

Generative Diffusion Models for Audio Inpainting



SAPIENZA
UNIVERSITÀ DI ROMA

Facoltà di Ingegneria dell'informazione, informatica e statistica
Engineering in Computer Science

Andrea Rodriguez
1834937

Advisor
Prof. Danilo Comminiello

Generative Diffusion Models for Audio Inpainting

1. Background

1. Diffusion models
2. Spectrograms
3. Audio generation and inpainting

2. State-of-the-Art for audio generation

1. AudioLDM
2. TANGO

3. Selected approaches for inpainting

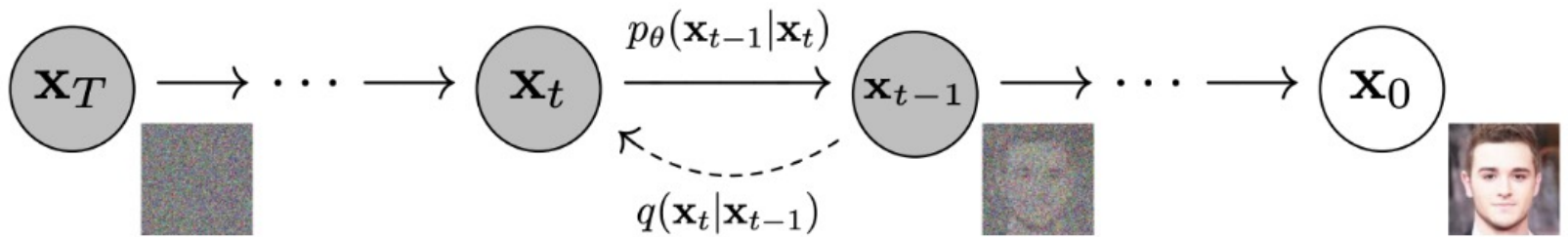
1. DDNM and DDNM⁺
2. RePaint and RePaint⁺

4. Results and Evaluation

5. Audio inpainting in communication scenarios

Diffusion models

- Main concepts**
- Iteratively transform an initial noise distribution into a target distribution
 - Generate high-quality samples and perform denoising and inpainting



Forward process: using a variance schedule, small amount of Gaussian noise is added to the sample for T steps

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

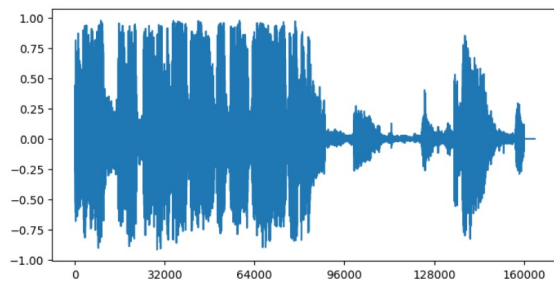
Reverse process: the noise added at each step of the forward process is predicted and removed from the sample

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$$

Spectrograms

A **spectrogram** is a visual representation that displays the frequency content of an audio signal over time

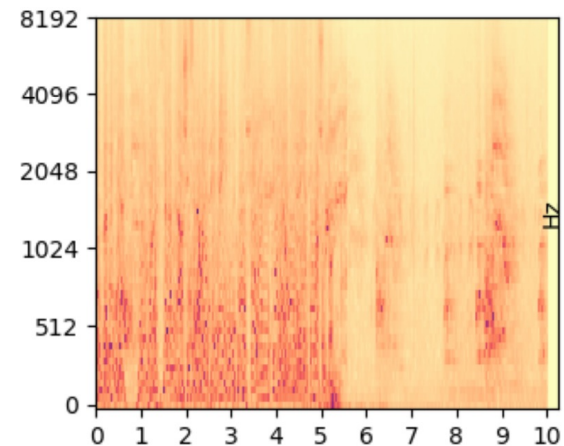
Vector 1 x 160000 (10 seconds)



(a) Audio wave

Sparsity → Computationally demanding and Risk of overfitting

Matrix 1024 x 64 (10 seconds)



(b) Spectrogram

Dense representation → Capture long term dependencies and Efficient training

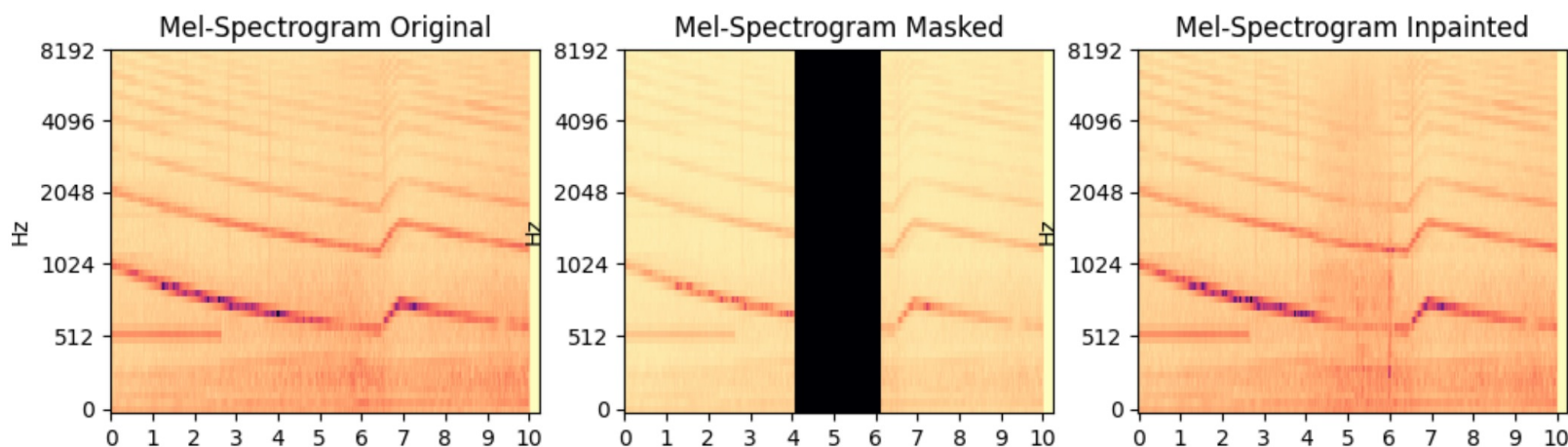
Audio generation and Audio inpainting

Audio Generation

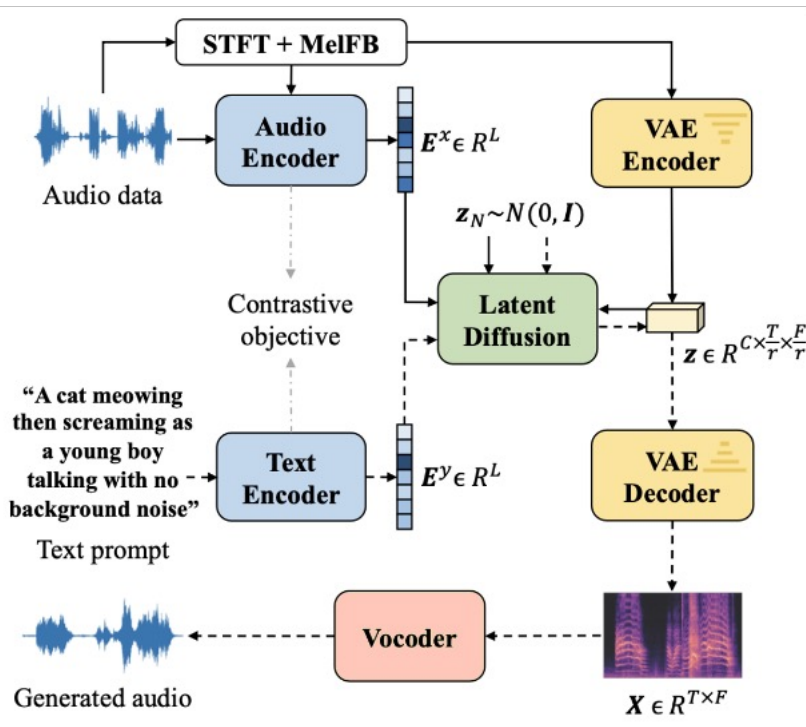
- Produce audio from textual descriptions that is indistinguishable from human-created or real-world audio
- Generate audio samples with similar characteristics to the training data, also showcasing innovative attributes

Audio Inpainting

- Reconstruct missing or corrupted portions of audio signals and restore the original audio content



State-of-the-Art: AudioLDM



CLAP (Contrastive Language-Audio Pretraining):

- encode audio descriptions and audio clips into a shared audio-text embedding space

- **VAE** (Variational Autoencoder): compress the spectrogram into a compact latent space

- **Latent Diffusion**: the conditioning information is integrated into the feature extraction process

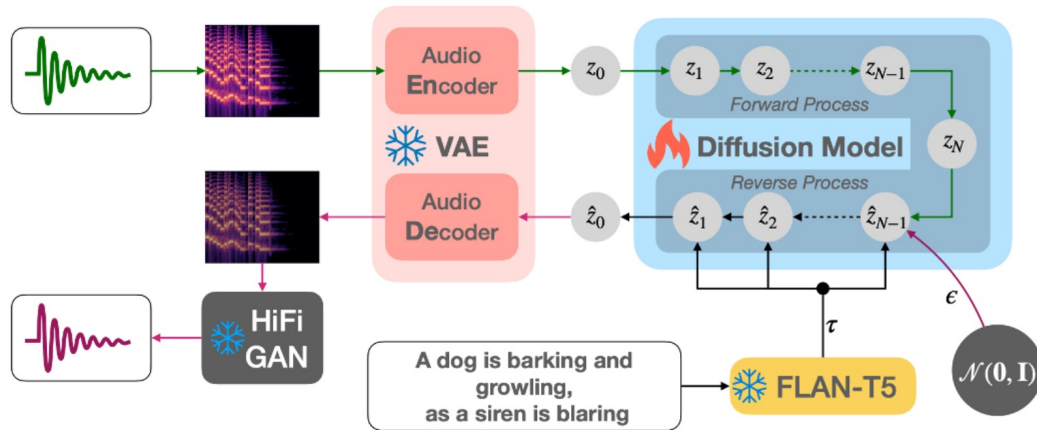
Vocoder: HiFi-GAN is employed to synthesize

- the audio waveform from the generated spectrogram

Inpainting: at each step of the reverse process, the

latent representation z_{n-1} is modified using the mask m , $z'_{n-1} = (1 - m) \odot z_{n-1} + m \odot z_{n-1}^{ob}$
the generated z_{n-1} and the known part of the audio z^{ob}

State-of-the-Art: Tango



FLAN-T5:

- Instruction-tuned LLM architecture used as text encoder
- Trained on a large-scale chain-of-thought and instruction-based dataset
- Can learn new tasks mimicking gradient descent through attention weights

TANGO achieves comparable results compared to AudioLDM while extremely reducing the training dataset and time

Inpainting:

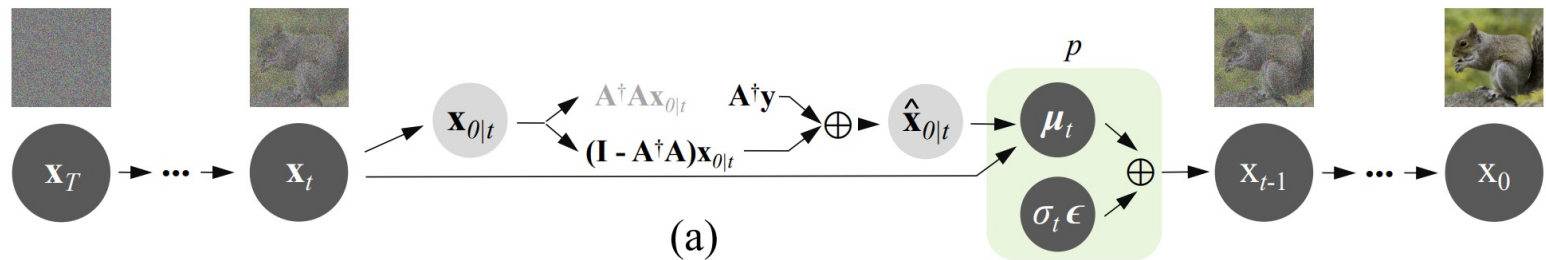
- x^{known} is sampled from the known part
- $x^{unknown}$ is sampled from the model

$$x_{t-1} = m \odot x_{t-1}^{known} + (1 - m) \odot x_{t-1}^{unknown}$$

The two components are combined using the mask m

Selected Approaches: DDNM and DDNM⁺

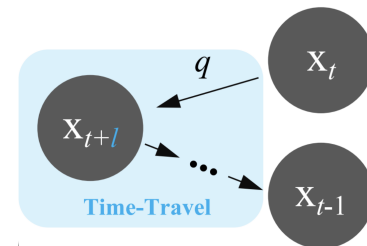
DDNM (Denoising Diffusion Null-Space Model) is a framework for image restoration which refines the null-space content during the reverse diffusion process to produce results that satisfy data consistency and realism



Task: reconstruct \hat{x} from y where $y = Ax$

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{x}_{0|t} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \mathcal{Z}_\theta(\mathbf{x}_t, t) \sqrt{1 - \bar{\alpha}_t})$
- 4: $\hat{\mathbf{x}}_{0|t} = \mathbf{A}^\dagger \mathbf{y} + (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_{0|t}$
- 5: $\mathbf{x}_{t-1} \sim p(\mathbf{x}_{t-1} | \mathbf{x}_t, \hat{\mathbf{x}}_{0|t})$
- 6: **return** \mathbf{x}_0

DDNM⁺: Through the time-travel trick we generate a better “past”, which in turn leads to a better “future”



Selected Approaches: RePaint and RePaint⁺

RePaint: sample the unmasked regions using the available information as conditioning during the diffusion process

$$x_{t-1}^{known} \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

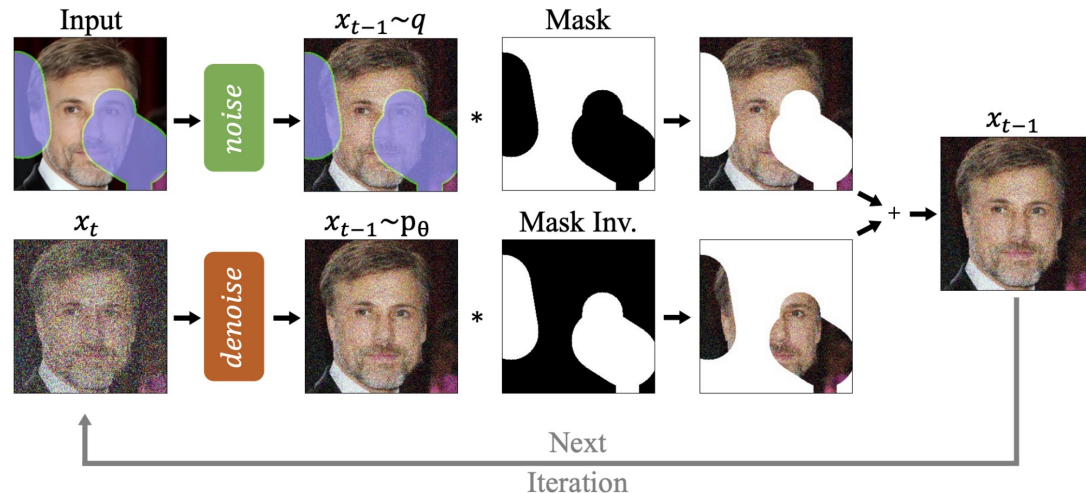
$$x_{t-1}^{unknown} \sim \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

$$x_{t-1} = m \odot x_{t-1}^{known} + (1 - m) \odot x_{t-1}^{unknown}$$

x_{t-1} is diffused back to x_t

$$x_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

the denoising step is performed again



The model has **more time** to effectively incorporate the provided information with the generated part

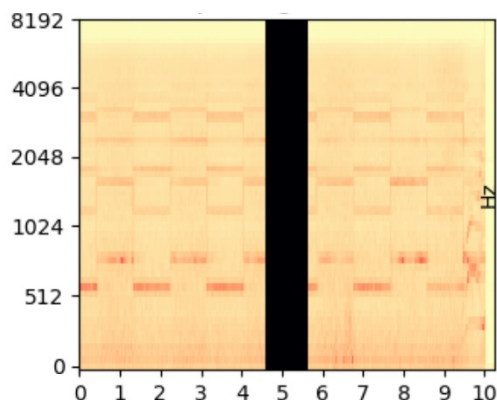
RePaint⁺: The latent representation is diffused back to multiple previous steps and then all of them are sequentially performed again

Audio samples and Inference details

Tests were done using 24 audio clips from the **AudioCaps** dataset



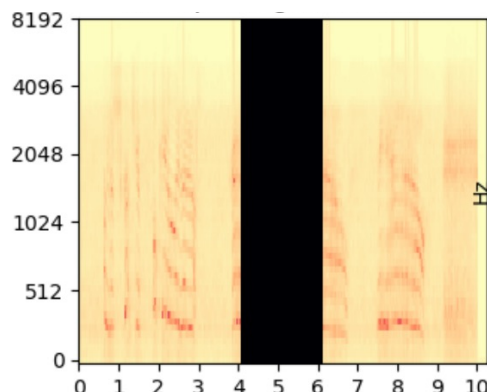
“Female and male are having conversation”



1 second gap



“An emergency vehicles’ siren with a brief male yell”



2 seconds gap



“Duck quacking repeatedly”

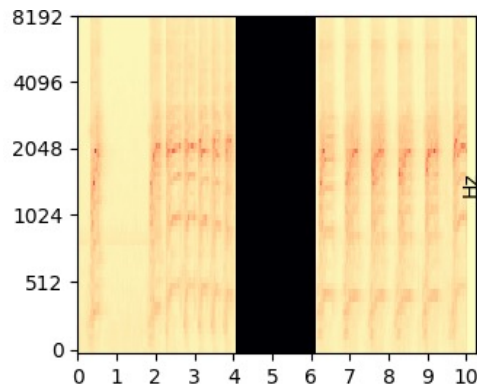
	1 Clip	Batch of 8 Clips
AudioLDM	1	4
Tango	10	40
Tango + DDNM	10	40
Tango + DDNM⁺	120	480
Tango + RePaint	100	400
Tango + RePaint⁺	100	400

Inference times in minutes using
one GPU NVIDIA Tesla T4

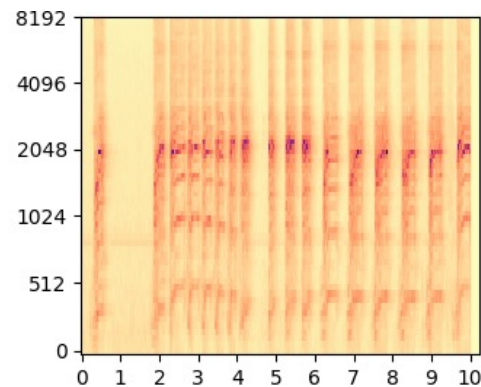
Listen to the audio clips at: <https://www.andrearodriguez.it/inpainting>

Results

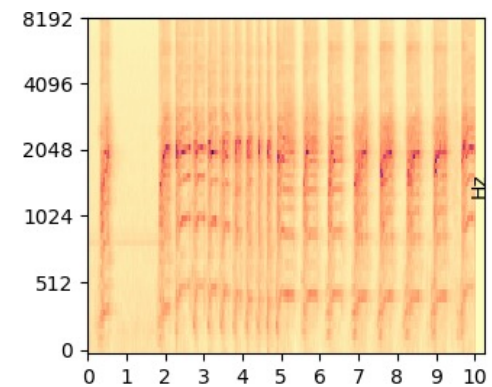
“Duck quacking repeatedly”



Masked

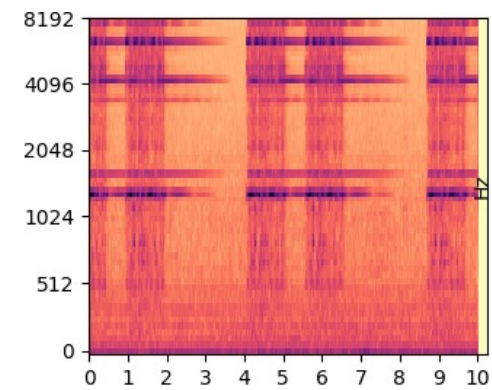
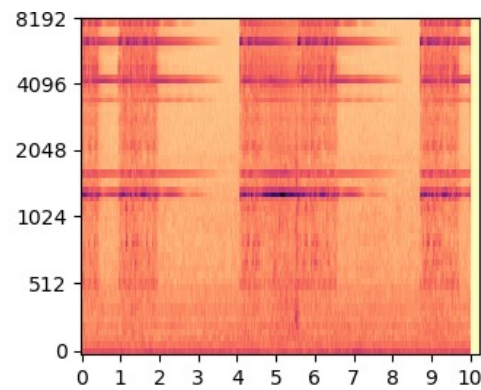
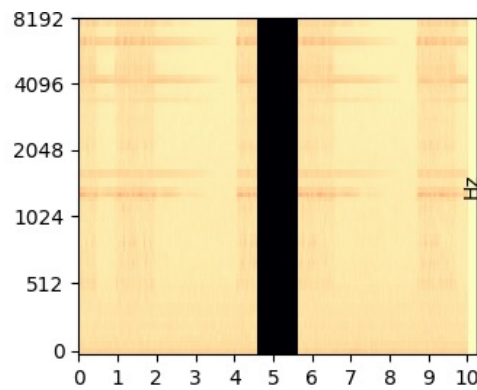


Inpainted using RePaint



Original

“A telephone ringing”



Metrics

- SDR = Signal Distortion Ratio
- SNR = Signal Noise Ratio
- PSNR = Peak Signal Noise Ratio
- SSIM = Structural Similarity Index Measure

1 second gap

4.5-5.5	AudioLDM	Tango	Tango DDNM	Tango DDNM ⁺	Tango RePaint	Tango RePaint ⁺
SDR	-3.27	5.48	5.03	5.47	5.96	4.97
SNR	-0.25	5.73	5.39	5.71	6.17	5.28
PSNR	39.46	43.35	42.44	43.22	44.08	42.38
SSIM	98.40	99.25	99.21	99.28	99.18	99.14

2 seconds gap

4-6	AudioLDM	Tango	Tango DDNM	Tango DDNM ⁺	Tango RePaint	Tango RePaint ⁺
SDR	-4.97	1.48	1.53	1.99	1.64	1.48
SNR	-1.45	2.82	2.90	2.21	2.71	2.02
PSNR	35.44	39.74	39.85	38.61	39.92	38.56
SSIM	97.84	98.34	98.46	98.45	98.59	98.56

* SSIM values are multiplied by 100

Audio inpainting in communication scenarios



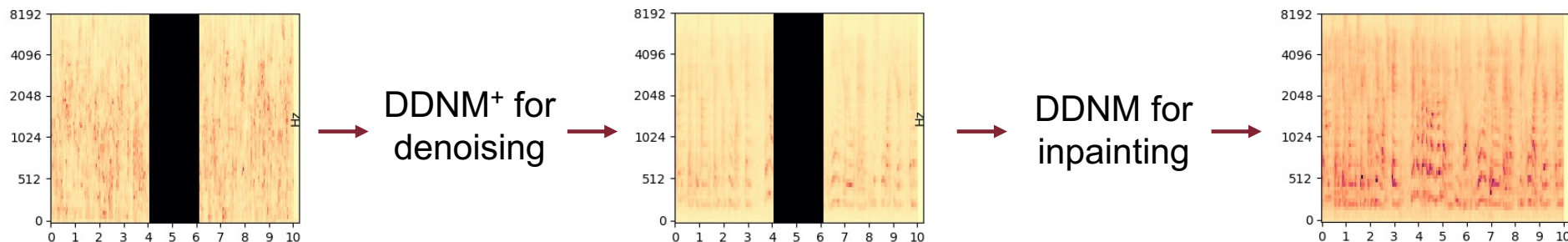
We tested the technique adding white **Gaussian noise** to some audio clips and then trying to remove it

We worked with **PSNR** 20 and 30

$$\hat{x}_{0|t} = \Sigma_t A^\dagger y + (I - \Sigma_t A^\dagger A) x_{0|t}$$

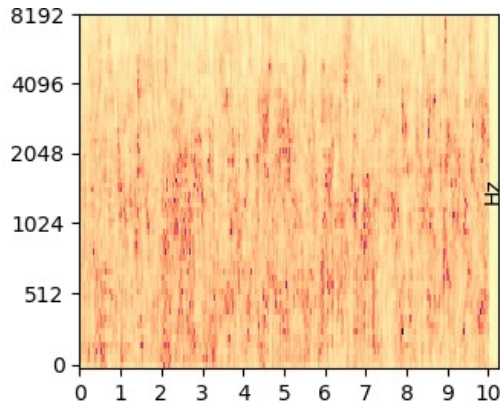
$$x_{t-1} = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \hat{x}_{0|t} + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \sqrt{\Phi_t} \epsilon$$

DDNM+: introduces a new parameter σ_y which represents the estimated noise level that we aim to eliminate from the audio signal

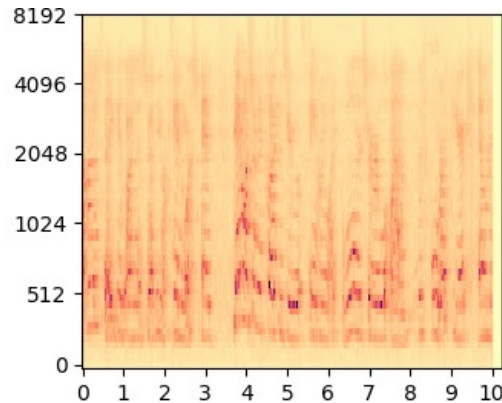


Results

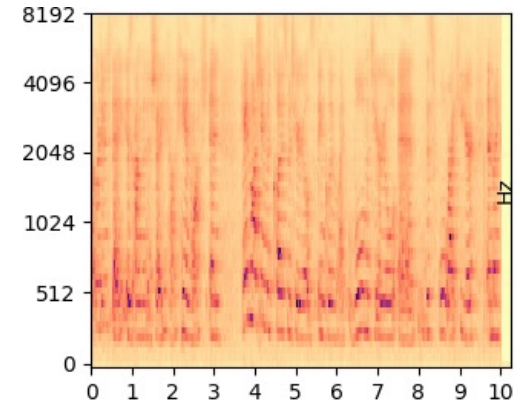
“An adult female is speaking in a quiet environment”



Noisy



Denoised



Original



SNR

PSNR 20	Clip 1	Clip 2	Clip 3	Clip 4
Noisy	-9.80	-9.11	-8.61	-10.10
Denoised	-2.47	-3.80	-2.64	-2.58

SNR

PSNR 30	Clip 1	Clip 2	Clip 3	Clip 4
Noisy	-3.53	-2.45	-3.54	-3.61
Denoised	-1.82	-1.14	-1.32	-2.20

Conclusions

- The proposed techniques achieve **better results** than the two baselines
- **RePaint** consistently achieves superior results compared to the other methodologies
- It is crucial to consider the **trade-off** between quality and inference time

and Future works

- Create an **automated communication system** which performs denoising and intelligently identifies and inpaints problematic segments
- **Remove conditioning** from text and perform inpainting based only on the known portion of audio

References

1. Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. 2020. arXiv: 2006.11239 [cs.LG].
2. Haohe Liu et al. AudioLDM: Text-to-Audio Generation with Latent Diffusion Models. 2023. arXiv: 2301.12503 [cs.SD].
3. Deepanway Ghosal et al. Text-to-Audio Generation using Instruction-Tuned LLM and Latent Diffusion Model. 2023. arXiv: 2304.13731 [eess.AS].
4. Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-Shot Image Restoration Using Denoising Diffusion Null-Space Model. 2022. arXiv: 2212.00490 [cs.CV].
5. Andreas Lugmayr et al. RePaint: Inpainting using Denoising Diffusion Probabilistic Models. 2022. arXiv: 2201.09865 [cs.CV].
6. Chris Dongjoo Kim et al. “AudioCaps: Generating Captions for Audios in The Wild”. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 119–132. doi: 10.18653/v1/N191011. url: <https://aclanthology.org/N19-1011>.