



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen



SCHOOL OF  
DATA SCIENCE  
數據科學學院

# A Comprehensive Guide to Amphion's Singing Voice Conversion

**Xueyao Zhang**

*The Chinese University of Hong Kong, Shenzhen*

2024/02

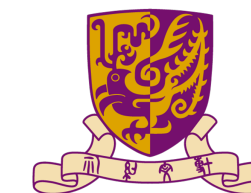
# About me



## Xueyao Zhang (张雪遥)

- ◆ **Second-year PhD student**, Supervised by Prof Zhizheng Wu  
School of Data Science, CUHK-Shenzhen  
Homepage: <https://www.zhangxueyao.com/>
- ◆ **Amphion v0.1's co-founder**  
Project: <https://github.com/open-mmlab/Amphion> (**3.5k stars**)
- ◆ **Research interest: "AI + Music"**, especially on:
  - Singing Voice Processing
  - Music Generation

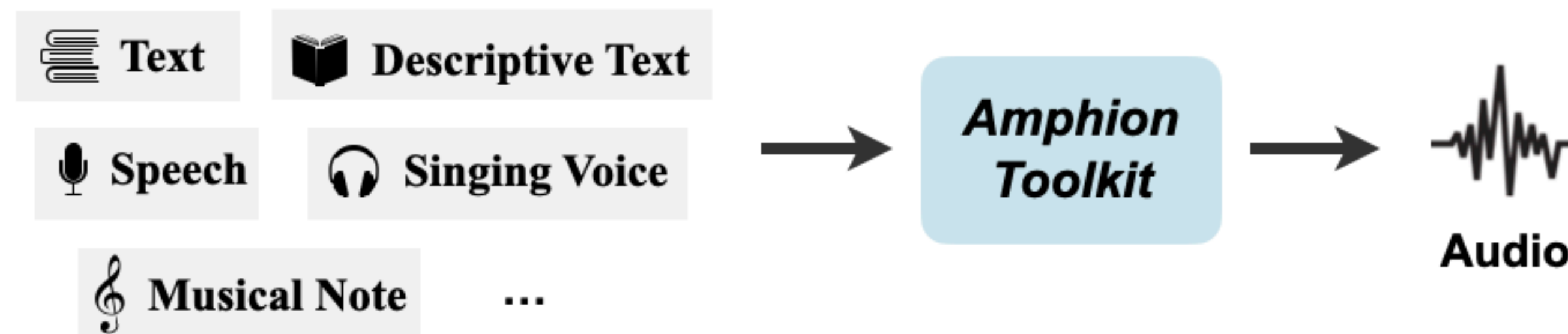
- 📎 **Amphion Technical Report:** <https://arxiv.org/abs/2312.09911>
- 👤 **Amphion GitHub:** <https://github.com/open-mmlab/Amphion>
- 🎯 **Amphion Demos/Models/Datasets:** <https://huggingface.co/amphion>





# About Amphion

- Support **reproducible research** and help **junior researchers and engineers** get started in the field of audio, music, and speech generation research and development.



**Our North-Star Objective:**  
**Any to Audio**

## Amphion: An Open-Source Audio, Music and Speech Generation Toolkit

Xueyao Zhang<sup>\*,1</sup>, Liumeng Xue<sup>\*,1</sup>, Yicheng Gu<sup>\*,1</sup>, Yuancheng Wang<sup>\*,1</sup>, Haorui He<sup>3</sup>,  
Chaoren Wang<sup>1</sup>, Xi Chen<sup>1</sup>, Zihao Fang<sup>1</sup>, Haopeng Chen<sup>1</sup>, Junan Zhang<sup>2</sup>, Tze Ying Tang<sup>1</sup>,  
Lexiao Zou<sup>3</sup>, Mingxuan Wang<sup>1</sup>, Jun Han<sup>1</sup>, Kai Chen<sup>2</sup>, Haizhou Li<sup>1</sup>, Zhizheng Wu<sup>†,1,2,3</sup>

<sup>1</sup>School of Data Science, The Chinese University of Hong Kong, Shenzhen

<sup>2</sup>Shanghai AI Lab

<sup>3</sup>Shenzhen Research Institute of Big Data

- TTS: Text to Speech (🚩 supported)
- SVS: Singing Voice Synthesis (👷 developing)
- VC: Voice Conversion (👷 developing)
- SVC: Singing Voice Conversion (🚩 supported)
- TTA: Text to Audio (🚩 supported)
- TTM: Text to Music (👷 developing)
- more...





*Amphion + Sora*



# Roadmap

---

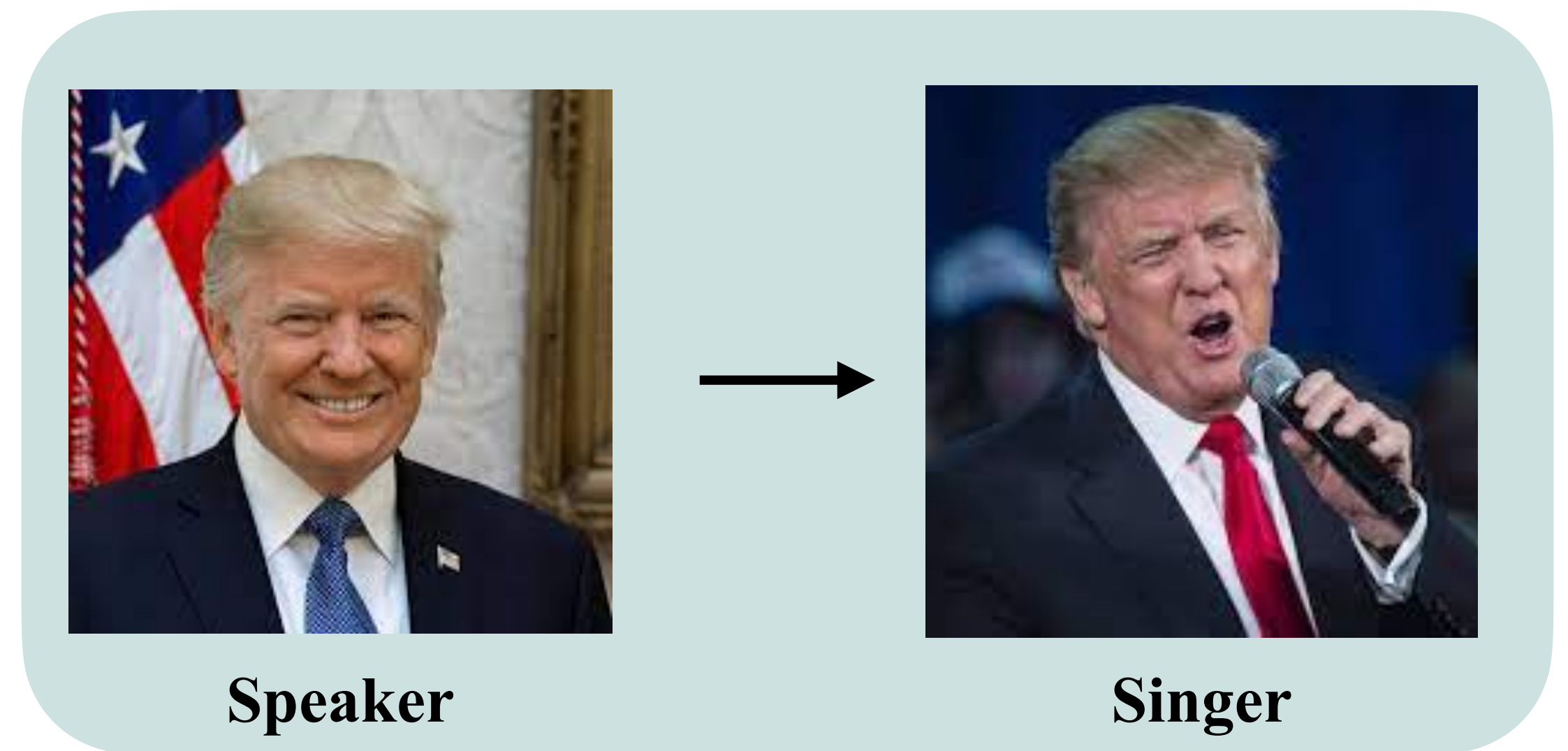
- **Singing Voice Conversion**
  - Definition, Classic Works, and Modern Pipeline
- **Singing Voice Conversion in Amphion**
  - Supported Model Architectures
  - Our research: *Leveraging Diverse Semantic-based Audio Pretrained Models for Singing Voice Conversion*
- **Amphion's Philosophy**
  - Unique strengths, Supported Features, and Visualization



# What is Singing Voice Conversion (SVC)?



*Inter-singer Conversion*



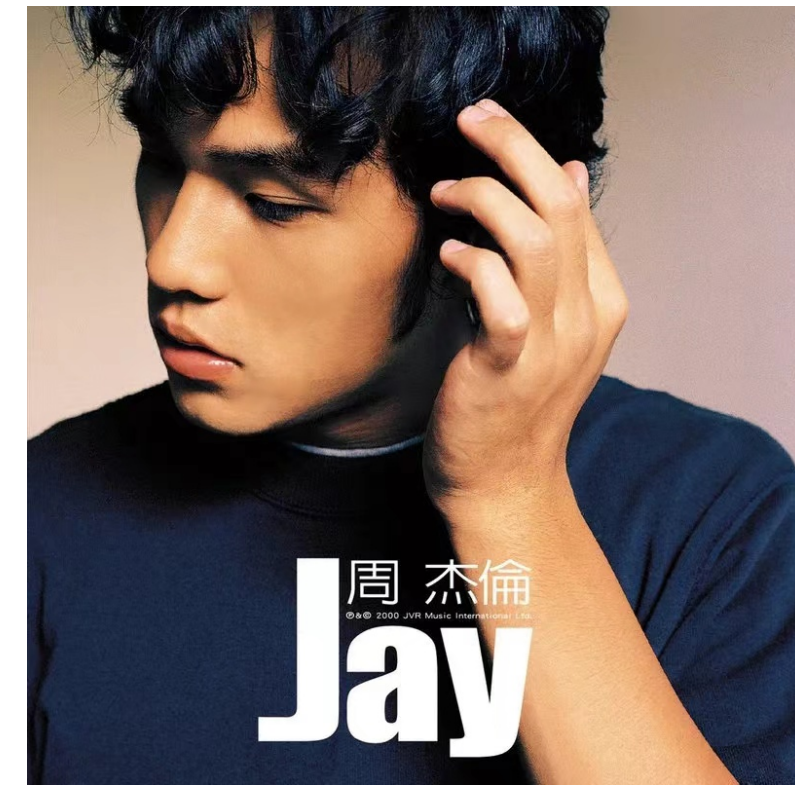
*Cross-domain Conversion*



*Intra-singer Conversion*



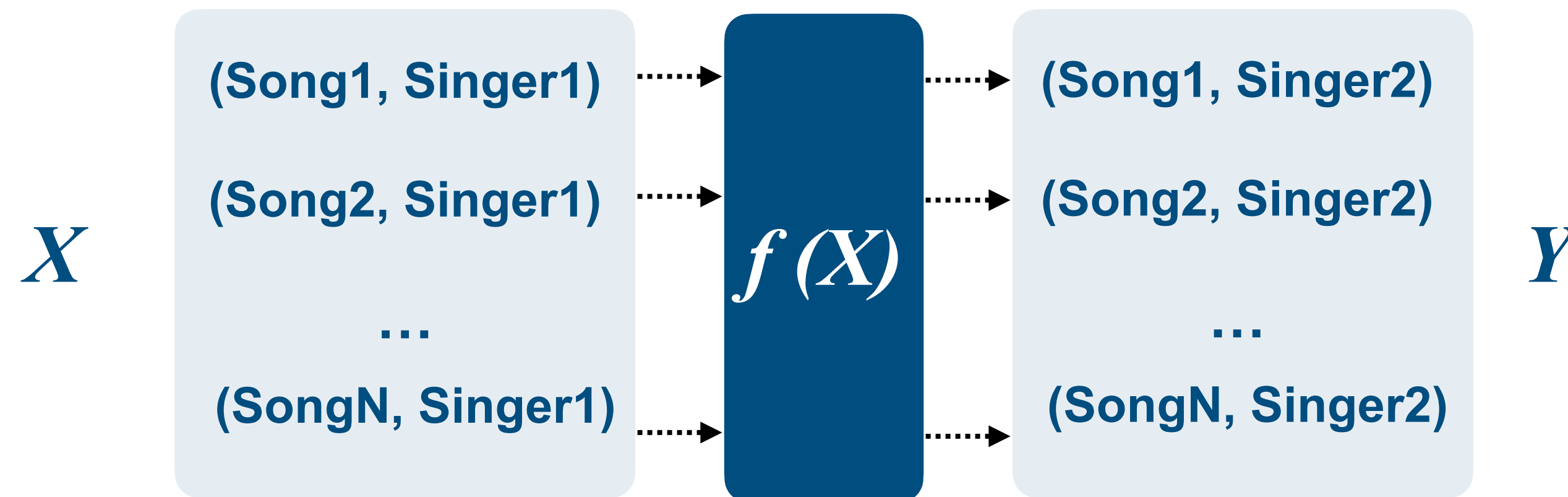
# Parallel Singing Voice Conversion



Professional Singer1



Professional Singer2

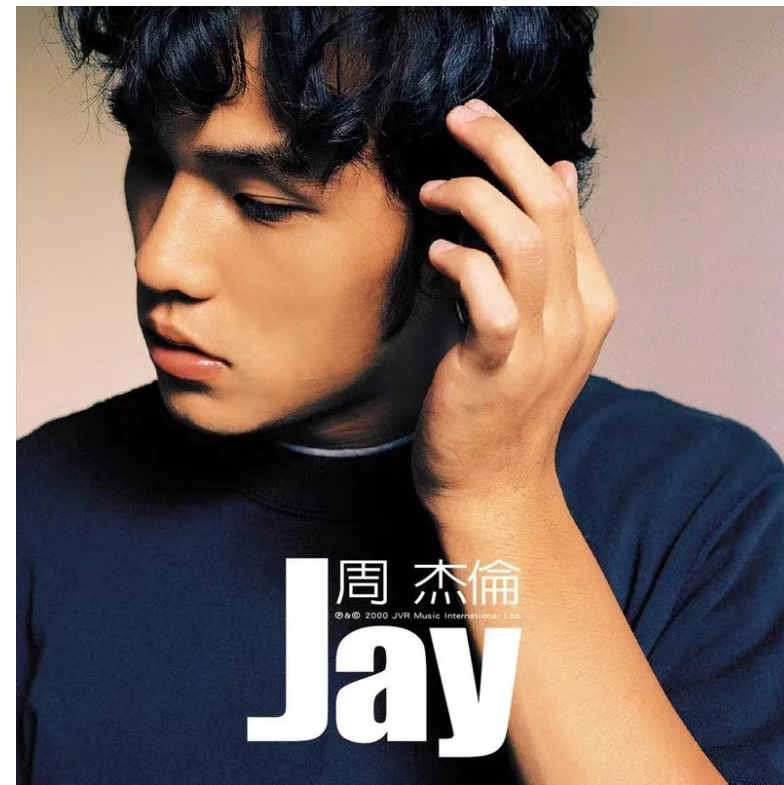


**Parallel corpus is hard to collect!**



# Non-Parallel Singing Voice Conversion

---



Professional Singer1



Professional Singer2

X

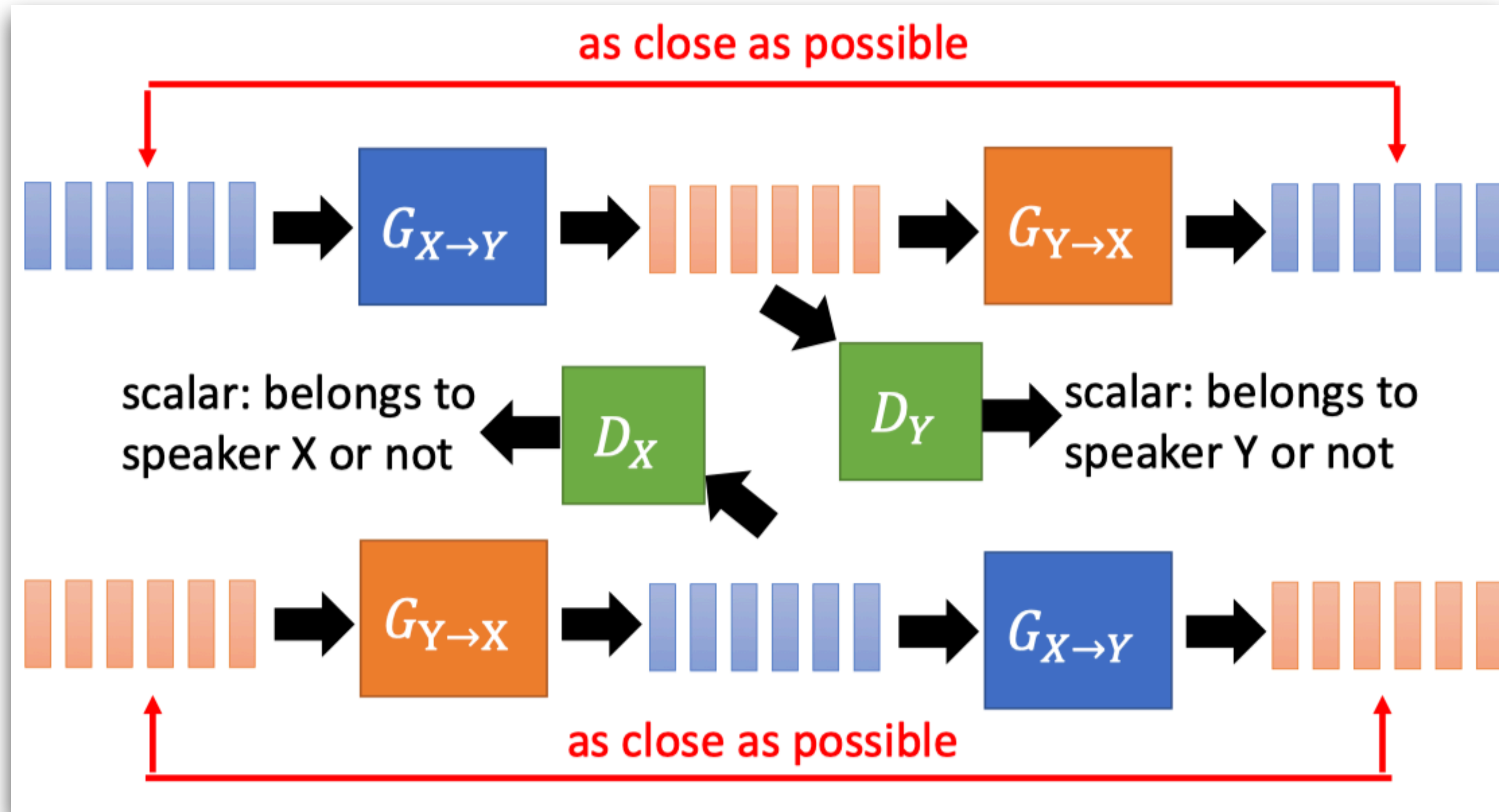
Singer1's Songs

Singer2's Songs

How to decouple the singer identity?



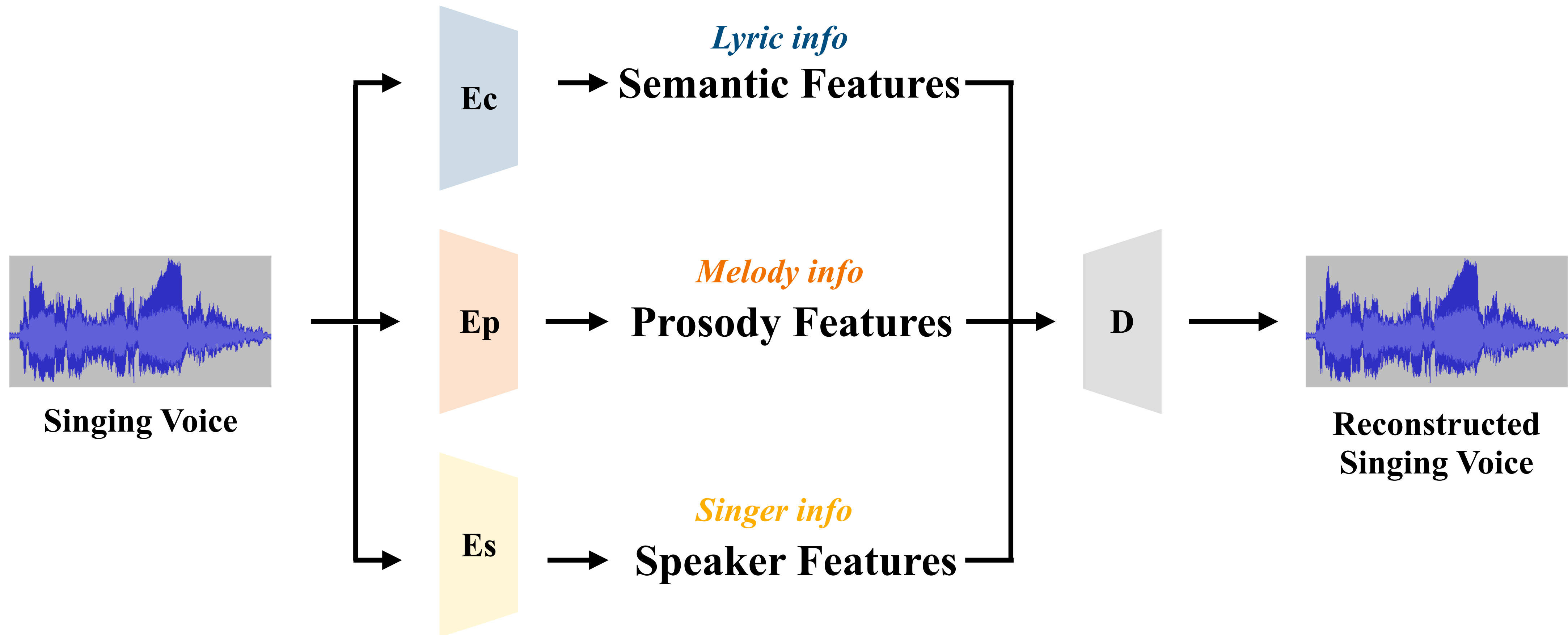
# Non-Parallel SVC: GAN School



Credit: Voice Conversion, Hung-yi Lee.



# Non-Parallel SVC: **Auto-Encoder School**

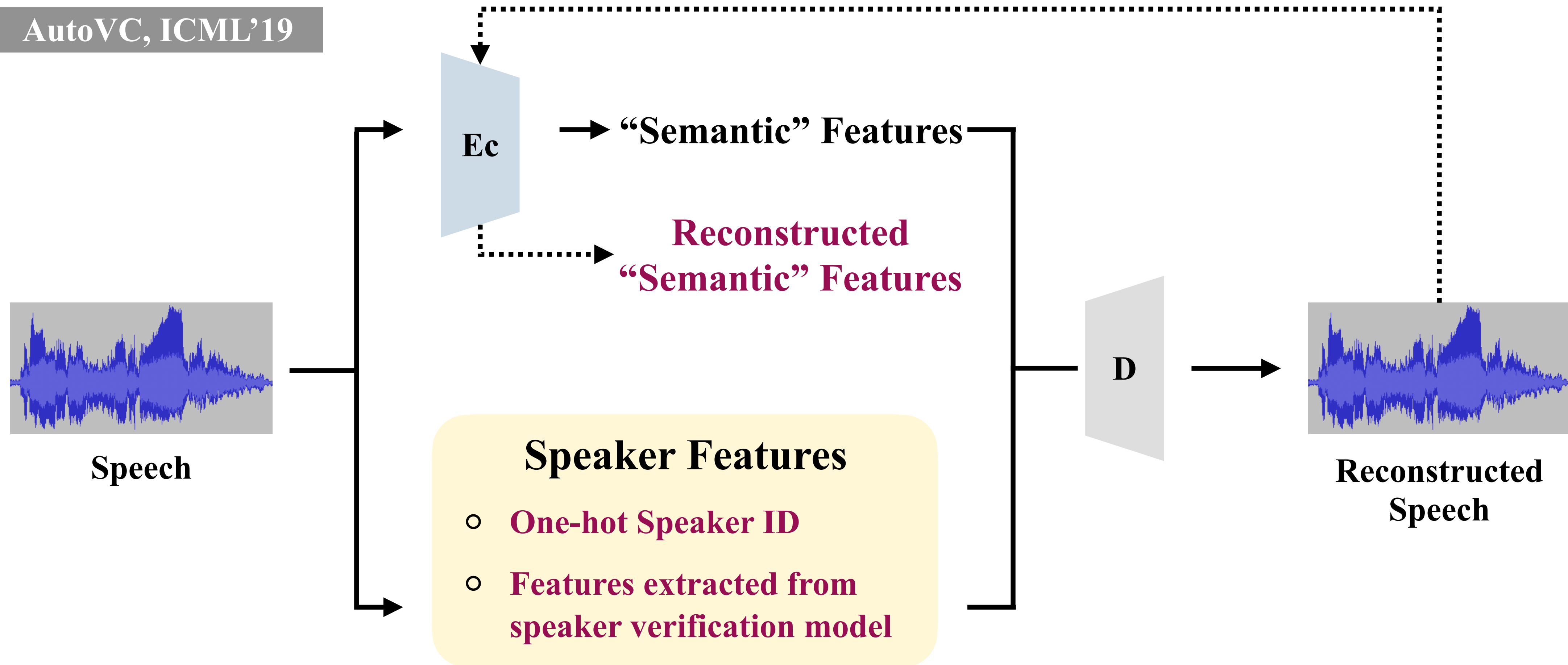


- How to ensure the disentanglement of different features?
- How to ensure there is enough information of each features?



# Auto-Encoder VC: The Early Researches

AutoVC, ICML'19

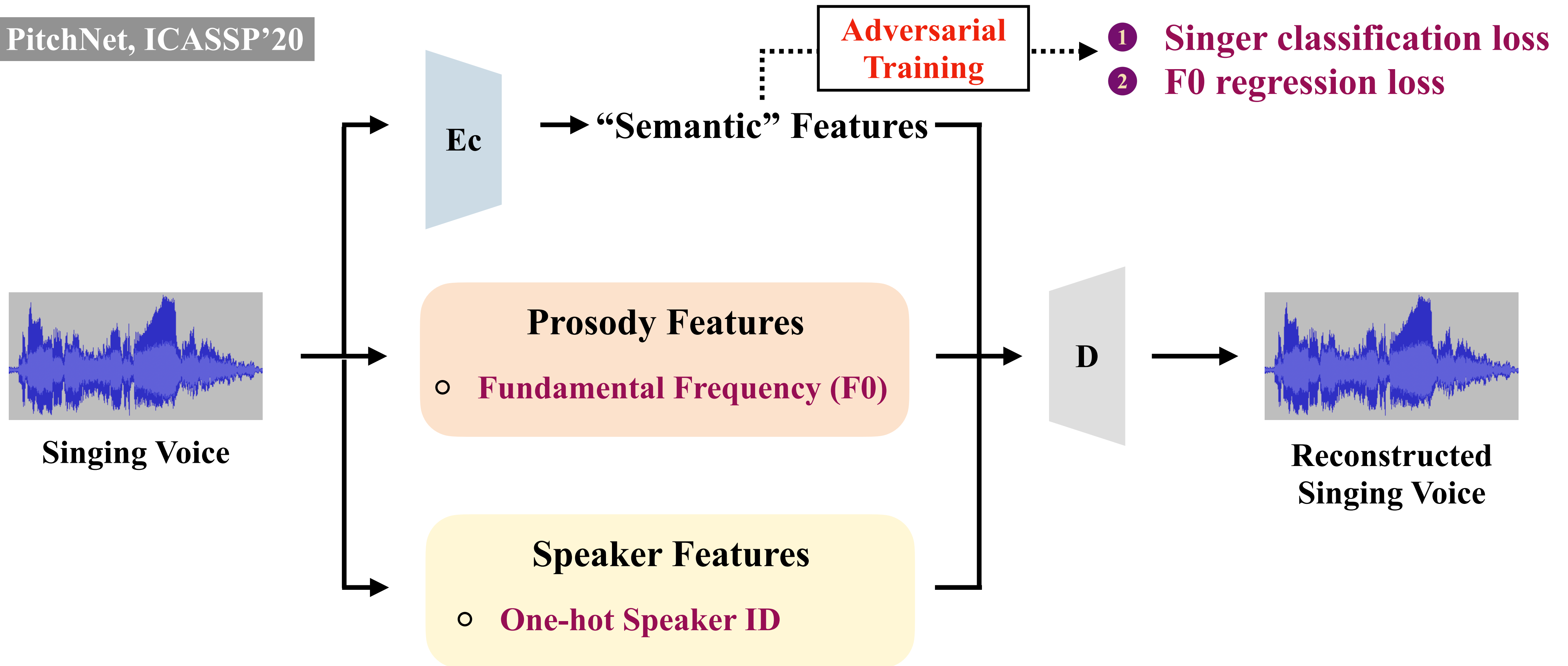


**AutoVC: “To carefully design the dimension of the *semantic* features”**



# Auto-Encoder SVC: The Early Researches

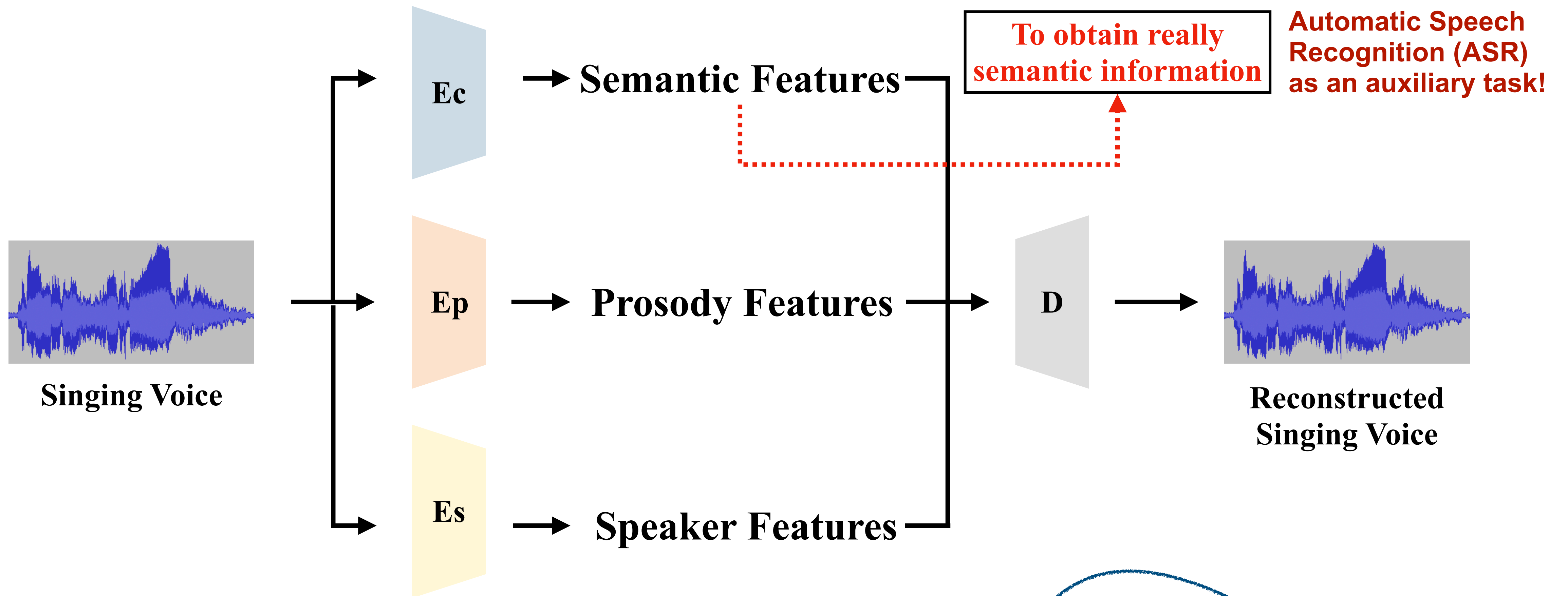
PitchNet, ICASSP'20



**PitchNet: “Adopt adversarial training to disentangle better”**



# (Review) Non-Parallel SVC: Auto-Encoder School



- How to ensure the disentanglement of different features?
- How to ensure there is enough information of each features?

🎯 Solved to some extent

🤔 How to address?



# Non-Parallel VC/SVC — a.k.a Recognition & Synthesis VC/SVC

---



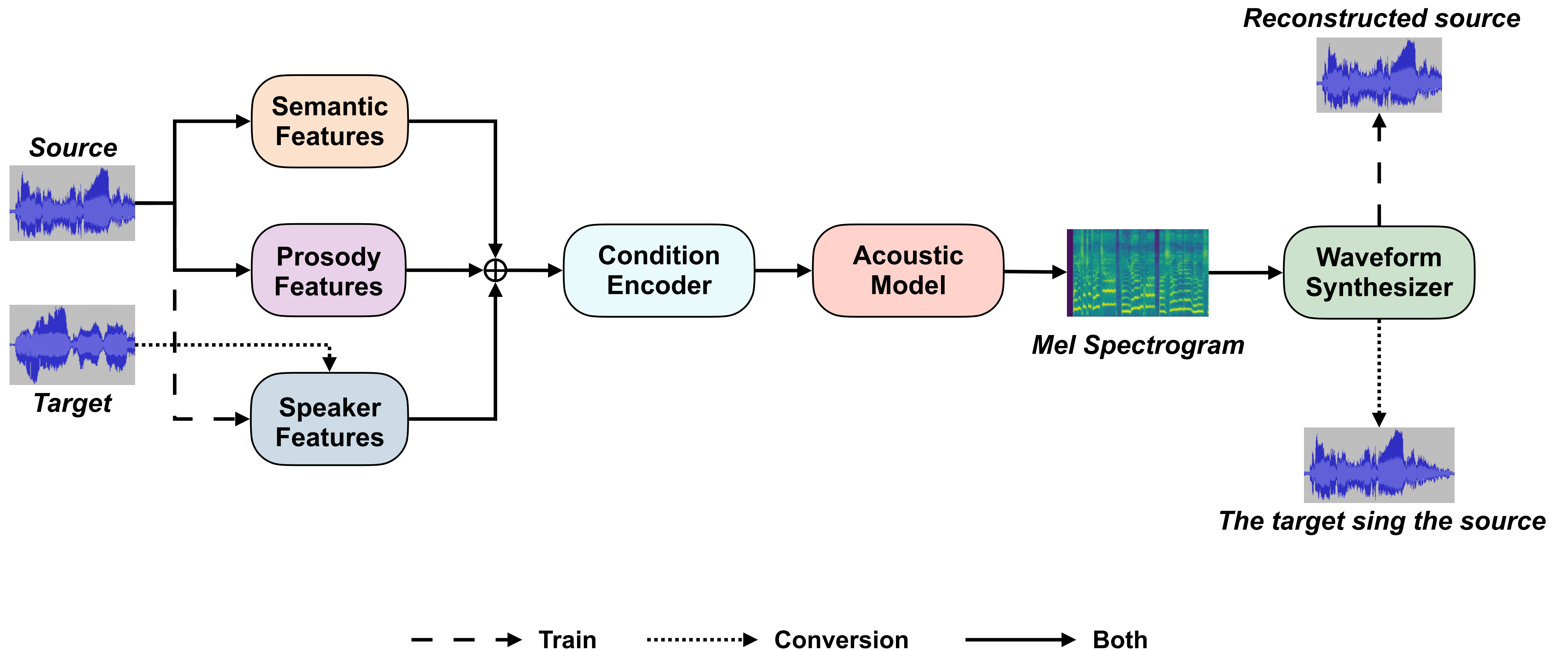
Using the intermediate output as “semantic-based” features

🤔 Why do we use the *dense semantic features* instead of the *symbolic text*?

- 1 There are errors for the recognized symbolic text.
- 2 It takes more time to obtain the symbolic text than just extracting dense features.
- 3 There are more acoustic information (such as pronunciation) in the dense features, which is better for improving the intelligibility of the synthesized voice.



# Modern Singing Voice Conversion Pipeline



# Roadmap

---

- **Singing Voice Conversion**
  - Definition, Classic Works, and Modern Pipeline
- **Singing Voice Conversion in Amphion**
  - Supported Model Architectures
  - Our research: *Leveraging Diverse Semantic-based Audio Pretrained Models for Singing Voice Conversion*
- **Amphion's Philosophy**
  - Unique strengths, Supported Features, and Visualization



# Amphion SVC: Supported Model Architectures

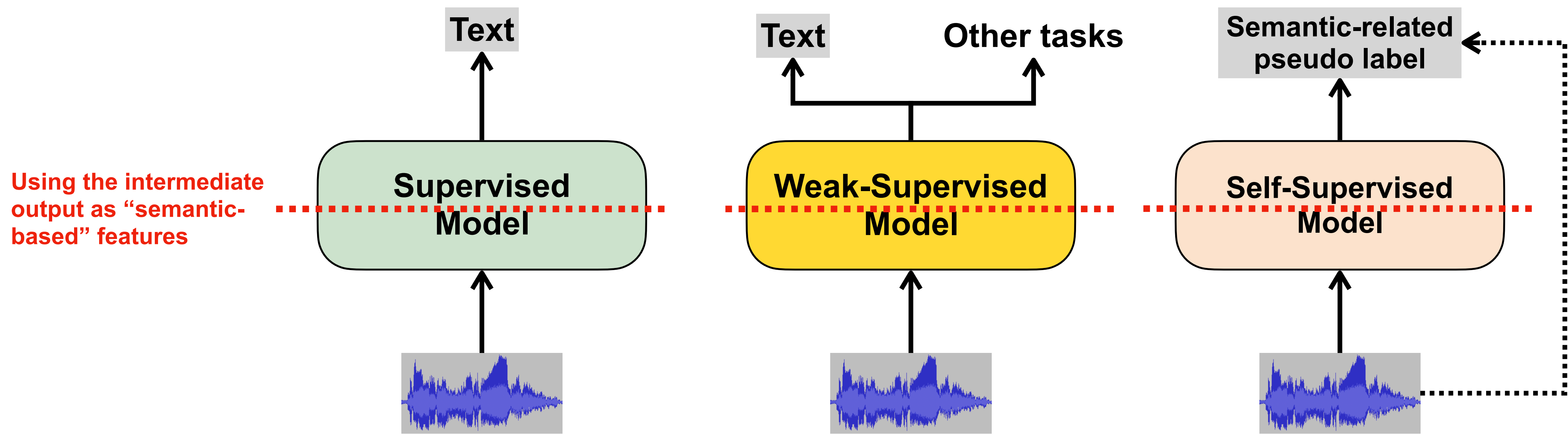
---

- **Semantic Features Extractor**
  - WeNet, Whisper, ContentVec
  - Joint Usage of Diverse Semantic Features Extractors
- **Prosody Features**
  - F0 and energy
- **Speaker Features**
  - One-hot Speaker ID
  - Features of Pretrained SV model
- **Acoustic Model**
  - Diffusion-based
  - Transformer-based
  - VAE- and Flow-based
- **Waveform Synthesizer**
  - GAN-based
  - Diffusion-based

# Semantic Features: Why Joint Usage of Multiple Extractors?

- **Background of Semantic-based Pretrained Models:**

- **Varied choices:** Classic ASR models, Whisper, HuBERT, Wav2Vec, WavLM, ContentVec...
- **High pretraining cost:** E.g., Whisper-large (1.5B parameters, 680k hours training data)

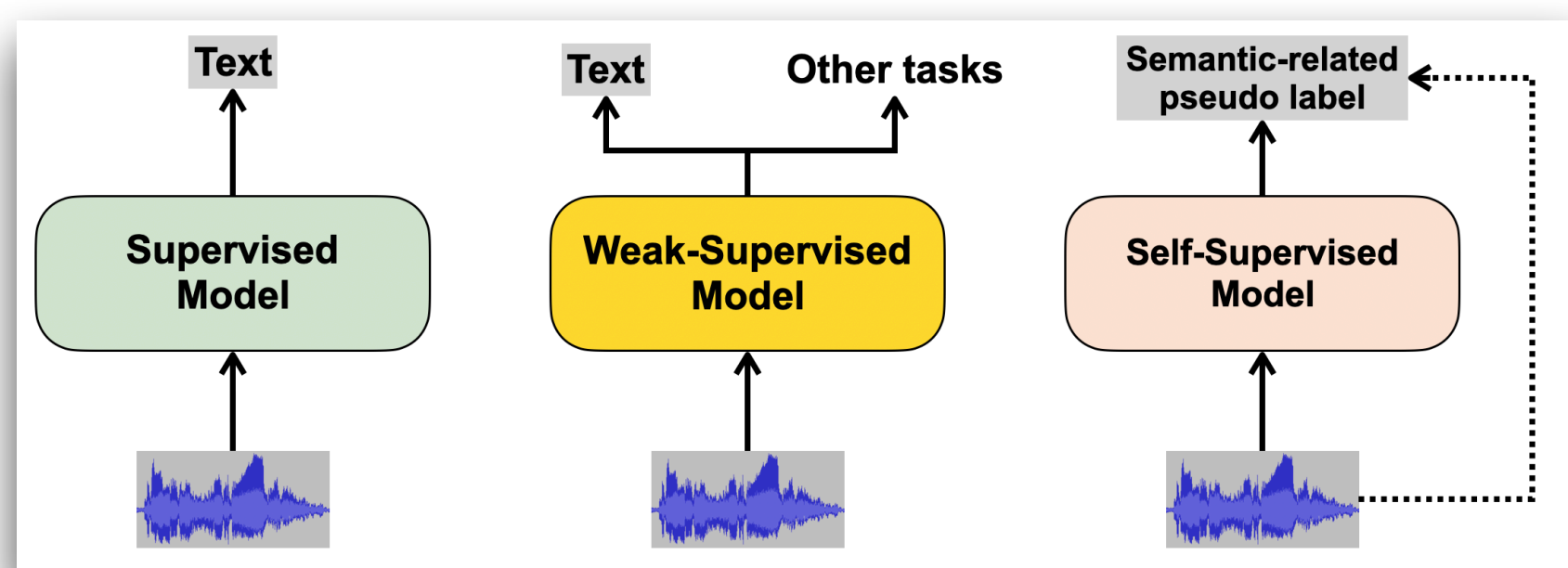




# Semantic Features: Why Joint Usage of Multiple Extractors?

- However, for the downstream tasks (such as SVC), the underlying knowledge of these models remains largely unknown:

Requirements of SVC	Capability of the Semantic-based Features
To model melody	<b>Whether could or not</b> remains unknown
To model lyrics	Could. But <b>exactly how much</b> remains unknown
To model auxiliary acoustic information	Could. But <b>whether the information is speaker-agnostic or not</b> remains unknown
To be robust for in-the-wild acoustic environment	<b>Whether is robust or not</b> remains unknown

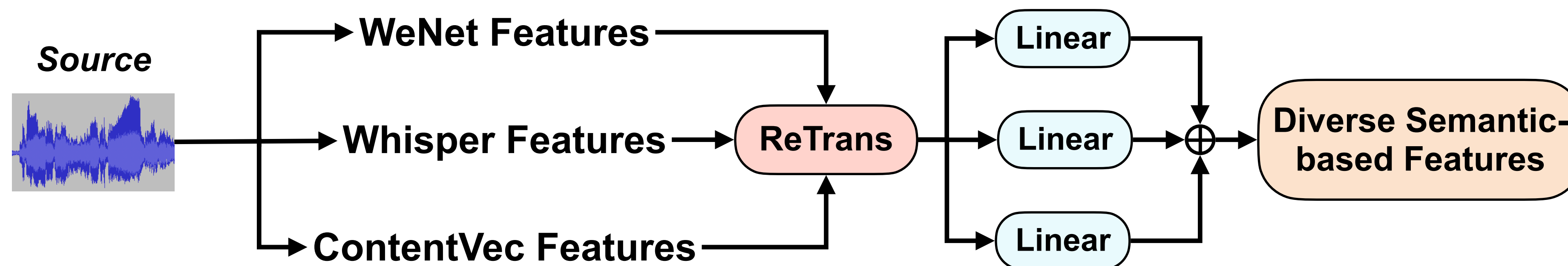


Different pretraining ways will yield different underlying knowledge



Could diverse pretrained models be complementary for SVC?

# Challenge: Time Resolution Mismatch of Multiple Semantic Features



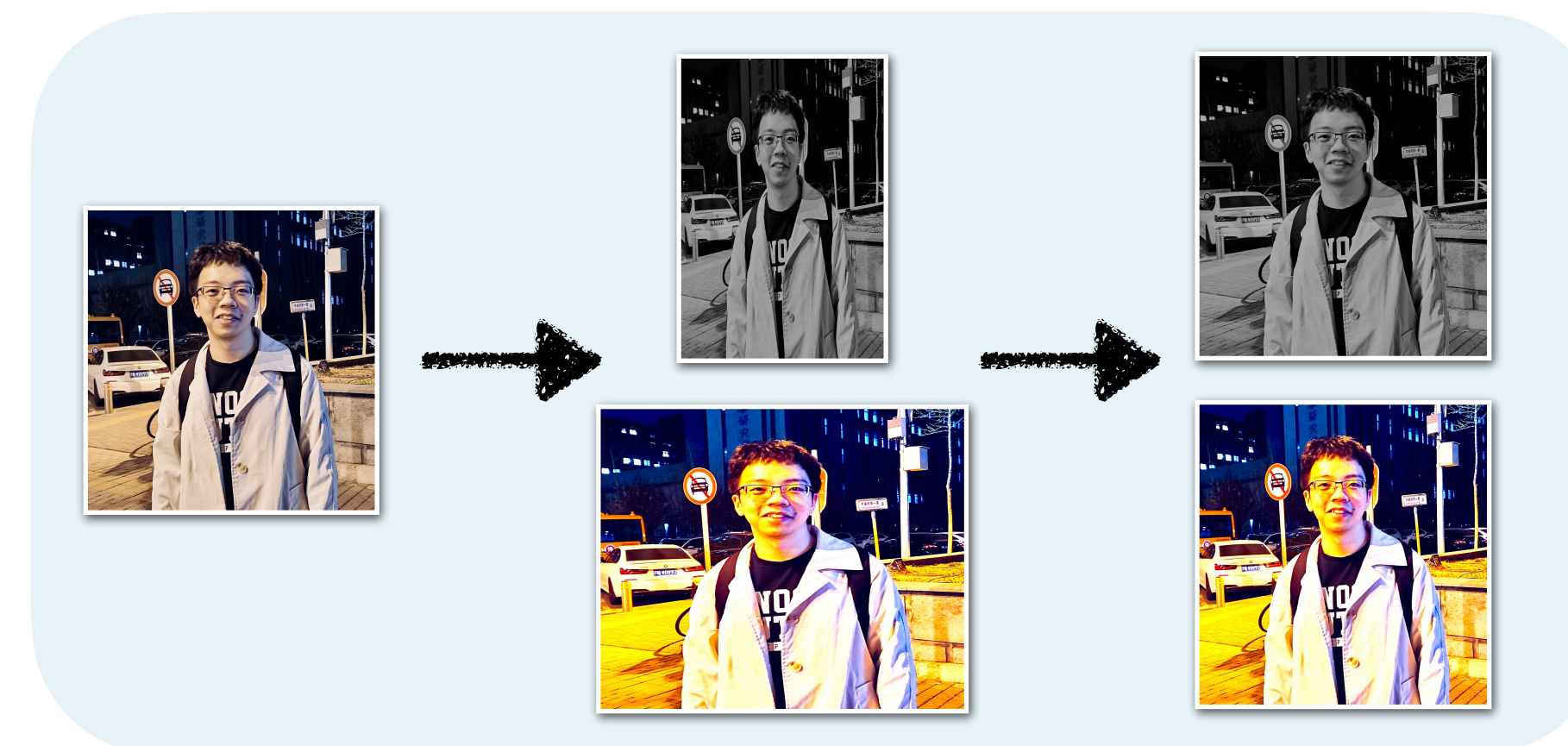
## Algorithm 1 ReTrans for features of audio pretrained models

**Input:**  $f$  – the features to be transformed whose frame rate is  $r_1$  Hz.

**Parameter:**  $r_2$  – The desired frame rate to transform

**Output:**  $f'$  – the transformed features whose frame rate is  $r_2$  Hz

- 1:  $c \leftarrow \text{gcd}(r_1, r_2)$   $\triangleright$  The greatest common divisor
- 2:  $f' \leftarrow \text{upsample}(f, r_2/c)$   $\triangleright$  Upsampling
- 3:  $f' \leftarrow \text{downsample}(f', r_1/c)$   $\triangleright$  Downsampling
- 4: **return**  $f'$



**High efficiency;  
No more training cost**



# Results: Using Only Semantic-based Features for SVC

Semantic-based Features	MCD (↓)	F0CORR (↑)	F0RMSE (↓)	CER (↓)	SIM (↑)
Ground Truth	0.000	1.000	0.0	12.9%	1.000
WeNet	10.324	0.203	423.4	38.2%	0.912
Whisper	8.229	0.524	297.3	18.9%	0.914
ContentVec	8.972	① 0.491	361.0	② 22.1%	③ <b>0.918</b>
WeNet + Whisper	8.345	0.540	284.2	16.8%	0.911
WeNet + ContentVec	8.870	0.525	329.5	19.9%	0.912
Whisper + ContentVec	<b>8.201</b>	0.548	279.6	16.9%	0.912
WeNet + Whisper + ContentVec	8.249	<b>0.572</b>	<b>278.5</b>	<b>16.1%</b>	0.913

① **To model melody:**

Whisper > ContentVec > WeNet

② **To model lyrics:**

Whisper > ContentVec > WeNet

③ **To be speaker-agnostic:**

All the three is good

\* Weak-supervised and self-supervised models is more robust for singing voice

\* Large-scale pretraining corpus is necessary

# Results: Complementary roles of Diverse Semantic-based Features

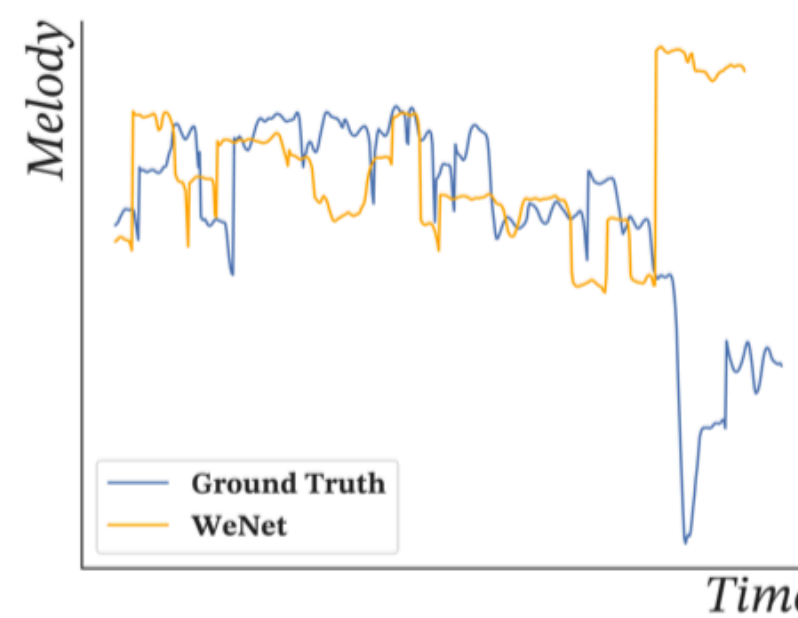
Semantic-based Features	MCD (↓)	F0CORR (↑)	F0RMSE (↓)	CER (↓)	SIM (↑)
Ground Truth	0.000	1.000	0.0	12.9%	1.000
WeNet	10.324	0.203	423.4	38.2%	0.912
Whisper	8.229	0.524	297.3	18.9%	0.914
ContentVec	8.972	0.491	361.0	22.1%	<b>0.918</b>
WeNet + Whisper	8.345	0.540	284.2	16.8%	0.911
WeNet + ContentVec	8.870	0.525	329.5	19.9%	0.912
Whisper + ContentVec	<b>8.201</b>	0.548	279.6	16.9%	0.912
WeNet + Whisper + ContentVec	8.249	<b>0.572</b>	<b>278.5</b>	<b>16.1%</b>	0.913

①

After  
Introducing F0

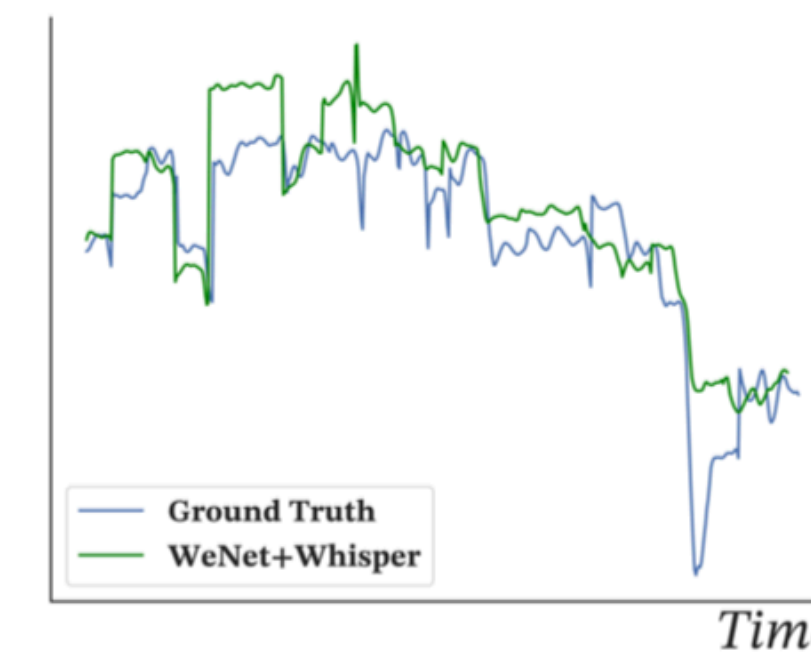
Using diverse semantic-based features:

- ① Most results are promoted stage by stage
- ② Introducing explicit melody modeling for SVC remains necessary

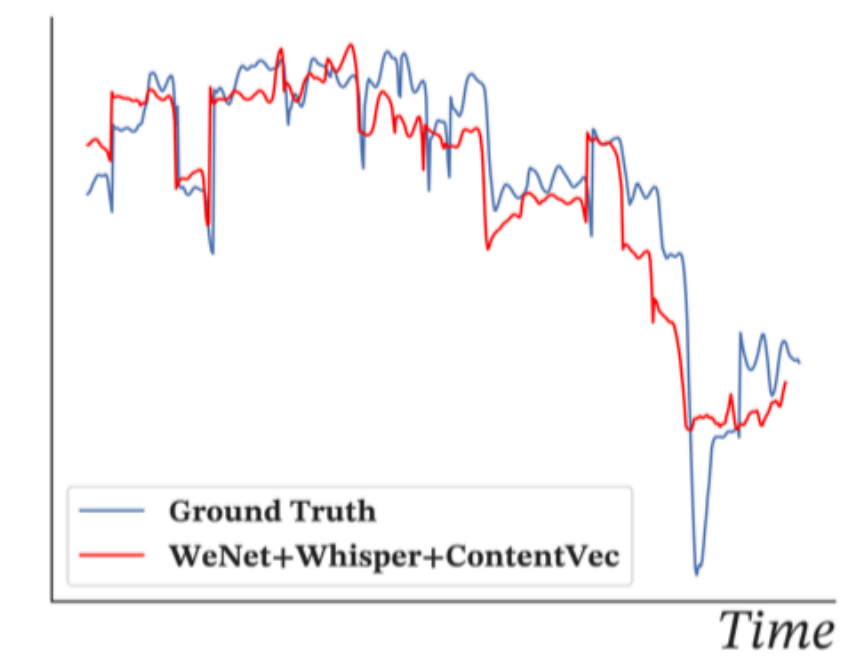


Source

WeNet



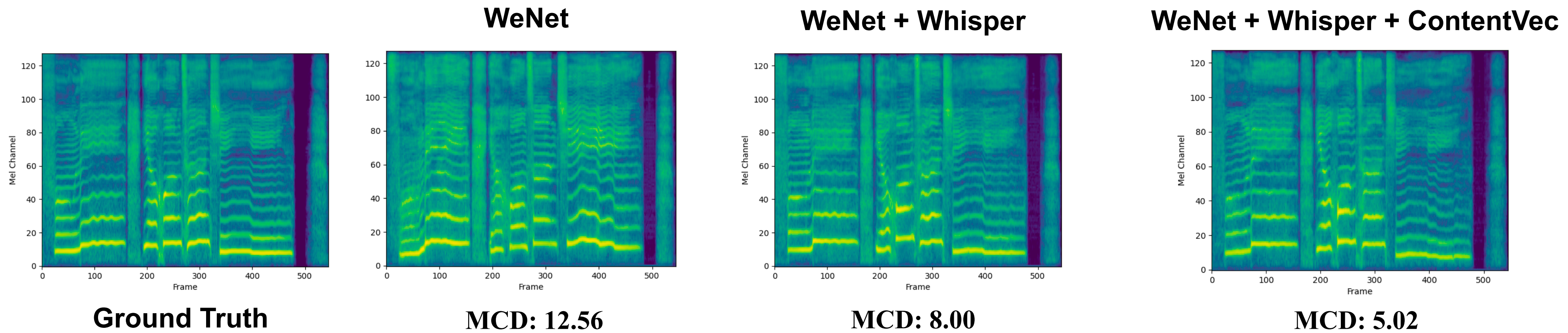
WeNet + Whisper



WeNet + Whisper + ContentVec



# Results: Complementary roles of Diverse Semantic-based Features



Spectrogram Reconstruction

The figure shows the transcription results for the same three models. Each result is presented in a box with the Chinese text, the Pinyin, and the CER percentage. The ground truth transcription is: 这只是刚刚入门接下来你还会会弹琴会写歌会双截棍. The WeNet model (CER: 81.8%) produces a completely incorrect transcription: 这只是当童话吗地下安利啊狗货看清单写报上节归. The WeNet + Whisper model (CER: 13.0%) produces a partially correct transcription: 这只是刚刚入门接下来你还会弹琴会写歌会双截棍. The WeNet + Whisper + ContentVec model (CER: 8.7%) produces a transcription that is almost identical to the ground truth: 就只是刚刚入门接下来你还会弹琴会写歌会双截棍.

**Ground Truth**      **WeNet**      **WeNet + Whisper**      **WeNet + Whisper + ContentVec**

CER: 81.8%      CER: 13.0%      CER: 8.7%

Intelligibility

# Results: SVC Framework based on Diverse Semantic-based Features Fusion

## Objective Evaluation

Base Model	Semantic-based Features	Recording Studio Setting				In-the-Wild Setting			
		F0CORR (↑)	F0RMSE (↓)	CER (↓)	SIM (↑)	F0CORR (↑)	F0RMSE (↓)	CER (↓)	SIM (↑)
TransformerSVC	WeNet	0.849	149.3	15.6%	0.878	0.871	210.0	40.0%	0.865
	WeNet + Whisper	0.924	77.2	14.9%	0.881	0.848	183.8	18.7%	0.867
	WeNet + Whisper + ContentVec	0.931	75.5	16.2%	0.883	0.857	186.7	23.3%	0.868
VitsSVC	WeNet	0.937	175.3	19.1%	0.890	0.919	91.3	57.7%	0.869
	WeNet + Whisper	0.945	144.4	17.8%	0.890	0.920	86.9	35.2%	0.869
	WeNet + Whisper + ContentVec	0.946	112.9	17.7%	0.886	0.921	79.5	32.3%	0.870
DiffWaveNetSVC	WeNet	0.936	55.5	15.8%	0.875	0.901	87.8	60.8%	0.855
	WeNet + Whisper	0.943	49.5	15.2%	0.884	0.921	73.6	21.1%	0.865
	WeNet + Whisper + ContentVec	0.940	55.2	15.7%	0.884	0.919	79.9	23.3%	0.867

## Subjective Evaluation for DiffWaveNetSVC

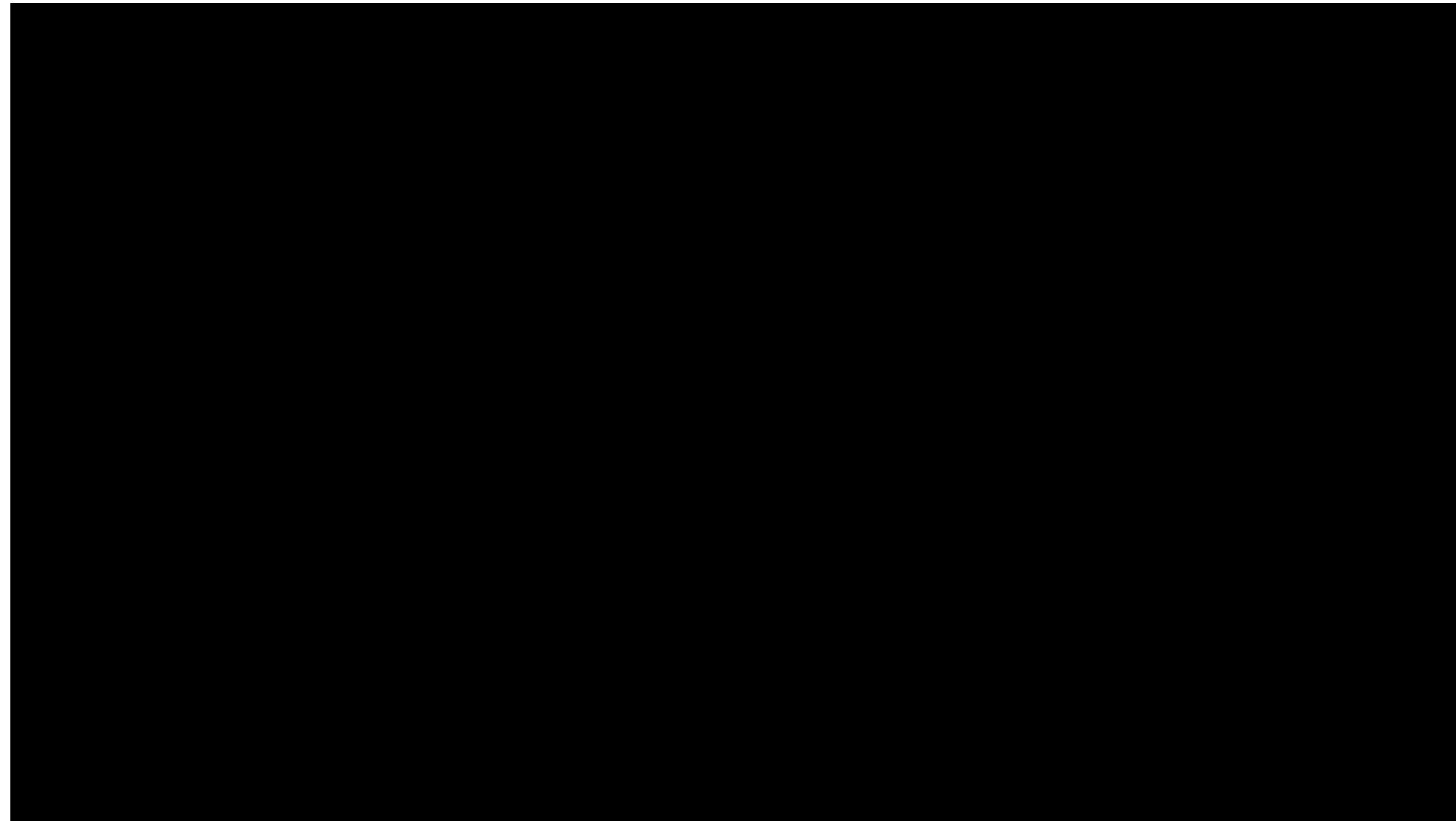
Semantic-based Features	Recording Studio Setting		In-the-Wild Setting	
	Naturalness (↑)	Similarity (↑)	Naturalness (↑)	Similarity (↑)
WeNet	2.72 ±0.22	2.64 ±0.21	2.85 ±0.21	2.34 ±0.20
WeNet + Whisper	4.02 ±0.18	3.13 ±0.17	3.70 ±0.18	2.86 ±0.23
WeNet + Whisper + ContentVec	4.14 ±0.19	3.25 ±0.18	3.71 ±0.18	2.82 ±0.23

The full scores of Naturalness and Similarity are 5 and 4

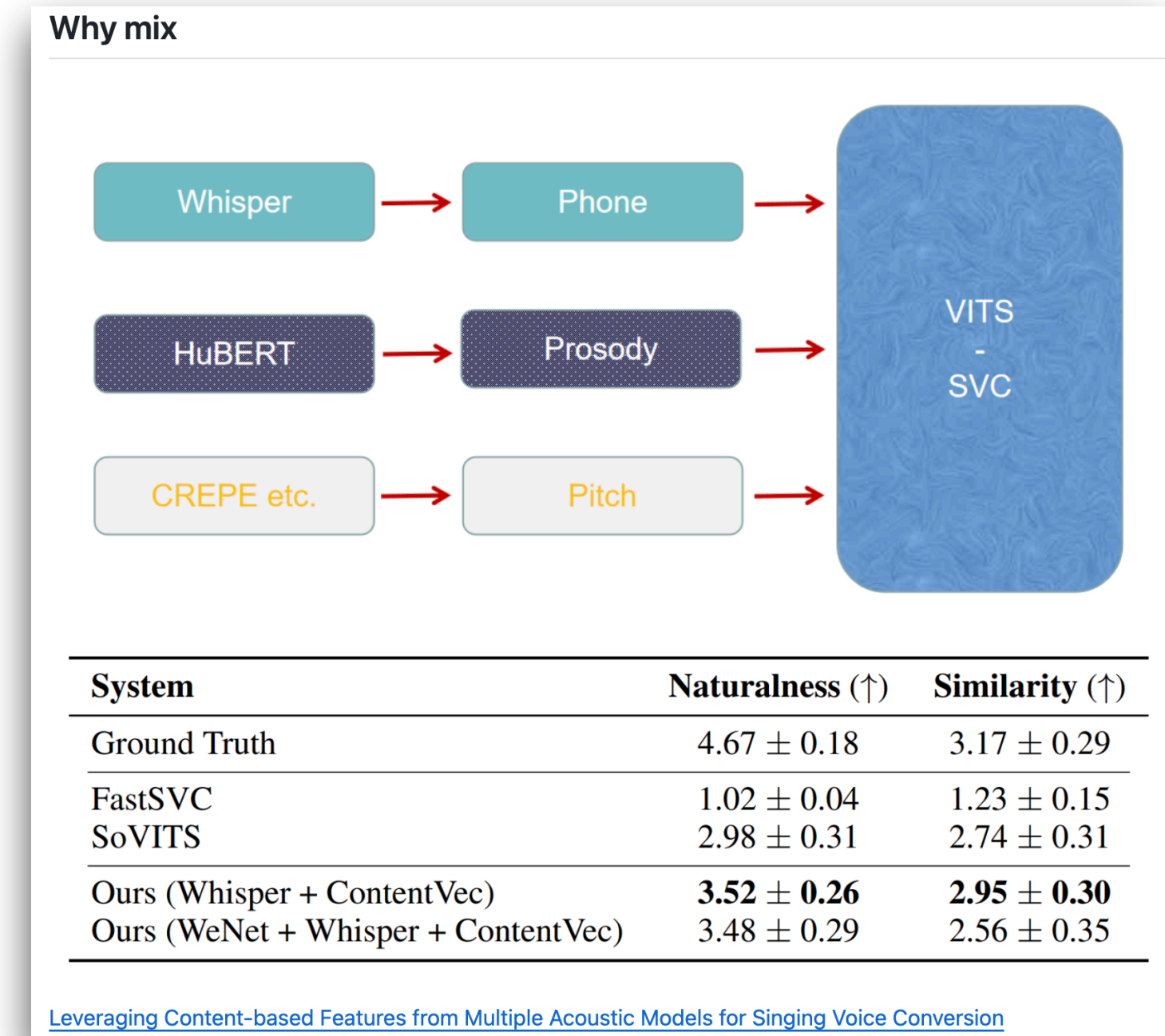
- ① **Generalization:** The idea of diverse semantic-based features fusion work for various base models in both settings.
- ② **Robustness:** for the more challenging in-the-wild setting, such solution is also effective.



# AI Singer Demo and Impact



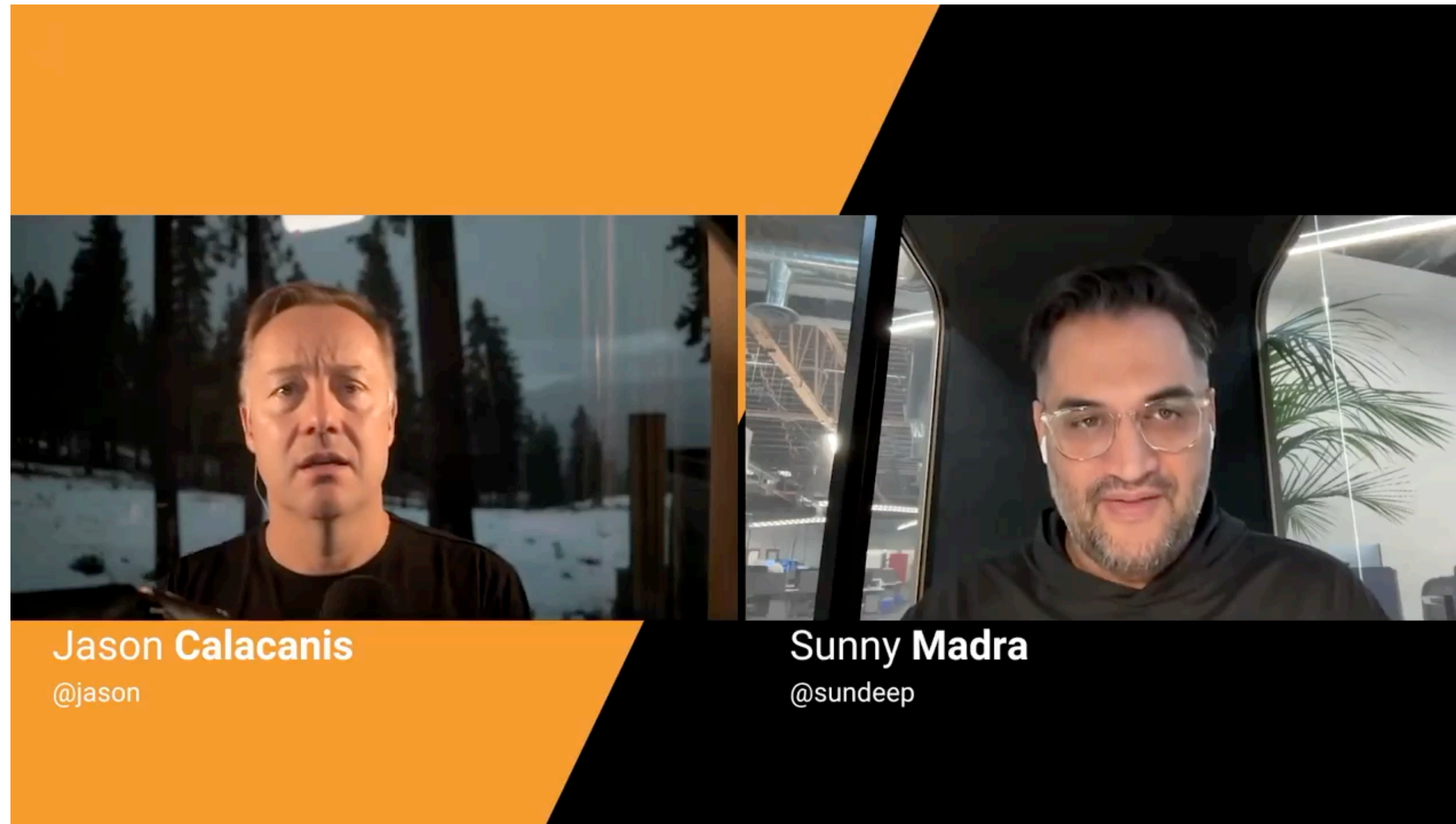
- ◆ Make Taylor Swift sing Mandarin song!



- ◆ Our idea of using multiple content features has been **borrowed and integrated into So-VITS-SVC 5.0** (Github over 2k stars)

# AI Singer Demo and Impact

---



- ◆ Highly positive comments from the market

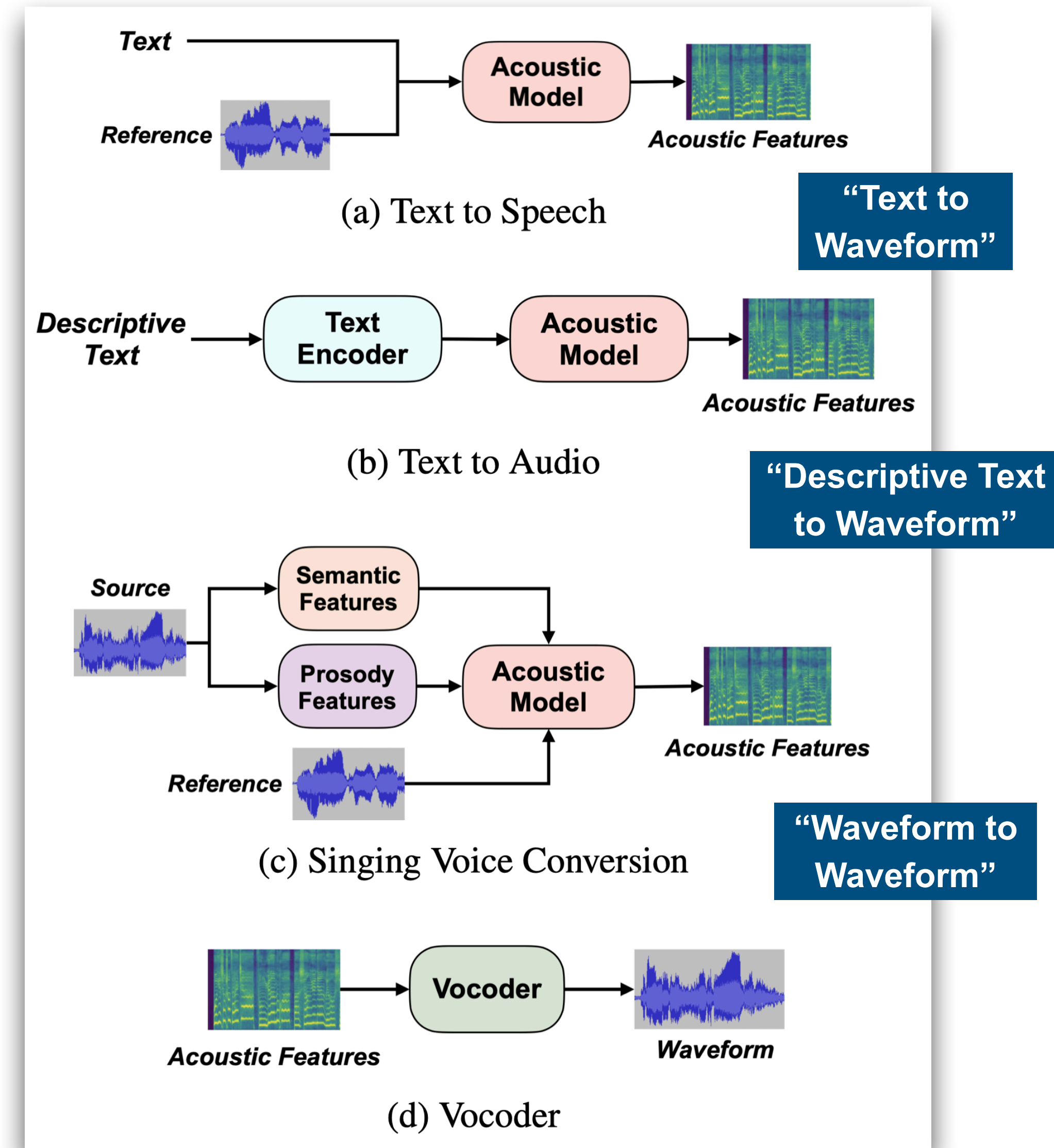
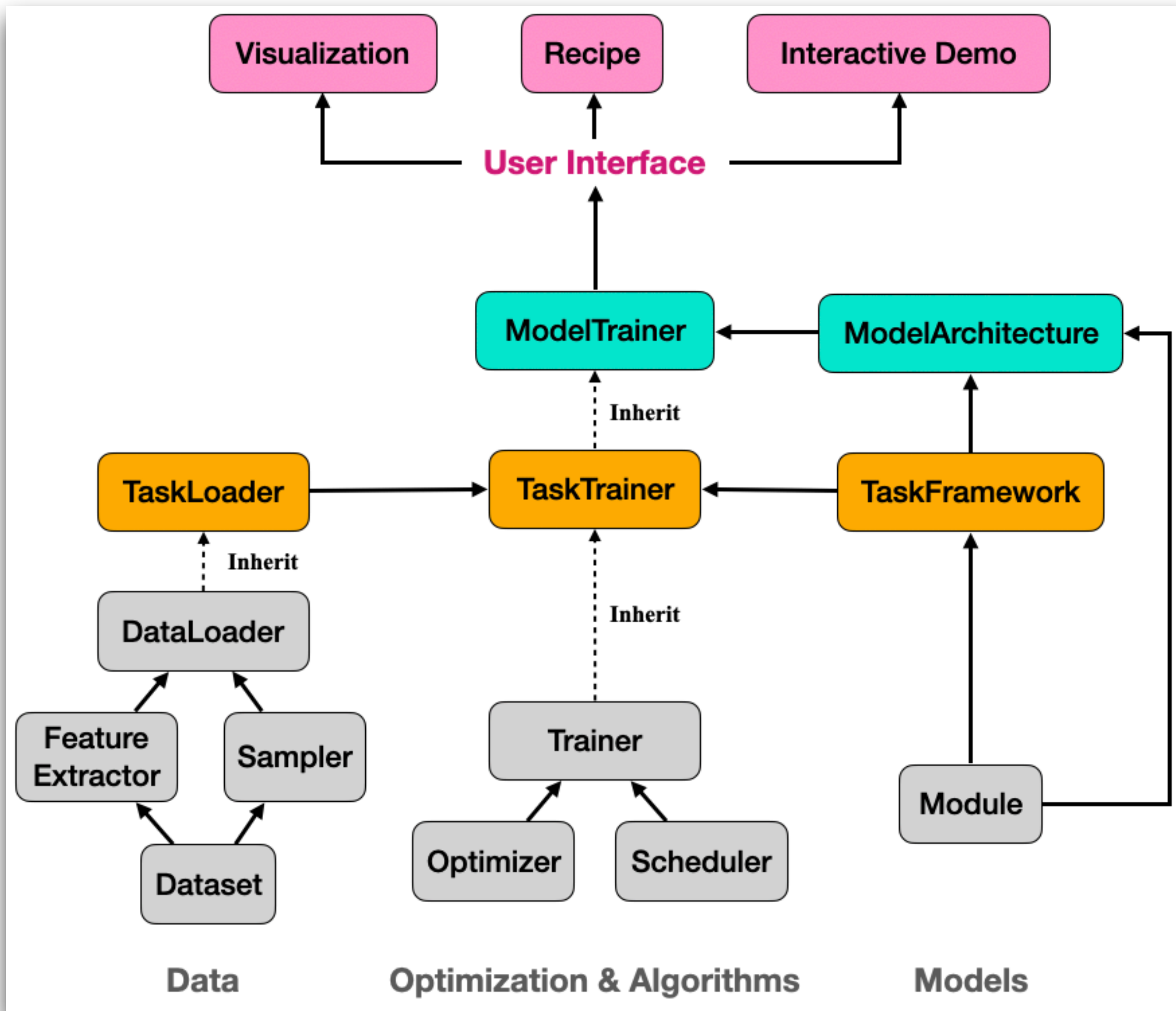


# Roadmap

---

- **Singing Voice Conversion**
  - Definition, Classic Works, and Modern Pipeline
- **Singing Voice Conversion in Amphion**
  - Supported Model Architectures
  - Our research: *Leveraging Diverse Semantic-based Audio Pretrained Models for Singing Voice Conversion*
- **Amphion's Philosophy**
  - Unique strengths, Supported Features, and Visualization

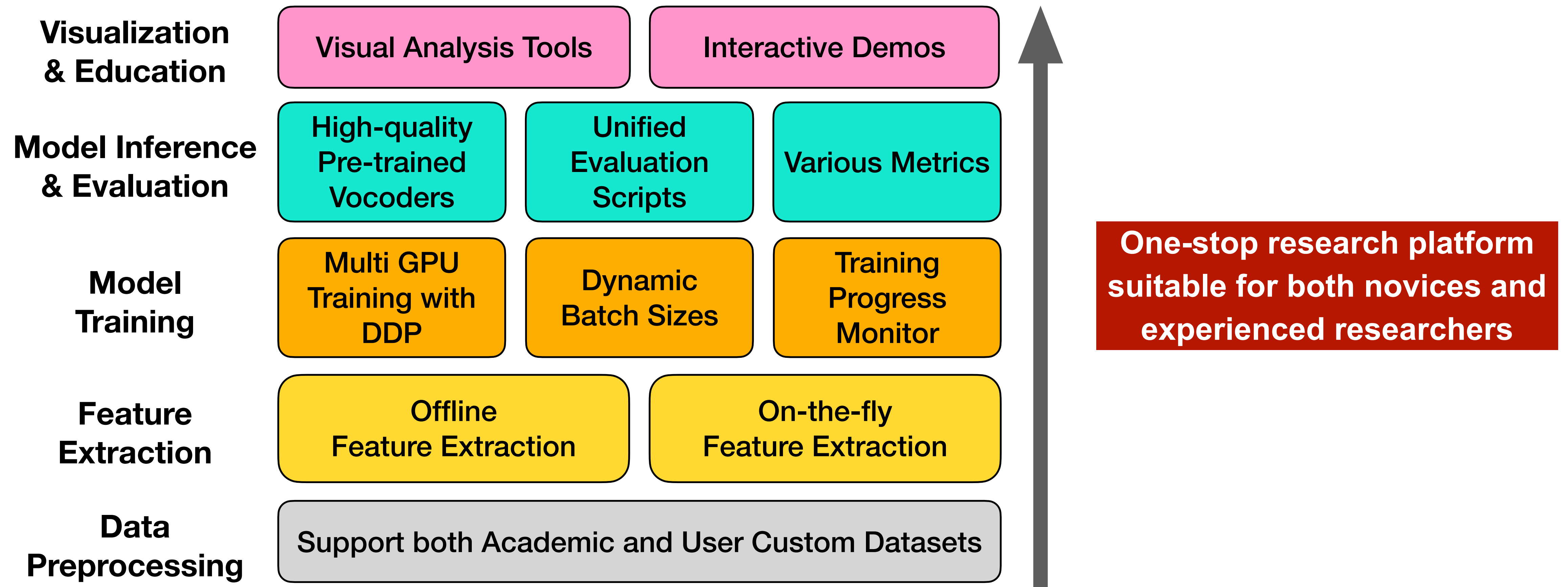
# Strength1: Unified Audio Generation Framework





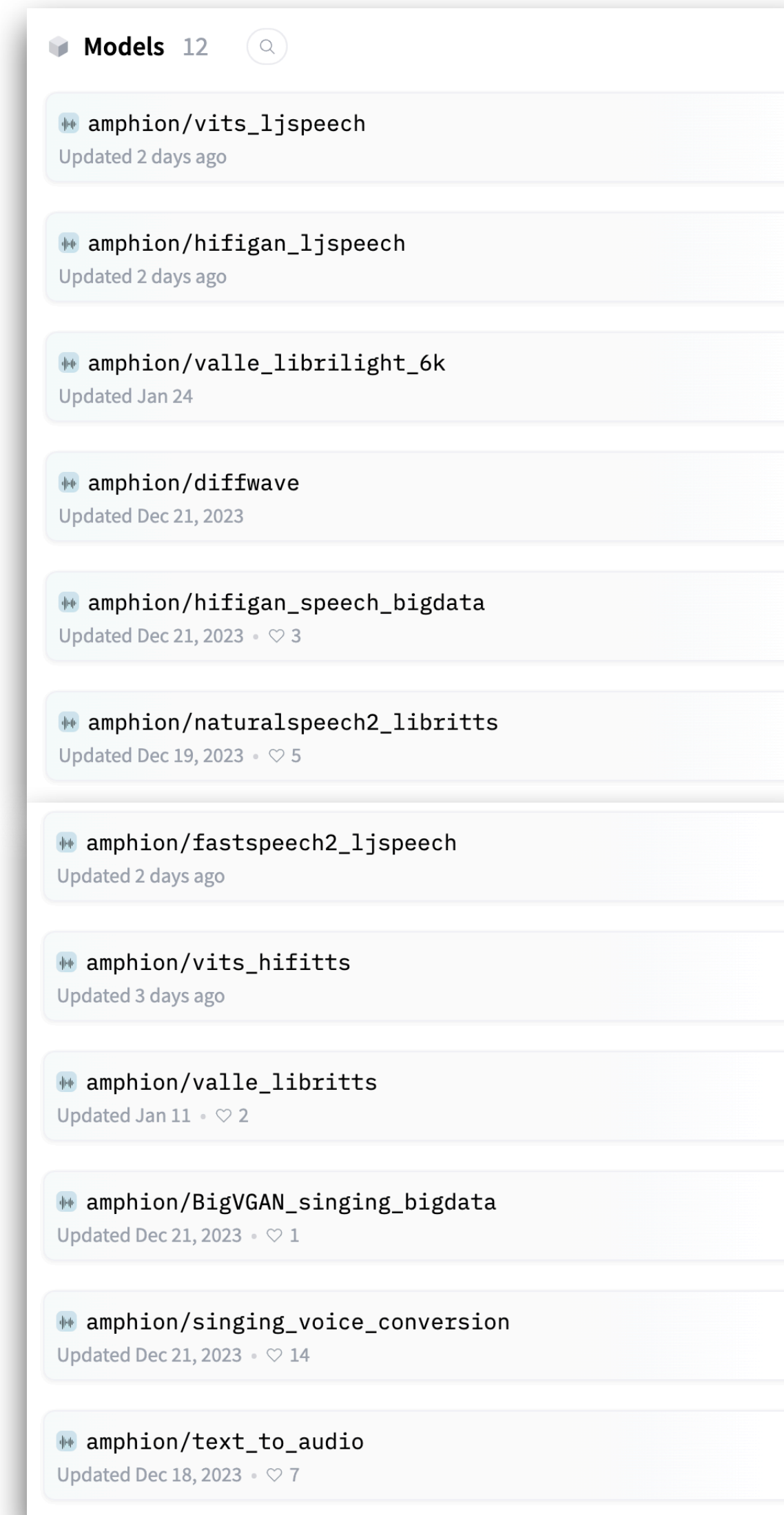
# Strength2: Beginner-friendly End-to-End Workflow

---



# Strength3: Open Pre-trained Models

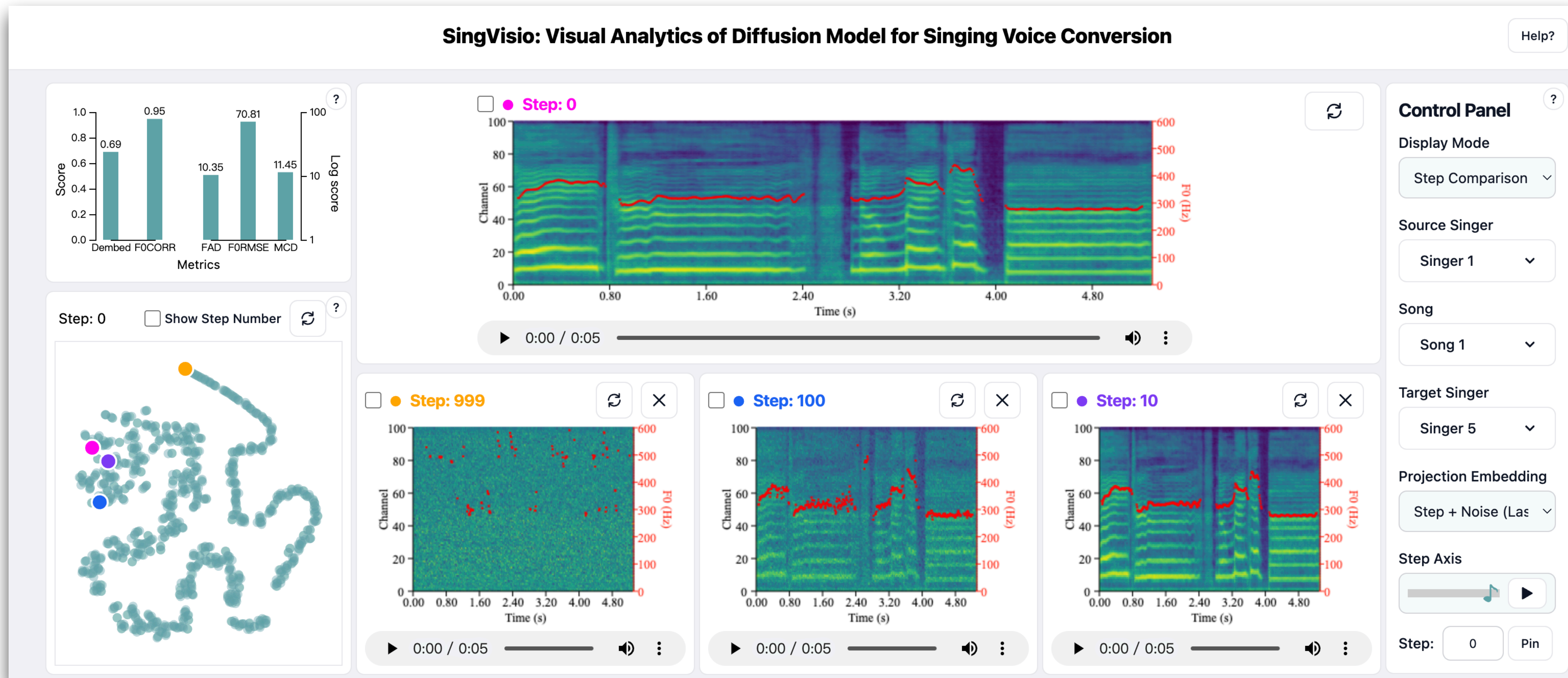
Release Criteria	Description
Model Metadata	Detail the model architecture and the number of parameters.
Training Datasets	List all the training corpus and their sources.
Training Configuration	Detail the training hyperparameters (like batch size, learning rate, and number of training steps) and the computational platform
Evaluation Results	Display the evaluation results and the performance comparison to other typical baselines.
Usage Instructions	Instruct how to inference and fine-tune based on the pre-trained model.
Interactive Demo	Provide an online interactive demo for users to explore.
License	Clear the licensing details including how the model can be utilized, shared, and modified.
Ethical Considerations	Address ethical considerations related to the model's application, focusing on privacy, consent, and bias, to encourage responsible usage.



**Supported  
Pretrained Models  
(Updating)**



# Strength4: Visualization and Interactivity



Liumeng Xue\*, Chaoren Wang\*, Mingxuan Wang, Xueyao Zhang, Jun Han, Zhizheng Wu. *SingVisio: Visual Analytics of Diffusion Model for Singing Voice Conversion*.



# THANKS



## Xueyao Zhang (张雪遥)

- ◆ **Second-year PhD student**, Supervised by Prof Zhizheng Wu  
School of Data Science, CUHK-Shenzhen  
Homepage: <https://www.zhangxueyao.com/>
- ◆ **Amphion v0.1's co-founder**  
Project: <https://github.com/open-mmlab/Amphion> (**3.5k stars**)
- ◆ **Research interest: "AI + Music"**, especially on:
  - Singing Voice Processing
  - Music Generation



**Amphion Official Account**

- 📎 **Amphion Technical Report:** <https://arxiv.org/abs/2312.09911>
- 👤 **Amphion GitHub:** <https://github.com/open-mmlab/Amphion>
- 🎧 **Amphion Demos/Models/Datasets:** <https://huggingface.co/amphion>



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen