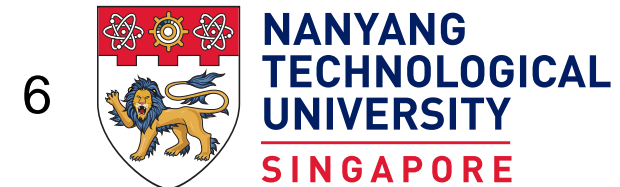


SongBsAb: A Dual Prevention Approach against Singing Voice Conversion based Illegal Song Covers

Guangke Chen¹, Yedi Zhang², Fu Song³
Ting Wang⁴, Xiaoning Du⁵, Yang Liu⁶



AI Generative Music

AI-based New Music Generation



MusicGPT 



AI-based Automated Song Cover

covers.ai



Jammable



Singify



Media.io

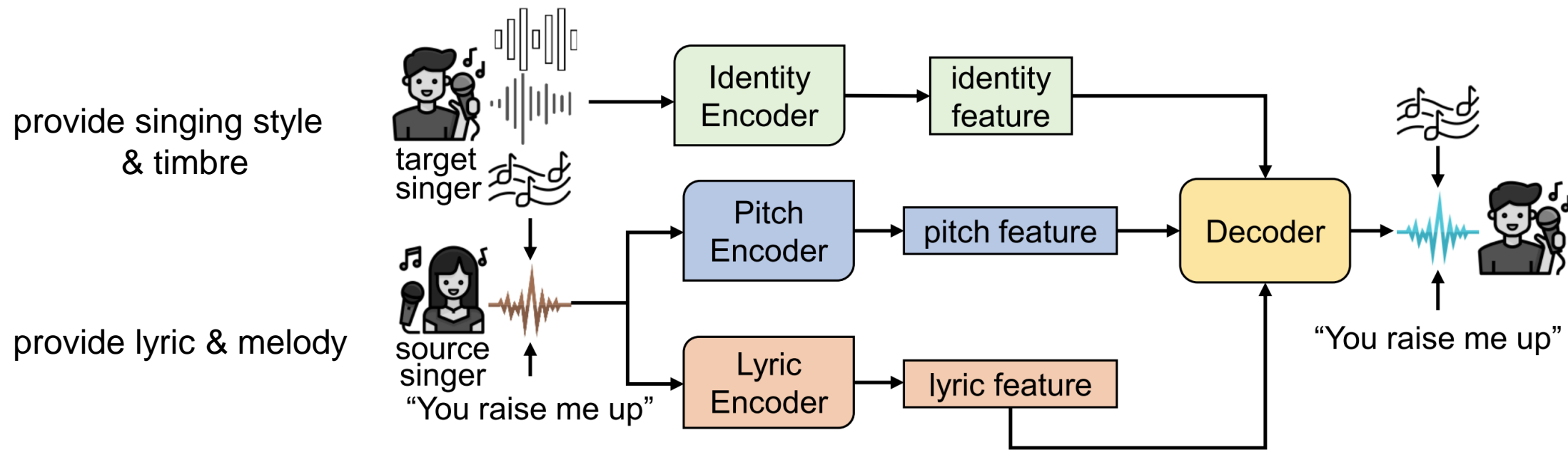
TopMedia 



Musicfy

AI-based Automated Song Cover by Singing Voice Conversion (SVC)

- transforms a song's vocal rendition from one singer to another's singing style and timbre while preserving the original lyrics and melody



SVC-based song covers

Singing Voice Conversion (SVC): Challenging Music Industry

- low entry barriers →
wide spread of AI-based song cover



Popular “AI Sun Yanzi” in China

An artificial intelligence-generated song, which mimics the voices of **Drake** and **The Weeknd** with terrifying accuracy, has been submitted for **Grammy** consideration.

■ The Impact on Music Industry:

- Infringement of singers' civil rights over voices & reputation
- Infringement of record companies' rights to release & distribute songs
- Infringement of the copyright of lyrics and melodies
- Erosion of singers' skill competitiveness (rely on for livelihood)
- Unfair competition faced by record companies



What should we do?

■ Reactive Detection:

reactive detection	SingFake-T02			low accuracy
	Method	Mixture	Vocals	
	AASIST	58.12	37.91	
	Spectrogram+ResNet	51.87	37.65	
	LFCC+ResNet	45.12	54.88	
	Wav2Vec2+AASIST	56.75	57.26	

Singing Voice Deepfake Detection Challenge [1]

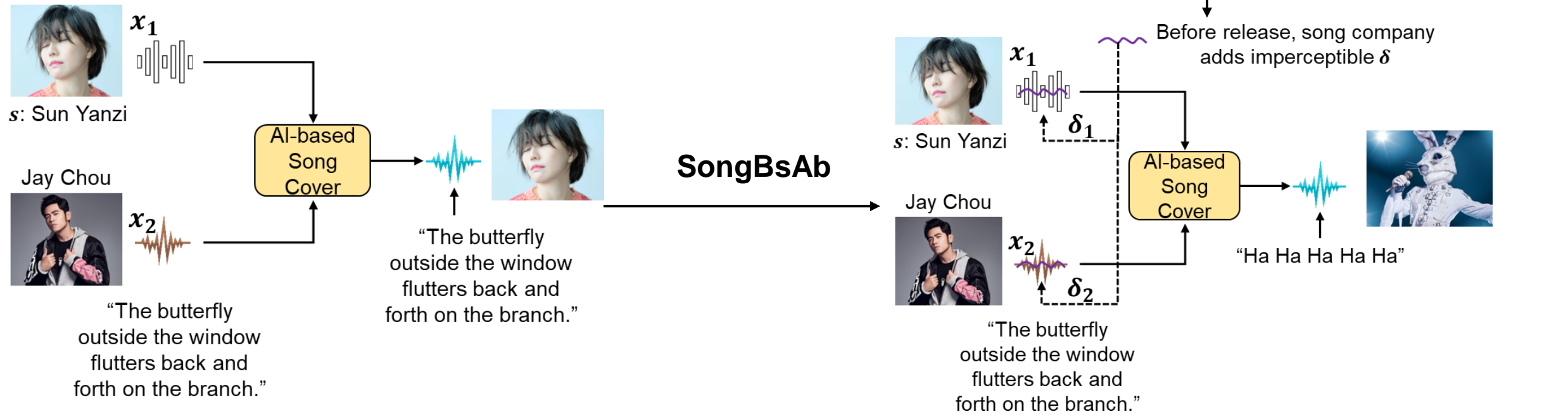


- Infringement already committed
- High quality, hard to detect
- Large number, inefficient to detect

[1] W. Huang, L. P. Violeta, S. Liu, J. Shi, Y. Yasuda, and T. Toda, “The singing voice conversion challenge 2023,” CoRR, vol. abs/2306.14422, 2023.

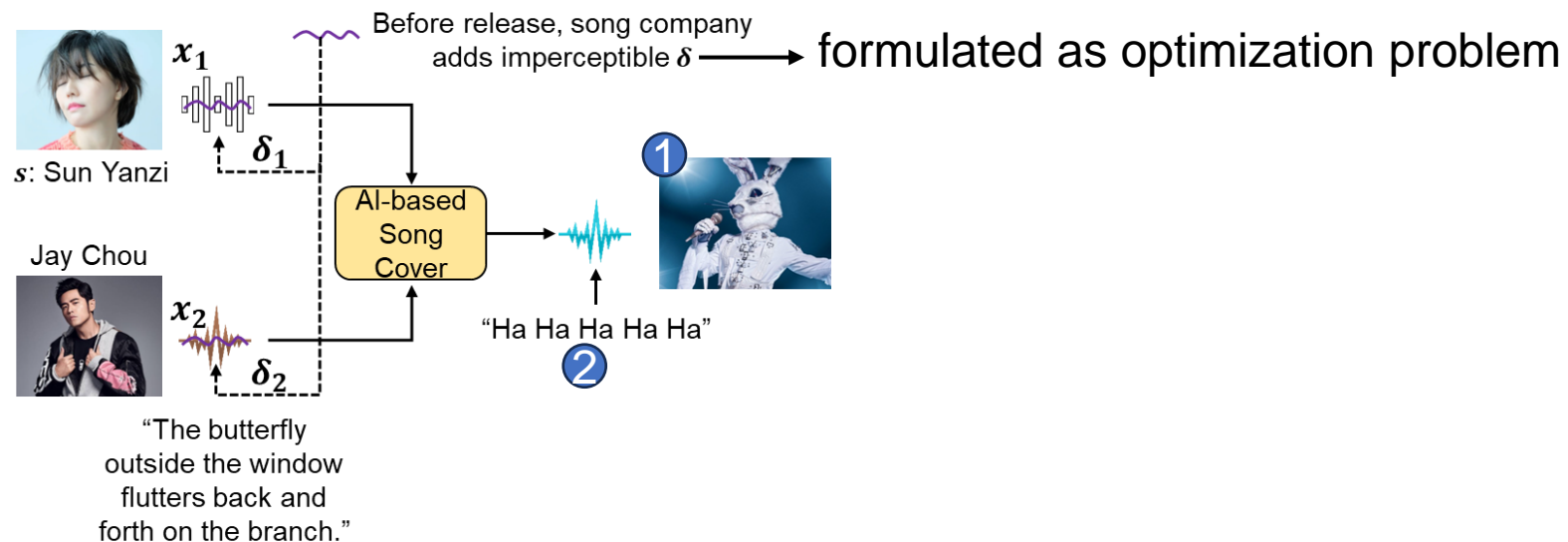
What should we do?

■ Proactive Prevention:



■ **Challenge:** do not know in advance if songs will be used for target or source songs

■ **Solution:** Dual Prevention



• ① Identity disruption

Gender-Transformation Loss

$$f_{ID} = \text{Distance}(\Theta(x + \delta), \Theta(s'))$$

Θ singer timbre extractor
 s' opposite-gender singers

• ② Lyric disruption

High/Low Hierarchy Loss

$$f_{LD} = \text{Distance}(\Phi_H(x + \delta), \Phi_H(x')) \\ + \text{Distance}(\Phi_L(x + \delta), \Phi_L(x'))$$

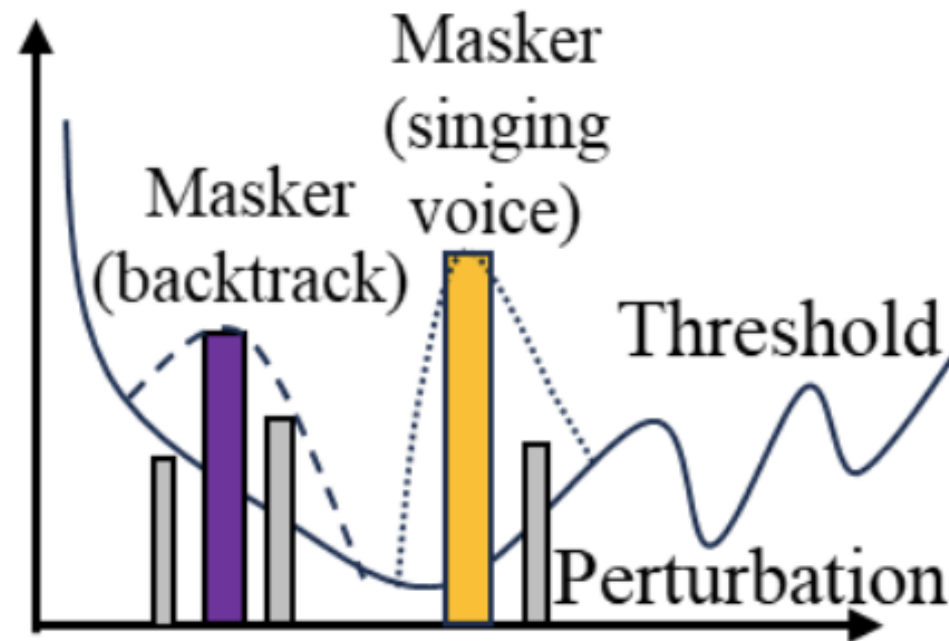
Φ_H high-level lyric features

Φ_L low-level lyric features

x' A song with different lyrics

$$\arg \min_{\delta} f_{ID} + f_{LD}$$

- **Challenge:** high quality requirements of songs
- **Solution:** backing track-refined simultaneous masking



the perturbation will not be audible as long as it is under one of the masking thresholds of the singing voice and the backing track

■ **Challenge:** Transferability to unknown SVC models exploited by adversaries

■ **Solution:** encoder ensemble & frame-level interaction reduction loss

- **encoder ensemble:** craft perturbation on multiple local white-box SVC models

- **frame-level interaction reduction loss:**

no perturbation
↓

perturbation interaction [2]:
$$\frac{\mathbb{E}_i(v(\Omega) + v(\emptyset) - v(\Omega \setminus \{i\}) - v(\{i\}))}{n - 1}$$

Ω : set of perturbation units
 \emptyset : no perturbation
 $\Omega/\{i\}$: only unit i not perturbed
 $\{i\}$: only unit i perturbed

interaction is negatively correlated with transferability \longrightarrow minimize the interaction loss at frame-level

[2] X. Wang, J. Ren, S. Lin, X. Zhu, Y. Wang, and Q. Zhang, “A unified approach to interpreting and boosting adversarial transferability,” in ICLR, 2021.

- **Few-shot SVC model:** Lora-SVC, Vits-SVC, Grad-SVC, NeuCo-SVC
- **Dataset:** OpenSinger (Chinese), NUS-48E (English)
- **Metric:** Identity Similarity with target singer; lyric Word Error Rate (WER)
- **Baseline:** Attack-VC [3], AntiFake [4]

[3] C. Huang, et al, “Defending your voice: Adversarial attack on voice conversion,” in SLT, 2021.

[4] Z. Yu, et al, “Antifake: Using adversarial audio to prevent unauthorized speech synthesis,” in CCS, 2023.

y Covered songs w/o prevention

\tilde{y} Covered songs w/ prevention

Dataset	SVC Model	Approach	Prevention Effectiveness			
			Identity Similarity ↓		Lyric WER (%) ↑	
			y	\tilde{y}	y	\tilde{y}
Open Singer	Lora-SVC	AntiFake		0.15		13.2
		AttackVC	0.54	0.55	13.9	13.9
		AttackVC-W		0.24		13.1
		SongBsAb		0.05		76.1
	Vits-SVC	AntiFake		0.15		15.2
		AttackVC	0.51	0.26	14.9	14.7
		AttackVC-W		0.09		90.4
		SongBsAb				
	Grad-SVC	AntiFake		0.17		31.4
		AttackVC	0.48	0.23	32.1	30.9
		AttackVC-W		0.11		103.6
		SongBsAb				
	NeuCo-SVC	AntiFake		0.33		20.8
		AttackVC	0.65	0.28	18.1	20.1
		AttackVC-W		0.22		86.5
		SongBsAb				

Dataset	SVC Model	Approach	Prevention Effectiveness			
			Identity Similarity ↓		Lyric WER (%) ↑	
			y	\tilde{y}	y	\tilde{y}
NUS-48E	Lora-SVC	AntiFake		0.22		22.0
		AttackVC	0.47	0.25	23.3	23.1
		AttackVC-W		0.12		79.9
		SongBsAb				
	Vits-SVC	AntiFake		0.19		19.4
		AttackVC	0.48	0.25	18.4	19.3
		AttackVC-W		0.12		78.4
		SongBsAb				
	Grad-SVC	AntiFake		0.24		41.2
		AttackVC	0.45	0.24	41.1	43.6
		AttackVC-W		0.16		94.5
		SongBsAb				
	NeuCo-SVC	AntiFake		0.24		22.7
		AttackVC	0.59	0.22	22.6	21.7
		AttackVC-W		0.16		76.6
		SongBsAb				



Input: target singer (style & timbre)



Input: source singer (lyric & melody)



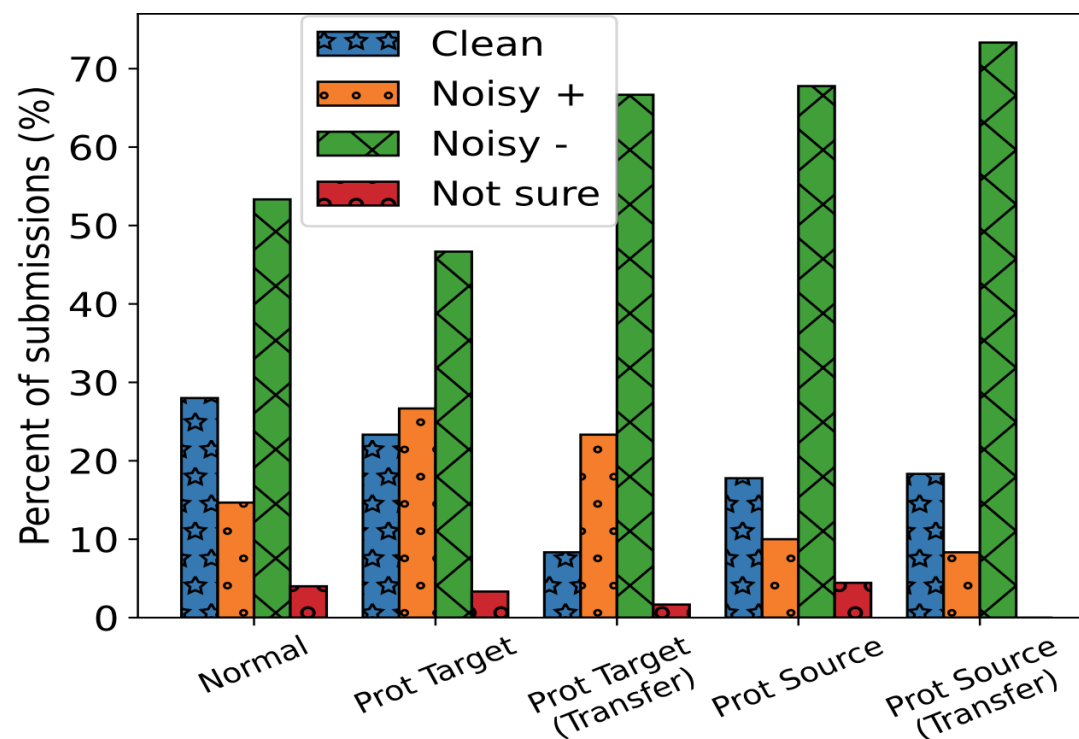
Output w/o SongBsAb



Output w/ SongBsAb

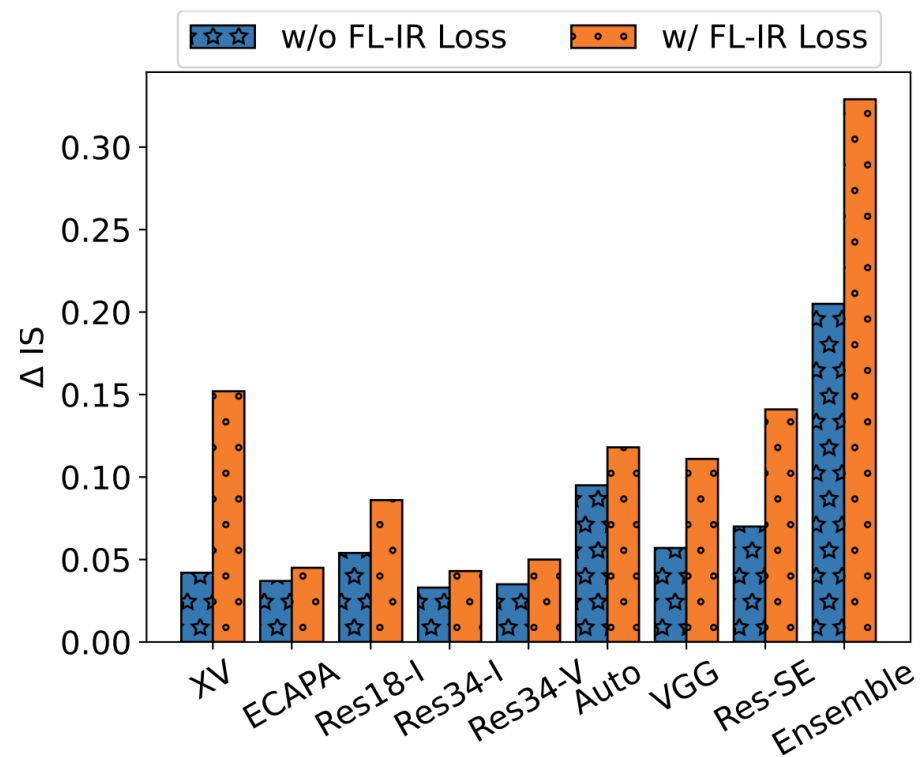
■ Impact on song quality and enjoyment experience

human study: if a given song contains any background noise and if so, how the noise influences their enjoyment of the song



“Noise +” and “Noise -” denote the answers “noisy w/ influence” and “noisy w/o influence”

■ Transferability

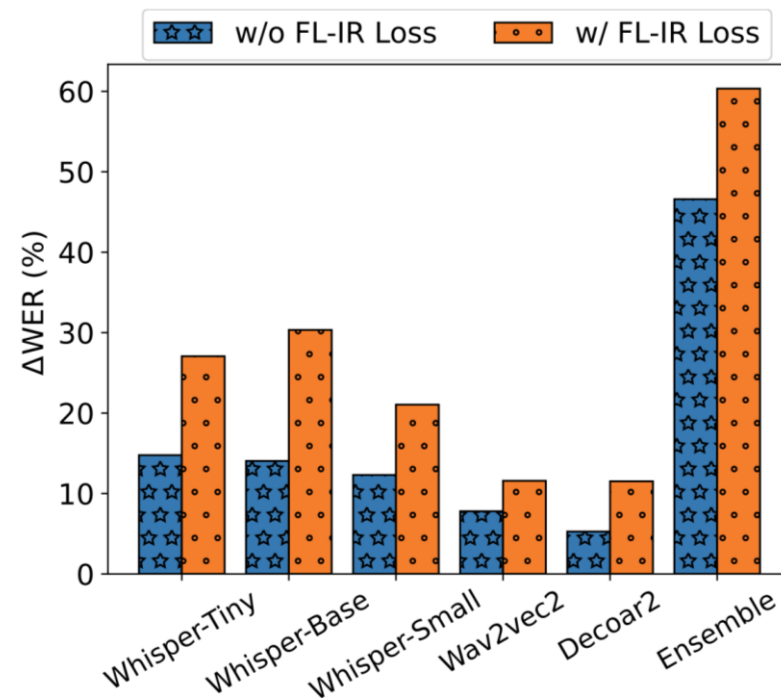


Identity Disruption

ΔIS : decrease of identity similarity

Ensemble: encoder ensemble

FL-IR: frame-level interaction reduction



Lyric Disruption

ΔWER : increase of lyric Word Error Rate

Ensemble: encoder ensemble

FL-IR: frame-level interaction reduction

Take away

- The first proactive prevention against singing voice conversion-based illegal song covers
- Dual prevention: identity & lyric disruption
- Backing track-refined simultaneous masking to preserve song quality
- Encoder ensemble & frame-level interaction reduction loss to enhance transferability
- Application: copyright & civil rights protection by record companies & singers

Website: <https://sites.google.com/view/songbsab>

Any Question?
Thanks!