

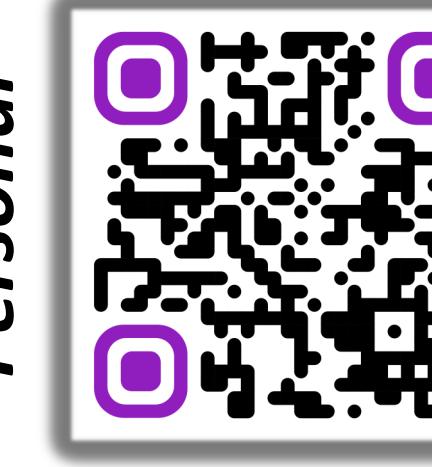
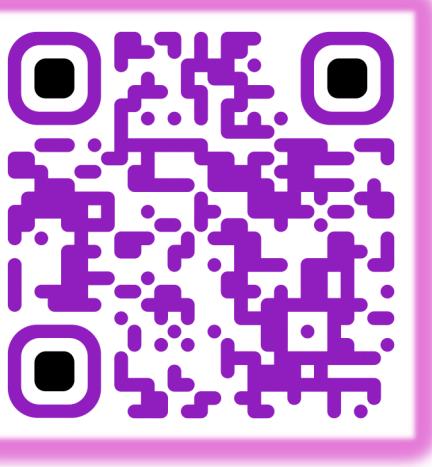
DITTO: Diffusion Inference-Time T-Optimization for Music Generation

Zachary Novack^{1,2}, Julian McAuley¹, Taylor Berg-Kirkpatrick¹, Nicholas J. Bryan²

¹UC – San Diego

²Adobe Research

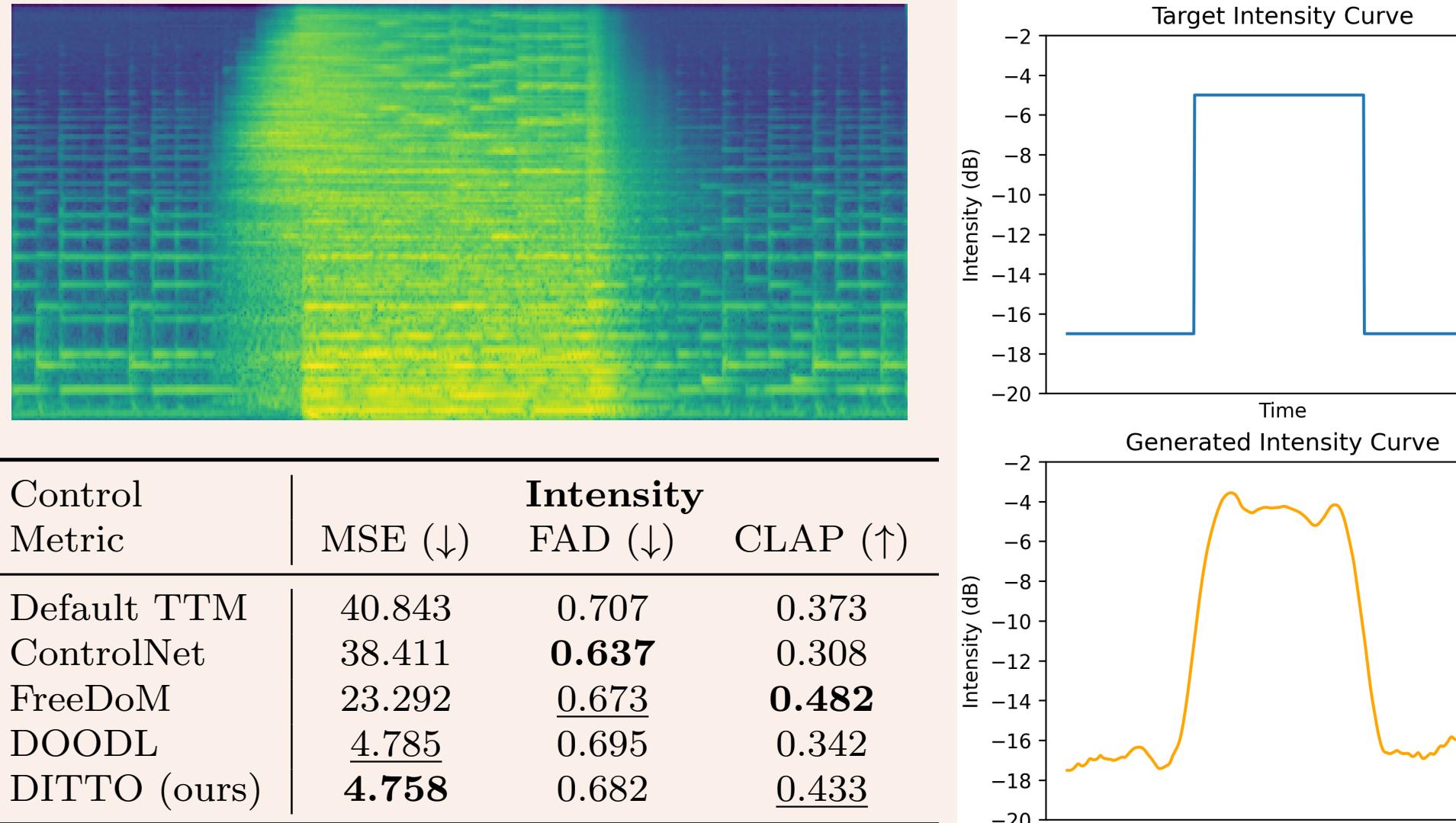
Project



Control Tasks

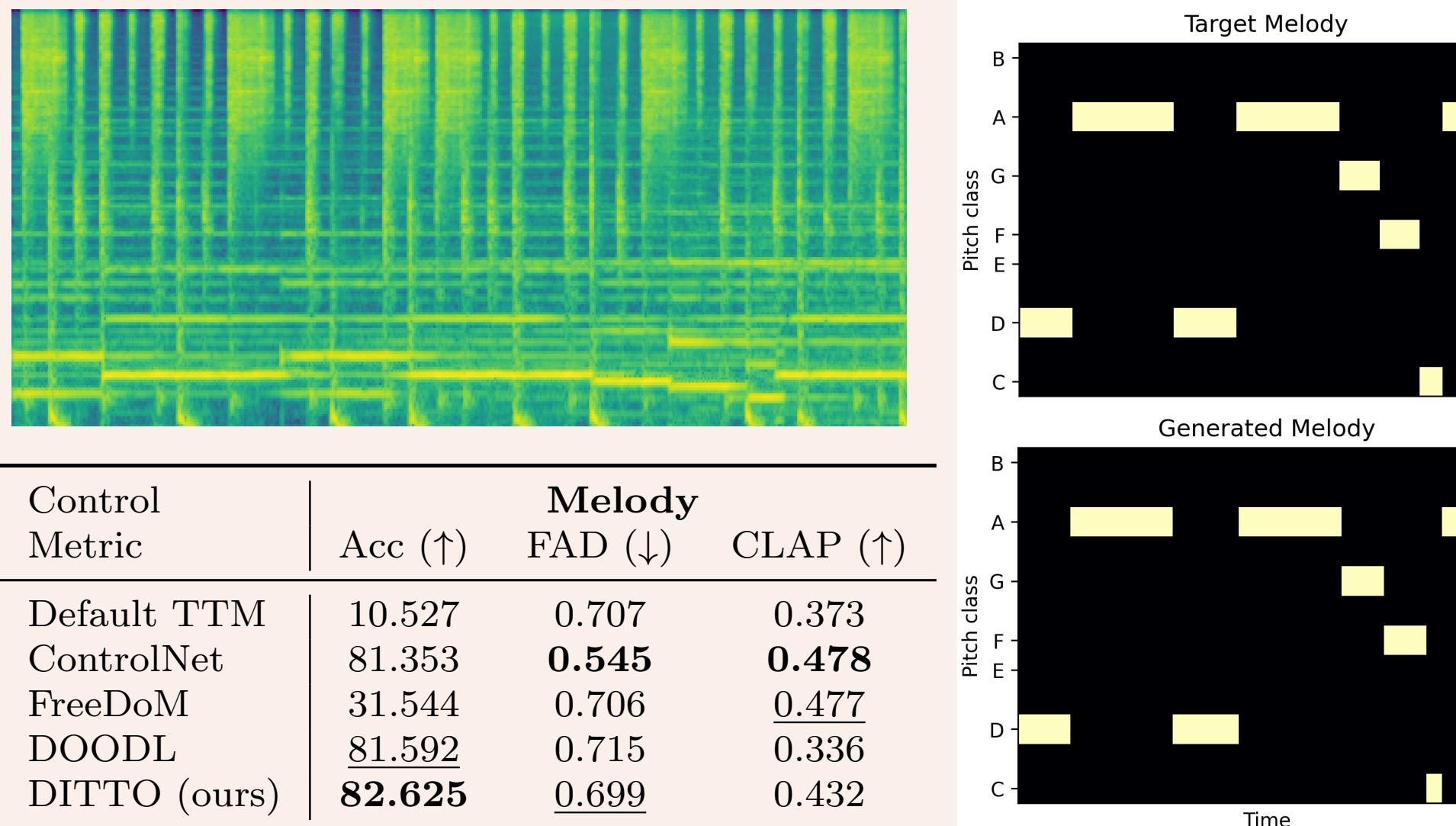
Intensity Control

$f(\cdot)$ = Smooth RMS Energy, \mathcal{L} = MSE



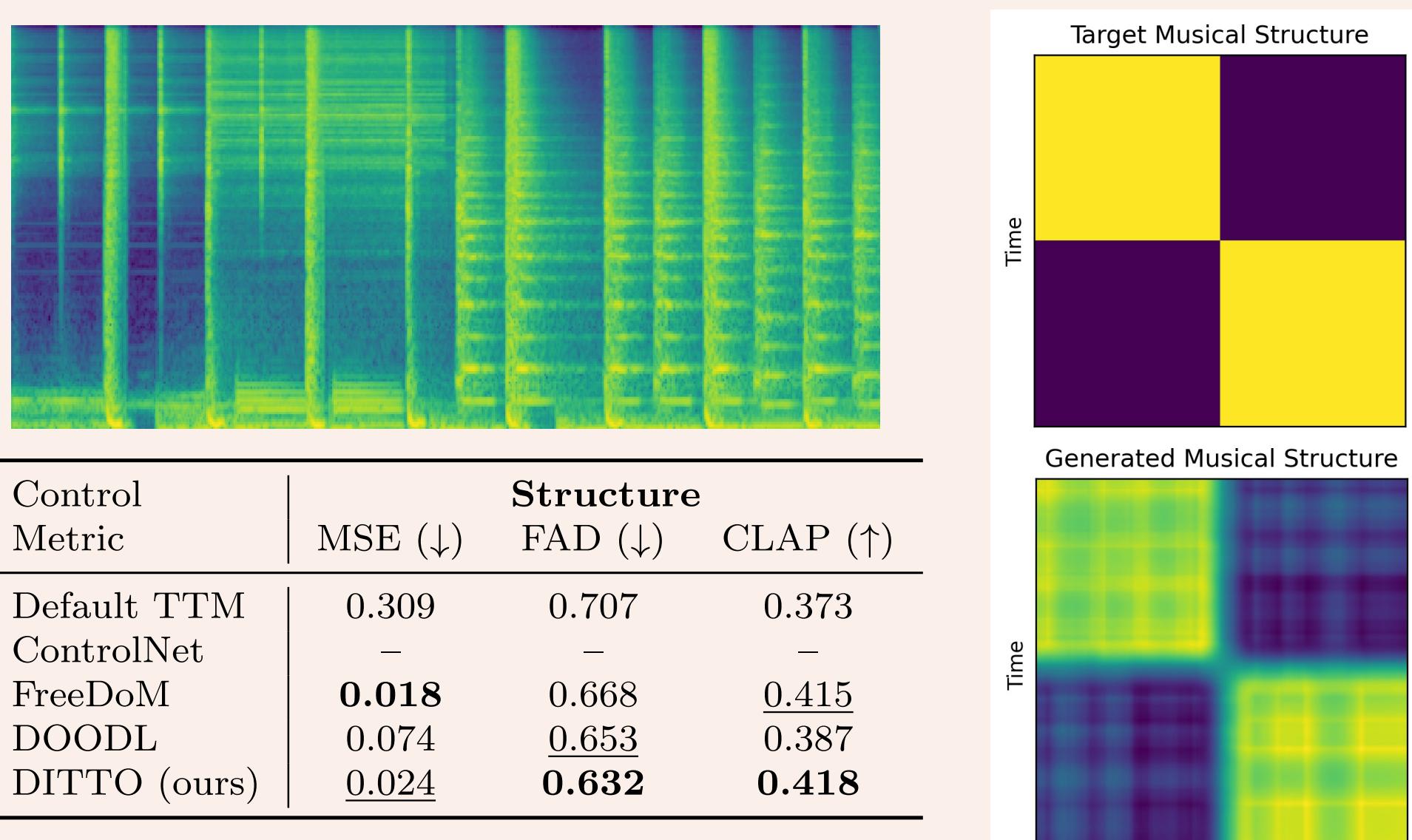
Melody Control

$f(\cdot)$ = Normalized Chromagram, \mathcal{L} = CELoss



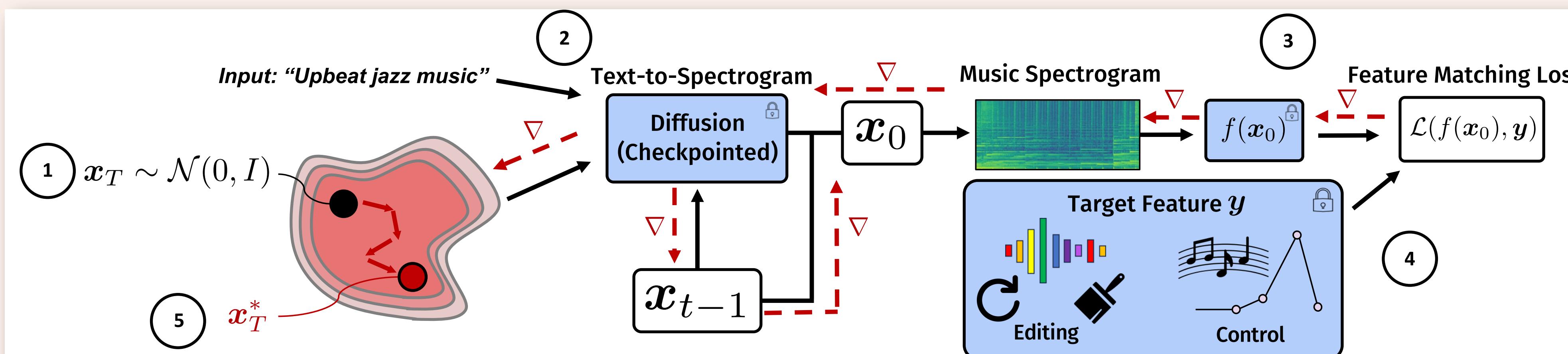
Structure Control

$f(\cdot)$ = Smooth MFCC Self-Similarity, \mathcal{L} = MSE



DITTO: Diffusion Inference-Time T-Optimization

"writing about music is like dancing about architecture"



TLDR: We optimize the **initial noise latent x_T** through diffusion sampling to match some feature of our generation $f(x_0)$ to the target feature y , w/**gradient checkpointing** for efficient memory usage

Method Comparison

Training-Based (Conditioning)

Music-ControlNet, JASCO

- ✓ Arbitrary controls
- ✓ Control/quality balance
- ✓ Fast @ inference time
- ✗ Large-scale training
- ✗ Paired/labeled control data
- ✗ Fixed controls @ training

DITTO (Optimization)

- ✓ Any (differentiable) control
- ✓ Architecture/sampler agnostic
- ✓ Zero training
- ✓ Moderate inference costs
- ✗ Approximate gradients
- ✗ Limited control in low SNR
- ✗ Bad at fine-grained controls

Training-Free (Guidance)

Classifier Guidance, FreeDoM

- ✓ Any (differentiable) control
- ✓ Zero training
- ✓ Moderate inference costs
- ✗ Approximate gradients
- ✗ Limited control in low SNR
- ✗ Bad at fine-grained controls

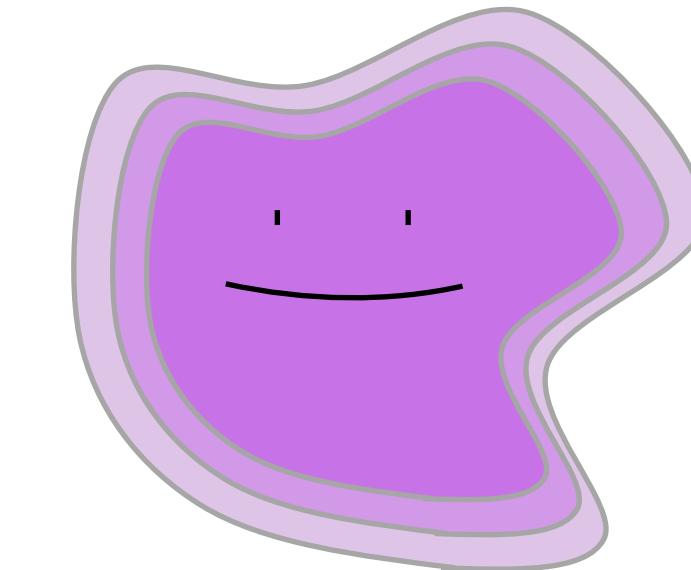
And More!

Further Analysis:

- Subject listening study
- Efficiency evaluation w/DOODL

Extra Use Cases:

- Correlation-based intensity control
- Multi-feature optimization
- Reference-free looping
- Musical structure transfer
- Optimized latent reuse
- Real-audio inversion



Editing Tasks

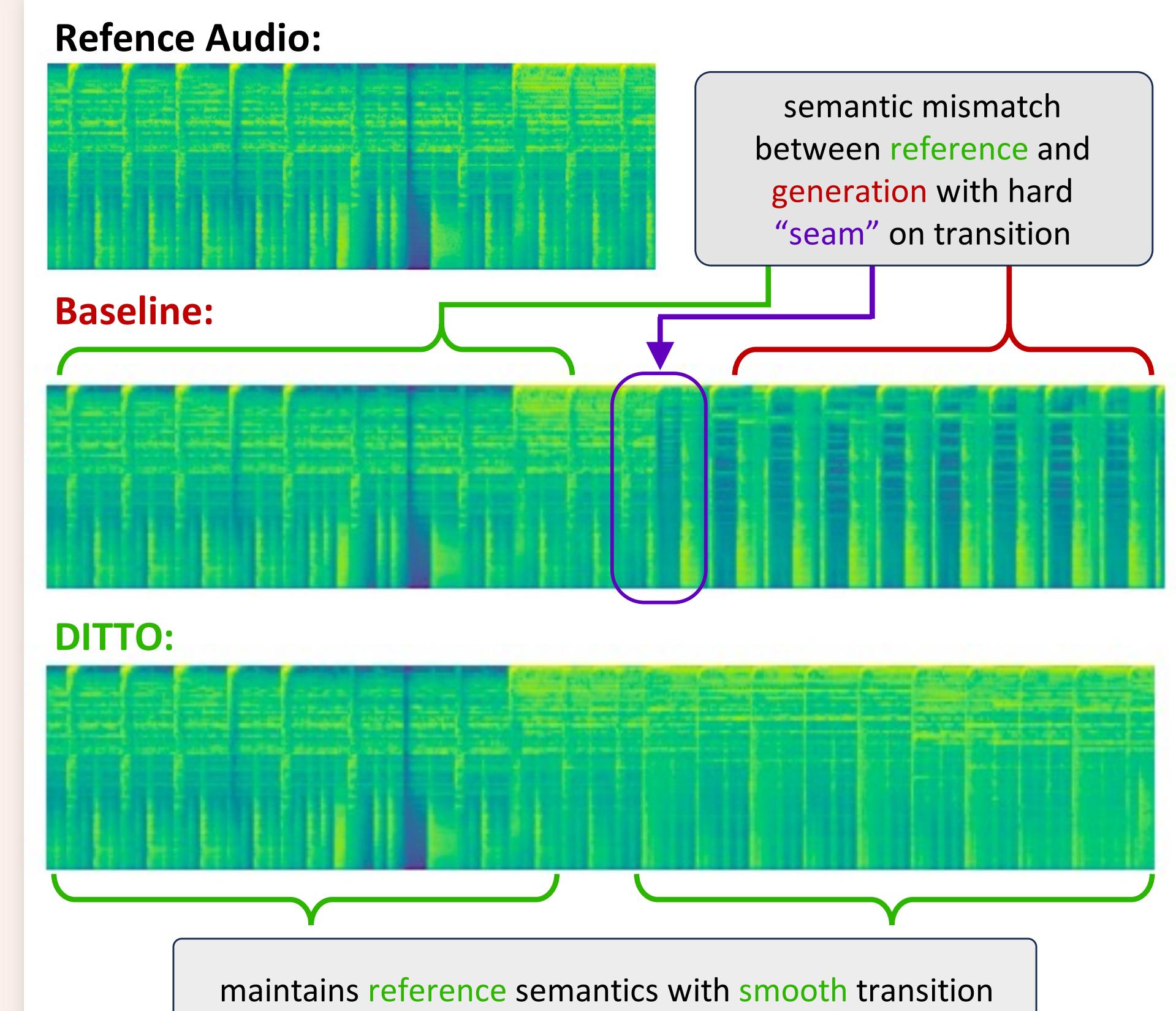
Outpainting

$f(\cdot) = 1$ (LR) Pixel Mask, $\mathcal{L} = \text{MSE}$

Looping

$f(\cdot) = 2$ (LR, RL) Pixel Masks, $\mathcal{L} = \text{MSE}$

Method	Outpainting			Inpainting			Looping
	$o = 1$	$o = 2$	$o = 3$	$o = 2$	$o = 3$	$o = 4$	
DOODL	0.719	0.707	0.700	0.696	0.693	0.688	0.750
Naive	0.722	0.716	0.712	0.707	0.705	0.697	0.753
MD	0.733	0.716	0.710	0.701	0.694	0.690	0.749
MD-50	0.718	0.714	0.705	0.711	0.708	0.701	0.752
GG	0.754	0.738	0.719	0.717	0.709	0.700	0.774
FreeDoM	0.726	0.723	0.715	0.719	0.709	0.704	0.758
DITTO (ours)	0.716	0.703	0.698	0.690	0.688	0.686	0.746



Algorithm

Algorithm 1 Diffusion Inference-Time T-Optimization (DITTO)

```

input :  $\epsilon_\theta$ , Sampler, sampling steps  $T$ , feature extractor  $f$ , loss  $\mathcal{L}$ , target feature  $y$ , starting latent  $\mathbf{x}_T$ , text conditioning  $c$ , optimization steps  $K$ , optimizer  $g$ .
1: // Run optimization
2: for  $i = 1$  to  $K$  do
3:    $\mathbf{x}_t \leftarrow \mathbf{x}_T$  // Initialize noise latents
4:   for  $t = T$  to 1 do // Diffusion sampling w/checkpointing
5:      $\mathbf{x}_{t-1} = \text{Checkpoint}(\text{Sampler}, \epsilon_\theta, \mathbf{x}_t, t, c)$ 
6:   end for
7:    $\hat{\mathbf{y}} = f(\mathbf{x}_0)$  // Extract features from generated output
8:    $\mathbf{x}_T \leftarrow \mathbf{x}_T - g(\nabla_{\mathbf{x}_T} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}))$  // Compute feature loss & backprop
9: end for
output :  $\mathbf{x}_0$ 

```