A DIFFUSION-INSPIRED TRAINING STRATEGY FOR SINGING **VOICE EXTRACTION IN THE WAVEFORM DOMAIN**

Genís Plaja-Roglans, Marius Miron and Xavier Serra

Universitat Pompeu Fabra, Music Technology Group, Barcelona, Spain



REVERSE PROCES

CONTEXT

SINGING VOICE EXTRACTION



DIFFUSION MODELS

DIFFUSION



tps://developer.nvidia.com/blog/improving-diffusion-models-as-an-alternative-to-gans-part



PROPOSED APPROACH



DIFFUSION PROCESS

• Given the parametrization in [1], we compute each diffusion step as: $\alpha_t = 1 - \beta_t$

$$q(x_t|x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} m \left\{ \bar{\alpha}_t = \prod_{s=1}^t \alpha_s \right\}$$

NOISE SCHEDULES

• We build our schedules following [1].

- We propose to use a deterministic signal to perturb the input signal for diffusion.
- By gradually converting a singing voice into its corresponding mixture, we train a small model to conduct the reverse operation.

REVERSE PROCESS

• We parametrize the reverse process as:

$$x_{t-1} = \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon(x_t, t)\right) \frac{1}{\sqrt{\alpha_t}}$$

this is an adaptation of the original formula in [1] that yields improved results.

- → Using 20 steps: $\beta_{20} = \{\beta_0 = 1^{-4}, \beta_T = 0.2\}$
- → Using 8 steps: $\beta_8 = \{\beta_0 = 1^{-4}, \beta_T = 0.5\}$

NETWORK

- We train a non-autoregressive WaveNet [2] to learn the reverse process.
- This network has been previously used for music source separation [3].
- We rely on the DiffWave [4] implementation.

WIENER FILTERING



TRAINING OBJECTIVE

$$\| \epsilon - \epsilon_{\theta} (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} m, t)) \|$$

We use the **accompaniment** as target!

EVALUATION

| | | | Sir | Singing Voice | | Accompaniment | | | 5 - | CI (9 | 5%) | | | 4.7 |
|--------------------------------|---|--------------|------|---------------|------|---------------|-------|-------|----------|---------|---------------------|---------------|------|-------|
| Model | Params | Diff. steps | SDR | SIR | SAR | SDR | SIR | SAR | | [1.5 | 7, 1.82] | | | |
| WaveNet | $\approx 3.3 M$ | - | 4.49 | 13.52 | 6.17 | 11.39 | 16.37 | 13.49 | ent) | [2.2 | 2, 2.50] | | | |
| Wave-U-Net | $\approx 10.2 \mathrm{M}$ | 2. 2 | 4.97 | 13.98 | 4.41 | 11.11 | 15.30 | 11.44 | celle 4 | [2.5 | 5, 2.86] | | | |
| ConvTas-Net | $\approx 8.75 \mathrm{M}$ | 3 — 1 | 6.43 | - | - | - | - | - | -EX | | 5, 2.88] 7 4 84] | 0.71 | 2 72 | |
| Demucs (v1) | - | - | 5.44 | - | - | - | - | - | d/5 ۳ | [4.0 | 2 35 | 2.71 | 2.72 | |
| Demucs (v2) | $\approx 450 \mathrm{M}$ | | 6.84 | - | - | - | - | - | l-Ba | | 2.55 I | | | |
| Ours (vocal) | $\approx 750 \mathrm{K}$ | 1 | 4.81 | 9.21 | 8.09 | - | - | - | u 2 - | 1.71 | | | | |
| Ours (vocal) | $\approx 750 \mathrm{K}$ | 8 | 5.63 | 10.55 | 8.86 | - | - | - | uati | I | | | | |
| Ours (vocal) | $\approx 750 \mathrm{K}$ | 20 | 5.59 | 10.78 | 8.89 | - | - | - | | | | | | |
| Ours (vocal) + Wiener | $\approx 750 \mathrm{K}$ | 20 | 5.66 | 11.60 | 8.49 | - | - | - | | | | | | |
| Ours (accomp) | $\approx 750 \mathrm{K}$ | 100 | - | - | - | 11.12 | 13.11 | 16.44 | | | | | | |
| Ours (vocal & accomp) + Wiener | $\approx 750\mathrm{K} + 750\mathrm{K}$ | 20 + 100 | 6.07 | 12.77 | 8.61 | 11.72 | 14.44 | 16.81 | 0 - | WaveNet | Wave-U-Ne | et Demucs(v2) | Ours | Groun |



J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models". In Proc. of the 33th Advances in Neural Information Processing Systems (NeurIPS), Online, pp. 6840–6851, 2020.

D. Rethage, J. Pons, and X. Serra, "A wavenet for speech denoising". In Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP), Paris, France, vol. 2018-April, pp. 5069–5073, 2018. [2]

F. Lluís, J. Pons, and X. Serra, "End-to-end music source separation: Is it possible in the waveform domain?". In Proc. of the Int. Speech Communication Association (INTERSPEECH), Graz, Austria, pp. 4619–4623, 2019 [3]

Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, "DiffWave: A Versatile Diffusion Model for Audio Synthesis". In Proc. of the 9th Int. Conf. on Learning Representations (ICLR), Vienna, Austria, 2021. [4]

This work was carried out under the projects Musical AI - PID2019-111403GB-100/AEI/10.13039/501100011033 and NextCore - RTC2019-007248-7 funded by the Spanish Ministerio de Ciencia, Innovación y Universidades (MCIU) and the Agencia Estatal de Investigación (AEI). We would like to thank the 40 participants that took the perceptual test for this work.







