

Conditional LSTM-GAN for Melody Generation from Lyrics

Yi Yu¹, Simon Canales²

¹Digital Content and Media Sciences Research Division, National Institute of Informatics, Tokyo

²Institut de génie électrique et électronique, École Polytechnique Fédérale de Lausanne, Switzerland
yiyu@nii.ac.jp, simon.canales@epfl.ch

Abstract—Melody generation from lyrics has been a challenging research issue in the field of artificial intelligence and music, which enables to learn and discover latent relationship between interesting lyrics and accompanying melody. Unfortunately, the limited availability of paired lyrics-melody dataset with alignment information has hindered the research progress. To address this problem, we create a large dataset consisting of 12,197 MIDI songs each with paired lyrics and melody alignment through leveraging different music sources where alignment relationship between syllables and music attributes is extracted. Most importantly, we propose a novel deep generative model, conditional Long Short-Term Memory - Generative Adversarial Network (LSTM-GAN) for melody generation from lyrics, which contains a deep LSTM generator and a deep LSTM discriminator both conditioned on lyrics. In particular, lyrics-conditioned melody and alignment relationship between syllables of given lyrics and notes of predicted melody are generated simultaneously. Experimental results have proved the effectiveness of our proposed lyrics-to-melody generative model, where plausible and tuneful sequences can be inferred from lyrics.

Index Terms—Lyrics-conditioned melody generation, conditional LSTM-GAN

I. INTRODUCTION

Music generation is also referred to as music composition with the process of creating or writing an original piece of music, which is one of human creative activities [1]. Without understanding music rules and concepts well, creating pleasing sounds is impossible. To learn these kinds of rules and concepts such as mathematical relationships between notes, timing, and melody, the earliest study of various music computational techniques related to Artificial Intelligence (AI) has emerged for music composition in the middle of 1950s [2]. Markov models as a representative method of machine learning have been applied to algorithmic composition [3]. However, due to the limited availability of paired lyrics-melody dataset with alignment information, research progress of lyrics-conditioned music generation has been obstructed.

With the development of available lyrics and melody dataset and deep neural networks, musical knowledge mining between lyrics and melody has gradually become possible [4],[5]. Melody [6] is a sequence of musical notes over time, in which each note is sounded with a particular pitch and duration. Generating a melody from lyrics is to predict a melodic sequence when given lyrics as a condition. Existing works, e.g., Markov models [7], random forests[8], and recurrent

neural network (RNN)[9], can generate lyrics-conditioned music melody. However, these methods cannot ensure that the distribution of generated data is consistent with that of real samples. Generative adversarial networks (GANs) proposed in [10] is a generative model which can generate data samples following a given distribution, and has achieved a great success in the generation of image, video, and text.

Inspired by the great success of GANs in various generative models in the area of computer vision and natural language processing, we propose a conditional LSTM-GAN model to compose lyrics-conditioned melody where a discriminator can help to ensure that generated melodies have the same distribution as real ones. To the best of our knowledge, this is the first study that conditional LSTM-GAN is proposed for melody generation from lyrics, which takes lyrics as additional context to instruct deep LSTM-based generator network and deep LSTM-based discriminator network. Our proposed generation framework has several significant contributions, as follows:

- i) A LSTM network is trained to learn a joint embedding in the syllable-level and word-level to capture syntactic structures of lyrics, which can represent semantic information of lyrics.
- ii) A conditional LSTM-GAN is optimized to generate discrete-valued sequences of music data by introducing a quantizer.
- iii) A large-scale paired lyrics-melody dataset with 12,197 MIDI songs is built to demonstrate that our proposed conditional LSTM-GAN can generate more pleasant and harmonious melody compared with baseline methods.

II. RELATED WORKS

Automatic music generation has experienced a significant change in computational techniques related to artificial intelligence and music [11], spanning from knowledge-based or rule-based methods to deep learning methods. The majority of traditional music generation methods are based on music knowledge representation, which is a natural way to solve the issue with some kind of composition rules [12]. Knowledge-based music rules are utilized to generate melodies when given the specified emotions by users [13]. Moreover, several statistical models [14] such as hidden Markov models, random walk, and stochastic sampling are discussed for music generation. For example, Jazz chord progressions are generated by Markov

model for music generation in [15] and statistical models are applied to music composition in [16].

With rapid advancement of neural networks, deep learning has been extended to the field of music generation. A hierarchical RNN for melody generation is proposed in [17], which includes three LSTM subnetworks. Beat profile and bar profile are exploited to represent rhythm features at two different time scales respectively. A neural network architecture [18] is suggested to compose polyphonic music with a manner of preserving translation invariance of dataset. Motivated by convolution that can obtain transposition-invariance and generate joint probability distributions over a musical sequence, two extended versions, Tied-parallel LSTM-NADE and bi-axial LSTM, are proposed to achieve better performance. A continuous RNN with adversarial training (C-RNN-GAN) [19] is proposed to compose MIDI classical music. RNN is considered to model sequences of MIDI data during adversarial learning. The generator is to transform random noise to MIDI sequences, while the discriminator is to distinguish the generated MIDI sequence from real ones.

Earliest work [20] for lyrics-conditioned melody generation is defined as generating a melody when given Japanese lyrics, patterns of music rhythms, and harmony sequences. Some constraints are determined to associate syllables with notes. Melody generation is realized by dynamic programming. In [21] the rhythmic patterns occurred in notes can be classified. Pitches that are most suitable for accompanying the lyrics are generated using n-gram models. Three stylistic categories such as nursery rhymes, folk songs, and rock songs are generated for given lyrics. A recently proposed ALYSIA songwriting system [8] is a lyrics-conditioned melody generation system based on exploiting a random forest model, which can predict the pitch and rhythm of notes to determine the accompaniments for lyrics. When given Chinese lyrics, melody and exact alignment are predicted in a lyrics-conditional melody composition framework [9], which is an end-to-end neural network model including RNN-based lyrics encoder, RNN-based context melody encoder, and a hierarchical RNN decoder. The authors create large-scale Chinese language lyrics-melody dataset to evaluate the proposed learning model.

Our work focuses on lyrics-conditioned melody generation using LSTM-based conditional GAN, which is quite distinct from existing works. A skip-gram model is trained to transform raw textual lyrics into syllable embedding vector which is taken as input together with noise vector for training a generator model. A discriminator model is trained to distinguish generated MIDI note sequences from real ones. A large English language lyrics-melody dataset is built to validate the effectiveness of our proposed recurrent conditional GAN for lyrics-to-melody generation.

III. PRELIMINARIES

Before we describe our melody generation algorithm, a brief introduction is given in the following to help understanding musical knowledge, the sequential alignment relationship

between lyrics and melody, and how to build the lyrics-melody music dataset.

A. Melody

Melody and lyrics provide complementary information in understanding a song with the richness of human beings' emotions, cultures, and activities. Melody, as a temporal sequence containing musical notes, plays an important role. A note contains two music attributes: pitch and duration. Pitches are perceptual properties of sounds that organize music by highness or lowness on a frequency-related scale, which can be played in patterns to create melodies [22]. Piano keys have MIDI numbers ranging from 21 to 108, which also represent the corresponding pitch numbers. For example, 'D5' and 'A4' can be respectively represented as 74 and 67 according to the mapping between notes and MIDI numbers. In music, duration [23] represents the length of time that a pitch or tone is sounded. Rests [24] are intervals of silence in pieces of music, marked by symbols indicating the length of the pause.

B. Lyrics

Lyrics as natural language represent music theme and story, which are a very important element for creating a meaningful impression of the music. An English syllable [25] is a unit of sound, which may be a word or a part of a word. According to timestamp information in MIDI files of music tracks, melodies and lyrics are synchronized together to parse the data and extract alignment information.

C. Alignment

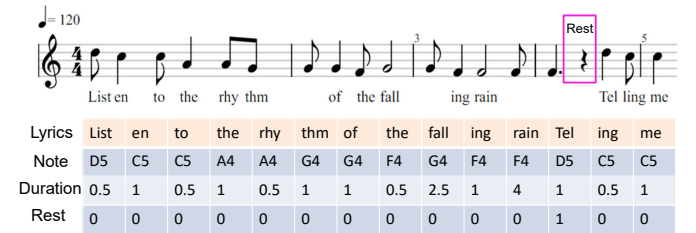


Fig. 1: An example of alignment between lyrics and melody.

An example data structure of the alignment between lyrics and melodies is shown in Fig. 1. Lyrics are divided to syllables. Each column represents one syllable with its corresponding note, note duration, and rest. With these music attributes, sheet music can be produced.

IV. MELODY GENERATION FROM LYRICS

Our proposed conditional LSTM-GAN for lyrics-to-melody generation model is shown in Fig. 2, which is an end-to-end generative learning conditioned with lyrics. A sequence of syllable embedding vectors concatenated with noisy vectors is taken as input of the generator network. The generated MIDI sequences together with the sequence of syllable embedding vectors are taken as input of the discriminator network, which aims to train a model for distinguishing generated MIDI note sequences from real ones. In addition, a tuning scheme

	RNN1 (generator)	RNN2 (discriminator)
Input	30 (random noise), 20 (syll. embedding)	3 (MIDI attributes), 20 (syll. embedding)
Layer 1	400, Fully-connected, ReLU	400, LSTM, tanh
Layer 2	400, LSTM, tanh	400, LSTM, tanh
Layer 3	400, LSTM, tanh	2 (real or fake), sigmoid
Layer 4 (output)	3 (MIDI attributes), fully-connected, linear	N/A

TABLE I: Configuration of the generator and discriminator

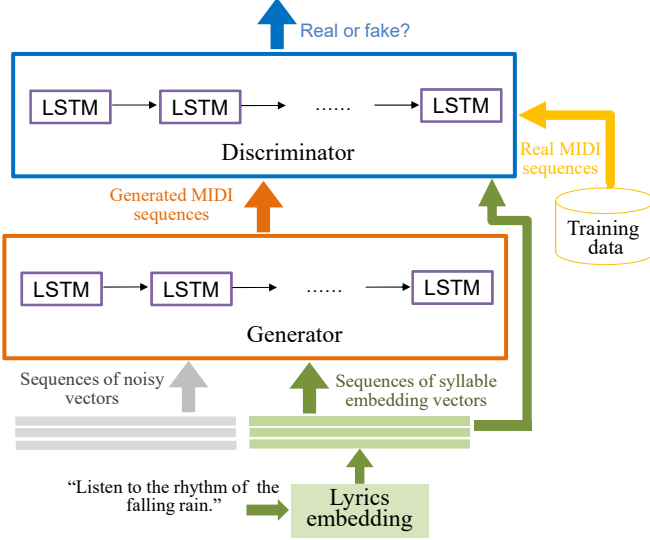


Fig. 2: Conditional LSTM-GAN for melody generation from melody.

is introduced to quantize MIDI numbers so as to generate melody sequence with discrete attributes. Both the generator and discriminator are unidirectional RNN networks with the configuration shown in Table I.

A. Problem formulation

Taking lyrics as input, our goal is to predict a melody sequentially aligned with the lyrics, MIDI numbers, note duration, and rest duration which can be synthesized with lyrics to generate a song. Our research problem can be formulated as follows: The syllables of lyrics as input are represented by a sequence $Y = (y^1, \dots, y^{|Y|})$. The melody as output is a sequence $X = (x^1, \dots, x^{|X|})$, where MIDI numbers, note duration, and rest duration are simultaneously predicted as the $x^i = \{x_{MIDI}^i, x_{dur}^i, x_{rest}^i\}$. Moreover, the time length of output sequence

$$\sum_{i=1}^{|X|} x_{dur}^i + x_{rest}^i \quad (1)$$

determines the length of synthesized song with lyrics.

B. Lyrics embedding

As the vocabulary of our lyrics data is large without any labels, we need to train an unsupervised model that can learn the context of any word or syllable. Specifically, we utilize our lyrics data to train a skip-gram model, which enables us

to obtain the vector representation that can encode linguistic regularities and patterns in the context of given lyrics.

We try to obtain embedding vectors of lyrics at both syllable-level and word-level. To this end, lyrics of each song are divided into sentences, each sentence is divided into words, and each word is further divided into syllables. Words $W = \{w_1, w_2, w_3, \dots, w_n\}$ are taken as tokens for training a word-level embedding model and syllables $S = \{s_1, s_2, \dots, s_m\}$ are taken as tokens for training a syllable-level embedding model. Then, we train each skip-gram model as a logistic regression with stochastic gradient descent as the optimizer, and the learning rate with an initial value 0.03 is gradually decayed every epoch until 0.0007. Context window spans 7 adjacent tokens and negative sampling distribution parameter is $\alpha = 0.75$. We train the models to respectively obtain the word-level and syllable-level embedding vectors of dimensions $V = 10$.

Denote $E_w(\cdot)$ and $E_s(\cdot)$ the word-level and syllable-level encoders respectively, and denote s a syllable from word w . Then, syllable embedding and word embedding are concatenated as follows:

$$\mathbf{y} = E_w(w) || E_s(s) = \mathbf{w} || \mathbf{s} \quad (2)$$

where $\mathbf{s} = E_s(s) \in \mathbb{R}^{10}$ and $\mathbf{w} = E_w(w) \in \mathbb{R}^{10}$ are the embedding of syllable s and word w , respectively, and the dimension of the overall embedding is $V = 20$.

C. Conditional LSTM-GAN model

In our work, an end-to-end deep generative model is proposed for lyrics-conditioned melody generation. LSTM is trained to learn semantic meaning and relationships between lyrics and melody sequences. Conditional GAN is trained to predict melody when given lyrics as input based on considering music alignment relationship between lyrics and melody.

1) *LSTM*: LSTM [26] networks are an extension to RNNs, which not only contain internal memory but also have capability of learning longer dependencies. An LSTM cell has three gates: input, forget, and output. These gates decide whether or not allow new input in, forget old information, and affect output at current time-step. In particular, at time-step t , the three states of the gates in an LSTM cell are given by:

$$i_t = \sigma(w_i[h_{t-1}, x_t] + b_i) \quad (3)$$

$$f_t = \sigma(w_f[h_{t-1}, x_t] + b_f) \quad (4)$$

$$o_t = \sigma(w_o[h_{t-1}, x_t] + b_o) \quad (5)$$

where i_t , f_t and o_t denote the input, forget, and output gates states, h_{t-1} is the output of the LSTM cell at previous time-step, w 's and b 's are weights and biases, x_t is the input of the LSTM cell, and $\sigma(\cdot)$ is the sigmoid function.

Then, the current output of the cell is computed by:

$$h_t = o_t \circ \tanh(c_t) \quad (6)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \quad (7)$$

$$\tilde{c}_t = \tanh(w_c[h_{t-1}, x_t] + b_c). \quad (8)$$

where \circ denotes the element-wise multiplication between vectors.

2) *GAN*: A GAN is proposed by Ian Goodfellow, et al. [10], which aims to train generative models by mitigating complex computation of approximating many probabilities. The general idea of GAN is to simultaneously train a generator $G(\cdot)$ and a discriminator $D(\cdot)$ with conflicting objectives. This method learns to predict new data with the same statistics as the training set. The generator $G(\cdot)$ tries to capture data distribution of training set. It takes a uniform noise vector \mathbf{z} as an input and outputs a vector $\tilde{\mathbf{x}} = G(\mathbf{z})$. In an adversarial way, the discriminator $D(\cdot)$ tries to identify samples produced by the generator from real ones. That is to say, $G(\cdot)$ and $D(\cdot)$ play the following two-player minimax game:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (9)$$

3) *Lyrics-conditional GAN*: Conditional GAN is proposed in [27], with the goal of instructing the generation process by conditioning the model with additional information, which motivates us to train a generative model for lyrics-conditioned melody generation. In this work, the generator and discriminator are conditioned on lyrics. The lyrics are encoded to a sequence of 20-dimensional embedding vectors.

The melody contains a sequence of music attributes $\mathbf{x}^{(i)}$, representing MIDI note, duration, and rest. Therefore, in the context of lyrics-conditioned melody generation, the input of the generator is the paired sequences of syllable embedding vectors $\mathbf{y}^{(i)}$ and uniform random vector $\mathbf{z}^{(i)}$ in $[0, 1]^k$, where $k = 30$. The generated music attributes, syllable embedding vectors $\mathbf{y}^{(i)}$, and real music attributes $\mathbf{x}^{(i)}$, are fed to the discriminator. $G(\mathbf{z}^{(i)}|\mathbf{y}^{(i)})$ is a sequence of triplets containing attributes $\hat{\mathbf{x}}^i = \{\hat{x}_{MIDI}^i, \hat{x}_{dur}^i, \hat{x}_{rest}^i\}$. Both the generator and the discriminator contain LSTM cells. The loss functions in the following are implemented to jointly train the generator and discriminator, where m is mini batch size.

$$L_G = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)}|\mathbf{y}^{(i)}))) \quad (10)$$

$$L_D = \frac{1}{m} \sum_{i=1}^m [-\log D(\mathbf{x}^{(i)}|\mathbf{y}^{(i)}) - \log(1 - D(G(\mathbf{z}^{(i)}|\mathbf{y}^{(i)})))] \quad (11)$$

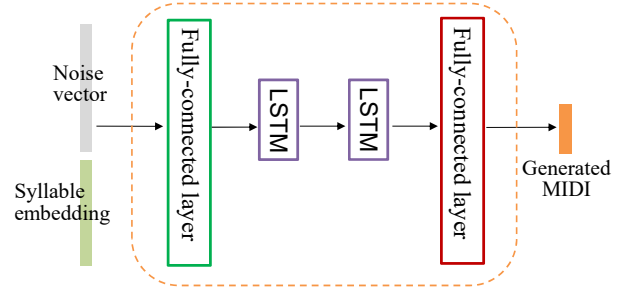


Fig. 3: Generator network for one MIDI note generation, conditioned with an encoded syllable $\mathbf{y} \in \mathbb{R}^{20}$, with an input random noise vector $\mathbf{z} \in \mathbb{R}^{30}$, and output MIDI attributes $\hat{\mathbf{x}} \in \mathbb{R}^3$.

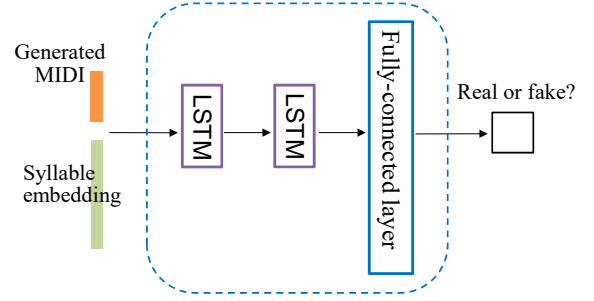


Fig. 4: Discriminator network for one MIDI note, conditioned with an encoded syllable $\mathbf{y} \in \mathbb{R}^{20}$, with the generated MIDI attributes $\hat{\mathbf{x}} \in \mathbb{R}^3$, and output the decision of real or fake.

4) *Generator network in Fig.3*: The generator is to learn the distribution of real samples, which is trained to increase the error rate of the discriminator. In this work, each melody sequence has 20 notes, which needs 20 LSTM cells to learn the sequential alignment between lyrics and melody. The first layer in the generator network uses ReLU (rectified linear unit). When given a 50-dimensional vector concatenated by an encoded syllable $\mathbf{y} \in \mathbb{R}^{20}$ and an input random noise vector $\mathbf{z} \in \mathbb{R}^{30}$, the output dimension of the first layer is scaled to 400-dimensional vector to fit the number of internal hidden units of the LSTM cells of the generator. We tried different amounts of LSTM layers and found that 2 layers are sufficient in learning alignment relationship between lyrics and melody. Then, the fourth linear layer produces a triplet of music attributes $\hat{\mathbf{x}}^i = \{\hat{x}_{MIDI}^i, \hat{x}_{dur}^i, \hat{x}_{rest}^i\}$.

The same LSTM cell is used for each syllable in the lyrics. The output with the triplet music attributes of previous LSTM cell is concatenated with current 20-dimensional syllable embedding, which are further fed to current LSTM cell. This procedure is repeated until the generator can succeed to fool the discriminator.

5) *Discriminator network in Fig.4*: The discriminator is to distinguish generated melody from real samples, which is trained by estimating the probability that a sample is from real training dataset rather than the generator. Since lyrics as context information are used as condition in the discriminator,

the generated triplet of music attributes concatenated with syllable embedding together as 23-dimensional vector are input to the first LSTM cell in the discriminator. The hidden size of the LSTM cell in the discriminator is also 400. The output 400-dimensional vector from the second LSTM layer is input to the third linear layer, followed by a sigmoid activation function which estimates the decision output of real or fake by a value in the range $[0, 1]$. With conditioning lyrics, the discriminator and generator are simultaneously learned until the training converges.

D. Tuning scheme

The continuous-valued sequence is output from the generator, which needs to be constrained to the underlying musical representation of discrete-valued MIDI attributes. Quantization of music attributes (MIDI number, note duration, and rest duration) are done during the generation in the experiments of validation and testing. Music attributes are constrained to their closest discrete values. The quantized music attributes are estimated to see if each generated sequence has a perfect scale consistency of melody. In particular, the most likely standard scales of music attributes this sequence belongs to is generated from syllable embedding in validation and testing datasets. The remaining out-of-tune notes are mapped to their closest in-tune music attributes.

V. LYRICS-MELODY DATA ACQUIREMENT

There is no aligned lyrics-melody music dataset publicly available for music generation. In this work, a large-scale music dataset with sequential alignment between lyrics and melody is created to investigate the feasibility of this research with deep conditional LSTM-GAN. We acquire lyrics-melody paired data from two sources based on considering melodies with enough English lyrics, where 7,998 MIDI files come from “LMD-full” MIDI Dataset [28] and 4,199 MIDI files come from reddit MIDI dataset [29]. Altogether there are 12,197 MIDI files in the dataset, which contain 789.34 hours of melodies. The average length of melodies is 3.88 minutes. This dataset is available on Github ¹.

1) *Data selection*: In our experiment, 20,934 unique syllables and 20,268 unique words from “LMD-full” MIDI and reddit MIDI dataset are used for training a skip-gram model to extract embedding vectors of lyrics. As for training the LSTM-GAN model, only paired lyrics-melody sequences in “LMD-full” dataset are used. In particular, if a MIDI file has more than 20 notes but less than 40 notes, one 20-note sequence is taken as our data sample; if a MIDI file has more than 40 notes, two 20-note sequences are taken as our data samples. Accordingly, 13,937 sequences each with 20 notes and 278,740 syllable-note pairs are acquired, which are used for training, validation, and testing datasets.

2) *Parsing MIDI file*: Triplets of music attributes with {MIDI Number, note duration, rest duration} are obtained by parsing each MIDI file from the LMD-full dataset which contains English lyrics. The parsing is made as follows:

- The beats-per-minute (BPM) value for each MIDI file is extracted.
- If a note has a corresponding English syllable, its MIDI Number is extracted. This value is taken as the first of music attributes in our melody representation.
- If a note has a corresponding English syllable, its note-on and note-off values are stored.
- From the note-on and note-off values, the note duration and rest duration attributes are calculated, using the formula

$$x = \phi\left(t \times \frac{\text{BPM}}{60}\right) \quad (12)$$

where x is an attribute of note k (either note duration or rest duration), t is a time in seconds ($t = \text{note-off}_k - \text{note-on}_k$ to calculate note duration and $t = \text{note-on}_k - \text{note-on}_{k-1}$ to calculate rest duration) and $\phi(\cdot)$ is an operator which constrains the value to the closest value in the set of values used to represent the corresponding attribute.

3) *Note duration*: A note duration means the length of time that a note is played. The association relationship between note and note duration used in this work is shown in Table II.

0.25	0.5	0.75	1	1.5	2	3	4	6	8	16	32

TABLE II: Relationship between note duration and note.

4) *Rest*: A rest means how long the silence in a piece of melody will last. The rest values and corresponding rest symbols are shown in Table III.

0	1	2	4	8	16	32
No rest						

TABLE III: Relationship between rest values and corresponding symbols.

5) *Distribution of music attributes*: The distribution of each attribute in our music dataset is respectively shown in Fig. 5, which indicates that most MIDI note numbers range from 60 to 80, quarter note is most frequently played, and *rest* = 0 appeared in most cases of melodies.

VI. EVALUATION

In this section, experimental setup, validation method, and experimental results are introduced to investigate the feasibility of our proposed conditional LSTM-GAN model.

A. Experimental setup

The entire dataset is split with a 0.8/0.1/0.1 proportion between training, validation and testing sets. The model is trained using mini-batch gradient descent for a total of 400

¹<https://github.com/yy11lab/Lyrics-Conditioned-Neural-Melody-Generation>

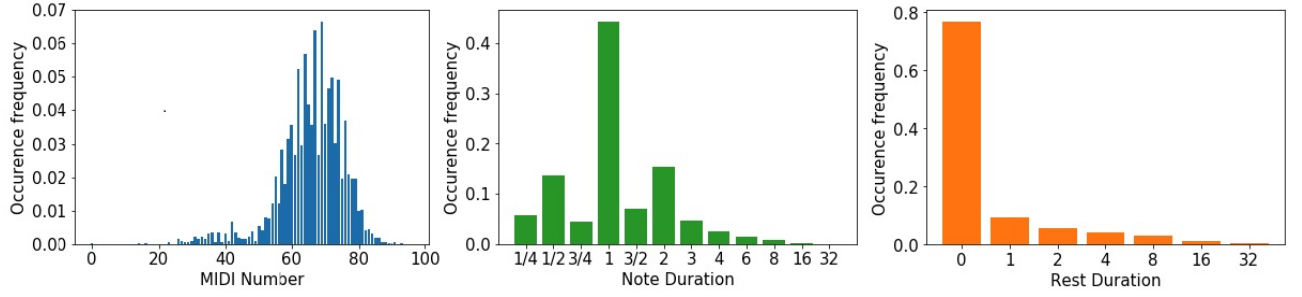


Fig. 5: Dataset distribution of music attributes

epochs. The learning rate starts at a value of 0.1, and gradually decreases.

During both validation and testing stages, the sequences of triplet continuous-valued attributes are first constrained to their closest discrete value. The candidate values for the MIDI Numbers is in the range $\{21, \dots, 108\}$. In addition, the quantized sequence is checked to see if it belongs to most likely scale, where the MIDI number of the out-of-tune notes is changed to the closest MIDI number which belongs to the most likely scale.

B. Validation using MMD

The validation is made using a Maximum Mean Discrepancy (MMD) [30] unbiased estimator. Given two sets of samples, MMD^2 takes a value between 0 and 1, indicating how likely the two sets of samples are coming from the same distribution (a value of 0 indicates that the two sets are sampled from the same distribution). Therefore, at the end of each training epoch, the MMD between the generated sequences and the validation sequences is calculated. The weights and biases from the model corresponding to the lowest MMD value are selected.

Let $X_m := \{x_1, \dots, x_m\}$ and $Y_n := \{y_1, \dots, y_n\}$ be two sets of independently and identically distributed (i.i.d) samples from P_x and P_y respectively. Then, an unbiased empirical estimate of MMD^2 is given by [31]:

$$\begin{aligned} \text{MMD}_u^2(\mathcal{F}, X_m, Y_n) = & \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) \\ & + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(y_i, y_j) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j). \end{aligned} \quad (13)$$

where \mathcal{F} is a reproducing kernel Hilbert space (RKHS), with the kernel function $k(x, x') := \langle \phi(x), \phi(x') \rangle$, and continuous feature mapping $\phi(x) \in \mathcal{F}$ for each x . We used $k(x, x') = \exp(-\|x - x'\|^2 / (2\sigma^2))$ as kernel function, with kernel bandwidth σ set such that $\|x - y\| / (2\sigma^2)$ equals 1 when the distance between x and y is the mean distance between points from both datasets X_m and Y_n [32].

C. Baseline model

The baseline model for the following experiments is inspired by [33]. Melodies of 20 notes are created by randomly sampling the testing set based on the dataset distribution for music attributes (i.e. the distribution shown in Fig. 5). Sequences generated by the baseline model are also judged to see if out-of-tune MIDI number needs to be changed. The MIDI numbers are restricted in the set $\{60, \dots, 80\}$, meaning that if a MIDI number lower than 60 is sampled from the MIDI numbers distribution, then it takes a value of 60, and similarly a MIDI number higher than 80 is set to 80.

D. Training stage analysis

The MIDI note number, spectrum, and relative amplitude of generated songs are investigated at 1, 51, and 351 epochs as shown in Fig. 6. The corresponding sheet score with alignment between lyrics and melodies at 1, 51, and 351 epochs are shown in Fig. 7. It is obvious that the generated melody gets better when the learning goes deep by increasing the number of epochs. Additionally, We ask volunteers to listen to the generated songs at different epochs. This also confirms the effectiveness of our deep conditional LSTM-GAN.

E. Music quantitative evaluation

Some quantitative measurements are designed to compare the melodies generated by both our proposed model and the baseline, which are shown in the following:

- MIDI numbers span: the difference between the highest MIDI number of the sequence and the lowest one.
- 3-MIDI number repetitions: a count of how many MIDI numbers 3-grams repetitions occur throughout the sequence.
- 2-MIDI number repetitions: a count of how many MIDI numbers 2-grams repetitions occur throughout the sequence.
- Number of unique MIDI numbers: a count of how many different MIDI numbers are present in the sequence.
- Number of notes without rest: a count of how many rest duration have a value of 0 throughout the sequence.
- Average rest value within a song: an averaged value of the rest duration attribute.
- Song length: the sum of all the note duration attributes and all the rest duration attributes of the sequence.

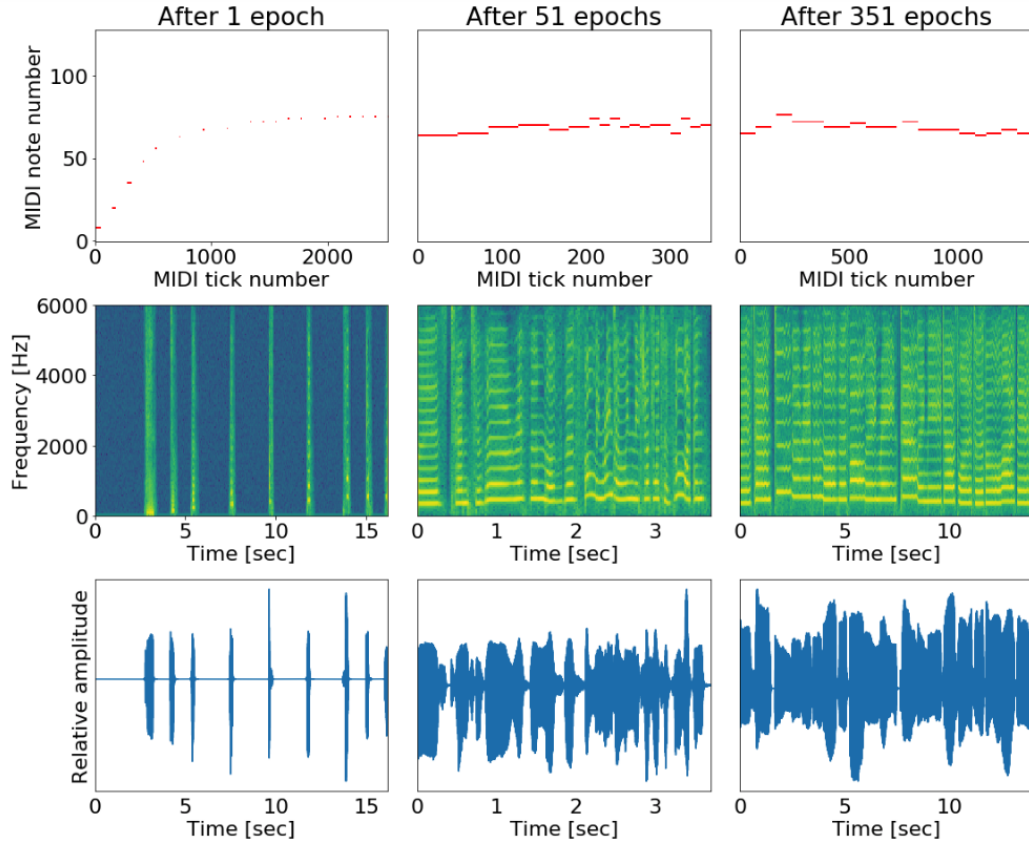


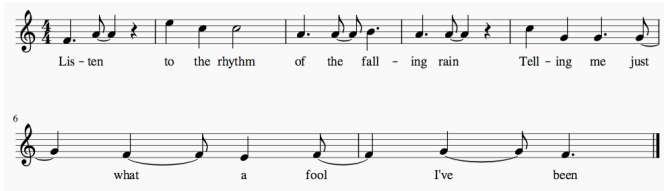
Fig. 6: Generated songs by generators trained for 1, 51 and 351 epochs respectively.



(a) Model trained for 1 epoch



(b) Model trained for 51 epochs



(c) Model trained for 351 epochs

Fig. 7: Different sheet music trained for 1, 51 and 351 epochs respectively.

Figure 8 demonstrates the evolution of each of these values averaged over 1,394 generated sequences (one sequence per testing set lyrics).

For pitch-related attributes, the proposed model outperforms the baseline in every aspect. The model tends to converge to a value which is relatively close to the ground truth (i.e. the average value from the dataset). However, the values of 2-MIDI numbers and 3-MIDI numbers repetitions converge to a value which is significantly lower to the corresponding measurement over the dataset.

For metrics which are related to temporal attributes (i.e. note duration and rest duration), the baseline is closer to the ground truth value. This is expected, since these metrics are nothing but an average of attributes which the baseline samples from the ground truth distribution. Hence, these values tend to the ground truth value as the number of generated baseline examples increases. Table IV shows the numerical values of the results.

	Ground truth	Cond-LSTM-GAN	Baseline
MIDI Numbers span	10.7	8.2	18.3
3-MIDI numbers repetitions	5.2	2.1	0.2
2-MIDI numbers repetitions	12.7	9.2	2.04
Number of unique MIDI numbers	6.0	5.3	9.2
Number of notes without rest	15.4	16.2	15.4
Average rest value within song	0.9	0.7	0.9
Song length	45.3	40.1	45.2

TABLE IV: Metrics evaluation of in-songs attributes

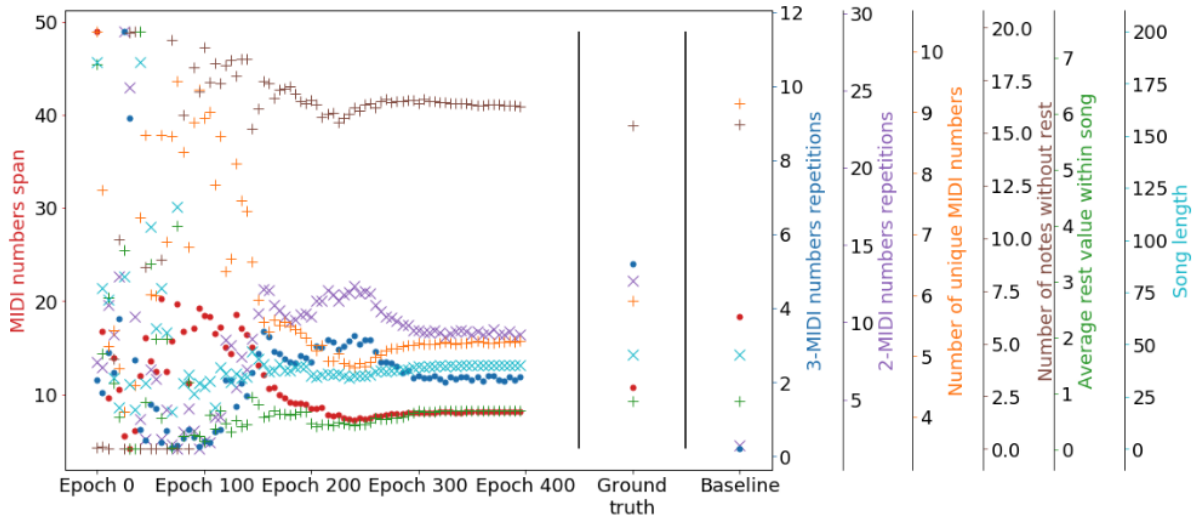


Fig. 8: Music-related measurements

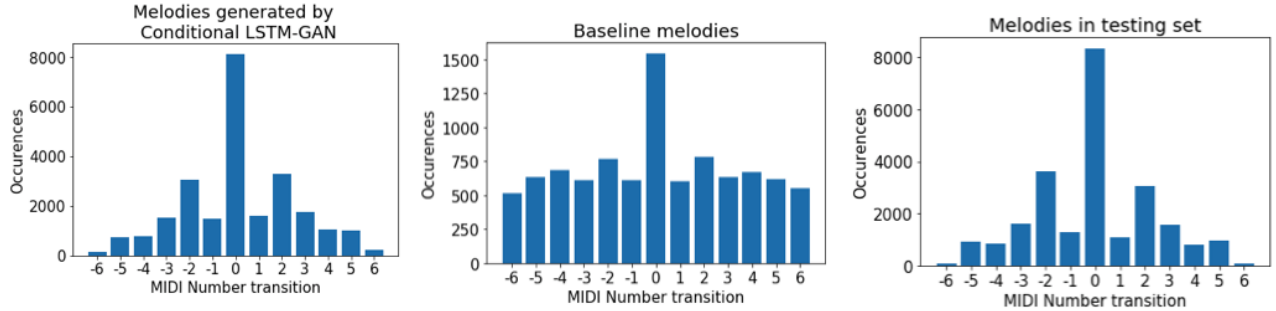


Fig. 9: Distribution of transitions

Another important attribute in music is the distribution of the transitions between MIDI numbers. Figure 9 shows the distributions of the transitions for the melodies generated by our model, the baseline model and the testing set. The melody generated by the proposed model approximates well the distribution of the human composed music. Although the testing set melody has a slightly higher transition from a note to a lower-pitched one, the melody generated by our model has more transitions from a note to a higher-pitched one.

F. Effect of lyrics conditioning

1) *MIDI numbers*: An example of the influence of lyrics on the generated MIDI numbers is shown in Fig. 10. Given two kinds of lyrics, 1,000 songs are generated for each lyrics, and the distribution of the generated MIDI numbers is estimated. In our example, the second lyrics leads to songs with lower MIDI numbers than the first one, which helps to deliver some semantic information latent in the lyrics.

2) *Note duration and rest duration*: In this experiment, an evaluation method to seeing if the lyrics conditioning has an effect on the generated note duration and rest duration is presented. The following focuses on note duration attribute, but the same is valid for rest duration as well.

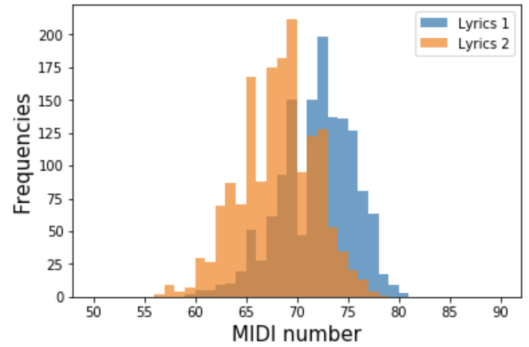


Fig. 10: Estimated distributions of music generated by two different lyrics. Lyrics 1 are taken from the song *It Must Have Been Love (Christmas for the Broken Hearted)* by Roxette (written by Per Gessle). Lyrics 2 are taken from *Love's Divine* by Seal (written by Mark Batson and Sealhenry Samuel).

Let $D \in \mathbb{R}^{N \times M}$ be a matrix composed of N note duration (or note rest) sequences of M note durations, where each sequence is taken from a different song from a N songs dataset used to train a generator G . Therefore, $D_{i,j}$ is the

j -th note duration of the i -th song. Let $G \in \mathbb{R}^{N \times M}$ be a matrix composed of note duration sequences generated by G by feeding the syllables sequences corresponding to each row of D . This means that $G_{i,j}$ is the j -th note duration of the sequence generated by G when the syllables corresponding to \mathbf{d}_i are fed to it, where \mathbf{d}_i denotes the i -th row of D .

Let $\text{randrow}(\cdot)$ be an operator which randomizes the order of the row of a matrix. Therefore $D_{rs} = \text{randrow}(D)$ can be seen as a matrix with correct in-sequence note duration order, but wrong song order when compared to D . $D_{rn} = (\text{randrow}(D^T))^T$ can be seen as a matrix for which the song order is the same as D 's, but the note duration sequences are randomized. Finally, $D_{rms} = \text{randrow}(D_{rn})$ can be seen as a matrix in which both the song and note order are randomized when compared to D . The subscripts rs , rn and rms denote "random songs", "random notes", and "random notes + songs" respectively. Since $\text{randrow}(\cdot)$ is a random operator, D_{rs} , D_{rn} , D_{rms} are matrices of random variables (random matrices).

In this experiment, $d = \frac{1}{NM} \|D - G\|_F$ (which is a real value) is compared to the distribution of the random variables $d_{rs} = \frac{1}{NM} \|D_{rs} - G\|_F$, $d_{rn} = \frac{1}{NM} \|D_{rn} - G\|_F$ and $d_{rms} = \frac{1}{NM} \|D_{rms} - G\|_F$, with $N = 1,394$ (number of songs in testing set) and $M = 20$. The experiment is made on the testing set.

Results are shown in Fig. 11 (note duration), and Fig. 12 (rest duration). The three distributions are estimated using 10,000 samples for each random variable. In each case, d is statistically lower than the mean value, indicating G learned useful correlation between syllable embeddings and note/rest durations.

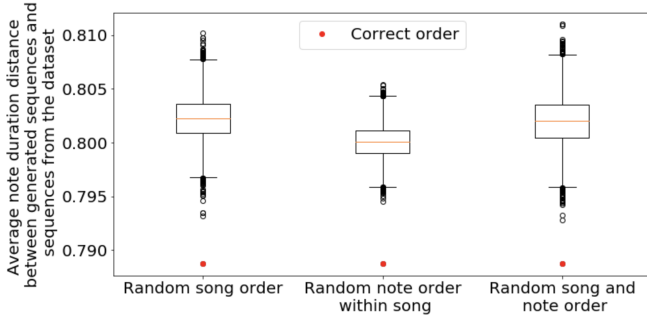


Fig. 11: Boxplots of the distributions of d_{rs} , d_{rn} and d_{rms} (For the note duration attribute). $d = 0.788$ is highlighted in red in each boxplot. Mean values are $\mu_{rs} = 0.802$, $\mu_{rn} = 0.801$ and $\mu_{rms} = 0.802$ respectively.

G. Subjective evaluation

4 different lyrics are randomly selected from ground truth dataset. Accordingly, 12 melodies are obtained by using baseline method, our model, and ground truth. All melodies are sung by a female voice produced by synthesizer V [34]. 4 male and 3 female subjects without knowing our research and any musical knowledge were invited to listen to these 12 melodies, where each melody with around 20 seconds is played 2 times in a random order. Three questions as metrics

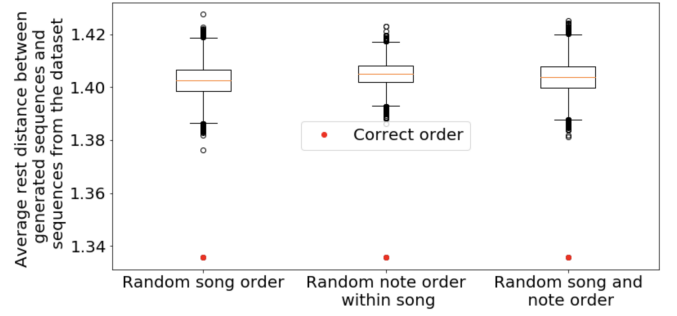


Fig. 12: Boxplots of the distributions of d_{rs} , d_{rn} and d_{rms} (for the rest duration attribute). $d = 1.336$ is highlighted in red in each boxplot. Mean values are $\mu_{rs} = 1.404$, $\mu_{rn} = 1.407$ and $\mu_{rms} = 1.404$ respectively.

are used for evaluation: how about the entire melody? how about the rhythm? and does the melody fit the lyrics well? The subjects are asked to give a score from 1 to 5 (where 1 corresponds to "very bad", 2 to "bad", 3 to "OK", 4 to "good" and 5 to "very good").

The first run is taken to enable subjects to get used to the type of melodies they were listening. The scores of evaluation metrics are respectively averaged based on listening results of baseline, our model, and ground truth on the second run.

Evaluation results are shown in Figure 13. It is obvious that the melodies generated by the proposed model are closer to the ones composed by humans than the baseline in each metric. The feedback from subjects indicates that relatively low scores of melody evaluation are generated which might be due to the limited capability of the synthesizer for high pitches. From these results of all three metrics, we also can find there still are the gaps between melodies generated by our model and ones from human composition, which tells us there is much space we can investigate to improve capability of neural melody generation.

VII. CONCLUSION AND FUTURE WORK

Melody generation from lyrics in music and AI is still unexplored well. Making use of deep learning techniques for melody generation is a very interesting research area, with the aim of understanding music creative activities of human. Several contributions are done in this work: i) the largest paired English lyrics-melody dataset is built to facilitate the learning of alignment relationship between lyrics and melody. This dataset is very useful for the area of melody generation. ii) a skip-gram model is trained to exact lyrics embedding vectors, which can be taken as a lyrics2vec model for English lyrics feature extraction. iii) Conditional LSTM-GAN is proposed to model sequential alignment relationship between lyrics and melody, followed by a tuning scheme that has the capability of constraining a continuous-valued sequence to the closest in-tune discrete-valued sequence. iv) Evaluation method of melody generation is suggested to demonstrate the effectiveness of our proposed deep generative model.

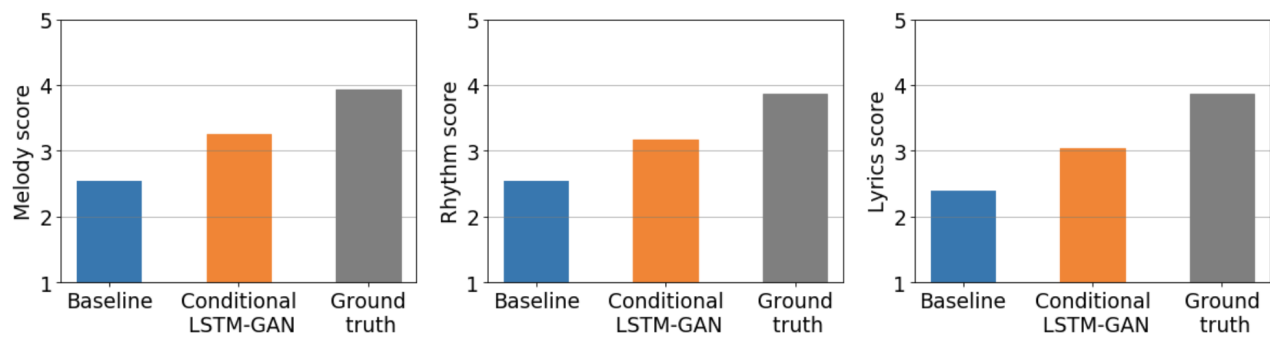


Fig. 13: Subjective evaluation results

Several interesting future works will be investigated as follows: i) how to compose melody with the sketch of uncomplete lyrics. ii) how to compose polyphonic melody with lyrics. iii) how to predict lyrics when given melody as a condition.

REFERENCES

- [1] G. A. Wiggins, "A preliminary framework for description, analysis and comparison of creative systems," *Journal of Knowledge Based Systems*, vol. 19, no. 7, pp. 449–458, 2006.
- [2] L. A. Hiller and L. M. Isaacson, "Musical composition with a high-speed digital computer," *Journal of the Audio Engineering Society*, vol. 6, no. 3, pp. 154–160, 1958.
- [3] D. Ponsford, G. Wiggins, and C. Mellish, "Statistical learning of harmonic movement," *Journal of New Music Research*, vol. 28, no. 2, pp. 150–177, 1999.
- [4] J. Briot and F. Pachet, "Music generation by deep learning - challenges and directions," *CoRR*, vol. abs/1712.04371, 2017. [Online]. Available: <http://arxiv.org/abs/1712.04371>
- [5] Y. Yu, S. Tang, F. Raposo, and L. Chen, "Deep cross-modal correlation learning for audio and lyrics in music retrieval," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 15, no. 1, pp. 1–15, 2019.
- [6] <https://en.wikipedia.org/wiki/Melody>.
- [7] M. Scirea, G. A. B. Barros, N. Shaker, and J. Togelius, "Smug: Scientific music generator," in *Sixth International Conference on Computational Creativity*, 2015, pp. 204–211.
- [8] M. Ackerman and D. Loker, "Algorithmic songwriting with ALYSIA," *CoRR*, vol. abs/1612.01058, 2016. [Online]. Available: <http://arxiv.org/abs/1612.01058>
- [9] H. Bao, S. Huang, F. Wei, L. Cui, Y. Wu, C. Tan, S. Piao, and M. Zhou, "Neural melody composition from lyrics," *CoRR*, vol. abs/1809.04318, 2018. [Online]. Available: <http://arxiv.org/abs/1809.04318>
- [10] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Networks," *arXiv e-prints*, 2014.
- [11] J. D. F. Rodriguez and F. J. Vico, "AI methods in algorithmic composition: A comprehensive survey," *CoRR*, vol. abs/1402.0585, 2014.
- [12] T. Anders and E. R. Miranda, "Constraint programming systems for modeling music theories and composition," *ACM Comput. Surv.*
- [13] M. Delgado, W. Fajardo, and M. Molina-Solana, "Inmamusys: Intelligent multiagent music system," *Expert Syst. Appl.*
- [14] D. Conklin, "Music generation from statistical models," in *SAISB Symposium on Artificial Intelligence and Creativity in the Arts and Sciences*, 2003, pp. 30–35.
- [15] A. Eigenfeldt and P. Pasquier, "Realtime generation of harmonic progressions using controlled markov selection," in *International Conference on Computational Creativity*, 2010, pp. 16–25.
- [16] D. Cope, *Computer Models of Musical Creativity*. The MIT Press, 2005.
- [17] J. Wu, C. Hu, Y. Wang, X. Hu, and J. Zhu, "A hierarchical recurrent neural network for symbolic melody generation," *CoRR*, vol. abs/1712.05274, 2017.
- [18] D. D. Johnson, "Generating polyphonic music using tied parallel networks," in *International Conference on Evolutionary and Biologically Inspired Music and Art*, 2017, pp. 128–143.
- [19] O. Mogren, "C-RNN-GAN: continuous recurrent neural networks with adversarial training," *CoRR*, vol. abs/1611.09904, 2016.
- [20] S. Fukayama, K. Nakatsuma, S. Sako, T. Nishimoto, and S. Sagayama, "Automatic song composition from the lyrics exploiting prosody of the Japanese language," in *International Conference of Sound and Music Computing*, 2010, pp. 299–302.
- [21] K. Monteith, T. R. Martinez, and D. Ventura, "Automatic generation of melodic accompaniments for lyrics," in *Proceedings of the Third International Conference on Computational Creativity*, 2012, pp. 87–94.
- [22] <http://www.musiccrashcourses.com/lessons/pitch.html>.
- [23] [Online]. Available: [https://en.wikipedia.org/wiki/Duration_\(music\)](https://en.wikipedia.org/wiki/Duration_(music))
- [24] [Online]. Available: [https://en.wikipedia.org/wiki/Rest_\(music\)](https://en.wikipedia.org/wiki/Rest_(music))
- [25] <https://en.wikipedia.org/wiki/Syllable>.
- [26] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [27] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *CoRR*, vol. abs/1411.1784, 2014. [Online]. Available: <http://arxiv.org/abs/1411.1784>
- [28] <https://colinraffel.com/projects/lmd/>.
- [29] <https://www.reddit.com/r/datasets/>.
- [30] A. Smola, A. Gretton, L. Song, and B. Schölkopf, "A hilbert space embedding for distributions," in *Algorithmic Learning Theory*, M. Hutter, R. A. Servedio, and E. Takimoto, Eds. Springer Berlin Heidelberg, 2007.
- [31] W. Bounliphone, E. Belilovsky, M. B. Blaschko, I. Antonoglou, and A. Gretton, "A Test of Relative Similarity For Model Selection in Generative Models," *arXiv e-prints*, 2015.
- [32] S. J. Reddi, A. Ramdas, B. Poczos, A. Singh, and L. Wasserman, 2014.
- [33] H.-P. Lee, J.-S. Fang, and W.-Y. Ma, "iComposer: An automatic songwriting system for Chinese popular music," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*.
- [34] <https://synthesizerv.com/en/>.