

Homework-3

Symbolic Music Generation

HW3 TA : 艾芯 (Ivy Ai)

Office hour: Tue. 14:00 ~ 15:30 @BL505

Outline

- Rules
- Timeline
- Overview
- Detailed Explanation
- Submission
- Scoring

Rules

- Don't cheat
- Use transformer-based generation model
- Can use pretrained model, but not suggested
- Can use public codes with citation in report
- Don't use extra data

Timeline

- W8 – 10/24 (Thursday): Announcement of HW3
- W11 – 11/13 (Wednesday 23:59pm): Deadline
 - Late submission: 1 day (-20%), 2 days (-40%), after that (-60%)

Overview

1. Learn to manipulate MIDI file and represent symbolic music as tokens.
2. Learn to train a transformer-based model for symbolic music generation.

Detailed Explanation

- Symbolic music generation
- Dataset
- Evaluation
- Required Task
- Optional Task
- How to start
- Warning & Training Tips

Symbolic music generation

- MIDI file manipulation: represent symbolic music as tokens.
- Train a transformer-based model to generate symbolic music.

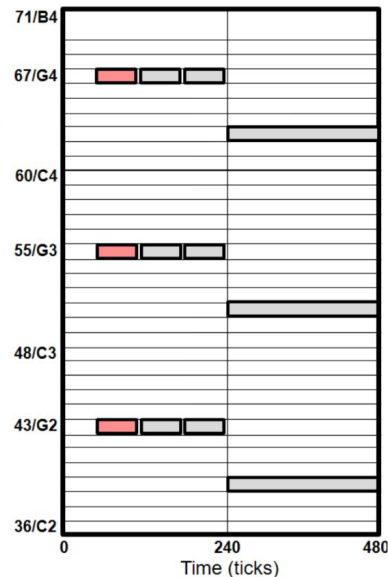
- Step (1): Task & Data
- Step (2): Token representation
- Step (3): Choose a model
- Step (4): Train the model till the loss is sufficiently low
- Step (5): Do inference
- Step (6): Listen to the generated music!
- Step (7): Do evaluation

MIDI file manipulation

By viewing the MIDI messages as “tokens”.

- A **text-like** representation of *.MIDI, using **MIDI messages** and **timestamps**
 - *MIDI note number* (0-127)
 - *Key velocity* (0-127): intensity of the sound
 - *MIDI channel* (different instruments)
 - *Time stamp*: how many *clock pulses or ticks* to wait before the command is executed

Time (Ticks)	Message	Channel	Note Number	Velocity
60	NOTE ON	1	67	100
0	NOTE ON	1	55	100
0	NOTE ON	2	43	100
55	NOTE OFF	1	67	0
0	NOTE OFF	1	55	0
0	NOTE OFF	2	43	0
5	NOTE ON	1	67	100
0	NOTE ON	1	55	100
0	NOTE ON	2	43	100
55	NOTE OFF	1	67	0
0	NOTE OFF	1	55	0
0	NOTE OFF	2	43	0
5	NOTE ON	1	67	100
0	NOTE ON	1	55	100
0	NOTE ON	2	43	100
55	NOTE OFF	1	67	0
0	NOTE OFF	1	55	0
0	NOTE OFF	2	43	0
5	NOTE ON	1	63	100
0	NOTE ON	2	51	100
0	NOTE ON	2	39	100
240	NOTE OFF	1	63	0
0	NOTE OFF	2	51	0
0	NOTE OFF	2	39	0



Music piano representation

There are many representation ways for music performance.

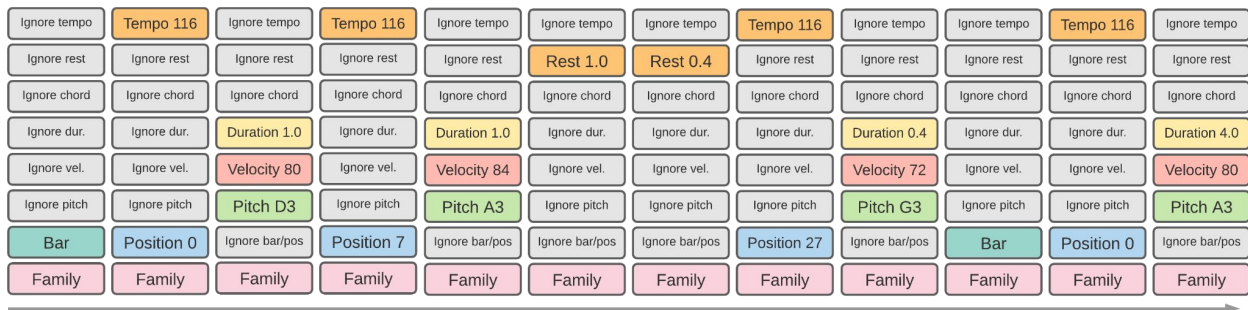
- MIDI-like
- REMI
- REMI+
- CP word
- etc...



MIDI-like



REMI

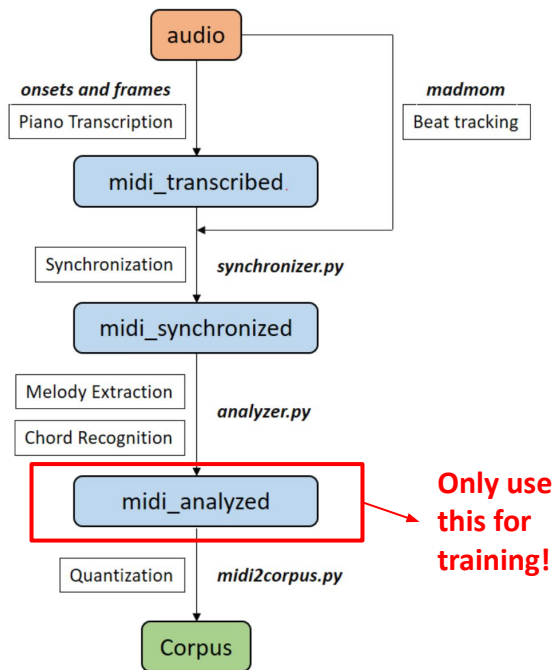


CP Word

Dataset: Pop1K7

- 1747 pop music piano performances (mid or midi file) transcribed from youtube audio.
 - Single track, 4 minutes in duration, totaling 108 hours.
 - 4/4 time signature (four beats per bar).
- You can use **whole dataset** for training, there is **no need to split** train/validation set for generation task.
- **Don't use other datasets** (for the evaluation score of your generation, the closer it is to the Pop1k7 dataset, the better)

data processing pipeline



Ref: Hsiao, Wen-Yi, et al. "Compound word transformer: Learning to compose full-song music over dynamic directed hypergraphs." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 35. No. 1. 2021.

Evaluation

[1] Huang, Yu-Siang, and Yi-Hsuan Yang. "Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions." *Proceedings of the 28th ACM international conference on multimedia*. 2020.

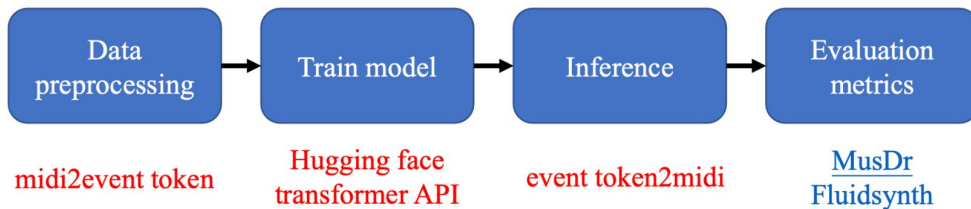
Objective [1] metrics for generation results

- **Pitch-Class Histogram Entropy (H)** : measures erraticity of pitch usage in shorter timescales (e.g., 1 or 4 bars). (**H4 is required**)
 - **Grooving Pattern Similarity (GS)** : measures consistency of rhythm across the entire piece. (**required**)
 - **Structureness Indicator (SI)** : detects presence of repeated structures within a specified range of timescale. (**optional**)
- Use [MusDr](#) package to calculate those metrics. (TA will provide sample code: *eval_metrics.py*)
- The **closer** the score of generation results compare to the original dataset, the better.

Task 1: Symbolic music generation

Train a transformer-based model **from scratch** to generate **32 bars** symbolic music.

- You may **randomly select** the prompt, treat it as a priming sequence for Transformer decoder to generate a continuation
 - For example: [Bar_None, Position_1/16, Tempo_class, Tempo_value]
- Model: You can use **any transformer-based** model.
- Either 1-stage generation or 2-stage generation is fine.



Task 1: Symbolic music generation

Report

1. Which token representation
2. Which transformer-based model
3. Implementation detail (model architecture, data augmentation, hyper-parameters...).
4. Implement at least **3 combinations** of your inference configurations. (it's quite influential!)
 - e.g. model loss, sampling method, top-k, temperature, etc. (see Lecture 7.)
5. For **each** combination, generate 20 mid/midi files (32 bars) to calculate their **average objective metrics results** (H4, GS).
6. Choose **one** combination, generate **20 mid/midi files (32 bars)** and convert into **wav files** as listening samples.

Task 1: Symbolic music generation

Example for objective result comparison

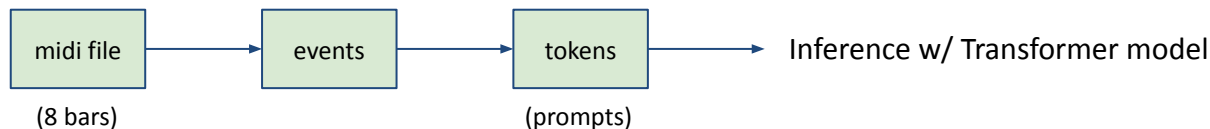
Model	representation	event	loss	top-k	temperature	H1	H4	GS	SI_short	SI_mid	SI_long
GPT2	REMI	w/ chord		5	1.2						
Transf ormer- XL	CP Word	w/o chord									
Real data											

Prompt files will be released
a few days later, focus on
Task 1 first!

Task 2: Symbolic music continuation

TA will provide 3 midi files (8 bars) as prompt, you need to generate their continuation for 24 bars. (Total: 8+24 bars)

- You can use checkpoint in Task 1, don't need to train another model.
- Inference pipeline:



- You can consider this task as a 「指定曲」 competition. Since everyone is given the same prompts, you can compare your generation results with others.

Task 2: Symbolic music continuation

Report

1. **Don't** need to report objective metrics results, it's not a matter in Task2.
2. For **each** prompt song, need to generate **1 mid/midi files (8+24 bars)** for **3 different inference configurations** and convert into **wav files** as listening samples.
3. Need to **specify the model, representation, inference configurations** you used for each continuation.
4. For **each** prompt song, need to **specify the one** you think that generates the best among different inference configurations.
5. We will choose some students to present their best results!

Optional Task

There are several optional tasks you can try. It's not required, but we look forward to see these in your report!

1. Different transformer model

- 1-stage generation, 2-stage generation, different model architecture, ...

2. Different token representation

- functional representation, chord extraction, ...

3. Difference tokenization strategy

- Byte Pair Encoding, Unigram, WordPiece, ...

4. Data augmentation

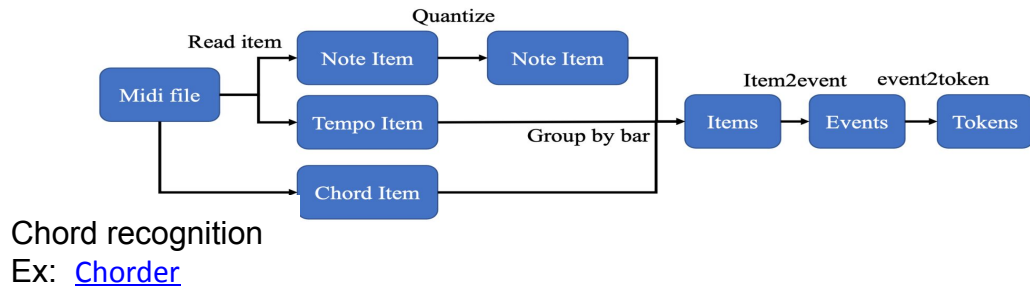
5. Conditional generation

- Given key, chords, lyrics, prime melody, theme, ...

How to start

Tokenization - Step-by-step tutorial

- TA will provide code (*tutorial.ipynb*) with some blank part (basis tutorial), you can follow the tutorial to learn and play with music tokenization. It may help your data pre-processing if you are not familiar with the MIDI processing pipeline.
- This tutorial is for [midi file to REMI events](#).



How to start

It's a powerful library but a little bit boring if you want to dig into midi file manipulation...

Tokenization - MidiTok [[Docs](#)], [[Github](#)]

- You also may use this package for MIDI file tokenization.
- It can tokenize symbolic music files (MIDI, abc), i.e. convert them into sequences of tokens ready to be fed to models such as Transformer, for any generation.
- Include: REMI, REMI+, MIDI-Like, TSD, Structured, CPWord, Octuple, MuMIDI, MMM, PerTok.
- Tokenizers can be trained with BPE, Unigram or WordPiece.
- Feature data augmentation (on MIDI level or token level).
 - e.g. increase velocities, durations of notes, shift pitches by octaves, ...

How to start

Tokenization - MidiTok [[Docs](#)], [[Github](#)]

- For example

```
# Our parameters
pitch_range = range(21, 109)
beat_res = {(0, 4): 8, (4, 12): 4}
num_velocities = 32
additional_tokens = {'Chord': True, 'Rest': True, 'Tempo': True,
                    'rest_range': (2, 8), # (half, 8 beats)
                    'num_tempos': 32, # num of tempo bins
                    'tempo_range': (40, 250), # (min, max)
                    'Program': False}

# Creates the tokenizer and loads a MIDI
tokenizer = REMI(pitch_range, beat_res, nb_velocities, additional_tokens) # REMI encoding

# tokenizer = CPWord(pitch_range, beat_res, nb_velocities, additional_tokens) # CP encoding

midi = MidiFile('Optimus-VIRTUOSO-RGA-Edition-Main-Sample.mid')

# Converts MIDI to tokens, and back to a MIDI
tokens = tokenizer.midi_to_tokens(midi)
```

How to start

Training - [Hugging Face Transformer](#)

- You can use Hugging Face Transformer API to for training and inference.
- It may be easier if you are not familiar with Transformer model implementation.
- Also it's easier to compare different models, training and inference configurations, ...
- Feel free to explore the documentation and source code to dig deeper.

```
config = GPTConfig(VOCAB_SIZE,
                  max_seq,
                  dim_feedforward=dim_feedforward,
                  n_layer=6,
                  n_head=8,
                  n_embd=512,
                  enable_rpr=True,
                  er_len=max_seq)
model = GPT(config).to(get_device())
```

A good start from MidiTok:

https://github.com/Natooz/MidiTok/blob/main/colab-notebooks/MidiTok_Full_Workflow_Tutorial.ipynb

How to start

Training - Open source code

You may consider using those open source code, it is helpful to learn the whole pipeline of symbolic music generation.

- 1-stage generation: [Pop Music Transformer](#)
- 2-stage generation: [Compose & Embellish](#)

TA also provides sample code for whole training & inference pipeline with some blank part.

Warning & Training Tips

Start this homework as soon as possible!

- It takes time to train a Transformer with a large dataset. (maybe several days for training to generate beautiful music)
- TA's experience (REMI + Transformer-XL)
 - 15~40 mins / epoch (GTX 1080 Ti) depending on training config
 - Model starts to converge when over about 100 epoch
- For those don't have good GPU, you may consider using smaller subset, or use smaller x_len (like 512) for model input, or use CP word to reduce sequence length.
- Write the whole pipeline and ensure that it can generate midi properly first. Then start training!
- At least 12-layers Transformer is recommended.
- For SI score, it is time-consuming to compute scape plot. Start early if you want to report it.
- Don't use other dataset, as it may cause the result tokens distribution far from Pop1K7.

Report

- Write with PPT or PPT-like format (16:9)
- Upload studentID_report.pdf (ex: r12345678_report.pdf)
- Please create a report that is clear and can be understood without the need for oral explanations.
- There is **no specific length requirement**, but it should clearly communicate the experiments conducted and their results. Approximately 10 pages is a suggested standard, but not a strict limitation.

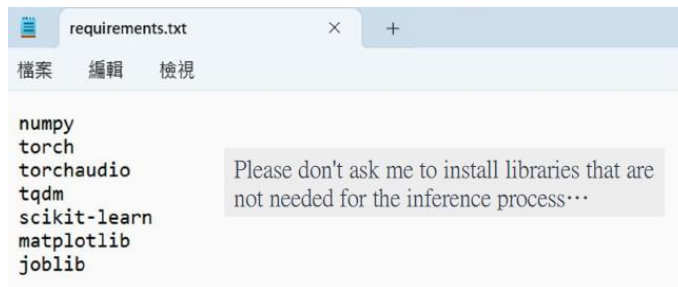
Report template

- **Cover page:** your name, student ID etc
- **Novelty highlight** (one page; optional): what's special about your work?
- **Methodology highlight** (one page): how did you make it? Or, list the attempts you have made
- **Result highlight** (one page): result on your test set
- **Findings highlight** (one page): main takeaways of your study
- **Details of your approach** (multi-pages): If you use open source code, you may want to read some of the associated paper(s) and summarize your understanding of the paper(s) (e.g., why it works)
- **Result analysis & discussion** (multi-pages)

Code

- Upload all your **source code and model** to a cloud drive, **open access permissions**, and then **upload the link to the NTU Cool assignment HW3_report in comments**, as well as include it on the first page of the report.
- You will need to upload **requirements.txt**
- I'll run : `pip install -r requirements.txt`

If you have used third-party programs that cannot be installed directly via 'pip install,' please write the URL and install method command by command on **your readme file**.



The screenshot shows a code editor window titled 'requirements.txt'. The file contains the following text:

```
numpy
torch
torchaudio
tqdm
scikit-learn
matplotlib
joblib
```

Below the code, there is a light gray box with the text: "Please don't ask me to install libraries that are not needed for the inference process..."

Code

- You will also need to upload **README.txt** or **README.pdf** to guide me on how to perform inference on your model. I will inference the generation from scratch, which is to make sure your generation results is not from others.
- The inference code should generate some listening samples.
- The inference process should allow me to set the output file path.

Submission file and details

1. Report (to NTU Cool)
 2. Readme file and requirements.txt (to your cloud drive)
 3. Code and one model checkpoint for inference (to your cloud drive)
 4. Generation results (to your cloud drive)
 - a. Task 1: 20 midi files + 20 wav files for one inference configuration
 - b. Task 2: 3 midi files + 3 wav files for each prompt song (specify the best one for each prompt song!)
- We will randomly select several classmates' code to run inference on your model and run the score on your results, so please ensure that the files you upload include trained model which can successfully execute the entire inference process.
 - **Don't** upload: training data, preprocessed data, others model, cache file

Scoring

- HW3 accounts for 15% of the total grade
- Report: 100%

ALL things you need to do before 11/13 23:59

- HW3_report
 - StudentID_report.pdf
- Cloud drive link
 - generation results
 - README.txt or README.pdf
 - requirements.txt
 - Codes and model to run inference
 - Others codes



When you encounter problem

1. Check out all course materials and announcement documents
2. Use the power of the internet and AI
3. Use **Discussions** on NTU COOL
4. Email me r12942156@ntu.edu.tw or come to office hour