# Text To Song Generative System

# Project Vivy

**Zhen Zhang**  **Coauthor**
Department of Electric and Computer Engineering  Affiliation
Virginia Polytechnic Institute and State University  Address
Blacksburg, VA 24061
zhenz@vt.edu  email

**Luis Vazquez Morales**  **Coauthor**
Department of Electric and Computer Engineering  Affiliation
Virginia Polytechnic Institute and State University  Address
Blacksburg, VA 24061
lvazquez@vt.edu  email

## Abstract

This paper proposes a novel text-to-song system. The system can generate a song that is entirely created, arranged, and sung by AI, based on user input such as lyrics and style. This research project focuses on the synthesis of singing voices, the text-to-singing-voice aspect, and develops a deep learning model named Vivy for this task. Vivy successfully synthesized the corresponding files as anticipated, and its generated results met experiment expectations. The system proposed in the paper has been preliminarily validated, making the design of text-to-song systems feasible and further advancing the application and development of AI in the field of art.

## 1  Introduction

In recent years, artificial intelligence has made significant strides in various domains, such as chatbot, self-driving, image generation, and music generation. However, in the field of art, especially in music, there is no generative model or multi-model system in research that can produce complete music from text. To accelerate the application of AI in the field of music, and to attempt to use AI to generate a complete piece of music, a text-to-song generative system is proposed in this paper.

In the project, the aim was to develop a complete text-to-song system where users can input lyrics, song style, BPM, and other information to obtain a fully AI-generated, complete piece of music with synthesized singing voice. This is to further accelerate the application of AI in the field of art. This novel project has not been implemented so far, thus was one of the reasons the project was chosen. Other reasons were to allow people who were previously unable to create music due to any lack of knowledge or skills can make the music they now desire. This would help accelerate new music and possible genres in the music industry due to the ease of making of music.

Due to the substantial complexity of this project and time constraints during the development process in this project, only the voice generation part of the project was able to be built. The rhythm generation and the creation of complete music will be implemented in the future.

## 2    Related Work

There has been no research paper so far that has tried to implement this project's proposed idea of generating a singing voice through a lyric input. Though there are projects some projects that have generated a singing voice, most significant being one is called Diffsinger that used Diffsinger (.ds) file that contains phoneme sequence, phoneme durations, and music score parameters. Generating the .ds files to create the singing voice with the lyrics that the user wants is complex, which is something that attempted to simplify.

## 3    Methodologies

This project involves training a text to singing model, named Vivy. This section will introduce the overall structure of the project, as well as the approach used to train Vivy.

### 3.1    Project Structure

The project system consists of two parts: one is the voice generation system, and the other is rhythm generation system. Users will input information such as lyrics and style, which will then be fed into different systems for processing.

Multiple models are utilized to perform each task. For the lyric processing part, GPT, a Natural Language Processing (NLP) model, is chosen for processing. In the text-to-singing part, Vivy is developed for generation, and its output is sent to Diffsinger, a voice synthesis model based on shallow diffusion, for voice synthesis. For the rhythm synthesis part, the plan is to use musicGen, a rhythm generation model from audioCraft. The results will be sent to a note extraction model to extract all notes for Vivy's training. Simultaneously, this rhythm will also be combined with the generated voice to compose a complete song.
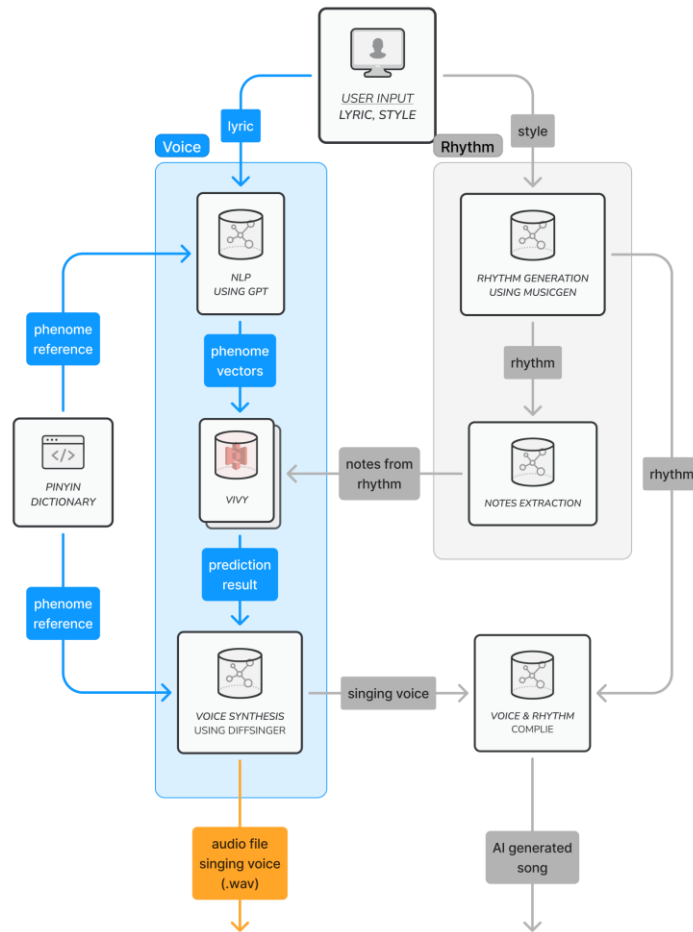
Fig 1. Project structure

### 3.1.1  Voice Generation

The user-input lyrics are processed through a NLP model. In this system, GPT is used for this purpose. Since the voice synthesis model employed is Diffsinger, A pre-trained model provided by them for voice synthesis. This pre-trained model is based on pinyin phonemes, hence the input for Diffsinger should be based on pinyin phonemes. Therefore, after feeding the Pinyin dictionary that is made by Opencpop into GPT, it first converts English into pinyin, then uses the pinyin-based English text to make predictions through Vivy. The results of these predictions are then further passed to Diffsinger for voice synthesis.

### 3.1.2  Rhythm Generation

The style information input by the user is sent to musicGen for rhythm generation. Then, a note extraction model is developed to analyze the notes within the rhythm and send these notes to Vivy for sound generation. The rhythm generation part has not been developed due to time constraints and will be the subject of further research in the future.

### 3.2  Data Collection

A total of 9.7 hours of singing text, 89315 entries were acquired, most of which came from Opencpop, with the remaining part from sample texts in Diffsinger. All texts are in the form of text (.txt) or Diffsinger (.ds) files, containing song phonemes, duration of each phoneme, vocal notes, and f0 time sequences. In the .ds files, they include all the elements.

91 In the txt files, there are phonemes, phoneme duration, and notes.
92     To train the model more efficiently, all the redundant parts in the .ds files, such as
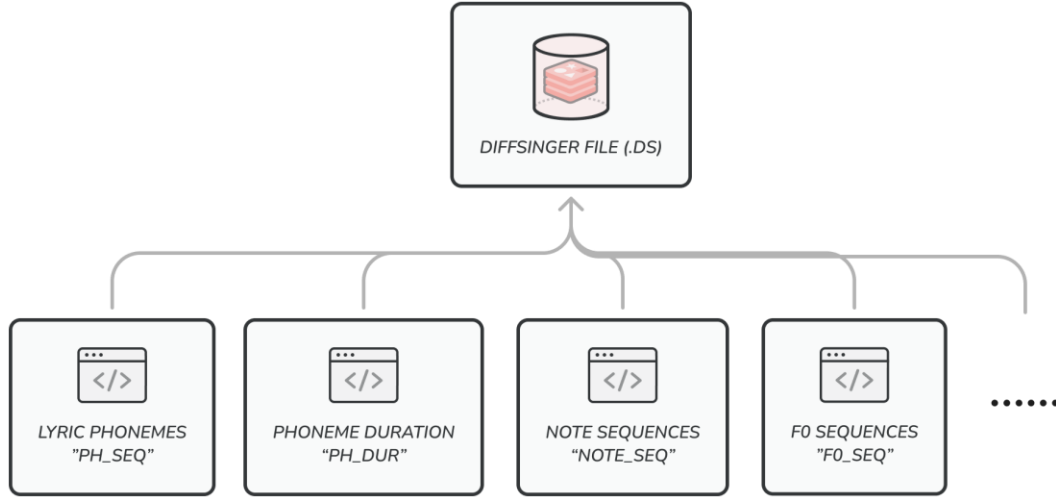93 note duration, energy, breath time, and other numerical values, have been removed.
94



95
96                         Fig 2. Diffsinger file structure
97
**3.3 Data Processing**

99     In this section, the process of handling all text data and normalizing
100 them for model training is explained. First, convert the text files into
101 Diffsinger files to facilitate batch processing later. Secondly, filter out all
102 irrelevant numerical values, retaining only those essential for sound
103 generation, such as ph_seq, which refers to all the phonemes of the lyrics;
104 ph_dur, indicating the duration of each phoneme; and note_seq, which
105 denotes the notes of each pronunciation.
106
**3.3.1   Text to Diffsinger File**

108     Since the training texts provided by OpenCpop primarily consist of text files and
109 raw audio files, and because these texts from OpenCpop make up over 80% of all training
110 samples, using their files for training can significantly improve the model's output. It is
111 essential to prioritize converting them into .ds files for subsequent batch processing. After
112 reading the text files, all texts are segmented and tagged, finally outputting them in the .ds
113 file format and incorporating them into the dataset.
114
**3.3.2   Data Filtering**

116     All .ds files are read, and any data unnecessary for audio generation is removed.
117 Additionally, this part also generates some statistical data to assist researchers in their analysis,
118 such as the total count of each value and the total duration of each phoneme duration.
119
**3.3.2   Normalization**

121     After importing the organized filtered dataset, the ph_seq and note_seq are first
122 converted into numerical values using integer encoding technology. Since notes in music
123 follow a regular pattern, it's important for the model to learn this pattern. Therefore, before
124 converting, a manually organized dictionary was prepared, in which all notes are sorted by
125 pitch from low to high. In this project, 54 notes ranging from "C#2/Db2" to "A#5" are used
126 for conversion. A special note, 'rest,' which represents a brief pause during singing, is assigned
127 the number 55.

128 Subsequently, after obtaining the encoded results, to ensure uniformity in the length
129 of all data, padding is applied to all data based on the length of the longest data in the set. This
130 padding is done using the mean value.

131
132 **3.4    Model Development**

133 Due to the time-sensitive nature of this project, several models adept at learning
134 sequences were selected for training. Two variants of Recurrent Neural Networks (RNN) were
135 used: Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU). Additionally, the
136 Transformer model was also employed for training. All models selected are simple standard
137 models, with a training epoch set to 1000.

138
139 **3.5    Testing and Evaluation**

140 To train the model, the training and validation data are randomly split in a 70/30 ratio,
141 using the validation data to check and analyze the results. Researchers hope the model will
142 exhibit creativity, so no specific expectations are set for the model's predicted results. However,
143 it is expected that the model should learn to pause at AP or SP, a type of rest symbol commonly
144 found in Diffsinger, typically occurring between sentences or at the end, thereby generating a
145 'rest'.

146
147 **4    Results**

148 The results of the project were a model called Vivy that generated prediction results from
149 phoneme vectors, a Pinyin Dictionary that converted English lyrics to Chinese phoneme, a
150 framework that connected the Diffsinger model with the previous results to form the input to singing
151 voice audio file, and three audio files generated by using three different model implementations.

152 The model implementation used were LSTM, GRU, and Transformer. Each model had
153 1000 epochs for the dataset to pass through the model. Each model was trained on the same computer,
154 which is equipped with a Ryzen 9 7000 series 16-core CPU, 64GB of RAM, and an NVIDIA RTX
155 4080 graphics card. The duration and size for each model are described in Table 1.

156

| Model Name | Training duration | Model size |
|------------|-------------------|------------|
| GRU | 44 min 40 s | 15.98 KB |
| LSTM | 47 min 33 s | 19.86 KB |
| Transformer | 1 hr 13 min 22 s | 23.47 KB |

157 Table 1. Prediction result

158

159 In this experiment, a line from Rick Astley's *Never Gonna Give You Up* was used: *Never*
160 *gonna give you up, never gonna let you down.* This lyric was tokenized by GPT, resulting in: 'AP n
161 ei f a g e n a j i f u y u a p u AP n ei f a g e n a l e y u d ao en AP'. The experiment aims to observe
162 if the model can generate a reasonable sequence of notes and durations for Diffsinger to synthesize
163 audio. Here, 'AP' represents a brief pause, and its corresponding note should be 'rest'

164 Since this study hopes that the audio generation model will be more creative, it did not set
165 too many expectations. However, the trained model should correctly identify that a rest note should
166 be output when an AP or SP symbol is encountered. After observing the outputs of three models,
167 only the GRU model successfully predicted a rest note in the aforementioned situations. Both LSTM
168 and Transformer models failed to make this prediction. Additionally, the duration of the GRU's
169 output is closer to that of normal singing.

170 Based on listening to the output audio files for each model, it was concluded that the audio

171 file from the GRU model gave the best result based on its tempo and flow throughout the audio file.
172 The audio files from the LSTM and Transformer had a slow tempo and felt somewhat to drag on.
173 These audio files were ranked based on the results shown in Table 2.

174

| Model Name | Duration | Is phoneme duration making sense? | Is AP note predict correctly? |
|---|---|---|---|
| GRU | 12s | Yes | Yes |
| LSTM | 16s | Yes | No |
| Transformer | 19s | No | No |

175 Table 2. Prediction result

176
## 177 4 Conclusion & Discussion

178 This project proposes a novel text-to-song generation system, providing a viable solution
179 for future developments in the field of art. By implementing a model that is able to generate
180 prediction results from phoneme vectors were able to produce a working project that allows for lyric
181 input to be used to produce a singing voice that allows normal people to generate a singing voice.
182 By testing different implementation of the model the results showed that a GRU implementation
183 produced the best results for a singing voice. Thus, a singing voice was able to be produced to show
184 proof of concept that the common person is able to produce songs that they were previously unable
185 to do so before.

186
### 187 4.1 Discussion

188 Due to time constraints on the project, the research team was only able to implement some
189 of their ideas, and it was not possible to deliver the complete project design within the limited time.
190 Additionally, due to the limited data resources collected, all training data was based on Chinese text,
191 making it more challenging to generate English lyrics. This issue could be addressed by manually
192 creating more English text data for more comprehensive training. To simplify the training process,
193 many important aspects were omitted, such as the duration of notes and the sequence of f0. If these
194 values are considered in further training, the results would be more pleasing to the ear.

195
### 196 4.2 Future Work

197 This project should continue to optimize the Vivy model to achieve better outputs, which
198 includes training with the vocal and rhythm parts of existing audio files. At the same time, it's
199 necessary to implement rhythm generation, extract the note sequence from the generated results,
200 and use it as input for Vivy to predict. This way, the rhythm and vocals can be more in sync.
201 Additionally, designing a user-friendly graphic user interface can further facilitate the general public
202 in using this system for creative purposes, thereby further promoting the application and
203 development of AI in the field of art.

204

205 Dataset, prediction examples and audio results can be found in project repository:
206 https://github.com/nikmomo/Song-Generation-Model-with-Vocal-Project-Vivy

207

## 208 References

209 [1] "Music generation," Papers With Code, https://paperswithcode.com/task/music-generation
210 (accessed Sep. 12, 2023).

211 [2] J.-P. Briot, G. Hadjeres, and F.-D. Pachet, "Deep learning techniques for Music Generation

212    -- A survey," arXiv.org, https://arxiv.org/abs/1709.01620 (accessed Sep. 12, 2023).

213    [3] J.-P. Briot and F. Pachet, "Deep learning for music generation: Challenges and directions
214    -       neural      computing      and      applications,"      SpringerLink,
215    https://link.springer.com/article/10.1007/s00521-018-3813-6 (accessed Sep. 12, 2023).

216    [4] J.-P. Briot, "From artificial neural networks to deep learning for music generation: History,
217    concepts    and    trends    -    neural    computing    and    applications,"    SpringerLink,
218    https://link.springer.com/article/10.1007/s00521-020-05399-0 (accessed Sep. 12, 2023).

219    [5] V. Kalingeri and S. Grandhe, "Music Generation with Deep Learning," arXiv.org,
220    https://arxiv.org/abs/1612.04928 (accessed Sep. 12, 2023).

221    [6] H. H. Mao, T. Shin and G. Cottrell, "DeepJ: Style-Specific Music Generation," 2018 IEEE
222    12th International Conference on Semantic Computing (ICSC), Laguna Hills, CA, USA, 2018,
223    pp. 377-382, doi: 10.1109/ICSC.2018.00077.

224    [7] Han, K., Xiao, A., Wu, E., Guo, J., Xu, C., & Wang, Y. (2021). Transformer in transformer.
225    Advances in Neural Information Processing Systems, 34, 15908-15919.

226    [8] Liu, J., Li, C., Ren, Y., Chen, F., & Zhao, Z. (2022, June). Diffsinger: Singing voice
227    synthesis via shallow diffusion mechanism. In Proceedings of the AAAI conference on
228    artificial intelligence (Vol. 36, No. 10, pp. 11020-11028).

229    [9] Copet, J., Kreuk, F., Gat, I., Remez, T., Kant, D., Synnaeve, G., ... & Défossez, A. (2023).
230    Simple and Controllable Music Generation. arXiv preprint arXiv:2306.05284.

231    [10] Wang, Y., Wang, X., Zhu, P., Wu, J., Li, H., Xue, H., ... & Bi, M. (2022). Opencpop: A
232    high-quality open source chinese popular song corpus for singing voice synthesis. arXiv
233    preprint arXiv:2201.07429.