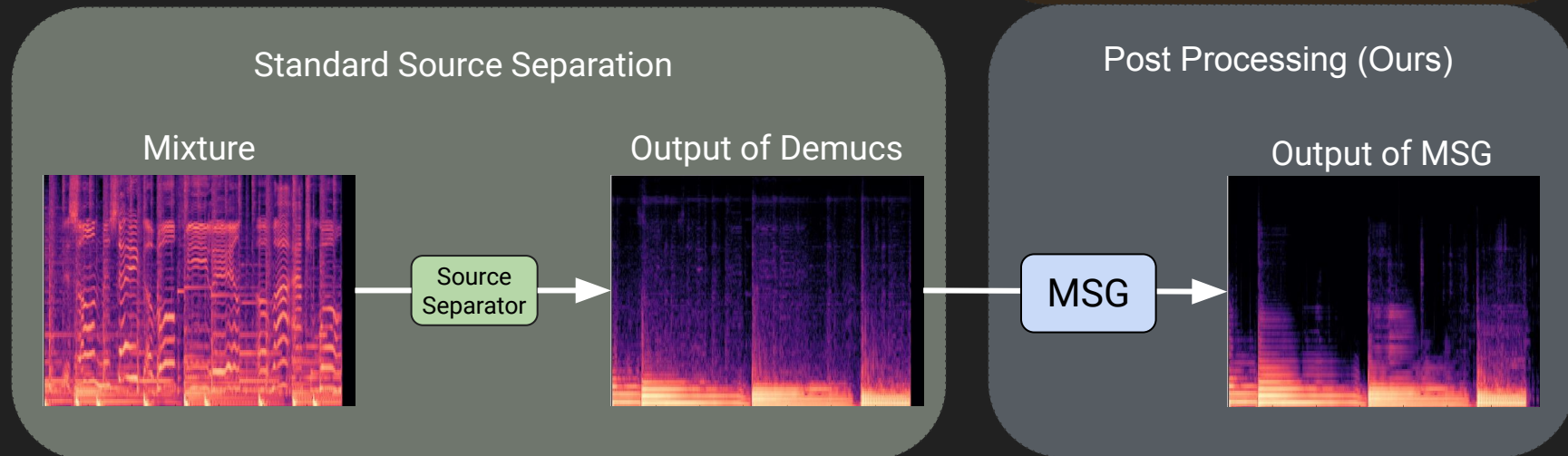# MSG: Make it Sound Good

Boaz Cogan, 🍜 Noah Schaffer, 🍜 Ethan Manilow, Bryan Pardo

🍜 = equal contribution
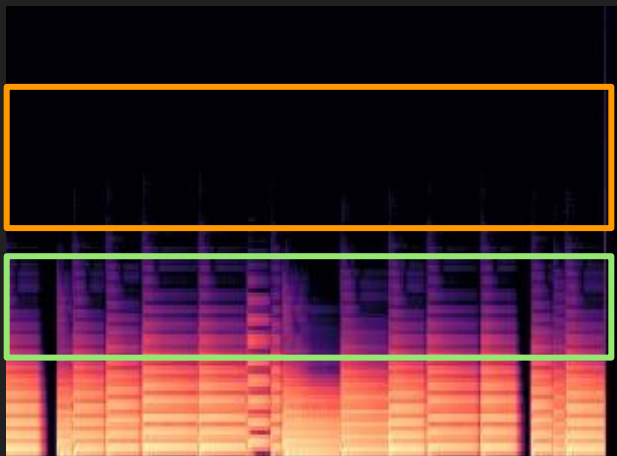
## Goal:

- Use a post-processor to <u>make source separation output sound better</u>
  - Reconstruct missing data
  - Reduce artifacts

Ground-Truth

Standard Source Separation

Mixture

Source Separator

Output of Demucs
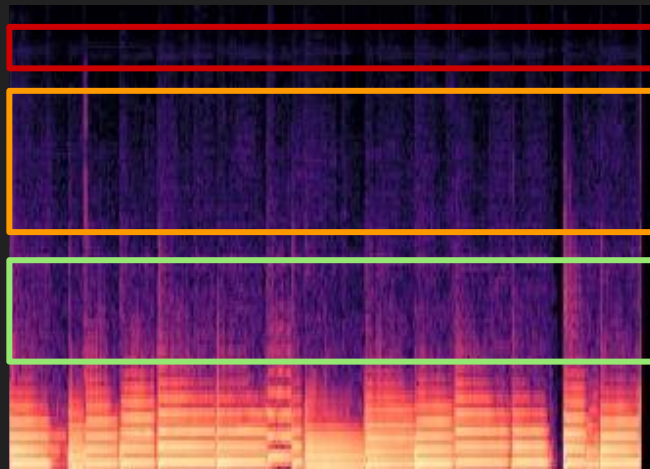
Post Processing (Ours)

MSG

Output of MSG

2

# SOTA Source Separation is Imperfect!

Ground Truth – Bass

Demucs – Bass Est.

- Source estimates that contain added noise and artifacts.
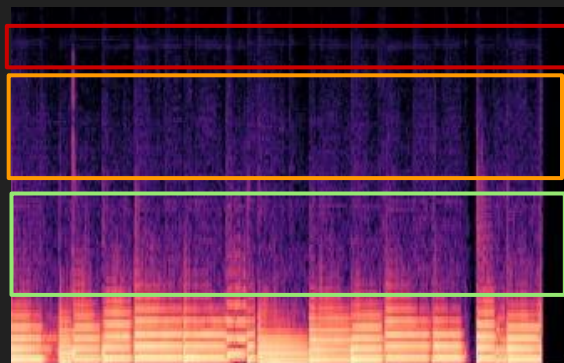
- How can we make this sound better?
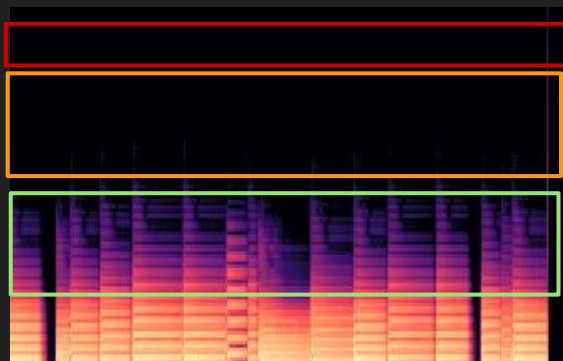
Missing Frequencies!

Added Noise!

Strange Artifacts!

# Clean up source separation to <u>M</u>ake it <u>S</u>ound <u>G</u>OOD!
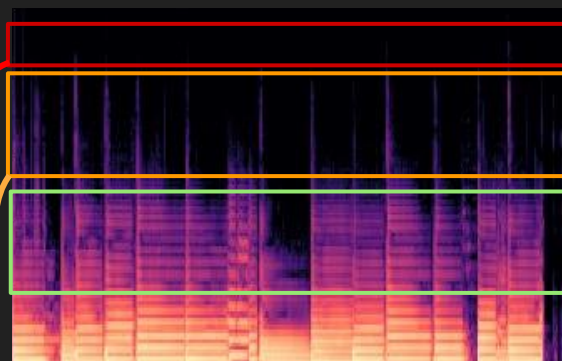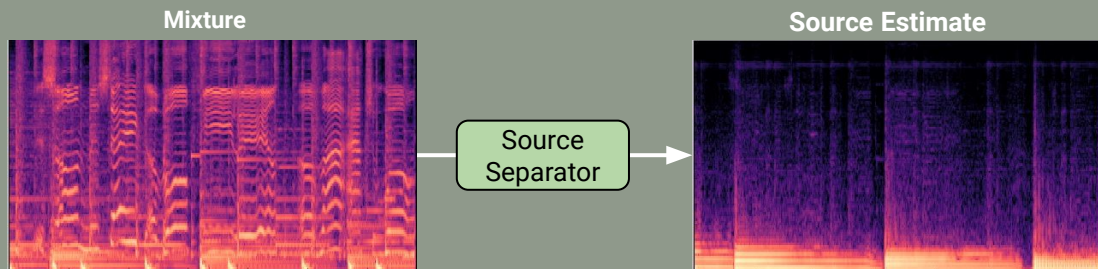
Demucs

Ground Truth

MSG (Ours!)

Strange artifact is gone!
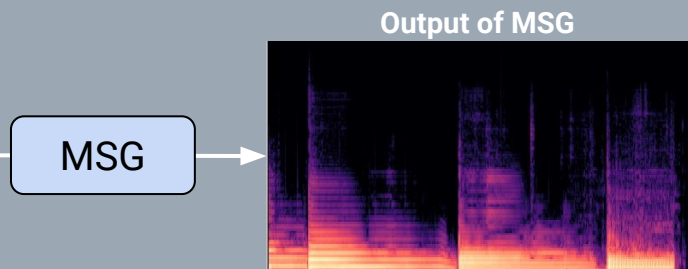
No noise!

Overtones are back!

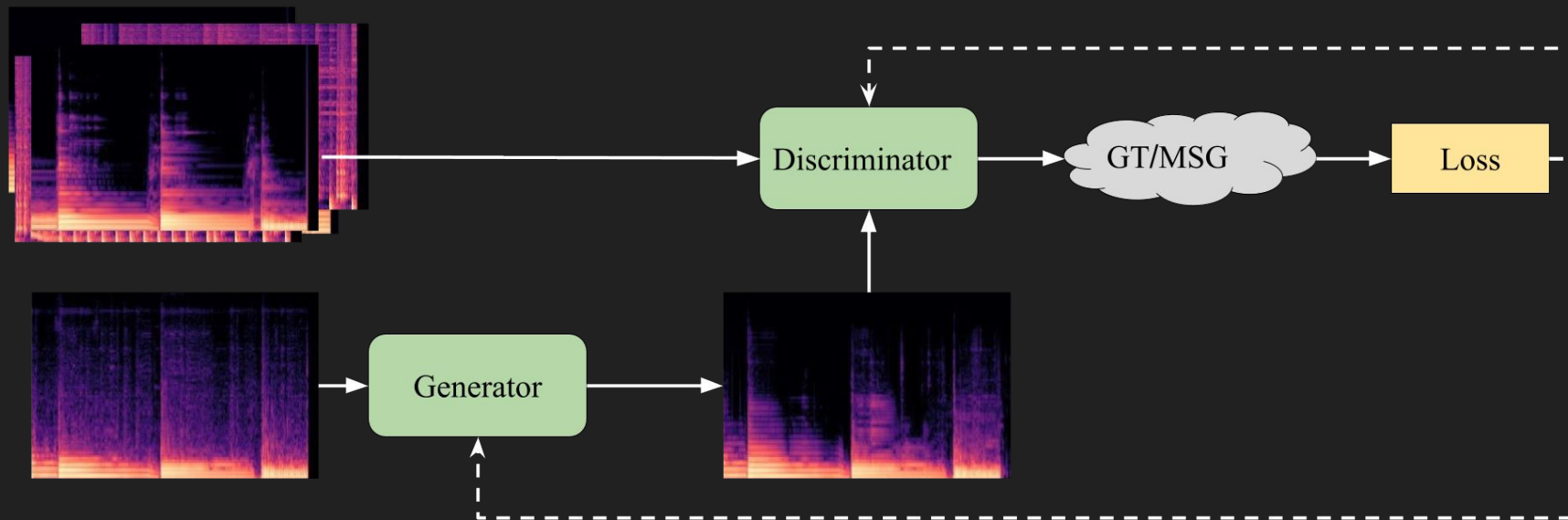# What's a good way to do that?



Standard Source Separation

Mixture

Source Separator

Source Estimate

Post Processing (Ours)

MSG

Output of MSG

- Treat source separator as a black box
- Use a post-processor to enhance its output

# Use a GAN!

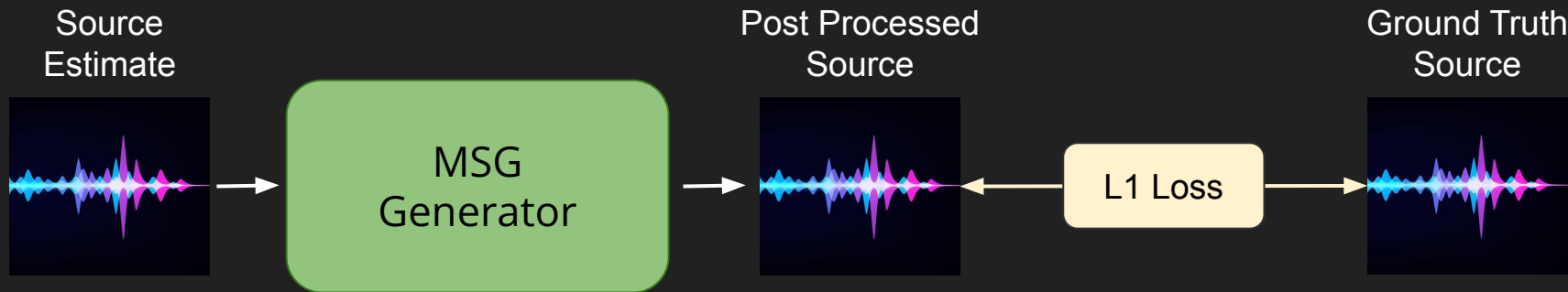Similar to speech denoising (e.g., HiFi-GAN [Su et. al.])
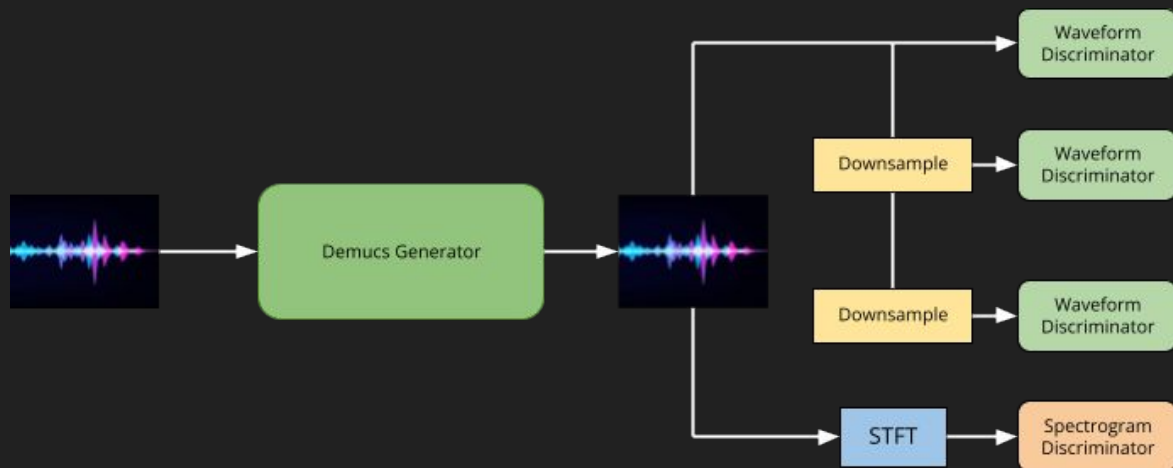
# MSG Training Procedure

1. Supervised pretraining using L1 waveform Loss

2. Add adversarial loss w/ discriminators

# Supervised Pretraining

- 50 epochs using L1 waveform Loss against GT sources

- More stable GAN training  →  No generator collapse

# Add Adversarial Loss



- 150 total training epochs

- Spectral discriminator weighting = combined waveform discriminator weighting

- Discriminator structure similar to HiFi-GAN [Su et. al. 2020]

# Objective Results – MUSDB18 test set

| Source | Demucs | Demucs + MSG (Ours) |
|--------|--------|---------------------|
| Bass | 6.91 | 6.78 |
| Drums | 7.25 | 7.03 |

Median SDR (dB)  –  ↑ Higher is better

*Slightly worse SDR!*
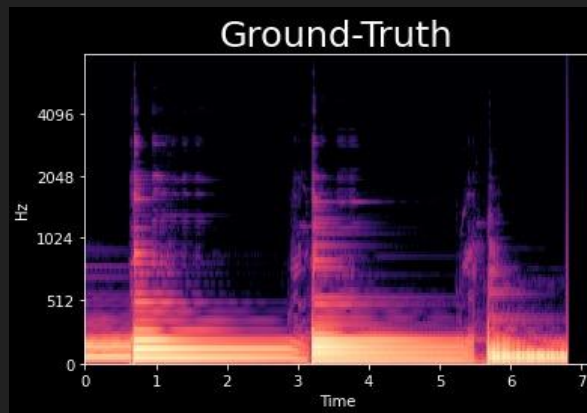
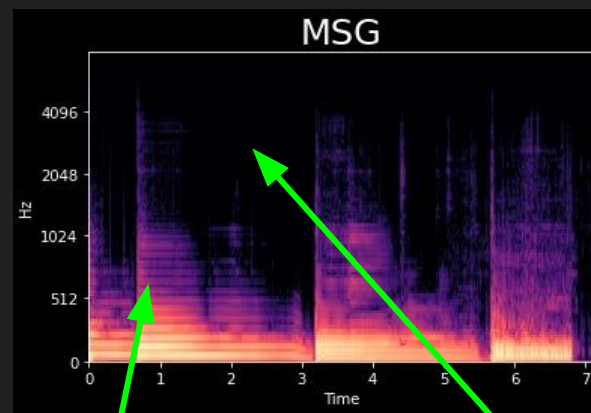*...how do they sound?*

# Illustrative Results — Bass



SDR: 7.40 · SDR: 6.88

Demucs · Ground-Truth · MSG

Missing Overtones · Added Noise · Reconstructed Overtones · Removed Noise

But which one *sounds* better?

# Failure case! – Drums



Demucs

SDR: 8.63
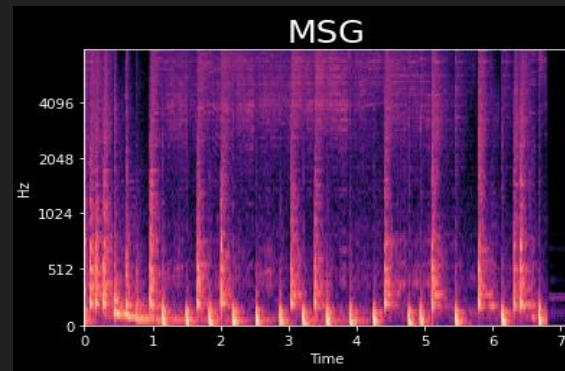
Ground-Truth

Can't reconstruct cymbal
frequencies if they're totally
removed

MSG

SDR: 8.54

# Discussion

- Mismatch between perceptual & SDR eval metrics → GANs get knocked on SDR

- GANs can add sounds to source estimates that hurt objective metrics...but they still might sound good.

- So what is the goal of source separation? Do well on SDR? Or to make it sound good?

# Next steps / Additional Ideas

- Listening studies!

- Can we also improve SDR?

- Condition MSG on the mix
  - Restore content that separation erases

- GAN Loss during separation training
  - Integrate MSG post-processor and separation

- Multi-source MSG models

# Full System Architecture