

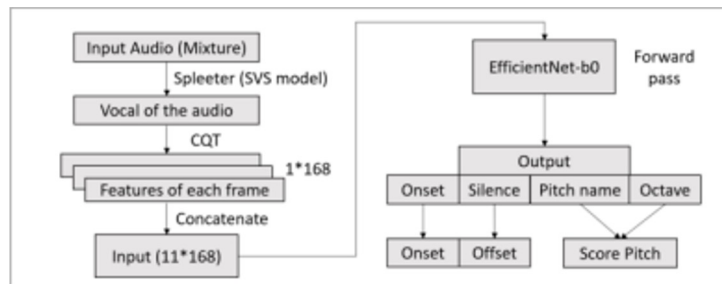


# Multi-modal Music Transcription System

# Automatic Melody Transcription

# AMT Assignment: Original Model Output

- **ONE** Nerual Network, **ONE** OutTensor to predict 4 targets



$$loss = onset\_loss + offset\_loss + octave\_loss + pitch\_loss \quad (weight?)$$

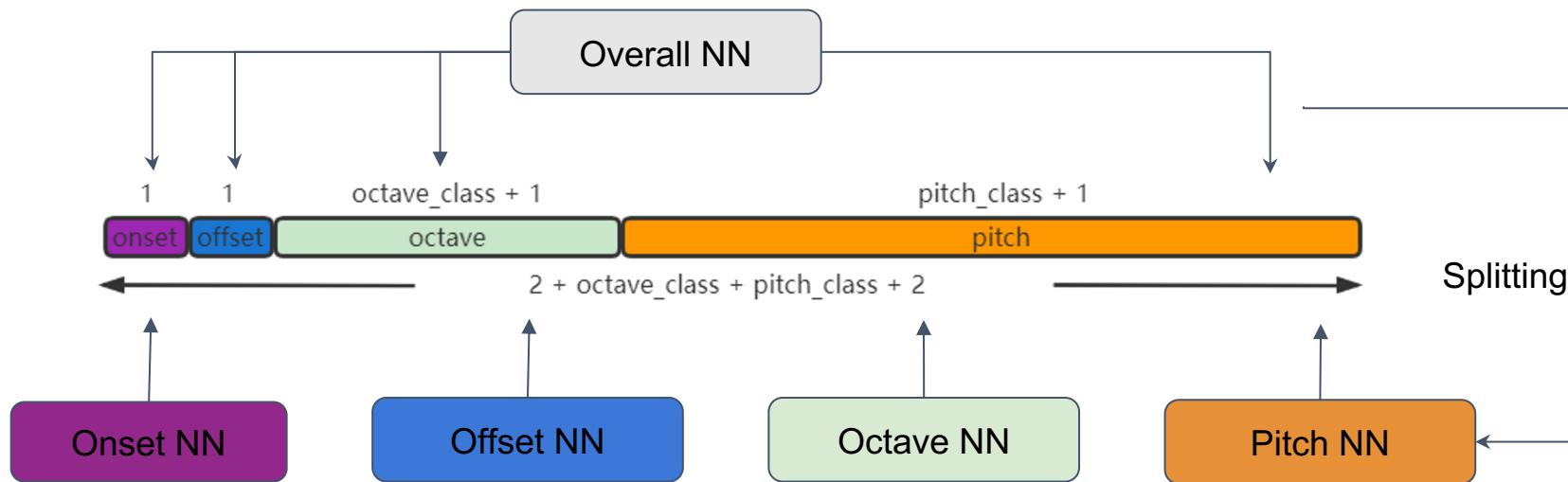
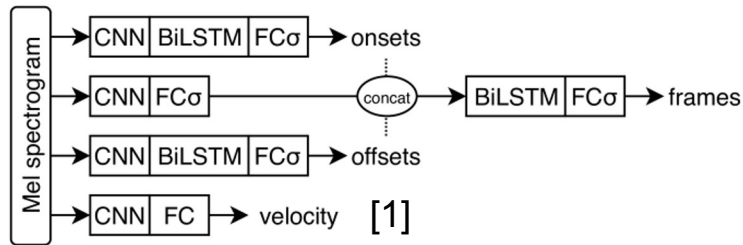
$$loss = k_1 * onset\_loss + k_2 * offset\_loss + k_3 * octave\_loss + k_4 * pitch\_loss$$

— Hard to **judge the best factor** ( $k_i$ ) of 4 different target

— Hard for model to predict 4 target **best at same time**

# Improvement: Model Splitting

**IDEA:** split different model to predict different metrics



[1] Github: [BShakhovsky/PolyphonicPianoTranscription](https://github.com/BShakhovsky/PolyphonicPianoTranscription): Recurrent Neural Network for generating piano MIDI-files from audio (MP3, WAV, etc.) ([github.com](https://github.com))

# Some improvement tricks

- *MIR-ST500 Dataset: only take annotation for vocal part ?*

Use extract tool (E.G. **demucs**) to extract vocal part !

- *Large-scale dataset, quick optimize and get worse soon ?*

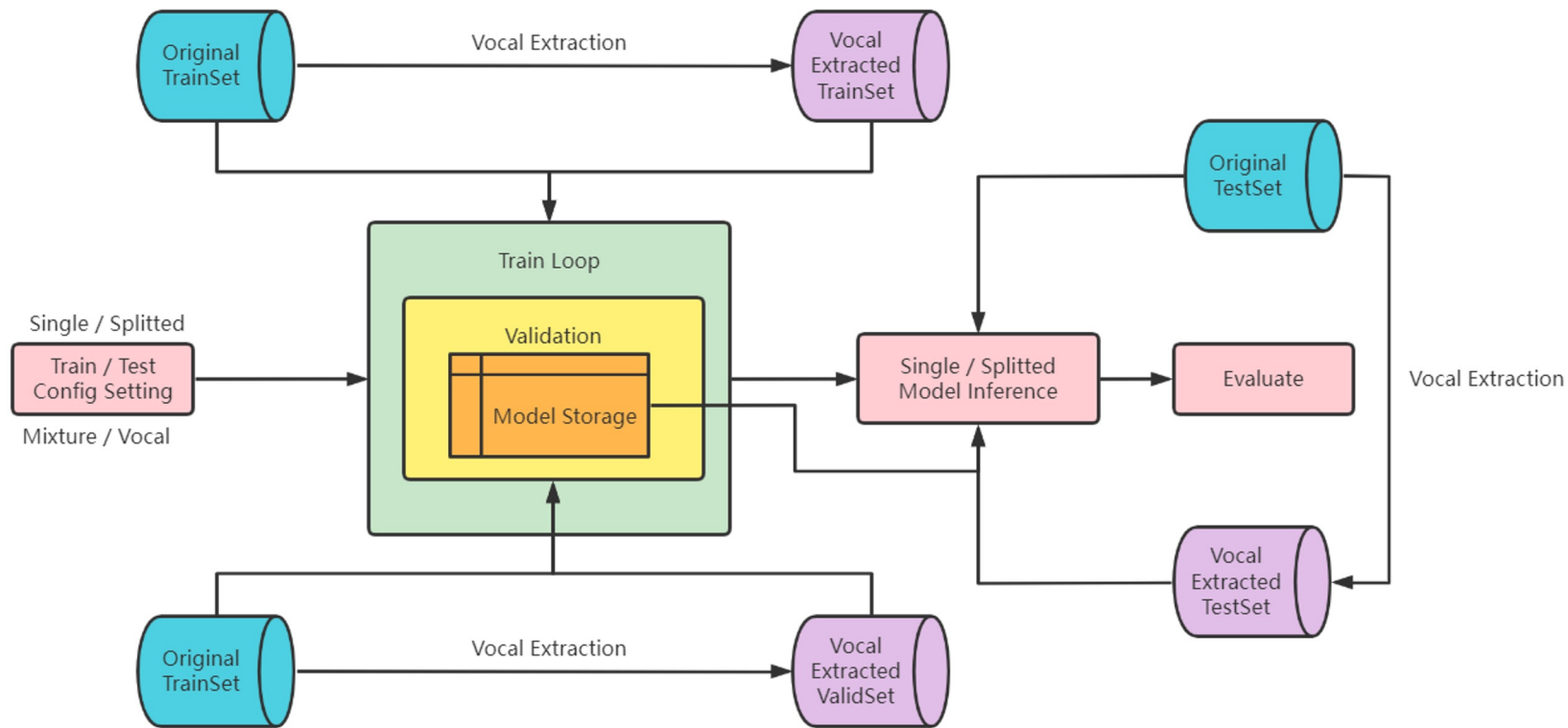
May skip some well-performed training model during epoch

Valid and save model after training **certain amount of data** !

For splitted model, **save best individually** and **combine** !

```
python inference.py --save_model_dir ./models7_vocal
--best_model_id_onset 59 --best_model_id_offset 70
--best_model_id_octave 19 --best_model_id_pitch 34
```

# Melody System Implementation



# Metrics Performance (TestSet)

EffNetb0  
F1-Score  
MIR-ST500  
[1]

ConPOff

0.4578

CONP

0.6663 ↓

CON

0.7544 ↓

Single Model + Mixture Audio  MIR-ST500	Metric	Precision	Recall	F1-Score
	ConPOff	0.334695	0.409107	0.366952
	CONP	0.569715	0.703664	0.627377
	CON	0.635402	0.784728	0.699628

Splitted Model + Mixture Audio  MIR-ST500	Metric	Precision	Recall	F1-Score
	ConPOff	0.357746 ↑	0.405023	0.378657
	CONP	0.605917 ↑	0.692694	0.644244 ↑
	CON	0.679555 ↑	0.775688	0.721938 ↑

Single Model + Vocal Audio  MIR-ST500	Metric	Precision	Recall	F1-Score
	ConPOff	0.402688	0.403332	0.402210
	CONP	0.693781	0.703302	0.697048 ↑
	CON	0.762988	0.772290	0.765943 ↑

Splitted Model + Vocal Audio  MIR-ST500	Metric	Precision	Recall	F1-Score
	ConPOff	0.402671	0.423197	0.411709 ↑
	CONP	0.674778	0.718016	0.694057 ↑
	CON	0.742230	0.788608	0.762811 ↑

- **Comparison Between Model Combination:** Better Precision (6%) (**Mixture**), obvious improvement when using mixture
- **Comparison Between Audio Input:** Significant improvement on **CONP**, **ConPOff** with vocal, better than EffNetb0
- **Splitted Model+Vocal Audio Structure:** Perform **best F1-Score** of **CONPOff** (0.4117, 12.2% higher than baseline) !



# Automatic Lyric Transcription



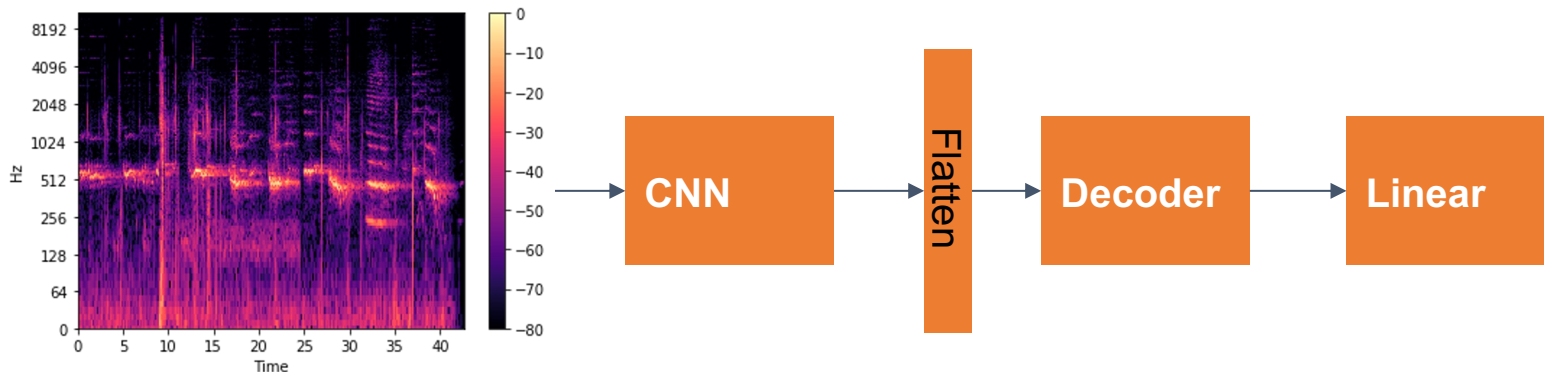
# ASR: Mel-Spectrogram Based

## 1. Mel-Spectrogram + CRDNN(CNN, RNN, DNN)

Spectrogram (raw feature) + CNN (feature extractor) + RNN (decoder) + DNN (linear)

## 1. Mel-Spectrogram + CNN + Transformer

Spectrogram (raw feature) + CNN (feature extractor) + Transformer (decoder) + DNN (linear)

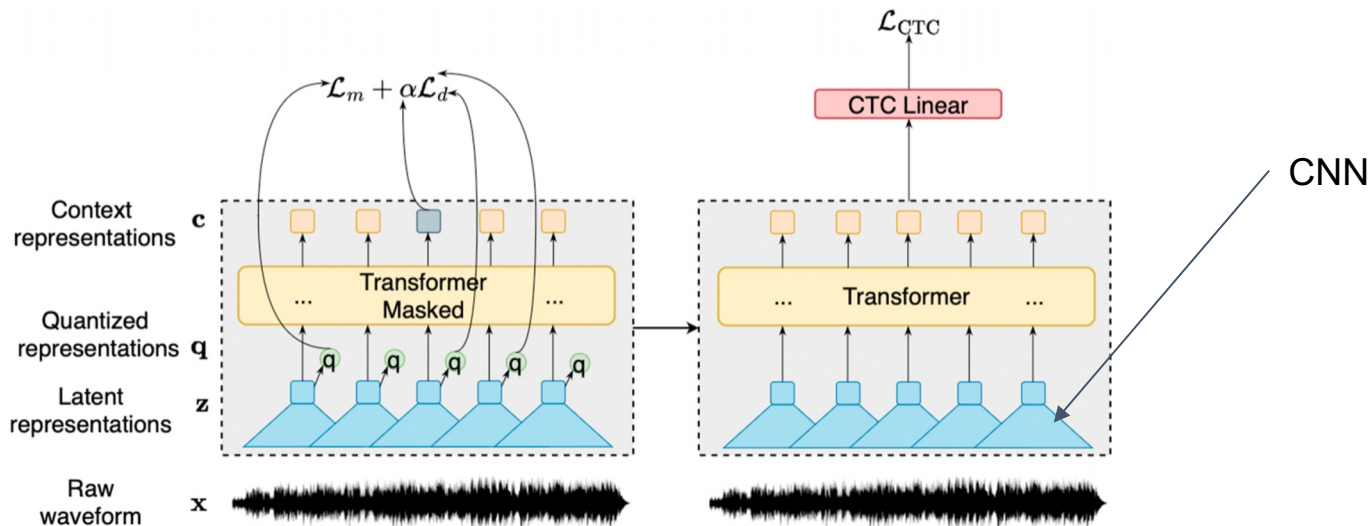


# ASR: End-to-End Wav2Vec 2.0

➤ Training wav2vec 2.0 has two stages:

Stage I: Self-supervised Contrastive Learning

Stage II: Supervised Fine-tuning

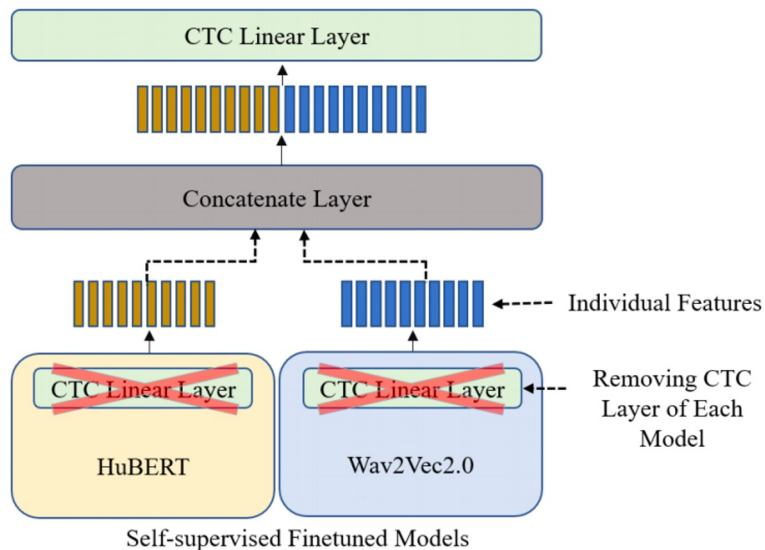


[1] Gu X, Ou L, Ong D, Wang Y. TRANSFER LEARNING OF WAV2VEC 2.0 FOR AUTOMATIC LYRIC TRANSCRIPTION[C]//Proceedings of the 30th ACM International Conference on Multimedia. 2022.

# Ensemble Methods on ASR

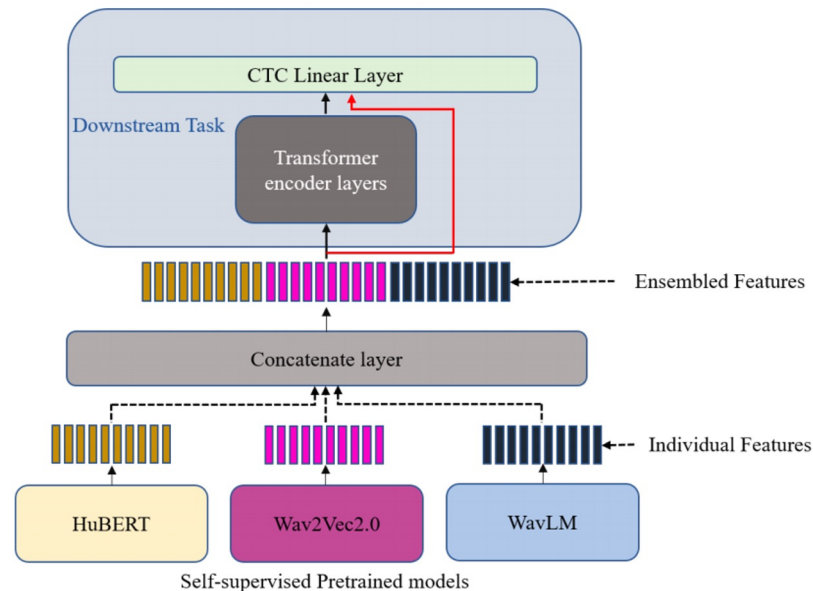
## 1. Model Ensemble

Concat before CTC Linear



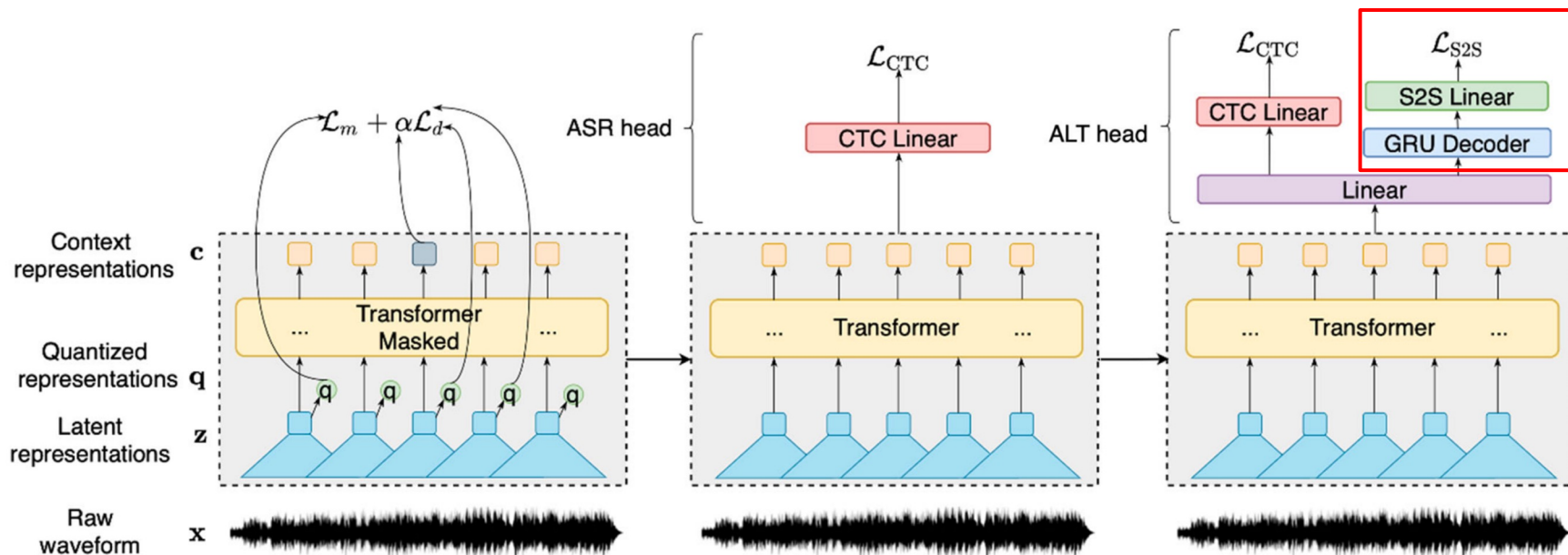
## 2. Feature Ensemble

Concat before Decoder



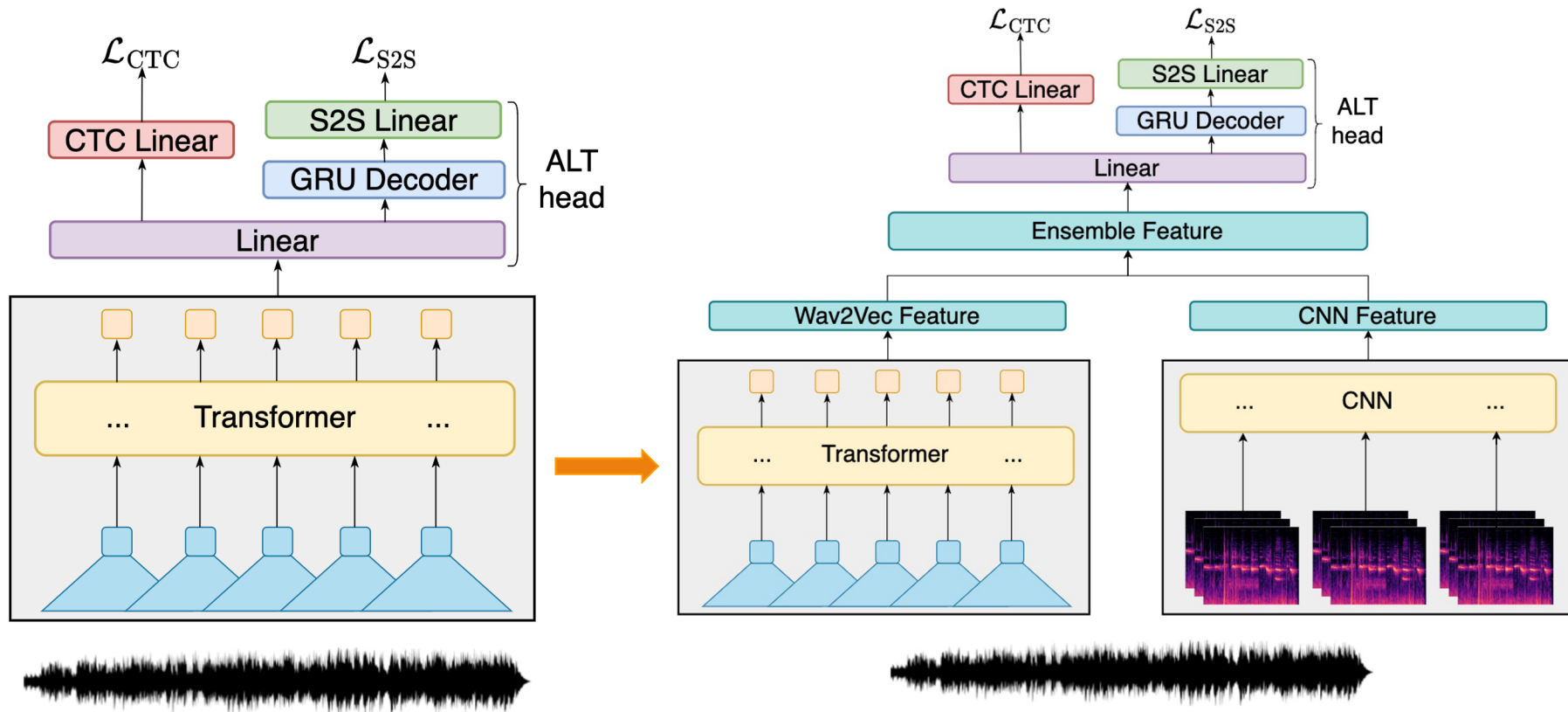
[1] A Arunkumar, Vrunda N Sukhadia, S. Umesh Investigation of Ensemble features of Self-Supervised Pretrained Models for Automatic Speech Recognition[C] International Speech Communication Association

# Transfer Learning to ALT



[1] Gu X, Ou L, Ong D, Wang Y. TRANSFER LEARNING OF WAV2VEC 2.0 FOR AUTOMATIC LYRIC TRANSCRIPTION[C]//Proceedings of the 30th ACM International Conference on Multimedia. 2022.

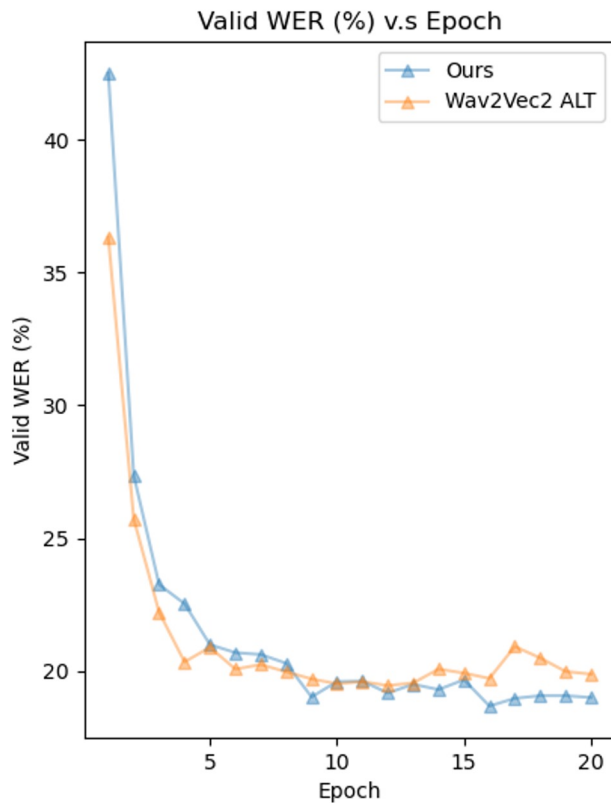
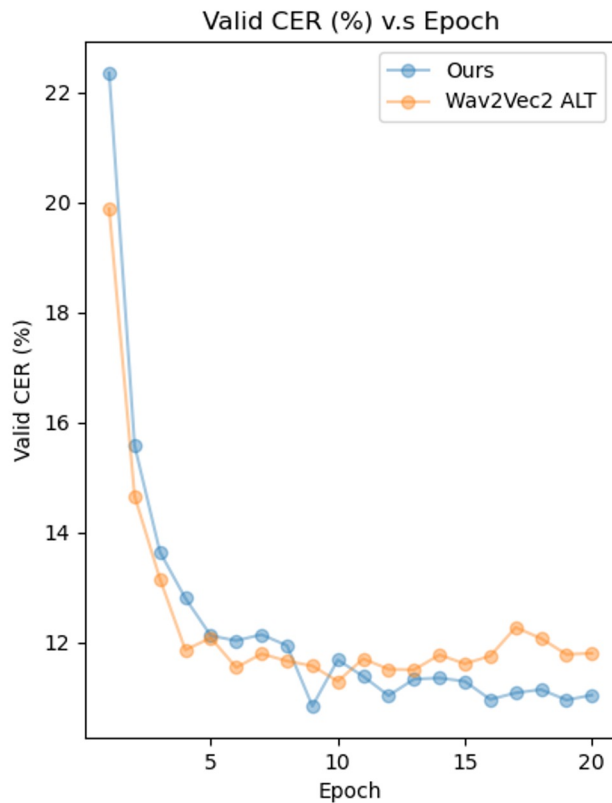
# Ensemble Features on ALT(Ours)



# Performance

Metric Model	Language Model	CER (%)	WER (%)	SER (%)
Wav2vec2+ALT head	RNNLM trained on DSing1	12.28	21.66	64.38
<b>Wav2Vec 2.0 + CNN +ALT Head(Ours)</b>		10.66	19.29	61.06
<b>Error Rate Decrease By</b>		<b>13.2%</b>	<b>10.9%</b>	<b>5.2%</b>
Wav2vec2+ALT head	RNNLM trained on DSing30 (10 times larger than Dsing1)	11.35	17.75	50.83
<b>Wav2Vec 2.0 + CNN +ALT Head(Ours)</b>		10.83	17.30	49.58
<b>Error Rate Decrease By</b>		<b>4.5%</b>	<b>2.5%</b>	<b>2.5%</b>

# Performance



# Model Parameters

ALT Model	Model Parameters
Wav2Vec 2.0	316.5M
ALT Head	89.4M
Wav2Vec 2.0 + ALT Head(Current SOTA)	405.9 M
Wav2Vec 2.0 + CNN + ALT Head(Ours)	408.7 M
<b>Model Increase By</b>	<b>0.5% (2.8M)</b>

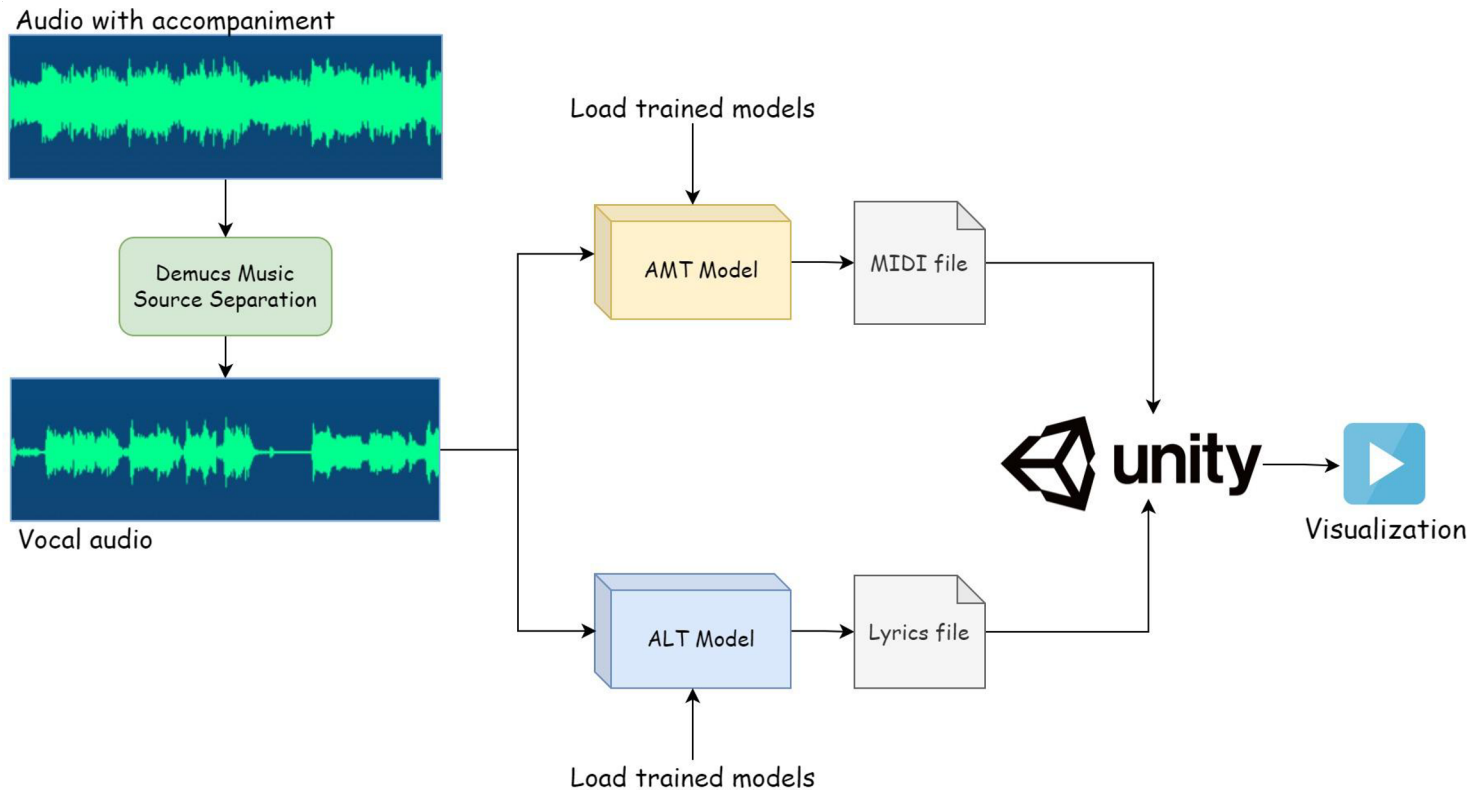
Our model improves the performance with only a slight increase of parameters!



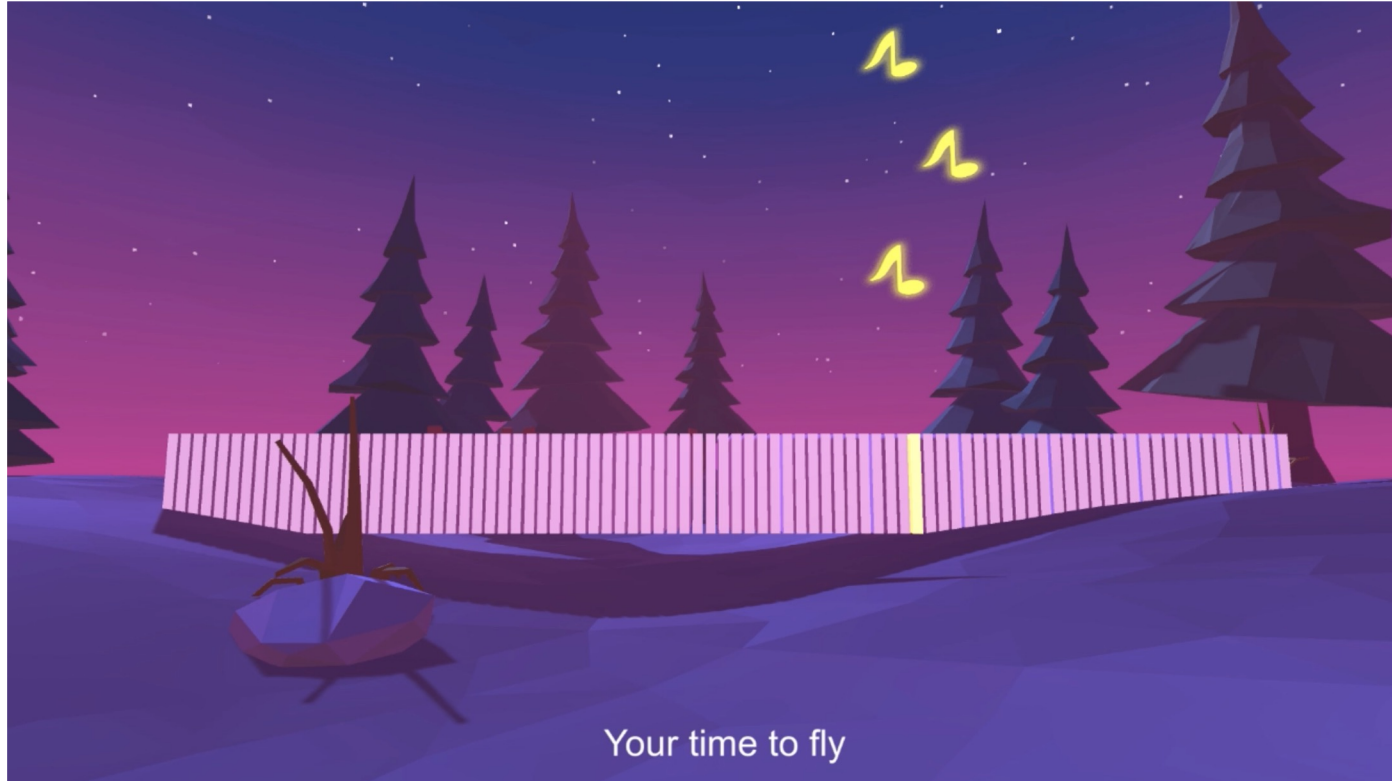


# Visuialization

# System Workflow



# Visualize pitch and lyrics



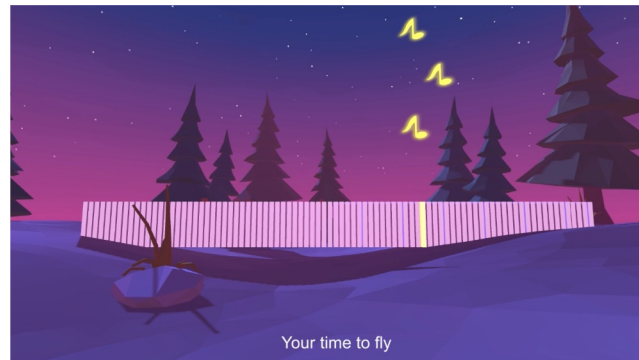
# Visualize pitch and lyrics

## Melody

- 88 keys
- Input: 2d array
- `[[start, end, pitch],...[start, end, pitch]]`
- generate a note when time to start

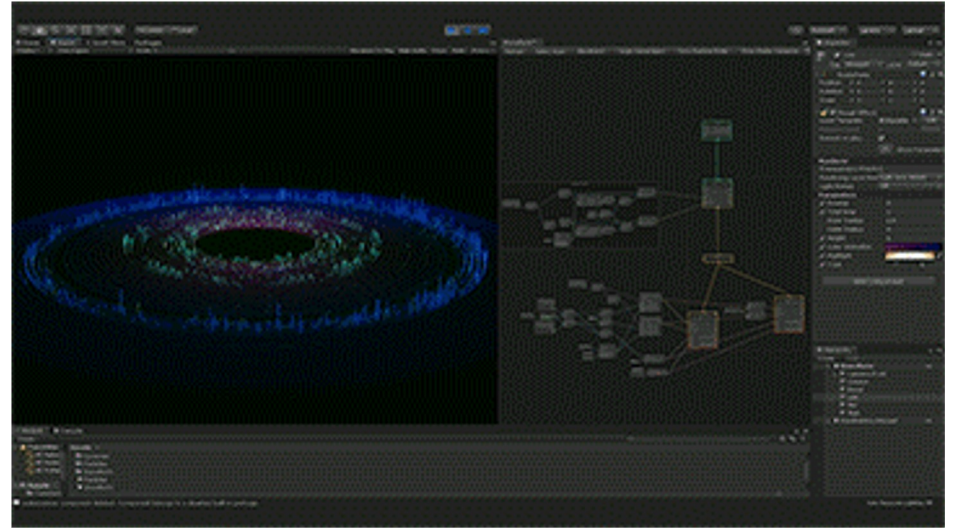
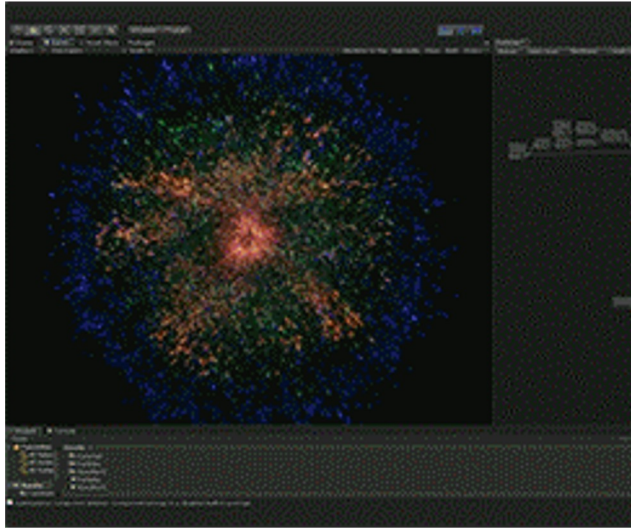
## Lyrics

- Input: 2d array
- `[[start, end, lyrics],...[start, end, lyrics]]`





# Audio Active visual



Customize particle system provided by Visual Effect Graphics in Unity



# Demonstration

Import Audio



So if you ever feel like giving up





**THANK YOU**