

Modeling the Latent Space of Symbolic String Quartet Music

Alex Kylo
akyllo@uw.edu
University of Washington
Bothell, WA, USA

ABSTRACT

This paper presents the results of a project to apply a recurrent variational autoencoder to the task of modeling a latent space for generating pieces of multi-instrument symbolic classical music through interpolation.

KEYWORDS

deep learning, recurrent neural networks, sequential models, music modeling, latent space modeling, generative modeling

1 INTRODUCTION

Machine learning models of music have interesting applications in music information retrieval and creative tools for musical artists and educators. Generative models can create accompaniments for music, blend and transfer styles between clips of music, and even generate entirely new music. Music is challenging to model because it is inherently high-dimensional, exhibiting a complex hierarchy of recurring patterns and long-range temporal dependencies, and because musical scores have multiple possible digital representations with distinct advantages and disadvantages.

Depending on the task, machine learning models of music may be trained on the audio signal of a musical performance, either in a time domain or a frequency domain representation, or they may be trained on a digital symbolic representation of music, the most common of which is MIDI (Musical Instrument Digital Interface) notation. MIDI is an encoding of music as streams of bytes in one or more tracks or channels, each representing a sequence of 128 possible pitch values (where 0 is the lowest pitch and 127 is the highest), along with timing, pressure and instrument identifier values. A generative symbolic music model can produce a symbolic score in MIDI format, which must be played by synthesizer software or by humans to produce an audio music performance. This project focuses on generative modeling of symbolic (MIDI) music to compose original musical scores by blending input scores through continuous latent space interpolation.

2 RELATED WORK

The state of the art in music generation still has a long way to go before it can consistently generate music scores or performances that would be enjoyable and popular for humans to listen to, but for this reason it is an area of significant opportunity, where a number of recent research projects have shown promising progress.

Google’s Magenta is an umbrella project for music deep learning research and development of software tools to expose these models for use by creative artists and students.

MusicVAE, part of the Magenta project, is a variational Long Short-Term Memory (LSTM) autoencoder for MIDI that incorporates a novel hierarchical structure using a “conductor” recurrent

layer in its decoder model to better capture structure at multiple levels and avoid the problem of “posterior/mode collapse” whereby a generative model learns to ignore its latent code and rely on autoregression [13]. This model is trained on 16-bar paragraphs of music and is capable of generating new melodies that blend two given melodies via latent space interpolation.

Another Magenta model called Music Transformer is a generative model that borrows its approach from the Natural Language Processing (NLP) domain, using a self-attention network to model MIDI music as a sequence of discrete tokens with relative positional dependencies [5]. The focus of this model is on learning long-term dependencies in music to produce longer clips of music with coherent structure. Music Transformer was trained on a dataset of Piano-e-competition performances [4] and its generated piano music received favorable qualitative (Likert scale) ratings from human listeners for its resemblance to human-composed music [5].

MuseGAN [3] is an application of Generative Adversarial Networks (GAN) to polyphonic MIDI music generation, trained on four-bar phrases of a multi-track pianoroll representation of rock songs from the Lakh Midi Dataset [11]. Like MusicVAE, MuseGAN includes a two-level generator that first samples latent codes at the phrase or bar level, then generates notes within the bars, to produce longer-term structural patterns.

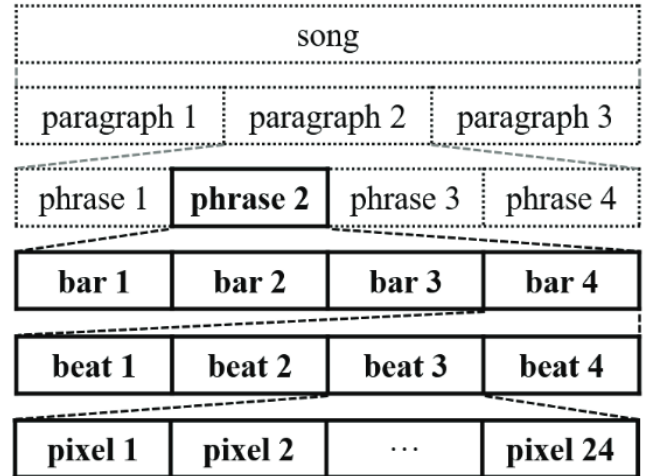


Figure 1: Diagram of the hierarchical structure of a musical composition, from the MuseGAN paper [3].

A major advantage of working with the symbolic representation of music is that it is of far lower dimensionality than the raw audio waveforms of a recorded performance, which makes it less computationally expensive. However, there are many stylistic aspects

of musical performance that are not captured by a symbolic representation, and may be specific to a particular performer, so the expressiveness of symbolic generative models is limited in comparison [7].

Other research has focused on modeling raw audio waveforms directly. WaveNet is a causal convolutional neural network for generating raw audio waveforms, developed by Google DeepMind, which achieves state of the art performance in generating natural sounding speech from text, but is also capable of generating short, realistic snippets of audio music [9]. Another model named SampleRNN generates raw audio waveforms using a three-tier hierarchy of gated recurrent units (GRU) to model recurrent structure at multiple temporal resolutions [8].

Prior work points out that the division between symbolic music notes and music performances, is analogous to the division between symbolic language and utterances in speech, which may inspire ideas for combining the two approaches [4]. A paper from Boston University describes an effort to combine the symbolic and waveform approaches to music modeling, by training an LSTM to learn melodic structure of different styles of music, then providing generations from this model as conditioning inputs to a WaveNet-based raw audio generator [7].

3 METHODS

3.1 Datasets

The dataset used for this research project is MusicNet, which is a collection of 330 freely licensed European classical music recordings with aligned MIDI scores [14]. The model is trained on 36 string quartets (four-part arrangements with two violins, one viola and one cello) by composers Bach, Beethoven, Dvorak, Haydn, Mozart, and Ravel.

3.2 Data Preprocessing

Several choices must be made in how to preprocess binary MIDI files into training examples for a neural network. There are multiple open-source Python packages that assist with the process of reading MIDI files from their binary on-disk representations into Python objects, such as pretty_midi [12], Pypianoroll [2] and music21 [1].

In order to accommodate polyphonic music, we each MIDI file is converted into a modified pianoroll representation, an example of which is visualized in Figure 2. In this modified representation, each instrument track is a sequence of integers from 0-129 representing the MIDI note pitch value being played at the given time step, where values 0-127 are pitches, 128 is a rest (silence) and 129 is sustain, indicating that the previously played pitch is held continuously. Because chords (multiple notes being played simultaneously by the same instrument) are relatively rare in string quartets, the lowest note of each chord is taken and the rest discarded, to reduce dimensionality. Because notes at the extreme low and high end are rare and are out of range of common instruments, the note pitch values are clipped to a narrower range and ordinal encoded. For the 36 string quartets found in the MusicNet database, this results in 66 distinct values, where 0-63 represent notes played by the violins, 64 is rest and 65 is sustain. Our software package also includes functions to convert this representation back into a Music21 Score object, which can then be written to a MIDI file.

Because songs are typically at least a few minutes long and of varying length, it will not be feasible to train with entire songs as examples, so we will crop songs into phrases of equal numbers of measures to use as training data.

The result of this preprocessing is that each training example is a 2D matrix of shape (time steps x instrument tracks) and stacking the training examples produces a 3D tensor. Information regarding which track corresponds to which musical instrument is lost and therefore must be maintained separately to reconstruct the MIDI representation. For string quartets this instrument ensemble is always [40, 40, 41, 42], which is an array of MIDI instrument codes representing two violins, one viola and one cello. Tempo information is also lost from this representation, so the outputs are normalized to the default tempo of 120 beats per minute.

As the model will produce one output per instrument track, it can incorporate only a fixed selection of instrument parts, similar to how MusicVAE models three-part (drum, bass and melody) [13] and MuseGAN models five-part (drum, bass, guitar, string, piano) arrangements. Therefore this model is purpose-built for string quartets, which are the most common type of arrangement found in the MusicNet dataset.

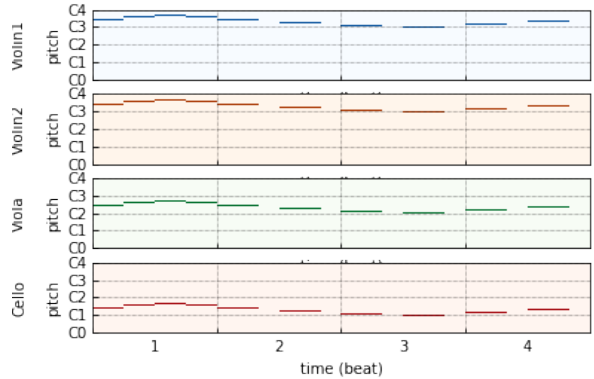


Figure 2: Pianoroll visualization of the first measure of Beethoven's Serioso String Quartet in F Minor



Figure 3: Sheet music of the first measure of Beethoven's Serioso String Quartet in F Minor

Data augmentation is also possible and we assess its impact on the results—the literature suggests augmentation via pitch shifting each training example up or down by up to six semitones, and increasing or reducing the speed by up to 10% in order to create additional training examples and reduce overfitting [10].

The multi-stream representation discussed in [6] addresses the issue of extreme class imbalance and sparsity in the pianoroll representation.

3.3 Model Design

The model is a recurrent variational autoencoder; both Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) cell types were tested.

The encoder model utilizes an embedding layer to encode the note pitch integers as vectors; embedding vector length is a tunable hyperparameter and we found the best results with it set to 8. Then two LSTM layers, with a dropout layer in between, convert the sequences into latent codes, which are modeled as a multidimensional Gaussian distribution.

The decoder model samples from the latent space z with the distribution parameterized by the μ and $\log(\sigma)$ values learned by the encoder.

TODO: model architecture plots

3.4 Model Evaluation

Evaluation of generative models is challenging because there is no equivalent of an accuracy metric like what is used in supervised learning. For autoencoder models we can measure how accurately the model can reconstruct its own inputs, but this does not tell us the quality of the interpolated examples. Generative models are typically evaluated using a combination of qualitative metrics whereby human judges rate the quality of the generated examples (essentially a Turing test), and quantitative metrics that assess the differences in the parametric distributions of generated and real examples. Yang and Lerch (2020) proposes a set of metrics informed by music theory, for probabilistically evaluating how similar the generations are to known sample distributions of real music [15]. These metrics include counts, ranges, histograms and transition matrices of pitches and note lengths, then the Kullback-Leibler divergence and overlapping area of the probability density functions are used to compare against known reference distributions per musical genre [15]. Due to the cost and time requirements associated with designing a human subjects experiment, we utilize this quantitative approach to the quality assessment of generated samples.

4 RESULTS

TODO: learning curve plot

5 DISCUSSION

REFERENCES

- [1] Michael Scott Cuthbert and Christopher Ariza. [n.d.]. music21: A Toolkit for Computer-Aided Musicology and Symbolic Music Data. ([n. d.]), 6.
- [2] Hao-Wen Dong and Wen-Yi Hsiao. 2018. Pypianoroll: Open Source Python Package for Handling Multitrack Pianoroll. (2018), 2.
- [3] Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. 2017. MuseGAN: Multi-track Sequential Generative Adversarial Networks for Symbolic Music Generation and Accompaniment. arXiv:1709.06298 [eess.AS]
- [4] Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. 2019. Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset. arXiv:1810.12247 [cs.SD]
- [5] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinulescu, and Douglas Eck. 2018. Music Transformer. (2018). arXiv:1809.04281 <http://arxiv.org/abs/1809.04281>
- [6] Harish Kumar and Balaraman Ravindran. 2019. Polyphonic Music Composition with LSTM Neural Networks and Reinforcement Learning. arXiv:1902.01973 [cs.SD]
- [7] Rachel Manzeili, Vijay Thakkar, Ali Siahkamari, and Brian Kulis. 2018. Conditioning Deep Generative Raw Audio Models for Structured Automatic Music. CoRR abs/1806.09905 (2018). arXiv:1806.09905 <http://arxiv.org/abs/1806.09905>
- [8] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo, Aaron Courville, and Yoshua Bengio. 2017. SampleRNN: An Unconditional End-to-End Neural Audio Generation Model. (2017). arXiv:1612.07837 <http://arxiv.org/abs/1612.07837>
- [9] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. (2016). arXiv:1609.03499 <http://arxiv.org/abs/1609.03499>
- [10] Sageev Oore, Ian Simon, Sander Dieleman, Douglas Eck, and Karen Simonyan. 2018. This Time with Feeling: Learning Expressive Musical Performance. (2018). arXiv:1808.03715 <http://arxiv.org/abs/1808.03715>
- [11] Colin Raffel. 2016. Learning-Based Methods for Comparing Sequences, with Applications to Audio-to-MIDI Alignment and Matching.
- [12] Colin Raffel and Daniel P.W. Ellis. 2018. Intuitive Analysis, Creation and Manipulation of MIDI Data With pretty_midi. (2018), 2.
- [13] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. 2018. A Hierarchical Latent Vector Model for Learning Long-Term Structure in Music. In *Proceedings of the 35th International Conference on Machine Learning*. 10.
- [14] John Thickstun, Zaid Harchaoui, and Sham Kakade. 2017. Learning Features of Music from Scratch. arXiv:1611.09827 [stat.ML]
- [15] Li-Chia Yang and Alexander Lerch. 2020. On the evaluation of generative models in music. 32, 9 (2020), 4773–4784. <https://doi.org/10.1007/s00521-018-3849-7>