

# SYNTHESIZING COMPOSITE HIERARCHICAL STRUCTURE FROM MUSIC DATASETS

Ilana Shapiro

UC San Diego

ilshapiro@ucsd.edu

## ABSTRACT

Music is an innately hierarchical system, comprising multiple semantic levels informed by music theory. Such levels include formal structure segmentation, disjoint motif repetition, and harmonic and melodic contour. Historically, researchers in the music information retrieval community have focused on developing analyses for single levels in this hierarchy. However, existing research has addressed neither (1) how to combine arbitrarily many levels of structure analyses into a single unified model and (2) how to extract a representative such structure from a corpus of music, rather than just a single piece. In this work, we propose a novel data structure called the *semantic temporal graph* that captures the semantic (i.e. hierarchical music theoretic) relationships between levels of the hierarchy, as well as the temporal relationships between adjacent-level analyses. Furthermore, given a corpus of such graphs derived from individual pieces, we introduce a method rooted in stochastic optimization to derive a representative *centroid* graph encoding the music dataset’s overall structure. We provide a qualitative evaluation of the semantic temporal graph [where we....], as well as a quantitative evaluation of the centroid graph [where we...].

## 1. INTRODUCTION

Music is both composed and comprehended with a hierarchical structure. Melodies comprising individual notes form the bottom of the hierarchy, followed by harmonic contour, rhythmic patterns, disjoint motifs, and finally large scale sections. Together, this composite hierarchy defines the overall structure of a piece [1].

Automatic identification of musical structure, also known as *music structure analysis* (MSA), is a major interest to both musicologists and the MIR community. Research thus far has focused on the automatic contiguous segmentation (both flat and hierarchical) of musical form [1–13], involving a boundary detection step followed by a segment labeling step. Motif detection [14–16], which looks for disjoint repeating musical patterns, has also

been investigated. More recently, researchers have studied harmonic [17], functional harmonic [18], and melodic [19–21] contour extraction. The techniques used are diverse, ranging from matrix factorization to deep learning. These tasks have all been proposed in annual competitions of the Music Information Retrieval eXchange (MIREX) [22–25], which provides a standard output format.

To our knowledge, all existing MSA research addresses a single aspect of the compositional hierarchy, such as motif extraction, or melodic contour. There is currently no notion of how reconcile differing levels of the hierarchy into a single, unified model of structure, even though their amalgamation is central to a piece’s compositional architecture and cohesive integrity. Indeed, Dai et al. have shown empirically that there are significant interactions between different structural levels [13]. In identifying the components necessary for integrating the levels, we find there are two central challenges: how to convey each level’s semantic, music theoretic level in the hierarchy, and how to encapsulate the temporal relationships between the results of the structural analyses at adjacent hierarchical levels.

Furthermore, prior MSA research has only addressed the problem of identifying structure in a single piece, and there is presently no methodology for describing the overall structure of a musical corpus, even across existing single-level analyses. The one exception is Oriol Nieto’s proposed technique for merging multiple segment boundary annotations [1], but this is intended to be used with multiple boundary detection algorithms over a single piece to alleviate the problem of subjectivity, and does not address the problem of reconciling differing labels or levels.

### 1.1 Contributions

To address the first gap, in Section 3 we introduce the *semantic temporal graph* (STG), a  $k$ -partite directed acyclic graph (DAG). In this novel data structure, the semantic, music theoretic levels of the compositional hierarchy are represented as levels in the  $k$ -partite structure, nodes represent structure labels that are the result of the relevant analysis, and edges between nodes of adjacent levels convey the temporal relationships between them. Each node has an associated time interval determined by the relevant MSA algorithm. A node must have one or two parents at the level above it: one if its associated time interval is a total subset of its parents, and two if its time interval begins in one parent and ends in the other. In order to easily parse the results of MSA algorithms into this data structure, the



standard MIREX format is consulted.

Importantly, the STG is incredibly flexible, and supports the representation of arbitrarily many layers and layer types. Furthermore, the STG is totally decoupled from any specific MSA algorithm or input format, meaning that the chosen MSA algorithm for any level can be easily swapped out, as long as its output adheres to the standard MIREX format. This is crucial as single-level MSA algorithms are constantly improving, and the STG must be adaptable enough to accommodate this.

Finally, to address the second gap, in Section 4 we examine the problem of finding a *centroid*, or most representative, graph given a corpus of STGs derived from individual pieces. We use the label-aware graph edit distance as the cost between two STGs. Given a set of STGs  $G$ , we seek to construct the STG  $g^*$  that minimizes this cost from  $g^*$  to every graph in  $G$ . This is a constraint satisfaction problem, but is NP-hard and intractable at this size. Thus, we must rely on approximation techniques, and utilize Markov Chain Monte Carlo methods, demonstrating how to use the Metropolis Hastings algorithm to infer an optimal solution and thus arrive at the centroid graph most descriptive of the entire corpus by construction.<sup>1</sup>

## 2. RELATED WORK

The majority of existing MSA algorithms focus on the contiguous formal segmentation task: boundary detection, followed by section labeling. Segmentation algorithms can be flat or hierarchical, where each level is an increasingly granular contiguous segmentation of the piece. Existing approaches generally utilize matrix factorization, formal grammars, or more recently, neural networks. Oriol Nieto’s Music Structure Analysis Framework (MSAF) toolkit [1] features most of the state-of-the art matrix factorization approaches, including ordinal linear discriminant analysis [8], convex nonnegative matrix factorization [5], checkerboard [2], spectral clustering [9], the Structural Features algorithm [3], 2D-Fourier Magnitude Coefficients [6], and the Variable Markov Oracle [4]. These methods use variants of self-similarity matrices (SSMs) for boundary detection, which are symmetric matrices storing pair-wise comparisons between a given set of features. To assign labels, methods such as Gaussian mixture models and nearest neighbor search are employed [1].

More recently, both supervised [26] and unsupervised [12, 27] deep learning approaches to segmentation have studied, as have graph- and grammar-based approaches. Hernandez-Olivan et al. use the graph-based G-PELT algorithm [10], while Dai et al. employ a graph-based approach informed by interactions between sections, melody, harmony and rhythm [13]. In the grammar realm, Finkensiep et al. attempt a unified model of structure by using a minimal context-free grammar to combine repetition with formal prototypes in a tree-based approach, but they still only operate within the contiguous segmentation domain [11].

In the motif extraction domain, algorithms search for disjoint, repeating, and possibly overlapping patterns in a piece. They generally fall into three categories: string-based approaches (e.g. the Variable Markov Oracle [15]), where music data is represented as a one-dimensional pitch sequence and repeated patterns are detected with sub-string matching; geometry-based approaches (e.g. Hsiao et al’s BPS-motif discovery algorithm [14]), where music data is represented as multidimensional point sets and *translatable subsets* identify repeating patterns; and feature-based (e.g. [16]) approaches, which learn features from music data, and retrieve patterns with clustering or classification of the features as repeated patterns.

Recent approaches in harmony identification are centered around neural networks, such as using multi-task learning with recurrent neural networks and long short term memory units to detect functional harmonic relationships in a piece [18], as well as developing transformer models that incorporate chord segmentation into the recognition process [17]. Until very recently, Justin Salamon’s Melodia algorithm was the state of the art in melody extraction. It comprises sinusoid extraction, salience functions, contour creation, and melody selection, and automatically estimates the fundamental frequency corresponding to the pitch of the predominant melodic line of a piece of polyphonic music. Since then, approaches have shifted to neural networks, such as lightweight deep bidirectional LSTM models [20] and transformers [21].

To our knowledge, little work has been done in developing a unified model of structure for either a single piece or across a corpus. The closest we are aware of is Halley Young’s *prototype graph*, a bipartite graph that represents musical form as a network of relationships between “prototype nodes” (specific musical elements) and the music they represent, as well as the music-theoretic relationships between musical spans and their respective nodes [28]. However, the prototype graph does not encode hierarchy, and it is limited to a single piece.

## 3. ABSTRACT REPRESENTATION

### 3.1 Semantic Temporal Graph

We seek a unified model of musical structure that captures the semantic, music theoretic levels of the compositional hierarchy, as well as the temporal relationships between them. The *semantic temporal graph*, or STG, is a novel data structure serving as a meta-representation of musical structure that unifies these objectives. The STG is  $k$ -partite directed acyclic graph (DAG). Each of the  $k$  layers encodes a level in the semantic hierarchy dictated by music theory. From top to bottom, this hierarchy consists of large scale form, motifs, rhythm, harmony, and finally melodies constituting specific note events [1]. Individual levels themselves can form sub-hierarchies of increasing granularity, as is commonly seen with the segmentation algorithms corresponding to large scale form.

In the STG, nodes encode labels from the relevant MSA algorithm, and contain associated time intervals given by

<sup>1</sup> Source code for this project can be found at [https://github.com/ilanashapiro/constraints\\_project](https://github.com/ilanashapiro/constraints_project)

the algorithm. At each level, these intervals are converted to indices by ordering their start times. Edges encode temporal relationships between nodes of adjacent levels. Specifically, for node  $n$  at level  $i$ ,  $n$  must have either one or two parents in level  $i - 1$ : one if its associated time interval is a total subset of its parent's, and two if its time interval begins in one parent's, and ends in the other's.

### 3.2 Building the Graph

In order to visualize the STG more clearly, we must first establish how it is built. In this paper, we build STGs that unify hierarchical formal segmentation and disjoint motif repetition, with plans in the near future to add layers for rhythm, harmony, and melody. In order to parse the results of any MSA algorithm into a format suitable for the STG, we consider the standard MIREX formats they adhere to.

The structural segmentation task was most recently proposed in the 2017 MIREX competition [22]. It defines the following standard output format:

```
<onset_time(sec)>\t<offset_time(sec)>\t<label>\n
<onset_time(sec)>\t<offset_time(sec)>\t<label>\n
...
```

where  $\backslash t$  denotes a tab and  $\backslash n$  denotes the end of line. The  $<$  and  $>$  characters are not included. Thus, an example output file could look like:

```
0.000 5.223 A
5.223 15.101 B
15.101 20.334 A
```

Hierarchical segmentations will have multiple of these lists separated by a newline.

MIREX 2017 also dictates the standard output format for motif detection algorithms [23]. Ontimes (in seconds) of notes in the pattern are in the left column, and MIDI note numbers are on the right. Each pattern occurrence is given before continuing to the next pattern, and occurrences can be differing lengths. An example output would resemble:

```
pattern1
occurrence1
7.00000, 45.00000
...
11.00000, 60.00000
occurrence2
31.00000, 57.00000
...
patternM
occurrence1
9.00000, 58.00000
...
occurrencem
100.00000, 62.00000
...
```

Given these standard MIREX formats, we parse this data into a format suitable for the STG. Each "chunk" of the data, separated by a newline, forms a level in the STG, and each line in a chunk corresponds to a node. Consider Figure 1, which demonstrates the parsing process for a segmentation hierarchy with two levels. We start with the MIREX segmentation output in step 1. In step 2, we number each segmentation level (i.e. text chunk) and extract the segment label and time interval from each line. This gives us the preliminary node label  $S\{\text{segment label}\}L\{\text{segment$

level number $\}I\{\text{time interval}\}$ . Finally, in step 3 we convert each time interval into an index corresponding to that line's position in the level, sorted by start time. In the case of segmentation, this matches the MIREX line order. We now have the final segmentation node label  $S\{\text{segment label}\}L\{\text{segment level number}\}N\{\text{label index within level}\}$ .

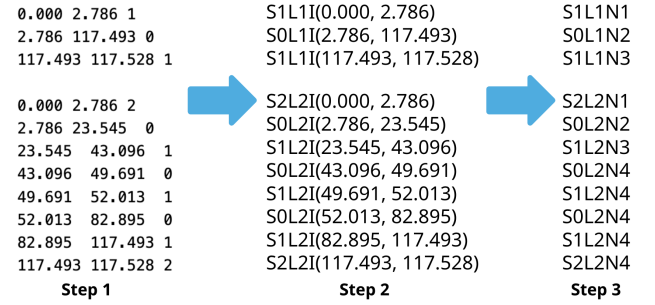


Figure 1. Parsing Formal Segmentation

We repeat this process for motifs as in Figure 2. Motif analyses are single level, i.e. they do not output sub-hierarchies like segmentation often does. Consider the MIREX motif output in step 1. Each pattern/motif occurrence's time interval is defined as the ontime of its first note to the ontime of its last. We extract pattern and occurrence numbers in step 2 to arrive at the preliminary node label  $P\{\text{pattern number}\}O\{\text{occurrence number}\}I\{\text{time interval}\}$ . We then order the preliminary labels by start time, as unlike in segmentation, pattern occurrences are not ordered temporally in MIREX format. We then assign indices to this ordering to arrive at the final motif node label  $P\{\text{pattern number}\}O\{\text{occurrence number}\}N\{\text{chronologic occurrence index}\}$ . This is reflected in step 3, where occurrence 2 of pattern 1 is assigned index 3 instead of 2.

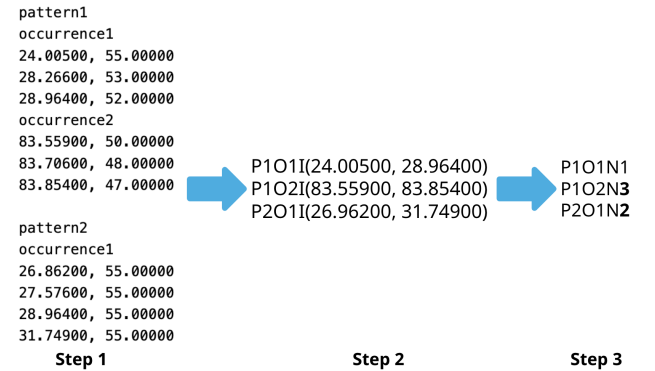


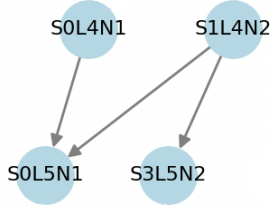
Figure 2. Parsing Motifs

In all cases, the time interval is stored as node metadata, and determines what that node's parent(s) will be.

### 3.3 STG Examples

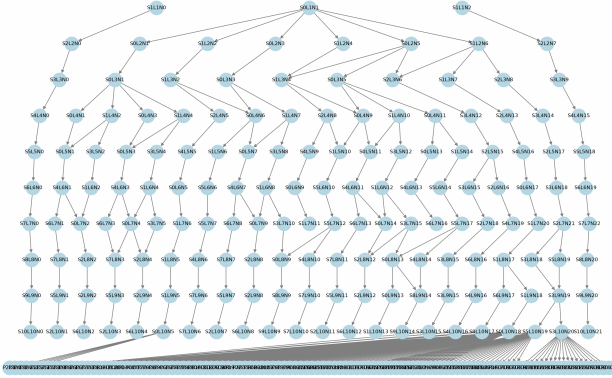
In the following examples, we use hierarchical spectral clustering from the MSAF toolkit for segmentation [1, 9], and Hsiao et al's BPS-motif discovery algorithm for motifs [14]. To better visualize the STG's structure, first con-

sider the following subgraph taken from an STG generated for the first movement exposition of the piano transcription of Beethoven’s Symphony No. 1, Op. 21 in Figure 3.<sup>2</sup> The nodes are taken from levels 4 and 5 of the segmentation sub-hierarchy. Node S3L5N2’s time interval is



**Figure 3.** STG Subgraph Example

a total subset of S1L4N2’s, while S0L5N1’s time interval starts in S0L4N1’s, but ends in S2L4N2’s. The entire STG generated for this piece, unifying hierarchical formal segmentation with disjoint motif detection, is shown in Figure 4. The segmentation nodes are easily visible, and the motif nodes form the crowded bottom layer.



**Figure 4.** STG Subgraph Example

## 4. SYNTHESIS

### 4.1 Distance Metric

In order to derive a representative "centroid" STG from a corpus of STGs, we must first define a distance metric between them. The distance between two STGs  $G_1$  and  $G_2$  is quantified by the *graph edit distance* (GED) between them, i.e. the minimum number of graph operations necessary to transform  $G_1$  into  $G_2$ . We define six valid graph operations: node and edge deletion, insertion, and substitution.

Calculating exact graph edit distance is an NP-hard problem, and thus intractable for graphs of these sizes. Fortunately, the Python package networkX offers a highly customizable approximation algorithm `optimize_edit_paths`<sup>3</sup> that efficiently computes the approximate GED and also returns the associated sequence of transforms in the edit path.

<sup>2</sup> MIDI file is from the LOP database, which can be found here: <https://qsdf.github.io/LOP/database>

<sup>3</sup> [https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.similarity.optimize\\_edit\\_paths.html](https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.similarity.optimize_edit_paths.html)

### 4.2 Stochastic Optimization

Consider the set  $G$  of STGs. With the GED given by `optimize_edit_paths` as our cost function  $c$ , we seek the centroid STG  $g^*$  that minimizes  $c$  over  $G$ . This is an NP-hard constraint satisfaction problem. Due to the graph sizes, the problem is intractable with deterministic SMT solvers like Z3. If we represent this an integer linear programming (ILP) problem, the size of each encoded graph would still blow up even the most efficient ILP solver. Thus, we turn to stochastic optimization.

Although this is still a work in progress, we introduce a framework for finding  $g^*$  using a Markov Chain Monte Carlo (MCMC) sampler. MCMC methods sample from an arbitrary probability distribution  $p$ . The equation for  $p$  is known, but it is unknown how to generate a random number from  $p$ . The idea behind MCMC is to define a Markov chain over possible  $x$  values, denoted  $\{x^{(0)}, x^{(1)}, \dots, x^{(N)}\}$  such that as  $n \rightarrow \infty$ ,  $x_n \sim p(x)$ .

The Metropolis Hastings algorithm is a popular MCMC algorithm that we adapt for the centroid graph problem. Its pseudocode is shown in Algorithm 1.

---

#### Algorithm 1 Metropolis-Hastings Algorithm

---

```

1: procedure METROPOLISHASTINGS( $p, q, x_0, N$ )
2:   Initialize  $x^{(0)} = x_0$ 
3:   for  $i = 0$  to  $N - 1$  do
4:     Propose  $t \sim q(x'|x^{(i)})$ 
5:     Apply  $t$  to  $x^{(i)}$  to obtain rewrite  $x'$ 
6:     Calculate acceptance ratio
        $\alpha = \min\left(1, \frac{p(x')q(x^{(i)}|x')}{p(x^{(i)})q(x'|x^{(i)})}\right)$ 
7:     Draw  $u \sim \text{Uniform}(0, 1)$ 
8:     if  $u < \alpha$  then
9:       Accept the proposal:  $x^{(i+1)} = x'$ 
10:    else
11:      Reject the proposal:  $x^{(i+1)} = x^{(i)}$ 
12:    end if
13:  end for
14:  return  $\{x^{(0)}, x^{(1)}, \dots, x^{(N)}\}$ 
15: end procedure

```

---

In more detail, given a probability distribution  $p$  we want to sample from, suppose that the current state of the Markov chain (i.e. the current rewrite/centroid graph) is  $x^{(i)}$ , and we want to generate  $x^{(i+1)}$ . We first generate proposal  $t$  from the proposal distribution  $q$ , which we know how to sample from. We apply  $t$  to  $x^{(i)}$  to obtain the new rewrite, or updated centroid graph,  $x'$ . Then, we execute the accept-reject step by computing the acceptance probability  $\alpha$  (line 5 of Algorithm 1). Consider the ratio  $\frac{q(x^{(i)}|x')}{q(x'|x^{(i)})}$  that appears in the  $\alpha$  equation.  $q(x'|x^{(i)})$  describes the probability that proposal  $t$  is chosen given the current state  $x^{(i)}$ .  $q(x^{(i)}|x')$ , however, is the probability that this proposal is *undone*, i.e. the probability that we choose the inverse of  $t$  to return from the new rewrite  $x'$  to the current rewrite  $x^{(i)}$ . If the proposal distribution is symmetric, these probabilities are equal. We accept the rewrite  $x'$  if the uniformly generated  $u$  falls in  $\alpha$ , setting



$x^{(i+1)} = x'$ , and reject otherwise. In the limit, we are guaranteed to arrive at the optimal rewrite  $x^{(N)}$  [29].

We adapt this algorithm to find the centroid graph. Our target distribution  $p$  is derived from the discrete cost function  $c$ , which recall is the label-aware GED given by `optimize_edit_paths`. In order to convert this into a continuous probability distribution, we use the technique suggested by Schkufza et al. [30] in Equation 1.

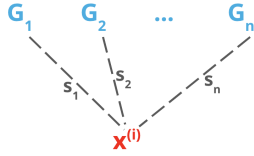
$$p(x') = \frac{1}{Z} \exp(-\beta \cdot c(x')) \quad (1)$$

where  $\beta$  is the temperature parameter, and  $Z$  is the normalizing constant. Computing  $Z$  is normally intractable, but since we take the ratio of  $p$  in  $\alpha$ , this becomes unnecessary as the  $Z$ 's cancel out. We define our  $c(x')$  in Equation 2.

$$c(x') = \sum_g^G \text{optimize\_edit\_paths}(x', g) \quad (2)$$

In other words,  $c(x')$  is the sum of the GED between itself and each STG  $g$  in the corpus.

We now consider the proposal distribution  $q$ , where we sample a potential transform  $t$  that will be applied to the current rewrite (i.e. centroid graph)  $x^{(i)}$ . An example  $t$  would be “add node S0L7N4”. In order to strategically define the space of possible transforms, we consider the other value `optimize_edit_paths` returns: the sequence of transforms constituting the edit path. Consider Figure 5. Each edit path  $s_j$  from  $x^{(i)}$  to a graph  $G_j$  in the corpus



**Figure 5.** Transform Paths from  $x^{(i)}$

consists of a sequence of transforms  $(t_{j_1}, t_{j_2}, \dots, t_{j_m})$  that changes  $x^{(i)}$  into  $G_j$ .

We collect all  $t$  values from the intersection of  $s_1, s_2, \dots, s_n$  to form the set  $T$ , which represents the space of transforms, or proposals, that we will sample from. In order to prioritize more likely transforms, we compute the count of each transform in  $T$ , and assign its count as its weight value. This results in the asymmetric proposal distribution  $q$ . It is likely that many transforms will not have inverses in  $T$ , since  $T$  only includes the transforms that occur in the edit paths. We solve this by manually adding the nonexistent transform inverses with counts of zero, and then perform additive smoothing to give them very small nonzero probabilities for the sake of computing  $\alpha$ . Our  $q$  is also *adaptive*, meaning it changes at each iteration.

There are some problems with this approach. We are still unsure whether the adaptive proposal distribution  $q$  undermines the Markov property of the chain, and therefore convergence of the algorithm. Furthermore, all MCMC algorithms are guaranteed convergence in the infinite limit, but we are unsure how to determine when the algorithm

sufficiently converges in practice. The former could be solved by discarding the custom (although theoretically more efficient) proposal distribution to make it symmetric and constant, but the latter remains unclear. Furthermore, there are two levels of approximation occurring here: one at the level of the GED itself, and another at the level of the MCMC process. This introduces a potentially unacceptable level of approximation in the result. For these reasons, we may shift our approach to other optimization strategies, such as a custom greedy edit distance algorithm.

## 5. CONCLUSIONS AND FUTURE WORK

In this paper we have presented the semantic temporal graph, a unified meta-representation of complete musical structure that encapsulates the both the intrinsic compositional hierarchy and the temporal relationships between elements of adjacent levels. Such a model is essential to comprehending the fusion of the fundamental components that constitute the piece’s core architecture, something cannot be done by considering each level in isolation. The STG thus is a vital contribution to the holistic cognition and analysis of musical form.

Furthermore, we have introduced a novel approach rooted in stochastic optimization for deriving an STG representative of a dataset of music. By defining the GED between two STGs as a cost function to minimize over the corpus, we are able to describe an application of the Metropolis Hastings algorithm allowing us to construct the “centroid” STG for the corpus that minimizes this cost by construction. To our knowledge, this is the inaugural presentation of a method for discerning the overarching structure of an entire dataset, rather than just individual pieces.

However, due to the uncertainties we have regarding the convergence and magnitude of approximation of the algorithm, we intend branch out and explore a custom edit distance algorithm that builds the centroid greedily. We also will add the remaining layers of the compositional hierarchy to the STG; namely, rhythmic contour, functional harmonic contour, and melodic contour, and derive the centroid using these more robust STGs. In addition, we plan to conduct a qualitative evaluation of the STG as a concept, as well as a quantitative evaluation of the centroid STG.

Finally, in future research we plan to use the STG as a system of constraints to induce structure in generative models, particularly transformers. Such models suffer from what is known as the *long sequence problem*, in that they cannot handle very long sequences as input since their self-attention operations exhibit quadratic time and space complexity [31]. This results in a lack of global relational cohesion in their outputs. By using the STG as a system of unified musical constraints, we will explore how to *limit*, rather than condition, the model to adhere to this structure. Similar work has been done by Young et al. [28], Wang et al. [32], and Dai et al. [33], though no one has ever attempted to limit, rather than condition, a model to a complete hierarchy of music theoretic constraints.

## 6. REFERENCES

- [1] O. Nieto, “Discovering structure in music: Automatic approaches and perceptual evaluations,” Ph.D. dissertation, New York University, 2015. [Online]. Available: <https://www.proquest.com/openview/09f67403121bcb7d2ee431985bf0568/1>
- [2] J. Foote, “Automatic audio segmentation using a measure of audio novelty,” in *2000 IEEE International Conference on Multimedia and Expo. ICME2000. Proceedings. Latest Advances in the Fast Changing World of Multimedia (Cat. No.00TH8532)*, vol. 1, 2000, pp. 452–455 vol.1.
- [3] J. Serrà, M. Müller, P. Grosche, and J. L. Arcos, “Unsupervised music structure annotation by time series structure features and segment similarity,” *IEEE Transactions on Multimedia*, vol. 16, no. 5, pp. 1229–1240, 2014.
- [4] C.-i. Wang and G. J. Mysore, “Structural segmentation with the variable markov oracle and boundary adjustment,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 291–295.
- [5] O. Nieto and T. Jehan, “Convex non-negative matrix factorization for automatic music structure identification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2013, Vancouver, BC, Canada, May 26-31, 2013*. IEEE, 2013, pp. 236–240. [Online]. Available: <https://doi.org/10.1109/ICASSP.2013.6637644>
- [6] O. Nieto and J. P. Bello, “Music segment similarity using 2d-fourier magnitude coefficients,” in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 664–668.
- [7] B. McFee, O. Nieto, M. M. Farbood, and J. P. Bello, “Evaluating hierarchical structure in music annotations,” *Frontiers in Psychology*, vol. 8, 2017. [Online]. Available: <https://www.frontiersin.org/journals/psychology/articles/10.3389/fpsyg.2017.01337>
- [8] B. McFee and D. P. W. Ellis, “Learning to segment songs with ordinal linear discriminant analysis,” in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2014, Florence, Italy, May 4-9, 2014*. IEEE, 2014, pp. 5197–5201. [Online]. Available: <https://doi.org/10.1109/ICASSP.2014.6854594>
- [9] B. McFee and D. Ellis, “Analyzing song structure with spectral clustering,” in *Proceedings of the 15th International Society for Music Information Retrieval Conference, ISMIR 2014, Taipei, Taiwan, October 27-31, 2014*, H. Wang, Y. Yang, and J. H. Lee, Eds., 2014, pp. 405–410. [Online]. Available: [http://www.terasoft.com.tw/conf/ismir2014/proceedings/T073\\_319\\_Paper.pdf](http://www.terasoft.com.tw/conf/ismir2014/proceedings/T073_319_Paper.pdf)
- [10] C. Hernandez-Olivan, S. R. Llamas, and J. R. Beltrán, “Symbolic music structure analysis with graph representations and changepoint detection methods,” *CoRR*, vol. abs/2303.13881, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2303.13881>
- [11] C. Finkensiep, M. Haeberle, F. Eisenbrand, M. Neuwirth, and M. Rohrmeier, “Repetition-structure inference with formal prototypes,” in *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, A. Sarti, F. Antonacci, M. Sandler, P. Bestagini, S. Dixon, B. Liang, G. Richard, and J. Pauwels, Eds., 2023, pp. 383–390. [Online]. Available: <https://doi.org/10.5281/zenodo.10265305>
- [12] M. Buisson, B. McFee, S. Essid, and H. C. Crayencour, “Learning multi-level representations for hierarchical music structure analysis,” in *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, P. Rao, H. A. Murthy, A. Srinivasamurthy, R. M. Bittner, R. C. Repetto, M. Goto, X. Serra, and M. Miron, Eds., 2022, pp. 591–597. [Online]. Available: <https://archives.ismir.net/ismir2022/paper/000071.pdf>
- [13] S. Dai, H. Zhang, and R. B. Dannenberg, “Automatic analysis and influence of hierarchical structure on melody, rhythm and harmony in popular music,” *CoRR*, vol. abs/2010.07518, 2020. [Online]. Available: <https://arxiv.org/abs/2010.07518>
- [14] Y. Hsiao, T. Hung, T. Chen, and L. Su, “Bps-motif: A dataset for repeated pattern discovery of polyphonic symbolic music,” in *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, A. Sarti, F. Antonacci, M. Sandler, P. Bestagini, S. Dixon, B. Liang, G. Richard, and J. Pauwels, Eds., 2023, pp. 281–288. [Online]. Available: <https://doi.org/10.5281/zenodo.10265277>
- [15] C. Wang, J. Hsu, and S. Dubnov, “Music pattern discovery with variable markov oracle: A unified approach to symbolic and audio representations,” in *Proceedings of the 16th International Society for Music Information Retrieval Conference, ISMIR 2015, Málaga, Spain, October 26-30, 2015*, M. Müller and F. Wiering, Eds., 2015, pp. 176–182. [Online]. Available: [http://ismir2015.uma.es/articles/78\\_Paper.pdf](http://ismir2015.uma.es/articles/78_Paper.pdf)
- [16] G. Velarde, D. Meredith, and T. Weyde, *A Wavelet-Based Approach to Pattern Discovery in Melodies*. Cham: Springer International Publishing, 2016, pp. 303–333. [Online]. Available: [https://doi.org/10.1007/978-3-319-25931-4\\_12](https://doi.org/10.1007/978-3-319-25931-4_12)
- [17] T. Chen and L. Su, “Harmony transformer: Incorporating chord segmentation into harmony recognition,”

- in *Proceedings of the 20th International Society for Music Information Retrieval Conference, ISMIR 2019, Delft, The Netherlands, November 4-8, 2019*, A. Flexer, G. Peeters, J. Urbano, and A. Volk, Eds., 2019, pp. 259–267. [Online]. Available: <http://archives.ismir.net/ismir2019/paper/000030.pdf>
- [18] —, “Functional harmony recognition of symbolic music data with multi-task recurrent neural networks,” in *Proceedings of the 19th International Society for Music Information Retrieval Conference, ISMIR 2018, Paris, France, September 23-27, 2018*, E. Gómez, X. Hu, E. Humphrey, and E. Benetos, Eds., 2018, pp. 90–97. [Online]. Available: [http://ismir2018.ircam.fr/doc/pdfs/178\\_Paper.pdf](http://ismir2018.ircam.fr/doc/pdfs/178_Paper.pdf)
- [19] J. Salamon, E. Gomez, D. P. W. Ellis, and G. Richard, “Melody extraction from polyphonic music signals: Approaches, applications, and challenges,” *IEEE Signal Processing Magazine*, vol. 31, no. 2, pp. 118–134, 2014.
- [20] K. Kosta, W. T. Lu, G. Medeot, and P. Chanquion, “A deep learning method for melody extraction from a polyphonic symbolic music representation,” in *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, P. Rao, H. A. Murthy, A. Srinivasamurthy, R. M. Bittner, R. C. Repetto, M. Goto, X. Serra, and M. Miron, Eds., 2022, pp. 757–763. [Online]. Available: <https://archives.ismir.net/ismir2022/paper/000091.pdf>
- [21] Y.-H. Chou, I.-C. Chen, C.-J. Chang, J. Ching, and Y.-H. Yang, “Midibert-piano: Large-scale pre-training for symbolic music understanding,” 2021.
- [22] M. McCallum. (2017) Structural Segmentation. Music Information Retrieval Evaluation eXchange (MIREX). [Online]. Available: [https://www.music-ir.org/mirex/wiki/2017:Structural\\_Segmentation](https://www.music-ir.org/mirex/wiki/2017:Structural_Segmentation)
- [23] T. Collins. (2017) Discovery of Repeated Themes & Sections. Music Information Retrieval Evaluation eXchange (MIREX). [Online]. Available: [https://www.music-ir.org/mirex/wiki/2017:Discovery\\_of\\_Repeated\\_Themes\\_%26\\_Sections](https://www.music-ir.org/mirex/wiki/2017:Discovery_of_Repeated_Themes_%26_Sections)
- [24] Sutashu. (2010) Harmonic Analysis. Music Information Retrieval Evaluation eXchange (MIREX). [Online]. Available: [https://www.music-ir.org/mirex/wiki/2010:Harmonic\\_Analysis](https://www.music-ir.org/mirex/wiki/2010:Harmonic_Analysis)
- [25] D. Evans. (2021) Audio Melody Extraction. Music Information Retrieval Evaluation eXchange (MIREX). [Online]. Available: [https://www.music-ir.org/mirex/wiki/2021:Audio\\_Melody\\_Extraction](https://www.music-ir.org/mirex/wiki/2021:Audio_Melody_Extraction)
- [26] J. Wang, J. B. L. Smith, W. T. Lu, and X. Song, “Supervised metric learning for music structure features,” in *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, J. H. Lee, A. Lerch, Z. Duan, J. Nam, P. Rao, P. van Kranenburg, and A. Srinivasamurthy, Eds., 2021, pp. 730–737. [Online]. Available: <https://archives.ismir.net/ismir2021/paper/000091.pdf>
- [27] M. C. McCallum, “Unsupervised learning of deep features for music segmentation,” *CoRR*, vol. abs/2108.12955, 2021. [Online]. Available: <https://arxiv.org/abs/2108.12955>
- [28] H. Young, M. Du, and O. Bastani, “Neurosymbolic deep generative models for sequence data with relational constraints,” in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 37 254–37 266. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/f13ceb1b94145aad0e54186373cc86d7-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/f13ceb1b94145aad0e54186373cc86d7-Paper-Conference.pdf)
- [29] D. Navarro. The Metropolis-Hastings algorithm. [Online]. Available: [https://blog.djnavarro.net/posts/2023-04-12\\_metropolis-hastings](https://blog.djnavarro.net/posts/2023-04-12_metropolis-hastings)
- [30] E. Schkufza, R. Sharma, and A. Aiken, “Stochastic program optimization,” *Commun. ACM*, vol. 59, no. 2, pp. 114–122, 2016. [Online]. Available: <https://doi.org/10.1145/2863701>
- [31] Y. Mao, M. Ester, and K. Li, “Iceformer: Accelerated inference with long-sequence transformers on CPUs,” in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=6RR3wU4mSZ>
- [32] Z. Wang, L. Min, and G. Xia, “Whole-song hierarchical generation of symbolic music using cascaded diffusion models,” in *The Twelfth International Conference on Learning Representations*, 2024. [Online]. Available: <https://openreview.net/forum?id=sn7CYWyavh>
- [33] S. Dai, Z. Jin, C. Gomes, and R. B. Dannenberg, “Controllable deep melody generation via hierarchical music structure representation,” *CoRR*, vol. abs/2109.00663, 2021. [Online]. Available: <https://arxiv.org/abs/2109.00663>