
A Differentiable Alignment Framework for Sequence-to-Sequence Modeling via Optimal Transport

Anonymous Author(s)

Affiliation

Address

email

Abstract

Accurate sequence-to-sequence (seq2seq) alignment is critical for applications like medical speech analysis and language learning tools relying on automatic speech recognition (ASR). State-of-the-art end-to-end (E2E) ASR systems, such as the Connectionist Temporal Classification (CTC) and transducer-based models, suffer from peaky behavior and alignment inaccuracies. In this paper, we propose a novel differentiable alignment framework based on one-dimensional optimal transport, enabling the model to learn a single alignment and perform ASR in an E2E manner. We introduce a pseudo-metric, called Sequence Optimal Transport Distance (SOTD), over the sequence space and discuss its theoretical properties. Based on the SOTD, we propose Optimal Temporal Transport Classification (OTTC) loss for ASR and contrast its behavior with CTC. Experimental results on the TIMIT, AMI, and LibriSpeech datasets show that our method considerably improves alignment performance compared to CTC and the more recently proposed Consistency-Regularized CTC, though with a trade-off in ASR performance. We believe this work opens new avenues for seq2seq alignment research, providing a solid foundation for further exploration and development within the community.

1 Introduction

Sequence-to-sequence (seq2seq) alignment is a fundamental challenge in automatic speech recognition (ASR), where, beyond text prediction, precise alignment of text to the corresponding speech is crucial for many applications. For example, in medical domain, accurate alignment helps speech and language pathologists pinpoint speech segments for analyzing pathological cues, such as stuttering or voice disorders. In real-time subtitling, precise alignment ensures that subtitles are synchronized with spoken words, which is crucial for live broadcasts and streaming content. In language learning tools, ASR systems use alignment to provide feedback on pronunciation and fluency, allowing learners to compare their speech to target pronunciations. In these ASR-driven applications, while word error rate (WER) is an important performance metric, frame-level and word-level alignment accuracy are equally important for improving the system’s applicability and responsiveness.

In the literature, two primary approaches to ASR have emerged, i.e., hybrid systems and end-to-end (E2E) models. In hybrid approaches, a deep neural network-hidden Markov model (DNN-HMM) [1, 2, 3, 4, 5, 6, 7] system is typically trained, where the DNN is optimized by minimizing cross-entropy loss on the forced alignments generated for each frame of audio embeddings from a hidden Markov model-Gaussian mixture model (HMM-GMM). One notable disadvantage of the hybrid approach is that the model cannot be optimized in an E2E manner, which may result in suboptimal performance [8]. More recently, E2E models for ASR have become very popular due to their superior performance. There are three popular approaches for training an E2E model: (i) attention-based

encoder-decoder (AED) models [9, 10, 11, 12], (ii) using Connectionist Temporal Classification (CTC) loss [13, 14], and (iii) neural Transducer-based models [15, 16, 17]. AED models use an encoder to convert the input audio sequence into a hidden representation. The decoder, typically auto-regressive, generates the output text sequence by attending to specific parts of the input through an attention mechanism, often referred to as soft alignment [18] between the audio and text sequences. This design, however, can make it challenging to obtain word-level timestamps and to do teacher-student training with soft labels. Training AED models also requires a comparatively large amount of data, which can be prohibitive in low-resource setups. In contrast to AED models, CTC and transducer-based models maximize the marginal probability of the correct sequence of tokens (transcript) over all possible valid alignments (paths), often referred to as hard alignment [18]. However, recent research has shown that only a few paths, which are dominated by blank labels, contribute meaningfully to the marginalization, leading to the well-known peaky behavior that can result in suboptimal ASR performance [19]. Unfortunately, it is not possible to directly identify these prominent paths, or those that do not disproportionately favor blank labels, in advance within E2E models. This observation serves as the main motivation of our work.

In this paper, we introduce the Optimal Temporal Transport Classification (OTTC) loss function, a novel approach to ASR where our model jointly learns temporal sequence alignment and audio frame classification. OTTC is derived from the Sequence Optimal Transport Distance (SOTD) framework, which is also introduced in this paper and defines a pseudo-metric for finite-length sequences. At the core of this framework is a novel, parameterized, and differentiable alignment model based on one-dimensional optimal transport, offering both simplicity and efficiency, with linear time and space complexity relative to the largest sequence size. This design allows OTTC to be fast and scalable, maximizing the probability of exactly one path, which, as we demonstrate, helps avoid the peaky behavior commonly seen in CTC based models.

To summarize, our contributions are the following:

- We propose a novel, parameterized, and differentiable seq2seq alignment model with linear complexity both in time and space.
- We introduce a new framework, i.e., SOTD, to compare finite-length sequences, examining its theoretical properties and providing guarantees on the existence and characteristics of a minimum.
- We derive a new loss function, i.e., OTTC, specifically designed for ASR tasks.
- Finally, we conduct proof-of-concept experiments on the TIMIT [20], AMI [21], and Librispeech [22] datasets, demonstrating that our method mitigates the peaky behavior, improves alignment performance, and achieves promising results in E2E ASR.

2 Related Work

CTC loss. The CTC criterion [13] is a versatile method for learning alignments between sequences. This versatility has led to its application across various seq2seq tasks [23, 24, 18, 25, 26, 27]. However, despite its widespread use, CTC has numerous limitations that impact its effectiveness in real-world applications. To address issues such as peaky behavior [19], label delay [28], and alignment drift [29], researchers have proposed various extensions. These extensions aim to refine the alignment process, ensuring better performance across diverse tasks. Delay-penalized CTC [30] and blank symbol regularization [31, 32, 33] attempt to mitigate label delay issues. Other works have tried to control alignment through teacher model spikes [34, 35] or external supervision [36, 37, 38], though this increases complexity. More recently, Bayes Risk CTC [28] offer customizable, E2E approaches to improve alignment without relying on external supervision. The latest advancement, Consistency-Regularized CTC (CR-CTC) [39], mitigates extreme peaky behavior by enforcing consistency between CTC distributions obtained from different augmented views of the same audio.

Transducer loss. The transducer loss was introduced to address the conditional independence assumption of CTC by incorporating a predictor network [15]. However, similarly to CTC, transducer models suffer from label delay and peaky behavior [40]. To mitigate these issues, several methods have been proposed, such as e.g., Pruned RNN-T [16], which prunes alignment paths before loss computation, FastEmit [40], which encourages faster symbol emission, delay-penalized transducers [41], which add a constant delay to all non-blank log-probabilities, and minimum latency training

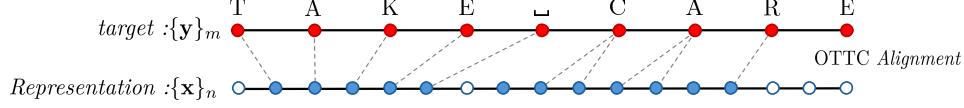


Figure 1: **Alignment between embeddings of frames and target sequence.** Red bullets represent the elements of the target sequence $\{y\}_m$, while the blue bullets indicate the frame embeddings $\{x\}_n$. In OTTC, the alignment guides the prediction model F in determining which frames should map to which labels. Additionally, the alignment model has the flexibility to leave some frames unaligned, as represented by the blue-and-white bullets, allowing those frames to be dropped during inference.

[42], which augments the transducer loss with the expected latency. Further extensions include CIFTransducer for efficient alignment [43], self-alignment techniques [44], and lightweight transducer models using CTC forced alignments [45].

Over the years, the CTC and transducer-based ASR models have achieved state-of-the-art performance. Despite numerous efforts to control alignments and apply path pruning, the fundamental formulation of marginalizing over all valid paths remains unchanged and directly or indirectly contributes to several of the aforementioned limitations. Instead of marginalizing over all valid paths as in CTC and transducer models, we propose a differential alignment framework based on optimal transport, which can jointly learn a single alignment and perform ASR in an E2E manner.

3 Problem Formulation

We define $\mathcal{U}_{\leq N}^d = \bigcup_{n \leq N} \mathcal{U}_n^d$ to be the set of all d -dimensional vector sequences of length at most N . Let us consider a distribution $\mathcal{D}_{\mathcal{U}_{\leq N}^d \times \mathcal{U}_{\leq N}^d}$ and pairs of sequences $(\{x_i\}_{i=1}^n, \{y_i\}_{i=1}^m)$ of length n and m drawn from $\mathcal{D}_{\mathcal{U}_{\leq N}^d \times \mathcal{U}_{\leq N}^d}$. For notational simplicity, the sequences of the pairs $(\{x_i\}_{i=1}^n, \{y_i\}_{i=1}^m)$ will be respectively denoted by $\{x\}_n$ and $\{y\}_m$ in the following. The goal in seq2seq tasks is to train a classifier that can accurately predict the target sequence $\{y\}_m$ from the input sequence $\{x\}_n$, enabling it to generalize to unseen examples. Typically, $n \neq m$, creating challenges for accurate prediction as there is no natural alignment between the two sequences. In this paper, we introduce a framework to address this class of problems, applying it specifically to the ASR domain. In this context, the first sequence $\{x\}_n$ represents an audio signal, where each vector $x_i \in \mathbb{R}^d$ corresponds to a time frame in the acoustic embedding space. The second sequence $\{y\}_m$ is the textual transcription of the audio, where each element y_i belongs to a predefined vocabulary $L = \{l_1, \dots, l_{|L|}\}$, such that $\{y\}_m \in L^m$, where L^m denotes the set of all m -length sequences formed from the vocabulary L .

4 Optimal Temporal Transport Classification

The core idea is to model the alignment between two sequences as a mapping to be learned along with the frame labels (see Figure 1). As the classification of audio frames improves, inferring the correct alignment becomes easier. Conversely, accurate alignments also improve frame classification. This mutual reinforcement between alignment and classification highlights the benefit of addressing both tasks simultaneously, contrasting with traditional hybrid models that treat them as separate tasks [1]. To achieve this, we propose the SOTD, a framework for constructing pseudo-metrics over the sequence space $\mathcal{U}_{\leq N}^d$, based on a differentiable, parameterized model that learns to align sequences. Using this framework, we derive the OTTC loss, which allows the model to learn both the alignment and the classification in a unified manner.

Notation. We denote $\llbracket 1, n \rrbracket = \{1, \dots, n\}$.

4.1 Preliminaries

Definition 1. Discrete monotonic alignment. Given two sequences $\{x\}_n$ and $\{y\}_m$, and a set of index pairs $\mathbf{A} \subset \llbracket 1, n \rrbracket \times \llbracket 1, m \rrbracket$ representing their alignment, we say that \mathbf{A} is a discrete monotonic alignment between the two sequences if:

- **Complete alignment of $\{y\}_m$:** Every element of $\{y\}_m$ is aligned, i.e.,

$$\forall j \in \llbracket 1, m \rrbracket, \exists k \in \llbracket 1, n \rrbracket, (k, j) \in \mathbf{A}.$$

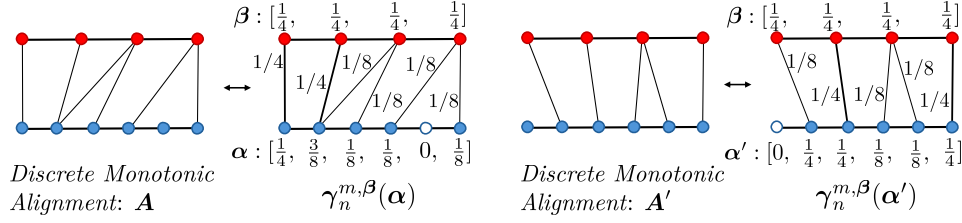


Figure 2: **Discrete monotonic alignment as 1D OT solution.** A discrete monotonic alignment represents a temporal alignment between two sequences (target on top, frame embeddings on bottom). It can be modeled by $\gamma_n^{m,\beta}$, as illustrated in the graph. The thickness of the links reflects the amount of mass $\gamma_n^{m,\beta}(\alpha)_{i,j}$ transported, with thicker links corresponding to higher mass.

127

- **Monotonicity:** The alignment is monotonic, meaning that for all $(i, j), (k, l) \in \mathbf{A}$

$$i \leq k \Rightarrow j \leq l.$$

128

Discrete monotonic alignments model the relationship between temporal sequences, such as those in ASR, by determining which frame should predict which target. The conditions imposed on the target sequence $\{y\}_m$ ensure that no target element is omitted, while the absence of similar constraints on the source sequence $\{x\}_n$ allows certain audio frames to be considered irrelevant and dropped (see Figure 2). The monotonicity condition preserves the temporal order, ensuring the sequential structure is maintained. In the following sections, we will develop a model capable of differentiating within the space of discrete monotonic alignments.

133

135 4.2 Differentiable Temporal Alignment with Optimal Transport

136

In the following, we introduce 1D OT and define our alignment model. Consider the 1D discrete distributions $\mu[\alpha, n]$ and $\nu[\beta, m]$ expressed as superpositions of δ measures, i.e., a distribution that is zero everywhere except at a single point, where it integrates to 1

137

$$\mu[\alpha, n] = \sum_{i=1}^n \alpha_i \delta_i \quad \text{and} \quad \nu[\beta, m] = \sum_{i=1}^m \beta_i \delta_i. \quad (1)$$

139

The bins of $\mu[\alpha, n]$ and $\nu[\beta, m]$ are $\llbracket 1, n \rrbracket$ and $\llbracket 1, m \rrbracket$, respectively, whereas the weights α_i and β_i are components of the vectors $\alpha \in \Delta^n$ and $\beta \in \Delta^m$, with Δ^n the simplex set defined as $\Delta^n = \{v \in \mathbb{R}^n | 0 \leq v_i \leq 1, \sum_{i=1}^n v_i = 1\} \subset \mathbb{R}^n$. OT theory provides an elegant and versatile framework for computing distances between distributions such as $\mu[\alpha, n]$ and $\nu[\beta, m]$, depending on the choice of the cost function [46] (chapter 2.4). One such distance is the 2-Wasserstein distance \mathcal{W}_2 , which measures the minimal cost of transporting the weight of one distribution to match the other. This distance is defined as

145

$$\mathcal{W}_2(\mu[\alpha, n], \nu[\beta, m]) = \min_{\gamma \in \Gamma^{\alpha, \beta}} \sum_{i,j=1}^{n,m} \gamma_{i,j} \|i - j\|_2^2, \quad (2)$$

146

where $\|i - j\|_2^2$ is the cost of moving weight from bin i to bin j and $\gamma_{i,j}$ is the amount of mass moved from i to j . The optimal coupling matrix γ^* is searched within the set of valid couplings $\Gamma^{\alpha, \beta}$, defined as

148

$$\Gamma^{\alpha, \beta} = \{\gamma \in \mathbb{R}_+^{n \times m} | \gamma \mathbf{1}_m = \alpha \text{ and } \gamma^T \mathbf{1}_n = \beta\}. \quad (3)$$

149

This constraint ensures that the coupling conserves mass, accurately redistributing all weights between the bins. A key property of optimal transport in 1D is its monotonicity [47]. Specifically, if there is mass transfer between bins i and j (i.e., $\gamma_{i,j}^* > 0$) and similarly between bins k and l (i.e., $\gamma_{k,l}^* > 0$), then it must hold that $i \leq k \Rightarrow j \leq l$. Consequently, when β has no zero components – meaning that every bin from ν is reached by the transport – the set $\{(i, j) \in \llbracket 1, n \rrbracket \times \llbracket 1, m \rrbracket | \gamma_{i,j}^* > 0\}$ satisfies the conditions of Definition 1, thereby forming a discrete monotonic alignment. This demonstrates that the optimal coupling can effectively model such alignments (see Figure 2).

156

Parameterized and differentiable temporal alignment. Given any sequences length n and m and β with no zero components, we can define the alignment function $\gamma_n^{m, \beta}$

157

$$\begin{aligned} \gamma_n^{m, \beta} : \Delta^n &\rightarrow \Gamma^{*, \beta}[n] \\ \alpha &\mapsto \gamma^* = \arg \min_{\gamma \in \Gamma} \mathcal{W}(\mu[\alpha, n], \nu[\beta, m]), \end{aligned} \quad (4)$$

where $\Gamma^{*,\beta}[n]$ is the space of all 1D transport solutions between $\mu[\alpha, n]$ and $\nu[\beta, m]$ for any α . Differently from β , α may have zero components, giving the model the flexibility to suppress certain bins, which acts similarly to a blank token in traditional models. In the context of ASR, α and β can be referred to as OT weights and label weights, respectively.

Lemma 1: The function $\alpha \mapsto \gamma_n^{m,\beta}(\alpha)$ is bijective from Δ^n to $\Gamma^{*,\beta}[n]$.

Proof. The proof can be found in Appendix A.2.1.

Proposition 1. Discrete Monotonic Alignment Approximation Equivalence. For any β that satisfies the condition above, any discrete set of alignments $A \subset \llbracket 1, n \rrbracket \times \llbracket 1, m \rrbracket$ between sequences of lengths n and m can be modeled by $\gamma_n^{m,\beta}$ through the appropriate selection of α , i.e.,

$$\forall A, \exists \alpha \in \Delta^n, (i, j) \in A \iff \gamma_n^{m,\beta}(\alpha)_{i,j} > 0. \quad (5)$$

Proof. The proof can be found in Appendix A.2.2.

Thus, we have defined a family of alignment functions $\gamma_n^{m,\beta}$ that are capable of modeling any discrete monotonic alignment, which can be chosen or adapted based on the specific task at hand. The computational cost of these alignment functions is low, as the bins are already sorted, eliminating the need for additional sorting. This results in linear complexity $O(\max(n, m))$ depending on the length of the longest sequence (see Algorithm A.1.1 in the Appendix). Furthermore, these alignments are differentiable, with $\gamma_n^{m,\beta}(\alpha)_{i,j}$ explicitly expressed in terms of α and β , allowing direct computation of the derivative $\frac{d\gamma_n^{m,\beta}(\alpha)_{i,j}}{d\alpha}$ via its analytical form.

4.2.1 Sequence-to-Sequence Distance

Here, we use the previously designed alignment functions to build a pseudo-metric over sets of sequences $\mathcal{U}_{\leq N}^d$.

Definition 1. Sequences Optimal Transport Distance (SOTD). Consider an n -length sequence $\{\mathbf{x}\}_n \in \mathcal{U}_{\leq N}^d$, an m -length sequence $\{\mathbf{y}\}_m \in \mathcal{U}_{\leq N}^d$, $p = \max(n, m)$, and $q = \min(n, m)$. Let $C : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$, be a differentiable positive cost function. Considering $r \in \mathbb{N}^*$ and a family of vectors $\{\beta\}_N = \{\beta_1 \in \Delta^1, \beta_2 \in \Delta^2, \dots, \beta_N \in \Delta^N\}$ without zero components, we define the SOTD \mathcal{S}_r as

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) = \min_{\alpha \in \Delta^n} \left(\sum_{i,j=1}^{n,m} \gamma_p^{q,\beta_q}(\alpha)_{i,j} \cdot C(\mathbf{x}_i, \mathbf{y}_j)^r \right)^{1/r}. \quad (6)$$

Note that β_q obviously depends on q , but could a priori depend on $\{\mathbf{x}\}_n$ and $\{\mathbf{y}\}_m$. To simplify the notation, we only denote its dependence on q . However, all the results in this section remain valid under such dependencies, as long as β_q components never becomes zero.

Proposition 2. Validity of the definition. SOTD is well-defined, meaning that a solution to the problem always exists, although it may not be unique.

Proof. The proof and the discussion about the non-unicity is conducted in Appendix A.2.3.

Proposition 3. SOTD is a Pseudo-Metric. If the cost matrix C is a metric on \mathbb{R}^d , then \mathcal{S}_r defines a pseudo-metric over the space sequences with at most N elements $\mathcal{U}_{\leq N}^d$.

Proof. The proof can be found in Appendix A.2.4.

Since \mathcal{S}_r is a pseudo-metric, there are sequences $\{\mathbf{x}\}_n \neq \{\mathbf{y}\}_m$ such that $\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) = 0$. The following proposition describes the conditions when this occurs.

Proposition 4. Non-Separation Condition. Let \mathcal{A} be the sequence aggregation operator which removes consecutive duplicates, i.e., $\mathcal{A}(\{\dots, \mathbf{x}, \mathbf{x}, \dots\}) = \{\dots, \mathbf{x}, \dots\}$. Let \mathcal{P}_α be the sequence pruning operator which removes any element \mathbf{x}_i from sequences corresponding to an $\alpha_i = 0$, i.e., $\mathcal{P}_\alpha(\{\dots, \mathbf{x}_{i-1}, \mathbf{x}_i, \mathbf{x}_{i+1}, \dots\}) = \{\dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots\}$ iff $\alpha_i = 0$. Further, let us consider $\{\mathbf{x}\}_n$ and $\{\mathbf{y}\}_m$ such that $\{\mathbf{x}\}_n \neq \{\mathbf{y}\}_m$. Without loss of generality, we assume that $n \geq m$. Then

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) = 0 \text{ iff } \mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)) = \mathcal{A}(\{\mathbf{y}\}_m), \quad (7)$$

where α^* is a minimum for which $\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) = 0$. It should be noted that this condition holds also when C is neither symmetric nor satisfies the triangular inequality, but is separated (like the cross-entropy for example). (*Proof.* See Appendix A.2.5.)

The consequence of the previous proposition is that we can learn a transformation through gradient descent using a trainable network F which maps input sequences $\{\mathbf{x}\}_n$ to target sequences $\{\mathbf{y}\}_m$ (with $n \geq m$) by solving the optimization problem

$$\min_F \mathcal{S}_r(F(\{\mathbf{x}\}_n), \{\mathbf{y}\}_m). \quad (8)$$

We are then guaranteed that a solution $F^*\{\mathbf{x}\}_n$ allows us to recover the sequence $\mathcal{A}(\{\mathbf{y}\}_m)$. In cases where retrieving repeated elements in $\{\mathbf{y}\}_m$ (e.g., double letters) is important, we can intersperse blank labels $\phi \notin L$ between repeated labels as follows: $\{\mathbf{y}\}_m = \{\dots, l_i, l_i, \dots\} \rightarrow \{\dots, l_i, \phi, l_i, \dots\}$.

Note on Dynamic Time Warping (DTW): A note on the distinction between our approach and DTW-based methods [48] can be found in Appendix A.4.

4.3 Application to ASR: OTTC Loss

In ASR, the target sequences $\{\mathbf{y}\}_m$ are d -dimensional one-hot encoding of elements from the set $L \cup \{\phi\}$, where ϕ is a blank label used to separate repeated labels. The encoder F predicts the label probabilities for each audio frame, such that

$$F(\{\mathbf{x}\}_n) = \{[p_{l_1}(\mathbf{x}_1), \dots, p_{l_{|L|+1}}(\mathbf{x}_1)]^T\}_{i=1}^n. \quad (9)$$

The alignment between $F(\{\mathbf{x}\}_n)$ and $\{\mathbf{y}\}_m$ is parameterized by $\alpha[\{\mathbf{x}\}_n, W] \in \Delta^n$, defined as

$$\alpha[\{\mathbf{x}\}_n, W] = \text{softmax}(W(\mathbf{x}_1), \dots, W(\mathbf{x}_n))^T \quad (10)$$

where W is a network that outputs a scalar for each frame \mathbf{x}_i . Using the framework built in Section 4.2.1 (with $r = 1$ and $C = C_e$, where C_e is the cross-entropy) to predict $\{\mathbf{y}\}_m$ from $\{\mathbf{x}\}_n$, we train both W and F by minimizing the OTTC objective

$$\mathcal{L}_{OTTC} = - \sum_{i,j=1}^{n,m} \gamma_n^{m,\beta_m}(\alpha[\{\mathbf{x}\}_n, W])_{i,j} \cdot \log p_{\mathbf{y}_j}(\mathbf{x}_i). \quad (11)$$

The choice of the cross-entropy C_e as the cost function arises naturally from the probabilistic encoding of the predicted output of F and the one-hot encoding of the target sequence. Additionally, since C_e is differentiable, it makes the OTTC loss differentiable with respect to F , while the differentiability of the OTTC with respect to W stems from the differentiability of γ_n^{m,β_m} with respect to its input $\alpha[\{\mathbf{x}\}_n, W]$. Thus, by following the gradient of this loss, we jointly learn both the alignment (via W) and the classification (via F).

Note: The notation $\gamma_n^{m,\beta}$ in Eq. 11 is valid in the context of ASR since $n \geq m$.

4.4 Link with CTC Loss

In this section, we link the CTC and the proposed OTTC losses. In the context of CTC, we denote by \mathcal{B} the mapping which reduces any sequences by deleting repeated vocabulary (similarly to the previously defined \mathcal{A} mapping in Proposition 5) and then deleting the blank token ϕ (e.g., $\mathcal{B}(\{GGOO\phi ODD\}) = \{GOOD\}$). The objective of CTC is to maximise the probability of all possible paths $\{\pi\}_n$ of length n through minimizing

$$-\log \sum_{\{\pi\}_n \in \mathcal{B}^{-1}(\{\mathbf{y}\}_m)} p(\{\pi\}_n) = -\log \sum_{\{\pi\}_n \in \mathcal{B}^{-1}(\{\mathbf{y}\}_m)} \prod_{i=1}^n p(\pi_i), \quad (12)$$

where $\{\pi\} \in L^n$ is an n -length sequence and $\mathcal{B}^{-1}(\{\mathbf{y}\}_m)$ is the set of all sequences collapsed by \mathcal{B} into $\{\mathbf{y}\}_m$.

Let us consider a path $\{\pi\}_n \in \mathcal{B}^{-1}(\{\mathbf{y}\}_m)$. Such a path can be seen as an alignment (see Figure 3), where $\{\mathbf{x}_i\}$ and $\{\mathbf{y}_j\}$ are aligned iff $\pi_i = \mathbf{y}_j$. By denoting \mathbf{A}_π as the corresponding discrete monotonic alignment, one can write

$$-\log p(\{\pi\}_n) = - \sum_{i=1}^n \log p_{\pi_i}(\mathbf{x}_i) = \sum_{\substack{i,j=1 \\ (i,j) \in \mathbf{A}_\pi}}^{n,m} C_e(\pi_j, \mathbf{y}_i) \stackrel{\exists \alpha \in \Delta^n}{=} \sum_{\substack{i,j=1 \\ \gamma_n^{m,\beta_m}(\alpha)_{i,j} > 0}}^{n,m} C_e(\pi_j, \mathbf{y}_i). \quad (13)$$

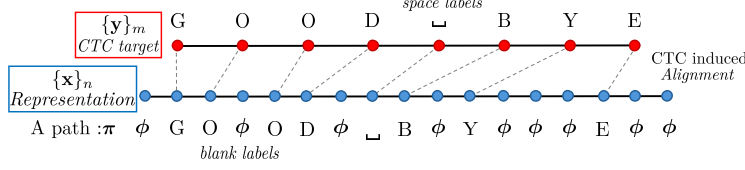


Figure 3: A **CTC alignment**. Here, we illustrate one of the valid alignments for CTC. The CTC loss maximizes the marginal probability over all such possible alignments.

Table 1: Alignment performance of the CTC-, CR-CTC-, and OTTC-based ASR models on the TIMIT and AMI datasets. [†]For TIMIT, we subtract the percentage of real silence, as it is available, unlike in AMI.

Model	TIMIT (Phoneme Level)			AMI (Word Level)		
	Peaky [†] (↓)	F1 Score (↑)	IDR (↑)	Peaky (↓)	F1 Score (↑)	IDR (↑)
CTC	53.51	88.77	26.98	81.93	83.94	16.75
CR-CTC	35.62	88.98	35.82	80.40	84.58	18.20
OTTC	0.76	89.27	76.72	54.75	84.81	42.84

with C_e representing the cross-entropy. The last equality arises from Proposition 1 and the fact that A_π represents a discrete monotonic alignment.

The continuous relaxation (i.e., making the problem continuous with respect to alignment) of the last term in this sequence of equalities results in $-\mathcal{L}_{OTTC}$. Therefore, OTTC can be seen as relaxation of the probability associated with a single path, enabling a differentiable path search mechanism. Essentially, OTTC optimization focuses on maximizing the probability of exactly one path, in contrast to CTC, which maximizes the probability across all valid paths.

Additionally, OTTC does not incentivize paths containing many blank tokens, unlike CTC. In CTC, the peaky behavior arises because maximizing the marginal probability over all valid paths can incentivize the model to assign more frames to the blank token [19]. In contrast, OTTC does not rely on a blank token to indicate that a frame i should not be classified (blank tokens are only used to separate consecutive tokens). Instead, the model simply sets the corresponding weight α_i to 0 (see Figure 2). This mechanism avoids the peaky behavior exhibited by CTC.

5 Experimental Setup

To demonstrate the viability of the proposed OTTC loss framework, we conduct several proof-of-concept experiments on the ASR task. To this end, we compare alignment quality and ASR performance using the proposed OTTC framework and existing CTC-based models. Note that an efficient batched implementation of OTTC along with the full code to reproduce our experimental results will be made publicly available.

Datasets. We conduct our experiments on popular open-source datasets, i.e., the TIMIT [20], AMI [21], and LibriSpeech [22]. TIMIT is a 5-hour English dataset with time-aligned transcriptions, including exact time-frame phoneme transcriptions, making it a standard benchmark for ASR and phoneme segmentation tasks. We report results on the standard eval set. AMI is an English spontaneous meeting speech corpus that serves as a good benchmark to evaluate our approach in a realistic conversational scenario, due to its spontaneous nature and prior use in alignment evaluation [49]. For our experiments on this dataset, we train models on the individual head microphone (IHM) split comprising 80 hours of audio, and report results on the official eval set. LibriSpeech is an English read-speech corpus derived from audiobooks, containing 1000 hours of data. It is a standard benchmark for reporting ASR results. For our experiments, we train models on the official 100-hour, 360-hour, and 960-hour splits, and report results on the two official test sets.

Baselines. We benchmark our performance against the standard CTC. To specifically compare alignment quality, particularly regarding the mitigation of the peaky behavior inherent in CTC-based models, we also include CR-CTC [39]. CR-CTC serves as a strong baseline, chosen for its established effectiveness against such peaky alignments.

Model architectures. We use the 300M parameter version of the well-known Wav2Vec2-large [50] as the base model for acoustic embeddings in all the experiments conducted in this work. The

Wav2Vec2 is a self-supervised model pre-trained on 60K hours of unlabeled English speech. For the baseline CTC-based models, we stack a dropout layer followed by a linear layer for logits prediction, termed the *logits prediction head*. For the proposed OTTC loss based model, we use a dropout and a linear layer (identical to the baseline) for logits prediction. In addition, as described in Section 4.3, we apply a dropout layer followed by two linear layers on top of the Wav2Vec2-large model for OT weight prediction, with a GeLU [51] non-linearity in between, termed the *OT weights prediction head*. Note that the output from the Wav2Vec2-large model is used as input for both the logit and OT weight prediction heads, and the entire model is trained using the OTTC loss.

Performance metrics. Alignment quality is assessed using three metrics: peaky behavior, starting frame accuracy, and Intersection Duration Ratio (IDR). Peaky behavior, a common characteristic of CTC-based models, refers to a large proportion of audio frames being assigned to blank or space symbols (non-alphabet symbols) [19]. To quantify this, we compute the average percentage of frames mapped to these symbols. Starting frame accuracy is evaluated using the F1 score, following the methodology proposed in [49]. It is important to note that this F1 score reflects only the correctness of the predicted token’s starting frame and does not fully capture alignment quality. To address this, we introduce IDR, which measures the overlap between predicted and reference word segments, normalized by the reference duration. This provides a finer-grained assessment of temporal alignment. These alignment metrics are computed only on the TIMIT and AMI datasets due to the lack of reliable ground-truth or forced-alignment annotations for LibriSpeech. On TIMIT, where ground-truth alignments are available, we assess alignment at the phoneme level. For AMI, which lacks ground-truth timestamps, we follow the forced-alignment approach in [49], but restrict evaluation to word-level timestamps, as they are generally more reliable than phoneme-, letter-, or subword-level annotations. Finally, ASR performance is evaluated using the WER on all considered databases.

Training details. In all our experiments, we use the AdamW optimizer [52] for training. For TIMIT and LibriSpeech, the initial learning rate is set to $lr = 2e^{-4}$, with a linear warm-up for the first 500 steps followed by a linear decay until the end of training. For AMI, the initial learning rate is set to $lr = 1.25e^{-3}$, with a linear warm-up during the first 10% of the steps, also followed by linear decay. We train all considered models for 40 epochs, reporting the test set WER at the final epoch. In our OTTC-based models, both the logits and OT weight prediction heads are trained for the first 30 epochs. During the final 10 epochs, the *OT weight prediction head* is fixed, while training continues on the *logits prediction head*. For experiments on the LibriSpeech (*resp.* TIMIT) dataset, we use character-level (*resp.* phoneme-level) tokens to encode text. Given the popularity of subword-based units for encoding text [53], we sought to observe the behavior of OTTC-based models when tokens are subword-based, where a token can contain more than one character. For the experiments on the AMI dataset, we use the SentencePiece tokenizer [54] to train subwords from the training text. Greedy decoding is used for all considered models to generate the hypothesis text.

Choice of label weights (β_q). To simplify the training setup for our OTTC-based models, we use a fixed and uniform β_q (see Sections 4.2 & 4.3), where the length q of β is equal to the total number of tokens in the text after augmenting with the blank (ϕ) label between repeating characters.

6 Results and Discussion

Alignment quality. We begin by analyzing the alignment performance of the models on the TIMIT and AMI datasets, with results shown in Table 1. Our proposed OTTC model consistently outperforms the CTC-based models across all alignment metrics on both datasets. A key observation is the significant difference in the percentage of frames assigned to non-alphabet symbols by the CTC-based models, highlighting the peaky behavior inherent in these models. Specifically, the baseline CTC-based models tend to assign a large proportion of frames to blank or space symbols, reflecting a misalignment in predicted word boundaries. In contrast, the OTTC model avoids this issue, preventing extreme peaky behavior observed in CTC-based models. While the OTTC model also outperforms the CTC-based models in F1 score, the margin of improvement is smaller. However, the IDR reveals a substantial advantage for OTTC, with a significant improvement over CTC and CR-CTC. This indicates that CTC-based models often either delay the prediction of word starts or assigns too few frames to non-blank symbols, reinforcing the peaky behavior. Additionally, the performance improvement on the AMI dataset is particularly significant, given its nature of meeting speech. This demonstrates how effectively the OTTC loss adapts to varying speaking rates, showcasing the robustness of our framework in learning alignments despite speech variability.

Table 2: Word Error Rate (WER%) comparison between the baseline CTC model and the proposed OTTC model on all considered datasets. Lower WER is better.

Model	TIMIT eval	AMI eval	100h-LibriSpeech		360h-LibriSpeech		960h-LibriSpeech	
			test-clean	test-other	test-clean	test-other	test-clean	test-other
CTC	8.38	11.75	3.36	7.36	2.77	6.58	2.20	5.23
OTTC	8.76	14.27	3.77	8.55	3.00	7.44	2.52	6.16

WER. ASR performance in terms of WER for the CTC model and the proposed OTTC model is depicted in Table 2 for all considered datasets. On the TIMIT dataset, the OTTC model shows a slightly higher WER compared to the CTC model, and while the performance gap is larger on the AMI dataset, it’s encouraging to observe consistent performance despite the varied nature of speech. On the LibriSpeech dataset, using the 100-hour training split, the OTTC model achieves a WER of 3.77% on test-clean. As we scale the training dataset (100h \rightarrow 360h \rightarrow 960h), we observe a monotonic improvement in WER for the proposed OTTC-based models, similarly to the CTC-based models. Although the WERs achieved by the OTTC-based models are typically higher than the CTC-based models, the presented results underscore the experimental validity of the SOTD as a metric and demonstrate that learning a single alignment can yield promising results in E2E ASR.

Qualitative alignment comparison. Apart from quantitative alignment comparison (Table 1), we show an alignment from the CTC- and OTTC-based models in Figure 4. For CTC, it can be seen that the best path aligns most frames to the blank token, resulting in peaky behavior [19]. In contrast, the OTTC model learns to align all frames to non-blank tokens. This effectively mitigates the peaky behavior observed in the CTC model. Note that OTTC allows dropping frames during alignment (see Section 4.4), however, in practice, we observed that only a few frames are dropped. For additional insights, we plot the evolution of the alignment for the OTTC model during the course of training in Figures 6 & 7. It is evident that the alignment learned early in the training process remains relatively stable as training progresses. The most notable changes occur at the extremities of the predicted label clusters. This observation led us to the decision to freeze the OT weight predictions for the final 10 epochs, otherwise, even subtle changes in alignment could adversely impact the logits predictions because same base model is shared for predicting both the logits and the alignment OT weights.

In summary, the presented results demonstrate that the proposed OTTC models achieve significant improvements in alignment performance, effectively mitigating the peaky behavior observed in CTC models. Although there is an increase in WER, the improvement in alignment accuracy indicates better temporal modeling. This enhanced alignment could benefit tasks that require precise timing information, such as speech segmentation, event detection, and applications in the medical domain, where accurate temporal alignment is crucial for tasks like clinical transcription or patient monitoring.

7 Conclusion and Future Work

Learning effective sequence-to-sequence mapping along with its corresponding alignment has diverse applications across various fields. Building upon our core idea of modeling the alignment between two sequences as a learnable mapping while simultaneously predicting the target sequence, we define a pseudo-metric known as the Sequence Optimal Transport Distance (SOTD) over sequences. Our formulation of SOTD enables the joint optimization of target sequence prediction and alignment, which is achieved through one-dimensional optimal transport. We theoretically show that the SOTD indeed defines a distance with guaranteed existence of a solution, though uniqueness is not assured. We then derive the Optimal Temporal Transport Classification (OTTC) loss for ASR where the task is to map acoustic frames to text. Experiments across multiple datasets demonstrate that our method significantly improves alignment performance while successfully avoiding the peaky behavior commonly observed in CTC-based models. Other sequence-to-sequence tasks could be investigated using the proposed framework, particularly those involving the alignment of multiple sequences, such as audio, video, and text.

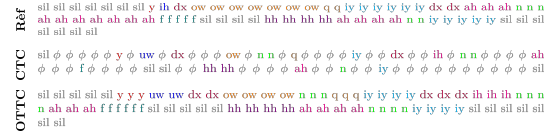


Figure 4: *CTC and OTTC alignments.* Phoneme-level transcription of CTC and OTTC, compared to a reference from TIMIT.

References

- [1] Nelson Morgan and Herve Bourlard. Continuous speech recognition using multilayer perceptrons with hidden markov models. In *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 413–416, Albuquerque, USA, Apr. 1990.
- [2] Herve A. Bourlard and Nelson Morgan. *Connectionist speech recognition: A hybrid approach*, volume 247. Springer Science & Business Media, 2012.
- [3] Steve Young. A review of large-vocabulary continuous-speech. *IEEE Signal Processing Magazine*, 13(5), Sept. 1996.
- [4] Daniel Povey. *Discriminative training for large vocabulary speech recognition*. PhD thesis, University of Cambridge, 2005.
- [5] Ossama Abdel-Hamid, Abdel-rahman Mohamed, Hui Jiang, and Gerald Penn. Applying convolutional neural networks concepts to hybrid NN-HMM model for speech recognition. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4277–4280, Kyoto, Japan, Mar. 2012.
- [6] Alex Graves, Navdeep Jaitly, and Abdel-rahman Mohamed. Hybrid speech recognition with deep bidirectional lstm. In *Proc. IEEE Workshop on Automatic Speech Recognition and Understanding*, pages 273–278, Olomouc, Czech Republic, Dec. 2013.
- [7] George E Dahl, Dong Yu, Li Deng, and Alex Acero. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):30–42, Jan. 2012.
- [8] A Hannun. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014.
- [9] William Chan, Navdeep Jaitly, Quoc V Le, and Oriol Vinyals. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4960–4964, Brisbane, Australia, Apr. 2015.
- [10] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *Proc. International Conference on Machine Learning*, pages 28492–28518, Honolulu, USA, July 2023.
- [11] Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R Hershey, and Tomoki Hayashi. Hybrid CTC/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253, Oct. 2017.
- [12] Rohit Prabhavalkar, Takaaki Hori, Tara N Sainath, Ralf Schlüter, and Shinji Watanabe. End-to-end speech recognition: A survey. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32, Oct. 2023.
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proc. International Conference on Machine learning*, pages 369–376, Pittsburgh, USA, June 2006.
- [14] Alex Graves and Navdeep Jaitly. Towards end-to-end speech recognition with recurrent neural networks. In *Proc. International Conference on Machine Learning*, pages 1764–1772, Beijing, China, June 2014.
- [15] Alex Graves. Sequence transduction with recurrent neural networks. *arXiv preprint arXiv:1211.3711*, 2012.
- [16] Fangjun Kuang, Liyong Guo, Wei Kang, Long Lin, Mingshuang Luo, Zengwei Yao, and Daniel Povey. Pruned RNN-T for fast, memory-efficient ASR training. In *Proc. Annual Conference of the International Speech Communication Association*, pages 2068–2072, Incheon, Korea, Sept. 2022.

- [17] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6645–6649, Vancouver, Canada, May 2013.
- [18] Brian Yan, Siddharth Dalmia, Yosuke Higuchi, Graham Neubig, Florian Metze, Alan W Black, and Shinji Watanabe. CTC alignments improve autoregressive translation. In *Proc. Conference of the European Chapter of the Association for Computational Linguistics*, pages 1623–1639, Dubrovnik, Croatia, May 2022.
- [19] Albert Zeyer, Ralf Schlüter, and Hermann Ney. Why does CTC result in peaky behavior? *arXiv preprint arXiv:2105.14849*, 2021.
- [20] John S Garofolo, Lori F Lamel, William M Fisher, David S Pallett, Nancy L Dahlgren, Victor Zue, and Jonathan G Fiscus. Timit acoustic-phonetic continuous speech corpus. (*No Title*), 1993.
- [21] Jean Carletta et al. The AMI meeting corpus: A pre-announcement. In *Proc. International Workshop on Machine Learning for Multimodal Interaction*, pages 28–39, Edinburgh, UK, July 2005.
- [22] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An ASR corpus based on public domain audio books. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 5206–5210, South Brisbane, Australia, Apr. 2015.
- [23] Yuchen Liu, Junnan Zhu, Jiajun Zhang, and Chengqing Zong. Bridging the modality gap for speech-to-text translation. *arXiv preprint arXiv:2010.14920*, 2020.
- [24] Shun-Po Chuang, Yung-Sung Chuang, Chih-Chiang Chang, and Hung-yi Lee. Investigating the reordering capability in CTC-based non-autoregressive end-to-end speech translation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1068–1077, Aug. 2021.
- [25] Jiatao Gu and Xiang Kong. Fully non-autoregressive neural machine translation: Tricks of the trade. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 120–133, Aug. 2021.
- [26] Alex Graves and Jürgen Schmidhuber. Offline handwriting recognition with multidimensional recurrent neural networks. In *Proc. Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2008.
- [27] Pavlo Molchanov, Xiaodong Yang, Shalini Gupta, Kihwan Kim, Stephen Tyree, and Jan Kautz. Online detection and classification of dynamic hand gestures with recurrent 3d convolutional neural networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, pages 4207–4215, Las Vegas, USA, June 2016.
- [28] Jinchuan Tian, Brian Yan, Jianwei Yu, Chao Weng, Dong Yu, and Shinji Watanabe. Bayes risk CTC: Controllable CTC alignment in sequence-to-sequence tasks. In *Proc. International Conference on Learning Representations*, Kigali, Rwanda, May 2023.
- [29] Haşim Sak, Andrew Senior, Kanishka Rao, Ozan Irsoy, Alex Graves, Françoise Beaufays, and Johan Schalkwyk. Learning acoustic frame labeling for speech recognition with recurrent neural networks. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4280–4284, South Brisbane, Australia, Apr. 2015.
- [30] Zengwei Yao, Wei Kang, Fangjun Kuang, Liyong Guo, Xiaoyu Yang, Yifan Yang, Long Lin, and Daniel Povey. Delay-penalized CTC implemented based on finite state transducer. In *Proc. Annual Conference of the International Speech Communication Association*, pages 1329–1333, Dublin, Ireland, Sept. 2023.
- [31] Yifan Yang, Xiaoyu Yang, Liyong Guo, Zengwei Yao, Wei Kang, Fangjun Kuang, Long Lin, Xie Chen, and Daniel Povey. Blank-regularized CTC for frame skipping in neural transducer. In *Proc. Annual Conference of the International Speech Communication Association*, pages 4409–4413, Dublin, Ireland, Sept. 2023.

- 475 [32] Zeyu Zhao and Peter Bell. Investigating sequence-level normalisation for CTC-Like End-to-End
476 ASR. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*,
477 pages 7792–7796, Singapore, May 2022.
- 478 [33] Théodore Bluche, Hermann Ney, Jérôme Louradour, and Christopher Kermorvant. Framewise
479 and CTC training of neural networks for handwriting recognition. In *Proc. International
480 Conference on Document Analysis and Recognition*, pages 81–85, Nancy, France, Aug. 2015.
- 481 [34] Shahram Ghorbani, Ahmet E. Bulut, and John H.L. Hansen. Advancing multi-accented LSTM-
482 CTC speech recognition using a domain specific student-teacher learning paradigm. In *Proc.
483 IEEE Spoken Language Technology Workshop*, pages 29–35, Athens, Greece, Dec. 2018.
- 484 [35] Gakuto Kurata and Kartik Audhkhasi. Guiding CTC posterior spike timings for improved pos-
485 terior fusion and knowledge distillation. In *Proc. Annual Conference of the International
486 Speech Communication Association*, pages 1616–1620, Graz, Austria, Sept. 2019.
- 487 [36] Albert Zeyer, André Merboldt, Ralf Schlüter, and Hermann Ney. A new training pipeline
488 for an improved neural transducer. In *Proc. Annual Conference of the International Speech
489 Communication Association*, pages 2812–2816, Shanghai, China, Sept. 2020.
- 490 [37] Andrew Senior, Haşim Sak, Félix de Chaumont Quitry, Tara Sainath, and Kanishka Rao.
491 Acoustic modelling with CD-CTC-SMBR LSTM RNNS. In *Proc. IEEE Workshop on Automatic
492 Speech Recognition and Understanding*, pages 604–609, Scottsdale, USA, Dec. 2015.
- 493 [38] Peter Plantinga and Eric Fosler-Lussier. Towards real-time mispronunciation detection in kids’
494 speech. In *Proc. IEEE Automatic Speech Recognition and Understanding Workshop*, pages
495 690–696, Singapore, Dec. 2019.
- 496 [39] Zengwei Yao, Wei Kang, Xiaoyu Yang, Fangjun Kuang, Liyong Guo, Han Zhu, Zengrui Jin,
497 Zhaoqing Li, Long Lin, and Daniel Povey. Cr-ctc: Consistency regularization on ctc for
498 improved speech recognition. *arXiv preprint arXiv:2410.05101*, 2024.
- 499 [40] Jiahui Yu et al. FastEmit: Low-latency streaming ASR with sequence-level emission regular-
500 ization. In *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*,
501 pages 6004–6008, Toronto, Canada, June 2021.
- 502 [41] Wei Kang, Zengwei Yao, Fangjun Kuang, Liyong Guo, Xiaoyu Yang, Long Lin, Piotr Żelasko,
503 and Daniel Povey. Delay-penalized transducer for low-latency streaming ASR. In *Proc. IEEE
504 International Conference on Acoustics, Speech and Signal Processing*, Rhodes Island, Greece,
505 June 2023.
- 506 [42] Yusuke Shinohara and Shinji Watanabe. Minimum latency training of sequence transducers
507 for streaming end-to-end speech recognition. In *Proc. Annual Conference of the International
508 Speech Communication Association*, pages 2098–2102, Incheon, Korea, Sept. 2022.
- 509 [43] Tian-Hao Zhang, Dinghao Zhou, Guiping Zhon, and Baoxiang Li. A novel CIF-based transducer
510 architecture for automatic speech recognition. In *Proc. IEEE International Conference on
511 Acoustics, Speech and Signal Processing*, Seoul, Republic of Korea, Apr. 2024.
- 512 [44] Jaeyoung Kim, Han Lu, Anshuman Tripathi, Qian Zhang, and Hasim Sak. Reducing streaming
513 ASR model delay with self alignment. pages 3440–3444, Aug. 2021.
- 514 [45] Genshun Wan, Mengzhi Wang, Tingzhi Mao, Hang Chen, and Zhongfu Ye. Lightweight
515 transducer based on frame-level criterion. In *Proc. Annual Conference of the International
516 Speech Communication Association*, pages 247–251, Kos, Greece, Sept. 2024.
- 517 [46] Gabriel Peyré and Marco Cuturi. Computational optimal transport: With applications to data
518 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- 519 [47] Gabriel Peyré. Numerical optimal transport and its applications. 2019.
- 520 [48] Fumitada Itakura. Minimum prediction residual principle applied to speech recognition. *IEEE
521 Transactions on Acoustics, Speech, and Signal Processing*, 23:154–158, Jan. 1975.

- 522 [49] Elena Rastorgueva, Vitaly Lavrukhin, and Boris Ginsburg. Nemo forced aligner and its appli-
523 cation to word alignment for subtitle generation. In *INTERSPEECH 2023*, pages 5257–5258,
524 2023.
- 525 [50] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0:
526 A framework for self-supervised learning of speech representations. *Advances in neural*
527 *information processing systems*, 33:12449–12460, 2020.
- 528 [51] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs). *arXiv preprint*
529 *arXiv:1606.08415*, 2016.
- 530 [52] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. International*
531 *Conference on Learning Representations*, New Orleans, USA, May 2019.
- 532 [53] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words
533 with subword units. In *Proc. Annual Meeting of the Association for Computational Linguistics*
534 *(Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, Aug. 2016.
- 535 [54] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword
536 tokenizer and detokenizer for neural text processing. In *Conference on Empirical Methods in*
537 *Natural Language Processing*, Brussels, Belgium, Oct. 2018.
- 538 [55] Marco Cuturi and Mathieu Blondel. Soft-DTW: A differentiable loss function for time-series.
539 In *Proc. International Conference on Machine Learning*, Sydney, Australia, Aug. 2018.
- 540 [56] Sanjay Haresh, Sateesh Kumar, Huseyin Coskun, Shahram N. Syed, Andrey Konin, M. Zeeshan
541 Zia, and Quoc-Huy Tran. Learning by aligning videos in time. In *Proc. IEEE/CVF Conference*
542 *on Computer Vision and Pattern Recognition*, pages 5544–5554, Nashville, USA, Nov. 2021.
- 543 [57] Amit Meghanani and Thomas Hain. LASER: Learning by aligning self-supervised repre-
544 sentations of speech for improving content-related tasks. In *Proc. Annual Conference of the*
545 *International Speech Communication Association*, Kos, Greece, Sept. 2024.
- 546 [58] Titouan Vayer, Romain Tavenard, Laetitia Chapel, Nicolas Courty, Rémi Flamary, and Yann
547 Soullard. Time series alignment with global invariances. *Transactions on Machine Learning*
548 *Research*, Oct. 2022.
- 549 [59] Feng Zhou and Fernando De la Torre. Canonical time warping for alignment of human behavior.
550 In *Proc. Neural Information Processing Systems*, Vancouver, Canada, Dec. 2009.

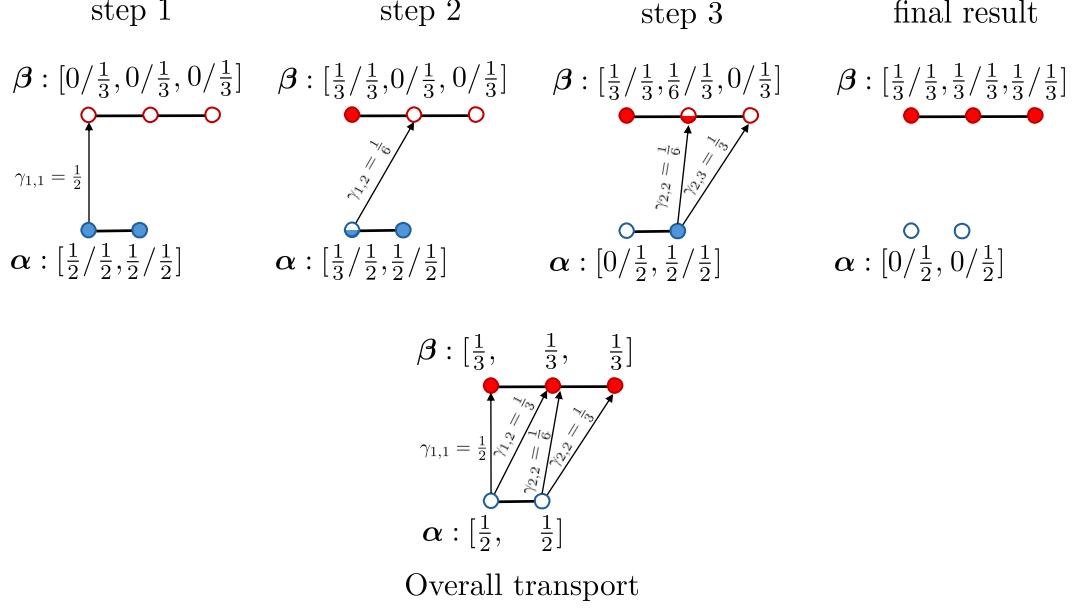


Figure 5: **1D OT transport computation.** Illustration of the optimal transport process, computed iteratively by transferring probability mass from the smallest bins to the largest.

A Appendix

A.1 Algorithm and Implementation Details

A.1.1 Alignment Computation

The algorithm to compute $\gamma_n^{m,\beta}$ is given in Algorithm 1. This algorithm computes the 1D optimal transport between $\mu[\alpha, n]$ and $\nu[\beta, m]$, exploiting the monotonicity of transport in this dimension. To do so the first step consist in sorting the bins which has the complexity $O(n \log n) + O(m \log m) = O(\max(n, m) \log \max(n, m))$. Then we transfer the probability mass from one distribution to another, moving from the smallest bins to the largest. A useful way to visualize this process is by imagining that the bins of μ each contain a pot with a volume of a_i filled with water, while the bins of ν each contain an empty pot with a volume of b_j . The goal is to fill the empty pots of ν using the water from the pots of μ . At any given step of the process, we always transfer water from the smallest non-empty pot of μ to the smallest non-full pot of ν . The volume of water transferred from i to j is denoted by $\gamma_{i,j}$. An example of this process is provided in Figure 5.

In the worst case, this process requires $O(n + m)$ comparisons. However, since the bins are already sorted in SOTD, the overall complexity remains $O(n + m) = O(\max(n, m))$. In practice, this algorithm is not directly used in this work, as we never compute optimal transport solely; it is provided here to illustrate that the dependencies of $\gamma_n^{m,\beta}$ on α are explicit, making it differentiable with respect to α . An efficient batched implementation version for computing SOTD will be released soon.

A.2 Properties of OTTC

Here can be found proof and more insight about the properties of SOTD, \mathcal{S}_r .

A.2.1 Lemma 1 : Bijectivity

Proof of Lemma 1. *Surjectivity:* The surjectivity come from definition of $\Gamma^{*,\beta}[n]$. *Injectivity:* Suppose $\gamma_n^{m,\beta}(\alpha) = \gamma_n^{m,\beta}(\sigma)$, so $\alpha = [\sum_{j=1}^m \gamma_n^{m,\beta}(\alpha)_{i,j}, \dots, \sum_{j=1}^m \gamma_n^{m,\beta}(\alpha)_{i,j}]^T =$

Algorithm 1 : Transport Computation - $\gamma_n^{m,\beta}(\alpha)$

Ensure: Compute $\gamma_n^{m,\beta}(\alpha)$.

Require: $\alpha \in \mathbb{R}^n$.

Set $\gamma \in \mathbb{R}^{n \times m} = \mathbf{0}_{n \times m}$.

Set $i, j = 0$.

while $T == \text{True}$ **do**

if $\alpha_i < \beta_j$ **then**

$\gamma_{i,j} = \beta_j - \alpha_i$

$i = i + 1$

if $i == n$ **then**

$T = \text{false}$

end if

$\beta_j = \beta_j - \alpha_i$

else

$\gamma_{i,j} = \alpha_i - \beta_j$

$j = j + 1$

if $j == m$ **then**

$T = \text{false}$

end if

$\alpha_i = \alpha_i - \beta_j$

end if

end while

Return $\gamma = 0$

575 $[\sum_{j=1}^m \gamma_n^{m,\beta}(\sigma)_{i,j}, \dots, \sum_{j=1}^m \gamma_n^{m,\beta}(\sigma)_{i,j}]^T = \sigma$ (because $\gamma_n^{m,\beta}(\alpha) \in \Gamma^{\alpha,\beta}$ and $\gamma_n^{m,\beta}(\sigma) \in$
576 $\Gamma^{\sigma,\beta}$), which conclude the proof.

577 **A.2.2 Proposition 1 : Discrete Monotonic Alignment Approximation Equivalence.**

578 **Proof of proposition 1.** Let's consider the following proposition $P(k)$:

$$P(k) : \exists \alpha^i \in \Delta^n, \forall i, \forall j \leq k, (i, j) \in \mathbf{A} \iff \gamma_n^{m,\beta}(\alpha^i)_{i,j} > 0. \quad (14)$$

579 **Initialisation** - $P(1)$. $P(1)$ is true. Consider the set $E_1 = \{j \in \llbracket 1, m \rrbracket \mid (1, j) \in \mathbf{A}\}$, which
580 can be written as $E_1 = \{1, 2, \dots, \max(E_1)\}$ since \mathbf{A} is a discrete monotonic alignment. Define
581 $\alpha^1 = [\sum_{j \in E_1} \beta_j, \dots]^T$, where the remaining coefficients are chosen to sum to 1.

582 Since the alignment $\gamma_n^{m,\beta}$ is computed monotonically (see Appendix A.1.1), $\gamma_n^{m,\beta}(\alpha^1)_{1,j} > 0$ if and
583 only if $\alpha_1^1 \leq \beta_1 + \dots + \beta_j$, which corresponds exactly to the set of indices $j \in E_1$, *i.e.*, the aligned
584 indices in \mathbf{A} . This proves $P(1)$.

585 **Heredity** - $P(k) \Rightarrow P(k+1)$. The proof follows similarly to $P(1)$. However two cases need to be
586 considered :

- 587 • When $(k+1, \max(E_k)) \in \mathbf{A}$, in this cases we must consider $E_{k+1} = \{j \in \llbracket 1, m \rrbracket \mid (k+1, j) \in \mathbf{A}\} = \{\max(E_k) = \min(E_{k+1}), \min(E_{k+1})+1, \dots, \max(E_{k+1})\}$ (because β has
588 no components) and define $\alpha^{k+1} = [\alpha_1^1, \dots, \alpha_k^k - \frac{\beta_{\max(E_k)}}{2}, \sum_{j \in E_{k+1}} \beta_j - \frac{\beta_{\max(E_k)}}{2}, \dots]^T$,
589 where the remaining parameters are chosen to sum to 1.
- 590 • When $(k+1, \max(E_k)) \notin \mathbf{A}$, we must consider $E_{k+1} = \{j \in \llbracket 1, m \rrbracket \mid (k+1, j) \in \mathbf{A}\} = \{\max(E_k) \neq \min(E_{k+1}), \min(E_{k+1}) + 1, \dots, \max(E_{k+1})\}$ (because β has no
591 components) and define $\alpha^{k+1} = [\alpha_1^1, \dots, \alpha_k^k, \sum_{j \in E_{k+1}} \beta_j, \dots]^T$, where the remaining
592 parameters are chosen to sum to 1.

By induction, the proposition holds for all n . Therefore, Proposition 1 (*i.e.*, $P(n)$) is true. An α verifying the condition is :

$$\alpha = [\alpha_1^1, \dots, \alpha_n^n]^T$$

595 A.2.3 Proposition 2 : Validity of SOTD definition

596 **Proof of proposition 2.** Since $\gamma_n^{m,\beta}$ is differentiable so continuous, it follows that $\alpha \mapsto$
 597 $\sum_{i,j=1}^{n,m} \gamma_n^{m,\beta}(\alpha)_{i,j} \cdot C(x_i, y_j)$ is continuous over Δ^n . Given that Δ^n is a compact set and ev-
 598 ery continuous function on a compact space is bounded and attains its bounds, the existence of an
 599 optimal solution α^* follows.

600 **Non-unicity of the solution.** The non unicity come from that if there is a solution α^* and two integer
 601 k, l such that $\gamma_n^{m,\beta}(\alpha^*)_{k,l} \geq \epsilon > 0$ and $\gamma_n^{m,\beta}(\alpha^*)_{k+1,l} \geq \epsilon > 0$ and $C(x_k, y_l) = C(x_{k+1}, y_l)$,
 602 therefore the transport $\hat{\gamma}$ such that :

- 603 • $\forall i \in \llbracket 1, n \rrbracket, j \in \llbracket 1, m \rrbracket, (i, j) \neq (k, l), \hat{\gamma}_{i,j} = \gamma_n^{m,\beta}(\alpha^*)_{i,j}.$
- 604 • $\hat{\gamma}_{k,l} = \gamma_n^{m,\beta}(\alpha^*)_{k,l} - \epsilon/2$
- 605 • $\hat{\gamma}_{k+1,l} = \gamma_n^{m,\beta}(\alpha^*)_{k+1,l} + \epsilon/2$

606 provide a distinct solution. Let's denote $\sigma = \{\gamma_n^{m,\beta}\}^{-1}(\hat{\gamma}_{i,j})$. First $\sigma \neq \alpha$ because $\sigma_k =$
 607 $\sum_{l=1}^m \hat{\gamma}_{k,l} = \sum_{l=1}^m \gamma_n^{m,\beta}(\alpha^*)_{k,l} - \epsilon/2 = \alpha_k^* - \epsilon/2$. Second, it's clear that $\sum_{i,j=1}^{n,m} \gamma_n^{m,\beta}(\alpha^*)_{i,j} \cdot$
 608 $C(x_i, y_j) = \sum_{i,j=1}^{n,m} \gamma_n^{m,\beta}(\sigma)_{i,j} \cdot C(x_i, y_j)$. Then σ is distinct solution.

609 A.2.4 Proposition 3 : SOTD is a pseudo Metric

610 **Proof of proposition 3. Pseudo-separation.** It's clear that $\mathcal{S}_r(\{x\}_n, \{x\}_n) = 0$, this value is
 611 attained for $\alpha^* = \beta_n$; where the corresponding alignment $\gamma_n^{n,\beta_n}(\alpha^*)$ corresponds to a one-to-one
 612 alignment. Since the two sequences are identical, all the costs are zero.

613 **Symmetry.** We have $\mathcal{S}_r(\{x\}_n, \{y\}_m) = \mathcal{S}_r(\{y\}_m, \{x\}_n)$ because the expression for \mathcal{S}_r in Eq. 6
 614 is symmetric. Specifically, because C is symmetric as it is a metric.

615 **Triangular inequality.** Consider three sequences $\{x\}_n, \{y\}_m$ and $\{z\}_o$. Let $p = \max(n, m)$,
 616 $q = \min(n, m), u = \max(m, o), v = \min(m, o)$. Define the optimal alignments $\gamma_p^{q,\beta_q}(\alpha^*)$ between
 617 $\{x\}_n$ and $\{y\}_m$; and $\gamma_u^{v,\beta_v}(\rho^*)$ between $\{y\}_m$ and $\{z\}_o$. $\forall i \in \llbracket 1, n \rrbracket, \forall j, k \in \llbracket 1, m \rrbracket, \forall l \in \llbracket 1, o \rrbracket$,
 618 we define :

$$\gamma_{i,j}^{xy} = \begin{cases} \gamma_p^{q,\beta_q}(\alpha^*)_{i,j} & \text{if } n \geq m \\ \gamma_p^{q,\beta_q}(\alpha^*)_{j,i} & \text{otherwise.} \end{cases} \quad (15)$$

$$\gamma_{k,l}^{yz} = \begin{cases} \gamma_u^{v,\beta_v}(\rho^*)_{k,l} & \text{if } k \geq l \\ \gamma_u^{v,\beta_v}(\rho^*)_{l,k} & \text{otherwise.} \end{cases} \quad (16)$$

$$\gamma_{j,k}^{yy} = \gamma_p^{q,\sigma^*}(\beta_q)_{j,k} \quad (17)$$

619 and we define :

$$b_j = \begin{cases} \sum_{i=1}^n \gamma_{i,j}^{xy} & \text{if } > 0 \\ 1 & \text{otherwise.} \end{cases} \quad (18)$$

$$c_k = \begin{cases} \sum_{l=1}^o \gamma_{k,l}^{yz} & \text{if } > 0 \\ 1 & \text{otherwise.} \end{cases} \quad (19)$$

620 So γ^{xy} is the optimal transport between $\mu[\alpha^*, p]$ and $\nu[\beta_q, q]$; γ^{yy} is the optimal transport between
 621 $\mu[\beta_q, q]$ and $\nu[\sigma^*, u]$ and γ^{yz} is the optimal transport between $\mu[\sigma^*, u]$ and $\nu[\beta_v, v]$, since in 1D
 622 optimal transport can be composed, the composition $\frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k}$ is an optimal transport between
 623 $\mu[\alpha^*, p]$ and $\nu[\beta_v, v]$. Therefore by bijectivity of $\gamma_{\max(p,v)}^{\min(p,v), \beta_{\min(p,v)}}$, there is a $\theta \in \mathbb{R}^{\max(p,v)}$ such
 624 that :

$$\gamma_{\max(p,v)}^{\min(p,v),\beta_{\min(p,v)}}(\theta) = \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \quad (20)$$

625 Thus, by the definition of $\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o)$:

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \gamma_{\max(p,v)}^{\min(p,v),\beta_{\min(p,v)}}(\theta) \cdot C(\mathbf{x}_i, \mathbf{z}_l)^r \right)^{1/r} \quad (21)$$

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \cdot C(\mathbf{x}_i, \mathbf{z}_l)^r \right)^{1/r} \quad (22)$$

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \cdot (C(\mathbf{x}_i, \mathbf{y}_j) + C(\mathbf{y}_j, \mathbf{y}_k) + C(\mathbf{y}_k, \mathbf{z}_l))^r \right)^{1/r} \quad (23)$$

626 Applying the Minkowski inequality:

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \cdot (C(\mathbf{x}_i, \mathbf{y}_j))^r \right)^{1/r} + \quad (24)$$

$$\left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \cdot (C(\mathbf{y}_j, \mathbf{y}_k))^r \right)^{1/r} + \quad (25)$$

$$\left(\sum_{i,l=1}^{n,o} \sum_{j,k=1}^{m,m} \frac{\gamma_{i,j}^{xy} \gamma_{j,k}^{yy} \gamma_{k,l}^{yz}}{b_j c_k} \cdot (C(\mathbf{y}_k, \mathbf{z}_l))^r \right)^{1/r} \quad (26)$$

627 Then :

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \left(\sum_{i,j=1}^{n,m} \gamma_{i,j}^{xy} \cdot C(\mathbf{x}_i, \mathbf{y}_j)^r \right)^{1/r} + \quad (27)$$

$$\left(\sum_{j,k=1}^{m,m} \gamma_{j,k}^{yy} \cdot C(\mathbf{y}_j, \mathbf{y}_k)^r \right)^{1/r} + \quad (28)$$

$$\left(\sum_{k,l=1}^{m,o} \gamma_{k,l}^{yz} \cdot C(\mathbf{y}_k, \mathbf{z}_l)^r \right)^{1/r} \quad (29)$$

628 By definition :

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) + \mathcal{S}_r(\{\mathbf{y}\}_m, \{\mathbf{y}\}_m) + \mathcal{S}_r(\{\mathbf{y}\}_m, \{\mathbf{z}\}_o) \quad (30)$$

629 So finally since $\mathcal{S}_r(\{\mathbf{y}\}_m, \{\mathbf{y}\}_m) = 0$, the triangular inequality holds :

$$\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{z}\}_o) \leq \mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) + \mathcal{S}_r(\{\mathbf{y}\}_m, \{\mathbf{z}\}_o). \quad (31)$$

630 This concludes the proof.

631 **Note:** If β 's depends on $\{\mathbf{x}\}_n$, $\{\mathbf{y}\}_m$ and $\{\mathbf{z}\}_m$, we need to introduce the appropriate γ^{zz} to
632 construct the composition in Equation 20, ensuring the proof remains valid.

633 A.2.5 Proposition 4 : Non-separation condition

634 *Proof.* Suppose $\mathcal{S}_r(\{\mathbf{x}\}_n, \{\mathbf{y}\}_m) = 0$, and $\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)) \neq \mathcal{A}(\{\mathbf{y}\}_n)$. So :

$$\sum_{i,j=1}^{n,m} \gamma_n^{m,\beta}(\alpha^*)_{i,j} \cdot C(\mathbf{x}_i, \mathbf{y}_j)^r = 0 \quad (32)$$

635 Let $\mathcal{A}_{\{\mathbf{x}\}_n}$ denote the aggregation operator on Δ^n , which groups indices where consecutive elements
 636 in $\{\mathbf{x}\}_n$ are identical (i.e., $\mathcal{A}([\dots, \alpha_i, \dots, \alpha_{i+k}, \dots]^T) = [\dots, \alpha_i + \dots + \alpha_{i+k}, \dots]^T$ iff $\mathbf{x}_i =$
 637 $\dots = \mathbf{x}_{i+k}$). By expanding the right term, we show that; $\forall \alpha \in \Delta^n$:

$$\sum_{i,j=1}^{n,m} \gamma_n^{m,\beta}(\alpha)_{i,j} \cdot C(\mathbf{x}_i, \mathbf{y}_j)^r = \sum_{i,j=1}^{n,m} \gamma_n^{m,\mathcal{A}_{\{\mathbf{y}\}_m}(\beta)}(\mathcal{A}_{\{\mathbf{x}\}_n}(\alpha))_{i,j} \cdot C(\mathcal{A}(\mathcal{P}_{\alpha}(\{\mathbf{x}\}_n)), \mathcal{A}(\{\mathbf{y}\}_n))^r \quad (33)$$

638 Therefore :

$$\sum_{i,j=1}^{n,m} \gamma_n^{m,\mathcal{A}_{\{\mathbf{y}\}_m}(\beta)}(\mathcal{A}_{\mathcal{P}_{\alpha}(\{\mathbf{x}\}_n)}(\alpha^*))_{i,j} \cdot C(\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)), \mathcal{A}(\{\mathbf{y}\}_n))^r = 0 \quad (34)$$

639 Since $\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)) \neq \mathcal{A}(\{\mathbf{y}\}_n)$ there is a $k \in \llbracket 1, m \rrbracket$ such that :

$$\forall k' < k, \mathcal{A}(\{\mathbf{x}\}_n)_{k'} = \mathcal{A}(\{\mathbf{y}\}_n)_{k'} \quad \text{and} \quad \mathcal{A}(\{\mathbf{x}\}_n)_k \neq \mathcal{A}(\{\mathbf{y}\}_n)_k \quad (35)$$

640 Because the optimal alignment is monotonous and lead to a 0 cost, necessarily :

$$\forall k' < k, \mathcal{A}_{\mathcal{P}_{\alpha}(\{\mathbf{x}\}_n)}(\alpha^*)_{k'} = \mathcal{A}_{\{\mathbf{y}\}_m}(\beta)_{k'} \quad (36)$$

641 which is the only way to have alignment between the k first element which led to 0 cost. Because
 642 of the monotonicity of $\gamma_n^{m,\mathcal{A}_{\{\mathbf{y}\}_m}(\beta)}(\mathcal{A}_{\mathcal{P}_{\alpha}(\{\mathbf{x}\}_n)}(\alpha^*))$ the next alignment (s, t) is between the next
 643 element with a non zeros weights for both sequences. Since β has non zero component and by the
 644 definition of \mathcal{P}_{α} , $s = k$ and $t = k$. Therefore the term $\gamma_n^{m,\mathcal{A}_{\{\mathbf{y}\}_m}(\beta)}(\mathcal{A}_{\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)}(\alpha^*))_{k,k}$ is non
 645 null and the term :

$$\gamma_n^{m,\mathcal{A}_{\{\mathbf{y}\}_m}(\beta)}(\mathcal{A}_{\mathcal{P}_{\alpha}(\{\mathbf{x}\}_n)}(\alpha^*)) C(\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)), \mathcal{A}(\{\mathbf{y}\}_n)_k)$$

646 belong to the sum in depicted in Eq. 34. So $C(\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)), \mathcal{A}(\{\mathbf{y}\}_n)_k) = 0$ i.e.,
 647 $\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)) = \mathcal{A}(\{\mathbf{y}\}_n)_k$ because C is separated. Here a contradiction so we can conclude
 648 that :

$$\mathcal{A}(\mathcal{P}_{\alpha^*}(\{\mathbf{x}\}_n)) = \mathcal{A}(\{\mathbf{y}\}_n)$$

649 .

650 A.3 Supplementary Experimental Insights

651 A.4 Note on Dynamic Time Warping (DTW)

652 It is important to highlight the distinction between our approach and DTW-based [48] alignment
 653 methods, particularly the differentiable variations such as soft-DTW [55]. These methods generally
 654 have quadratic complexity [55], making them significantly more computationally expensive than ours.
 655 Furthermore, in DTW-based methods, the alignment emerges as a consequence of the sequences

themselves. When the function F is powerful, the model can collapse by generating a sequence $F(\{\mathbf{x}\}_n)$ that induces a trivial alignment [56] (see Appendix A.4.1, where we conducted experiments using soft-DTW for ASR to illustrate this). To mitigate this issue, regularization losses [56, 57] or constraints on the capacity of F [58, 59] are commonly introduced. However, using regularization losses lacks theoretical guarantees and introduces additional hyperparameters. Furthermore, constraining the capacity of F , although more theoretically sound, makes tasks requiring powerful encoders on large datasets impractical. In contrast, our method decouples the computation of the alignment from the transformation function F , offering more flexibility to the model as well as built-in temporal alignment constraints and theoretical guarantees against collapse.

A.4.1 Ablation Studies

This section explores the effects of various design choices and configurations on the performance of the proposed OTTC framework and provides additional insights on its comparison to soft-DTW.

Training with single-path alignment from CTC. A relevant question that arises is whether the gap between the OTTC and CTC models arises from the use of a single alignment in OTTC rather than marginalizing over all possible alignments. To investigate this, we conducted a comparison with a single-path alignment approach. Specifically, we first obtained the best path (forced alignment using the Viterbi algorithm) from a trained CTC-based model on the same dataset. A new model was then trained to learn this single best path using Cross-Entropy. On the 360-hour LibriSpeech setup with Wav2Vec2-large as the pre-trained model, this single-path approach achieved a WER of 7.04% on the test-clean set and 13.03% on the test-other set. In contrast, under the same setup, the OTTC model achieved considerably better results, with a WER of 3.00% on test-clean and 7.44% on test-other (see Table 2). These findings indicate that the OTTC model is effective with learning a single alignment, which may be sufficient for achieving competitive ASR performance.

Fixed OT weights prediction (α). We conducted an additional ablation experiment where we replaced the learnable *OT weight prediction head* with fixed and uniform OT weights (α). This approach removes the model’s ability to search for the best path, assigning instead a frame to the same label during training. Consequently, the model loses the localization of the text-tokens in the audio. For this experiment, we used the 360-hour LibriSpeech setup with Wav2Vec2-large as the pre-trained model. The results show a WER of 3.51% on test-clean, compared to 2.77% for CTC and 3.00% for OTTC with learnable OT weights. On test-other, the WER was 8.24%, compared to 6.58% for CTC and 7.44% for OTTC with learnable OT weights. These results demonstrate that while using fixed OT weights leads to a slight degradation in performance, the localization property is completely lost, highlighting the importance of learnable OT weights for preserving both performance and localization in the OTTC model.

Impact of freezing OT weights prediction head across epochs. In our investigations so far, we arbitrarily selected the number of epochs for which the *OT weights prediction head* (α predictor) remained frozen (see Section 6), as a hyperparameter without any tuning. To further understand its impact, we conducted additional experiments on the 360h-LibriSpeech setup using the Wav2Vec2-large model while freezing the *OT weights prediction head* for the last 5 and 15 epochs. When frozen for the last 5 epochs, we achieve a WER of 3.01%, whereas when frozen for the last 15 epochs, the WER is 3.10%. As shown in the Table 2, freezing the OT head for the last 10 epochs results in a WER of 3.00%. Based on these results, it appears that the model’s performance doesn’t change considerably when the model is trained for a few more epochs after freezing the alignment part of the OTTC model.

Oracle experiment. We believe that the proposed OTTC framework has the potential to outperform CTC models by making β learnable with suitable constraints or by optimizing the choice of static β . To illustrate this potential, we conduct an oracle experiment where we first force-align audio frames and text tokens using a CTC-based model trained on the same data. This alignment is then used to calculate the β values. For example, given the target sentence YES and the best valid path from the Viterbi algorithm ($\phi Y \phi \phi EES$), we re-labeled it to $(\phi Y \phi ES)$ and set $\beta = [1/7, 1/7, 2/7, 2/7, 1/7]$. This approach enabled OTTC to learn a uniform distribution for α , mimicking CTC’s highest probability path. As a result, in both the 100h-LibriSpeech and 360h-LibriSpeech setups, the OTTC model converged much faster and matched the performance of CTC. This experiment underscores the critical role of β , suggesting that a better strategy for its selection or training will lead to further improvements.

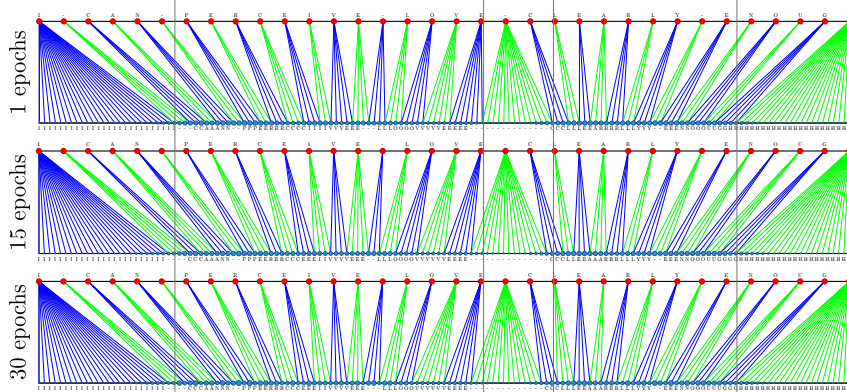


Figure 6: **Evolution of alignment in the OTTC model during the course of training.** The red bullets represent elements of the target sequence $\{y\}_m$, while the blue bullets indicate the predicted OT weights for each frame. The size of the blue bullets is proportional to the predicted OT weight.

Comments on soft-DTW. In soft-DTW, only the first and last elements of sequences are guaranteed to align, while all in-between frames or targets may be ignored; *i.e.*, there is no guarantee that soft-DTW will yield a discrete monotonic alignment. A “powerful” transformation F can map \mathbf{x} to $F(\mathbf{x})$ in such a way that soft-DTW ignores the in-between transformed frames ($F(\mathbf{x})$) and targets (\mathbf{y}), which we refer to as a collapse (Section 4.2.1). This is why transformations learned through sequence comparison are typically constrained (e.g., to geometric transformations like rotations) [58]. Since transformer architectures are powerful, they are susceptible to collapse as demonstrated by the following experiment we conducted using soft-DTW as the loss function. On the 360h-LibriSpeech setup with Wav2Vec2-large model, the best WER achieved using soft-DTW is 39.43%. In comparison, CTC yields 2.77% whereas the proposed OTTC yields 3.00%. A key advantage of our method is that, by construction, such a collapse is not possible.

A.4.2 Alignment Analysis

Temporal evolution of alignment. An example of the evolution of the alignment in the OTTC model during training for 40 epochs without freezing *OT weights prediction head* is shown in Figure 7. Note that during the initial phase of training, there is significant left/right movement of boundary frames for all groups. As training progresses, the movement typically stabilizes to around 1-2 frames. While this can be considered “relatively stable” in terms of alignment, the classification loss (*i.e.*, cross-entropy) in the OTTC framework is still considerably affected by these changes. This change of the loss is what impacts the final performance and the performance difference between freezing or not-freezing the alignments.

B Limitations

The primary limitation of the current work is the observed trade-off between significantly improved alignment quality and a higher WER in ASR tasks compared to CTC. Further research is necessary to bridge this transcription accuracy gap. Additionally, the framework’s performance, particularly the quality of learned alignments and ASR accuracy, can be sensitive to the configuration of label weights (β_q). The current use of fixed, uniform weights is a simplification, and developing strategies to learn β_q or devise more adaptive approaches without encountering degenerate solutions or overly complex training dynamics remains an area for future exploration. Finally, while the SOTD framework and OTTC loss show promise, their empirical validation and necessary adaptations have been primarily focused on ASR, with extensive investigation for a broader range of sequence-to-sequence tasks still required.

C Broader Impacts

Our work has the potential to positively impact several application areas. Improved temporal alignment can benefit domains such as medical speech analysis (e.g., detecting pathological cues), language

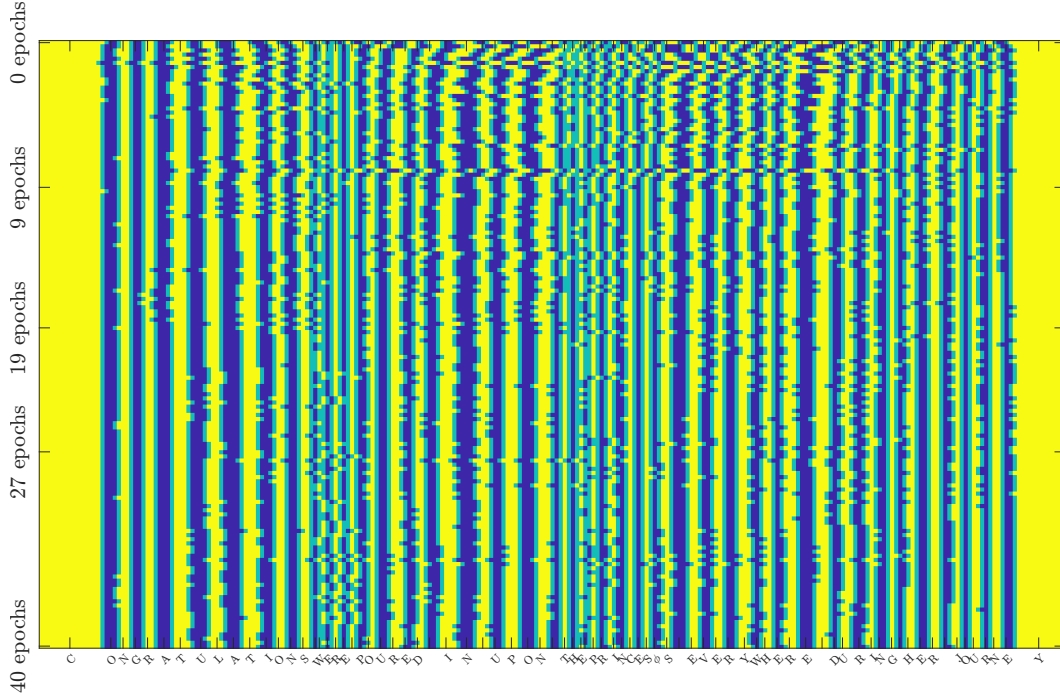


Figure 7: *Alignment evolution in the OTTC model during training for 40 epochs without freezing OT weights prediction head (α predictor).* On the x -axis, each pixel corresponds to one audio frame, while the y -axis represents the epoch. Frames grouped by tokens are shown in alternating colors (yellow and dark blue), with the boundaries of each group highlighted in light blue/green. One can note that during the initial phase of training, there is significant left/right movement of boundary frames for all groups. As training progresses, the movement typically stabilizes to around 1-2 frames.

745 learning tools (e.g., pronunciation feedback), and real-time captioning systems (e.g., enhanced syn-
 746 chronization for accessibility). The proposed methodology also advances sequence modeling by
 747 introducing a more interpretable alignment mechanism.

748 However, responsible deployment remains essential. The current trade-offs in transcription accuracy
 749 must be carefully considered before applying this approach in high-stakes scenarios. Additionally,
 750 as with all ASR technologies, there is a risk of biased performance across different demographic
 751 groups or speaking styles. Future work should address these concerns by incorporating fairness and
 752 robustness considerations. The interpretability gained from a single, learned alignment path may also
 753 support transparency and error analysis.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: Our proposed OTTC loss significantly outperforms CTC and CR-CTC baselines on alignment metrics.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Please see Section B.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: Please see detailed proofs for all propositions in Section A.1.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Please see experiment section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: The code will be open-sourced after publication.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Please see experiment section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Not applicable for relevant experiments on automatic speech recognition.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [No]

Justification: The authors have not provided this information at the current stage.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The paper raises no ethical concerns and complies with the NeurIPS guidelines. No high-risk applications or data are involved, and fairness and privacy considerations are acknowledged where applicable.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The authors include a “Broader Impacts” section that discusses potential benefits and risks in the appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The models released do not pose potential for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All datasets and code used are properly cited, and licenses are referenced. There is no evidence of license violations.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: Any new assets introduced include usage instructions, and documentation appears sufficient for independent use.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: There are no experiments related to crowdsourcing and research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: There is no research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1067 **16. Declaration of LLM usage**
1068 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1069 non-standard component of the core methods in this research? Note that if the LLM is used
1070 only for writing, editing, or formatting purposes and does not impact the core methodology,
1071 scientific rigorousness, or originality of the research, declaration is not required.
1072 Answer: [NA]
1073 Justification: [NA]
1074 Guidelines:
1075 • The answer NA means that the core method development in this research does not
1076 involve LLMs as any important, original, or non-standard components.
1077 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1078 for what should or should not be described.