

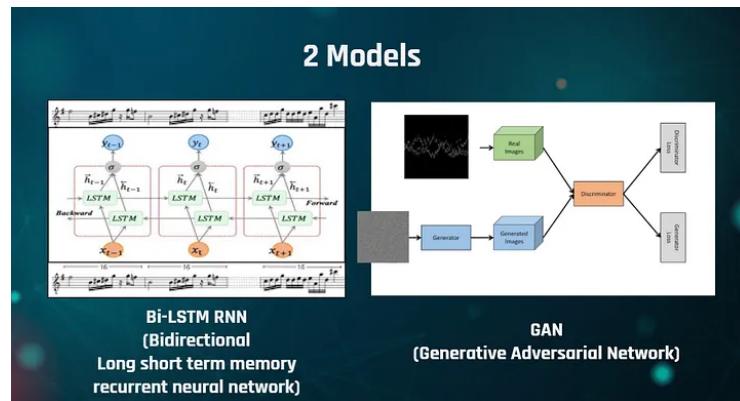
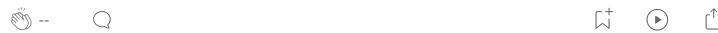
AI Classical Music Composer — Bi-LSTM & CNN-GAN

How to generate classical music for composers, musicians, or even non-specialists without prior knowledge of classical music theories and background.



Jonathan C.T. Kuo · Follow

Published in Analytics Vidhya · 27 min read · Oct 25, 2021



2 proposed AI models

Note: The source code of this project can be accessed [here](#).

1. Abstract

Artificial Intelligence could bring music composition to another level with limitless possibilities as an assistant for human musicians or an AI musician itself. Living in a digital era, classical music plays a dominant role in commercial films, movie trailers, game soundtracks, and more. However, there are no existing works that generate classical music in different eras. To fill this gap, this project proposes an AI music generator for classical music. So that it would be possible for composers, musicians, or even non-specialists without prior knowledge of classical music theories and backgrounds could quickly compose classical music according to their

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

Sign up for free

Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

Try for \$5/month

According to specific probabilistic rules, a Markov Chain is a system that changes its state from one to another in fixed timesteps. When applying the Markov Chain to music generation, each note in the training set is assigned a unique state. The model will then generate new states (notes), learning from the past states sequentially. The most straightforward design of a Markov Chain is to limit the model to use a single previous state to predict another. It is reasonable and feasible to apply Markov Chain for assisting real-time improvisational performance like jazz as musicians usually play spontaneously.

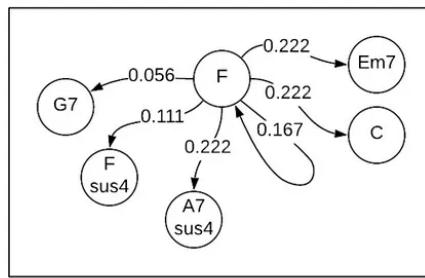


Figure 2-1 An example of weighted next possible notes. The next note after 'F' could be Em7, C, A7, F, G7 or itself: F.

However, when Markov Chains generated music from an existing musical corpus with full compositions, the results were not satisfying. It showed weird arrangements of notes and chords that were unmusical (Moorer, 1972).

2.1.2 Recurrent Neural Network

A Recurrent Neural Network (RNN) is a neural network where the current input is fed in with output from the previously hidden units. Each node of an RNN has a hidden state which allows the network to memorize a sequence of outcomes from the past. This feature of memorization means that RNN can be applied to tasks like Natural Language Processing (NLP), Time Series Pattern Predictions, and more. Since the traditional RNNs have drawbacks of not being able to memorize long sequences of data and suffering from problems of exploding and vanishing gradient, most studies applied Long Short-Term Memory (LSTM, a kind of RNN) as an alternative to generate longer sequential information like text, speech, and music. Typically, an LSTM consists of a forget gate that removes useless information, an input gate that adds useful information together, and finally an output gate that extracts valuable information as an output of the current cell. Eck and Schmidhuber (2002) demonstrated that an LSTM model was able to compose novel blues music in styles that were similar to its trained (input) data.



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

F:9		F:9 F:9 F:9 F:9 D:min7
D:min7	G:9	
C:maj	F:9	
C:maj		C:maj

Figure 2-2 An example of 4-bar text-based LSTM training data. The right-hand side represents the score from the left with removal of the bars.

Boulanger-Lewandowski et al. (2012) suggested a combined approach that improved their previous study in generating polyphonic music by introducing the RNN-RBM model. Each RNN hidden unit is combined with an RBM. An RBM or Restricted Boltzmann Machine is a bidirectional connected and stochastic (or generative) system with visible nodes and a hidden node where all visible nodes are connected, but the hidden nodes do not. Each state of the RNN unit gets both inputs from its previous state and the RBM observation vector, and this allows the model to better predict the coming notes by learning from the prior timestep of notes. However, the Magenta project (Shiebler, 2016) commented that this model had difficulty generating rhythmic and melodic music for longer timesteps.

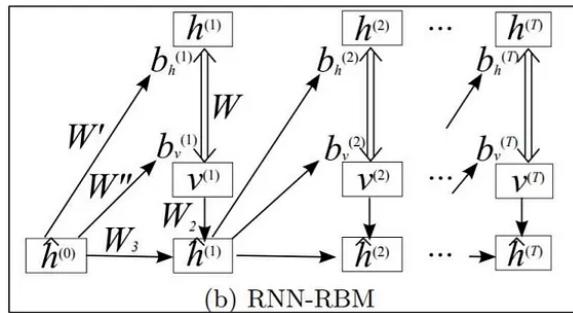


Figure 2-3 The structure of RNN-RBM model. each RNN hidden unit (lower-half of the diagram) is combined with an RBM (upper-half of the diagram)

2.1.3 Generative adversarial network

A Generative Adversarial Network (GAN) (Goodfellow et al., 2014) is a type of neural network that produces plausible data from scratch. There are two systems in a GAN – the generator and the discriminator. The generator tries to fool the discriminator by producing fake data while the discriminator distinguishes real and fake samples. By competing, the discriminator penalizes the generator for generating fake data, and the generator gets better at creating new instances that look real. A study by Dong et al. (2017) adopted Convolutional GAN with Wasserstein Loss to generate four bars of multi-track piano-rolls. The input data size is set as 96 (timesteps) times 84



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

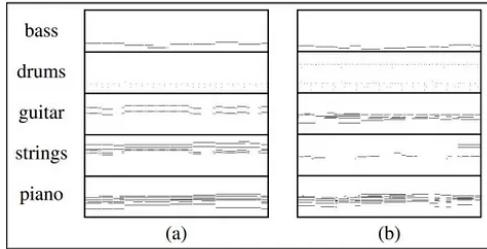


Figure 2-4 An example of the training data for the Convolutional GAN model. There are two instances (a & b) of music fragments of four bars with five tracks.

2.2 Proposed Methods

This project will evaluate two neural network architectures — the Deep Bidirectional LSTM model and a GAN model with convolutional layers. It is hypothesized that the Deep Bidirectional LSTM model could improve shortcomings of existing studies by predicting musical structures with longer timesteps and more meaningful melodies. On the other hand, the GAN models try to solve the LSTM models' problem of getting stuck in a loop (generating the same musical structure) by developing music from scratch (random noise). Both proposed models are supposed to tackle the drawbacks of current studies, and all results will be compared and evaluated by appropriate evaluation metrics — the Fréchet Inception Distance (FID) score and nearest neighbor search — and human survey.

3. Design, Solution, and System

3.1 Deep Bidirectional LSTM RNN

3.1.1 Deep neural networks

It is believed that a deeper neural network is more capable of analyzing and processing thorough characteristics of information than a shallow one. Typically, a shallow neural network contains only one or two hidden layer(s), while a deep neural network consists of more than two hidden layers. Therefore, to let the proposed model have a better understanding of input data, a deep neural network is adopted in this project.

3.1.2 LSTM Architecture

Each LSTM layer contains blocks of LSTM cells with multiple gates to memorize sequences of information. The cell states transfer valuable information to other cells and multiple gates with unique functions to keep the LSTM network learning and memorizing information. The below figure illustrates the architecture of an LSTM cell.



Sign up to discover human stories that deepen your understanding of the world.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

There are three types of gates in a LSTM cell: a forget gate, an input gate, and an output gate.

First, the forget gate will abandon useless information before passing them to future cells. The formula of the forget gate is as follows:

$$f_t = \sigma(W_f * [h_{t-1}, x_t] + b_f)$$

Second, useful information will be decided by the input gate (i_t) and a candidate cell state (\tilde{C}_t) will also be created for regulating the current cell state. The equations are as follows:

$$i_t = \sigma(W_i * [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_c * [h_{t-1}, x_t] + b_c)$$

Now, the output of the forget gate, input gate and the candidate cell state are ready, an updated cell state can be derived as follows:

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Finally, the output gate will determine the hidden state (h_t). The hidden state is calculated as follows:

$$h_t = \sigma(W_o * [h_{t-1}, x_t] + b_o) * \tanh(C_t)$$

The above equations can be written in python pseudocode as follows:

```
1  def LSTM(c_t_prev, h_t_prev, input):
2      combine = h_t_prev + c_t_prev
3      f_t = forgetGate(combine)
4      i_t = inputGate(combine)
5      c_t_candidate = tanh(combine)
6      c_t = c_t_prev * f_t + c_t_candidate * i_t
7      h_t = outputGate(combine) * tanh(c_t)
8
9  c_t = [array of c_t]
10 h_t = [array of h_t]
11
12 for input in allInput:
13     c_t, h_t = LSTM(c_t, h_t, input)
```

Figure 3-2 The algorithm to calculate LSTM's next hidden state

3.1.3 Bidirectional LSTM

A bidirectional LSTM model is a system containing two LSTM layers that pass-through data in opposite directions. Each output from the system receives information from both the hidden states of the backward and forward layers. Comparing with unidirectional LSTMs, bidirectional LSTMs not only learn information from the past but from the future as well. This feature allows the model to predict sequential information more accurately. From the image below, we can see that it is an acyclic graph. The output for each timestep can be derived as follows:



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

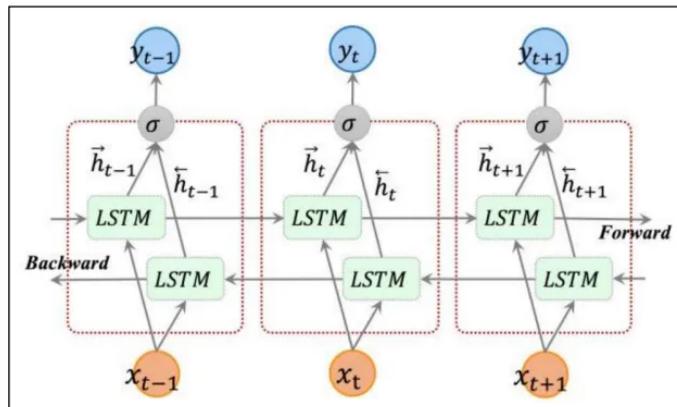


Figure 3-3 An illustration of Bi-LSTM architecture

3.1.3 A combined approach

A Deep Bidirectional LSTM RNN is a model of multiple bidirectional LSTM layers stacked together. It is hypothesized that combining deep neural networks and the feature of bidirectional layers will perform better than methods done by previous studies. Therefore, the first model in this project will be based on this design.

The model contains three bidirectional LSTM layers with 128 neurons each, followed by two fully connected layers to modify the input vectors' dimension and adapt the desired output shape. The initial plan for the activation function and the optimizer is to use Rectified Linear Unit (ReLU) and RMSProp since they are most widely used for LSTM models. However, other optimizers and activation functions will be further tested during the implementation steps.

3.2 CNN based GAN

A GAN is made up of a generator and a discriminator with any neural network architecture that is best for classifying the training data. It is mostly used for generating images with CNN. However, other neural networks like LSTM or RNN have also been used to predict sequential data. In general, the training process for GAN can be divided into two parts – training of the discriminator and training of the generator. First, randomly generated data will be fed into the generator.

The randomly generated input size does not necessarily need to be the same as the desired generated output. Since the dimension of the input data can be changed in the generator, usually a smaller size than the dimension of the output data will be used. After the generator has produced a fake (generated)



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

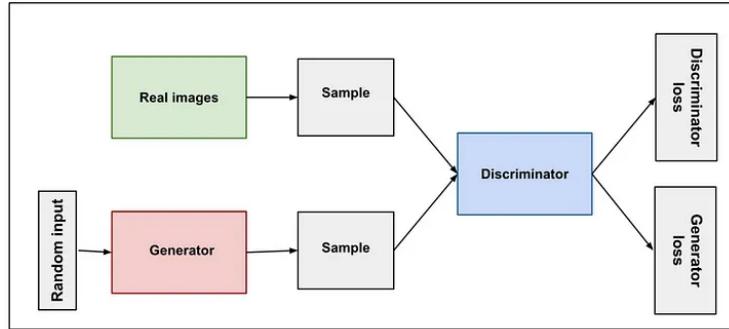


Figure 3-4 The main architecture of a GAN model

This model utilizes convolutional layers as the main neural networks for both the discriminator and the generator. Since a midi file is a file with different frequencies that progress with time, each midi file can be transformed to a 2-D array where the x-axis denotes the timesteps, and the y-axis represents different pitches of notes. In this case, we can treat the 2-D arrays as images to be trained in a convolutional neural network-based GAN model easily.

The generator consists of a fully connected layer with many neurons followed by four convolutional layers that gradually modify the shape of the data vectors into the desired data shape. The final output shape of the generator will be the shape of the generated music. Here, we set the output shape to be 1000 times 84, where 1000 represents the number of timesteps in the music 17 and 84 stands for the number of pitches. It is worth noting that the number of pitches is set to be 84 in order to save computational cost. Since the notes and chords of almost all pieces fall between in midi index 24 to 108, which contains in total 84 notes, it is better to truncate the unnecessary 0s to speed up the training process. After the training and generating process, the generated pieces will then be added two arrays of 0s to become arrays of size of 1000 by 128. The kernel size for each convolutional layer is 5 by 5, with strides of either (2, 1), (1, 2), or (2, 2) to fit the final image shape. The last convolutional 2D transpose layer uses tanh as the activation function to determine the final output of each neuron.

The discriminator contains two convolutional layers and one fully connected layer. Each convolutional layer has a kernel size of 5 by 5 and a stride of (2, 2), which is similar as the generator's settings.

3.3 One-hot encoding on Midi files

One-hot encoding is the process of transforming categorical information



Sign up to discover human stories that deepen your understanding of the world.

- Free**
- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

The Fréchet Inception Distance (FID) is designed for evaluating the generated results (fake data) produced by the GAN models with the real data. The equation to calculate the evaluation score is shown below.

$$FID \text{ score} = d^2 = \|\mu_1 - \mu_2\|^2 + \text{tr}(C_1 + C_2 - 2\sqrt{C_1 C_2})$$

Here μ_1 and μ_2 denotes the mean of all feature vectors of the true and fake data, respectively. Similarly, C_1 and C_2 denotes the covariance matrix of all feature vectors of the true and fake data, and tr denotes the trace linear algebra operation. A lower FID score means the generated results are more similar to the real data, which is desired in most cases. On the other hand, a higher FID score signifies that the generated results are less like the real data.

3.4.2 Pitch Histogram

Note consistency and note diversity can be visualized by plotting the pitch histogram. A pitch histogram is plotted by summing up each unique frequency count of notes and chords in a histogram style. It is believed that it is an effective measure to classify different genres of music. In this project, it can be used as a tool to distinguish the eras among all generated pieces of music by all models.

3.4.3 Nearest Neighbor search

A way to check whether there is any similarity between a generated piece and the pieces from the training set is by searching the nearest neighbor of that generated piece. By calculating the root mean square error between that generated piece and all pieces from the training set, we should retrieve the closest piece by searching for the lowest error.

4 Methodology and Implementation

4.1 Midi files Pre-processing

This project uses the Music21 and Pypianoroll python libraries to extract information from all midi files into notes and chord sequences. Their pitches and octaves correspond to each time step. Since this project's scope only covers a single instrument's generation, the training data will only contain pieces of music played by the piano.

All training data is acquired from the GiantMIDI-Piano repository on the ByteDance GitHub repository. According to their documentation, there are 10854 midi files of classical music available. However, duplicates and empty files were found in the acquired dataset after the data cleaning and categorization process. The final dataset contains 546 pieces of baroque



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

```
[ 'A2',
  'C#3',
  'E3',
  'A3',
  'G#3',
  'F#3',
  'E3',
  'G#3',
  'B3',
  'D4',
  'C#4',
  'B3',
  'A3',
```

Figure 4-1: An excerpt of unique notes and chords set from the baroque music dataset.

Since LSTM models learn from given sequences and their next prediction(s), an additional process is needed to build series of sequences and their corresponding output (next forecast). This project sets the sequence length as 50, meaning that the LSTM network learns to predict the next note given its previous 50 notes or chords. First, the program will store the first 50 notes or chords, then transform them into integers according to the unique note-chord dictionary built in the previous into the input sequence array. The next corresponding note or chord of the last sequence 20 is then stored in the output array. Retrieving all sequences of notes and chords, each element in the output array is normalized between 0 and 1 for the later training process's ease.

first sequence	second sequence	third sequence
[154]	[168]	[186]
[168]	[186]	[155]
[186]	[155]	[199]
[155]	[199]	[191]
[199]	[191]	[186]
[191]	[186]	[199]
[186]	[199]	[164]
[199]	[164]	[178]
[164]	[178]	[169]
[178]	[169]	[164]
[169]	[164]	[155]
[164]	[155]	[169]
[155]	[169]	[187]
[169]	[187]	[156]
[187]	[156]	[200]
[156]	[200]	[192]
[200]	[192]	[187]
[192]	[187]	[200]
[187]	[200]	[165]
[200]	[165]	[179]
[165]	[179]	[170]
[179]	[170]	[165]
[170]	[165]	[156]
[165]	[156]	[165]
[156]	[165]	[170]
next input: 165	next input: 170	next input: 179

Figure 4-2: An example of sequence of notes and chords stored in arrays and their corresponding outputs in integers format.



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

Suppose a note or chord occurs in any timestep, the array of that particular timestep store a ‘1’. On the other hand, empty notes are denoted by a ‘-1’.

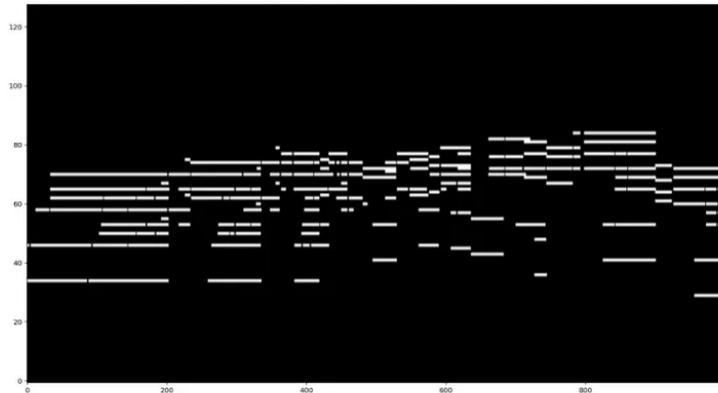


Figure 4-4: An example of a piece of music represented in 2-D array image format (*Valse carnavalesque*, Op.73 by Chaminade, Cécile). Here, white dots or lines indicate a presence of notes.

The final step is to store all images into a NumPy array with a shape of (x, y, z) where x denotes the number of unique pieces of music, y indicates timesteps of each training piece, and z means the range of notes. In this project, the array shape setting is (x, 1000, 84) for simplicity since the objective is to assess the CNN-GAN model’s capability to generate structured music and meaningful melodies instead of generating long pieces that usually require excessive computing power to finish the task. In standard Midi expression, there are 128 notes in total. However, by observing all images from the training set, most of them only have notes span from minimum C1 (index 24) to maximum C7 (index 108). This means that notes below index 24 and above index 108 can be ignored in that they are empty notes which are not conducive for the latter training process. Hence, a trimming process is done for all images before training. It results in the final input data shape of (x, 1000, 84), where x denotes the number of unique pieces of music.



Figure 4-5: Trimming the image above to have a range of notes from C1 to C7.



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

batch normalization is added. The first fully connected layer contains 128 nodes with a rectified linear activation function (ReLU). The final layer contains x nodes where x is determined by the distinct notes or chords from the training corpus. Since it is known that if a note or chord is not within the training corpus, it will not show up or be predicted in the future by the neural network, we can safely implement the final fully connected layer in this manner to save memory and computational time. The model is compiled with a loss function of cross-entropy loss and RMSprop for its optimizer.

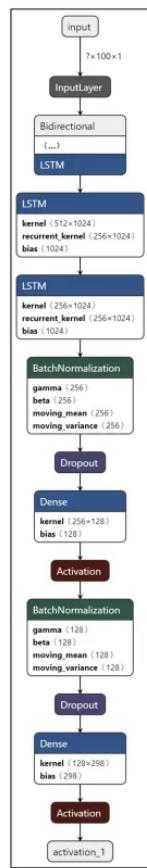


Figure 4-6 The Bi-LSTM network structure

Similar to the training step, the prediction network possesses the same architecture as the training model. The only difference is that the optimal weights from the last epoch of the training process are added to the neural network for sequence prediction.

Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.



Figure 4-7: generating 2 notes in 2 iterations. The first row (sequence) is the initial sequence determined by the random index. Here, B3 is the first note being generated, and D4 being the second. The process goes on until reaching the maximum number of notes (iterations) specified.

4.2.2 CNN based GAN

This CNN-GAN model consists of 2 neural networks — a discriminator and a generator. Both networks adopt cross-entropy loss as their loss function by computing the difference between predicted labels and true labels. Since the generator is trying to fool the discriminator by generating fake images that make it the discriminator challenging to discern the authenticity, the generator's loss function compares the discriminator's decisions on fake images to an array of ones. As for the discriminator, there are two steps (two losses). It compares real images with an array of ones and fake images with an array of zeros. Combining both losses yields the total loss of the discriminator's loss.

```
def generator_loss(fake_output):
    return cross_entropy(tf.ones_like(fake_output), fake_output)
```

Figure 4-8: the implementation of the generator's loss

```
def discriminator_loss(real_output, fake_output):
    real_loss = cross_entropy(tf.ones_like(real_output), real_output)
    fake_loss = cross_entropy(tf.zeros_like(fake_output), fake_output)
    total_loss = real_loss + fake_loss
    return total_loss
```

Figure 4-9: the implementation of the discriminator's loss

The generator's network structure begins with introducing a small random seed as input to a fully connected layer with a much larger output size. After reshaping the output data, four transposed convolutional layers are stacked together to upsample the data. Each convolutional layer has a kernel size of 5 (both height and width) and a stride of (2, 1) or (1, 2), or (2, 2) to upsample the data to meet the same shape of the input data which is (1000, 84). The paddings are the same to preserve the original size of the data. Finally, except the last layer uses tanh as an activation function, all layers are



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

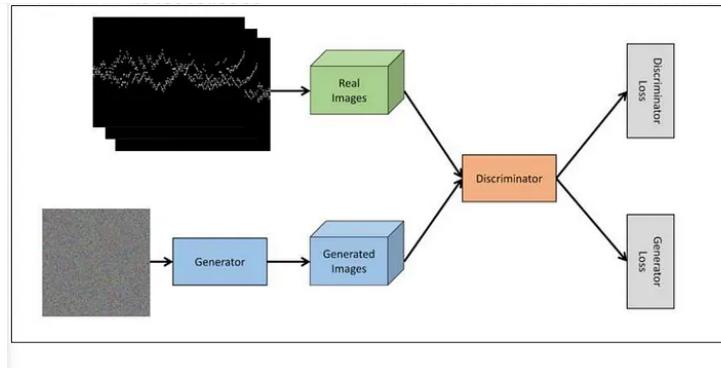


Figure 4-10: The CNN-GAN model structure (The Generator is fed in with some noise)

4.2 Data Conversion

For the Bi-LSTM model, after a new sequence of numbers is generated, all elements from the sequence will first be converted to the representation of either notes or chords accordingly by referencing the unique set of notes and chords from the training set, then be converted to either notes or chords in Midi format. Here, a limit of the note number and chords need to be specified, and so does the tempo. Take modernist music as an example: a maximum of 300 notes and an offset of 0.3 between notes are good settings. Here, the offset controls the tempo of the music. When the number is larger, the tempo is slower, and vice versa.

For the CNN-GAN model, the final step is to convert the array with the size of (1000, 84) back to the size of (1000, 128), the acceptable data format by the Pypianoroll library, then convert it to a midi file. The predicted array concatenates an empty array with a size of (1000, 24), and the combined array is concatenated by the other empty array of size (1000, 20).

4.3 Evaluation

4.3.1 Pitch histogram

The purpose of using pitch histograms to analyze generated music is that this graph is an efficient way to visualize note similarities and diversities from a given piece of music. By counting the most frequent notes occurring in a piece, it is common to determine the music. However, suppose most notes fall in only a few (1 or 2) categories of note. In that case, it may indicate a failure of the generation of music since almost most pieces have a specific diversity of distribution of notes.

Pitch Class Histogram



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

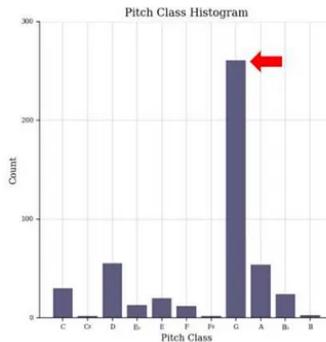


Figure 4-12: An example considered as a failure. It is a generated piece in baroque style with a few training epochs (20 epochs). This indicates that the model did not learn much from the training process and could only produce music with repeated notes. The dominant key is pointed by the red arrow.

4.3.2 FID

There are two required components to calculate the FID scores — the set of real images and the generated image. For each genre, all music from the training set belongs to the ‘real’ image set, and the generated image belongs to the ‘generated’ image set. The features that will be extracted for both the ‘real’ and ‘generated’ sets are part of the images themselves, and only a segment of each piece will be extracted. The segments are images with a shape of (400, 48), where 400 denotes timesteps and 48 means the range of notes. The reason to set the range of notes to 48 is that from observation, most notes and chords fall in this range, and it saves a lot of computational costs to retrieve features from each piece of music from the training set. Next, the mean of all feature vectors and the covariance matrix of all feature vectors for the ‘real’ and ‘generated’ set are first calculated. The FID score can then be retrieved from the below equation.

$$FID \text{ score} = d^2 = \|\mu_1 - \mu_2\|^2 + \text{tr}(C_1 + C_2 - 2\sqrt{C_1 C_2})$$

```

1 ▼ def calculate_fid(real, generated):
2     # calculate mean and covariance statistics
3     mu1, sigma1 = real.mean(axis=0), cov(real, rowvar=False)
4     mu2, sigma2 = generated.mean(axis=0), cov(generated, rowvar=False)
5     # calculate sum squared difference between means
6     ssdiff = numpy.sum((mu1 - mu2)**2.0)
7     # calculate sqrt of product between cov
8     covmean = sqrtm(sigma1.dot(sigma2))
9     # check and correct imaginary numbers from sqrt
10    if iscomplexobj(covmean):
11        covmean = covmean.real
12    # calculate score
13    fid = ssdiff + trace(sigma1 + sigma2 - 2.0 * covmean)
14    return fid

```

 Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

► Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

One piece from each genre and each category ('real' an 'generated') is chosen, and only an excerpt (between 15 to 30 seconds) of that piece will be available to the respondents since letting the respondents listen to the whole piece in a survey will make them fatigued and therefore deteriorate the survey result. In each survey, the number of 'real' and 'generated' music is balanced and ordered randomly for fairness. The below snippet is an example of a question.

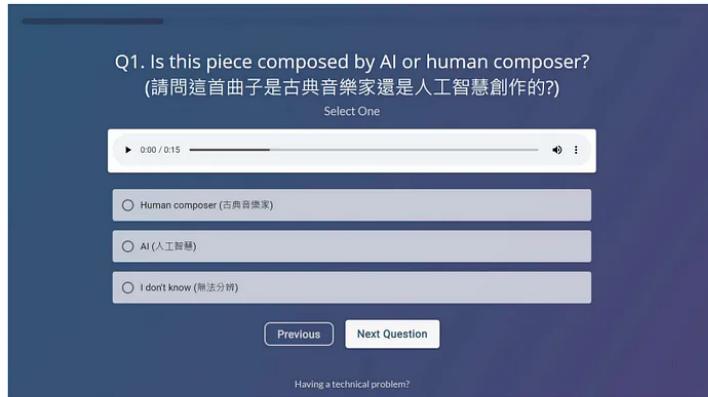


Figure 4-14 A snippet of one of the survey questions.

5 Results and Evaluation

5.1 Deep Bidirectional LSTM RNN

5.1.1 Note Diversity & Analysis

Below are four examples of generated music in different styles. Beginning with the Baroque, Classical styles and followed by the Romantic and Modernist styles. The results shown below demonstrate the Bi-LSTM model's capabilities to compose novel music passages by learning from the given training set.

- An example of a generated Baroque style music



Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

Music in this era often uses simple harmonic melodies to fulfill a balanced musical structure. This kind of composing idea was often found in Beethoven's music. The music generated below demonstrates a repetition of similar motifs (musical phrases or structures) that are waltz-like, which corresponds to one of the most common music composition ideas in the classical era. In my opinion, this very generated piece is one of the most successful one that resembles music from Beethoven.

Figure 5-2 Classical-I-Bi-LSTM

- An example of generated Romantic style music

According to the Music 21 library analysis, the expected key for the below piece of music is C sharp minor with a confidence of 70.2%, which corresponds precisely to the pitch histogram below.

Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

```

timeSignature = base_midi.getTimeSignatures()[0]
music_analysis = base_midi.analyze('key')
print("Music time signature: {0}/{1}".format(timeSignature.beatCount, timeSignature.denominator))
print("Expected music key: {0}".format(music_analysis))
print("Music key confidence: {0}".format(music_analysis.correlationCoefficient))
print("Other music key alternatives:")
for analysis in music_analysis.alternateInterpretations:
    if (analysis.correlationCoefficient > 0.5):
        print(analysis)

Music time signature: 4/4
Expected music key: c# minor
Music key confidence: 0.7027908535567045
Other music key alternatives:
C# major

```

Figure 5-4 Romantic-1-Bi-LSTM. An analysis on the piece shown above. It tries to predict its key with a confidence rate.

- An example of generated Modernist style music



Figure 5-5 Modernist-1-Bi-LSTM

- Pitch histogram of four generated pieces

From the below figures, we can observe that all pieces have diverse distributions of pitches. No single pitch noticeably dominates the whole pieces. Thus, it can be concluded that the Bi-LSTM model handles the note diversity well and does not tend to memorize partial information from the training set.

Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

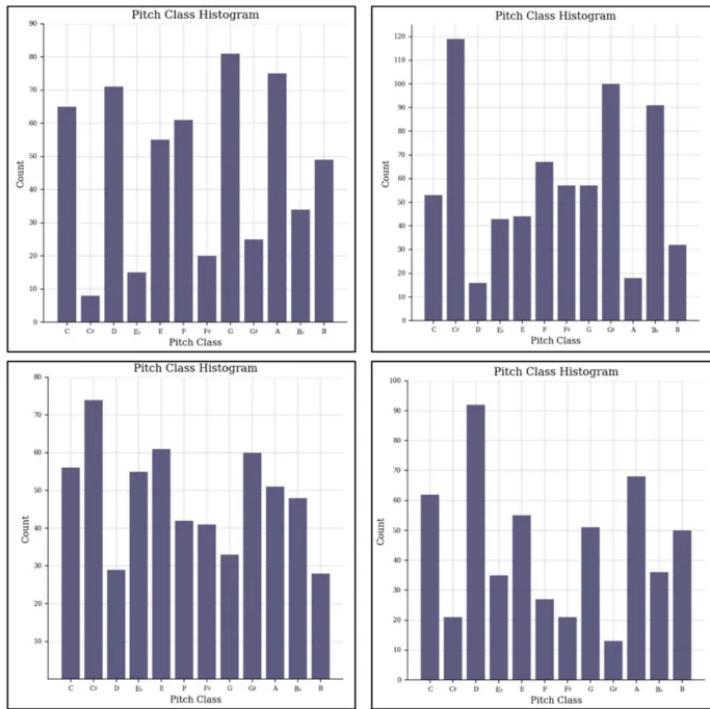
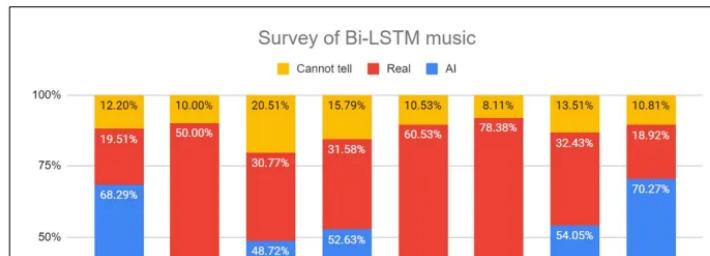


Figure 5-6 Four pitch histograms of the generated pieces from the Bi-LSTM models. Clockwise from top-left: Baroque-1-Bi-LSTM, Classical-1-Bi-LSTM, Modernist-1-Bi-LSTM, and Romantic-1-Bi-LSTM.

5.1.2 Survey

From the stacked column chart below, it is clear that over half of the respondents have difficulty in telling whether the given piece of music is composed by AI or human. Among all human-composed pieces, only half of which receives over fifty percentages of confidence indicating them to be composed by the AI. The results are the same for all pieces composed by the AI. Though some respondents regarded 'real' music as composed by AI, the statistics showing that most respondents picked 'Real' when they are given to listen to a piece composed by the AI. This indicates that the proposed Bi-LSTM model have the capability to trick humans' judgment.



Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

(image), which represents the initial states of the GAN training process, and the training set. Since the lower the FID scores, the similar the pieces (images) are to the training set, we can observe that each generated piece has a far lower FID score than the random one that corresponds to it. Therefore, we can say that the generated pieces look similar to those in the training set, which means that the GAN model mimics music from the training set well.

Major things that might affect the FID score could be the number of pieces in a particular training set, the piece itself, and the features that are being extracted for the FID evaluation. Also, unlike other image generating tasks which normally use the pretrained InceptionV3 model to evaluate generated image qualities, this project does not adopt this model as part of the FID score calculation due to the project's nature. The generated pieces are not exact the same as traditional images that have RGB color channels and have association with ImageNet labels. Hence, in this case, I consider the FID score as a reference of the model performance, instead of a hard standard to determine the goodness of a generated piece of music.

	Baroque-1-GAN	Classical-1-GAN	Romantic-1-GAN	Modernist-1-GAN
FID	2437.382	3596.256	4408.329	5939.203
FID (random)	20947.692	21013.249	20440.701	20610.990

5.2.2 Note Diversity & Analysis

Below are four examples of generated music in different styles. Beginning with the Baroque, Classical styles and followed by the Romantic and Modernist styles. The results shown below demonstrate the GAN model's capabilities to compose novel musical notes/ chords and structures by learning from the given training set.

- An example of a generated Baroque style music



 Medium

Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.

both the upper and the lower stave possess their melody, or at least the lower staff acts as the companion for the upper one. In contrast, the Bi-LSTM model does not showcase this kind of feature. Instead, it combines all chords or notes coinciding as a whole and represents them as a single chord. Though this might not simplify the musical structures, the GAN model might be a better option for generating more complex musical structures and producing more readable music sheets.

A way to check whether there is any similarity between a generated piece and the pieces from the training set is by searching the nearest neighbor of that generated piece. By calculating the root mean square error between Baroque-GAN-1 and all piece from the training set, we get the closest piece by searching for the lowest error. Below is a comparison between Baroque-GAN-1 (the generated piece) and the closest piece being found by the previously mentioned calculation. Observing the melody patterns of both pieces, we could say that both pieces have similar tendencies in developing continuous arpeggios. It is a good sign of showing the GAN model's learning capability, and not entirely copying the musical phrases from the training set.

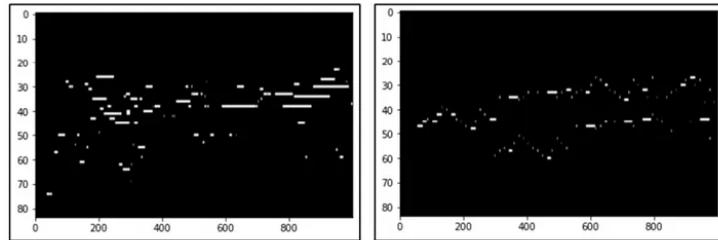


Figure 5-92 The nearest neighbor of Baroque-I-GAN. Left: Baroque-I-GAN; Right: a piece of music from the training set.

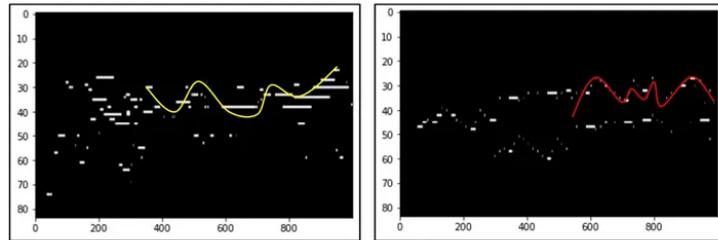


Figure 5-3 Both pieces have similar tendency in developing arpeggios going up and down continuously.

- An example of a generated Classical style music



Sign up to discover human stories that deepen your understanding of the world.

Free

- ✓ Distraction-free reading. No ads.
- ✓ Organize your knowledge with lists and highlights.
- ✓ Tell your story. Find your audience.

♦ Membership

- ✓ Access the best member-only stories.
- ✓ Support independent authors.
- ✓ Listen to audio narrations.
- ✓ Read offline.
- ✓ Join the Partner Program and earn for your writing.