

## Context

Differentiable Digital Signal Processing (DDSP): training neural networks to estimate parameters of signal models (e.g sinusoidal frequency, amplitudes) using synthesis and reconstruction.

### Main takeaways

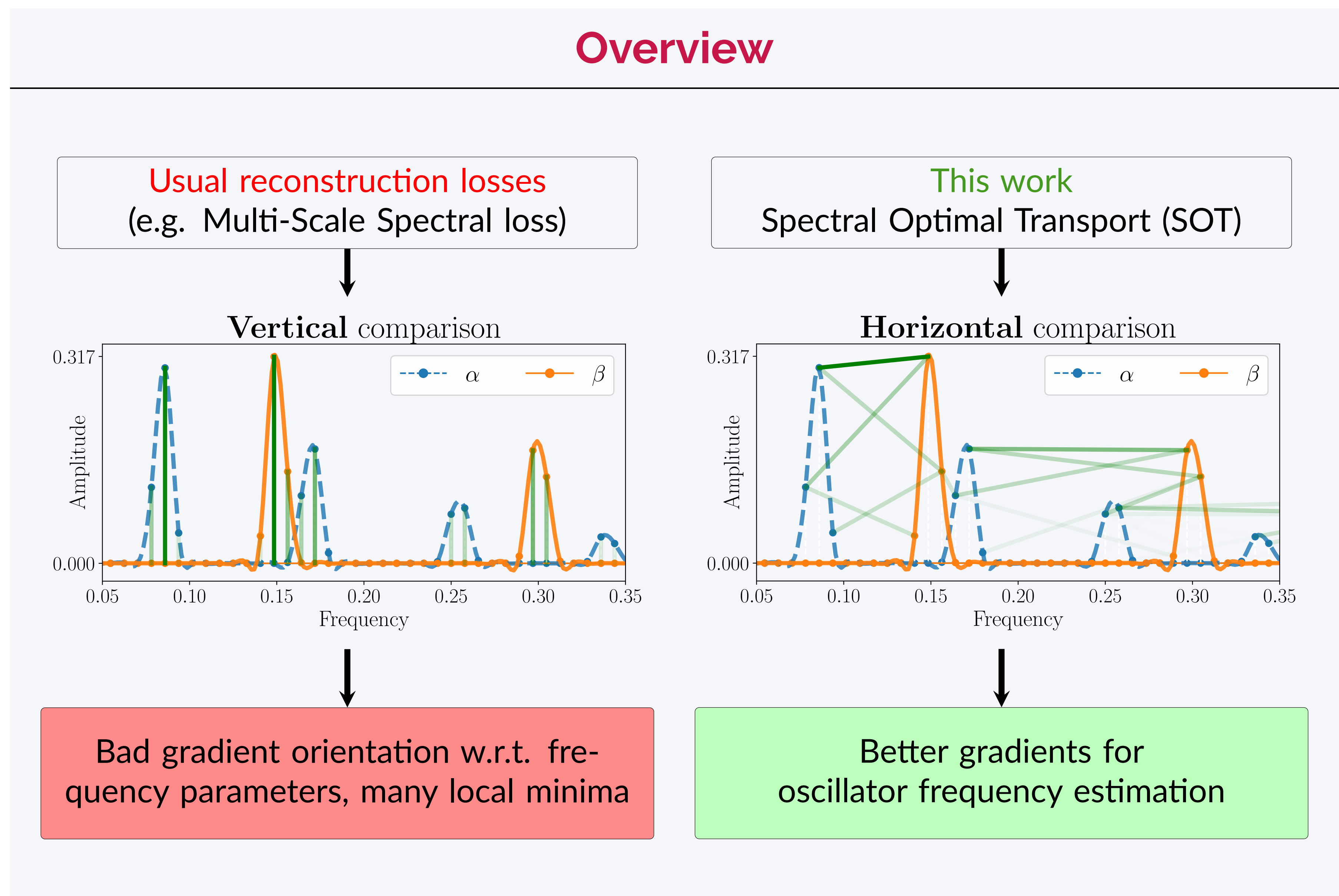
Spectral Optimal Transport (SOT) compares audio measuring frequency displacement of spectral frames

✓ Improves pitch accuracy and reconstruction error when estimating jointly the  $f_0$  and amplitudes of a DDSP harmonic synthesizer (no supervision)

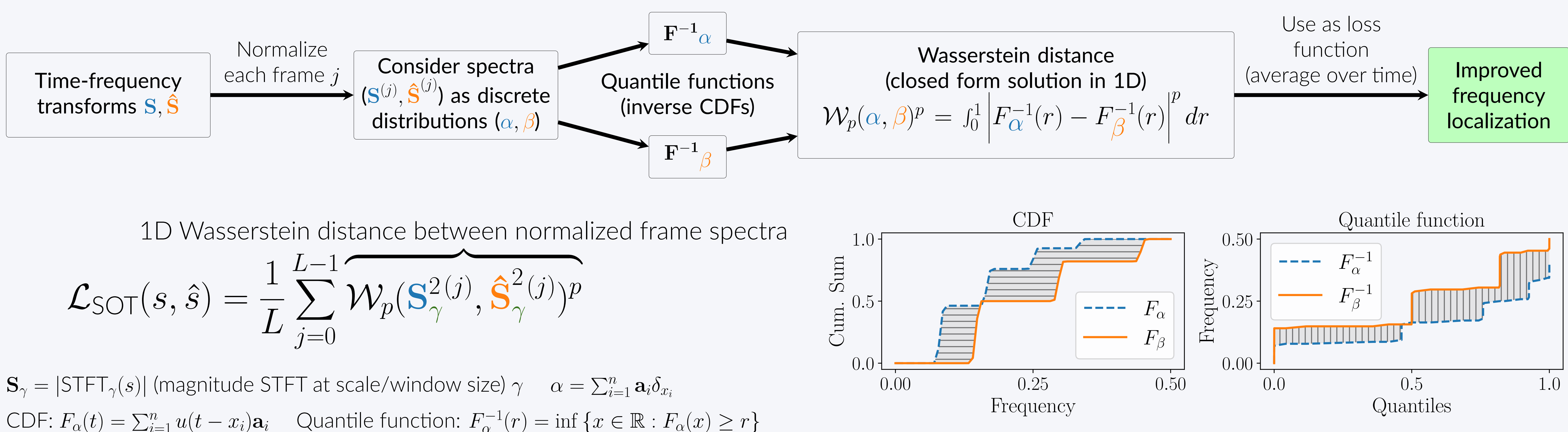
! Sensitive to spectrum normalization and leakage



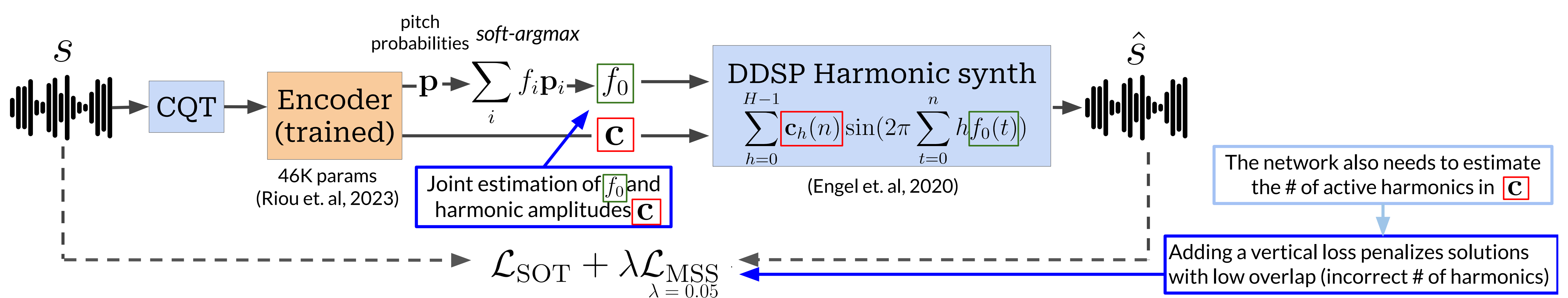
## Overview



## How does Spectral Optimal Transport work ?



## Unsupervised harmonic parameter estimation autoencoder



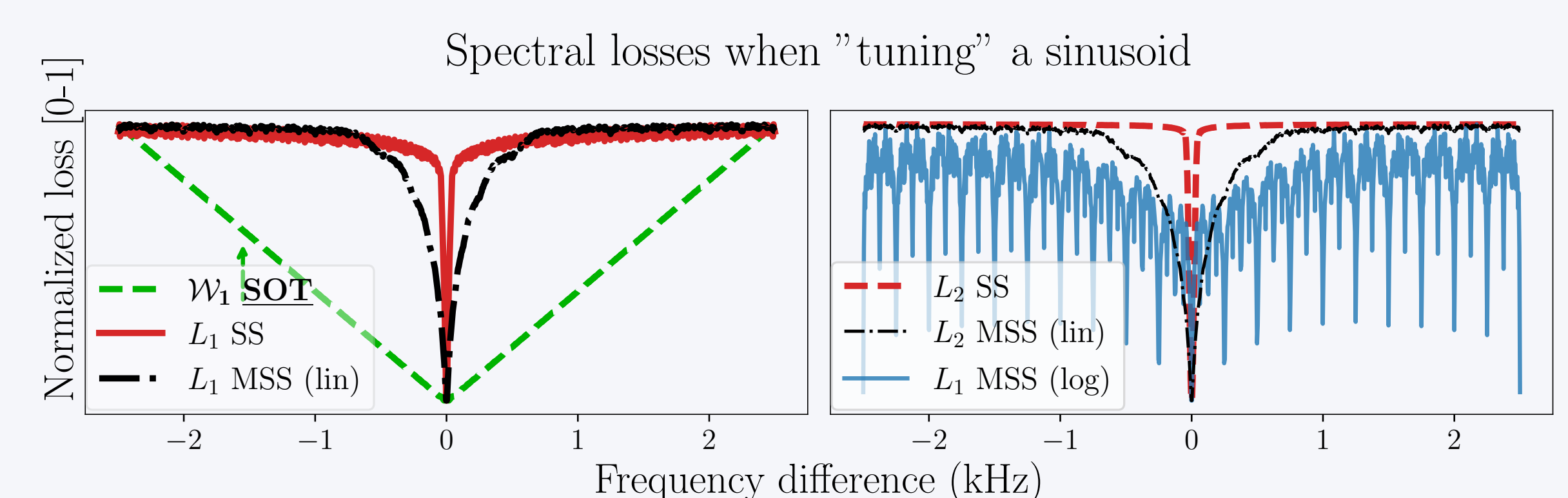
## Results

Baseline: Multi-Scale spectral loss:  $\mathcal{L}_{\text{MSS}}(s, \hat{s}) = \sum_{\gamma \in \Gamma} \left( \|\mathbf{S}_\gamma - \hat{\mathbf{S}}_\gamma\|_1 + \left| \log(\mathbf{S}_\gamma) - \log(\hat{\mathbf{S}}_\gamma) \right| \right)$

Synthetic dataset varying  $f_0$ , harmonic amplitudes, and # of harmonics ([1-8])

Evaluation on pitch estimation and reconstruction metrics

$\mathcal{L}$	Variations	Mean/Median (STD) test metrics (5 runs)						
		$\gamma$ ( $\mathcal{L}_{\text{SOT}}$ )	$\Gamma$ ( $\mathcal{L}_{\text{lin}}$ )	LogF	$f_{\text{cut}}$	LSD [dB] ↓	RPA [%] ↑	RCA [%] ↑
LIN	-	$\Gamma_0$	-	-	-	46.4 / 58.0 (21.4)	20.2 / 0.2 (44.6)	26.9 / 3.9 (42.7)
MSS	-	$\Gamma_0$	-	-	-	80.5 / 82.6 (15.1)	1.4 / 0.1 (2.7)	4.0 / 3.2 (4.5)
SOT	2048	$\Gamma_0$	×	✓	✓	23.5 / 24.5 (3.5)	75.0 / 99.7 (43.2)	99.2 / 99.8 (1.6)
SOT	512	$\Gamma_0$	×	✓	✓	40.5 / 26.6 (23.5)	42.9 / 63.6 (39.4)	62.3 / 75.2 (42.6)
SOT	512	$\Gamma_0$	✓	✓	✓	25.9 / 25.0 (2.5)	55.4 / 63.7 (36.1)	86.8 / 95.6 (16.2)
SOT	2048	$\Gamma_0$	×	×	×	70.6 / 77.6 (31.8)	23.7 / 20.0 (30.3)	46.0 / 45.0 (36.4)
SOT	2048	{512}	×	✓	✓	97.9 / 101.1 (32.5)	14.1 / 4.7 (25.5)	28.6 / 11.6 (32.6)



Vertical ( $L_1$ ,  $L_2$ ) converge smoothly **only** when close to global min.  
 ✓ Horizontal SOT has good gradient orientation

- ✗ High sensitivity to initialization (specially MSS baseline)
- ✓ SOT improves on MSS
- ✓ Larger window size, logarithmic frequency scaling and frequency cutoff → improved metrics
- ! Uncertainty in the number of harmonics → tricky optimization