

Judah Daniels

Inferring Harmony from Free Polyphony

Computer Science Tripos – Part II

Clare College

July, 2023

Declaration of originality

I, Judah Daniels of Clare College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose. I am content for my dissertation to be made available to the students and staff of the University.

Signed Judah Daniels

Date April 22, 2023

Proforma

Candidate Number: **2200D**
Project Title: **Inferring Harmony from Free Polyphony**
Examination: **Computer Science Tripos – Part II, July, 2023**
Word Count: **5724¹**
Code Line Count: **2272²**
Project Originator: **Christoph Finkensiep**
Supervisor: **Dr Peter Harrison**

Original Aims of the Project

Work Completed

Special Difficulties

There were no special difficulties encountered in this project

¹This word count was computed by `texcount main.tex`

²This code line count was computed by using `cloc`

Contents

1	Introduction	9
1.1	Related Work	9
1.2	Achievements	10
2	Preparation	11
2.1	Probabilistic Inference	11
2.1.1	Tools of Probabilistic Inference	11
2.1.2	Inference with Latent variables	12
2.2	Overview of Approach	12
2.3	Inferring Labels	13
2.4	Inferring Structure	14
2.4.1	Voices	14
2.4.2	Primitive Protovoice Operations	16
2.4.3	Outer Structure	16
2.5	Parsing	18
2.6	Heuristic Search Algorithms	19
2.7	Starting Point	20
2.7.1	Relevant courses and experience	20
2.7.2	Existing codebase	20
2.8	Requirements Analysis	20
2.9	Software Engineering Techniques	21
2.9.1	Development model	21
2.9.2	Languages, libraries and tools	22
2.9.3	Hardware, version control and backup	22
3	Implementation	25
3.1	ProtoVoice Harmony Model	25
3.2	End to End Pipeline	26

3.3	Protovoice Parser	27
3.4	Harmony Model	32
3.4.1	Baseline algorithms	33
3.5	Extension Implementation	34
3.5.1	Extending the Harmony Model	34
3.5.2	Heuristic Design	35
3.5.3	Heuristic Search	40
3.5.4	Beam Search	40
3.5.5	Dual Beam Search	41
3.5.6	Stochastic Dual Beam Search	41
3.6	Choosing hyper parameters	42
3.7	Testing	42
3.7.1	Unit Tests	42
3.7.2	Qualitative Tests	42
3.8	Repository Overview:	42
4	Evaluation	45
4.1	Metrics	45
4.2	Harmony Model	45
4.2.1	Accuracy	45
4.2.2	Sensitivity Analysis	46
4.2.3	Qualitative Analysis	46
4.3	Baseline Algorithms	46
4.4	Extension Algorithms	46
4.4.1	Accuracy	46
4.4.2	Scalability	46
4.4.3	Qualitative Analysis	46
4.5	Limitations	46
5	Conclusions	49
5.1	Achievements	49
5.2	Lessons learned	49
5.3	Future Work	49
	Bibliography	49
	A Additional Information	55

B	Project Proposal	57
B.1	Abstract	2
B.2	Substance and Structure	2
B.2.1	Core: Search	2
B.2.2	Core: Evaluation	3
B.2.3	Extension	3
B.2.4	Overview	3
B.3	Starting Point	4
B.4	Success Criteria	4
B.5	Timetable	6
B.6	Resources	7
B.7	Supervisor Information	7

Acknowledgements

I'd like to thank Christoph Finkensiep and Peter Harrison for being awesome.

Introduction

Most of western tonal music can be described using a sequence of chords, representing a higher level harmonic structure of a piece. Automatic Chord Estimation (ACE) is the task of inferring the sequence of chords for a given piece from symbolic data. There is a small, finite set of chord types, but each chord can be expressed in a practically infinite number of ways. Given a score S (a symbolic representation of a piece of music), we wish to infer the sequence of underlying chord types $L = l_0, \dots, l_n$.

$$\hat{L} = \arg \max_L p(L|S) \quad (1.1)$$

Automatic Chord Estimation has both theoretical and practical applications. Analysis of music often starts with the manual labeling of each chord, which is a time consuming and cognitively demanding expert task. Sequences of chords provide compact representations for use in analysis, music identification and music similarity finding. More broadly speaking, any system that involves the understanding of written tonal music will benefit from chord estimation.

The paper *Modeling and Inferring Protovoice Structure in Free Polyphony* describes a generative model that expresses a piece of music as a result of the recursive application of a few primitive operations on notes (units of sound). The combination of these operations results in a hierarchical structure which encodes explicit relations between all the notes in a piece. This process can be reversed by *parsing* a piece of music according to the grammar described by the model, resulting in a list of possible *derivations*, that is, a list of all the ways the piece of music could have been produced through the recursive application of these simple rules. Each of these derivations capture an explanation of the underlying note-level structure of a piece, which can in turn be used to infer a list of descriptive chord labels.

While the original model could in theory be used to generate harmonic annotations, an exhaustive parse strategy would be prohibitively time-consuming in practice for any but the shortest musical extracts; one half measure can have over 100,000 valid derivations [8].

1.1 Related Work

Automatic chord estimation systems first emerged in in the 60's, utilising handcrafted grammar/rule-based systems [27] [48], followed by the development of optimisation algorithms in the early 2000s [32]. In more recent years, supervised learning approaches have have risen in popularity, exploiting large datasets and improved compute power [31] [28] [25].

The protovoice model offers a novel approach that unifies three aspects of tonal music analysis that are traditionally considered independently. Voice-leading refers how notes relate to each other *sequentially* through time, while harmony emerges from the simultaneous interaction between notes. Note *function* how notes relate to each other through recursive *functional dependencies*. Previous models have been developed alongside inference algorithms to perform automatic chord estimation that consider these dimensions of musical structure separately [?] [48], but the protovoice is the first to model free polyphony, one of the least restrictive forms of Western tonal music, and does so at a granular level.

1.2 Achievements

This was an ambitious project; I met all the Success Criteria and completed the extension tasks. I show that the protovoice model can be used to effectively annotate pieces with chord labels, and these results provide a promising foundation for the model being developed further as a sophisticated tool for the automated analysis of western total music. The PVHM has been made **open source** to accommodate future research in the area.

Preparation

2.1 Probabilistic Inference

Probabilistic inference is a statistical inference paradigm which reduces to representing the *joint distribution* of all the random variables in a system, such that we can compute any probability of interest in that system. This allows us to *infer* information about *latent*(hidden) variables based on *observed* evidence. Let X and Z denote the observed and latent(unobserved) variables respectively, and let Y denote the random variable to be inferred. By formulating a model of the *joint distribution* $P(X, Y, Z)$, the variable Y can be inferred from the evidence X by solving $\arg \max_Y P(Y|X)$.

2.1.1 Tools of Probabilistic Inference

There are two methods used to find the most likely hypothesis Y :

- The **maximum a priori estimate** maximises the *conditional likelihood* of the hypothesis given the evidence, given by $P(Y|X)$. This requires us to know the *prior* $P(Y)$, the distribution of the hypothesis, beforehand. The prior can be learned from labeled data.

$$\begin{aligned}\hat{Y} &= \arg \max_Y P(Y|X) \\ &= \arg \max_Y P(X|Y)P(Y)\end{aligned}\tag{2.1}$$

- The **maximum likelihood estimate** maximises the conditional likelihood of the observed evidence given the hypothesis, given by $P(X|Y)$. This method allows us to capture *uncertainty* in the prediction, given by $P(X|\hat{Y})$.

$$\hat{Y} = \arg \max_Y P(X|Y)\tag{2.2}$$

Definition 2.1.1 (Marginalisation).

$$P(A) = \sum_B P(A|B) P(B)\tag{2.3}$$

Definition 2.1.2 (Chain Rule of Probability).

$$P(A, B) = P(A|B) P(B)\tag{2.4}$$

2.1.2 Inference with Latent variables

To incorporate the latent variables Z , we use *marginalisation* and the *chain rule of probability* to show that:

$$\begin{aligned} P(Y|X) &= \sum_Z P(Y, Z|X) \\ &= \sum_Z P(Y|X, Z) P(Z|X) \end{aligned} \quad (2.5)$$

Given an abstract model $P(Z, X)$, the sum over all values of Z can be avoided by finding $\hat{Z} = \arg \max_Z P(Z|X)$, the most likely value of Z given X , using MLE or MAP estimation.

We thus update the prior $P(Z|X)$ as follows:

$$P(Z|X) = \begin{cases} 1 & \text{if } Z = \hat{Z} \\ 0 & \text{otherwise} \end{cases}$$

This gives us an approximation of $P(Y|X)$:

$$P(Y|X) \approx P(Y|X, \hat{Z}) \quad (2.6)$$

2.2 Overview of Approach

Probabilistic programming is a programming paradigm that makes use of model definitions and statistical inference algorithms to compute the conditional distribution of inputs that could have given rise to an observed output.

In the context of ACE, the **input** is the underlying sequence of chord labels L_0, L_1, \dots, L_n , and the score is the **observed output** S . In this sense, ACE can be solved by finding the most likely sequence of labels for the given surface, described by the equation:

$$\hat{L} = \arg \max_L P(L|S) \quad (2.7)$$

The difficulty arises from the complexity and prohibitively large number of the **latent variables** ϕ ; in reality, we need to maximise $\sum_{\phi} P(L|S, \phi) P(\phi|S)$.

$$\hat{L} = \arg \max_L \sum_{\phi} P(L|S, \phi) P(\phi|S) \quad (2.8)$$

The set of latent variables ϕ is **practically infinite**. These include the author's compositional conception, their musical conception, cognitive phenomena experienced by listeners, shared experience distilled into music theory, musical trends/culture and notational conventions. This cannot be solved analytically, but \hat{L} can be approximated using models that encode domain specific knowledge about both the music generation and labelling processes, combined with efficient inference algorithms.

Approximating Conditional Likelihood with MLE

Now consider the set of latent RVs ϕ as the union of two *disjoint sets*, ϕ' and ϕ'' , where ϕ' is the subset of ϕ which we can *infer* given L and S , and $\phi'' = \phi \setminus \phi'$. Assuming $P(\phi'|L)$ and $P(\phi''|L)$ are independent, it follows:

$$\begin{aligned} P(L|S) &= \sum_{\phi'} \sum_{\phi''} P(L|S, \phi', \phi'') P(\phi', \phi''|S) \\ &= \sum_{\phi'} \sum_{\phi''} P(L|S, \phi', \phi'') P(\phi'|S) P(\phi''|S) \end{aligned} \quad (2.9)$$

The uniform distribution is used to model $P(\phi''|S)$ as we have no prior knowledge of those latent variables. Using the technique described in Equation 2.6, we find $\hat{\phi}'$, thus giving:

$$\hat{L} \approx \arg \max_L P(L|S, \hat{\phi}') \quad (2.10)$$

2.3 Inferring Labels

The task of inferring harmony (ACE) poses **three main challenges**:

- **Segmentation:** Segmentation is the task of dividing a score into discrete regions described by different chord labels. For example, Figure 2.1 shows each segment separated by a dashed grey line. As these segments are typically not given a priori, both the segmentation and harmony needs to be inferred. Performing the joint task of segmentation and labelling is beyond the scope of this project, however, so it is assumed that the segmentation is given.
- **Ambiguity:** The notes in a given segment may not be sufficient to determine the chord label. For example, the slice containing notes D and B in Figure 2.1 could be a realisation of a Bm triad, but the context of the neighbouring slices as well as the functional dependencies of the notes makes it clear that the label should be a G. Furthermore, ambiguity is often used with much license to create artistic interest.

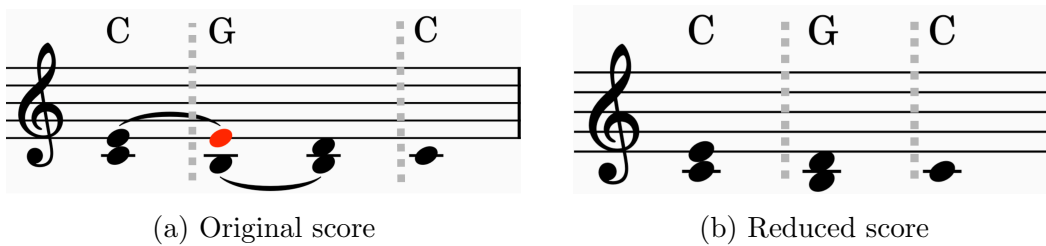


Figure 2.1: Applying a reduction to remove a non-chord-tone

- **Non chord-tones:** Any given segment will typically consist of mixture of *chord-tones* and *non chord-tones*. The chord tones directly define the harmony, so to perform ACE,

an algorithm will need to distinguish between the two. Previous approaches have involved modelling generation as a noisy process, such that the non chord-tones are considered as noise [43]. The protovoice model allows the non-tone tones to be *explicitly* explained away as being generated as *ornamentations*, or decorations of chord-tones. In Figure 2.1a, the non chord-tone is denoted in red. By applying a *reduction* that removes the neighbour note, E, we result in a single slice in that segment which only contains chord-tones(Figure 2.1b), thus describing the chord label more clearly, shown in Figure 2.1b.

1

2.4 Inferring Structure

The protovoice model is a formal generative model which represents a piece of music as a graph where each note is a node, and the edges between nodes represent relations between notes. The model is primarily concerned with the analysis of Western Classical music, although its expressiveness and generality means it could be applied to different musical styles including jazz or some popular western music [8].

2.4.1 Voices

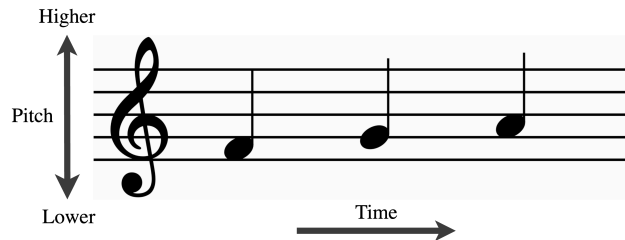


Figure 2.2: An example of music notation showing an ascending stepwise sequence.

The input is called a *score*, a symbolic abstraction of a piece of music based on a *2-dimensional axis*. The marks on on score represents notes, with the *pitch* of the note a function of its position on the vertical axis ², and the notes' placement in time represented by its position on the horizontal axis.

The notion of a *voice* is crucial for the understanding of the protovoice model. A voice typically refers to a single melodic line (sequence of notes) that is part of a musical composition. The

¹I'm considering adding a small section describing some relevant distributions. Dirchelet, Beta, Multinomial, categorical? It took a while for me to get my head around the fact that we used probability distributions of parameters of other probability distributions, e.g. sampling from Dirichlet for parameters of a multinomial or categorical.

²This is a simplification as there are other factors that determine the pitch, such as the key signature, accidentals and intonation.

term is derived from its use in choral music, such as J.S Bach’s four-voice chorales, which consist of 4 sung melodic lines. The term voice is used more generally however, as the melodic lines do not need to be sung or voice-like in character and can be performed by any melodic instrument.

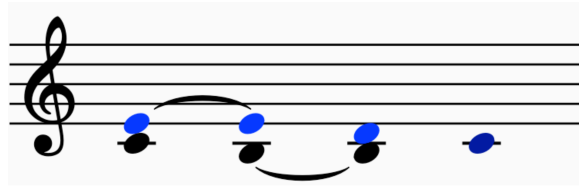


Figure 2.3: A short cadential phrase with two voices.

Polyphony refers to a piece of music that contains more than one voice. Typically, polyphonic music will have a set number of voices throughout the piece, but *free polyphony* refers to music wherein the number of concurrent voices is *arbitrary*; voices can start and stop freely and combine with one another throughout the piece.

There are three types of relations between notes that form the basis for the protovoice model:

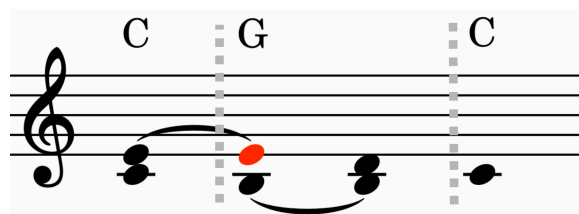


Figure 2.4: Chord labels shown with segment boundaries.

- **Horizontal:** As music is perceived in time, natural sequential relations arise between subsequent notes, in fact, a total order can be defined on the notes of a piece of music based on their positions on the horizontal axis.
- **Vertical:** This refers to the pitch axis on the score. The vertical arrangement of the notes determine the emergent harmony when they are perceived simultaneously. In most cases, multiple *voices* heard together will lead to an emerging sequence of harmonies, which can be described using chord labels.

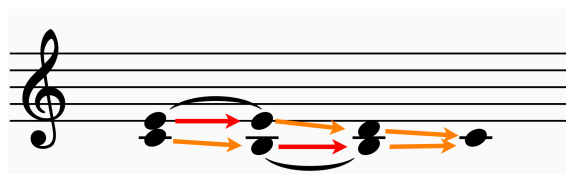


Figure 2.5: Functional dependencies between notes.

- **Functional:** Functional relations refer to the purpose or function of a note relative to another note. These functions can include repetitions, where both notes have the same *pitch*, or ornaments/neighbour notes, where the child note is a step (single unit of pitch) away from the parent note. These relations can be applied recursively, giving rise to a network of dependencies.

2.4.2 Primitive Protovoice Operations

Protovoices are represented as a directed graph G with one vertex for each note, a vertex each for the beginning (\bowtie) and end (\bowtie) of the piece, and edges that indicate connections between notes. A note is an instance of a pitch, defined as $p \in \mathcal{P}$ where \mathcal{P} is the set of pitches, i.e the set of vertical positions on the score³. A *protovoice* is a path within this graph. The protovoice model is characterised by stepwise generative operations on notes. A *repetition* inserts a child note of the same pitch on either side of the parent note, a *neighbour* inserts a child note that is a step away from the parent, and finally *passing* notes connect two protovoices.

The generation of a piece begins with the empty piece ($\bowtie \rightarrow \bowtie$) and is defined by the recursive application of operations. These operations relate child notes to one or two *parent* notes.

Single-sided operations are represented by attaching a new *child* note with an edge connected to a parent note:

$$\begin{aligned}
 x &\Longrightarrow n \rightarrow x && \text{left-neighbor} \\
 x &\Longrightarrow x \rightarrow n && \text{right-neighbor} \\
 x &\Longrightarrow x' \rightarrow x && \text{repeat-before} \\
 x &\Longrightarrow x' \rightarrow x && \text{repeat-after}
 \end{aligned} \tag{2.11}$$

Double-sided operations are represented by *edge replacement*.

$$\begin{aligned}
 \bowtie \rightarrow \bowtie &\Longrightarrow \bowtie \rightarrow x \rightarrow \bowtie && \text{root-note} \\
 x \rightarrow y &\Longrightarrow x \rightarrow x' \rightarrow y && \text{repeat-after'} \\
 x \rightarrow y &\Longrightarrow x \rightarrow y' \rightarrow y && \text{repeat-before'} \\
 x_1 \rightarrow x_1 &\Longrightarrow x \rightarrow n \rightarrow x_2 && \text{full-neighbor} \\
 x_1 \rightarrow x_1 &\Longrightarrow x \rightarrow' \rightarrow x_2 && \text{full-repeat}
 \end{aligned} \tag{2.12}$$

Finally, *passing* operations fill passing edges with notes from the left or right until the edge is fully stepwise.

$$\begin{aligned}
 x \rightarrow y &\Longrightarrow x \rightarrow p \rightarrow y && \text{passing-left} \\
 x \rightarrow y &\Longrightarrow x \rightarrow p \rightarrow y && \text{passing-right} \\
 x \rightarrow y &\Longrightarrow x \rightarrow p \rightarrow y && \text{passing-final}
 \end{aligned} \tag{2.13}$$

2.4.3 Outer Structure

The inner structure provided by protovoices captures the sequential and functional organisation of notes, but does not capture when notes are simultaneous. To model simultaneity of notes, *slices* are introduced, representing a group of notes heard together, and *transitions* which contain the protovoice edges between notes in the two neighbouring slices.

As slices and transitions contain notes and edges respectively, slices and transitions are called *outer structure*, and the notes and edges contained therein are called *inner structure*.

³See Appendix C for a detailed explanation of the pitch representation used in this project

Definition 2.4.1 (Multiset). A *multiset* is a set that allows multiple instances for each of its elements, formally defined as an ordered pair (A, m) where A is the *underlying set* of the multiset, and $m : A \rightarrow \mathbb{Z}^+$ gives the *multiplicity*, such that the number of occurrences of a in (A, m) is given by $m(a)$.

Definition 2.4.2 (Slice). A *slice* $s \in \mathcal{S}$ is defined as a multiset of pitches (\mathcal{P}, m) .

$$s = \{p_1^{m(p_1)}, \dots, p_n^{m(p_n)}\} \quad (2.14)$$

Definition 2.4.3 (Transition). A *transition* $t \in \mathcal{T}$ relates two adjacent slices, s_l and s_r , with a configuration of edges $e = (e_{reg}, e_{pass})$.

$$t = (s_l, e, s_r) \quad (2.15)$$

The slices and transitions form a graph given slices as nodes and transitions as edges (containing inner edges). However, as transitions only relate sequentially adjacent slices, the outer structure is in fact a *path graph* and can thus be represented as an alternating list slices and transitions.

Definition 2.4.4 (Path Graph). A *path graph*, referred to as a **path**, is a graph which can be represented as an alternating sequence of elements from two sets A and B , defined inductively as:

$$\begin{aligned} P &= a & a &\in A \\ P &= abP & a &\in A, b \in B \end{aligned} \quad (2.16)$$

Definition 2.4.5 (Outer Structure).

The *outer structure* is thus defined as a **path**:

$$P = t_1 s_1 t_2 s_2 \dots t_n \quad t_i \in \mathcal{T}, s_i \in \mathcal{S}, i \in 1 \dots n \quad (2.17)$$

The outer structure is generated by applying three operations described as production rules recursively:

- A **split** replaces a parent transition t by inserting a new slice s' and two surrounding transitions t_l and t_r . Each of the edges in t can have by one or more inner operations applied to it. The edges generated by these inner operations can either be discarded, or kept to form the new edges of t_l and t_r .

$$t \rightarrow t'_l s' t'_r \quad (2.18)$$

- A **spread** replaces a parent slice s by distributing its notes to two child slices s'_l and s'_r . The three non-terminals on the left hand side is what makes the protovoice grammar *context-sensitive*.

$$t_l s t_r \rightarrow t'_l s'_l t'_m s'_r t'_r \quad (2.19)$$

- A **freeze** marks a transition as terminal, such that the edges within the transition can no longer have operations applied to it.

$$t \rightarrow \underline{t} \quad (2.20)$$

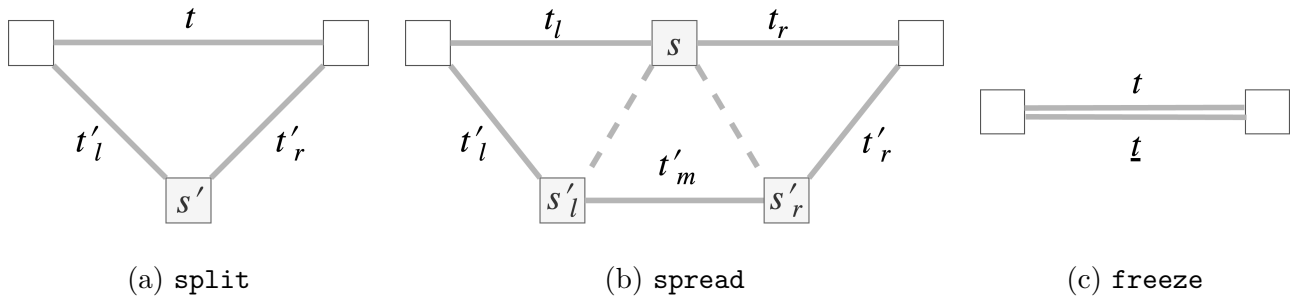


Figure 2.6: The three operations on outer structure.

Figure 2.6 shows a diagrammatic representation of each of these operations. The original slices and transitions are shown at the top of each diagram, while the lower part shows the generated structure after each outer operation is applied. The generation of a piece begins with the empty transition t_0 , followed by a **split** operation, which generates the root slice of the piece.

Figure 2.7 gives an example derivation of the short phrase shown in Figure 2.3. In Figure 2.7a, we can see that 4 outer operations have been applied to generate the surface. The spread operation introduces a passing edge between e and c, shown as the dashed line in Figure 2.7b. Note that the vertically aligned notes in Figure 2.7b correspond to the notes in the surface slices in Figure 2.7a, and the final score shown on the right of Figure 2.7c.

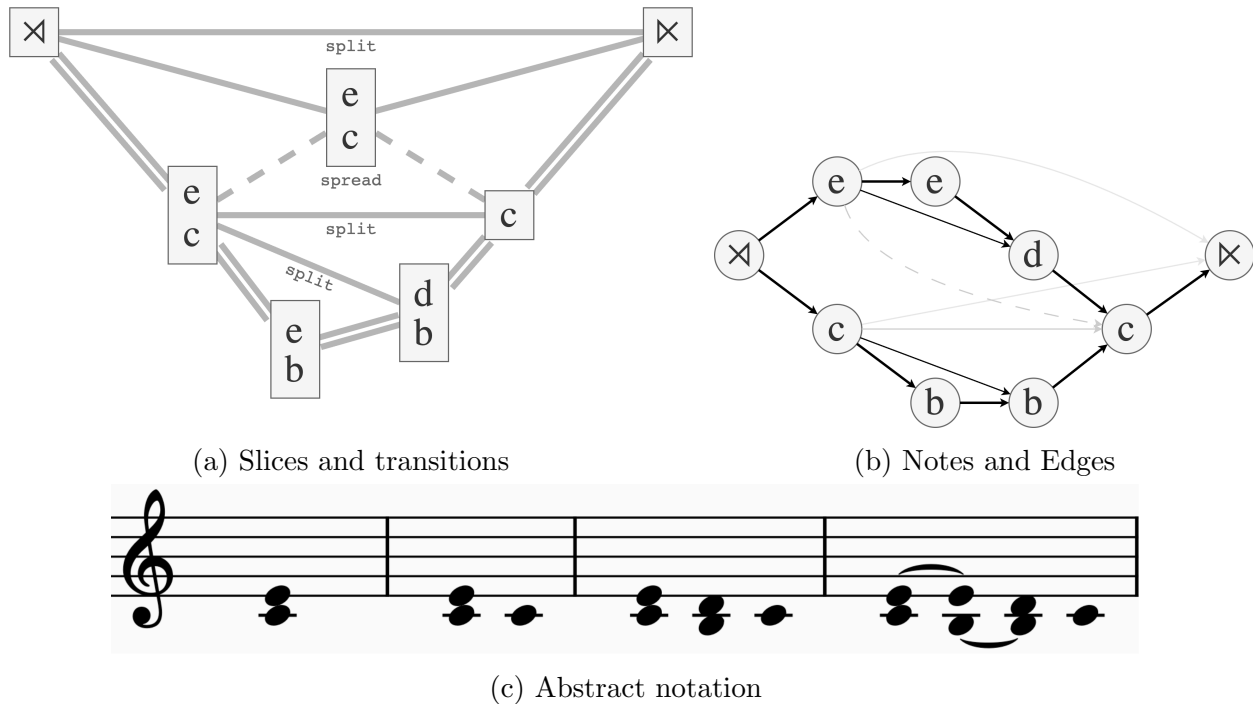


Figure 2.7: An example derivation of the phrase shown in Figure 2.3.

2.5 Parsing

An *ambiguous grammar* is a context-free grammar for which a string can have multiple derivations. In this case, naive parsing results in a *parse forest*, a set of possible derivations of a

particular string. In many cases, such an exhaustive parse is not possible as a result of a combinatorial explosion. With highly ambiguous grammars, such as those describing natural language, or the protovoice model, only a small subset of derivations produced by a grammar are *plausible*. For a parse to be useful, it should yield a plausible derivation.

In the case of the protovoice model, a *plausible* derivation is one which is consistent with a valid musical *interpretation*. Although this is non-trivial, derivation plausibility can be approximately represented by a probability distribution, analogously to probabilistic context free grammars (PCFGs).

Definition 2.5.1 (Derivation plausibility). The plausibility of a derivation is given by the product of the probabilities of each of the production rules used. Given a derivation D consisting of N operations, its plausibility is defined:

$$p(D) = \prod_i^N p(d_i | d_0, \dots, d_{i-1}) \quad (2.21)$$

Where d_0 is the first step in the generative direction. Assuming we can calculate $p(D_i | S)$, we can find the most plausible derivation by taking the maximum likelihood estimate (MLE) of the distribution

$$\hat{D} = \arg \max_D p(D | S) \quad (2.22)$$

This presents **two key problems** to be solved:

- **Calculating $p(d_i | d_{>i})$.** Production probabilities can be viewed as parameters of the model; a common approach with PCFGs is to learn these parameters using machine learning. However, as the protovoice model is not context-free and the volume of data required is not available an alternative method must be found. Furthermore, this would require a top down parse, which is not computationally viable.
- **Combinatorial explosion:** Even if $p(D | S)$ could be solved analytically, the computation would be prohibited by the large branching factor.

2.6 Heuristic Search Algorithms

Heuristic search algorithms are introduced in 1B *Artificial Intelligence*, so it is therefore assumed the reader has knowledge of the heuristic search paradigm.

Parsing can be defined as a search problem wherein the goal is to find a valid derivation of the input according to a grammar. The search space consists of partial derivations, and begins with the input. The space is defined by a **directed acyclic graph** (DAG), with edges corresponding to the inverse of the generative rules. The naive parsing algorithm corresponds to an *exhaustive search* strategy. This would theoretically allow us to find the most plausible derivation (assuming we can rank them in some way), but is computationally infeasible.

Best-first search is a heuristic search algorithm that selects the most *promising* node to expand based on a heuristic evaluation function. In general, the heuristic function $h(n)$ depends

on the description of n , the description of the goal, information gathered by the search up to that point, and most importantly domain specific knowledge [33]. Best-first suffers from memory problems if the number of possible next states are very large.

Beam search is an optimisation of best-first search that serves to reduce its memory requirements by limiting number of best states stored as candidates to expand, dependent on the *beam width*. Beam search is greedy algorithm, so it does not necessarily produce the optimal solution, but trades optimality for improved complexity. Beam search suffers from lack of diversity as the beam tends to get dominated by expansions of a low cost state, causing the search to get stuck in local optima. []

2.7 Starting Point

2.7.1 Relevant courses and experience

Haskell I was introduced to Haskell during an internship during the summer before starting this project (July to August 2022). As a result, I had 2 months of experience with the language beforehand. I chose to use Haskell in order to further familiarise myself with the language, and because of its amenability to parsing algorithms. Furthermore, the functional paradigm and rich type system lends itself to modular software development.

Python I have experience coding in Python from personal projects as well as the 1A *Scientific Computing Practical Course*. This was chosen to made use of powerful existing libraries for handling large datasets and running experiments.

IB Two modules from 1B provide a foundation for this project. *Formal Models of Language* introduces the ideas and terminology used in the protovoice model, and *Artificial Intelligence* provides a background for classical search algorithms as well as some of the probabilistic frameworks used in the project.

2.7.2 Existing codebase

This project was built on a fork of the pre-existing *protovoices-haskell* GitHub repository. This repository contains custom data structures and types describing the protovoice model, allowing interoperability with other projects. I also make use of learned parameters from the implementation of the paper *Bayesian Model of Extended Chord Profiles* [9], as learning these parameters would be beyond the scope of this project.

2.8 Requirements Analysis

The Success Criteria are given in the Project Proposal. During the preparation phase, the Success Criteria were refined in light of increased clarity from reading the literature related to

the project.

This project will be deemed a success given it achieves:

- An implementation of a parser for the protovoice model which outputs possible derivations of a piece of music.
- A harmony module that can take the output of the search algorithm and quantitatively evaluate its accuracy against the ground truth annotations.
- A baseline search algorithm that outputs a reduction of a piece using the parser in order to predict chord labels.
- Extension: Develop and evaluate one or more search algorithms that use additional heuristics to inform the search.

Table 2.1: Project Deliverables

ID	Deliverable	Priority	Risk
core1	Harmony Model	High	Low
core2	End to End Pipeline	High	Medium
core3	Protovoice Parser	High	High
base1	Baseline Reduction	High	Low
base2	Random Search	High	Low
ext1	Heuristic Design	Medium	High
ext2	Greedy Search	Medium	Medium
ext3	Heuristic Search	Low	Very High

There is a high general risk in this project due to the fact that the protovoice model is a result of an unpublished thesis, and thus may have flaws that have not been discovered. Table 2.1 shows a list of project deliverables with associated priorities and risk, denoted qualitatively. The highest risk task is designing and implementing the protovoice parser, but as there exists an implementation, it should be possible. The heuristic design task is expansive, requiring creativity, research and iterative development, hence the risk is high. The baseline inference deliverable comprises a standard method for ACE, so there is minimal risk. The design and optimisation of search algorithms is also a substantial task, which has the risk of sinking a practically infinite amount of time.

2.9 Software Engineering Techniques

2.9.1 Development model

Based on the risk analysis (Table 2.1), a plan was created describing which modules to implement in which order, with a list of milestones on a 2 week basis. Notion was used to maintain a list of core tasks and corresponding subtasks with associated priorities, facilitating the selection of the next tasks to work on. The development strategy chosen drew from the Agile

methodology, involving two-week long sprints with regular re-evaluations of the plan informed by experimental data and testing. GitHub’s continuous integration features were used to run a test suite on the repository after every commit.

2.9.2 Languages, libraries and tools

Table 2.2 shows a justified list of the key languages, libraries and tools used in the project. The licensing agreements for all the tools used in the project were determined and analysed. For the most part, these are all permissive licenses, guaranteeing freedom to use, modify and redistribute as well as permitting proprietary derivative works.

2.9.3 Hardware, version control and backup

The code was developed using Vim for Haskell and Visual Studio Code for Python notebook development, on my personal laptop (16’ MacBook Pro 2022, M1 Max, 32GB). Experiments were run on a server provided by the EPFL Digital Cognitive Musicology Lab (Dell PowerEdge R740XD Server, 2x Xeon Gold 625R, 768GB), using Jupyter notebooks to conduct the evaluation. I used GitHub for all my notes, development and dissertation writing. Finally, this dissertation was written in Vim with VimTeX.

Table 2.2: Languages, libraries and tools

Tool	Purpose	Justification	License
<i>Languages</i>			
Haskell	Main language used for the core, baseline and extension implementations	Protovoice model implementation is in Haskell. Functional and amenable to parser development.	GHCL
Python	Secondary language for experiments and analysis	Powerful library ecosystem for running experiments and creating plots	PSFL
<i>Libraries</i>			
Musicology	Haskell Library with data-types for pitches	Contains a robust implementation of spelled pitch classes, which would be tedious to reimplement.	BSD-3.0
Timeit	Lighweight wrapper to show the used CPU time of a monadic computation	This is used to time the runtime of the algorithms as part of analysis	BSD-3.0
Dimcat	Python library: Digital Musicology Corpus Analysis Toolkit	This library was written to work with the datasets used in this project	GPL-3.0
Numpy	Python library used for preprocessing and analysis	Powerful standard library that is used in conjunction with Seaborn to run analysis and visualise data	BSD-3.0
Pandas	Python library for preprocessing and analysis	This is a standard library for data manipulation and processing	BSD-3.0
Seaborn	Python data visualisation library used for analysis	Creates high quality graphs and charts	BSD-3.0
<i>Tools</i>			
Docker	Containerised software service used to run repeatable experiments	The use of many libraries creates a web of dependencies to be resolved. Ensures the code will last and can be executed on different devices	Free/Paid
Git	Version Control, Continuous Integration	Provides natural backups and allows for reverts to previous commits if necessary	GPL-3.0
GitHub	Hosting source code	Free, reliable hosting	GPL-3.0
GHC	Compiling and profiling.	This is the standard Haskell compiler.	BSD-3.0
Stack	Haskell building and testing	Creates reliable builds, and includes a powerful testing framework.	BSD-3.0
Undotree	Vim Plugin: stores all past actions as a tree	Solves the problem of linear undo history being lost. Protects code between commits.	BSD-3.0
MuseScore	Music notation software	The raw inputs are in the MusicXML format, which is used by MuseScore 3	GPL-3.0
PAT	Protovoice Annotation Tool, Used to view protovoice derivations on a web browser	The protovoice derivations are huge and very complex, so it's vital to have a viewing tool for use in analysis and iterative development	GPL-3.0

Implementation

3.1 ProtoVoice Harmony Model

The result of this project is the **ProtoVoice Harmony Model**, a **novel** inference system that uses the protovoice model to infer harmony from free polyphony.

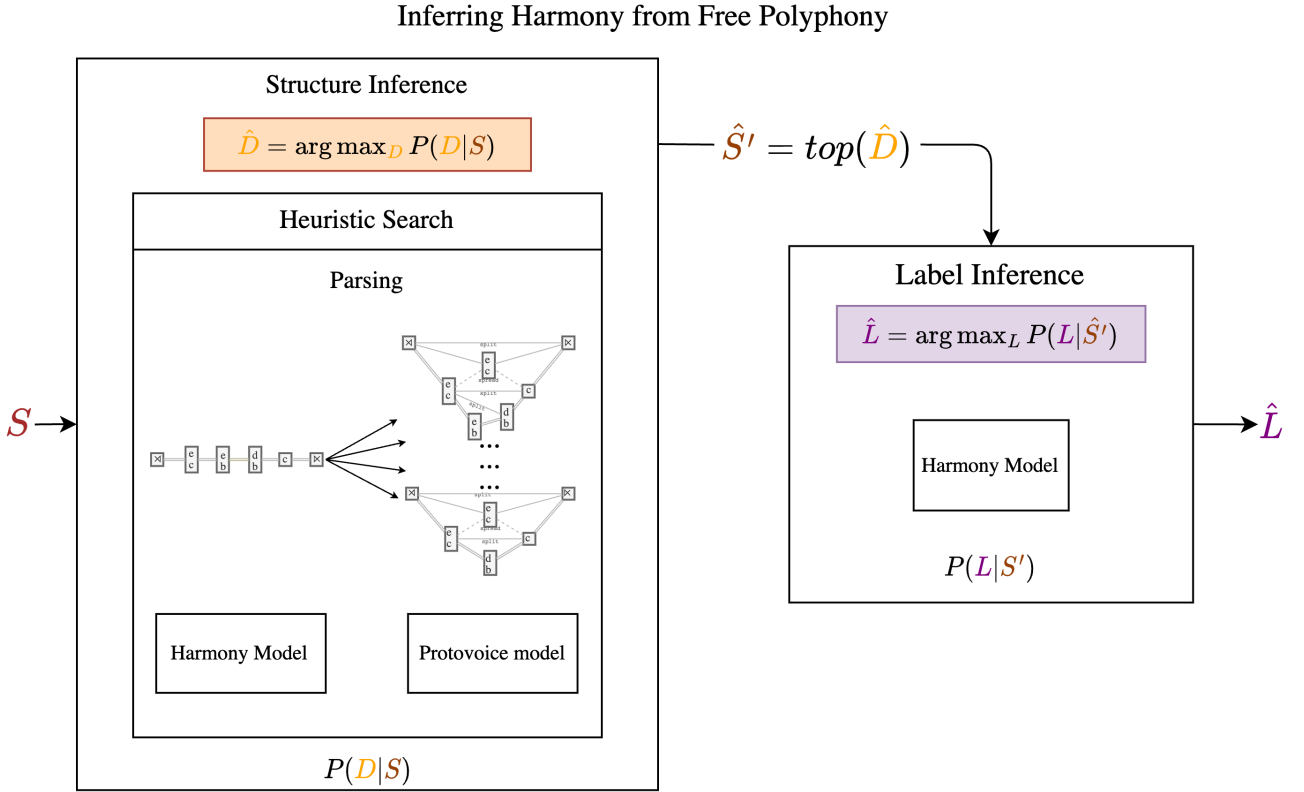


Figure 3.1: Overview of the inference process.

The PVHM is described by the following equation:

$$\hat{L} = \arg \max_L \left(P(S|L, \hat{S}', \hat{D}) \cdot P(L) \right) \quad (3.1)$$

Where \hat{D} denotes the inferred *protovoice derivation* of the score S , obtained using a modified beam search algorithm to parse S , guided by a heuristic \mathcal{H}_r . \hat{S} is the associated reduction of the surface, given by $\text{top}(\hat{D})$.

In the context of a protovoice derivation, the *surface* refers to the input, the original score to be parsed, and the result of the parse is the *top* of the derivation. The core of the PVHM is finding the most likely reduction $\mathcal{R}(S)$ of the surface S by inferring the most likely protovoice structure \hat{D} . The baseline algorithm achieves through a random search through the space of valid partial derivations. The extension algorithms use an informed heuristic H_0 to guess $P(D|S)$ in order to guide the search, and explore the trade-offs between exploration and exploitation with different search algorithms. Given the inferred derivation, the corresponding reduced surface is extracted with $top : \mathbf{D} \rightarrow \mathbf{S}'$. Finally, the reduced surface $\mathcal{R}(S)$ is used to solve $\hat{L} = \arg \max_L P(L|S, \mathcal{R}(S)) \cdot P(L)$, finding the most likely sequence of chord labels \hat{L} .

Development plan

This system is modular by design, leading to a natural progression of complexity throughout the implementation; the baseline models provide a subset of the functionality required for the full system. For example, the harmony model is a prerequisite for the first baseline, but is also used in all subsequent models.¹

- The Baseline Reduction computes a simplified $\mathcal{R}(S)$ using random sampling.
- The Random Search uses $\mathcal{R}(S) = top(D_0)$, where D_0 is the derivation found by randomly choosing each parse operation.
- The first extension involves developing an informed heuristic \mathcal{H}_0
- The Greedy Search uses heuristic \mathcal{H}_0 to approximate $\arg \max_D P(D|S)$.
- The Beam Search uses \mathcal{H}_0 to implement a search parameterised by the beam width β
- The Dual Beam Search implements a custom search algorithm using \mathcal{H}_0 , parameterised by α and β , which control operation dependent beam widths.

3.2 End to End Pipeline

The first step was to implement a full end-to-end pipeline from piece to chord label predictions. This was achieved by writing `preprocess.py`, `runExperiment.py` and `analysis.ipynb`. The repository was then containerised using *Docker*, so that experiments could be run and analysed using a remote server with automated dependency resolution. Experiments were executed on a server, initialised via *ssh*, and *tmux* was used to maintain the server environment across connections, with automated scripts to pull the latest changes from GitHub. A *jupyter-lab* environment was set up on the server, and scripts written to fetch the results of the latest experiments and produce plots and tables, allowing analysis to be conducted through a browser.

¹This could be complemented by a diagram showing the progression

Dataset selection

The datasets used contain bodies of work by a single composer as shown below, these were selected to exhibit a range of different styles. Each dataset contains an entire set work by a composer consisting of between 15 – 50 annotated scores, each with varying lengths. Each dataset was split into training, validation, and test sets, using a 60:20:20 split with *stratified sampling* in order to maintain a balanced representation of the different composers. In order to choose hyperparameters, the training and validation sets from each of the five datasets were combined into a single training pool, for use in cross-validation. For the final evaluation, the performance of all developed algorithms will be evaluated on the hold-out test set (20%). This is to provide an unbiased estimate of each algorithm’s performance on **unseen** data.

3.3 Protovoice Parser

The Protovoice Parser is a bottom-up parser for the protovoice model that enumerates all the valid protovoice derivations that produce a given score $S \rightarrow [D]$. Recall that within the context of the protovoice model, the original score is called the *surface*, and the surface is reduced by parsing, removing non-chord-tones. A derivation is described as a sequence of rule applications in *left-most derivation order*, applied to a fully reduced surface, S' .

Formally, a derivation D is defined as a pair (top, ops) , where the *surface*, S , is derived by starting with the fully reduced top and applying each operation in ops in order.

It is informative to consider the *generation order* and *parse order* of a protovoice derivation. In generation order, we begin with the reduced surface top , consisting of only chord-tones, with a pointer at the left-most transition and apply each **split** or **spread** operation to the two left-most non-terminal transitions in the path graph, generating new slices and transitions. The **freeze** operation marks a transition as terminal thus shifts the pointer to the right. Operations are applied until the pointer is right-most, and all transitions are terminal, resulting in the fully elaborated *surface* and a derivation $D = (top, ops)$, where ops is the sequence of operations applied, $ops = [d_N \dots d_0]$, resulting in the *surface* $= d_0 \circ \dots \circ d_N(top)$.

Parsing is the inverse of generation: the parse begins with the elaborated surface S consisting of only frozen transitions with a pointer at the right-most frozen transition. During the parse, we can either **unfreeze** the transition at the pointer, shifting the pointer to the left, or apply a **unsplit** or **unspread** reduction to the two transitions to the right of the pointer. Once the pointer is left-most, all transitions are unfrozen and there is a *single reduced slice* for each chord segment at the top of the derivation, the resulting derivation is $D = (top, ops)$, where the reduced surface is $S' = top = (d_N^{-1} \circ \dots \circ d_0^{-1})(surface)$.

This is a form of *shift-reduce* parser, making one pass across the surface, right to left. The state of the parse is represented by **parseState** $= (top, ops)$ and implements **expand**: **parseState** \rightarrow [**parseState**] such that we can use search algorithms to conduct the parse.

In the initial parse state, **psFrozen**, the surface is represented as a **path** from the end to the beginning of the piece, with the right-most transition at its head. The **start**(\times) and **stop**(\times) symbols are not explicitly stored in the path, but are implied.

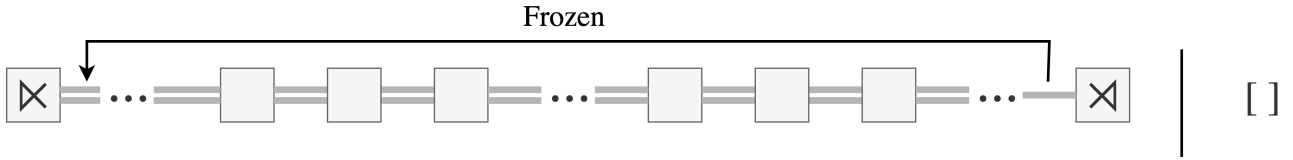


Figure 3.2: pSFrozen

All transitions are initialised as frozen, indicated by the double lines in Figure 3.2. The **path** begins on the right. In this initial state the only option is to unfreeze the rightmost transition, moving to the **pSSemiOpen** state.

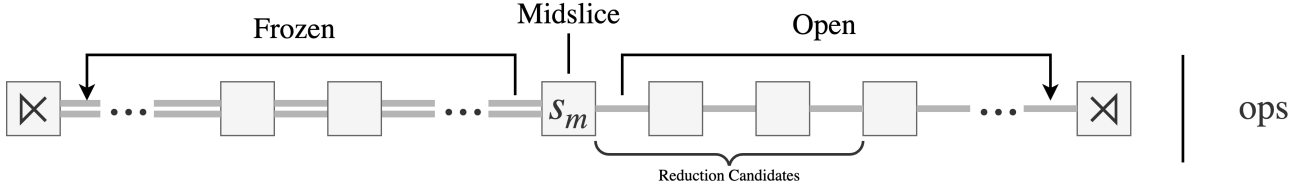


Figure 3.3: pSSemiOpen

The **pSSemiOpen** represents the majority of the parse. The pointer is represented by **midSlice**, and points to the rightmost frozen transition. The **open** path contains all unfrozen transitions (denoted by a single line) from the right of the pointer to the end of the piece, and the **frozen** path contains all frozen transitions from the pointer to the start of the piece.

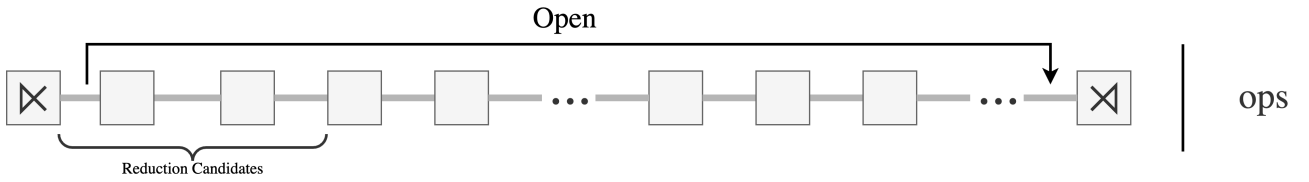


Figure 3.4: pSOpen

In the final state, **pSOpen**, the pointer is at the beginning of the piece, and comprises a single path **open** containing only unfrozen transitions from the beginning to the end of the piece.

Each state also stores a list of operations, *ops* that when applied to the current reduction, results in the original surface. This is empty at the beginning of the parse; for each reduction applied, the corresponding generative operation *op* is consed to the list: $ops' := op : ops$.

This is not the only way to parse according to the protovoice model; the protovoice reduction rules can be applied to any pair of open transitions. Allowing reductions at any point along the surface could allow for different parse strategies to be employed that may have advantages. I choose not to pursue this due to the associated combinatorial explosion, and the benefits of being able to represent a protovoice derivation as a sequence of left-most reductions.

Parsing Operations

The **protoVoiceEvaluator** [10] provides methods that are used to enumerate the possible operations for each reduction type ². It also includes an implementation of a greedy parser, from which ideas were adapted and expanded upon to create the protovoice parser. This allows

the parser to consider only the outer structure of slices and transitions while parsing without exposing the internal notes and edges.

$$\begin{aligned}
 \text{evalUnfreeze} : & \quad t \rightarrow [(t', op)] \\
 \text{evalUnsplit} : & \quad (t_l, s_l, t_m) \rightarrow [(t_{top}, op)] \\
 \text{unSpreadLeft} : & \quad (s_l, t_l) \rightarrow s_{top} \rightarrow [t_{topl}] \\
 \text{unSpreadRight} : & \quad (s_r, t_r) \rightarrow s_{top} \rightarrow [t_{topr}] \\
 \text{unSpreadMiddle} : & \quad (s_l, t_{top}, s_r) \rightarrow \text{Maybe } (s_{top}, op)
 \end{aligned} \tag{3.2}$$

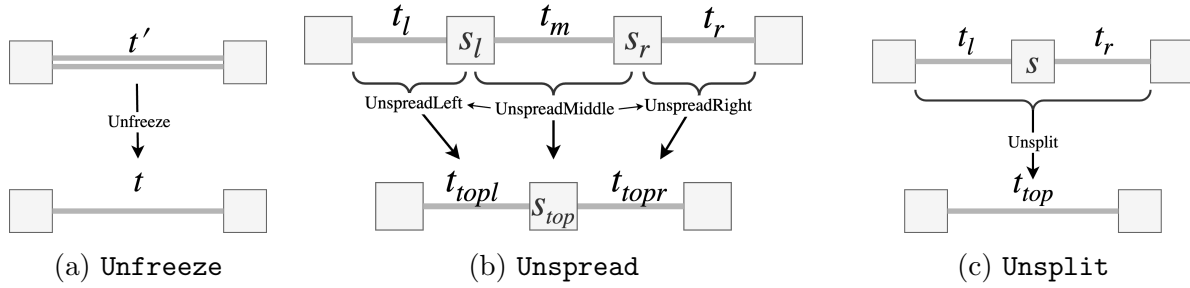


Figure 3.5: Reduction operations.

Boundary handling

Recall that the goal of the parse is to reduce the original surface such that there is one slice per segment, containing only chord-tones. In order to achieve this, constraints are imposed on the reduction operations dependent on adjacent segment boundaries.

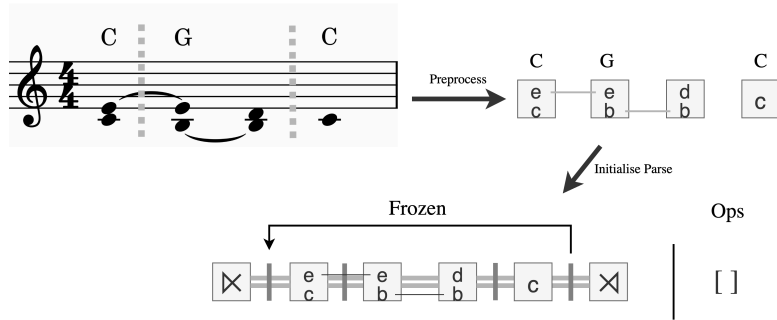


Figure 3.6: Parse Initialisation

Each transition has a boolean *boundary* value which indicates if the transition is across a segment boundary, denoted by the dashed vertical line in Figure 3.6. Let $B_t : t \rightarrow \{\mathbf{True}, \mathbf{False}\}$, such that B_f, B_l, B_m , and B_r denote the boundary values of t_f, t_l, t_m and t_r respectively. As an **unspread** operation ($t_l s_l t_m s_r t_r \rightarrow t_{topl} s_{top} t_{topr}$) merges two slices s_l and s_r , removing t_m , the constraint $\neg B_m$ is imposed to prevent two segments from merging. Similarly, an **unsplit** operation (eg. $t_l s_l t_m \rightarrow t_{top}$) combines two transitions, thus we impose the constraint $\neg(B_l \wedge B_m)$.

²Note that the notation used here for functions combines type definitions with variable names. For example, the function **evalUnsplit** has type $\text{evalUnsplit} :: (\text{transition}, \text{slice}, \text{transition}) \rightarrow [(\text{transition}, \text{operation})]$, but type names (e.g. *transition*) have been substituted by variable names (e.g. t_l, t_r, t_{top}) for expository purposes.

Finally, the `unfreeze` operation shifts the context to the left. If a boundary transition is unfrozen, no more reduction operations can be applied to its right, hence it is a necessary condition that the segment has been fully reduced. Thus the constraint imposed on an `unfreeze` operation is $\neg B_f \vee (B_l \wedge B_m \wedge B_r)$. Note that there are some configurations that are *unreachable*, such as $B_f \wedge (B_l \wedge B_m \wedge \neg B_r)$, thus the `unfreeze` constraint could be simplified to $\neg B_f \vee (B_l \wedge B_m)$, but I have chosen not to as it would obfuscate the semantics, and the cost of an additional logic check is negligible.

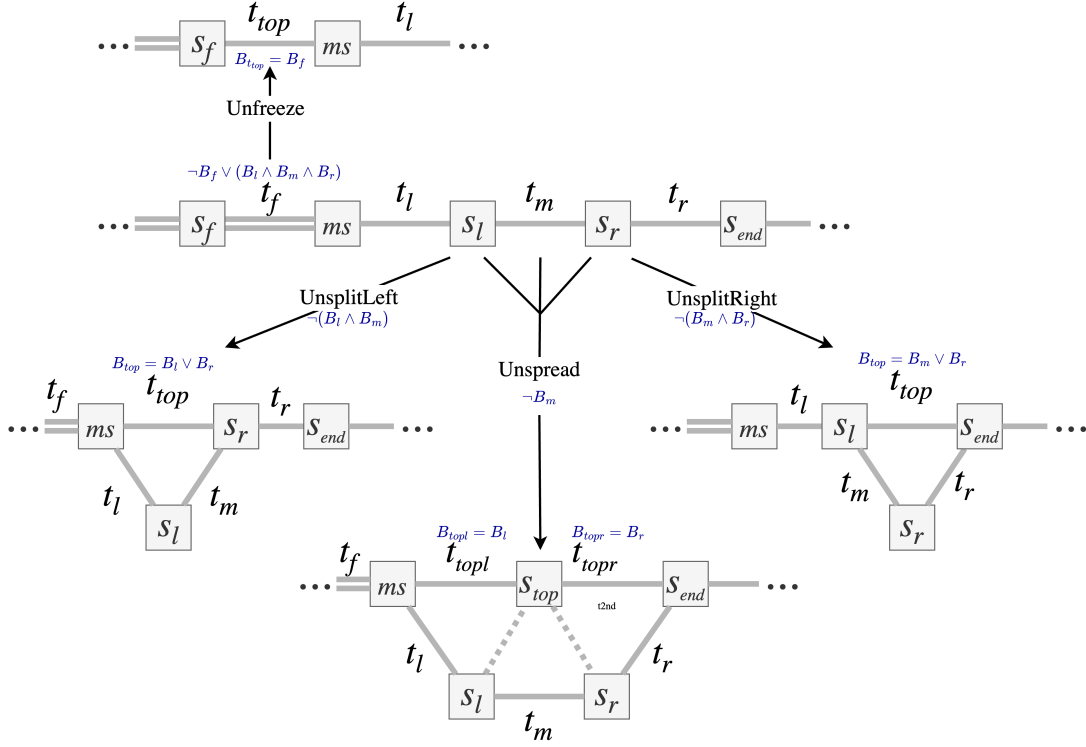


Figure 3.7: Boundary reduction conditions and propagation

Figure 3.7 provides an illustration of these constraints. The boundary must also be propagated to parent transitions following a reduction, defined by the union of each child transitions' boundary values for a given reduction.

The `evalUnfreeze` and `evalUnsplit` functions allow the associated generative operation, `freeze` and `split` to be determined trivially, but the possible `spread` operations need to be derived using `unSpreadLeft`, `unSpreadRight`, and `unSpreadMiddle`. This is achieved using the *list monad*, wherein multiple branches of sequential computation appear to be executed simultaneously, returning the results of all possible computation branches in a list.

Algorithm 1 Enumerate unspread reductions

```

1: function COLLECTUNSPREADS( $s_l$   $t_l$   $s_m$   $t_r$   $s_r$ )
2:   using ListMonad do
3:     ( $s_{top}, op$ )  $\leftarrow$  maybeToList $ evalUnspreadMiddle( $s_l, t_m, s_r$ )
4:      $t_{topl}$   $\leftarrow$  evalUnspreadLeft ( $s_l, t_l$ )  $s_{top}$ 
5:      $t_{topr}$   $\leftarrow$  evalUnspreadRight ( $s_r, t_r$ )  $s_{top}$ 
6:     return ( $l_{top}, s_{top}, r_{top}$ )
7: end function

```

Finally, given that an **unspread** reduction involves two transitions, and an **unsplit** reduction involves only one, it is important to ensure the parse pointer is such that there are two open transitions to its right wherever possible. This avoids the case of applying alternating **unsplit** and **unfreeze** reductions, as there is only zero or one open transition within the current segment, and **unspread** reductions require two. This is not important for the sake of a parse, which enumerates all derivations, but is important later on when designing heuristic search algorithms.³

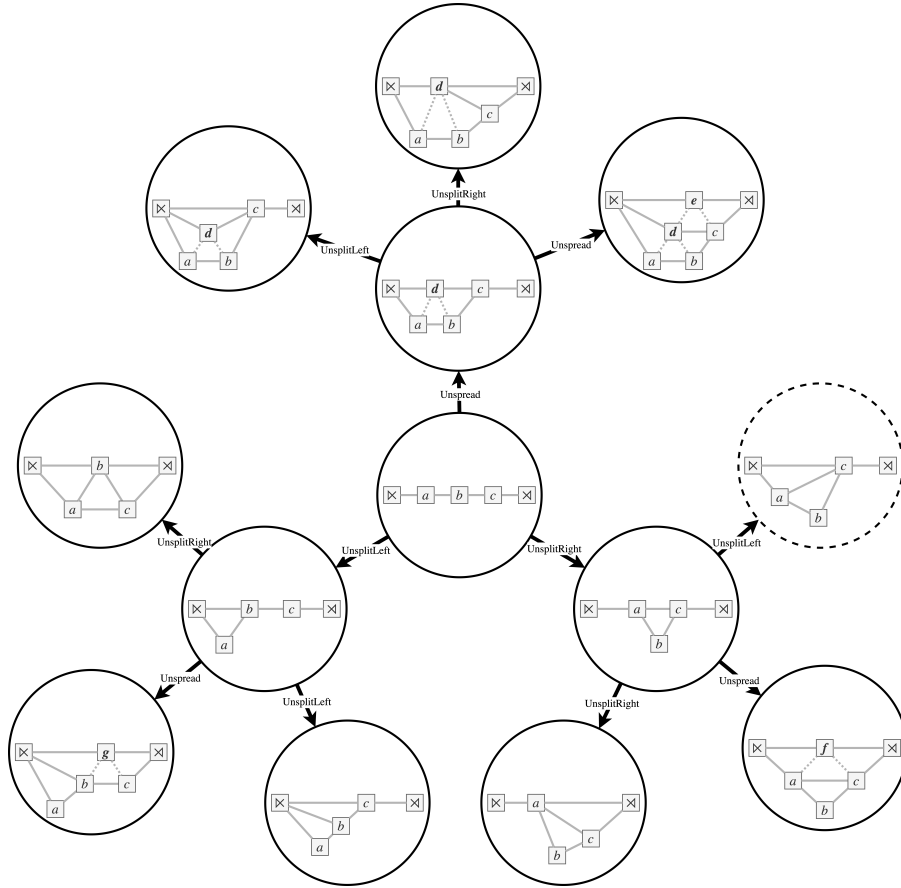


Figure 3.8: Parse state transitions for an open segment with three slices.

Figure 3.8 shows the 9 possible outer configurations of slices after two reduction steps, shown for a single segment reduction consisting of three slices and four transitions. One of these configurations, **UnsplitLeft** \circ **UnsplitRight**, is redundant as it results in the same configuration as **UnsplitLeft** \circ **UnsplitLeft**, so it is prohibited.

Note that while the number of possible outer structure reductions grows exponentially, each reduction consists of a set of inner operations on the notes and edges contained within each slice and transition, which also grows exponentially with the number of possible configurations of inner operations in each slice, an exponential function of the number of pairs of neighbor/repetitions. The number of pairs of neighbors/repetitions is also exponential with respect to the number of notes in a slice. Thus as the parsing runtime grows exponentially with respect to three input variables, it has **3-EXP**, or *triply exponential* complexity. These variables, however, are not entirely independent, and convention means that some of these dimensions of complexity are limited for most pieces. It is still the case that the combinatorial blowup in the parse dimension means its impossible to enumerate all parses for *any* piece, and there are cases

³I guess I'll move this to the heuristic design section?

when storing the each the of different **unsplit** operations at a single step in memory is not possible.

The parse states define a **directed acyclic graph** where the parseStates are nodes and the parse operations are edges. Traversing this graph results in either reach a valid derivation or dead state.

3.4 Harmony Model

The next step was to devise and implement a harmony model $P(L, S')$, describing the relationship between chord-tones and chord labels. This provides a framework to evaluate all developed algorithms. Log probabilities are used throughout.

Recall that input is a sequence of slices, groups of notes, that are only chord-tones. Each note has an associated pitch $p \in \mathcal{P}$. A chord label is defined as a tuple of *root note* pitch and chord-quality $\mathcal{L} = (\mathcal{P}, \mathcal{C})$, where the chord qualities are taken from the DCML chord label annotation standard as used in the datasets: $\mathcal{C} = \{M, m, Mm7, om, o7, mm7, \%7, MM7, +, Ger, It, Fr, mM7, +7\}$.

Given a slice $s = \{p_1^{m(p_1)}, \dots, p_n^{m(p_n)}\}$ where $p_i \in \mathcal{P}$ and $n = |\mathcal{P}|$, the goal is to find the most plausible chord label:

$$\hat{l} = \arg \max_l P(l|s) \quad (3.3)$$

We choose to find $P(l|x)$ rather than $P(s|l)$ as it provides a *confidence* in the prediction, given by $P(\hat{l}|s)$. The generation of a slice s of size n from a label l is modelled as follows:

$$s \sim \text{Multinomial}(n, |\mathcal{P}|, \Phi_l) \quad (3.4)$$

$$\text{where } \phi_l^{(p)} = P(\text{Pitch} = p | \text{Label} = l)$$

The slice vector is defined $\mathbf{v}(s) = [m(p_1) \ m(p_2) \ \dots \ m(p_k)]$ where $k = |\mathcal{P}|$. Then $P(l|s)$ is given by: ⁴

$$\begin{aligned} P(l|s) &= P(s|l) \cdot P(l) \\ &= f(\mathbf{v}(s), \phi_l) \cdot P(l) \end{aligned} \quad (3.5)$$

where f is the multinomial probability density function

$$f(x_1, \dots, x_k; p_1, \dots, p_k) = \frac{\Gamma\left(\sum_i x_i + 1\right)}{\prod_i \Gamma(x_i + 1)} \prod_{i=1}^k p_i^{x_i}$$

⁴Struggling with naming here. Using P for pitches conflicts the P for probability.

The log probability $\log P(l|s)$ is given by

$$\begin{aligned}
 \log P(l|s) &= \log P(s|l) + \log P(l) \\
 &= \dots \\
 &= \log \left(\Gamma \left(\sum_i x_i + 1 \right) \right) + \sum (x_i \log p_i^{x_i} - \log \Gamma(x_i + 1)) + \log P(l)
 \end{aligned} \tag{3.6}$$

So to find the most likely label l given slice s take:

$$\hat{l} = \arg \max_l \left(\log \left(\Gamma \left(\sum_i x_i + 1 \right) \right) + \sum (x_i \log p_i - \log \Gamma(x_i + 1)) + \log P(l) \right) \tag{3.7}$$

The full derivation has been omitted for the sake of brevity. Given that $|\mathcal{L}| = |\mathcal{P}| \times |\mathcal{C}| = 29 \times 12 = 348$, the equation can be solved directly.

3.4.1 Baseline algorithms

The goal of protovoice parser is to find a derivation $D = (top, ops)$ which results in a top sequence of slices top consisting of only chord-tones.

The simplest possible baseline is one that generates slices randomly for each segment, without considering the piece at all. This is the **RandomSample** algorithm. While this is naive, it provides a useful reference for evaluation metrics.

As a second baseline, consider the simplest form of reduction, removing notes by randomly guessing if they are chord tones or non-chord-tones. This is the **RandomReduction** algorithm. It is equivalent to combine all the slices in a given segment and take a random sample without replacement of n notes from the combined slice. The number of notes is sampled from a Poisson distribution, $n \sim \text{Poisson}(\lambda)$, where λ is learned from data.

Next, the core algorithm is the **RandomWalk** algorithm using the protovoice parser. This algorithm takes a random walk through the parse states graph. The protovoice model does not make judgements based on harmony or note function explicitly, but by comparing this algorithm to the other baselines, it can be determined whether the constraints enforced by using a protovoice derivation has any effect on the quality of the reduction.

The `ParseAlgo` type-class provides an interface for running these algorithms. The result is wrapped in a `Maybe` as it is possible for the search to get stuck, and is in the `IO` monad to allow non-determinism.

```

class (Show algo) => ParseAlgo algo where
  runParse :: algo -> AlgoInput -> IO (Maybe AlgoResult)

```

Listing 3.1: Algorithm type-class

3.5 Extension Implementation

Let us take a step back and consider what has achieved up to this point. Recall the goal is to find $\arg \max_L P(L|S)$.

By finding a protovoice derivation of the piece, $D = (top, ops)$ and corresponding reduced surface $top \in \mathcal{R}(S)$, the harmony model $P(S, L)$ is used to infer the chord labels. These additional variables are introduced formally by marginalising:

$$\begin{aligned} \hat{L} &= \arg \max_L P(L|S) \\ &= \arg \max_L P(S|L) \cdot P(L) \\ &= \arg \max_L \sum_{D, \mathcal{R}(S)} (P(S|D) \cdot P(D|\mathcal{R}(S)) \cdot P(\mathcal{R}(S)|L)) \cdot P(L) \end{aligned} \quad (3.8)$$

It is not possible to enumerate D in the general case, so the optimisation is approximated by finding the most plausible derivation $\hat{D} = (top, ops)$.

$$\hat{L} = \arg \max_L P(L|top) \quad (3.9)$$

So far, the protovoice parser has been used to find an arbitrary derivation, but this derivation may not be plausible as it assumes random decisions made at each step. To find the most likely derivation $\hat{D} = (top, ops)$, consider abstract relation between derivations and surfaces described by $P(D, S) = P(D) \cdot P(S|D)$. The surface is given so $P(S|D) = 1$, so it is necessary to maximise $P(D)$.

The probability of the derivation $P(D)$ is defined $\prod_i P(d_i|d_0, \dots, d_{i-1})$ where d_0 is the first step in the generative direction and d_N is the last step, or the first step in the parsing direction. The plausibility of the derivation given the surface S from the generative direction is thus defined $P(D|S) = \prod_i P(d_i|d_{>i}, S)$. This cannot be calculated as it requires enumerating all derivations.

An approximation is made that is akin to the Markov assumption, estimating the plausibility of the reduction so far through a first order approximation of each term. For each step d_k , the plausibility of the reduction up to that point is $P(d_k \dots d_N|S) = P(d_{>k+1}|S) \cdot P(d_k|d_{>k+1}, S)$.

The plausibility of each derivation step is approximated by considering only local information: the current top surface and the reduction operation applied. The heuristic is a function $\mathcal{H}_0 : ((\text{ParseState}, \text{Cost}), \text{Op}) \rightarrow \text{Cost}$. Recall that $\text{ParseState} = (top, ops)$, where $top = (d_0 \circ \dots \circ d_N)^{-1}(surface)$.

$$H_0((d_{>i}, c), d_i) = c \cdot P() \quad (3.10)$$

3.5.1 Extending the Harmony Model

In order to design the heuristic \mathcal{H}_0 , the harmony model needs to be extended to incorporate both chord-tones and non-chord-tones that will be present during the search.

Predicting chordness

The goal of the search is to reduce the piece in to a sequence of slices containing only chord-tones. In order to do this, it is helpful to have a heuristic that indicates how close we are to this state in order to guide the search. *Chordness* is defined as the likelihood that a given slice could have been generated by the most likely label \hat{l} .

$$\text{Chordness}(s) = P(s|\hat{l})P(\hat{l}) \quad (3.11)$$

where $\hat{l} = \arg \max_l P(l|s)$ as above.

An alternative approach would be to marginalise over all chords find the "expected" chordness.⁵

Ornamentation likelihood

The original model $P(\mathcal{R}(S), L)$ is extended to account for both chord-tones and non-chord-tones, $P(S, L)$, by introducing a *mixture model*. Notes that are not chord-tones are referred to as *ornaments* as music theory states that they often have a *function* with respect to an adjacent chord-tone, ornamenting that chord-tone. This phenomena is captured in the protovoice model through *neighbour* notes. Here, a neighbour note and ornamentation are considered synonymous.

It is assumed that each note in a chord l is generated as either a chord-tone or ornament according to a Bernoulli distribution $t \sim \text{Bernoulli}(\theta_l)$. Then each note is generated from a categorical distribution $n \sim \text{Categorical}(n, \Phi_l^{(p)})$ where Φ is either for the chord tone distribution or ornament distribution depending on t ⁶

3.5.2 Heuristic Design

In order to calculate the cost of a single reduction, we consider the plausibility of the corresponding operation being applied in the generative direction. Recall that the protovoice model comprises an inner and outer structure. The outer structure provides an abstraction of the inner structure, allowing the parsing algorithm to be designed using operations applied to slices and transitions. In order to evaluate the plausibility of operations, the details of the operation needs to be considered.

Recall that the core of the protovoice model is a set of generative operations on notes.

Any operation consists of a set of regular edges, which are either neighbour

Step 1: Design heuristic to be as accurate as possible. I.e the extreme is to consider every possible parse, but for a single piece there can be over $10^{10^{10^{10^{10}}}}$ different parses. We consider 1 step at a time at first - this still results in needing to choose an operation out of upwards of 30,000,000 options for just a single step.

⁵Need to add a clever reason why I didn't choose to do so, or drop this sentence

⁶todo formalise

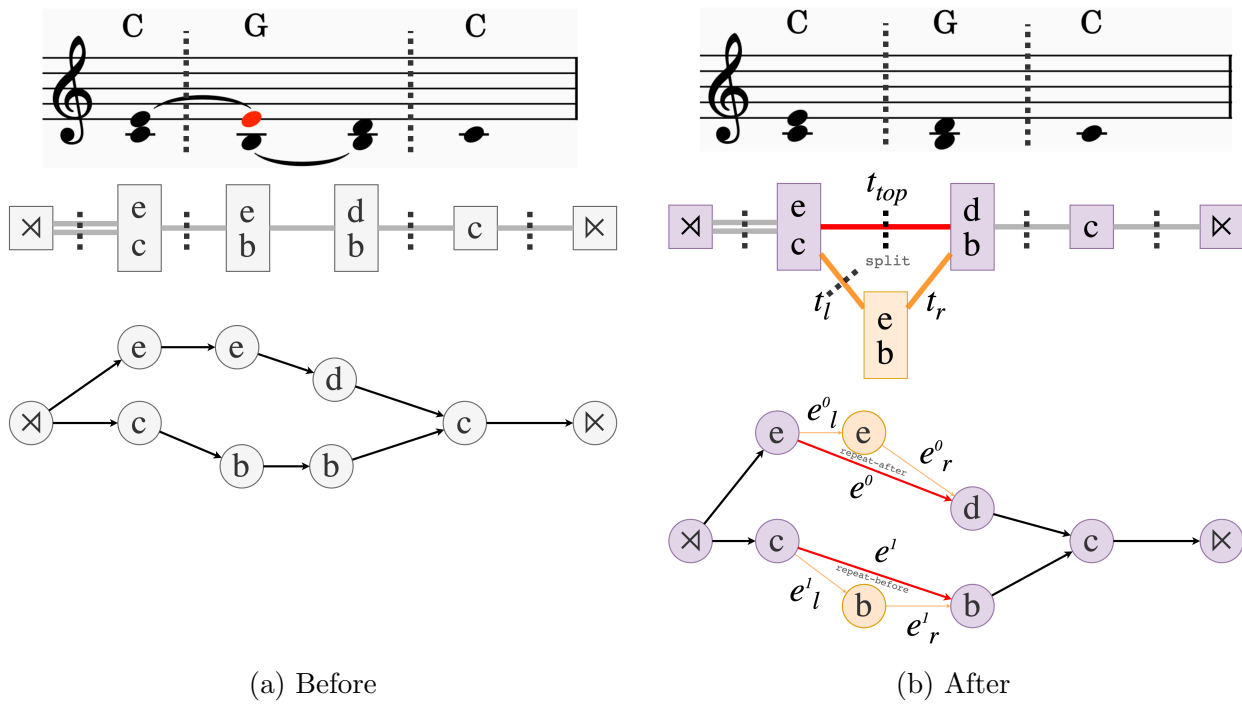


Figure 3.9: Single Reduction Step

First the full piece heuristic parse

The goal of the heuristic is to find the probability of a reduction.

Here the `ProbVectorSingle` function a vector for the.

Algorithm 2 Heuristic function \mathcal{H}_0

```

1: function EVALUATEOPERATION( $State, Op$ )
2:   if  $Op$  is a Split then
3:     return EVALUATESPLIT( $State, Op$ )
4:   else if  $Op$  is a Spread then
5:     return EVALUATESPREAD( $State, Op$ )
6:   else if  $Op$  is a Freeze then
7:     return LOG(1) ▷ Likelihood of 1
8:   end if
9: end function
10:
11: function EVALUATESPREAD( $State, Op$ )
12:    $s_l, s_r \leftarrow \text{PARENTSLICES}(State, Op)$ 
13:    $E \leftarrow \{\text{REGULAREDGES}(Op) \cup \text{PASSINGEDGES}(Op)\}$ 
14:    $LogLikelihoods \leftarrow \{\text{EVALUATEEDGE}(e, s_l, s_r) : e \in E\}$ 
15:   return AVERAGE( $LogLikelihoods$ )
16: end function
17:
18: function EVALUATESPLIT( $State, Op$ ) ▷ How much detail? keepEdges?
19:    $s_l, s_r \leftarrow \text{PARENTSLICES}(State, Op)$ 
20:    $e_{reg} \leftarrow \text{REGULAREDGES}(Op)$ 
21:    $e_{pass} \leftarrow \text{PASSINGEDGES}(Op)$ 
22:    $e_{single} \leftarrow \text{LEFTEGES}(Op) \cup \text{RIGHTEDGES}(Op)$ 
23:    $e_{right} \leftarrow \text{REGULAREDGES}(Op)$ 
24:    $LogLikelihoods \leftarrow \{\text{EVALUATEEDGE}(e, s_l, s_r) : e \in \{e_{reg} \cup e_{pass}\}\}$ 
25:   return AVERAGE( $LogLikelihoods$ )
26: end function
27:
28: function EVALUATEEDGE( $e, s_l, s_r$ )
29:    $n \leftarrow$  child note introduced in  $e$ 
30:   if  $e$  elaborates two parents,  $p_l$  and  $p_r$  then
31:      $\vec{\phi} \leftarrow \text{PROBVECTORDOUBLE}(n, p_l, p_r, s_l, s_r)$ 
32:   else if  $e$  elaborates single left parent  $p$  then
33:      $\vec{\phi} \leftarrow \text{PROBVECTORSINGLE}(n, p_l, s_l)$ 
34:   else if  $e$  elaborates a right parent  $p$  then
35:      $\vec{\phi} \leftarrow \text{PROBVECTORSINGLE}(n, p_r, s_r)$ 
36:   end if
37:   return CATEGORICALLOGPDF( $\vec{\phi}, n$ )
38: end function
39:

```

Scoring Unsplit Operations

Consider the Split rule :

$$t \rightarrow t'_l \ s' \ t'_r$$

During a split, each edge in the transition and each node in an adjacent slice can be elaborated by one or more inner operations. These new edges can be discarded or kept to form the new edge of t'_l and t'_r .

The notes in the child slice s can either have edges connected to the left neighboring slice or right neighbouring slice, or both. In other words each note in the child slice was generated as an elaboration of the left or right parent, or both in the case of a double sided operation.

The plausibility of fa

previous note, subsequent note, both, or repetition of prev note, subsequent note etc. So we consider the chord tone profiles of the involved slices.

We first guess the chord type each parent slice.

$$\theta_l = \underset{c \in C}{\operatorname{argmax}} P(s_l | c) \quad , \quad \theta_r = \underset{c \in C}{\operatorname{argmax}} P(s_r | c)'$$

We now consider each edge individually, considering their likelihoods based on the probabilistic model of harmony along with theoretical assumptions.

Single Sided Operations

- Right Neighbour (Left Neighbour analogously)

$$x \implies x \rightarrow n \quad , x, n \in P$$

$$x \sim \text{Categorical}(\sigma_{ct}^{\theta_l})$$

$$n \sim \text{Categorical}(\sigma_{or}^{\theta_r})$$

Find

$$P(x, n \mid \theta_l)$$

- Right Repeat (Left Repeat analogously)

$$x \implies x \rightarrow x \quad , x \in P$$

$$x \sim \text{Categorical}(\sigma_{ct}^{\theta_l})$$

Find

$$P(x \mid \theta_l)$$

Two Sided Operations

- Root Note: This operation is only done once in the original model. In our case we do not need to consider due to segment boundaries.
- Full Repeat:

$$x \implies x \rightarrow n, x, n \in P$$

$$x \sim \text{Categorical}(\sigma_{ct}^{\theta_l})$$

$$n \sim \text{Categorical}(\sigma_{or}^{\theta_r})$$

Find

$$P(x, n \mid \theta_l)$$

- Left Repeat of Right:

$$x \rightarrow y \implies x \rightarrow y' \rightarrow y$$

$$y \sim \text{Categorical}(\sigma_{ct}^{\theta_l})$$

Find

$$P(y \mid \theta_l)$$

- Full Neighbour:

$$x_1 \rightarrow x_2 \implies x_1 \rightarrow n \rightarrow x_2, x \in P$$

Find

$$P(\mid \theta_l, \theta_r)$$

Scoring Unspread Operations

Consider the Spread rule :

$$t_l s_r \rightarrow t'_l s_l t'_m s_r t'_r$$

We make the assumption that s , s_l , & s_r are all realisations of the same chord. This lines up with the music theorretical basis for this operation in the model(justify).

Thus we find the most likely chord (optional extension: marginalise over all chords)

$$\theta = \underset{c \in C}{argmax} P(s|c)$$

When then measure the extent to which the parent slices match this chord.

$$p(s_l, s_r \mid \theta)$$

We can calculate $p(s_l \mid \theta)$ and $p(s_r \mid \theta)$ using the multinomial distribution probability density function as described in the preparation chapter.

Scoring Unfreeze Operations

We assign 0 cost to unfreeze operations. This means we need to be careful about ensure that we don't just unfreeze the entire piece immediately. Careful construction of the search algorithm can ensure this. More later.

Full state evaluation

Figure 3.10: What does this do?

```

1: function FINDMYSTERIES( $S, k$ )
2:   if  $k \leq 1$  then
3:     return  $\{\}$                                 ▷ The empty set  $\emptyset$ 
4:   end if
5:    $n \leftarrow |S|$                                 ▷ Cardinality of  $S$ 
6:    $k \leftarrow \text{MIN}(k, n + 1)$                     ▷ Since maximum of  $n$  mysteries
7:    $r \leftarrow \lfloor \lfloor k/2 \rfloor (\frac{n+1}{k}) \rfloor$         ▷ 1-indexed rank of the pivot
8:    $p \leftarrow \text{MEDIANOFMEDIANSSELECT}(S, r)$       ▷ element of rank  $r$ 
9:    $S_1 \leftarrow \{x \in S : x < p\}$ 
10:   $S_2 \leftarrow \{x \in S : x > p\}$ 
11:  return  $\text{FINDMYSTERIES}(S_1, \lfloor k/2 \rfloor) \cup \{p\} \cup \text{FINDMYSTERIES}(S_2, \lceil k/2 \rceil)$ 
12: end function

```

We need to combine all of these in a fair way. Also the distinction between splits and spreads need to be considred, as they are different operations, the calculations of likelihood may cause an imbalance. All likelihoods are stored in log space.

3.5.3 Heuristic Search

Given the heuristic \mathcal{H}_0 , the first option is to use a greedy search. We keep unfreeze operations as an option.

We find that the number of operations causes a combinatorial blowup, so the search doesn't end.

Problem of very large slices.

Segment by segment heuristic parse - avoids the problem, but is slightly hacky. Can we incorporate our knowledge regarding the relative proportion of chord tones and ornaments. Should we allow duplicates of notes in slices? Perhaps we should favour spreads more.

Always consider a certain number of slices and spreads.

3.5.4 Beam Search

Step 2: Relax the heuristic search in order to reduce runtime/ lower complexity.

Algorithm 3 Greedy Search

```

initialise:  $state \leftarrow initialState$ 
while  $state$  is not a goal state do
   $state \leftarrow \arg \min_{s \in \text{EXPAND}(state)} H(s)$ 
   $freezes, spreads, splits \leftarrow \text{Split } nextStates \text{ by operation type}$ 
   $open \leftarrow \text{Take best } \beta \text{ best states from } freezes \cup spreads \cup splits$ 
end while
return Best state in open

```

In the case that there are 85,000,000 options, perhaps we should sample the options rather than evaluating all of them.

This version of heuristic search should be able to parse full pieces (hopefully), so can be used to compare with the baselines on an entire corpus.

Beam of size n , with 1 for a freeze, k for spread, $n-k-1$

3.5.5 Dual Beam Search

It isn't clear how to determine how we should balance the costs of unspread and unsplit op

Algorithm 4 MultiBeam Search

```

hyper-parameters:  $\beta \leftarrow BeamWidth, \gamma \leftarrow ReservoirSize$ 
initialise:  $open \leftarrow (initialState, 0)$ 
while  $open$  does not contain any goal states do
   $nextStates \leftarrow \bigcup_{(s,c) \in open} \{(s', c' + H(s')) : (s', c') \in \text{EXPAND}((s, c))\}$ 
   $freezes, spreads, splits \leftarrow \text{Split } nextStates \text{ by operation type}$ 
   $open \leftarrow \text{Take best } \beta \text{ best states from } freezes \cup spreads \cup splits$ 
end while
return Best state in open

```

3.5.6 Stochastic Dual Beam Search

We propose sampling from the options, increasing the proportion that we ignore dependent on the number of options.

Here H assigns a cost for moving from state s to s' . This is defined as the negated log likelihood in \mathcal{H} .

Algorithm 5 Reservoir MultiBeam Search

hyper-parameters: $\beta \leftarrow \text{BeamWidth}, \gamma \leftarrow \text{ReservoirSize}$
initialise: $\text{open} \leftarrow (\text{initialState}, 0)$
while open does not contain any goal states **do**
 $\text{nextStates} \leftarrow \bigcup_{(s,c) \in \text{open}} \{(s', c' + H(s')) : (s', c') \in \text{EXPAND}((s, c))\}$
 $\text{freezes}, \text{spreads}, \text{splits} \leftarrow \text{Split } \text{nextStates} \text{ by operation type}$
 $\text{unFreezes} \leftarrow \bigcup_{s \in \text{freezes}} (\text{RESERVOIRSAMPLE}(\gamma, s))$
 $\text{unSpreads} \leftarrow \bigcup_{s \in \text{spreads}} (\text{RESERVOIRSAMPLE}(\gamma, s))$
 $\text{unSplits} \leftarrow \bigcup_{s \in \text{splits}} (\text{RESERVOIRSAMPLE}(\gamma, s))$
 $\text{open} \leftarrow \text{Take best } \beta \text{ best states from } \text{unFreezes} \cup \text{unSpreads} \cup \text{unSplits}$
end while
return Best state in open

3.6 Choosing hyper parameters

3.7 Testing

3.7.1 Unit Tests

3.7.2 Qualitative Tests

Throughout the design of the heuristic H_0 and the implementation of different search algorithms, a few segments we used as recurring test examples.

3.8 Repository Overview:

Repository Justification

The repository has been split into four main folders, with the addition of `Main.hs` which serves as an interface between the python experiment code and the algorithms developed in Haskell.

- Firstly, the `src/Core/` folder contains all the core code, including the implementation of the parsing and inference functions using the probabilistic model of harmony, as well as some helper code for file handling.
- The `experiments/` folder contains all the python code that is used for this project. The experiments consist of three stages, as described by the three main files: `preprocess.py`, `experiments.py` and `analysis.py`. Splitting these stages up prevents wasteful computation, as all the pre-processing can be done just once, while experiments are run on the processed data iteratively alongside algorithm development.
- The `src/Algorithms/` folder contains all the parsing algorithms including the baseline and extension search algorithms. Having all the algorithms contained in one module allows

experiments to be run using any selection of algorithms and input data, facilitating the evaluation process.

- Finally, the `test/` folder contains unit tests and end-to-end tests for use in Continuous Integration.

Figure 3.11 illustrates how these modules are connected.

Table 3.1: Repository Overview

File/Folder	Description	LOC
protovoices-haskell/	Root directory	2272
├──src/		
│ ├──HeuristicParser.hs, HeuristicSearch.hs	Core Implementation (Section x)	470
│ ├──RandomChoiceSearch.hs, RandomSampleParser.hs	Baseline Implemetation (Section x)	121
│ ├──Heuristics.hs, PBHModel.hs	Extension Implementation (Section x)	383
│ ├──FileHandling.hs	Utilities	188
│ └──...		
├──app/		431
│ └──MainFullParse.hs	Entry Point	
├──harmonic-inference		
├──experiments/	Running Experiments	115
│ ├──preprocess.ipynb		
│ ├──experiments.ipynb		
│ ├──analysis.ipynb		
│ ├──dcml_params.json		
│ └──inputs/		611
└──test/	Unit Tests (Section x)	

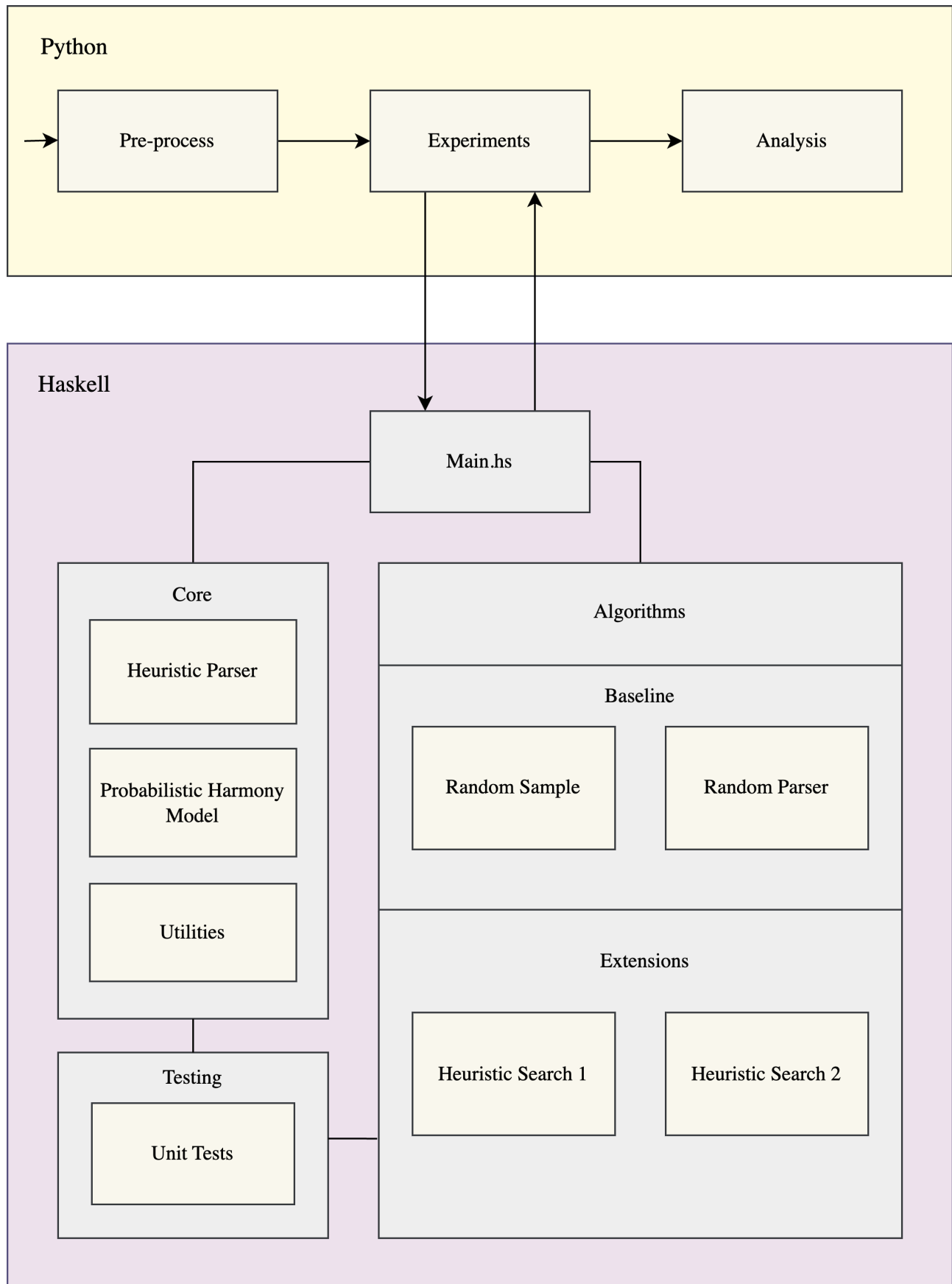


Figure 3.11: Block diagram of project components

Evaluation

In this chapter, I provide qualitative and quantitative evaluations of the work completed. I then provide and interpret evidence to show that the success criteria were met.

The main questions to answer are as follows:

- *Can the proto-voice model be used to accurately infer chord labels?*
- *Can the proto-voice model be used to practically infer chord labels?*
- *How well my heuristic search algorithms infer chord labels?*

4.1 Metrics

Ambiguity. which metrics exist? discussion. Accuracy and likelihood.

4.2 Harmony Model

4.2.1 Accuracy

When evaluating using the protovoice model: we assume that we result in only chord tones for each segment. Thus we use the chord tone probabilities to evaluate the prediction.

When just using a random sample, we have to assume that there is a mixture model of chord tones and ornaments. We use the learnt parameters to determine the distribution.

These two measures of likelihoods are comparable as they are drawn from the same distributions.

We also need to infer chord labels. We can simply choose the chord that is most likely according to our model.

This gives us two key metrics, likelihood and accuracy.

Could also use a more sophisticated notion of accuracy, using a chord similarity function [16]. The `mir_eval` package provides a plethora of metrics to compare chord label predictions [37].

4.2.2 Sensitivity Analysis

4.2.3 Qualitative Analysis

4.3 Baseline Algorithms

Reminder of what the search is actually doing.

Things to note

- The fact that segmentation is known ahead of time provides a great deal of information [13]
- So we use comparisons between the random sample from each segment algorithm and the random parse algorithm to see if the use of the grammar provides an advantage over just sampling the notes directly, without looking at relations between notes.
- Then we want a heuristic search algorithm that considers each option exhaustively and finds the best local option. This is too computationally expensive to be used for whole pieces.
- Given there can be millions of possible next states in the search, we need to look at different strategies to avoid searching through them all. E.g just sample states.
- Sensitivity Analysis for the heuristic search is useful for the evaluation. Explore how robust it is to handcrafted attacks/ different types of passages.
- Could evaluate by segments instead of pieces.

4.4 Extension Algorithms

4.4.1 Accuracy

4.4.2 Scalability

4.4.3 Qualitative Analysis

I've shown that i met success criteria x by this analysiss. etc.

4.5 Limitations

I'm solving:

$$\hat{L} = \arg \max_L (p(S|L)) \quad (4.1)$$

But I could be solving:

$$\hat{L} = \arg \max_L (p(S|L)p(L)) \quad (4.2)$$

In which case:

To compute the conditional probability $p(L|S)$, we use Bayes' theorem:

$$p(L|S) = \frac{p(S|L) p(L)}{P(S)} \quad (4.3)$$

Finding the most likely sequence of labels is found using:

$$\arg \max_L \left(\underbrace{p(S|L)}_{\text{likelihood}} \underbrace{p(L)}_{\text{prior}} \right) \quad (4.4)$$

The prior probability of a chord sequence $p(L)$ can be learned from a labeled dataset of chord sequences, and the likelihood can be found using a probabilistic harmony model. The likelihood $p(S|L)$ can be found using

This would be better, but was beyond the scope of the project.

Conclusions

In this chapter, I first discuss the success achieved by the project then offer a reflection on lessons learned. Finally, I consider the directions in which there is potential for future work.

5.1 Achievements

5.2 Lessons learned

5.3 Future Work

Bibliography

- [1] CHEN, T.-P., AND SU, L. Functional Harmony Recognition of Symbolic Music Data with Multi-task Recurrent Neural Networks, Sept. 2018.
- [2] CHEN, T.-P., AND SU, L. Harmony Transformer: Incorporating Chord Segmentation into Harmony Recognition. In *International Society for Music Information Retrieval Conference* (2019).
- [3] COHEN, D. E. “The Imperfect Seeks Its Perfection”: Harmonic Progression, Directed Motion, and Aristotelian Physics. *Music Theory Spectrum* 23, 2 (Oct. 2001), 139–169.
- [4] DAVIDSON-PILON, C. Bayesian Methods for Hackers, Mar. 2023.
- [5] DENG, J., AND KWOK, Y.-K. Large vocabulary automatic chord estimation using bidirectional long short-term memory recurrent neural network with even chance training. *Journal of New Music Research* 47, 1 (Jan. 2018), 53–67.
- [6] EBCIOĞLU, K. An Expert System for Harmonizing Four-Part Chorales. *Computer Music Journal* 12, 3 (1988), 43–51.
- [7] FEISTHAUER, L., BIGO, L., GIRAUD, M., AND LEVÉ, F. Estimating keys and modulations in musical pieces. In *Sound and Music Computing Conference (SMC 2020)* (June 2020).
- [8] FINKENSIEP, C. *The Structure of Free Polyphony*. PhD thesis, EPFL, Lausanne, 2023.
- [9] FINKENSIEP, C., ERICSON, P., KLASSMANN, S., AND ROHRMEIER, M. Chord Types and Ornamentation - A Bayesian Model of Extended Chord Profiles. *Open Research Europe* (Apr. 2023).
- [10] FINKENSIEP, C., AND ROHRMEIER, M. Modeling and Inferring Proto-voice Structure in Free Polyphony. In *Proceedings of the 22nd ISMIR Conference* (Online, Nov. 2021).
- [11] GJERDINGEN, R. Cognitive Foundations of Musical Pitch Carol L. Krumhansl. *Music Perception: An Interdisciplinary Journal* 9 (July 1992), 476–492.
- [12] GOODMAN, J. Semiring Parsing. *Computational Linguistics* 25, 4 (1999), 573–606.
- [13] GOTHAM, M., KLEINERTZ, R., WEISS, C., MÜLLER, M., AND KLAUK, S. What if the ‘When’ Implies the ‘What’?: Human harmonic analysis datasets clarify the relative role of the separate steps in automatic tonal analysis. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021* (2021), J. H. Lee, A. Lerch, Z. Duan, J. Nam, P. Rao, P. van Kranenburg, and A. Srinivasamurthy, Eds., pp. 229–236.

- [14] GRANROTH-WILDING, M. Harmonic Analysis of Music Using Combinatory Categorical Grammar. *The University Of Edinburgh* (2013).
- [15] HARASIM, D. Harmonic Syntax in Time: Rhythm Improves Grammatical Models of Harmony. In *Proceedings of the 20th ISMIR Conference* (2019), T. J. O'Donnell and M. A. Rohrmeier, Eds., ISMIR.
- [16] HUMPHREY, E. J., AND BELLO, J. P. Four Timely Insights on Automatic Chord Estimation.
- [17] JOHANNES, H. Ms3 - Parsing MuseScore 3. <https://github.com/johentsch/ms3>, 2021.
- [18] KIDNEY, D. O., AND WU, N. Algebras for weighted search. *Proceedings of the ACM on Programming Languages* 5, ICFP (Aug. 2021), 1–30.
- [19] KOOPS, H. V., DE HAAS, W. B., BRANSEN, J., AND VOLK, A. Automatic chord label personalization through deep learning of shared harmonic interval profiles. *Neural Computing and Applications* 32, 4 (Feb. 2020), 929–939.
- [20] KOOPS, H. V., DE HAAS, W. B., BURGOWNE, J. A., BRANSEN, J., KENT-MULLER, A., AND VOLK, A. Annotator subjectivity in harmony annotations of popular music. *Journal of New Music Research* 48, 3 (May 2019), 232–252.
- [21] KORZENIOWSKI, F., AND WIDMER, G. Improved Chord Recognition by Combining Duration and Harmonic Language Models, Aug. 2018.
- [22] KRUMHANSL, C. L., AND KESSLER, E. J. Tracing the dynamic changes in perceived tonal organization in a spatial representation of musical keys. *Psychological Review* 89 (1982), 334–368.
- [23] LERDAHL, F., AND JACKENDOFF, R. *A Generative Theory of Tonal Music*, repr. ed. MIT Press, Cambridge, Mass., 2010.
- [24] MARSDEN, A. Schenkerian Analysis by Computer: A Proof of Concept. *Journal of New Music Research* 39, 3 (Sept. 2010), 269–289.
- [25] MASADA, K., AND BUNESCU, R. Chord Recognition in Symbolic Music: A Segmental CRF Model, Segment-Level Features, and Comparative Evaluations on Classical and Popular Music, Oct. 2018.
- [26] MAUCH, M., MULLENSIEFEN, D., DIXON, S., AND WIGGINS, G. Can Statistical Language Models be used for the Analysis of Harmonic Progressions?
- [27] MAXWELL, H. J. An expert system for harmonizing analysis of tonal music. In *Understanding Music with AI: Perspectives on Music Cognition*. MIT Press, Cambridge, MA, USA, Aug. 1992, pp. 334–353.
- [28] MCLEOD, A., AND ROHRMEIER, M. A Modular System for the Harmonic Analysis of Musical Scores using a Large Vocabulary. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021* (2021), J. H. Lee, A. Lerch, Z. Duan, J. Nam, P. Rao, P. van Kranenburg, and A. Srinivasamurthy, Eds., pp. 435–442.

- [29] MEARNS, L. The Computational Analysis of Harmony in Western Art Music.
- [30] NEUWIRTH, M., HARASIM, D., MOSS, F. C., AND ROHRMEIER, M. The Annotated Beethoven Corpus (ABC): A Dataset of Harmonic Analyses of All Beethoven String Quartets. *Frontiers in Digital Humanities* 5 (July 2018), 16.
- [31] NI, Y., MCVICAR, M., SANTOS-RODRIGUEZ, R., AND DE BIE, T. An end-to-end machine learning system for harmonic analysis of music, July 2011.
- [32] PARDO, B., AND BIRMINGHAM, W. P. Algorithms for Chordal Analysis. *Computer Music Journal* 26, 2 (June 2002), 27–49.
- [33] PEARL, J. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. The Addison-Wesley Series in Artificial Intelligence. Addison-Wesley Pub. Co, Reading, Mass, 1984.
- [34] PICKENS, J. *Harmony Modeling for Polyphonic Music Retrieval*. PhD thesis, University of Massachusetts, 2004.
- [35] QUINN, I. Are Pitch-Class Profiles Really “Key for Key”? *Zeitschrift der Gesellschaft für Musiktheorie [Journal of the German-Speaking Society of Music Theory]* 7 (Jan. 2010).
- [36] RADICIONI, D. P., AND ESPOSITO, R. BREVE: An HMPerception-Based Chord Recognition System. In *Advances in Music Information Retrieval*, J. Kacprzyk, Z. W. Raś, and A. A. Wiecezorkowska, Eds., vol. 274. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010, pp. 143–164.
- [37] RAFFEL, C., MCFEE, B., HUMPHREY, E. J., SALAMON, J., NIETO, O., LIANG, D., AND ELLIS, D. P. W. Mir_eval: A Transparent Implementation of Common MIR Metrics. In *Proceedings of the 15th International Conference on Music Information Retrieval* (2014).
- [38] RAPHAEL, C., AND STODDARD, J. Functional Harmonic Analysis Using Probabilistic Models. *Computer Music Journal* 28, 3 (Sept. 2004), 45–52.
- [39] ROCHER, T., ROBINE, M., HANNA, P., AND STRANDH, R. Dynamic Chord Analysis for Symbolic Music.
- [40] SAPP, C. Visual hierarchical key analysis. *Computers in Entertainment* 3 (Oct. 2005), 1–19.
- [41] SAPP, C. 6 Computational Chord-Root Identification in Symbolic Musical Data: Rationale, Methods, and Applications.
- [42] SIDOROV, K., JONES, A., AND MARSHALL, D. MUSIC ANALYSIS AS A SMALLEST GRAMMAR PROBLEM.
- [43] TEMPERLEY, D. An Algorithm for Harmonic Analysis. *Music Perception* 15, 1 (Oct. 1997), 31–68.
- [44] TEMPERLEY, D. A Bayesian Approach to Key-Finding. In *Music and Artificial Intelligence*, G. Goos, J. Hartmanis, J. van Leeuwen, C. Anagnostopoulou, M. Ferrand, and A. Smaill, Eds., vol. 2445. Springer Berlin Heidelberg, Berlin, Heidelberg, 2002, pp. 195–206.

- [45] TEMPERLEY, D. A Bayesian Approach to Key-Finding. In *Music and Artificial Intelligence* (Berlin, Heidelberg, 2002), C. Anagnostopoulou, M. Ferrand, and A. Smaill, Eds., Lecture Notes in Computer Science, Springer, pp. 195–206.
- [46] TEMPERLEY, D. A Unified Probabilistic Model for Polyphonic Music Analysis. *Journal of New Music Research* 38, 1 (Mar. 2009), 3–18.
- [47] VOLK, A. Improving Audio Chord Estimation by Alignment and Integration of Crowd-Sourced Symbolic Music. 141–155.
- [48] WINOGRAD, T. Linguistics and the Computer Analysis of Tonal Harmony. *Journal of Music Theory* 12, 1 (1968), 2–49.

Additional Information

Project Proposal

Inferring Harmony from Free Polyphony

Judah Daniels

April 22, 2023

DOS: Prof. Larry Paulson

Crsid: jasd6

College: Clare College

B.1 Abstract

A piece of music can be described using a sequence of chords, representing a higher level harmonic structure of a piece. There is a small, finite set of chord types, but each chord can be realised on the musical surface in a practically infinite number of ways. Given a score, we wish to infer the underlying chord types.

The paper *Modeling and Inferring Proto-voice Structure in Free Polyphony* describes a generative model that encodes the recursive and hierarchical dependencies between notes, giving rise to a grammar-like hierarchical system [10]. This proto-voice model can be used to reduce a piece into a hierarchical structure which encodes an understanding of the tonal/harmonic relations of a piece.

Christoph Finkensiep suggests in his paper that the proto-voice model may be an effective way to infer higher level latent entities, such as harmonies or voice leading schemata. Thus in this project I will ask the question: is this parsing model an effective way to annotate harmonies? By ‘effective’ we are referring to two things:

- Accuracy: can the model successfully emulate how experts annotate harmonic progressions in musical passages?
- Practicality: can the model be used to do this within a reasonable time frame?

While the original model could in theory be used to generate harmonic annotations, its exhaustive search strategy would be prohibitively time-consuming in practice for any but the shortest musical extracts; one half measure can have over 100,000 valid derivations [10]. My approach will be to explore the use of heuristic search algorithms to solve this problem.

B.2 Substance and Structure

B.2.1 Core: Search

The core of this project is essentially a search problem characterised as follows:

- The state space S is the set of all possible partial reductions of a piece along with each reduction step that has been done so far.
- We have an initial state $s_o \in S$, which is the empty reduction, corresponding to the unreduced surface of the piece. The score is represented as a sequence of slices grouping notes that sound simultaneously. We are also given the segmentation of the original chord labels that we wish to retrieve.
- We have a set of actions, A modelled by a function $action : A \times S \rightarrow S$. These actions correspond to a single reduction step.
 - The reduction steps are the inverses of the operations defined by the generative proto-voice model.

- Finally we have a goal test, $goal : S \rightarrow \{true, false\}$ which is true iff the partial reduction s has exactly one slice per segment of the input.
 - This means the partial reduction s contains a sequence of slices which start and end positions corresponding to the segmentation of the piece.
- At the first stage, this will be implemented using a random graph search algorithm, picking each action randomly, according to precomputed distributions.

B.2.2 Core: Evaluation

The second core task is to create an evaluation module that iterates over the test dataset, and evaluates the partial reduction computed by the search algorithm above. This will be done by comparing the outputs to ground truth annotations from the Annotated Beethoven Corpus.

In order to do this I will make use of the statistical harmony model from Finkensiep’s thesis, *The Structure of Free Polyphony* [8]. This model provides a way of mapping between the slices that the algorithm generates and the chords in the ground truth. This can be used to empirically measure how closely the slices match the expert annotations.

B.2.3 Extension

Once the base search implementation and evaluation module have been completed, the search problem will be tackled by heuristic search methods, with different heuristics to be trialled and evaluated against each other. The heuristics will make use of the chord profiles from Finkensiep’s statistical harmony model discussed above. These profiles relate note choices to the underlying harmony. Hence the heuristics may include:

- How the chord types relate to the pitches used.
- How the chord types relate which notes are used as ornamentation, and the degree of ornamentation.
- Contextual information about neighboring slices

B.2.4 Overview

The main work packages are as follows:

Preliminary Reading – Familiarise myself with the proto-voice model, and read up on similar models and their implementations. Study heuristic search algorithms.

Dataset Preparation – Pre-process the Annotated Beethoven Corpus into a suitable representation for my algorithm.

Basic Search – Implement a basic random search algorithm that takes in surface and segmentations, and outputting the sequence of slices matching the segmentations.

Evaluation Module – Implement an evaluation module to evaluate the output from the search algorithm.

End-to-end pipeline – Implement a full pipeline from the data to the evaluation that can be used to compare different reductions.

Heuristic Design – Extension – Trial different heuristics and evaluate their performance against each other.

Dissertation – I intend to work on the dissertation throughout the duration of the project. I will then focus on completing and polishing the project upon completion.

B.3 Starting Point

The following describes existing code and languages that will be used for this project:

Haskell – I will be using Haskell for this project as it is used in the proto-voice implementation. It must be noted that my experience with Haskell is limited, as I was first introduced to it via an internship this summer (July to August 2022).

Python – Python will be used for data handling. I have experience coding in Python.

Prior Research - Over the summer I have been reading the literature on computational models of music, as well as various parsing algorithms such as semi-ring parsing [12], and the CYK algorithm, which is used in the implementation of the proto-voice model.

Protovoices-Haskell – The paper *Modeling and Inferring Proto-Voice Structure in Free Polyphony* [10] includes an implementation of the proto-voice model in Haskell. A fork of this repository will form the basis of my project. This repository includes a parsing module which will be used to perform the actions in the search space of partial reductions. There is a module that can exhaustively enumerate reductions of a piece, but this is infeasible in practice due to the blowup of the derivation forest.

MS3 – This is a library for parsing MuseScore Files and manipulating labels [17], which I will use as part of the data processing pipeline.

ABC – The *Annotated Beethoven Corpus* [30] contains analyses of all Beethoven string quartets composed between 1800 and 1826, encoded in a human and machine readable format. This will be used as a dataset for this project.

B.4 Success Criteria

This project will be deemed a success if I complete the following tasks:

- Develop a baseline search algorithm that uses the proto-voice model to output a partial reduction of a piece of music up to the chord labels.
- Create an evaluation module that can take the output of the search algorithm and quantitatively evaluate its accuracy against the ground truth annotations by providing a score based on a statistical harmony model.
- Extension: Develop one or more search algorithms that use additional heuristics to inform the search, and compare the accuracy with the baseline algorithm.

B.5 Timetable

Time frame	Work	Evidence
Michaelmas (Oct 4 to Dec 2)		
Oct 14 to Oct 24	<i>Oct 14</i> : Final proposal deadline. Preparation work: familiarise myself with the dataset and the proto-voice model implementation. Work on manipulating reductions using the proto-voice parser provided by the paper.	None
Oct 24 to Nov 7	Dataset preparation and handling.	Plot useful metrics about the dataset using Haskell
Nov 7 to Nov 21	Random Search implementation	None
Nov 21 to Dec 5	Evaluation Module. Continue with search implementation.	Evaluate a manually created derivation and plot results
Vacation (Dec 3 to Jan 16)		
Dec 5 to Dec 11	Evaluate performance of random search. Begin to work on extensions	Plot results
Dec 10 to Dec 21	Trial different heuristics. Implement an end-to-end pipeline from input to evaluation.	None
Dec 21 to Dec 27	None	None
Dec 27 to Jan 10	Continue trialing and evaluating heuristics	<i>Fulfill success criterion: At least one heuristic technique gives better performance than random search.</i>
Lent (Jan 17 to Mar 17)		
Jan 4 to Jan 20	Buffer Period to help keep on track	None
Jan 20 to Feb 3	<i>Feb 3</i> : Progress Report Deadline. Write progress report and prepare presentation. Write draft <i>Evaluation</i> chapter	Progress Report (approx. 1 page)
Feb 3 to Feb 17	Prepare presentation.	<i>Feb 8 – 15</i> : Progress Report presentation
Feb 17 to Mar 3	<i>Feb 17</i> : How to write a Dissertation briefing. Write draft Introduction and Preparation chapters. Incorporate feedback on Evaluation chapter.	Send draft Introduction and Preparation chapter to supervisor
Mar 3 to Mar 17	Write draft Implementation chapters. Incorporate feedback on Introduction and Preparation chapters.	Send draft Implementation chapters to Supervisor
Vacation (Mar 18 to		

B.6 Resources

I plan to use my own laptop for development: MacBook Pro 16-inch, M1 Max, 32GB Ram, 1TB SSD, 24-core GPU.

All code will be stored on a GitHub repository, which will guarantee protection from data loss. I will easily be able to switch to using university provided computers upon hardware/software failure.

The project will be built upon work that has been done in the DCML (Digital cognitive musicology lab) based in EPFL. The files are in their Github repository, and I have been granted permission to access their in-house datasets of score annotations, as well as software packages which are used to handle the data.

B.7 Supervisor Information

Peter Harrison, head of Centre for Music and Science at Cambridge, has agreed to supervise me for this. We have agreed on a timetable for supervisions for this year. I am also working with Christoph Finkensiep, a PHD student at the DCML, and originator of the proto-voice model. Professor Larry Paulson has agreed to be the representative university teaching officer.

Bibliography

- [1] CHEN, T.-P., AND SU, L. Functional Harmony Recognition of Symbolic Music Data with Multi-task Recurrent Neural Networks, Sept. 2018.
- [2] CHEN, T.-P., AND SU, L. Harmony Transformer: Incorporating Chord Segmentation into Harmony Recognition. In *International Society for Music Information Retrieval Conference* (2019).
- [3] COHEN, D. E. “The Imperfect Seeks Its Perfection”: Harmonic Progression, Directed Motion, and Aristotelian Physics. *Music Theory Spectrum* 23, 2 (Oct. 2001), 139–169.
- [4] DAVIDSON-PILON, C. Bayesian Methods for Hackers, Mar. 2023.
- [5] DENG, J., AND KWOK, Y.-K. Large vocabulary automatic chord estimation using bidirectional long short-term memory recurrent neural network with even chance training. *Journal of New Music Research* 47, 1 (Jan. 2018), 53–67.
- [6] EBCIOĞLU, K. An Expert System for Harmonizing Four-Part Chorales. *Computer Music Journal* 12, 3 (1988), 43–51.
- [7] FEISTHAUER, L., BIGO, L., GIRAUD, M., AND LEVÉ, F. Estimating keys and modulations in musical pieces. In *Sound and Music Computing Conference (SMC 2020)* (June 2020).
- [8] FINKENSIEP, C. *The Structure of Free Polyphony*. PhD thesis, EPFL, Lausanne, 2023.
- [9] FINKENSIEP, C., ERICSON, P., KLASSMANN, S., AND ROHRMEIER, M. Chord Types and Ornamentation - A Bayesian Model of Extended Chord Profiles. *Open Research Europe* (Apr. 2023).
- [10] FINKENSIEP, C., AND ROHRMEIER, M. Modeling and Inferring Proto-voice Structure in Free Polyphony. In *Proceedings of the 22nd ISMIR Conference* (Online, Nov. 2021).
- [11] GJERDINGEN, R. Cognitive Foundations of Musical Pitch Carol L. Krumhansl. *Music Perception: An Interdisciplinary Journal* 9 (July 1992), 476–492.
- [12] GOODMAN, J. Semiring Parsing. *Computational Linguistics* 25, 4 (1999), 573–606.
- [13] GOTHAM, M., KLEINERTZ, R., WEISS, C., MÜLLER, M., AND KLAUK, S. What if the ‘When’ Implies the ‘What’?: Human harmonic analysis datasets clarify the relative role of the separate steps in automatic tonal analysis. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021* (2021), J. H. Lee, A. Lerch, Z. Duan, J. Nam, P. Rao, P. van Kranenburg, and A. Srinivasamurthy, Eds., pp. 229–236.

- [14] GRANROTH-WILDING, M. Harmonic Analysis of Music Using Combinatory Categorical Grammar. *The University Of Edinburgh* (2013).
- [15] HARASIM, D. Harmonic Syntax in Time: Rhythm Improves Grammatical Models of Harmony. In *Proceedings of the 20th ISMIR Conference* (2019), T. J. O'Donnell and M. A. Rohrmeier, Eds., ISMIR.
- [16] HUMPHREY, E. J., AND BELLO, J. P. Four Timely Insights on Automatic Chord Estimation.
- [17] JOHANNES, H. Ms3 - Parsing MuseScore 3. <https://github.com/johentsch/ms3>, 2021.
- [18] KIDNEY, D. O., AND WU, N. Algebras for weighted search. *Proceedings of the ACM on Programming Languages* 5, ICFP (Aug. 2021), 1–30.
- [19] KOOPS, H. V., DE HAAS, W. B., BRANSEN, J., AND VOLK, A. Automatic chord label personalization through deep learning of shared harmonic interval profiles. *Neural Computing and Applications* 32, 4 (Feb. 2020), 929–939.
- [20] KOOPS, H. V., DE HAAS, W. B., BURGOYNE, J. A., BRANSEN, J., KENT-MULLER, A., AND VOLK, A. Annotator subjectivity in harmony annotations of popular music. *Journal of New Music Research* 48, 3 (May 2019), 232–252.
- [21] KORZENIOWSKI, F., AND WIDMER, G. Improved Chord Recognition by Combining Duration and Harmonic Language Models, Aug. 2018.
- [22] KRUMHANSL, C. L., AND KESSLER, E. J. Tracing the dynamic changes in perceived tonal organization in a spatial representation of musical keys. *Psychological Review* 89 (1982), 334–368.
- [23] LERDAHL, F., AND JACKENDOFF, R. *A Generative Theory of Tonal Music*, repr. ed. MIT Press, Cambridge, Mass., 2010.
- [24] MARSDEN, A. Schenkerian Analysis by Computer: A Proof of Concept. *Journal of New Music Research* 39, 3 (Sept. 2010), 269–289.
- [25] MASADA, K., AND BUNESCU, R. Chord Recognition in Symbolic Music: A Segmental CRF Model, Segment-Level Features, and Comparative Evaluations on Classical and Popular Music, Oct. 2018.
- [26] MAUCH, M., MULLENSIEFEN, D., DIXON, S., AND WIGGINS, G. Can Statistical Language Models be used for the Analysis of Harmonic Progressions?
- [27] MAXWELL, H. J. An expert system for harmonizing analysis of tonal music. In *Understanding Music with AI: Perspectives on Music Cognition*. MIT Press, Cambridge, MA, USA, Aug. 1992, pp. 334–353.
- [28] MCLEOD, A., AND ROHRMEIER, M. A Modular System for the Harmonic Analysis of Musical Scores using a Large Vocabulary. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021* (2021), J. H. Lee, A. Lerch, Z. Duan, J. Nam, P. Rao, P. van Kranenburg, and A. Srinivasamurthy, Eds., pp. 435–442.

- [29] MEARNS, L. The Computational Analysis of Harmony in Western Art Music.
- [30] NEUWIRTH, M., HARASIM, D., MOSS, F. C., AND ROHRMEIER, M. The Annotated Beethoven Corpus (ABC): A Dataset of Harmonic Analyses of All Beethoven String Quartets. *Frontiers in Digital Humanities* 5 (July 2018), 16.
- [31] NI, Y., MCVICAR, M., SANTOS-RODRIGUEZ, R., AND DE BIE, T. An end-to-end machine learning system for harmonic analysis of music, July 2011.
- [32] PARDO, B., AND BIRMINGHAM, W. P. Algorithms for Chordal Analysis. *Computer Music Journal* 26, 2 (June 2002), 27–49.
- [33] PEARL, J. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. The Addison-Wesley Series in Artificial Intelligence. Addison-Wesley Pub. Co, Reading, Mass, 1984.
- [34] PICKENS, J. *Harmony Modeling for Polyphonic Music Retrieval*. PhD thesis, University of Massachusetts, 2004.
- [35] QUINN, I. Are Pitch-Class Profiles Really “Key for Key”? *Zeitschrift der Gesellschaft für Musiktheorie [Journal of the German-Speaking Society of Music Theory]* 7 (Jan. 2010).
- [36] RADICIONI, D. P., AND ESPOSITO, R. BREVE: An HMPerception-Based Chord Recognition System. In *Advances in Music Information Retrieval*, J. Kacprzyk, Z. W. Raś, and A. A. Wieczorkowska, Eds., vol. 274. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010, pp. 143–164.
- [37] RAFFEL, C., MCFEE, B., HUMPHREY, E. J., SALAMON, J., NIETO, O., LIANG, D., AND ELLIS, D. P. W. Mir_eval: A Transparent Implementation of Common MIR Metrics. In *Proceedings of the 15th International Conference on Music Information Retrieval* (2014).
- [38] RAPHAEL, C., AND STODDARD, J. Functional Harmonic Analysis Using Probabilistic Models. *Computer Music Journal* 28, 3 (Sept. 2004), 45–52.
- [39] ROCHER, T., ROBINE, M., HANNA, P., AND STRANDH, R. Dynamic Chord Analysis for Symbolic Music.
- [40] SAPP, C. Visual hierarchical key analysis. *Computers in Entertainment* 3 (Oct. 2005), 1–19.
- [41] SAPP, C. 6 Computational Chord-Root Identification in Symbolic Musical Data: Rationale, Methods, and Applications.
- [42] SIDOROV, K., JONES, A., AND MARSHALL, D. MUSIC ANALYSIS AS A SMALLEST GRAMMAR PROBLEM.
- [43] TEMPERLEY, D. An Algorithm for Harmonic Analysis. *Music Perception* 15, 1 (Oct. 1997), 31–68.
- [44] TEMPERLEY, D. A Bayesian Approach to Key-Finding. In *Music and Artificial Intelligence*, G. Goos, J. Hartmanis, J. van Leeuwen, C. Anagnostopoulou, M. Ferrand, and A. Smaill, Eds., vol. 2445. Springer Berlin Heidelberg, Berlin, Heidelberg, 2002, pp. 195–206.

- [45] TEMPERLEY, D. A Bayesian Approach to Key-Finding. In *Music and Artificial Intelligence* (Berlin, Heidelberg, 2002), C. Anagnostopoulou, M. Ferrand, and A. Smaill, Eds., Lecture Notes in Computer Science, Springer, pp. 195–206.
- [46] TEMPERLEY, D. A Unified Probabilistic Model for Polyphonic Music Analysis. *Journal of New Music Research* 38, 1 (Mar. 2009), 3–18.
- [47] VOLK, A. Improving Audio Chord Estimation by Alignment and Integration of Crowd-Sourced Symbolic Music. 141–155.
- [48] WINOGRAD, T. Linguistics and the Computer Analysis of Tonal Harmony. *Journal of Music Theory* 12, 1 (1968), 2–49.