

Emulating the Style of Johannes Brahms: Neural Network-Based Generation of Romantic Piano Music

James Doherty

Technological University Dublin
Central Quad, Grangegorman Lower, Dublin 7, D07ADY7
Dublin, Ireland
e-mail: j_doherty1992@hotmail.com

Brendan Tierney

Technological University Dublin
Central Quad, Grangegorman Lower, Dublin 7, D07ADY7
Dublin, Ireland
e-mail: brendan.tierney@tudublin.ie

Abstract – Neural network architectures currently are only able to employ music generation tasks to similar levels of human composers when the music is at a basic compositional standard. Research has shown that various neural network models struggle with the complex motifs and harmonic structures of Western Classical Music. This study aims to determine if various data preprocessing and augmentation techniques can train a neural network model to generate pieces of piano music to a similar level of musicality and emotion as Romantic Period composer Johannes Brahms. Quantitative experimentation involving Music Information Retrieval was conducted as well as a quantitative survey with respondents consisting of only professional musicians, composers and conductors. Analysis of the results from the quantitative experiment demonstrated that the Transformer models using Self-Attention, Relative Self-Attention and Local Windowed Attention generated statistically indistinguishable results to the original piano works of Brahms in terms of duration distribution, pitch class distribution, rhythmic variation and pulse clarity. Analysis of the quantitative survey discovered that participants struggled to distinguish between the pieces generated by Brahms and the models, with two generated pieces being at majority mistaken as one of Brahms' own works. The results indicate that various data preprocessing and augmentation methods do have an impact on model accuracy and for numerous musical characteristics, can be statistically indistinguishable to Brahms through quantitative experiment and survey. Further research is needed to identify other techniques to improve the musical characteristics of entropy and global energy in the generated pieces so that they are at a statistically comparable level to Brahms. The paper discovers that transformer models using varying attention mechanisms were able to generate longer sequences of music up to 2 minutes long which contained the composite motifs and harmonic structures of romantic period piano music.

Keywords–artificial Intelligence; music generation; neural network architecture; Brahms.

I. INTRODUCTION

The intention of this project is to generate piano music in the style of classical music composer Johannes Brahms by training a Recurrent Neural Network (RNN), a Long Short-Term Memory (LSTM) based RNN, various Transformer models with different attention mechanisms, and a Perceiver AR model [1]. These models will be trained with a pre-processed and augmented dataset containing MIDI files of Brahms' piano works. To deem the success of the project, the best musical pieces generated from the neural network models must show statistical indistinguishability from the Brahms corpus in various musical variables using Music Information Retrieval (MIR). Along with this, the pieces must also be mistaken by professional musicians, composers, and conductors as one of Brahms' own works through a quantitative survey. Although there have been many examples of AI models generating music in the style of particular composers, no models have been created to generate the work of Brahms. The lack of a detailed computational analysis of Brahms shows a gap in the study of romantic period composers, which Brahms was a key figure of [2]. According to studies taken, computer-generated music has traditionally only sounded human-like when short excerpts were created and struggled with the complex motifs and harmonic cadences of romantic period piano music. This could be down to them being poor at handling higher-level musical structures due to the models only learning how to play the next note according to the previous. Zheng et al. (2022) stated that their model was limited due to it learning how to play the next note according to the previous ones however it did not have any knowledge on the structure of music [3]. In their paper, Child et al. developed a sparse transformer and stated that it was able to extract complex patterns from sequences up to 30 times longer than possible previously [4]. Likewise, Hawthorne et al. developed a Perceiver AR model which had the ability to efficiently handle longer sequences with an improved memory efficiency [5]. After listening to the

examples from the papers, the generated pieces still consisted of basic harmonic and rhythmic structures and struggled with the complex motifs and harmonic cadences of romantic period piano music. The most common quantitative experiment was surveys. It appears that simple questions were asked to a small group of test subjects who were mainly computer science students and participants who had little to no music compositional knowledge. A more realistic survey would involve professional musicians, composers and conductors who specialise in the subject and would give a critical and educated opinion on the generated pieces compared to the composers that they have great knowledge about.

The importance of the research problem not only addresses AI's ability to generate music, but also highlights potential significance on how music could be composed in the future, particularly for those with no previous musical knowledge [6]. The research assumes that neural network models can already be trained to learn the general characteristics and patterns of various musical composers. To help with producing optimal results, the selected MIDI files for the dataset were exact replications of Brahms' piano music without errors and inconsistencies.

Based on the issues raised above, this study focused on understanding what configurations of the neural network models performed best when tasked with producing music in the style of Brahms. Therefore, the research question is:

To what extent can the accuracy of various Neural Network Models, trained with Long Short-Term Memory and numerous Attention mechanisms, be significantly improved by augmenting MIDI files containing the compositional works of Johannes Brahms with an augmentation pipeline to generate pieces of music that are mistaken by professional musicians, composers and conductors as one of Brahms' own works?

Due to conclusions from studies taken and the general state of knowledge at the time of beginning the project, the null hypothesis for this research project is;

H0: Neural network models cannot generate piano works to the same level of musicality and emotion as Brahms. Due to this, generated pieces will not be statistically indistinguishable through music information retrieval or mistaken as a work of Brahms by professional musicians, composers and conductors through a quantitative survey using Likert Scales.

Through the utilisation of an augmentation pipeline to expand the MIDI dataset containing the compositional piano works of Johannes Brahms, more musical variations could be created including transposition, rhythmical and note durations. In addition to this, various preprocessing techniques including track splitting, quantisation and normalisation could help make the MIDI files more readable for the models. This provides an alternate hypothesis;

H1: If an augmentation pipeline is utilised to expand a MIDI dataset of pre-processed files containing the piano works of Johannes Brahms, then various neural network

models trained with Long Short-Term Memory and numerous Attention mechanisms could generate pieces of music that is statistically indistinguishable from Brahms and could be mistaken as one of Brahms' own piano works by professional musicians, composers and conductors through a quantitative survey using Likert Scales and various Independent-Samples T-Tests and Hotelling's T^2 Tests being implemented to determine whether the p -value is > 0.05 to support the conclusion of statistical indistinguishability

This paper contains a total of four sections. Section II describes the experiment design, methodology and how the dataset was prepared and pre-processed. Along with this, the neural network models obtained for the project and MIR functions will be explained as well as ethical considerations. Section III analyses and evaluates the results of the quantitative experiment and survey to determine if the experiments provide evidence that the null hypothesis is incorrect. Section IV summarises what has been learnt and proposes recommendations and adjustments for future studies.

II. DESIGN AND METHODOLOGY

The research project was carried out through three stages. Firstly, data was collected, pre-processed, and augmented. The dataset contained 67 MIDI files of Johannes Brahms' piano works and was obtained for offline manipulation from *Classical Archives* and *MIDIworld*. MIDI (Musical Instrument Digital Interface) is seen as one of the most important tools for both musicians and producers as it represents a method to store and transmit musical data. They contain *Note-on* and *Note-off* messages to indicate the beginning and end of each note and delivers information about various musical characteristics in musical notes, including pitch, duration and dynamics.

An augmentation pipeline was employed to create variations in melodies, tempo, rhythm, and transpositions. This was followed by obtaining and training various neural network models with the augmented dataset. The models were analysed for training accuracy and loss to determine the best configurations. After training, the best generated examples from each model were analysed through various MIR functions using *MIDItoolbox* and *MIRToolbox* to determine the best performing models [7][8]. Finally, the best models generated pieces of music that were evaluated against the pieces of Brahms' repertoire for statistical equivalence. Along with this, the same generated pieces were used in a quantitative listener survey to test professional musicians, composers and conductors on whether they could differentiate between the generated pieces against the Brahms original. All quantitative experiments involved evaluation to test the research hypothesis through Independent-Samples T-Tests and Hotelling's T^2 Tests.

A. Preparation and Preprocessing of Dataset

The preparation and preprocessing of the dataset involved numerous steps to optimise the potential of training the neural network models. Track splitting involved dividing the MIDI tracks into smaller segments of 30 second clips to make the tracks more digestible for the models in training [9]. The conversion of multi-track MIDI files into a single monophonic track allowed for further simplicity in the files. Normalisation was performed to ensure that note velocities of the files were standard to provide consistent values for training and evaluation tasks. Finally, all the MIDI files were quantized to semiquavers to adjust the timings of notes and align them with the correct timing to ensure consistency in rhythm and therefore making it easier for the models to digest them [10]. Figure 1 displays an example of Brahms' *Variations on a Hungarian Song in D major, Op.21, No. 11* being quantised to semiquavers, also known as sixteenth notes.

Depending on the model being used, there were several packages utilised to extract information from MIDI files:

- **Pyfluidsynth:** Python package that enables audio playback of MIDI files in a Google Colab setting [11].
- **pretty_midi:** Package that contains various functions and classes for manipulating MIDI data so that it could be easily modified for information extraction. Figure 2 shows an example of `pretty_midi` extracting variables from each note [12].

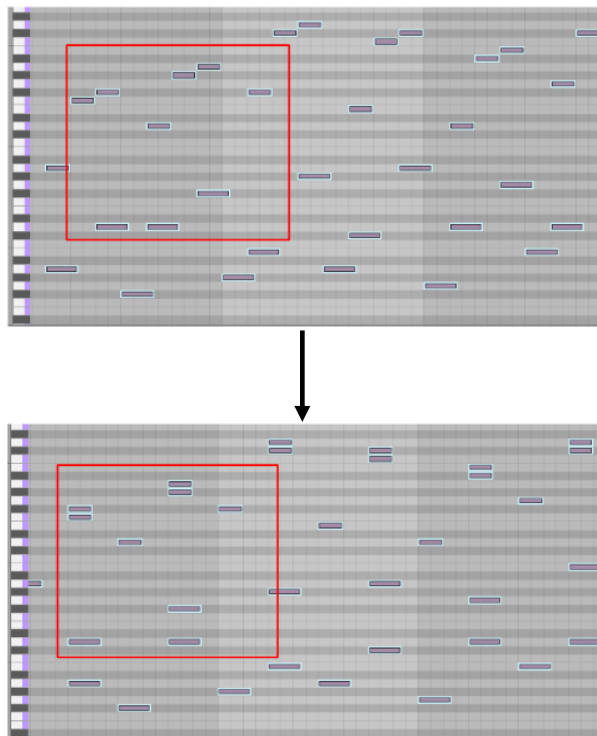


Figure 1. Example of Sixteenth note Quantisation

```
for i, note in enumerate(instrument.notes[:10]):
    note_name = pretty_midi.note_number_to_name(note.pitch)
    duration = note.end - note.start
    print(f'{i}: pitch={note.pitch}, note_name={note_name}, '
          f'duration={duration:.4f}')

0: pitch=79, note_name=G5, duration=0.2487
1: pitch=67, note_name=G4, duration=0.2708
2: pitch=67, note_name=G4, duration=0.1589
3: pitch=79, note_name=G5, duration=0.3060
4: pitch=67, note_name=G4, duration=0.2201
5: pitch=79, note_name=G5, duration=0.3346
6: pitch=64, note_name=E4, duration=0.0977
7: pitch=67, note_name=G4, duration=0.2943
8: pitch=79, note_name=G5, duration=0.3958
9: pitch=72, note_name=C5, duration=0.2005
```

Figure 2. Example of `pretty_midi` extracting variables from each note

- **NoteSequences:** Numerical representations of music notes including pitch number representation, along with note start and end time [13].

An augmentation pipeline was utilised to create variations in melodies, tempo, rhythm and transpositions. Several techniques were used to increase the dataset size. Time-stretching was applied to make each MIDI file 5% faster or slower. Another method was to transpose each of the MIDI files so that the pitches would be raised or lowered by a third [9]. Doing these increased the dataset by 500% with a total of 445 tracks. The augmented dataset was divided into training and testing sets which were defined at an 90%/10% split with repeated Montecarlo sampling (100 times). This same split was used for all the models to conduct a fair experiment. In comparison with prior studies, data preparation and preprocessing was influenced by previous research, which was collected, adapted and integrated into this paper.

B. Neural Network Models

Numerous Neural Network Architectures were obtained for offline manipulation and trained with the augmented Brahms MIDI dataset. A Recurrent Neural Network was acquired from TensorFlow [12]. A RNN was described by IBM as “a type of artificial neural network which used sequential data or time series data” [14]. Sequential data was utilised to predict the next output based off previous elements in the sequence. For this model, the training dataset was created through extraction of notes from the MIDI files with three values being taken from each note; pitch, step and duration. The model was trained on batches of note sequences with each example consisting of a sequence of notes as the input tokens and the subsequent notes served as the target output. This meant that the model was trained to predict the next note in the sequence containing 100 notes as the input, a common practice used in text classification. Similar to the training dataset, the model contained pitch, step and duration outputs. For the step and duration outputs, a loss function was constructed based on *mean squared error* equation to encourage the

model to output non-negative values. Equation (1) shows the formula where y_i represents the n^{th} token value and \hat{y}_i being the predicted value from the model. n represented the number of observations [15].

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n} \quad (1)$$

Evaluating the model showed that *pitch loss* was much greater than *duration loss* and *step loss*. This result indicated that predicting the next note had proven to be a more difficult task when compared to the other two. The model was trained with 22 epochs and testing recorded a training loss of 0.550 and accuracy of 0.628. RNNs suffer from vanishing gradient problem. With the neural network using the gradient decent algorithm to update the weight, the gradients therefore decreased in growth the further down the layers the network progressed.

A solution to this problem is the use of LSTM which utilises gating mechanisms to control the movement of information and gradients to allow for the network to learn and maintain information over longer sequences. [16] An LSTM-based RNN was obtained from Huang et al. [17]. The RNN utilised event sequence encoding with the following representing the event types:

- **NOTE_ON:** The start of a musical note
- **NOTE_OFF:** The end of a musical note
- **TIME_SHIFT:** Change of time (bpm)
- **VELOCITY:** Change of attack on notes

The model supported expressive timing with *time_shift* events allowing it to control the tempo at increments between 10ms and 1 second. *Velocity* events allowed for a dynamic to be set at the current *note_on* event which helped build more natural tension and timing. For the model to be able to read the Brahms MIDI dataset, they were converted to *NoteSequences*, a numerical representation of musical notes including pitch number representation, along with the note start and end times [13]. The model was trained for a total of 10002 checkpoints accumulating to 17 hours of training. Although LSTM provides a solution to vanishing gradient problem. They still suffer from gradient explosion which occurs when gradients become too large, causing a model to become unstable and unable to learn from training data. LSTMs struggle to memorise earlier information when very long sequences were being processed due to its limited context window. As the forget gate became more prominent, older information was replaced by new data.

Various Transformer models containing different attention mechanisms were obtained from Project Los Angeles [18]. Proposed by Google researchers Vaswani et al., transformer models do not rely on recurring processing of data. Instead, they operate on an attention mechanism [19]. Attention allows for neural networks to concentrate on particular parts of the input. Transformer models contain

two main components, an encoder, which produces a sequence of hidden states from an input, and a decoder, which takes the hidden states from the encoder and generates an output. The attention mechanism enables the decoder to access the encoder's hidden states meaning that the decoder can learn long-range dependencies in the input. This feature allows for the models to effectively handle larger amounts of data [20]. When compared to Simon & Oore's (2017) *Performance RNN* model, Huang et al. (2018) stated that while the LSTM-based model was able to generate credible music for very short samples, there was an apparent lack of long-term structure. However, their model *Music Transformer* could generate music with a consistent style that created multiple phrases out of a single motif. The attention mechanisms obtained for the transformer models tested were:

- **Self-Attention** – Processes inputs in the same sequence, enabling the model to capture dependencies within the input [21].
- **Relative Attention** – The relative position of tokens is considered based on similarity of other tokens in the sequence [22].
- **Local Windowed Attention** – Restricts attention to a fixed window of tokens, enabling the model to focus on nearby information [23].
- **Relative Self-Attention** – Combination of Self and Relative attention allowing the model to focus on relevant information based on the positional relationships of tokens [17].
- **Sparse Attention** – Attends to a subset of tokens instead of the entire sequence to improve efficiency while retaining information [24].

Traditional transformers with self-attention grow quadratically with the sequence length. This form of scaling could restrict the model's ability to work with longer sequences due to it requiring lots of memory and computational power. Sparse attention addresses this issue with the introduction of sparse factorisations of the attention matrix and therefore attending to a reduced subset of elements strategically chosen by the model. This enables the mechanism to still capture necessary information whilst the model performance is greatly increased [24]. Equation (2) displays the reformulation of the self-attention mechanism for sparse attention.

$$O(n^2) \longrightarrow O(n\sqrt{n}) \quad (2)$$

To train the transformer models, the MIDI files had to be processed with the TMIDIX MIDI processor provided by Tegrity code which converted the MIDI files to a score representation so that it could be understood by the model [25].

Pieces were generated in the style of OpenAI's *MuseNet* workflow where a premier seed MIDI file would be input and treated as the beginning of the generated piece. Doing so gave information to the model of the tempo, style and dynamics of the piece. A continuation generator was then implemented to append a chosen number of tokens and the multiple generations were batched to append the best results. This provided a degree of outcomes and the model was tested with variations on temperature, generation tokens and memory tokens. According to Huang & Yang (2020), amongst the great number of deep learning-based models that had been proposed for automatic music composition, the Transformer "stood out as the prominent approach for generating expressive classical piano performance with a coherent structure of up to one minute" [26].

A Perceiver Autoregressive (AR) model was also obtained from Project Los Angeles. Seen as an improvement to the Transformer model using latent array to distinguish the size of inputs and outputs, allowed for the model to effortlessly handle longer sequences with an improved memory efficiency [27]. Perceiver AR expands from the foundations of the original Perceiver Model with the addition of autoregressive capabilities, which was the ability to generate outputs sequentially based off previously generated outputs. This, along with the use of cross-attention to encode inputs into latent space, gave it the ability to train up to 100,000 elements. Perceiver AR contains 1024 latents and 24 self-attention layers. Due to each new note being predicted based off the preceding full sequence, the model can generate pieces of music containing high levels of melodic, harmonic and rhythmic consistency. The processing of MIDI files and generation of music was done so in a similar style to that of the transformer models.

C. Quantitative Experiment & Survey

A quantitative experiment was conducted to evaluate and compare the generated pieces of music from each of the models using Music Information Retrieval (MIR). MIR is the extraction of significant audio features from music data to be analysed through various algorithms and techniques. Raw audio signals were extracted and analysed from Waveform Audio Files (.wav) including pitch, harmony, rhythm and timbre that acquired an expressive description that is machine processible. MIR was used to test the neural network model's ability to replicate the musical characteristics and motifs of Brahms' piano works. Various MIR variables were utilised to gather various quantitative data from the models to test against not only each other, but the original works of Brahms. These variables included:

- **Entropy** – The measure of uncertainty or unpredictability
- **Duration Distribution** – The statistical analysis of note durations as well as silence
- **Pitch Class Distribution** – The evaluation of frequency of musical pitches

- **Mean Roughness** – The measure of dissonance or clashing sounds
- **Global Energy** – The total amount of energy held within a waveform. Calculated using the Root Mean Square (RMS) (3), each audio signal sample's amplitude x_i was squared with the mean of the square values being calculated to which the square root of that mean is taken. N represented the total number of samples in the signal [28].
- **Normalised Pairwise Variability Index (nPVI)** – The analysis of variability between successive durations to draw comparisons between rhythmic structures in music. Equation (4) shows the formula for nPVI, with d^k representing the duration of the interval k . m represented the number of intervals. The use of the formula standardised the difference between successive intervals in relation to their average length [29].
- **Pulsation Clarity** – The strength of the rhythmic pulse in a piece of music

$$RMS = \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2} \quad (3)$$

$$nPVI = \frac{100}{m-1} \times \sum_{k=1}^{m-1} \left| \frac{d_k - d_{k+1}}{\frac{d_k + d_{k+1}}{2}} \right| \quad (4)$$

A quantitative survey was also conducted and contained a total of 10 questions all containing the question "Rate the likelihood that this piece was composed by Johannes Brahms as opposed to being generated by AI" In random order, 5 pieces contained the works of Brahms and 5 were generated by the models. Answers consisted of a Likert Scale of 5 values ranging from *Definitely generated by AI* to *Definitely generated by Brahms*. Members of the Irish Defence Forces School of Music, RTÉ Concert Orchestra and National Symphony Orchestra were recruited for the survey to provide professional expertise in the subject. Participants were selected based on their extensive experience in classical music, including familiarity with Brahms' works, having performed his pieces in the past.

Using IBM's SPSS software, quantitative values from the MIR functions *MIRToolbox* and *MIDIToolbox* were tested to obtain p-values. Various Independent-Samples T-Tests and Hotelling's T^2 tests were implemented to evaluate musical variables to determine if there was statistical indistinguishability between Brahms' piano works and the generated pieces from the AI models.

D. Ethical Considerations

Several ethical considerations were adhered to in order to correctly conduct research including copyright issues and collection of data from survey participants. A total of 67 pieces of music were obtained for offline manipulation for the MIDI dataset. According to German Federal Law Gazette, copyright protection for musical and artistic works expired 70 years after the death of the creator. With Brahms passing away in 1897, his compositions therefore resided in the public domain. The use of MIDI files also prevented any issues with performers rights as recordings of Brahms' works were not being used. All the neural network models obtained for testing were open-source and ran under the Apache 2.0 License. MIR tasks were performed using the *MATLAB* functions *MIRToolbox* and *MIDIToolbox*. To utilise the tools, *MATLAB* had to be downloaded free of charge under the GNU General Public License. Ethical considerations were vital when collecting data from personnel for the quantitative survey. Participation for the survey was voluntary and those who chose to partake were informed of the purpose of the study. No personal information was required from participants and the confidentiality of participants was guaranteed from the designer of the survey. The results of the survey were not tampered with and therefore were accurate.

III. EVALUATION

Overfitting issues were observed during training, with the LSTM and Perceiver AR models initially replicating the training material excessively instead of generating unique motifs. To mitigate this, the dataset was further augmented by transposing musical phrases and altering rhythms through quantisation to introduce additional variations in tempo and phrasing. Experimenting with different temperature values for each model allowed for control over the balance between simplicity and randomisation. With lower temperatures producing more basic outputs that resembled the training data, while higher temperatures encouraged greater diversity but could lead to compositions with less structure. Table 1 shows the model performance metrics stating that the transformer model with relative self-attention scored the best training loss and accuracy with scores of 0.015 and 0.995 respectively. The recurrent neural network scored worst in terms of training accuracy with a score of 0.628.

Through quantitative experimentation and a survey, the neural network models utilised for the project were evaluated extensively in a numerical form and through professional human judgement in order to confirm or refute the research hypothesis that neural network models could generate pieces of music with statistically indistinguishable musical characteristics and emotion as Brahms. To conduct a fair experiment, each model had to generate a two-minute-long piece which contained 300 prime tokens (30 seconds) from the beginning of six of Brahms' piano works. These generated pieces were then compared with the first two

TABLE 1. TRAINING LOSS AND ACCURACY SCORES

Model	Training Loss	Training Accuracy
Recurrent Neural Network	0.550	0.628
RNN with Long Short-Term Memory	0.154	0.983
Transformer w/ Self-Attention	0.737	0.939
Transformer w/ Relative Attention	0.447	0.985
Transformer w/ Relative Self-Attention	0.015	0.995
Transformer w/ Local Windowed Attention	0.015	0.993
Sparse Transformer	0.612	0.803
Perceiver AR	0.825	0.774

minutes of the original Brahms piece to evaluate the evolution of the generated pieces and determine their ability to maintain the style and structure of Brahms. MIR evaluation concluded that the Transformer models with self-attention, local windowed attention and relative global attention performed best in generating music most similar to the Brahms original.

Several statistical tests were conducted on the best performing models to obtain p-values to test the research hypothesis. The variables Entropy, nPVI, Global Energy, Mean Roughness and Pulsation Clarity were tested with an Independent-Samples t-test and the variables Duration Distribution and Pitch Class Distribution being tested with a Hotelling's T^2 test. Based on guidance from Leard Statistics, the choice of test for each variable depended on whether the variable required a joint test. The independent-samples t-test was not suitable in those cases, whereas Hotelling's T^2 test was. Testing concluded that no statistically significant difference was found with most of the variables therefore supporting the alternative hypothesis that neural network models can produce music that is indistinguishable from Brahms. However, the variables Entropy and Global Energy were deemed to contain statistical significance within them stating that further work must be done to improve complexity, uncertainty and energy to a similar level to Brahms. Using *MIRToolbox*, the waveforms of the generated pieces were evaluated.

The entropy of distributions was evaluated to measure uncertainty or unpredictability in succeeding phrases. Table 2 shows how to models performed in terms of entropy from using 300 prime tokens of Brahms' *2 Rhapsodies No.1* to determine which model scored most like the original piece. Results showed that the Transformer model with Local Windowed Attention scored most similar with a score of 0.9214 compared to the original piece's 0.9115.

The transformer model with Local Windowed Attention was then further tested against the other Brahms original tracks to obtain more Entropy values. Table 3 shows the comparison of values with each piece.

With all this data, an independent-samples t-test was performed to determine if there was statistical significance between the Brahms and Transformer model's generated music. First, a Levene's Test for Equality of Variances was conducted to determine if equal variances should be

TABLE 2. ENTROPY VALUES BETWEEN BRAHMS AND AI MODELS

Model	Entropy	Difference
Brahms 2 Rhapsodies No. 1 – Original	0.9115	/
RNN	0.8802	0.0313
LSTM	0.8899	0.0216
Transformer w/ Self-Attention	0.9247	0.0132
Transformer w/ Relative Attention	0.9007	0.0108
Transformer w/ Local Windowed Attention	0.9214	0.0099
Transformer w/ Relative Global Attention	0.9241	0.0126
Sparse Transformer	0.9310	0.0195
Perceiver Transformer	0.9241	0.0126

TABLE 3. ENTROPY VALUES BETWEEN BRAHMS AND TRANSFORMER MODEL WITH LOCAL WINDOWED ATTENTION

Piece of Music	Type of output	Entropy
2 Rhapsodies No. 1	Brahms	0.9115
	Transformer w/ Local Windowed Attention	0.9214
2 Rhapsodies No. 2	Brahms	0.9140
	Transformer w/ Local Windowed Attention	0.9214
3 Intermezzi No.1	Brahms	0.9135
	Transformer w/ Local Windowed Attention	0.9168
4 Ballads No. 4	Brahms	0.9094
	Transformer w/ Local Windowed Attention	0.9167
7 Fantasias No. 7	Brahms	0.9135
	Transformer w/ Local Windowed Attention	0.9192
Hungarian Dances No. 2	Brahms	0.9163
	Transformer w/ Local Windowed Attention	0.9133

assumed for the t-test. With a scored p-value of 0.401 being much greater than 0.05, equal variances were assumed. Since both groups containing Brahms' original works and generated pieces from the transformer models were being evaluated, a two-sided p-value was observed. The independent-samples t-test stated a value of 0.012 indicating that there was statistical significance between the Brahms and AI generated music groups. The audio waveforms for the Brahms and AI generated pieces were analysed using *MIRToolbox* to discover the reason for the statistical significance. An onset curve depicted a greater variance in the detection of successive notes in the AI generated pieces therefore giving a higher entropy value.

Figure 3 shows the brightness curve of Brahms' 2 *Rhapsodies No. 1* alongside the generated piece from the transformer model with local windowed attention. The generated piece again displays a greater variance of frequencies, resulting in a higher entropy score. The entropy results show that the generated pieces provided too much variation and complexity in succeeding phrases when compared to the Brahms originals.

Statistical significance was also recorded in Global Energy between the Brahms pieces and the best scoring model in

this category, the Transformer model with Self-Attention. In this case, the Levene's Test for Equality of Variances calculated a p-value of 0.0002 meaning that equal variances should not be assumed. Therefore, the two-sided p-value of 0.035 was provided, indicating there was statistical significance between Brahms and AI.

The difference in global energy is depicted in Figure 4 where the temporal evolution curve reveals a much greater variance in the generated piece compared to Brahms'

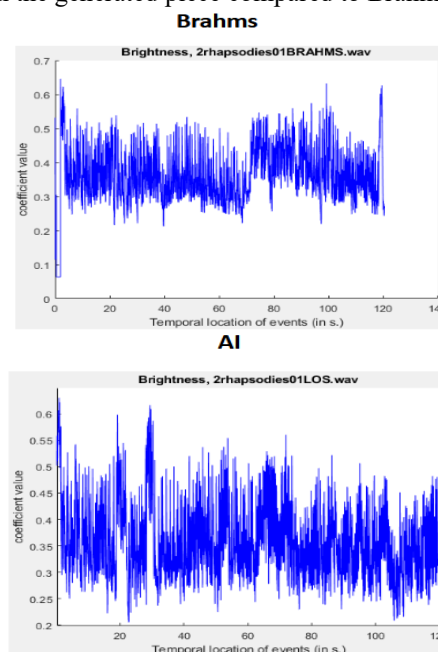


Figure 3. Brightness Curve between Brahms and AI

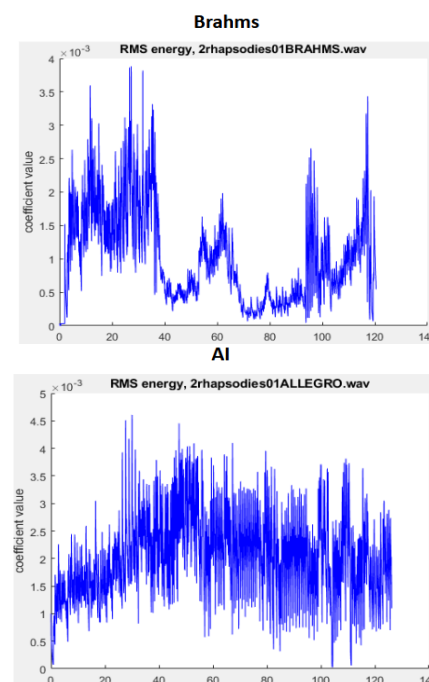


Figure 4. Global Energy Between Brahms and AI

original work. While Brahms' piece maintains a steady flow of increase and decrease of tension, the generated piece has much greater variations in timbre and harmonics throughout. The evaluation of these waveforms determined that the entropy and global energy values may have been much higher than Brahms' works due to the *MuseNet* inspired workflow the transformer models undertook causing the pieces to be generated in blocks and therefore sound unnatural.

Although there was statistical significance recognised for global energy, the independent-samples t-test did not recognise any difference in sensory dissonance. With Brahms being considered a more conservative composer who involved little dissonance, or clashing sounds, in his music, it was important to replicate this. Figure 5 shows that the generated piece composed the same methods as Brahms by answering to dissonance with a harmonic resolution. This can be seen with the sudden peaks in the coefficient value being resolved to lower roughness values.

Figure 6 displays the Pitch Class Distribution in the form of a box plot for Brahms' *3 Intermezzi No.1* and the generated piece from the transformer model with local windowed attention. With the piece being in the key of D# Major, which has an enharmonic equivalent of C minor, in order to stay within the key signature, the majority of the notes must be from the following triads:

- D# Major – D#, G, A#
- C Minor – C, D#, G

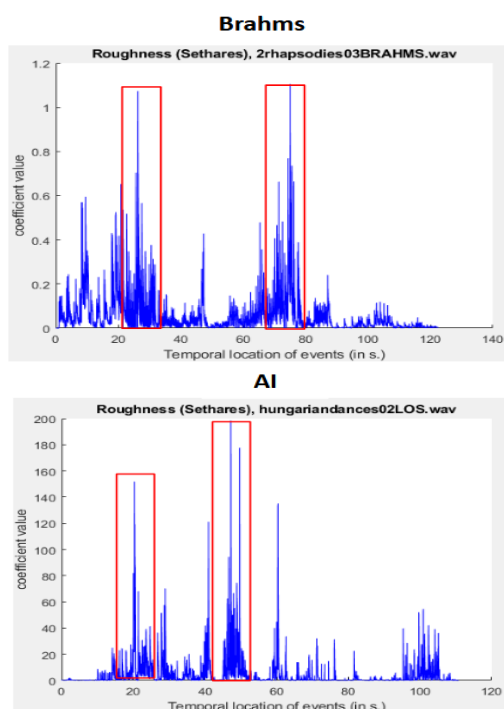


Figure 5. Sensory Dissonance between Brahms and AI

The pitch class distribution showed that both the original Brahms piece and AI generated kept within these guidelines, with particular emphasis being placed on the notes D# and G, as they feature in the triads of both D# major and C minor. Although the transformer model focused on the notes within the two triads to ensure consistency in key signature, it was apparent that the model was reluctant to incorporate accidentals to further add colour to the piece. The use of extended chords was a key factor of the romantic period in which Brahms lived in and it was an era that bridged the gap between classical and modern music. *MIRToolbox* was able to identify both the Brahms and AI pieces to be in the key of D# major. And with the function *mirmode*, it identified that the generated piece scored a higher probability of being in a major key than the Brahms original.

The chromagram at Figure 7 visualises the pitch distribution throughout the generated piece, showing that the D# major triad was used frequently throughout, meaning that the AI piece kept firmly within the major key. Although it was a positive that the generated piece was able to keep within the key signature for longer generated sequences, it does show an inability to evolve melodically into different harmonics.

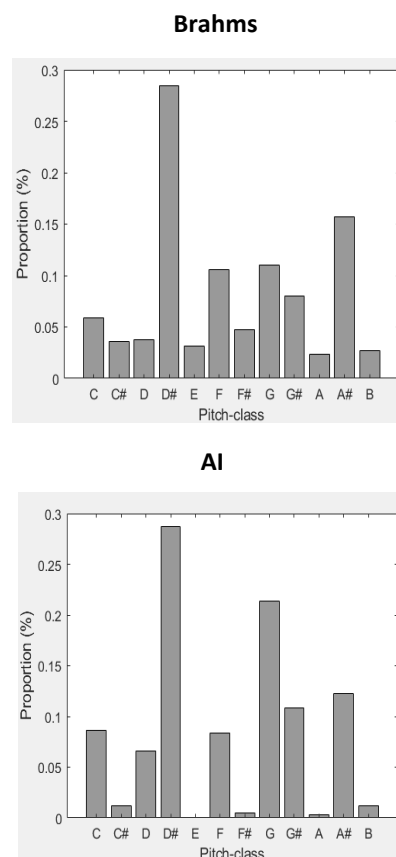


Figure 6. Pitch Class Distribution between Brahms and AI

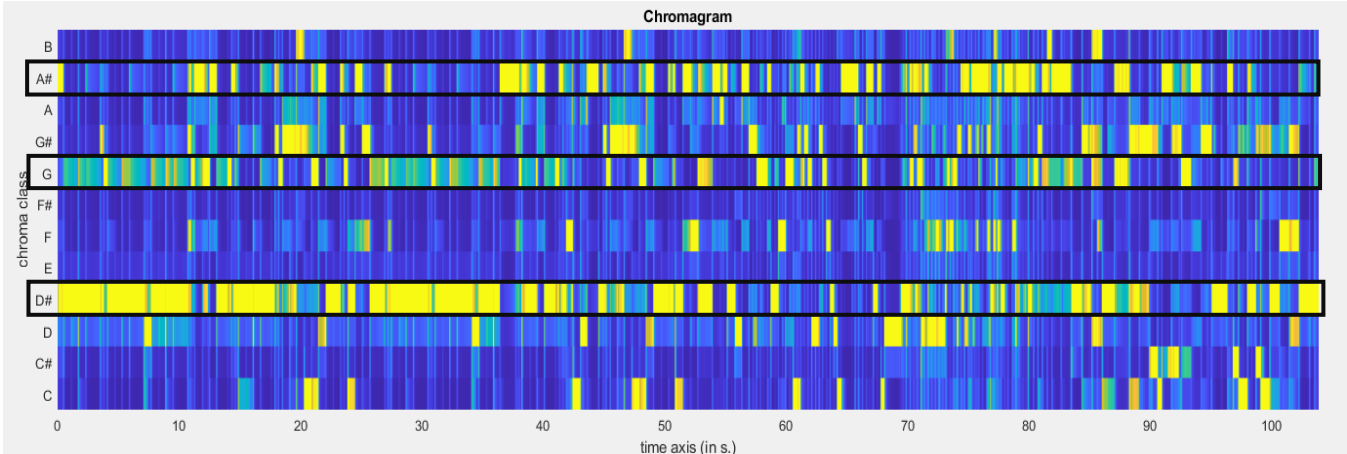


Figure 7. Chromagram showing the pitch class frequency throughout the generated piece by transformer model with local windowed attention

Figure 8 shows the Pulsation Clarity between Brahms and AI. The independent-samples t-test did not recognise any statistical difference between the two therefore the generated music from the transformer model with relative self-attention produced indistinguishable pulse clarity characteristics from Brahms. Pulse clarity did not necessarily mean that the piece must be perfectly in tempo. With the Romantic period being an expressive music era,

compositions often involved many tempo pushes and pulls, also known as *rubato*. Therefore, it was important for the AI model to generate the same pulsation characteristics as Brahms.

The figure shows how the transformer model generated similar pulsation and rhythmic characteristics to the original Brahms pieces with the higher values indicating a strong pulse followed by dips in coefficient value before regaining its original value. This could have been a sign of the music taking *rubato* into effect.

A quantitative survey was also conducted to obtain human evaluation on 30 second clips of generated pieces from the models against the original works of Brahms. A total of 56 participants featuring only professional musicians, composers and conductors, displayed difficulty in recognising distinction between the Brahms and generated pieces provided, with the majority incorrectly identifying two of the generated pieces as one of Brahms’ own works.

Table 4 shows the percentages for each of the questions alongside whether the track was that of Brahms or AI. The total score was also calculated with the number of responses per answer being multiplied dependant on how similar to Brahms it was scored with *Definitely generated by AI* scoring 1 and *Definitely generated by Brahms* scoring 5.

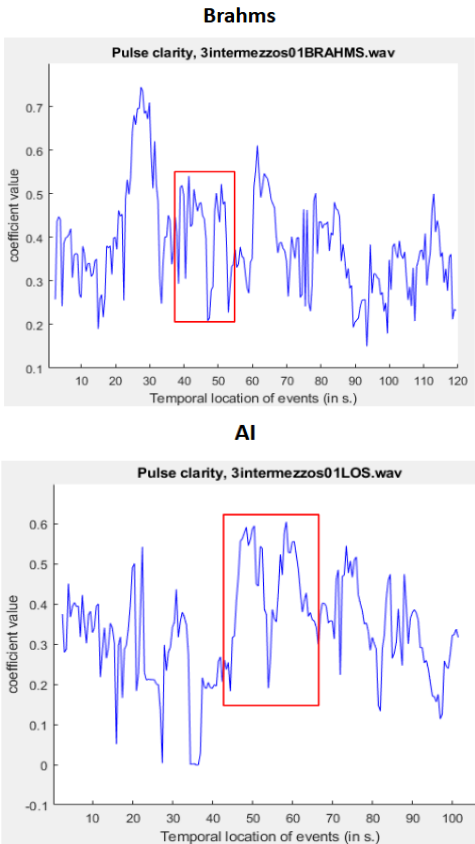


Figure 8. Pulse Clarity between Brahms and AI

TABLE 4. QUANTITATIVE SURVEY RESULTS

	AI	Probably AI	Unsure	Probably Brahms	Brahms	Total Score
Brahms	3.64%	5.45%	3.64%	36.36%	50.91%	237
Brahms	12.50%	25.00%	3.57%	37.50%	21.43%	185
Brahms	0%	23.21%	7.14%	39.29%	30.36%	211
Brahms	17.86%	44.64%	3.57%	28.57%	5.36%	145
Brahms	7.27%	7.27%	18.18%	41.82%	25.45%	207
AI	23.21%	32.14%	12.50%	23.21%	8.93%	147
AI	10.71%	19.64%	19.64%	28.57%	21.43%	185
AI	5.36%	41.07%	14.29%	32.14%	7.14%	165
AI	14.55%	34.55%	5.45%	32.73%	12.73%	165
AI	7.27%	30.91%	7.27%	36.36%	18.18%	183

This scoring design was intentional so that uncertainty was treated as a reward for the AI models as they still had not been identified as not Brahms by professionals.

An Independent-Samples T-Test was conducted which stated that there was no statistically significant difference between the Brahms and generated pieces therefore supporting the alternative hypothesis by that neural network models trained with an augmented and pre-processed dataset could generate music at the level of musicality and emotion as of Brahms so much that through a quantitative survey the difference could not be identified by professional musicians, composers and conductors.

IV. CONCLUSION AND FUTURE WORK

The work conducted in this paper differed to previous studies as it trained various neural network models with a dataset containing the piano works of Brahms, an important figure of the Romantic Period in Classical Music. By doing this, a gap in the research was addressed by analysing a very important composer in an era of classical music where harmonic and rhythmic structures began to diverge from the traditional aspects of Renaissance and Classical Period music while also bridging the gap between traditionalism and modernism. While the literature review stated a gap in the research that previous models struggled with complex motifs and harmonic cadences of romantic period piano music. This paper found that statistical testing in various musical categories stated there was not statistical significance between the Brahms and AI generated pieces. A quantitative survey containing participants who were educated in the subject mistook two of the neural network models generated pieces as Brahms' own works suggesting the model's ability to generate pieces of music with the complexity in rhythmic and harmonic characteristics of Brahms.

The results from the research carried out suggest that transformer models with self-attention, relative self-attention and local windowed attention were able to generate various characteristics to a statistically indistinguishable level from Brahms with a dataset of his piano works by utilising an augmentation pipeline and various preprocessing techniques. A couple of musical characteristics however proved to be statistically significant to Brahms, these being entropy and global energy. This concludes that while the transformer models are able to replicate a vast amount of Brahms' compositional traits, it still falls behind in reproducing the rhythmical and harmonical complexities, uncertainties and global energy levels of Brahms' works. The possibility of increasing the dataset to the orchestral, ensemble and choral works of Brahms could greatly increase the abilities of generated music from just solo piano works. This also could adhere to limitations regarding a small dataset and therefore improve accuracies in entropy and global energy from the generated pieces.

While participants noted difficulty in distinguishing between the Brahms and AI pieces, some commented that the use of MIDI files made all the music sound robotic and therefore made it even more challenging to differentiate between the two. Future work could focus on converting the generated pieces into musical notation and have a professional pianist perform them. This would enable an experiment to assess the generated music on an acoustic piano, the instrument it was originally intended for. The present study included generated music that used a short input sequence derived from the Brahms dataset. Future work will focus exclusively on cold-start generation to ensure all evaluated pieces represent novel compositional output.

REFERENCES

- [1] J. Doherty and B. Tierney, "The Generation of Piano Music in the Style of Johannes Brahms Using Neural Network Architectures", *In Proc of the International Conference on Creative Content Technologies*, no. 17, pp. 7-12, IARIA, 2025.
- [2] J. D. Fernández and F. Vico, "AI Methods in Algorithmic Composition: A Comprehensive Survey", *Journal of Artificial Intelligence Research*, no. 48, pp. 513-582, 2013. <https://doi.org/10.48550/arXiv.1402.0585>
- [3] K. Zheng, R. Meng, C. Zheng, X. Li, J. Sang, J. Cai, and J. Wang, "EmotionBox: a music-element-driven emotional music generation system using Recurrent Neural Network", *Frontiers in Psychology*, no. 13, pp. 1-14, 2021. <https://doi.org/10.3389/fpsyg.2022.841926>
- [4] R. Child, S. Gray, A. Radford and I. Sutskever, "Generating Long Sequences with Sparse Transformers", *PsyArXiv*, no. 1, pp. 1-10, 2019. <https://doi.org/10.48550/arXiv.1904.10509>
- [5] C. Hawthorne, A. Jaegle, C. Cangea, S. Borgeaud, C. Nash, M. Malinowski, S. Dieleman, O. Vinyals, M. Botvinick, I. Simon, H. Sheahan, N. Zeghidour, J. B. Alayrac, J. Carreira, and J. Engel, "General-purpose, long-context autoregressive modelling with Perceiver AR", *In Proc of the International Conference on Machine Learning*, no. 39, pp. 1-24, 2022. <https://doi.org/10.48550/arXiv.2202.07765>
- [6] E. Deruty, M. Grachten, S. Lattner, J. Nistal and C. Aouameur, "On the Development and Practice of AI technology for Contemporary Popular Music Production", *Transactions of the International Society for Music Information Retrieval*, no. 5, pp. 35-49, 2022. <https://doi.org/10.5334/tismir.100>
- [7] O. Lartillot, P. Toiviainen, P. Saari and T. Eerola, "MIRtoolbox", 2007. Accessed: Mar. 15, 2024. [Online]. Available: <https://www.jyu.fi/hytk/fi/laitokset/mutku/en/research/materials/mirtoolbox>
- [8] T. Eerola and P. Toiviainen, "MIDI Toolbox: MATLAB Tools for Music Research", 2004. Accessed: Apr. 22, 2024. [Online]. Available: <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=6e06906ca1ba0bf0ac8f2cb1a929f3be95ceadfa>
- [9] S. Oore, I. Simon, S. Dieleman, D. Eck and K. Simonyan, "This Time With Feeling: Learning Expressive Musical Performance", *Neural Computing and Applications*, no. 32, pp. 995-967, 2018. <https://doi.org/10.1007/s00521-018-3758-9>
- [10] J. A. Franklin, "Jazz Melody Generation from Recurrent Network Learning of Several Human Melodies", *In Proc of the International Florida Artificial Intelligence Research*

- Society Conference*, no. 18, pp. 57-62, 2005. <https://cdn.aaai.org/FLAIRS/2005/Flairs05-010.pdf>
- [11] FluidSynth, “pyFluidSynth”, 2014. Accessed: Apr. 14, 2024. [Online]. Available: <https://github.com/nwhitehead/pyfluidsynth>
- [12] TensorFlow, “Generate music with an RNN”, n.d. [Online]. Available: https://www.tensorflow.org/tutorials/audio/music_generation
- [13] C. Leja, “Making Music With Machine Learning: Introduction to Magenta Music.js,” 2020. Accessed May 27, 2024. [Online]. Available: <https://christopher-leja.medium.com/making-music-with-machine-learning-introduction-to-magenta-music-js-30b4c9dc6279>
- [14] IBM, “What are recurrent neural networks?”, n.d. Accessed: Mar. 5, 2024. [Online]. Available: <https://ibm.com/topics/recurrent-neural-networks>
- [15] J. Frost, “Mean Squared Error (MSE),” n.d. Accessed: May 11, 2024. [Online]. Available: <https://statisticsbyjim.com/regression/mean-squared-error-mse/>
- [16] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory”, *Neural Computation*, no. 9, pp. 1735-1780, 1997. <https://doi.org/10.1162/neco.1997.9.8.1735J>.
- [17] C. Z. A. Huang, A. Vaswani, J. Uszkoreit, N. Shazeer, I. Simon, C. Hawthorne, A. M. Dai, M. D. Hoffman, M. Dinculescu and D. Eck, “Music Transformer: Generating Music with Long-Term Structure”, *In Proc of the International Conference on Learning Representations*, no. 17, pp. 1-14, 2018. <https://doi.org/10.48550/arXiv.1809.04281>
- [18] A. Sigalov, “Project Los Angeles”, 2019. Accessed: Apr. 3, 2024. [Online]. Available: <https://github.com/asigalov61>
- [19] A. K. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser and I. Polosukhin, “Attention is All you Need”, *Advances in Neural Information Processing Systems*, no. 31, pp. 1-15, 2017. <https://doi.org/10.48550/arXiv.1706.03762>
- [20] S. Shankar, “Understanding Google’s “Attention is all you Need” Paper and its Groundbreaking Impact”, n.d. Accessed: Mar. 22, 2024. [Online]. Available: <https://alokshankar.medium.com/understanding-googles-attention-is-all-you-need-paper-and-its-groundbreaking-impact-c5237043540a>
- [21] A. K. Huang, P.A. Szerlip, M. E. Norton, T.A. Brindle, Z. Merritt and K.O. Stanley, “Visualizing music self-attention”, *In Proc of the Conference on Neural Information Processing Systems*, no. 32, pp. 1-5, 2018. <https://openreview.net/pdf?id=ryfxVNEajm>
- [22] P. Shaw, J. Uszkoreit and A. Vaswani, “Self-Attention with Relative Positional Representations,” *In Proc of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, no. 2, pp. 464-468, 2018. <https://doi.org/10.18653/v1/N18-2074>
- [23] Z. Liu, Y. Lin, H. Hu, Y. Wei, Z. Zhang, S. Lin and B. Guo, “Swim Transformer: Hierarchical Vision Transformer using Shifted Windows,” *In Proc of the IEEE/CVF International Conference on Computer Vision*, no.23, pp. 9992-10002, 2021. <https://doi.org/10.48550/arXiv.2103.14030>
- [24] OpenAI, “Generative modelling with sparse transformers,” 2019. Accessed: May 16, 2024. [Online]. Available: <https://openai.com/index/sparse-transformer>
- [25] asigalov61. “TMIDIX.py,” n.d. Accessed: May 18, 2024. [Online]. Available: <https://github.com/asigalov61/Giant-Music-Transformer/blob/main/TMIDIX.py>
- [26] S. Y. Huang and Y. H. Yang, “Pop Music Transformer: Beat-based Modeling and Generation of Expressive Pop Piano Compositions,” *In Proc of the ACM International Conference on Multimedia*, no. 28, pp. 1180-1188, 2020. <https://doi.org/10.1145/3394171.3413671>
- [27] A. Jaegle, F. Gimeno, A. Brock, A. Zisserman, O. Vinyals and J. Carreira, “Perceiver: General Perception with Iterative Attention,” *In Proc of the International Conference on Machine Learning*, no. 38, pp. 1-43, 2021. <https://doi.org/10.48550/arXiv.2103.03206>
- [28] M. Rayden, “What is RMS in audio world?,” 2024. Accessed: June 2, 2024. [Online]. Available: <https://majormixing.com/what-is-rms-in-audio-world/>
- [29] A. D. Patel and J. R. Daniele, “An empirical comparison of rhythm in language and music,” *Cognition*, no. 87, pp. 35-45, 2002. [https://doi.org/10.1016/S0010-0277\(02\)00187-7](https://doi.org/10.1016/S0010-0277(02)00187-7)