

Time-Frequency Audio Similarity using Optimal Transport

Linda Fabiani¹, Filip Elvander¹, Sebastian J. Schlecht²

¹ Structured and Stochastic modeling group, Aalto University

*² Multimedia Communications and Signal Processing, Friedrich-
Alexander-Universität Erlangen-Nurnberg*

Topics

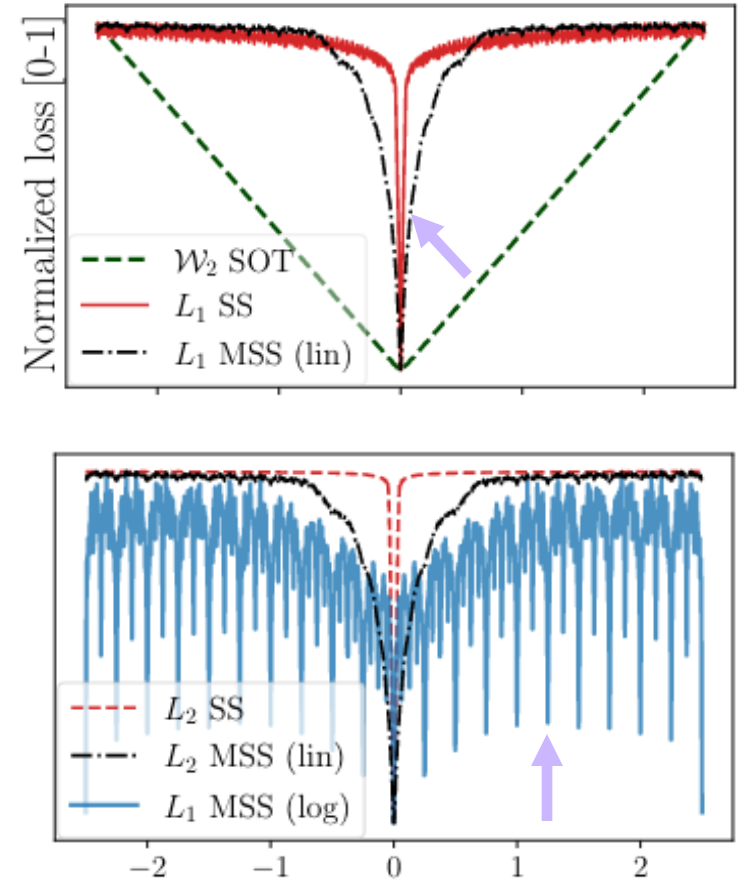
1. Problem motivation
2. Background on optimal transport
3. Optimal transport for audio
4. Applications and results
5. Conclusions and future works

1. Problem motivation

Typical audio-to-audio losses (l_1 , l_2 , Multi-scale Spectral loss) show limited performances when used in neural audio synthesis tasks



- ✗ **rapid growth** for small shifts and **saturate** as soon as the compared signals no longer have shared time-frequency support
- ✗ Presence of **local minima** that can impact the gradients propagation

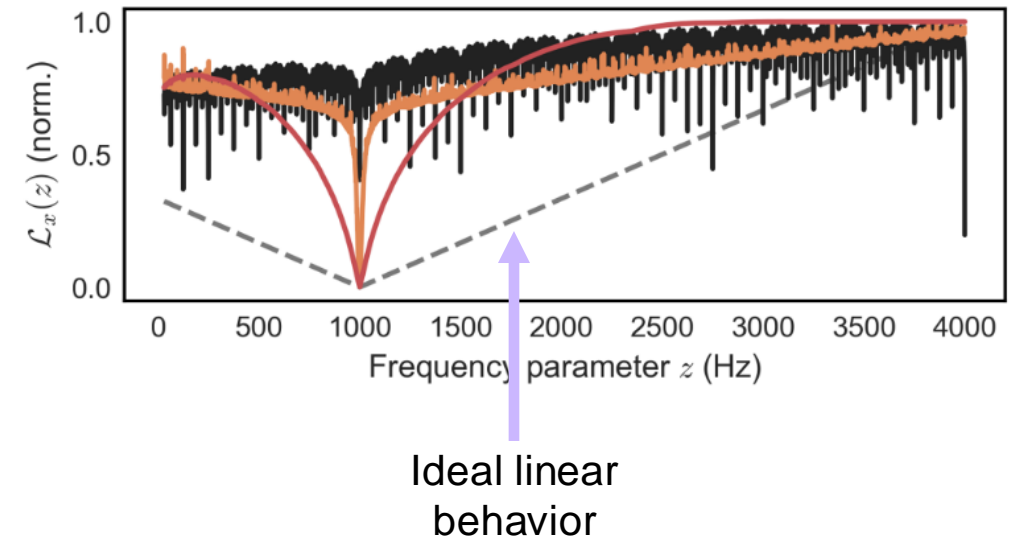


1. Problem motivation


Typical audio-to-audio losses (l_1 , l_2 , Multi-scale Spectral loss) show limited performances when used in neural audio synthesis tasks



- ✓ New loss function for comparing audio signals in **time** and **frequency simultaneously** to produce clean and informative gradients
- ✓ Use optimal transport to ensure robustness with respect to the **geometric space** of the signals' spectral power

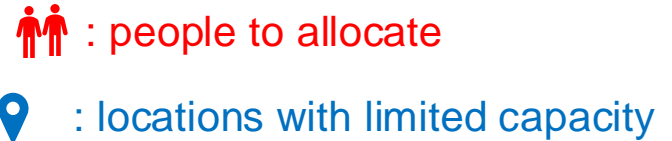


A!

$\mathbf{a} \in \mathbb{R}^N$
 $\mathbf{b} \in \mathbb{R}^M$  Two continuous or discrete distributions (histograms)

$$\mathbf{C} \in \mathbb{R}^{N \times M} \longrightarrow c_{ij} = c(x_i, y_j)$$

Cost of transport



2. Background on optimal transport

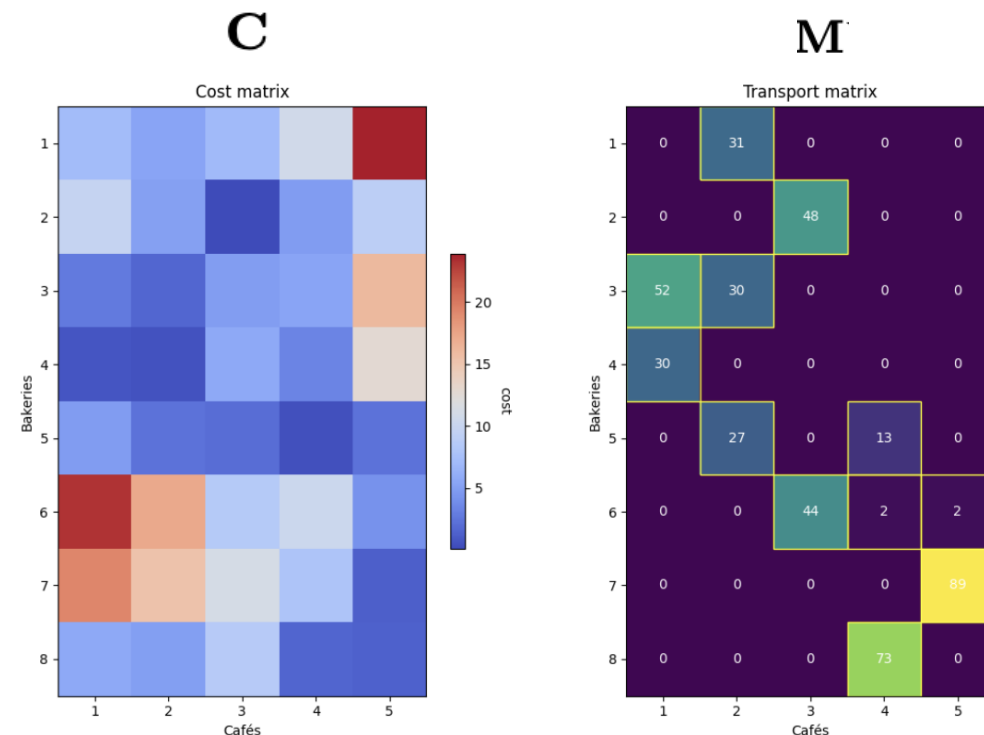
$$\begin{aligned} & \underset{\mathbf{M} \geq 0}{\text{minimize}} && \langle \mathbf{C}, \mathbf{M} \rangle \\ & \text{s.t.} && \mathbf{M} \mathbf{1}_N = \mathbf{a}, \mathbf{M}^T \mathbf{1}_M = \mathbf{b}. \end{aligned}$$

mass conservation constraints

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) = \{ \mathbf{M} \in \mathbb{R}_+^{N \times N} : \mathbf{M} \mathbf{1}_N = \mathbf{a}, \mathbf{M}^T \mathbf{1}_N = \mathbf{b} \}$$

\mathbf{M}^* = optimal transport **plan** out of all possible pairings between \mathbf{a} and \mathbf{b}

✓ Convex objective function and constraints \longrightarrow solve as a **linear program** (Sinkhorn, entropy regularized)



$$\Psi(\mathbf{a}, \mathbf{b}) \triangleq \langle \mathbf{C}, \mathbf{M}^* \rangle$$

3. Optimal transport for spectrograms: Sinkhorn-based approach

Mass: Normalized spectral energy

Cost: between pairs of **time** and **frequency** bins

$$c((t_1, \omega_1), (t_2, \omega_2)) = \underbrace{(t_1 - t_2)^2}_{C_\tau} + \underbrace{(\omega_1 - \omega_2)^2}_{C_\Omega}$$

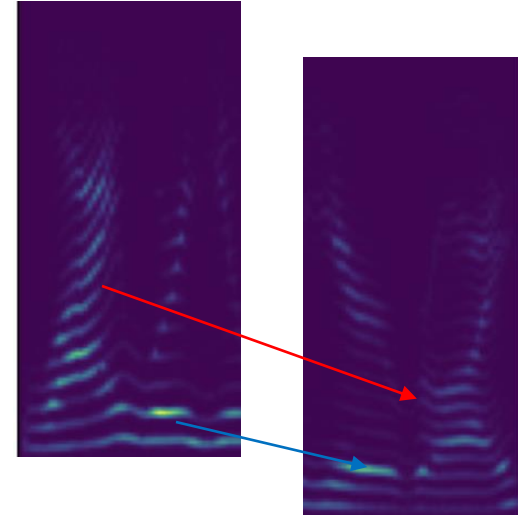


Construct the cost matrix related to transport in **time-frequency domain** using the spectrograms in vectorized form:

$$\mathbf{C} = k \cdot \mathbf{C}_\tau \otimes \mathbf{1}_{|\Omega| \times |\Omega|} + \mathbf{1}_{|\mathcal{T}| \times |\mathcal{T}|} \otimes \mathbf{C}_\Omega$$



The **parameter** k is used to balance the time and frequency costs magnitudes



$$\mathcal{T} = \{t_1, t_2, \dots, t_{|\mathcal{T}|}\}$$

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_{|\Omega|}\}$$

3. Optimal transport for spectrograms: Sliced Wasserstein approach

- Helpful to reduce the **dimensionality** of the problem by projecting the distributions onto subspaces of lower dimensions
- Solve multiple 1D OT problems and average their distances

New "position" vector

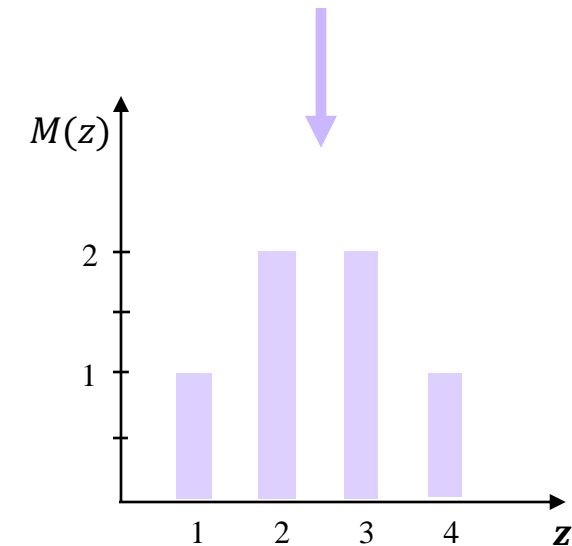
$$\mathbf{z}^i = k \cdot t_\ell \cdot \theta_i^1 + \omega_m \cdot \theta_i^2$$

with $\varphi_i = \frac{2\pi i}{N}$ for $i = 0, 1, 2, \dots, N-1$

$\rightarrow M(z)$

distribution of mass corresponding to the unique values of z

$\omega_1 = 1$	1	2	2	1
$\omega_2 = 2$				
$\omega_3 = 3$				
	$t_1 = 1$	$t_2 = 2$	$t_3 = 3$	$t_4 = 4$



4. Experiments – Sinusoids

Data: time-limited sinusoid sampled at 8 kHz, with central frequency 2 kHz and 2 seconds duration

Comparing OT-based losses with l_1 spectral loss

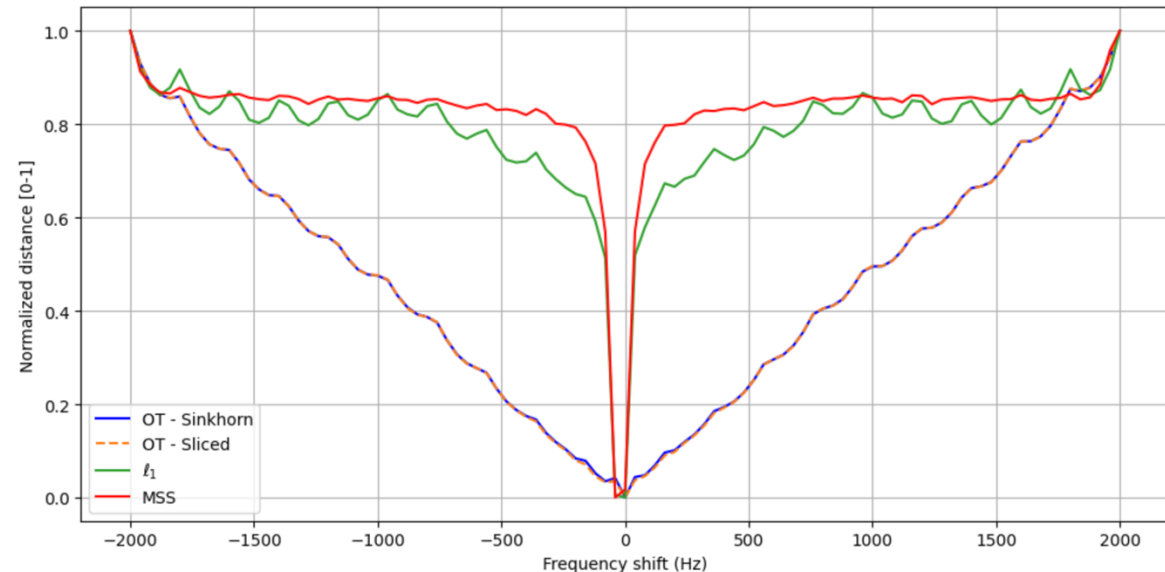
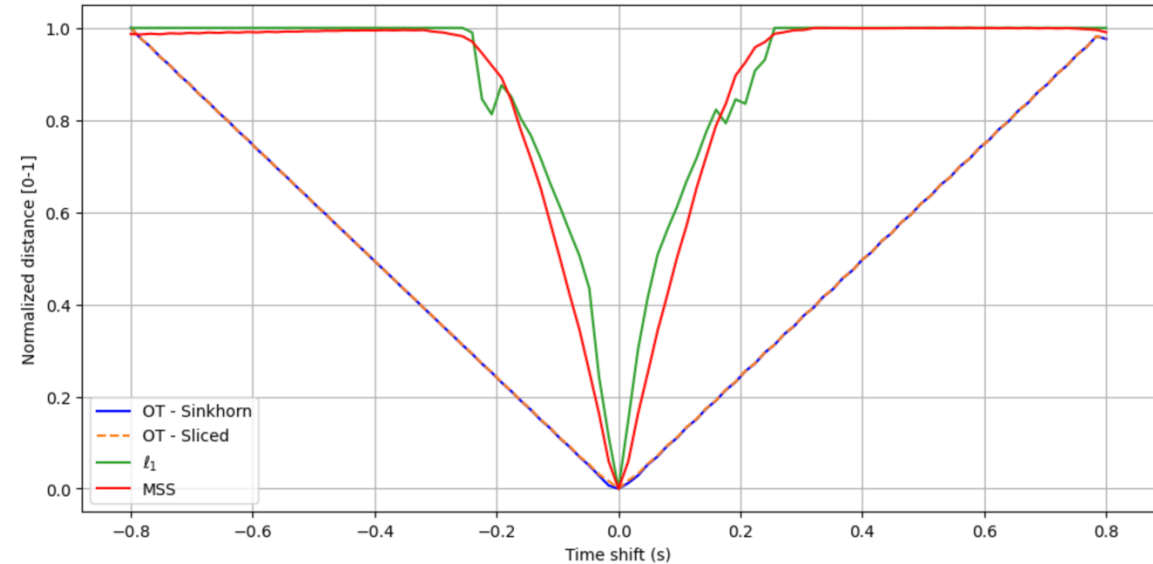
$$\mathcal{L}_1(\mathbf{S}, \hat{\mathbf{S}}) = \|\mathbf{S} - \hat{\mathbf{S}}\|_1$$

and **Multi-Scale Spectral loss**

$$\mathcal{L}_{MSS}(\mathbf{S}_\gamma, \hat{\mathbf{S}}_\gamma) = \frac{1}{|\Gamma|} \sum_{\gamma \in \Gamma} \|\mathbf{S}_\gamma - \hat{\mathbf{S}}_\gamma\|_1 + \|\log(\mathbf{S}_\gamma) - \log(\hat{\mathbf{S}}_\gamma)\|$$

A!

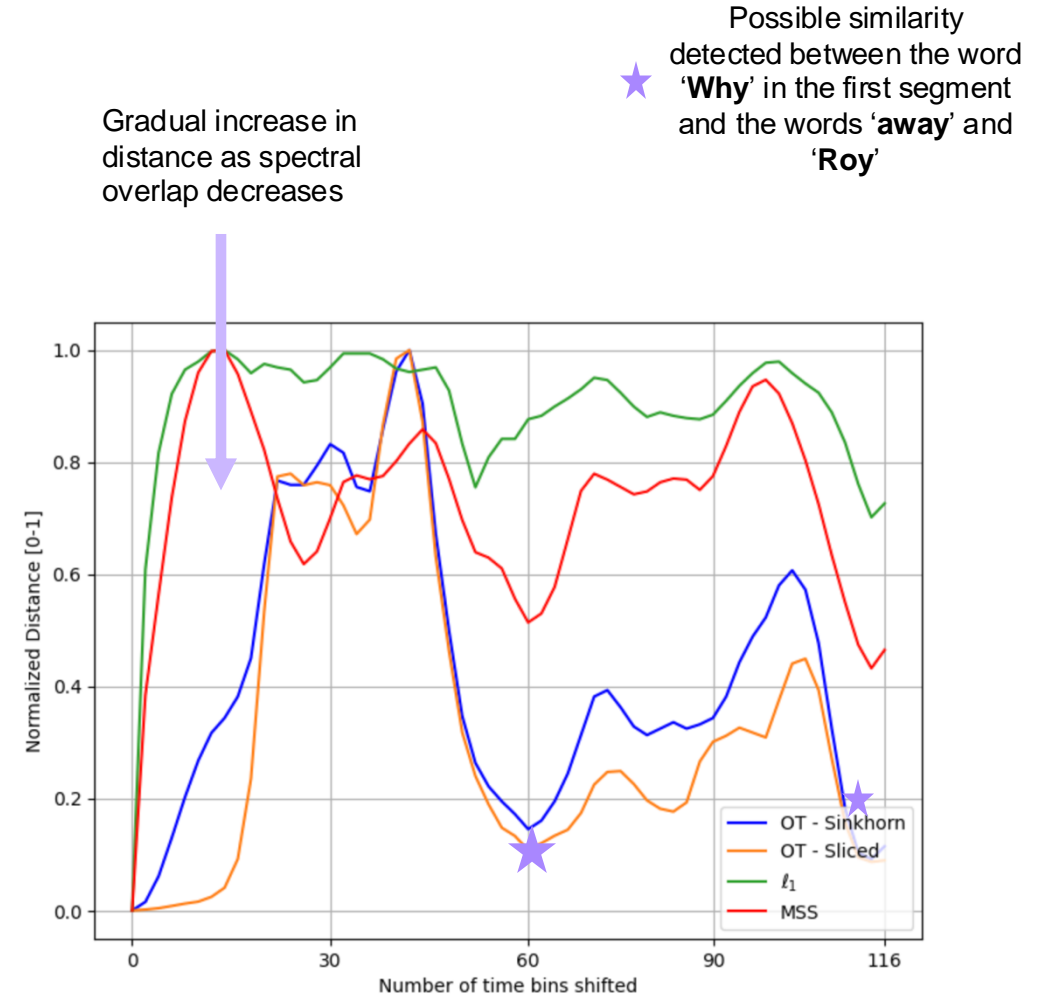
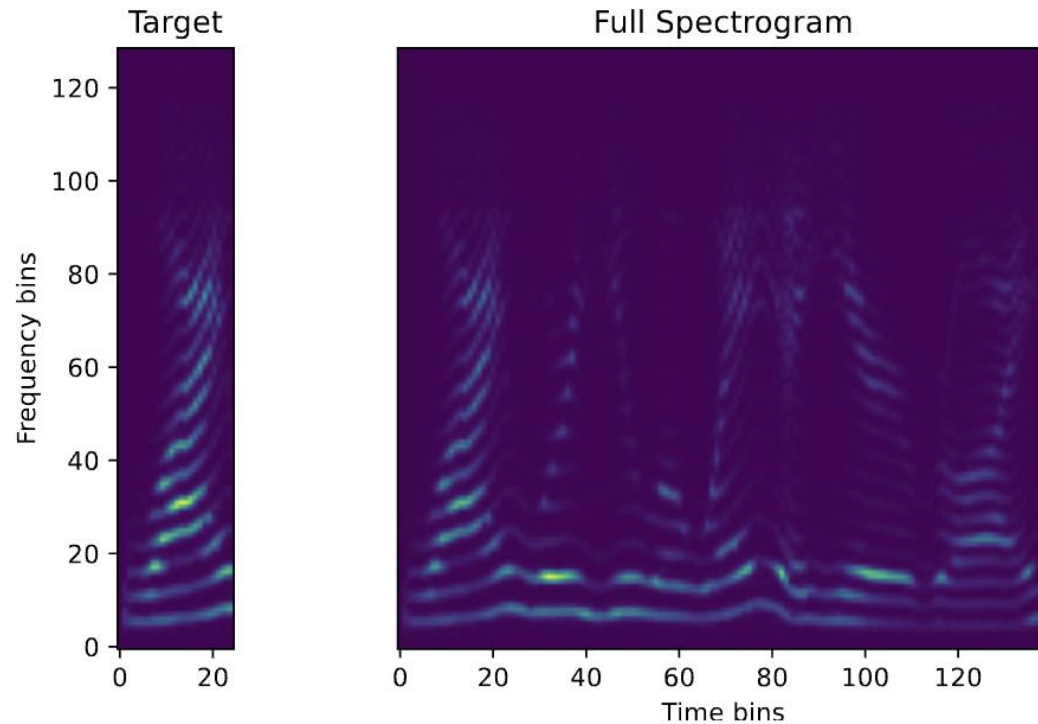
Time shifts



Frequency shifts

4. Experiments – Speech signal

Data: real speech signal corresponding to the phrase "Why were you away a year, Roy?" compared through segments of length 24 bins



A!

Conclusions - Future works

- ✓ OT methods work with **magnitude** of the spectrograms → Extend them to account for the **phase** of complex valued signals
- ✓ OT losses show state-of-the-art performances for simple **DDSP reconstruction task** → explore more advanced DDSP tasks (timbre transfer, singing voice synthesis)
- 💡 Optimize distance calculation to improve computational expenses
- 💡 Include topics from psychoacoustics to the metric and validate with **listening experiments**

A!

Thank you!