

# CompoNet: Toward Incorporating Human Perspective in Automatic Music Generation Using Deep Learning

SeyyedPooya HekmatiAthar

North Carolina Agricultural and Technical State University

Mohd Anwar (✉ [manwar@ncat.edu](mailto:manwar@ncat.edu))

North Carolina Agricultural and Technical State University

---

## Research Article

**Keywords:** CompoNet, Human, Deep Learning

**Posted Date:** September 7th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-812168/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# 1 CompoNet: Toward Incorporating Human 2 Perspective in Automatic Music Generation using 3 Deep Learning

4 SeyyedPooya HekmatiAthar<sup>1</sup> and Mohd Anwar<sup>1,\*</sup>

5 <sup>1</sup>North Carolina A&T State Univeristy, Department of Computer Science, Greensboro, 27411, United States

6 \*Corresponding Author: Mohd Anwar (manwar@ncat.edu)

## 7 ABSTRACT

8 The art nature of music makes it difficult, if not impossible, to extract solid rules from composed pieces and express them mathematically. This has led to the lack of utilization of music expert knowledge in the AI literature for automation of music composition. In this study, we employ *intervals*, which are the building blocks of music, to represent musical data closer to human composers' perspectives. Based on intervals, we developed and trained OrchNet which translates musical data into and from numerical vector representation. Another model called CompoNet is developed and trained to generate music. Using *intervals* and a novel monitor-and-inject mechanism, we address the two main limitations of the literature: lack of orchestration and lack of long-term memory. The music generated by CompoNet is evaluated by Turing Test: whether human judges can tell the difference between the music pieces composed by humans versus. generated by our system. The Turing test results were compared using Mann-Whitney U Test, and there was no statistically significant difference between human-composed music versus what our system has generated.

9  
10 Automatic Music Generation (AMG) started in 1957 by Ljaren Hiller and Leonard Isaacson. They employed Markov chains  
11 and rule-based logic and developed software capable of composing a string quartet<sup>1</sup>. In early 1960s, Iannis Xenakis generated  
12 numerical patterns using computers, which he later transcribed into sheet music<sup>2</sup>. After a period of silence in the literature which  
13 was probably due to AI winter, in early 2000s, Blackwell created an improvisational musical system by applying multi-swarms<sup>3</sup>.  
14 Concurrently, Eck et al. employed Deep learning for the first time in AMG. They used two Long Short-Term Memory (LSTM)  
15 models to create music<sup>4</sup>. Schulze and Merwe proposed a model which can handle two-instruments harmony. This model,  
16 which was based on Hidden Markov Model generated pieces for piano with melodies on the right hand and simple chord  
17 progressions on the left hand<sup>5</sup>. Kosta et al. achieved chord sequence generation based on feeding observed examples into an  
18 unsupervised multi-stage framework<sup>6</sup>. Colombo et al. composed melody using Recurrent Neural Network (RNN) architectures<sup>7</sup>.  
19 Similar to Colombo's work, Sturm et al. employed RNN to composite music in digital format<sup>8</sup>. Jaques et al. combined RNN  
20 and reinforcement learning to develop music generation models<sup>9</sup>. Dieleman et al. proposed multiple deep architectures for  
21 generating music in raw audio format<sup>10,11</sup>. Huang et al. applied deep sequential generative models with self-attention to  
22 generate structured compositions in synthesizing keyboard compositions<sup>12</sup>. Hadjeres et al. handled monophonic AMG based  
23 on Variational Auto-Encoder (VAE)<sup>13</sup>. Inspired by the exciting progress in image and video related tasks, some researchers  
24 have employed Generative Adversarial Networks (GAN)<sup>14</sup> to extract temporal features and applied them to AMG. Yang et  
25 al. proposed midinet which is a convolution GAN<sup>15</sup> with specified conditions in 1D and 2D during the generation process.  
26 Midinet effectively solves the problem of lack of context association in AMG. Considering the temporal structure of music,  
27 C-RNN-GAN and SeqGAN combined RNN and GAN to generate music<sup>16,17</sup>. SeqGan combined reinforcement learning with  
28 GAN for music generation<sup>17,18</sup>. SeqGAN tried to extract temporal features by RNN; however, SeqGAN is not suitable for  
29 multi-instrument symbolic AMG in pianoroll format because pianoroll for multi-instruments symbolic music is represented as  
30 high-dimensional data, and SeqGAN is difficult to train with high-dimensional data. Dong et al. proposed MusGan which is a  
31 multi-track symbolic AMG model based on a hybrid GAN architecture<sup>19,20</sup>. Table 1 summarizes the most relevant works in the  
32 literature.

## 33 The Gaps in the Existing Literature

34 The existing literature leaves two major challenges/limitations unanswered: lack of orchestration and lack of long-term memory.

**Table 1.** Summary of most relevant works in the literature

Reference	Model Name	Dataset	Feature Engineering
Guan <i>et al.</i> <sup>21</sup>	DMB-GAN	Lakh MIDI Dataset <sup>22</sup>	Convolution - Self Attention
Hadjeres <i>et al.</i> <sup>13</sup>	GLSR-VAE	music21 <sup>9</sup>	RNN
Eck & Schmidhuber <sup>4</sup>	-	<sup>23</sup>	LSTM
Strum <i>et al.</i> <sup>8</sup>	char-rnn	<sup>24</sup>	LSTM
Strum <i>et al.</i> <sup>8</sup>	folk-rnn	<sup>24</sup>	LSTM
Jaques <i>et al.</i> <sup>9</sup>	Note RNN	Not given	RNN
Yu <i>et al.</i> <sup>17</sup>	SeqGAN	Nottingham Dataset <sup>25</sup>	Convolution
Yang <i>et al.</i> <sup>15</sup>	MidiNet	TheoryTab <sup>26</sup>	Convolution
Mogren <sup>16</sup>	C-RNN-GAN	Not Given	RNN
Dong <i>et al.</i> <sup>19,20</sup>	MuseGAN	Lakh MIDI Dataset <sup>22</sup>	Convolution

**Lack of Orchestration**

Numbers in music theory should be interpreted differently than in the algebra. This leads to extra challenges when it comes to musical data representation. In other words, since numerical patterns in music follow different rules than patterns in other domains, existing pattern recognition techniques fail to correctly and efficiently extract the patterns in a given piece of music. Since common practices in pattern recognition need to put extra effort to compensate for music-specific differences, the number of instruments and ensembles needs to be lowered.

**Lack of Long-Term Memory**

Existing methodologies can be categorized into generative and memory-based models. While memory-based models, unlike generative models, do have memory capabilities, they inherit LSTM's most famous limitation: lack of long term memory. It has been observed that in sequence generation based on LSTM, an object, name, or idea gets introduced in the beginning, but it vanishes as the sequence gets longer.

**Methods**

In this paper, we aim to address the aforementioned challenges/limitations by employing the embedding mechanism based on music theory<sup>27</sup>. The embedded data is then given to a Stacked AutoEncoder called Orchestration Network (OrchNet). OrchNet specializes in translating musical data to and from vector representation. In the latent space of OrchNet, Composer Network (CompoNet) handles generation while ensuring long-term dependencies are maintained using a novel monitor-and-inject mechanism.

**Pre-Processing Data****Dataset**

Classical pieces are chosen as a use case for this paper. A MIDI dataset of 3,920 pieces from 144 classical composers is available on <https://www.kaggle.com/blanderbuss/midi-classic-music/28>. There is no documentation on how this dataset is gathered; therefore, it needs to be validated and cleaned.

**Data Validation**

The first step was to check file formats. There were some compressed files (.zip) in the dataset. The content of these files are extracted and checked. Then, a number of files (n=100) were randomly selected and verified by a human music analyst. This was done by comparing the content of the files with the available recordings of the same piece. In addition to these files, some famous pieces were selected and verified. The list of these pieces is provided in Table 2.

**Data Cleaning**

Data cleaning was done using Pypianoroll<sup>29</sup> and pretty\_midi<sup>30</sup> packages. The following criteria were followed for excluding a file:

1. if either Pypianoroll or pretty\_midi cannot handle the file and throw an error,
2. if the file has more than one time signature change event,
3. if the file contains less than 10 bars (the number of bars was calculated as:  $\frac{\text{active length}/\text{beat resolution}/\text{time signature}}{\text{numerator}}$ ),

**Table 2.** List of verified pieces

<b>Composer</b>	<b>Piece Name</b>
J.S. Bach	Ave Maria
	Bwv 565 Toccata and Fugue In Dm
	Bwv 772 Invention n01
	Bwv 1041 Violin Concerto n1 3mov
	Unfinished Fugue
Beethoven	Bwv 1053 Harpsichord Concerto n2 1 mov
	Fur Elise
	String Quartet no 1 op 18 1 mov
	Egmont Overture
	Sonata No.14 Op 27 Moonlight Sonata
Brahms	Symphony No. 9
	Two rondos for piano
Chopin	Waltz No. 1
	Hungarian Dance No.5
Fauré	Etude No. 1
	Nocturne Opus 62 No.1
Haydn	Spleen
	Romance sans paroles
Pachelbel	Symphony No. 94 "Surprise"
	Symphony No. 96 "Miracle"
Liszt	werde munter mein gemüte
	Canon Fantasy
Mozart	Hungarian Rhapsody No. 2
	Dante Sonata
	A piece for Piano, K. 176
	Piano Concerto No. 18, K. 456, "Paradis"
	Symphony No. 25 K. 183
Paganini	Piano Concerto No. 26, K. 537 "Coronation"
	Sonata for Piano Duet in C major, K.19d
Rachmaninoff	Symphony No. 40 K. 550
	Witches Dance
R. Korsakov	Moto Perpetuo
	Piano Concerto No.2 Op.18
Schubert	Prelude in A minor Op. 32 No. 8
	Flight Of the Bumblebee
J. Strauss	Aria of Sadko
	Erlkönig . D328 Op. 1
Wagner	Symphony D. 944 No. 9 "Great"
	The Blue Danube Op.314
Tchaikovsky	Ride of the Valkyries
	Waltz of the Flowers
Vivaldi	Swan Lake
	The Four Seasons
	Violin Concerto Op. 9 No. 2 RV. 345

70 4. if the numerator of the first time signature event is 1.

71 As a result, 653 files were excluded due to the first criterion while the other criteria reduced the number of files from 3,267 to  
72 2,416.

### 73 **Data Embedding**

74 The Music Embedding package<sup>27</sup> was used for data embedding. This package accepts inputs in the pianoroll format; therefore,  
75 Pypianoroll package is used to convert MIDI files into pianorolls and then input them to the embedder of the Music Embedding  
76 package. The embedder converts pianoroll into sequences of intervals. In music theory, an interval is defined as the difference  
77 in pitch between two notes. The results of this embedding is samples with five dimensions:

- 78 • Order: the ordinal number in the interval,
- 79 • Type: the type or mode of the interval,
- 80 • Direction: is interval ascending or descending,
- 81 • Octave: indicates the octave difference of the two notes of the interval,
- 82 • RLE: the value of RLE encoding, i.e. note duration.

### 83 **OrchNet**

#### 84 **Input and Output**

85 The proposed OrchNet has music theory built into it and translates music data to and from vector representation. OrchNet,  
86 which is a Stacked AutoEncoder, receives data embedded with Music Embedding package. *Order*, *type*, and *direction* are  
87 categorical while *octave* and *RLE* are numerical; therefore, it is reasonable to have one-hot encoding for *Order*, *type*, and  
88 *direction* and integer encoding for *octave* and *RLE*. However, since OrchNet is an AutoEncoder, its input and output must be  
89 identical; and having different encoding schemes in the output makes the model extremely difficult to train as each encoding  
90 scheme requires a specific loss function and combining different loss functions is not always possible. For this reason, *octave*  
91 and *RLE* are also encoded with one-hot encoding. Hence, each time step of the data is encoded as:

- 92 • Order: a one-hot vector of 7,
- 93 • Type: a one-hot vector of 5,
- 94 • Direction: a one-hot vector of 2,
- 95 • Octave: a one-hot vector of 4,
- 96 • RLE: a one-hot vector of 72.

97 This encoding, which requires 90 bits, allows encoding all seven ordinal possibilities of an interval, all five type possibilities of  
98 an interval, both ascending and descending directions, and 72 different possibilities for note duration; while allowing to have  
99 identical loss functions, hence an easier and more stable training routine.

100 The above-mentioned 90-bit encoding represents a single interval. Yet, music is the sequence of intervals. Per music theory,  
101 it is reasonable to break a piece into bars. However, each bar can have a different time signature and different number of notes  
102 and intervals. To mitigate this variable-length issue, padding and truncating mechanisms can be used. Therefore, each data  
103 sample for OrchNet is a 64 by 90 tensor; if a data sample has more than 64 time steps, it is truncated to 64; while the vast  
104 majority of data samples, which have less than 64 time steps, are padded with zeros to make them 64 by 90. The decision for  
105 having 64 time steps in each data sample is made because (1) the chosen number needs to be a power of two due to (GPU)  
106 hardware considerations, (2) 32 and smaller numbers are too aggressive and they truncate lots of data, while 128 is too much  
107 and pads lots of zeros which in turn increases sparsity in the data. 64 is the sweet spot, it truncates only 2% of samples while  
108 keeps the number of padded zeros manageable.

#### 109 **Model Architecture**

110 Since the data has temporal dependencies, OrchNet is built based on LSTM blocks. As explained earlier, the input is a 64 by 90  
111 tensor in which the sequence is padded with zeros to make its length consistent with other sequences. Therefore, OrchNet starts  
112 with a masking layer which return a mask of *True* and *False* values which indicates which time steps are real data and which of  
113 them are the padded zeros. An LSTM layer consisting 64 units follows the masking layer. This layer has a dropout value of 0.2  
114 to prevent over-fitting. Finally, a fully connected layer with 128 neurons is added to generate a vector of 128 as the latent space  
115 representation of the input.

116 The decoder starts with *repeat vector* operation to convert the 1D latent space into 2D which allows the following 64-unit  
117 LSTM layer reproduce the input. A fully connected layer with 90 neurons prepares the data to flow into the five final fully  
118 connected branches, each of which representing a dimension in the data. All branches have categorical cross-entropy loss  
119 function.

120 A music piece can be embedded in three different ways: *melodic*, *harmonic*, and *with-respect-to-bar*. After applying  
121 data validation and cleaning procedures, *melodic*, *harmonic*, and *with-respect-to-bar* embeddings were applied on the dataset.  
122 Then, the aforementioned model, which is illustrated in Fig. 1, was trained as the building block of OrchNet, irrespective of  
123 embedding type, i.e., all three embeddings were used equally as the training dataset. Then, three copies of the trained model  
124 were obtained and each of them were re-trained separately on a specific type of embedding.

125 The trained models can translate a single track (instrument) into vector representation. However, this is not enough for an  
126 orchestral piece. To address this, another AutoEncoder was trained to combine the vector representations of different tracks. As  
127 shown in Fig. 2, this fully connected model accepts up to 16 tracks. The number of tracks in the model is arbitrarily decided to  
128 be 16 as this number (1) satisfies hardware acceleration constraints, (2) does not lead to hard-to-handle sparsity in the data, and  
129 (3) covers 95% of the dataset without losing any data.

130 The building block of OrchNet has a latent space of 128. Combining the latent spaces of 16 copies of this model, each of  
131 which having three different embedding types, results in a vector of  $3 \times 16 \times 128 = 6144$ . Passing through two fully connected  
132 layers with 2048 and 1024 neurons respectively, OrchNet encodes each bar of an up-to-16 tracks piece of music into a vector of  
133 1024.

## 134 **CompoNet**

### 135 **LSTM**

136 CompoNet works in the latent space of OrchNet and generates sequences which can be decoded with OrchNet and presented as  
137 symbolic music. The core component of this sequence generation is an LSTM model as shown in Fig. 4. The model starts with  
138 two layers of LSTM units arranged in the opposition direction of each other (bidirectional), each consisting of 128 units. The  
139 output of these two layers is then given to a fully connected layer of 1024 neurons to generate the output. A dropout of 0.2 is  
140 used as a regularization mechanism.

141 This model works based on a sliding window mechanism which is illustrated in Fig. 3. In this mechanism, the model  
142 receives 16 samples as the input and generates the next one. The generated sample is added to the sequence and the window  
143 moves one step forward to include the newly generated sample and exclude the oldest sample. The limitations of this solution  
144 are:

- 145 • it is not self-sufficient, the first 16 samples are needed from external sources,
- 146 • it suffers from the long-term memory problem which was mentioned earlier.

### 147 **GAN**

148 The model shown in Fig. 4 requires the first 16 samples to be given from another source. While this source can be an existing  
149 piece of music, it is also possible to employ generative models to generate the required 16 samples. As illustrated in Fig. 5, a  
150 model based on GAN architecture can be developed and trained out of OrchNet Encoder, a Generator, and a Discriminator. The  
151 generator receives a vector of 2048 which is filled with random values. Using a fully connected layer of 2048 neurons and two  
152 layers of 1024 LSTM units (bidirectional), it converts the random input into a 16 by 1024 vector which matches the size of 16  
153 samples required by the LSTM model.

154 The input to the discriminator comes from either the generator or from OrchNet Encoder encoding 16 consecutive bars.  
155 The discriminator performs a binary classification task to determine the source of a given input. The training routine starts  
156 with training the discriminator. At this stage, since generator is producing random output, discriminator easily learns how to  
157 differentiate between inputs from OrchNet Encoder and the generator. The value of loss function is then back-propagated  
158 through discriminator to calculate the loss of the generator. This leads to training the generator to produce more realistic outputs.  
159 The cycle repeats and discriminator becomes better at differentiating the source and generator becomes better at producing  
160 outputs that it becomes hard for the discriminator to distinguish. At the end, the generator is trained enough to generate the  
161 required 16 samples similar to the existing samples in the dataset from a given random vector.

### 162 **The Monitor-and-Inject Mechanism**

163 The proposed architecture of CompoNet has two major issues:

- 164 • being based on LSTM, it suffers from lack of long-term memory,
- 165 • it follows a fixed sequential approach which makes it unsuitable to generate different styles and forms.

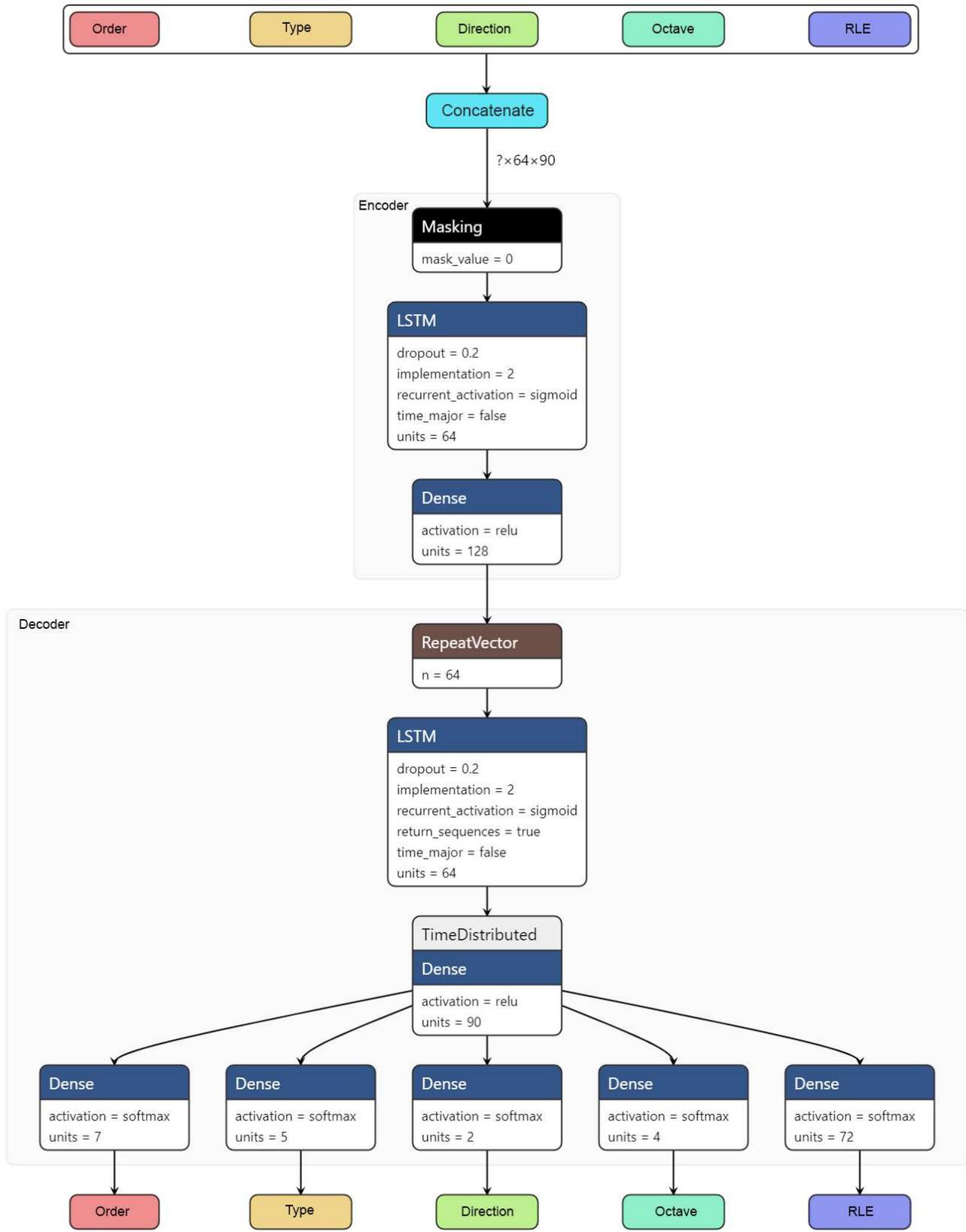


Figure 1. The building block of OrchNet

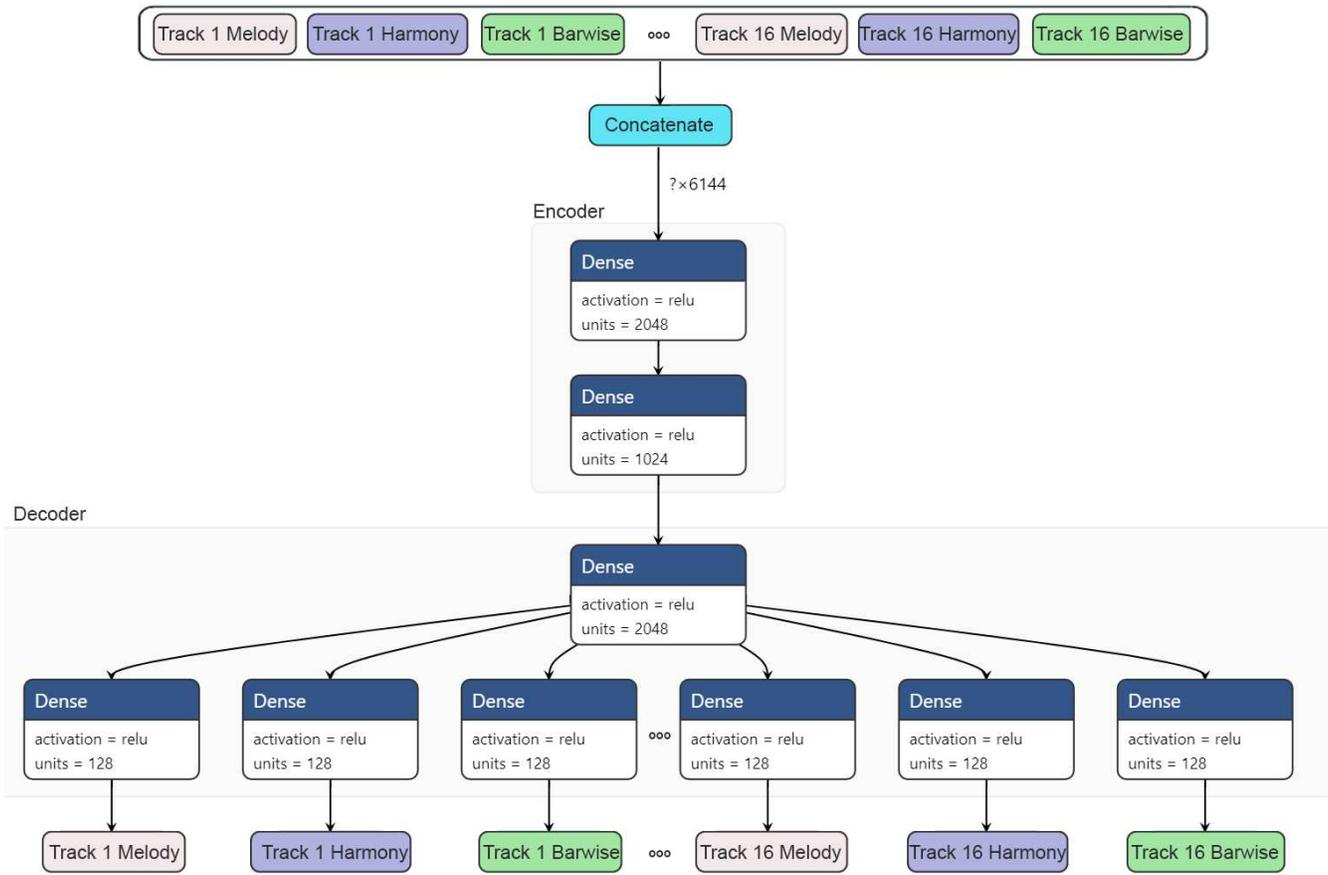
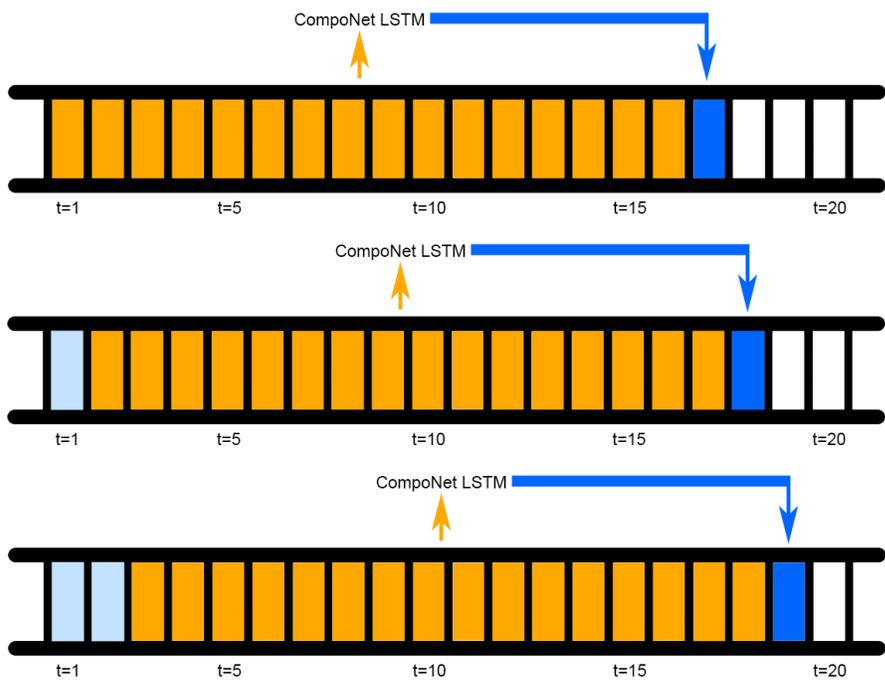


Figure 2. The full stack of OrchNet



window.png

Figure 3. The sliding window mechanism in CompoNet

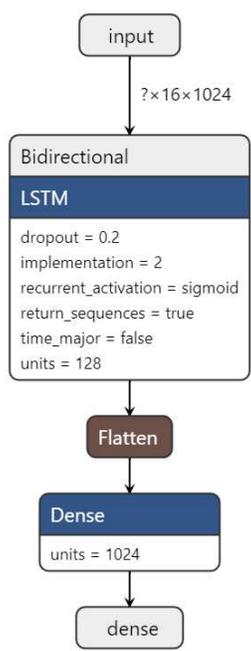


Figure 4. The LSTM model in CompoNet

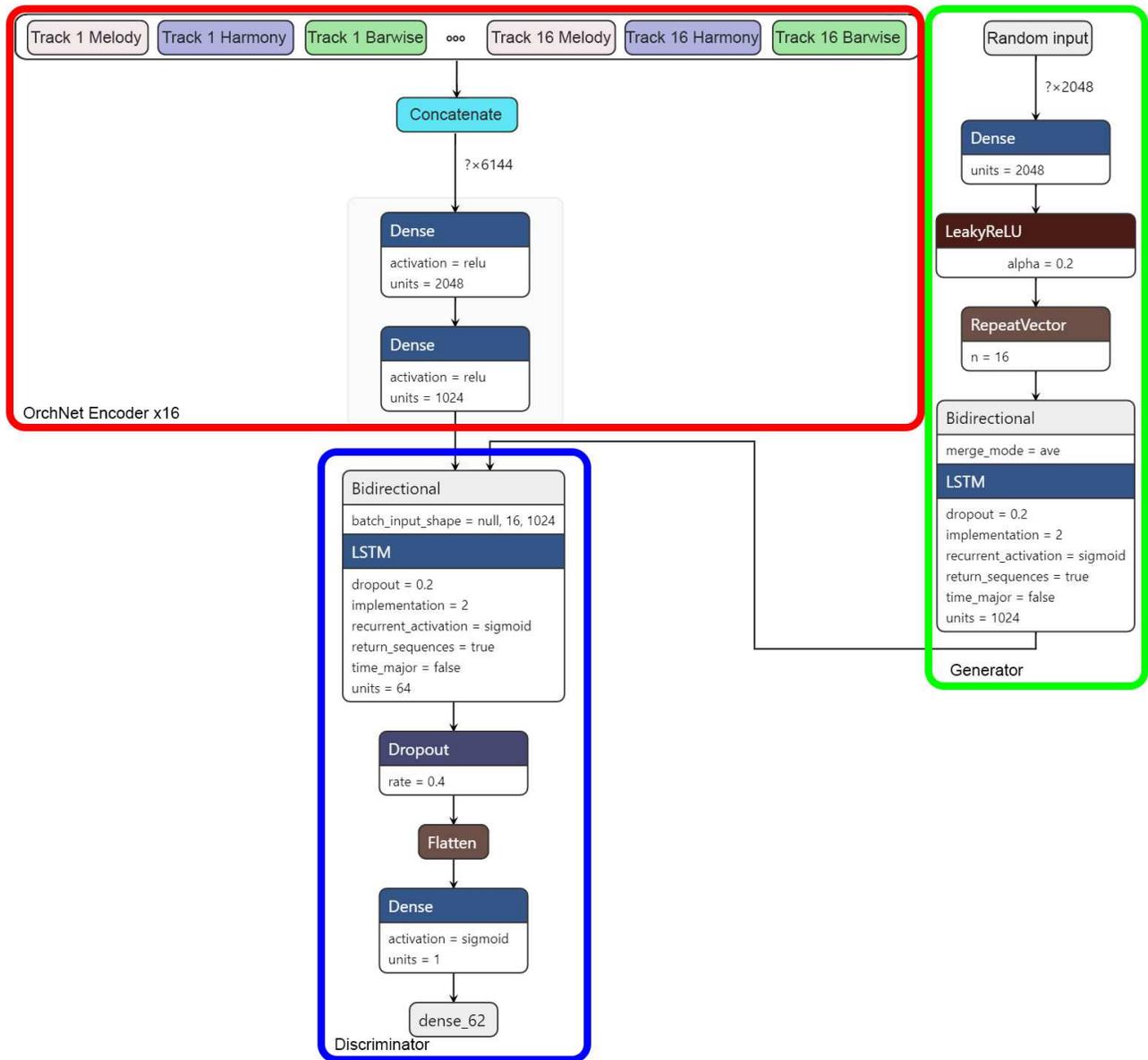


Figure 5. The GAN model in CompoNet

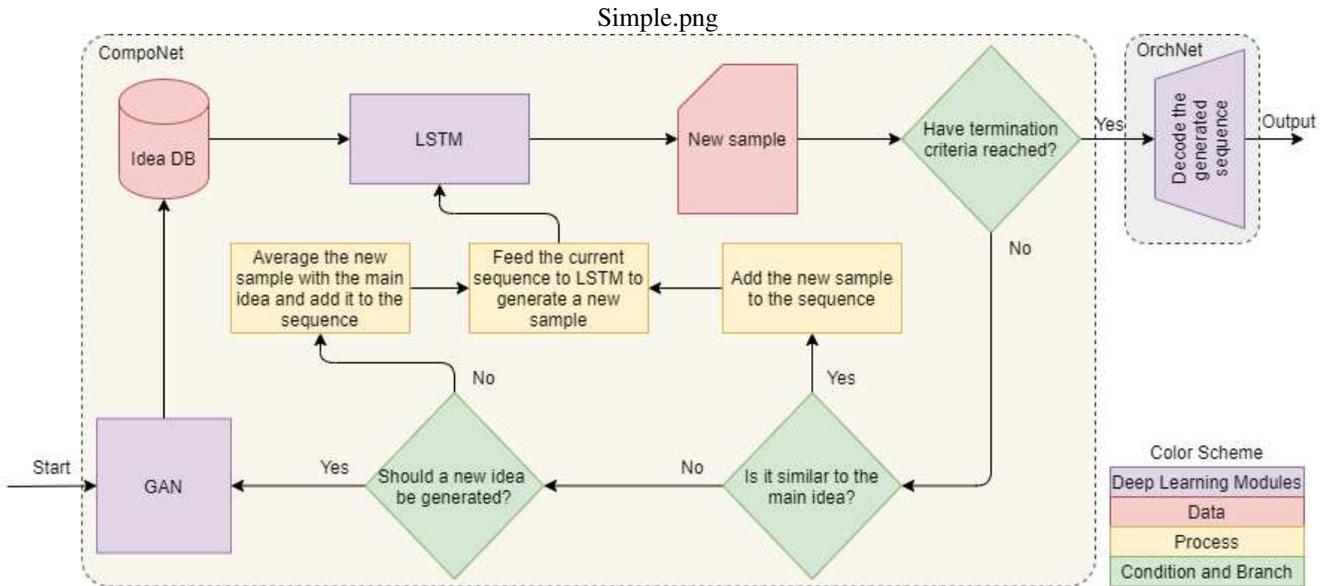


Figure 6. An example flowchart of the monitor-and-inject mechanism for maintaining long-term consistency and adopting different forms and styles

166 Both of these issue can be addressed by monitoring the output of CompoNet and altering it on the fly as needed. In other words,  
 167 since CompoNet has a subject/idea generation engine (the GAN model) and a model which maintains the short-term temporal  
 168 dependencies (the LSTM model), these two components can be combined in virtually infinite number of ways to generate  
 169 various styles and form such as *Sonata*, *Fugue*, *Nocturne* and so on, as well as new creative ways.

170 Fig. 6 illustrates an example of how the two components of CompoNet can be combined. While this example does not  
 171 implement any predefined form or style of music and it has some level of randomness built into it, it ensures long-term  
 172 dependencies are maintained. The process of generating a new piece of music starts by the GAN model generating the initial  
 173 seed. The seed is stored in the idea database, then it is given to the LSTM model to generate a new sample. The monitoring  
 174 process begins by *virtually* adding the new sample to the existing sequence. Then, the similarity of the newest 16 elements of the  
 175 generated sequence is compared to the initial seed from the GAN. While the impact and application of using different metrics  
 176 for measuring similarity is a topic for further investigation, here, Mean-Squared Error (MSE) is used. If the similarity meets a  
 177 certain threshold value, then the new sample will be added to the sequence and the cycle repeats. Otherwise, a decision has to  
 178 be made to either generate a new seed or reuse a previous one. While different logics can be implemented to determine how  
 179 previously stored ideas can be reused, here the simple mechanism of roulette wheel is used to give both options of generating  
 180 new seed and reusing the last seed equal chances. This is important because if the balance between the two options is not  
 181 maintained, then the output either loses the long-term consistency or generates a repetitive sequence. Either way, if an idea  
 182 needs to be injected to the sequence, it can be either replaced with the existing sequence or be averaged with it. Here, the  
 183 averaging option is chosen because it has less impact on short-term dependencies.

## 184 Implementation

185 Python is the main language used to implement this work. Music Embedding, Pypianoroll, and pretty\_midi packages were used  
 186 for data-handling tasks while Keras and TensorFlow were used to define and train the Deep Learning models.

187 A considerable number of models had to be trained and evaluated for this work. Therefore, it was crucial to exploit hardware  
 188 accelerators to be able to deliver results in a timely manner. For this reason, all hyper-parameters were set to powers of two if it  
 189 was possible. Google Colab was chosen as the platform due to its availability and ease of use. Google Colab offers virtual  
 190 machines (runtime) on Google Cloud Platform (GCP). Each runtime runs on Ubuntu OS and uses Jupyter Notebook as an  
 191 interface. The free version of Google Colab offers 13 GB of RAM along with access to K80 GPUs. Some models of this work  
 192 have millions of parameters and the dataset contains hundreds of thousands tensors. Therefore, purchasing the Pro version  
 193 of Google Colab was necessary. The Pro version of Google Colab provides access to P100 GPUs and 26 GB of RAM. Still,  
 194 at some models, there was memory limitations on the number of parameters. Also, for some data operations, software-level  
 195 techniques had to be implemented to mitigate the limited memory problem.

## 196 **Post-Processing Results**

197 CompoNet generates symbolic music, i.e., it is not an end-to-end approach and the output requires post-processing before being  
198 ready to be heard. While different tasks can be implemented in the post-processing, the following were implemented in the  
199 current work:

### 200 **Trimming**

201 OrchNet generates 64 time steps per each sample, the trailing time steps which are the result of zeros padded in the input need  
202 to be trimmed. These trailing time steps share the same probability value. This characteristic was utilized by traversing the  
203 generated sequence backwards and monitoring the generated probabilities. The sequence was trimmed at the point where the  
204 first change in the probabilities was observed.

### 205 **Sampling**

206 The output branches of OrchNet use softmax activation function. In a classification problem, this function generates the  
207 probability of each class being selected. Here, instead of choosing the option with the highest probability, all options are  
208 sampled with respect to their probabilities.

### 209 **Key and Scale Matching**

210 The dataset was generated from MIDI files. While MIDI standard is capable of storing key signature and scale information, it is  
211 not mandatory. Since many of the MIDI files in the dataset did not contain this information and it is not possible to accurately  
212 determine the scale and key signature of a piece automatically, OrchNet and CompoNet are not designed to contain and learn  
213 these information. Therefore, the generated music needs to be matched to a scale and key. This was done by checking all 12  
214 major scales and counting the number of notes that do not follow the scale pattern (alteration). The scale with the minimum  
215 number of alterations was chosen and alterations were matched to the closest neutral note (Tonic, Subdominant, Dominant)  
216 in the scale. The sequence was trimmed at a note (or chord) matching the scale's Tonic to create a cadence sensation. The  
217 minor scales were not considered because each minor scale is relative to a major scale and they share the same key signature;  
218 subsequently, they share the same notes. Although a major scale is different from its relative minor, for the purpose of scale  
219 matching they can be used interchangeably.

### 220 **Note Quantification**

221 Since 72 different possibilities are considered for note duration, the output does not automatically snap to standard note duration  
222 such as whole or black notes. It is possible to have a note which is only one pixel short of a standard note duration. Similarly,  
223 it is possible to have a note with only one pixel. In these cases it is reasonable to apply note quantification to enhance both  
224 rhythm and readability of the generated score. It is worth mentioning that the 72 possibilities were considered for the input data  
225 because different pieces have different rhythmic patterns and there is no one-fits-all approach to mitigate this.

### 226 **Assigning Instrument**

227 One advantage of generating music in symbolic format rather than audio format is that each track can be assigned to different  
228 instruments and ensembles to find the most suitable timbre for the generated sequence. This was done considering the available  
229 options and aesthetic aspects.

### 230 **Octave and Range Matching**

231 OrchNet has 4 bits to encode the octave of an interval. Due to the employed sampling mechanism, it is possible to have intervals  
232 spanning over three octaves. However, the assigned instrument might have a limited range and not be able to perform the  
233 generated note. In such scenario, the out-of-range notes were mapped to identical note which is in the feasible range.

### 234 **Audiolization**

235 The generated symbolic music needs to be converted to audio format to be heard. This was done using *Sibelius Music Notation*  
236 *Software 6* and its *Essential Sound Set*<sup>31</sup>.

## 237 **Turing Test**

238 Evaluation of results is one the major challenges in AMG<sup>32</sup>. This is because it is not possible to develop metrics to assess the  
239 quality of generated results in an objective manner. To mitigate this, Turing test can be utilized to evaluate the results if human  
240 music listeners can tell the difference between the pieces generated by our system and human music composers (shown in Table  
241 3).

**Table 3.** List of pieces used in the Turing test

Composer	Piece Name
Albeniz	Torre Bermeja Serenata from Piezas caracteristicas
Alkan	Etude Op 35 No 5
Bacewitz	Sonata No 4 for Piano Violin Mov 4
Coleridge-Taylor	La Tarantelle Fretillante
Faure	Dolly suite
Liszt	Valse Oubliee No 1
Meyerbeer	Grand March
Satie	Gnossienne No.3
Shostakovich	Ten short piano pieces No 1 Recitative
Wagner	Traume

### 242 **Test Design**

243 To perform the Turing test, 10 pieces were generated by CompoNet to be compared against 10 randomly-selected pieces from  
244 the dataset. To keep the test fair and reasonable, the following steps were taken:

- 245 • To ensure respondents are judging the pieces rather than answering the question from their memories, famous pieces  
246 were removed from consideration.
- 247 • To ensure the sequence of pieces is not impacting the results, they were put in a random sequence.
- 248 • To ensure the quality of audio is not impacting the results, all 20 pieces were audiolized with the same engine.
- 249 • To increase the likelihood of respondent completely listening to the piece before making a decision, all pieces were kept  
250 under one minute.

251 For each piece, respondents were asked "*How do you evaluate this piece of music?*" and they were asked to choose one the  
252 following options:

- 253 • A. I am confident it is composed by a human.
- 254 • B. I think it might be composed by a human.
- 255 • C. I am not sure.
- 256 • D. I think it might be composed by AI Software.
- 257 • E. I am confident it is composed by AI Software.

### 258 **Conducting the Test**

259 *Questionpro* website was used as the platform of conducting the test. The link to the test was published via author's social  
260 media, asking for volunteers to take the test. The link was visited 260 times, the test was started 169 times, and 88 respondents  
261 completed the test. This study was approved by the institutional review board of North Carolina A&T State University. No  
262 personal information was collected from the study participants. Informed consent was obtained from each participant prior to  
263 their participation in the study. The study was fully compliant to the IRB guidelines.

### 264 **Demography of Respondents**

265 Before beginning the test, respondents were asked to self-identify their age and gender groups voluntarily. Fig. 7 and Fig. 8  
266 illustrate the age and gender distributions of respondents respectively. Additionally, *Questionpro* automatically provides  
267 geographical distribution of respondents which is provided in Table 4 and Fig. 9. Based on these information, the respondents  
268 of this test include a diverse range of age, gender, and geographical location, which reduces potential biases in the results.

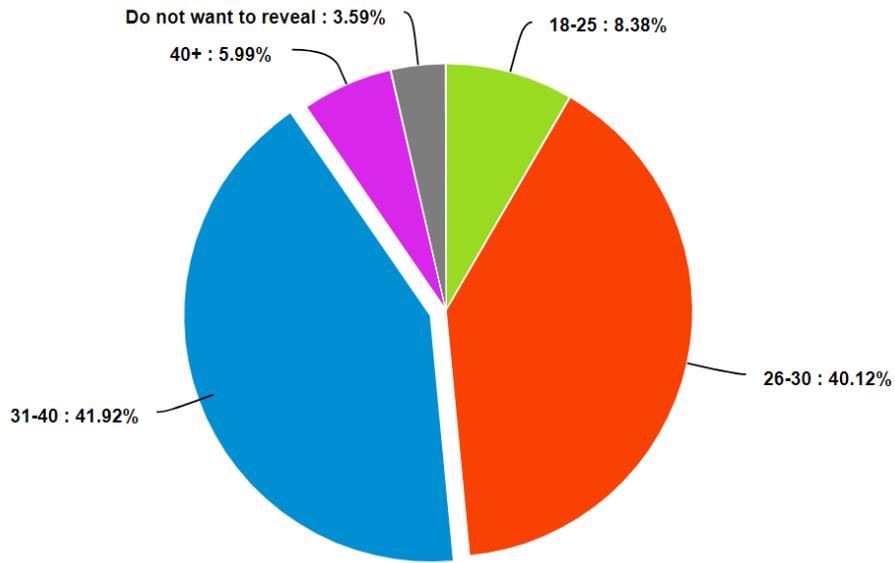


Figure 7. Age distribution of all respondents

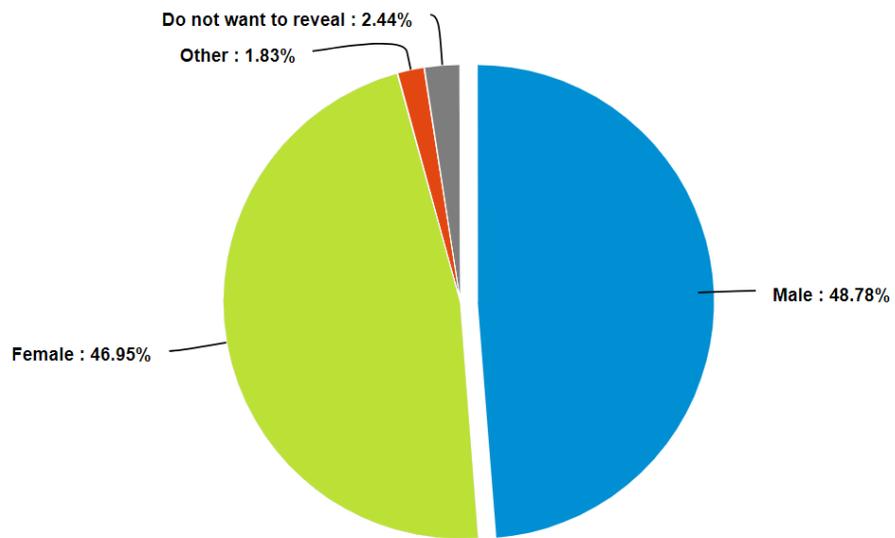


Figure 8. Gender distribution of all respondents



Figure 9. Geographical distribution of all respondents

Table 4. Geographical distribution of all respondents

Country	Responses
United States	65.19%
Iran	23.76%
France	2.76%
Germany	1.66%
Canada	1.10%
United Kingdom	1.10%
Italy	1.10%
Netherlands	1.10%
Romania	0.55%
Switzerland	0.55%
Austria	0.55%
Australia	0.55%
Total	100.00%

**Table 5.** Mapping responses to numerical values

Response	Value
I am confident it is composed by a human.	1
I think it might be composed by a human.	0.5
I am not sure	0
I think it might be composed by AI Software.	-0.5
I am confident it is composed by AI Software.	-1

## Test Results

A scoring mechanism was used to compare the pieces of music in the test. The score of each piece was calculated as:

$$s = \frac{\sum_{i=1}^n v_i}{n} \quad (1)$$

where  $s$  is the score of the piece,  $i$  is the number of respondent,  $n$  is the total number of respondents, and  $v_i$  is the value of the  $i_{th}$  respondent's response for the piece based on Table 5. The calculated score is in the range of  $[-1, 1]$ , where  $-1$  means all respondents have confidently identified the piece synthesized by CompoNet, while  $1$  shows that all respondents were confident the piece was composed by human.

Based on Table 6, which contains the score of each piece calculated from the result of respondents who completed the test, both the highest and lowest scores belong to CompoNet outputs. Both human-composed pieces and CompoNet outputs are almost equally present in the upper and lower halves of Table 6. The statistical characteristics of CompoNet and pieces composed by human are provided in Table 7. Since the obtained scores do not follow a normal distribution and the number of samples in each group is 10, the non-parametric Mann-Whitney U test was used to check the validity of the null hypothesis:

$H_0$ ; There is a difference between the ranks of the pieces generated by CompoNet and composed by human.

versus the alternative hypothesis:

$H_1$ ; There is no difference between the ranks of the pieces generated by CompoNet and composed by human.

The result of Mann-Whitney test is as follows:  $U(N_{CompoNet} = 10, N_{Human} = 10) = 33.5, z\text{-score} = -1.209941, p\text{-value} = 0.226301 (p(x \leq Z) = 0.113151)$ . The test statistic  $Z$  is in the 95% critical value accepted range:  $[-1.9600 : 1.9600]$ . The  $U\text{-value}$  satisfies the critical value of  $U$  at  $p < 0.05$  which is 23. Therefore, the result is not significant at  $p < 0.05$ . The  $p\text{-value}$  equals 0.226301,  $(p(x \leq Z) = 0.113151)$ . This means that if we would accept  $H_0$ , the chance of type II error (accepting a false null hypothesis) would be too high: 0.2263 (22.63%). In other words, the difference between the randomly selected output of CompoNet and pieces composed by human is not high enough to be statistically significant. This indicates that CompoNet output is indistinguishable from pieces which are composed by human and CompoNet can mimic classic composers whose works are present in the dataset.

## References

- Hiller, L. A. & Isaacson, L. M. *Experimental Music; Composition with an Electronic Computer* (Greenwood Publishing Group Inc., USA, 1979).
- Xenakis, I. *Formalized Music: Thought and Mathematics in Composition*. Harmonologia series (Pendragon Press, 1992).
- Blackwell, T. Swarm music: improvised music with multi-swarms. *Artif. Intell. Simul. Behav. Univ. Wales* **10**, 142–158 (2003).
- Eck, D. & Schmidhuber, J. A first look at music composition using lstm recurrent neural networks. *Istituto Dalle Molle Di Studi Sull Intelligenza Artif.* **103**, 48 (2002).
- Van Der Merwe, A. & Schulze, W. Music generation with markov models. *IEEE MultiMedia* **18**, 78–85, [10.1109/MMUL.2010.44](https://doi.org/10.1109/MMUL.2010.44) (2011).
- Kosta, K., Marchini, M. & Purwins, H. Unsupervised chord-sequence generation from an audio example. In *ISMIR*, 481–486 (2012).
- Colombo, F., Muscinelli, S. P., Seeholzer, A., Brea, J. & Gerstner, W. Algorithmic composition of melodies with deep recurrent neural networks. *arXiv preprint arXiv:1606.07251* (2016).

**Table 6.** Score of each piece based on the Turing test

Piece	Composer	Score
8	CompoNet	0.3920
La Tarantelle Fretillante	Coleridge-Taylor	0.3391
Gnossienne No.3	Satie	0.3352
9	CompoNet	0.1591
4	CompoNet	0.1080
Sonata No 4 for Piano Violin Mov 4	Bacewitz	0.0852
Valse Oubliee No 1	Liszt	0.0398
Dolly suite	Faure	0.0341
Etude Op 35 No 5	Alkan	0.0000
10	CompoNet	0.0000
Grand March	Meyerbeer	-0.0568
Traume	Wagner	-0.0739
1	CompoNet	-0.0795
Torre Bermeja Serenata from Piezas caracteristicas	Albeniz	-0.1092
6	CompoNet	-0.1477
Ten short piano pieces No 1 Recitative	Shostakovich	-0.2102
5	CompoNet	-0.2386
3	CompoNet	-0.2644
7	CompoNet	-0.2670
2	CompoNet	-0.2784

**Table 7.** Statistical characteristics of the performed Turing test

	Median	Mean	Standard Deviation
CompoNet	-0.1136	-0.0616	0.2136
Human	0.0170	0.0383	0.1695

- 304 8. Sturm, B. L., Santos, J. F., Ben-Tal, O. & Korshunova, I. Music transcription modelling and composition using deep  
305 learning. *arXiv preprint arXiv:1604.08723* (2016).
- 306 9. Jaques, N. *et al.* Tuning recurrent neural networks with reinforcement learning (2017).
- 307 10. Dieleman, S., van den Oord, A. & Simonyan, K. The challenge of realistic music generation: modelling raw audio at scale.  
308 In *Advances in Neural Information Processing Systems*, 7989–7999 (2018).
- 309 11. Van den Oord, A. *et al.* Wavenet: A generative model for raw audio. arxiv 2016. *arXiv preprint arXiv:1609.03499* .
- 310 12. Huang, C. A. *et al.* An improved relative self-attention mechanism for transformer with application to music generation.  
311 *CoRR abs/1809.04281* (2018). [1809.04281](https://arxiv.org/abs/1809.04281).
- 312 13. Hadjeres, G., Nielsen, F. & Pachet, F. Glsr-vae: Geodesic latent space regularization for variational autoencoder  
313 architectures. 1–7, [10.1109/SSCI.2017.8280895](https://arxiv.org/abs/1710.1109) (2017).
- 314 14. Goodfellow, I. *et al.* Generative adversarial nets. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D. &  
315 Weinberger, K. Q. (eds.) *Advances in Neural Information Processing Systems 27*, 2672–2680 (Curran Associates, Inc.,  
316 2014).
- 317 15. Yang, L.-C., Chou, S.-Y. & Yang, Y.-H. Midinet: A convolutional generative adversarial network for symbolic-domain  
318 music generation. *arXiv preprint arXiv:1703.10847* (2017).
- 319 16. Mogren, O. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*  
320 (2016).
- 321 17. Yu, L., Zhang, W., Wang, J. & Yu, Y. Seqgan: Sequence generative adversarial nets with policy gradient. In *Thirty-first*  
322 *AAAI conference on artificial intelligence* (2017).
- 323 18. Guimaraes, G. L., Sanchez-Lengeling, B., Outeiral, C., Farias, P. L. C. & Aspuru-Guzik, A. Objective-reinforced generative  
324 adversarial networks (organ) for sequence generation models. *arXiv preprint arXiv:1705.10843* (2017).
- 325 19. Dong, H.-W., Hsiao, W.-Y., Yang, L.-C. & Yang, Y.-H. Musegan: Multi-track sequential generative adversarial networks  
326 for symbolic music generation and accompaniment. In *Thirty-Second AAAI Conference on Artificial Intelligence* (2018).
- 327 20. Dong, H.-W. & Yang, Y.-H. Convolutional generative adversarial networks with binary neurons for polyphonic music  
328 generation. *arXiv preprint arXiv:1804.09399* (2018).
- 329 21. Guan, F., Yu, C. & Yang, S. A gan model with self-attention mechanism to generate multi-instruments symbolic music. In  
330 *2019 International Joint Conference on Neural Networks (IJCNN)*, 1–6 (2019).
- 331 22. Raffel, C. *Learning-based methods for comparing sequences, with applications to audio-to-midi alignment and matching*.  
332 Ph.D. thesis, Columbia University (2016).
- 333 23. Schmidhuber, J. Composing music with lstm recurrent networks - blues improvisation, <http://people.idsia.ch/~juergen/blues>  
334 (2002).
- 335 24. The session data <https://github.com/adactio/thesession-data> (2020).
- 336 25. Nottingham dataset, <http://www-labs.iro.umontreal.ca/~lisa/deep/data> (2012).
- 337 26. Theory tab, <https://www.hooktheory.com/theorytab>.
- 338 27. HekmatiAthar, S. & Anwar, M. Music embedding: A tool for incorporating music theory into computational music  
339 applications (2021). [2104.11880](https://arxiv.org/abs/2104.11880).
- 340 28. Fedorak, B. Classical music of various composers in midi format, <https://www.kaggle.com/blanderbuss/midi-classic-music/>  
341 (2018).
- 342 29. Hao-Wen Dong, Y.-H. Y., Wen-Yi Hsiao. Pypianoroll: Open source python package for handling multitrack pianorolls  
343 (2018).
- 344 30. Colin Raffel, D. P. W. E. Intuitive analysis, creation and manipulation of midi data with pretty\_midi (2014).
- 345 31. Avid. Sibelius essential sound set (2021).
- 346 32. Liebman, E. & Stone, P. Artificial musical intelligence: A survey. *arXiv preprint arXiv:2006.10553* (2020).

## 347 Author contributions statement

348 S.H. and M.A. designed the study. S.H. implemented the algorithm and conducted the study. S.H. and M.A. analyzed the  
349 results. S.H. wrote the first draft. M.A. supervised the research and contributed to the writing of the manuscript.

350 **Competing interests**

351 The authors declare no competing interests.