

# Diffusion-based speech enhancement in matched and mismatched conditions using a Heun-based sampler

P. Gonzalez<sup>1</sup> Z.-H. Tan<sup>2</sup> J. Østergaard<sup>2</sup> J. Jensen<sup>2</sup> T. S. Alstrøm<sup>3</sup> T. May<sup>1</sup>

<sup>1</sup>Department of Health Technology  
Technical University of Denmark

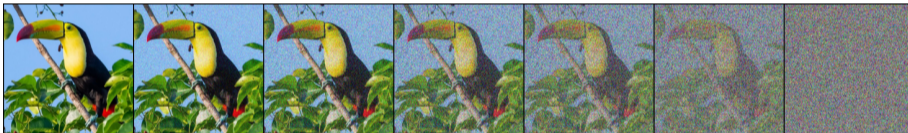
<sup>2</sup>Department of Electronic Systems  
Aalborg University

<sup>3</sup>Department of Applied Mathematics and Computer Science  
Technical University of Denmark

April 16th, 2024

# Introduction – Diffusion models

- Diffusion models are a new class of generative models that have recently been applied to speech enhancement<sup>1,2</sup>
- They iteratively add noise to the training data and learn to undo this process



<sup>1</sup>Y.-J. Lu *et al.*, "Conditional diffusion probabilistic model for speech enhancement," *ICASSP*, 2022

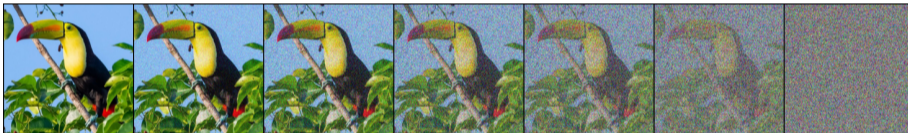
<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Diffusion models

- Diffusion models are a new class of generative models that have recently been applied to speech enhancement<sup>1,2</sup>
- They iteratively add noise to the training data and learn to undo this process

Forward process

$$\longrightarrow dx_t = f(t, x_t) dt + g(t) d\omega_t \longrightarrow$$



<sup>1</sup>Y.-J. Lu *et al.*, "Conditional diffusion probabilistic model for speech enhancement," *ICASSP*, 2022

<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Diffusion models

- Diffusion models are a new class of generative models that have recently been applied to speech enhancement<sup>1,2</sup>
- They iteratively add noise to the training data and learn to undo this process

Forward process

$$\longrightarrow dx_t = f(t, x_t) dt + g(t) d\omega_t \longrightarrow$$



- **Noise schedule:** defined by the hyperparameters  $f(t, x_t)$  and  $g(t)$

<sup>1</sup>Y.-J. Lu *et al.*, "Conditional diffusion probabilistic model for speech enhancement," *ICASSP, 2022*

<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Diffusion models

- Diffusion models are a new class of generative models that have recently been applied to speech enhancement<sup>1,2</sup>
- They iteratively add noise to the training data and learn to undo this process

Forward process

$$\longrightarrow dx_t = f(t, x_t) dt + g(t) d\omega_t \longrightarrow$$



$$\longleftarrow dx_t = [f(t, x_t) - g(t)^2 \nabla_{x_t} \log p_t(x_t)] dt + g(t) d\omega_t \longleftarrow$$

Reverse process

- **Noise schedule:** defined by the hyperparameters  $f(t, x_t)$  and  $g(t)$

<sup>1</sup>Y.-J. Lu *et al.*, "Conditional diffusion probabilistic model for speech enhancement," *ICASSP, 2022*

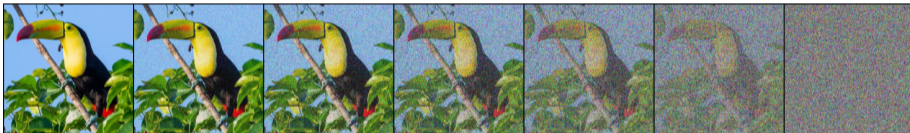
<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Diffusion models

- Diffusion models are a new class of generative models that have recently been applied to speech enhancement<sup>1,2</sup>
- They iteratively add noise to the training data and learn to undo this process

Forward process

$$\longrightarrow dx_t = f(t, x_t) dt + g(t) d\omega_t \longrightarrow$$



$$\longleftarrow dx_t = [f(t, x_t) - g(t)^2 \nabla_{x_t} \log p_t(x_t)] dt + g(t) d\omega_t \longleftarrow$$

Reverse process

- **Noise schedule:** defined by the hyperparameters  $f(t, x_t)$  and  $g(t)$
- **Sampler:** numerical method used to integrate the reverse process

<sup>1</sup>Y.-J. Lu *et al.*, "Conditional diffusion probabilistic model for speech enhancement," *ICASSP, 2022*

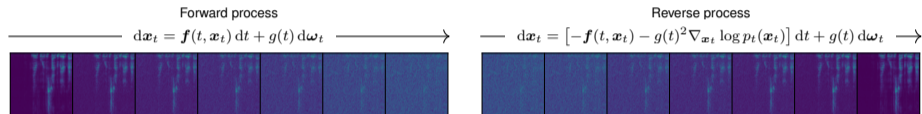
<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Adaptation to speech enhancement

- Output image ↔ Clean spectrogram  
Text prompt ↔ Noisy spectrogram

# Introduction – Adaptation to speech enhancement

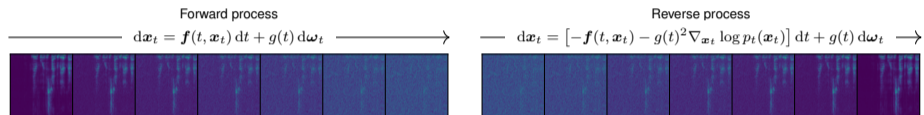
- Output image  $\longleftrightarrow$  Clean spectrogram  
 Text prompt  $\longleftrightarrow$  Noisy spectrogram





# Introduction – Adaptation to speech enhancement

- Output image  $\longleftrightarrow$  Clean spectrogram  
Text prompt  $\longleftrightarrow$  Noisy spectrogram
- Provoke a drift towards the noisy spectrogram<sup>12</sup>

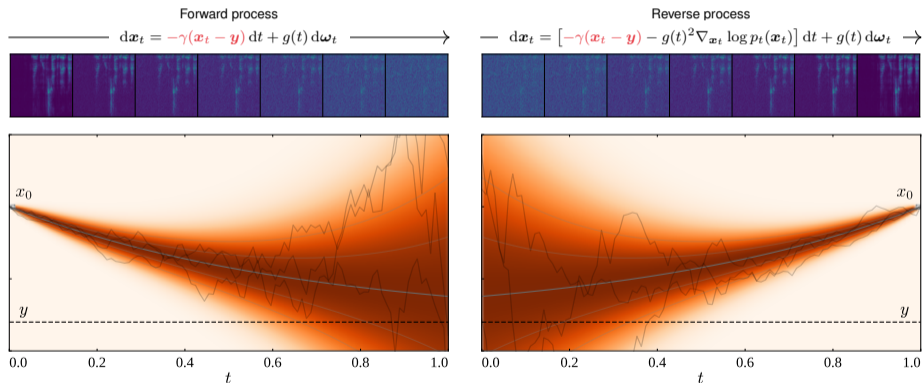


<sup>1</sup>S. Welker *et al.*, "Speech enhancement with score-based generative models in the complex STFT domain," *INTERSPEECH*, 2022

<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Introduction – Adaptation to speech enhancement

- Output image  $\longleftrightarrow$  Clean spectrogram  
Text prompt  $\longleftrightarrow$  Noisy spectrogram
- Provoke a drift towards the noisy spectrogram<sup>12</sup>



<sup>1</sup>S. Welker *et al.*, "Speech enhancement with score-based generative models in the complex STFT domain," *INTERSPEECH*, 2022

<sup>2</sup>J. Richter *et al.*, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, 2023

# Motivation

- **Goal:** apply findings from image generation literature<sup>1</sup> to improve speech enhancement performance
  - Neural network preconditioning based on first principles
  - Second-order Heun-based sampler

<sup>1</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS, 2022*

# Motivation

- **Goal:** apply findings from image generation literature<sup>1</sup> to improve speech enhancement performance
  - Neural network preconditioning based on first principles
  - Second-order Heun-based sampler
- **Problem:** the drift coefficient  $-\gamma(\mathbf{x}_t - \mathbf{y})$  cannot be written in the form  $f(t)\mathbf{x}_t$

<sup>1</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS, 2022*

# Motivation

- **Goal:** apply findings from image generation literature<sup>1</sup> to improve speech enhancement performance
  - Neural network preconditioning based on first principles
  - Second-order Heun-based sampler
- **Problem:** the drift coefficient  $-\gamma(\mathbf{x}_t - \mathbf{y})$  cannot be written in the form  $f(t)\mathbf{x}_t$
- **Solution:** change of variable  $\mathbf{n}_t = \mathbf{x}_t - \mathbf{y}$

$$d\mathbf{x}_t = -\gamma(\mathbf{x}_t - \mathbf{y}) dt + g(t) d\boldsymbol{\omega}_t$$

$$\Downarrow$$

$$d\mathbf{n}_t = -\gamma\mathbf{n}_t dt + g(t) d\boldsymbol{\omega}_t$$

Or more generally,

$$d\mathbf{n}_t = f(t)\mathbf{n}_t dt + g(t) d\boldsymbol{\omega}_t$$

<sup>1</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

# Motivation

- **Goal:** apply findings from image generation literature<sup>1</sup> to improve speech enhancement performance
  - Neural network preconditioning based on first principles
  - Second-order Heun-based sampler
- **Problem:** the drift coefficient  $-\gamma(\mathbf{x}_t - \mathbf{y})$  cannot be written in the form  $f(t)\mathbf{x}_t$
- **Solution:** change of variable  $\mathbf{n}_t = \mathbf{x}_t - \mathbf{y}$

$$d\mathbf{x}_t = -\gamma(\mathbf{x}_t - \mathbf{y}) dt + g(t) d\omega_t$$

$$\Downarrow$$

$$d\mathbf{n}_t = -\gamma\mathbf{n}_t dt + g(t) d\omega_t$$

Or more generally,

$$d\mathbf{n}_t = f(t)\mathbf{n}_t dt + g(t) d\omega_t$$

- This allows to write

$$p(\mathbf{n}_t | \mathbf{n}_0) = \mathcal{N}(\mathbf{n}_t; s(t)\mathbf{n}_0, s(t)^2\sigma(t)^2\mathbf{I})$$

where

$$s(t) = \exp \int_0^t f(\xi) d\xi \quad \text{and} \quad \sigma(t)^2 = \int_0^t \frac{g(\xi)^2}{s(\xi)^2} d\xi$$

<sup>1</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS, 2022*

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP, 2023*

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML, 2023*

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS, 2022*

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR, 2021*

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021



# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP, 2023*

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML, 2023*

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS, 2022*

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR, 2021*

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**  
**Noise:**  
**Room:**

TIMIT  
TAU  
Surrey

LibriSpeech  
NOISEX  
ASH

WSJ  
ICRA  
BRAS

Clarity  
DEMAND  
CATT

VCTK  
ARTE  
AVIL

Training  
Testing

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

## Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK	Training Testing
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE	
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL	

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK	Training Testing
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE	
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL	

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK	Training Testing
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE	
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL	

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**  
**Noise:**  
**Room:**

TIMIT	LibriSpeech	WSJ	Clarity
TAU	NOISEX	ICRA	DEMAND
Surrey	ASH	BRAS	CATT

VCTK
ARTE
AVIL

Training  
Testing

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK	Training Testing
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE	
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL	

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

## Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**  
**Noise:**  
**Room:**

TIMIT	LibriSpeech	WSJ	Clarity
TAU	NOISEX	ICRA	DEMAND
Surrey	ASH	BRAS	CATT

VCTK
ARTE
AVIL

Training  
Testing

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021



## Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**  
**Noise:**  
**Room:**

TIMIT	LibriSpeech	WSJ
TAU	NOISEX	ICRA
Surrey	ASH	BRAS

Clarity
DEMAND
CATT

VCTK
ARTE
AVIL

Training  
Testing

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**

TIMIT

LibriSpeech

WSJ

Clarity

VCTK

**Noise:**

TAU

NOISEX

ICRA

DEMAND

ARTE

**Room:**

Surrey

ASH

BRAS

CATT

AVIL

Training  
Testing

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

## Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

<b>Speech:</b>	TIMIT	LibriSpeech	WSJ	Clarity	VCTK	Training Testing
<b>Noise:</b>	TAU	NOISEX	ICRA	DEMAND	ARTE	
<b>Room:</b>	Surrey	ASH	BRAS	CATT	AVIL	

<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

# Experimental setup

- Modifications to SGMSE+M<sup>1</sup>:
  - Shifted-cosine noise schedule<sup>2</sup>
  - Neural network preconditioning<sup>3</sup>
  - Heun-based sampler<sup>3</sup> instead of PC sampler<sup>4</sup>
- Evaluation in matched and mismatched conditions:
  - Noisy speech is generated using 5 speech corpora, 5 noise databases and 5 room impulse response databases
  - $N = 1$  or  $N = 4$  databases are used for the training condition and the remaining databases are used for the mismatched condition
  - Results are averaged across 5 different combinations

**Speech:**

TIMIT

LibriSpeech

WSJ

Clarity

VCTK

**Noise:**

TAU

NOISEX

ICRA

DEMAND

ARTE

**Room:**

Surrey

ASH

BRAS

CATT

AVIL

Training  
Testing

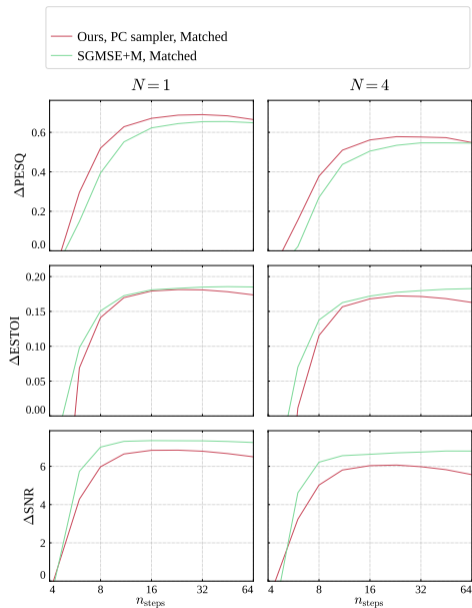
<sup>1</sup>J.-M. Lemerrier *et al.*, "Analysing diffusion-based generative approaches versus discriminative approaches for speech restoration," *ICASSP*, 2023

<sup>2</sup>E. Hoogeboom *et al.*, "simple diffusion: End-to-end diffusion for high resolution images," *ICML*, 2023

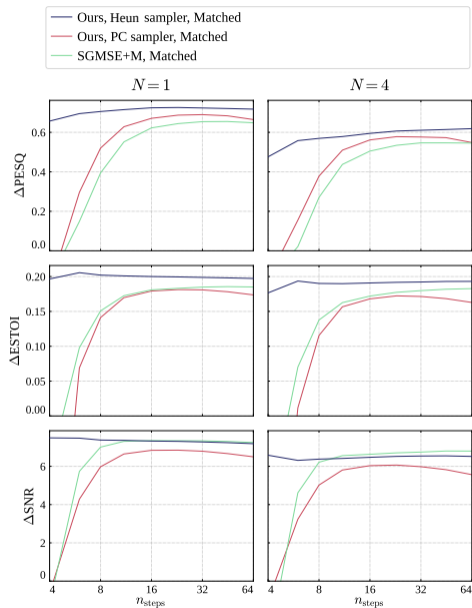
<sup>3</sup>T. Karras *et al.*, "Elucidating the design space of diffusion-based generative models," *NeurIPS*, 2022

<sup>4</sup>Y. Song *et al.*, "Score-based generative modeling through stochastic differential equations," *ICLR*, 2021

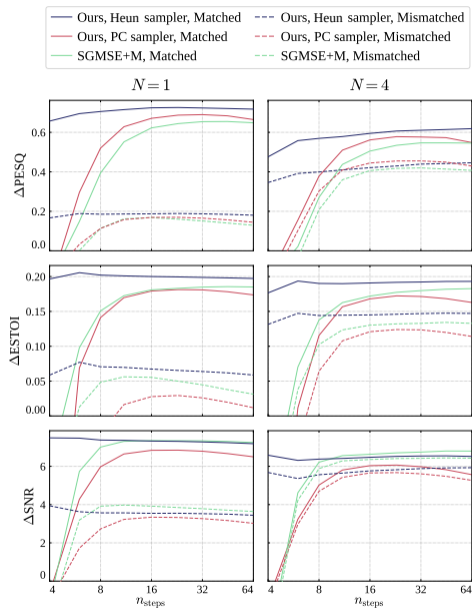
# Results



# Results



# Results



# Results

		$\Delta$ PESQ	$\Delta$ ESTOI	$\Delta$ SNR
Matched	Conv-TasNet	0.63	0.19	<b>8.58</b>
	DCCRN	0.41	0.12	7.11
	MANNER	0.68	0.17	7.24
	SGMSE+M	0.65	0.18	7.33
	Ours	<b>0.72</b>	<b>0.20</b>	7.27
Mismatched	Conv-TasNet	0.12	0.01	3.24
	DCCRN	0.11	0.02	3.25
	MANNER	0.17	0.04	3.40
	SGMSE+M	0.15	0.04	<b>3.79</b>
	Ours	<b>0.19</b>	<b>0.06</b>	3.52

(a)  $N = 1$

		$\Delta$ PESQ	$\Delta$ ESTOI	$\Delta$ SNR
Matched	Conv-TasNet	0.44	0.16	<b>7.38</b>
	DCCRN	0.34	0.12	6.52
	MANNER	0.52	0.15	6.23
	SGMSE+M	0.55	0.18	6.74
	Ours	<b>0.61</b>	<b>0.19</b>	6.54
Mismatched	Conv-TasNet	0.29	0.09	5.77
	DCCRN	0.23	0.07	5.42
	MANNER	0.38	0.10	5.61
	SGMSE+M	0.42	0.13	<b>6.42</b>
	Ours	<b>0.44</b>	<b>0.15</b>	5.88

(b)  $N = 4$

**Table:** Average  $\Delta$ PESQ,  $\Delta$ ESTOI and  $\Delta$ SNR scores in matched and mismatched conditions when training with  $N = 1$  (a) or  $N = 4$  (b) speech corpora, noise databases and BRIR databases. SGMSE+M and Ours use  $n_{steps} = 32$ .



# Conclusion

- The drift towards the noisy speech previously proposed for diffusion-based speech enhancement makes it difficult to apply recent advances from image generation literature
- To overcome this, the diffusion process is reformulated using a change of variable
- A different neural network preconditioning and noise schedule only had a small effect on performance
- The Heun-based sampler substantially improved the performance at few sampling steps
- All systems substantially benefited from training with multiple corpora in mismatched conditions

Thank you!