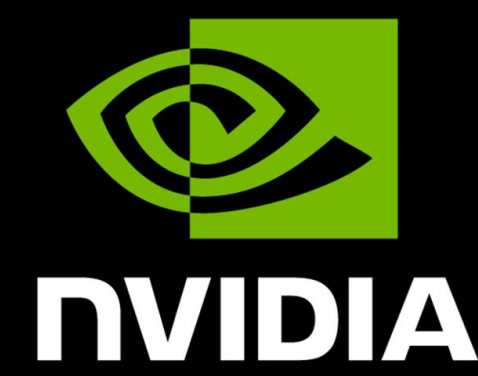


Score Distillation Sampling for Audio: Source Separation, Synthesis, and Beyond

Jessie Richter-Powell, Antonio Torralba, Jonathan Lorraine



ICML
International Conference
On Machine Learning



Introduction

- We adapt Score Distillation Sampling (SDS), originally for 3D generation [1], to audio models.

With one large, pretrained audio diffusion model...

...the Audio-SDS update improves performance...

...on diverse audio tasks.

source separation

physical impact synthesis

tuning FM synthesizers

... and more

Our Method

- High-level idea for the SDS update:

Finds parameters θ that render audio \mathbf{x} the diffusion model finds probable.

Rendered Audio:

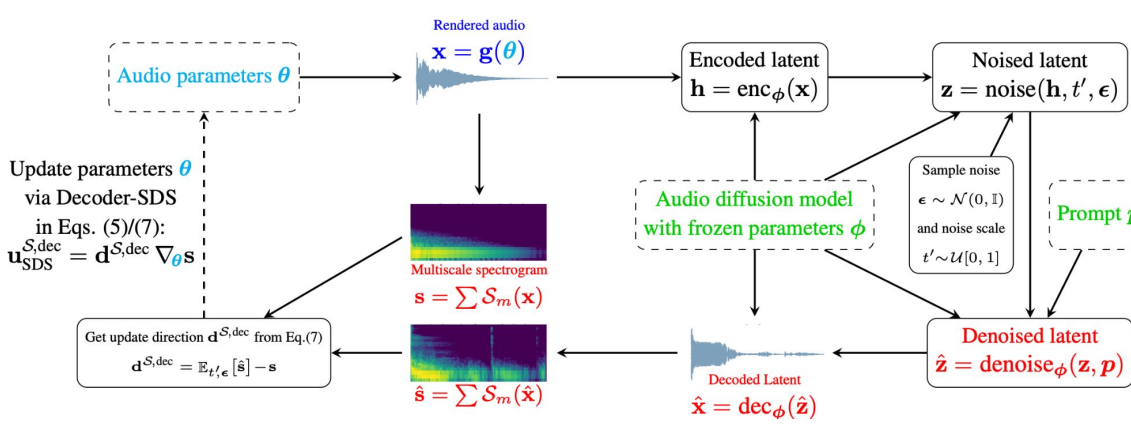
$$\mathbf{x} = \mathbf{g}(\theta, \mathbf{c})$$

SDS update (roughly):

$$\mathbf{u}_{\text{SDS}}^{\text{dec}} = (\mathbb{E}[\hat{\mathbf{x}}] - \mathbf{x}) \nabla_{\theta} \mathbf{x}$$

- All the gory details for SDS and our audio-variant:

$$\text{Full-SDS Update: } \mathbf{u}_{\text{SDS}}(\theta; \mathbf{p}) = \mathbb{E}_{t', \epsilon, \mathbf{c}} [\omega(t') (\hat{\epsilon}_{\phi}(\mathbf{z}(\theta, \mathbf{c}), t', \mathbf{p}) - \epsilon) \nabla_{\theta} \mathbf{z}(\theta, \mathbf{c})]$$



Our SDS-variant with spectrogram emphasis:

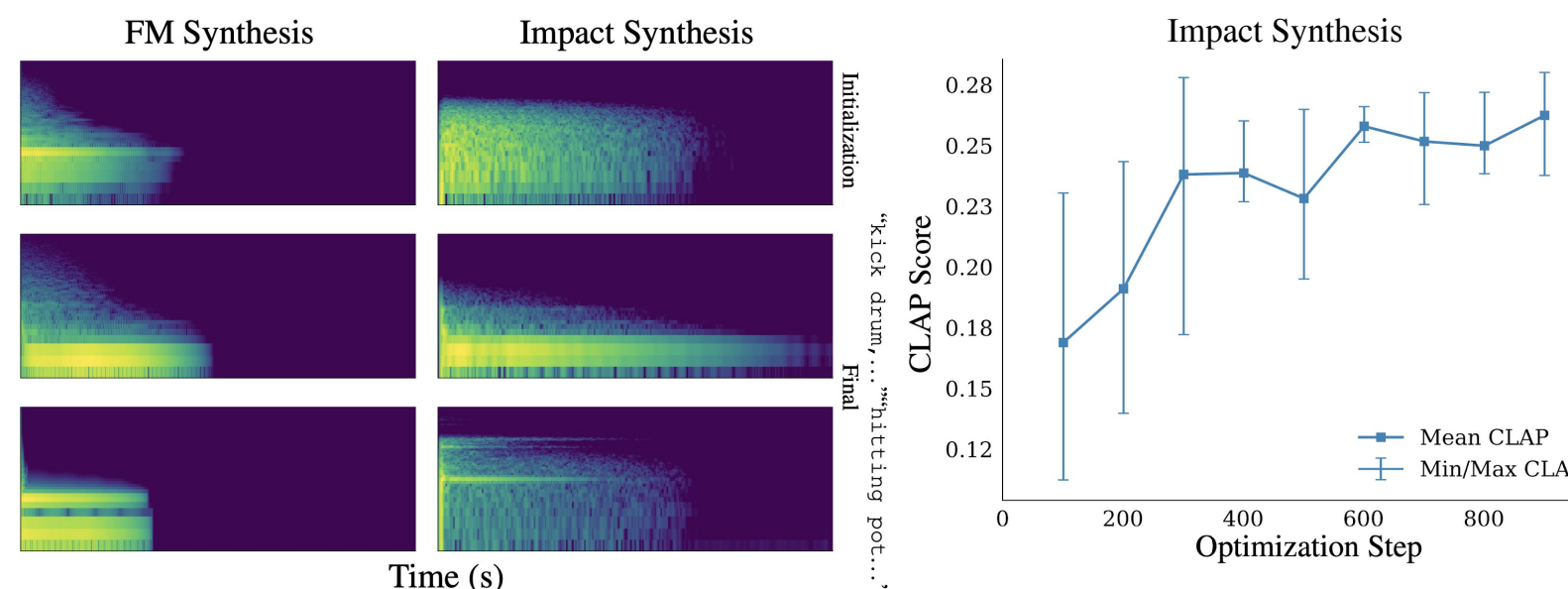
$$\mathbf{u}_{\text{SDS}}^{\text{S, dec}} = (\mathbb{E}[\hat{\mathbf{s}}] - \mathbf{s}) \nabla_{\theta} \mathbf{s}$$

Our Proposed Tasks

Task	FM Synthesis	Impact Synthesis	Source Separation
Use-Case	Toy setup, to generate FM synthesizer settings consistent with prompts like “kick drum, bass, reverb”	Generate impacts consistent with prompts like “hitting pot with wooden spoon”	Prompt-conditional source separation for a source audio, such as separating out a “sax ...” and “cars ...” from a music recording on a road
Optimizable Parameters θ	Envelope params., fundamental freq., FM Matrix $\theta = \{\alpha_v, \delta_v, \omega_v\}_{v=1}^V, \mathbf{A}$	Frequency, damping, amplitude of sinusoids $\theta = \{\lambda_n, d_n, a_n, \tilde{\lambda}_n, \tilde{d}_n, \tilde{a}_n\}_{n=1}^N$	Latent / raw audio for each source $\theta = \{\mathbf{x}_k\}_{k=1}^K$
Rendering Function $\mathbf{g}(\theta)$	Sine oscillators modulate each other's frequency $f_v(t) = \max(0, \min(t/(\alpha_v + 10^{-5}), \exp((\alpha_v - t)/\delta_v^2))(\delta_v - t - \alpha_v)/\delta_v)$ $\mathbf{u}_v[t] = \sin(t \cdot \omega_v + \langle \mathbf{A}_v, \mathbf{u}[t-1] \rangle) f_v(t)$ $\mathbf{g}(\theta)[t] = \langle \mathbf{A}_v, \mathbf{u}_{t-1} \rangle$	Convolution of impact, object and reverb impulse $\mathbf{I}_{\text{obj}}^{\theta}[t] = \sum_{n=1}^N \tilde{a}_n \exp(-\tilde{d}_n t) \cos(\tilde{\lambda}_n t)$ $\mathbf{I}_{\text{reverb}}^{\theta}[t] = \sum_{n=1}^N \tilde{a}_n \exp(-\tilde{d}_n t) F(\mathcal{W}(t), \tilde{\lambda}_n)$ $\mathbf{I}_{\text{impact}} = \text{Delta function}$ $\mathbf{g}(\theta) = \mathbf{I}_{\text{impact}} * \mathbf{I}_{\text{obj}}^{\theta} * \mathbf{I}_{\text{reverb}}^{\theta}$	Simply the sum of audio over sources $\mathbf{g}(\theta) = \sum_{k=1}^K \mathbf{x}_k$
Parameter Update	Audio-SDS	Audio-SDS	Reconstruction Loss Gradient + $\gamma \cdot \text{Audio-SDS}$
Visualization			

Experimental Results

- We assess quality with audio-prompt alignment via CLAP
- FM Synthesis works for simple prompts
- Impact synthesis improves CLAP on impact-oriented prompts
- Ablations provided in the paper



Experimental Results

- Prompt-guided Source Separation:** A user takes mixed audio and chooses a set of prompts for channels to split it.
- We improve SDR to ground truth and CLAP when ground truth unavailable.
- Fully-automatic pipeline:** We take a video, caption it with a model, and provide that to an LLM-assistant suggesting source decompositions.

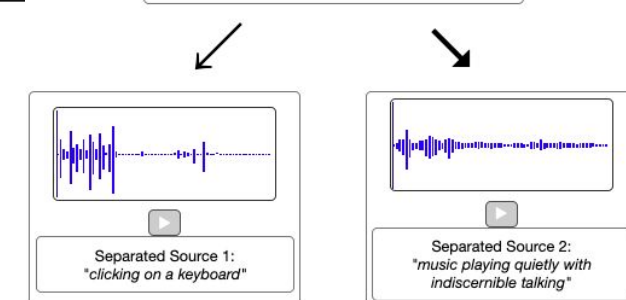


Give the target audio from video to an audio captioning model, such as AudioCaps.

Provide an LLM-assistant, such as ChatGPT, with the audio caption (here, “Someone is clicking on a keyboard and talking”) and a task description (e.g., “...suggest different prompts for sources given the audio caption...”).

LLM Output:
...
• Channel 1 Prompt: “clicking on a keyboard”
• Channel 2 Prompt: “music playing quietly with indiscernible talking”
...

Run Audio-SDS source separation using the prompts for each channel.



Mixture	SDRs for reconstructed (source ₁ , source ₂ , mixture)	
	Initialization	Us
Traffic + Sax	(−0.7, −5.2, 3.2)	(8.5, 8.1, 13.1)
Bongo + Waves	(1.2, −6.4, 3.9)	(1.2, −2.1, 8.7)
Pipes + Glass	(−4.2, 0.0, 2.8)	(1.5, 6.5, 7.7)
Clock + Bongo	(−15.8, 10.4, 11.8)	(−10.4, 3.6, 13.9)
Wind + Pipes	(−6.1, 1.9, 5.8)	(−0.5, 8.6, 7.3)
Mixture	CLAPs for reconstructed source ₁ and source ₂ , then SDR for mixture	
	Initialization	Us
Traffic + Sax	(0.17, 0.02, 3.2)	(0.2, 0.05, 13.1)
Bongo + Waves	(0.15, 0.14, 3.9)	(0.16, 0.08, 8.7)
Pipes + Glass	(0.25, 0.27, 2.8)	(0.21, 0.3, 7.7)
Clock + Bongo	(0.22, 0.24, 11.8)	(0.30, 0.34, 13.9)
Wind + Pipes	(0.25, 0.06, 5.8)	(0.25, 0.09, 7.3)

Our View of the Future

- Limitations to improve on:
 - Clip-Length Budget:** We optimized on ≤ 10 s clips; minute-scale audio may have artifacts or blow up memory.
 - Audio-Model Bias:** We use Stable Audio Open, so when this struggles, e.g., on rare instruments, speech, so do we.
- Future:** Easily use one video + audio diffusion model with SDS-style updates for various tasks — impacts, lighting, cloth, fluids, and more.

Links

- Webpage in QR-code: research.nvidia.com/labs/toronto-ai/Audio-SDS/
- [1] Poole, Ben, et al. "DreamFusion: Text-to-3D using 2D Diffusion." The Eleventh International Conference on Learning Representations.