

# FINALLY

## Fast and Universal Speech with Studio-like Quality

Nicholas Babaev\*\*, Kirill Tamogashev\*\*, Azat Saginbaev\*, Ivan Shchekotov\*,  
Hanbin Bae\*, Hosang Sung\*, WonJun Lee\*, Hoon-Young Cho\*, Pavel Andreev\*\*

\*Equal Contribution \*Samsung Research

# Mode Collapse and Speech Enhancement

Speech Enhancement aims to restore noisy or degraded sound into the clean one

## **Common probabilistic formulation:**

Learn a distribution of clean signals given a noisy one.

$y \sim p_{clean}(y|x)$ ,  $x$  – noisy signal,  $y$  – clean signal

# Mode Collapse and Speech Enhancement

Speech Enhancement aims to restore noisy or degraded sound into the clean one

## **Common probabilistic formulation:**

Learn a distribution of clean signals given a noisy one.

$y \sim p_{clean}(y|x)$ ,  $x$  – noisy signal,  $y$  – clean signal

## **Proposed probabilistic formulation:**

Learn the most likely clean signal given a noisy one.

$y = \operatorname{argmax}_y p_{clean}(y|x)$ ,  $x$  – noisy signal,  $y$  – clean signal

# Mode Collapse and Speech Enhancement

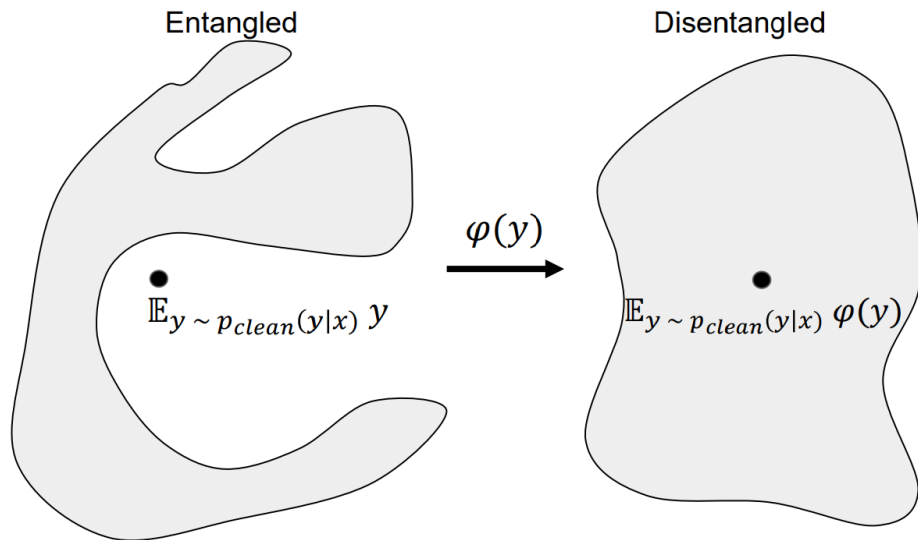
We show that GANs are a natural choice for predicting the main mode of the conditional clean speech distribution  $p_{clean}(y|x)$

**Proposition 1.** *Let  $p_{clean}(y|x) > 0$  be a finite and Lipschitz continuous density function with a unique global maximum and  $p_g^\xi(y|x) = \xi^n / 2^n \cdot \mathbf{1}_{y-g_\theta(x) \in [-1/\xi, 1/\xi]^n}$ , then*

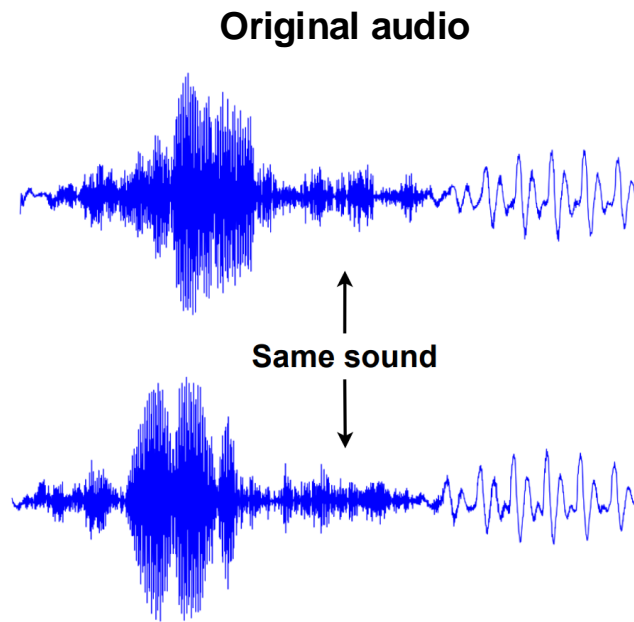
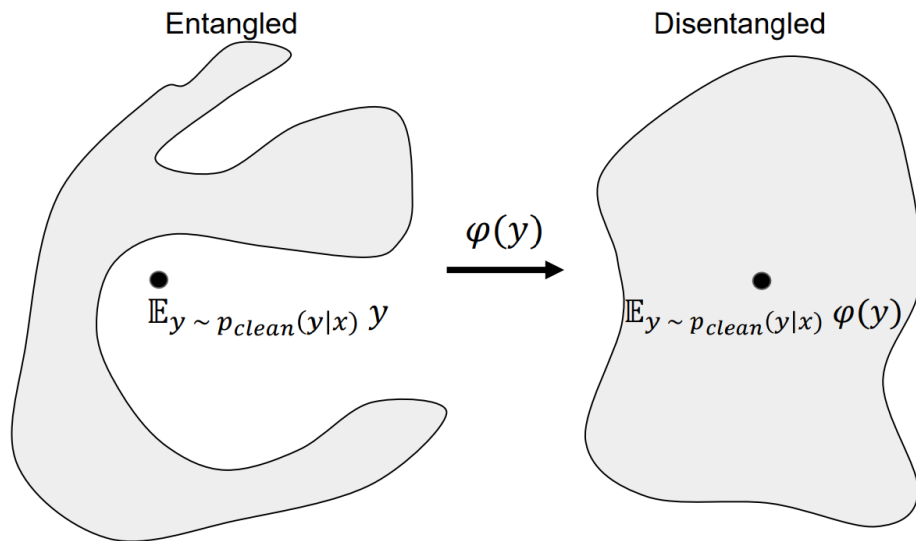
$$\lim_{\xi \rightarrow +\infty} \arg \min_{g_\theta(x)} \chi_{Pearson}^2(p_g^\xi || (p_{clean} + p_g^\xi)/2) = \arg \max_y p_{clean}(y|x)$$

However, **adversarial training is unstable**, therefore, **auxiliary regression losses** are needed to push generator close to the desirable solution.

# Practical Aspects: an issue with regression losses



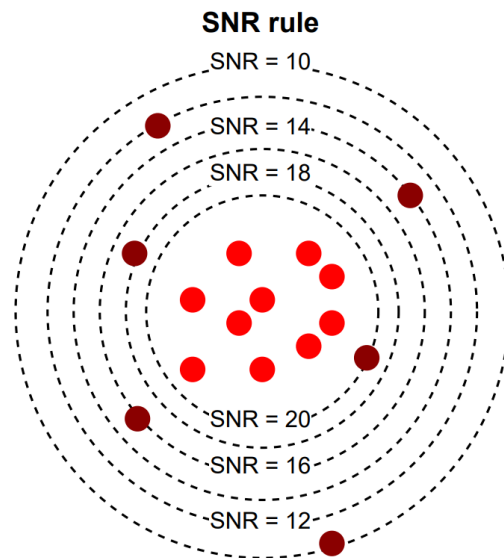
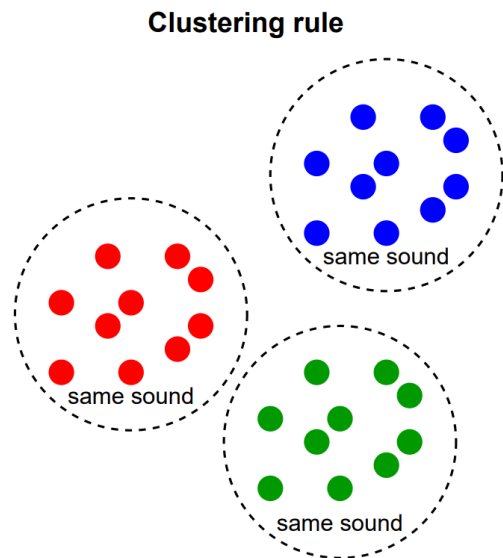
# Practical Aspects: an issue with regression losses



Audio, reconstructed with HiFi-GAN vocoder

# Practical Aspects: rules for an ideal regression loss

1. **Clustering rule:** same speech sounds should be closer to each other, while different speech sounds should be separated.
2. **SNR rule:** the higher the noise level on the noisy signal, the further it should be from its clean version.



# Practical Aspects: choice of the mapping $\phi$

Feature space	Rand score ( $\uparrow$ ) (Clustering rule)	Negative correlation ( $\uparrow$ ) (SNR rule)	MOS ( $\uparrow$ ) (Vocoding)
Waveform	$0.00 \pm 0.00$	$0.31 \pm 0.02$	Failed
Spectrogram	$0.00 \pm 0.00$	$0.08 \pm 0.03$	$1.78 \pm 0.08$
Wav2Vec 2.0	$0.25 \pm 0.03$	$0.19 \pm 0.03$	$1.65 \pm 0.08$
Wav2Vec 2.0-conv	$0.94 \pm 0.01$	$0.78 \pm 0.02$	$2.23 \pm 0.09$
WavLM	$0.46 \pm 0.05$	$0.46 \pm 0.03$	$1.71 \pm 0.07$
WavLM-conv	<b><math>0.96 \pm 0.01</math></b>	<b><math>0.89 \pm 0.02</math></b>	<b><math>3.27 \pm 0.10</math></b>
EnCodec	$0.55 \pm 0.03$	$0.67 \pm 0.03$	$1.80 \pm 0.08$
CDPAM	$0.00 \pm 0.00$	$0.17 \pm 0.03$	Failed

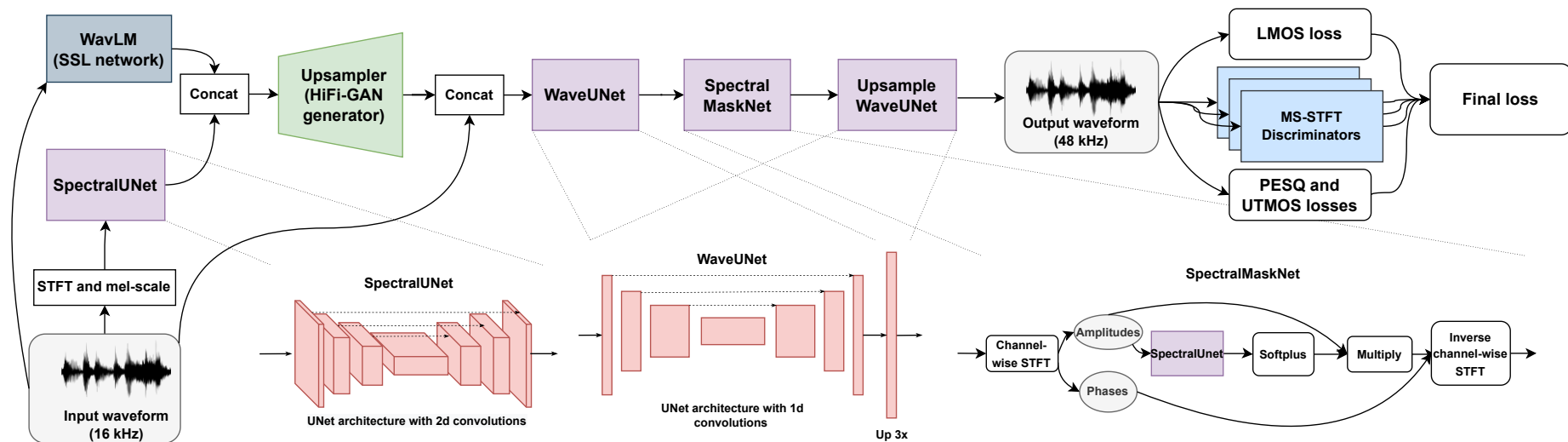


# Practical Aspects: choice of the mapping $\phi$

Feature space	Rand score ( $\uparrow$ ) (Clustering rule)	Negative correlation ( $\uparrow$ ) (SNR rule)	MOS ( $\uparrow$ ) (Vocoding)
Waveform	$0.00 \pm 0.00$	$0.31 \pm 0.02$	Failed
Spectrogram	$0.00 \pm 0.00$	$0.08 \pm 0.03$	$1.78 \pm 0.08$
Wav2Vec 2.0	$0.25 \pm 0.03$	$0.19 \pm 0.03$	$1.65 \pm 0.08$
Wav2Vec 2.0-conv	$0.94 \pm 0.01$	$0.78 \pm 0.02$	$2.23 \pm 0.09$
WavLM	$0.46 \pm 0.05$	$0.46 \pm 0.03$	$1.71 \pm 0.07$
WavLM-conv	<b><math>0.96 \pm 0.01</math></b>	<b><math>0.89 \pm 0.02</math></b>	<b><math>3.27 \pm 0.10</math></b>
EnCodec	$0.55 \pm 0.03$	$0.67 \pm 0.03$	$1.80 \pm 0.08$
CDPAM	$0.00 \pm 0.00$	$0.17 \pm 0.03$	Failed

$$\mathcal{L}_{\text{LMOS}}(\theta) = \mathbb{E}_{x, y \sim p(x, y)} \left[ 100 \cdot \|\phi(y) - \phi(g_\theta(x))\|_2^2 + \left| \|\text{STFT}(y)\| - \|\text{STFT}(g_\theta(x))\| \right| \right]$$

# FINALLY Architecture



# Training in 3 stages

1. Training in 16kHz without Upsample WaveUNet and with regression losses only.

$$\mathcal{L}_{\text{gen}}(\theta) = \underbrace{\lambda_{\text{LMOS}} \cdot \mathcal{L}_{\text{LMOS}}(\theta)}_{\text{1st stage (16 kHz)}}$$

# Training in 3 stages

1. Training in 16kHz without Upsample WaveUNet and with regression losses only
2. Adversarial training in 16 kHz without Upsample WaveUNet with GAN loss and features matching loss.

$$\mathcal{L}_{\text{gen}}(\theta) = \underbrace{\lambda_{\text{LMOS}} \cdot \mathcal{L}_{\text{LMOS}}(\theta)}_{\text{1st stage (16 kHz)}} + \overbrace{\lambda_{\text{GAN}} \cdot \mathcal{L}_{\text{GAN-gen}}(\theta) + \lambda_{\text{FM}} \cdot \mathcal{L}_{\text{FM}}(\theta)}^{\text{2nd stage (16 kHz)}}$$

$$\mathcal{L}_{\text{disc}}(\varphi_i) = \mathcal{L}_{\text{GAN-disc}}(\varphi_i), \quad i = 1, \dots, k.$$

# Training in 3 stages

1. Training in 16kHz without Upsample WaveUNet and with regression losses only.
2. Adversarial training in 16 kHz without Upsample WaveUNet with GAN loss and features matching loss.
3. Adversarial Training with Upsample WaveUNet in 48 kHz using all previous losses and also Human Feedback losses.

$$\mathcal{L}_{\text{gen}}(\theta) = \underbrace{\lambda_{\text{LMOS}} \cdot \mathcal{L}_{\text{LMOS}}(\theta)}_{\text{1st stage (16 kHz)}} + \underbrace{\lambda_{\text{GAN}} \cdot \mathcal{L}_{\text{GAN-gen}}(\theta) + \lambda_{\text{FM}} \cdot \mathcal{L}_{\text{FM}}(\theta)}_{\text{2nd stage (16 kHz)}} + \lambda_{\text{HF}} \cdot \mathcal{L}_{\text{HF}}(\theta),$$

3rd stage (48 kHz)

$$\mathcal{L}_{\text{disc}}(\varphi_i) = \mathcal{L}_{\text{GAN-disc}}(\varphi_i), \quad i = 1, \dots, k.$$

# Evaluations and Results

Dataset	Model	MOS ( $\uparrow$ )	RTF ( $\downarrow$ )
VoxCeleb	Input	$3.46 \pm 0.07$	-
	HiFi-GAN-2 (by Adobe)	$4.47 \pm 0.05$	0.5
	Ours	<b><math>4.63 \pm 0.04</math></b>	<b>0.03</b>
UNIVERSE (validation set)	Input	$2.87 \pm 0.05$	-
	UNIVERSE (by Dolby)	$4.10 \pm 0.07$	0.5
	Ours	<b><math>4.23 \pm 0.07</math></b>	<b>0.03</b>
LibriTTS	Input	$3.59 \pm 0.07$	-
	MIIPHER (by Google)	$4.18 \pm 0.06$	N/A
	Ours	<b><math>4.54 \pm 0.05</math></b>	<b>0.03</b>

# Evaluations and Results

**VoxCeleb (HiFi-GAN-2 validation set, real data)**

Model	MOS ( $\uparrow$ )	UTMOS ( $\uparrow$ )	WV-MOS ( $\uparrow$ )	DNSMOS ( $\uparrow$ )	-	RTF ( $\downarrow$ )
Input	$3.46 \pm 0.07$	$2.76 \pm 0.13$	$2.90 \pm 0.16$	$2.72 \pm 0.11$	-	-
VoiceFixer	$3.41 \pm 0.07$	$2.60 \pm 0.09$	$2.79 \pm 0.09$	$3.08 \pm 0.06$	-	0.02
DEMUCS	$3.79 \pm 0.07$	$3.51 \pm 0.08$	$3.72 \pm 0.08$	$3.27 \pm 0.04$	-	0.08
STORM	$3.75 \pm 0.06$	$3.29 \pm 0.08$	$3.54 \pm 0.09$	$3.17 \pm 0.04$	-	1.05
BBED	$3.97 \pm 0.06$	$3.30 \pm 0.10$	$3.47 \pm 0.08$	$3.23 \pm 0.04$	-	0.43
HiFi-GAN-2	$4.47 \pm 0.05$	$3.67 \pm 0.09$	<b><math>3.96 \pm 0.06</math></b>	<b><math>3.32 \pm 0.03</math></b>	-	0.50
Ours	<b><math>4.63 \pm 0.04</math></b>	<b><math>4.05 \pm 0.07</math></b>	<b><math>3.98 \pm 0.06</math></b>	<b><math>3.31 \pm 0.04</math></b>	-	<b>0.03</b>

Model	MOS ( $\uparrow$ )	UTMOS ( $\uparrow$ )	WV-MOS ( $\uparrow$ )	DNSMOS ( $\uparrow$ )	PESQ ( $\uparrow$ )	STOI ( $\uparrow$ )	SI-SDR ( $\uparrow$ )	WER ( $\downarrow$ )
Input	$3.18 \pm 0.07$	$3.06 \pm 0.14$	$2.99 \pm 0.24$	$2.53 \pm 0.10$	$1.98 \pm 0.17$	$0.92 \pm 0.01$	$8.4 \pm 1.2$	$0.09 \pm 0.03$
MetricGAN+	$3.75 \pm 0.06$	$3.62 \pm 0.09$	$3.89 \pm 0.10$	$2.95 \pm 0.05$	$3.14 \pm 0.10$	$0.93 \pm 0.01$	$8.6 \pm 0.7$	$0.10 \pm 0.04$
DEMUCS	$3.95 \pm 0.06$	$3.95 \pm 0.05$	$4.37 \pm 0.06$	$3.14 \pm 0.04$	$3.04 \pm 0.12$	<b><math>0.95 \pm 0.01</math></b>	$18.5 \pm 0.6$	<b><math>0.07 \pm 0.03</math></b>
HiFi++	$4.08 \pm 0.05$	$3.89 \pm 0.06$	$4.36 \pm 0.06$	$3.10 \pm 0.04$	$2.90 \pm 0.12$	<b><math>0.95 \pm 0.01</math></b>	$17.9 \pm 0.6$	$0.08 \pm 0.03$
HiFi-GAN-2	$4.13 \pm 0.05$	$3.99 \pm 0.05$	$4.26 \pm 0.05$	$3.12 \pm 0.05$	$3.14 \pm 0.12$	<b><math>0.95 \pm 0.01</math></b>	$18.6 \pm 0.6$	<b><math>0.07 \pm 0.03</math></b>
DB-AIAT	$4.22 \pm 0.05$	$4.02 \pm 0.05$	$4.38 \pm 0.06$	$3.18 \pm 0.04$	<b><math>3.26 \pm 0.12</math></b>	<b><math>0.96 \pm 0.01</math></b>	<b><math>19.3 \pm 0.8</math></b>	<b><math>0.07 \pm 0.03</math></b>
Ours (16 kHz)	<b><math>4.41 \pm 0.04</math></b>	<b><math>4.32 \pm 0.02</math></b>	<b><math>4.87 \pm 0.05</math></b>	<b><math>3.22 \pm 0.04</math></b>	$2.94 \pm 0.10$	$0.92 \pm 0.01$	$4.6 \pm 0.3$	<b><math>0.07 \pm 0.03</math></b>
Ours (48 kHz)	<b><math>4.66 \pm 0.04</math></b>	<b><math>4.32 \pm 0.02</math></b>	<b><math>4.87 \pm 0.05</math></b>	<b><math>3.22 \pm 0.04</math></b>	$2.94 \pm 0.10$	$0.92 \pm 0.01$	$4.6 \pm 0.3$	<b><math>0.07 \pm 0.03</math></b>
GT (16 kHz)	$4.26 \pm 0.05$	$4.07 \pm 0.04$	$4.52 \pm 0.04$	$3.16 \pm 0.04$	-	-	-	-
GT (48 kHz)	$4.56 \pm 0.03$	$4.07 \pm 0.04$	$4.52 \pm 0.04$	$3.16 \pm 0.04$	-	-	-	-

For more information, please, visit our [demo](#).