

Lecture
**Selected Topics in Deep Learning
for Audio, Speech, and Music Processing**

Connectionist Temporal Classification (CTC) Loss

Frank Zalkow and Meinard Müller

International Audio Laboratories Erlangen
frank.zalkow@audiolabs-erlangen.de

21.06.2021

Lecturers

Frank Zalkow

- Music Informatics and Musicology (B.A., University of Music Karlsruhe)
- Music Informatics (M.A., University of Music Karlsruhe)
- Ph.D. student in music information retrieval (FAU, supervisor: Meinard Müller)



Meinard Müller

- Mathematics (Diplom/Master, Bonn University)
- Computer Science (Ph.D., Bonn University)
- Information Retrieval (Habilitation, Bonn University)
- Combinatorics (Postdoc, Keio University, Japan)
- Senior Researcher (Max-Planck Institute, Saarland)
- Professor: Semantic Audio Processing (FAU)



Overview

1. Introduction
2. CTC Loss Computation
3. Applications
4. Outlook and Further Notes

Introduction

- Connectionist Temporal Classification
- Graves, Fernández, Gomez, and Schmidhuber. *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. ICML 2006. [[ACM](#)]
- “Temporal Classification”: Labelling un-segmented data sequences
- “Connectionist”: Refers to the use of deep learning

Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks

Alex Graves¹
Santiago Fernández¹
Faustino Gomez¹
Jürgen Schmidhuber^{1,2}

ALEX@IDSIA.CH
SANTIAGO@IDSIA.CH
TINO@IDSIA.CH
JUERGEN@IDSIA.CH

¹ Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA), Galleria 2, 6928 Manno-Lugano, Switzerland

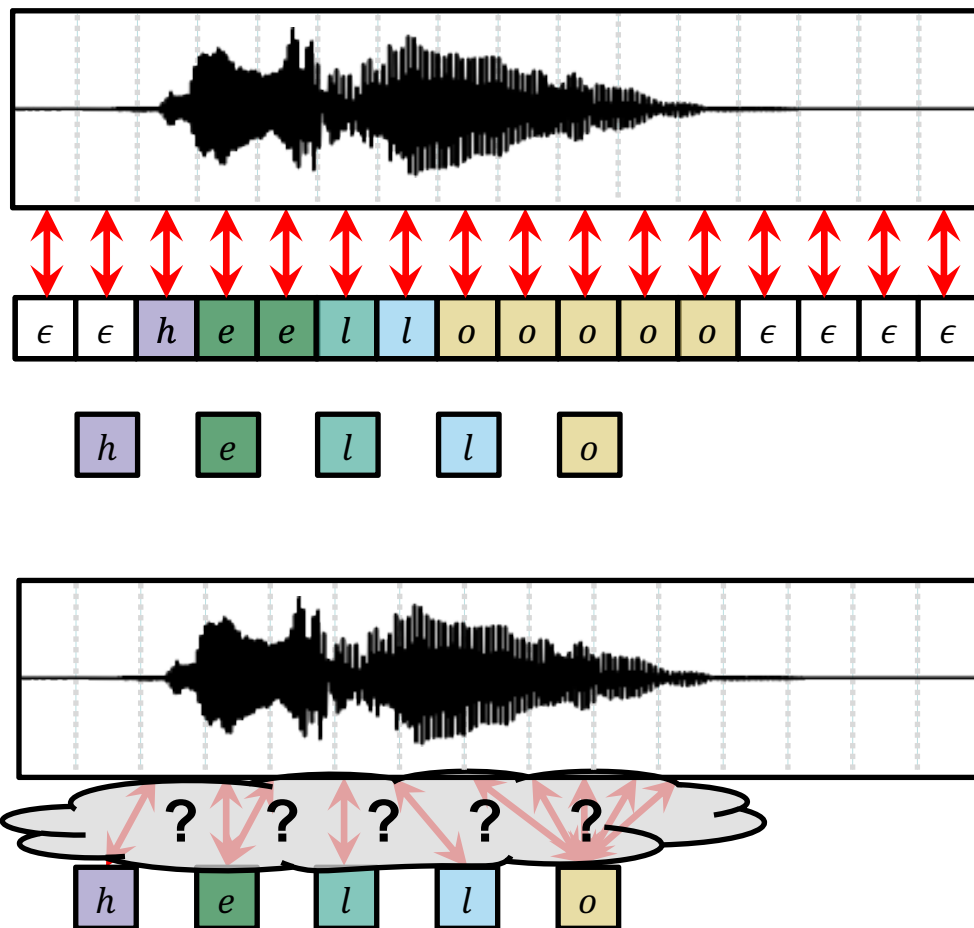
² Technische Universität München (TUM), Boltzmannstr. 3, 85748 Garching, Munich, Germany

Introduction



Training data in speech recognition

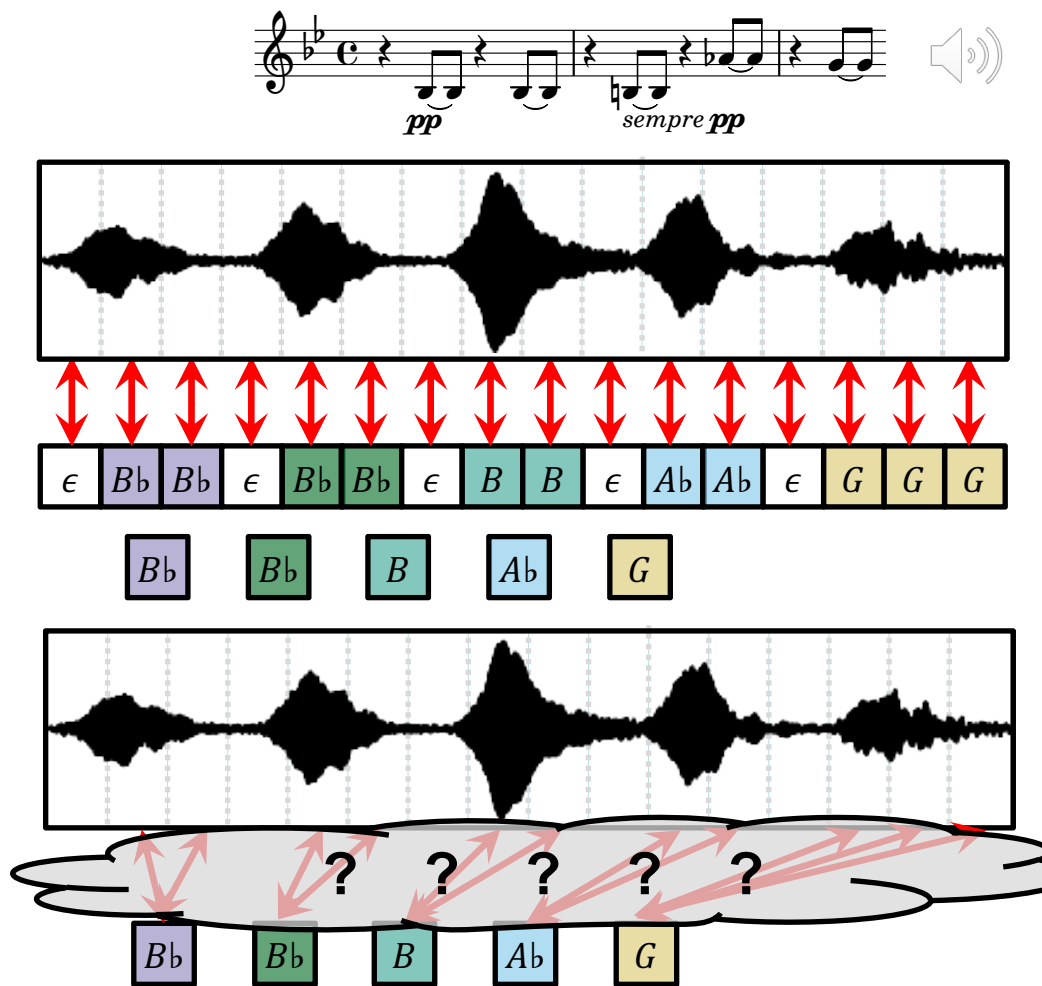
- Strongly aligned training data
 - Character annotations (labels) for each time step
 - Can be used for training in a standard classification setup
 - Tedious to annotate
- Weakly aligned training data
 - Globally corresponding character sequence without local alignment
 - Cannot be used for training in a standard classification setup
 - Easier to annotate
- Aim of CTC: Employing weakly aligned data for training
- Useful for many applications



Introduction

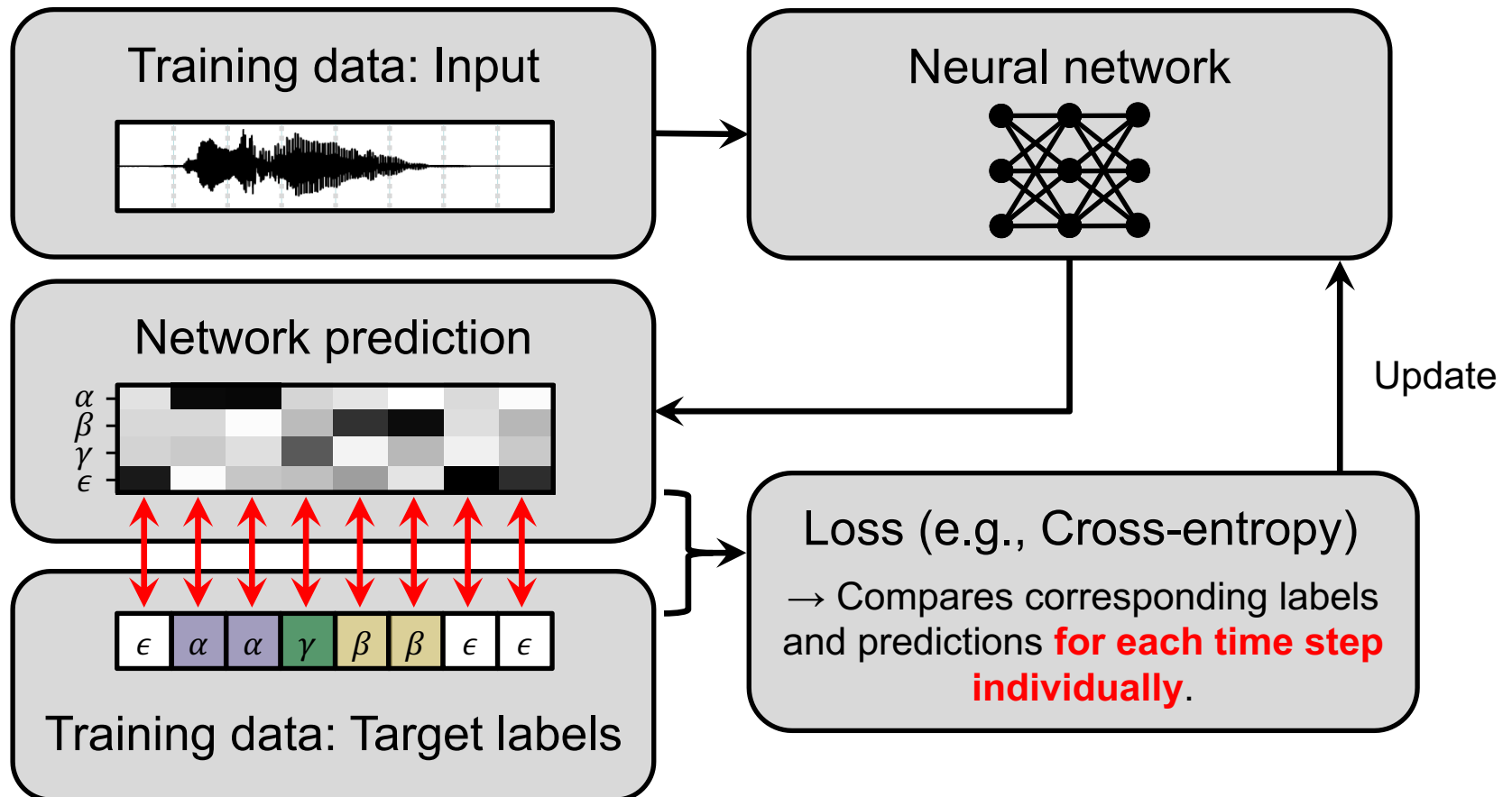
Training data in theme-based music retrieval

- Strongly aligned training data
 - Pitch/chroma annotations (labels) for each time step
 - Can be used for training in a standard classification setup
 - Tedious to annotate
- Weakly aligned training data
 - Globally corresponding pitch/chroma sequence without local alignment
 - Cannot be used for training in a standard classification setup
 - Easier to annotate
- Aim of CTC: Employing weakly aligned data for training
- Useful for many applications



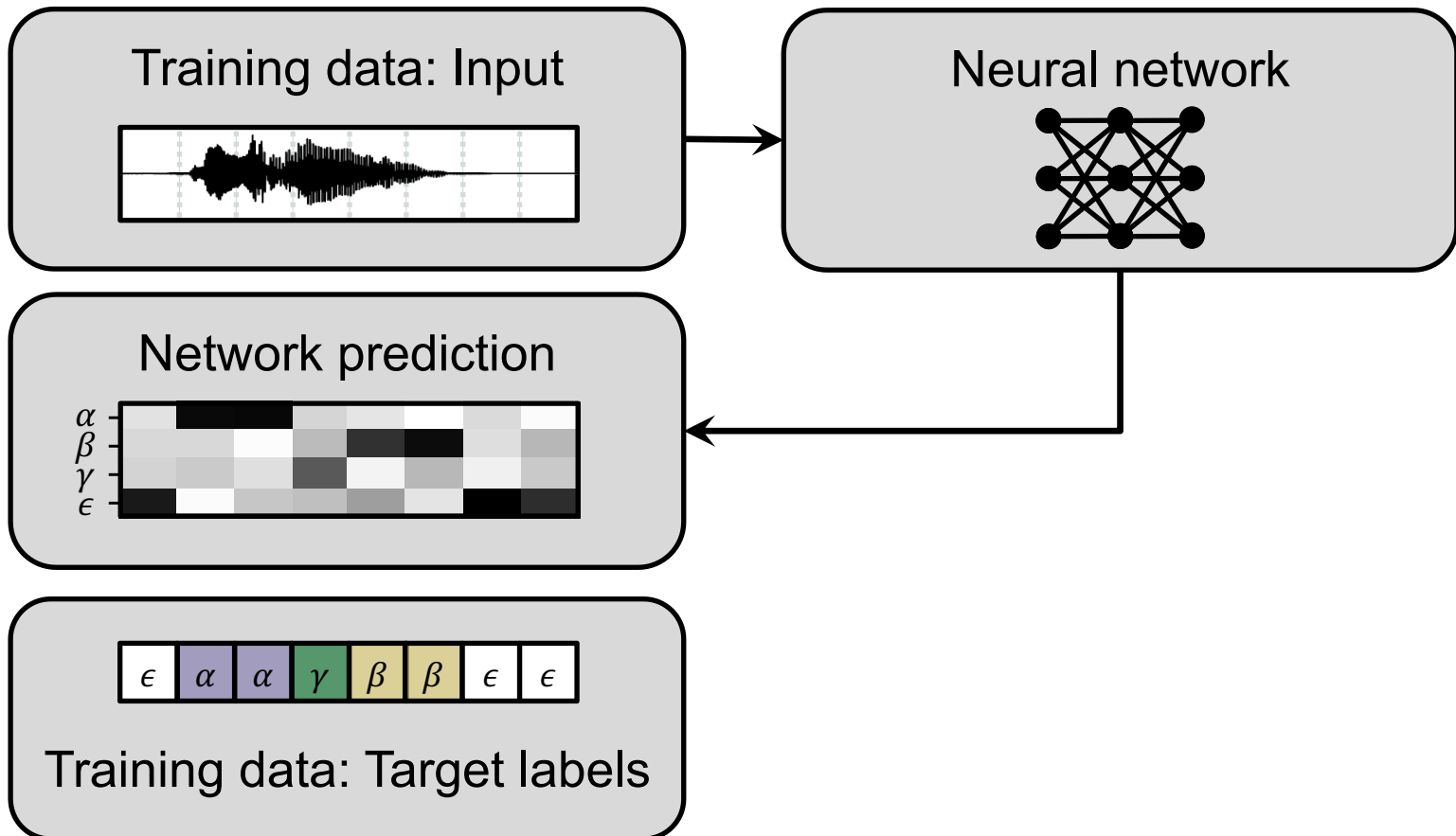
Introduction

Standard deep learning setup: Strongly aligned training data



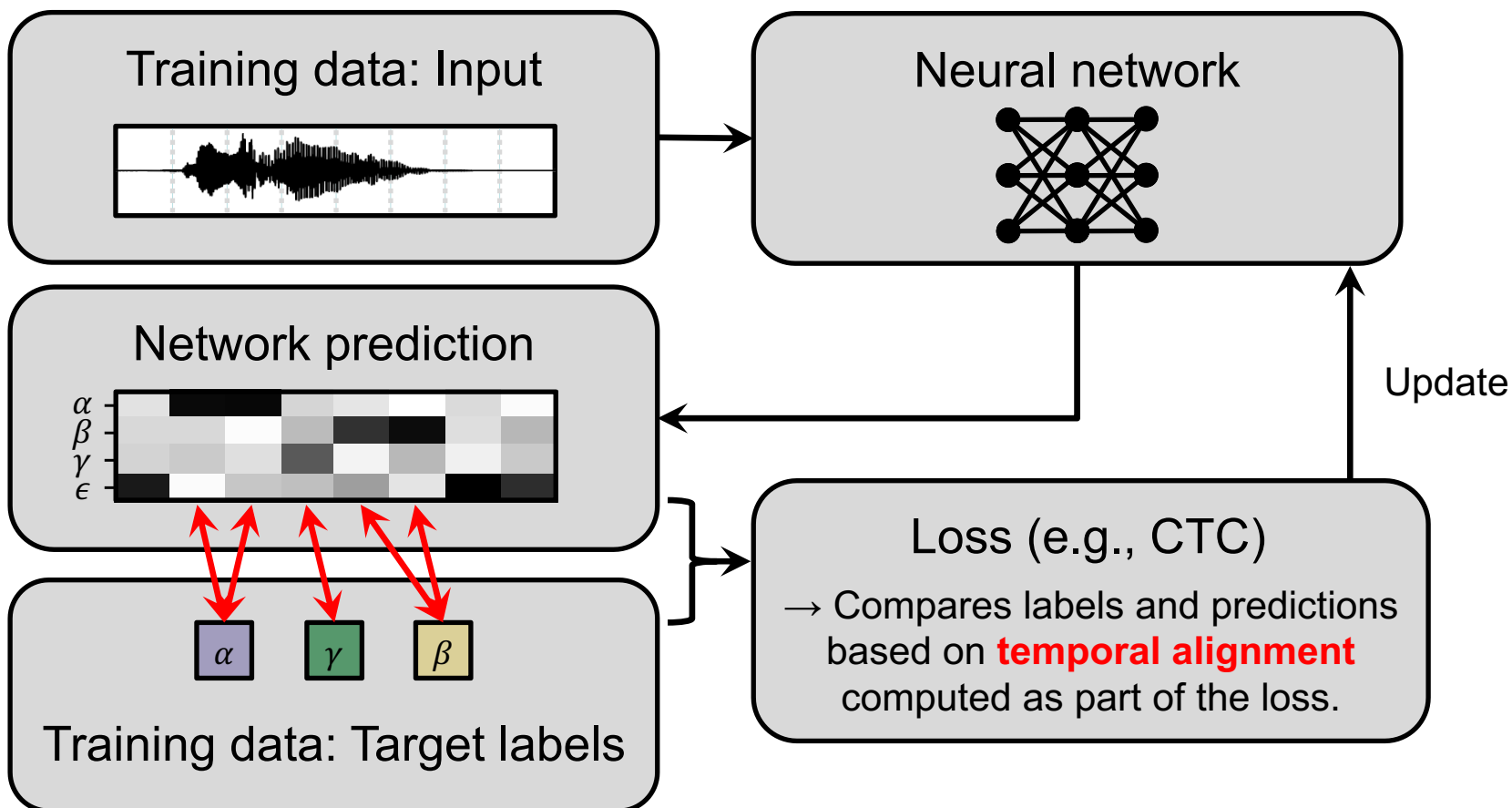
Introduction

Standard deep learning setup: Strongly aligned training data



Introduction

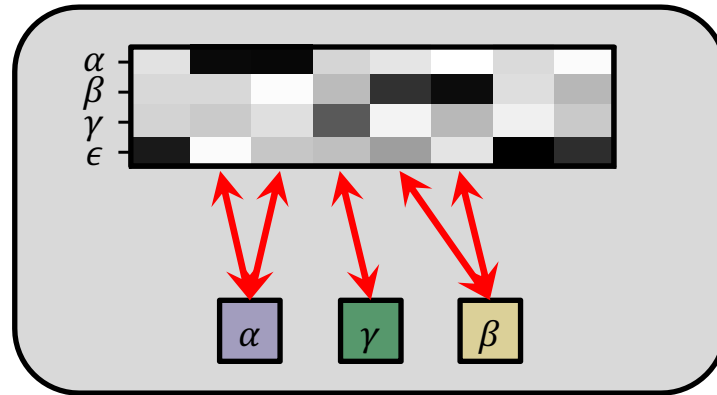
Non-standard deep learning setup: Weakly aligned training data



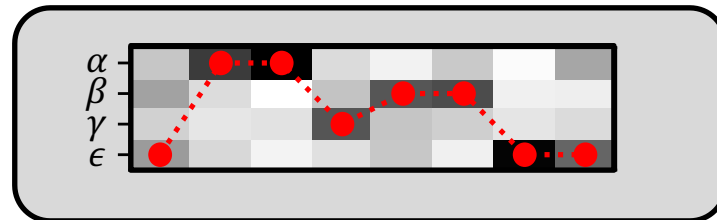
Introduction

Alignment Representations

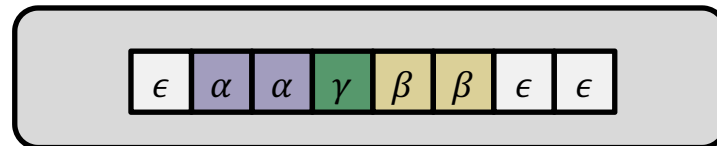
“Arrow” representation



“Point” representation



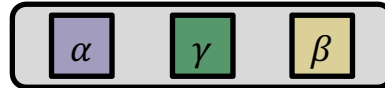
“Unfolded” representation



CTC Loss Computation: Intuition

- Alphabet $\mathbb{A} = \{\alpha, \beta, \gamma\}$

- Label sequence $Y = (\alpha, \gamma, \beta)$



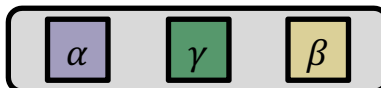
- Network output $f_{\theta}(\mathbf{X}) =$



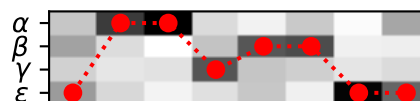
- Alignment A is “expansion” of Y to the temporal length of $f_{\theta}(\mathbf{X})$ (possibly consecutive duplicates and blank symbols ϵ)

CTC Loss Computation: Intuition

- Alphabet $\mathbb{A} = \{\alpha, \beta, \gamma\}$
- Label sequence $Y = (\alpha, \gamma, \beta)$



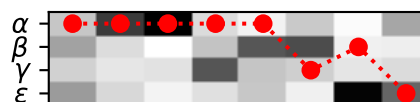
- Naive idea: “Hard” alignment
(Related: Viterbi decoding)



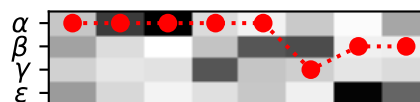
$$P(\epsilon, \alpha, \alpha, \gamma, \beta, \beta, \epsilon, \epsilon) \approx 0.015$$

- Not suitable for gradient-descent-based training (not differentiable)!

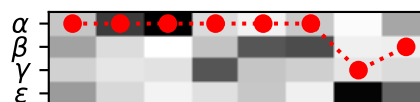
- Therefore: “Soft” alignment
(Related: Forward algorithm)



$$P(\alpha, \alpha, \alpha, \alpha, \alpha, \gamma, \beta) \approx 7.14 \cdot 10^{-7}$$

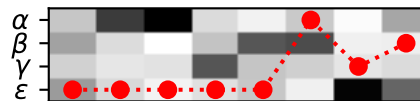


$$P(\alpha, \alpha, \alpha, \alpha, \alpha, \gamma, \beta, \beta) \approx 3.98 \cdot 10^{-7}$$



$$P(\alpha, \alpha, \alpha, \alpha, \alpha, \gamma, \beta, \epsilon) \approx 3.23 \cdot 10^{-6}$$

⋮



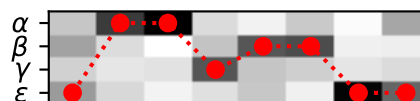
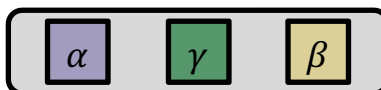
$$P(\epsilon, \epsilon, \epsilon, \epsilon, \epsilon, \alpha, \gamma, \beta) \approx 5.82 \cdot 10^{-6}$$

⋮

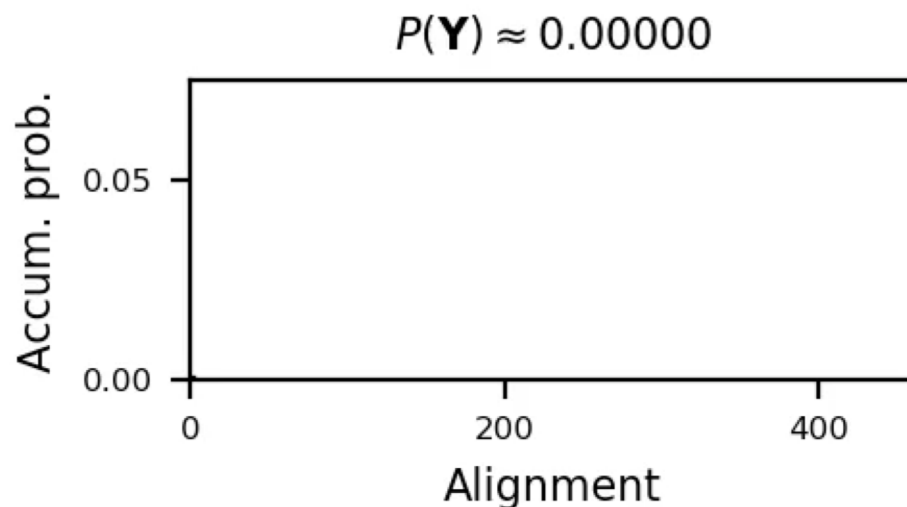
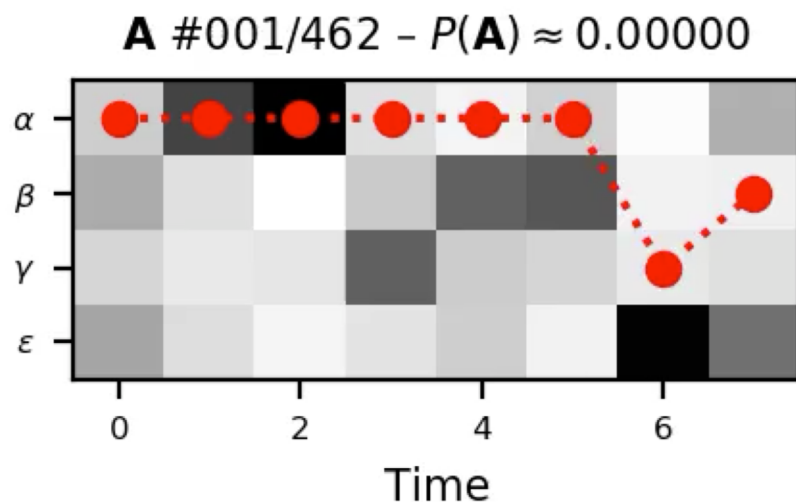
$$\sum_{\mathbf{A}} P(\mathbf{A}) \approx 0.069$$

CTC Loss Computation: Intuition

- Alphabet $\mathbb{A} = \{\alpha, \beta, \gamma\}$
- Label sequence $\mathbf{Y} = (\alpha, \gamma, \beta)$
- Naive idea: “Hard” alignment
(Related: Viterbi decoding)
- Not suitable for gradient-descent-based training (not differentiable)!
- Therefore: “Soft” alignment
(Related: Forward algorithm)



$$P(\epsilon, \alpha, \alpha, \gamma, \beta, \beta, \epsilon, \epsilon) \approx 0.015$$



CTC Loss Computation: Formal Description

- Input feature sequence (length N , elements $\mathbf{x}_n \in \mathbb{R}^D$, frame n , dimensionality $D \in \mathbb{N}$)

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

Examples



CTC Loss Computation: Formal Description

- Input feature sequence (length N , elements $\mathbf{x}_n \in \mathbb{R}^D$, frame n , dimensionality $D \in \mathbb{N}$)

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

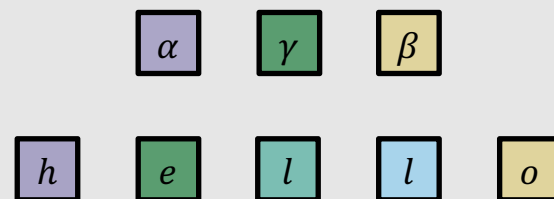
$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

$$\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$$

- Label sequence (length M , elements $\mathbf{y}_m \in \mathbb{A}$, index m , alphabet \mathbb{A})

$$\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$$

Examples



CTC Loss Computation: Formal Description

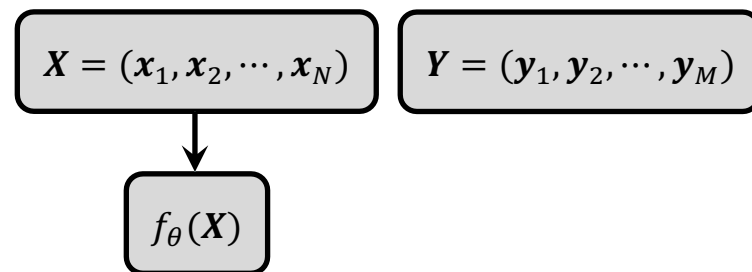
- Input feature sequence (length N , elements $\mathbf{x}_n \in \mathbb{R}^D$, frame n , dimensionality $D \in \mathbb{N}$)

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

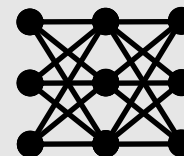
- Label sequence (length M , elements $\mathbf{y}_m \in \mathbb{A}$, index m , alphabet \mathbb{A})

$$\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$$

- Neural network f_θ



Examples



CTC Loss Computation: Formal Description

- Input feature sequence (length N , elements $\mathbf{x}_n \in \mathbb{R}^D$, frame n , dimensionality $D \in \mathbb{N}$)

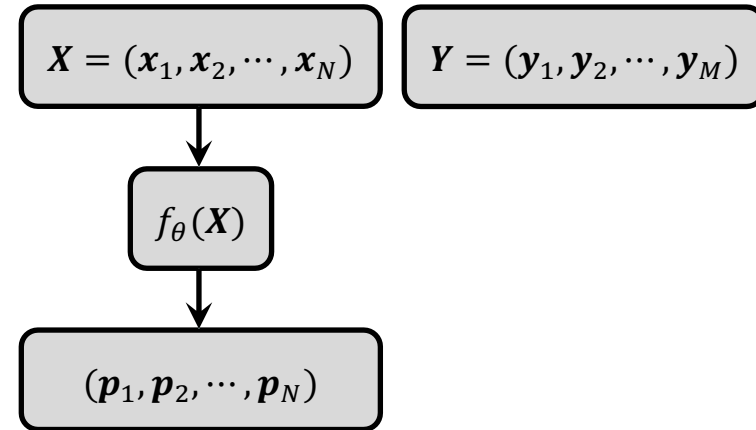
$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

- Label sequence (length M , elements $\mathbf{y}_m \in \mathbb{A}$, index m , alphabet \mathbb{A})

$$\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$$

- Neural network f_θ
- Network output: probability distribution $\mathbf{p}_n: \mathbb{A}' \rightarrow [0,1]$ ($\mathbb{A}' = \mathbb{A} \cup \{\epsilon\}$)

$$f_\theta(\mathbf{X}) = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$$



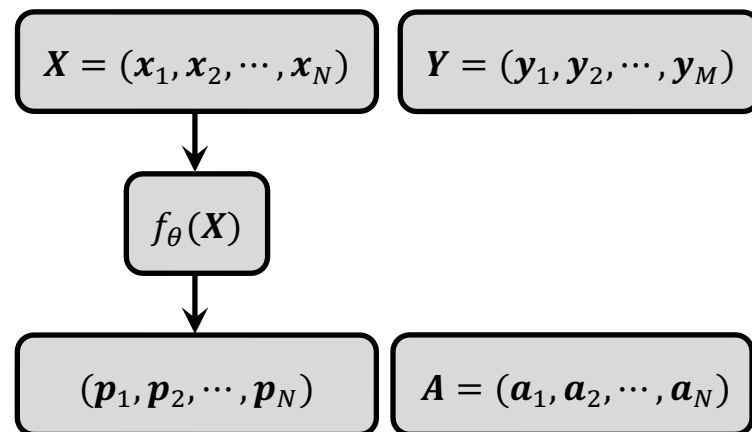
Examples



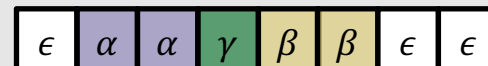
CTC Loss Computation: Formal Description

- Alignment of label sequence Y to feature sequence X yields sequence (length N , elements $\mathbf{a}_n \in \mathbb{A}'$)

$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N)$$



Examples



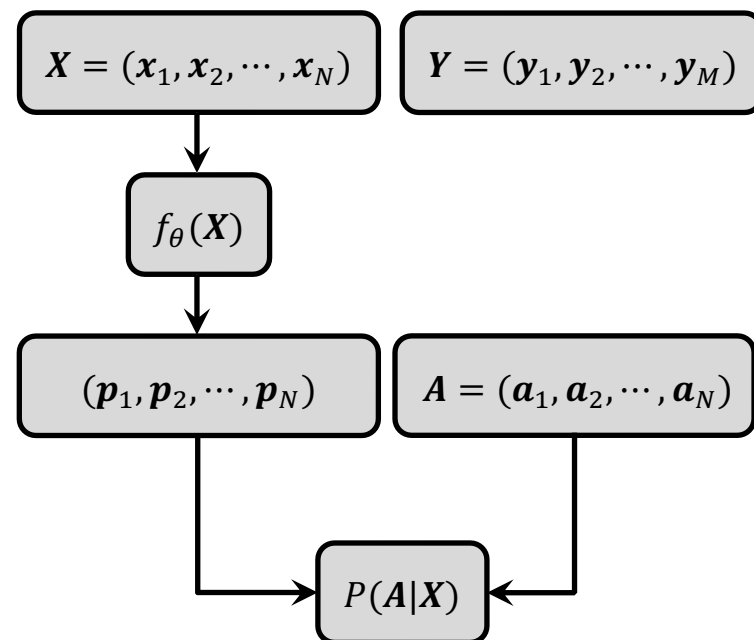
CTC Loss Computation: Formal Description

- Alignment of label sequence Y to feature sequence X yields sequence (length N , elements $\mathbf{a}_n \in \mathbb{A}'$)

$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N)$$

- Probability of alignment

$$P(\mathbf{A}|\mathbf{X}) = \prod_{n=1}^N p_n(\mathbf{a}_n)$$



Examples

$$P(\epsilon \alpha \alpha \gamma \beta \beta \epsilon \epsilon \mid \text{audio waveform}) \approx 0.043$$

CTC Loss Computation: Formal Description

- Alignment of label sequence Y to feature sequence X yields sequence (length N , elements $\mathbf{a}_n \in \mathbb{A}'$)

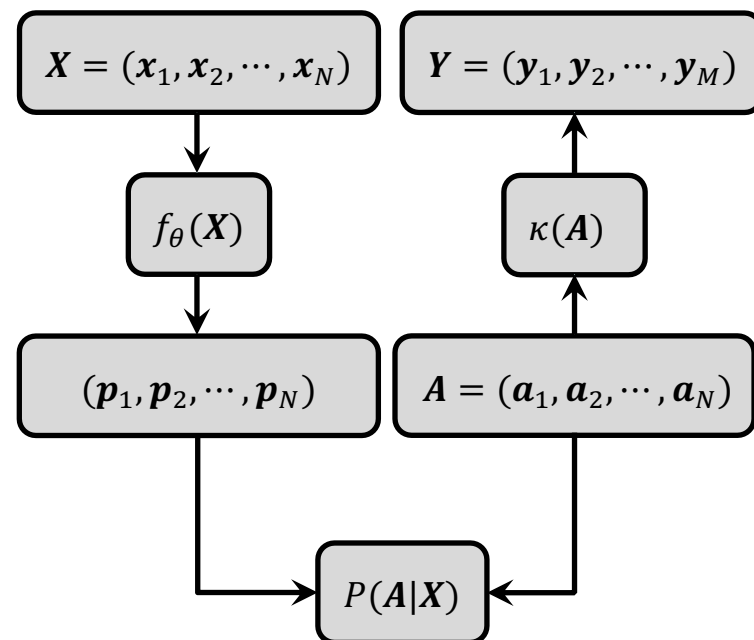
$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N)$$

- Probability of alignment

$$P(\mathbf{A}|\mathbf{X}) = \prod_{n=1}^N p_n(\mathbf{a}_n)$$

- Alignment \mathbf{A} can be reduced to label sequence Y by composition of two functions:

- Removing consecutive duplicates κ'
- Removing blank symbols κ''
- $\kappa = \kappa'' \circ \kappa'$
- $\kappa(\mathbf{A}) = \kappa''(\kappa'(\mathbf{A})) = Y$



Examples

$$\mathbf{A} = (\alpha, \alpha, \beta, \beta, \beta)$$

$$\kappa'(\mathbf{A}) = \kappa(\mathbf{A}) = (\alpha, \beta)$$

$$\mathbf{A} = (\alpha, \epsilon, \alpha, \beta, \epsilon, \beta, \beta)$$

$$\kappa'(\mathbf{A}) = (\alpha, \epsilon, \alpha, \beta, \epsilon, \beta)$$

$$\kappa''(\kappa'(\mathbf{A})) = \kappa(\mathbf{A}) = (\alpha, \alpha, \beta, \beta)$$

CTC Loss Computation: Formal Description

- Many different alignments \mathbf{A} map to the same label sequence \mathbf{Y} (“valid alignments”)

- Set of all valid alignments

$$\mathbb{K}_{\mathbf{X}, \mathbf{Y}} = \{\mathbf{A} \in (\mathbb{A}')^N : \kappa(\mathbf{A}) = \mathbf{Y}\}$$

- Probability of label sequence

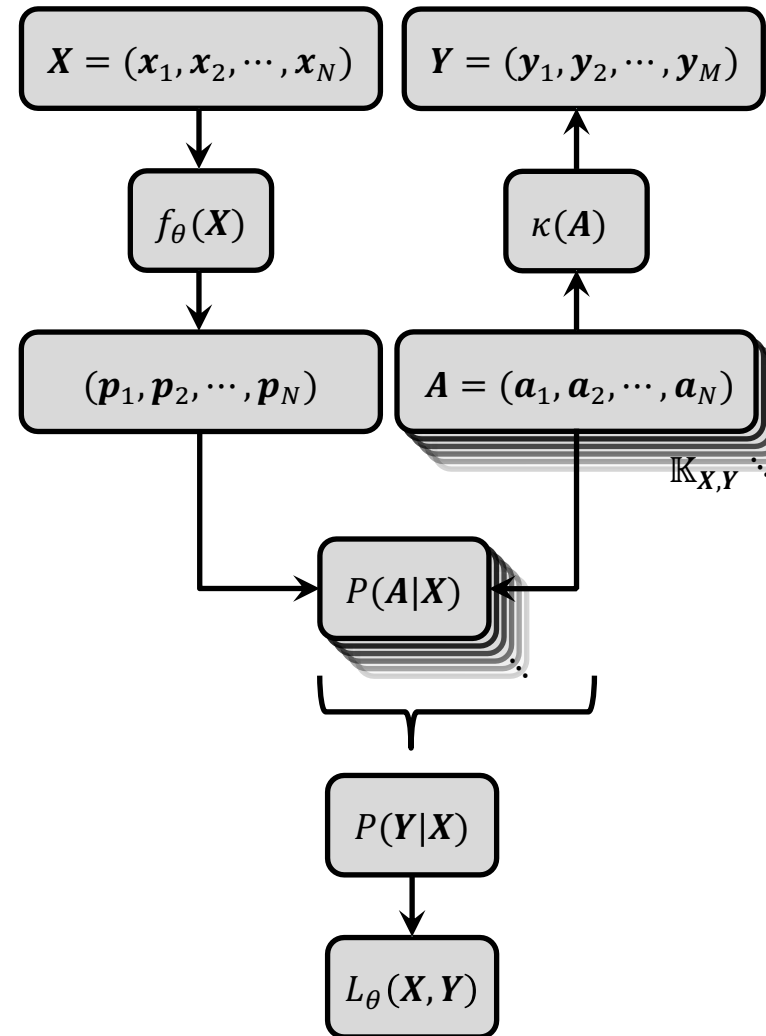
$$P(\mathbf{Y}|\mathbf{X}) = \sum_{\mathbf{A} \in \mathbb{K}_{\mathbf{X}, \mathbf{Y}}} P(\mathbf{A}|\mathbf{X})$$

- CTC loss

$$L_{\theta}(\mathbf{X}, \mathbf{Y}) = -\log P(\mathbf{Y}|\mathbf{X})$$

- Problem: Combinatorial explosion in the cardinality of the set $\mathbb{K}_{\mathbf{X}, \mathbf{Y}}$

- Solution: Compute $P(\mathbf{Y}|\mathbf{X})$ by dynamic programming



CTC Loss Computation: Formal Description

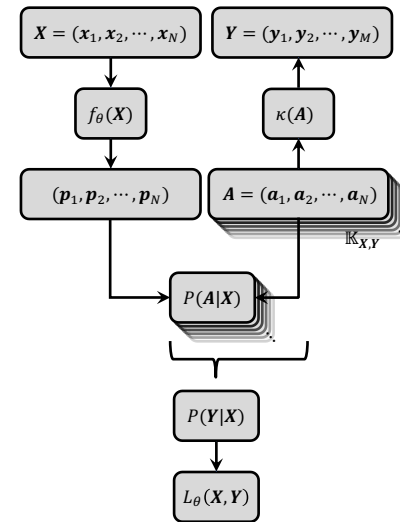
- Modified label sequence (length $2M + 1$, elements $y_m \in \mathbb{A}'$ for $m \in \{1, \dots, 2M + 1\}$)

$$\mathbf{Z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{2M+1})$$

obtained by inserting blank symbols ϵ into \mathbf{Y} at the beginning, the end, and between each two consecutive elements

- Idea: Align network output $(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$ to modified label sequence $(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{2M+1})$, but allow for skipping the inserted blank symbols ϵ (except where it is used to separate consecutive duplicates)
- Set of all alignments that correspond to \mathbf{Z}

$$\mathbb{K}'_{\mathbf{X}, \mathbf{Z}} = \{\mathbf{A} \in (\mathbb{A}')^N : \kappa'(\mathbf{A}) = \mathbf{Z}\}$$



Examples

$$Y = (\alpha, \beta)$$

$$Z = (\epsilon, \alpha, \epsilon, \beta, \epsilon)$$

$$Y = (\alpha, \alpha, \beta, \beta, \beta)$$

$$Z = (\epsilon, \alpha, \epsilon, \alpha, \epsilon, \beta, \epsilon, \beta, \epsilon, \beta, \epsilon)$$

CTC Loss Computation: Formal Description

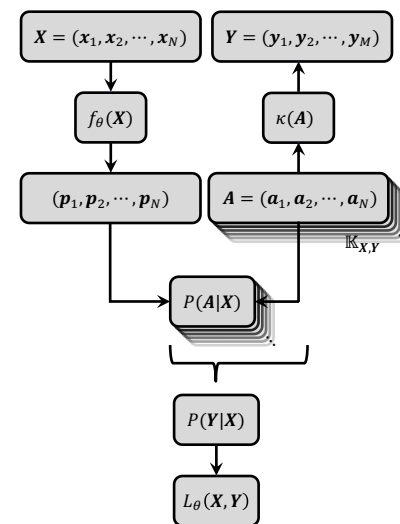
- Forward matrix

$$D(m, n) = \sum_{\substack{A \in \mathbb{K}'_{X,Z} \\ \kappa'(A(1:n)) = Z(1:m)}} \prod_{i=1}^n p_i(a_i)$$

of $n \in \{1, \dots, N\}$ and $m \in \{1, \dots, 2M + 1\}$,
encoding the probability that the first n time
steps correspond to the first m elements of the
modified label sequence Z

- Forward matrix contains probability of label
sequence

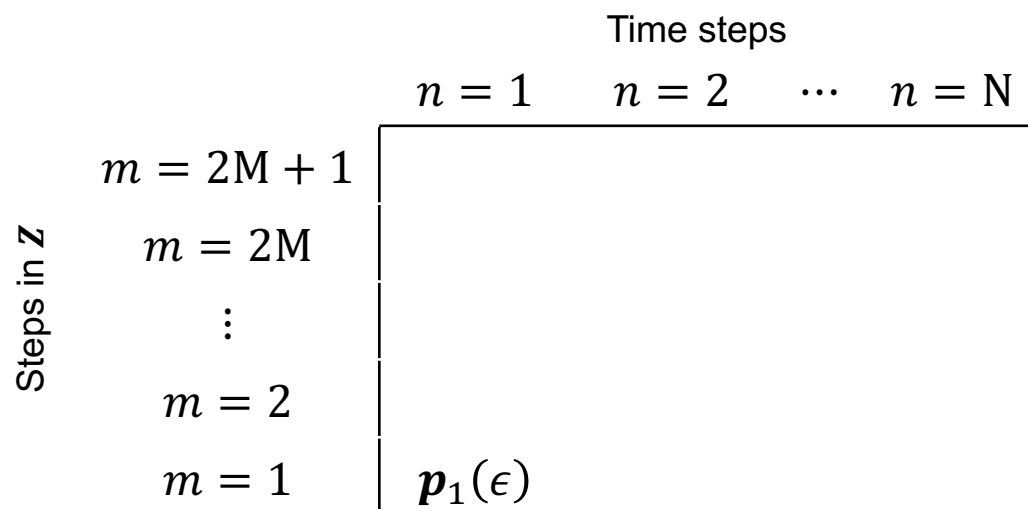
$$P(Y|X) = D(2M + 1, N) + D(2M, N)$$



CTC Loss Computation: Formal Description

- Initialization

$$D(1,1) = \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon)$$

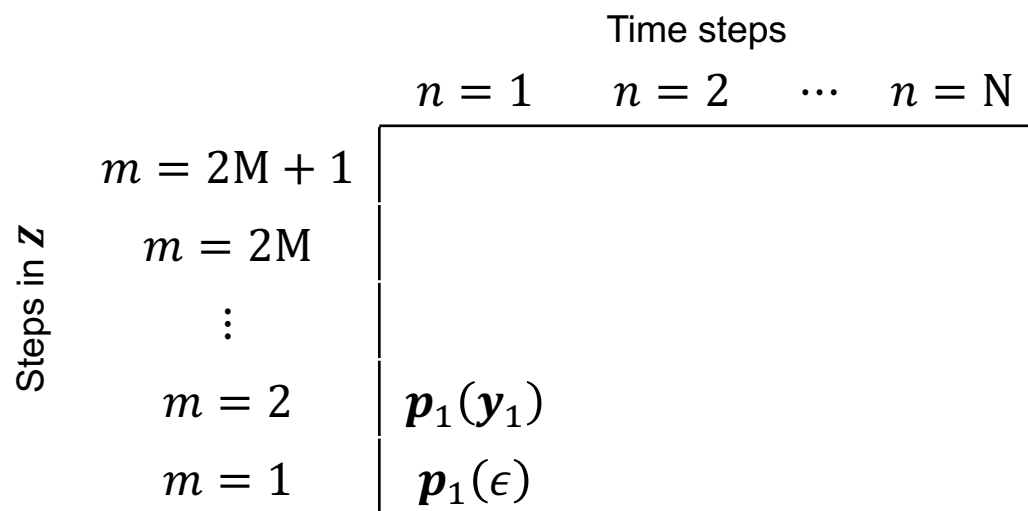


CTC Loss Computation: Formal Description

- Initialization

$$D(1,1) = p_1(\mathbf{z}_1) = p_1(\epsilon)$$

$$D(2,1) = p_1(\mathbf{z}_2) = p_1(\mathbf{y}_1)$$



CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned} D(1,1) &= \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon) \\ D(2,1) &= \mathbf{p}_1(\mathbf{z}_2) = \mathbf{p}_1(\mathbf{y}_1) \\ D(m,1) &= 0, \quad \forall m > 2 \end{aligned}$$

		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in Z	$m = 2M + 1$	0			
	$m = 2M$	0			
	\vdots	0			
	$m = 2$	$\mathbf{p}_1(\mathbf{y}_1)$			
	$m = 1$	$\mathbf{p}_1(\epsilon)$			

CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned}
 \mathbf{D}(1,1) &= \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon) \\
 \mathbf{D}(2,1) &= \mathbf{p}_1(\mathbf{z}_2) = \mathbf{p}_1(\mathbf{y}_1) \\
 \mathbf{D}(m, 1) &= 0, \quad \forall m > 2
 \end{aligned}$$

- Recursion

$$\mathbf{D}(m, n) = \mathbf{p}_n(\mathbf{z}_m) \cdot \begin{cases} (\mathbf{D}(m, n-1) + \mathbf{D}(m-1, n-1)), & \text{if } \mathbf{z}_m = \epsilon \end{cases}$$

		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in \mathbf{Z}	$m = 2M + 1$	0			
	$m = 2M$	0			
	\vdots	0			
	$m = 2$	$\mathbf{p}_1(\mathbf{y}_1)$			
	$m = 1$	$\mathbf{p}_1(\epsilon)$			

Example

$$\mathbf{Z} = (\dots, \epsilon, \alpha, \epsilon, \alpha, \epsilon, \beta, \epsilon, \dots)$$

CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned}
 D(1,1) &= \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon) \\
 D(2,1) &= \mathbf{p}_1(\mathbf{z}_2) = \mathbf{p}_1(\mathbf{y}_1) \\
 D(m, 1) &= 0, \quad \forall m > 2
 \end{aligned}$$

- Recursion

$$D(m, n) = \mathbf{p}_n(\mathbf{z}_m) \cdot \begin{cases} (D(m, n-1) + D(m-1, n-1)), & \text{if } \mathbf{z}_m = \epsilon \\ (D(m, n-1) + D(m-1, n)), & \text{if } \mathbf{z}_{m-2} = \mathbf{z}_m \end{cases}$$

		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in \mathbf{Z}	$m = 2M + 1$	0			
	$m = 2M$	0			
	\vdots	0			
	$m = 2$	$\mathbf{p}_1(\mathbf{y}_1)$			
	$m = 1$	$\mathbf{p}_1(\epsilon)$			

Example

$$\mathbf{Z} = (\dots, \epsilon, \alpha, \epsilon, \alpha, \epsilon, \beta, \epsilon, \dots)$$

CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned}
 D(1,1) &= \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon) \\
 D(2,1) &= \mathbf{p}_1(\mathbf{z}_2) = \mathbf{p}_1(\mathbf{y}_1) \\
 D(m, 1) &= 0, \quad \forall m > 2
 \end{aligned}$$

- Recursion

$$D(m, n) = \mathbf{p}_n(\mathbf{z}_m) \cdot \begin{cases} (D(m, n-1) + D(m-1, n-1)), & \text{if } \mathbf{z}_m = \epsilon \\ (D(m, n-1) + D(m-1, n-1)), & \text{if } \mathbf{z}_{m-2} = \mathbf{z}_m \\ (D(m, n-1) + D(m-1, n-1) + D(m-2, n-1)), & \text{otherwise} \end{cases}$$

		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in \mathbf{Z}	$m = 2M + 1$	0			
	$m = 2M$	0			
	\vdots	0			
	$m = 2$	$\mathbf{p}_1(\mathbf{y}_1)$			
	$m = 1$	$\mathbf{p}_1(\epsilon)$			

Example

$$\mathbf{Z} = (\dots, \epsilon, \alpha, \epsilon, \alpha, \epsilon, \beta, \epsilon, \dots)$$

CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned}
 D(1,1) &= \mathbf{p}_1(\mathbf{z}_1) = \mathbf{p}_1(\epsilon) \\
 D(2,1) &= \mathbf{p}_1(\mathbf{z}_2) = \mathbf{p}_1(\mathbf{y}_1) \\
 D(m, 1) &= 0, \quad \forall m > 2
 \end{aligned}$$

- Recursion

$$D(m, n) = \mathbf{p}_n(\mathbf{z}_m) \cdot \begin{cases} (D(m, n-1) + D(m-1, n-1)), & \text{if } \mathbf{z}_m = \epsilon \\ (D(m, n-1) + D(m-1, n-1)), & \text{if } \mathbf{z}_{m-2} = \mathbf{z}_m \\ (D(m, n-1) + D(m-1, n-1) + D(m-2, n-1)), & \text{otherwise} \end{cases}$$

		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in Z	$m = 2M + 1$	0	\dots	\dots	$P(Y X)$
	$m = 2M$	0	\dots	\dots	
	\vdots	0	\dots	\dots	\dots
	$m = 2$	$\mathbf{p}_1(\mathbf{y}_1)$	\dots	\dots	\dots
	$m = 1$	$\mathbf{p}_1(\epsilon)$	\dots	\dots	\dots

CTC Loss Computation: Formal Description

- Initialization

$$\begin{aligned}
 D(1,1) &= p_1(z_1) = p_1(\epsilon) \\
 D(2,1) &= p_1(z_2) = p_1(y_1) \\
 D(m,1) &= 0, \quad \forall m > 2
 \end{aligned}$$

- Recursion

$$D(m,n) = p_n(z_m) \cdot \begin{cases} (D(m,n-1) + D(m-1,n-1)), & \text{if } z_m = \epsilon \\ (D(m,n-1) + D(m-1,n-1)), & \text{if } z_{m-2} = z_m \\ (D(m,n-1) + D(m-1,n-1) + D(m-2,n-1)), & \text{otherwise} \end{cases}$$

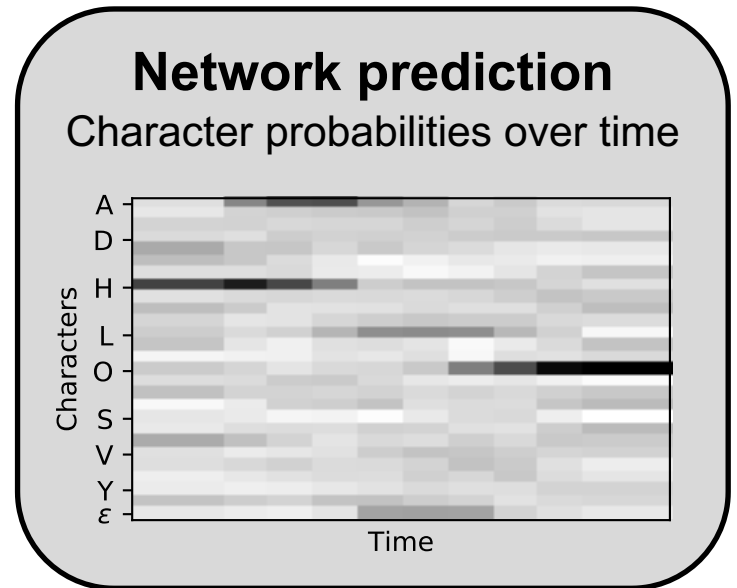
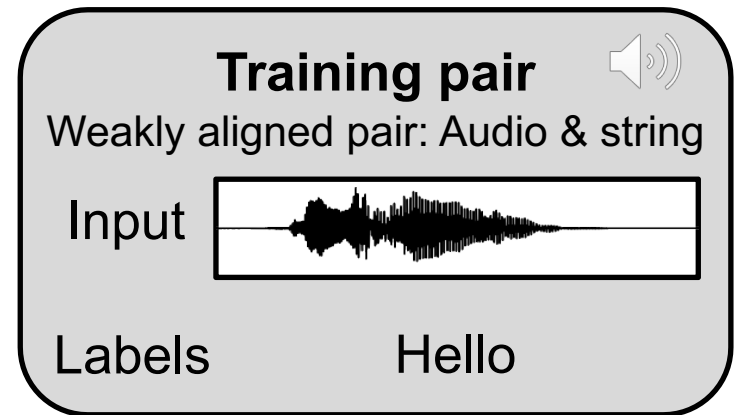
		Time steps			
		$n = 1$	$n = 2$	\dots	$n = N$
Steps in Z	$m = 2M + 1$	0	\dots	\dots	$P(Y X)$
	$m = 2M$	0	\dots	\dots	
	\vdots	0	\dots	\dots	\dots
	$m = 2$	$p_1(y_1)$	\dots	\dots	\dots
	$m = 1$	$p_1(\epsilon)$	\dots	\dots	\dots

Difference to hard alignment: Sum instead of max!

Applications

Application: Speech Recognition

- Task: Estimate character labels from waveform
- CTC-based training: Using weakly aligned pairs to train DNN for computing character probabilities¹
- Classes: 26 characters, space, and blank symbol
- Approach: Finding most probable character sequence for given character probabilities, e.g., using beam search
- CTC is a core technology used in today's speech recognizing systems, e.g., in the Google App²

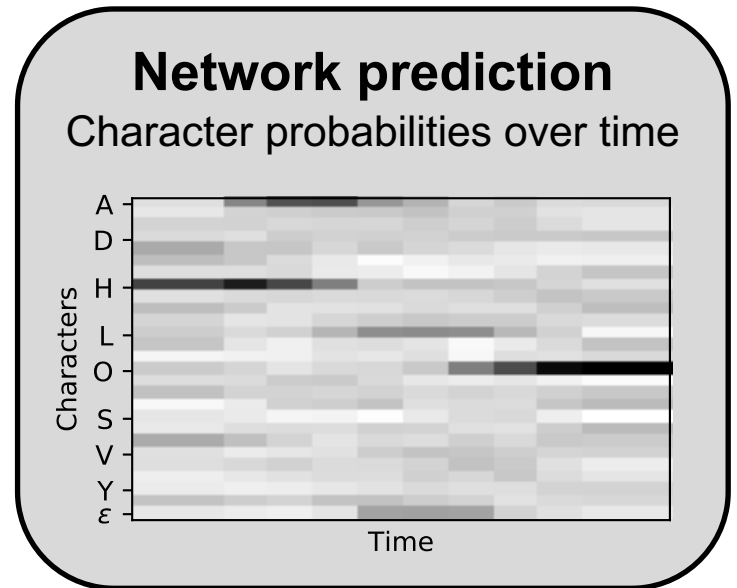
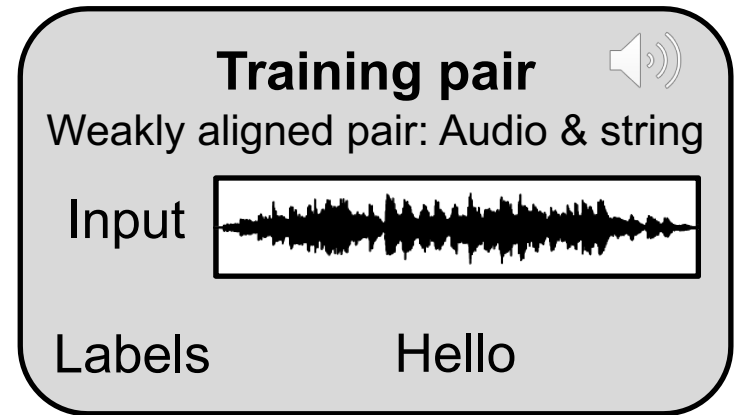


¹ Graves et al. *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. ICML 2006. [\[ACM\]](#)

² Sak et al. *Google Voice Search: Faster and More Accurate*. Google AI Blog, 2015 (<https://ai.googleblog.com/2015/09/google-voice-search-faster-and-more.html>).


Application: Lyrics Alignment

- Task: Align character labels from lyrics to music recording¹
- CTC-based training: Using weakly aligned pairs to train DNN for computing character probabilities (as before)²
- Classes: 26 characters, space, and blank symbol (as before)
- But: Now acoustic model needs to ignore non-singing-voice components of input representation
- Approach: Finding most probable alignment for given lyrics (dynamic programming similar to Viterbi)
- Beyond lyrics alignment, lyrics transcription is a challenging problem¹



¹ Stoller et al. *End-to-end Lyrics Alignment for Polyphonic Music Using an Audio-To-Character Recognition Model*. ICASSP 2019. [\[IEEE\]](#)

² Graves et al. *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. ICML 2006. [\[ACM\]](#)



Demo: Synced
lyrics feature in the
Apple Music App

Song: *Jingle Bells*
by Pentatonix

Source: Twitter,
@PTXofficial
[https://twitter.com/PTXofficial/
status/1337556809553301506](https://twitter.com/PTXofficial/status/1337556809553301506)

Sorry for the untimely
Christmas music!

Application: Theme-Based Music Retrieval

- Task: Given a symbolically encoded musical theme, find music recordings, where theme is being played^{1–3}
- Challenges due to differences in:¹
 - Modality (symbolic vs. audio)
 - Tuning
 - Transposition
 - Tempo
 - Degree of polyphony

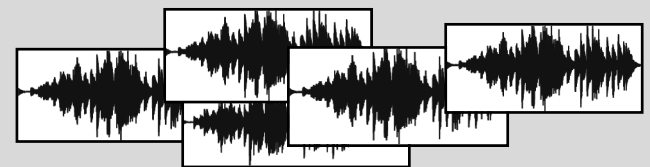
Query

Symbolically encoded monophonic musical theme



Database

Audio recordings of polyphonic music



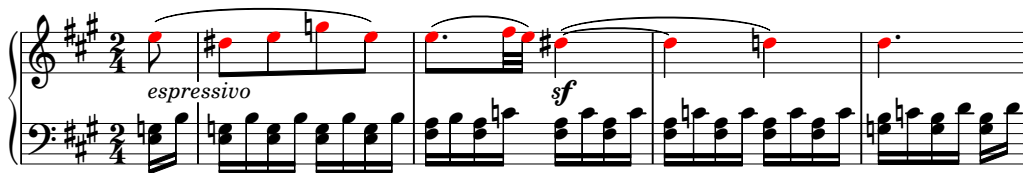
¹ Balke et al. *Retrieving Audio Recordings Using Musical Themes*. ICASSP 2016. [\[IEEE\]](#)

² Zalkow et al. *Evaluating Saliency Representations for Cross-Modal Retrieval of Western Classical Music Recordings*. ICASSP 2019. [\[IEEE\]](#)

³ Zalkow and Müller. *Using Weakly Aligned Score–Audio Pairs to Train Deep Chroma Models for Cross-Modal Music Retrieval*. ISMIR 2020. [\[Zenodo\]](#)

Application: Theme-Based Music Retrieval

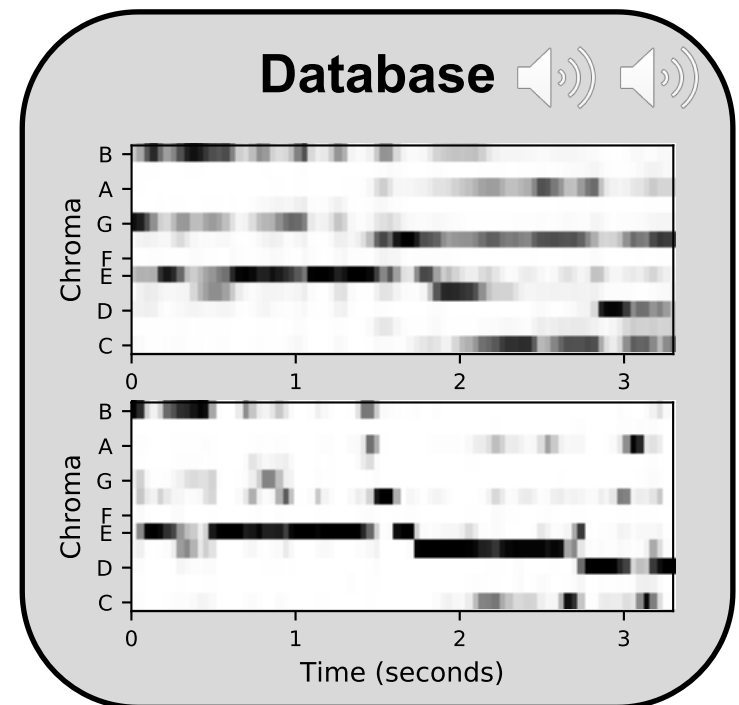
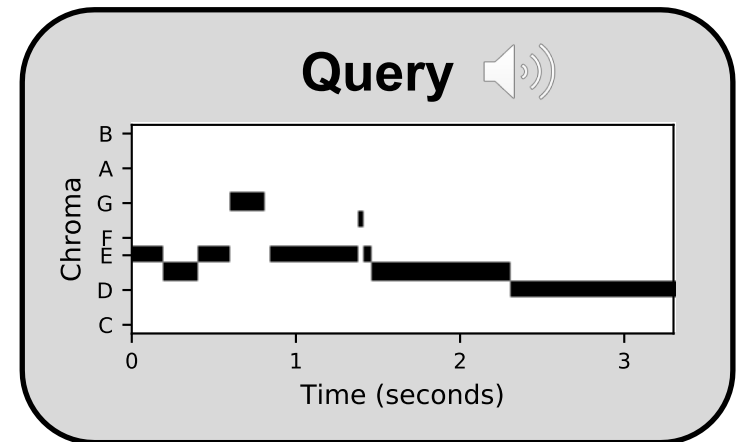
- Retrieval procedure based on chroma features and sequence alignment algorithm (subsequence dynamic time warping)¹⁻³
- Standard chroma features capture the full spectral content (influenced by theme and accompaniment)
- Learning enhanced chroma features with CTC loss (mainly influenced by theme and not by accompaniment)³



¹ Balke et al. *Retrieving Audio Recordings Using Musical Themes*. ICASSP 2016. [\[IEEE\]](#)

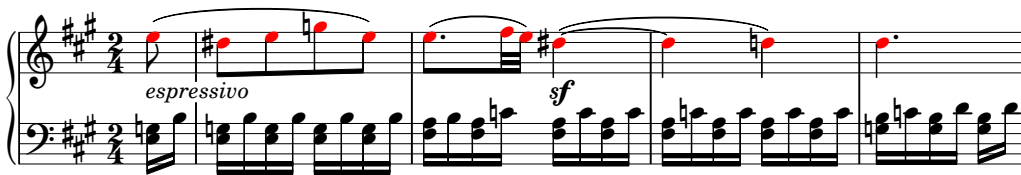
² Zalkow et al. *Evaluating Saliency Representations for Cross-Modal Retrieval of Western Classical Music Recordings*. ICASSP 2019. [\[IEEE\]](#)

³ Zalkow and Müller. *Using Weakly Aligned Score–Audio Pairs to Train Deep Chroma Models for Cross-Modal Music Retrieval*. ISMIR 2020. [\[Zenodo\]](#)



Application: Theme-Based Music Retrieval

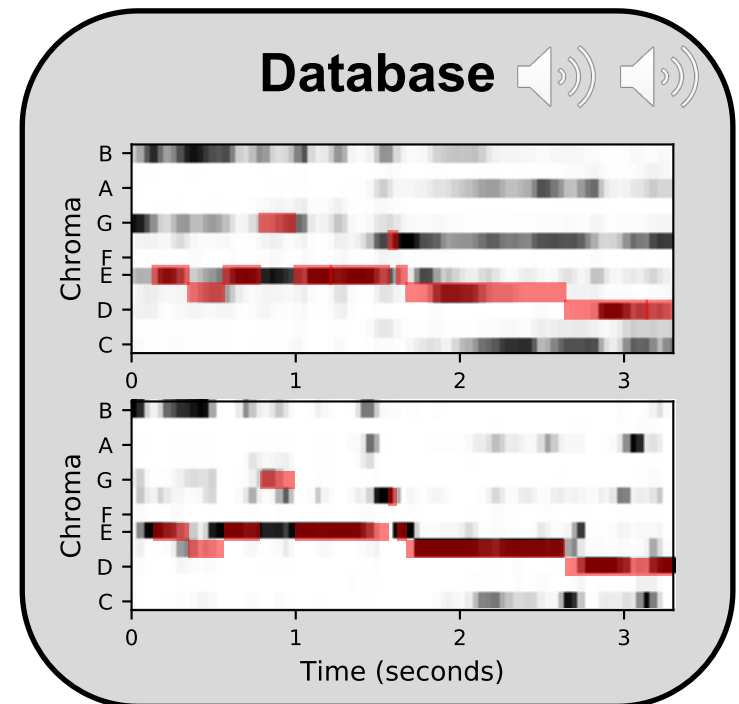
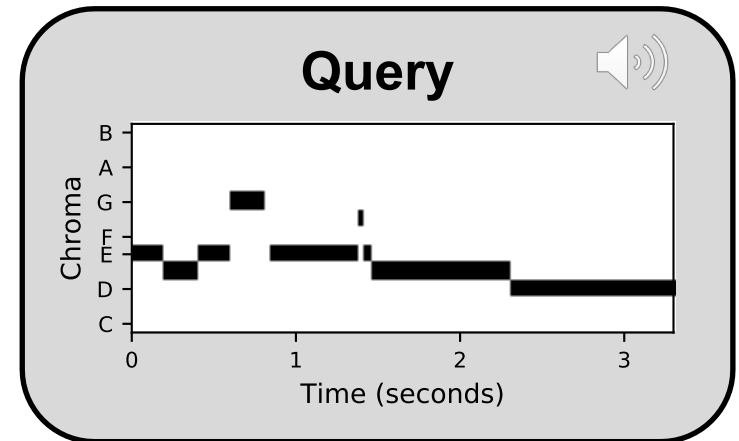
- Retrieval procedure based on chroma features and sequence alignment algorithm (subsequence dynamic time warping)¹⁻³
- Standard chroma features capture the full spectral content (influenced by theme and accompaniment)
- Learning enhanced chroma features with CTC loss (mainly influenced by theme and not by accompaniment)³



¹ Balke et al. *Retrieving Audio Recordings Using Musical Themes*. ICASSP 2016. [\[IEEE\]](#)

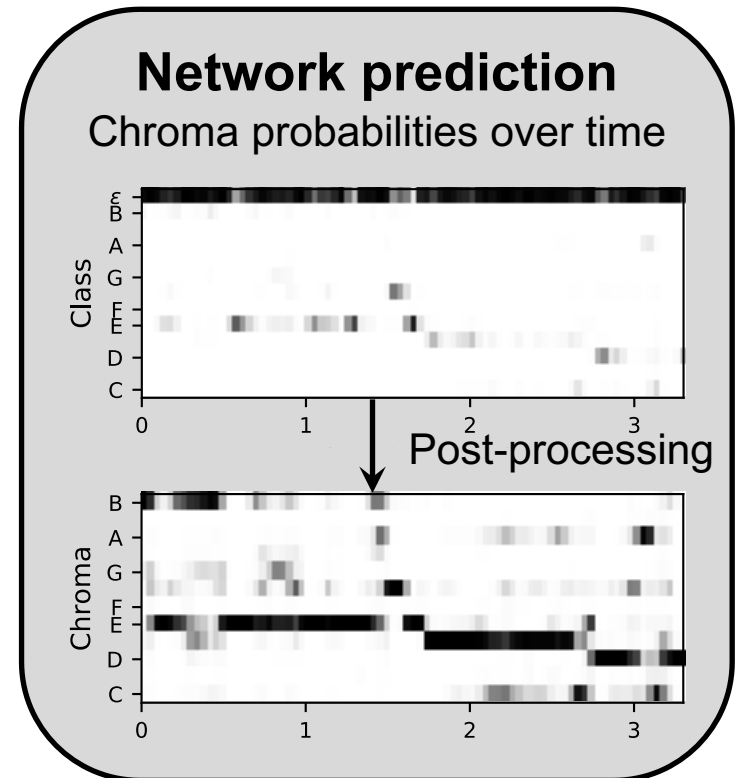
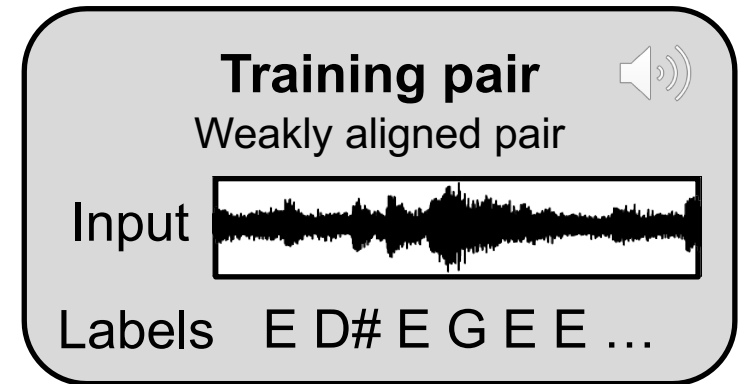
² Zalkow et al. *Evaluating Saliency Representations for Cross-Modal Retrieval of Western Classical Music Recordings*. ICASSP 2019. [\[IEEE\]](#)

³ Zalkow and Müller. *Using Weakly Aligned Score–Audio Pairs to Train Deep Chroma Models for Cross-Modal Music Retrieval*. ISMIR 2020. [\[Zenodo\]](#)



Application: Theme-Based Music Retrieval

- Task: Learn chroma representation that represents musical themes³
- CTC-based training: Using weakly aligned score–audio pairs to train DNN for computing chroma probabilities
- Classes: 12 chroma labels and blank symbol
- Observation: Blank-probabilities are active most of the time
- Approach: Post-processing of network prediction (remove blank probabilities and ℓ^2 -normalize each column)



¹ Balke et al. *Retrieving Audio Recordings Using Musical Themes*. ICASSP 2016. [\[IEEE\]](#)

² Zalkow et al. *Evaluating Saliency Representations for Cross-Modal Retrieval of Western Classical Music Recordings*. ICASSP 2019. [\[IEEE\]](#)

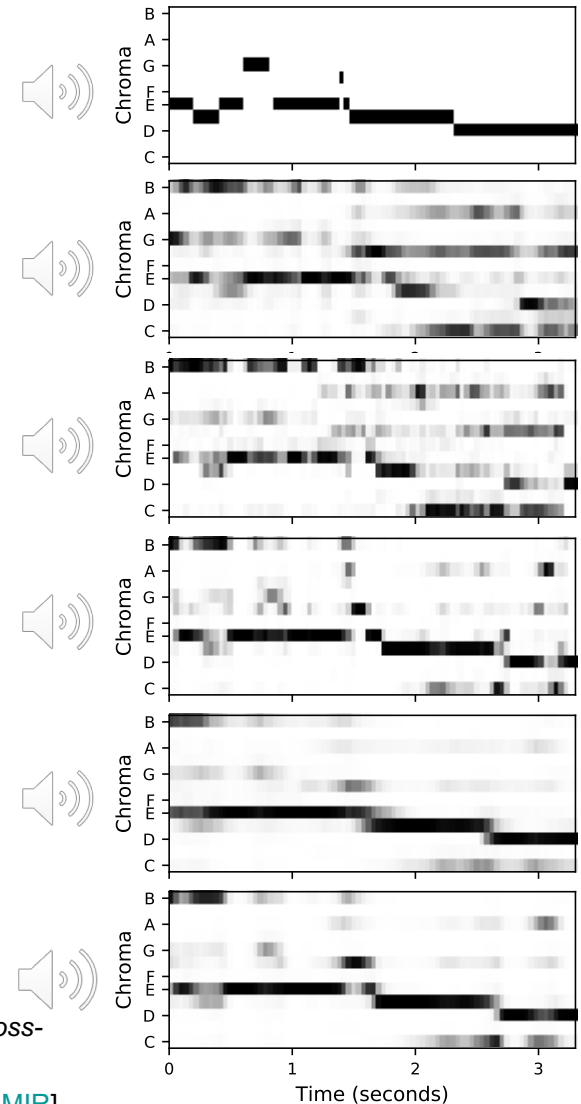
³ Zalkow and Müller. *Using Weakly Aligned Score–Audio Pairs to Train Deep Chroma Models for Cross-Modal Music Retrieval*. ISMIR 2020. [\[Zenodo\]](#)

Application: Theme-Based Music Retrieval

- Retrieval results¹ for dataset² with 2067 musical themes



Chroma Variant	Top-1	Top-10
Standard chroma features	0.561	0.723
Enhanced chroma features (baseline)	0.824	0.861
DNN-based chroma features (CTC)	0.867	0.942
DNN-based chroma features (linear scaling)	0.829	0.914
DNN-based chroma features (strong alignment)	0.882	0.939



¹ Zalkow and Müller. *Using Weakly Aligned Score–Audio Pairs to Train Deep Chroma Models for Cross-Modal Music Retrieval*. ISMIR 2020. [\[Zenodo\]](#)

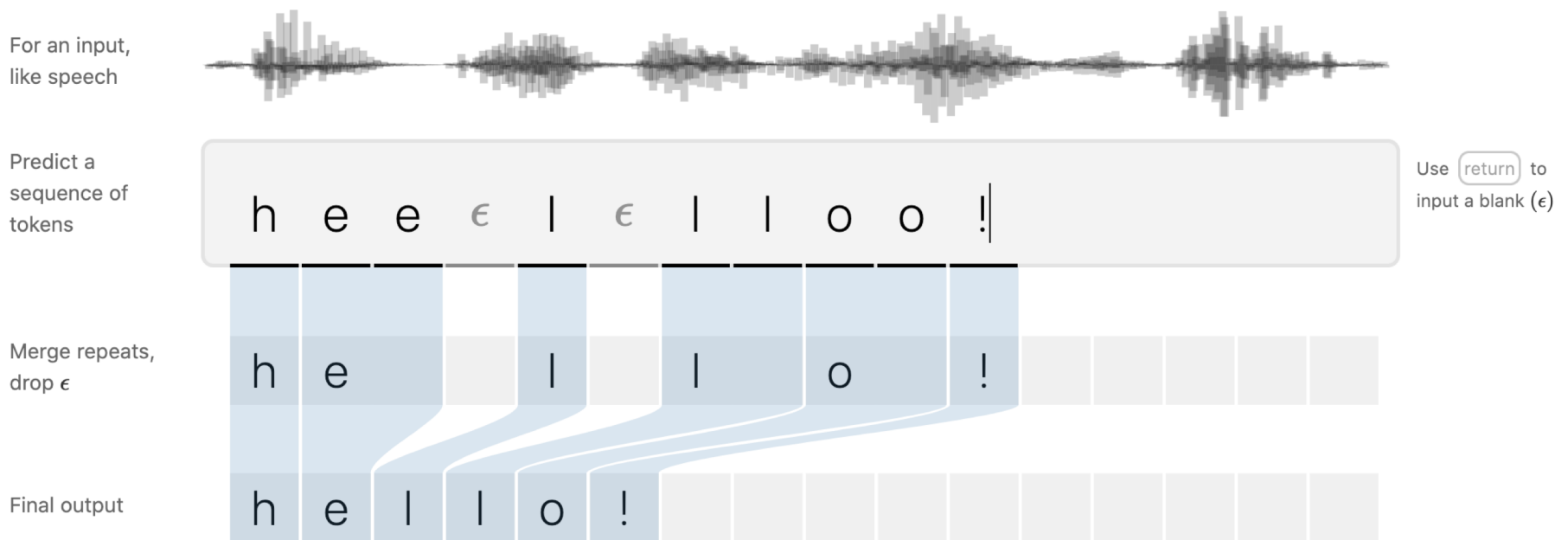
² Zalkow et al. *MTD: A Multimodal Dataset of Musical Themes for MIR Research*. TISMIR 2020. [\[TISMIR\]](#)

³ Bosch and Gómez. *Melody Extraction Based on a Source-Filter Model Using Pitch Contour Selection*. SMC 2016. [\[UPF\]](#)

Outlook and Further Notes

Outlook and Further Notes

- Good review of loss computation by Hannun.^{1,2}



¹ Hannun. *Sequence Modeling with CTC*. Distill, 2017 (<https://distill.pub/2017/ctc/>)

² Hannun. *Transcribing Real-valued Sequences with Deep Neural Networks*. PhD Thesis, Stanford University, 2018. [[Stanford](#)]

Outlook and Further Notes

- CTC is suited for gradient-based training because it accumulates all possible alignments (i.e., summation over the set $\mathbb{K}_{X,Y}$). This procedure is different from related alignment algorithms such as dynamic time warping (DTW), where the optimal alignment is computed (corresponding to a maximization over $\mathbb{K}_{X,Y}$). A variant of DTW adapted for the usage with neural networks is known as soft-DTW.¹

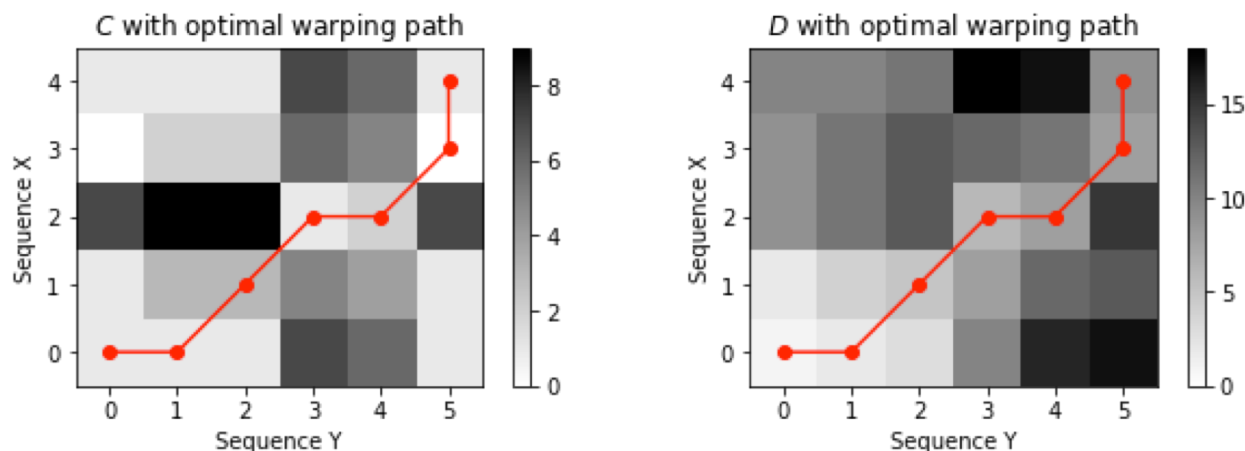


Figure from FMP²

¹ Cuturi and Blondel: *Soft-DTW: A Differentiable Loss Function for Time-Series*, ICML 2017. [\[PMLR\]](#)

² Müller and Zalkow. *FMP Notebooks: Educational Material for Teaching and Learning Fundamentals of Music Processing*. ISMIR 2019. [\[Zenodo\]](#)

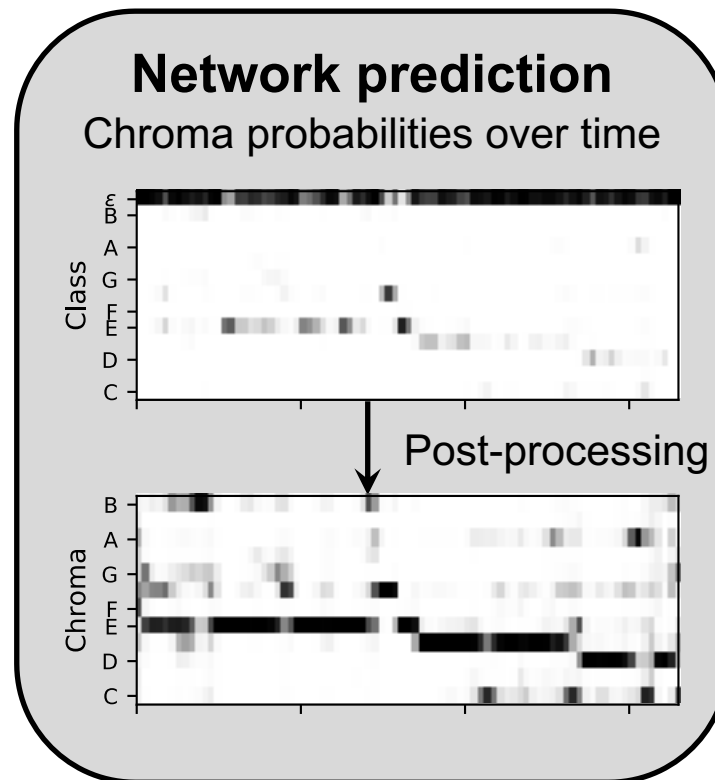
Outlook and Further Notes

- A CTC-based network models the dependencies between input and output. However, the dependencies between individual output elements are not modeled explicitly. One may use an external “language model” if needed. An extension of CTC for jointly modeling both input–output and output–output dependencies is denoted as Sequence Transduction.¹

¹ Graves: *Sequence Transduction with Recurrent Neural Networks*. ICML Workshop on Representation Learning 2012. [[arXiv](#)]

Outlook and Further Notes

- Many studies show that the probabilities for the blank symbol are often dominant in the output of CTC-based networks (“spiky problem”). One may modify the CTC algorithm to compensate for that.¹



¹ Li and Wang: Reinterpreting CTC Training as Iterative Fitting. Pattern Recognition, 2020. [[ScienceDirect](#), [arXiv](#)]

Outlook and Further Notes

- In its basic form, CTC is only able to model one label at a time (music application: “monophonic”). There are also extensions of CTC to multi-label problems (music application: “polyphonic”).¹

¹ Wignington et al. *Multi-Label Connectionist Temporal Classification*. ICDAR 2019. [\[IEEE\]](#)