# AI Music Composition

Xu Tan
Microsoft Research Asia
xuta@microsoft.com

# Self-introduction

- Xu Tan (谭旭)

- Senior Researcher @ Machine Learning Group, Microsoft Research Asia

- Research interests: deep learning and its applications on NLP/Speech/Music
  - Text to speech
  - Automatic speech recognition
  - Music understanding and generation
  - Neural machine translation
  - Language/speech pre-training

- Homepage:  https://www.microsoft.com/en-us/research/people/xuta/

- Speech related research: https://speechresearch.github.io/

# Background

- Pipeline of music composition
  - Song Writing (Lyric/Melody) → Accompaniment/Arrangement → Instrumental Recording → Vocal Recoding → Mixing

- General pipeline
  - Score Generation → Performance Generation → Sound Generation

- How deep learning can help?
  - Music is not only about art, also logic/rule/theory!
  - Data, model, and computation

  - Score/Performance generation → Language generation
  - Sound generation → Speech synthesis

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation

- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding

- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# Song writing

- Melody and lyric generation
    - Lack of paired melody and lyric data
    - The connection between melody and lyric is weak
        - Unlike other tasks: Automatic Speech Recognition, Text to Speech, Neural Machine Translation
        - Needs large amount of paired data
        - Or motivate us to find connections from other aspects

- How to model the connections
    - Learning: SongMASS
    - knowledge based on rhythm/structure: StructMelody
    - Combine them together: ongoing

# SongMASS: Automatic Song Writing with Masked Sequence to Sequence Pre-training, AAAI 2021

- Background
  - Lyric-to-melody and melody-to-lyric generation are two important tasks for song writing
  - Lyric and melody are weakly coupled, but strictly aligned



**Paired Aligned Data :**

| Lyric | Another | | | | day | has | gone | I'm | | still | alone | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pitch | R | G3 | E4 | D4 | C4 | B3 | C4 | R | E4 | C4 | B3 | C4 |
| Duration | $\frac{7}{16}$ | $\frac{1}{16}$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{3}{16}$ | $\frac{1}{16}$ | $\frac{5}{16}$ | $\frac{1}{4}$ | $\frac{1}{8}$ | $\frac{3}{16}$ | $\frac{1}{16}$ | $\frac{5}{16}$ |

# SongMASS

- Background
  - Lack of training data
    - The two domains are weak coupled, need a lot of data to build the relationship
    - A lot of unpaired data available on the web
    - Previous works only use supervised data from training, the quality is limited

  - **Solution**
    - **Adapt masked sequence to sequence pre-training (MASS) on song writing for both tasks**

# SongMASS

- Background
  - Lyric and melody alignment
    - For each word/syllable, which note to align? How many notes to align?



《再见二丁目》
作词：林夕
作曲：于逸尧
演唱：杨千嬅

《开始懂了》
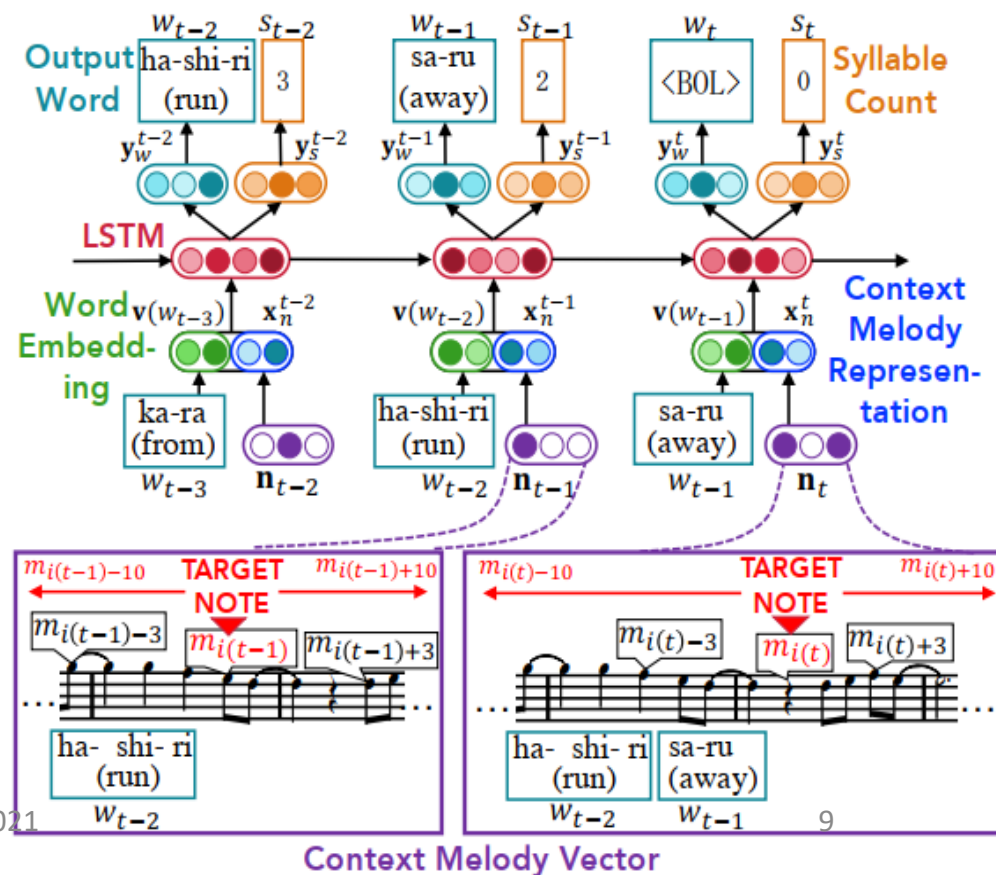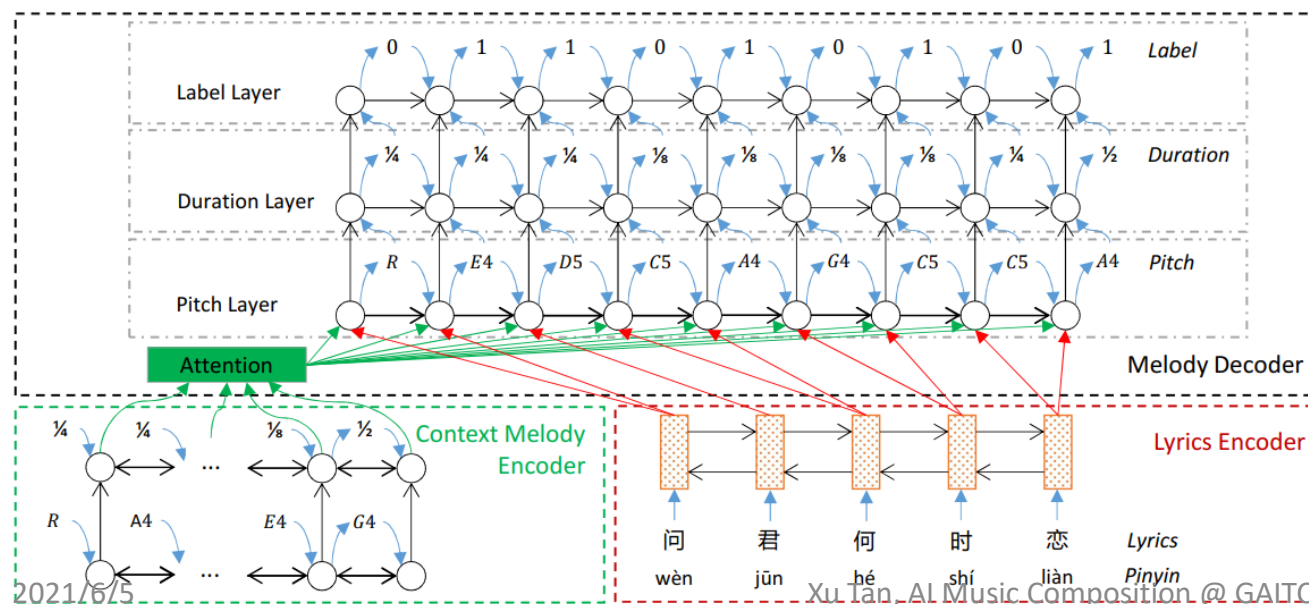作词：姚若龙
作曲：李偲菘
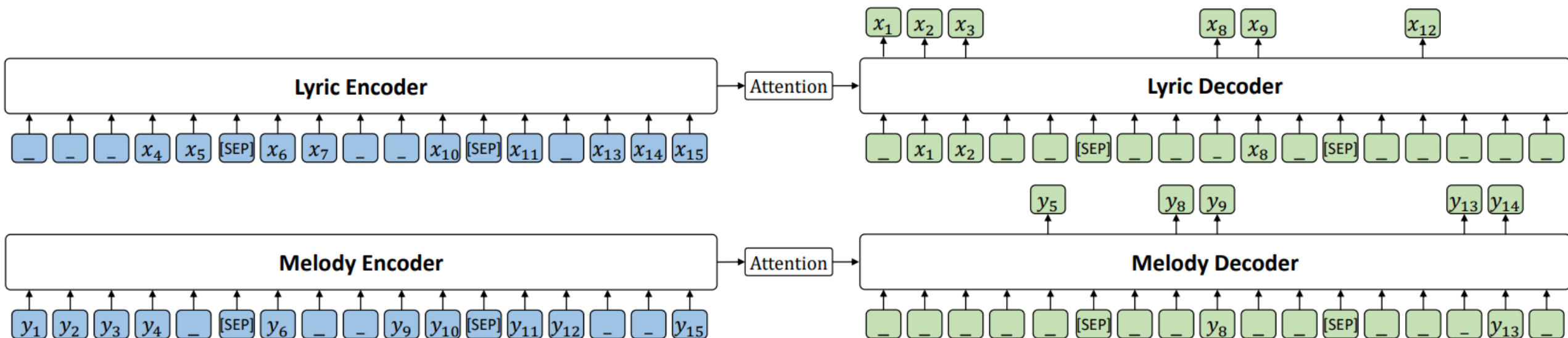演唱：孙燕姿

# SongMASS

- Background
  - Lyric and melody alignment
    - For each word/syllable, which note to align? How many notes to align?
    - Previous works
      - Decide if switch to next word when predicting notes (lyri
      - Predict how many syllable in predicting word, to decide

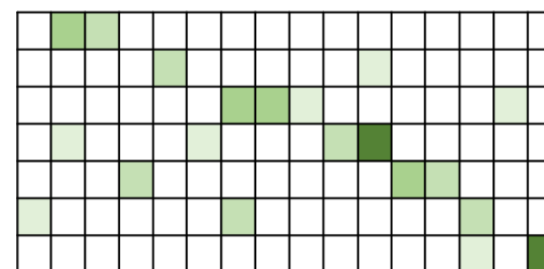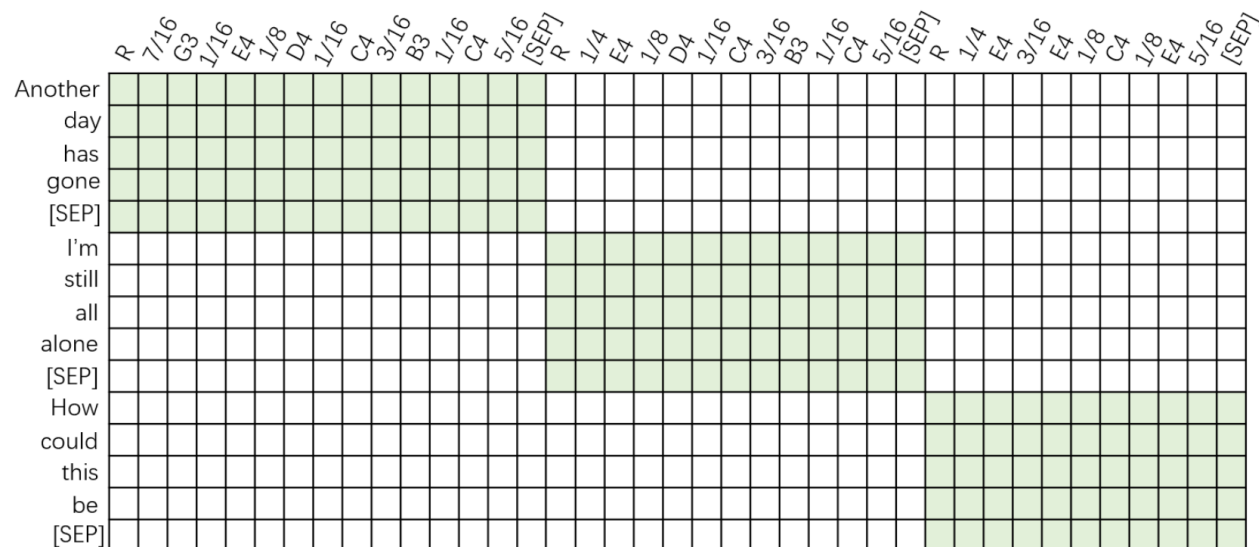Xu Tan, AI Music Composition @ GAITC 2021

# SongMASS

- MASS pre-training
  - Document-level MASS, mask each a segment in each sentence and predict all segments in the target
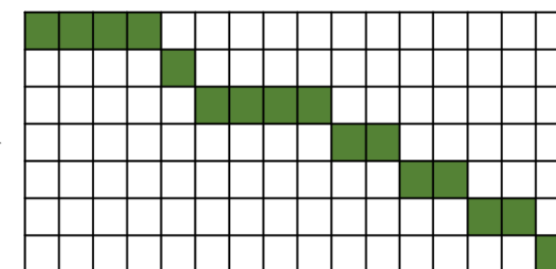  - Separate encoder and decoder, add supervised loss to guide the pre-training



Xu Tan, AI Music Composition @ GAITC 2021

# SongMASS

- Lyric and melody alignment
  - Sentence-level and token-level alignment
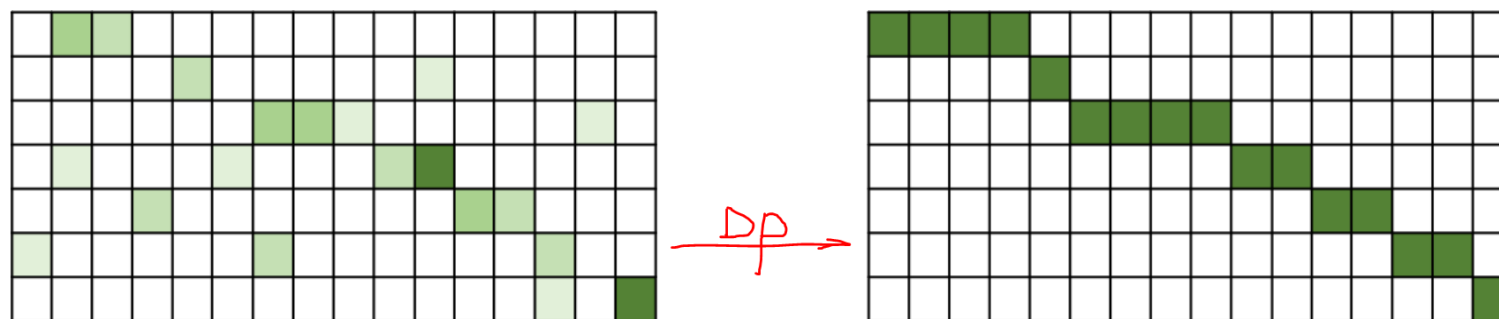  - During training, attention constraint

# SongMASS

- Lyric and melody alignment
  - Sentence-level and token-level alignment
  - During training, attention constraint
  - During inference
    - Sentence-level: SEP token
    - Token-level: Dynamic programming



**Algorithm 2** DP for Duration Extraction

1: **Input**: Alignment matrix $A \in \mathbb{R}^{\mathcal{T} \times \mathcal{S}}$
2: **Output**: Phoneme duration $D \in \mathbb{R}^{\mathcal{T}}$
3: **Initialize**: Initialize reward matrix $O \in \mathbb{R}^{\mathcal{T} \times \mathcal{S}}$ with zero matrix. Initialize the prefix sum matrix $C \in \mathbb{R}^{\mathcal{T} \times \mathcal{S}}$ to the prefix sum of each row of $A$, that is, $C_{i,j} = \sum_{k=0}^{j} [A]_{i,k}$. Initialize all elements in the splitting boundary matrix $B_m \in \mathbb{R}^{\mathcal{T} \times \mathcal{S}}$ to zero.
4: **for** each $j \in [0, \mathcal{S})$ **do**
5:     $[O]_{0,j} = [C]_{0,j}$
6: **end for**
7: **for** each $i \in [1, \mathcal{T})$ **do**
8:     **for** each $j \in [0, \mathcal{S})$ **do**
9:         **for** each $k \in [0, \mathcal{S})$ **do**
10:             $O_{new} = [O]_{i-1,k} + [C]_{i,j} - [C]_{i,k}$
11:             **if** $O_{new} > [O]_{i,j}$ **then**
12:                 $[O]_{i,j} = O_{new}$
13:                 $[B_m]_{i,j} = k$
14:             **end if**
15:         **end for**
16:     **end for**
17: **end for**
18: $P = \mathcal{S} - 1$
19: **for** each $i \in [\mathcal{T} - 1, 0]$ **do**
20:     $[D]_i = P - [B_m]_{i,P}$
21:     $P = [B_m]_{i,P}$
22: **end for**
23: **return** $D$

# SongMASS

- Experiments
  - Datasets
    - Unpaired data: total 362,237 song lyrics, 65,000 song melodies
    - Paired data: LMD, 7998 songs
  - Data preprocessing
    - Pitch normalized to C major or A minor
    - Duration normalized to 1/16 note
    - Lyrics: BPE sequence
    - Melody: pitch, duration, pitch, duration, ...
  - Metrics
    - Objective
      - Pitch distribution (PD), duration distribution (DD), Melody Distance (MD), Alignment similarity (AS), Perplexity (PPL)
    - Subjective
      - Lyric: Listenability, Grammaticality, Meaning, Quality. Melody: Emotion, Rhythm, Quality

# SongMASS

- Experiments
  - Results in objective evaluation

|  | Lyric-to-Melody | | | | Melody-to-Lyric |
|  | PD (%) ↑ | DD (%) ↑ | MD ↓ | PPL ↓ | PPL ↓ |
|---|---|---|---|---|---|
| Baseline | 38.20 | 52.00 | 2.92 | 3.27 | 37.50 |
| **SongMASS** | 57.00 | 65.90 | 2.28 | 2.41 | 14.66 |
| − pre-training | 43.50 | 57.00 | 2.79 | 3.72 | 45.10 |
| − separate encoder-decoder | 55.00 | 64.80 | 2.32 | 2.53 | 15.57 |
| − supervised loss | 47.20 | 53.60 | 3.29 | 2.92 | 27.50 |
| − alignment | 56.10 | 65.20 | 2.36 | 2.07 | 8.54 |

  - Results in subjective evaluation

| Metric | Baseline | SongMASS |
|---|---|---|
| *Lyric* | | |
| Listenability | 1.67 ± 0.62 | 2.00 ± 0.65 |
| Grammaticality | 3.00 ± 0.76 | 3.27 ± 0.59 |
| Meaning | 2.20 ± 0.68 | 3.20 ± 0.68 |
| Quality | 2.27 ± 0.46 | 3.00 ± 0.38 |
| *Melody* | | |
| Emotion | 2.40 ± 1.06 | 3.53 ± 0.64 |
| Rhythm | 2.33 ± 1.18 | 2.87 ± 0.74 |
| Quality | 2.33 ± 1.05 | 2.93 ± 0.70 |

# SongMASS

- Experiments
  - Study on the alignment constraints

| | L2M Acc ↑ | M2L Acc ↑ |
|---|---|---|
| **SongMASS** | 62.6 | 45.4 |
| - TC | 62.1 | 44.8 |
| - SC | 56.2 | 44.0 |
| - TC - SC | 55.3 | 43.8 |
| - TC - SC - PT | 48.3 | 37.1 |
| - DP | 15.7 | 11.3 |

# SongMASS

- Demo
  - https://speechresearch.github.io/songmass/

Xu Tan, AI Music Composition @ GAITC 2021

```
1 3 5 3  2      1    6  1
you have loved lots of girls
1  1    7      6    5 3 6
in the sweet long ago
1  -    1   7  6    5     3  6
and each one has meant heaven to you
3   5 5 3 2 1     6     1
you have      vowed your affection
1  1    7   6  5 3
to each one in turn
3  3    5     3  2    1 6 1
and have sworn to them be  true
6 6 6 5    5 3     2   1
you   have kissed the moon
1    1   7    7     6 5 3
while the world seemed in  tune
6    3    3   5  3    2 1    2
then left her to hunt a new game
1    3 5 3    2     1 6   1
does it  ever occur to you later
1  2 1 3
my boy
1 2 1 3 2 1 3   2
that      doing the
6 6     5      5 3 2 1  |
i wonder kissing her    now
6 1    1      2 1 3
wonder teaching her
1    2    1    3   -
wonder looking into her eyes
1      6     -    1
breathing sighs telling lies
1 1    7    6   5 3 6
i wonder buying the wine
1  1   7    6 5    3 -    6
```

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# StructMelody

- Background
  - Lyric and melody is weakly correlated
  - Data hungry but low-resource
  - However, lyric and melody has its own structures

- Solution
  - Lyric → Structure, Structure → Melody
  - Lyric → Structure': learned based on supervised data
  - Structure'' → Melody: self-supervised learning from music data
  - Close the gap between Structure' and Structure''

# StructMelody

- Structure: Rhythm, Beat, Bar, Chord, Form
- How to get lyric-structure data



Xu Tan, AI Music Composition @ CAAI 2021

# StructMelody

- Experiment results
  - 古诗词：《春晓》
    - 春眠不觉晓，处处闻啼鸟。
    - 夜来风雨声，花落知多少。

  - 散文诗：《童话》
    - 我给你们讲
    - 一位森林仙女
    - 她的样子和你们一样的
    - 她是一位女河神的妹妹
    - 她的衣裳多么离奇
    - 那是用露水和月光的薄纱做的
    - 这位仙女
    - 在树叶里面正要睡去
    - 活像这个时候的你们

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# DeepRapper: Neural Rap Generation with Rhyme and Rhythm Modeling, ACL 2021

- Explore a new lyric-melody relationship: Rap

- Rap is a musical form of vocal delivery that incorporates "rhyme, rhythmic speech, and street vernacular"
  - Originated in America in the 1970s
  - Popular in the world especially in young people

- Hip-Hop
  - 1970s originated from New York, young people in African-American and Latino
  - Street culture
  - Four elements in Hip-Hop
    - DJ (Disc Jockey)
    - Rap (MC)
    - Street Dance (B-Boy)
    - Graffiti

# DeepRapper

- Lyric with Rhyme and Rhythm, and sing out
  - Rhyme and Rhythm (beat) is important
  - Rap cares more about beat/duration, rather than pitch (melody)
- However, previous works on rap generation only consider rhyme, but ignores rhythm
  - How they control rhyme? Use Rhyme list. Complicated and not learned end-to-end
  - No rhythm/beat information, cannot be directly used!

# DeepRapper

- Generated results
  - N Rhyme: single, double, multiple
  - 下苦功 练武功 变武松
  - Diversity in rhyme

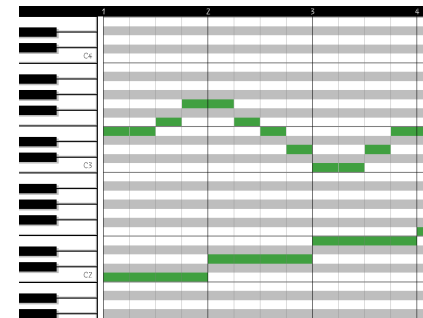- Demo
  - https://deeprapper.github.io/

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
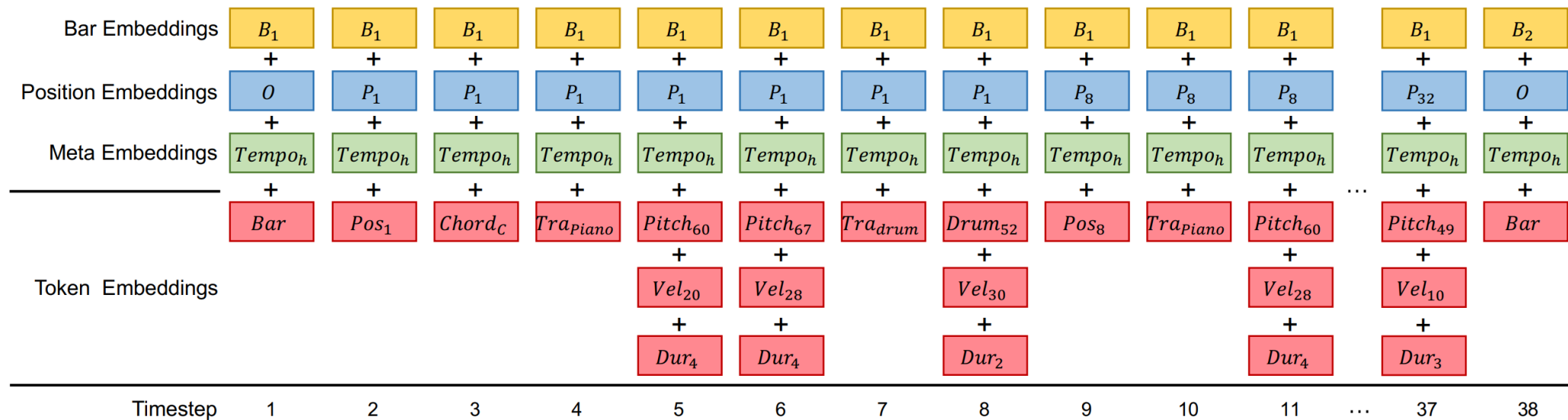  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# PopMAG: Pop Music Accompaniment Generation, ACM MM 2020

- Music accompaniment generation/arrangement are challenging
  - Multi-track generation: Lead, Chord → Drum, Bass, Guitar, Piano, String
  - Arrangement: ensure the harmony between tracks

- Previous works
  - Pianoroll: MuseGAN, MIDI-Sandwich
    - Generate as image, suffers from data sparsity
  - Multi-track MIDI: Xiaoice Band, LakhNES
    - Cannot ensure the dependency in the same step
  - There are no explicitly dependency among tracks

# PopMAG

- MUlti-track MIDI representation (MuMIDI)
  - enables simultaneous multi-track generation in a single sequence
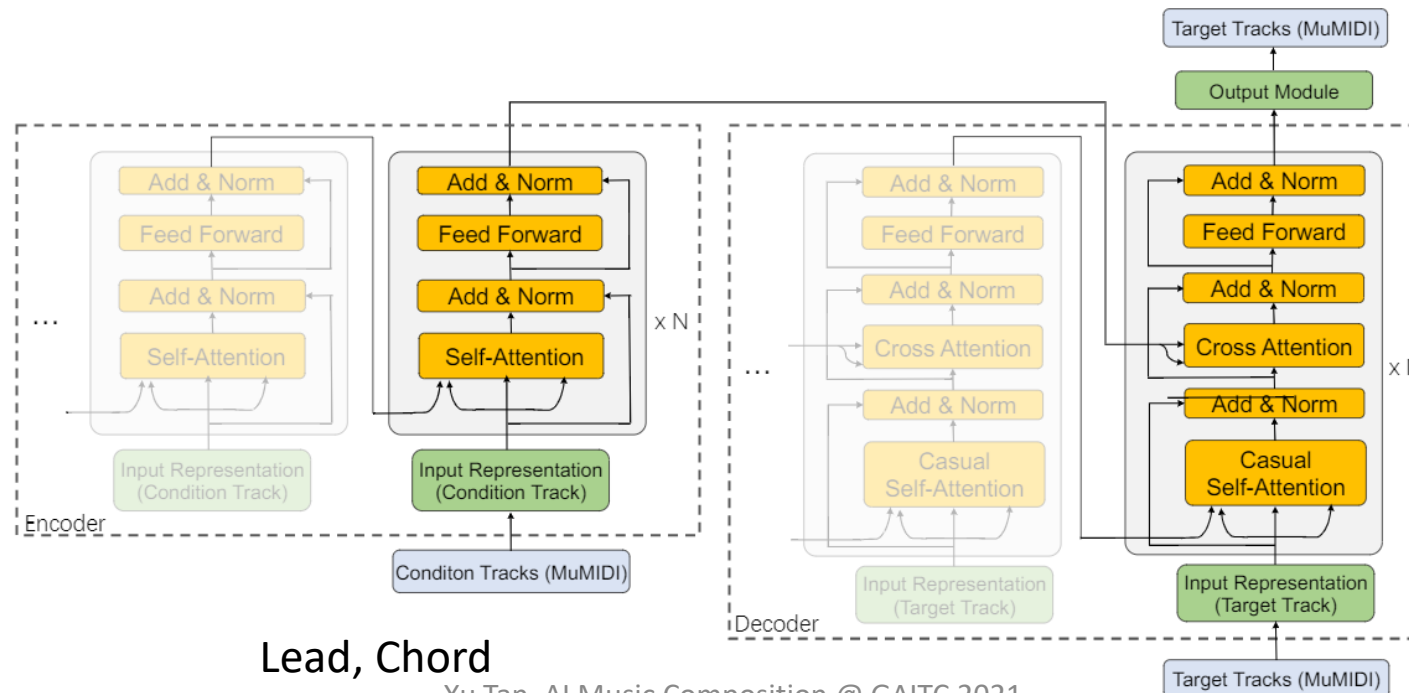  - explicitly models the dependency of the notes from different tracks



**Bar**: <Bar> token,  **Position**: 32 position (1/32),   **Chord**: 12 chord root * 7 types = 84 chords
**Track**: Lead, Chord, Drum, Bass, Guitar, Piano, String,   **Note**: Pitch, Duration, Velocity

# PopMAG

- MuMIDI sequence is long and challenging for long-term music modeling
  - Shorten the sequence length: modeling multiple note attributes (e.g., pitch, duration, velocity) in one step
  - Introduce long-term context as memory



Lead, Chord

Drum, Bass, Guitar, Piano, String

# PopMAG

- Experiments
  - Dataset
    - Lakh MIDI
    - FreeMIDI
    - An internal Chinese Pop MIDI (CPMD)

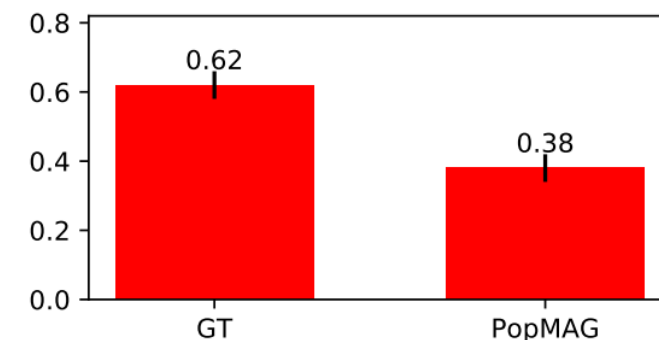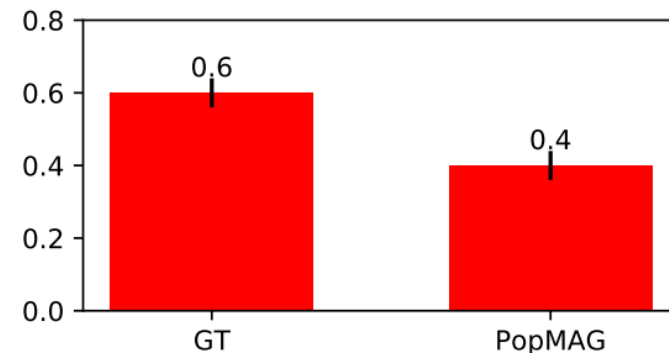| Dataset | #Musical Pieces | #Bars | Duration (hours) |
|---------|-----------------|-------|------------------|
| LMD | 21916 | 372339 | 255.13 |
| FreeMidi | 5691 | 92825 | 52.32 |
| CPMD | 5344 | 94170 | 54.12 |

Melody        Melody+ Generated Accompaniment



**(a) Preference scores on LMD.**



**(b) Preference scores on FreeMidi.**



**(c) Preference scores on CPMD.**

Xu Tan, AI Music Composition @ GAITC 2021
https://speechresearch.github.io/popmag/

# Arrangement

- Horizonal axis (time): music form, chord progression
- Vertical axis (harmony): texture (Melody, Harmony, Base, Rhythm, Noise）

| Music Form: verse-chorus | Intro: 4 | Verse: 16 | Chorus: 16 | Interlude: 4 | Verse: 8 | Chorus: 16 | Outro: 6 |
|---|---|---|---|---|---|---|---|
| Melody | | Sequence | Syncopation | | | Strengthen | Slow |
| Harmony | Guitar | Guitar | Piano | | | | |
| Base | | | Bass | | | | |
| Rhythm | | | Drum | | | | |
| Noise | Sea Wave | | | | | | |

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# MusicBERT: Symbolic Music Understanding with Large-Scale Pre-Training, ACL 2021

- Understanding music is important for generation
  - Emotion recognition
  - Genre classification
  - Melody/accompaniment extraction
  - Structure analysis

- Previous works on music understanding
  - PiRhDy, ACM MM 2020 best paper, contextual word embedding
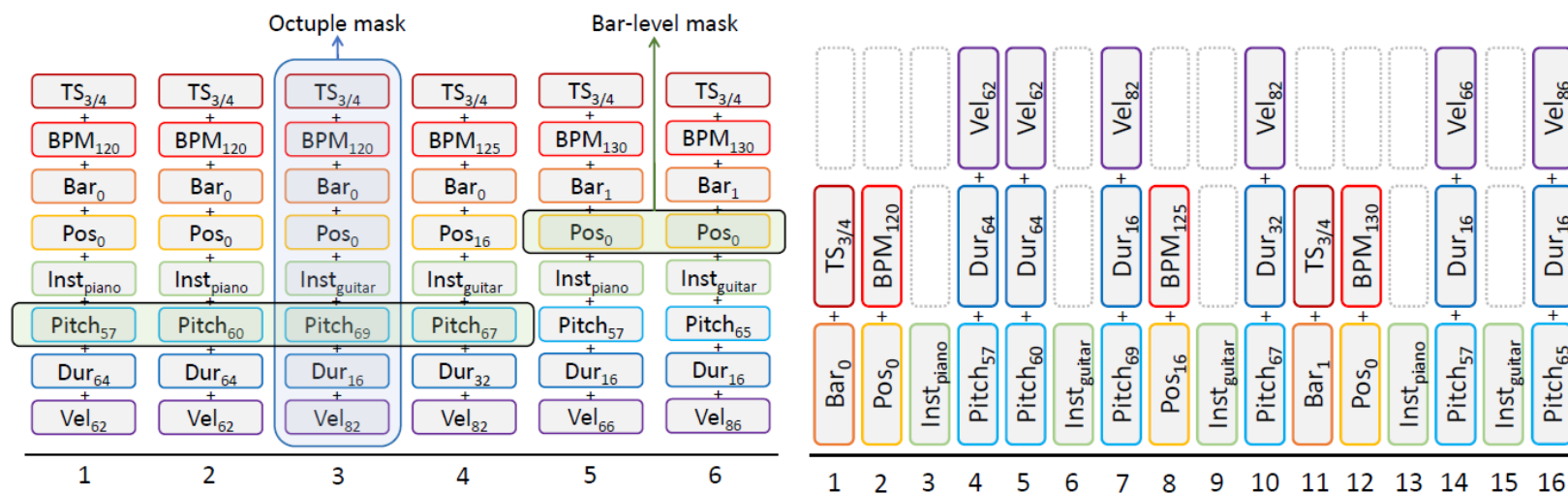  - Shallow model, too much complicated design with music knowledge

# MusicBERT

- Dataset construction: Million MIDI Dataset (MMD)
  - Crawled from various MIDI and sheet music websites
  - 1.5 million songs after deduplication and cleaning (10x larger than LMD)

| Dataset | Songs | Notes (Millions) |
|---|---|---|
| MAESTRO | 1,184 | 6 |
| GiantMIDI-Piano | 10,854 | 39 |
| LMD | 148,403 | 535 |
| MMD | **1,524,557** | **2,075** |

- Data representation: OctupleMIDI
  - Compound token: (Bar_1, TimeSig_4/4, Pos_35, Tempo_120, Piano, Pitch_64, Dur_12, Vel_38)
  - Supports changing tempo and time signature
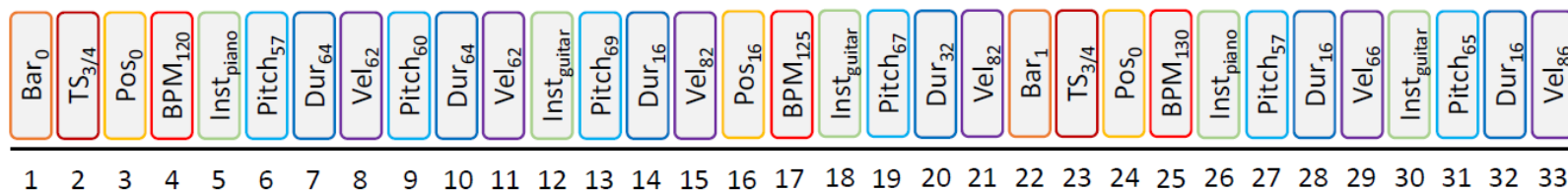  - Shorter length compared to REMI and MuMIDI in PopMAG

# MusicBERT

- OctupleMIDI representation



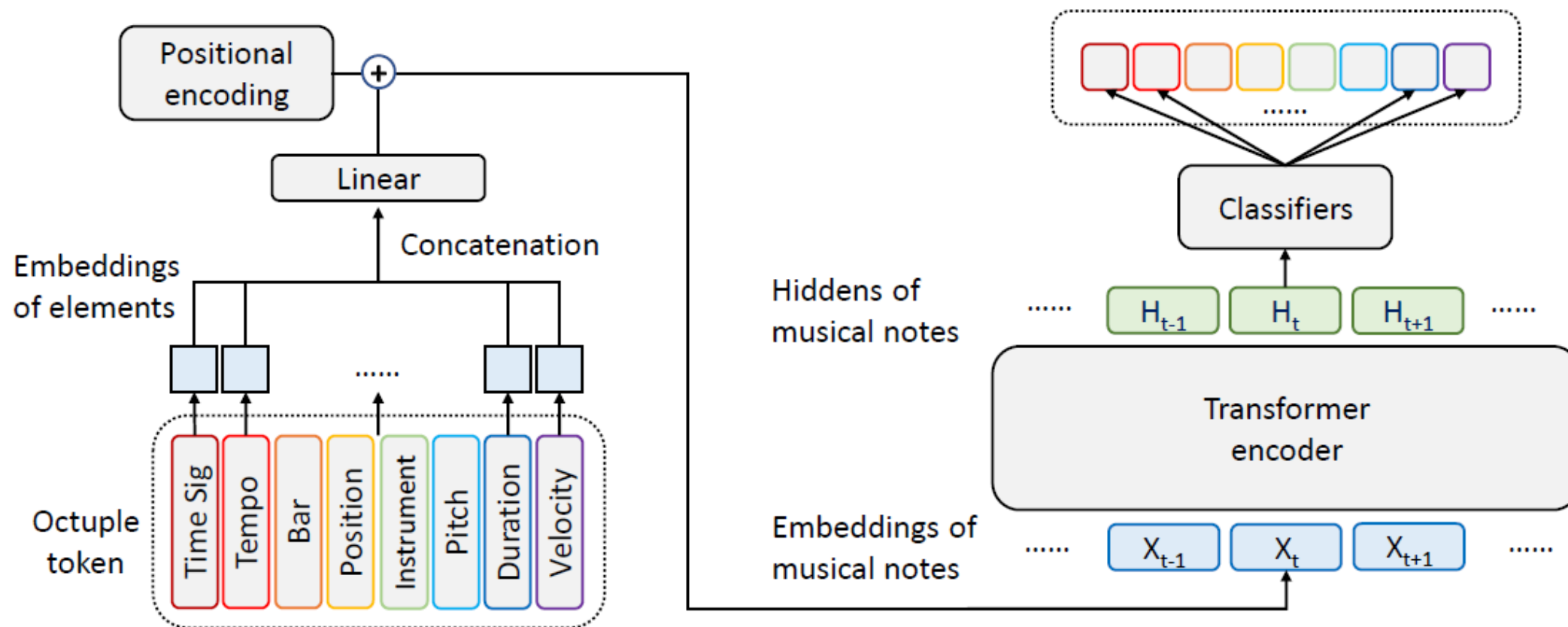(a) OctupleMIDI encoding.

(b) CP-Like encoding.

(c) REMI-Like encoding.

| Encoding | OctupleMIDI | CP-like | REMI-like |
|---|---|---|---|
| Tokens | **3607** | 6906 | 15679 |

Xu Tan, AI Music Composition @ GAITC 2021

# MusicBERT

- Model structure

# MusicBERT

- Experiments
  - Melody completion
    - Two sequences classification
  - Accompaniment completion
    - Melody and accompaniment sequences classification
  - Genre classification
    - Single sentence classification

| Model | Melody Completion | | | | | Accompaniment Suggestion | | | | | Classification | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MAP | HITS @1 | HITS @5 | HITS @10 | HITS @25 | MAP | HITS @1 | HITS @5 | HITS @20 | HITS @25 | Genre F1 | Style F1 |
| melody2vec$_F$ | 0.646 | 0.578 | 0.717 | 0.774 | 0.867 | - | - | - | - | - | 0.649 | 0.299 |
| melody2vec$_B$ | 0.641 | 0.571 | 0.712 | 0.772 | 0.866 | - | - | - | - | - | 0.647 | 0.293 |
| tonnetz | 0.683 | 0.545 | 0.865 | 0.946 | 0.993 | 0.423 | 0.101 | 0.407 | 0.628 | 0.897 | 0.627 | 0.253 |
| pianoroll | 0.762 | 0.645 | 0.916 | 0.967 | 0.995 | 0.567 | 0.166 | 0.541 | 0.720 | 0.921 | 0.640 | 0.365 |
| PiRhDy$_{GH}$ | 0.858 | 0.775 | 0.966 | 0.988 | 0.999 | 0.651 | 0.211 | 0.625 | 0.812 | 0.965 | 0.663 | 0.448 |
| PiRhDy$_{GM}$ | 0.971 | 0.950 | 0.995 | 0.998 | 0.999 | 0.567 | 0.184 | 0.540 | 0.718 | 0.919 | 0.668 | 0.471 |
| MusicBERT$_{small}$ | 0.979 | 0.966 | 0.995 | 0.998 | **1.000** | 0.920 | 0.325 | 0.834 | 0.991 | 0.996 | 0.762 | 0.604 |
| MusicBERT$_{base}$ | **0.984** | **0.973** | **0.997** | **0.999** | **1.000** | **0.945** | **0.333** | **0.856** | **0.995** | **0.998** | **0.784** | **0.651** |

# MusicBERT

- Experiments
  - Ablation studies

| Encoding | Melody | Accom. | Genre | Style |
|---|---|---|---|---|
| CP-like | 96.6 | 88.0 | 0.750 | 0.594 |
| REMI-like | 96.7 | 88.4 | 0.734 | 0.562 |
| OctupleMIDI | **96.9** | **88.7** | **0.762** | **0.604** |

| Mask | Melody | Accom. | Genre | Style |
|---|---|---|---|---|
| Random | 96.7 | 88.1 | 0.753 | 0.602 |
| Octuple | 96.7 | 88.1 | 0.751 | 0.606 |
| Bar | **97.0** | 88.1 | **0.766** | **0.610** |

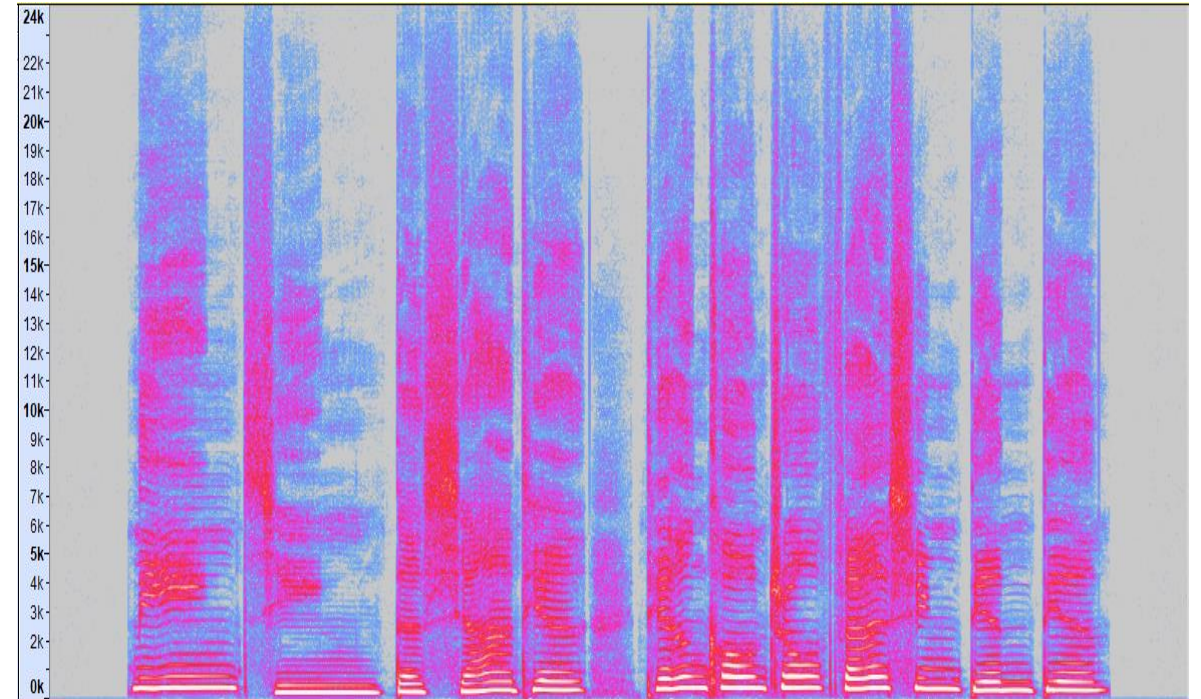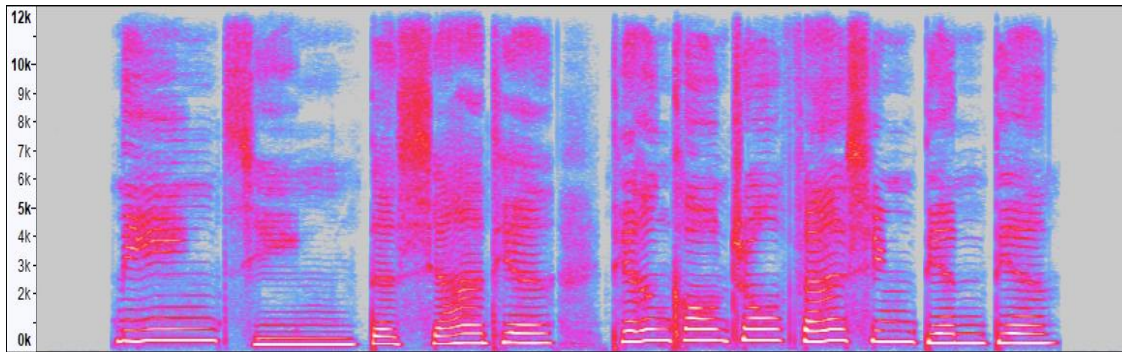| Model | Melody | Accom. | Genre | Style |
|---|---|---|---|---|
| No pre-train | 93.7 | 77.4 | 0.677 | 0.450 |
| MusicBERT | **96.9** | **88.7** | **0.762** | **0.604** |

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- **Singing voice synthesis**
  - **HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis**
  - **XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)**

# HiFiSinger: Towards High-Fidelity Neural Singing Voice Synthesis

- Compared with speaking voice, singing voice need high-fidelity to convey expressiveness and emotion

- How to ensure high-fidelity?  High sampling rate
  - Speaking voice in TTS: 16KHz or 24KHz
  - Human can perceive frequency 20~20K
  - According to Nyquist-Shannon frequency, 16KHz or 24KHz can convey 8KHz or 12KHz frequency

- Increase to 48KHz, can convey 24KHz frequency, fully satisfy human ear

- Challenges of 48KHz
  - 48KHz vs 24KHz, wide frequency cause challenges to acoustic model
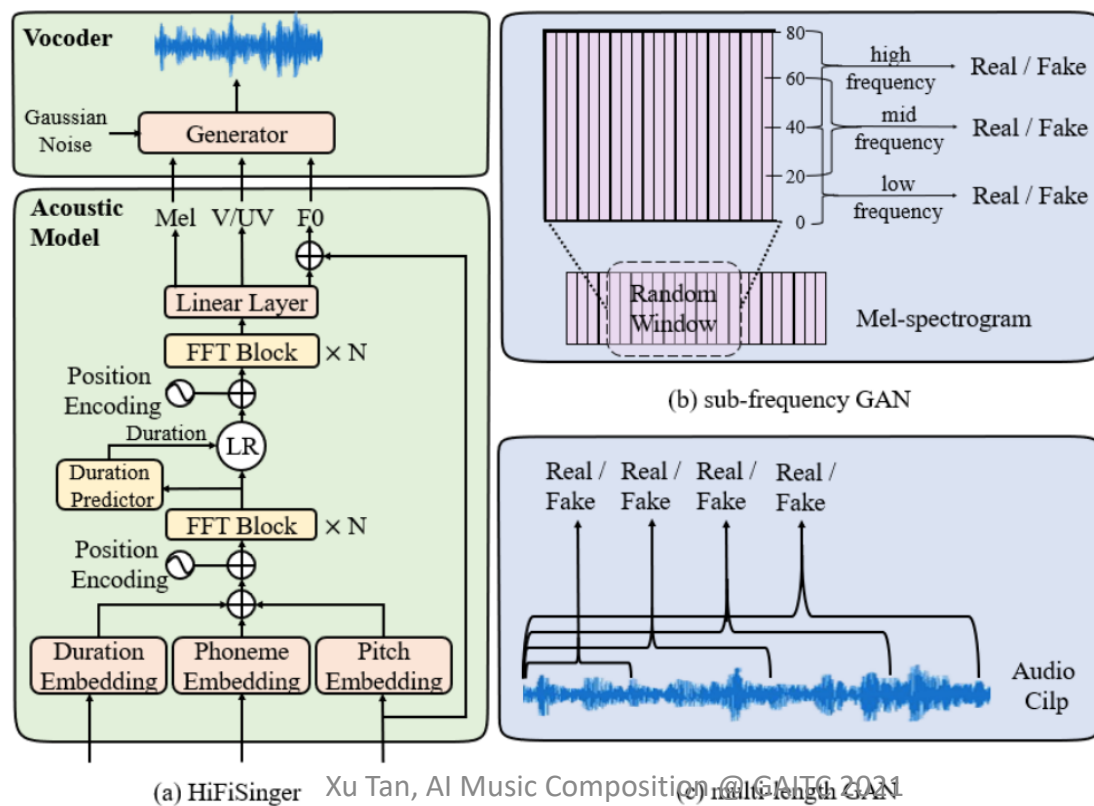  - 48KHz, 1s has 48000 waveform points, cause challenges to vocoder
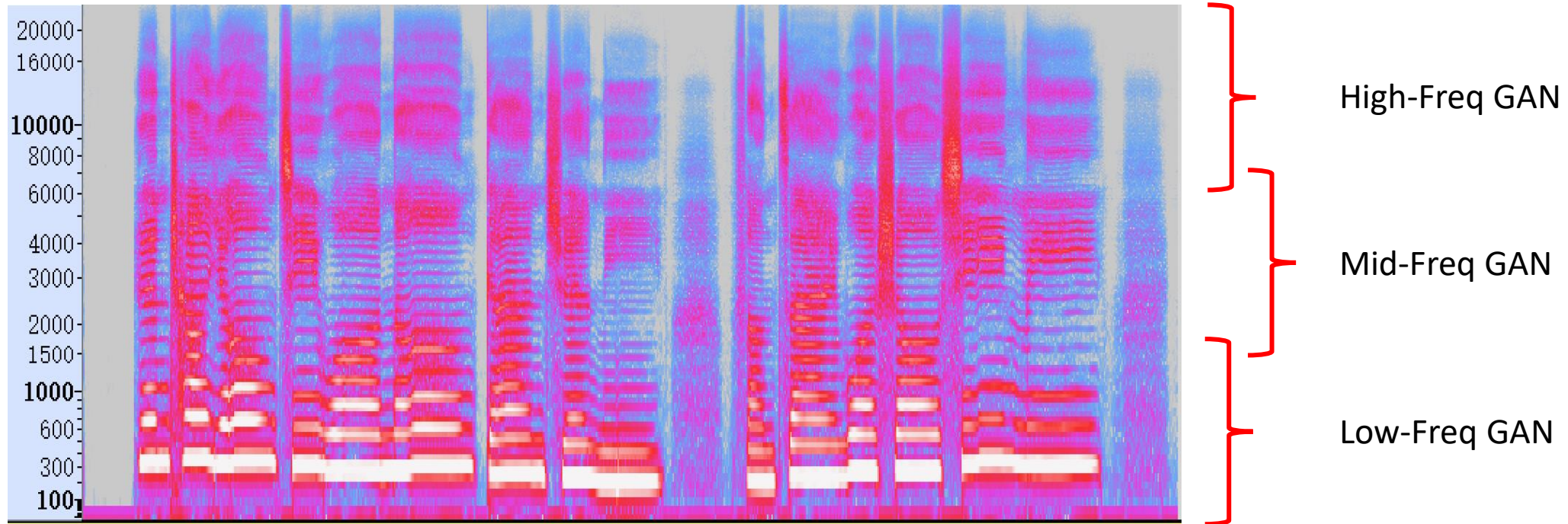
# HiFiSinger

- Demo voice

# HiFiSinger

- Model pipeline
  - Acoustic model: lyric + score → mel-spectrogram
  - Vocoder: mel-spectrogram → waveform

Xu Tan, AI Music Composition @IGAITC 2021

# HiFiSinger

- Sub-frequency GAN
  - Use different GAN focus on different frequencies
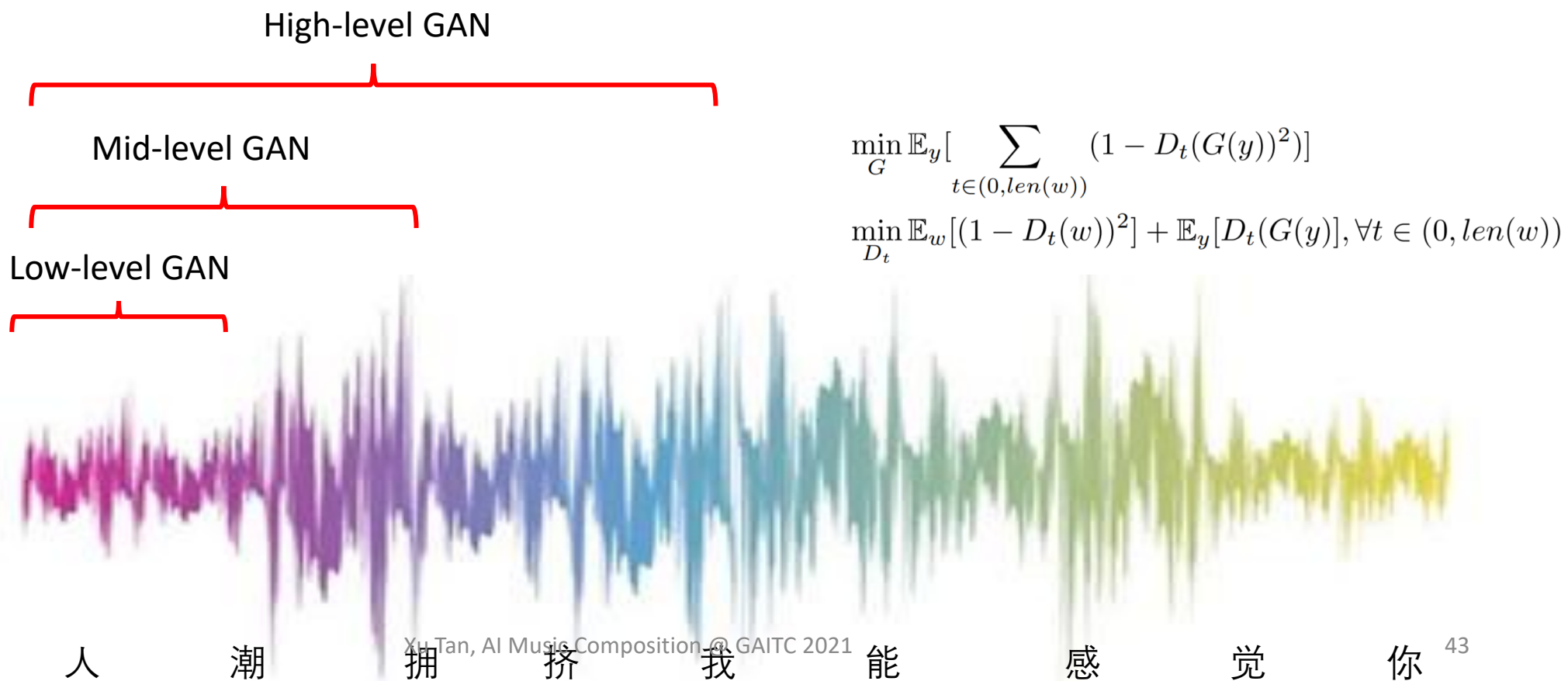


$$\min_{G} \mathbb{E}_x \Big[ \sum_{f \in \{low, mid, high\}} (1 - D_f(G(x))^2) \Big]$$

$$\min_{D_f} \mathbb{E}_y [(1 - D_f(y))^2] + \mathbb{E}_x [D_f(G(x)], \forall f \in \{low, mid, high\}$$

Xu Tan, AI Music Composition @ GAITC 2021

# HiFiSinger

- Multi-length GAN
  - Use different GAN focus on different time resolution

High-level GAN

Mid-level GAN

Low-level GAN

$$\min_{G} \mathbb{E}_y \Big[ \sum_{t \in (0, len(w))} (1 - D_t(G(y))^2) \Big]$$

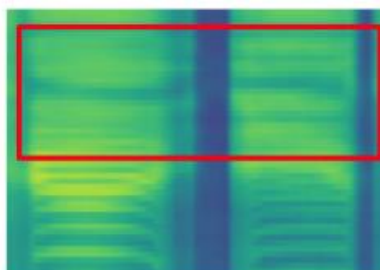$$\min_{D_t} \mathbb{E}_w [(1 - D_t(w))^2] + \mathbb{E}_y [D_t(G(y)], \forall t \in (0, len(w))$$

人　潮　拥　挤　我　能　感　觉　你

# HiFiSinger

- Systematic improvements
  - Hop size/window size tradeoff
  - Pitch/UV
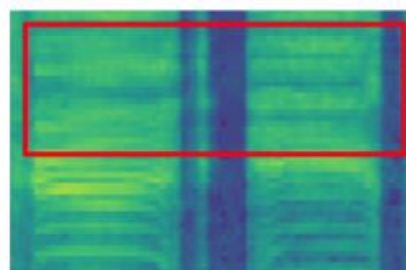  - Increase receptive field
  - Use long audio clips
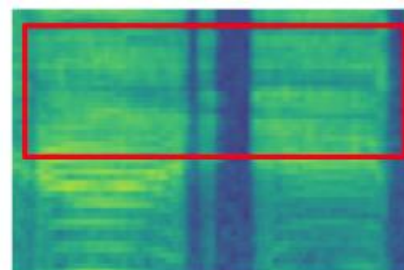
# HiFiSinger

- Experiments
  - Audio quality

| Method | MOS |
|---|---|
| Recording | $4.03 \pm 0.06$ |
| Recording (24kHz) | $3.70 \pm 0.08$ |
| XiaoiceSing (Lu et al., 2020) | $2.93 \pm 0.06$ |
| Baseline (24kHz) | $3.32 \pm 0.09$ |
| Baseline (24kHz upsample) | $3.38 \pm 0.08$ |
| Baseline | $3.44 \pm 0.08$ |
| HiFiSinger (24kHz) | $3.47 \pm 0.06$ |
| HiFiSinger | $3.76 \pm 0.06$ |

  - Ablation study

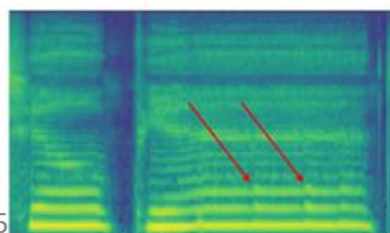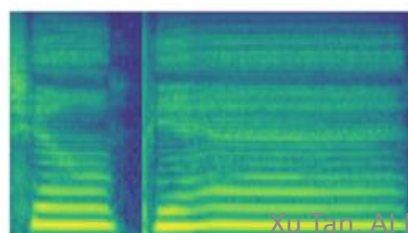

(a) HiFiSinger w/o SF-GAN     (b) HiFiSinger     (c) Ground truth

(a) HiFiSinger w/o ML-GAN     (b) HiFiSinger     (c) Ground truth

https://speechresearch.github.io/hifisinger/

# Our research work

- Song writing
  - SongMASS (AAAI 2021), for lyric and melody generation
  - StructMelody (ongoing), for melody generation
  - DeepRapper (ACL 2021), for lyric and rhythm generation
- Arrangement
  - PopMAG (ACM MM 2020), for accompaniment
  - MusicBERT (ACL 2021), for music structure understanding
- Singing voice synthesis
  - HiFiSinger (arXiv 2020), for high-fidelity singing voice synthesis
  - XiaoiceSing (INTERSPEECH 2020), DeepSinger (KDD 2020)

# Research challenges

- Music structure
  - Clear theme and self-repetitive structure （Motif → Sequence）
  - Music form: rondo, variation, sonata, ternary, verse-chorus, Chinese
  - Arrangement: harmony, orchestration

- Emotion and Style
  - How to recognize emotion and style
  - How to control the emotion and style in generation

- Interaction
  - Retain a certain level of creative freedom when composing music with AI

- Originality
  - How to ensure innovation, instead of fitting data distribution

# Thank You!

Xu Tan (谭旭)
Senior Researcher @ Microsoft Research Asia
xuta@microsoft.com

https://www.microsoft.com/en-us/research/people/xuta/
https://www.microsoft.com/en-us/research/project/ai-music/