

AUDIO LABS

Loss Functions Matter

Three Case Studies in Informed Loss Design

Meinard Müller

International Audio Laboratories Erlangen meinard.mueller@audiolabs-erlangen.de

Lecture Series "Musical Informatics"

Linz, November 19, 2025





Meinard Müller

- Mathematics (Diplom/Master, 1997) Computer Science (PhD, 2001) Information Retrieval (Habilitation, 2007)
- Senior Researcher (2007-2012)
- Professor Semantic Audio Processing (since 2012)
- Former President of the International Society for Music Information Retrieval (MIR)
- IEEE Fellow for contributions to Music Signal Processing

© AudioLabs, 2025





mpii





International Audio Laboratories Erlangen



- Integrated Circuits IIS
- Largest Fraunhofer institute with > 1000 members
- Applied research for sensor, audio, and media technology









- Friedrich-Alexander Universität Erlangen-Nürnberg (FAU)
- One of Germany's largest universities with ≈ 40,000 students
- Strong Technical Faculty



International Audio Laboratories Erlangen





International Audio Laboratories Erlangen



© AudioLabs, 2025 Loss Functions Matter 5



Meinard Müller: Research Group

- Ben Maman
- Simon Schwär
- Johannes Zeitler
- Peter Meier
- Illi Berendes Vlora Arifi-Müller

Sebastian Strahl





- Ching-Yu Chiu (Sunny)
- Yigitcan Özer
- Michael Krause Christof Weiß
- Sebastian Rosenzy
- Frank Zalkow
- Jonathan Driedger

Hendrik Schreiber

Christian Dittmar

Stefan Balke









© AudioLabs, 2025



Meinard Müller: Research Group

- Ben Maman
- Simon Schwär
- Sebastian Strahl Uli Berendes
- Johannes Zeitler Vlora Arifi-Müller Peter Meier
 - Stefan Balke



- Ching-Yu Chiu (Sunny)
- Yigitcan Özer
- Michael Krause Christof Weiß
- Sebastian Rosenzweig

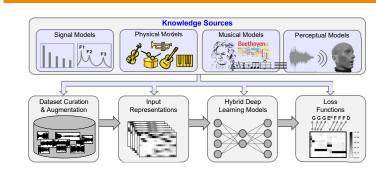


- Stefan Balke
- Jonathan Driedger Thomas Prätzlich



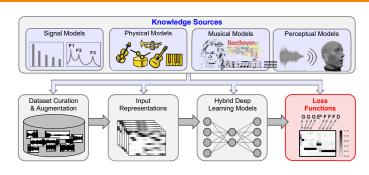






Richard, Lostanlen, Yang, Müller: Model-Based Deep Learning for Music Information Research: Leveraging Diverse Knowledge Sources to Enhance Explainability, Controllability, and Resource Efficiency. IEEE Signal Processing Magazine, 41(6): 51–59, 2024





Richard, Lostanlen, Yang, Müller: Model-Based Deep Learning for Music Information Research: Leveraging Diverse Knowledge Sources to Enhance Explainability, Controllability, and Resource Efficiency. IEEE Signal Processing Magazine, 41(6): 51–59, 2024.

LABS

Overview

- Multi-Scale Spectral Loss Knowledge Source: Signal Representations
- Hierarchical Classification Loss Knowledge Source: Musical Hierarchies
- Differentiable Alignment Loss Knowledge Source: Temporal Coherence



Simon Schwär



Michael Krause



Loss Functions Matter



Overview

- Multi-Scale Spectral Loss Knowledge Source: Signal Representations
- Hierarchical Classification Loss
- Differentiable Alignment Loss







- Literature

 Turian, Henry: I'm sorry for your loss: Spectrally-based audio distances are bad at pitch. Proc. Adv. Neural Inf. Process. Syst., 2020.

 Hayes, Saltis, Fazekas: Sinusoidal frequency estimation by gradient descent. Proc. ICASSP. 2023.

 Torres, Peeters, Richard: Unsupervised Harmonic Parameter Estimation Using DDSP and Spectral Optimal Transport. Proc. ICASSP, 2024.

 Schwär, Müller: Multi-Scale Spectral Loss Revisited. IEEE Signal Processing Letters, 30: 1712–1716, 2023.

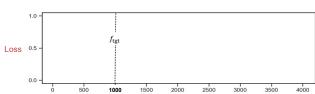
© AudioLabs, 2025

Example Scenario: Sinusoidal Frequency Estimation

Sinusoid with target frequency: $f_{\rm tgt} = 1000~{\rm Hz}$







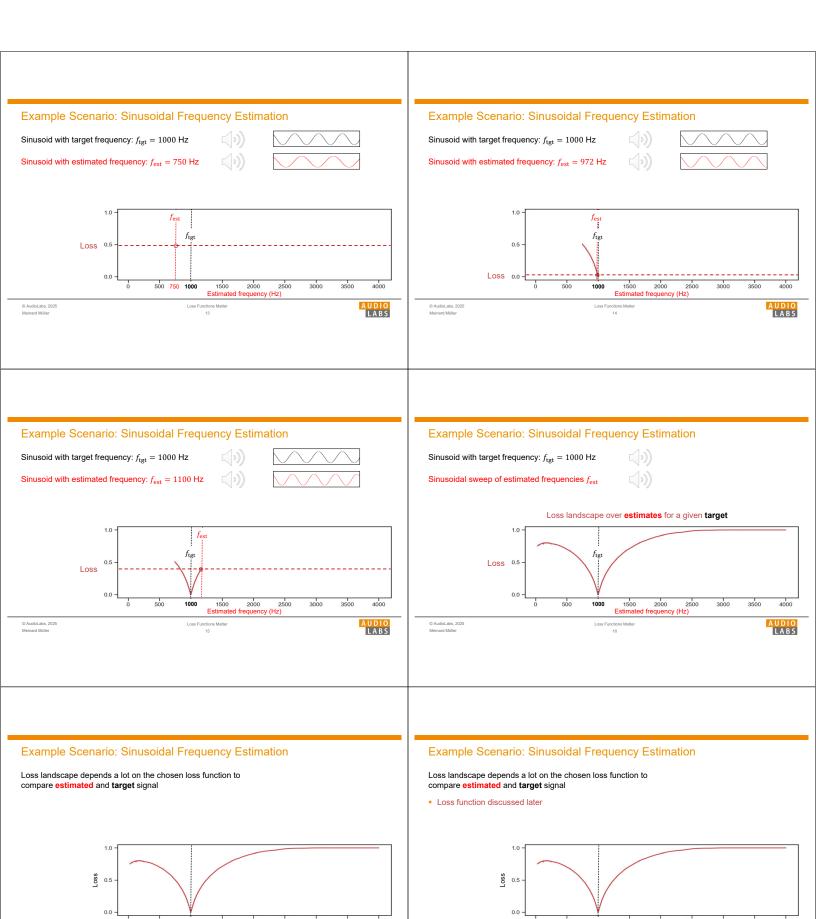
Loss Functions Matter 12

AUDIO LABS

© AudioLabs, 2025

Loss Functions Matter





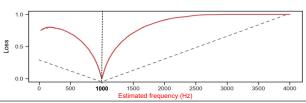
LABS

AUDIO LABS

Example Scenario: Sinusoidal Frequency Estimation

Loss landscape depends a lot on the chosen loss function to compare **estimated** and **target** signal

- Loss function discussed later
- Ideal convex loss



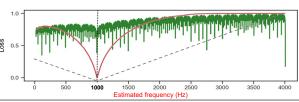
© AudioLabs, 2025

Functions Matter 19 AUDIO LABS

Example Scenario: Sinusoidal Frequency Estimation

Loss landscape depends a lot on the chosen loss function to compare **estimated** and **target** signal

- Loss function discussed later
- Ideal convex loss
- Multi-Scale Spectral (MSS) loss with standard settings



Window Type

© AudioLabs, 2025

unctions Matter

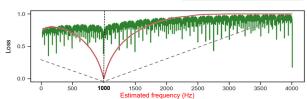
LABS

Example Scenario: Sinusoidal Frequency Estimation

Loss landscape depends a lot on the chosen loss function to compare **estimated** and **target** signal

- Loss function discussed later
- Ideal convex loss
- Multi-Scale Spectral (MSS) loss with standard settings

The MSS loss is what we widely use in audio processing (e.g., DDSP)



© AudioLabs, 2025 Meinard Müller

21

AUDIO

Multi-Scale Spectral Loss

- x input signal
- N window size
- H hop size
- w window function
- p compression function
- d distance function
- \mathcal{N} set of window sizes
- ullet ${\mathcal P}$ set of compression function

 $\mathbf{Spectrum} \qquad \qquad \mathcal{Y}_{w,N,p}(m,k) \ = \ p \left(\left| \sum_{n=0}^{N-1} x[n+mH]w[n] \exp\left(\frac{-i2\pi kn}{N}\right) \right| \right)$

MSS loss

 $\mathcal{L}_{\text{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w, N, p}, \hat{\mathcal{Y}}_{w, N, p})$

© AudioLabs, 2025 Meinard Müller oss Functions Matte 22 AUDIO

Multi-Scale Spectral Loss

۰	х	input signal	
i	N	window size	

H hop size

H HOP SIZE
 W Window function

p compression function

d distance function

lacksquare $\mathcal N$ set of window sizes

ullet ${\cal P}$ set of compression function

$$\mathcal{Y}_{w,N,p}(m,k) = p\left(\left|\sum_{n=0}^{N-1} x \right| n + mH\right] w[n] \exp\left(\frac{-i2\pi kn}{N}\right) \left|\right)$$

MSS loss

$$\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w, N, p}, \hat{\mathcal{Y}}_{w, N, p})$$

© AudioLabs, 2025 Meinard Müller AUDIO

 $d(\mathcal{Y}, \hat{\mathcal{Y}}) = ||\mathcal{Y} - \hat{\mathcal{Y}}||_1$ $d(\mathcal{Y}, \hat{\mathcal{Y}}) = ||\mathcal{Y} - \hat{\mathcal{Y}}||_2^2$

Multi-Scale Spectral Loss

x	input signal	
N	window size	

H hop sizew window function

p compression function

d distance function

• \mathcal{N} set of window sizes

lacksquare set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) = p \left(\sum_{n=0}^{N-1} x[n+nH] w[n] \exp\left(\frac{-i2\pi kn}{N}\right) \right)$$

MSS loss

$$\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w, N, p}, \hat{\mathcal{Y}}_{w, N, p})$$

© AudioLabs, 2025 Meinard Müller Loss Functions Matt

AUDIO

 $d(\mathcal{Y}, \hat{\mathcal{Y}}) = ||\mathcal{Y} - \hat{\mathcal{Y}}||_1$ $d(\mathcal{Y}, \hat{\mathcal{Y}}) = ||\mathcal{Y} - \hat{\mathcal{Y}}||_2^2$

Multi-Scale Spectral Loss

- input signal
- N window size

 H hop size

• w window function

• p compression function

• d distance function N set of window sizes

P set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) \ = p\Biggl(\biggl| \sum_{n=0}^{N-1} x[n+mH \boxed{w} n] \exp\left(\frac{-i2\pi kn}{N} \right) \biggr| \Biggr)$$

MSS loss

$$\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w,N,p}, \hat{\mathcal{Y}}_{w,N,p})$$

© AudioLabs, 2025

Window Type

Window Type

Magnitude Compression



$$\begin{split} \mathcal{N} &= \{67, 127, 257, 509, 1021, 205\\ \mathcal{P} &= \{x\}\\ \mathcal{P} &= \{\log(x+\varepsilon)\}, \ \varepsilon = 10^{-7}\\ \mathcal{P} &= \{\log(x+\gamma)\}, \ \gamma = 1\\ \mathcal{P} &= \{20\log_{10}(x+\gamma)\}, \ \gamma = 1\\ \mathcal{P} &= \{2\log\log_{10}(x+\gamma)\}, \ \varepsilon = 10^{-7}\\ \mathcal{P} &= \{x, \log(x+\varepsilon)\}, \ \varepsilon = 10^{-7}\\ \mathcal{d}(\mathcal{Y}, \hat{\mathcal{Y}}) &= \|\mathcal{Y} - \hat{\mathcal{Y}}\|_1 \end{split}$$

Description Rectangular window

Flat v_{e} $\mathcal{N} = \{64\}$ $\mathcal{N} = \{512\}$ $\mathcal{N} = \{2018\}$ $\mathcal{N} = \{2018\}$ $\mathcal{N} = \{64, 122, 257, 509, 1021, 2053\}$ $\mathcal{P} = \{x\}$ $\mathcal{P} = \{x\}$ $\mathcal{P} = \{\log(x + \epsilon)\}, \epsilon = 10^{-7}$ $\mathcal{P} = \{\log(x + \epsilon)\}, \epsilon = 10^{-7}$ $\mathcal{P} = \{x\}, \log(x + \epsilon)\}, \epsilon = 10^{-7}$ $\mathcal{P} = \{x\}, \log(x + \epsilon)\}, \epsilon = 10^{-7}$ $\mathcal{P} = \{x\}, \log(x + \epsilon)\}, \epsilon = 10^{-7}$ $\mathcal{P} = \{x\}, \log(x + \epsilon)\}, \epsilon = 10^{-7}$

Multi-Scale Spectral Loss

- input signal
- N window size
- *H* hop size
- window function • w
- p compression function
- distance function
- set of window sizes
- set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) = p \left(\left| \sum_{n=0}^{N-1} x[n+mH]w[n] \exp\left(\frac{-i2\pi kn}{N}\right) \right| \right)$$

Window Type

Window Type

Flat Top window $N = \{64\}$

 $\mathcal{P} = \{x\}$ $\mathcal{P} = \{\log(x + \varepsilon)\}, \ \varepsilon = 10^{-7}$ $\mathcal{P} = \{\log(1 + \gamma x)\}, \ \gamma = 1$ $\mathcal{P} = \{20\log_{10}(x + \varepsilon)\}, \ \varepsilon = 1$

 $d(\mathcal{Y}, \hat{\mathcal{Y}}) = \|\mathcal{Y} - \hat{\mathcal{Y}}\|$ $d(\mathcal{Y}, \hat{\mathcal{Y}}) = \|\mathcal{Y} - \hat{\mathcal{Y}}\|$

MSS loss

$$\mathcal{L}_{\text{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w,N,p}, \hat{\mathcal{Y}}_{w,N,p})$$

© AudioLabs, 2025



Rectangular window Hann window Flat Top window $\mathcal{N} = \{64\}$ $\mathcal{N} = \{512\}$ $\mathcal{N} = \{2048\}$ $\mathcal{N} = \{64, 128, 256, 512, 1024, 2048\}$ $\mathcal{N} = \{64, 128, 257, 509, 1021, 2053\}$ $\mathcal{P} = \{r\}$

$$\begin{split} &N = \{\alpha, 1st, sor, sor\} \\ &P = \{x\} \\ &P = \{\log(x + \varepsilon)\}, \, \varepsilon = 10^{-7} \\ &P = \{\log(x + \varepsilon)\}, \, \varepsilon = 10^{-7} \\ &P = \{\log(x + \varepsilon)\}, \, \varepsilon = 10^{-7} \\ &P = \{20\log_0(x + \varepsilon)\}, \, \varepsilon = 10^{-7} \\ &d(\mathcal{Y}, \hat{\mathcal{Y}}) = \|\mathcal{Y} - \hat{\mathcal{Y}}\|_1 \\ &d(\mathcal{Y}, \hat{\mathcal{Y}}) = \|\mathcal{Y} - \hat{\mathcal{Y}}\|_2^2 \end{split}$$

Multi-Scale Spectral Loss

- input signal
- N window size
- H hop size
- window function
- compression function
- distance function d
- set of window sizes N
- . P set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) = p\left(\left|\sum_{n=0}^{N-1} x[n+mH]w[n] \exp\left(\frac{-i2\pi kn}{N}\right)\right|\right)$$

MSS loss

$$\mathcal{L}_{\text{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} \boxed{d(\mathcal{V}_{w,N,p}, \hat{\mathcal{Y}}_{w,N,p})}$$

Loss Functions Matter



Multi-Scale Spectral Loss

- x input signal
 - window size
- hop size H
- w window function
- compression function • p
- distance function d
- N set of window sizes
 - set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) = p\left(\left|\sum_{n=0}^{N-1} x[n+mH]w[n] \exp\left(\frac{-i2\pi kn}{N}\right)\right|\right)$$

MSS loss

$$\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w, N, p}, \hat{\mathcal{Y}}_{w, N, p})$$



Multi-Scale Spectral Loss

- input signal window size
- H hop size
- window function • w
- p compression function
- d distance function N set of window sizes
- set of compression function

Spectrum

$$\mathcal{Y}_{w,N,p}(m,k) = p\left(\left|\sum_{n=0}^{N-1} x[n+mH]w[n] \exp\left(\frac{-i2\pi kn}{N}\right)\right|\right)$$

MSS loss

$$\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w,N,p}, \hat{\mathcal{Y}}_{w,N,p})$$

© AudioLabs, 2025

AUDIO LABS

Multi-Scale Spectral Loss

- input signal
- window size
- *H* hop size window function

• w

- p compression function
- d distance function
- N set of window sizes
- P set of compression function

	$\left(N - $

 $\mathcal{Y}_{w,N,p}(m,k) = p\left(\left|\sum_{n=0}^{N-1} x[n+mH]w[n]\exp\left(\frac{-i2\pi kn}{N}\right)\right|\right)$

 $\mathcal{L}_{\mathrm{MSS}}(x, \hat{x}) := \sum_{N \in \mathcal{N}} \sum_{p \in \mathcal{P}} d(\mathcal{Y}_{w,N,p}, \hat{\mathcal{Y}}_{w,N,p})$ MSS loss

© AudioLabs, 2025

Spectrum

MSS loss with

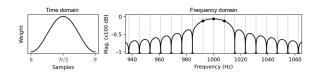
standard settings:

(WH, S4, C4, D1)

AUDIO

Spectrum-Based Distance

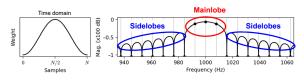
Hann window



• Input signal: Sinusoid with frequency f = 1000 Hz

AUDIO

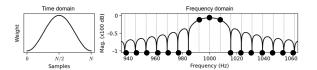
Spectrum-Based Distance



- Input signal: Sinusoid with frequency f = 1000 Hz
- $\, \blacksquare \,$ STFT \rightarrow Spectral leakage due to windowing

Spectrum-Based Distance

Hann window

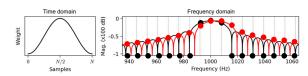


- Input signal: Sinusoid with frequency f = 1000 Hz
- $\, \bullet \,$ STFT \rightarrow Spectral leakage due to windowing
- Discrete STFT \rightarrow Frequency grid

AUDIO

Spectrum-Based Distance

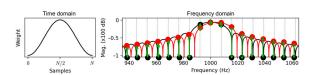
Hann window



- Input signal: Sinusoid with frequency f = 1000 Hz
- $\, \blacksquare \,$ STFT \rightarrow Spectral leakage due to windowing
- Discrete STFT \rightarrow Frequency grid
- Second signal: Sinusoid with frequency f = 1003.9 Hz

Spectrum-Based Distance

Hann window



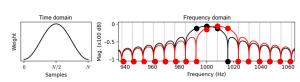
- Input signal: Sinusoid with frequency $f=1000~{
 m Hz}$
- STFT → Spectral leakage due to windowing
- Discrete STFT → Frequency grid
- Second signal: Sinusoid with frequency f = 1003.9 Hz

Distance depends on

- Grid sampling
- Mainlobe & sidelobes
- Window type
- STFT parameters

Spectrum-Based Distance

Hann window



- Input signal: Sinusoid with frequency $f=1000~{
 m Hz}$
- STFT → Spectral leakage due to windowing
- Discrete STFT \rightarrow Frequency grid
- Second signal: Sinusoid with frequency f = 1007.8 Hz

Distance depends on Grid sampling

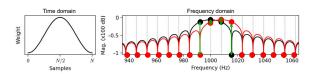
- Mainlobe & sidelobes

- Window type STFT parameters

AUDIO LABS

Spectrum-Based Distance

Hann window



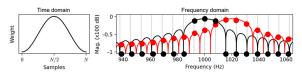
- Input signal: Sinusoid with frequency f = 1000 Hz
- STFT \rightarrow Spectral leakage due to windowing
- Discrete STFT → Frequency grid
- Second signal: Sinusoid with frequency f = 1007.8 Hz

Distance depends on

- Grid samplingMainlobe & sidelobes
- Window type
- STFT parameters

AUDIO

Spectrum-Based Distance



- Input signal: Sinusoid with frequency f = 1000 Hz
- $\, \blacksquare \,$ STFT \rightarrow Spectral leakage due to windowing
- Discrete STFT → Frequency grid
- Second signal: Sinusoid with frequency f = 1020 Hz

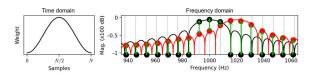
Distance depends on

- Grid sampling
 Mainlobe & sidelobes
- Window type
- STFT parameters



Spectrum-Based Distance

Hann window



- Input signal: Sinusoid with frequency f = 1000 Hz
- $\, \bullet \,$ STFT \rightarrow Spectral leakage due to windowing
- $\blacksquare \ \, \mathsf{Discrete} \ \mathsf{STFT} \to \textbf{Frequency} \ \textbf{grid}$
- Second signal: Sinusoid with frequency f = 1020 Hz

Distance depends on

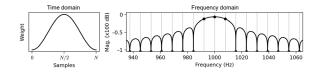
- Grid sampling Mainlobe & sidelobes
- Window typeSTFT parameters

Loss Functions Matter



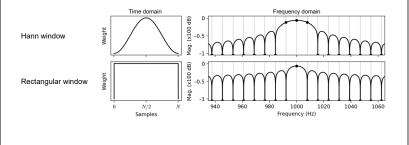
Dependency: Window Type

Hann window



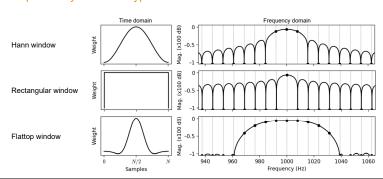


Dependency: Window Type

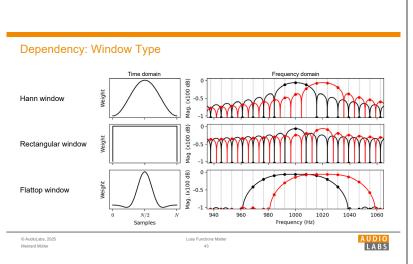


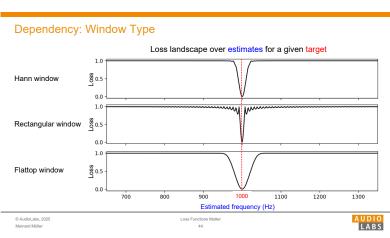
AUDIO

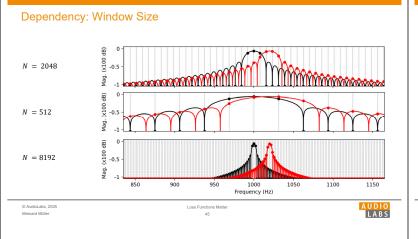
Dependency: Window Type

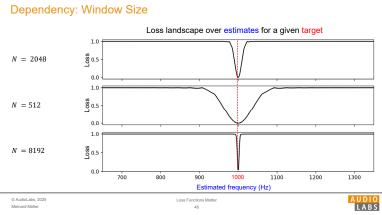


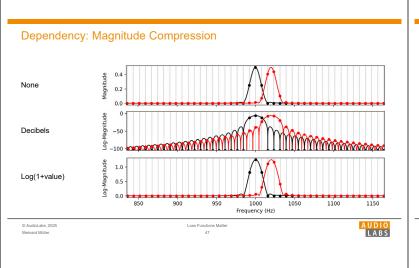
AUDIO LABS

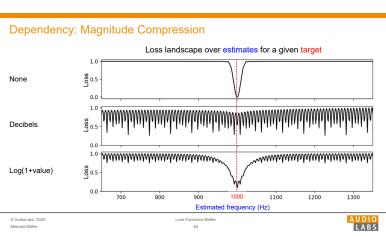


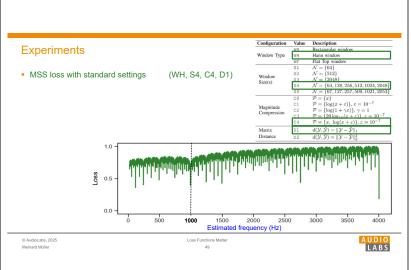




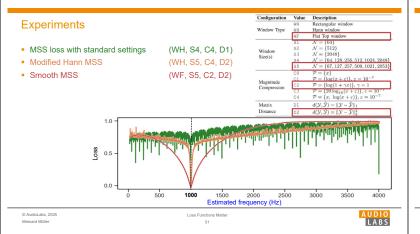








Experiments | Configuration | Value | Description | Window Type | Size | Rectangular staindow | Window Type | Window |

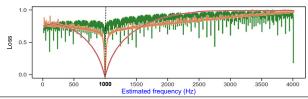


Experiments

GRA (Gradient-Sign Ranking Accuracy)

- Measures how often the loss gradient points in the correct direction.
- Step size distinguishes local gradient behavior from global trend.

Configuration	GRA			
Step Size	0.3 ct.	3 ct.	30 ct.	300 ct.
Standard MSS	0.523	0.529	0.573	0.775
Modified Hann MSS	0.613	0.635	0.708	0.923
Smooth MSS	0.999	0.993	0.952	0.860



AudioLabs, 2025 Loss Functions Matternation Matternation

AUDIO LABS

Overview

- Multi-Scale Spectral Loss
 Knowledge Source: Signal Representations
- Hierarchical Classification Loss
 Knowledge Source: Musical Hierarchies
- Differentiable Alignment Loss
 Knowledge Source: Temporal Coherence







use

© AudioLabs, 2025

Wagner Ring Dataset

Tetralogy (four operas)

Opera Das Rheingold Die Walküre Siegfried Götterdämmerung

Literature

Silla, Freitas: A survey of hierarchical classification across different application domains. Data Mining and Knowledge Discovery, 22(1-29: 31-72, 2011.
Wehrmann, Cerri, Barros: Hierarchical multi-label classification networks. Proc. ICML, 2018.

Krause, Müller: Hierarchical Classification for Singing Activity, Gender, and Type in Complex Music Recordings. Proc. ICASSP, 2022.

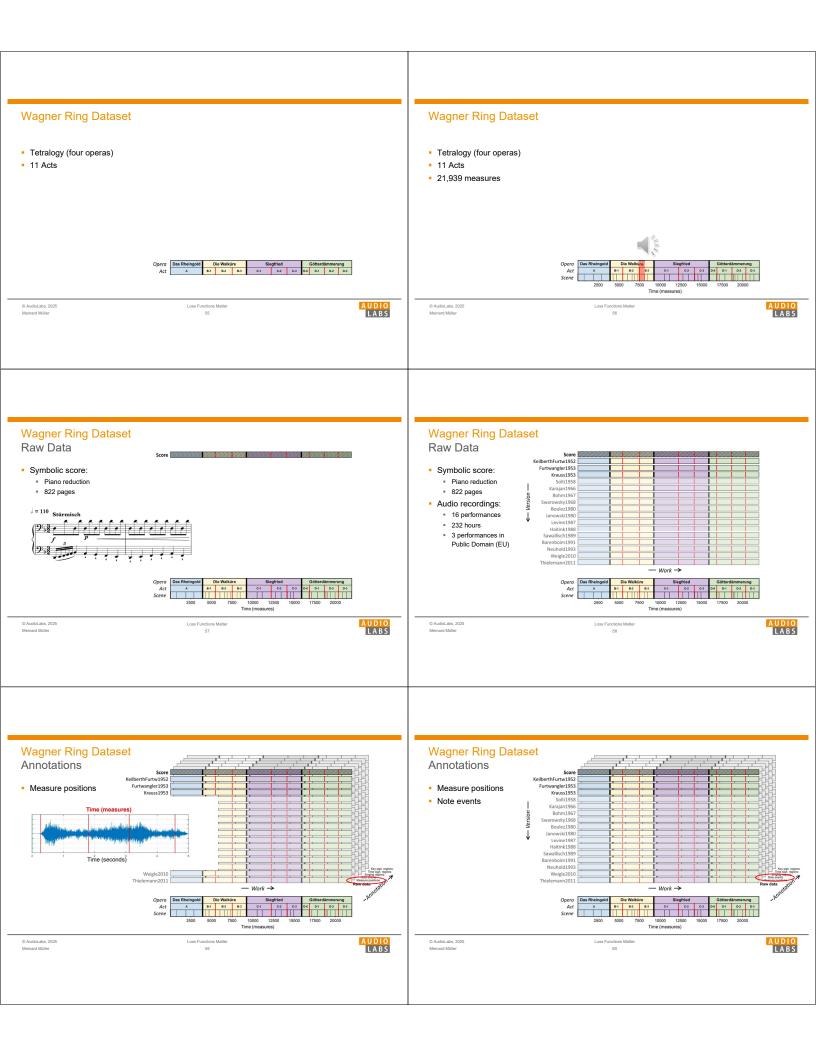
Krause, Müller: Hierarchical Classification for Instrument Activity Detection in Orchestral Music Recordings. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 31: 2567–2578, 2023.

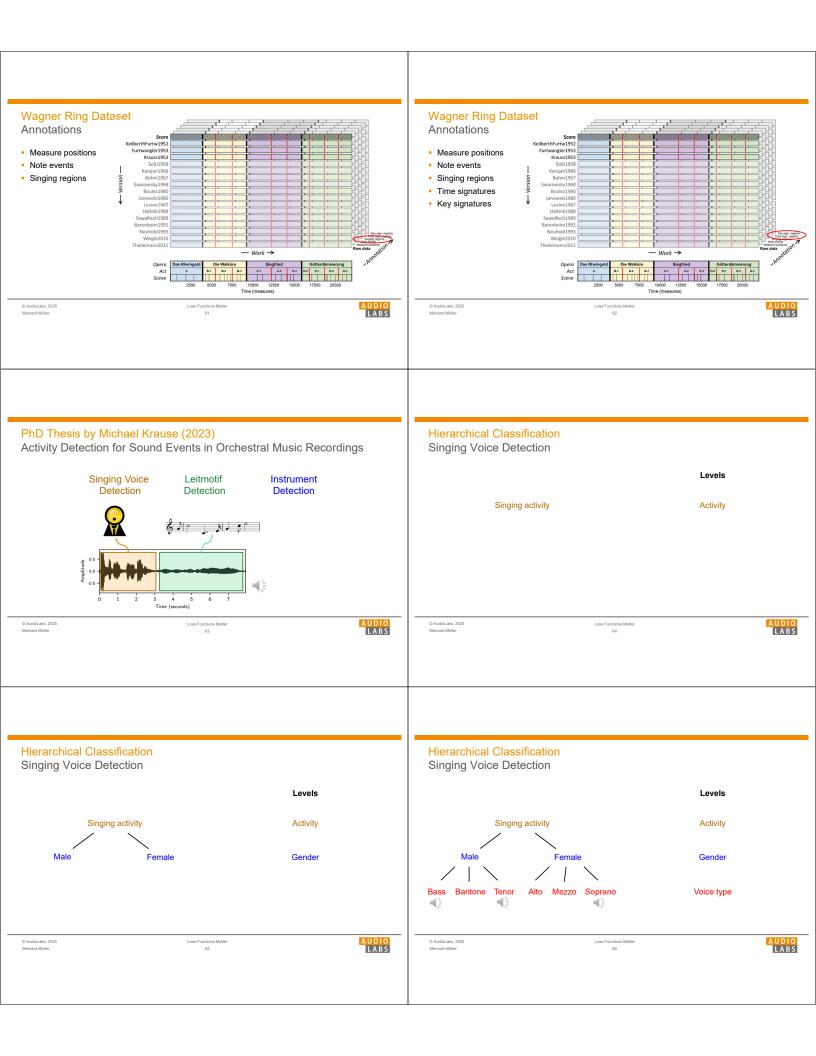
Weiß, Arifi-Müller, Krause, Zalkow, Klauk, Kleinertz, Müller: Wagner Ring Dataset: A Complex Opera Scenario for Music Processing and Computational Musicology. Transaction of the International Society for Music Information Retrieval (TISMIR), 6(1): 135–149, 2023.

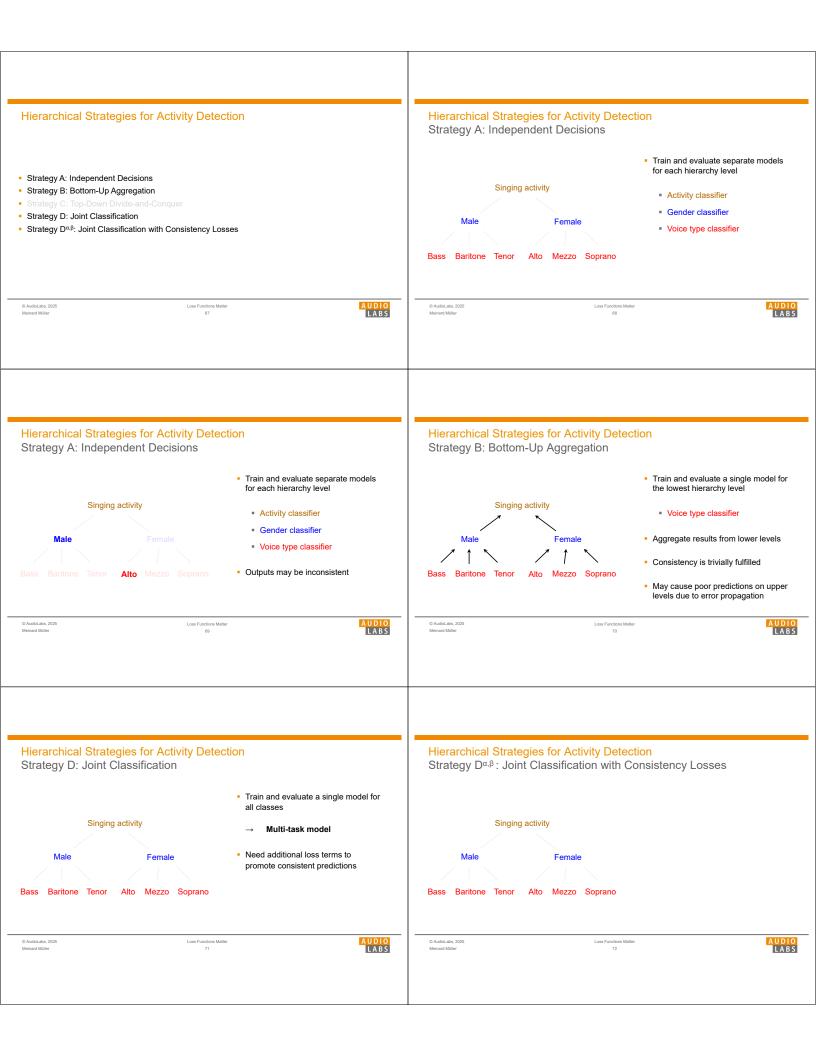


Loss Functions Matter 54









Hierarchical Strategies for Activity Detection

Strategy $D^{\alpha,\beta}$: Joint Classification with Consistency Losses



- Notation:
 - c: a class
 - p_c: probability of c

Hierarchical Strategies for Activity Detection

Strategy $D^{\alpha,\beta}$: Joint Classification with Consistency Losses



- Notation:
 - c: a class
 - p_c: probability of c
 - c↓: child classes of c

AUDIO

AUDIO

Hierarchical Strategies for Activity Detection

Strategy $D^{\alpha,\beta}$: Joint Classification with Consistency Losses



- Notation:
 - c: a class
 - p_c: probability of c
 - c↓: child classes of c
- For **bottom-up** consistency, minimize

 $\sum_{c' \in c \downarrow} \max\{0, p_{c'} - p_c\}^2$

 p_c should be at least as high as any $p_{c'}$

 \rightarrow penalty for every $p_{c'} > p_c$

Loss Functions Matter

Hierarchical Strategies for Activity Detection

Strategy $D^{\alpha,\beta}$: Joint Classification with Consistency Losses



- Notation:
 - c: a class
 - p_c: probability of c
 - c↓: child classes of c
- For top-down consistency, minimize

$$\max\{0, p_c - \max_{c' \in c\downarrow} p_{c'}\}^2$$

 p_c should not be above largest $p_{c'}$

Hierarchical Strategies for Activity Detection

Strategy D^{α,β}: Joint Classification with Consistency Losses

Bottom-up loss term:

$$\mathcal{L}_{\uparrow} = \frac{1}{|\mathbf{C} \setminus \mathbf{C}^H|} \sum_{h=2}^H \sum_{c \in \mathbf{C}^h} \sum_{c' \in c, \downarrow} \max\{0, \rho_{c'} - \rho_c\}^2$$

Top-down loss term:

$$\mathcal{L}_{\downarrow} = \frac{1}{|\mathbf{C} \setminus \mathbf{C}^1|} \sum_{h=2}^{H} \sum_{c \in \mathbf{C}^h} \max\{0, p_c - \max_{c' \in c_{\downarrow}} p_{c'}\}^2$$

Joint loss term:

$$\mathcal{L} = \mathcal{L}_{\mathsf{BCE}} + \alpha \mathcal{L}_{\downarrow} + \beta \mathcal{L}_{\uparrow}$$

Loss Functions Matter 77

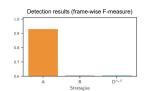
Notation C All classes

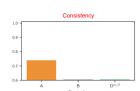
- Ch Classes at level h H Number of levels
- C ↓ Children of c

AUDIO

Pc Probability for c

Results: Female Singing





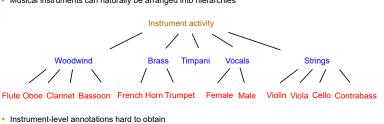
Consistency $\mathcal{I}_c^{\mathrm{Est}}$ Frames predicted as c $\mathcal{I}_{c\downarrow}^{\mathrm{Est}}$ Frames predicted as child of c $\gamma_c = \frac{|\mathcal{I}_c^{\mathrm{Est}} \cap \mathcal{I}_{c\downarrow}^{\mathrm{Est}}|}{|\mathcal{I}_c^{\mathrm{Est}} \cup \mathcal{I}_{c\downarrow}^{\mathrm{Est}}|}$

Strategy A (Independent Decisions) yields good but inconsistent results

© AudioLabs, 2025

© AudioLabs, 2025

Results: Female Singing Results: Female Singing Consistency $\mathcal{I}_c^{\mathrm{Est}}$ Frames predicted as cDetection results (frame-wise F-measure) Consistency Detection results (frame-wise F-measure) $\mathcal{I}_{c\downarrow}^{\mathrm{Est}}$ Frames predicted as child of c $|\mathcal{I}_c^{\mathrm{Est}} \cap \mathcal{I}_{c\perp}^{\mathrm{Est}}|$ $|\mathcal{I}_{c}^{\mathrm{Est}} \cup \mathcal{I}_{c}^{\mathrm{Est}}|$ Strategy A (Independent Decisions) yields good but inconsistent results Strategy A (Independent Decisions) yields good but inconsistent results Strategy B (Bottom-Up Aggregation) gives worse but consistent results Strategy B (Bottom-Up Aggregation) gives worse but consistent results - Strategy $D^{\alpha,\beta}$ (Joint with Consistency Losses) provides good trade-off AUDIO Scenario: Hierarchical Instrument Classification Overview Multi-Scale Spectral Loss Musical instruments can naturally be arranged into hierarchies



Loss Functions Matter 81

Instrument-level annotations hard to obtain



- Hierarchical Classification Loss
- Differentiable Alignment Loss Knowledge Source: Temporal Coherence







Consistency

 $\mathcal{I}_c^{\mathrm{Est}}$ Frames predicted as c

 $\mathcal{I}_{c\downarrow}^{\mathrm{Est}}$ Frames predicted as child of c

 $|\mathcal{I}_{c}^{\mathrm{Est}} \cap \mathcal{I}_{c}^{\mathrm{Est}}|$ $\gamma_c = \frac{|\mathcal{I}_c| \cdot \cdot \cdot \mathcal{I}_{c\downarrow}|}{|\mathcal{I}_c^{\text{Est}} \cup \mathcal{I}_{c\downarrow}^{\text{Est}}|}$

- Literature

 Cuturi, Blondei: Soft-DTW: A Differentiable Loss Function for Time-Series. ICML, 2017.

 Blondel, Mensch, Vert: Differentiable Divergences Between Time Series. AISTATS, 2021.

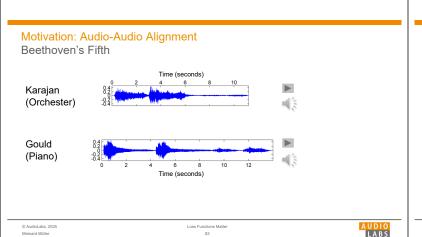
 Krause, WellS, Muller: Soft Dynamic Time Warping For Multi Pitch Estimation And Beyond. Proc. ICASSP, 2023.

 Zeitler, Deniffel, Krause, Müller: Stabilizing Training with Soft Dynamic Time Warping: A Case Study for Pitch Class Estimation with Weakly Aligned Targets. Proc. ISMR, 2023.

 Zeitler, Krause, Müller: Soft Dynamic Time Warping with Variable Step Weights. Proc. ICASSP, 2024.

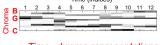
© AudioLabs, 2025





Motivation: Audio-Audio Alignment Beethoven's Fifth

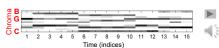
Karajan (Orchester)



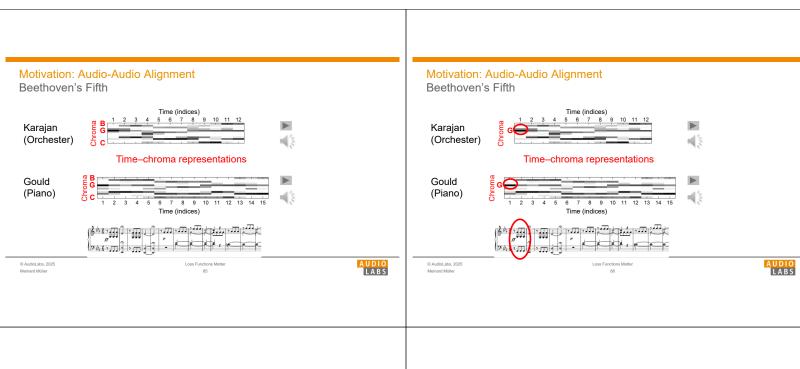


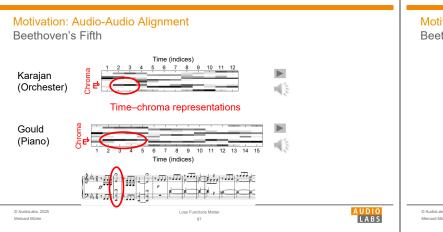
Time-chroma representations

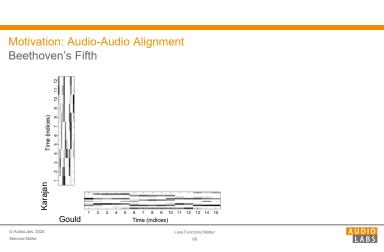
Gould (Piano)

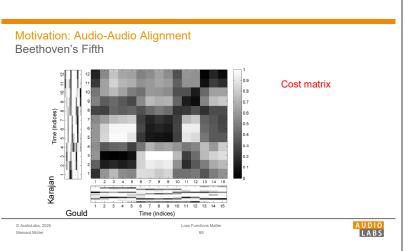


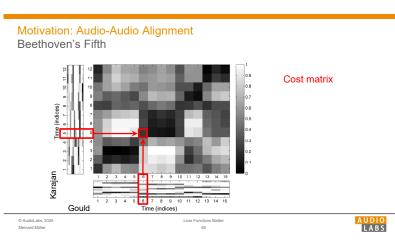
© AudioLabs, 2025





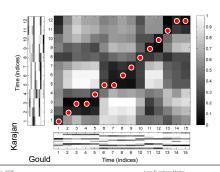






Motivation: Audio-Audio Alignment

Beethoven's Fifth



Cost-minimizing

warping path

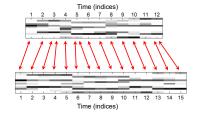
AUDIO

Motivation: Audio-Audio Alignment

Beethoven's Fifth



Gould (Piano)



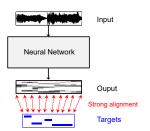
Cost-minimizing warping path

→ Strong alignment

© AudioLabs, 2025



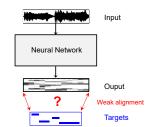
Feature Learning



- Task: Learn audio features using a neural network
- Loss: Binary cross-entropy
 - framewise loss
 - requires strongly aligned targets
 - hard to obtain



Feature Learning



- Task: Learn audio features using a neural network
- Loss: Binary cross-entropy
 - framewise loss
 - requires strongly aligned targets
 - hard to obtain
- Alignment as part of loss function
 - requires only weakly aligned targets
 - needs to be differentiable
- Problem: DTW is not differentiable → Soft DTW

Loss Functions Matter



Dynamic Time Warping (DTW)

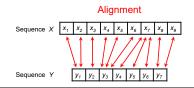
$$X := (x_1, x_2, \dots, x_N)$$

$$Y:=(y_1,y_2,\ldots,y_M)$$

$$x_n, y_m \in \mathcal{F}, n \in [1:N], m \in [1:M]$$

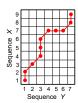
= Feature space

© AudioLabs, 2025



Alignment matrix $A \in \left\{0,1\right\}^{N \times M}$

Set of all possible alignment matrices $\mathcal{A}_{N,M} \subset \{0,1\}^{N \times M}$



AUDIO LABS

Dynamic Time Warping (DTW)

$$X := (x_1, x_2, \dots, x_N)$$

$$Y:=(y_1,y_2,\ldots,y_M)$$

$$x_n, y_m \in \mathcal{F}, n \in [1:N], m \in [1:M]$$

 \mathcal{F} = Feature space

Alignment matrix $A \in \left\{0,1\right\}^{N \times M}$

Set of all possible alignment matrices

 $\mathcal{A}_{N,M} \subset \{0,1\}^{N \times M}$

 $c: \mathcal{F} \times \mathcal{F} \to \mathbb{R}_{\geq 0}$ Cost measure:

 $C \in \mathbb{R}^{N \times M} \quad \text{with} \quad C(n,m) := c(x_n,y_m)$ Cost matrix:

Cost of alignment: $\langle A, C \rangle$

DTW cost: $DTW(C) = \min (\{ \langle A, C \rangle \mid A \in \mathcal{A}_{N,M} \})$ Optimal alignment: $A^* = \operatorname{argmin} (\{\langle A, C \rangle \mid A \in \mathcal{A}_{N,M}\})$

© AudioLabs, 2025 Meinard Müller

Dynamic Time Warping (DTW)

DTW cost:

$$\mathrm{DTW}(C) = \min \left(\left\{ \left\langle A, C \right\rangle \ | \ A \in \mathcal{A}_{N,M} \right\} \right)$$

Efficient computation via Bellman's recursion in O(NM)

$$D(n,m) = \min\{D(n-1,m), D(n,m-1), D(n,m)\} + C(n,m)$$

for n>1 and m>1 and suitable initialization.

$$DTW(C) = D(N, M)$$

- Problem: DTW(C) is not differentiable with regard to C
- Idea: Replace min-function by a smooth version

$$\min^{\gamma} (\mathcal{S}) = -\gamma \log \sum\nolimits_{s \in \mathcal{S}} \exp \left(-s/\gamma \right)$$

for set $\, \mathcal{S} \subset \mathbb{R} \,$ and temperature parameter $\, \gamma \in \mathbb{R} \,$



Soft-DTW
Cuturi, Blondel: Soft-DTW: A
Differentiable Loss Function
for Time-Series. ICML, 2017

Soft Dynamic Time Warping (SDTW)

SDTW cost: $SDTW^{\gamma}(C) = min^{\gamma} (\{\langle A, C \rangle \mid A \in \mathcal{A}_{N,M}\})$

Efficient computation via Bellman's recursion in O(NM) still works:

$$D^{\gamma}(n,m) = \min^{\gamma} \{ D^{\gamma}(n-1,m), D^{\gamma}(n,m-1), D^{\gamma}(n,m) \} + C(n,m)$$

for n>1 and m>1 and suitable initialization.

$$SDTW^{\gamma}(C) = D^{\gamma}(N, M)$$

- Limit case: $SDTW^{\gamma}(C) \xrightarrow{\gamma \to 0} DTW(C)$
- SDTW(C) is differentiable with regard to C
- Questions:
 - How does the gradient look like?
 - Can it be computed efficiently?
 - How does SDTW generalize the alignment concept?



Soft Dynamic Time Warping (SDTW)

SDTW cost: $\mathrm{SDTW}^{\gamma}(C) = \min^{\gamma} \left(\left\{ \langle A, C \rangle \mid A \in \mathcal{A}_{N,M} \right\} \right)$

• Define $p^{\gamma}(C)$ as the following "probability" distribution over $\mathcal{A}_{N,M}$:

$$p^{\gamma}(C)_{A} = \frac{\exp\left(-\langle A, C \rangle / \gamma\right)}{\sum_{A' \in \mathcal{A}_{N,M}} \exp\left(-\langle A', C \rangle / \gamma\right)} \qquad \text{ for } A \in \mathcal{A}_{N,M}$$

- The expected alignment with respect to $p^{\gamma}(C)$ is given by:

$$E^{\gamma}(C) = \sum\nolimits_{A \in \mathcal{A}_{N,M}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$$

The gradient is given by:

$$\nabla_C \mathrm{SDTW}^{\gamma}(C) = E^{\gamma}(C)$$

The gradient can be computed efficiently in O(NM) via a recursive algorithm.

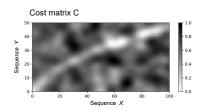
Loss Functions Matter

AUDIO LABS

Soft Dynamic Time Warping (SDTW)

Expected alignment : $E^{\gamma}(C) = \sum_{A \in \mathcal{A}_{N,M}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ

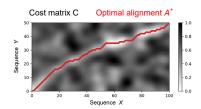


AUDIO

Soft Dynamic Time Warping (SDTW)

Expected alignment : $E^{\gamma}(C) = \sum_{A \in \mathcal{A}_{NM}} p^{\gamma}(C)_A A \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ

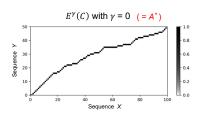


AUDIO LABS

Soft Dynamic Time Warping (SDTW)

Expected alignment : $E^{\gamma}(C) = \sum_{A \in \mathcal{A}_{NM}} p^{\gamma}(C)_A A \in \mathbb{R}^{N \times M}$

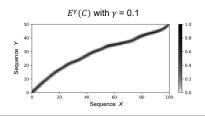
- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ



Soft Dynamic Time Warping (SDTW)

 $\text{Expected alignment}: \quad E^{\gamma}(C) = \sum\nolimits_{A \in \mathcal{A}_{N,M}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ

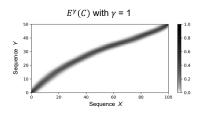


AUDIO

Soft Dynamic Time Warping (SDTW)

 $\text{Expected alignment}: \quad E^{\gamma}(C) = \sum\nolimits_{A \in \mathcal{A}_{N,M}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ

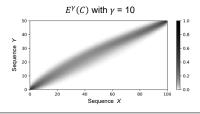


AUDIO LABS

Soft Dynamic Time Warping (SDTW)

 $\text{Expected alignment}: \quad E^{\gamma}(C) = \sum\nolimits_{A \in \mathcal{A}_{N,M}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ

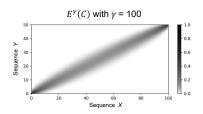


AUDIO

Soft Dynamic Time Warping (SDTW)

 $\text{Expected alignment}: \quad E^{\gamma}(C) = \sum\nolimits_{A \in \mathcal{A}_{NM}} p^{\gamma}(C)_A A \quad \in \mathbb{R}^{N \times M}$

- Can be interpreted as a smoothed version of an alignment
- Degree of smoothing depends on temperature parameter γ



AUDIO

Soft Dynamic Time Warping (SDTW)

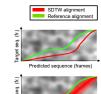
Conclusions

- Direct generalization of DTW (replacing min by smooth variant)
- Gradient is given by expected alignment
- Fast forward algorithm: O(NM)
- Fast gradient computation: O(NM)
- SDTW yields a (typically) poor lower bound for DTW
- Can be used as loss function to learn from weakly aligned sequences

AUDIO LABS

Soft Dynamic Time Warping (SDTW) Stabilizing Training

- Standard SDTW often unstable Unstable training in early stages
 - Degenerate output alignment
- Hyperparameter adjustment
 - · High temperature to smooth alignments
 - Temperature annealing
- Diagonal prior
- Modified step size condition





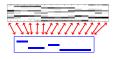
Soft Dynamic Time Warping (SDTW)

Representation Learning

- Symmetric application
 - Learn representation of both sequences
 - Needs a contrastive loss term
- Assymmetric application

 - Use fixed (e.g., binary) encoding of target
 Learn representation of only one sequences
 - No contrastive loss term need
- Simulation of CTC-loss using SDTW possible
- Many DTW variants also possible for SDTW





Conclusions

- Multi-Scale Spectral Loss Knowledge Source: Signal Representations
- Hierarchical Classification Loss Knowledge Source: Musical Hierarchies
- Differentiable Alignment Loss Knowledge Source: Temporal Coherence











Conclusions

- Multi-Scale Spectral Loss Knowledge Source: Signal Representations
- Hierarchical Classification Loss Knowledge Source: Musical Hierarchies
- Differentiable Alignment Loss Knowledge Source: Temporal Coherence











Müller, Zeitler: **2025 ISMIR Tutorial**Differentiable Alignment Techniques for Music
Processing: Techniques and Applications

Loss Functions Matter

