofon, který si chce do not přepsat své oblíbené jazzové sólo, výstup algoritmu mu dá užitečnou informaci o tom, jaký tón zní v jakou chvíli. Z této reprezentace už pak hráči zbývá nalezené tóny projít a zapsat je do notové osnovy.

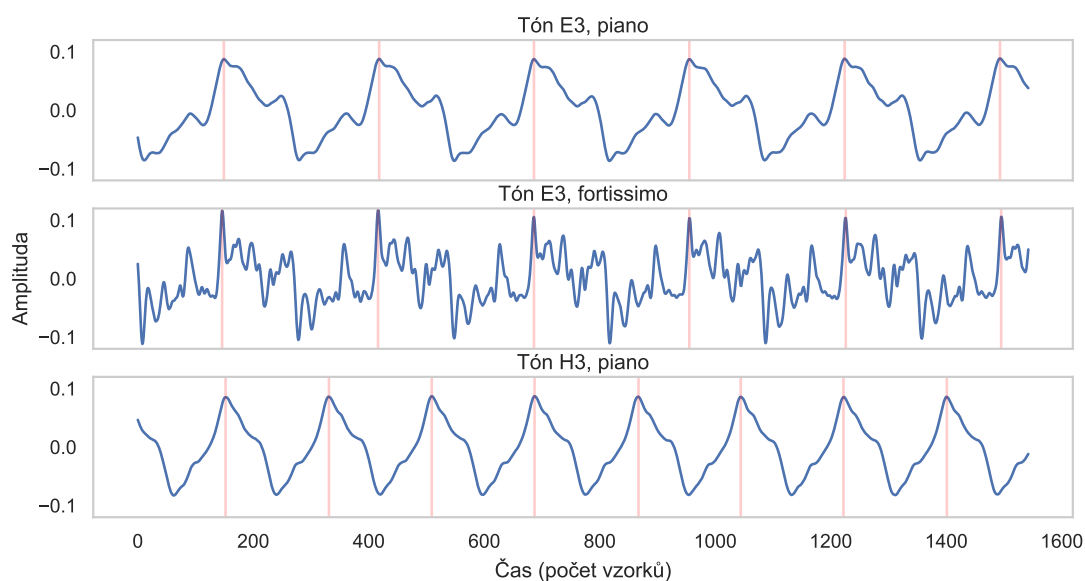
1.1 Analýza hudebního signálu

Proč je ale extrakce melodie otevřený problém? Příbuzná úloha, která spočívá v přepisu nahrávky jednoho izolovaného nástroje, je v podstatě vyřešena (Mauch a Dixon, 2014); proč se tato úloha po přidání hudebního doprovodu stává výrazně obtížnější? Pro vysvětlení zásadního problému, který se s přepisem nahrávky pojí, musíme nejdříve přiblížit vůbec povahu zvuku a možnosti jeho zkoumání.

Naše zkušenost se zvukem probíhá primárně skrze sluch. Teprve na hlasitém koncertu však člověk pocítí, že zvuk je ve své fyzikální podstatě změna tlaku vzduchu, putující od zdroje k posluchači. Díky sluchu z těchto vibrací dokážeme oddělit jednotlivé zdroje a identifikovat v nich i velmi jemné rozdíly. Ačkoli jde o subjektivní vjemy, zvuky lze částečně rozřadit podle toho, jak snadno v nich rozeznáme nějakou konkrétní výšku.

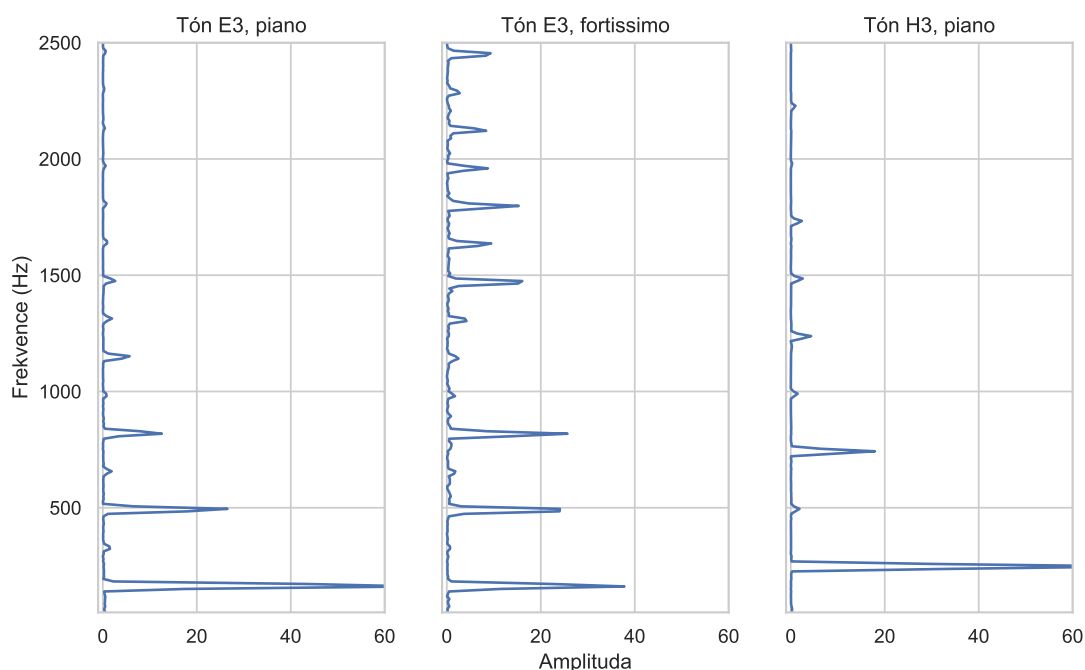
Čtenář této práce si nyní může postupně vybavit: hrající violoncello, odbíjení kostelního zvonu, cinknutí příboru, štěkot psa, plynutí potoka, šelest listí stromů, trhání papíru, tlesknutí a prasknutí balónku.

Se ztrácející se zřetelností výšky nejprve přijdeme o možnost zpívat společně se zdrojem zvuku v harmonii a posléze i o možnost si představit „vyšší“ a „nižší“ instance toho samého zvuku (jak zní vysoké a nízké prasknutí balónku?). To, co mají první z uvedených příkladů společné, je výrazná a stabilní periodicitu jejich signálu — daný tlakový průběh se opakuje v čase. Díky sluchu tuto periodicitu interpretujeme jako výšku, přičemž různé výšky se od sebe liší frekvencí, se kterou se signál opakuje. Hudební nástroje jsou jedním ze zdrojů těchto pravidelných vibrací, jejichž frekvenci lze zpravidla měnit (pomocí klapky, pohybu prstu po struně, atd.). Hlas nástroje však není charakteristický pouze svou výškou, nýbrž i barvou. Ta je určena podobou signálu v rámci jedné periody.



Obrázek 1.2: Zvuk klarinetu, tóny s různou výškou a dynamikou, 35 milisekund signálu se vzorkovací frekvencí 44 100 Hz. **nějak uvést zdroj zvuku** https://www.philharmonia.co.uk/explore/sound_samples/clarinet?p=3

Na obrázku 1.2 můžeme srovnat tři tóny hrané klarinetem, první dva mají stejnou výšku, jsou ale zahráné s různou dynamikou. Jejich vizuální rozdíl částečně odpovídá i rozdílu v barvě tónu, první tón má příjemný, měkký zvuk, druhý je výraznější a hrubší. Třetí tón se od zbylých liší svou výškou, což lze pozorovat na kratší periodě signálu, která je na obrázku 1.2 vyznačená svislými úsečkami.



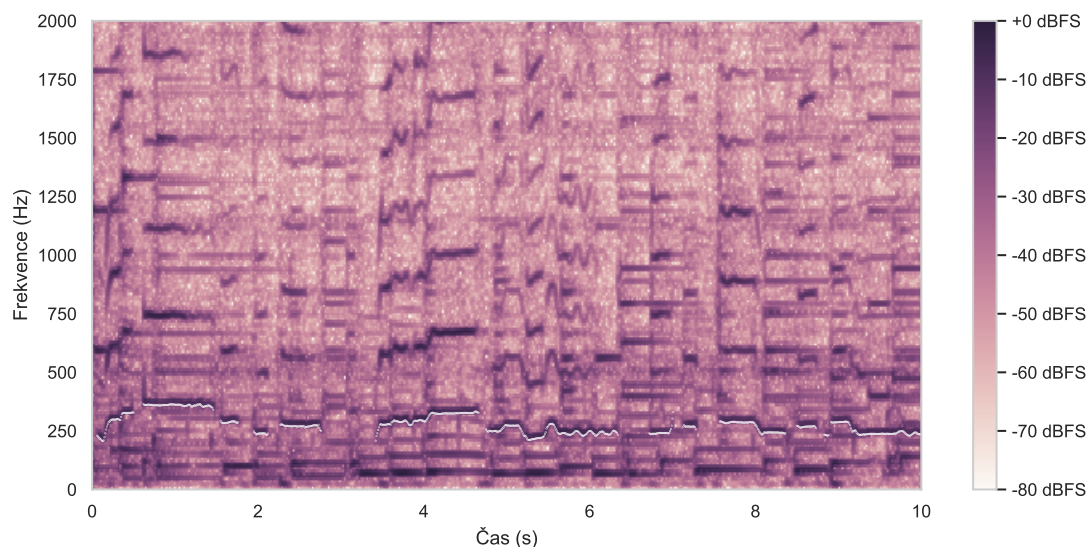
Obrázek 1.3: Zvuk klarinetu, absolutní hodnota výstupu Fourierovy transformace signálu délky 4096 s oknem typu Hamming.

Jedním ze způsobů analýzy zvukového signálu je pomocí Fourierovy transformace (DFT). Základní myšlenkou je, že na signál lze hledět jako na vážený součet jednodušších signálů. Podobně, jako když se barvy na obrazovce míchají ze tří základních, libovolný zvuk můžeme smíchat ze sady sinusoid. Výslednou kombinaci všech frekvencí, ze kterých se zvuk skládá, označujeme *zvukové spektrum*. Na obrázku 1.3 vidíme část výsledku Fourierovy transformace zvuků klarinetu z předchozího příkladu. To zásadní, co na spektru tónu můžeme pozorovat, je jeho podstata jakožto součet *harmonických složek*. Tón, kterému posluchač přisoudí výšku f_0 , se zpravidla skládá ze součtu sinusoid, jejichž frekvence je celočíselným násobkem základní frekvence f_0 (jinak také nazývaná *fundamentální frekvence*). Například tedy tón E3 se na obrázku 1.3 skládá ze složek o frekvenci 165 Hz, 330 Hz, 495 Hz, ..., zároveň intenzita těchto harmonických frekvencí určuje barvu hlasu.

Ukazuje se, že práce s touto reprezentací zvuku je pro analýzu signálu užitečnější, než práce s nezpracovaným signálem. Ze spektrální reprezentace je například na první pohled zřejmý vztah fundamentálních frekvencí porovnávaných signálů, který odpovídá lidské intuici o výšce zvuků — tón H3 je na obrázku 1.3 opravdu „výše“ než tón E3. Díky spektrální analýze lze také pozorovat charakteristiky hlasů různých nástrojů. Pro hlas klarinetu platí, že liché harmonické frekvence jsou mnohem výraznější než sudé (na obrázku 1.3 vypadají sudé harmonické jako malé vrcholky mezi výraznými lichými), naopak například lidský zpěv je charak-

teristický výraznějšími sudými harmonickými složkami. Další výhodou je možnost hledání rozdílů v barvě tónů — na spektru vidíme, že vyšší harmonické jsou u tónu hraném fortissimo mnohem výraznější než u tónu hraném piano. Tyto vyšší frekvence způsobují zmiňovanou hrubost tónu.

Harmonická struktura, která je vlastní lidskému hlasu a téměř všem zvukům hudebních nástrojů, je zásadní pro metody extrakce melodie. Je to vlastnost, která zvuky potenciálně nesoucí melodii odlišuje od bubnového doprovodu, od šumu nebo od jiných nemelodických rušení. Díky ní se také můžeme pokoušet rozložit souzvuk různých vysokých tónů na jejich původní, čisté signály.



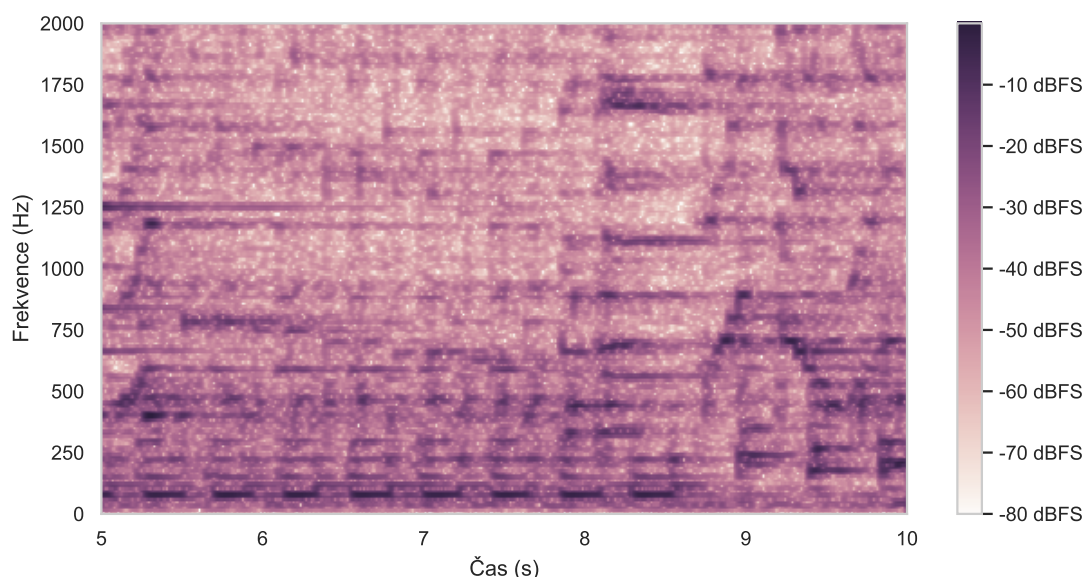
Obrázek 1.4: Spektrogram zpěvu s doprovodem piana, basy a perkusí; zpívaná melodie je vyznačena bílým obrysem.

Obrázek 1.4 vznikl pomocí opakované Fourierovy transformace, která byla aplikována na po sobě jdoucí, krátké časové úseky vstupní nahrávky, přičemž intenzity frekvenčních složek v každém časovém okamžiku zvuku jsou nyní znázorněny odstínem barvy. Časově-frekvenční reprezentaci signálu nazýváme obecně *spektrogram*, a jeho výpočet je prvním krokem většiny metod pro extrakci melodie.

Na spektrogramu na obrázku 1.4 můžeme pozorovat harmonické struktury tónů — vyznačenou konturu hlasu, která se na frekvenční ose pohybuje volněji, a pak klavírní a basový doprovod, charakteristický frekvenční stabilitou a v čase slábnoucí amplitudou. Na tomto jednoduchém příkladu lze melodii zpozorovat poměrně snadno, nese ji velmi výrazný, v poměru k doprovodu nejsilnější hlas. Lze na něm však prezentovat první ze základních problémů extrakce melodie.

Frekvence tónů, ze kterých se skládá hudební skladba, jsou uspořádány do stupnic, které definují pevně dané poměry (hudební *interval*), ve kterých se tyto tóny ve skladbě mohou vyskytovat. Principem libozvučnosti jsou však takové intervaly, které způsobují, že harmonické frekvence jednotlivých tónů se překrývají a ve výsledné směsi pak není zřejmé, zda-li daná harmonická frekvence patří k jednomu, či více hlasů. Hudební doprovod, pro lidské ucho znějící „pod melodií“, tedy často svými harmonickými frekvencemi zasahuje do melodie samotné, což je zřejmé ze spektrální analýzy.

Dekompozice signálu na jednotlivé znějící hlasy, která je pro člověka přirozená podobně, jako porozumění řeči v rušné kavárně, se kvůli této harmonické povaze tónů a intervalů stává pro algoritmy extrakce melodie obtížným problémem. To, co pro nás činí hudbu zajímavou pro poslech, ji činí obtížně analyzovatelnou pro počítač.



Obrázek 1.5: Spektrogram orchestrální skladby s obtížně detekovatelnou melodií. Čtenář se může pokusit ve spektrogramu nalézt melodii, na obrázku 1.6 nalezneme řešení.

1.2 Definice melodie

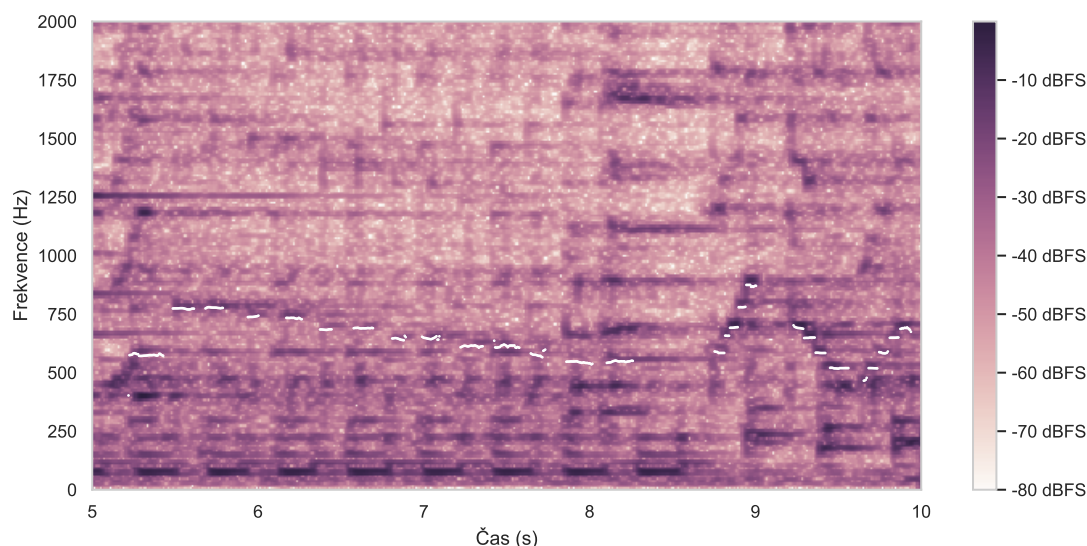
Rozpoznání melodie v rámci hudební skladby je pro většinu posluchačů intuitivní schopností, která je součástí prožitku poslechu hudby, a která jejímu poslechu vůbec dává smysl. Ačkoli je melodie tedy termín, který je subjektivně jasný, formální, obecně přijímanou muzikologickou definicí, která by se zpětně neodkazovala k posluchači, nemá.

Z tohoto důvodu si výzkumné týmy zabývající se automatickou transkripcí melodie volí pragmaticky spíše užší definice melodie, se kterými se v jejich kontextu pracuje lépe. Práce Goto a Hayamizu (1999), která je považována za jednu z prvních prací v oboru, chápe melodii jako „konturu fundamentální frekvence sestávající se z nejsilnějších tónů hrajících v omezeném frekvenčním rozsahu“. Práce se tedy omezuje na poměrně úzké chápání melodie, obecně se totiž tóny melodie jistě mohou vyskytovat i mimo autory specifikovaný frekvenční rozsah a také nemusí být vždy v poměru s doprovodem nejhlasitější složkou signálu. Z technického hlediska však umožnila autorům implementaci algoritmu běžícího v reálném čase, který poskytoval sémanticky bohatý popis vstupních nahrávek. Navazující články již pracují s volnějšími definicemi, které lépe reflektují podstatu melodie.

Kompromisem mezi subjektivní a praktickou definicí se na dlouhou dobu stala „extrakce základní frekvence hlavního, neměnného, melodického hlasu“. Ačkoli

melodii v reálném hudebním materiálu obvykle nese více hlasů, které se v hraní střídají (například píseň se zpěvem a kytarovým sólem), v letech 2005 – 2015 se v soutěži MIREX (mezinárodní soutěž pro metody řešící MIR úlohy) provádí evaluace pouze nad krátkými výňatky, kde tato definice není omezující. Přestože se může na první pohled zdát, že tato definice pouze uměle zjednodušuje celou úlohu, její formulace vede k rozvoji nových a zajímavých přístupů, které se sice formálně soustředí právě na extrakci pouze jednoho neměnného hlasu, ale ve výsledku překvapivě dobře fungují i na složitější skladby. Příkladem nového směru může být extrakce melodie pomocí modelování hudebního záznamu jako součtu signálu jednoho hlasu a doprovodu (práce Durrieu a David (2010) nebo Bosch a Gómez (2016)) s přesahem do příbuzné úlohy oddělení hlasů (source separation). Některé práce se zaměřují ještě konkrétněji na separaci lidského zpěvu a doprovodu (Hsu a Jang (2010), Ikemiya a kol. (2016)). Nově se také objevují práce, které „hlavní“ melodický hlas neinterpretují nutně jako „nejsilnější“ a k jeho rozlišení využívají dalších rysů, jako je barva, vibrato nebo délka not. Například Salamon a Gomez (2012) využívají těchto rysů pro finální výběr mezi extrahovanými kandidátními konturami.

Ve svém přehledovém článku Salamon a kol. (2014) dochází k závěru že výzkum začal v letech 2009–2012 stagnovat, nová data jsou proto pro další vývoj oboru zásadní. Výrazným posunem v rámci MIR komunity proto bylo zveřejnění nových datasetů MedleyDB (Bittner a kol., 2014) a ORCHSET (Bosch a kol., 2016), oba obsahují data, ve kterých již melodii nenese pouze jeden hlas po celou dobu skladby. V porovnání s do té doby dostupnými daty jde také o mnohem rozmanitější kolekce a v případě MedleyDB jde o první volně dostupný dataset, ve kterém se objevují celé skladby, nikoli pouze výňatky.



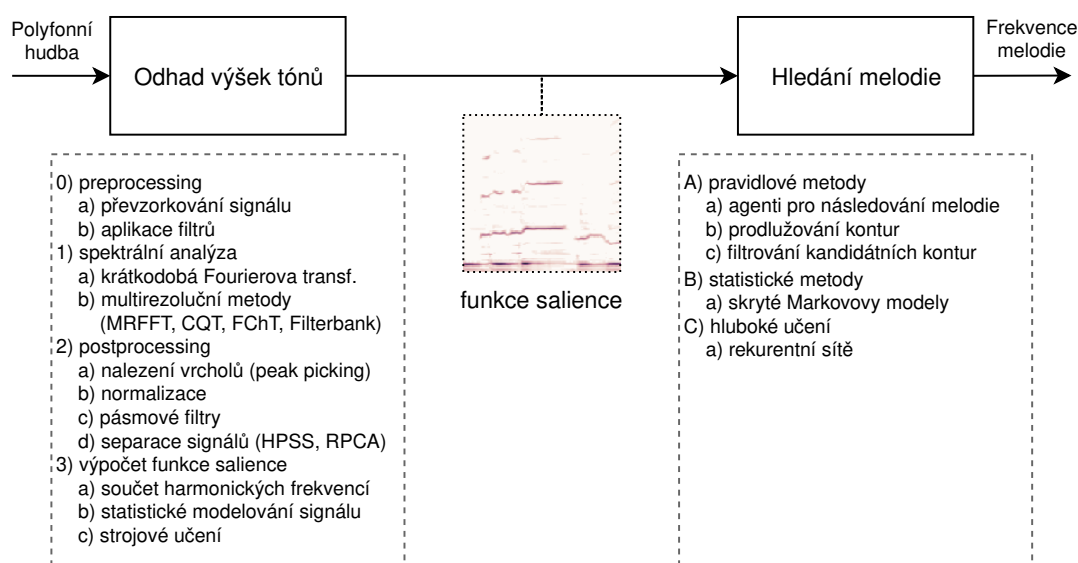
Obrázek 1.6: Spektrogram orchestrální skladby s vyznačenou melodií.

Bosch a kol. (2016) pro práci na datasetu ORCHSET definuje melodii jako „jednohlasou sekvenci tónů, kterou bude posluchač nejspíše reprodukovat, pokud jej požádáme o zapískání či zabroukání příslušné skladby“ (na základě článku Poliner a kol. (2007)), pro sestavení kolekce dat proto opravdu využívá skupiny posluchačů, které po poslechu krátkých ukázek orchestrální hudby následně žádá

o přezpívání melodie. V případě MedleyDB se na anotacích melodie podílí skupina profesionálních hudebníků a vznikají tři různé přepisy melodie s různě volnými formulacemi definice melodie.

Celkový směr výzkumu je tak ve výsledku velmi podmíněn dostupnými daty. Ta tvoří jakýsi protipól k ryze technickým a objektivním cílům algoritmických metod. Nadějí je, že tato dialektika vývoje algoritmů a práce na datech vyústí jednak v metody extrakce, které věrně zachycují podstatu subjektivního prožitku porozumění hudbě, a jednak v celkovém důsledku také snad v lepší porozumění pojmu melodie obecně.

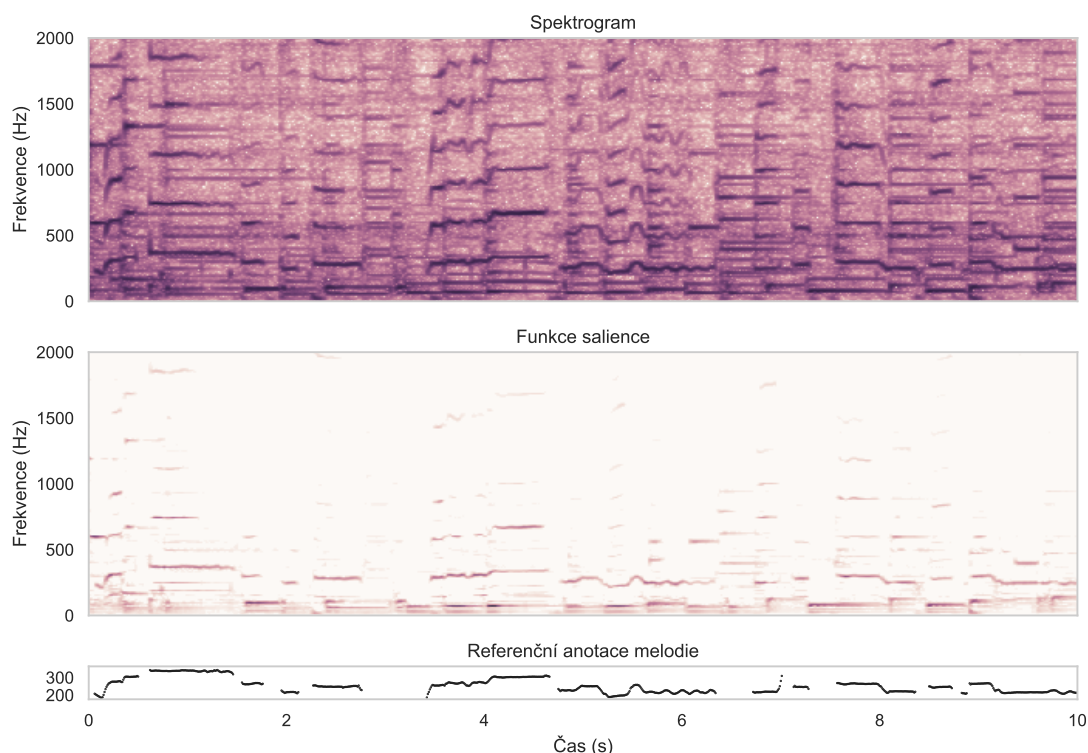
1.3 Metody extrakce melodie



Obrázek 1.7: Diagram obvyklého návrhu metod pro extrakci melodie.

Základním a společným přístupem k problému extrakce melodie je dekompozice na podproblémy odhadu výšek znějících hlasů v signálu a následného výběru melodické linie z této reprezentace znějících výšek. Jednotlivé metody se pak liší ve způsobech řešení těchto podproblémů, mimo jiné také v míře abstrakce od původního signálu, které při zpracovávání dosahují — zatímco některé přístupy po celou dobu pracují pouze se signálem jako takovým a důmyslnými způsoby ho transformují tak, aby získaly co nejpřesnější odhad výšky melodie, jiné metody při výpočtu vytváří symbolický popis jednotlivých not a následně i celých frází a melodii pak hledají v tomto vysokoúrovňovém popisu nahrávky. Přehled části používaných přístupů, blíže popsanych v kapitole Související práce, můžeme vidět na diagramu 1.7.

Prvním krokem všech existujících metod pro extrakci melodie je převod zvukového signálu do frekvenční domény, ať už pomocí zmiňované STFT nebo použitím jiných metod vyvinutých pro analýzu harmonického signálu (Multi-Resolution Fast Fourier Transform, Constant-Q Transform, Fan Chirp Transform a další, blíže popisované v sekci 2.1.1). Jednou z hlavních výhod těchto dalších metod je logaritmická osa frekvence, díky které je snadné pracovat s harmonickými poměry a hudebními intervaly nezávislé na výšce tónu. Abychom tuto vlastnost



Obrázek 1.8: Příklad výstupu výpočtu salienční funkce pomocí váženého sčítání harmonických frekvencí. Ačkoli je zpěv velmi zvýrazněn a salienční funkce na většině nahrávky dobře zachycuje výšku znějící melodie, doprovod kolem čtvrté sekundy nahrávky má vyšší hodnotu než zpěv, což neodpovídá lidskému vnímání zpěvu jakožto nejdůležitější složky signálu.

ilustrovali, uvedeme příklad — tóny A4 a E5 jsou vzdálené o kvintu (na klavíru od tónu A4 musíme postupně zmáčknout 7 bílých a černých klapek, abychom se dostali k tónu E5), tóny A5 a E6 jsou také vzdálené o kvintu. Rozdíl frekvencí těchto tónů je však 220 Hz mezi první dvojicí a 440 Hz mezi druhou dvojicí, tudíž vzdálenost frekvencí daného intervalu závisí na tónu, od kterého se tento interval počítá. Protože je hudební interval jistý *poměr* mezi dvěma tóny, použitím logaritmické osy frekvence se tyto poměry budou jevit jako absolutní rozdíly ($\log n f_0 = \log n + \log f_0$).

Zpracováním spektrogramu vstupu pak vzniká tzv. *funkce salience*, která každé znějící frekvenci v signálu přiřazuje jisté ohodnocení, které vyjadřuje poměrnou důležitost dané frekvence k ostatním slyšitelným složkám. Funkce salience je tedy v jistém slova smyslu speciální frekvenčně-časová reprezentace, která podává zejména informace o znějící melodii, nikoli o celkové kompozici signálu. Na obrázku 1.8 můžeme srovnat funkci salience vypočtenou pomocí váženého součtu harmonických frekvencí se vstupním spektrogramem.

Druhým krokem je pak výběr melodie na základě funkce salience. Triviálním řešením je výběr takových frekvencí, kterým funkce salience přiřazuje nejvyšší ohodnocení. Problémem tohoto přístupu je však to, že jakmile signál obsahuje více podobně ohodnocených tónů, výstup tohoto řešení má tendenci mezi těmito kandidáty často „přeskakovat“. Algoritmy pro extrakci melodie proto volí různé pokročilé metody vyhlazování, případně metody hledání nejpravděpodobnějšího

průchodu posloupností stavů (například pomocí Viterbiho algoritmu).

1.4 Hluboké učení

Motivací pro použití metod strojového učení je překonání limitů člověkem navržených, rigidních, pravidlových systémů. Cílem je automatické nalezení optimálního postupu pro řešení úlohy na základě množství dat, ve kterých strojové učení dokáže nalézt a využít jejich pravidelností tak, aby zvolený model dokázal tyto vlastnosti zobecnit a nalezený postup na trénovací množině poté přinášel co nejlepší výsledky na datových instancích, které během trénování k dispozici nebyly. V našem případě pak po metodě založené na strojovém učení požadujeme, aby na základě příkladů z trénovací množiny vytvořila funkci salience.

Výhodou tohoto přístupu je, že o vstupních datech nemusíme dělat žádné předpoklady. Vzniklá metoda pak může v praxi zohledňovat více faktorů ovlivňujících přítomnost melodie, jako je její barva, frekvenční modulace (vibrato, glissando) nebo jiné vzájemné srovnání současně znějících tónů. Na základě trénovacích příkladů může být tato metoda robustnější vůči většímu spektru barev hlasů nástrojů — zatímco některé metody pro extrakci melodie často uvažují signály s postupně se snižujícím podílem harmonických frekvencí, opravdové signály hudebních nástrojů často tento předpoklad nesplňují (viz obrázek 1.3).

První pokus o využití těchto metod představili Poliner a Ellis (2005), vstupní signál transformovali pomocí krátkodobé Fourierovy transformace a část spektra po jednoduché normalizaci použili jako vstupní data pro metodu podpůrných vektorů (SVM). Jejich metoda měla své limity, výstup byl kvantizován na úroveň jednoho půltónu a tudíž metoda nedokázala dobře postihnout například vibrata. I přesto však tým dosáhl srovnatelných výsledků s ostatními metodami.

Po roce 2005 jakékoli pokusy o aplikaci strojového učení ustávají a na nové metody se čeká až do roku 2016. Jedním z důvodů byl jistě nedostatek dat, tuto situaci zlepšil například dataset MedleyDB (Bittner a kol., 2014) nebo dnes již zaniklý iKala (Chan a kol., 2015). Zájem o strojové učení však znovu stoupá, také díky úspěšnému využití hlubokých neuronových sítí napříč ostatními obory. Na konferenci International Society of Music Information Retrieval (ISMIR) 2016¹ objevují dva články týmů Kum a kol. (2016) a Rigaud a Radenen (2016), založené právě na hlubokém učení. V roce 2017 publikuje své metody Bittner a kol. (2017) (ISMIR), Balke a kol. (2017) (ICASSP²), následující rok přináší metody D Basaran, S Essid (2018) (ISMIR), Bittner a kol. (2018). Všechny zmíněné popisujeme v kapitole Související práce. V oboru lze tedy od roku 2016 vidět velmi výrazný trend právě směrem k hlubokému učení, a stejný směr je patrný i v příbuzných úlohách přepisu hudby. Tým z laboratoře Google Brain dokázal výrazně zlepšit přepis klavírních skladeb pomocí kombinace konvoluční a rekurentní architektury (Hawthorne a kol., 2017). Neuronové sítě také zlepšují výsledky na poli oddělení signálů (Stoller a kol., 2018).

V této práci se pokusíme navázat na zmiňované práce a otestovat nové architektury hlubokých neuronových sítí pro úlohu extrakce melodie, zejména pak pro hledání nových způsobů výpočtu funkce salience. Předpokládáme, že čtenář

¹International Society for Music Information Retrieval Conference

²International Conference on Acoustics, Speech, and Signal Processing

práce má příslušné znalosti o metodách hlubokého učení, zejména pak o konvolučních neuronových sítích. Pro potřebný základ se odkazujeme na knihy Goodfellow a kol. (2016) a Bishop (2006)³

1.5 Přínosy práce

V práci navrhujeme tři nové metody pro výpočet funkce salience založené na hlubokém učení. Pro návrh architektur se inspirujeme pracemi Kim a kol. (2018), van den Oord a kol. (2016a) a Bittner a kol. (2017), které se zabývají využitím hlubokého učení pro zpracování zvukového signálu. Architektury představované v článcích upravujeme pro účel extrakce melodie a systematickým prohledáváním nastavení hyperparametrů a úprav topologie sítí tyto modely dále optimalizujeme. V případě architektury vycházející z práce Bittner a kol. (2017) navrhujeme nový způsob úpravy vnitřních vrstev, díky kterému je síť schopna lépe zachytit harmonické závislosti tónů přítomných v signálu. Úpravu nazýváme *harmonická transformace* a tuto novou architekturu pak *Harmonic Convolutional Neural Network* (HCNN).

Architektury srovnáváme s vybranými state-of-the-art metodami pro extrakci melodie a v případě poslední navrhované architektury HCNN prokazujeme její účinnost překonáním dosavadních výsledků na většině veřejně dostupných datasetů k úloze. Následně provádíme kvalitativní analýzu porovnáním výstupů HCNN se stávajícími metodami a navrhujeme východiska pro navazující práce.

1.6 Struktura práce

V následující kapitole Související práce shrnujeme dosavadní přístupy existujících metod extrakce melodie a vybíráme tři, které v této práci používáme pro srovnání výsledků. V kapitole Datasets uvádíme kompletní seznam všech veřejně dostupných dat k úloze, porozumění jejich obsahu nám dovolí interpretovat výsledky experimentů. Kapitola Metody evaluace obsahuje informace o způsobech porovnávání výsledků jednotlivých metod z teoretického i praktického hlediska. Popis nových přístupů k úloze, spolu s experimenty, které ke konkrétní podobě těchto přístupů vedly, popisujeme v kapitole Experimenty. A konečně srovnání těchto nových metod s existujícími nalezneme v kapitole Výsledky.

³Případně na skvělé přednášky Deep Learning vedené M. Strakou. Video záznam dostupný na: <https://ufal.mff.cuni.cz/courses/npf1114>

2. Související práce

Pokusy o vytvoření automatické metody pro kompletní transkripci hudby se podle Poliner a kol. (2007) objevují již od sedmdesátých let, z důvodu značné obtížnosti této úlohy, která strmě roste s každým dalším přidaným hlasem ve zkoumaném signálu, však dodnes jedná o otevřený problém. Z tohoto důvodu se od devadesátých let objevují práce, které se pokouší alespoň o částečný automatický popis některých muzikálních aspektů skladeb.

Jednou z prvních je práce týmu Goto a Hayamizu (1999), který se záměrně omezuje na identifikaci jedné, nejhlasitější, spojitě křivky fundamentní frekvence hlasu (F0) v omezeném frekvenčním rozsahu. Vzniklé transkripce pak sice nejsou kompletní, na druhou stranu je jejich získání výpočetně nenáročné a přitom poskytují sémanticky bohatý popis nahrávek, který je poměrně často shodný s melodií. Ustanovením úlohy pojmenované jako „Predominant-F0 Estimation“ (Pre-FEst) byly položeny základy pro vznik navazujících prací a soutěží zabývajících se automatickým přepisem melodie.

Největší rozkvět v oboru začal od roku 2004. Uspořádáním první soutěže pro porovnání systémů pro automatický popis hudby v rámci konference ISMIR (ISMIR 2004 Audio Description Contest), se ustanovily priority, formalizovaly podmínky evaluace a byly sestaveny první kolekce dat pro testování algoritmů (Downie a kol. (2010)). Soutěž se v následujícím roku přerodila do samostatné každoroční události, v níž soutěží stále více týmů v rostoucím počtu úloh.

V této kapitole představíme existující metody a společné přístupy k řešení úlohy extrakce melodie. Nejprve úlohu dekomponujeme na podúlohy a pro tyto podúlohy uvedeme příklady existujících metod. V závěru kapitoly pak provádíme kvantitativní srovnání existujících metod na základě dat ze soutěže MIREX, abychom vybrali nejnadhlejší metody, které replikovat a na něž v této práci navázat.

2.1 Průzkum existujících metod

Jen do soutěže MIREX se od roku 2005 přihlásilo 45 týmů s 62 různými metodami pro extrakci melodie, s různou mírou přesnosti přepisu. Mezi přístupy k tomuto problému tedy existuje velká rozmanitost, jejíž kompletní popis přesahuje rámec této práce. Zaměříme se proto na společné rysy a celkové trendy v oboru.

Shrnující práce od Poliner a kol. (2007) a Salamon a kol. (2014) se při charakterizaci systémů pro transkripci odkazují na příbuznou úlohu odhadu fundamentální frekvence monofonní nahrávky. Algoritmy pro monofonní tracking na základě vstupního signálu $x_{mono}(t)$ počítají *funkci salience* $S_{x_{mono}}(f_{\tau}, \tau)$ pro každý krátký časový okamžik (okno) τ a frekvenci f_{τ} . Výsledkem této funkce je relativní ohodnocení (příp. pravděpodobnost) jednotlivých frekvencí obsažených ve vstupním signálu, které značí, zda-li je daná frekvence fundamentální frekvencí znějícího hlasu.

Pro zvýšení odolnosti vůči šumu, přeslechu, dozvuku a jiným vlivům, které zhoršují kvalitu odhadu salience, se využívá také spojitosti fundamentální frekvence. Pro zajištění kontinuity extrahovaných frekvenčních kontur se zohledňuje také faktor temporálních závislostí $C(\mathbf{f})$, jejímž vstupem je kandidátní kontura

\mathbf{f} a výstupem je ohodnocení této celé kontury na její spojitost. Například tato funkce může penalizovat odhady, ve kterých výstupní F0 často přeskakuje o oktavu, což je u skutečného signálu nepravděpodobné a naopak jde o častou chybu při výpočtu salienční funkce signálů se silnými sudými harmonickými frekvencemi.

Výstupem monofonního trackingu je posloupnost frekvencí s maximální saliencí a spojitostí, tedy posloupnost frekvencí, které jsou nejlépe ohodnocenými kandidáty na fundamentální frekvenci a zároveň tato celá sekvence má také vysoké ohodnocení spojitosti.

$$\hat{\mathbf{f}} = \operatorname{argmax}_{\mathbf{f}} \left[\sum_{\tau} S_{x_{mono}}(f_{\tau}, \tau) + C(\mathbf{f}) \right]$$

Přejdeme-li k úloze extrakce melodie, obecně se vstupní polyfonní signál $x(t) = x_m(t) + x_d(t)$ skládá ze směsi melodického hlasu $x_m(t)$ a hudebního doprovodu $x_d(t)$, cílem metod pro extrakci je z pohledu přepisu fundamentální frekvence zvýšení robustnosti algoritmu vůči mnohem výraznějšímu druhu šumu - hudebnímu doprovodu $x_d(t)$. Výstupem našeho systému tedy bude posloupnost odhadů frekvence v každém časovém okně vstupního signálu, reprezentovaná vektorem $\hat{\mathbf{f}}$:

$$\hat{\mathbf{f}} = \operatorname{argmax}_{\mathbf{f}} \left[\sum_{\tau} S'_x(f_{\tau}, \tau) + C'(\mathbf{f}) \right]$$

kde f_{τ} je frekvence na pozici τ ve vektoru \mathbf{f} . $S'_x(f_{\tau}, \tau)$ je upravená funkce salience, která při výpočtu zohledňuje vliv doprovodu a složka $C'(\mathbf{f})$ představuje ohodnocení celého průběhu melodie.

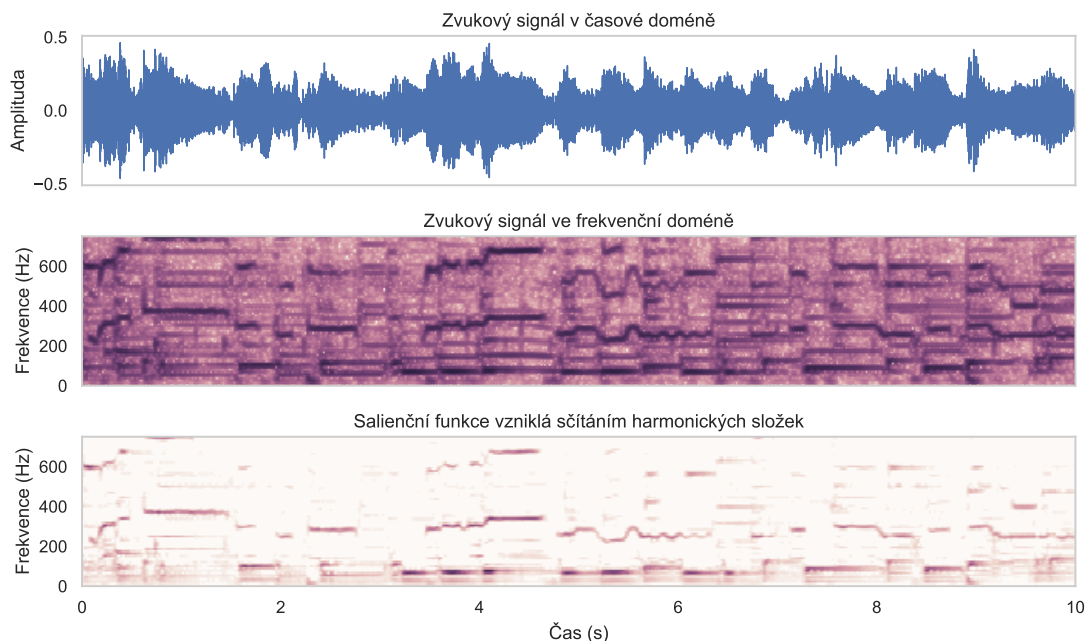
Spolu s odhadem frekvencí by také měl systém na výstupu určit úseky, ve kterých v nahrávce melodie zní a kdy nikoli. K výstupu tedy patří také vektor $\hat{\mathbf{v}}$, se stejným počtem složek jako $\hat{\mathbf{f}}$, který indikuje přítomnost melodie v každém časovém okně τ .

Většina existujících metod sdílí podobnou základní strukturu při řešení extrakce, která se zakládá na popsané formalizaci. Prvním krokem je transformace zvuku do frekvenční domény a následný odhad znějících výšek tónů v polyfonním signálu (výpočet funkce salience), druhým krokem je pak zpracování těchto odhadů a výběr melodie (tedy zpřesnění výsledné $\hat{\mathbf{f}}$ pomocí $C'(\mathbf{f})$). Přístupy k řešení těchto dvou kroků již s konkrétními příklady nastíníme v dalších sekcích.

2.1.1 Spektrální analýza

Zvuk hraného tónu na melodickém nástroji je z fyzikálního pohledu periodická změna tlaku vzduchu. Perioda tohoto signálu se nazývá fundamentální frekvence (označujeme F0) a zpravidla je tento signál složen ze součtu řady sinusoid, jejichž frekvence jsou celočíselným násobkem fundamentální frekvence. V čase mění se amplitudy těchto *harmonických frekvencí* udávají hlasitost a barvu hlasu, výška první harmonické frekvence (tj. výška fundamentální frekvence) pak ve většině případů odpovídá posluchačem vnímané výšce tónu.

Prvním krokem metod pracujících s hudebním signálem je proto provedení spektrální analýzy, jde o převod zvuku do frekvenční reprezentace, která odhaluje tyto harmonické struktury tónů a umožňuje s nimi dále pracovat.



Obrázek 2.1: Znázornění kroků spektrální analýzy a výpočtu funkce salience.

Krátkodobá Fourierova transformace

Ačkoliv je přístup k spektrální analýze více, nejpřímochařejší a podle Dressler (2016) nejčastěji používaná metoda je *krátkodobá Fourierova transformace* (STFT). Jejím principem je rozdělení vstupního signálu na množinu překrývajících se oken konstantní délky a výpočet Fourierovy transformace těchto krátkých zvukových úseků. Komplexní výsledek transformace umocníme a získáme tzv. výkonové spektrum signálu, které obsahuje informaci o poměrech energie frekvencí, ze kterých se signál v okně skládá. Spektrogram $X(f, \tau)$ vypočteme v čase τ a frekvenční složce f jako:

$$X(f, \tau) = \left| \sum_{n=-\infty}^{\infty} x(n)w(n - \tau)e^{-2\pi i \cdot n \cdot f} \right|^2$$

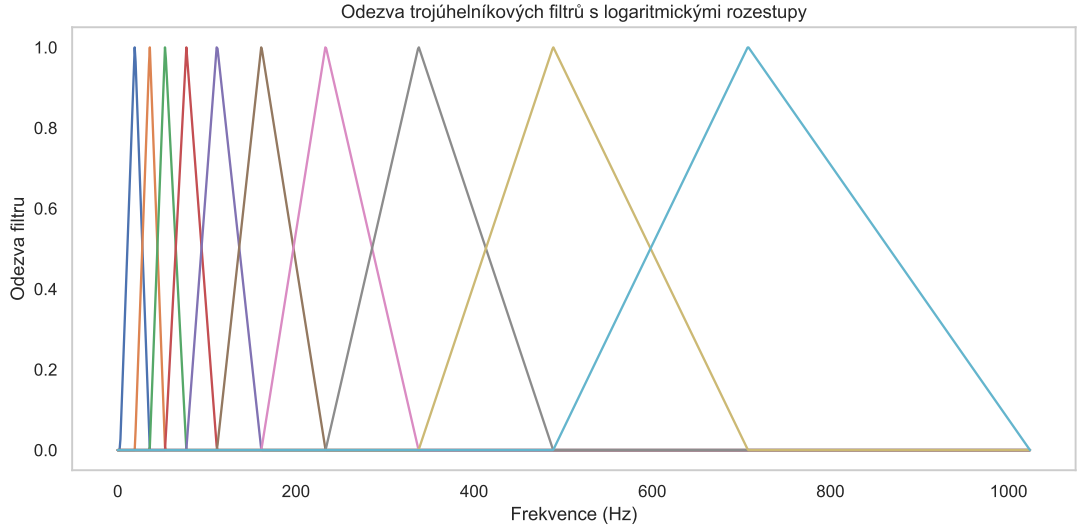
Pro diskretní vstupní signál $x(n)$ a okénkovou funkci $w(n)$. V této práci budeme pro výpočty spektrogramů používat Hannovo okno, které omezuje tzv. prořezávání ve spektru. **POPSAT HANNOVO okno. citace? proč omezuje přesřekování?**

$$w_{hann}(n) = \begin{cases} \cos^2 \frac{\pi n}{N}, & \text{pro } -\frac{N}{2} \leq n \leq \frac{N}{2}, \\ 0, & \text{jinak.} \end{cases}$$

Kde N je požadovaná velikost okna STFT transformace (kvůli obvyklé praxi výpočtu STFT pomocí algoritmu rychlé Fourierovy transformace (FFT) se hodnota N volí z $N \in \{512, 1024, 2048, \dots\}$).

Logaritmická osa spektrogramu

Výsledkem STFT je rozložení signálu na jednoduché frekvenční složky (sinusoidy) s konstantně vzdálenými frekvencemi. Jinými slovy frekvenční osa spektrogramu vytvořeném metodou STFT je lineární. Jak jsme již nastínili v úvodu,



Obrázek 2.2: Příklad trojúhelníkových filtrů pro transformaci frekvenční domény na logaritmickou škálu.

povaha hudebních intervalů a harmonických struktur tónů spočívá v tom, že téměř všechny periodické signály se v hudební skladbě vyskytují ve vzájemných poměrech (v případě intervalů v poměrech $2^{\frac{n}{12}}$ a v případě harmonických frekvencí v celočíselných). K tónu, jehož základní frekvence je rovna 440 Hz, patří také harmonické složky s frekvencemi 880 Hz, 1320 Hz, \dots , tedy absolutní vzdálenosti na spektrogramu STFT mezi frekvencemi harmonických složek jsou závislé na výšce základní frekvence. Toto způsobuje obtíže při analýze signálů, jelikož všechny uvažované frekvenční rozdíly jsou relativní. Častým druhem zpracování STFT je proto převod frekvenční osy na logaritmickou, na té pak platí pro vzdálenost libovolné fundamentální frekvence f_0 a její libovolné h -té harmonické frekvence $f_h = f_0 \cdot h$:

$$\log f_h - \log f_0 = \log (h \cdot f_0) - \log f_0 = \log h + \log f_0 - \log f_0 = \log h = \text{const.}$$

Tedy absolutní vzdálenosti uvnitř harmonické struktury tónů na spektrogramu s logaritmickou osou frekvence zůstávají konstantní nezávisle na výšce fundamentální frekvence. Tento přepočít se obvykle provádí pomocí banky filtrů s trojúhelníkovou odezvou, pro transformovaný signál $X(f, \tau)$ s N frekvenčními složkami spočteme nový spektrogram s logaritmickou osou následovně:

$$X_{\log}(\omega, \tau) = \sum_{f=-\infty}^{\infty} g_{\omega}(f) X(f, \tau)$$

Přičemž $g_{\omega}(f)$ značí odezvu trojúhelníkového filtru pro výstupní frekvenční složku ω .

Multirezoluční transformace

Na libovolnou metodu převodu diskretizovaného signálu na frekvenční doménu se vztahuje Gaborův limit, který popisuje závislost přesnosti lokalizace signálu ve frekvenční a časové doméně (Gabor a Member, 1945). Volbou délky vstupního

okna transformace zpřesňujeme buď frekvenční nebo časové rozlišení výsledné spektrální reprezentace. Zvolíme-li krátké vstupní okno, zvyšujeme časové rozlišení (krátké okno lépe zachycuje rychlé změny průběhu signálu), avšak ztrácíme přesnost na frekvenční ose, opačný vztah platí pro volbu delšího okna.

Tato limitace je markantní zejména pokud STFT používáme pro hudební data. , u vyšších tónů jsou proto vzdálenosti mezi frekvencemi signálů větší než u nižších tónů. Frekvenční rozlišení STFT je však konstantní na celém výstupním frekvenčním rozsahu a volba velikosti okna transformace zajistí dobrý poměr frekvenčního a časového rozlišení jen pro část rozsahu. Ve výsledku je pak buď pro vyšší frekvence okno příliš velké (zbytečně detailní frekvenční rozlišení na úkor časového rozlišení) a nebo naopak pro nižší frekvence je okno nedostačující (rozlišení frekvence nemusí být pro basy ani na úrovni půltónů).

Z tohoto důvodu existují vedle STFT i další metody, jejichž cílem je nabídnout lepší kompromis frekvenčně-časového rozlišení v kontextu melodických dat. Goto a Hayamizu (1999) používají MRFFT (Multi-Resolution Fast Fourier Transform), principem je opakovaný downsampling signálu (převzorkování na nižší vzorkovací frekvenci) a aplikace Fourierovy transformace na každý vzniklý signál; s každou iterací spektrum obsahuje čím dál podrobnější informace o nižších frekvencích, protože vyšší frekvence se při downsamplingu ztratí. Brown (1990) popsala metodu Constant-Q Transform (CQT), která spočívá v použití proměnné délky okna Fourierovy transformace pro výstupní frekvenční pásma, která rovnoměrně pokrývají logaritmickou osu frekvence. Cancela a kol. (2010) kombinuje CQT a Chirp Z-transform, jedná se o zobecnění diskrétní Fourierovy transformace, ve které je signál rozložen na tzv. lineární čerpy — sinusoidy s proměnnou frekvencí. Díky tomu transformace dokáže s lepším rozlišením zachytit signály, které rychle mění výšku tónu, v hudbě například vibrato. Paiva a kol. (2004) napodobují mechanismy lidského sluchu pomocí banky pásmových filtrů (Cochleagram) s logaritmicky rozmístěnými mezními frekvencemi a sumy autokorelací na jednotlivých frekvenčních pásmech signálu (Summary correlogram).

I přes uvedené důvody se Salamon a kol. (2014) a Dressler (2016) domnívají, že metoda zpracování signálu příliš neovlivňuje výslednou přesnost algoritmů pro přepis melodie. Tvrzení dokládají jednak celkovým srovnáním výsledků metod ze všech ročníků soutěže MIREX a jednak neochvějnou převahou využití krátkodobé Fourierovy transformace, jakožto efektivní a dostačující metody pro spektrální analýzu.

Postprocessing spektrogramu

Po převodu signálu na frekvenční reprezentaci následuje u většiny metod některý druh úpravy celého spektrogramu, předcházející samotnému výpočtu *funke salience*. Výsledkem tohoto kroku může být potlačení šumu a nemelodických částí signálu, zpřesnění informace o výšce znějících frekvencí nebo normalizace či jiná úprava amplitud spektrogramu.

Nejčastější úpravou je nalezení lokálních maxim (vrcholů). Potlačením ne-maximálních oblastí se zbavíme velkého množství nemelodických složek signálu, přitom informaci o těch melodických neztratíme. Výhodou práce s množinou vrcholů je, že jejich frekvenci lze na základě spektra dále zpřesnit pomocí parabolické interpolace (Rao a Rao (2010)) a nebo využitím úhlové frekvence (Salamon a Gomez (2012), Dressler (2009)).

Jiným druhem úpravy jsou různé způsoby normalizace, ať jednoduché aplikace logaritmu na jednotlivé hodnoty spektrogramu (Cancela (2008), Bittner a kol. (2017)) nebo složitější strategie normalizace, které se aplikují na celá výstupní okna krátkodobé Fourierovy transformace (spectral whitening, Ryyänen a Klapuri (2008)), či které používají pohyblivé průměry nebo jinou metodu, beroucí v potaz širší zvukový kontext. Cílem normalizace je zvýraznění slabších harmonických frekvencí a potlačení celopásmových zvuků (například perkusí). Principiálně podobným krokem je aplikace pásmového filtru (Goto a Hayamizu (1999)) pro zvýraznění frekvencí obsahující melodii. Případně využití psychoakustických filtrů modelující lidské vnímání hlasitosti (Salamon a Gomez (2012), Ikemiya a kol. (2016)). Ze signálu lze také oddělit melodické nástroje a perkusivní doprovod pomocí metod separace signálů (source separation). Používanými metodami jsou například Harmonic/Percussive Sound Separation (HPSS) (Tachibana a kol. (2010)) nebo Robust principal component analysis (RPCA) (Ikemiya a kol. (2016)).

2.1.2 Funkce salience

Salience tónu vyjadřuje míru důležitosti či nápadnosti ke svému hudebnímu okolí. Nejvíce ji ovlivňuje hlasitost v poměru ke zbylým znějícím hlasům, vliv má ale také řada dalších charakteristik hraní. Dressler (2016) mezi příklady uvádí například frekvenční modulaci, jako je vibrato nebo glissando, zejména oproti frekvenčně stálému hudebnímu doprovodu (piano, kytara). Velký vliv má samozřejmě také i barva hlasu. Lidský zpěv nebo obecně zvuky se silnějšími vyššími alikvótními frekvencemi snadněji upoutají pozornost. V případě vícehlasu mají obecně posluchači potíže rozeznat výšky tónů uvnitř souzvuku. Pokud má posluchač přiřadit pocítovanou výšku tónu akordu, obvykle volí nejvyšší či nejnižší ze znějících frekvencí.

Výstupem funkce salience je ohodnocení každé výšky tónu v každém časovém okamžiku nahrávky, které co nejlépe odpovídá v relativních poměrech výše popsané zvukové salienci. Jelikož neexistují žádné studie, které by se zabývaly měřením a kvantifikací toho, co člověk považuje za salientní v hudbě, nelze posoudit, jak dobře výsledky obvyklých způsobů výpočtu funkce salience korelují s mírou, ze které vychází. Bittner (2018) se však domnívá, že odhad bude velmi hrubý, většina postupů totiž do výpočtu zahrnuje pouze hlasitost hlasu, a tedy vynechává řadu jiných důležitých faktorů, které salienci ovlivňují.

Přístupy k výpočtu by se daly zařadit do tří kategorií - sčítání harmonických frekvencí, odhad parametrů modelujících vstup, a metody strojového učení.

Sčítání harmonických frekvencí

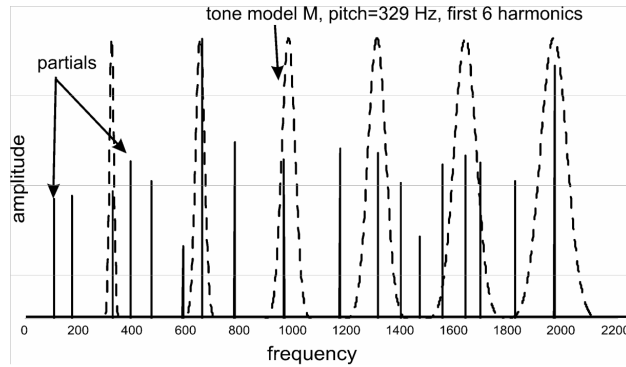
Metody založené na sčítání harmonických frekvencí jsou principiálně nejjednodušší skupinou. Vychází z práce Hermes (1988), jejich podstatou je využití harmonické struktury zvuku tónů. Ohodnocení frekvence f_τ získáme váženou sumou amplitud všech jejích harmonických frekvencí $h \cdot f_\tau$. Pro spektrogram $X(f, \tau)$ signálu x , funkci vah $g(f_\tau, h)$ a N_h počet zahrnutých harmonických frekvencí v sumě:

$$S_x(f_\tau, \tau) = \sum_{h=1}^{N_h} g(f_\tau, h) |X(h \cdot f, \tau)|$$

Pro ilustraci uvažme jednohlasý harmonický signál s fundamentální frekvencí f^* . Hodnoty spektrogramu $X(f, \tau)$ tedy budou vyšší kolem frekvencí $H_{f^*} = \{1 \cdot f^*, 2 \cdot f^*, 3 \cdot f^*, \dots\}$ a jinde nulové. Funkce $S_x(f, \tau)$ bude tedy pro $f \notin H_{f^*}$ nulová a v f^* bude nabývat globálního maxima. Příklad výstupu této metody pro polyfonní nahrávku je na obrázku 2.1.

Dressler (2011) tuto metodu vylepšuje zpracováváním dvojic vrcholů vstupního spektrogramu, její výsledný salienční spektrogram obsahuje méně kandidátů na fundamentální frekvenci. Cancela (2008) se pokouší zmenšovat chybné hodnoty salienční funkce pomocí vyhodnocení subharmonických frekvencí.

Statistické modelování signálu



Obrázek 2.3: Ilustrace modelu tónu spolu se signálem, převzato z Marolt (2004)

Jiným přístupem k počítání funkce salience, který používá Goto a Hayamizu (1999), je modelování okna spektrogramu váženým součtem harmonických struktur (modelů tónů). Snažíme se vrcholy ve spektru rozdělit mezi různě silně znějící tóny tak, aby v součtu co nejlépe odpovídaly měřeným intenzitám. Přístup se jinými slovy snaží zjistit, jaké tóny musely v danou chvíli znít, aby vzniklo dané spektrum.

Vstupem metody je okno normalizovaného spektrogramu $p_X^{(t)}(x)$ v čase t . Okno se pokusíme modelovat jako hustotu pravděpodobnosti $p(x; \theta^{(t)})$ vzniklou váženou směsí modelů všech možných tónů melodie v definovaném rozsahu frekvencí v intervalu $[F_l, F_h]$. Hustotu pravděpodobnosti jednoho z tónů s fundamentální frekvencí F označíme jako $p(x|F)$ (obrázek 2.3 ukazuje jeden z možných modelů tónu), a jako $w^{(t)}(F)$ označíme váhu, kterou model tónu $p(x|F)$ přispívá do celkové smíšené hustoty pravděpodobnosti. Pak $p(x; \theta^{(t)})$ definujeme jako:

$$p(x; \theta^{(t)}) = \int_{F_l}^{F_h} w^{(t)}(F) p(x|F) dF$$

$$\theta^{(t)} = \{ w^{(t)}(F) \mid F_l < F < F_h \}$$

Cílem pak je nalezení takových parametrů $\theta^{(t)}$, aby model $p(x; \theta^{(t)})$ dobře popisoval pozorování $p_X^{(t)}(x)$. K tomu Goto a Hayamizu (1999) využívá Expectation-Maximization (EM) algoritmus. Výsledné parametry $\theta^{(t)}$ jsou pak hodnoty salienční funkce.

Metody strojového učení

K výpočtu funkce salience můžeme využít také metody strojového učení. První práci využívající těchto metod představují Poliner a Ellis (2005), kteří úlohu formulují jako klasifikační. Vstupní okno signálu jejich metoda klasifikuje do jednotlivých tříd tónů melodie, výstup je tedy kvantizován na rozlišení jednoho půltónu. Pro tento účel využívají metodu podpůrných vektorů (SVM) přičemž vstupními příznaky je část spektrogramu signálu, konkrétně vektor s 256 složkami, který získávají pomocí krátkodobé Fourierovy transformace podvzorkovaného signálu. Metoda dosahovala průměrných výsledků v soutěžích MIREX 2005 a 2006. Strojové učení se v oboru příliš neuchytilo, z důvodu nedostupnosti dat, ale možná také i kvůli limitacím, které Poliner a Ellis (2005) ve své práci prezentují.

V roce 2016 na tyto nedostatky odpovídá práce Kum a kol. (2016). Augmentací dat a zaměřením se na odhad výšky zpěvu místo melodie úspěšně překovávají problém nedostatku dat. Místo klasifikátoru SVM používají hluboké neuronové sítě (3 skryté plně propojené vrstvy) a problém kvantizace na úroveň půltónu řeší natrénováním tří nezávislých sítí s různě jemným rozlišením výstupu, které pak spojí v jeden výstup. V jejich případě platí, že jemná síť má sice lepší výstupní rozlišení, celkově má ale horší přesnost, naopak je tomu u sítě s hrubým rozlišením. Proto jejich výsledky kombinují a dosahují tehdejších state-of-the-art výsledků na vokálních datech. Podobný postup v tomtéž roce podniká i tým Rigaud a Radenen (2016), podobně jako v práci Kum a kol. (2016) se tým omezuje na odhad výšky zpěvu, používá augmentaci dat, přechází na hlubokou neuronovou síť (2 skryté plně propojené vrstvy) a používá jemnější rozlišení výsledného vektoru pravděpodobnostního rozdělení znějící výšky. Rozdílem je ale předzpracování dat pomocí dekompozice na harmonické a perkusivní složky pomocí HPSS. Práce Balke a kol. (2017) používá neuronovou síť s jednou skrytou vrstvou pro anotaci jazzových sol ze spektrogramu s logaritmickou osou frekvence.

Bittner a kol. (2017) představuje první pokus o použití hlubokých sítí na úlohu extrakce melodie bez zaměření na zpěv. Úlohu formuluje jako odstranění šumu obrázku, cílem je ze vstupního spektrogramu pomocí hluboké sítě s konvolučními vrstvami vytvořit salienční funkci. Vstupní a výstupní data mají tedy stejné měřítko, výstup však v ideálním případě obsahuje pouze informace o fundamentálních frekvencích znějící melodie. Velmi přínosný je také popis vstupní spektrální reprezentace HCQT, která spočívá ve výpočtu několika CQT spektrogramů (kanálů), jejichž počáteční frekvence jsou vzdálené v harmonických poměrech, tudíž (protože CQT spektrogram používá logaritmickou osu frekvence) všechny související harmonické složky na celém spektru jsou na těchto spektrogramech zarovnané nad sebou v ose kanálů. Tento koncept je hlouběji popsán v kapitole Experimenty v sekci HCNN, kde tento nápad aplikujeme nejen na vstupní spektrogram, ale nově také na propojení celé neuronové sítě. Dalším důležitým přínosem práce je úprava reprezentace cílové salienční funkce pomocí rozostření hranic cílové výšky tónu. Zatímco předchozí metody trénovali síť jako diskrétní klasifikaci s jednou správnou výstupní třídou, v této práci je cíl trénování gaussian se střední hodnotou uprostřed výšky tónu.

Práce Bittner a kol. (2017) dokázala překonat state-of-the-art metody pouhým výpočtem salienční funkce. Metoda tedy úplně přeskakuje vyhlazování odhadů v čase a používá pouze nejzákladnější metody pro detekci melodie pomocí práhování. Zároveň však její salienční funkce pro odhad jednoho okna délky ≈ 11 ms

zpracovává přibližně 150 ms vstupního okna, tedy v porovnání s existujícími metodami její salienční funkce zpracovává velmi široký kontext, proto není vyhlazování jejich výsledků nezbytné.

D Basaran, S Essid (2018) na této práci buduje a jednak prozkoumává možnosti použití jiné vstupní spektrální reprezentace (založené na práci Durrieu a kol. (2011)) a jednak do architektury sítě zabudovává rekurentní neuronové sítě, které zajišťují zmiňované vyhlazování odhadů. Jeho metoda překonává výsledky Bittner a kol. (2017), nevýhodou jeho přístupu je však hrubý výstup s rozlišením na půltóny.

V rámci soutěže MIREX také přibývají od roku 2016 nové metody založené na hlubokých sítích, ty se však buď stále zaměřují pouze na extrakci zpěvu (například Su (2018)) případně k nim nelze dohledat související článek (v roce 2016 metody účastníka Zhe-Cheng Fan a v roce 2018 metoda od týmu Sanguen Kum, Juhan Nam).

2.1.3 Hledání melodie

Po výpočtu funkce salience $S_x(f_\tau, \tau)$ máme k dispozici odhady fundamentálních frekvencí v signálu, z těchto ohodnocení pak musíme vybrat výslednou konturu melodie. V závislosti na způsobu výpočtu funkce salience tyto ohodnocení frekvencí více či méně odpovídají jejich důležitosti v signálu. Triviálním řešením zpracování těchto hodnot by bylo vybrat frekvence s maximální salienčí pro každé časové okno $\hat{f}_\tau = \operatorname{argmax}_{f_\tau} S_x(f_\tau, \tau)$, tento jednoduchý přístup však mnoho metod nevolí, protože jeho výstup má u složitějších skladeb tendenci „přeskakovat“ mezi doprovodem a melodií.

Obecně se dají přístupy rozdělit na pravidlové metody a statistické metody, dále se pak metody liší v tom, zda na základě salienční funkce vytváří abstraktnější popisy obsahu - ať už na úrovni jednotlivých celistvých konturů, tónů nebo celých frází.

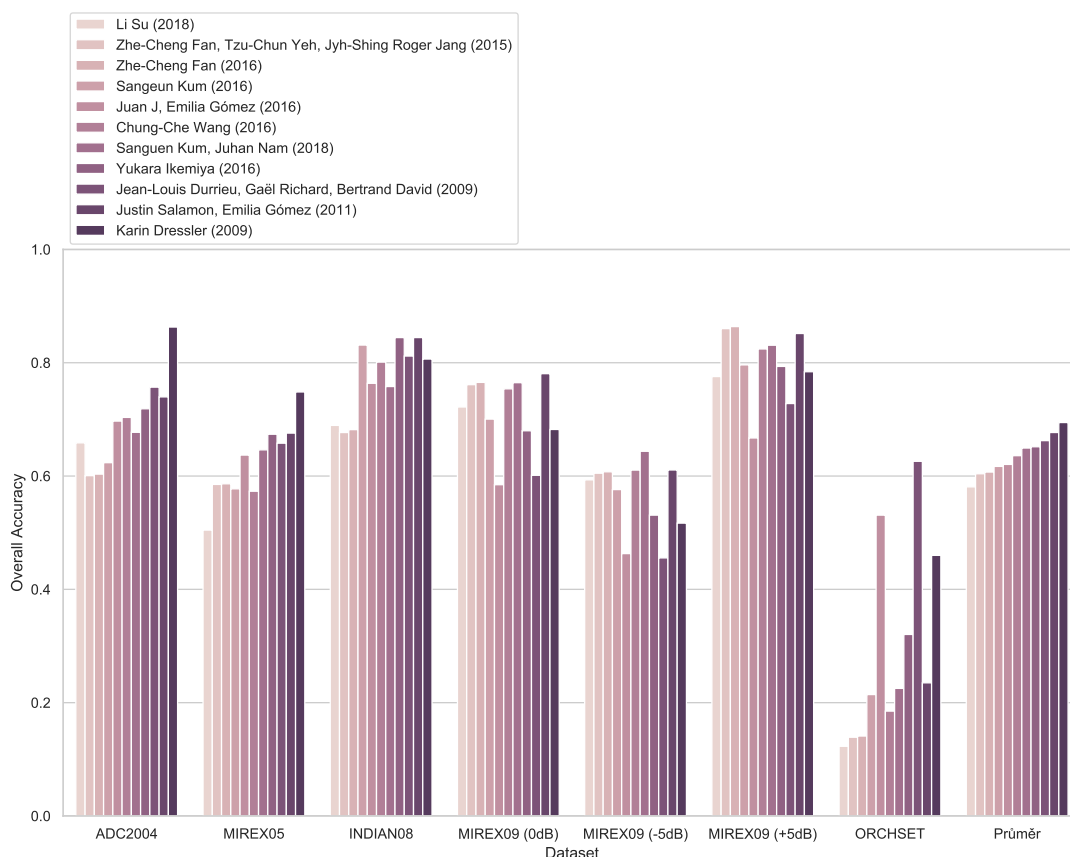
Goto a Hayamizu (1999) pro sledování melodie používá množinu „agentů“, pohybujících se v čase po výstupu salienční funkce a na základě předem definovaných pravidel jejich pohyb zaručuje kontinuitu výstupní fundamentální frekvence. Podobné sledování kontur v čase na základě salienční funkce využívá i Dressler (2009). Jinou pravidlovou metodou je opakované nalezení globálního maxima, jeho iterativní prodlužování v obou směrech časové osy a následné vymazání této nalezené kontury z výstupu salienční funkce, čímž dovolíme nalezení nového globálního maxima (Cancela (2008), Salamon a Gomez (2012)). Z extrahovaných kontur následně můžeme vybrat ty, které splňují kritéria pro melodické kontury.

Jiným přístupem k hledání melodie je použití statistických metod, jako je například modelování pomocí skrytých Markovových modelů. Na salienční funkci tyto metody pohlíží jako na sérii pozorování a pomocí hledání nejpravděpodobnější cesty skrz stavy modelu s vhodně nastavenými či z dat získanými pravděpodobnostmi přechodu získávají vyhlazenou konturu melodie. Tyto modely mohou být velmi komplexní a můžou zahrnovat modely průběhu not (Ryynänen a Klapuri, 2008) nebo naopak velmi minimalistické, zahrnující pouze stavy pro jednotlivé tóny (Yeh a kol. (2012))

Přítomnost melodie (voicing)

Důležitou součástí algoritmů pro extrakci melodie je detekce melodie v signálu. Většina metod tento krok provádí na konci vyhodnocování pomocí pevně nastaveného či dynamického práhování, jiné metody detekci melodie řeší filtrováním melodických kontur (Salamon a Gomez (2012)). V případě statistických metod je stav neznějící melodie často přímo zabudován do statistického modelu (Ryynänen a Klapuri, 2008). Některé metody také používají klasifikační metody strojového učení (Rigaud a Radenen, 2016).

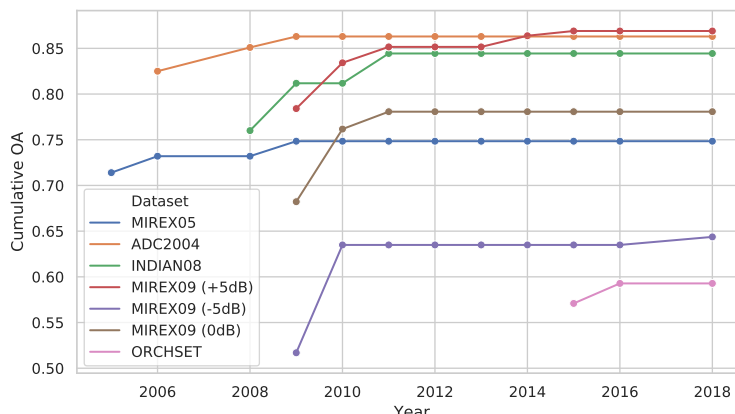
2.2 Srovnání existujících metod



Obrázek 2.4: Výsledky metod v soutěži MIREX v letech 2015-2018 s vybranými metodami ze starších ročníků.

Pro celkové kvantitativní srovnání metod jsme zpracovali výsledky všech ročníků soutěže MIREX. Tuto soutěž a používané datasety blíže popisujeme v kapitole Metody evaluace. V této sekci prezentujeme shrnutí výsledků metod, pro které existují výsledky na všech evaluačních datasetech soutěže MIREX, vybíráme tedy převážně z metod od roku 2015, kdy byl vydán doposud nejnovější evaluační dataset ORCHSET. Díky práci Bosch a Gómez (2014), který starší metody spustil na svém datasetu ORCHSET, můžeme ke srovnání přidat také výsledky metod Dressler (2009), Salamon a Gomez (2012) a Durrieu a David (2010). Celkové srovnání nalezneme v tabulce 2.4. Na základě těchto výsledků vybíráme metodu

Salamon a Gomez (2012), jejíž implementace je volně dostupná a spolu s prací Dressler (2009) dosahuje v průměru na datasetech MIREX nejlepších výsledků. V této práci ji lze považovat také jako zástupce metod, které nejsou založeny na strojovém učení.¹



Obrázek 2.5: Stagnující vývoj metod pro extrakci melodie.

Na základě grafu 2.4 také vidíme, že největší variabilitu mají výsledky na datasetu ORCHSET, zde mají metody velký prostor pro zlepšení, naopak u datasetů INDIAN08 a variant MIREX09 je zřejmá jistá hranice kterou je pro metody obtížné překonat. Salamon a kol. (2014) ve svém přehledovém článku dochází k závěru, že vývoj metod extrakce melodie začal od roku 2009 stagnovat, na obrázku 2.5 znázorňujeme maximální dosaženou celkovou přesnost metod na jednotlivých datasetech od počátku soutěže MIREX. Bohužel musíme konstatovat, že v rámci soutěže MIREX stagnace pokračuje doposud, přitom výzkum metod stále pokračuje a zejména díky strojovému učení se obor posouvá. V rámci MIREXu však nebyly vyhodnoceny nové stěžejní metody Bittner a kol. (2017) a D Basaran, S Essid (2018) a zájem o soutěž v této kategorii postupně upadá (v roce 2017 nesoutěžily žádné týmy, v roce 2018 pouze dva). Důvodem může být právě nedostatečný prostor pro zlepšení kvůli nedostatku nových, zajímavých dat.

2.2.1 Replikace výsledků

Pro srovnání metod představovaných v této práci spouštíme metody Salamon a Gomez (2012), Bittner a kol. (2017) a D Basaran, S Essid (2018) na testovacích množinách. Všechny tři metody mají volně dostupnou implementaci, první ve formě VAMP plug-inu,² zbylé jsou implementovány v jazyce Python a používají standardní knihovny určené pro hluboké učení.³⁴ Výsledky těchto metod uvádíme v kapitole Výsledky.

Poznamenáme, že implementace algoritmu Salamon a Gomez (2012) existují dvě, druhá v rámci knihovny Essentia,⁵ tato implementace však v porovnání s

¹Pro metodu Dressler (2009) implementace zveřejněná není.

²<https://www.upf.edu/web/mtg/melodia>

³<https://github.com/rabitt/ismir2017-deepsalience/>

⁴https://github.com/dogacbasaran/ismir2018_dominant_melody_estimation

⁵<https://essentia.upf.edu/documentation/>

implementací VAMP podávala výrazně horší výsledky napříč datasety. V kapitole Výsledky proto používáme implementaci VAMP.

Pokusili jsme se také o replikaci výsledků Bosch a Gómez (2014), bohužel se nám ale kvůli nekompatibilitě mezi verzemi knihoven nepodařilo tuto 5 let starou metodu spustit

3. Datasets

Nedostupnost dostatečného množství dat pro automatickou transkripci melodie představuje zejména pro metody strojového učení značný problém. Zatímco pro vzdáleně příbuznou úlohu automatického přepisu mluveného slova existuje tisíce hodin nahrávek (například dataset LibriSpeech, který vznikl na základě audioknih), největší dataset s přepsanou melodickou linkou MedleyDB má celkovou délku pod šest hodin. Do roku 2014, kdy MedleyDB vznikl, existovaly datasety, které byly buď rozmanité, ale příliš krátké (ADC04, MIREX05, INDIAN08) nebo naopak celkově větší, ale žánrově a hudebně homogenní (MIREX09, MIR1K, RWC). V roce 2015 byl vydán dataset Orchset, který obsahuje 23 minut výňatků z orchestrálních skladeb různých období. Za dataset pro extrakci melodie se také dá považovat Weimar Jazz Database, který je sice primárně zaměřený na využití v muzikologii, nicméně obsahuje přes 450 přepsaných jazzových sól. Novinkou z roku 2017 je vydání datasetu MDB-melody-synth, který byl automaticky vygenerován základě vstupní vícestopé hudby (převzaté z MedleyDB), existuje tedy naděje, že současný korpus pro přepis melodie by se mohl v budoucnu rozšířit o velkou část automaticky přesyntetizovaných, veřejně dostupných vícestopých nahrávek.

Co se týče blízké úlohy transkripce hudby, velikost největších datasetů se pohybuje v řádu desítek hodin, tudíž jde stále o omezené kolekce. Mezi největší se řadí MusicNet (orchestrální, 34 hodin), MAPS (klavír, 18 hodin), MDB-mf0-synth (multižánrový, 4,7 hodin), GuitarSet (kytara, 3 hodiny) a URMP (komorní orchestr, 1,3 hodiny). I když jde o úlohu, která je lépe definovaná (na rozdíl od extrakce melodie zde nehraje roli subjektivita volby hlavního hlasu), s použitím polyfonních nástrojů vyvstává problém náročné ruční anotace.

Vytváření nových datasetů je obecně velmi pracné a nákladné. Obvyklý postup totiž zahrnuje buď kompletní ruční přepis nahrávky nebo alespoň ruční opravu výstupu automatického přepisu jednohlasých nahrávek, přičemž tuto práci odvedou kvalitně pouze zaškolení hudebníci. Každá vzniklá anotace se také musí překontrolovat, a to nejlépe jiným hudebníkem. Dalším problémem je vůbec identifikace melodie - jelikož je určení hlavní melodické linie subjektivní, musí se na výsledné anotaci shodnout co nejvíce posluchačů. Ve výsledku se proto do datasetů buď vybírají takové nahrávky, které nejsou sporné, nebo na každé anotaci pracuje celý tým, který melodii společně určí. S tím souvisí také zavedení a pečlivé dodržování anotační politiky u komplexnějších skladeb (například orchestrálních), kde může melodii nést více hlasů zároveň současně či střídajíc se. Také množství výchozích dat pro vznik datasetů není velké. Jednak musí být skladby šiřitelné, pokud má být dataset volně dostupný a jednak by k nim měly být dostupné *audio stopy* (nahrávky samostatných hlasů), ze kterých je smíchán finální mix, jelikož ruční anotace finálního mixu je mnohem náročnější než anotace oddělených stop.

Existence dostatečně velkých datasetů je obecně vzato zásadním předpokladem pro využití metod strojového učení pomocí hlubokých neuronových sítí, zejména pak pro netriviální úlohy, jakou je například přepis melodie, jelikož dovoluje zvětšení celkové kapacity modelu, aniž by docházelo k přeučení. Také pro evaluaci metod, například i v soutěži MIREX, jsou potřeba takové datasety, které dobře reprezentují reálná data, přitom dataset MedleyDB vznikl mimo jiné z dů-

vodu, že stávající datasety nestačily ani pro účel evaluace.

V následující sekci uvádíme přehled veřejně dostupných dat a jejich společnou strukturu, po této sekci následuje podrobnější popis jednotlivých datasetů.

3.1 Struktura dostupných dat a jejich přehled

Dataset, který chceme použít pro řešení úlohy extrakci melodie, musí obsahovat soubory se zvukem a k nim příslušící anotace melodie. Standardním formátem zvukových souborů je jedno- nebo vícekanálový formát WAVE, se vzorkovací frekvencí 44 100 Hz. Anotace melodie je uložena jako CSV soubor s dvěma sloupci — časem a frekvencí. Výška melodie je tedy určena její fundamentální frekvencí a je specifikována pro každý časový okamžik v nahrávce. Délka anotačního okna je standardně 10 ms, případně $\frac{256}{44100} \doteq 5.8$ ms. Nepřítomnost melodie se označuje hodnotou 0.

Výjimkou je dataset ORCHSET, který neobsahuje přesné anotace fundamentální frekvence melodie, ale pouze frekvence not. Tedy frekvence nejsou spojitě, nýbrž jsou omezené na přesnost jednoho půltónu (V tabulce 3.1 je tato informace zohledněna řádkem MIDI melodie). Důležitou poznámkou je, že zde nejde o diskretizaci původní spojitě křivky, ale opravdu jde o anotaci not, tedy pokud melodii nese nástroj hrající vibrato a svou výškou se dostane nad rozsah jednoho půltónu, v anotaci tato skutečnost není zaznamenána.

Formát CSV pro zápis anotací melodie, který byl zaveden v rámci soutěže MIREX, dodržují všechny dostupné datasety a ačkoli existují pokusy o změnu tohoto formátu (Humphrey a kol., 2014), MIREX formát je natolik jednoduchý a prozatím dostačující, že k přechodu na sofistikovanější formáty zatím nedošlo. Pro ilustraci přikládáme část referenční anotace ženského zpěvu, v anotaci se vyskytuje krátká pomlka mezi znějícími tóny. Grafické znázornění referenční anotace melodie celé nahrávky můžeme nalézt v úvodu na obrázku 1.1.

8.568	381.349
8.574	379.959
8.580	378.229
8.586	376.067
8.591	372.236
8.597	369.793
8.603	0.000
8.609	0.000
8.615	0.000
8.620	0.000
8.626	0.000
8.632	0.000
8.638	352.272
8.644	338.922

Datasety však mohou obsahovat více informací či audio souborů. Užitečné jsou například přiložené audio stopy, ze kterých je vytvořena výsledná píseň (mix),

informace o všech znějících výškách (Multi-F0) nebo notách (MIDI), o melodické prioritě jednotlivých audio stop nebo o instrumentaci skladby.

V tabulce 3.1 nalezneme přehledné shrnutí obsahu všech dostupných datasetů.

	MedleyDB	Orchset	ADC04	MIREX05 train	MDB-synth	WJAZZD	MIR-1K	RWC
Audio	Ano	Ano	Ano	Ano	Ano	Ne ¹	Ano	Ano ²
F0 melodie	Ano	Ne	Ano	Ano	Ano	Ano	Ano	Ano
MIDI melodie	Ne	Ano	Ne	Ne	Ne	Ano	Ne	Ne
Audio stopy	Ano ³	Ne	Ne	Ne	Ano	Ne	Ano ⁴	Ne
Multi-F0	Ne ⁵	Ne	Ne	Ne	Ano	Ne	Ne	Ne
MIDI	Ne	Ne	Ne	Ne	Ne	Ne	Ne	Ano
Priorita stop	Ano	Ne	Ne	Ne	Ano	Ne	Ne	Ne
Informace o instrumentaci	Ano	Ano	Ne	Ne	Ano	Ano	Ne	Částečné
Celková délka	7.3 h ⁶	23.4 m	6.1 m	6.5 m	3.19 h	8.85 h	2.22 h	—
Poměr znějící melodie	60.9%	93.69%	85.7%	63.1%	50.4%	62.8%	—	—
Počet nahrávek	122 ⁷	64	20	13	65	299	1000	315
Webová stránka	⁸	⁹	¹⁰	¹¹	¹²	¹³	¹⁴	¹⁵
Žánr	mnoho- žánrový	klasika	pop,jazz, opera,midi	pop, midi	mnoho- žánrový	jazz	karaoke	pop, jazz klasika
Účel v této práci	Trénování Validace Testování	Testování	Testování	Testování	Testování	Testování	Žádný	Žádný

Tabulka 3.1: Souhrnná tabulka se základními informacemi o veřejně dostupných datasetech.

3.2 MedleyDB

MedleyDB je žánrově rozmanitý dataset obsahující 122 nahrávek, k 108 z nich je dostupná anotace melodie. Kromě té dataset obsahuje také metadata o všech písních s informacemi o žánru a instrumentaci. S celkovou délkou 7.3 hodiny jde o nejdelší volně dostupný dataset, který obsahuje více žánrů hudby. O rozmanitosti svědčí i to, že se v datasetu vyskytuje řada nástrojů mimoevropského původu, a že jen přibližně polovina písní obsahuje zpěv. Na rozdíl od ostatních datasetů jsou nahrávky ve většině případů celé písně, tedy nejde pouze o krátké výňatky, a ke každé jsou poskytnuty audiostopy, ze kterých je vytvořen výsledný mix. Na základě diskuze, kterou shrnujeme v kapitole o definici melodie, autoři datasetu Bittner a kol. (2014) poskytují tři verze anotací, na základě různých obecných

¹Autoři audio poskytují neveřejně pro výzkumné účely

²Přístup k datasetu je zpoplatněn

³Část stop obsahuje přeslech ostatních nástrojů, informace o přeslechu je součástí metadat každé skladby.

⁴Oddělený zpěv a karaoke doprovod

⁵Je dostupný přepis všech znějících melodií, viz Definice 3 v sekci MedleyDB.

⁶5.59 h s anotací melodie

⁷108 s anotací melodie

⁸<https://medleydb.weebly.com/>

⁹<https://www.upf.edu/web/mtg/orchset>

¹⁰http://ismir2004.ismir.net/melody_contest/results.html

¹¹<https://labrosa.ee.columbia.edu/projects/melody/>

¹²<http://synthdatasets.weebly.com/mdb-melody-synth.html>

¹³<https://jazzomat.hfm-weimar.de/>

¹⁴<https://sites.google.com/site/unvoicedsoundseparation/mir-1k>

¹⁵<https://staff.aist.go.jp/m.goto/RWC-MDB/>

definice:

1. Základní frekvence nejvýraznějšího melodického hlasu, jehož zdroj zůstává po dobu nahrávky neměnný.¹⁶
2. Základní frekvence nejvýraznějšího melodického hlasu, jehož zdroje se mohou měnit.
3. Základní frekvence všech melodických hlasů, potenciálně pocházejících z více zdrojů.

Ačkoli třetí definice umožňuje, aby v anotaci znělo více melodických linek zároveň, v datasetu se nejedná o kompletní přepis nahrávek (použitelný pro úlohu multi-f0 estimation, tedy pro úplný přepis všech fundamentálních frekvencí znějících tónů), ten autoři neposkytují.

Dataset vznikl obvyklou cestou ruční anotace. Ze shromážděného vícestopého materiálu byly vybrány stopy s potenciálním výskytem melodie, stopy s přeslechem byly předzpracovány pomocí algoritmu pro oddělení hlasu a doprovodu (source-separation) s ručně doladěnými parametry pro každou jednotlivou stopu, následně byla na monofonní stopy spuštěna metoda pYIN pro odhad výšky v monofonních datech (pitch tracker) a výsledné automaticky získané anotace opravilo a vzájemně zkontrolovalo pět anotátorů s hudebním vzděláním.

3.3 Orchset

Dataset vytvořený týmem Bosch a kol. (2016) orientovaný na orchestrální repertoár pocházející z různých historických období včetně 20. století. Obsahuje 64 výňatků délky od 10 do 32 sekund. Výňatky byly vybírány tak, aby obsahovaly zřejmou melodii, dataset tedy obsahuje v porovnání málo pasáží bez melodie (6% z celkové délky). Vzhledem k komplexitě uvažovaných žánrů autoři vycházejí z kombinace rozšířené definice melodie podle Bittner a kol. (2014) a definice Poliner a kol. (2007). Melodii ve výňatcích proto zpravidla nese více hudebních nástrojů (nebo celých sekcí), které se v průběhu střídají, případně mohou části hrát společně v rozdílných oktávách (nebo jiných intervalech, tvoříce tak harmonický doprovod).

Pro zjištění melodie se v takto vrstveném materiálu autoři uchylují k úplnému základu definice melodie (Poliner a kol. (2007)) a nechávají si skupinou čtyř posluchačů výňatky přezpívat. Tato hrubá data pak autoři sumarizují a odebírají z datasetu ty výňatky, na jejichž melodii se posluchači neshodli. Přezpívané tóny bylo nutné ručně opravit, aby načasováním přesně seděly na výňatek. Lidský hlas také samozřejmě nemá rozsah plného orchestru, proto bylo dalším krokem transponovat anotace tak, aby zněly ve správných oktávách. Zde se opět může vyskytnout problém subjektivity, pokud melodii hraje dva různé nástroje, pouze v jiných oktávách, pak je sporné, který nástroj označit jako hlavní, a v některých případech taková otázka ani nedává příliš smysl. Částečným řešením je zvolit libovolnou anotační politiku a tu konzistentně dodržovat (žádná společná

¹⁶Tato definice je shodná pro evaluační datasety používané v soutěži MIREX, s výjimkou Orchsetu

v komunitě MIR neexistuje), v případě Orchsetu byla snaha minimalizovat skoky v melodické kontuře, což zároveň respektuje obecné pozorování, že v melodii se vyskytují mnohem častěji malé skoky mezi tóny (nejčastěji prima a malá/velká sekunda) než větší. Tedy například pokud pasáží hrané ve dvou různých oktávách předcházela pasáž hraná v jedné, anotace obou pasáží lze transponovat do společné oktávy tak, abychom na rozhraní těchto pasáží minimalizovali skok v anotaci.

Dataset obsahuje pouze hrubé anotace tónů melodie, nikoli přesnou základní frekvenci nástroje, který v danou chvíli melodii hraje. Článek o tomto rozhodnutí příliš nediskutuje, vychází ale opět logicky z volby dat. U orchestrálních dat je tento abstraktnější pojem melodie mnohem méně sporný. Pokud hraje melodii sekce nástrojů v unisonu, přesná základní frekvence není dobře definovaná, jelikož se základní frekvence znějících hlasů vzájemně překrývají.

3.4 MIREX datasety

Datasety MIREX05 train a MIR-1K byly vydány jako trénovací data v rámci soutěže MIREX. Jde o malé množství dat, MIREX05 train se skládá z několika anotovaných populárních skladeb a několika syntetizovaných písní z MIDI sborů, MIR-1K obsahuje 1000 úryvků zpěvu s karaoke doprovodem. První dataset používáme jako testovací, druhý vzhledem k dostupnosti jiných, rozmanitějších testovacích dat nepoužíváme.

Dataset ADC2004, použitý ve stejnojmenné soutěži, která předcházela vzniku MIREXu, byl po konci soutěže zveřejněn včetně testovací množiny, stále je však využíván jako jeden z testovacích datasetů v soutěži MIREX. Celý dataset proto také používáme jako testovací. Obsahuje 20 výňatků ze žánrů popu, jazzu a opery a dále pak 4 syntetické skladby.

3.5 Weimar Jazz Database

Weimar Jazz Database (práce německého týmu Pfeiderer a kol.) obsahuje přes 450 transkripce jazzových sol ze všech období vývoje jazzu. Data původně zamýšlená pro muzikologické studie využívající statistické metody ale lze využít i pro potřeby extrakce melodie, jelikož uvažované nahrávky spadají zřejmě pod nejrestriktivnější definici melodie (definici používanou v soutěži MIREX) - melodii nese jistě právě jeden, sólový nástroj, a po celou dobu výňatku je jistě nejvýraznější. Výběr sólových nástrojů se omezuje pouze na jednohlasé, jelikož ruční anotace vícehlasých je příliš obtížná. Hlavním problémem při využívání je restriktivní licence, která platí na nahrávky, tudíž zdrojové audio, na základě kterého anotace vznikaly, není veřejně přístupné.

Dataset v této práci používáme pouze pro testování metod, nikoliv trénování.

3.6 MDB-synth

Hlavním přínosem práce Salamon a kol. (2017) je navržení způsobu anotace základní frekvence monofonních audiostop takovým způsobem, že výsledná dvojice

zvukové stopy a anotace nevyžaduje další manuální kontrolu. Anotace monofonní stopy probíhá ve dvou krocích: nejprve získáme pomocí libovolné metody přepisu jednohlasu křivku základní frekvence. Poté na základě této křivky, která může obsahovat chybně anotované úseky, syntetizujeme novou stopu, která zachovává barvu nahrávky, ale výšku tónu určuje právě tato automaticky získaná anotace. Díky tomu je pak přesnost anotace pro tuto novou, syntetickou nahrávku stoprocentní, přitom (v ideálním případě) neztrácí charakteristiky původní nahrávky.

Pro vytváření datasetu je toto významné zjednodušení, protože tím algoritmus odstraňuje časově nejnáročnější část práce — ruční kontrolu anotací audiostop. Pokud by se ukázalo, že syntéza významně neubírá na kvalitě dat, použitím navrhované metody by mohlo vzniknout velké množství nových dat (například repozitář Open Multitrack Testbed obsahuje stovky vícestopých nahrávek, které by bylo možné využít). Autoři v článku provádí kvantitativní analýzu pomocí srovnání state-of-the-art algoritmů pro extrakci melodie a prokazují, že výsledky těchto metod na syntetických datech se významně neliší od výsledků na původních, tím je podle autorů potvrzená možnost použití dat jak pro trénování tak pro evaluaci nových metod.

Metoda má ale bohužel svá omezení. Mezi ty zásadní patří, že se dá aplikovat pouze na stopy, které obsahují monofonní signál, vstupní data tedy nesmí obsahovat přeslech a nahrávaný nástroj může hrát pouze jednohlas. V důsledku tedy nelze zpracovat například klavír či kytara, které hrají zpravidla vícehlas. Pro generování datasetu určeného pro přepis melodie tato limitace není zásadní, jelikož melodii často hraje jeden hlas a doprovod může být vícehlasý. Problémem je spíše generování datasetů pro úlohu kompletního přepisu (multi-f0 estimation).

Bohužel také k článku není zveřejněná kompletní refereční implementace algoritmu, tudíž algoritmus nelze snadno spustit na nových datech. Ve výsledku je tudíž největším praktickým přínosem nová sada syntetických dat pro úlohy přepisu melodie, přepisu basy, přepisu jednohlasu a kompletního přepisu. Každý z datasetů určených pro zmíněné úlohy obsahuje destíky nahrávek. Vícestopá data použitá pro syntézu byla převzata z MedleyDB, tudíž nové datasety nerozšiřují celkový hudební záběr, pouze zpřesňují již existující.

Dataset v této práci používáme pouze pro testování metod, nikoliv trénování.

3.7 Dataset RWC

Dataset RWC (práce Goto a kol. (2002)) je první vzniklá rozsáhlá kolekce dat určená pro úlohy Music Information Retrieval, mezi které v tomto případě patří i extrakce melodie. Dataset obsahuje 100 popových, 50 orchestrálních a 50 jazzových skladeb. Přístup k datasetu RWC je však zpoplatněn, proto ho v této práci nepoužíváme.

4. Metody evaluace

V této sekci prezentujeme zavedené a standardní způsoby vzájemného kvantitativního i kvalitativního porovnávání výsledků metod pro extrakci melodie, které vychází z postupů evaluace v soutěži MIREX. V první sekci uvádíme soutěž MIREX jako takovou, následně definujeme používané metriky a v závěru kapitoly zmiňujeme některé praktické metriky používané při trénování modelů v této práci.

4.1 MIREX

Soutěž MIREX (Music Information Retrieval Evaluation eXchange) probíhá již od roku 2005 a v MIR komunitě zastává hlavní postavení jakožto každoroční událost pro nezávislé, objektivní srovnání state-of-the-art metod a algoritmů pro řešení širokého spektra úloh souvisejících se zpracováním hudebních dat. Mezi tyto úlohy patří například rozpoznání žánru, odhad tempa, odhad akordů, identifikace coveru a samozřejmě také extrakce melodie.

Na rozdíl od jiných úloh, kde debata o zvolení nejvhodnějších objektivních metrik pro porovnávání metod stále probíhá, metriky pro extrakci melodie se ustanovily již v prvním ročníku (na základě dřívějších zkušeností) a zůstaly neměnné dodnes (Raffel a kol., 2014). Naopak data použitá pro testování se postupně kumulují a dnes soutěž probíhá již s řadou datasetů (ADC04, MIREX05, MIREX08, MIREX09, ORCHSET15). V tabulce 4.1 uvádíme přehled těchto evaluačních datasetů.

4.1.1 Formát výstupu v soutěži MIREX

Obvyklý formát výstupu algoritmů je CSV soubor se dvěma sloupci. První sloupec obsahuje pravidelné časové značky, druhý sloupec pak odhad základní frekvence melodie. Některé algoritmy uvádí i odhady výšky základní frekvence mimo detekovanou melodii (může jít například o doprovod, který zní i po hlavním melodickém hlasu). Aby tyto odhady byly odlišené od odhadů hlavní melodie, jsou uvedeny v záporných hodnotách. Díky tomu pak lze nezávisle vyhodnotit přesnost *odhadu výšky* a *detekce melodie*. Odhad výšky se vyhodnocuje podle absolutní hodnoty frekvence ve všech časových oknech, ke kterým existuje anotace, detekce melodie pak na všech hodnotách vyšších než 0.

4.2 Trénovací, validační a testovací množina

Z dostupných dat, které pro úlohu máme k dispozici, musíme vyhradit množiny pro trénování, validaci a testování, aby byly metody porovnatelné jak mezi sebou, tak se stávajícími state-of-the-art metodami. Pro trénování se jeví jako nejvhodnější dataset MedleyDB, jednak pro svou délku a jednak pro žánrovou rozmanitost, proto je použit pro všechny popsání experimenty. Rozdělení datasetu na trénovací, validační a testovací množiny jsme převzali z práce Bittner a kol. (2017), díky tomu je korektní naše výsledky na testovací množině MedleyDB přímo porovnávat s výsledky uvedenými v článku. Tento postup využití

Dataset	Popis
ADC2004	20 výňatků délky kolem 20 sekund, žánry pop, jazz, opera. Obsahuje živé a syntetické MIDI nahrávky. Celková délka 369 sekund. Dataset je kompletně zveřejněn.
MIREX05	25 výňatků délky 10-40 sekund, žánry rock, R&B, pop, jazz a sólové piano. Obsahuje živé a syntetické MIDI nahrávky. Celková délka 686 sekund. K datasetu je dostupná trénovací množina MIREX05train
INDIAN08	Čtyři minutové výňatky tradičního zpěvu z oblasti severní Indie. Ke každému výňatku existují dvě verze s různě hlasitým doprovodem, celkově se tedy dataset skládá z osmi nahrávek. Celková délka 501 sekund. Dataset je neveřejný
MIREX09 (0 dB)	374 čínských karaoke písní. Poměr hlasitosti zpěvu k doprovodu je 0dB. Celková délka 10 020 sekund. K datasetu je dostupná trénovací množina MIR1K
MIREX09 (-5 dB)	Stejné výňatky jako v případě MIREX09 (0 dB), poměr hlasitosti zpěvu k doprovodu je však -5dB. Celková délka 10 020 sekund.
MIREX09 (+5 dB)	Stejné výňatky jako v případě MIREX09 (0 dB), poměr hlasitosti zpěvu k doprovodu je však +5dB. Celková délka 10 020 sekund.
ORCHSET15	64 orchestrálních výňatků z různých období vývoje klasické hudby. Anotace melodie je uvedena s přesností jednoho půltónu, Celková délka 1404 sekund. Dataset je kompletně zveřejněn.

Tabulka 4.1: Přehled evaluačních datasetů v soutěži MIREX, částečně převzato z práce Salamon a kol. (2014).

existujícího rozdělení dat používá ve své práci také D Basaran, S Essid (2018), i s jeho výsledky proto naši práci můžeme porovnávat.

Data byla rozdělena na trénovací, validační a testovací množinu tak, aby skladby od jednoho interpreta náležely právě do jedné z množin. Využili jsme existujícího rozdělení používaného v článcích Bittner a kol. (2017) a D Basaran, S Essid (2018), validační i testovací výsledky jsou tedy díky tomu porovnatelné s výsledky uvedenými v článcích.

Dalším užitečným zdrojem dat je dataset MDB-melody-synth, který je syntetizován z vícestopých nahrávek MedleyDB. Proto je vhodné použít stejné rozdělení dat, jaké se používá pro MedleyDB, protože data budou silně korelovaná s původními nahrávkami a tak bychom mohli dojít k chybným, uměle vysokým výsledkům, pokud bychom množiny změnili. Jelikož dataset neobsahuje veškerá data, ale pouze jejich podmnožinu, i v experimentech používané rozdělení dat obsahuje pouze podmnožinu z původního rozdělení datasetu MedleyDB.

Posledním velkým datasetem používaným v práci, je Weimar Jazz Database. Původním záměrem bylo dataset použít i pro trénování a validaci metod proto jsme jej rozdělili na datové množiny, nakonec byla však použita pouze testovací množina. Při rozdělování datasetu jsme postupovali podle práce Bittner a kol.

(2017), data jsme rozdělili na tři části. Skladby jsou rozděleny do částí podle interpretů tak, aby se každý interpret vyskytoval právě v jedné části datasetu. Toto omezení na podmnožiny Bittner a kol. (2017) nediskutuje, lze však doložit (práce Sturm (2013)), že pro úlohu rozpoznání žánru metody založené na strojovém učení vykazují po trénování a validaci na datech bez tohoto filtru výrazně lepší výsledky než stejné metody spuštěné na roztríděných datech, takové zlepšení výkonu je ale jistě umělým důsledkem špatné volby trénovací množiny.

Délky a počty nahrávek datových množin používaných v této práci shrnujeme v tabulce 4.2, přesný výčet skladeb použitých v každé množině je v příloze práce.

Dataset	Množina	Počet nahrávek	Délka
MedleyDB	train	67	3.06 h
	valid	15	0.85 h
	test	27	1.71 h
MDB-melody-synth	train	44	1.88 h
	valid	8	0.46 h
	test	14	0.87 h
WJazzD	train	169	5.49 h
	valid	54	1.25 h
	test	74	2.08 h

Tabulka 4.2: Přehled trénovacích, validačních a testovacích množin použitých v práci.

Dále k testování používáme datasety ADC04, MIREX05train a ORCHSET, blíže popsané v kapitole Datasety. Protože jsou datasety ADC04 a ORCHSET jsou v práci použity pouze jako testovací data, můžeme výsledky přímo srovnávat s odpovídajícími žebříčky úlohy Melody Extraction v soutěži MIREX.

4.3 Metriky

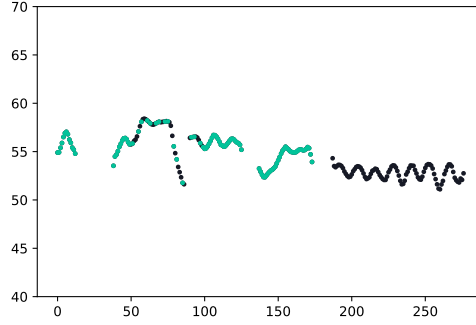
Celková kvalita výstupu metody pro extrakci melodie je určena její schopností odhadnout správně výšku tónu hrající melodie a také rozpoznat části skladby, které melodii neobsahují. Jelikož jsou tyto podúlohy na sobě nezávislé, standardní sada metrik zahrnuje jak celkové vyhodnocení přesnosti, tak dílčí vyhodnocení pro *odhad výšky* a *detekci melodie*.

V práci pracujeme s metrikami vypočítanými na validačních a testovacích sadách, ty se počítají jako průměr výsledků na jednotlivých skladbách v sadě.

4.3.1 Definice metrik

Většina metrik je definována na základě porovnávání jednotlivých anotačních oken — tedy typicky srovnáním odhadovaných a pravdivých výšek melodie po konstantních časových skocích. Datasety používané pro vyhodnocování v soutěži MIREX používají časový skok délky 10 ms. V definicích budu vycházet ze značení v práci Salamon a kol. (2014).

Označme vektor odhadovaných základních frekvencí \mathbf{f} a cílový vektor \mathbf{f}^* , složka f_τ je buď rovna hodnotě f_0 melodie, nebo 0, pokud v daném čase melodie nezní. Obdobně zavedme vektor indikátorů \mathbf{v} , jehož prvek na pozici τ je roven $v_\tau = 1$, pokud je v daném časovém okamžiku detekována melodie a $v_\tau = 0$ v opačném případě. Podobným způsobem zavedeme i vektor cílových indikátorů melodického hlasu \mathbf{v}^* a také vektor indikátorů absence melodie $\bar{v}_\tau = 1 - v_\tau$.

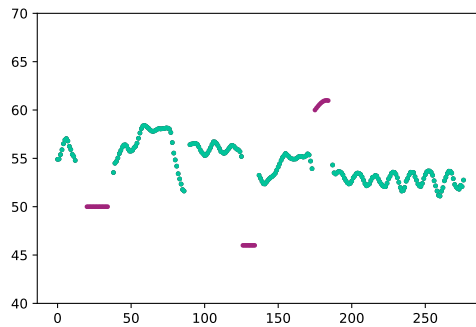


Obrázek 4.1: Příklady chyb špatné detekce melodie negativně ovlivňující metriku VR.

Voicing Recall rate (Úplnost detekce)

Poměr počtu časových oken, které byly správně označené jakožto obsahující melodii, a počtu časových oken doopravdy obsahujících melodii podle anotace.

$$\text{VR}(\mathbf{v}, \mathbf{v}^*) = \frac{\sum_{\tau} v_{\tau} v_{\tau}^*}{\sum_{\tau} v_{\tau}^*}$$

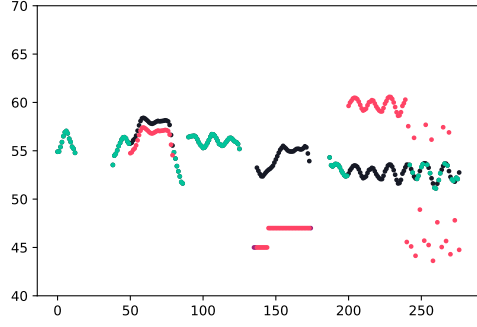


Obrázek 4.2: Příklady chyb špatné detekce melodie zvyšující metriku VFA.

Voicing False Alarm rate (Nesprávné detekce)

Poměr počtu časových oken, které byly nesprávně označené jako obsahující melodii, a počtu časových oken, které doopravdy melodii neobsahují. Pro interpretaci této metriky platí, že nižší hodnota je lepší, vyšší hodnota je horší.

$$\text{FA}(\mathbf{v}, \mathbf{v}^*) = \frac{\sum_{\tau} v_{\tau} \bar{v}_{\tau}^*}{\sum_{\tau} \bar{v}_{\tau}^*}$$



Obrázek 4.3: Příklady chyb ovlivňujících metriku RPA a RCA.

Raw Pitch Accuracy (Přesnost odhadu výšky)

Poměr správně odhadnutých tónů k celkovému počtu oken, které obsahují melodii. Výška správně určeného tónu se může lišit až o jeden půltón.

$$\text{RPA}(\mathbf{f}, \mathbf{f}^*) = \frac{\sum_{\tau} v_{\tau}^* v_{\tau} \mathcal{T}[\mathcal{M}(f_{\tau}) - \mathcal{M}(f_{\tau}^*)]}{\sum_{\tau} v_{\tau}^*}$$

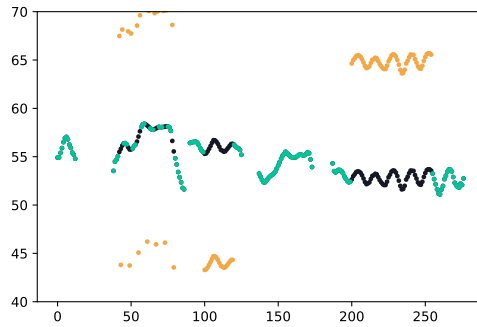
kde \mathcal{T} je prahová funkce

$$\mathcal{T}[a] = \begin{cases} 1 & \text{pro } |a| \leq 0.5 \\ 0 & \text{jinak} \end{cases}$$

Dodáme, že práh 0.5 je v některých případech použití vhodné snížit na restriktivnější hodnoty, jako je 0.25 nebo dokonce 0.1. Jedná se zejména o použití metriky pro úlohu odhadu výšky jednohlasu (jako například v práci Kim a kol. (2018)).

a \mathcal{M} je funkce zobrazující frekvenci f na reálné číslo počtu půltónů od nějakého referenčního tónu f_{ref} (například od 440 Hz, tedy komorního A4).

$$\mathcal{M}(f) = 12 \log_2\left(\frac{f}{f_{\text{ref}}}\right)$$



Obrázek 4.4: Příklady chyb „o oktávu“ ovlivňujících pouze metriku RPA, nikoli metriku RCA.

Raw Chroma Accuracy (Přesnost odhadu výšky nezávisle na oktávě)

Počítá se podobně jako *Přesnost odhadu tónu*, výstupní a cílové tóny jsou však mapovány na společnou oktávu. Metrika tedy ignoruje chyby odhadu způsobené špatným určením oktávy tónu.

$$\text{RCA}(\mathbf{f}, \mathbf{f}^*) = \frac{\sum_{\tau} v_{\tau}^* v_{\tau} \mathcal{T}[\langle \mathcal{M}(f_{\tau}) - \mathcal{M}(f_{\tau}^*) \rangle_{12}]}{\sum_{\tau} v_{\tau}^*}$$

Nezávislost na oktávě zajistíme pomocí zobrazení rozdílu cílového a výstupního tónu na společnou oktávu.

$$\langle a \rangle_{12} = a - 12 \lfloor \frac{a}{12} + 0.5 \rfloor$$

Overall Accuracy (Celková přesnost)

Celková přesnost měří výkon algoritmu jak v odhadu melodie tak v detekci melodie. Počítá se jako podíl správně odhadnutých oken a celkového počtu oken.

$$\text{OA}(\mathbf{f}, \mathbf{f}^*) = \frac{\sum_{\tau} v_{\tau}^* v_{\tau} \mathcal{T}[\mathcal{M}(f_{\tau}) - \mathcal{M}(f_{\tau}^*)] + \bar{v}_{\tau}^* \bar{v}_{\tau}}{L}$$

Poznámka k definicím metrik

Definice RPA, RCA a OA zde uvedené se mírně liší od výchozích v práci Salamon a kol. (2014), jejich přímá implementace podle vzorce totiž vede kvůli nedostatečně dobře zadanému vektoru frekvencí \mathbf{f} k chybě, která se týká zejména metriky RCA. Ta v původním znění definice chybně zahrnovala jako správné tóny ty, které algoritmus odhadl jako nulové (tedy neznějící) a zároveň jejich pravdivá hodnota byla po zobrazení na jednu společnou oktávu blízká nule (tedy původní tón byl blízký nějakému násobku referenční frekvence). Kvůli zobrazení na společnou oktávu se stanou „neznělé nulové odhady“ a tóny blízké referenčním frekvencím nerozlišitelné a byly nesprávně považované za korektní.¹

4.3.2 Další metriky

Protože princip vnitřního fungování neuronových sítí často není zřejmý, je užitečné mít co nejvíce různých indikátorů, abychom měli při porovnávání jednotlivých modelů alespoň podrobnou informaci, v jakých ohledech se síť zlepšuje nebo zhoršuje. Pro tento účel jsem při práci implementoval další metriky, které při hledání architektur sítí pomáhaly.

Chroma Overall Accuracy

Počítá se obdobně jako Overall Accuracy, ale tóny jsou mapovány na společnou oktávu.

¹Tato chyba byla přítomná i v nejpoužívanější, veřejné implementaci MIR metrik *mir_eval*. V praxi rozdíl chybné a opravené hodnoty této metriky na datasetu MedleyDB mohla dosahovat až sedmi procentních bodů, na repozitáři hostovaném na serveru Github jsme již spolu s autory chybu odstranili (odkaz na Github issue: https://github.com/craffel/mir_eval/issues/311). Opravný patch bude zahrnut do další verze balíku. Výsledky v této práci jsou počítány s opravenou verzí.

Raw Harmonic Accuracy

Metrika počítá odhadovaný tón jako správný, pokud se trefil do některé z harmonických frekvencí tónu. Protože je harmonických frekvencí teoreticky nekonečné množství, parametrem metriky je do jakého celočíselného násobku se ještě odhad počítá.

$$\text{RHA}(\mathbf{f}, \mathbf{f}^*, n) = \frac{\sum_{k=1}^n \sum_{\tau} v_{\tau}^* v_{\tau} \mathcal{T}[\mathcal{M}(f_{\tau}) - \mathcal{M}(kf_{\tau}^*)]}{\sum_{\tau} v_{\tau}^*}$$

Matice záměn not

Pro podrobnější souhrnný přehled četností chyb se pro klasifikační úlohy používá matice záměn. Sloupce označují správné noty, řádky odhadované. Buňka na pozici (x, y) má pak hodnotu podle četnosti odhadu noty y místo správné noty x .

Histogram vzdáleností odhadu

Histogram hodnot rozdílů $\mathcal{M}(\mathbf{f}) - \mathcal{M}(\mathbf{f}^*)$, tedy histogram vzdáleností odhadů výšky tónů od správné hodnoty.

Kvalitativní příklady

Modely byly při práci vyhodnocovány na několikaminutových množinách výňatků z validačních a testovacích dat. Metodika výběru spočívala v poslechu nahrávek a ručním výběru zajímavých hudebních jevů. Několik ukázek bylo také vybráno na základě seřazení nahrávek podle úspěšnosti přepisu stávajícími metodami a výběrem výňatků z tohoto seznamu nejproblematictějších příkladů.

5. Experimenty

V této kapitole prezentujeme návrhy nových architektur hlubokých neuronových sítí pro extrakci melodie a jejich výsledky na validační množině. Modely byly trénované nad datasetem MedleyDB. Data byla rozdělena na trénovací, validační a testovací množinu tak, aby skladby od jednoho interpreta náležely právě do jedné z množin. Využili jsme existujícího rozdělení používaného v článcích Bittner a kol. (2017) a D Basaran, S Essid (2018), validační i testovací výsledky jsou tedy díky tomu porovnatelné s výsledky uvedenými v článcích.

Navrhované sítě jsou koncipovány jako nové způsoby výpočtu funkce salience. Cílem je ze vstupního okna signálu získat ohodnocení znějících tónů, přičemž v ideálním případě nejvyšší ohodnocení bude mít nejvýraznější tón — tedy ten, který nese melodii. Naopak tóny doprovodu a frekvenční složky perkusí by měly být v této výstupní reprezentaci co možná nejvíce upozaděny. Zdůrazníme, že těžištěm této práce je odhad výšky tónů; pro detekci melodie používáme pouze jednoduchou metodu práhování. Po výpočtu funkce salience také nepoužíváme žádný modul pro vyhlazování predikcí v čase, pouze vybíráme maximální hodnotu výstupu v každém časovém okně. Tímto zaměřením pouze na výpočet funkce salience je práce podobná práci Bittner a kol. (2017).

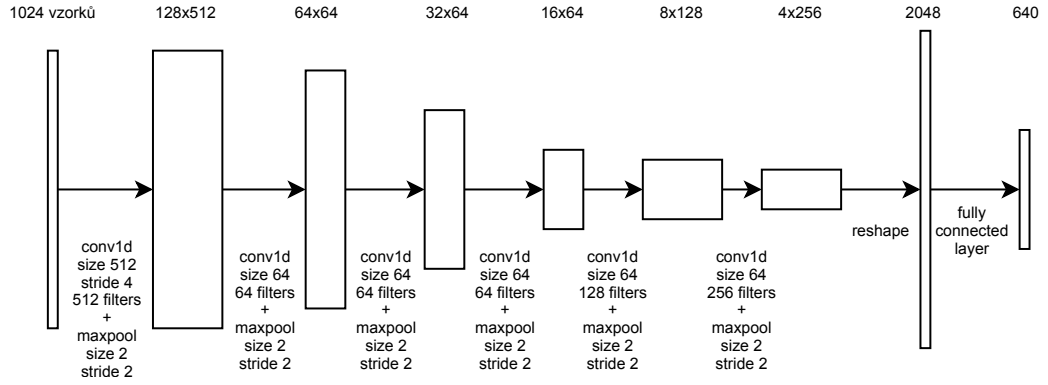
Pro výpočet salienční funkce porovnáváme tři různé architektury inspirované pracemi Kim a kol. (2018), van den Oord a kol. (2016a) a Bittner a kol. (2017), tyto architektury se pokoušíme systematickým hledáním přizpůsobit pro účel extrakce melodie pomocí vhodných úprav jak v jejich topologii tak nastavením hyperparametrů. První práce popisuje architekturu CREPE, původně navrženou pro sledování výšky tónů v jednohlasých nahrávkách. Architektura WaveNet popsaná týmem van den Oord a kol. (2016a) je složena z dilatovaných konvolucí a původně určena pro generování lidské řeči, autoři však zmiňují její možné použití i pro přepis řeči. Bittner a kol. (2017) používá speciální vstupní reprezentaci signálu, kterou zpracovává pomocí hluboké konvoluční sítě, její práce se zabývá kompletním přepisem nahrávek i přepisem melodie.

5.1 Architektura CREPE

První sada experimentů se zakládá na architektuře popsané v článku od Kim a kol. (2018) použité pro sledování jednohlasu. Jak blíže popisujeme v kapitole Související práce, cílem sledování jednohlasu je určit konturu základní frekvence melodického nástroje v jednohlasé nahrávce. Tato nahrávka se zpravidla skládá ze směsi čistého signálu hlasu a šumu v pozadí. Pokud však rozšíříme pojem šumu v pozadí tak, aby zahrnoval i melodický doprovod, pak dostáváme polyfonní signál, tedy vstupní signál pro metody extrakce melodie.

Jinými slovy je sledování jednohlasu speciálním případem extrakce melodie a tudíž přinejmenším stojí za zkoušku pokusit se tuto architekturu pro extrakci využít. Mimo to jednohlasé stopy často obsahují přeslech ostatních nástrojů, pokud nahrávka vznikala při společném hraní ve studiu, tudíž by model trénovaný na vícehlasých mixech mohl být robustní vůči tomuto druhu rušení.

Architektura CREPE se skládá ze šesti konvolučních a pooling vrstev, pro regularizaci používá batch normalization (Ioffe a Szegedy (2015)) a dropout (Sut-



Obrázek 5.1: Diagram architektury CREPE, multiplikační koeficient 16x.

skever a kol. (2014)) po každé konvoluční vrstvě, jako nelineární aktivace je použita funkce ReLU. Po konvolucích následuje výstupní plně propojená vrstva, jako finální aktivační funkce je použita sigmoida. Vstupem modelu je okno o velikosti 1024 vzorků jednonákanalového audio signálu, převzorkovaného na 16 kHz. Před první konvolucí je signál normalizován tak, aby každé jednotlivé vstupní okno mělo střední hodnotu 0 a směrodatnou odchylku 1 — aplikováním normalizace o signálu neztratíme pro nás důležité informace, protože výška tónu nezávisí na absolutním posunu signálu nebo na jeho amplitudě. Naopak tento krok zrychluje trénování, protože síť se nemusí učit být vůči těmto rozdílům v datech invariantní. Celková struktura a podrobnější popis modelu je naznačen na obrázku.

Výsledný vektor o 640 složkách aproximuje pravděpodobnostní rozdělení výšky základní frekvence uprostřed vstupního okna, přičemž tento vektor pokrývá rozsah od noty C_{-1} po G_9 , mezi dvěma sousedními predikovanými výškami je tudíž vzdálenost 20 centů. Výšky tónů v centech označíme $\zeta_1, \zeta_2, \dots, \zeta_{640}$. Rozsah výstupního vektoru tedy bezpečně pokrývá obvyklé rozsahy hlasů hudebních nástrojů a na jednu notu připadá 5 složek (tónů) výsledného vektoru.

$$\zeta(f) = 1200 \log_2 \frac{f}{f_{\text{ref}}}$$

Pro trénování modelu potřebujeme také cílové diskrétní pravděpodobnostní rozdělení základní frekvence tónu. Jako cílovou pravděpodobnostní funkci použijeme normální rozdělení se střední hodnotou v bodě cílové základní frekvence $\zeta(f_{\text{ref}})$ a se směrodatnou odchylkou 25 centů. Toto rozdělení diskretizujeme tak, aby měl cílový vektor stejné dimenze jako odhadovaný.

$$y_i = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\zeta_i - \zeta_{\text{ref}})^2}{2\sigma^2}\right)$$

Převod z pravděpodobnostní reprezentace výstupního vektoru na konkrétní hodnotu výšky noty provedeme pomocí výpočtu střední hodnoty výstupní distribuce. Jelikož by při výpočtu střední hodnoty ale hodnotu výsledné výšky tónu ovlivňoval i doprovod, který se na výstupním vektoru objevuje, počítáme střední hodnotu pouze z okolí maxima výstupu. Tím zajistíme, že získáme střední hodnotu gaussiánu náležícímu pouze jednomu tónu.

$$\hat{c} = \sum_{i, |\hat{c}_i - \hat{c}_m| < 50} \hat{y}_i \hat{c}_i / \sum_{i, |\hat{c}_i - \hat{c}_m| < 50} \hat{y}_i, m = \operatorname{argmax}_i(\hat{y}_i)$$

Optimalizovaná ztrátová funkce modelu (loss funkce) $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$ se počítá jako cross-entropy mezi vektorem cílových pravděpodobností \mathbf{y} a výstupním vektorem $\hat{\mathbf{y}}$.

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{i=1}^{640} (-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i))$$

Optimalizace probíhá pomocí algoritmu Adam (Kingma a Ba, 2014) s parametrem learning rate 0.0002.

	1.	2.	3.	4.	5.	6.	Celk. parametrů
CREPE 4x	128	16	16	16	32	64	558 240
CREPE 8x	256	32	32	32	64	128	177 1200
CREPE 16x	512	64	64	64	128	256	6 163 200

Tabulka 5.1: Počty filtrů v konvolučních vrstvách v architektuře CREPE v závislosti na multiplikačním koeficientu.

V následujících experimentech se zaměříme nejprve na replikaci výsledků pro sledování jednohlasu, následně v základním nastavení model spustíme i pro melodická data. Poté prozkoumáme nejvhodnější nastavení cílové reprezentace referenční anotace melodie, na základě které se síť učí. Tyto experimenty budou užitečné následně i pro další testované architektury, jelikož způsob reprezentace melodie využíváme ve všech experimentech práce stejný. V předposledním experimentu prozkoumáváme vliv zvětšení kontextu, který má síť k dispozici pro predikci výšky tónu. V závěru se pokusíme vyřešit autory zmíněný problém architektury, který se týká odhadu tónů, které mají vysokou frekvenci.

5.1.1 Replikace výsledků CREPE

Abychom ověřili správnost implementace architektury CREPE pro sledování jednohlasu, spustíme model na syntetických, jednohlasých datech MDB-stem-synth, která byla zveřejněna spolu s článkem od Salamon a kol. (2017).

Na rozdíl od článku Kim a kol. (2018), ve kterém autoři používají pro celkové vyhodnocení architektury pětinasobnou křížovou validaci, jsme použili pouze jednu trénovací a testovací množinu. Zásadní rozdíly mezi implementacemi modelu jsme na základě článku a veřejně dostupného kódu neobjevili.

Po jedné epoše trénování model dosáhl na testovací množině 98.6% přesnosti odhadu výšky. Kim a kol. (2018) uvádí přesnost modelu 97%. V jejich případě jde o průměrný výsledek pěti nezávislých běhů trénování a testování na různě rozdělených datových množinách. Rozdíl mezi dosaženými přesnostmi tedy přičítáme odlišné evaluační strategii. Přehled výsledků je uveden v tabulce 5.2¹

¹Při replikaci experimentu jsme narazili na důležitost správného promíchání dat. Framework Tensorflow použitý pro trénování promíchává data vždy pomocí bufferu pevné velikosti pro dvojice vstupů a cílových výstupů. V praxi je však potřeba buď nastavit buffer na velikost větší

Metrika	Práh	Průměrná hodnota	Hodnota Kim a kol. (2018)
RCA	50 centů	0.988	0.970
RPA	50 centů	0.986	0.967
RPA	25 centů	0.975	0.953
RPA	10 centů	0.937	0.909

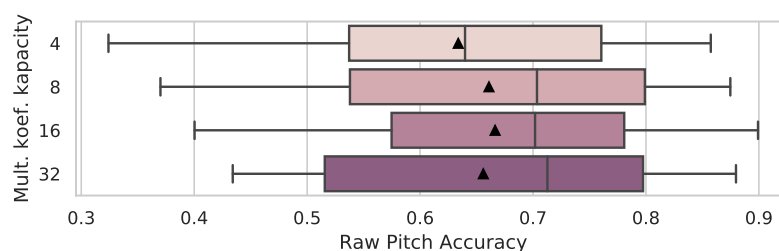
Tabulka 5.2: Výsledky pokusu o replikaci architektury CREPE pro sledování jednohlasu. Přesnosti nejsou přímo srovnatelné kvůli různým evaluačním strategiím.

5.1.2 CREPE pro extrakci melodie

Jako první experiment extrakce melodie z polyfonních dat spustíme nezměněnou architekturu CREPE, v následujících experimentech se tuto baseline pokusíme překonat. Abychom urychlili trénování následujících experimentů, přesnost určíme pro sítě s různou kapacitou. Pokud se výsledky při různých kapacitách nebudou příliš lišit, můžeme experimenty provádět s architekturou s nižší kapacitou a tím snížit trénovací čas. Kapacity upravíme pomocí multiplikačního koeficientu počtu filtrů u všech konvolučních vrstev, počty filtrů jsou uvedeny v tabulce 5.1.

Mult. koef. kapacity	RPA	RCA
4	0.634	0.753
8	0.661	0.766
16	0.666	0.771
32	0.656	0.753

Tabulka 5.3: Výsledky experimentu s nezměněnou architekturou CREPE spuštěnou pro extrakci melodie. Přesnosti uvádíme pro různé kapacity modelu.



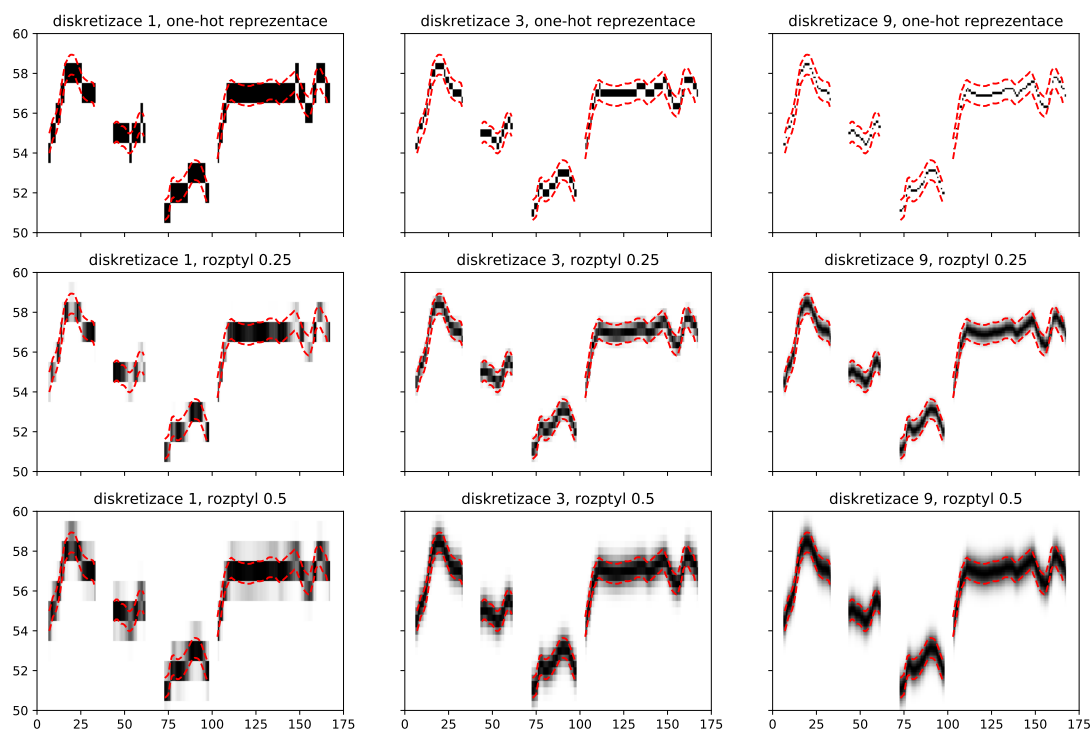
Obrázek 5.2: Výsledky experimentu s nezměněnou architekturou CREPE spuštěnou pro extrakci melodie.

V tabulce 5.3 a na obrázku 5.2 uvádíme výsledky tohoto experimentu. Z validačních výsledků po 200 000 iteracích (přibližně 6 epoch) vidíme, že se výsledek modelů CREPE 8x a CREPE 16x liší řádově o desetiny procentních bodů. Přitom

než je celková velikost datasetu, a nebo implementovat vlastní míchání přes všechna dostupná data. Při nedostatečně promíchaných datech totiž trénovací dávka (batch) není reprezentativní pro celý dataset, ale pouze pro jeho podmnožinu, což se negativně projevuje kolísající validační přesností modelu.

model s větší kapacitou se trénuje o 35% delší dobu. Proto pro většinu následujících srovnávání zvolíme architektury s multiplikačním koeficientem 8, modely s dobrými výsledky případně přetrénujeme s vyšší kapacitou.

5.1.3 Vliv rozlišení diskretizace výšky noty

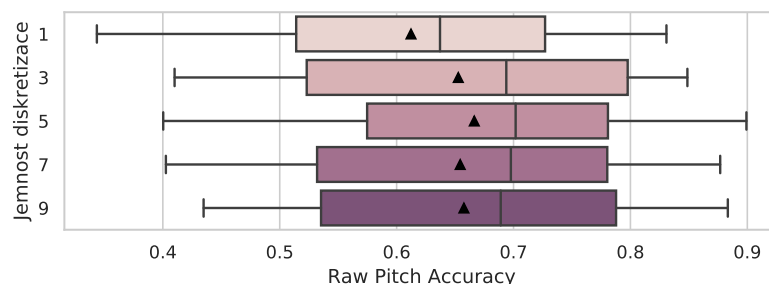


Obrázek 5.3: Ukázky cílové reprezentace použité pro učení modelů. Pružovaná linka značí dolní a horní hranici korektního odhadu výšky melodie.

Otestujeme nastavení granularity výstupního vektoru. V článku Kim a kol. (2018) se totiž důvod volby pěti frekvencí na půltón nediskutuje. Intuitivně by však mělo vyšší rozlišení spíše pomáhat, důvodem je, že nástroje a zejména lidský hlas se často při hraní odchylojí od přesných, definovaných frekvencí hraných not a vyšší rozlišení tyto odchylky může lépe zachytit. Ve výsledku by pak síť s jemnějším výstupem měla dělat méně chyb, u kterých se skutečná a výstupní hodnota liší o jeden půltón. Na obrázku 5.3 pro lepší představu uvádíme příklad cílových reprezentací s různým nastavením rozlišení diskretizace.

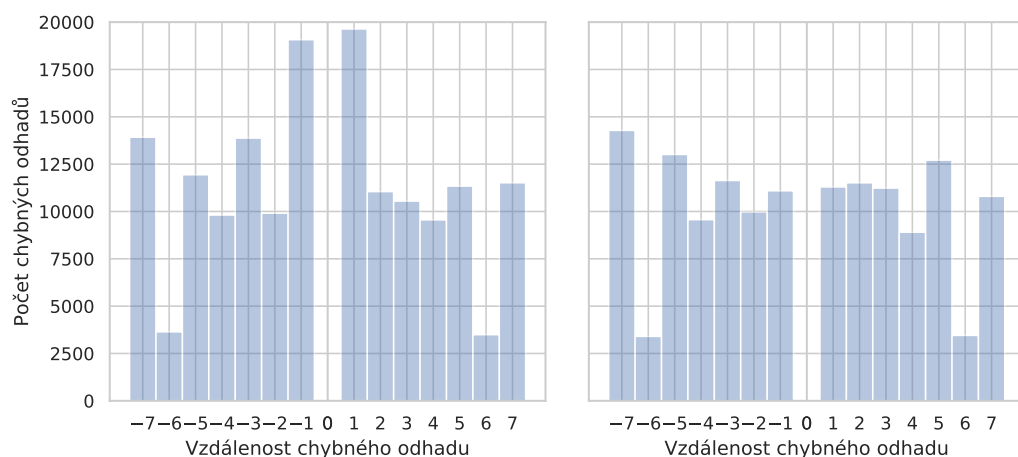
Jemnost diskretizace	RPA	RCA
1	0.612	0.711
3	0.653	0.760
5	0.666	0.771
7	0.654	0.763
9	0.658	0.760

Tabulka 5.4: Architektura CREPE s různou jemností diskretizace.



Obrázek 5.4: Architektura CREPE s různou jemností diskretizace.

Jak můžeme pozorovat na výsledných hodnotách, jemná granularita výstupu zlepšuje přesnost sítě. Abychom ověřili domněnku, že vyšší rozlišení pomáhá zmenšit počet chyb o půltón, vytvoříme histogram vzdáleností cílového a odhadovaného tónu. V tomto histogramu by pak měl být zřetelný pokles v příslušných třídách. Podle histogramu se počet chyb o půltón mezi zkoumanými modely liší téměř o polovinu, zlepšení tohoto druhu chyb je tedy podstatné.



Obrázek 5.5: Histogramy vzdálenosti chybného odhadu, výstup prvního modelu má rozlišení 50 centů, výstup druhého 10 centů.

5.1.4 Vliv rozptylu cílové pravděpodobnostní distribuce výšky noty

Podle Bittner a kol. (2017) pomáhá cílová distribuce s vyšším rozptylem snížit penalizaci sítě za téměř korektní odhady výšek tónů. Mimo to u dostupných dat často nejsou anotace naprosto perfektní, jisté rozostření hranice anotace tudíž pomáhá i v případě nepřesné cílové anotace, síť pak není tolik penalizována za svou případnou správnou odpověď.

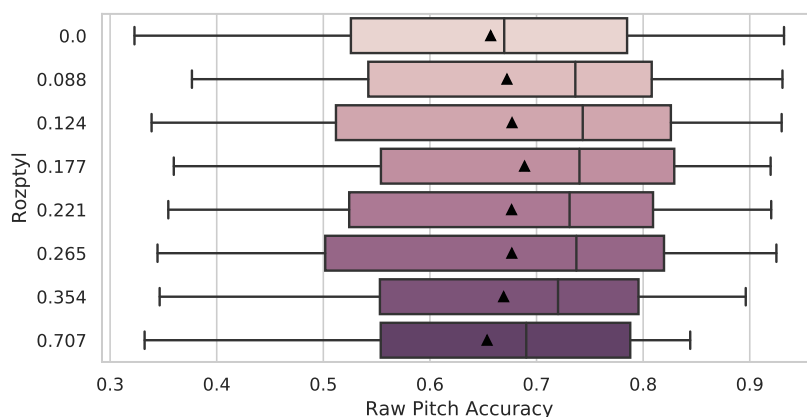
V článku se však nediskutuje, proč bylo zvoleno nastavení směrodatné odchylky na 20 centů. Kim a kol. (2018) používá odchylku 25 centů a není na první pohled zřejmé, jaká je optimální hodnota. Příliš vysoký rozptyl způsobí, že síť bude tolerovat více chyb o půltón, příliš nízký rozptyl naopak penalizuje i za téměř správné odhady. Intuitivně se nejlepší nastavení pravděpodobně bude po-

hybovat kolem používaných 25 centů, jelikož to je hranice chybné klasifikace. Na obrázku 5.3 pro lepší představu uvádíme příklad cílových reprezentací s různým nastavením rozptylu.

Testovaná síť má vstupní okno široké 4096 vzorků, používá multiplikátor kapacity 16x a vstup zpracovává 6 různě širokými konvolučními vrstvami (viz experiment 5.1.6).

Rozptyl	RPA	RCA
0.000	0.657	0.759
0.088	0.672	0.775
0.124	0.677	0.773
0.177	0.689	0.784
0.221	0.677	0.771
0.265	0.677	0.770
0.354	0.669	0.773
0.707	0.654	0.757

Tabulka 5.5: Architektura CREPE, vliv rozptylu cílové distribuce.



Obrázek 5.6: Architektura CREPE, vliv rozptylu cílové distribuce.

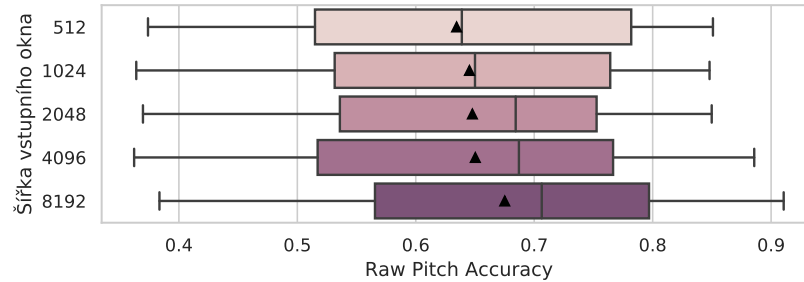
Z experimentů vyplývá, že optimální směrodatná odchylka se pohybuje kolem hodnoty 0.177, tedy níže než v porovnávaných pracích. Spolu s výsledky ze sekce Vliv rozlišení diskretizace výšky noty, ve které ověřujeme nejlepší nastavení frekvenčního rozlišení, tedy docházíme k optimálnímu nastavení cílové reprezentace melodie, které budeme používat i v následujících experimentech se zbylými architekturami.

5.1.5 Vliv šířky vstupního okna

Architektura CREPE byla navržena pro monopitch tracking, dá se předpokládat, že jelikož je v monofonních nahrávkách oproti polyfonním daleko méně (melodického) šumu, není pro určení výšky tónu potřeba větší kontext než použitých 1024 vzorků (při vzorkovací frekvenci 16kHz toto odpovídá 64 milisekundám audia). To ale nemusí platit pro složitější signály, kde by síť mohla z delšího kontextu těžit. Otestujeme tedy vliv většího vstupního okna na výslednou přesnost.

Šířka vstupního okna	RPA	RCA
512 (32 ms)	0.634	0.748
1024 (64 ms)	0.645	0.763
2048 (128 ms)	0.648	0.760
4096 (256 ms)	0.650	0.762
8192 (512 ms)	0.675	0.775

Tabulka 5.6: Architektura CREPE, vliv šířky vstupního okna.



Obrázek 5.7: Architektura CREPE, vliv rozptylu cílové distribuce.

Na obrázku 5.6 vidíme výsledky experimentu, mezi přesnostmi modelů s šířkou okna 1024 až 4096 vzorků nevidíme příliš veliké zlepšení, výsledek okna 8192 je však skokově lepší. Tento výrazný rozdíl spíše přičítáme jistě variabilitě výsledků při trénování neuronových sítí, není totiž zřejmé, jaká vlastnost vstupních dat by mohla způsobit tento skokový rozdíl.

5.1.6 Vliv násobného rozlišení první konvoluční vrstvy

Podle Kim a kol. (2018) se přesnost CREPE snižuje s výškou tónu. Autoři si tuto skutečnost vysvětlují neschopností modelu generalizovat na barvy a výšky tónů neobsažených v trénovací množině, generalizaci by ale mohla pomoci také úprava modelu. Protože k rozpoznání vyšších frekvencí stačí méně vzorků než pro rozpoznání nižších, mohli bychom se pokusit upravit první konvoluční vrstvu sítě, která tento úkol rozpoznávání frekvencí ve vstupu zastává, a rozdělit ji na množinu různě širokých konvolucí, jejichž kanály následně sloučíme zpět do jednotné vrstvy. To by mohlo mít za následek, že rozpoznávání vysokých tónů budou zastávat užší konvoluce a jejich kernel bude jednodušší a obecnější, než když tuto funkci zastávají zbytečně široké kernely, u kterých je možné, že ve svých vahách obsahují redundantní informace.

První vrstvu s kernelem s 256 filtry (tj. počet filtrů první vrstvy s multiplikátorem 8x, viz první experiment) jsme rozdělili na více různě širokých kernelů s menším počtem filtrů, tak aby kapacita sítě zůstala přibližně stejná a sítě byly porovnatelné. Experiment jsme provedli na síti se vstupním oknem 2048 vzorků a multiplikátorem kapacity 8.

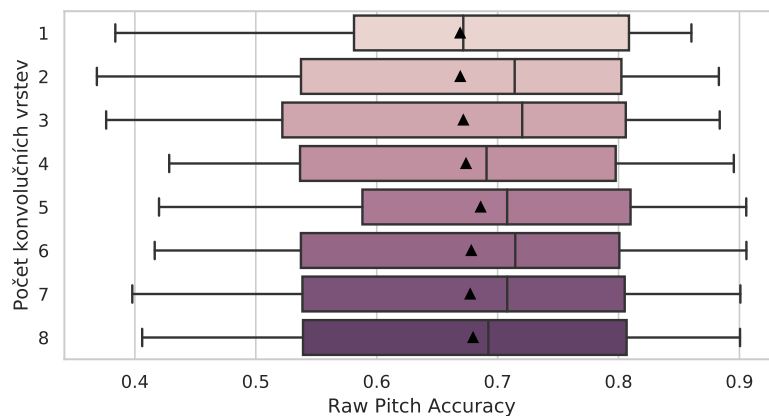
Zlepšení výsledků se pohybuje v řádu desetin procentních bodů, nejvíce patrné je v případě pěti různě širokých konvolučních vrstev, kde dosahuje 1.3 procentního bodu. Analýzou výsledků přesnosti podle výšky noty se mi nepodařilo prokázat domněnku, že by konvoluce s více rozlišeními pomáhala u odhadu not vyšších

Šířka vrstev Počet vrstev	512	256	128	64	32	16	8	4	Počet parametrů
1	256								2098880
2	128	128							2066112
3	85	85	85						2041918
4	64	64	64	64					2029248
5	51	51	51	51	51				2016350
6	42	42	42	42	42	42			2001944
7	36	36	36	36	36	36	36		1996184
8	32	32	32	32	32	32	32	32	2000448

Tabulka 5.7: Počet filtrů prvních vrstev multirezoluční vstupní konvoluční vrstvy v architektuře CREPE.

Počet konvolučních vrstev	RPA	RCA
1	0.669	0.779
2	0.669	0.773
3	0.672	0.773
4	0.674	0.778
5	0.686	0.781
6	0.678	0.780
7	0.677	0.779
8	0.680	0.778

Tabulka 5.8: Architektura CREPE, vliv multirezoluční vstupní konvoluční vrstvy.



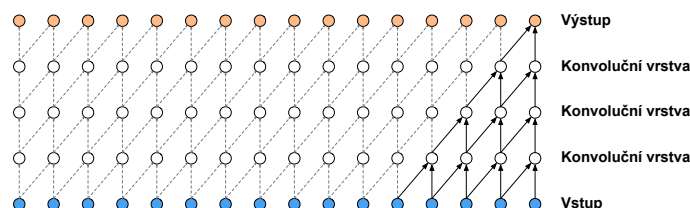
Obrázek 5.8: Architektura CREPE, vliv multirezoluční vstupní konvoluční vrstvy.

frekvencí. Její přínos je malý a projevuje se na většině frekvenčních pásem.

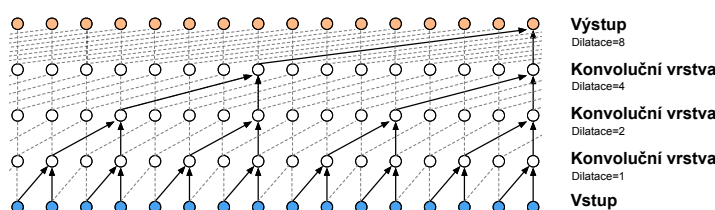
5.2 Architektura WaveNet

Generativní model WaveNet popsaný týmem van den Oord a kol. (2016a) je architektura navržená pro generování zvukového signálu. Autoři však v článku zmiňují, že se architekturu pokusili využít i pro převod mluvené řeči na text (da-

taset TIMIT), a podařilo se jim dosáhnout výsledků srovnatelných se state-of-the-art. Architektura spočívá ve vrstvení dilatovaných konvolucí s rozšiřujícím se rozsahem. Díky exponenciálně rostoucím dilatacím se také exponenciálně zvětšuje receptivní pole jednotlivých konvolučních vrstev. Díky této vlastnosti pak například stačí pro pokrytí 1024 vzorků vstupu pouze 9 vrstev s šířkou kernelu 2 a dilatacemi 1,2,4,8 ... 512. Pokud bychom stejného receptivního pole chtěli dosáhnout pomocí obvyklých konvolucí počet potřebných vrstev by byl lineární vzhledem k šířce pole. Vrstvení konvolucí je porovnáno na obrázcích 5.9 a 5.10. Síť tedy velmi snadno pokryje široký kontext, což je vlastnost, která je pro zpracování zvukového signálu užitečná.



Obrázek 5.9: Vrstvení obyčejných konvolucí s lineárně rozšiřovaným dosahem, obrázek převzat z van den Oord a kol. (2016a).



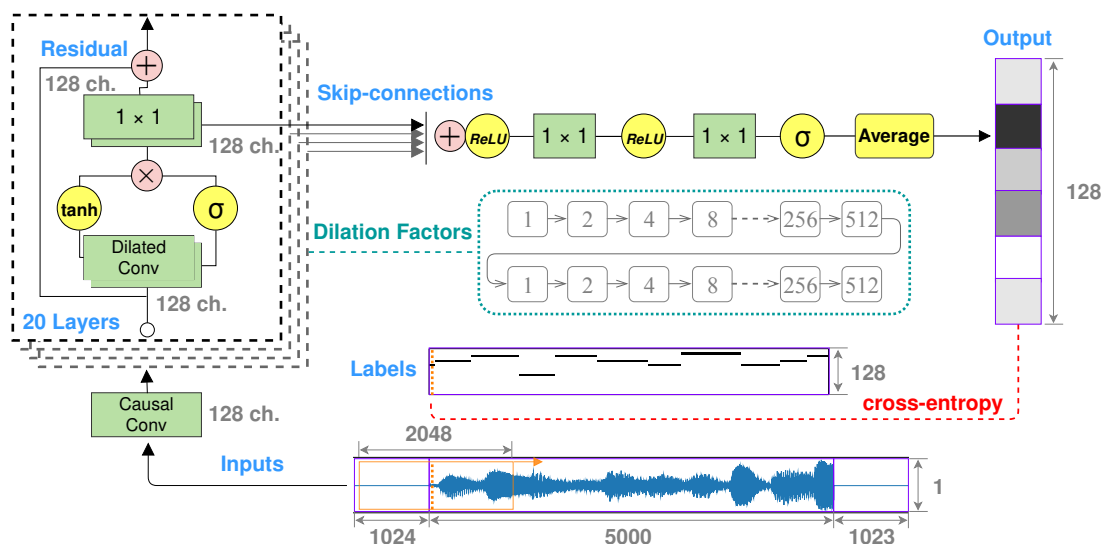
Obrázek 5.10: Vrstvení dilatovaných konvolucí s exponenciálně rozšiřovaným dosahem, obrázek převzat z van den Oord a kol. (2016a).

Síť se pro Music Information Retrieval úlohy od svého zveřejnění příliš neuchytila. Její použití se v oblasti hudby se omezuje na generativní úlohy (Hawthorne a kol. (2018), Yang a kol. (2017), Engel a kol. (2017) a další), případně pro source-separation (Stoller a kol., 2018). Jediný publikovaný pokus s použitím architektury WaveNet pro příbuznou úlohu kompletního automatického přepisu podnikli Martak a kol. (2018) s použitím datasetu MusicNet. Jejich model však netestovali na standardních evaluačních datasetech ze soutěže MIREX, tudíž není zřejmé, jakých výsledků v porovnání s existujícími metodami autoři dosáhli.

V dalších sekcích se pokusíme nejprve spustit model Martak a kol. (2018) na melodická data. Následně se pokoušíme tento výchozí výsledek zlepšit pomocí úpravy architektury WaveNet. Upravujeme celkovou kapacitu sítě pomocí nastavení počtu filtrů v dilatačních blocích. Dále pak měníme šířku kontextu, který je síti pro odhad jednoho bodu výšky tónu k dispozici, pomocí nastavení počtu dilatačních vrstev, bloků a velikosti šířky kernelů dilatací. Prozkoumáme také předzpracování vstupu do dilatačních bloků a způsob zpracování jejich výstupu.

5.2.1 Baseline na základě Martak a kol. (2018)

Pro extrakci melodie využijeme jako výchozí model upravenou architekturu od Martak a kol. (2018), jejíž struktura je naznačena na obrázku 5.11. Vstupem



Obrázek 5.11: Architektura WaveNet upravená pro kompletní přepis skladeb, upraveno na základě Martak a kol. (2018).

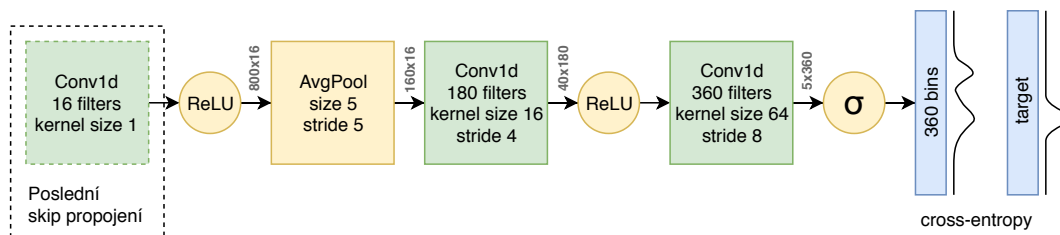
je okno velikosti 4894 vzorků audia převzorkovaného na 16 kHz, toto okno nejprve zpracujeme standardní konvoluční vrstvou s šířkou kernelu 2 a 128 výstupními kanály, takto zpracovaný signál dále prochází dvěma bloky po desíti vrstvách, které obsahují dilatované konvoluce. Vrstvy obsahují dvě dilatované konvoluce, jednu s aktivací hyperbolickým tangens a jednu s aktivací funkcí sigmoid. Výstupy těchto konvolucí jsou po složkách vynásobeny, díky čemuž konvoluce s aktivací sigmoid funguje jako nastavitelná propust signálu (Gated activation unit, van den Oord a kol. (2016b)), poté je výstup opět zpracován dvěma konvolucemi, tentokrát se šířkou kernelu 1x1. Výstup první sečteme s původním vstupem celé vrstvy (jedná se o residuální propojení poprvé popsané v práci He a kol. (2015)) a zpracujeme ho nadcházející vrstvou. Výstup druhé konvoluce (autoři tuto cestu nazývají *skip propojení*) zařadíme mezi výstupy skip propojení všech ostatních vrstev.

Pro výpočet odhadů výšky tónu sečteme všechna skip propojení, tento součet zpracujeme dvěma konvolucemi, druhá z nich má na výstupu 128 kanálů zpracovaných sigmoidou, což odpovídá rozsahu odhadovaných not. Jelikož všechny popsané transformace zachovávají šířku vstupu, dostáváme po této konvoluci odhady výšek not pro každý vstupní vzorek zvuku, taková anotace je zbytečně podrobná, tudíž výsledné anotace podvzorkujeme pomocí pooling vrstvy počítající průměr hodnot.

Po 20 000 iteracích s velikostí dávky 20 a parametrem learning rate 0.001 síť dosahuje na validačních datech přesnost odhadu tónu 0.583 (RPA) a přesnost odhadu tónu nezávisle na oktávě 0.692 (RCA). V porovnání s výsledky architektury CREPE jsou tyto výsledky výrazně nižší, krom toho učení sítě trvá velmi dlouho (14 minut na 1000 iterací), zejména kvůli topologické komplexitě modelu. Síť také v této konfiguraci vykazuje známky přeučení. Pokusíme se tedy zrychlit trénování a zabránit přeučení použitím pouze jednoho dilatačního bloku a snížením počtu filtrů dilatačních a skip propojení ze 128 na 16. Také snížíme velikost trénovací dávky na 8 pro další zrychlení učení. Nová síť obsahuje 10 vrstev s dilatacemi (1, 2, 4, 8, ..., 512), celkový počet trénovatelných parametrů se z 2 009 728 snížil

na 34 736. Tato síť dosahuje RPA 0.598 a RCA 0.696 po 100 000 iteracích při rychlosti trénování 30 sekund na 1000 iterací. Síť tedy dosahuje stejných výsledků ve výrazně kratším čase.

Z předchozích experimentů na architektuře CREPE víme, že jemnější výstupní reprezentace pomáhá snížit počet chyb o půltón, upravíme tedy síť tak, aby na jeden půltón připadalo 5 výstupních složek vektoru, zmenšíme výstupní frekvenční rozsah a hodnoty cílového vektoru změníme z ostré predikce konkrétního tónu na „rozmlžený“ gaussian se směrodatnou odchylkou 18 centů (pro podrobnější popis viz 5.1). Síť po této úpravě dosahuje RPA 0.629 a RCA 0.731 po 100 000 iteracích.



Obrázek 5.12: Úprava posledních vrstev WaveNet architektury.

Předběžným hledáním výchozí architektury pro nadcházející sadu experimentů jsme došli k následujícím úpravám. Dilatované konvoluce ve všech vrstvách mají šířku 3 místo 2, jejich vstup tedy závisí na všech okolních vzorcích, nikoli jen na předchozích, jak je naznačeno na obrázku 5.10. Dále jsme odstranili konvoluční vrstvu, která zpracovává vstup, ten je tedy přímo zpracován dilatačním blokem. Ze skip propojení se zpracovává pouze poslední, nedochází ke sčítání všech, což je úprava, kterou pro přepis řeči používají původní autoři článku WaveNet van den Oord a kol. (2016a). Tento výstup je následně zpracován pooling vrstvou a dvěma konvolucemi se skoky (stride), viz obrázek 5.12. Tato síť dosahuje RPA 0.655 a RCA 0.759 po 100 000 iteracích a je základem pro následující experimenty.

5.2.2 Vliv počtu filtrů dilatačních vrstev a skip propojení

Jedním ze zásadních faktorů ovlivňujících kapacitu sítě je počet filtrů dilatačních vrstev a skip propojení. Tyto kanály nesou informaci o vzorku a jeho okolí, v závislosti na vrstvě dilatace. Výstup poslední vrstvy má tedy délku počtu vstupních vzorků 2848 a 16 kanálů. Každý výstupní vzorek nese informaci o svém okolí délky 2047. Podobně jako v případě architektury CREPE nalezneme vhodnou kapacitu sítě tak, aby nedocházelo k přeučení (overfitting) ani k nedoučení (underfitting) na trénovací množině.

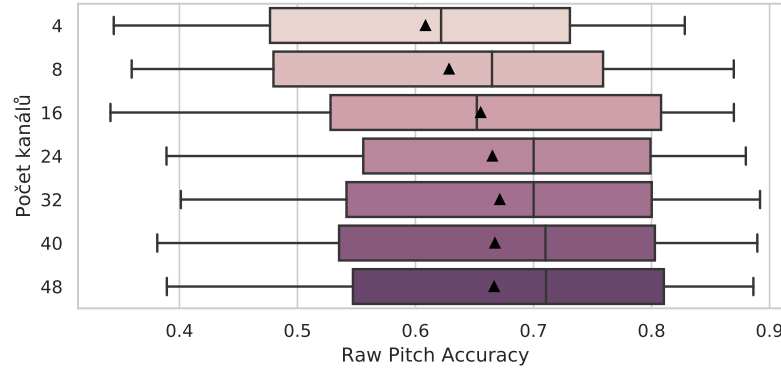
Na základě výsledků volíme pro další experimenty nastavení 16 filtrů pro dilatované konvoluce a skip propojení. Při nastavení 24 a více filtrů dosažená přesnost sítě stagnuje, 16 filtrů je kompromisem z hlediska rychlosti trénování.

5.2.3 Systematické prohledávání počtu dilatačních vrstev a bloků

Velikost zpracovávaného kontextu lze v architektuře WaveNet ovlivnit třemi různými hyperparametry. Jde o počet vrstev v bloku n_{layers} , počet dilatačních

Počet kanálů	RPA	RCA
4	0.609	0.714
8	0.628	0.739
16	0.655	0.759
24	0.665	0.764
32	0.671	0.771
40	0.667	0.766
48	0.667	0.764

Tabulka 5.9: Architektura WaveNet, vliv počtu filtrů dilatačních vrstev a skip propojení.



Obrázek 5.13: Architektura WaveNet, vliv počtu filtrů dilatačních vrstev a skip propojení.

bloků poskládaných nad sebou n_{stacks} a šířka kernelu dilatací n_{width} . Přesně lze dosah vypočítat jako

$$\text{receptive_field} = (n_{\text{width}} - 1) \cdot \left(\sum_d^{n_{\text{layers}}} 2^{(d-1)} \right) \cdot n_{\text{stacks}} + 1$$

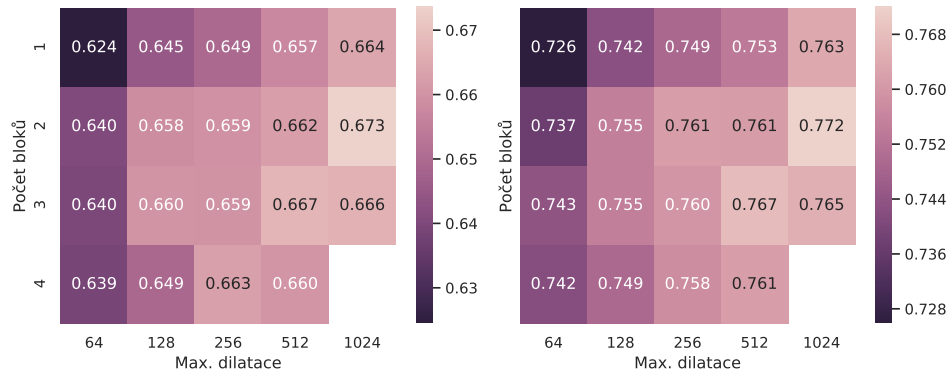
Dosud jsme testovali síť s jedním blokem o deseti vrstvách s dilatacemi $(1, 2, 4, 8, 16, \dots, 512)$ s šířkou konvolucí 3, její dosah je tedy 2047 vzorků. Systematickým prohledáváním prozkoumáme vliv delšího kontextu měněného pomocí n_{layers} a n_{stacks} . Při vyhodnocování experimentu je nutné vzít v úvahu, že přidání nové vrstvy, případně celého nového bloku, zároveň zvyšuje kapacitu modelu, což také ovlivňuje výslednou přesnost modelu. Šířky kontextu a kapacity jsou uvedené v tabulce 5.10.

V architektuře modelu jsme pro tento experiment provedli změnu ve zpracování skip propojení. V této sadě se nebere poslední skip propojení, nýbrž všechna se spojí v ose kanálů a slouží jako vstup pro následující pooling vrstvu a konvoluci. Z toho důvodu má model mnohem více parametrů, na druhou stranu je zajištěno, že poslední vrstvy mohou využít informace ze všech skip spojení k odhadu výšek. Na základě srovnání tréninkové a validační ztrátové funkce však modely nevykazují známky přeučení.

Na základě obrázku 5.14 lze dojít k pozorování, že přidávání počtu vrstev zvyšuje výslednou přesnost víceméně vždy, to neplatí o počtu bloků, u kterých

Max. dilatace	64	128	256	512	1024
Počet bloků					
1 (dosah)	255	511	1023	2047	4095
1 (kapacita)	4 483 644	4 531 836	4 580 028	4 628 220	4 676 412
2 (dosah)	509	1021	2045	4093	8189
2 (kapacita)	4 820 988	4 917 372	5 013 756	5 110 140	5 206 524
3 (dosah)	763	1531	3067	6139	12 283
3 (kapacita)	5 158 332	5 302 908	5 447 484	5 592 060	5 736 636
4 (dosah)	1017	2041	4089	8185	—
4 (kapacita)	5 495 676	5 688 444	5 881 212	6 073 980	—

Tabulka 5.10: Architektura WaveNet, dosah a kapacita v závislosti na dilatačních počtu vrstev a bloků.



Obrázek 5.14: Architektura WaveNet, systematické prohledávání počtu dilatačních vrstev a bloků, vlevo hodnoty RPA, vpravo RCA.

se zdá, že nejvhodnější počet je 2. Pokud se zaměříme na srovnání výsledků s podobným dosahem, zdá se že výsledky jsou poměrně podobné, zejména pro širší dosahy, pro kratší se zdá lepší využít spíše více vrstev než více bloků.

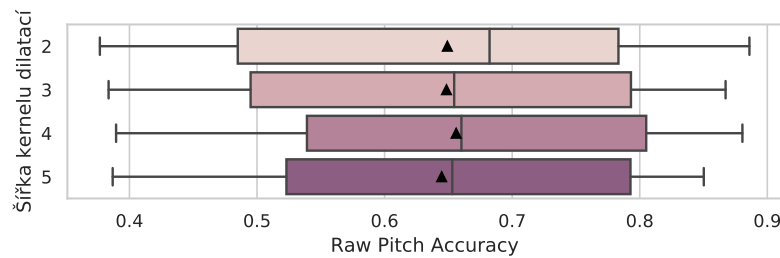
5.2.4 Vliv velikosti šířky kernelu dilatací

Jiným způsobem, jak zvýšit velikost zpracovávaného kontextu, je volba velikosti šířky kernelu dilatací n_{width} . Provedeme čtyři experimenty s různou šířkou. Šířka 2 je zvolena v původním článku z toho důvodu, že konvoluce se tím stávají „kauzální“ — jejich výstup závisí pouze na vzorcích před nimi (viz obrázky 5.9, 5.10). To je výhodné pro generativní úlohy, u kterých chceme, aby hodnota nového generovaného vzorku v audio signálu nezávisela na budoucích, zatím neexistujících vzorcích. U klasifikačních úloh takto omezení nejsme, tudíž můžeme s nastavením experimentovat.

Z tabulky 5.11 a grafu 5.15 vyplývá, že tato síť při změnách šířky kernelu dilatace svou výslednou přesnost na validačních datech nemění.

Šířka kernelu dilatací	RPA	RCA
2	0.649	0.754
3	0.648	0.755
4	0.656	0.759
5	0.645	0.756

Tabulka 5.11: Architektura WaveNet, vliv velikosti šířky kernelu dilatací.



Obrázek 5.15: Architektura WaveNet, vliv velikosti šířky kernelu dilatací.

5.2.5 Vliv výstupní transformace skip propojení

Volba transformace skip propojení se v článku týmu van den Oord a kol. (2016a) nediskutuje, Martak a kol. (2018) také přejímají součet všech skip výstupů. Vyzkoušíme proto dvě další možnosti práce se skip propojeními - výběr posledního skip propojení a dále spojení všech skip propojení do mnohakanálového výstupu. V případě výběru posledního propojení či součtu jde o transformace, které lze vyjádřit jako speciální případ konvoluce nad spojenými skip propojeními. Dalo by se tedy říct, že výběr a součet jsou v tomto případě speciálními případy spojených výstupů. Testovaná možnost výběr poslední vrstvy skip propojení naopak přináší jinou výhodu - pro výpočet výsledku sítě nejsou potřeba zbylá skip propojení, což urychluje trénování.

Transf. skip vrstev	RPA	RCA
Konkatenace	0.654	0.754
Poslední	0.651	0.754
Součet	0.646	0.753

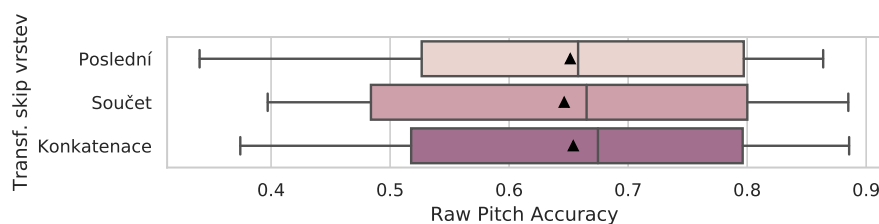
Tabulka 5.12: Architektura WaveNet, vliv výstupní transformace skip propojení.

Všechny tři transformace vedou ke stejné přesnosti. Síť tedy netěží z větších možností práce s dřívejšími vrstvami a nejdůležitější informace se objevují až na vrchu bloků dilatovaných konvolucí.

5.2.6 Vliv velikosti první konvoluce

Tým van den Oord a kol. (2016a) používá pro předzpracování zvukového signálu před vstupem do dilatovaných bloků obvyklou konvoluci, nespecifikuje však její šířku. Veřejná implementace architektury WaveNet pracuje s šířkou 32 ².

²https://github.com/ibab/tensorflow-wavenet/blob/master/wavenet_params.json

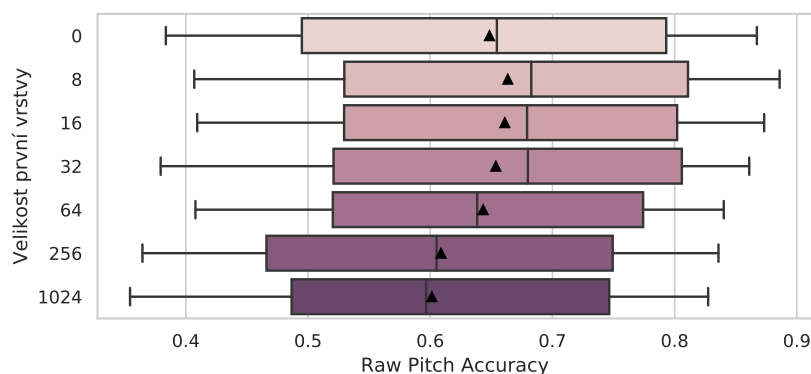


Obrázek 5.16: Architektura WaveNet, vliv výstupní transformace skip propojení.

Martak a kol. (2018) používají šířku 2, čímž fakticky jen duplikují první dilatační vrstvu. První vrstva může sloužit jako banka filtrů, podobně jako první vrstva v architektuře CREPE. Aby však pro tento účel mohla fungovat plnohodnotně, musela by být široká minimálně 512 vzorků, aby mohla zachytit periodu i té nejnižší frekvence ve výstupním rozsahu. Pomocí sady experimentů se pokusíme zjistit, zda má první konvoluční vrstva pozitivní vliv na výslednou přesnost.

Velikost první vrstvy	RPA	RCA
0	0.648	0.755
8	0.663	0.762
16	0.661	0.765
32	0.654	0.754
64	0.643	0.756
256	0.609	0.735
1024	0.601	0.734

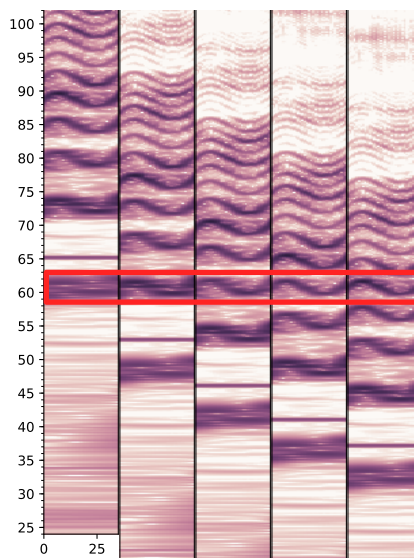
Tabulka 5.13: Architektura WaveNet, vliv velikosti první konvoluce.



Obrázek 5.17: Architektura WaveNet, vliv velikosti první konvoluce.

Sítě, které mají možnost využít první konvoluce jako banky filtrů, dosahují výrazně horší přesnosti, než sítě s menší konvolucí. Přeskočení tohoto předzpracování však také není nejlepším nastavením pro tuto architekturu. Vhodná šířka první konvoluce se pohybuje kolem 8 vstupních vzorků, do takové šířky se vejde jedna perioda signálu frekvence 2000 Hz, jinými slovy vrstva jistě nemůže zastávat ani z části funkci banky filtrů. Mohla by však popisovat sklon křivky vstupních dat, ze kterých se v dalších vrstvách dilatovaných bloků složí již celý frekvenční popis vstupu.

5.3 Architektura HCNN



Obrázek 5.18: Znázornění harmonických závislostí na spektrogramu s logaritmickou osou frekvence.

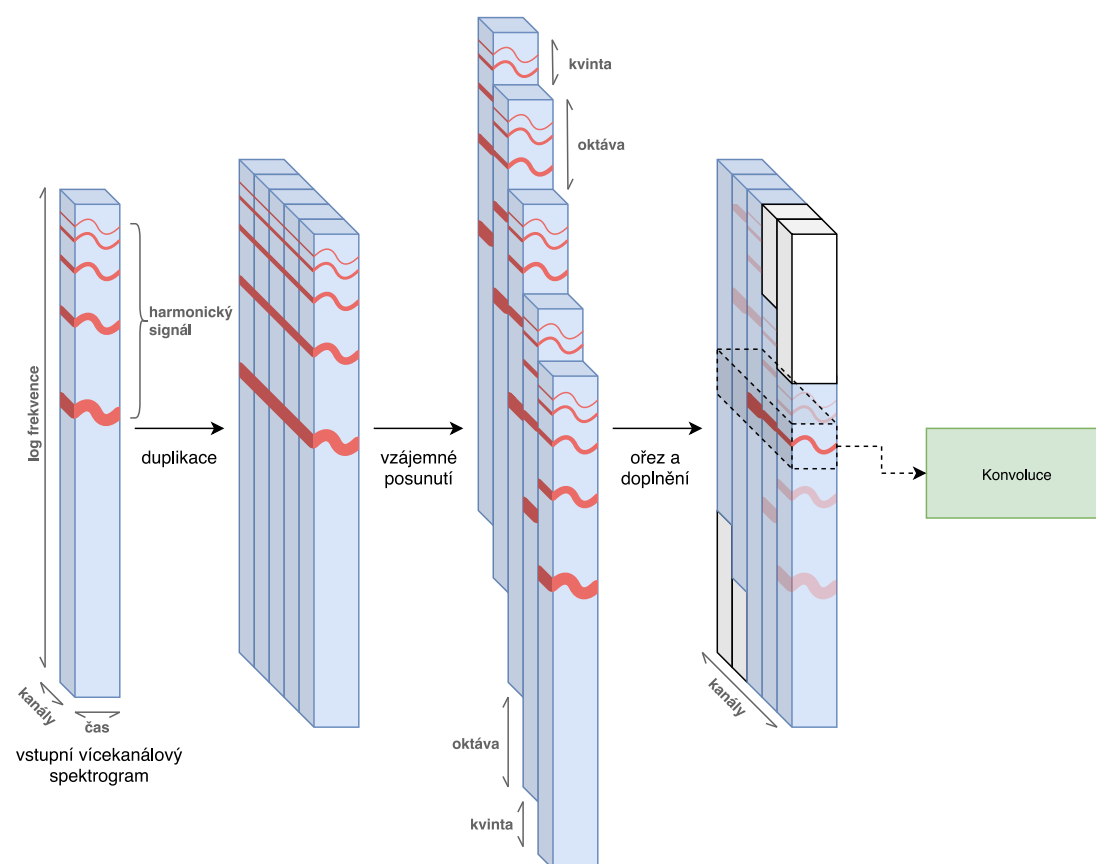
V práci prezentujeme také novou architekturu pro zpracování harmonicky strukturovaných zvukových signálů kterou nazýváme *Harmonic Convolutional Neural Network*. Inspirujeme se prací Bittner a kol. (2017), která problém výpočtu salienční funkce zasazuje do rámce úlohy odstranění šumu v obrázku pomocí konvolučních sítí. Její architektura vstupní spektrogram zpracovává několika konvolučními vrstvami, výsledkem pak je nový spektrogram s odstraněným „šumem hudebního doprovodu“.

Dalším přínosem její práce je popis vstupní spektrální reprezentace HCQT. CQT spektrogramy mají logaritmickou osu frekvence, tato reprezentace má tu vlastnost, že dvě libovolné harmonické frekvence jednoho tónu mají na spektrogramu mezi sebou konstantní vzdálenost, nezávislou na fundamentální frekvenci. HCQT spočívá ve výpočtu několika CQT spektrogramů, jejichž frekvenční rozsahy jsou vzájemně posunuty právě o tyto konstantní vzdálenosti. Pokud pak tyto spektrogramy poskládáme za sebe v nové ose kolmé na osu frekvence a času, harmonické frekvence se překryjí (viz obrázek 5.18). Tento vícekanálový spektrogram pak může sloužit jako vstup pro konvoluční vrstvu, která má při výpočtu jedné frekvenční složky k dispozici také informace o souvisejících harmonických složkách.

Z této myšlenky úpravy vstupní reprezentace pro rozšiřování receptivního pole konvoluce na související harmonické složky vychází naše architektura HCNN. Sít je koncipována jako obvyklé vrstvení konvolučních vrstev, každému konvolučnímu bloku však předchází úprava vstupu takovým způsobem, aby do výpočtu nadcházející konvoluce byly zahrnuty i potenciální harmonické složky dané frekvence. Naše metoda tedy nepracuje pouze s výpočtem různých CQT spektrogramů na vstupu jako práce Bittner a kol. (2017), upravujeme totiž mezivýsledky všech vrstev v síti a naše transformace spočívá v pouhém posunutí, nikoli přepočítávání. Na obrázku 5.19 tuto úpravu vstupu znázorňujeme; výstup předchozí konvoluční vrstvy duplikujeme a vzájemně posuneme o pevně dané počty půltónů tak, aby

harmonicky související složky byly umístěny „nad sebou“ v dimenzi kanálů. Následně výstupy v této dimenzi spojíme a chybějící složky doplníme nulami.

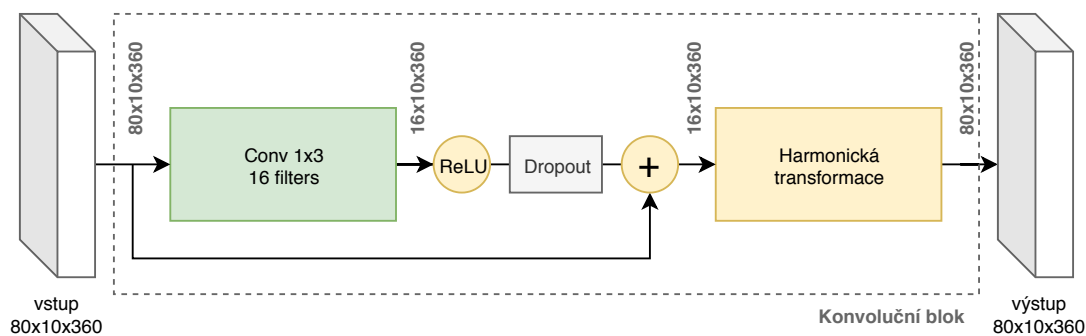
Tato úprava přináší dvě výhody. Jednak síť může využít relevantní harmonické informace v každé vrstvě a ne pouze na vstupu, což vede k lepší přesnosti modelů. Také ale můžeme díky harmonické transformaci velmi výrazně zmenšit samotné konvoluční vrstvy. Zatímco konvoluční vrstvy v architektuře Bittner a kol. (2017) používají velké kernely (zejména pak ve frekvenční ose) a na výstupu vrací velké množství kanálů, aby byly schopná zachytit vzdálené frekvenční závislosti, díky struktuře HCNN takto velké konvoluční vrstvy nejsou potřeba. Ve výsledku naše testované architektury obsahují řádově desetkrát méně parametrů a trénují se v rámci desítek až stovek minut.



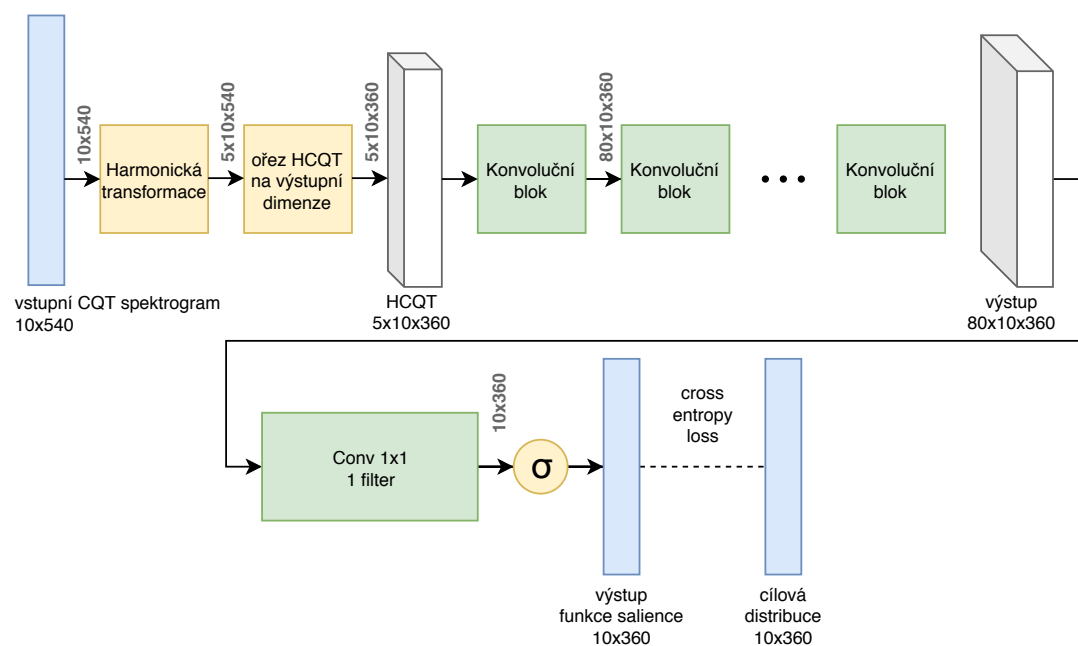
Obrázek 5.19: Diagram transformace vstupu konvoluční vrstvy pro zachycení harmonických souvislostí.

Zbývající prvky architektury již čerpají ze zavedených postupů. Diagram 5.20 znázorňuje strukturu jednoho konvolučního bloku; po každé konvoluční vrstvě s aktivací funkcí ReLU následuje dropout vrstva a kolem konvoluce vedeme residuální propojení pro umožnění trénování hlubokých architektur. Celková architektura (obrázek 5.21) pak spočívá v harmonické transformaci vstupní CQT (čímž efektivně vytvoříme reprezentaci odpovídající HCQT bez přepočítávání každého CQT zvlášť), aplikaci několika konvolučních bloků a aplikaci konvoluce 1x1 pro transformaci výstupu zpět na jeden dvoudimenzionální výstup.

Vstupní CQT počítáme z původních audio souborů se vzorkovací frekvencí 44 100 Hz, používáme implementaci algoritmu z balíku `librosa` a parametry transformace jsou: počáteční frekvence (f_{min}) = 32.7 (Tón C1), počet složek na



Obrázek 5.20: Diagram konvolučního bloku architektury HCNN.



Obrázek 5.21: Diagram celkového propojení architektury HCNN.

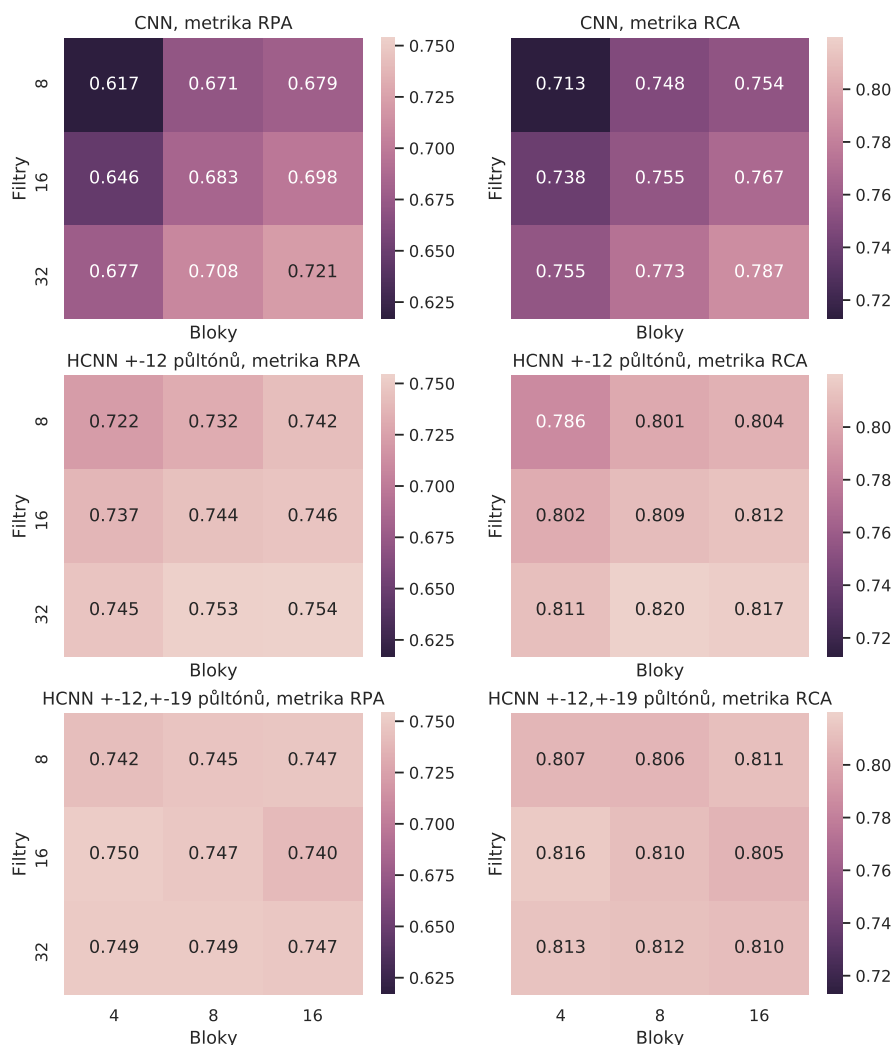
oktávu (`bins_per_octave`) = 60 (vychází 5 složek na půltón, jako v předchozích sekcích), počet složek (`n_bins`) = 540 (vychází na 9 oktáv, po první harmonické transformaci však spektrum ořízneme na výšku 360, zbylé složky jsou přítomny, aby zaplnily chybějící hodnoty při harmonických posunech), `hop_size=256` (tudíž jedno anotační okno MedleyDB odpovídá jednomu sloupci vstupního spektru). Dále pak hodnoty vstupu převedeme na logaritmickou škálu dBFS (funkce `librosa.core.amplitude_to_db`).

Cílový výstup počítáme v souladu s implementacemi v předchozích sekcích, tedy použijeme normální rozdělení se střední hodnotou rovnou výšce znějící melodie a rozptylem 18 centů.

V následujících experimentech nejprve prokážeme užitečnost aplikace harmonické transformace uvnitř sítí a nalezneme vhodnou hloubku a kapacitu modelů, dále prezentujeme jednu z limitací vstupní reprezentace práce Bittner a kol. (2017) a pokusíme se navrhnout vlastní rozšíření vstupní reprezentace. V závěru se pokusíme rozšířit síť zpracováváný kontext a zběžně otestujeme i možnost použití augmentace vstupních spektrů.

5.3.1 Harmonické transformace

Abychom prokázali pozitivní vliv harmonické transformace, natrénujeme tři sady modelů. Jednu sadu bez jakýchkoliv transformací, sadu s transformací s posunutím o ± 12 půltónů (zachycení druhé harmonické složky) a třetí sadu s posunutím o ± 12 a o ± 19 půltónů (zachycení druhé a třetí harmonické složky). Vstupem všech sítí bude HCQT (implementovaný efektivně pomocí harmonických posunů), aby byl posouzen vliv transformací uvnitř sítě, nikoli také na vstupu.

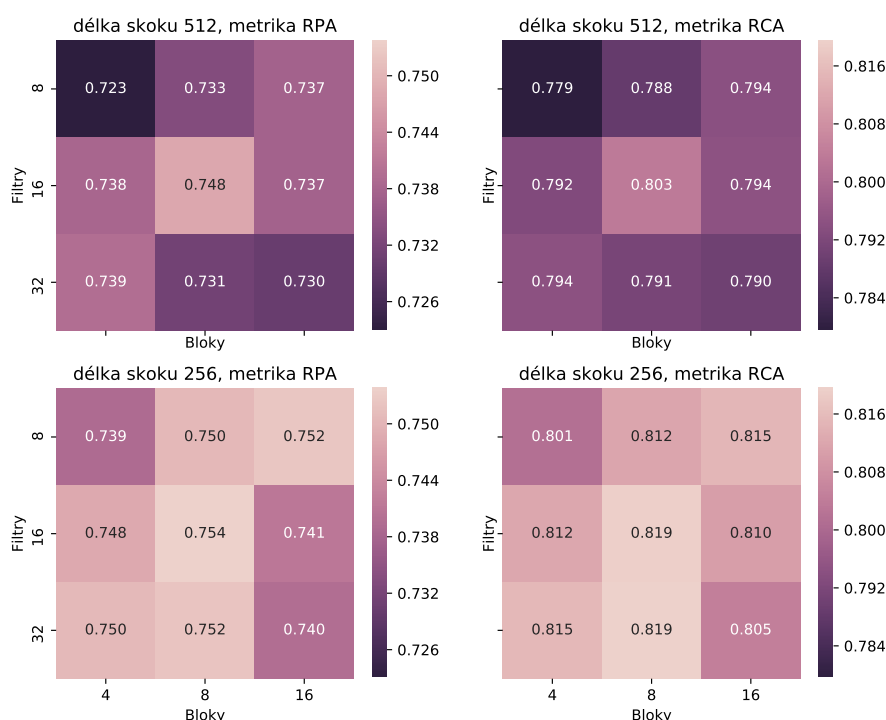


Obrázek 5.22: Architektura HCN, Detekce melodie, vliv počáteční pooling vrstvy.

Na základě porovnání všech kombinací počtů filtrů a bloků na obrázku 5.22 je jisté, že použití harmonické transformace uvnitř konvolučních sítí pomáhá zvyšovat přesnost. U všech testovaných šířek a hloubek sítí došlo k několikaprocentnímu nárůstu přesnosti přepisu. Ukazuje se, že k tomuto nárůstu došlo přidáním pouze dvou posunutí — posunutí o oktávu v ose frekvence nahoru a dolů. Přidáním další zachycené harmonické frekvence sice pomůžeme zvýšit přesnost modelů s menší kapacitou, modelům s vyšší kapacitou již úprava nepomáhá.

5.3.2 Parametr `hop_size`

Metoda Bittner a kol. (2017) pracuje se spektrogramy CQT vypočítanými s velikostí skoku ≈ 11 ms. Data MedleyDB, která v práci používáme, jsou však anotována s velikostí skoku ≈ 5.8 ms. To znamená, že při trénování musíme anotace v trénovací množině MedleyDB podvzorkovat a podobně při vyhodnocování výsledků sítě musíme výstup převzorkovat zpět na původní časové rozlišení datasetu. V obou případech převzorkování můžeme ztratit informace, které se projeví v horší výsledné přesnosti. Je také otázkou, zda-li je nejlepší postup, pokud síť budeme trénovat na neoriginálních, převzorkovaných datech. Na druhou stranu výhodou většího skoku je větší receptivní pole konvoluce, kernel o šířce větší než 1 má na podvzorkovaných datech dvojnásobný dosah. Tento efekt lze nicméně napodobit nastavením dvojnásobné dilatace konvoluce. Na experimentech s kernely šířky 1 posoudíme dopad převzorkování bez této zmiňované výhody.



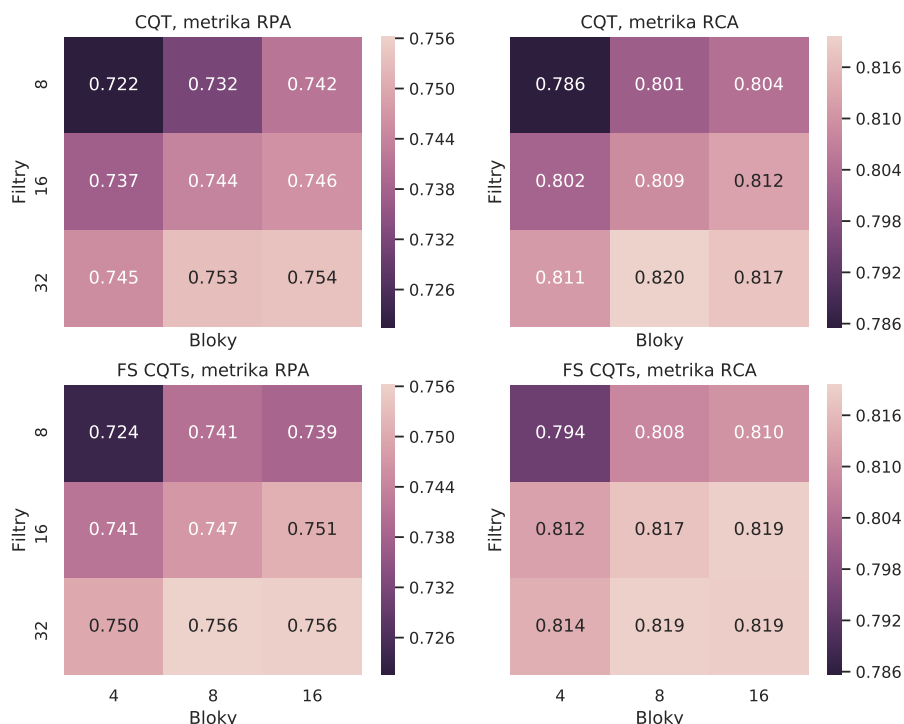
Obrázek 5.23: Architektura HCNN, Vliv velikosti skoku pro výpočet vstupního spektrogramu

Výsledek prezentovaný na obrázku 5.23 je potenciálně využitelný i pro metodu Bittner a kol. (2017). Je možné, že pokud se původní architektura přetrénuje na datech s menším `hop_size`, výsledky sítě se také zvýší, jako v našem případě.

5.3.3 Vícekanálový vstup CQT

Při výpočtu CQT spektrogramu můžeme specifikovat parametr `filter_scale`, což je koeficient ovlivňující velikost okna transformace. S menší velikostí okna se pojí lepší časové rozlišení, s větší naopak lepší frekvenční. Otestujeme proto sítě, které mají na vstupu k dispozici více variant CQT spektrogramu. Další možností úpravy vstupního spektrogramu je druh práhování nastavitelný parametrem `top_db` ve funkci `librosa.core.amplitude_to_db`, prozkoumáme i tato

nastavení. Vstupní vícekanálový spektrogram se proto bude skládat ze tří spektrogramů s dvojicemi parametrů `filter_scale`, `top_db`: (0.5, 60), (1, 80), (2, 100).



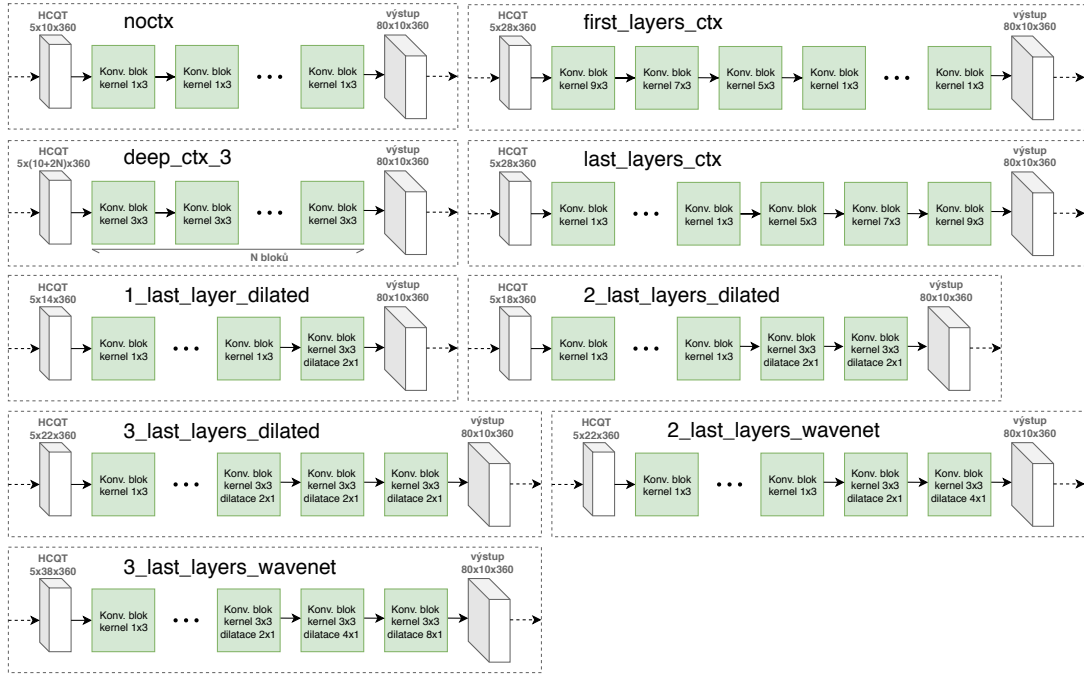
Obrázek 5.24: Architektura HCNN, Vliv vstupního vícekanálového CQT

Výsledné zlepšení se pohybuje v desetinách procentního bodu u většiny testovaných konfigurací. Na druhou stranu vstupní reprezentace je nezávislá na architektuře zbytku modelu, tudíž je zde naděje, že toto zlepšení bude nezávislé na zvolené architektuře. Použití tohoto vstupu by tedy mohlo vylepšit i výsledky nejlepších objevených architektur.

5.3.4 Kontext

Doposud jsme trénovali síť, které pro výpočet salience uvažovaly pouze jedno okno vstupního CQT. Tyto síť pro odhad výšky měly k dispozici pouze okno šířky ≈ 5.8 ms. Následující sadou experimentů se pokusíme prozkoumat možnosti zvětšení tohoto kontextu. Všechny úpravy spočívají ve změnách konvolučních vrstev v blocích architektury (viz obrázek 5.25), konkrétně nastavujeme velikost konvoluce z 1×3 na jiné hodnoty. Také prozkoumáváme možnost využití dilatovaných konvolucí, jejichž princip vysvětlujeme v sekci Architektura WaveNet, které umožňují snadno zvětšovat zpracováváný kontext pomocí menšího počtu vrstev než s použitím obvyklých konvolucí.

Protože v předchozích experimentech dosahujeme vysoké úspěšnosti i bez použití kontextu navíc, ve většině testovaných architekturách zachováváme část této původní architektury a rozšíření receptivního pole zajišťují pouze poslední vrstvy. Testujeme tedy řadu různých změn architektur, pro každou variantu budeme trénovat verzi s celkově čtyřmi konvolučními bloky a s celkově osmi. Také každou variantu přetrénujeme ve verzi s 8 a s 16 konvolučními filtry v každém bloku.

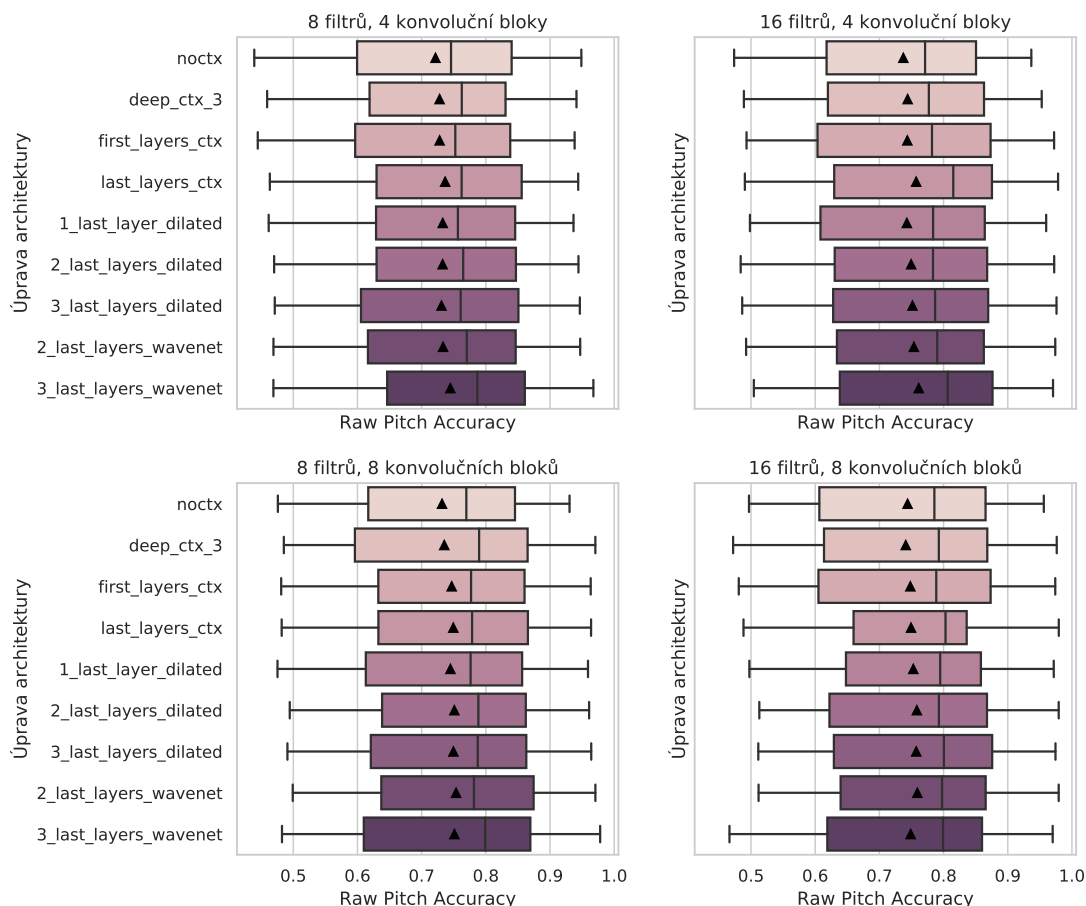


Obrázek 5.25: Úpravy architektury HCNN z obrázku 5.21 pro rozšíření uvažovaného kontextu.

Celkově tedy budeme trénovat 4 varianty od každé architektury. Všechny sítě aplikují harmonickou transformaci ± 12 v každém konvolučním bloku, vstupní reprezentací je HCQT s `hop_size=256`. Testované architektury mají následující struktury:

- `noctx`: síť bez kontextu navíc, celkový kontext pro výpočet jednoho odhadu výšky je 5.8 ms.
- `deep_ctx_3`: síť s konvolucemi velikosti 3x3. Celkový kontext $2N \cdot 5.8$ ms na jeden odhad je tedy dán hloubkou sítě N . Konkrétně tedy kontext vychází 46.4 ms pro verzi se čtyřmi konvolučními bloky a 92.8 ms pro verzi s osmi.
- `first_layers_ctx`: síť s konvolucemi velikosti 9x3, 7x3 a 5x3, které zpracovávají vstupní spektrogram, následované konvolucemi 1x3, které již kontext nerozšiřují. Celkový kontext vychází $(8 + 6 + 4) \cdot 5.8$ ms = 104.4 ms.
- `last_layers_ctx`: síť zpracovávající vstup konvolucemi 1x3, následované řadou vrstev s konvolucemi 9x3, 7x3 a 5x3. Celkový kontext 104.4 ms.
- `1_last_layer_dilated`: síť zpracovávající vstup konvolucemi 1x3, následované jednou konvoluční vrstvou 3x3 s dilatací 2 v dimenzi času. Celkový kontext 23.2 ms.
- `2_last_layers_dilated`: síť zpracovávající vstup konvolucemi 1x3, následované dvěma konvolučními vrstvami 3x3 s dilatací 2 v dimenzi času. Celkový kontext 46.4 ms.
- `3_last_layers_dilated`: síť zpracovávající vstup konvolucemi 1x3, následované třemi konvolučními vrstvami 3x3 s dilatací 2 v dimenzi času. Celkový kontext 69.6 ms.

- `2_last_layers_wavenet`: síť zpracovávající vstup konvolucemi 1x3, následované dvěma konvolučními vrstvami 3x3 s dilatacemi 2 a 4 v dimenzi času. Celkový kontext 69.6 ms.
- `3_last_layers_wavenet`: síť zpracovávající vstup konvolucemi 1x3, následované třemi konvolučními vrstvami 3x3 s dilatacemi 2, 4 a 8 v dimenzi času. Celkový kontext 162.4 ms.



Obrázek 5.26: Architektura HCNNet, Vliv úpravy architektury ovlivňující receptivní pole modelu.

Poznamenáme také, že šířka vstupních spektrogramů v ose času je pro každý experiment rozšířena tak, aby konvoluční bloky měly k dispozici všechny potřebné informace pro výpočet hodnot na výstupním oknu. Například pokud tedy síť pro jeden bod anotace využívá také tři předcházející a tři nadcházející okna spektrogramu, rozšíříme vstupní spektrogram z obou stran o tři okna.

Výsledky uvádíme na obrázku 5.26, kvůli většímu rozměru tabulku s výsledky přesouváme do přílohy 6.5. Prvním důležitým pozorováním je, že přidání kontextu oproti kontrolní síti bez kontextu výsledky sítě zlepšuje téměř ve všech případech. Sítě tedy jsou schopné přidání kontextu využít k lepším odhadům. Dále vidíme, že síť `deep_ctx_3` si vede nejhůře ze všech testovaných, tato síť na rozdíl od ostatních neobsahuje žádné konvoluce 1x3, můžeme se tedy domnívat, že konvoluce 3x3 na datech v této architektuře nefungují natolik dobře, jako jejich užší varianty. Je však také možné, že informace o kontextu je pro síť obtížné využít,

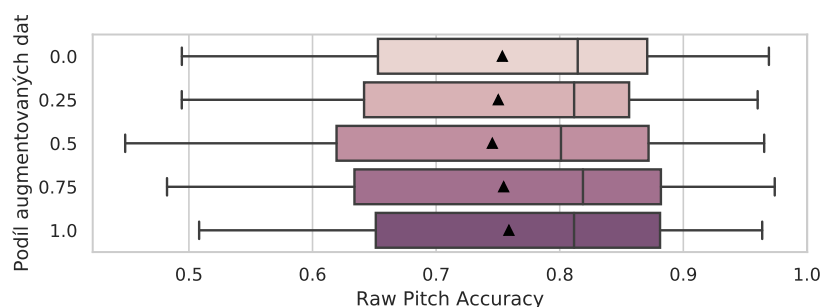
pokud tato informace má projít dlouhou kaskádou úzkých konvolucí, na které se architektura `deep_ctx_3` zakládá.

Nejlepší výsledek podává mělká architektura se čtyřmi bloky, jde o konfiguraci s 16 filtry, 4 konvolučními bloky a architekturou `3_last_layers_wavenet`, tato architektura tedy obsahuje pouze jednu konvoluční vrstvu šířky 1×3 a pak následují dilatované konvoluce 3×3 .

5.3.5 Augmentace vstupních spektrogramů

Jak shrnujeme v kapitole Související práce, metody, které používají hlubokého učení k extrakci melodie, také často prozkoumávají možnosti augmentace vstupních dat, aby snížily míru přeučení svých modelů. Prozkoumáme proto zběžně novou techniku augmentace nazvanou SpecAugment, která byla představena v článku Park a kol. (2019) a je určena pro úpravu vstupních časově-frekvenčních reprezentací. Augmentace je v principu velmi jednoduchá, spočívá ve vynechání náhodně zvolených frekvenčních pásem a časových úseků na vstupu sítě. Díky tomu se trénovaný model stává odolnější vůči neúplným a zkresleným příkladům. Na síť je také následně možné nahlížet i perspektivou generativních modelů, ve které se model snaží v místě chybějících informací predikovat příklad výstupu z distribuce podmíněné vstupním zvukovým okolím. Jinou motivací také může být paralela s lidským vnímáním, lidský sluch je totiž vůči vynechání frekvenčních pásem také odolný, zvuk po aplikaci pásmového filtru sice zní pozměněně, nicméně kvality jako je výška tónu si zachovává. Vzhledem ke stáří tohoto článku představujeme první pokus o aplikaci této augmentace na melodická data. Augmentaci testujeme na nejlepší síti z minulé sekce, tedy na architektuře `3_last_layers_wavenet` s 16 filtry a 4 konvolučními bloky. Navíc síti předkládáme vícekanálový vstup popsáný v sekci Vícekanálový vstup CQT.

Z článku SpecAugment implementujeme vynechávání frekvenčních pásem a časových úseků. Náhodnou část vstupních spektrogramů transformujeme ve frekvenční doméně překrytím dvou frekvenčních pásem šířky až 5.4 půltónu, v časové doméně pak vynecháváme jeden úsek délky až 29 ms.



Obrázek 5.27: Architektura HCNN, vliv použití SpecAugment na vstupní spektrogramy.

Výsledky nalezneme v tabulce 5.14 a na obrázku 5.27. Pro tuto síť se rozdíl mezi výsledky sítě bez augmentace vstupních dat a se vstupem pouze augmentovaných dat pohybuje kolem půl procentního bodu. Vzhledem k výsledkům prezentovaným v článku Park a kol. (2019) se však domníváme, že augmentace má

Podíl augmentovaných dat	RPA	RCA
0.0	0.754	0.818
0.25	0.750	0.824
0.5	0.745	0.820
0.75	0.755	0.828
1.0	0.759	0.824

Tabulka 5.14: Architektura HCNN, vliv použití SpecAugment na vstupní spektrogramy.

vyšší potenciál, zejména pak pro větší sítě s delším kontextem, které mají větší tendence k přeučení.

6. Výsledky

V této kapitole shrnujeme výsledky vybraných architektur z kapitoly Experimenty na testovacích datech a provádíme kvantitativní a kvalitativní srovnání se state-of-the-art systémy pro extrakci melodie, představené v kapitole Související práce.

6.1 Výběr testovaných modelů

Pro srovnání jsme vybrali nejlepší natrénované modely z kapitoly Experimenty od každé testované architektury. Každou vybranou architekturu krátce popíšeme a v tabulce uvedeme nalezené nastavení hyperparametrů.

Všechny testované modely mají společnou reprezentaci cílového výstupu při trénování. Rozlišení diskretizace výšky noty je nastaveno na 5 hodnot na půltón, rozptyl distribuce výšky noty je nastaven na 0.18, na základě experimentů v sekci Architektura CREPE.

6.1.1 Architektura CREPE

Zvolený model dosáhl na validačních datech přesnosti odhadu výšky 0.682 a přesnosti odhadu výšky nezávisle na oktávě 0.783. Počet trénovatelných parametrů tohoto modelu je 2 016 350. Proběhlo 360 000 iterací trénování. Trénování probíhalo přes dvě hodiny. Parametry vybrané architektury uvádíme v tabulce 6.1.

6.1.2 Architektura WaveNet

Zvolený model dosáhl na validačních datech přesnosti odhadu výšky 0.673 a přesnosti odhadu výšky nezávisle na oktávě 0.768. Počet trénovatelných parametrů tohoto modelu je 5 206 524. Proběhlo 360 000 iterací trénování. Trénování probíhalo téměř pět hodin. Parametry vybrané architektury uvádíme v tabulce 6.2.

6.1.3 Architektura HCNN

Pro úspěšnost modelů na validačních datech jsme se rozhodli porovnat dvě různé architektury HCNN. Architekturu s minimálním kontextem 5.8 ms, dále nazývanou HCNN noctx, a architekturu, která uvažuje širší kontext, dále jen HCNN.

HCNN noctx

Zvolený model dosáhl na validačních datech přesnosti odhadu výšky 0.751 a přesnosti odhadu výšky nezávisle na oktávě 0.813. Počet trénovatelných parametrů tohoto modelu je 27 857. Proběhlo 100 000 iterací trénování. Trénování probíhalo 22 minut. Parametry vybrané architektury uvádíme v tabulce 6.3.

Prohledávané parametry		Ostatní parametry	
Parametr	Hodnota	Parametr	Hodnota
Multiplikační koef. kapacity	8	Velikost dávky	32
Šířka vstupního okna	2048	Iterace trénování	360000
Násobné rozlišení první vrstvy	6 vrstev	Learning rate	0.001

Tabulka 6.1: Nastavení architektury a hyperparametrů pro testovanou architekturu CREPE.

Prohledávané parametry		Ostatní parametry	
Parametr	Hodnota	Parametr	Hodnota
Velikost první konvoluce	0	Iterace trénování	100000
Počet filtrů	16	Velikost dávky	8
Počet bloků	2	Learning rate	0.001
Maximální dilatace	1024		
Transformace skip propojení	concat		

Tabulka 6.2: Nastavení architektury a hyperparametrů pro testovanou architekturu WaveNet.

Prohledávané parametry		Ostatní parametry	
Parametr	Hodnota	Parametr	Hodnota
Parametr hop_size	256	Iterace trénování	100000
Kontext	noctx	Velikost dávky	16
Vstupní reprezentace	HCQT	Learning rate	0.001
Filtry, Bloky	16, 8	Dropout	0.3
Harmonická trans.	$\pm 12, \pm 19$		

Tabulka 6.3: Nastavení architektury a hyperparametrů pro testovanou architekturu HCNN noctx.

Prohledávané parametry		Ostatní parametry	
Parametr	Hodnota	Parametr	Hodnota
Parametr hop_size	256	Iterace trénování	100000
Kontext	3_last_layers _wavenet	Velikost dávky	8
Vstupní reprezentace	Vícekan. HCQT	Learning rate	0.001
Filtry, Bloky	16, 4	Dropout	0.3
Harmonická trans.	± 12		
Pst. augmentace	0.75		

Tabulka 6.4: Nastavení architektury a hyperparametrů pro testovanou architekturu HCNN.

HCNN

Zvolený model dosáhl na validačních datech přesnosti odhadu výšky 0.755 a přesnosti odhadu výšky nezávisle na oktávě 0.828. Počet trénovatelných parametrů tohoto modelu je 23 153. Proběhlo 100 000 iterací trénování. Trénování probíhalo 75 minut. Parametry vybrané architektury uvádíme v tabulce 6.4.

6.1.4 Detekce melodie

Pro detekci melodie jsme použili techniku práhování. Pro odhad přítomnosti melodie nalezneme maximální hodnotu funkce salience v daném okamžiku — pokud tato hodnota přesahuje jistý práh, daný časový okamžik se uvažuje jako obsahující melodii. Konkrétní nastavení práhu jsme určili na základě validačních dat pro každou metodu zvlášť. Tato metoda práhování je použita například také v práci Bittner a kol. (2017).

6.2 Kvantitativní srovnání

V tabulkách 6.5, 6.6 a 6.7 prezentujeme výsledky nově představovaných metod v porovnání s velmi silnými baseline metodami pro extrakci melodie. Jak uvádíme v kapitole Související práce, práce Salamon a Gomez (2012) dosahuje spolu s prací Dressler (2009) v průměru nejlepších výsledků v soutěži MIREX. Práce Bittner a kol. (2017) a D Basaran, S Essid (2018) představují metody, které dosahují nejlepších výsledků na prozatím nejrozsáhlejší datasetu MedleyDB.

Pro srovnání metod používáme datasety ADC04, MIREX05train a ORCH-SET, které jsme v práci vyhradili pouze pro testování. Také používáme testovací množiny datasetů MedleyDB a MDB-melody-synth, převzaté z prací Bittner a kol. (2017) a D Basaran, S Essid (2018), pro dataset WJazzD používáme vlastní testovací množinu. Metodika výběru příkladů do testovací množiny WJazzD je popsána v kapitole Datasety, všechny množiny jsou výčtem popsány v elektronické příloze.

Protože se v práci zaměřujeme zejména na odhad výšky, v tabulkách uvádíme standardní metriky celkové přesnosti (OA), přesnosti odhadu výšky (RPA) a přesnosti odhadu výšky nezávisle na oktávě (RCA). Metriky pro evaluaci detekce melodie (úplnost detekce a nesprávné detekce) přikládáme do přílohy 6.6. Připomínáme však, že metrika celkové přesnosti (OA) zahrnuje vyhodnocení odhadu výšky i vyhodnocení detekce melodie. K vyhodnocování používáme knihovnu `mir_eval`, která poskytuje transparentní a standardizovaný způsob výpočtu metrik pro úlohy oboru Music Information Retrieval.

6.2.1 Popis výsledků

Metody HCNN a HCNN noctx překonávají srovnávané algoritmy v metrikách celkové přesnosti (OA), přesnosti odhadu výšky (RPA) a přesnosti odhadu výšky nezávisle na oktávě (RCA) na datasetech ADC04, MIREX05train a WJazzD. Metoda HCNN pak překonává všechny srovnávané přístupy i na datasetu MedleyDB. Na obrázku 6.1 porovnáváme rozdělení dosažených výsledků na datasetu MedleyDB, na uvedeném krabicovém grafu lze navíc porovnat rozptyl výsledků.

Metoda	ADC04	MDB-m-s test	MIREX05 train.	MDB test	ORCH- SET	WJazzD test
Salamon	0.714	0.527	0.715	0.519	0.235	0.667
Bittner	0.716	0.633	0.702	0.611	0.407	0.692
Basaran	0.669	0.689	0.734	0.640	0.483	0.700
CREPE	0.590	0.562	0.652	0.502	0.248	0.671
WaveNet	0.681	0.528	0.649	0.503	0.256	0.648
HCNN noctx	0.737	0.626	0.723	0.635	0.439	0.715
HCNN	0.726	0.661	0.755	0.652	0.459	0.725

Tabulka 6.5: Výsledky celkové přesnosti (Overall Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.

Metoda	ADC04	MDB-m-s test	MIREX05 train.	MDB test	ORCH- SET	WJazzD test
Salamon	0.767	0.514	0.761	0.526	0.281	0.693
Bittner	0.814	0.606	0.807	0.670	0.519	0.774
Basaran	0.793	0.733	0.798	0.706	0.635	0.767
CREPE	0.794	0.550	0.779	0.616	0.408	0.782
WaveNet	0.796	0.528	0.792	0.595	0.345	0.759
HCNN noctx	0.827	0.647	0.833	0.701	0.511	0.805
HCNN	0.841	0.654	0.851	0.715	0.535	0.806

Tabulka 6.6: Výsledky přesnosti odhadu výšky (Raw Pitch Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.

Metoda	ADC04	MDB-m-s test	MIREX05 train.	MDB test	ORCH- SET	WJazzD test
Salamon	0.807	0.639	0.805	0.659	0.568	0.757
Bittner	0.855	0.666	0.824	0.735	0.694	0.785
Basaran	0.820	0.766	0.807	0.757	0.776	0.776
CREPE	0.851	0.617	0.810	0.714	0.607	0.808
WaveNet	0.843	0.597	0.828	0.703	0.564	0.793
HCNN noctx	0.862	0.699	0.845	0.767	0.683	0.821
HCNN	0.880	0.716	0.863	0.781	0.732	0.820

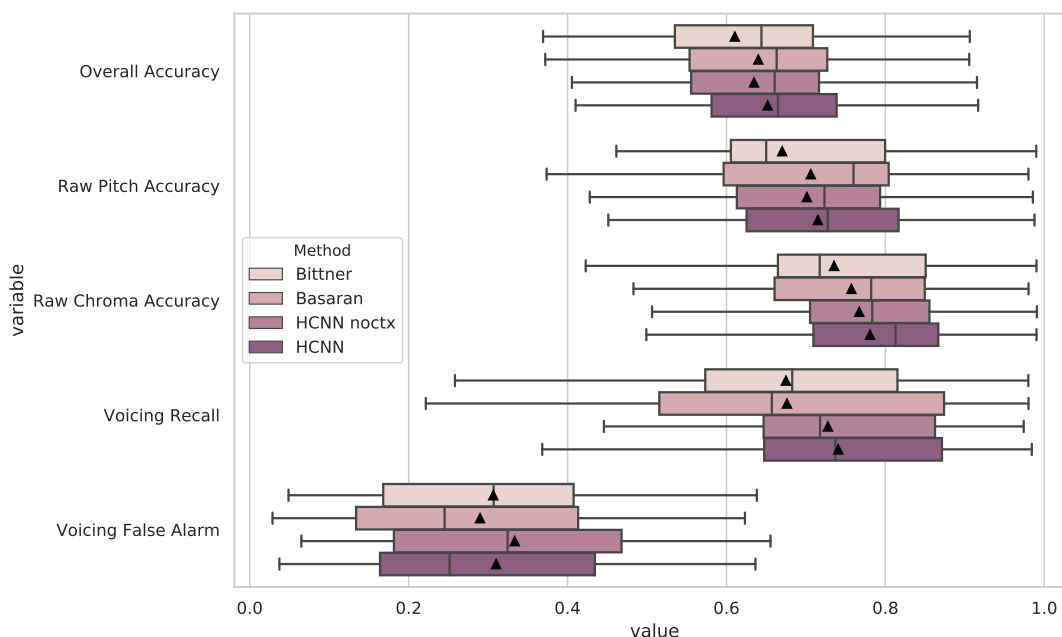
Tabulka 6.7: Výsledky přesnosti odhadu výšky nezávisle na oktávě (Raw Chroma Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.

V metrice RCA metody HCNN dosahují menší variability na sadě testovacích příkladů. Na zbylých datasetech MDB-melody-synth a ORCHSET překonává architektura HCNN ve všech uvažovaných metrikách pouze práce Salamon a Gomez (2012) a Bittner a kol. (2017).

Co se týče zbylých architektur CREPE a WaveNet, v metrikách přesnosti odhadu výšky (RPA) a přesnosti odhadu výšky nezávisle na oktávě (RCA) na téměř všech testovacích datasetech překonávají metodu Salamon a Gomez (2012), která není založena na strojovém učení. Výsledky v porovnání s HCNN, Bittner a kol. (2017) a D Basaran, S Essid (2018) jsou však až na výjimky nižší.

Podle očekávání na základě výsledků ze soutěže MIREX se na datasetu ORCHSET výsledky algoritmů liší nejvíce, v některých případech až o desítky procent. Jak popisujeme v kapitole Datasets, dataset je složen z orchestrálních nahrávek a kvůli vysokému stupni polyfonie a rozmanitým kombinacím barev nástrojů se jedná pro metody extrakce melodie o velmi náročný materiál. Naopak nejblíže, zejména v metrikách RPA a RCA, jsou si výsledky na datasetech ADC04, MIREX05train a WJazzD. Výňatky v těchto datasetech často obsahují velmi zřetelnou melodii a v porovnání s datasetem MedleyDB je jejich hudební obsah žánrově homogenní.

Vzhledem k dosaženým výsledkům, jednoduchosti sítí a rychlému trénování považujeme návrh architektury HCNN jako nadějný pro další rozšiřování. Abychom mohli uvedené výsledky interpretovat, provedeme nejprve také kvalitativní srovnání algoritmů.

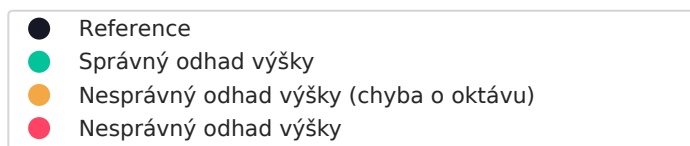


Obrázek 6.1: Výsledky nejúspěšnějších metod na datasetu MedleyDB

6.3 Kvalitativní srovnání

Na základě kvantitativního vyhodnocení vybíráme metody Bittnerové a Basarana pro podrobnější srovnání na jednotlivých příkladech. Z testovaných architek-

tur pak vybíráme obě varianty HCNN. V následujících srovnáních se soustředíme na odhad výšky, proto v obrázcích zobrazujeme pouze části odhadů, ve kterých podle referenční anotace melodie zní. Obrázky jsou bez tohoto zjednodušení příliš komplikované a v práci detekci melodie řešíme pouze okrajově. Metodika výběru kvalitativních příkladů spočívala v hledání skladeb, ve kterých se odhady jednotlivých algoritmů vzájemně nejvíce lišily s nadějí, že právě tyto příklady budou nejlépe ilustrovat limity porovnávaných metod. Vybíráme ale také příklady, které jsou napříč metodami pro odhad melodie obtížné a také ukázkou snadno analyzovatelného vstupu. Legenda barev použitých ve všech následujících obrázcích je vysvětlena na obrázku 6.2 Pro srovnání vybíráme příklady, ve kterých se výstupy sítí nejvíce lišily.



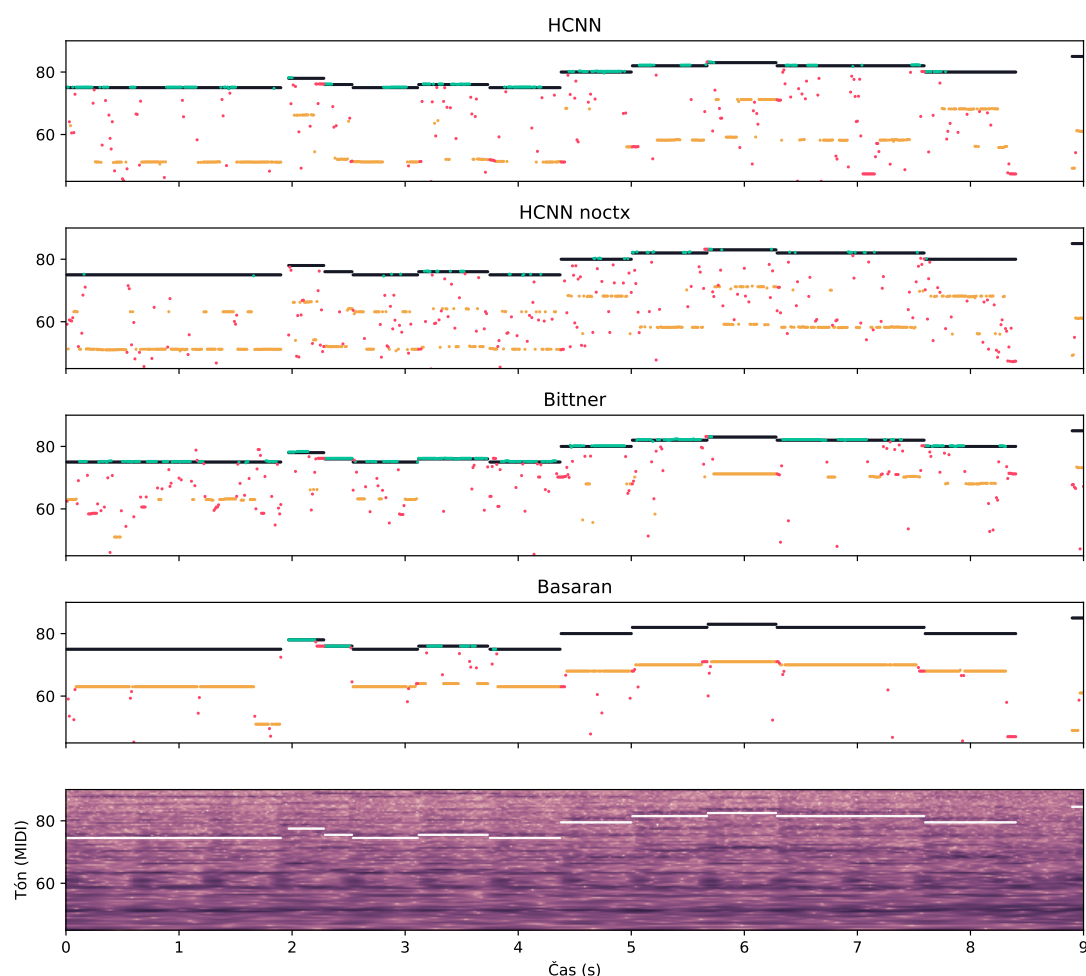
Obrázek 6.2: Legenda pro následující kvalitativní srovnání.



Obrázek 6.3: Příklad s vysokou úspěšností přepisu `mirex05_train01` z datasetu MIREX05train, se kterým je čtenář seznámen z úvodu práce.

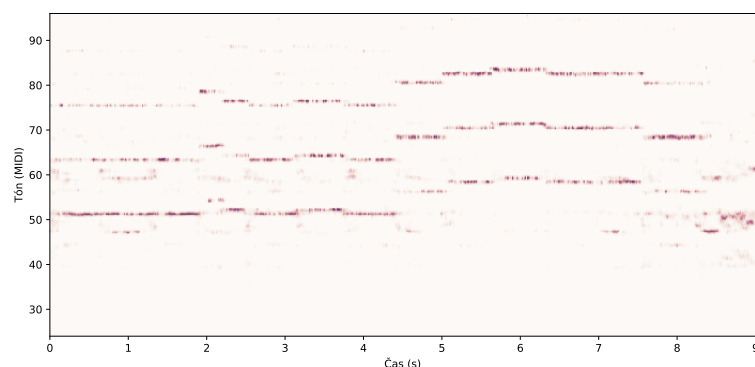
Na obrázku 6.3 můžeme vidět výsledky metod spuštěných na popové nahrávce, ve které melodii nese hlas zpěvačky. Napříč metodami je přepis této nahrávky, kterou uvádíme v úvodu práce, velmi spolehlivý, přesnost odhadu výšky se pohybuje mezi 0.86 a 0.89. Protože v následujících srovnáních ukazujeme zejména chyby přepisu a metody porovnáváme na obtížných příkladech, nechceme, aby

si po přečtení sekce čtenář odnesl, že existující metody přepisu nefungují. Proto uvádíme tento příklad jako pozitivní ukázkou toho, že na obvyklých vstupních datech všechny porovnávané metody fungují velmi dobře.



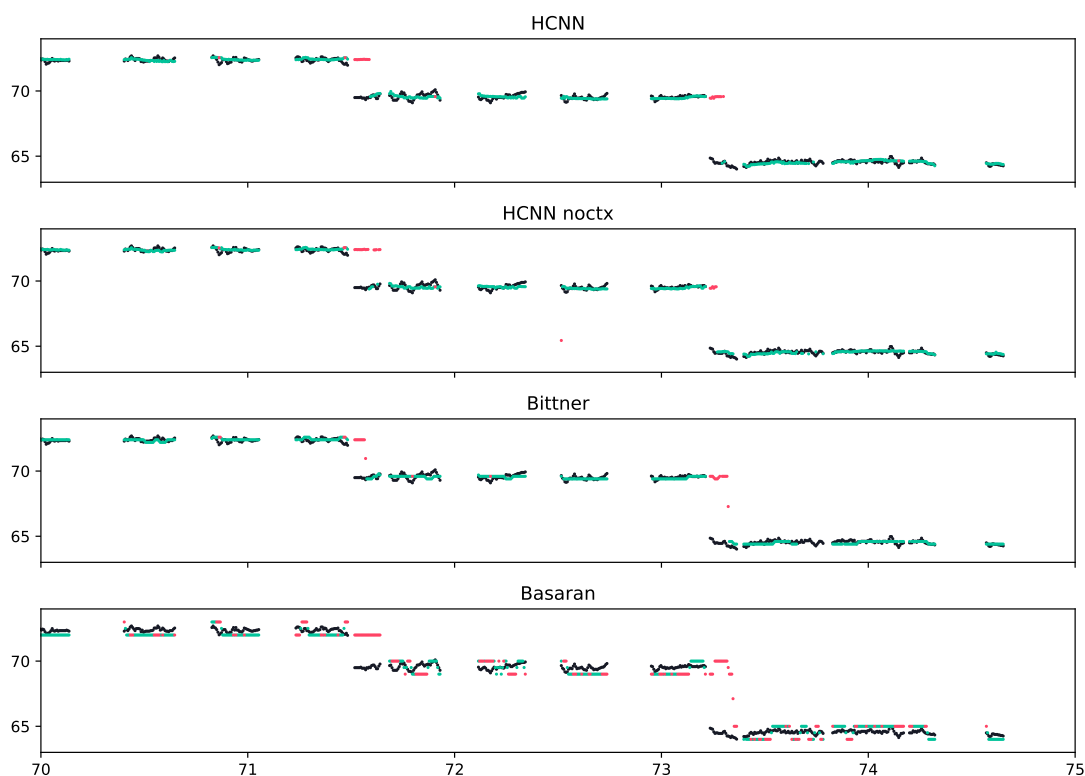
Obrázek 6.4: Výstup metod na testovacím souboru Musorgski-Ravel-PicturesExhibition-ex6 z datasetu ORCHSET.

Největší slabinou sítě HCNN noctx se stala podle očekávání časová kontinuita odhadů. Protože tato síť pro odhad výšek uvažuje vždy pouze 5.8 ms vstupu a po vytvoření funkce salience na odhady tónů neaplikujeme žádné způsoby vyhlazování, odhady jednotlivých časových oken na sebe nenavazují. To se nejvíce jeví jako zásadní problém v případech, kdy ve skladbě melodii nese více hlasů v souzvuku (viz výstup 6.3, 6.7). U některých orchestrálních skladeb však vzniká problém, například pokud melodii nese zároveň sekce smyčců a dechů v různých oktávách. Jak můžeme vidět na výstupu algoritmů 6.4, HCNN noctx pak „přeskakuje“ mezi oktávami. Problém je také dobře vidět na výstupní funkci salience 6.5, na které vidíme tři totožné kontury posunuté o oktávu. Mírné zlepšení tohoto problému vidíme na výstupech metod HCNN a Bittner a kol. (2017), které sice také nijak výsledek salienční funkce nezpracovávají, na druhou stranu pro její výpočet uvažují delší okna délky 162 ms v případě HCNN a 150 ms v případě metody Bittner. Na obrázku 6.4 vidíme, že množství odhadů těchto metod, je často chybný pouze kvůli nesprávně určené oktávě — díky většímu kontextu může me-



Obrázek 6.5: Výstupní salience metody HCNN na testovacím souboru Musorgski-Ravel-PicturesExhibition-ex6 z datasetu ORCHSET.

toda vybrat v čase navazující odhady a proto tyto výstupy obsahují méně velmi krátkých chybných úseků. Metoda týmu D Basaran, S Essid (2018) odhad výšky tónů vyhlazuje pomocí rekurentní architektury GRU, jejich výstup proto obsahuje nejméně skoků, jelikož metoda uvažuje celý kontext skladby, nikoli jen okno omezené délky. Použití rekurentní sítě díky tomu dovoluje zachytit ještě dlouhodobější závislosti a výstupní kontura pak často obsahuje nejmenší množství velmi krátkých, chybných skoků mimo hlavní melodii, které se na obrázku 6.4 hojně vyskytují u metod bez vyhlazování. Na obrázku 6.4 proto vidíme, že metoda D Basaran, S Essid (2018) se drží při odhadu jedné oktávy a přesto, že je tato oktáva zvolena špatně, výsledný přepis je koherentní.



Obrázek 6.6: Detail přepisu metod na testovacím souboru CannonballAdderley_SoWhat z datasetu WJazzD.

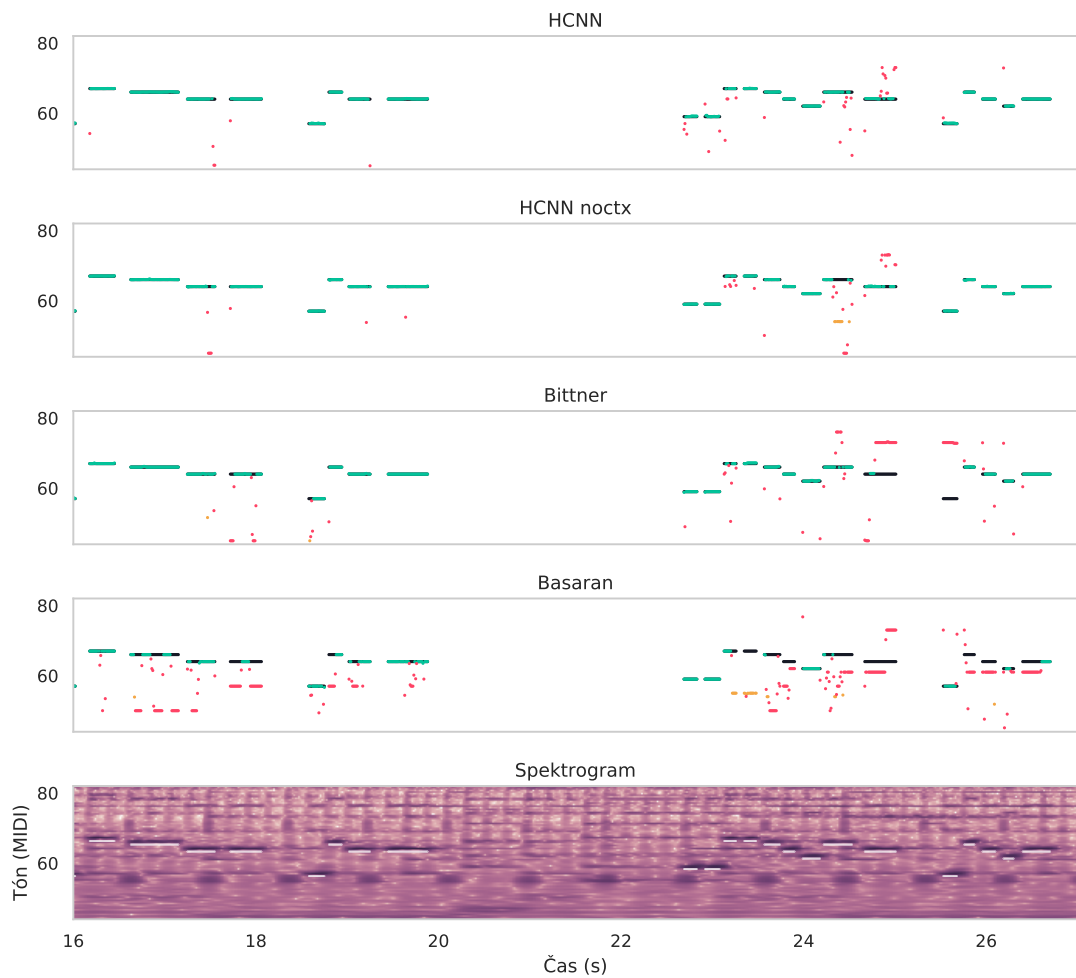
Basaranova metoda pro tuto koherenci výsledných kontur však obětovala frekvenční přesnost znějících výšek tónů, frekvenční rozlišení této metody je totiž na úrovni jednoho půltónu. Jak jsme již prezentovali v kapitole Experimenty, výstup metod kvantizovaný na půltóny obsahuje množství chyb navíc, jelikož často selhává v zachycení frekvenčních modulací. Na obrázku 6.6 vidíme další limitaci takového výstupu — pokud je obsah skladby laděný podle jiného referenčního tónu než jaký byl použit pro trénování sítě, znějící tóny vycházejí výškou „mezi“ výstupní složky. Z tabulky 6.8 je pak zřejmé, že kvůli této kvantizaci síť nedosahuje srovnatelných výsledků, přestože na jiných, žánrově shodných datech, které jsou laděny na správný referenční tón, podává kompetitivní výsledky. Limitace se proto týká zejména jazzových nahrávek pocházejících z období před rokem 1955, před zavedením referenčního tónu A4=440Hz ve standardu ISO16.

Metrika (Metoda)	CannonballAdderley_SoWhat
RPA (HCNN)	0.850
RPA (HCNN noctx)	0.848
RPA (Bittner)	0.828
RPA (Basaran)	0.653

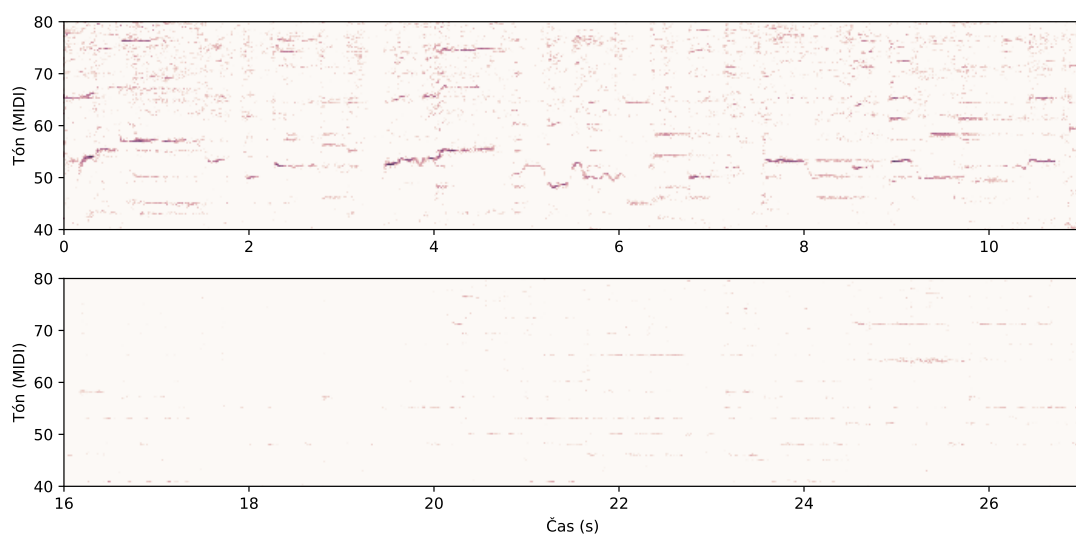
Tabulka 6.8: Přesnost metod na testovacím souboru CannonballAdderley_SoWhat z datasetu WJazzD.

Dalším problémem Basaranovy metody je zhoršená schopnost extrakce na syntetických datech. Příklady `midi2REF`, `midi3REF`, `train10REF` z datasetů ADC04 a MIREX05 jsou syntetizovány na základě MIDI pomocí základních zvukových fontů, nahrávky proto zní velmi uměle. Jak vidíme na obrázku 6.7, zatímco metody Bittnerové a HCNN si s touto syntetickou barvou hlasu dokáží poradit, výstup Basaranovy metody obsahuje šum a skoky k doprovázejícím nástrojům. Příčinou může být jiná vstupní reprezentace signálu, která je založena na práci Durrieu a David (2010) a spočívá na modelování hlavního hlasu pomocí zdroje a filtrů. Na obrázku 6.8 srovnáváme tuto reprezentaci pro vstupní signál s lidským zpěvem (nahore, `train01`) a pro signál se syntetickou flétnou (dole, `train10`). Je zřejmé, že zatímco lidský zpěv tato reprezentace dokáže zachytit, syntetický hlas na reprezentaci téměř zachycen není.

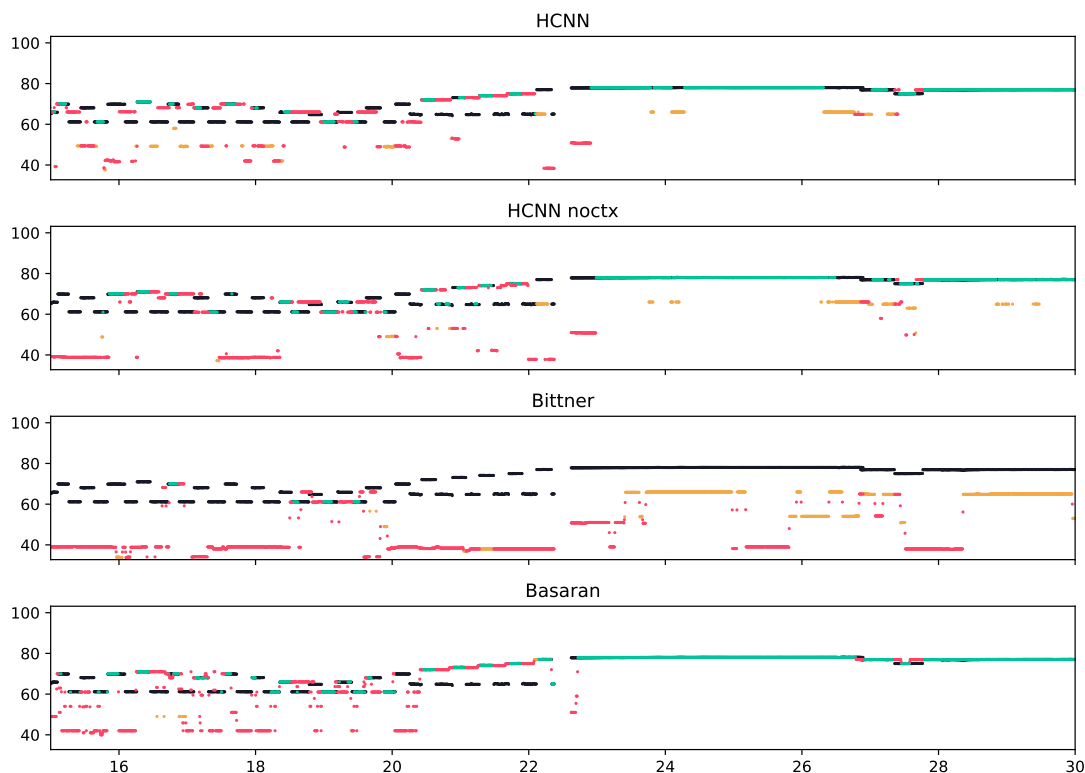
Pokud srovnáme metody HCNN a metodu Bittner, rozdílem v predikcích jsou zejména jiné priority, které přiřazují barvám hlasů. Lze tudíž nalézt mnoho příkladů, kde HCNN zároveň přepisuje melodii a nesprávně místy přeskakuje k nástrojům v doprovodu, zatímco metoda Bittnerové na stejném příkladu tuto chybu nedělá, podobně však existují i opačné příklady. Příkladem, ve kterém se tyto metody nejvíce rozcházejí, je soubor `MatthewEntwistle_FairerHopes` z kolekce MedleyDB, ve kterém melodii hraje harfa. Zvuk harfy se však nevyskytuje v množině trénovacích dat, rozdílem proto je, že zatímco metoda HCNN zvládá alespoň částečně generalizovat i na tuto dosud neslyšenou barvu hlasu, metoda Bittnerové tyto tóny úplně ignoruje a přepisuje doprovod pod harfou (viz obrázek 6.9).



Obrázek 6.7: Výstup metod na testovacím souboru `train10` z datasetu MIREX05.



Obrázek 6.8: Srovnání vstupní frekvenčně-časové reprezentace \mathbf{H}^{F_0} Basaranovy metody testovacích souborů `train01` a `train10` z datasetu MIREX05.



Obrázek 6.9: Výstup metod na testovacím souboru MatthewEntwistle_FairerHopes z datasetu MedleyDB.

6.4 Interpretace výsledků

Zásadní výhodou metody Basaran oproti HCNN a Bittner je zpracování výsledků funkce salience pomocí rekuretní sítě. Tento rozdíl spolu s jinou výchozí časově-frekvenční reprezentací signálu jeho metodu zvýhodňuje zejména na datasetu ORCHSET, kde jeho metoda v metrice RPA dosahuje o deset procentních bodů lepších výsledků. Pro metodu Basaran je na tomto datasetu také výhodné, že jeho referenční anotace mají půltónové frekvenční rozlišení, tedy stejné, jako výstup této metody. Tudíž jeho metoda při použití tohoto hrubého rozlišení nijak netratí. Na zbylých datasetech se však nižší frekvenční rozlišení projevuje více a metodu pravděpodobně spíše znevýhodňuje. Přesto jsou jeho predikce, zejména pak na složitějších vstupních datech, často koherentnější, obsahují méně šumu.

Výsledky HCNN a Bittner jsou si oproti výsledkům metody Basaran mnohem podobnější, ačkoliv je mezi nimi větší procentuální rozdíl. Při kvalitativním vyhodnocování jsme nenalezli příliš mnoho příkladů, na kterém by se výstupy metod výrazně lišily. Mezi sítěmi je však řada podobností, zejména stejná vstupní reprezentace, přibližně stejné velký zpracovávaný kontext, přeskočení vyhlazování výstupu ale i celková struktura sítě. Kvantitativní rozdíly na všech datasetech tedy přičítáme spíše lépe naučeným barvám nástrojů a jejich priorit v celkovém mixu. Podobnost výsledků ilustrujeme také výpočtem korelace, zatímco výsledky metriky RPA pro metody Bittner a HCNN mezi sebou mají korelaci 0.932, mezi Basaran a HCNN vychází nižší korelace 0.736.

Závěr

V práci jsme nastínili principy dosavadních přístupů k extrakci melodie, uvedli výčet veřejně dostupných dat a vysvětlili způsoby evaluace. Na těchto základech jsme pak prezentovali experimentální výsledky nových architektur určených pro výpočet funkce salience s důrazem na odhad výšky tónů v nahrávkách, jejichž návrhy čerpají z příbuzných oborů a úloh. Následně jsme tyto architektury porovnali s vybranými state-of-the-art metodami. Ze tří navrhovaných dosáhla architektura HCNN na většině veřejně dostupných datasetech nejlepších výsledků. Architektura HCNN spočívá v duplikaci a vzájemném posunu výstupů vnitřních konvolučních vrstev takovým způsobem, aby následujícím konvolučním vrstvám umožnila lépe zachytit harmonické složky znějících tónů v signálu.

Východisek pro navazující práce je více. Jako nejsnadnější se jeví přidání libovolné techniky vyhlazování výstupů. V práci se zaměřujeme na výpočet funkce salience, výstup této funkce však dále nezpracováváme, obvyklým krokem metod extrakce melodie je tuto reprezentaci dále používat buď pro detekci kandidátních kontur melodie a z těch pak vybírat výslednou melodii nebo použít statistický model, který výstup vyhladí, další možností je použití rekurentních neuronových sítí, které by dovolilo současné trénování HCNN a tohoto modulu zahrnujícího kontext. Díky tomu by se tento modul mohl naučit v čase zohledňovat nejen výšky tónů, ale i barvu nástroje (čímž by se metoda začala podobat existujícím, úspěšným metodám modelující hlavní hlas v nahrávce, Durrieu a kol. (2011)).

Dále je jistě možné se pokusit o zlepšení detekce melodie, ať už implementací složitějších systémů práhování nebo dokonce samostatného modulu pro detekci melodie. Při práci jsme zběžně testovali možnost použití nezávislé konvoluční sítě určené pouze pro detekci, jeho přesnost detekce však byla nižší, než pouhé práhování, výsledky jsme proto do práce již nezahrnovali.

Je také možné upravovat architekturu HCNN — pokoušet se dále zvětšovat kontext, nebo se naopak pokoušet ještě vylepšit základní bezkontextový model. Je možné experimentovat s vyššími konvolucemi (na frekvenční ose) pro zachycení závislostí mezi jednotlivými znějícími tóny, dále pak s méně často aplikovanou harmonickou transformací, která by dovolila trénování širších modelů, nebo například aplikace SpecAugment techniky nejen na vstupu sítě ale i uvnitř modelu jako speciální druh dropout regularizace.

Techniky augmentace dat jsou jistě také směr, kterým se dá ubírat, protože se s ní spojuje možnost zvyšování kapacity modelu bez přeučení. Přidávání šumu do nahrávek je augmentační evergreen, který výsledky zlepšil o několik desetin procentního bodu, žádná práce však nešla s augmentacemi příliš daleko. Možností je přitom více — například používání rozmanitých audio efektů z programů pro skládání hudby nebo mixování exponenciálně velkého množství trénovacích dat pomocí jednotlivých monofonních stop z datasetů

Jinou možností pokusu o řešení nedostatku dat je pomocí syntézy. Realistická syntéza dala například vzniknout novému klavírnímu datasetu MAESTRO (Hawthorne a kol., 2018), který svou délkou řádově předčil dosavadní datasety. Se zlepšujícími se generativními modely pro syntézu libovolných barev tónů (Engel a kol., 2019) se dá uvažovat o vzniku podobného datasetu i pro kompletní přepis nahrávek, a také pro přepis melodie.

Experimentovat lze samozřejmě i se vstupní reprezentací, například využitím zmiňované Fan Chirp Transform (Cancela a kol. (2010)), která nabízí velmi dobré rozlišení frekvenčně modulovaných signálů, autokorelace výstupů banky filtrů (Paiva a kol. (2006)), která je zatím nejbližší biologickým modelům lidského sluchu nebo využitím informace o fázi pro zpřesnění frekvenčních odhadů a amplitud vrcholků spektra (Instantaneous Frequency, Dressler (2009)).

Poslední možný směr, který předložíme je zmírnění vlivu malého množství dat pomocí multitask learningu. Díky přidání datasetů příbuzných úloh by bylo možné opět navýšit kapacitu modelů. Jako nejbližší se jeví úlohy kompletního přepisu hudby, tento směr již částečně prozkoumala Bittner a kol. (2018), případně oddělení signálů (source separation), o kterém publikované články nejsou.

Seznam použité literatury

- BALKE, S., DITTMAR, C., ABESSER, J. a MULLER, M. (2017). Data-driven solo voice enhancement for jazz music retrieval. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, pages 196–200. ISSN 15206149. doi: 10.1109/ICASSP.2017.7952145. URL <https://www.audiolabs-erlangen.de/content/05-fau/assistant/balke/01-publications/2017{ }BalkeDAM{ }SoloVoiceEnhancement{ }ICASSP.pdf>.
- BISHOP, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg. ISBN 0387310738.
- BITTNER, R., SALAMON, J., TIERNEY, M., MAUCH, M., CANNAM, C. a BELLO, J. (2014). MedleyDB: A multitrack dataset for annotation - intensive mir research. *International Society for Music Information Retrieval Conference*.
- BITTNER, R. M. (2018). Data-Driven Fundamental Frequency Estimation.
- BITTNER, R. M., MCFEE, B., SALAMON, J., LI, P. a BELLO, J. P. (2017). Deep Saliency Representations for F0 Estimation in Polyphonic Music. *Ismir*, pages 23–27. URL <https://bmcfee.github.io/papers/ismir2017{ }saliency.pdf>.
- BITTNER, R. M., MCFEE, B. a BELLO, J. P. (2018). Multitask Learning for Fundamental Frequency Estimation in Music. Technical report. URL <http://arxiv.org/abs/1809.00381>.
- BOSCH, J. J. a GÓMEZ, E. (2014). Melody Extraction in Symphonic Classical Music: a Comparative Study of Mutual Agreement Between Humans and Algorithms. *Proceedings of the Conference on Interdisciplinary Musicology*. URL <http://phenicx.upf.edu/system/files/publications/cim14{ }submission{ }114{ }ready.pdf>.
- BOSCH, J. J. a GÓMEZ, E. (2016). Melody Extraction Based on a Source-Filter Model Using Pitch Contour Selection. *13th Conference on Sound and Music Computing*, pages 67–74.
- BOSCH, J. J., MARXER, R. a GÓMEZ, E. (2016). Evaluation and combination of pitch estimation methods for melody extraction in symphonic classical music. *Journal of New Music Research*, **45** (2), 101–117. ISSN 17445027. doi: 10.1080/09298215.2016.1182191. URL <https://repositori.upf.edu/bitstream/handle/10230/26985/Bosch{ }NewMusic{ }Eval.pdf?sequence=1{ }isAllowed=y>.
- BROWN, J. C. (1990). Calculation of a constant Q spectral transform. *The Journal of the Acoustical Society of America*, **89**(1), 425–434. ISSN 0001-4966. doi: 10.1121/1.400476. URL <http://academics.wellesley.edu/Physics/brown/pubs/cq1stPaper.pdf>.

- CANCELA, P. (2008). Tracking melody in polyphonic audio. *4th Music Information Retrieval Evaluation eXchange*. URL <https://pdfs.semanticscholar.org/226e/a7b870fd229149fae2e6c8b15d8a3d4f9bb8.pdf>.
- CANCELA, P., LÓPEZ, E. a ROCAMORA, M. (2010). FAN CHIRP TRANSFORM FOR MUSIC REPRESENTATION. *Audio*, (1), 1–8. URL <https://iie.fing.edu.uy/publicaciones/2010/CLR10/CLR10.pdf>.
- CHAN, T. S., YEH, T. C., FAN, Z. C., CHEN, H. W., SU, L., YANG, Y. H. a JANG, R. (2015). Vocal activity informed singing voice separation with the iKala dataset. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, **2015-Augus**, 718–722. ISSN 15206149. doi: 10.1109/ICASSP.2015.7178063.
- D BASARAN, S ESSID, G. P. (2018). Main melody extraction with source-filter nmf and crnn. *Ismir*, pages 82–89. URL <http://ismir2018.ircam.fr/doc/pdfs/273{ }Paper.pdf>.
- DOWNIE, J. S., EHMANN, A. F., BAY, M. a JONES, M. C. (2010). eXchange : Some Observations and Insights. *Advances in Music Information Retrieval*, pages 93–115. URL <https://pdfs.semanticscholar.org/3836/15607f345a8cfbc5e884eaa5ca04f3c4b139.pdf>.
- DRESSLER, K. (2009). Audio melody extraction for mirex 2009. *Evaluation eXchange (MIREX)*, pages 1–3. URL <http://www.idmt.fraunhofer.de/content/dam/idmt/de/Dokumente/Produktflyer/QbH/white{ }paper{ }fraunhofer{ }idmt{ }audio{ }melody{ }extraction{ }mirex{ }2009.pdf>.
- DRESSLER, K. (2011). Pitch estimation by the pair-wise evaluation of spectral peaks. *42nd AES Conference*, pages 1–10. URL <http://www.aes.org/e-lib/browse.cfm?elib=15960>.
- DRESSLER, K. (2016). Automatic Transcription of the Melody from Polyphonic Music. URL <https://www.db-thueringen.de/servlets/MCRFileNodeServlet/dbt{ }derivate{ }00038847/ilm1-2017000136.pdf>.
- DURRIEU, J.-L. a DAVID, B. (2010). Source / Filter Model for Unsupervised Main Melody. **18**(3), 1–12. URL <https://www.irit.fr/{ }Cedric.Fevotte/publications/journals/ieee{ }asl{ }voice{ }extrac.pdf>.
- DURRIEU, J. L., DAVID, B. a RICHARD, G. (2011). A musically motivated mid-level representation for pitch estimation and musical audio source separation. *IEEE Journal on Selected Topics in Signal Processing*, **5**(6), 1180–1191. ISSN 19324553. doi: 10.1109/JSTSP.2011.2158801. URL <http://durrieu.ch/publis/durrieuDavidRichard{ }musicallyMotivatedRepresentation{ }JSTSP2011.pdf>.

- ENGEL, J., RESNICK, C., ROBERTS, A., DIELEMAN, S., ECK, D., SIMONYAN, K. a NOROUZI, M. (2017). Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders. URL <http://arxiv.org/abs/1704.01279>.
- ENGEL, J., AGRAWAL, K. K., CHEN, S., GULRAJANI, I., DONAHUE, C. a ROBERTS, A. (2019). GANSynth: Adversarial Neural Audio Synthesis. pages 1–17. URL <http://arxiv.org/abs/1902.08710>.
- GABOR, B. D. a MEMBER, A. (1945). THEORY OF COMMUNICATION * Part 1 . THE ANALYSIS OF INFORMATION. *Journal of the Institution of Electrical Engineers-Part III: Radio and Communication Engineering*, **93** (1946), 429–441.
- GOODFELLOW, I., BENGIO, Y. a COURVILLE, A. (2016). *Deep learning*. MIT press.
- GOTO, M. a HAYAMIZU, S. (1999). A real-time music scene description system: Detecting melody and bass lines in audio signals. *IJCAI-99 Workshop on Computational Auditory Scene Analysis*, (August), 31–40. URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.33.1085{&}rep=rep1{&}type=pdf>.
- GOTO, M., HASHIGUCHI, H., NISHIMURA, T. a OKA, R. (2002). RWC Music Database: Popular, Classical, and Jazz Music Databases. *International Conference on Music Information Retrieval*, (October), 287–288. URL <https://staff.aist.go.jp/m.goto/PAPER/ISMIR2002goto.pdf>.
- HAWTHORNE, C., ELSSEN, E., SONG, J., ROBERTS, A., SIMON, I., RAFFEL, C., ENGEL, J., OORE, S. a ECK, D. (2017). Onsets and Frames: Dual-Objective Piano Transcription. ISSN 09259902. doi: 10.1007/s10844-013-0258-3. URL <https://arxiv.org/pdf/1710.11153.pdf><http://arxiv.org/abs/1710.11153>.
- HAWTHORNE, C., STASYUK, A., ROBERTS, A., SIMON, I., HUANG, C.-Z. A., DIELEMAN, S., ELSSEN, E., ENGEL, J. a ECK, D. (2018). Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset. pages 1–12. URL <http://arxiv.org/abs/1810.12247>.
- HE, K., ZHANG, X., REN, S. a SUN, J. (2015). Deep Residual Learning for Image Recognition. URL <http://arxiv.org/abs/1512.03385>.
- HERMES, D. J. (1988). Measurement of pitch by subharmonic summation. *The Journal of the Acoustical Society of America*, **83**(1), 257–264. ISSN 0001-4966. doi: 10.1121/1.396427. URL <http://asa.scitation.org/doi/10.1121/1.396427>.
- HSU, C.-L. a JANG, J.-S. R. (2010). Singing Pitch Extraction by Voice Vibrato/Tremolo Estimation and Instrument Partial Deletion. *11th Int. Soc. for Music Info. Retrieval Conf.*, (Ismir), 525–530. URL <http://citeseerx.ist.psu.edu/viewdoc/download;jsessionid=9BE69D930B1BEADFC07E4534BA4E7E0?doi=10.1.1.231.4599{&}rep=rep1{&}type=pdf>.

- HUMPHREY, E. J., SALAMON, J., NIETO, O., FORSYTH, J., BITTNER, R. M. a BELLO, J. P. (2014). JAMS: A JSON Annotated Music Specification for Reproducible MIR Research. *Proceedings of the 15th International Society for Music Information Retrieval Conference*, (September), 591–596.
- IKEMIYA, Y., ITOYAMA, K. a YOSHII, K. (2016). Singing Voice Separation and Vocal F0 Estimation Based on Mutual Combination of Robust Principal Component Analysis and Subharmonic Summation. *IEEE/ACM Transactions on Audio Speech and Language Processing*, **24**(11), 2084–2095. ISSN 23299290. doi: 10.1109/TASLP.2016.2577879.
- IOFFE, S. a SZEGEDY, C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. URL <http://arxiv.org/abs/1502.03167>.
- KIM, J. W., SALAMON, J., LI, P. a BELLO, J. P. (2018). Crepe: A Convolutional Representation for Pitch Estimation. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, **2018-April**, 161–165. ISSN 15206149. doi: 10.1109/ICASSP.2018.8461329. URL <https://arxiv.org/pdf/1802.06182.pdf>.
- KINGMA, D. P. a BA, J. (2014). Adam: A Method for Stochastic Optimization. pages 1–15. URL <http://arxiv.org/abs/1412.6980>.
- KUM, S., OH, C. a NAM, J. (2016). Melody Extraction on Vocal Segments Using Multi-Column Deep Neural Networks. *Proceedings of the 17th International Society for Music Information Retrieval Conference, (ISMIR)*, (August), 819–825. URL https://wp.nyu.edu/ismir2016/wp-content/uploads/sites/2294/2016/07/119_{_}Paper.pdf.
- MAROLT, M. (2004). On finding melodic lines in audio recordings. *Proc. of the 7th Int. Conference on Digital Audio Effects*, pages 5–9. URL https://ant-s4.unibw-hamburg.de/dafx/paper-archive/2004/P_{_}217.PDF.
- MARTAK, L. S., SAJGALIK, M. a BENESOVA, W. (2018). Polyphonic Note Transcription of Time-Domain Audio Signal with Deep WaveNet Architecture. *International Conference on Systems, Signals, and Image Processing*, **2018-June**, 2–6. ISSN 21578702. doi: 10.1109/IWSSIP.2018.8439708. URL https://vgg.fiit.stuba.sk/wp-uploads/2018/09/iwssip2018_{_}wavenet.pdf.
- MAUCH, M. a DIXON, S. (2014). PYIN: A fundamental frequency estimator using probabilistic threshold distributions. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, **1**(1), 659–663. ISSN 15206149. doi: 10.1109/ICASSP.2014.6853678. URL <https://www.eecs.qmul.ac.uk/{~}simond/pub/2014/MauchDixon-PYIN-ICASSP2014.pdf>.
- PAIVA, R. P., MENDES, T. a CARDOSO, A. (2004). An Auditory Model Based Approach for Melody Detection in Polyphonic Musical Recordings. (May 2014), 21–40. doi: 10.1007/978-3-540-31807-1_2.

- PAIVA, R. P., MENDES, T. a CARDOSO, A. (2006). Melody Detection in Polyphonic Musical Signals: Exploiting Perceptual Rules, Note Saliency, and Melodic Smoothness. *Computer Music Journal*, **30**(4), 80–98. ISSN 0148-9267. doi: 10.1162/comj.2006.30.4.80.
- PARK, D. S., CHAN, W., ZHANG, Y., CHIU, C.-C., ZOPH, B., CUBUK, E. D. a LE, Q. V. (2019). SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition. URL <http://arxiv.org/abs/1904.08779>.
- PFLEIDERER, M., FRIELER, K., ABESSER, J., ZADDACH, W.-G. a BURKHART, B. *Inside the Jazzomat*. ISBN 9783959831246. URL <http://schott-campus.com/wp-content/uploads/2017/11/inside{ }the{ }jazzomat{ }final{ }rev{ }oa4.pdf>.
- POLINER, G. E. a ELLIS, D. P. W. (2005). A Classification Approach to Melody Transcription. *Proc. 6th Int. Conf. on Music Inform. Retrieval*, (1), 161–166.
- POLINER, G. E., ELLIS, D. P. W., EHMANN, A. F., GÓMEZ, E., STREICH, S. a ONG, B. (2007). Melody transcription from music audio: Approaches and evaluation. *IEEE Transactions on Audio, Speech and Language Processing*, **15**(4), 1247–1256. ISSN 15587916. doi: 10.1109/TASL.2006.889797. URL <https://academiccommons.columbia.edu/doi/10.7916/D8NC69RK/download>.
- RAFFEL, C., MCFEE, B., HUMPHREY, E. J., SALAMON, J., NIETO, O., LIANG, D. a ELLIS, D. P. W. (2014). mir_eval: A Transparent Implementation of Common MIR Metrics. *Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014)*, pages 367–372. URL <https://bmcfree.github.io/papers/ismir2014{ }mireval.pdf>.
- RAO, V. a RAO, P. (2010). Vocal melody extraction in the presence of pitched accompaniment in polyphonic music. *IEEE Transactions on Audio, Speech and Language Processing*, **18**(8), 2145–2154. ISSN 15587916. doi: 10.1109/TASL.2010.2042124. URL <https://www.ee.iitb.ac.in/course/{ }daplab/publications/PrePrintForWebsite.pdf>.
- RIGAUD, F. a RADENEN, M. (2016). Singing Voice Melody Transcription Using Deep Neural Networks. *Ismir*, pages 737–743. URL <https://s18798.pcdn.co/ismir2016/wp-content/uploads/sites/2294/2016/07/163{ }Paper.pdf>.
- RYYNÄNEN, M. P. a KLAURI, A. P. (2008). Automatic Transcription of Melody, Bass Line, and Chords in Polyphonic Music. *Computer Music Journal*, **32**(3), 72–86. ISSN 0148-9267. doi: 10.1162/comj.2008.32.3.72.
- SALAMON, J. a GÓMEZ, E. (2012). Melody extraction from polyphonic music signals using pitch contour characteristics. *IEEE Transactions on Audio, Speech and Language Processing*, **20**(6), 1759–1770. ISSN 15587916. doi: 10.1109/TASL.2012.2188515.
- SALAMON, J., GÓMEZ, E., ELLIS, D. P. a RICHARD, G. (2014). Melody extraction from polyphonic music signals: Approaches, applications, and challenges. *IEEE Signal Processing Magazine*, **31**

- (2), 118–134. ISSN 10535888. doi: 10.1109/MSP.2013.2271648. URL http://www.justinsalamon.com/uploads/4/3/9/4/4394963/salamon_{ }gomez_{ }ellis_{ }richard_{ }melodyextractionreview_{ }ieeespm_{ }2013.pdf.
- SALAMON, J., BITTNER, R. M., BONADA, J., BOSCH, J. J., GOMEZ, E. a JUAN PABLO BELLO (2017). An Analysis/Synthesis Framework for Automatic F0 Annotation of Multitrack Datasets. *Proceedings of the International Society for Music Information Retrieval {(ISMIR)} Conference*, pages 71–78.
- STOLLER, D., EWERT, S. a DIXON, S. (2018). Wave-U-Net: A Multi-Scale Neural Network for End-to-End Audio Source Separation. pages 334–340. ISSN 13514180. doi: arXiv:1806.03185v1. URL <http://arxiv.org/abs/1806.03185>.
- STURM, B. L. (2013). Classification accuracy is not enough. *Journal of Intelligent Information Systems*, **41**(3), 371–406. ISSN 0925-9902. doi: 10.1007/s10844-013-0250-y. URL <https://link.springer.com/content/pdf/10.1007/{%}2Fs10844-013-0250-y.pdf>.
- SU, L. (2018). Vocal melody extraction using patch-based CNN. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, **2018-April**, 371–375. ISSN 15206149. doi: 10.1109/ICASSP.2018.8462420. URL <https://arxiv.org/pdf/1804.09202.pdf>.
- SUTSKEVER, I., HINTON, G., KRIZHEVSKY, A. a SALAKHUTDINOV, R. R. (2014). Dropout : A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, **15**, 1929–1958.
- TACHIBANA, H., ONO, T., ONO, N. a SAGAYAMA, S. (2010). Melody line estimation in homophonic music audio signals based on temporal-variability of melodic source. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, (May 2014), 425–428. ISSN 15206149. doi: 10.1109/ICASSP.2010.5495764.
- VAN DEN OORD, A., DIELEMAN, S., ZEN, H., SIMONYAN, K., VINYALS, O., GRAVES, A., KALCHBRENNER, N., SENIOR, A. a KAVUKCUOGLU, K. (2016a). WaveNet: A Generative Model for Raw Audio. pages 1–15. URL <http://arxiv.org/abs/1609.03499>.
- VAN DEN OORD, A., KALCHBRENNER, N., VINYALS, O., ESPEHOLT, L., GRAVES, A. a KAVUKCUOGLU, K. (2016b). Conditional Image Generation with PixelCNN Decoders. URL <http://arxiv.org/abs/1606.05328>.
- YANG, L.-C., CHOU, S.-Y. a YANG, Y.-H. (2017). MidiNet: A Convolutional Generative Adversarial Network for Symbolic-domain Music Generation. URL <http://arxiv.org/abs/1703.10847>.
- YEH, T. C., WU, M. J., JANG, J. S. R., CHANG, W. L. a LIAO, I. B. (2012). A hybrid approach to singing pitch extraction based on trend estimation and hidden Markov models. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, (August),

457–460. ISSN 15206149. doi: 10.1109/ICASSP.2012.6287915. URL https://www.researchgate.net/profile/Jyh-Shing-Jang/publication/261113746_A_hybrid_approach_to_singing_pitch_extraction_base_links/55d5446f08ae6788fa352c8e/A-hybrid-approach-to-singing-pitch-extraction

Seznam obrázků

1.1	Příklad vstupu a výstupu metody pro extrakci melodie. přidat zvukový příklad do přílohy	3
1.2	Zvuk klarinetu, tóny s různou výškou a dynamikou, 35 milisekund signálu se vzorkovací frekvencí 44 100 Hz. nějak uvést zdroj zvuku https://www.philharmonia.co.uk/explore/sound_samples/clarinet?p=3	4
1.3	Zvuk klarinetu, absolutní hodnota výstupu Fourierovy transformace signálu délky 4096 s oknem typu Hamming.	5
1.4	Spektrogram zpěvu s doprovodem piana, basy a perkusí; zpívaná melodie je vyznačena bílým obrysem.	6
1.5	Spektrogram orchestrální skladby s obtížně detekovatelnou melodií. Čtenář se může pokusit ve spektrogramu nalézt melodii, na obrázku 1.6 nalezneme řešení.	7
1.6	Spektrogram orchestrální skladby s vyznačenou melodií.	8
1.7	Diagram obvyklého návrhu metod pro extrakci melodie.	9
1.8	Příklad výstupu výpočtu salienční funkce pomocí váženého sčítání harmonických frekvencí. Ačkoli je zpěv velmi zvýrazněn a salienční funkce na většině nahrávky dobře zachycuje výšku znějící melodie, doprovod kolem čtvrté sekundy nahrávky má vyšší hodnotu než zpěv, což neodpovídá lidskému vnímání zpěvu jakožto nejdůležitější složky signálu.	10
2.1	Znázornění kroků spektrální analýzy a výpočtu funkce salience.	15
2.2	Příklad trojúhelníkových filtrů pro transformaci frekvenční domény na logaritmickou škálu.	16
2.3	Ilustrace modelu tónu spolu se signálem, převzato z Marolt (2004)	19
2.4	Výsledky metod v soutěži MIREX v letech 2015-2018 s vybranými metodami ze starších ročníků.	22
2.5	Stagnující vývoj metod pro extrakci melodie.	23
4.1	Příklady chyb špatné detekce melodie negativně ovlivňující metriku VR.	34
4.2	Příklady chyb špatné detekce melodie zvyšující metriku VFA.	34
4.3	Příklady chyb ovlivňujících metriku RPA a RCA.	35
4.4	Příklady chyb „o oktávu“ ovlivňujících pouze metriku RPA, nikoli metriku RCA.	35
5.1	Diagram architektury CREPE, multiplikační koeficient 16x.	39
5.2	Výsledky experimentu s nezměněnou architekturou CREPE spuštěnou pro extrakci melodie.	41
5.3	Ukázky cílové reprezentace použité pro učení modelů. Pruhoaná linka značí dolní a horní hranici korektního odhadu výšky melodie.	42
5.4	Architektura CREPE s různou jemností diskretizace.	43
5.5	Histogramy vzdálenosti chybného odhadu, výstup prvního modelu má rozlišení 50 centů, výstup druhého 10 centů.	43
5.6	Architektura CREPE, vliv rozptylu cílové distribuce.	44

5.7	Architektura CREPE, vliv rozptylu cílové distribuce.	45
5.8	Architektura CREPE, vliv multirezoluční vstupní konvoluční vrstvy.	46
5.9	Vrstvení obyčejných konvolucí s lineárně rozšiřovaným dosahem, obrázek převzat z van den Oord a kol. (2016a).	47
5.10	Vrstvení dilatovaných konvolucí s exponenciálně rozšiřovaným dosahem, obrázek převzat z van den Oord a kol. (2016a).	47
5.11	Architektura WaveNet upravená pro kompletní přepis skladeb, upraveno na základě Martak a kol. (2018).	48
5.12	Úprava posledních vrstev WaveNet architektury.	49
5.13	Architektura WaveNet, vliv počtu filtrů dilatačních vrstev a skip propojení.	50
5.14	Architektura WaveNet, systematické prohledávání počtu dilatačních vrstev a bloků, vlevo hodnoty RPA, vpravo RCA.	51
5.15	Architektura WaveNet, vliv velikosti šířky kernelu dilatací.	52
5.16	Architektura WaveNet, vliv výstupní transformace skip propojení.	53
5.17	Architektura WaveNet, vliv velikosti první konvoluce.	53
5.18	Znázornění harmonických závislostí na spektrogramu s logaritmickou osou frekvence.	54
5.19	Diagram transformace vstupu konvoluční vrstvy pro zachycení harmonických souvislostí.	55
5.20	Diagram konvolučního bloku architektury HCNN.	56
5.21	Diagram celkového propojení architektury HCNN.	56
5.22	Architektura HCNN, Detekce melodie, vliv počáteční pooling vrstvy.	57
5.23	Architektura HCNN, Vliv velikosti skoku pro výpočet vstupního spektrogramu	58
5.24	Architektura HCNN, Vliv vstupního vícekanálového CQT	59
5.25	Úpravy architektury HCNN z obrázku 5.21 pro rozšíření uvažovaného kontextu.	60
5.26	Architektura HCNN, Vliv úpravy architektury ovlivňující receptivní pole modelu.	61
5.27	Architektura HCNN, vliv použití SpecAugment na vstupní spektrogramy.	62
6.1	Výsledky nejúspěšnějších metod na datasetu MedleyDB	68
6.2	Legenda pro následující kvalitativní srovnání.	69
6.3	Příklad s vysokou úspěšností přepisu <code>mirex05_train01</code> z datasetu MIREX05train, se kterým je čtenář seznámen z úvodu práce.	69
6.4	Výstup metod na testovacím souboru <code>Musorgski-Ravel-PicturesExhibition-ex6</code> z datasetu ORCHSET.	70
6.5	Výstupní salience metody HCNN na testovacím souboru <code>Musorgski-Ravel-PicturesExhibition-ex6</code> z datasetu ORCHSET.	71
6.6	Detail přepisu metod na testovacím souboru <code>CannonballAdderley_SoWhat</code> z datasetu WJazzD.	71
6.7	Výstup metod na testovacím souboru <code>train10</code> z datasetu MIREX05.	73
6.8	Srovnání vstupní frekvenčně-časové reprezentace \mathbf{H}^{F_0} Basaranovy metody testovacích souborů <code>train01</code> a <code>train10</code> z datasetu MIREX05.	73

6.9	Výstup metod na testovacím souboru MatthewEntwistle_FairerHopes z datasetu MedleyDB.	74
-----	---	----

I

Seznam tabulek

3.1	Souhrnná tabulka se základními informacemi o veřejně dostupných datasetech.	27
4.1	Přehled evaluačních datasetů v soutěži MIREX, částečně převzato z práce Salamon a kol. (2014).	32
4.2	Přehled trénovacích, validačních a testovacích množin použitých v práci.	33
5.1	Počty filtrů v konvolučních vrstvách v architektuře CREPE v závislosti na multiplikačním koeficientu.	40
5.2	Výsledky pokusu o replikaci architektury CREPE pro sledování jednohlasu. Přesnosti nejsou přímo srovnatelné kvůli různým evaluačním strategiím.	41
5.3	Výsledky experimentu s nezměněnou architekturou CREPE spuštěnou pro extrakci melodie. Přesnosti uvádíme pro různé kapacity modelu.	41
5.4	Architektura CREPE s různou jemností diskretizace.	42
5.5	Architektura CREPE, vliv rozptylu cílové distribuce.	44
5.6	Architektura CREPE, vliv šířky vstupního okna.	45
5.7	Počet filtrů prvních vrstev multirezoluční vstupní konvoluční vrstvy v architektuře CREPE.	46
5.8	Architektura CREPE, vliv multirezoluční vstupní konvoluční vrstvy.	46
5.9	Architektura WaveNet, vliv počtu filtrů dilatačních vrstev a skip propojení.	50
5.10	Architektura WaveNet, dosah a kapacita v závislosti na dilatačních počtu vrstev a bloků.	51
5.11	Architektura WaveNet, vliv velikosti šířky kernelu dilatací.	52
5.12	Architektura WaveNet, vliv výstupní transformace skip propojení.	52
5.13	Architektura WaveNet, vliv velikosti první konvoluce.	53
5.14	Architektura HCNN, vliv použití SpecAugment na vstupní spektrogramy.	63
6.1	Nastavení architektury a hyperparametrů pro testovanou architekturu CREPE.	65
6.2	Nastavení architektury a hyperparametrů pro testovanou architekturu WaveNet.	65
6.3	Nastavení architektury a hyperparametrů pro testovanou architekturu HCNN noctx.	65
6.4	Nastavení architektury a hyperparametrů pro testovanou architekturu HCNN.	65
6.5	Výsledky celkové přesnosti (Overall Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.	67
6.6	Výsledky přesnosti odhadu výšky (Raw Pitch Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.	67

6.7	Výsledky přesnosti odhadu výšky nezávisle na oktávě (Raw Chroma Accuracy). Vyznačené výsledky jsou pro daný dataset nejvyšší z porovnávaných v rámci daného datasetu.	67
6.8	Přesnost metod na testovacím souboru CannonballAdderley_SoWhat z datasetu WJazzD.	72
6.9	Architektura HCNN, Vliv úpravy architektury ovlivňující receptivní pole modelu.	91
6.10	Výsledky úplnosti detekce (Voicing Recall).	92
6.11	Výsledky nesprávné detekce (Voicing False Alarm). Nižší hodnota je lepší.	92

Seznam použitých zkratek

Přílohy

- 6.5 Architektura HCNN, Vliv úpravy architektury ovlivňující receptivní pole modelu
- 6.6 Výsledky detekce melodie testovaných metod

Konfigurace bloků	Úprava architektury	RPA	RCA
8 filtrů, 4 konv. bloky	noctx	0.722	0.786
	deep_ctx_3	0.728	0.796
	first_layers_ctx	0.728	0.795
	last_layers_ctx	0.737	0.800
	1_last_layer_dilated	0.733	0.797
	2_last_layers_dilated	0.733	0.796
	3_last_layers_dilated	0.731	0.796
	2_last_layers_wavenet	0.733	0.798
	3_last_layers_wavenet	0.745	0.804
16 filtrů, 4 konv. bloky	noctx	0.737	0.802
	deep_ctx_3	0.744	0.806
	first_layers_ctx	0.744	0.804
	last_layers_ctx	0.757	0.817
	1_last_layer_dilated	0.743	0.802
	2_last_layers_dilated	0.749	0.814
	3_last_layers_dilated	0.752	0.813
	2_last_layers_wavenet	0.754	0.813
	3_last_layers_wavenet	0.761	0.819
8 filtrů, 8 konv. bloků	noctx	0.732	0.801
	deep_ctx_3	0.735	0.796
	first_layers_ctx	0.747	0.814
	last_layers_ctx	0.749	0.811
	1_last_layer_dilated	0.745	0.806
	2_last_layers_dilated	0.751	0.814
	3_last_layers_dilated	0.749	0.806
	2_last_layers_wavenet	0.754	0.814
	3_last_layers_wavenet	0.751	0.803
16 filtrů, 8 konv. bloků	noctx	0.744	0.809
	deep_ctx_3	0.741	0.815
	first_layers_ctx	0.748	0.817
	last_layers_ctx	0.749	0.819
	1_last_layer_dilated	0.753	0.817
	2_last_layers_dilated	0.758	0.819
	3_last_layers_dilated	0.757	0.817
	2_last_layers_wavenet	0.759	0.819
	3_last_layers_wavenet	0.749	0.809

Tabulka 6.9: Architektura HCNN, Vliv úpravy architektury ovlivňující receptivní pole modelu.

Metoda	ADC04	MDB-m-s test	MIREX05 train.	MDB test	ORCH- SET	WJazzD test
Salamon	0.774	0.729	0.841	0.705	0.603	0.794
Bittner	0.796	0.638	0.796	0.675	0.614	0.846
Basaran	0.732	0.704	0.713	0.676	0.605	0.841
CREPE	0.584	0.431	0.576	0.449	0.326	0.680
WaveNet	0.765	0.595	0.747	0.618	0.494	0.784
HCNN noctx	0.806	0.684	0.824	0.728	0.729	0.880
HCNN	0.794	0.692	0.836	0.741	0.721	0.872

Tabulka 6.10: Výsledky úplnosti detekce (Voicing Recall).

Metoda	ADC04	MDB-m-s test	MIREX05 train.	MDB test	ORCH- SET	WJazzD test
Salamon	0.103	0.394	0.263	0.300	0.385	0.271
Bittner	0.278	0.273	0.308	0.306	0.490	0.333
Basaran	0.188	0.271	0.160	0.290	0.407	0.274
CREPE	0.178	0.252	0.171	0.243	0.235	0.213
WaveNet	0.311	0.383	0.387	0.397	0.426	0.370
HCNN noctx	0.246	0.312	0.336	0.333	0.535	0.339
HCNN	0.222	0.258	0.278	0.310	0.511	0.300

Tabulka 6.11: Výsledky nesprávné detekce (Voicing False Alarm). Nižší hodnota je lepší.