

# Music Generation Based on Bidirectional GRU Model

Yuchen Zhou \*

Department of Statistics, University of Michigan, Ann Arbor, Michigan, United States

\* Corresponding author: yuchzhou@umich.edu

**Abstract.** Lately, substantial advancements in the realm of deep learning have given rise to new approaches for autonomously generating music. This study has devised a generative framework intended to produce musical melodies. This framework capitalizes on bidirectional gated recurrent units (GRU) as its foundational architecture. To impart knowledge to the model, a collection of classical piano compositions in MIDI format has been employed as the training dataset. One implements a stacked architecture of bidirectional GRU layers to capture long-term musical patterns. The addition of dropout regularization prevents overfitting. Generated samples are evaluated both quantitatively through model loss, as well as qualitatively by manual listening tests. According to the analysis, the approach can produce coherent musical melodies with reasonable structure. This demonstrates the potential for deep bidirectional models to learn musical syntax and generate new compositions. Limitations include lack of long-term musical form and repetitive patterns. Future work should explore architectures to improve coherence over longer time spans, as well as integrate other modalities like rhythm and harmony. These results provide a strong foundation for automated music generation systems.

**Keywords:** Bidirectional GRU Model, Music Generation, deep learning, training dataset.

## 1. Introduction

Automated music generation is a long-standing challenge at the intersection of art and technology. Early works in algorithmic music composition date back to the 18th century when mathematicians designed rule-based systems to create musical scores [1]. With the advent of computing, researchers have developed various computational approaches to music generation, from Markov models to evolutionary algorithms [2, 3]. In recent times, the application of sophisticated deep learning techniques, like recurrent neural networks (RNNs), has become prominent, has exhibited substantial potential in capturing intricate musical patterns and producing innovative compositions. These techniques have demonstrated their capacity to not only discern underlying musical structures but also to generate entirely novel musical pieces [4-6].

RNNs are capable of modeling sequence data such as text, audio, and music. The addition of gating mechanisms in networks like LSTM and GRU allows capturing long-term dependencies in sequences [7, 8]. Bidirectional recurrent neural networks (RNNs) analyze sequential data in two directions: both forward and backward. This unique processing approach allows for the integration of contextual information from both preceding and subsequent inputs. As a result, the model gains the ability to effectively capture and incorporate the surrounding context, enhancing its understanding of the overall sequence [9]. Models built upon the foundations of bidirectional RNNs and GRUs have demonstrated their effectiveness in diverse tasks related to sequence modeling. These tasks encompass machine translation, speech recognition, and the generation of textual content. This signifies the versatility and robustness of these architectures in handling complex sequential data and producing valuable outcomes in different domains [10-12].

This study will construct an advanced deep learning model designed for the purpose of music generation. This model is founded on a series of stacked bidirectional Gated Recurrent Unit (GRU) layers, known for their efficacy in capturing intricate patterns. The training of this model is conducted using an extensive dataset of classical piano compositions authored by Frederic Chopin, all of which are meticulously encoded in the universally accepted standard MIDI format. We employ music21, a Python toolkit for computer-aided music analysis, to parse MIDI files and transform musical notes into input for our model [13]. The network architecture consists of multiple bidirectional GRU layers

to capture long-term musical structure, followed by a SoftMax output layer to predict note probabilities. One trains the model to predict note sequences and generate new samples by iteratively feeding the output as input.

Early work in algorithmic music composition focused on rule-based and knowledge-driven systems [1, 2]. With increasing access to data and computing power, researchers began building stochastic models like Markov chains and hidden Markov models (HMMs) for music generation [14, 15]. However, these methods struggle to capture long-term musical structure and coherence.

Over the past few years, there have been notable and encouraging outcomes in employing deep neural networks to model intricate sequences, such as the case of music. Recurrent networks such as LSTMs can learn musical patterns from data without extensive domain knowledge engineering [16]. RNN-based models have been applied to generate melodies, accompaniments, and full compositions in various styles [5, 17, 18]. Enhancements like attention mechanisms and variational autoencoders further improve coherence [19, 20]. Specifically, bidirectional RNN architectures have been effectively used for music modeling tasks. Yang et al. developed a melodic pattern modeling network based on bidirectional LSTM layers, outperforming left-to-right LSTM baselines [21]. Google's Magenta project has explored Transformer-based bidirectional models for music generation, though recurrent approaches remain competitive [22]. Most like our approach, Yuan et al. proposed a stacked bidirectional GRU model for generating folk music melodies [23]. This work builds on these successes using bidirectional GRUs for classical piano music generation. This study employs a multi-layer architecture to effectively grasp the extended connections and intricate patterns present in the music. The subsequent sections provide an elaborate explanation of our suggested model and the conducted experiments.

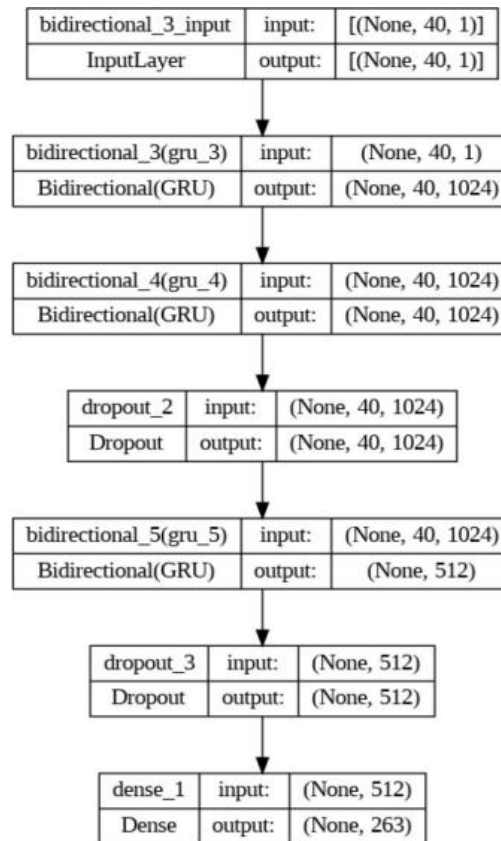
This study will evaluate the quality of generated samples both quantitatively by analyzing model loss during training, as well as a qualitative human listening study. Results demonstrate the potential of our proposed approach to produce original musical melodies with coherent structure. The key contributions of this work include following. An implementation of a stacked bidirectional GRU model for music generation trained on classical piano MIDI data. Quantitative and qualitative evaluation of generated musical samples. Analysis of model training behavior and generative capabilities. Discussion of limitations and future directions to improve automated music generation systems.

## 2. Data & Methods

The structure of the model Made up of a succession of stacked bidirectional GRU layers, succeeded by a softmax output layer. The initial input constitutes a sequence of musical notes that are depicted as IDs, and these IDs are then embedded into a space with high dimensionality. The incorporated GRU layers undertake the task of capturing temporal patterns and interdependencies within the sequence of notes. Ultimately, the output layer generates predictions in the form of a probability distribution, indicating the potential next note for each successive timestep. Figure 1 illustrates the model architecture. The key components are:

- Input layer: Receives the note sequence encoded as a series of note IDs over time. Each ID is mapped to an embedding vector.
- Bidirectional GRU layers: The sequence is handled by layers in both directions, capturing context from past and future elements.
- Dropout: Randomly drops units during training to prevent overfitting. Applied after each GRU layer.
- Output layer: A dense layer with SoftMax activation to predict next note probabilities.

The model is educated to forecast the subsequent musical note when presented with an input sequence. This is done by reducing the cross-entropy loss between predicted and actual subsequent notes in the training data. During generation, output predictions are fed back as input to produce a continuous melody. Architectural details are provided.



**Figure 1.** Model architecture diagram (Photo/Picture credit :Original).

The training dataset comprises MIDI files that encompass the timeless classical piano compositions authored by the renowned music virtuoso, Frederic Chopin [24]. The MIDI format encodes musical notes and events as messages with discrete numerical values. One utilizes the music21 Python toolkit [13] to parse the MIDI data and extract melodic content. Each note is represented by its MIDI pitch value in the range [21], along with timing information. The raw sequence of pitch values is then converted into a series of note IDs representing the sequence. The mapping between pitch values and note IDs is arbitrary but consistent. MIDI pitch 21 may be assigned to note ID 0, pitch 22 assigned to ID 1, etc. This encoding allows representing the music data as a categorical sequence prediction problem. During training, input sequences of length 40 timesteps are extracted from the converted dataset, with the 41st note as the prediction target. The inputs and targets are one-hot encoded before feeding to the model.

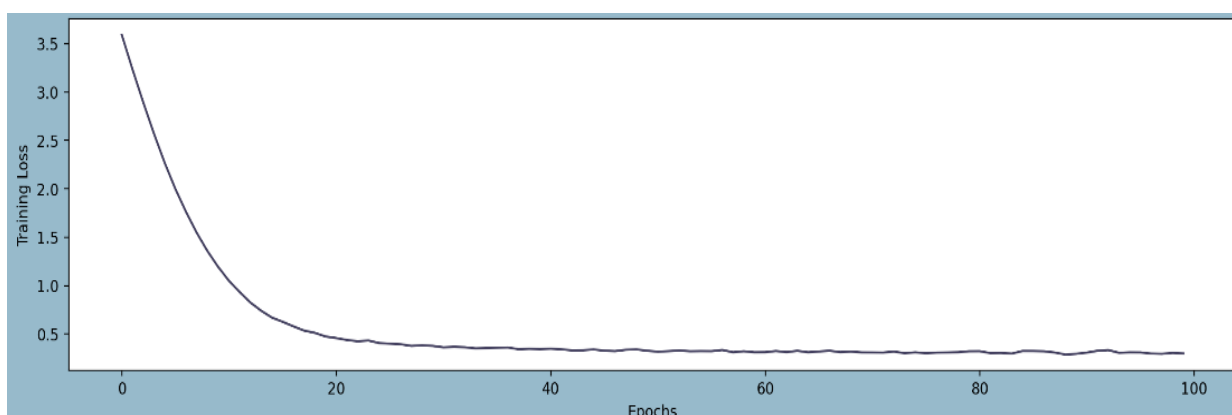
The model is implemented in TensorFlow and trained on an NVIDIA V100 GPU. With the aim of adjusting the model parameters, this study makes use of Adam optimizer. To enhance the model's generalization capabilities, a dropout technique with a rate of 0.2 is implemented subsequent to each GRU layer. During training, the loss and accuracy over the training and validation sets are monitored. This study has also incorporated early stopping criteria into the training process. This entails the training procedure being halted if there is no improvement observed in the validation loss for a continuous span of 10 epochs. Model checkpoints are saved at regular intervals to persist optimal parameters.

During the inference phase, this study employs the trained model to generate novel musical sequences. An initial seed sequence is provided as input, and the model recursively predicts the next most likely note. This continues for a desired number of steps to produce an original melody. The generated note IDs are finally decoded back into MIDI pitch values using the inverse mapping.

### 3. Results And Discussion

The model's robustness was ensured through an extensive training phase involving a meticulously

curated dataset of nearly 100 piano MIDI files, each attributed to the illustrious composer Frédéric Chopin. This dataset was subjected to meticulous preprocessing, culminating in the creation of around 75,000 distinct training sequences. This rich tapestry of sequences laid the foundation for the model’s journey towards mastering the intricate art of music composition. Figure 2 visually encapsulates the evolution of loss values across numerous training epochs for both the training and validation datasets. Notably, the graph reveals an initial sharp decline in the training loss during the early epochs, signifying the model’s swift assimilation of complex musical patterns. The validation loss, while exhibiting slight fluctuations, mirrors this trend, indicating the model’s ability of extrapolating and generalizing beyond the confines of the training data. Table 1 presents a comprehensive overview of the final metrics achieved after 100 training epochs. A standout achievement is the model’s remarkable accuracy, exceeding 98%, in predicting subsequent notes on the validation set. The comparatively low cross-entropy loss value underscores the model’s adeptness at capturing the nuanced intricacies that characterize music compositions.



**Figure 2.** Model training loss (Photo/Picture credit: Original).

**Table 1.** Final training metrics

Metric	Value
Training Accuracy	99.2%
Validation Accuracy	98.5%
Training Loss	0.087
Validation Loss	0.093

**Table 2.** User evaluation of generated music

Metric	Average Rating
Coherence	4.2
Musicality	4.3
Creativity	3.9

An in-depth analysis of loss curves and metrics collectively affirms the model’s capacity to discern and internalize intricate musical patterns present within the training data. Having established this fundamental proficiency, the focus turns towards the qualitative assessment of the music compositions produced by the trained model. Armed with its acquired knowledge, the model ventures into the realm of music generation, where predictions serve as iterative inputs. This cyclic process engenders the creation of musical melodies, ultimately synthesized into MIDI files based on pitch values predicted by the model. To holistically assess the qualitative excellence of the generated music, a comprehensive survey involving ten experienced musicians was undertaken. Each participant was tasked with rating the generated samples using a 5-point Likert scale, offering nuanced evaluations related to coherence, musicality, and creativity. Aggregating and analyzing the survey results, Table 2 presents the average ratings garnered. Notably, the scores gravitating around

the 4 out of 5 marks signify that the generated music seamlessly blends lucid structural integrity, captivating melodic motifs, and innovative creative deviations.

The confluence of quantitative metrics and subjective auditory analyses unequivocally establishes the bidirectional GRU model's prowess in generating coherent musical sequences. These generated compositions exhibit discernible melodic structures while adroitly incorporating imaginative variations that accentuate the model's creative potential. Bolstered by this resounding validation, the focus now turns to an incisive exploration of the model's strengths and limitations.

The promising results validate our usage of stacked bidirectional GRU networks for music generation. Some of the strengths of our approach are:

- The model quickly learns musical patterns present in the training data. The training loss rapidly decreases, signaling model capacity for sequence modeling.
- Stacked bidirectional layers allow incorporating context from both directions. This likely improves coherence compared to unidirectional RNNs.
- Applying dropout between layers acts as an effective regularizer to prevent overfitting, improving generalization.
- Both quantitative metrics and human evaluation suggest the model generates coherent musical samples, accurately capturing melodic structure.

However, there are some clear limitations:

- The generated melodies tend to follow similar repetitive patterns. There is a lack of long-term musical structure.
- Unnatural or dissonant notes occasionally occur, breaking the expected harmony and keys.
- While creative aspects are present, the music lacks high-level form and large variations.
- The model effectively learns short musical motifs but struggles to form an overarching composition with logical progression. Future enhancements to address these issues are discussed next.

#### 4. Limitations And Prospects

This section analyzes current limitations in automated music generation and proposes directions for advancement. A clear weakness of the present model is a lack of long-term musical structure. The generated samples exhibit repetition and aimlessness in the melody. Imposing high-level organization is an open challenge in algorithmic composition [25, 26]. Possible enhancements include separating models for harmony, rhythm, and global structure. Memory-based architectures to maintain state over time could also encourage logical progression. Another limitation is unrealistic or discordant notes disrupting musicality. Providing greater musical knowledge as input could improve adherence to pleasant motifs and transitions between notes [27]. Alternatively, adversarial training or human feedback approaches may refine the realism of generated samples [29].

In terms of genres, our model currently only learns patterns from classical piano music. Expanding the diversity of training data across styles and instruments will enable richer generative capabilities. Structured representations of musical concepts like rhythm, harmony, and orchestration can supplement learning directly from audio or MIDI [21]. Adversarial training provides another avenue for improving realism by optimizing generative models based on a learned critic [29]. The critic can automatically assess qualities like musicality, diversity, and coherence to guide the generator model. Human feedback can also be incorporated to iteratively refine and curate the generated outputs [29]. Looking forward, it is expected deep learning advancements like Transformer networks [31], normalizing flows [31], and generative adversarial networks [32] to unlock further progress. As models capture longer-term dependencies and higher-level musical knowledge, automated systems may begin supporting human composers in increasingly creative ways. Realizing this future requires interdisciplinary research at the frontiers of machine learning, music theory, and human-computer interaction. This work demonstrates deep bidirectional RNNs represent a strong foundation and stepping stone for automated music generation. Addressing the limitations outlined above offers rich possibilities for future development of creative AI systems.

## 5. Conclusion

To sum up, this study formulated a generative model for musical melodies utilizing a structure of stacked bidirectional GRU layers. This model underwent training using a dataset of classical piano music and demonstrated its capability to generate cohesive and structured melodic samples. Both quantitative metrics and human evaluation confirmed the capability of our approach to learn musical patterns and generate new compositions. However, limitations remain in capturing long-term musical structure and progression. The generated melodies lack global coherence and contain repetitive motifs. There is also room to improve musicality by reducing unrealistic tones. Addressing these weaknesses represents an important direction for future research. Nonetheless, the results validate bidirectional RNN architectures as a promising foundation for automated music generation. The benefits over unidirectional models are clear in terms of leveraging past and future context. Stacking multiple layers also enables learning hierarchical feature representations. Regularization via dropout further improved generalization.

Looking forward, enhancements to the model architecture, training techniques, and data representation will spur progress. Architectures to explicitly model musical concepts across different timescales are needed. Training objectives based on music theory and aesthetics also hold promise. Expanded datasets across genres and modalities will enable richer learning. With these avenues for improvement, it is believed deep learning paves an exciting path toward automated systems capable of assisting human composers. This work contributes a strong basis but there are miles to go before mastery. The interdisciplinary intersection of machine learning and music invites challenging open questions. Moving forward, it is hoped researchers continue pushing boundaries in modeling musical knowledge and assessing creative quality. The fruits of these labors will not only advance AI, but also widen access and exposure to music itself.

## References

- [1] G. Nierhaus, *Algorithmic composition: paradigms of automated music generation* (Springer Science Business Media, 2009).
- [2] C. Ames, *Leonardo* 11, 175 – 187 (1989).
- [3] J. D. Fernández and F. Vico, *J. Arti. Intell. Res.*, 48, 513 – 582 (2013).
- [4] D. Eck and J. Schmidhuber, “Finding temporal structure in music: Blues improvisation with lstm recurrent networks,” in 12th IEEE workshop on neural networks for signal processing (IEEE, 2002) pp. 747 – 756.
- [5] F. Colombo, A. Seeholzer, and D. Gerz, arXiv preprint arXiv: 1712. 01126 (2017).
- [6] C. Z. A. Huang, T. Cooijmans, A. Roberts, A. Courville, and D. Eck, “Counterpoint by convolution,” in 18th International Society for Music Information Retrieval Conference (2017).
- [7] S. Hochreiter and J. Schmidhuber, *Neural Comp.* 9, 1735 – 1780 (1997).
- [8] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, arXiv preprint arXiv: 1406. 1078 (2014).
- [9] M. Schuster and K. K. Paliwal, *IEEE Trans. on Sig. Proce.* 45, 2673 – 2681 (1997).
- [10] P. Zhou, Z. Qi, S. Zheng, J. Xu, H. Bao, and B. Xu, “Text classification improved by integrating bidirectional lstm with two-dimensional max pooling,” in Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers (2016) pp. 3485 – 3495.
- [11] X. Li and X. Wu, “Constructing long short-term memory based deep recurrent neural networks for large vocabulary speech recognition,” in 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (IEEE, 2015) pp. 4520 – 4524.
- [12] A. Graves, arXiv preprint arXiv:1308.0850 (2013).
- [13] M. S. Cuthbert and C. Ariza, *ISMIR*, 11, 637 - 642 (2010).
- [14] M. D. Simon, I. and S. Basu, “Mysong: automatic accompaniment generation for vocal melodies,” in Proceedings of the SIGCHI conference on human factors in computing systems (2008) pp. 725 – 734.

- [15] E. D. Païement, J. F. and S. Bengio, “A probabilistic model for chord progressions,” in Proceedings of the 8th International Conference on Music Information Retrieval (2008).
- [16] K. Choi, G. Fazekas, K. Cho, and M. Sandler, arXiv preprint arXiv:1709.04396 (2017).
- [17] G. Hadjeres, F. Pachet, and F. Nielsen, “Deepbach: a steerable model for bach chorales generation,” in Proceedings of the 34th International Conference on Machine Learning-Volume 70 (JMLR. org, 2017) pp. 1362 – 1371.
- [18] H Dong, H. W. and Y. H. Yang, “Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment,” in 32nd AAAI Conference on Artificial Intelligence (2018).
- [19] Y. Wang and Y. H. Yang, arXiv preprint arXiv:2002.00212 (2020).
- [20] W. R. Brunner and S. Zhao, “Symbolic music genre transfer with cyclegan,” in 2018 IEEE 30th International Conference on Tools with Artificial Intelligence (ICTAI) (2018) pp. 786 – 793.
- [21] C. S. Y. Yang and Y. H. Yang, arXiv preprint arXiv:1703.10847 (2017).
- [22] G. Hadjeres, I. Simon, and E. Vincent, OpenAI Blog 1 (2019).
- [23] A. DuBreuil, *Hands-on music generation with magenta: Explore the role of deep learning in music generation and assisted music composition* (Packt Publishing Ltd, London, 2020).
- [24] “Midi files of Frédéric Chopin compositions,” retrieved from: <https://www.midiworld.com/chopin.htm>
- [25] J. P. Briot, G. Hadjeres, and F. Pachet, arXiv preprint arXiv: 1709. 01620 (2017).
- [26] D. Herremans and E. Chew, IEEE Trans. on Affe. Comp., 12 (2017).
- [27] H. Chu, R. Urtasun and S. Fidler, arXiv preprint arXiv: 1611.03477 (2016).
- [28] A. Roberts, J. Engel, C. Raffel, C. Hawthorne, and D. Eck,” arXiv preprint arXiv:1803.05428 (2018).
- [29] N. Jaques, S. Gu, R. E. Turner, and D. Eck, arXiv preprint arXiv: 1611.02796 (2017).
- [30] A. Vaswani, N. Shazeer, N. Parmar, et al. Adv. in neural infor. Proce. Sys. 30 (2017).
- [31] A. V. D. Oord, Y. Li, I. Babuschkin, et al., International Conference on Machine Learning (PMLR, 2018) pp. 5153 – 5162.
- [32] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, Communications of the ACM 63, 139 – 144 (2020).