

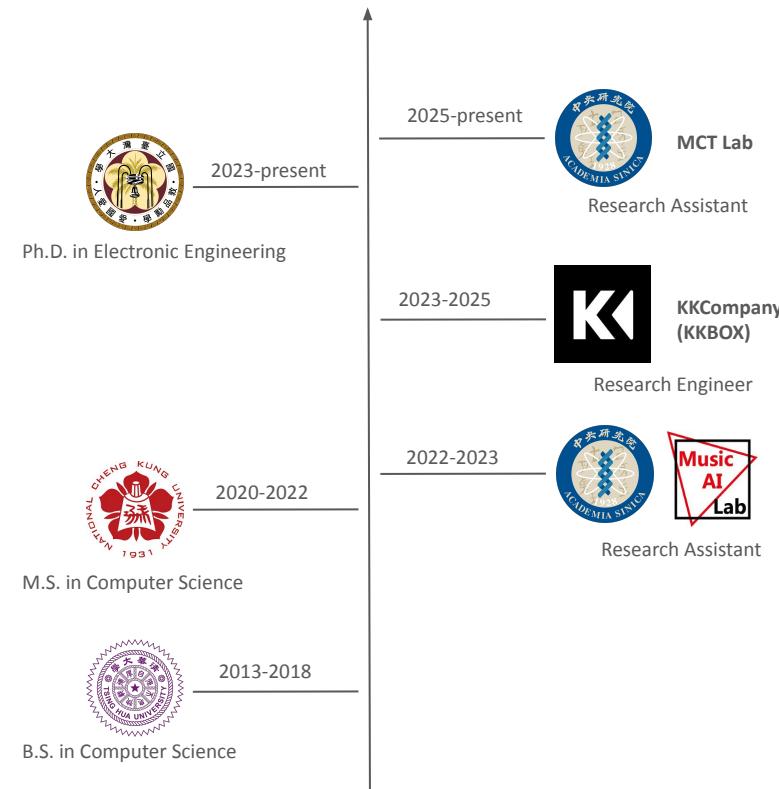
Music Adaptation Arrangement & Cover Generation

Chih-Pin Tan 譚至斌

National Taiwan University

About me

- Chih-Pin (CP) Tan
- Email: tanchihpin0517@gmail.com
- Research Topic
 - Music generation & arrangement
 - Cover song generation



Introduction

**Have you ever heard a song that sounds so familiar
to another song?**

Example

- papaya - Listen to me <-> 王心凌 - 愛你

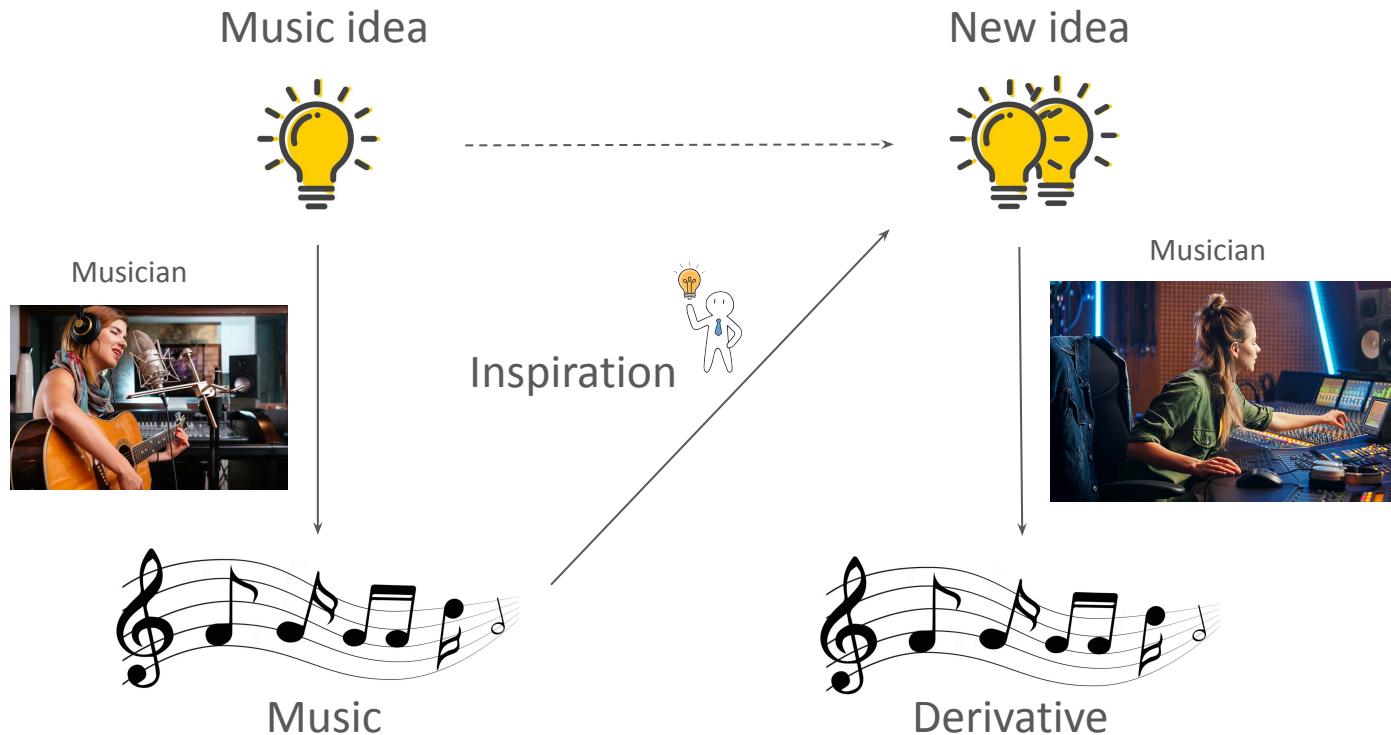


Song adaptation

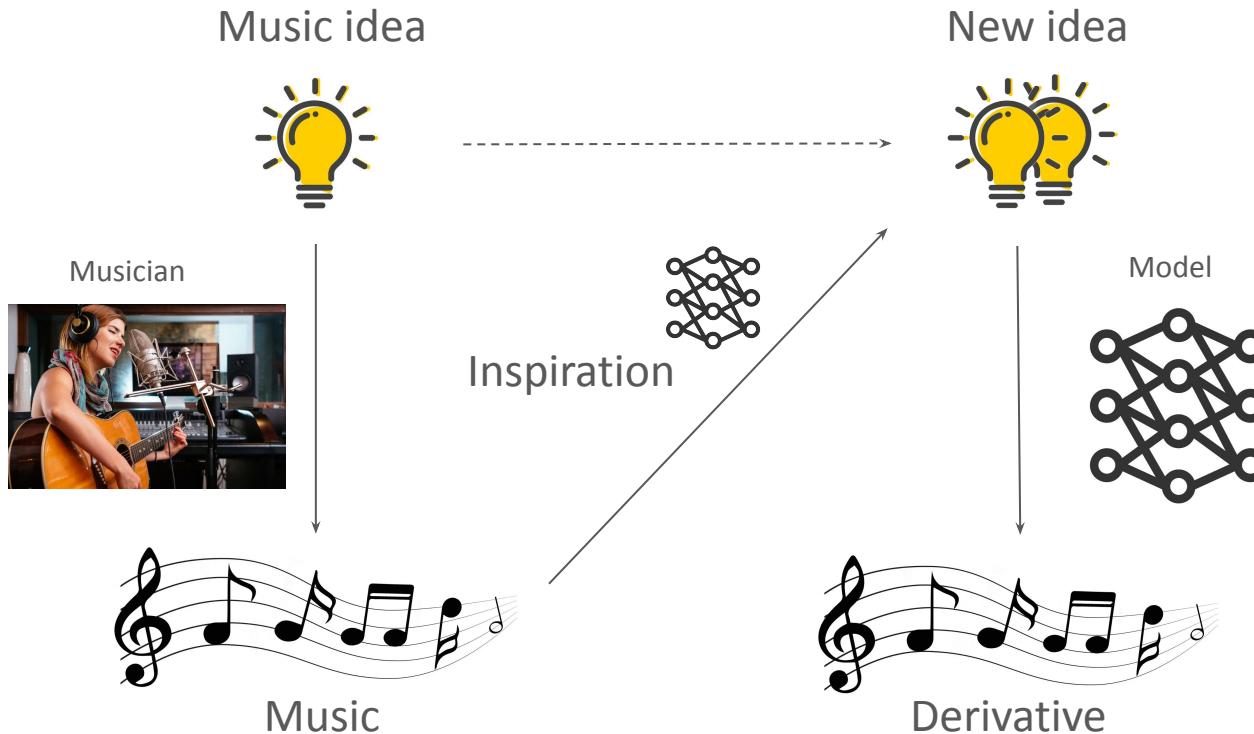
An adaptation in music, which falls under the category of *Derivative Work*, refers to the process of altering or modifying an existing musical piece to create a new version.

This creative process can take various forms, such as changing the **lyrics, melody, rhythm, or arrangement** of the original composition. It can also include *translating a song into another language or reshaping it to match a different musical genre or style*.

Create a music adaptation



Create a music adaptation with deep learning



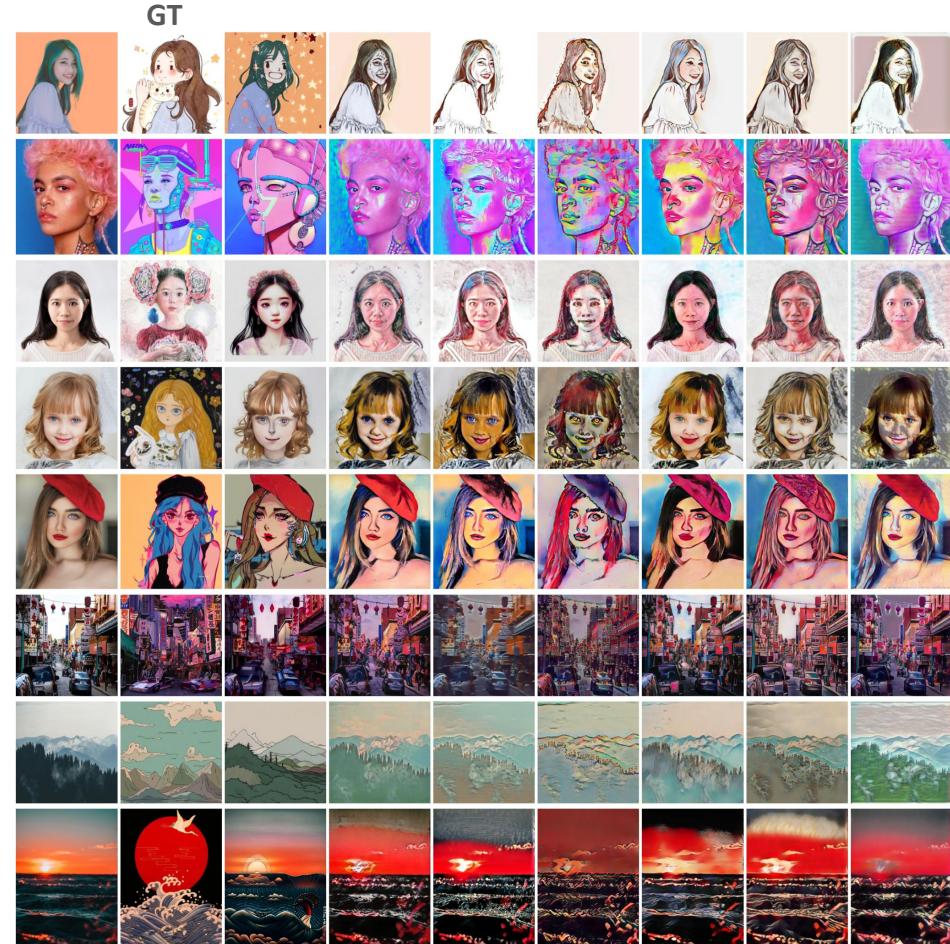
Challenges in cover song generation

- Data acquisition and preparation
 - Difficulty in collecting paired data.
 - Specifically, ensuring proper data alignment.
- Preserving the core musical idea
 - Defining the essential musical idea.
 - Establishing a measurable definition for "sounds similar".

Music Style Transfer

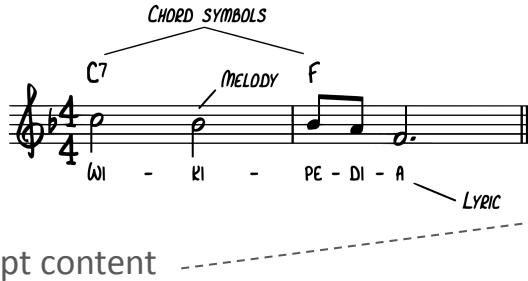
Image style transfer

We take the **edges** of the content image and dress them up in the texture of the style image



What's content and style of music?

- Content
 - Melody
 - Chord
- Style
 - Everything except content



Genre



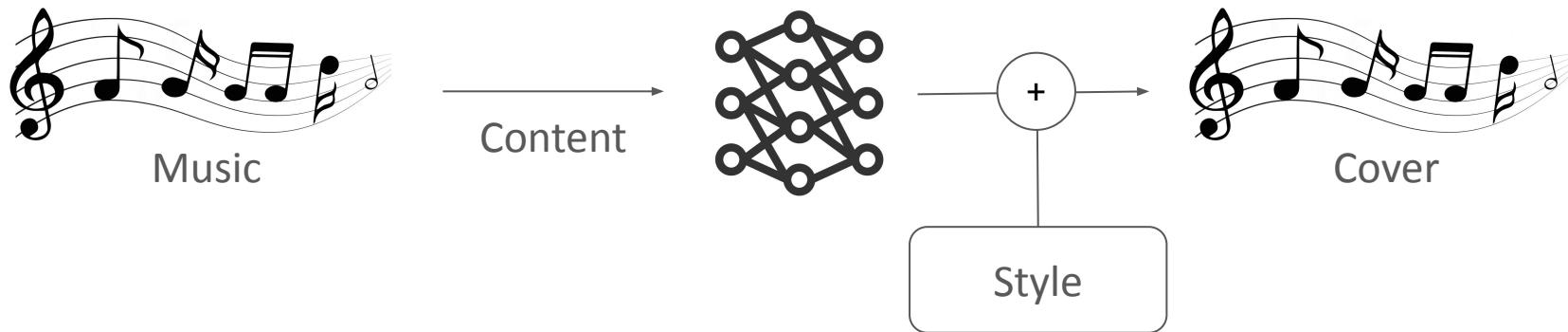
Instrument



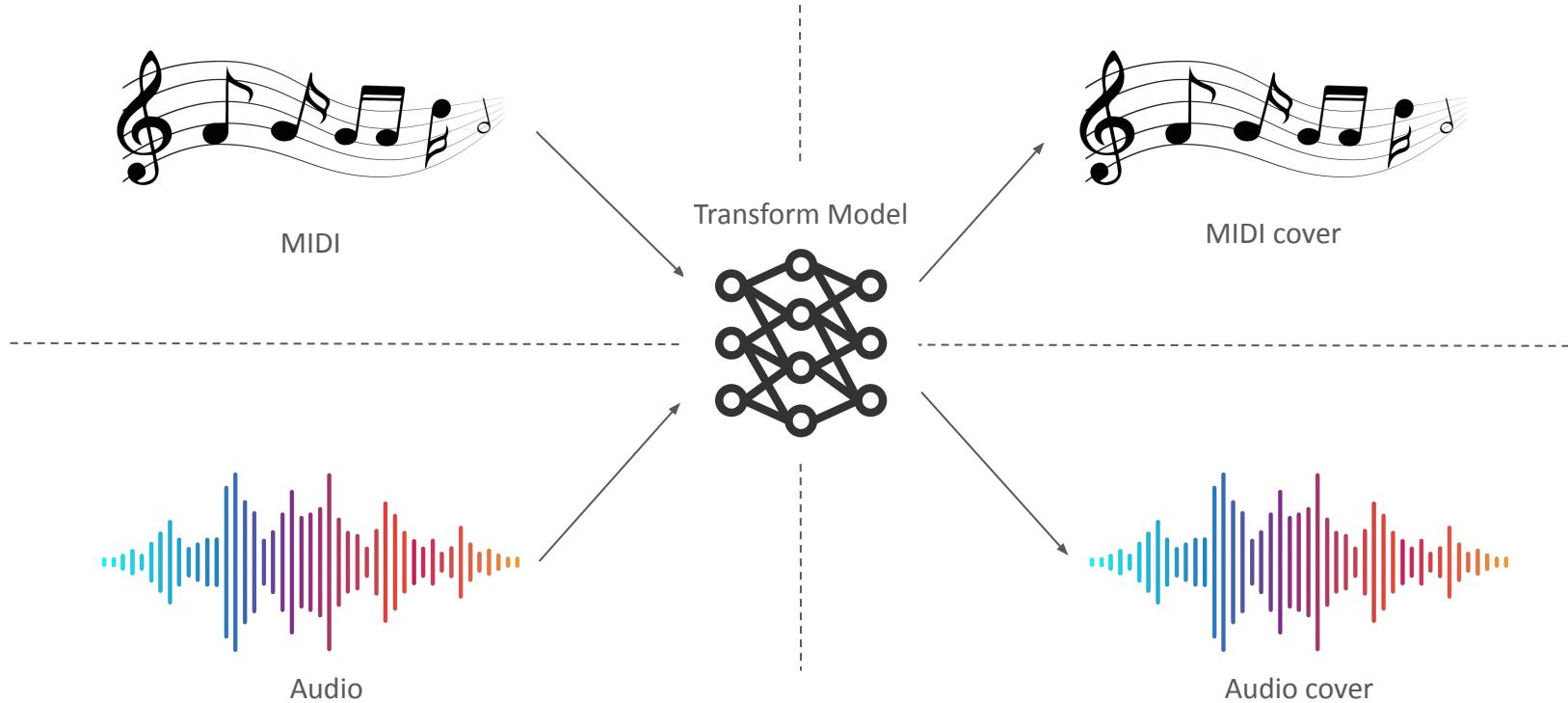
Arrangement

Problem formulation

- Cover Song Generation is formulated as a problem of ***Content-Preserving Style Transfer***
 - Extracting the content (core **melody**, **harmony**, and **structure**) from the original source track.
 - Replacing or re-rendering the expressive and temporal style (rhythm, timbre, articulation, dynamics) of the content with that of a target style (e.g., genre or artist).



Cross-modal music transformation



Music Arrangement

Music production workflow

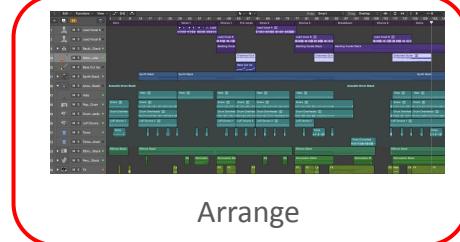
The process of creating a final musical piece involves a sequence of creative and technical stages, transforming an initial idea into a polished, publishable track.



Compose



Record



Arrange



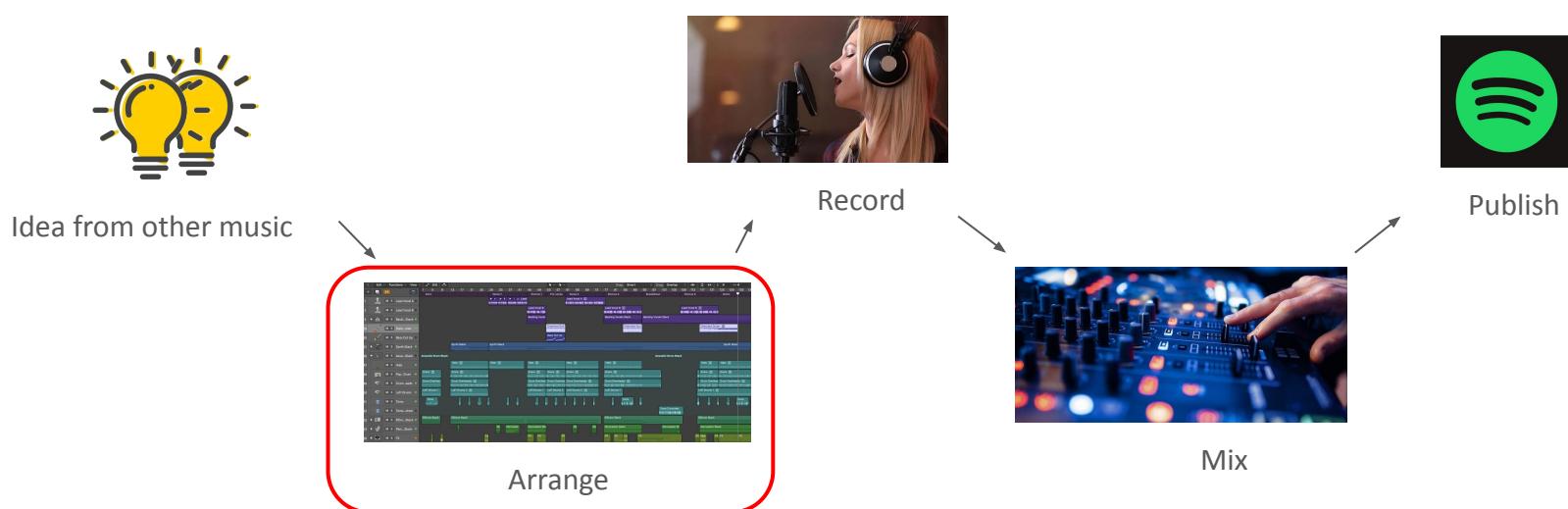
Mix



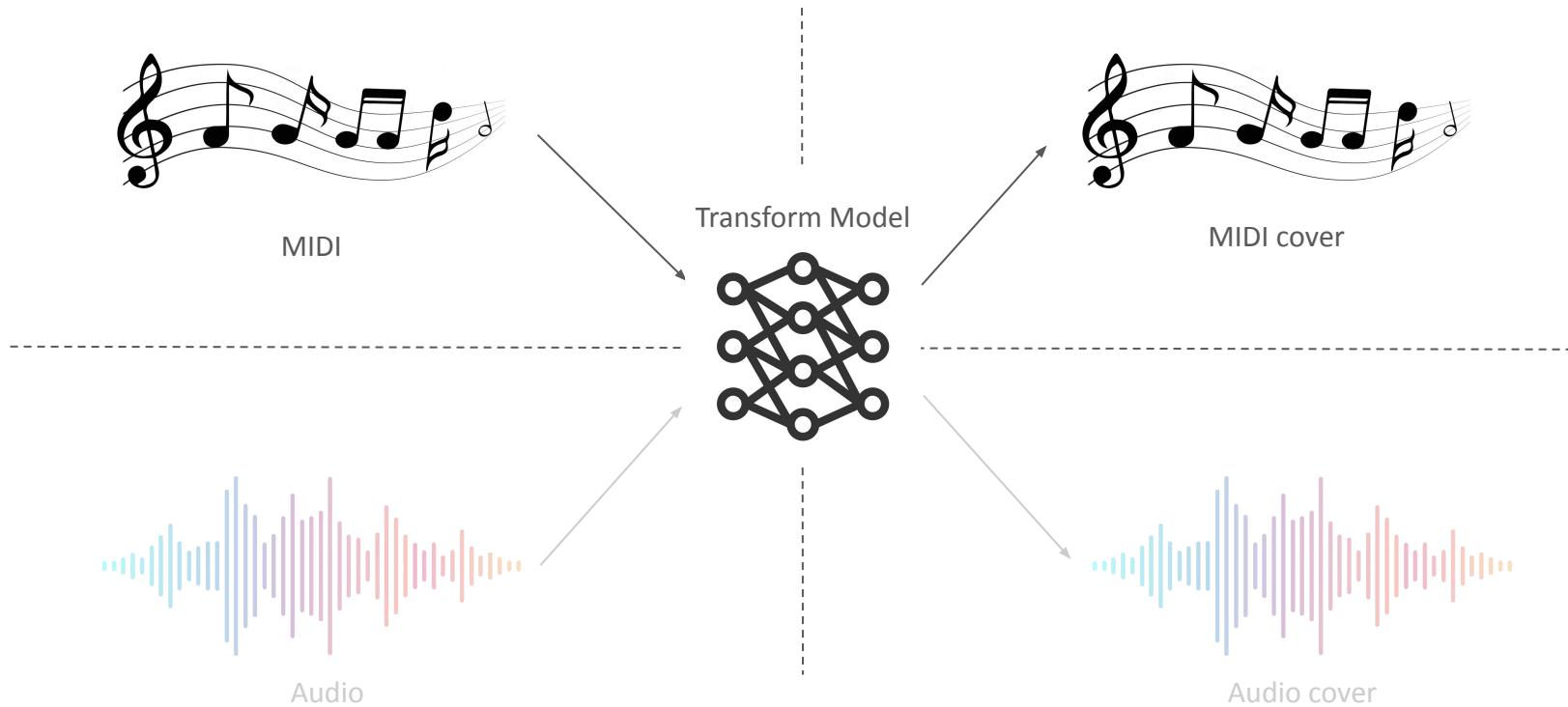
Publish

Music adaptation workflow

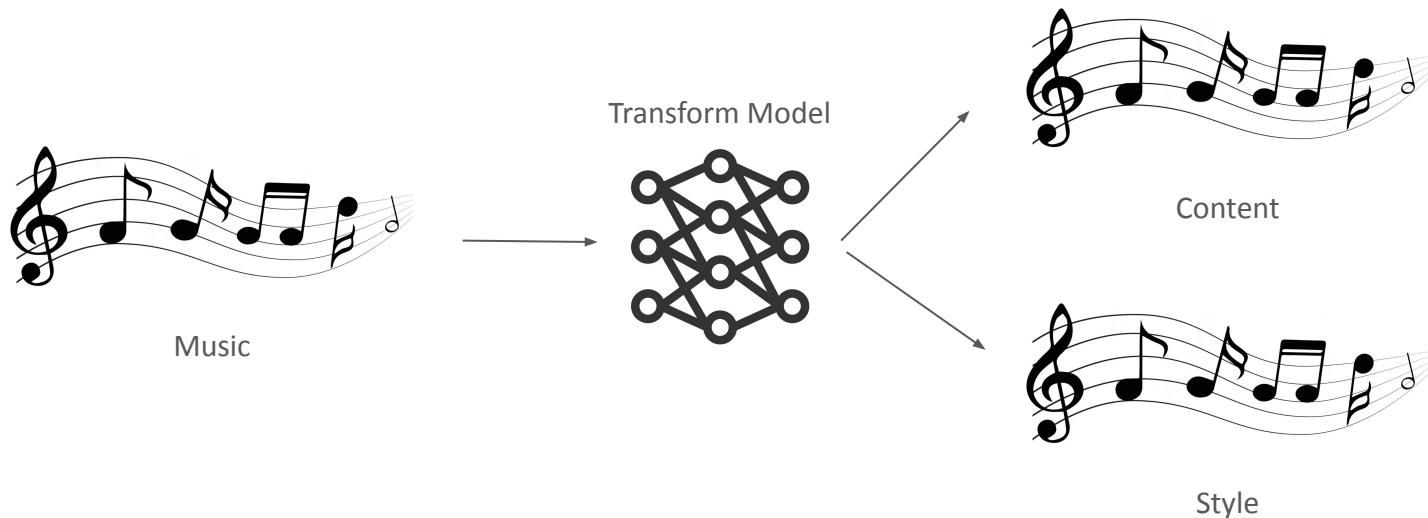
This workflow represents the creative process of developing a new piece of music *inspired by or derived from* existing material, rather than starting from an original composition.



MIDI to MIDI

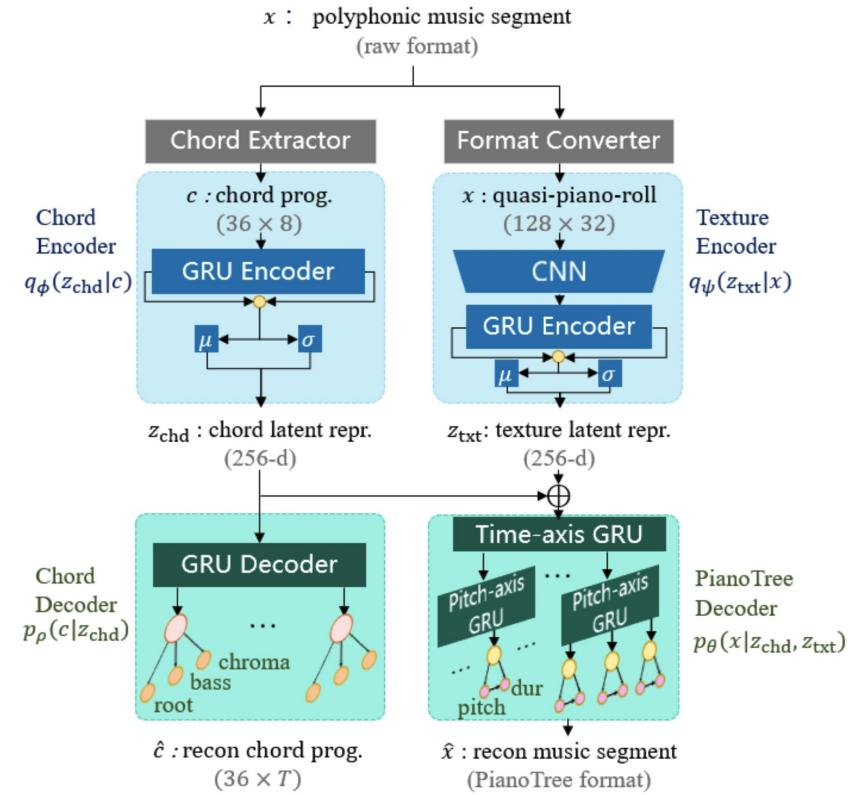


Task: disentangle music content

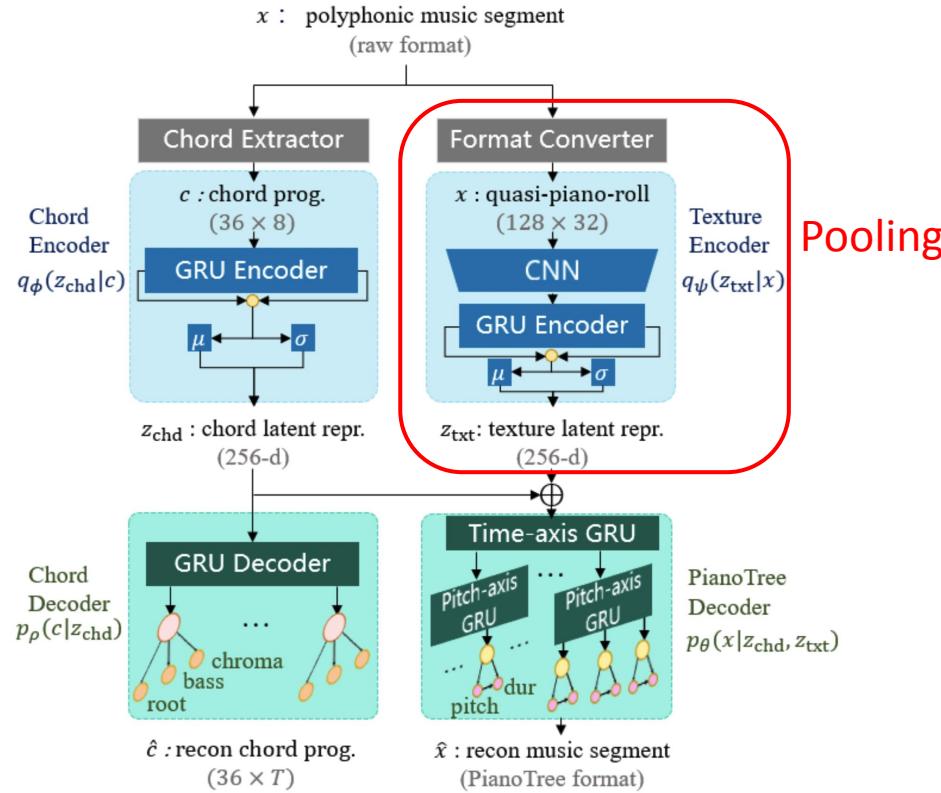


Controllable polyphonic music generation

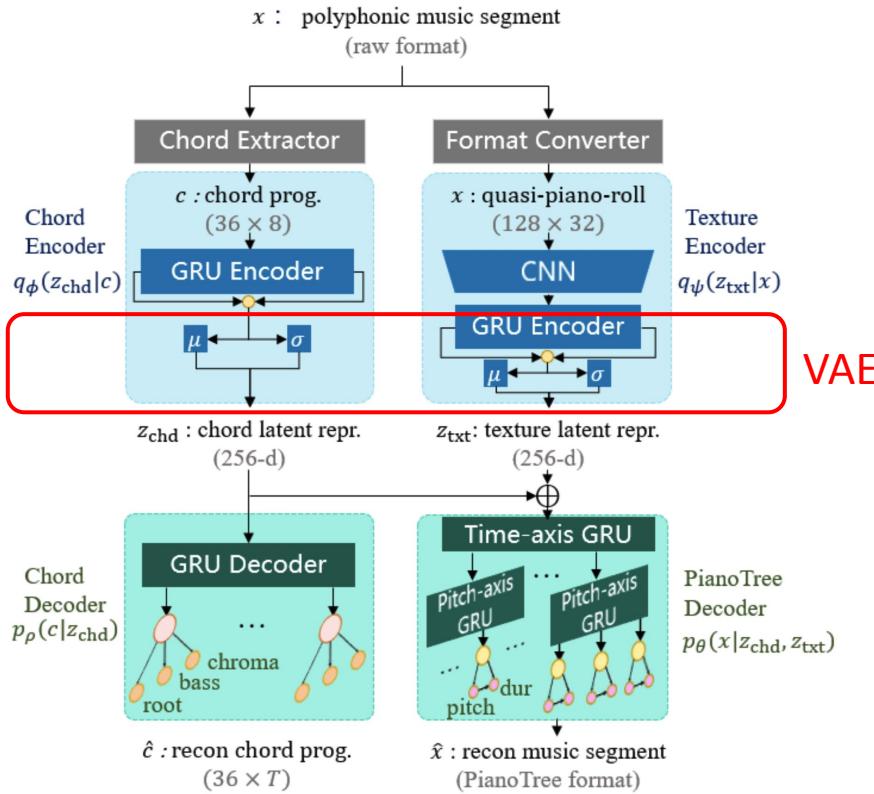
- Disentangle chord (content) from music
 - Melody is not considered as content



Controllable polyphonic music generation



Controllable polyphonic music generation



Controllable polyphonic music generation (VAE)



6 6 6 6 6 6 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
q 4 4 4 2 2 2 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0
q 2 2 2 2 2 2 2 5 5 6 0 0 0 0 0 0 0 0 0 0 2
q 4 2 2 2 2 2 2 3 3 5 5 6 0 0 0 0 0 0 0 2
q 9 4 2 2 2 2 2 3 3 3 3 5 5 5 5 5 5 5 3 2
q 9 9 4 2 2 2 2 3 3 3 3 3 5 5 5 5 5 5 3 2
q 9 9 9 4 2 2 2 3 3 3 3 3 5 5 5 5 5 5 3 2
q 9 9 9 9 9 8 3 3 3 3 3 3 5 5 8 8 2
q 9 9 9 9 9 8 3 3 3 3 3 3 8 8 8 8 2
q 9 9 9 9 9 8 8 8 8 8 8 8 8 8 8 8 2
q 9 9 9 9 9 8 8 8 8 8 8 8 8 8 8 8 2
q 9 9 9 9 9 8 8 8 8 8 8 0 0 6 6 0 0 5 5 2
q 9 9 9 9 9 8 8 8 8 8 8 6 6 6 6 6 6 5 5 2
q 9 9 9 9 9 9 9 9 9 9 9 6 6 6 6 6 6 6 5 5 2
q 9 4 4 4 9 9 9 9 9 9 9 6 6 6 6 6 6 6 5 5 2
q 9 4 4 4 9 9 9 9 9 9 9 6 6 6 6 6 6 6 6 2 2
q 9 4 4 4 9 9 9 9 9 9 9 6 6 6 6 6 6 6 6 2 2
q 9 9 9 9 9 9 9 9 9 9 9 1 1 1 1 1 1 1 1 1 1
q 9 9 9 9 9 9 9 9 9 9 9 1 1 1 1 1 1 1 1 1 1
q 7 7 7 7 7 7 7 1 1 1 1 1 1 1 1 1 1 1 1 1

Controllable polyphonic music generation

However, the melodic fidelity of the generated output is poor.

C \sharp minor, rock style Right-hand: Syncopated alto melody throughout

Left hand: 8th note arpeggio

Change to 16th note arpeggio

Two-voice melody

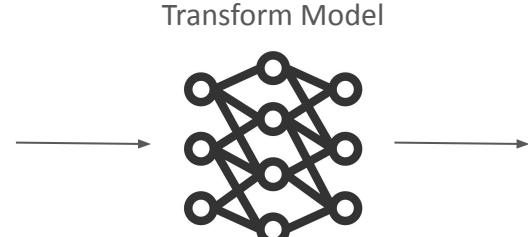
2 characteristic cut-offs

(a) A real piece.

Task: piano accompaniment generation



Melody + Chord

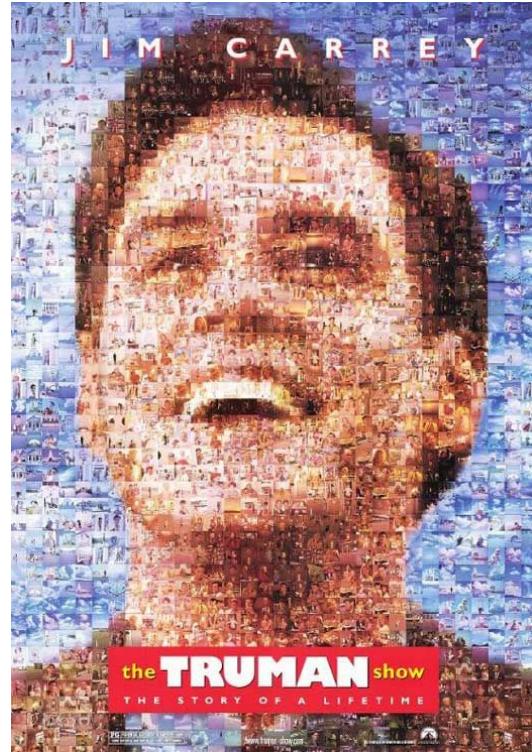


A piano score enclosed in a blue border. The top staff is for the treble clef part, and the bottom staff is for the bass clef part. The score consists of ten measures. Measure 1 starts with a dynamic 'mp'. Measures 2-4 show a continuation of the melody with some eighth-note patterns. Measures 5-8 show more eighth-note patterns, with measure 6 featuring a dynamic 'p'. Measures 9-10 conclude the piece. A small piano icon is positioned above the score.

Piano accompaniment

AccoMontage (Accompaniment Montage)

"A montage is a film editing technique in which a series of short shots are sequenced to condense space, time, and information." - Wikipedia



AccoMontage (POP909)

Melody

Melody2

Accompaniment

Structure

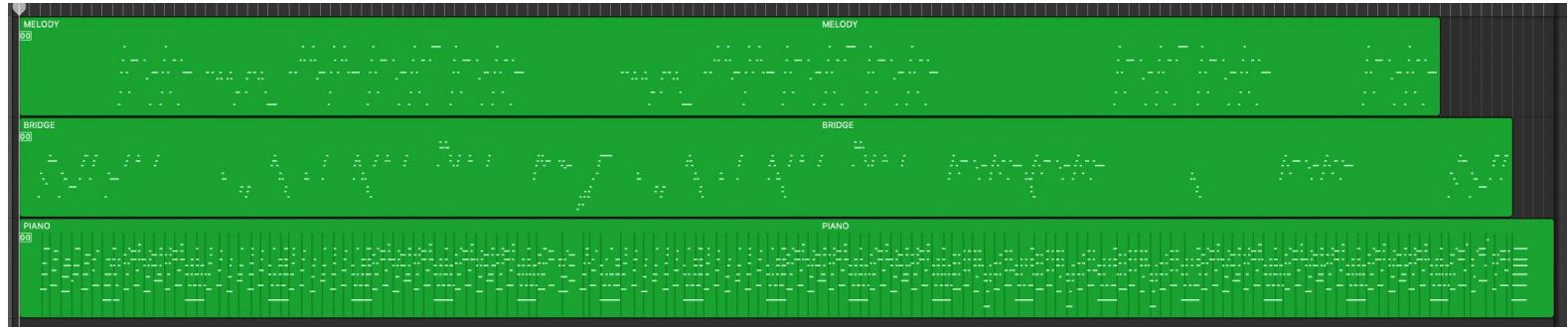
A1

A2

B1

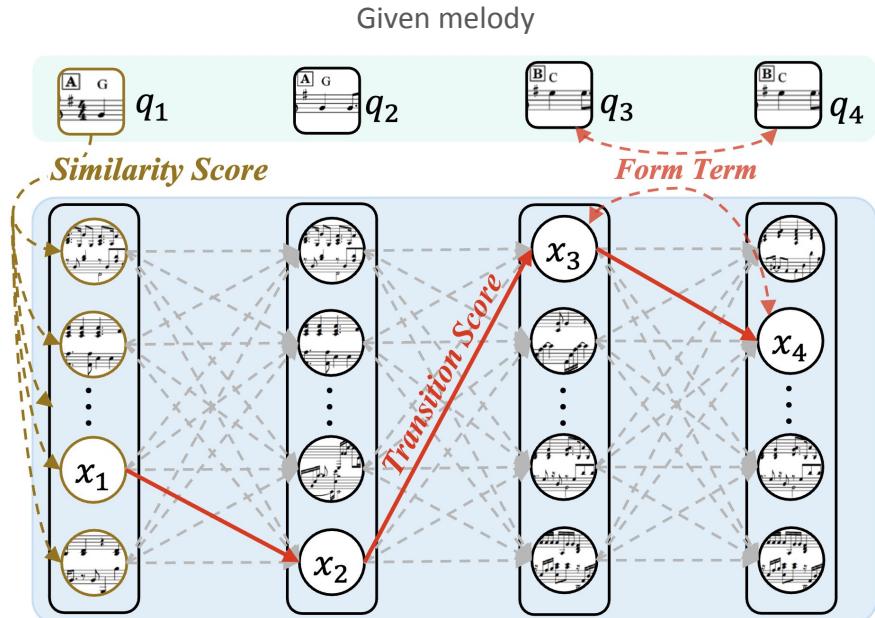
B2

.....



AccoMontage

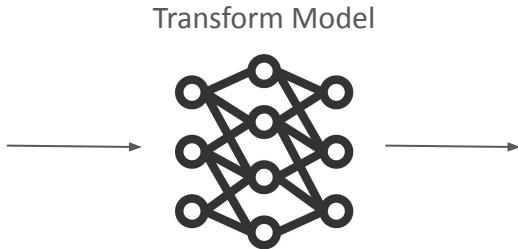
- Optimization
 - **Chord/Similarity Score:** Measures the fitness between the target melody query and the candidate accompaniment phrase
 - **Transition Score:** Measures the smoothness and musicality of the connection between two successive accompaniment phrases
 - **Form Term:** Enforces constraints on the overall musical structure and arrangement (e.g., song sections A, B, C).



Task: multitrack arrangement



Melody + Chord



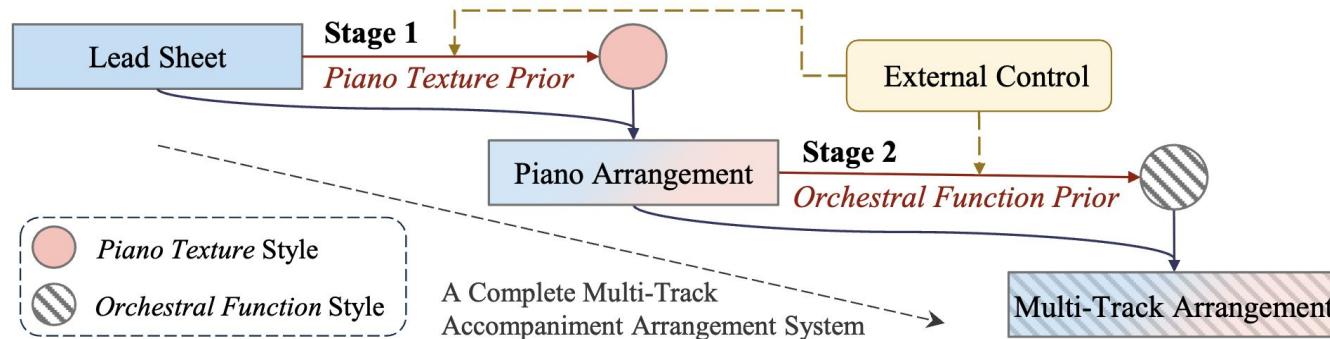
Transform Model



Multitrack arrangement

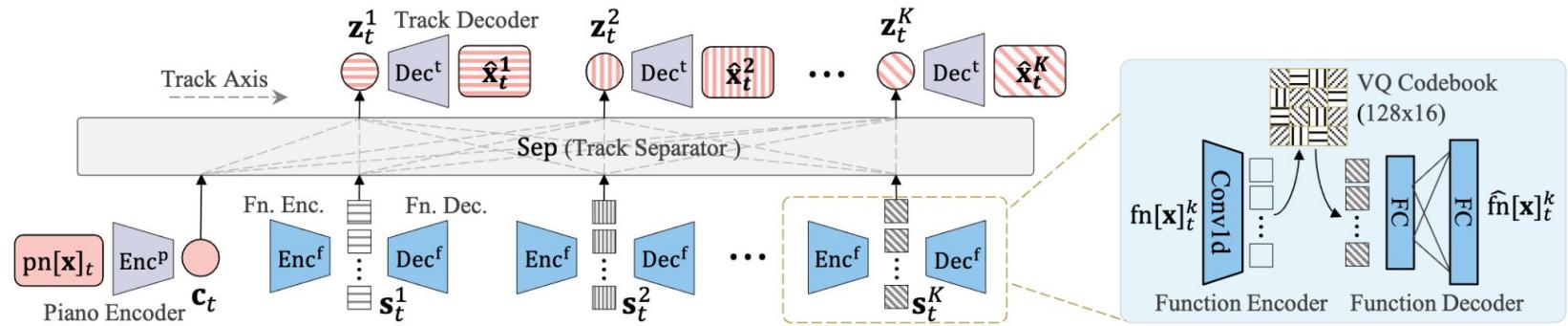
Multi-Track accompaniment arrangement

Concept 1: Arrange multi-track accompaniments from a simple lead sheet.



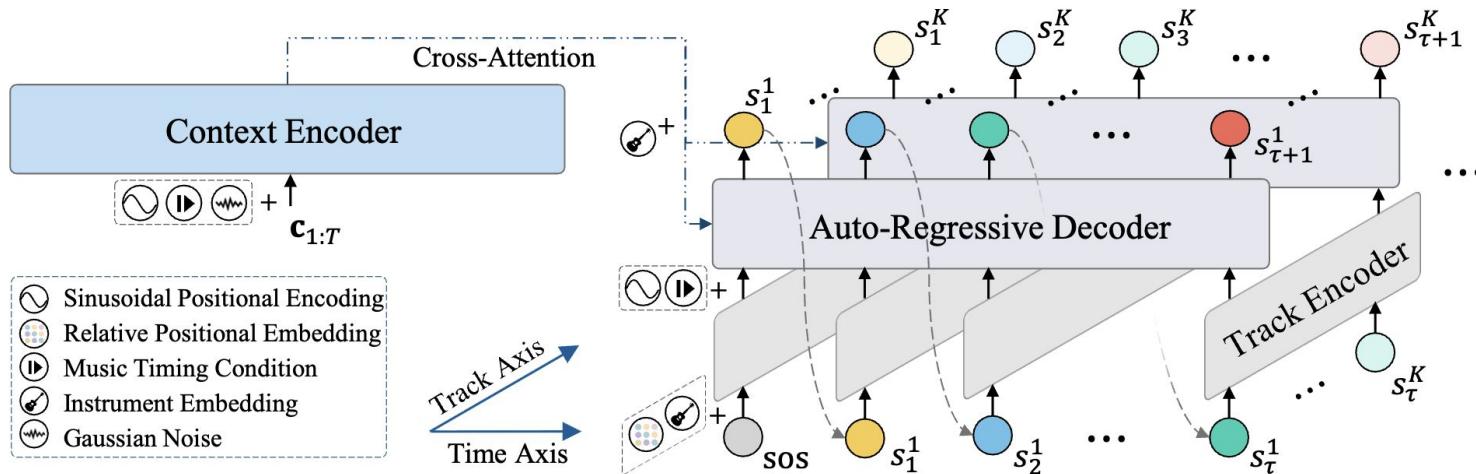
Multi-Track accompaniment arrangement

Concept 2: Generate each track by deriving the musical content from the **reduced piano arrangement** and assigning its instrumental role using a **quantized orchestral function**.

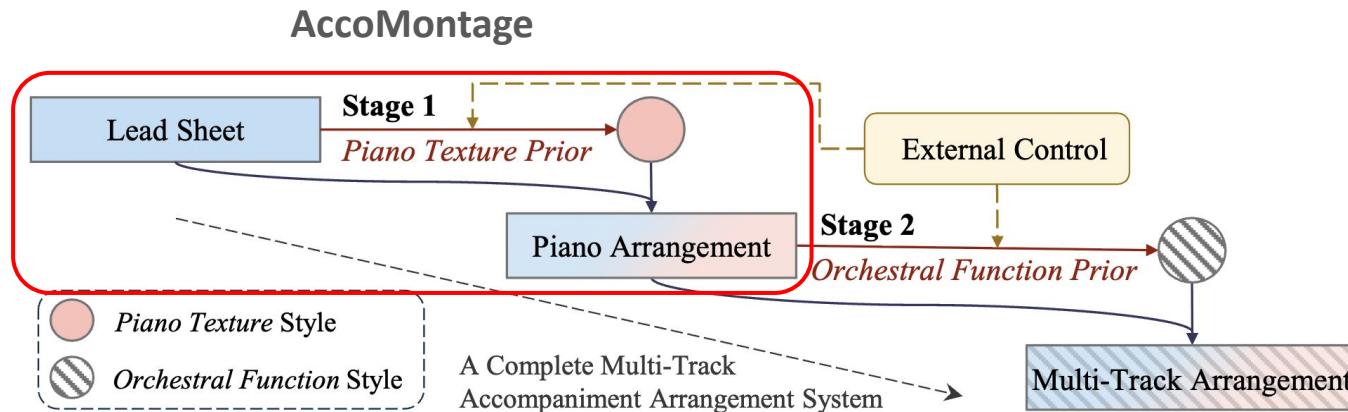


Multi-Track accompaniment arrangement

Model: Generate a complete orchestral score via the automatic orchestration of a piano reduction.



Multi-Track accompaniment arrangement



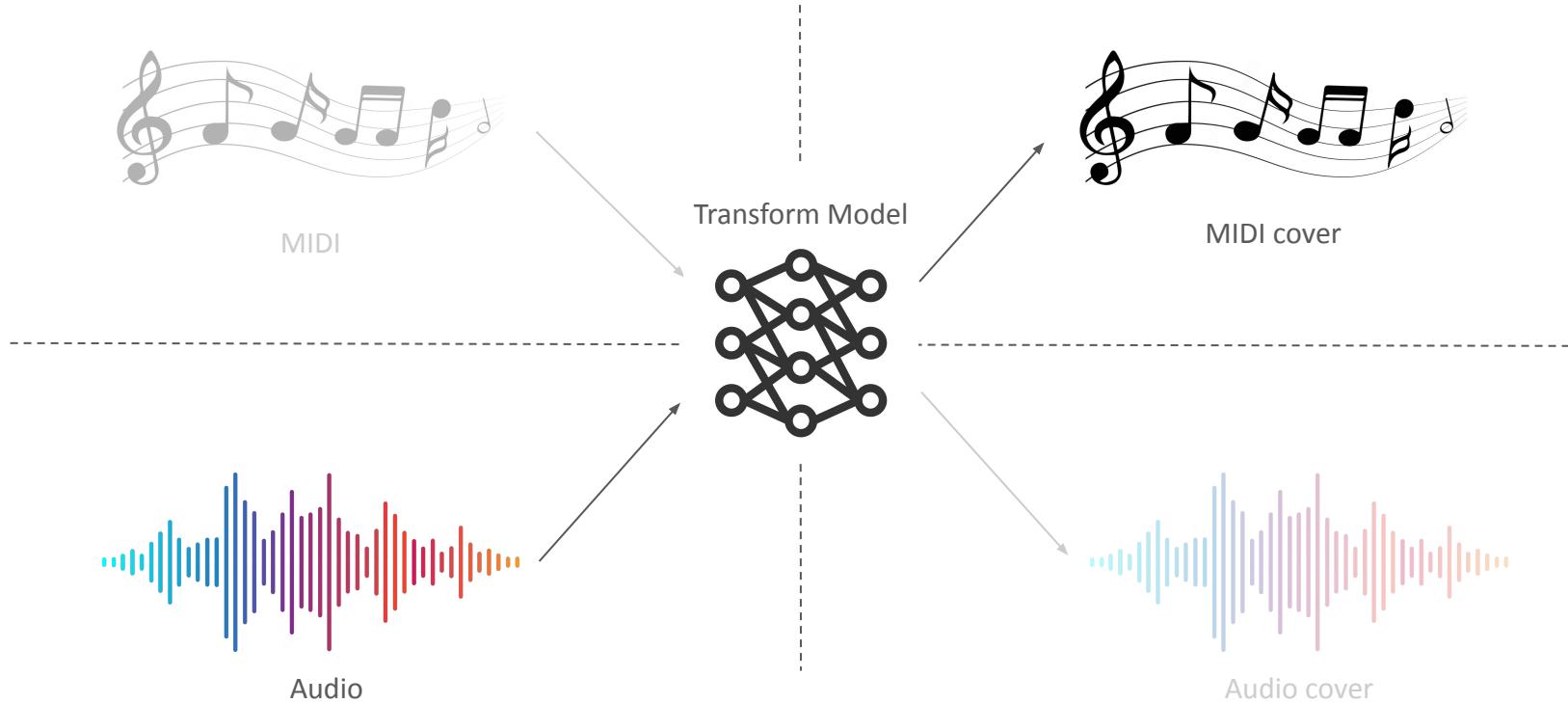
We didn't mention how to get content ...

- Pop909 is a *perfect* dataset
 - It may not fully represent real-world musical complexity
- Melody extraction from MIDI is easier than from raw audio.



Cover (Adaptation) Generation

Audio to MIDI

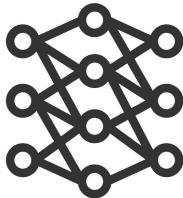


Task: piano cover generation



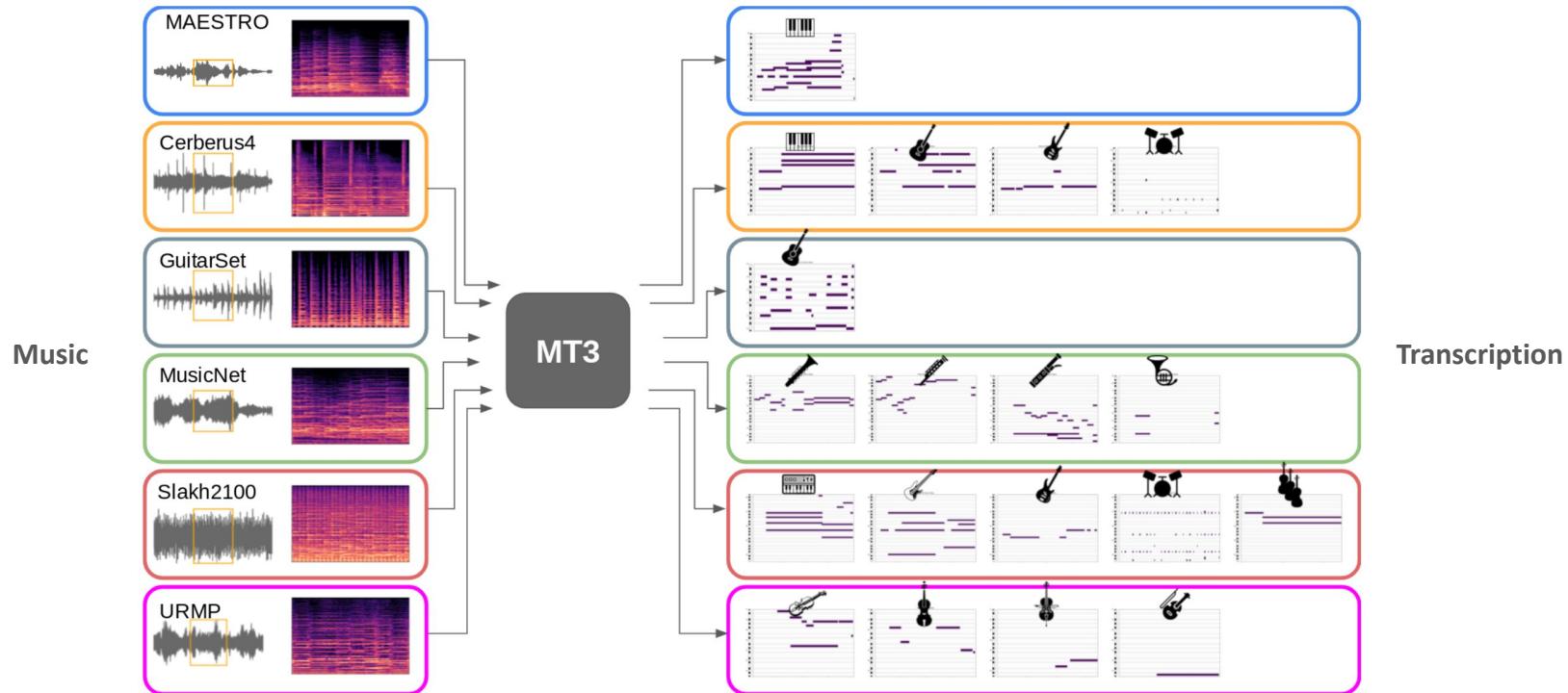
Song audio

Transform Model

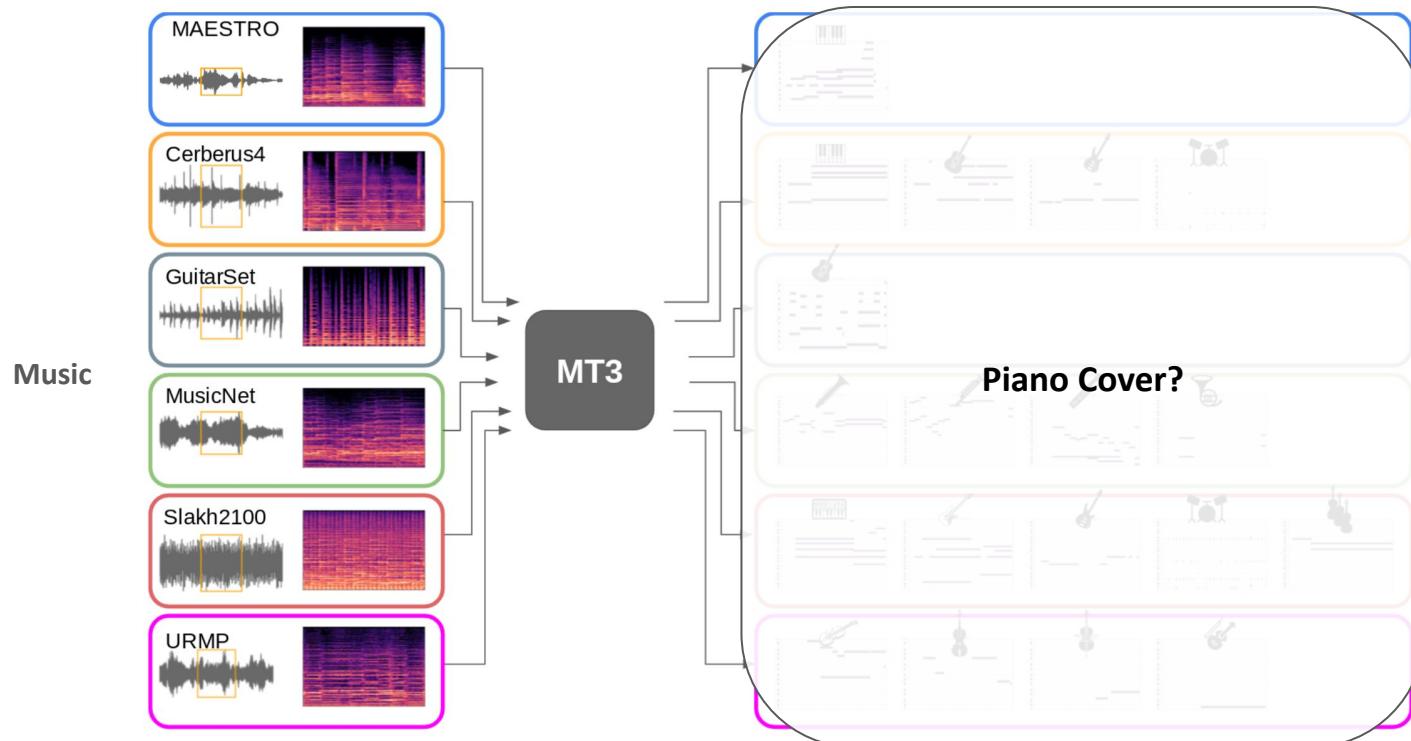


Piano cover

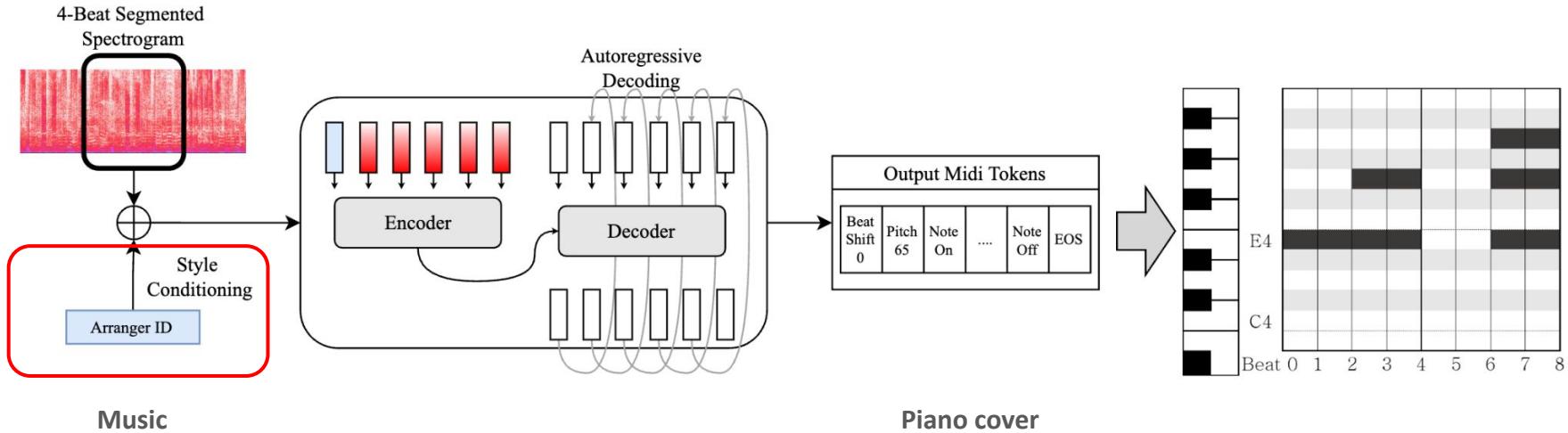
Pop2Piano (MT3)



Pop2Piano (MT3)



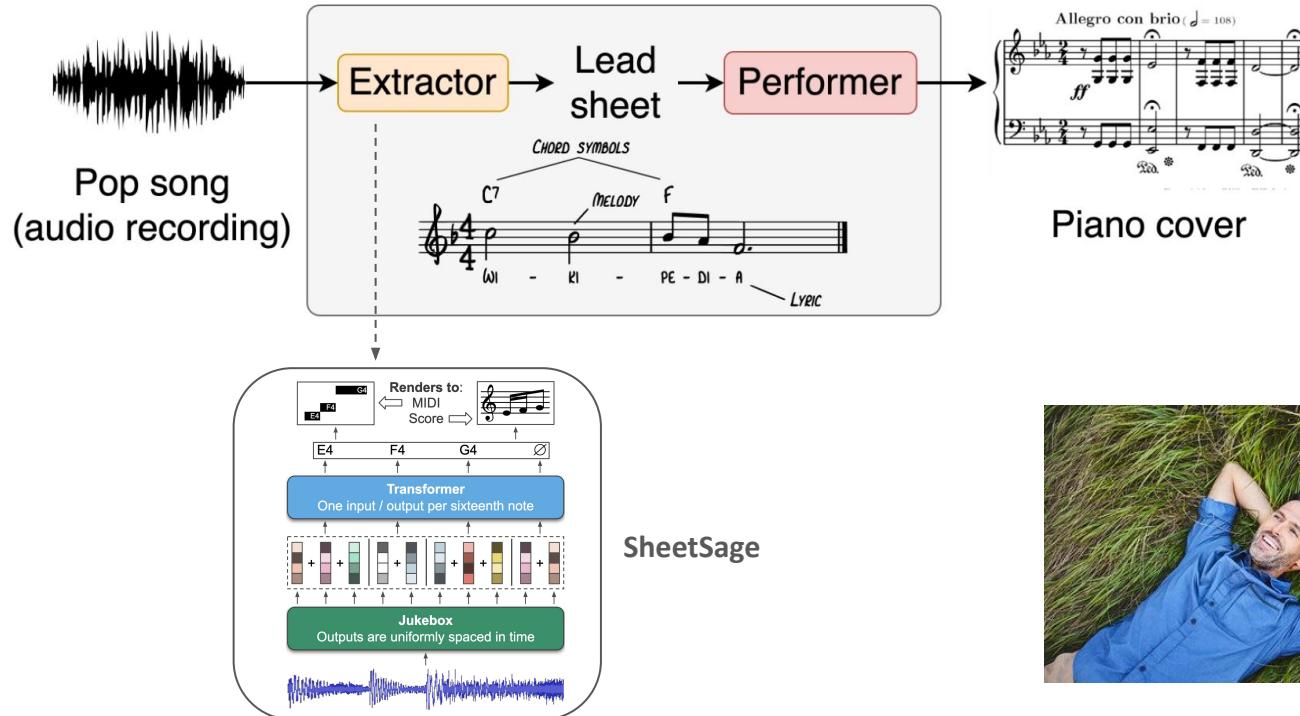
Pop2Piano



Rethink: cover generation = extraction + generation



PiCoGen1



Content preservation

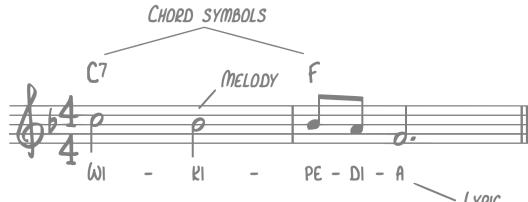
The Question: Should the core elements (melody, harmony, rhythm) be preserved exactly when generating a cover?

- No. We aim for **perceptual similarity**, not strict note-for-note identity. The content should be **adaptable** to the new style.



What's content and style of music?

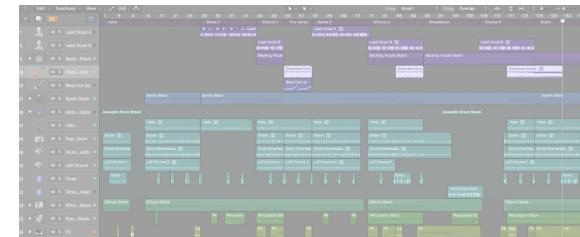
- Content
 - Melody
 - Chord
- Style
 - *Everything except content*



Genre



Instrument



Arrangement

Preservation of content and style

The Challenge of Content: to what extent must the **original melody** be altered (or preserved) to fit a target?

- Specifically: how do we define the "**essential content**" of the melody?
- The detail: which notes are **mandatory** for recognition, and which can be changed for embellishment?

The Challenge of Style: beyond the core melody and harmony, which **other musical attributes** should be preserved in the cover song?



What is the model?

Model = Deep Network

遇到問題用 deep learning 「硬 train 一發」就對了

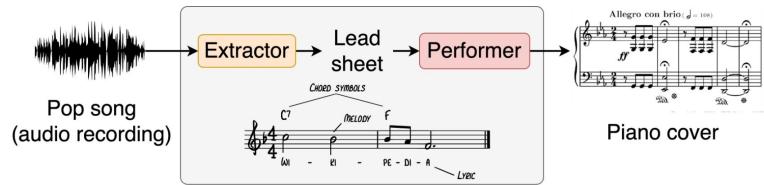


How about data-driven approach?

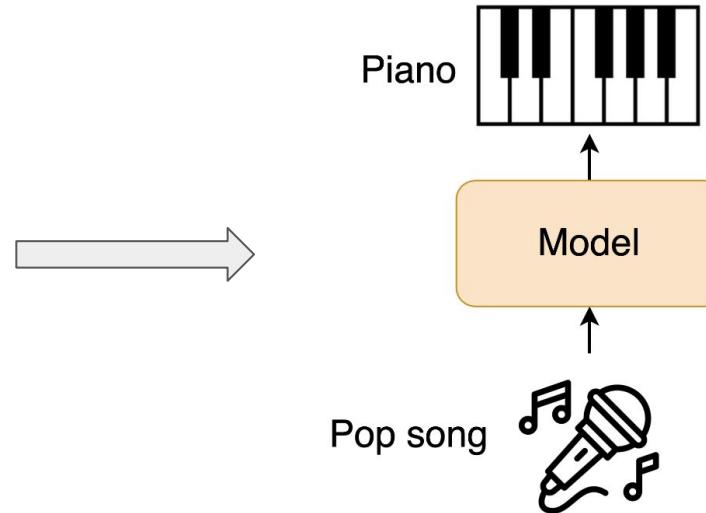
".....不知道什麼是'硬train一發'的，我試著解釋一下什麼是'硬train一發'。這個，'硬train一發'啊...它是一個非常神奇的東西：它是一種信念，它是一種夢想，它是一種浪漫，它是人類最原始的衝動，它是**恆古以來人類的目標**....." - by 李宏毅

PiCoGen2

We aim to build a complete, end-to-end system based on the PiCoGen1.



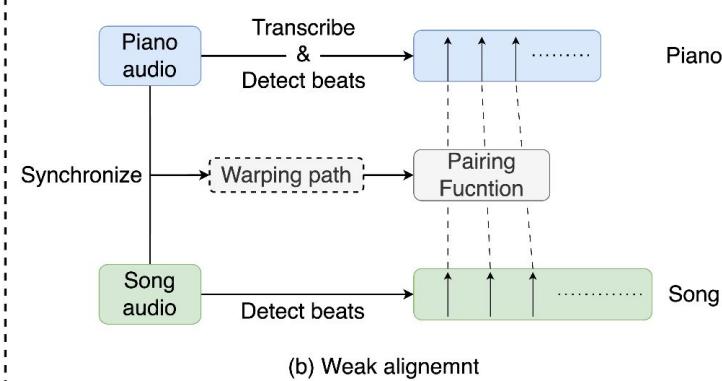
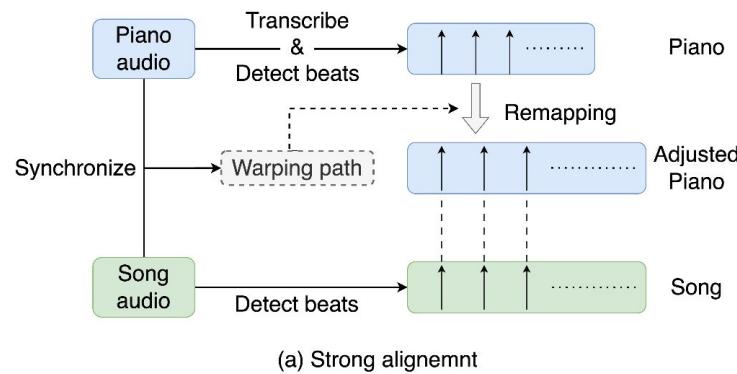
PiCoGen1



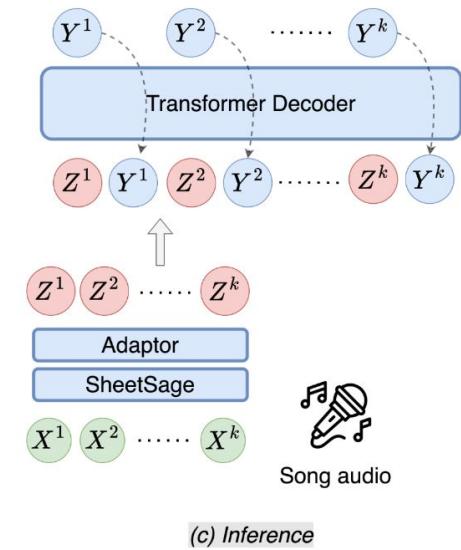
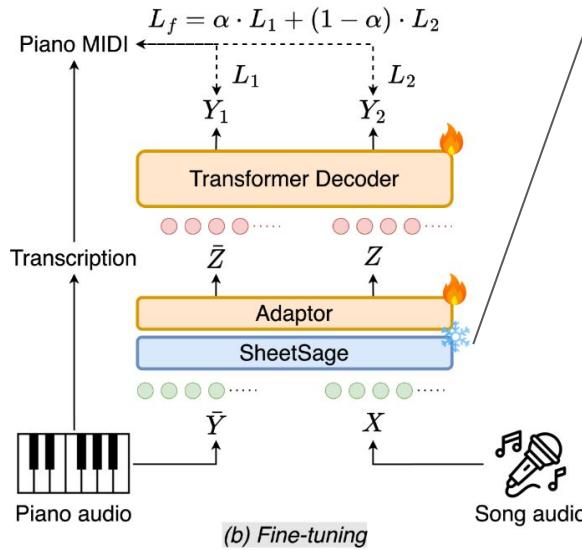
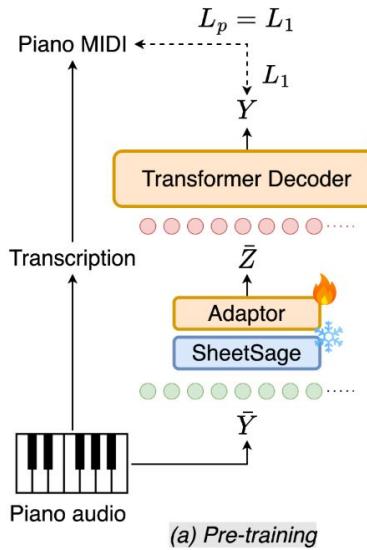
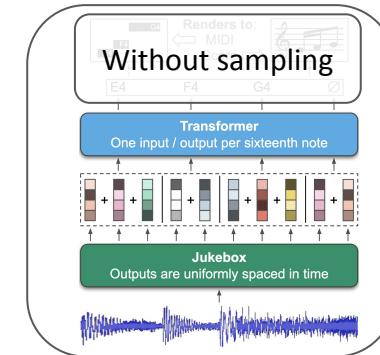
PiCoGen2

PiCoGen2 - data

Training data is synchronized based on the musical **beat** (quantized time), rather than raw **time** values.



PiCoGen2 - training & inference



PiCoGen2 - demo



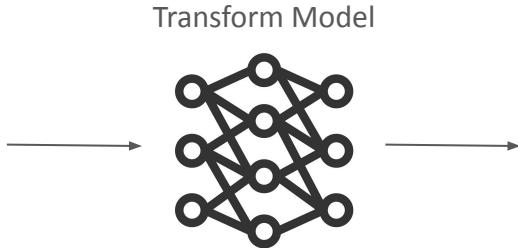
PiCoGen2

Piano cover generation with transfer learning approach and weakly aligned data

Task: multitrack cover



Song audio



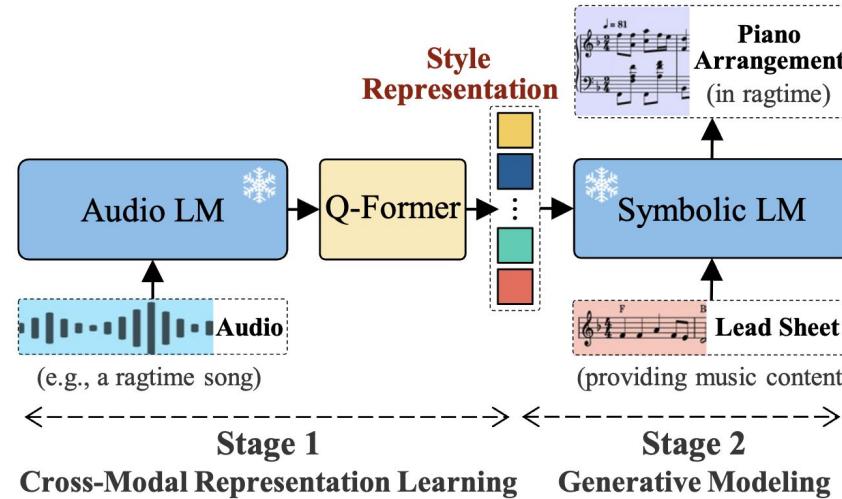
Transform Model



Multitrack cover

BOSSA

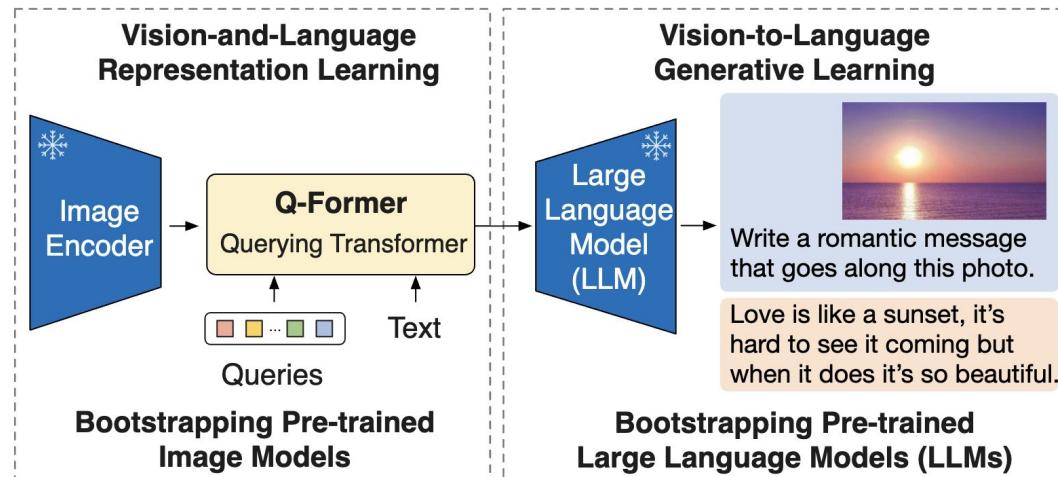
Learn a **style representation** (e.g., "ragtime") from an audio example and then use that compact style vector to condition a **symbolic music generation model**



BOSSA (BLIP-2)

BLIP-2 bootstraps vision-language pre-training from off-the-shelf frozen pretrained image encoders and frozen large language models.

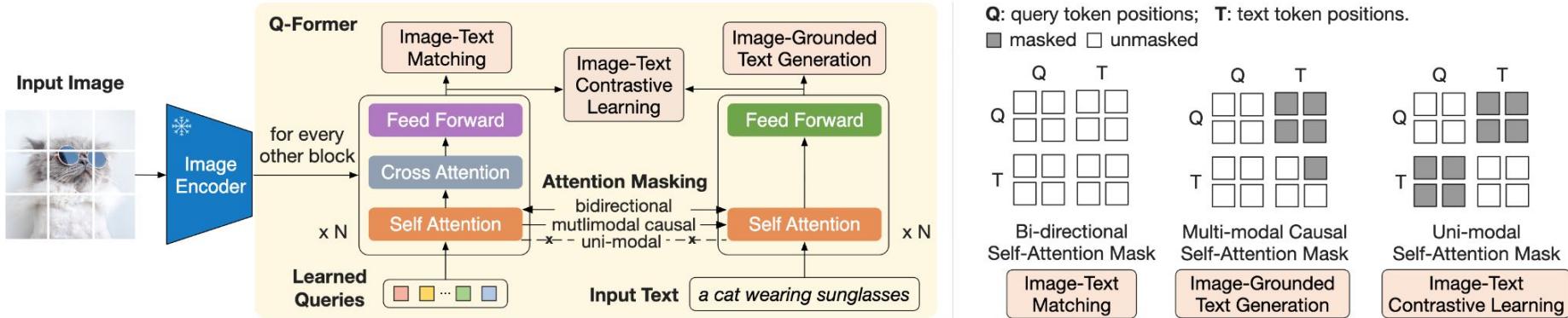
BLIP-2 bridges the modality gap with a lightweight **Querying Transformer**, which is pretrained in two stages.



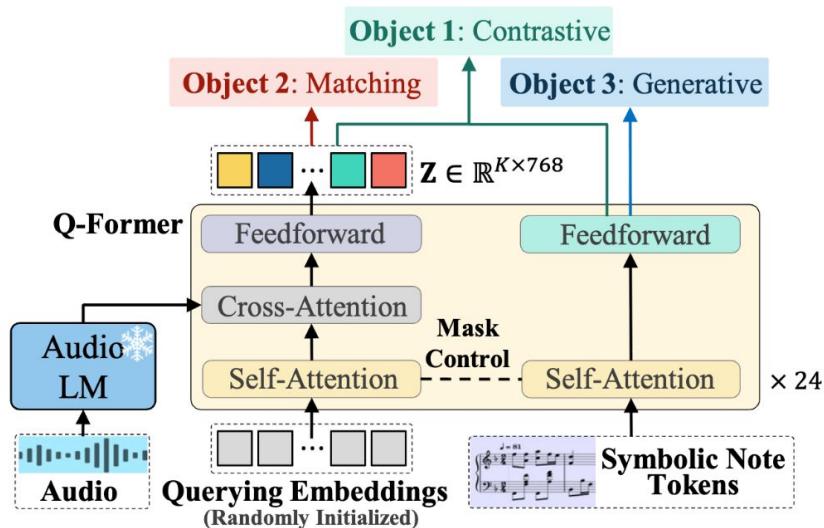
BOSSA (BLIP-2)

Optimize three objectives which enforce the queries (a set of learnable embeddings) to extract visual representation most relevant to the text.

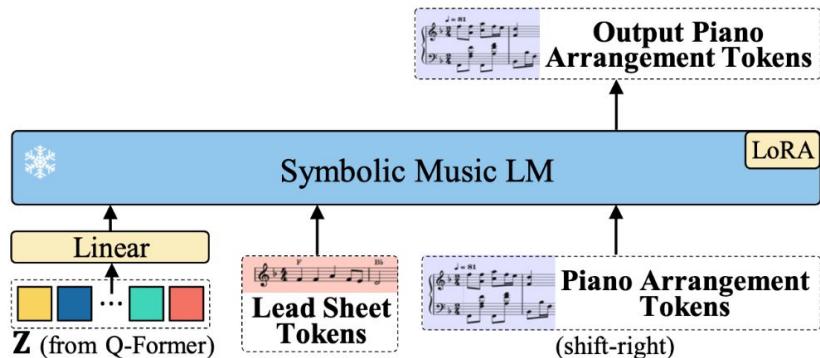
The self-attention masking strategy for each objective to control query-text interaction.



BOSSA

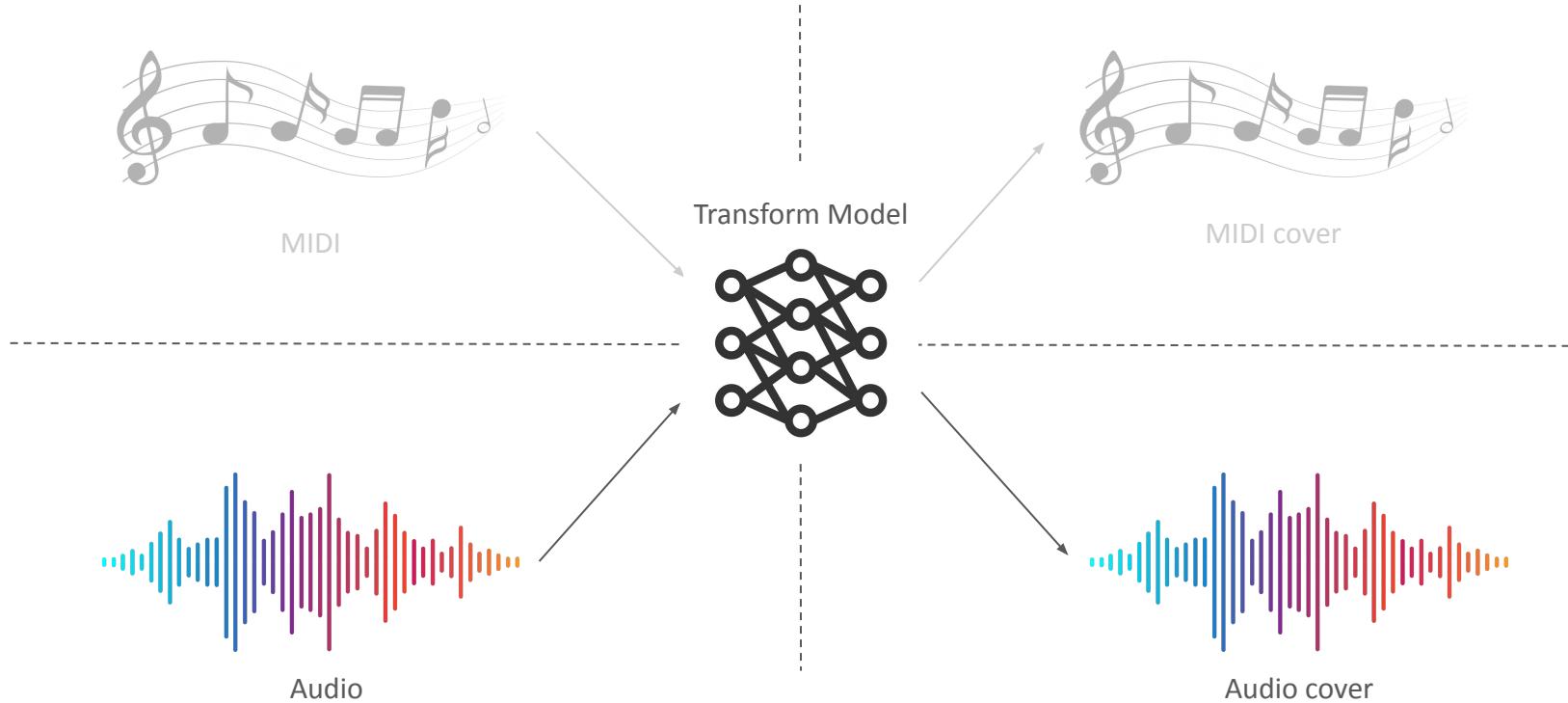


(a) Stage 1: Cross-modal learning with Q-Former.

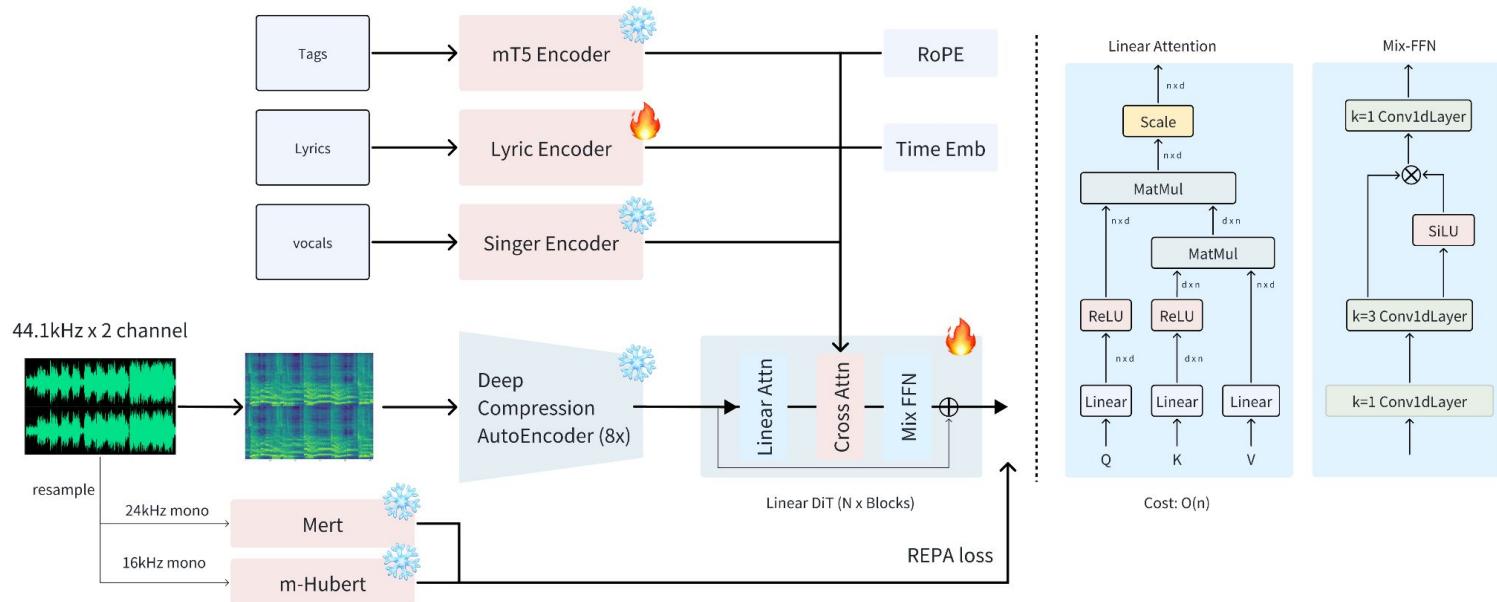


(b) Stage 2: Audio-to-symbolic arrangement.

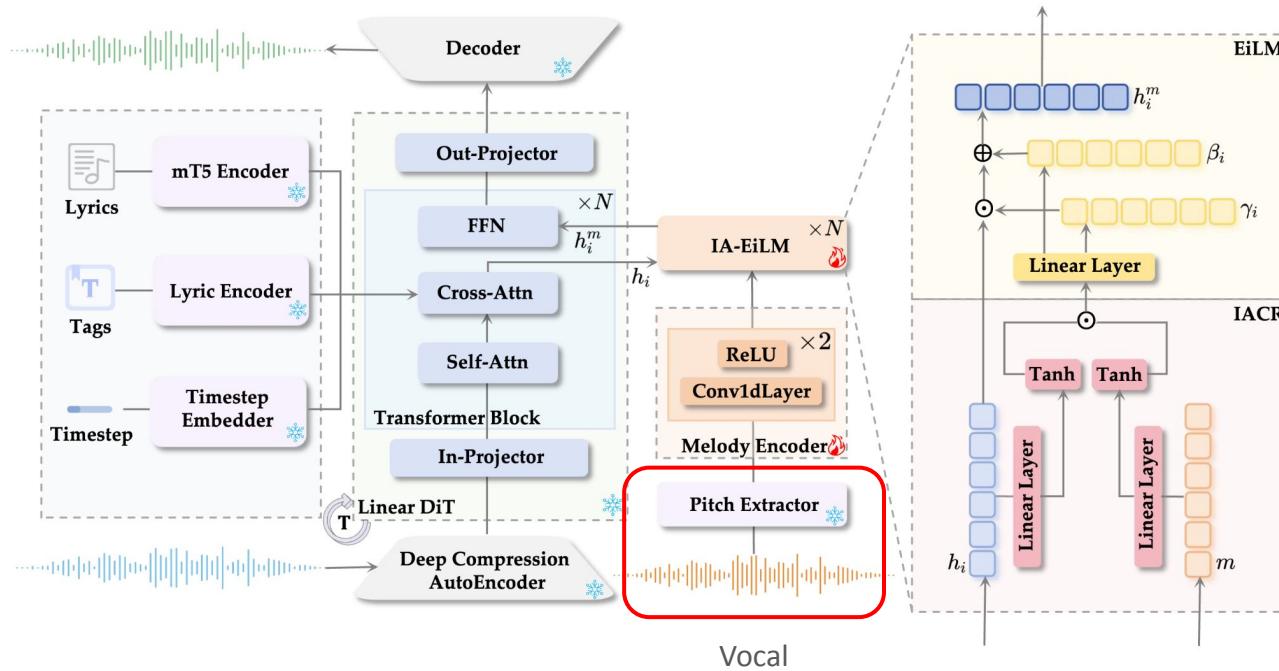
Audio to Audio



SongEcho (Ace-Step)



SongEcho



Evaluation - what a mess

What happened?

- The musical concept of "style" is not well-defined, leading to inconsistent model objectives.
- Evaluation metrics suffer from strong **personal biases**.
 - A high score on technical similarity is **not equivalent** to a good musical cover (a cover needs interpretation, not cloning).
 - Similarly, high overall quality does **not guarantee** a good cover (it needs to sound like the original and the new style).
- We currently rely on metrics borrowed from general music generation, which are **fundamentally unsuited for cover generation**.



Evaluation of cover song generation

There is still much room for improvement



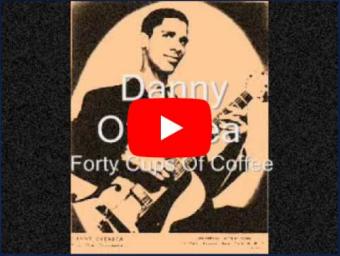
Resource

SecondHandSongs

SecondHandSongs EXPLORE DISCUSS PARTICIPATE PLAY SIGN IN / AD-FREE

Search Find cover songs, artists and more GO DETAILED SEARCH

Discover The Original



Danny
Overbea
Forty Cups Of Coffee

VS

SHUFFLE

40 Cups of Coffee
by Danny Overbea with King
Kolax and His Orchestra

Listen to The Cover



DECCA
RECORDS
RECORDED BY DECCA RECORDS INC., NEW YORK
RECORDED AND PRODUCED BY BILL HALEY
AND HIS COMETS
RECORDED NO. 9-30214
(100703) [3-10]
FORTY CUPS OF COFFEE
DANNY OVERBEA
BILL HALEY
AND HIS COMETS

Forty Cups of Coffee
by Bill Haley and His Comets

HookTheory

The image shows a screenshot of the Hookpad Musical Sketchpad software. At the top, there's a banner with the text "Create amazing music WITH HOOKPAD MUSICAL SKETCHPAD". Below the banner are two buttons: "Shop Hookpad" and "Learn More". The main interface features a toolbar at the top with various icons for Play, Record, Click, Mixer, Preview, Meter, Key, Tempo, Band, Lyrics, Stable, and Piano. The "Key" section is set to C Major. The "Tempo" section shows 120 BPM. On the left, there's a "Duration: 2" section with a 4/4 time signature and a "Chords in C major" section showing chords I, ii⁷, iii⁷, IV⁷, V⁷, vi⁷, and vii⁷. The main workspace shows a 5-measure staff with chords I, V, V, vi, IV, and V⁷ over the notes C, G, G, am, F, and G7 respectively. To the right of the workspace are several panels: "Chord Properties" (Type: Triad, Inversion: None, Options: sus2, add9, sus4, add11, add7, no3, no5), "Secondary" (None, Borrow From: N/A), and "Bass Sets" (Magic, Popular, Search, Progressions, Bass Sets). The bottom of the screen shows a navigation bar with icons for Home, Lessons, Practice, Theory, and Tools.

Dataset

- [POP909](#)
- [Live Orchestral Piano](#)
- [LargeSHS](#)
- [Discog-VI](#)
- [Pop2Piano](#)

**If you are interested in doing anything related to
cover song generation, feel free to contact me**