

# Discriminative Learning in Speech Recognition

Xiaodong He and Li Deng  
{xiaohe, deng}@microsoft.com

October 2007

Technical Report  
MSR-TR-2007-129

Microsoft Research  
Microsoft Corporation  
One Microsoft Way  
Redmond, WA 98052  
<http://www.research.microsoft.com>

## I. INTRODUCTION

Discriminative learning has become a major theme in recent statistical signal processing and pattern recognition research including practically all areas of speech and language processing, e.g., [9][10][13][22][29][35][44][45][48][50]. In particular, much of the striking progress in large scale automatic speech recognition over the past few years has been attributed to the successful development and applications of discriminative learning, e.g., [35][38][48][49]. A key to understanding the speech process is the dynamic characterization of its sequential or variable-length pattern. Two central issues in the development of discriminative learning methods for sequential pattern recognition are: 1) construction of the objective function for optimization; and 2) actual optimization techniques. There have been a wide variety of methods reported in the literature related to both of these issues (e.g., [9][18][22][29][35][38][42][46][50][53][55][59]). However, their relationships have not been adequately understood. Because of the practical and theoretical importance of this problem, there is a pressing need for a unified account of the numerous discriminative learning techniques in the literature. This article is aimed to fulfill this need while providing insights into the discriminative learning framework for sequential pattern classification and recognition. We intend to address the issue of how various discriminative learning techniques are related to and distinguished from each other, and what may be a deeper underlying scheme that can unify various ostensibly different techniques. Although the unifying review provided in this article is on a general class of pattern recognition problems associated with sequential characteristics, we will focus most of the discussions on those related to speech recognition and to the hidden Markov model (HMM) [11][51][57]. We note that the HMM as well as the various forms of discriminative learning have been used in many signal processing related areas beyond speech; e.g., in bioinformatics[6][16], in computational genetics [56], in text and image classification/recognition [33][63][66], in video object classification [64], in natural language processing [8][10], and in tele-robotics [65]. It is our hope that the unifying review and insights provided in this article will foster more principled and successful applications of discriminative learning in a wide range of signal processing disciplines, speech processing or otherwise.

In addition to presenting an extensive account of the basic ideas behind approaches and methods in discriminative learning, in this article we also desire to position our treatment of related algorithms in a wider context of learning and building statistical classifiers from a more general context of machine learning. Generative and discriminative approaches are two main paradigms for designing and learning statistical classifiers/recognition. Generative recognizers rely on a learned model of the joint probability distribution of the observed features and the corresponding class membership. They use this joint-probability model to perform the decision making task based on the posterior probability of the class computed by Bayes rule [12][51][67]. In contrast, discriminative classifiers/recognition directly employ the class posterior probability (or the related discriminant function), exemplified by the argument that “one should solve the (classification/recognition) problem directly and never solve a more general problem as an intermediate step” [58]. This recognizer design philosophy is the basis of a wide range of popular machine learning methods including support vector machine [58], conditional random field [32][45], and maximum entropy Markov models [19][34], etc., where the “intermediate step” of estimating the joint distribution has been avoided. For example, in the recently proposed

structured classification approach [19] [32][34][45] in machine learning and speech recognition, some well known deficiencies of the HMM are addressed by applying “direct” discriminative learning, replacing the need for a probabilistic generative model by a set of flexibly selected, overlapping “features”. Since the conditioning is made on the feature sequence and these “features” can be designed with long-contextual-span properties, the conditional-independence assumption made in the HMM is conceptually alleviated -- provided that proper “features” can be constructed. How to design such features is a challenging research direction and it becomes a critical factor for the potential success of the structured discriminative approach, which departs from the “generative” component or joint distribution. On the other hand, local features can be much more easily designed that are appropriate for the generative approach and many effective local features have been established (e.g., cepstra, filter-bank outputs, etc. [11][51] for speech recognition). Despite the complexity of estimating joint distributions when the sole purpose is discrimination, the generative approach has important advantages of facilitating knowledge incorporation and of conceptually straightforward analyses of recognizer’s components and their interactions.

Analyses of the capabilities and limitations associated with the two general machine learning paradigms discussed above lead to a practical pattern recognition framework being pursued here. That is, we attempt to establish a simplistic joint-distribution or generative model, with the complexity lower than what is required to accurately “generate” samples from the true distribution. In order to make such low-complexity generative models discriminate well, it requires parameter learning methods that are discriminative in nature to overcome the limitation in the simplistic model structure. This is in contrast to the generative approach of fitting the intra-class data as conventional maximum likelihood (ML) based methods intend to accomplish. This type of practical framework has been applied to and guiding much of the recent work in speech recognition research, where HMMs are used as the low-complexity joint distribution for the local acoustic feature sequences of speech and the corresponding underlying linguistic label sequences (sentences, words, or phones). Popular discriminative parameter learning techniques for HMMs are 1) maximum mutual information (MMI) [7][18][21][40][41][42][59][62]; 2) minimum classification error (MCE) [1][9][24][28][29][35][36][38][50][53][55], and 3) minimum phone error (MPE) and closely related minimum word error (MWE) [13][46][47][48][49].

In addition to providing a general overview on the above classes of techniques, this article has a special focus on three key areas in discriminative learning. First, it provides a unifying view of the three major discriminative learning objective functions, MMI, MCE, and MPE/MWE, for classifier parameter optimization, from which insights to the relationships among them are derived. We concentrate on a unified objective function that gives rise to various special cases associated with different levels of performance optimization for pattern recognition tasks --- including performance optimization levels of super-string unit, string unit, and sub-string unit. Second, we describe an efficient approach of parameter estimation in classifier design that unifies the optimization techniques for discriminative learning. This approach for parameter estimation in discriminative learning is based on the optimization framework of growth transformation (GT) (see its detailed introduction in Section IV). We show in a step-by-step fashion that this approach leads to unified parameter estimation formulas and it is scalable for large pattern recognition tasks. The third area is the algorithmic properties of the MCE and MPE/MWE based learning methods under the parameter estimation framework of growth transformation for sequential pattern recognition using HMMs..

The organization of this paper is as follows. In Section II, we provide an introduction to discriminative learning criteria of MMI, MCE and MPE/MWE. In Section III, we show that under certain assumptions, objective functions from MMI, MCE, and MPE/MWE criteria can be

formulated and unified to a rational-function form. From that, relations among MMI, MCE, and MPE/MWE criteria are studied. In Section IV, we provide an overview of the growth transformation (GT) and Extended Baum-Welch (EBW) algorithm based parameter optimization framework. Based on the unified rational-function approach of Section III, we show that the same GT based parameter estimation framework can be applied to objective functions from the discriminative learning criteria of MMI, MCE, and MPE/MWE in sequential pattern recognition. In Section V, we derive the GT based parameter optimization formulation of discriminative learning for the classifier design of discrete HMMs. In Section VI, we present the GT based parameter optimization formulation of discriminative learning for the classifier design of continuous-density HMMs (CDHMMs). In these studies, some familiarities of HMMs are assumed, such as those described in standard textbooks (e.g., [51][57]). Comments and related work are discussed in Section VII and the paper is summarized in Section VIII. For expository purposes, additional technical details are provided in the Appendixes.

## II. DISCRIMINATIVE LEARNING CRITERIA OF MMI, MCE AND MPE/MWE

MMI (maximum mutual information), MCE (minimum classification error), and MPE/MWE (minimum phone error/minimum word error) are the three most popular discriminative learning criteria in speech and language processing, which are the main subject of this paper.

Although the discussion of the discriminative classifier design in this paper has a focus on speech and language processing, they are equally applicable to other similar sequential pattern recognition problems such as handwriting recognition. References made in this paper to words, phones, strings, etc. are for the purpose of showing that the sequential dynamic pattern recognition problem can be based on different levels of recognition units. Moreover, the classifier in sequential pattern recognition can be constructed based on recognizing each pattern (or recognition unit) in isolation. If it can take advantage of the sequential correlation, the classifier can also be constructed based on recognizing a string of patterns (or a string of recognition units), e.g. phrases, word strings, sentences, etc. This flexibility in classifier design for sequential pattern recognition has been a fertile field of research, and many approaches have been developed [23][29][48].

To set up the stage, we denote by  $\Lambda$  the set of classifier parameters that needs to be estimated during the classifier design. For instance in speech and language processing, a (generative) joint distribution of observing a data sequence  $X$  given the corresponding labeled word sequence  $S$  can be written as follows:

$$p(X, S | \Lambda) = p(X | S, \Lambda) P(S) \quad (1)$$

In this notation, it is assumed that the parameters in the “language model”  $P(S)$  are not subject to optimization. Given a set of training data, we denote by  $R$  the total number of training tokens. In this paper, we focus on supervised learning, where each training token consists of an observation data sequence:  $X_r = x_{r,1}, \dots, x_{r,T_r}$ , and its correctly labeled (e.g., word) pattern sequence:  $S_r = W_{r,1}, \dots, W_{r,N_r}$ , with  $W_{r,i}$  being the  $i$ -th word in word sequence  $S_r$ . We use a lower case variable  $s_r$  to denote all possible pattern sequences that can be used to label the  $r$ -th token, including the correctly labeled sequence  $S_r$  and other sequences.

### A. Maximum Mutual Information (MMI)

In the MMI-based classifier design, the goal of classifier parameter estimation is to maximize the mutual information  $I(X, S)$  between data  $X$  and their corresponding labels/symbols  $S$ . From the information theory perspective, mutual information provides a measure of the amount of information gained, or the amount of uncertainty reduced, regarding  $S$  after seeing  $X$ . The MMI

criterion is well established in information theory. It possesses good theoretical properties, and it is different from the criterion of maximum likelihood (ML) used in generative model based learning. Quantitatively, mutual information  $I(X,S)$  is defined as

$$I(X,S) = \sum_{X,S} p(X,S) \log \frac{p(X,S)}{p(X)p(S)} = \sum_{X,S} p(X,S) \log \frac{p(S|X)}{p(S)} = H(S) - H(S|X) \quad (2)$$

where  $H(S) = -\sum_S p(S) \log p(S)$  is the entropy of  $S$ , and  $H(S|X)$  is the conditional entropy

given data  $X$ :  $H(S|X) = -\sum_{X,S} p(X,S) \log p(S|X)$ . When  $p(S|X)$  is based on model  $\Lambda$ , we

have

$$H(S|X) = -\sum_{X,S} p(X,S) \log p(S|X, \Lambda) \quad (3)$$

Assume that the parameters in  $P(S)$  (“language model”) and hence  $H(S)$  is not subject to optimization. Consequently, maximizing mutual information of (2) becomes equivalent to minimizing  $H(S|X)$  of (3) on the training data. When the tokens in the training data are drawn from an i.i.d. distribution,  $H(S|X)$  is given by

$$H(S|X) = -\frac{1}{R} \sum_{r=1}^R \log p(S_r | X_r, \Lambda) = -\frac{1}{R} \sum_{r=1}^R \log \frac{p(X_r, S_r | \Lambda)}{p(X_r)}.$$

Therefore, parameter optimization of MMI based discriminative learning is to maximize the following objective function:

$$O_{MMI}(\Lambda) = \sum_{r=1}^R \log \frac{p(X_r, S_r | \Lambda)}{P(X_r)} = \sum_{r=1}^R \log \frac{p(X_r, S_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)}, \quad (4)$$

where  $P(s_r)$  is the “language model” probability of pattern sequence  $s_r$ .

The objective function  $O_{MMI}$  of (4) is a sum of logarithms. For comparisons with other discriminative training criteria in following sections, we construct the monotonically increasing function of exponentiation for (4). This gives

$$\mathcal{O}_{MMI}(\Lambda) = \exp[O_{MMI}(\Lambda)] = \prod_{r=1}^R \frac{p(X_r, S_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)} \quad (5)$$

It should be noted that  $\mathcal{O}_{MMI}$  and  $O_{MMI}$  have the same set of maximum points, because maximum points are invariant to monotonically increasing transforms. For comparisons with other discriminative training criteria, we rewrite each factor in (5) as

$$\frac{p(X_r, S_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)} = 1 - \sum_{s_r \neq S_r} P(s_r | X_r, \Lambda) = 1 - \sum_{s_r} (1 - \delta(s_r, S_r)) P(s_r | X_r, \Lambda). \quad (6)$$

We define (6) as the model-based expected utility for token  $X_r$ , which equals one minus the model-based expected loss for that token.

### B. Minimum Classification Error (MCE)

The MCE based classifier design is a discriminant-function-based approach to pattern recognition [1][28][29]. The decision rule of the classifier is treated as comparisons among a set of discriminant functions, and parameter estimation involves minimizing the expected loss incurred when these decision rules are applied to the classifier. The loss function in MCE-based discriminative learning is constructed in such a way that the recognition error rate of the classifier

is embedded in a smooth functional form, and minimizing the expected loss of the classifier has a direct relation to classifier error rate reduction.

The objective (loss) function in MCE-based discriminative learning can be constructed from likelihood-based generative models through the following steps. For each training token  $X_r$ , the set of discriminant functions  $\{g_{s_r}\}$  is given as

$$g_{s_r}(X_r; \Lambda) = \log p(X_r, s_r | \Lambda),$$

which is the log joint probability of data  $X_r$  and the pattern sequence (string)  $s_r$  given model  $\Lambda$ . The decision rule of the classifier/recognizer is defined as

$$C(X_r) = s_r^* \quad \text{iff} \quad s_r^* = \arg \max_{s_r} g_{s_r}(X_r; \Lambda). \quad (7)$$

In practice, the  $N$  most confusable competing strings,  $s_{r,1}, \dots, s_{r,N}$ , against the correct string  $S_r$  are considered in MCE-based discriminative learning, where each of these  $N$ -best strings can be defined inductively by

$$\begin{aligned} s_{r,1} &= \arg \max_{s_r: s_r \neq S_r} \log p(X_r, s_r | \Lambda) \\ s_{r,i} &= \arg \max_{s_r: s_r \neq S_r, s_r \neq s_{r,1}, \dots, s_{r,i-1}} \log p(X_r, s_r | \Lambda), \quad i=2, \dots, N \end{aligned} \quad (8)$$

and  $\Lambda$  is the current model parameter set of the classifier.

Then, a misclassification measure  $d_r(X_r, \Lambda)$  can be defined to approximate the performance of the decision rule for each training token  $X_r$ , i.e.,  $d_r(X_r, \Lambda) \geq 0$  implies misclassification and  $d_r(X_r, \Lambda) < 0$  implies correct classification. In particular, such a misclassification measure can be defined by

$$d_r(X_r, \Lambda) = -g_{s_r}(X_r; \Lambda) + G_{S_r}(X_r; \Lambda) \quad (9)$$

where  $G_{S_r}(X_r; \Lambda)$  is a function that represents the scores from the incorrect competing strings and  $g_{s_r}(X_r; \Lambda)$  is the discriminant function for the correct string  $S_r$ .

In the case of 1-best string MCE approach ( $N=1$ ), only the most confusable incorrect string  $s_{r,1}$  is considered as the competitor where  $G_{S_r}(X_r; \Lambda)$  becomes

$$G_{S_r}(X_r; \Lambda) = g_{s_{r,1}}(X_r; \Lambda). \quad (10)$$

However, for the general case where  $N > 1$ , different definitions of  $G_{S_r}(X_r; \Lambda)$  can be used. One popular definition takes the following form [29]:

$$G_{S_r}(X_r; \Lambda) = \log \left\{ \frac{1}{N} \sum_{i=1}^N p^\eta(X_r, s_{r,i} | \Lambda) \right\}^{\frac{1}{\eta}} \quad (11)$$

Another popular form of  $g_{s_r}(X_r; \Lambda)$  and  $G_{S_r}(X_r; \Lambda)$  (the latter has similar effects to (11) and was used in [55]) is

$$\begin{cases} g_{s_r}(X_r; \Lambda) = \log p^\eta(X_r, s_r | \Lambda) \\ G_{S_r}(X_r; \Lambda) = \log \sum_{i=1}^N p^\eta(X_r, s_{r,i} | \Lambda) \end{cases} \quad (12)$$

where  $\eta$  is a scaling factor for joint probability  $p(X_r, s_r | \Lambda)$ . In this paper, we adopt  $G_{S_r}(X_r; \Lambda)$  with the form of (12) and set  $\eta = 1$  for mathematical tractability reasons. (The  $\eta \neq 1$  case will be discussed in Appendix II.)

Given the misclassification measure, the loss function can be defined for each training token  $r$ , and it is usually defined through a sigmoid function as originally proposed in [28][29]:

$$l_r(d_r(X_r, \Lambda)) = \frac{1}{1 + e^{-\alpha d_r(X_r, \Lambda)}} \quad (13)$$

where  $\alpha > 0$  is the slope of the sigmoid function, often determined empirically. As presented in [25] (pp. 156), we also use  $\alpha=1$  for simplifications in the exposition of this paper. (More discussions of  $\alpha$  in empirical studies are included in Appendix II). It should be noted that the loss function of (13) approximates the zero-one classification error count in a smooth functional form.

Given the set of all possible pattern sequences  $\{s_r\} = \{S_r, s_{r,1}, \dots, s_{r,N}\}$  associated with observation data  $X_r$ , and with  $\eta=1$  and  $\alpha=1$ , we substitute (12) into (13) and rewrite the loss function for the training token  $X_r$  as

$$l_r(d_r(X_r, \Lambda)) = \frac{\sum_{s_r, s_r \neq S_r} p(X_r, s_r | \Lambda)}{\sum_{s_r, s_r \neq S_r} p(X_r, s_r | \Lambda) + p(X_r, S_r | \Lambda)} = \frac{\sum_{s_r, s_r \neq S_r} p(X_r, s_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)}. \quad (14)$$

Correspondingly, we can define the utility function as one minus the loss function, i.e.,

$$u_r(d_r(X_r, \Lambda)) = 1 - l_r(d_r(X_r, \Lambda)). \quad (15)$$

Then, the goal in the MCE-based discriminative learning becomes minimization of the expected loss over the entire training data

$$L_{MCE}(\Lambda) = \frac{1}{R} \sum_{r=1}^R l_r(d_r(X_r, \Lambda)). \quad (16)$$

Obviously, minimizing  $L_{MCE}(\Lambda)$  in (16) is equivalent to maximizing the following MCE objective function:

$$O_{MCE}(\Lambda) = R(1 - L_{MCE}(\Lambda)) = \sum_{r=1}^R u_r(d_r(X_r, \Lambda)) = \sum_{r=1}^R \frac{p(X_r, S_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)} \quad (17)$$

It is noteworthy that the summation in (17) for combining utilities of all string tokens for MCE forms a sharp contrast to the MMI case as in (5) where a multiplication of utility functions is constructed for pooling all string tokens.

### C. Minimum “Phone” or “Word” Errors (MPE/MWE)

MPE/MWE is another approach to discriminative learning. It was originally developed in [46][48] and has demonstrated quite effective performance improvement in speech recognition. In contrast to MMI and MCE described earlier that are typically aimed at large segments of pattern sequences (e.g., at string or even super-string level obtained by concatenating multiple pattern strings in sequence), MPE aims at the performance optimization at the sub-string pattern level. In speech recognition, a pattern string usually corresponds to a sentence which consists of a sequence of words, and a sub-string as a constituent of the sentence can be words or phones (subwords).

The MPE objective function that needs to be maximized is defined as

$$O_{MPE}(\Lambda) = \sum_{r=1}^R \frac{\sum_{s_r} p(X_r, s_r | \Lambda) A(s_r, S_r)}{\sum_{s_r} p(X_r, s_r | \Lambda)} \quad (18)$$

where  $A(s_r, S_r)$  is the raw phone (sub-string) accuracy count in the sentence string  $s_r$  (proposed originally in [46][48]). The raw phone accuracy count  $A(s_r, S_r)$  is defined as the total phone (sub-

string) count in the reference string  $S_r$  minus the sum of insertion, deletion and substitution errors of  $s_r$  computed based on  $S_r$ .

The MPE criterion (18) equals the model-based expectation of the raw phone accuracy count over the entire training set. This relation can be seen more clearly by rewriting (18) as

$$O_{MPE}(\Lambda) = \sum_{r=1}^R \sum_{s_r} P(s_r | X_r, \Lambda) A(s_r, S_r)$$

where  $p(s_r | X_r, \Lambda) = \frac{p(X_r, s_r | \Lambda)}{p(X_r | \Lambda)} = \frac{p(X_r, s_r | \Lambda)}{\sum_{s_r} p(X_r, s_r | \Lambda)}$  is the model-based posterior probability.

The concept of the raw phone accuracy count  $A(s_r, S_r)$  in (18) can be generalized to define raw sub-string accuracy count. In particular, the raw word accuracy count  $A_l(s_r, S_r)$  can be defined in the same fashion as the total word (sub-string) count in the reference string  $S_r$  minus the sum of insertion, deletion and substitution errors of  $s_r$  computed based on  $S_r$ . Based on raw word accuracy count  $A_l(s_r, S_r)$ , we have the equivalent definition of the MWE criterion:

$$O_{MWE}(\Lambda) = \sum_{r=1}^R \frac{\sum_{s_r} p(X_r, s_r | \Lambda) A_l(s_r, S_r)}{\sum_{s_r} p(X_r, s_r | \Lambda)} \quad (19)$$

and therefore, in this paper, we merge these two approaches into one MPE/MWE category.

#### D. Discussions

Here we provide brief discussions on key differences among the three criteria presented in this section. At the single-token level, the MMI criterion uses a model-based expected utility of (6) while the MCE criterion uses an classifier-dependent smoothed empirical utility defined by (9), (13), and (15). Likewise, the MPE/MWE criterion also uses a model-based expected utility, but the utility is computed at the sub-string level; e.g., at the phone or word level. We note that for mathematical tractability reasons, in this paper, a specific misclassification measure (12) is used for MCE. As a consequence, the smoothed empirical utility (15) takes the same form as (6) (though they are derived from different motivations). This can be directly seen by substituting (14) to (15).

At the multiple-token level, by comparing (5), (17), (18), and (19), it is clear that MMI training maximizes a *product* of model-based expected utilities of training tokens, while MCE training maximizes a *summation* of smoothed empirical utilities over all training tokens and MPE/MWE training maximizes a *summation* of model-based expected utilities (computed on sub-string units). The difference between the product and the summation forms of the utilities differentiates MMI from MCE/MPE/MWE. This difference causes difficulties in extending the original GT/EBW formulas proposed for MMI to other criteria (See more detailed discussions on this point in Section 7.2, pp. 92 of [48].)

In the following sections, we will show how this difference is reflected in our unified criterion.

### III. THE COMMON RATIONAL-FUNCTION FORM FOR OBJECTIVE FUNCTIONS OF MMI, MCE, AND MPE/MWE

In this section, we show that the objective functions in discriminative learning based on the MMI, MCE and MPE/MWE criteria can be mapped to a canonical rational-function form where the denominator function is constrained to be positive valued. This canonical rational-function form has the benefit of offering insights into the relationships among MMI, MCE, and MPE/MWE based classifiers. In addition, it facilitates the development of a unified classifier parameter

optimization framework for applying MMI, MCE, and MPE/MWE objective functions in sequential pattern recognition tasks.

### A. Rational-Function Form for the Objective Function of MMI

Based on (5), the canonical rational-function form for MMI objective function can be constructed as:

$$\mathcal{O}_{MMI}(\Lambda) = \frac{p(X_1 \dots X_R, S_1 \dots S_R | \Lambda)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} = \frac{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda) C_{MMI}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \quad (20)$$

where

$$C_{MMI}(s_1 \dots s_R) = \prod_{r=1}^R \delta(s_r, S_r) \quad (21)$$

is a quantity that depends only on the sentence sequence  $s_1, \dots, s_R$ , and  $\delta(s_r, S_r)$  is the Kronecker delta function, i.e.,  $\delta(s_r, S_r) = \begin{cases} 1 & \text{if } s_r = S_r \\ 0 & \text{otherwise} \end{cases}$ . In (20), the first step uses the common assumption that different training tokens are independent of each other.

The MMI objective function is aimed at improving the conditional likelihood on the entire training data set instead of on each individual string (token). It can be viewed as a discriminative performance measure at the “super-string” level of all training data  $s_1, \dots, s_R$ , where  $C_{MMI}(s_1, \dots, s_R)$  can be interpreted as the binary function (as “accuracy count”) of the “super-string”  $s_1, \dots, s_R$ , which takes value one if the super-string  $s_1, \dots, s_R$  is correct and zero otherwise.

### B. Rational-Function Form for the Objective Function of MCE

Unlike the MMI case where the rational-function form can be obtained through a simple exponential transformation, the objective function of MCE as given in (17) is a sum of rational functions rather than a rational function in itself (i.e., a ratio of two polynomials). This creates the problem of making the objective function of MCE amenable to the parameter optimization framework of growth transform (GT). Consequently, the objective function of MCE is usually optimized using the generalized probabilistic descent (GPD) [9][28][29] algorithm or other gradient-based methods [37][38]. Despite the popularity and many successful applications, the gradient descent based sequential learning using GPD has two main drawbacks. First, it is a sample-by-sample learning algorithm. Algorithmically, it is difficult for GPD to parallelize the parameter learning process, which is critical for large scale tasks. Second, it is not a monotone learning algorithm and it does not have a monotone learning function to determine the stopping point of the discriminative learning. Recently, applying other batch-mode gradient-based optimization methods, including batch and semi-batch probabilistic descent, Quickprop, and Rprop, to MCE training have been proposed, and improved recognition results are reported [37][38]. However, monotone convergence of these methods has not been established.

In this paper, we take a different approach that makes the objective function for MCE-based discriminative learning directly suitable for GT-based parameter optimization. The scalability and monotone convergence learning properties of GT have the advantage being fast and stable. In order to realize this advantage, we need to re-formulate the MCE objective function and derive a canonical rational-function form for the objective function of MCE. The canonical rational-function form of MCE derived in this process has an additional benefit of unifying the MCE

objective function with MMI and MPE/MWE ones, upon which their differences and similarities can be studied.

The derivation of the rational-function form for the objective function of MCE is as follows:

$$\begin{aligned}
O_{MCE}(\Lambda) &= \sum_{r=1}^R \frac{\sum_{s_r} p(X_r, s_r | \Lambda) \delta(s_r, S_r)}{\sum_{s_r} p(X_r, s_r | \Lambda)} \tag{22} \\
&= \frac{\sum_{s_1} p(X_1, s_1 | \Lambda) \delta(s_1, S_1)}{\sum_{s_1} p(X_1, s_1 | \Lambda)} + \frac{\sum_{s_2} p(X_2, s_2 | \Lambda) \delta(s_2, S_2)}{\sum_{s_2} p(X_2, s_2 | \Lambda)} \\
&\quad + \frac{\sum_{s_3} p(X_3, s_3 | \Lambda) \delta(s_3, S_3)}{\sum_{s_3} p(X_3, s_3 | \Lambda)} + \dots + \frac{\sum_{s_R} p(X_R, s_R | \Lambda) \delta(s_R, S_R)}{\sum_{s_R} p(X_R, s_R | \Lambda)} \\
&= \frac{\sum_{s_1} \sum_{s_2} p(X_1, s_1 | \Lambda) p(X_2, s_2 | \Lambda) [\delta(s_1, S_1) + \delta(s_2, S_2)]}{\sum_{s_1} \sum_{s_2} p(X_1, s_1 | \Lambda) p(X_2, s_2 | \Lambda)} + O_3 + \dots + O_R \\
&= \frac{\sum_{s_1, s_2} p(X_1, X_2, s_1, s_2 | \Lambda) [C_{MCE}(s_1, s_2)]}{\sum_{s_1, s_2} p(X_1, X_2, s_1, s_2 | \Lambda)} + O_3 + \dots + O_R \\
&= \frac{\sum_{s_1, s_2, s_3} p(X_1, X_2, X_3, s_1, s_2, s_3 | \Lambda) [C_{MCE}(s_1, s_2, s_3)]}{\sum_{s_1, s_2, s_3} p(X_1, X_2, X_3, s_1, s_2, s_3 | \Lambda)} + O_4 + \dots + O_R \\
&= \frac{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda) C_{MCE}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \tag{23}
\end{aligned}$$

where  $C_{MCE}(s_1 \dots s_R) = \sum_{r=1}^R \delta(s_r, S_r)$ .  $C_{MCE}(s_1, \dots, s_R)$  can be interpreted as the string accuracy count

for  $s_1, \dots, s_R$ , which takes an integer value between zero and  $R$  as the number of correct strings in  $s_1, \dots, s_R$ . As it will be further elaborated, the rational-function form (23) for the MCE objective function will play a pivotal role in our study of MCE-based discriminative learning.

### C. Rational-Function Form for the Objective Functions of MPE/MWE

Similar to MCE, the MPE/MWE objective function is also a sum of multiple (instead of a single) rational functions, and hence it is difficult to derive GT formulas as discussed in [48] (pp. 92). In order to bypass this issue, a method of optimizing MPE/MWE objective functions based on a heuristic weak-sense auxiliary function (WSAF) was developed in [46][48]. In this paper, we show how to rigorously overcome the above difficulty by re-formulating the MPE/MWE objective function and by deriving a canonical rational-function form for the MPE/MWE objective function. This makes the parameter optimization in MPE/MWE-based discriminative learning directly amendable to the GT-based parameter estimation framework. It provides a unified parameter estimation framework with guaranteed monotone convergence properties

which are lacking in other alternative methods such as gradient-based and WSAF-based approaches.

An important finding is that the same method used to derive the rational-function form (23) for the MCE objective function can be applied directly to derive the rational-function form for MPE/MWE objective functions as defined in (18) and (19). Note that (18) and (19) are in the same form as (22), except that  $\delta(s_r, S_r)$  is replaced by  $A(s_r, S_r)$  or  $A_l(s_r, S_r)$ . Same derivation steps for the objective function of MCE can be applied here and rational-function forms for MPE/MWE are given as follows:

$$O_{MPE}(\Lambda) = \frac{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda) C_{MPE}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \quad (24)$$

where  $C_{MPE}(s_1 \dots s_R) = \sum_{r=1}^R A(s_r, S_r)$ , and

$$O_{MWE}(\Lambda) = \frac{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda) C_{MWE}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \quad (25)$$

where  $C_{MWE}(s_1 \dots s_R) = \sum_{r=1}^R A_l(s_r, S_r)$ .

$C_{MPE}(s_1, \dots, s_R)$  or  $C_{MWE}(s_1, \dots, s_R)$  can be interpreted as the raw phone or word (sub-string unit) accuracy count within the ‘‘super string’’  $s_1, \dots, s_R$ . Its upper-limit value is the total number of phones or words in the full training data (i.e., the correct super-string  $S_1, \dots, S_R$ ). However, the actual value can become negative, e.g., if there are too many insertion errors. Correspondingly,  $O_{MPE}(\Lambda)$  and  $O_{MWE}(\Lambda)$  can be interpreted as the model-based average raw phone or word accuracy count of the full training data set, respectively.

#### D. Comments and Discussions

The main result in this section is that all three discriminative learning objective functions, MMI, MCE, and MPE/MWE, can be formulated in a unified canonical rational-function form as follows:

$$O(\Lambda) = \frac{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda) \cdot C_{DT}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \quad (26)$$

where the summation over  $s=s_1 \dots s_R$  in (26) denotes all possible labeled sequences (both correct and incorrect ones) for all  $R$  training tokens. As it will be further elaborated, this huge number of possible strings can be drastically reduced in practical implementations.

In (26),  $X_1 \dots X_R$  denotes the collection of all observation data sequences (strings) in all  $R$  training tokens, which we also call a ‘‘super string’’ after concatenating them into one single string.  $p_\Lambda(X_1 \dots X_R, s_1 \dots s_R)$  is the joint distribution of the super-string data  $X_1 \dots X_R$  and its possible label sequence  $s_1 \dots s_R$ . MMI, MCE, and MPE/MWE are differentiated in (26) through the criterion-dependent weighting factors  $C_{MMI}(s_1 \dots s_R)$ ,  $C_{MCE}(s_1 \dots s_R)$ , and  $C_{MPE}(s_1 \dots s_R)$ , respectively. An important property is:  $C_{DT}(s_1 \dots s_R)$  is dependent only on the labeled sequence  $s_1 \dots s_R$ , and it is independent of the parameter set  $\Lambda$  to be optimized.

The rational-function formulation (26) for MMI, MCE, and MPE/MWE objective functions serves two main purposes. First, it unifies the objective functions for MMI, MCE, and

MPE/MWE in a canonical rational-function form upon which the relations among different discriminative learning criteria can be studied and their properties be compared. This provides insights into the various approaches in discriminative learning. Second, the unified objective function (26) overcomes the main obstacle for applying the GT-based parameter optimization framework in discriminative learning. It leads to a scalable and common parameter estimation framework for discriminative learning which is highly efficient and has well-founded algorithmic convergence properties. All these properties have been among the major concerns in the past when applying discriminative learning to sequential pattern recognition.

As presented in this section, the key difference in the rational-function form of MMI, MCE, and MPE/MWE criteria is the weighting factor in the numerator of (26), where  $C_{DT}(s_1 \dots s_R)$  as a generic weighting factor depends on what discriminative training (DT) criterion is being applied.

For example, for MMI,  $C_{DT}(s_1 \dots s_R) = \prod_{r=1}^R \delta(s_r, S_r)$ , and for MPE,  $C_{DT}(s_1 \dots s_R) = \sum_{r=1}^R A(s_r, S_r)$ . In

the case of MCE with general  $N$ -best competitors where  $N > 1$ ,  $C_{DT}(s_1 \dots s_R) = \sum_{r=1}^R \delta(s_r, S_r)$ , and for

one-best MCE ( $N=1$ ),  $s_r$  belongs to only the subset  $\{S_r, s_{r,1}\}$ . From the canonical rational-function form (26), direct comparisons can be made on the objective functions of MMI, MCE and MPE/MWE. Table 1 tabulates the relation among these discriminative objective functions. As discussed in [48], MPE/MWE has an important difference from MCE and MMI in that the weighting given by the MPE/MWE criterion to an incorrect string (sentence token) depends on the number of wrong sub-strings (e.g., wrong phones or words) within the string. MCE and MMI make a binary distinction based on whether the entire sentence string is correct or not, which may not be a good fit if the goal is to reduce the sub-string errors (e.g., word errors in speech recognition). This distinction can be clearly seen by comparing the sum of the binary function

$C_{DT}(s_1 \dots s_R) = \sum_{r=1}^R \delta(s_r, S_r)$  for MCE and the sum of non-binary functions

$C_{DT}(s_1 \dots s_R) = \sum_{r=1}^R A(s_r, S_r)$  for MPE/MWE. This key difference gives rise to the distinction of the

sub-string level versus the string level recognition performance optimization in MPE/MWE and MCE. Further, the product instead of summation form of the binary function associated with

MMI, i.e.,  $C_{DT}(s_1 \dots s_R) = \prod_{r=1}^R \delta(s_r, S_r)$ , makes it clear that MMI achieves performance

optimization at the super-string level, e.g., the joint product of Kronecker delta functions becomes zero if any sentence token is incorrect. Therefore, all summation terms in the numerator of (26) are zero except for the one corresponding to the correct label/transcription sequence. This criterion is apparently less desirable than MCE or MPE/MWE, as has been observed extensively in speech recognition experiments [35] [46][47][48].

Objective Functions	$C_{DT}(s_r)$	$C_{DT}(s_1 \dots s_R)$	Label Sequence Set Used in DT
MCE (N-best)	$\delta(s_r, S_r)$	$\sum_{r=1}^R C_{DT}(s_r)$	$\{S_r, s_{r,1}, \dots, s_{r,N}\}$
MCE (one-best)	$\delta(s_r, S_r)$	$\sum_{r=1}^R C_{DT}(s_r)$	$\{S_r, s_{r,1}\}$
MPE	$A(s_r, S_r)$	$\sum_{r=1}^R C_{DT}(s_r)$	all possible label sequences
MWE	$A_l(s_r, S_r)$	$\sum_{r=1}^R C_{DT}(s_r)$	all possible label sequences
MMI	$\delta(s_r, S_r)$	$\prod_{r=1}^R C_{DT}(s_r)$	all possible label sequences

Table 1:  $C_{DT}(s_1 \dots s_R)$  in the unified rational-function form for MMI, MCE, and MPE/MWE objective functions. The set of “competing token candidates” distinguishes  $N$ -best and one-best versions of the MCE. Note that the overall  $C_{DT}(s_1 \dots s_R)$  is constructed from its constituents  $C_{DT}(s_r)$ ’s in individual string tokens by either summation (for MCE, MPE/MWE) or product (for MMI).

Another insight from the unified form of the objective function (26) is that in the special case of having only one sentence token (i.e.,  $R=1$ ) in the training data and when the sentence contains only one phone, then all three MMI, MCE, and MPE/MWE criteria become identical. This is obvious because in this case  $C_{DT}(s_1 \dots s_R)$  becomes identical. The difference surfaces only when the training set consists of multiple sentence tokens. With multiple training tokens, the difference lies mainly in the  $\Lambda$ -independent weighing factor  $C_{DT}(s_1 \dots s_R)$  (as well as in the set of competitor strings) while the general rational-function form (26) for the three criteria remains unchanged.

Although we intend to derive the GT-based parameter optimization framework for the three types of objective functions of MMI, MCE, and MPE/MWE in sequential pattern recognition, it should be noted that the unified objective function (26) can provide a critical foundation to derive other parameter optimization methods in discriminative learning. For example, recently Jebara [26][27] proposed a parameter optimization method for rational functions as an alternative to the GT method. This method is based on the reverse Jensen inequality, upon which an elegant solution for HMMs with exponential-family densities is constructed [26].

#### IV. OPTIMIZING RATIONAL FUNCTIONS BY GROWTH TRANSFORMATION

GT-based parameter optimization refers to a family of batch-mode, iterative optimization schemes that “grow” the value of the objective function upon each iteration. That is, the new set of model parameter  $\Lambda$  is estimated from the current model parameter set  $\Lambda'$  through a transformation  $\Lambda = T(\Lambda')$  with the property that the target objective function “grows” in its value  $O(\Lambda) > O(\Lambda')$  unless  $\Lambda = \Lambda'$ . One particular algorithm of this type of optimization techniques is called Extended Baum-Welch (EBW) algorithm when HMM parameters are estimated. GT/EBW algorithm was initially developed for the homogeneous polynomial by Baum and his colleagues [3][4]. It was later extended to optimizing non-homogeneous rational functions as reported in [18]. EBW algorithm became popular for its successful application in MMI-based discriminative

training of discrete HMMs [18]. It was later extended and applied to MMI-based discriminative training of CDHMMs [2][21][42][60][62].

The importance of GT/EBW algorithm lies in its monotone convergence properties, its algorithmic effectiveness and scalability for parallel execution, and its closed-form parameter updating formulas for large-scale optimization problems. The unified parameter optimization framework of GT also alleviates the need for other heuristics, e.g., tuning the parameter-dependent learning rate as in some other methods [29][53].

Let  $G(\Lambda)$  and  $H(\Lambda)$  be two real-valued functions on the parameter set  $\Lambda$ , and the denominator function  $H(\Lambda)$  is positive valued. The goal of GT based parameter optimization is to find an optimal  $\Lambda$  that maximizes the objective function  $O(\Lambda)$  which is a rational function of the following form:

$$O(\Lambda) = \frac{G(\Lambda)}{H(\Lambda)}. \quad (27)$$

For example,  $O(\Lambda)$  can be one of the rational functions of (20), (23), (24) and (25) for the MMI, MCE, and MPE/MWE objective functions, respectively, or the general rational-function (26). In the general case of (26), we have

$$G(\Lambda) = \sum_s p(X, s | \Lambda) C(s), \text{ and } H(\Lambda) = \sum_s p(X, s | \Lambda) \quad (28)$$

where we use short-hand notation  $s=s_1 \dots s_R$  to denote the labeled sequences of all  $R$  training tokens/sentences, and  $X=X_1 \dots X_R$ , to denote the observation data sequences for all  $R$  training tokens.

#### Primary Auxiliary Function:

As originally proposed in [18], for the objective function (27), the GT-based optimization algorithm will constructs an auxiliary function of the following form:

$$F(\Lambda; \Lambda') = G(\Lambda) - O(\Lambda')H(\Lambda) + D \quad (29)$$

where  $D$  is a quantity independent of the parameter set, and  $\Lambda$  is the model parameter set to be estimated by applying GT to another model parameter set  $\Lambda'$ . The GT algorithm starts from the (initial) parameter set  $\Lambda'$  (e.g., obtained using maximum likelihood training). Then, it updates the parameter set from  $\Lambda'$  to  $\Lambda$  by maximizing the auxiliary function  $F(\Lambda; \Lambda')$ , and the process iterates until convergence is reached. Maximizing the auxiliary function  $F(\Lambda; \Lambda')$  can often be more feasible than directly maximizing the original rational function  $O(\Lambda)$ . The important property of GT-based parameter optimization is that as long as  $D$  is a quantity not relevant to the parameter set  $\Lambda$ , an increase of  $F(\Lambda; \Lambda')$  guarantees an increase of  $O(\Lambda)$ . This can be seen clearly from the following derivation.

Substituting  $\Lambda = \Lambda'$  into (29), we have

$$F(\Lambda'; \Lambda') = \underset{=0}{G(\Lambda') - O(\Lambda')H(\Lambda')} + D = D$$

Hence,

$$\begin{aligned} F(\Lambda; \Lambda') - F(\Lambda'; \Lambda') &= F(\Lambda; \Lambda') - D = G(\Lambda) - O(\Lambda')H(\Lambda) \\ &= H(\Lambda) \left( \frac{G(\Lambda)}{H(\Lambda)} - O(\Lambda') \right) = H(\Lambda)(O(\Lambda) - O(\Lambda')) \end{aligned}$$

Since  $H(\Lambda)$  is positive, we have  $O(\Lambda) - O(\Lambda') > 0$  on the right hand side, as long as  $F(\Lambda; \Lambda') - F(\Lambda'; \Lambda') > 0$  on the left hand side.

Secondary Auxiliary Function:

However,  $F(\Lambda; \Lambda')$  may still be too difficult to optimize directly, and a secondary auxiliary function can be constructed and optimized based on the previous auxiliary function  $F(\Lambda; \Lambda')$ . As proposed in [17], this secondary auxiliary function in GT-based parameter estimation can have the following form:

$$V(\Lambda; \Lambda') = \sum_s \sum_q \sum_{\chi} f(\chi, q, s, \Lambda') \log f(\chi, q, s, \Lambda) \quad (30)$$

where  $f(\chi, q, s, \Lambda)$  is a positive valued function which is constructed with discrete arguments<sup>1</sup> of  $\chi, q, s$  and which is related to the primary auxiliary function  $F(\Lambda; \Lambda')$  according to

$$F(\Lambda; \Lambda') = \sum_s \sum_q \sum_{\chi} f(\chi, q, s, \Lambda) \quad (31)$$

By applying the Jensen's inequality to the concave log function, it is easy to prove (proof omitted here) that an increase in the auxiliary function  $V(\Lambda; \Lambda')$  guarantees an increase in  $\log F(\Lambda; \Lambda')$ . Since logarithm is a monotonically increasing function, this implies an increase of  $F(\Lambda; \Lambda')$  and hence an increase of the original objective function  $O(\Lambda)$ .

**V. DISCRIMINATIVE LEARNING FOR DISCRETE HMMs BASED ON THE GT FRAMEWORK**

The GT/EBW-based discriminative learning for discrete HMMs needs to estimate the model parameters --  $\Lambda = \{\{a_{i,j}\}, \{b_i(k)\}\}$  consisting of the state transition and emitting probabilities. We derive the parameter optimization formula that “grows” the generic discriminative objective function  $O(\Lambda)$  in the form of (26) which covers MMI, MCE and MPE/MWE as special cases. The discriminative function  $O(\Lambda)$  is difficult to optimize directly. However, since it is a rational function, it is amenable to the GT/EBW-based parameter estimation framework. We can construct the auxiliary function  $F$  and then construct the secondary auxiliary function  $V$  based on  $F$ . We describe how to optimize  $V(\Lambda; \Lambda')$ , leading to the GT-based parameter estimation formulas for all three types of discriminative criteria: MMI, MCE, and MPE/MWE. This approach is applicable to any other discriminative criteria as long as the objective functions can be represented in a rational-function form of (26).

For the discrete HMM, the observation space is quantized by some discrete codebook. In this case,  $X = X_1 \dots X_R$  is a concatenation of all training tokens, and each training token  $X_r$  consists of a sequence of discrete indices obtained by mapping the time sequence of observations for  $r$ -th token to a discrete index sequence with each element  $x_{r,t} \in [1, 2, \dots, K]$ , where  $K$  is the size of the codebook index set and  $x_{r,t}$  is the index of the cell that the observation of the  $t$ -th frame in  $r$ -th token is quantized to.

**A. Constructing the Primary Auxiliary Function  $F(\Lambda; \Lambda')$**

Substituting (28) into (29), we obtain the following auxiliary function

---

<sup>1</sup> Examples of the arguments  $\chi, q, s$  are the discrete acoustic observation, the HMM state sequence, and the label sequence, respectively, in a discrete-HMM-based sequential classifier.

$$\begin{aligned}
F(\Lambda; \Lambda') &= \sum_s p(X, s | \Lambda) C(s) - O(\Lambda') \sum_s p(X, s | \Lambda) + D \\
&= \sum_s p(X, s | \Lambda) [C(s) - O(\Lambda')] + D = \sum_s \sum_q p(X, q, s | \Lambda) [C(s) - O(\Lambda')] + D
\end{aligned} \tag{32}$$

where  $q$  is an HMM state sequence, and  $s=s_1 \dots s_R$  is the ‘‘super’’ label sequence for all  $R$  training tokens (including correct or incorrect sentences). The main terms in the auxiliary function  $F(\Lambda; \Lambda')$  above can be interpreted as the average deviation of the accuracy count.

### B. Constructing the Secondary Auxiliary Function $V(\Lambda; \Lambda')$

Since  $p(s)$  depends on the language model and is irrelevant for optimizing  $\Lambda$ , we have  $p(X, q, s | \Lambda) = p(s) \cdot p(X, q | s, \Lambda)$ , and

$$\begin{aligned}
F(\Lambda; \Lambda') &= \sum_s \sum_q [C(s) - O(\Lambda')] p(s) p(X, q | s, \Lambda) + D \\
&= \sum_s \sum_q \sum_{\chi} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda)
\end{aligned} \tag{33}$$

where

$$\Gamma(\Lambda') = \delta(\chi, X) p(s) [C(s) - O(\Lambda')] \tag{34}$$

and  $D = \sum_s d(s)$  is a quantity independent of parameter set  $\Lambda$ . In (34),  $\delta(\chi, X)$  is the Kronecker delta function, where  $\chi$  represents the entire discrete data space where  $X$  belongs. Using ideas in [21], the summation over this data space is introduced here for satisfying the requirement in (29) and (33) that constant  $D$  be parameter independent. That is, in (33),  $\sum_s \sum_q \sum_{\chi} d(s) p(\chi, q | s, \Lambda) = \sum_s d(s) = D$  is a  $\Lambda$ -independent constant. Importantly, while the full sum is  $\Lambda$ -independent, each constituent  $d(s) p(\chi, q | s, \Lambda)$  is a  $\Lambda$ -dependent quantity in order to account for the possibility that the corresponding term  $\Gamma(\Lambda') p(\chi, q | s, \Lambda)$  may be negative. We elaborate this point below.

To construct the secondary auxiliary function for (30) based on function (33), we first identify from (33) that

$$f(\chi, q, s, \Lambda) = [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda)$$

according to (31). To ensure that  $f(\chi, q, s, \Lambda)$  above is positive,  $d(s)$  should be selected to be sufficiently large so that  $\Gamma(\Lambda') + d(s) > 0$  (note  $p(\chi, q | s, \Lambda)$  in (33) is non-negative). Then, using (30), we have

$$\begin{aligned}
V(\Lambda; \Lambda') &= \sum_q \sum_s \sum_{\chi} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') \log \left\{ \underbrace{[\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda)}_{\text{optimization-independent}} \right\} \\
&= \sum_q \sum_s \sum_{\chi} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') \log p(\chi, q | s, \Lambda) + \text{Const.} \\
&= \sum_q \sum_s p(X, q, s | \Lambda') (C(s) - O(\Lambda')) \log p(X, q | s, \Lambda) \\
&\quad + \sum_q \sum_s \sum_{\chi} d(s) p(\chi, q | s, \Lambda') \log p(\chi, q | s, \Lambda) + \text{Const.}
\end{aligned} \tag{35}$$

The auxiliary function (35) is easier to optimize than (33), because the new logarithm  $\log p(X, q | s, \Lambda)$  introduced in (35) (which is absent in (33)) can lead to significant simplification of  $V(\Lambda; \Lambda')$  which we outline below.

### C. Simplifying the Secondary Auxiliary Function $V(\Lambda; \Lambda')$

We first ignore optimization-independent constant in (35), and divide  $V(\Lambda; \Lambda')$  by another optimization-independent quantity,  $p(X/\Lambda')$ , in order to convert the joint probability  $p(X, q, s | \Lambda')$  to the posterior probability  $p(q, s | X, \Lambda') = p(s | X, \Lambda')p(q | X, s, \Lambda')$ . We then obtain an equivalent auxiliary function of

$$U(\Lambda; \Lambda') = \sum_q \sum_s p(s | X, \Lambda') p(q | X, s, \Lambda') (C(s) - O(\Lambda')) \log p(X, q | s, \Lambda) \\ + \sum_q \sum_s \sum_{\chi} d'(s) p(\chi, q | s, \Lambda') \log p(\chi, q | s, \Lambda) \quad (36)$$

$$\text{where } d'(s) = d(s) / p(X | \Lambda'). \quad (37)$$

Since  $X$  depends only on the HMM state sequence  $q$ , we have  $p(X, q | s, \Lambda) = p(q | s, \Lambda) \cdot p(X | q, \Lambda)$ . Therefore,  $U(\Lambda; \Lambda')$  can be further decomposed to four terms as follows:

$$U(\Lambda; \Lambda') = \sum_q \sum_s p(s | X, \Lambda') p(q | X, s, \Lambda') (C(s) - O(\Lambda')) \log p(X | q, \Lambda) \\ + \sum_q \sum_s \sum_{\chi} d'(s) p(\chi, q | s, \Lambda') \log p(\chi | q, \Lambda) \quad \text{term-II} \\ + \sum_q \sum_s p(s | X, \Lambda') p(q | X, s, \Lambda') (C(s) - O(\Lambda')) \log p(q | s, \Lambda) \\ + \sum_q \sum_s \sum_{\chi} d'(s) p(\chi, q | s, \Lambda') \log p(q | s, \Lambda) \quad \text{term-IV} \quad (38)$$

In this case,  $X = X_1 \dots X_R$  aggregates all training data with  $R$  independent sentence tokens. For each token  $X_r = x_{r,1}, \dots, x_{r,T_r}$ , the observation vector  $x_{r,t}$  is independent of each other and it depends only on the HMM state at time  $t$ . Hence,  $\log p(X | q, \Lambda)$  can be decomposed, enabling simplification of both *term-I* and *term-II* in (38). To simplify *term-III* and *term-IV* in (38), we decompose  $\log p(q | s, \Lambda)$  based on the property of the first-order HMM that state at time  $t$  depends only on state at time  $t-1$ . We now elaborate on the simplification of each of these four terms.

For *term-I*, we first define

$$\gamma_{i,r,s_r}(t) = \sum_{q, q_{r,t}=i} p(q | X, s, \Lambda') = p(q_{r,t} = i | X, s, \Lambda') = p(q_{r,t} = i | X_r, s_r, \Lambda') \quad (39)$$

The last equality comes from the fact that sentence tokens in the training set are independent of each other.  $\gamma_{i,r,s_r}(t)$  is the occupation probability of state  $i$  at time  $t$ , given the label sequence  $s_r$ , and observation sequence  $X_r$ , which can be obtained through an efficient forward-backward algorithm [51]. Using the definition of (39) and assuming the HMM state index is from 1 to  $I$ , we have

$$\text{term-I} = \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_q p(q | X, s, \Lambda') \sum_{r=1}^R \sum_{t=1}^{T_r} \log p(x_{r,t} | q_{r,t}, \Lambda)$$

$$\begin{aligned}
&= \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{q, q_{r,t}=i} p(q | X, s, \Lambda') \log p(x_{r,t} | q_{r,t} = i, \Lambda) \\
&= \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \gamma_{i,r,s_r}(t) \log p(x_{r,t} | q_{r,t} = i, \Lambda) \tag{40}
\end{aligned}$$

The simplification process for the second term in (38) is as follows. Using notations  $\mathbb{Q} = q_{1,1}, \dots, q_{r,t-1}, q_{r,t+1}, \dots, q_{R,T_R}$  and  $\mathbb{X} = \mathcal{X}_{1,1}, \dots, \mathcal{X}_{r,t-1}, \mathcal{X}_{r,t+1}, \dots, \mathcal{X}_{R,T_R}$ , we have

$$\begin{aligned}
\text{term-II} &= \sum_s d'(s) \sum_{q_{1,1}, \dots, q_{R,T_R}} \sum_{\mathcal{X}_{1,1}, \dots, \mathcal{X}_{R,T_R}} p(\mathcal{X}_{1,1}, \dots, \mathcal{X}_{R,T_R}, q_{1,1}, \dots, q_{R,T_R} | s, \Lambda') \sum_{r=1}^R \sum_{t=1}^{T_r} \log p(\mathcal{X}_{r,t} | q_{r,t}, \Lambda) \\
&= \sum_s d'(s) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q_{r,t}} \sum_{\mathcal{X}_{r,t}} p(\mathcal{X}_{r,t}, q_{r,t} | s, \Lambda') \sum_{\mathbb{X}} \sum_{\mathbb{Q}} p(\mathbb{X}, \mathbb{Q} | \mathcal{X}_{r,t}, q_{r,t}, s, \Lambda') \log p(\mathcal{X}_{r,t} | q_{r,t}, \Lambda) \\
&= \sum_s d'(s) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q_{r,t}} \sum_{\mathcal{X}_{r,t}} p(\mathcal{X}_{r,t}, q_{r,t} | s, \Lambda') \log p(\mathcal{X}_{r,t} | q_{r,t}, \Lambda) \\
&= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{\mathcal{X}_{r,t}} \sum_s d'(s) p(q_{r,t} = i | s, \Lambda') p(\mathcal{X}_{r,t} | q_{r,t} = i, \Lambda') \log p(\mathcal{X}_{r,t} | q_{r,t} = i, \Lambda) \\
&= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t, i) \sum_{\mathcal{X}_{r,t}} p(\mathcal{X}_{r,t} | q_{r,t} = i; \Lambda') \log p(\mathcal{X}_{r,t} | q_{r,t} = i; \Lambda) \tag{41}
\end{aligned}$$

$$\text{where } d(r, t, i) = \sum_s d'(s) p(q_{r,t} = i | s, \Lambda'). \tag{42}$$

To simplify *term-III* in (38), we first define

$$\xi_{i,j,r,s_r}(t) = \sum_{q: q_{r,t-1}=i, q_{r,t}=j} p(q | X, s, \Lambda') = p(q_{r,t-1} = i, q_{r,t} = j | X, s, \Lambda') = p(q_{r,t-1} = i, q_{r,t} = j | X_r, s_r, \Lambda') \tag{43}$$

which is the posterior probability of staying at state  $i$  at time  $t-1$  and staying at state  $j$  at time  $t$ , given the labeled sequence  $s_r$  and the observation sequence  $X_r$ . This posterior probability can be computed using an efficient forward-backward algorithm [51]. Further,  $p(q | s, \Lambda)$  can be decomposed as follows:

$$p(q | s, \Lambda) = \prod_{r=1}^R p(q_{r,1}, \dots, q_{r,T_r} | s_r, \Lambda) = \prod_{r=1}^R \prod_{t=1}^{T_r} a_{q_{r,t-1}, q_{r,t}}.$$

This leads to the following simplification:

$$\begin{aligned}
\text{term-III} &= \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_q p(q | X, s, \Lambda') \sum_{r=1}^R \sum_{t=1}^{T_r} \log a_{q_{r,t-1}, q_{r,t}} \\
&= \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{j=1}^I \sum_{q, q_{r,t-1}=i, q_{r,t}=j} p(q | X, s, \Lambda') \log a_{i,j} \\
&= \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{j=1}^I \xi_{i,j,r,s_r}(t) \log a_{i,j} \tag{44}
\end{aligned}$$

and

$$\text{term-IV} = \sum_s d'(s) \sum_q \sum_{\mathcal{X}} p(\mathcal{X}, q | s, \Lambda') \sum_{r=1}^R \sum_{t=1}^{T_r} \log a_{q_{r,t-1}, q_{r,t}}$$

$$\begin{aligned}
&= \sum_s d'(s) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_q \sum_{\mathcal{X}} p(\mathcal{X}, q | s, \Lambda') \log a_{q_{r,t-1}, q_{r,t}} \\
&= \sum_s d'(s) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{q_{r,t-1}} \sum_{q_{r,t}} p(q_{r,t-1}, q_{r,t} | s, \Lambda') \log a_{q_{r,t-1}, q_{r,t}} \\
&= \sum_s d'(s) \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{j=1}^I p(q_{r,t-1} = i | s, \Lambda') p(q_{r,t} = j | q_{r,t-1} = i, s, \Lambda') \log a_{i,j} \\
&= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t-1, i) \sum_{j=1}^I a'_{i,j} \log a_{i,j} \tag{45}
\end{aligned}$$

where  $a'_{i,j} = p(q_{r,t} = j | q_{r,t-1} = i, s, \Lambda')$  is the transition probability from the previous GT iteration.

Substituting (40), (41), (44) and (45) into (38), and denoting the emitting probability by  $b_i(x_{r,t}) = p(x_{r,t} | q_{r,t} = i, \Lambda)$  and  $b'_i(x_{r,t}) = p(x_{r,t} | q_{r,t} = i, \Lambda')$ , we obtain the decomposed and simplified objective function:

$$U(\Lambda; \Lambda') = U_1(\Lambda; \Lambda') + U_2(\Lambda; \Lambda') \tag{46}$$

where

$$\begin{aligned}
U_1(\Lambda; \Lambda') &= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) \log b_i(x_{r,t}) \\
&\quad + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t, i) \sum_{\mathcal{X}_{r,t}} b'_i(\mathcal{X}_{r,t}) \log b_i(\mathcal{X}_{r,t}) \tag{47}
\end{aligned}$$

$$\begin{aligned}
U_2(\Lambda; \Lambda') &= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_{j=1}^I \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \xi_{i,j,r,s_r}(t) \log a_{i,j} \\
&\quad + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t-1, i) \sum_{j=1}^I a'_{i,j} \log a_{i,j} \tag{48}
\end{aligned}$$

In (46),  $U_1(\Lambda; \Lambda')$  is relevant only to optimizing the emitting probability  $b_i(k)$ , and  $U_2(\Lambda; \Lambda')$  is relevant only to optimizing the transition probability  $a_{i,j}$ .

#### D. Establishing Growth Transformation by Optimizing the Auxiliary Function $U(\Lambda; \Lambda')$

In order to optimize the discrete distribution  $b_i(k) = p(x_{r,t} = k | q_{r,t} = i, \Lambda)$ ,  $k=1, 2, \dots, K$ , where the constraint  $\sum_{k=1}^K b_i(k) = 1$  is imposed, we apply the Lagrange multiplier method by constructing

$$W_1(\Lambda; \Lambda') = U_1(\Lambda; \Lambda') + \sum_{i=1}^I \lambda_i \left( \sum_{k=1}^K b_i(k) - 1 \right) \tag{49}$$

Setting  $\frac{\partial W_1(\Lambda; \Lambda')}{\partial \lambda_i} = 0$  and  $\frac{\partial W_1(\Lambda; \Lambda')}{\partial b_i(k)} = 0$ ,  $k=1, \dots, K$ , we have the following  $K+1$  equations:

$$\begin{aligned}
&\sum_{k=1}^K b_i(k) - 1 = 0 \\
0 &= \lambda_i b_i(k) + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{s \text{ s.t. } x_{r,t}=k} p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r, t, i) b'_i(k), \quad k=1, \dots, K.
\end{aligned}$$

where  $b_i(k)$  is multiplied on both sides. Solving for  $b_i(k)$ , we obtain the re-estimation formula:

$$b_i(k) = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{s \text{ s.t. } x_{r,t}=k} p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) + b'_i(k) \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t,i)}{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t,i)} \quad (50)$$

We now define

$$D_i = \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t,i) \quad (51)$$

$$\Delta\gamma(i,r,t) = \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) \quad (52)$$

and re-write (50) as

$$b_i(k) = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i,r,t) + b'_i(k) D_i}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i,r,t) + D_i} \quad (53)$$

In order to optimize transition probability  $a_{i,j}$ , with constraint  $\sum_{j=1}^I a_{i,j} = 1$ , we apply the Lagrange multiplier method by constructing

$$W_2(\Lambda; \Lambda') = U_2(\Lambda; \Lambda') + \sum_{i=1}^I \lambda_i \left( \sum_{j=1}^I a_{i,j} - 1 \right) \quad (54)$$

Setting  $\frac{\partial W_2(\Lambda; \Lambda')}{\partial \lambda_i} = 0$  and  $\frac{\partial W_2(\Lambda; \Lambda')}{\partial a_{i,j}} = 0, j=1, \dots, I$ , we have the following  $I+1$  equations:

$$\sum_{j=1}^I a_{i,j} - 1 = 0$$

$$0 = \lambda_i a_{i,j} + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \xi_{i,j,r,s_r}(t) + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t-1,i) a'_{i,j}, j=1, \dots, I.$$

Note that  $\sum_{j=1}^I \xi_{i,j,r,s_r}(t) = \gamma_{i,r,s_r}(t)$ . By solving  $a_{i,j}$ , we obtain the re-estimation formula with a standard procedure (used for deriving the EM estimate of transition probabilities [11]):

$$a_{i,j} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \xi_{i,j,r,s_r}(t) + a'_{i,j} \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t-1,i)}{\sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t-1,i)} \quad (55)$$

Now we define

$$\mathcal{D}_i = \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t-1,i) \quad (56)$$

$$\Delta\xi(i,j,r,t) = \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \xi_{i,j,r,s_r}(t) \quad (57)$$

and together with (52), we re-write (55) as

$$a_{i,j} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta \xi(i, j, r, t) + a'_{i,j} \mathcal{D}_i}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta \gamma(i, r, t) + \mathcal{D}_i} \quad (58)$$

The parameter re-estimation formulas (53) and (58) are unified across MMI, MCE, and MPE/MWE. What distinguishes among MMI, MCE, and MPE/MWE is the different weighing term  $\Delta \gamma(i, r, t)$  in (52) and  $\Delta \xi(i, j, r, t)$  in (57) due to the different  $C(s)$  contained in the unified objective function. Details for computing  $\Delta \gamma(i, r, t)$  for MMI, and MCE, and MPE/MWE are included in Appendix I.

### E. Setting Constant $D_i$

Values of constant  $D_i$  in (53) and  $\mathcal{D}_i$  in (56) determine the stability and convergence speed of the above GT/EBW algorithm. From (51), (42), and (37), we have

$$D_i = \sum_{r=1}^R \sum_{t=1}^{T_r} d(r, t, i) = \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s d'(s) p(q_{r,t} = i | s, \Lambda') = \frac{1}{p(X | \Lambda')} \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s d(s) p(q_{r,t} = i | s, \Lambda') \quad (59)$$

The theoretical basis for setting  $D_i$  to ensure that (53) and (58) are growth transformations is the requirement described in (33) that  $d(s)$  of (59) be sufficiently large so that  $\Gamma(\Lambda') + d(s) > 0$ . From (34),

$$\Gamma(\Lambda') = \delta(\chi, X) p(s) [C(s) - O(\Lambda')] = \begin{cases} p(s) [C(s) - O(\Lambda')] & \text{if } \chi = X \\ 0 & \text{otherwise} \end{cases}$$

Therefore,  $d(s) > \max\{0, -p(s)[C(s) - O(\Lambda')]\}$ . This gives

$$D_i > \frac{1}{p(X | \Lambda')} \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s \max\{0, p(s)[O(\Lambda') - C(s)]\} p(q_{r,t} = i | s, \Lambda') \quad (60)$$

Similarly, we can derive that

$$\mathcal{D}_i > \frac{1}{p(X | \Lambda')} \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s \max\{0, p(s)[O(\Lambda') - C(s)]\} p(q_{r,t-1} = i | s, \Lambda'). \quad (61)$$

In practice,  $D_i$  and  $\mathcal{D}_i$  given by (60) and (61) have often been found to be over conservative and unnecessarily large, causing slower convergence than those obtained through some empirical methods. We will not discuss such heuristics in this review, and would like to point out that this is still an interesting research problem and to refer the readers to the studies and discussions in [18] [42][43][48][55][60][62].

## VI. DISCRIMINATIVE LEARNING FOR CONTINUOUS DENSITY HMMs

For continuous-density HMMs (CDHMMs), the observation space is not quantized. In this case,  $X = X_1 \dots X_R$  is a concatenation of all training tokens, and each training token  $X_r$  consists of a sequence of continuous random variables. The formulation (26) applies to discriminative learning for CDHMMs. In particular,  $\chi$  in previous equations (30) and (31) is a continuous variable and hence the summation over domain  $\chi$  is changed to integration over  $\chi$ . That is, (30) is modified to

$$V(\Lambda; \Lambda') = \sum_s \sum_q \int_{\chi} f(\chi, q, s, \Lambda') \log f(\chi, q, s, \Lambda) d\chi \quad (62)$$

where the integrand  $f(\chi, q, s, \Lambda)$  is defined by

$$F(\Lambda; \Lambda') = \sum_s \sum_q \int_{\chi} f(\chi, q, s, \Lambda) d\chi \quad (63)$$

Correspondingly,

$$\begin{aligned} F(\Lambda; \Lambda') &= \sum_s \sum_q [C(s) - O(\Lambda')] p(s) p(X, q | s, \Lambda) + D \\ &= \sum_s \sum_q \int_{\chi} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda) d\chi \end{aligned} \quad (64)$$

where

$$f(\chi, q, s, \Lambda) = [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda) \quad (65)$$

and

$$\Gamma(\Lambda') = \delta(\chi, X) p(s) [C(s) - O(\Lambda')] \quad (66)$$

with  $\delta(\chi, X)$  in (66) being the Dirac delta function<sup>2</sup>. After a similar derivation as in the preceding section, it can be shown that the transition probability estimation formula (58) stays the same as the discrete HMM case. But for the emitting probability, (47) is changed to

$$\begin{aligned} U_1(\Lambda; \Lambda') &= \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i, r, s_r}(t) \log b_i(x_{r,t}) \\ &\quad + \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{i=1}^I d(r, t, i) \int_{\chi_{r,t}} b'_i(\chi_{r,t}) \log b_i(\chi_{r,t}) d\chi_{r,t} \end{aligned} \quad (67)$$

#### A. GT-Based Parameter Estimation for Gaussian Density CDHMM

We first derive the GT-based parameter estimation formulas for the CDHMM with Gaussian distributions and then generalize them to the case of mixture-Gaussian distributions in the subsequent subsection. For the CDHMM with Gaussian distributions, the observation probability density function  $b_i(x_{r,t})$  in (67) becomes a Gaussian distribution taking the following form:

$$b_i(x_{r,t}) \propto \frac{1}{|\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (x_{r,t} - \mu_i)^T \Sigma_i^{-1} (x_{r,t} - \mu_i) \right]. \quad (68)$$

where  $(\mu_i, \Sigma_i), i=1, 2, \dots, I$  are the mean vector and covariance matrix of the Gaussian component at state  $i$ .

In order to solve for  $\mu_i$  and  $\Sigma_i$ , based on (67), we set

$$\frac{\partial U_1(\Lambda; \Lambda')}{\partial \mu_i} = 0; \quad \text{and} \quad \frac{\partial U_1(\Lambda; \Lambda')}{\partial \Sigma_i} = 0.$$

This gives:

---

<sup>2</sup> Due to the Dirac delta function, Jensen's inequality no longer applies to the secondary auxiliary function. However, in section VI.B we will show that (62) is still a valid auxiliary function.

$$0 = \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) \Sigma_i^{-1} (x_{r,t} - \mu_i) + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t,i) \Sigma_i^{-1} \int_{\mathcal{X}_{r,t}} b'_i(\mathcal{X}_{r,t}) (\mathcal{X}_{r,t} - \mu_i) d\mathcal{X}_{r,t} \quad (69)$$

$$0 = \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,r,s_r}(t) \left[ \Sigma_i^{-1} - \Sigma_i^{-1} (x_{r,t} - \mu_i) (x_{r,t} - \mu_i)^T \Sigma_i^{-1} \right] + \sum_{r=1}^R \sum_{t=1}^{T_r} d(r,t,i) \int_{\mathcal{X}_{r,t}} b'_i(\mathcal{X}_{r,t}) \left[ \Sigma_i^{-1} - \Sigma_i^{-1} (\mathcal{X}_{r,t} - \mu_i) (\mathcal{X}_{r,t} - \mu_i)^T \Sigma_i^{-1} \right] d\mathcal{X}_{r,t} \quad (70)$$

For a Gaussian distribution  $b'_i(\mathcal{X}_{r,t}) = p(\mathcal{X}_{r,t} | q_{r,t} = i; \Lambda')$ , we have

$$\begin{aligned} \int_{\mathcal{X}_{r,t}} b'_i(\mathcal{X}_{r,t}) d\mathcal{X}_{r,t} &= 1, \\ \int_{\mathcal{X}_{r,t}} \mathcal{X}_{r,t} \cdot b'_i(\mathcal{X}_{r,t}) d\mathcal{X}_{r,t} &= \mu'_i, \\ \int_{\mathcal{X}_{r,t}} (\mathcal{X}_{r,t} - \mu'_i) (\mathcal{X}_{r,t} - \mu'_i)^T \cdot b'_i(\mathcal{X}_{r,t}) d\mathcal{X}_{r,t} &= \Sigma'_i. \end{aligned}$$

Hence integrals in (69) and (70) have closed-form results. Next, we left-multiply both sides of (69) by  $\Sigma_i$ , and left- and right-multiply both sides of (70) by  $\Sigma_i$ . Finally, solving  $\mu_i$  and  $\Sigma_i$  gives the ‘‘GT’’ formulas of

$$\mu_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i,r,t) x_i + D_i \mu'_i}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i,r,t) + D_i}, \quad (71)$$

$$\Sigma_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \left[ \Delta\gamma(i,r,t) (x_i - \mu_i) (x_i - \mu_i)^T \right] + D_i \Sigma'_i + D_i (\mu_i - \mu'_i) (\mu_i - \mu'_i)^T}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i,r,t) + D_i} \quad (72)$$

where  $\Delta\gamma(i,r,t)$  is defined in (52) and  $D_i$  defined in (51).

Just as in the discrete HMM case, (71) and (72) are based on the generic discriminative objective function  $O(\Lambda)$  in the form of (26), which covers MMI, MCE and MPE/MWE as special cases. This leads to unified, GT-based parameter estimation formulas for MMI, MCE, MPE/MWE, as well as for any other discriminative objective functions that can be mapped into the rational-function form (26). Moreover,  $\Delta\gamma(i,r,t)$  in (71) and (72) is defined in the same way as (52) in the discrete-HMM case — differing only in  $C(s)$  for MMI, MCE, and MPE/MWE, respectively, as will be illustrated further in Appendix I.

### B. Setting Constant $D_i$ for CDHMM

Based on Jensen’s inequality, the theoretical basis for setting an appropriate constant  $D_i$  to ensure that (71) and (72) are growth transformation is the requirement specified in (33), where  $d(s)$  in (59) needs to be sufficiently large to ensure that for any string  $s$  and any observation sequence  $\mathcal{X}$ ,  $\Gamma(\Lambda') + d(s) > 0$ , where  $\Gamma(\Lambda') = \delta(\mathcal{X}, X) p(s) [C(s) - O(\Lambda')]$  is defined in (34). However, for

CDHMM,  $\delta(\chi, X)$  is the Dirac delta function, which is a distribution with its density function value unbounded at the centre point, i.e.,  $\delta(\chi, X) = +\infty$  when  $\chi = X$ . Therefore, for a string  $s$  such that  $C(s) - O(\Lambda') < 0$ ,  $\Gamma(\Lambda')|_{\chi=X} = -\infty$ . Under this condition, it is impossible to find a bounded  $d(s)$  that ensures  $\Gamma(\Lambda') + d(s) > 0$  and hence Jensen's inequality may not apply. Note that this problem does not occur in the discrete HMM case, because in that case  $\delta(\chi, X)$  is a Kronecker delta function taking only a finite value of either zero or one.

The above-mentioned difficulty for CDHMMs can be overcome and the same derivation can still be used, if it can be shown that there exists a sufficiently large but still bounded constant  $D$  so that  $V(\Lambda; \Lambda')$  of (62), with the integrand defined by (65), is still a valid auxiliary function of  $F(\Lambda; \Lambda')$ ; i.e., an increase of the value of  $V(\Lambda; \Lambda')$  can guarantee an increase of the value of  $F(\Lambda; \Lambda')$ . Such a proof was developed in the recent work of [2] for GT-based MMI training for CDHMMs, and it holds for our common rational-function discriminative training criterion as well<sup>3</sup>, an important topic which we will outline in Appendix IV. Therefore, a bounded  $D_i$  exists according to (59), where  $d(s)$  satisfies (124) in Appendix IV.

Although a sufficiently large  $D_i$  guarantees monotone convergence of the GT-based iterative estimation formulas, i.e., (53)(58) for the discrete HMM and (71)(72) for the CDHMM, the value of  $D_i$  from the monotone convergence proof is a very loose upper bound and it can be too large for a reasonable convergence speed. In practice,  $D_i$  is often empirically set to achieve compromised training performance.

Empirical setting of  $D_i$  has been extensively studied from the day when EBW was proposed. In early days, only one global constant  $D$  was used for all parameters [18][42]. Later research discovered on the empirical basis that for CDHMM, a useful lower bound on (non-global)  $D_i$  is the value satisfying the constraint that newly estimated variances remain positive [43]. In [60][61], this constraint was further explored, leading to some quadratic inequalities upon which the lower bound of  $D_i$  can be solved. Most recently, in [55], constant  $D_i$  was further bounded by an extra condition that the denominators in re-estimation formulas remain non-singular.

In [62], the use of Gaussian-specific  $D_i$  was reported to give further improved convergence speed. For MMI, the Gaussian-specific constant  $D_i$  was set empirically to be the maximum of i) two times of the value necessary to ensure positive variances, i.e.,  $2 \cdot D_{\min}$ ; and ii) a global constant  $E$  multiplied by the denominator occupancy; e.g.,  $E \cdot \gamma_i^{den}$ . Specifically, for MMI in the

work of [62],  $\gamma_i^{den} = \sum_{r=1}^R \sum_{t=1}^{T_r} \gamma_{i,r}^{den}(t) = \sum_{r=1}^R \sum_{t=1}^{T_r} \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,s_r}(t)$ . For MPE reported in

[46][47][48], the empirical setting of  $D_i$  was the same as MMI, i.e.,  $D_i = \max\{2 \cdot D_{\min}, E \cdot \gamma_i^{den}\}$  except that the computation of the denominator occupancy became

$\gamma_i^{den} = \sum_{r=1}^R \sum_{t=1}^{T_r} \max(0, -\Delta \gamma(i, r, t))$ . In addition, these new parameters were further smoothed with

the ML estimate of parameters (which was called I-smoothing).

For MCE, the empirical setting of  $\gamma_i^{den}$  as  $\sum_{r=1}^R \sum_{t=1}^{T_r} p(S_r | X_r, \Lambda') \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,s_r}(t)$  was

developed in the recent work of [24][68]. It was based on the consideration that MCE and MMI are equivalent in the special case of having only one utterance in the training set. This setting was

---

<sup>3</sup> The authors would like to thank an anonymous reviewer for pointing us to the work of [2].

experimentally verified with strong recognition results as reported in [24][68]. Further discussions and comparisons of different settings of empirical  $D_i$  can be found in [18] [24][42][43] [48] [55][61] [62] .

### C. Parameter Estimation for Gaussian Mixture CDHMM

The model parameter estimation formulas for a Gaussian-mixture HMM are similar to those for a Gaussian HMM discussed earlier. For a Gaussian-mixture HMM, the continuous observation density function  $b_i(x_{r,t})$  for state  $i$  has the following form:

$$b_i(x_{r,t}) = \sum_{l=1}^L w_{i,l} N(x_{r,t} | \mu_{i,l}, \Sigma_{i,l}) \quad (73)$$

where  $b_i(x_{r,t})$  is a mixture of  $L$  Gaussian components,  $N(x_{r,t} | \mu_{i,l}, \Sigma_{i,l})$  is the  $l$ -th Gaussian mixture component that takes the same form as (68),  $w_{i,l}$  is a positive weight of the  $l$ -th Gaussian component, and  $\sum_{l=1, \dots, L} w_{i,l} = 1$ . Compared with a Gaussian HMM, there is an additional hidden component, the Gaussian component index sequence  $l$ . The hidden sequence  $l$  can be accommodated in (62) by the same way that we exploited to treat the hidden state sequence  $q$ . Then after similar derivation steps, we can obtain the following parameter estimation formulas:

$$\mu_{i,l} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, l, r, t) x_t + D_{i,l} \mu'_{i,l}}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, l, r, t) + D_{i,l}}, \quad (74)$$

$$\Sigma_{i,l} = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} [\Delta\gamma(i, l, r, t) (x_t - \mu_{i,l})(x_t - \mu_{i,l})^T] + D_{i,l} \Sigma'_{i,l} + D_{i,l} (\mu_{i,l} - \mu'_{i,l})(\mu_{i,l} - \mu'_{i,l})^T}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, l, r, t) + D_{i,l}} \quad (75)$$

where  $D_{i,l}$  and  $\Delta\gamma(i, l, r, t)$  are defined in a similar way to (51) and (52), i.e.,

$$D_{i,l} = \sum_{r=1}^R \sum_{t=1}^{T_r} d(r, t, i, l) \quad (76)$$

$$\Delta\gamma(i, l, r, t) = \sum_s p(s | X, \Lambda') (C(s) - O(\Lambda')) \gamma_{i,l,r,s}(t) \quad (77)$$

and  $\gamma_{i,l,r,s}(t) = p(q_{r,t} = i, l | X_r, s_r, \Lambda')$  is the occupation probability of Gaussian mixture component  $l$  of state  $i$ , at time  $t$  in the  $r$ -th utterance. Accordingly, the empirical setting of  $D_{i,l}$  takes similar forms as discussed in the previous section, except that  $\Delta\gamma(i, l, r, t)$  and  $\gamma_{i,l,r,s}(t)$  will be used instead. Estimation of the mixture component weights  $w_{i,l}$  is similar to the discrete HMM estimation case, and will not be described here.

## VII. RELATED WORK AND DISCUSSIONS

### A. Relation to other approaches

In recent papers [35][55], an approach was proposed to unify a number of discriminative learning methods including MMI, MPE, and MPE/MWE (the earlier paper [55] did not include

MPE/MWE). Functional similarities and differences among MMI, MCE, and MPE/MWE criteria were noted and discussed in [35][55]. In this paper, the proposed framework takes an additional step of unifying these criteria in a canonical rational-function form (26), and GT-based discriminative learning is applied to this rational-function which includes MMI, MCE and MPE/MWE criteria as special cases. This is significant from two perspectives. First, it provides a more precise and direct insight into the fundamental relations among MMI, MCE and MPE/MWE criteria at the objective function level based on the common rational-function form (26). Second, it enables a unified GT-based parameter optimization framework that applies to MMI, MCE, MPE/MWE and other discriminative criteria, as long as their objective functions can be represented by (26).

The proposed framework in [35] was based on the objective function of the following form (rewritten using the mathematical notations adopted in this paper for easy comparisons):

$$O(\Lambda) = \frac{1}{R} \sum_{r=1}^R f \left( \frac{1}{\eta} \log \frac{\sum_{s_r} p^\eta(X_r, s_r | \Lambda) C_{DT}(s_r)}{\sum_{s_r \in M_r} p^\eta(X_r, s_r | \Lambda)} \right) \quad (78)$$

where  $C_{DT}(s_r)$  takes the same value as in our Table 1. The choices of the smoothing function  $f(z)$ , the competing word sequences  $M_r$ , and the weight value  $\eta$  in (78) are provided in Table 2 for the different types of DT criteria. In Table 2,  $q$  is the slope of a sigmoid smoothing function.

Criteria	Smoothing Function $f(z)$	Alternative Word Sequences $M_r$	$\eta$
MCE ( $N$ -best)	$-1/[1 + \exp(2qz)]$	$\{s_r\}$ excluding $S_r$	$\geq 1$
MCE (one-best)	$-1/[1 + \exp(2qz)]$	$\{s_{r,l}\}$	N/A
MPE/MWE	$\exp(z)$	all possible label sequence $\{s_r\}$	1
MMI	$z$	all possible label sequence $\{s_r\}$	1

Table 2. Choices of the smoothing function  $f(z)$ , alternative word sequences  $M_r$ , and exponent weight  $\eta$  in (78) for various types of DT criteria. This is modified from the original table in [55].

Eq. (78) indicates that different discriminative criteria can have a similar form of kernel and differ by the criterion-dependent smoothing function  $f(z)$  that modulates the kernel, where the objective function is a sum of smoothing functions. Eq. (78) is a generic description of the objective functions of MMI, MCE, and MPE/MWE. However, it is not in a general form of a rational function (defined as a ratio of two polynomial functions) due to the presence of the nonlinear function  $f(z)$ . The important distinction of product vs. summation of utility functions among these criteria (as described in section II.D) is not explicitly addressed. In the approach presented in this paper, we address this issue directly and show that the objective functions from MMI, MCE, and MPE/MWE criteria can have a definitive rational-function form (26), and for each discriminative criterion, the objective function differs only by a model-independent quantity  $C_{DT}(s_1 \dots s_R)$ .

Furthermore, as shown in Table 2, since  $f(z)$  is a nonlinear function for the MPE/MWE and MCE criteria, the original GT solution [18], while directly applicable to MMI with  $f(z)$  being an identity function and  $z$  being the logarithm of a rational function (since sum of log becomes log of product), is not directly applicable to the objective functions of the MPE/MWE and MCE criteria (since the sum stays when  $f(z)$  is nonlinear). In order to circumvent this difficulty, the theorem described in [30] is applied. In [30], the original objective function is approximated by a Taylor series expansion. Then, via a similar approach to that of [18], the GT-based parameter optimization may be applied to the partial sum of the Taylor series expansion, which is a

polynomial with a finite degree. This forms the theoretical basis of the earlier GT-based methods for MCE and MPE/MWE [35][55]. However, the positive growth of the partial sum depends on the degree of that partial sum (see more detailed discussions on this point in [18]), and it vanishes when the degree goes to infinity. It may vanish even faster than the error of Taylor series approximation does. Therefore, it has not been definitively shown that the re-estimation formula ensures true growth of the value of the objective function with iteration.

In contrast, the unified rational-function approach described in this paper departs from the work of [35][55]. It is free from the Taylor series approximation and it shows that the objective functions for the MMI, MCE, and MPE/MWE criteria have a common definitive rational-function form (26). Therefore, the GT-based parameter optimization framework can be directly applied to (26) in a constructive way.<sup>4</sup> Moreover, the unified representation of the discriminative objective functions developed in this paper opens a way to apply other rational-function based optimization methods (e.g., the method based on the reverse Jensen inequality [26]) to MMI, MCE, and MPE/MWE-based classifier design. Using the structure of the rational function, we expect that all desirable algorithmic properties of the parameter optimization procedures presented in this paper can be established and justified.

### B. Relation to gradient-based optimization

The relation between the GT/EBW methods and gradient-based methods has been studied in the literature (e.g., [2][54][55]). In addition to the critical difference in the convergence properties, the learning speed of GT/EBW-based updating formula (71) is comparable to a quadratic Newton update; i.e., it can be formulated as a gradient ascent with the step size that approximates inverse Hessian  $H$  of the objective function. Let us take the *mean* vector estimation as an example for the objective function of the form (26) in the case of CDHMM. The gradient of  $O(\Lambda)$  w.r.t.  $\mu_i$  can be shown to be

$$\nabla_{\mu_i} O(\Lambda) |_{\Lambda=\Lambda'} = \Sigma_i'^{-1} \sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t)(x_t - \mu_i') \quad (79)$$

On the other hand, we can rewrite the GT formula of (71) into the following equivalent form:

$$\begin{aligned} \mu_i &= \mu_i' + \frac{1}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) + D_i} \cdot \sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t)(x_t - \mu_i') \\ &= \mu_i' + \frac{1}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) + D_i} \Sigma_i' \cdot \nabla_{\mu_i} O(\Lambda) |_{\Lambda=\Lambda'} \end{aligned} \quad (80)$$

Consider the quadratic Newton update, where the Hessian  $H_i$  for  $\mu_i$  can be approximated by the following equation after dropping the dependency of  $\mu_i$  with  $\Delta\gamma(i, r, t)$ :

$$H_i = \nabla_{\mu_i}^2 O(\Lambda) |_{\Lambda=\Lambda'} \approx -\Sigma_i'^{-1} \sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t)$$

Therefore, the updating formula of GT in (71) can be further re-written to

---

<sup>4</sup> Note that the approach we take in this paper is based on the work in [2][21] rather than on the work of [3][18].

$$\mu_i \approx \mu'_i - \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t)}{\sum_{r=1}^R \sum_{t=1}^{T_r} \Delta\gamma(i, r, t) + D_i} H_i^{-1} \nabla_{\mu_i} O(\Lambda) |_{\Lambda=\Lambda'} \quad (81)$$

which approximates the quadratic Newton update  $\mu_i = \mu'_i - \alpha \cdot H_i^{-1} \nabla_{\mu_i} O(\Lambda) |_{\Lambda=\Lambda'}$  and usually gives a faster learning speed than the simple gradient-based search.

Other popular and effective gradient-based methods exist for optimizing discriminative training criteria [38][31][15][22][37]. For instance, Quickprop [17] is a batch-mode, second-order optimization method that approximates Newton’s optimization, with the help of heuristics to determine the proper update step size. Rprop [52], which stands for “Resilient back-propagation”, is another batch-mode optimization method, which performs dynamic scaling of the update step size for each parameter based on different kinds of heuristics. In [38], a comprehensive study of gradient-based optimization methods for MCE training, including batch and semi-batch probabilistic descent (PD), QuickProp, and RProp, is given for large vocabulary speech recognition tasks. It was shown that the MCE criterion can be optimized by using these gradient-based methods, and improved recognition accuracies were reported. Furthermore, there exist other gradient-based methods such as BFGS and conjugate gradient search [5][14]. Although both of these methods are more complicated to implement for large scale discriminative training tasks, they are superior to other gradient-descent techniques in terms of the convergence properties. Readers are referred to [5][53] for further discussions.

In contrast to the popular gradient-based methods discussed above, we can view a class of optimization methods with a re-estimation style, including Expectation-Maximization (EM) algorithm and EBW algorithm, as GT-based methods in a broad sense. The GT-based methods are designed for the objective functions with special, rational-functional forms, and the GT algorithm can ensure rigorous monotone growth of the value of the objective functions iteratively. From this perspective, on the one hand, GT-based methods are less general than gradient-based ones. On the other hand, they give desirable monotone convergence in training. Further, although GT-based parameter re-estimation formulas may be rewritten into gradient-based forms, the step sizes are specifically derived so that monotone convergence is guaranteed. This critical property differentiates them from general gradient-based methods.

The advanced gradient-based methods discussed above, such as batch and semi-batch PD, QuickProp, RProp, BFGS, and conjugate gradient, are alternatives to the GT/EBW-method for optimizing discriminative training criteria. Although theoretically the GT/EBW method has the desirable monotone convergence property, empirical setting of  $D$  is used in practice to speed up training with the trade-off for monotone convergence. This makes rigorous comparisons between GT/EBW-based and advanced gradient-based methods difficult. In the literature, experimental results of both types of methods have been reported on various speech recognition tasks [22][24][31][55][43].

Algorithmic convergence of parameter estimation is a central issue for classifier design using discriminative training criteria. Search for more powerful discriminative criteria and optimization methods in classifier design remains an area of active and on-going research. It is our hope that the unified rational-function based objective function representation reviewed in this paper can provide additional structural formulation and can motivate the development of new learning algorithms to improve the discriminative power of sequential pattern classifiers and recognizers.

### VIII. SUMMARY

In this paper, we studied the objective functions of MMI, MCE and MPE/MWE for discriminative learning in sequential pattern recognition. We presented an approach that unifies the objective functions of MMI, MCE and MPE/MWE in a common rational-function form of (26). The exact structure of the rational-function form for each discriminative criterion was derived and studied. While the rational-function form of MMI has been known in the past, we provided the theoretical proof that the similar rational-function form exists for the objective functions of MCE and MPE/MWE. Moreover, we showed that the rational function forms for objective functions of MMI, MCE and MPE/MWE differ in the constant weighting factors  $C_{DT}(s_1 \dots s_R)$  and these weighting factors depend only on the labeled sequence  $s_1 \dots s_R$ , and are independent of the parameter set  $\Lambda$  to be optimized.

The derived rational-function form for MMI, MCE and MPE/MWE allows the GT/EBW-based parameter optimization framework to be applied directly in discriminative learning. In the past, lack of the appropriate rational-function form was a difficulty for MCE and MPE/MWE, because without this form, the GT/EBW-based parameter optimization framework cannot be directly applied. Based on the unified rational-function form, in a tutorial style, we derived the GT/EBW-based parameter optimization formulas for both discrete HMMs and CDHMMs in discriminative learning using MMI, MCE, and MPE/MWE criteria.

This paper is motivated by the striking success of the MMI, MCE, and MPE/MWE-based discriminative criteria in sequential pattern recognition and particularly in speech and language processing. Yet there was a lack of common understanding of the inter-relation among these techniques, despite the relatively long history of MMI (since 1987 [7]), MCE (since 1992 [28]), and MPE/MWE (since 2002 [46]). Due to the complexity of these techniques and the lack of a common underlying theoretical theme and structure, disparate discriminative learning procedures were developed and parameter optimization has become a major issue. The main goal of this paper is to provide an underlying foundation for MMI, MCE, and MPE/MWE at the objective function level to facilitate the development of new parameter optimization techniques and to incorporate other pattern recognition concepts, e.g., discriminative margins [68], with the current discriminative learning paradigm.

#### APPENDIX I: COMPUTING $\Delta\gamma(i, r, t)$ IN THE GROWTH TRANSFORMATION FORMULAS

In (52), computing  $\Delta\gamma(i, r, t)$  involves summation over all possible super-string label sequences  $s = s_1 \dots s_R$ . The number of training tokens (sentence strings),  $R$ , is usually very large. Hence, the summation over  $s$  needs to be decomposed and simplified. To proceed, we use the notations of  $s' = s_1 \dots s_{r-1}$ ,  $s'' = s_{r+1} \dots s_R$ ,  $X' = X_1 \dots X_{r-1}$ , and  $X'' = X_{r+1} \dots X_R$ . Then, from (52), we have,

$$\Delta\gamma(i, r, t) = \sum_{s_r} p(s_r | X_r, \Lambda') \left[ \sum_{s'} \sum_{s''} p(s', s'' | X', X''; \Lambda') (C(s', s_r, s'') - O(\Lambda')) \right] \gamma_{i, r, s_r}(t) \quad (82)$$

where factor  $\Psi$  is the average deviation of the accuracy count for the given string  $s_r$ . The remaining steps in simplifying the computation of  $\Delta\gamma(i, r, t)$  will be separate for MMI and MCE/MPE/MWE because the parameter-independent accuracy count function  $C(s)$  for them takes the product and summation form, respectively (as shown in Table 1 ).

### A. Product form of $C(s)$ (for MMI)

For MMI, we have  $C(s) = C(s_1, \dots, s_R) = \prod_{r=1}^R C(s_r) = \prod_{r=1}^R \delta(s_r, S_r)$  in a product form. Using  $C(s', s_r, s'') = C(s_r) \cdot C(s', s'')$ , we simplify factor  $\Psi$  in (82) to

$$\begin{aligned} \Psi &= C(s_r) \cdot \sum_{s'} \sum_{s''} p(s', s'' | X', X''; \Lambda') C(s', s'') - O(\Lambda') \\ &= O(\Lambda') \left( \frac{C(s_r) \cdot \sum_{s'} \sum_{s''} p(s', s'' | X', X''; \Lambda') C(s', s'')}{O(\Lambda')} - 1 \right) \end{aligned} \quad (83)$$

The idea behind the above steps is to make use of the product form of the  $C(s)$  function for canceling out common factors in both  $O(\Lambda')$  and  $C(s)$  functions. To proceed, we now factorize  $O(\Lambda')$  as follows:

$$\begin{aligned} O(\Lambda') &= \frac{\left[ \sum_{s_r} p(s_r, X_r | \Lambda') C(s_r) \right] \left[ \sum_{s'} \sum_{s''} p(s', s'' | X', X'' | \Lambda') C(s', s'') \right]}{[p(X_r | \Lambda')] [p(X', X'' | \Lambda')]} \\ &= p(S_r | X_r, \Lambda') \cdot \sum_{s'} \sum_{s''} p(s', s'' | X', X''; \Lambda') C(s', s'') \end{aligned}$$

where the last step uses  $C(s_r) = \delta(s_r, S_r)$ . Substituting this to (83) then gives the simplification of

$$\Psi = O(\Lambda') \left( \frac{C(s_r)}{p(S_r | X_r, \Lambda')} - 1 \right) \quad (84)$$

Substituting (84) to (82) and using  $C(s_r) = \delta(s_r, S_r)$  again, we obtain

$$\Delta\gamma(i, r, t) = O(\Lambda') \left[ \gamma_{i,r,S_r}(t) - \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,s_r}(t) \right] \quad (85)$$

In the re-estimation formulas (71) and (72), if we divide both the numerator and denominator by  $O(\Lambda')$ ,  $\Delta\gamma(i, r, t)$  in (85) can take a simplified form of

$$\Delta\mathcal{Y}(i, r, t) = \left[ \gamma_{i,r,S_r}(t) - \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,s_r}(t) \right] = \gamma_{i,r}^{num}(t) - \gamma_{i,r}^{den}(t) \quad (86)$$

The corresponding constant  $D_i$  in the re-estimation formulas (71) and (72) then becomes

$$\mathcal{D}_i = D_i / O(\Lambda'). \quad (87)$$

Substituting the above into (71) and (72), we have the GT/EBW formulas for MMI

$$\mu_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} [\gamma_{i,r}^{num}(t) - \gamma_{i,r}^{den}(t)] x_t + \mathcal{D}_i \mu'_i}{\sum_{r=1}^R \sum_{t=1}^{T_r} [\gamma_{i,r}^{num}(t) - \gamma_{i,r}^{den}(t)] + \mathcal{D}_i}, \quad (88)$$

$$\Sigma_i = \frac{\sum_{r=1}^R \sum_{t=1}^{T_r} [\gamma_{i,r}^{num}(t) - \gamma_{i,r}^{den}(t)] (x_t - \mu_i)(x_t - \mu_i)^T + \mathcal{D}_i \Sigma'_i + \mathcal{D}_i (\mu_i - \mu'_i)(\mu_i - \mu'_i)^T}{\sum_{r=1}^R \sum_{t=1}^{T_r} [\gamma_{i,r}^{num}(t) - \gamma_{i,r}^{den}(t)] + \mathcal{D}_i} \quad (89)$$

This gives the classical GT/EBW-based MMI re-estimation formulas described in [42][62].

Eqn. (85) or (86) gives an  $N$ -best string based solution to computing  $\Delta\gamma(i, r, t)$ . This is illustrated by the string-level summation over  $s_r$  (i.e., the label sequence for token  $r$ , including both correct and incorrect strings). For  $N$ -best string-level discriminative training, the summation over  $s_r$  in (85) or (86) amounts to going through all  $N$ -best string hypotheses and is computationally inexpensive when  $N$  is relatively small (e.g.,  $N$  in the order of thousands as typical for most  $N$ -best experiments).

When a lattice instead of an explicit  $N$ -best list is provided for competing hypotheses in discriminative training, in theory, (85) or (86) can be applied just as for the  $N$ -best string based solution already discussed. This is because a lattice is nothing more than a compact representation of  $N$ -best strings. However, since  $N$  in this equivalent “ $N$ -best list” would be huge (in the order of billions or more [67]), more efficient techniques for dealing with the summation over  $s_r$  in computing (85) or (86) will be needed. Readers are referred to Appendix III for details of such computation.

### B. Summation form of $C(s)$ (MCE and MPE/MWE)

Different from MMI, for MCE and MPE/MWE, we have  $C(s) = C(s_1, \dots, s_R) = \sum_{r=1}^R C(s_r)$ , or  $C(s', s_r, s'') = C(s_r) + C(s', s'')$ . That is, the  $C$  function is in a summation instead of a product form. This changes the simplification steps for factor  $\Psi$  of (82) as follows:

$$\begin{aligned} \Psi &= \sum_{s'_r} \sum_{s''_r} p_{\Lambda'}(s'_r, s''_r | X', X'') C(s_r) + \sum_{s'_r} \sum_{s''_r} p(s'_r, s''_r | X', X''; \Lambda') C(s'_r, s''_r) - O(\Lambda') \\ &= C(s_r) + \sum_{s'_r} \sum_{s''_r} p(s'_r, s''_r | X', X''; \Lambda') C(s'_r, s''_r) - O(\Lambda') \end{aligned} \quad (90)$$

The idea behind the above steps is to make use of the summation form of the  $C(s)$  function for subtracting out the common terms in the  $O(\Lambda')$  function. To achieve this, we decompose  $O(\Lambda')$ , based on its original non-rational form, e.g., (22) or (18) and (19), as follows:

$$\begin{aligned} O(\Lambda') &= \frac{\sum_{s_r} p(X_r, s_r | \Lambda') C(s_r)}{\sum_{s_r} p(X_r, s_r | \Lambda')} + \sum_{i=1, i \neq r}^R \frac{\sum_{s_i} p(X_i, s_i | \Lambda') C(s_i)}{\sum_{s_i} p(X_i, s_i | \Lambda')} \\ &= \frac{\sum_{s_r} p(X_r, s_r | \Lambda') C(s_r)}{\sum_{s_r} p(X_r, s_r | \Lambda')} + \sum_{s'_r, s''_r} p(s'_r, s''_r | X', X'' | \Lambda') C(s'_r, s''_r) \end{aligned}$$

The second term above cancels out the same term in (90), leading to the simplification of

$$\Psi = C(s_r) - \frac{\sum_{s_r} p(s_r, X_r | \Lambda') C(s_r)}{\sum_{s_r} p(s_r, X_r | \Lambda')} \quad (91)$$

Now, substituting (91) back to (82), we obtain

$$\Delta\gamma(i, r, t) = \sum_{s_r} p(s_r | X_r, \Lambda') \left( C(s_r) - \frac{\sum_{s_r} p(X_r, s_r | \Lambda') C(s_r)}{\sum_{s_r} p(X_r, s_r | \Lambda')} \right) \gamma_{i, r, s_r}(t) \quad (92)$$

For MCE that has  $C(s_r) = \delta(s_r, S_r)$ , the above equation can be further simplified as:

$$\Delta\gamma(i, r, t) = p(S_r | X_r, \Lambda') \left[ \gamma_{i, r, S_r}(t) - \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i, r, s_r}(t) \right] \quad (93)$$

Again, if a lattice instead of an  $N$ -best list is provided for discriminative learning, a huge number of terms in the summation over  $s_r$  in (92) would be encountered. In order to keep the computation manageable, efficient computation of (92) based on the lattice is needed, which we describe in Appendix III.

## APPENDIX II: TWO EMPIRICAL ISSUES IN MCE IMPLEMENTATION

In this Appendix, we discuss two empirical issues in MCE implementation that were raised in Section II.B. First, in (12), if we use the exponent scale factor  $\eta \neq 1$ , we can obtain the following result corresponding to (14) :

$$l_r(d_r(X_r, \Lambda)) = \frac{\sum_{s_r, s_r \neq S_r} p^\eta(X_r, s_r | \Lambda)}{\sum_{s_r} p^\eta(X_r, s_r | \Lambda)}$$

The corresponding result to (17) then becomes

$$O_{MCE}(\Lambda) = \sum_{r=1}^R \frac{p^\eta(X_r, S_r | \Lambda)}{\sum_{s_r} p^\eta(X_r, s_r | \Lambda)}$$

which can be reformulated into a rational function using the same steps as in Section III.B:

$$O_{MCE}(\Lambda) = \frac{\sum_{s_1 \dots s_R} p^\eta(X_1 \dots X_R, s_1 \dots s_R | \Lambda) C_{MCE}(s_1 \dots s_R)}{\sum_{s_1 \dots s_R} p^\eta(X_1 \dots X_R, s_1 \dots s_R | \Lambda)} \quad (94)$$

The remaining derivations in Sections V and VI will no longer follow strictly for the more general and practical case of (94). In the MCE implementation as in our experiments reported in [23], however, we modify (92) for computing  $\Delta\gamma(i, r, t)$  in the following way in order to include the effects of the exponent scale factor:

$$\Delta\gamma(i, r, t) = \sum_{s_r} \mathbb{h}(s_r | X_r, \Lambda') \left( C(s_r) - \sum_{s_r} \mathbb{h}(X_r | s_r, \Lambda') C(s_r) \right) \gamma_{i, r, s_r}(t) \quad (95)$$

where  $\mathbb{h}(s_r | X_r, \Lambda')$  is the generalized posterior probability of  $s_r$ , which can be computed as

$$\mathbb{h}(s_r | X_r, \Lambda') = \frac{p^\eta(X_r, s_r | \Lambda')}{\sum_{s_r} p^\eta(X_r, s_r | \Lambda')} \quad (96)$$

The second empirical MCE implementation issue concerns the use of  $\alpha \neq 1$  in (13). For 1-best MCE,  $\alpha$  acts as  $\eta$ , or we equivalently set  $\eta = \alpha$ , and  $\alpha = 1$ . Then we can compute

$\Delta\gamma(i, r, t)$  according to (95). For N-best MCE ( $N > 1$ ), given the discriminant function defined in (12) and the sigmoid function defined in (13), we have the following result (which is corresponding to (14)):

$$l_r(d_r(X_r, \Lambda)) = \frac{\left( \sum_{s_r, s_r \neq S_r} p^\eta(X_r, s_r | \Lambda) \right)^\alpha}{p^{\eta \cdot \alpha}(X_r, S_r | \Lambda) + \left( \sum_{s_r, s_r \neq S_r} p^\eta(X_r, s_r | \Lambda) \right)^\alpha} \quad (97)$$

Now,  $\alpha$  is applied outside of the summation of scaled joint probabilities over all competing strings, making the rigorous computation intractable. In our practical MCE implementation, we

instead use  $\sum_{s_r, s_r \neq S_r} p^{\eta \cdot \alpha}(X_r, s_r | \Lambda)$  to approximate  $\left( \sum_{s_r, s_r \neq S_r} p^\eta(X_r, s_r | \Lambda) \right)^\alpha$ . This approximation

(which is exact when  $\eta$  approaches infinity) makes it equivalent to setting the new “ $\eta$ ” as  $\alpha \cdot \eta$ , and setting new  $\alpha = 1$ . Then, again, we can compute  $\Delta\gamma(i, r, t)$  according to (95). It should be noted that, with this approximation, the computation for the lattice-based MCE does not require removing the correct word string  $S_r$  from the lattice. This contrasts with the solution in [55][35] where the removal was necessary without using the approximation, making it more difficult to implement in practice.

The above two empirical solutions have been implemented successfully in our speech recognition system, yielding strong practical results (published in [24][68]) that validate the solutions.

### APPENDIX III: COMPUTING $\Delta\gamma(i, r, t)$ WHEN USING LATTICES

A lattice, as illustrated in Fig. 1, is a compact representation of a large list of strings. It is an acyclic directed graph consisting of a number of nodes (nine in Fig. 1 as a highly simplified example) and a set of directed arcs each connecting two nodes. In Fig. 1, each node corresponds to a time stamp and each arc corresponds to a sub-string unit (e.g., a word in a sentence). A string in the lattice contains multiple arcs. A typical arc is shown as  $q$  in Fig. 1. Two time stamps,  $b_q$  and  $e_q$ , are associated with each arc, providing an estimate of the segment boundaries for the substring. For a time slice  $t$  within the arc segment  $q$ , we have  $b_q \leq t \leq e_q$ .

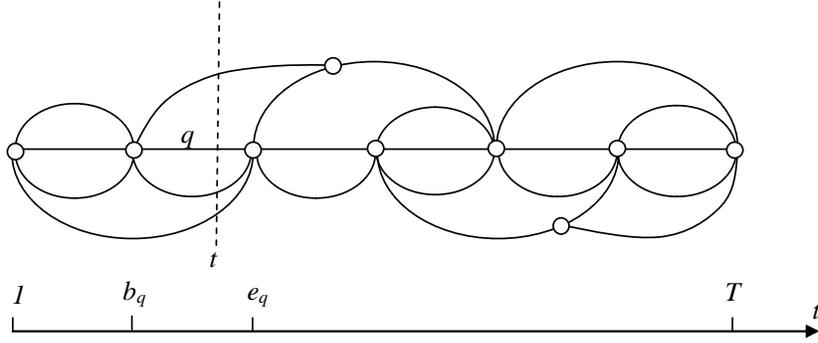


Fig. 1. A graphical illustration of a lattice, where  $q$  represents an arc in the lattice and  $t$  represents a time slice. The time span of arc  $q$  is  $b_q \leq t \leq e_q$  and that for the entire lattice is  $L \leq t \leq T$ . In this simple example, the total number of arcs ( $q$ ) is 21, which is substantially lower than the total number of paths ( $s_r$ ) of 420. The essence of the decomposition of occupation probability introduced in the text (Eq. (98)) is to enable fast computation by reducing the number of terms in summation over  $s_r$  to that over  $q$ .

We will show below that (86) and (92) can both be computed efficiently by a forward-backward algorithm. First, given the lattice in Fig. 1 and  $s_r$  as an arbitrary path in that lattice, we will show the occupancy given the entire string  $s_r$  can be computed as the occupancy given the local arc  $q$ , where arc  $q$  belongs to  $s_r$ . i.e.,

$$\gamma_{i,r,q}(t) = \gamma_{i,r,s_r}(t) \quad \text{when } b_q \leq t \leq e_q \quad (98)$$

To see this, let  $s_r$  be composed of three sub-strings:  $s'_r, q, s''_r$ , and correspondingly the observation sequence  $X_r$  is composed of three sub-sequences:  $X'_r, X_q, X''_r$ . Then the right hand side of (98) can be analyzed to be

$$\begin{aligned} \gamma_{i,r,s_r}(t: b_q \leq t \leq e_q) &= p(q_{r,t:b_q \leq t \leq e_q} = i \mid X_r, s_r, \Lambda') \\ &= p(q_{r,t:b_q \leq t \leq e_q} = i \mid X'_r, X_q, X''_r, s'_r, q, s''_r, \Lambda') \\ &= p(q_{r,t:b_q \leq t \leq e_q} = i \mid X_q, q, \Lambda') \\ &= \gamma_{i,r,q}(t: b_q \leq t \leq e_q) \end{aligned}$$

which is the left hand side of (98). The third step holds because usually the HMM of  $s_r$  is formed by concatenating phone specific HMMs, so the states in different arcs are belong to different HMMs, and are independent with each other, i.e., given arc  $q$ , its first HMM state  $q_{r,b_q}$  is independent of its preceding state  $q_{r,b_q-1}$ .

The essence of (98) is to decouple the dependency on the local arc  $q$  from the entire string  $s_r$ . This enables drastic simplification of the computation in (86) and (92), which we discuss below for three separate cases.

### A. Computing $\Delta\gamma(i, r, t)$ for MMI involving lattice approximation

The principal computation burden in (86) is the huge number ( $N$ ) of summation terms for  $s_r$  for the equivalent  $N$ -best list of a lattice in the following quantity in (86):

$$Y = \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,s_r}(t) \quad (99)$$

Using (98), we can significantly reduce the computation by the following simplification:

$$\begin{aligned} Y &= \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i,r,q}(t) = \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot p(q | X_r, \Lambda') = \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \frac{p(q, X_r | \Lambda')}{p(X_r | \Lambda')} \end{aligned} \quad (100)$$

Note the number of summation terms for  $q$  in (100) after the approximation is substantially smaller than that for  $s_r$  before the approximation. The key quantities in (100) can be efficiently computed as follows (proof omitted):

$$p(q, X_r | \Lambda') = \alpha(q) \beta(q); \quad (101)$$

$$p(X_r | \Lambda') = \sum_{q: q \in \{\text{ending arcs}\}} p(q, X_r | \Lambda') = \sum_{q: q \in \{\text{ending arcs}\}} \alpha(q) \quad (102)$$

where the ‘‘forward’’ and ‘‘backward’’ probabilities are defined by

$$\alpha(q) \quad p(q, X'_r(q), X_r(q) | \Lambda'); \quad (103)$$

$$\beta(q) \quad p(X''_r(q) | q, \Lambda'). \quad (104)$$

In (103),  $X'_r(q)$  denotes the  $r$ -th training token’s partial observation sequence preceding arc  $q$ , i.e., during  $1 \leq t < b_q$ .  $X_r(q)$  is the observation sequence bounded by arc  $q$  with  $b_q \leq t \leq e_q$ .  $X''_r(q)$  in (104) denotes the partial observation sequence succeeding arc  $q$ , or during  $e_q < t \leq T_r$ .  $\alpha(q)$  is the probability that lattice is at arc  $q$  during time  $b_q \leq t \leq e_q$ , and having generated partial observation  $X'_r(q)$  plus  $X_r(q)$ , i.e.,  $x_{r,1}, \dots, x_{r,e_q}$ .  $\beta(q)$  is the probability of generating partial observation  $X''_r(q)$  given that the lattice is at arc  $q$  at time  $t=e_q$ .

For each arc  $q$  in the lattice,  $\alpha(q)$  and  $\beta(q)$  can be computed by the following efficient forward and backward recursions, respectively (proofs omitted):

$$\alpha(q) = \sum_{\{p: p \text{ precedes } q\}} P(q | p, \Lambda') p(X_r(q) | q, \Lambda') \alpha(p) \quad (105)$$

and

$$\beta(q) = \sum_{\{v: v \text{ succeeds } q\}} P(v | q, \Lambda') p(X_r(v) | v, \Lambda') \beta(v) \quad (106)$$

where in (105),  $\{p: p \text{ precedes } q\}$  is the collection of all arcs  $p$  that directly connects to  $q$  in the lattice. Similarly,  $\{v: v \text{ succeeds } q\}$  in (106) is the collection of all arcs  $v$  that directly connect to  $q$  in the lattice.  $\alpha(q)$  is initialized at the starting arc  $q_0$  by  $\alpha(q_0) = \pi(q_0) p(X_r(q_0) | q_0, \Lambda')$ , and  $\beta(q)$  initialized at the ending arc  $q_E$  by  $\beta(q_E) = 1$ .

The recursive computation of  $\alpha(q)$  and  $\beta(q)$  is illustrated in Fig. 2. There is direct analogy between this forward and backward probability computation over the sub-lattice illustrated here and that for the standard HMM over time [51][11]. In Fig. 2, the arc  $q$  under consideration is analogous to the HMM state occupied at current time frame  $t$  in describing the HMM’s forward-

backward algorithm, the set of arcs  $\{p : p \text{ precedes } q\}$  is analogous to all states in HMM at frame  $t-1$ , the set  $\{v : v \text{ succeeds } q\}$  is analogous to all states in HMM at frame  $t+1$ .  $X_r'(q)$  plays the role of the sequence of observation vectors from 1 to  $t-1$ , and  $X_r''(q)$  plays the role of the sequence of observation vectors from  $t+1$  to the end.  $P(q/p, \Lambda')$  is analogous to the HMM's transition probability (and its value is available from the lattice as the phone or word's "bigram language model" score).  $p(X_r(q)/q, \Lambda')$  is analogous to the HMM's emission probability (and its value is available from the lattice as the "acoustic model" score for arc  $q$ ). Given these analogies, the forward and backward probability computation for (105) and (106) as illustrated in Fig. 2 becomes identical to that for the standard HMM (as illustrated in Figs. 6.5 and 6.6 of the textbook [51]).

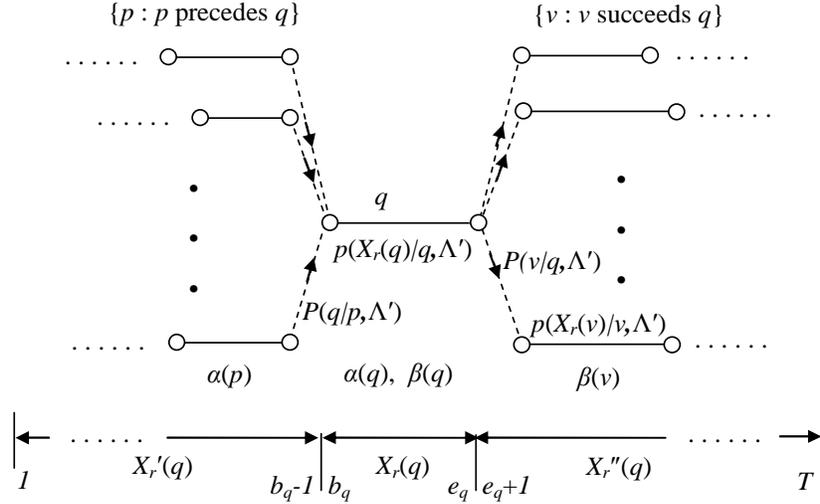


Fig. 2. Illustrations of the sub-lattice that contains arc  $q$  and of the computation of the forward and backward  $\alpha(q)$  and  $\beta(q)$  based on the sub-lattice. Each solid line represents an arc in the lattice, and each dashed line represents the direct connection between two arcs (i.e.,  $b_q-1=e_p$ ).

### B. Computing $\Delta\gamma(i, r, t)$ for MPE/MWE involving lattice approximation

We now describe how the computation burden in (92) due to the huge number of summation terms for  $s_r$  can be drastically reduced for the MPE/MWE case. It should be pointed out that (92) is a unified form for both MCE and MPE/MWE. However, due to the different properties of  $C(s_r)$  (i.e., MCE has each term as the Kronecker delta function, but not so for MPE/MWE), the lattice-based computation of (92) for MCE and MPE/MWE becomes different.

Consider a particular string token  $s_r$  that consists of a sequence of sub-tokens or sub-strings. For MCE,  $C(s_r) = \delta(s_r, S_r)$ , and hence if any of the sub-tokens is incorrect, the entire token is incorrect also. On the other hand, for MPE,  $C(s_r) = A(s_r, S_r)$ , which is the raw phone (sub-string) accuracy count in the sentence string  $s_r$ . Therefore, we have a sum of raw phone (sub-string) accuracy counts of all sub-tokens; i.e., for  $s_r = s_{r,1}, \dots, s_{r,N_r}$ , we have  $C(s_r) = \sum_{i=1}^{N_r} C(s_{r,i})$ , where  $C(s_{r,i})$  is the raw accuracy count of the sub-token  $s_{r,i}$ . Readers are referred to [48] for the computation of  $C(s_{r,i})$  for sub-token  $s_{r,i}$  in the lattice.

In this subsection, we discuss the lattice-based computation of (86) for MPE/MWE. (The lattice-based MCE will be discussed in the next subsection.)

To proceed, we define

$$\bar{C}_r = \frac{\sum_{s_r} p(X_r, s_r | \Lambda') C(s_r)}{\sum_{s_r} p(X_r, s_r | \Lambda')} \quad (107)$$

which is the average accuracy count of utterance  $r$ , given the observation sequence  $(X_r)$  and the lattice that represents all possible strings  $s_r$ .

Then, we make use of (98) to simplify (92) as follows:

$$\begin{aligned} \Delta\gamma(i, r, t) &= \sum_{s_r} p(s_r | X_r, \Lambda') (C(s_r) - \bar{C}_r) \gamma_{i,r,s_r}(t) \\ &= \sum_{q:t \in [b_q, e_q]} \sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') (C(s_r) - \bar{C}_r) \gamma_{i,r,q}(t) \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \left[ \sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') C(s_r) - \bar{C}_r \cdot \sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') \right] \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot \left[ p(q | X_r, \Lambda') \cdot \frac{\sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') C(s_r)}{p(q | X_r, \Lambda')} - \bar{C}_r \cdot p(q | X_r, \Lambda') \right] \\ &= \sum_{q:t \in [b_q, e_q]} \gamma_{i,r,q}(t) \cdot p(q | X_r, \Lambda') \cdot [\bar{C}_r(q) - \bar{C}_r] \end{aligned} \quad (108)$$

where  $p(q | X_r, \Lambda') = \sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') = \frac{p(q, X_r | \Lambda')}{p(X_r | \Lambda')}$  is computed in the same way as in (101)

and (102). In (108), we define

$$\bar{C}_r(q) = \frac{\sum_{s_r: q \in s_r} p(s_r | X_r, \Lambda') C(s_r)}{p(q | X_r, \Lambda')} = \frac{\sum_{s_r: q \in s_r} p(s_r, X_r | \Lambda') C(s_r)}{\sum_{v_r: q \in v_r} p(v_r, X_r | \Lambda')} \quad (109)$$

which is the average accuracy count of the utterance  $r$ , given observation sequence  $X_r$  and the sub-lattice that represents all strings  $s_r$  containing arc  $q$ .

The difficulty of computing  $\bar{C}_r(q)$  and  $\bar{C}_r$  in (108) lies in the very large number of terms in the summation over  $s_r: q \in s_r$  and over  $s_r$ , respectively. To efficiently compute  $\bar{C}_r(q)$  and  $\bar{C}_r$ , we now further define the following two additional ‘‘forward’’ and ‘‘backward’’ variables for each arc  $q$  (following [48]):

$$\varphi(q) = \frac{\sum_{\{s': s' \text{ precedes } q\}} p(s', q, X'_r(q), X_r(q) | \Lambda') C(s', q)}{\sum_{\{s': s' \text{ precedes } q\}} p(s', q, X'_r(q), X_r(q) | \Lambda')} \quad (110)$$

and

$$\psi(q) = \frac{\sum_{\{s'': s'' \text{ succeeds } q\}} p(s'', X''_r(q) | q, \Lambda') C(s'')}{\sum_{\{s'': s'' \text{ succeeds } q\}} p(s'', X''_r(q) | q, \Lambda')} \quad (111)$$

In (110),  $\varphi(q)$  is the weighted average accuracy count of the sub-lattice that represents all partial paths  $(s',q)$  ending inclusively in  $q$ , with the partial observation sequence  $X'_r(q) \cup X_r(q)$  (i.e.,  $x_{r,1}, \dots, x_{r,e_q}$ ). In (111),  $\psi(q)$  is the weighted average accuracy count of the sub-lattice that represents all partial paths  $s''$  that succeeds  $q$ , with the partial observation sequence  $X''_r(q)$ . Fig. 3 illustrates the sub-lattice that represents all  $s_r$  that contains arc  $q$ , together with all the relevant quantities for defining  $\varphi(q)$  and  $\psi(q)$  based on the sub-lattice. To show these quantities in defining  $\varphi(q)$ , we denote the accuracy count as  $C(s',q)$  for a given partial path  $(s',q)$  encircled by the dotted line to the left of Fig. 3. We denote the weight associated with this partial path as  $p(s',q, X'_r(q), X_r(q)|\Lambda')$ . The quantities defining  $\psi(q)$  are illustrated to the right of Fig. 3, including the partial path  $(s'')$  that is to the future of arc  $q$ , the accuracy count  $C(s'')$  associated with this path, and the associated weight of  $p(s'', X''_r(q)|\Lambda')$ .

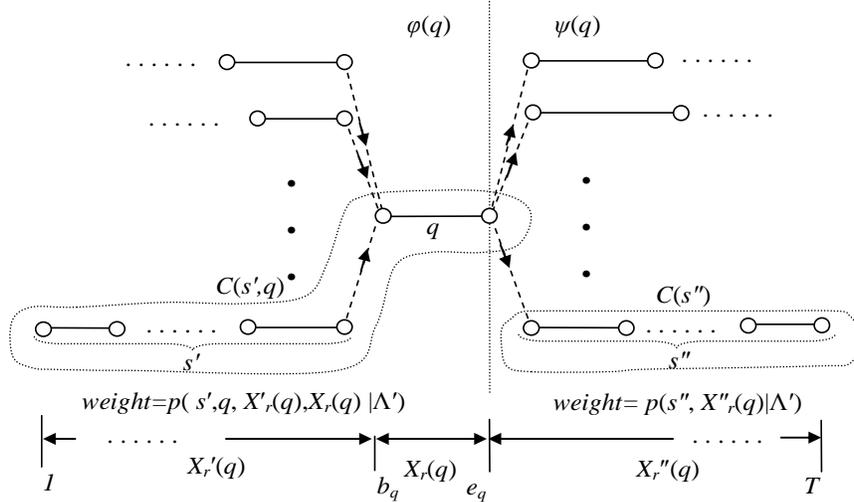


Fig. 3. Illustrations of the sub-lattice that contains arc  $q$ , and of the probability weights that define  $\varphi(q)$  of (110) and  $\psi(q)$  of (111) based on the sub-lattice.

We now describe how to compute  $\varphi(q)$  defined in (110) and  $\psi(q)$  defined in (111) efficiently for each arc  $q$  in the lattice. For  $\varphi(q)$ , we use the following efficient “forward” recursion (proof omitted):

$$\begin{aligned}
\varphi(q) &= \frac{\sum_{\{p: p \text{ precedes } q\}} P(q|p, \Lambda') p(X_r(q)|q, \Lambda') \alpha(p) [\varphi(p) + C(q)]}{\sum_{\{p: p \text{ precedes } q\}} P(q|p, \Lambda') p(X_r(q)|q, \Lambda') \alpha(p)} \\
&= \frac{\sum_{\{p: p \text{ precedes } q\}} P(q|p, \Lambda') \alpha(p) \varphi(p)}{\sum_{\{p: p \text{ precedes } q\}} P(q|p, \Lambda') \alpha(p)} + C(q)
\end{aligned} \tag{112}$$

where  $\varphi(q)$  is initialized for each starting arc  $q_0$  by  $\varphi(q_0) = C(q_0)$ , which is the raw phone or word accuracy for  $q_0$ . For  $\psi(q)$ , we use the following efficient “backward” recursion (proof omitted):

$$\psi(q) = \frac{\sum_{\{v: v \text{ succeeds } q\}} p(X_r(v) | v, \Lambda') P(v | q, \Lambda') \beta(v) [C(v) + \psi(v)]}{\sum_{\{v: v \text{ succeeds } q\}} p(X_r(v) | v, \Lambda') P(v | q, \Lambda') \beta(v)} \quad (113)$$

where  $\psi(q)$  is initialized for each ending arc  $q_E$  by  $\varphi(q_E) = 0$ .

The recursive computation of  $\varphi(q)$  in (112) is illustrated in Fig. 4. Given the partial observation sequence  $x_{r,1}, \dots, x_{r,e_q}$ ,  $[\varphi(p) + C(q)]$  is the mean accuracy count of the sub-lattice that represents all partial paths that pass  $p$  and end with  $q$ . These paths are marked by the dotted line in Fig. 4.  $\varphi(q)$  is a weighted sum and the weighted associated with each path passing arc  $p$  is  $\alpha(p)P(q|p, \Lambda')p(X_r(q)|q, \Lambda')$ , where each of the three factors is associated with each corresponding portion that makes up the path. The three factors are placed in the corresponding portions on the path in Fig. 4. The weighted average of  $[\varphi(p) + C(q)]$  over all arcs  $p$  (directly preceding  $q$ ) using the three-factor weight above gives the recursive form of  $\varphi(q)$  shown in the first line of (112). The second line of (112) removes some redundant computation and has been implemented in practice.

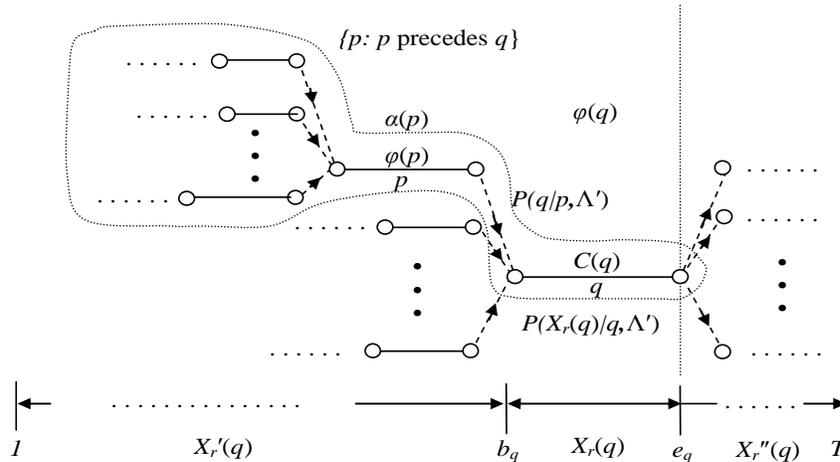


Fig. 4. Illustrations of the sub-lattice containing arc  $q$  and of the recursive  $\varphi(q)$  computation based on the sub-lattice. Each solid line represents an arc in the sub-lattice, and each dashed line represents the transition between two arcs. The dotted line encircles all partial paths that pass  $p$  and end with  $q$ .

The recursive computation of  $\psi(q)$  in (113) can be similarly interpreted as the weighted average of the accuracy count  $C(v) + \psi(v)$  for all arcs  $v$  directly following  $q$ .

Now given that both  $\varphi(q)$  and  $\psi(q)$  are computed, and assuming that arc  $q$  depends only on the arcs directly preceding it and succeeding it, we can use (110) and (111) to directly prove that  $\bar{C}_r(q) = \varphi(q) + \psi(q)$ . (114)

as one of the two quantities required to compute  $\Delta\gamma(i, r, t)$  in (108). The interpretation of (114) is offered using Fig. 3. By definition,  $\bar{C}_r(q)$  is the average accuracy count for utterance  $r$  over the sub-lattice shown in Fig. 3 that contains arc  $q$ . This count can be decomposed into two parts. The first part is the “forward” average accuracy count of the left part of the sub-lattice in Fig. 3 for the utterance from  $t=1$  to  $e_q$ , which is  $\varphi(q)$ . The second part is the “backward” average accuracy count of the right part of the sub-lattice for the utterance from  $t=e_q+1$  to  $T_r$ , which is  $\psi(q)$ .

The second quantity,  $\bar{C}_r$ , required to compute  $\Delta\gamma(i, r, t)$  in (108) can be proved to be

$$\bar{C}_r = \frac{\sum_{q:q \in \{\text{ending arcs}\}} \varphi(q)\alpha(q)}{\sum_{q:q \in \{\text{ending arcs}\}} \alpha(q)} \quad (115)$$

The interpretation of (115) is as follows. Let arc  $q$  be an ending arc in the lattice. And recall  $\varphi(q)$  is the average accuracy count of utterance  $r$  given the sub-lattice that represents all  $s_r$  containing (sub-lattice-ending) arc  $q$ , and  $\alpha(q)$  is the weight of this sub-lattice. Therefore,  $\bar{C}_r$ , which is defined in (107) as the average accuracy count for the entire lattice, becomes a weighted sum of the average accuracy counts of all sub-lattices as shown in (115).

This completes the description of the computation of  $\Delta\gamma(i, r, t)$  in (108).

### C. Computing $\Delta\gamma(i, r, t)$ for MCE involving lattice approximation

Finally, we discuss using (98) to compute  $\Delta\gamma(i, r, t)$  of (92) for MCE. As we mentioned earlier, while (92) is unified between MPE and MCE, the specific form of  $C(s_r) = \delta(s_r, S_r)$  in MCE permits special simplification of  $\Delta\gamma(i, r, t)$  of (92) for MCE. The simplification steps, followed by the use of (98), lead to

$$\begin{aligned} \Delta\gamma(i, r, t) &= \sum_{s_r} p(s_r | X_r, \Lambda') [C(s_r) - p(S_r | X_r, \Lambda')] \gamma_{i, r, s_r}(t) \\ &= p(S_r | X_r, \Lambda') \gamma_{i, r, S_r}(t) - p(S_r | X_r, \Lambda') \left[ \sum_{s_r} p(s_r | X_r, \Lambda') \gamma_{i, r, s_r}(t) \right] \\ &= p(S_r | X_r, \Lambda') \gamma_{i, r, S_r}(t) - p(S_r | X_r, \Lambda') \left[ \sum_{q:q \in [b_q, e_q]} \gamma_{i, r, q}(t) \cdot \sum_{s_r:q \in s_r} p(s_r | X_r, \Lambda') \right] \\ &= p(S_r | X_r, \Lambda') \gamma_{i, r, S_r}(t) - p(S_r | X_r, \Lambda') \left[ \sum_{q:q \in [b_q, e_q]} \gamma_{i, r, q}(t) \cdot p(q | X_r, \Lambda') \right] \\ &= p(S_r | X_r, \Lambda') \left[ \frac{\gamma_{i, r, S_r}(t)}{\gamma_{i, r}^{\text{num}}(t)} - \sum_{q:q \in [b_q, e_q]} \frac{\gamma_{i, r, q}(t) \cdot p(q | X_r, \Lambda')}{\gamma_{i, r}^{\text{den}}(t)} \right] \end{aligned} \quad (116)$$

The last line shows a striking similarity between the lattice-based MCE and MMI. In (116),  $p(q | X_r, \Lambda') = \frac{p(q, X_r | \Lambda')}{p(X_r | \Lambda')}$  is computed by (101) and (102) for the numerator and denominator,

respectively. Also in (116), we have  $p(S_r | X_r, \Lambda') = \frac{p(X_r | S_r, \Lambda') p(S_r | \Lambda')}{p(X_r | \Lambda')}$ , where correct

string  $S_r$  is known. Hence  $\gamma_{i, r, S_r}(t)$  and  $p(X_r | S_r, \Lambda')$  in (116) can be efficiently computed by the standard forward-backward algorithm for the HMM [51]. Finally, for the computation of  $p(S_r | \Lambda')$  and  $p(X_r | \Lambda')$ , we use the language model and  $\sum_{q:q \in \{\text{ending arcs}\}} \alpha(q)$ , respectively.

Note that the computation for the lattice-based MCE provided in (116) does not require removing the correct word string  $S_r$  from the lattice.

#### APPENDIX IV: EXISTENCE PROOF OF FINITE $D$ IN GT UPDATES FOR CDHMM

As shown in Section VI.B, Jensen's inequality based optimization cannot be applied directly to Gaussian CDHMM since the value  $D_i$  in the GT update formulas (71) and (72) may be infinite, making the algorithm's convergence infinitely slow. In this Appendix, we follow the insight provided in [2] to prove that there exist finite values of  $D_i$  which make the GT update formulas (71) and (72) practical for all MMI, MCE, and MPE/MWE.

To proceed the proof, we substitute (65) into (62) and obtain

$$V(\Lambda; \Lambda') = \sum_q \sum_s \int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') \log p(\chi, q | s, \Lambda) + Const. \quad (117)$$

We prove below that for CDHMM, given a sufficiently large but bounded (i.e., finite) constant  $D$ ,

$$F(\Lambda; \Lambda') - F(\Lambda'; \Lambda') \geq V(\Lambda; \Lambda') - V(\Lambda'; \Lambda') \quad (118)$$

First, we define

$$\Delta_D = [F(\Lambda; \Lambda') - F(\Lambda'; \Lambda')] - [V(\Lambda; \Lambda') - V(\Lambda'; \Lambda')] \quad (119)$$

and will show that  $\Delta_D \geq 0$  for any parameter set  $\Lambda$ . Substituting (64) and (62) into (119), we obtain

$$\begin{aligned} \Delta_D &= [F(\Lambda; \Lambda') - F(\Lambda'; \Lambda')] - [V(\Lambda; \Lambda') - V(\Lambda'; \Lambda')] \\ &= \sum_q \sum_s \int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] [p(\chi, q | s, \Lambda) - p(\chi, q | s, \Lambda')] d\chi \\ &\quad - \sum_q \sum_s \int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') [\log p(\chi, q | s, \Lambda) - \log p(\chi, q | s, \Lambda')] d\chi \\ &= \sum_q \sum_s \int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') \left[ \frac{p(\chi, q | s, \Lambda)}{p(\chi, q | s, \Lambda')} - 1 - \log \frac{p(\chi, q | s, \Lambda)}{p(\chi, q | s, \Lambda')} \right] d\chi \\ &= \sum_q \sum_s \int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi \end{aligned} \quad (120)$$

$$\text{where } H(\chi, q, s, \Lambda, \Lambda') = \left[ \frac{p(\chi, q | s, \Lambda)}{p(\chi, q | s, \Lambda')} - 1 \right] - \log \left[ \left[ \frac{p(\chi, q | s, \Lambda)}{p(\chi, q | s, \Lambda')} - 1 \right] + 1 \right].$$

Then, we need to show that there exists a bounded  $d(s)$  that ensures the summand of  $\Delta_D$  in (120) be non-negative. To proceed, we expand the summand to

$$\begin{aligned} &\int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi \\ &= p(s) [C(s) - O(\Lambda')] p(X, q | s, \Lambda') H(X, q, s, \Lambda, \Lambda') + d(s) \int_{\mathcal{X}} p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi \end{aligned} \quad (121)$$

We now use the following key theorem from [2]: If  $f(X, \Lambda)$  is non-negative and analytic for  $X \in \mathcal{X}$  and  $\Lambda \in \Omega$ , where  $\mathcal{X}$  and  $\Omega$  are the data space and model space, respectively, then there is a  $\Lambda$ -independent constant  $K > 0$  such that

$$\int_{\mathcal{X}} f(\chi, \Lambda) d\chi \geq K f(X, \Lambda) \quad (122)$$

for any valid model  $\Lambda$ . (Readers are referred to [2] for a rigorous proof.)

Define  $f(X, \Lambda) = p(X, q | s, \Lambda') H(X, q, s, \Lambda, \Lambda')$ . Here  $f(X, \Lambda)$  is non-negative and analytic since both  $p(X, q | s, \Lambda')$  and  $H(X, q, s, \Lambda, \Lambda')$  are non-negative and analytic (for CDHMM). Using (122), we have

$$\int_{\mathcal{X}} p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi \geq K p(X, q | s, \Lambda') H(X, q, s, \Lambda, \Lambda') \quad (123)$$

Now we construct non-negative  $d(s)$  as follows:

$$d(s) = \begin{cases} 0 & \text{if } C(s) \geq O(\Lambda') \\ \frac{1}{K} p(s) (O(\Lambda') - C(s)) & \text{if } C(s) < O(\Lambda') \end{cases} \quad (124)$$

Then, (123) becomes

$$d(s) \int_{\mathcal{X}} p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi > -p(s) [C(s) - O(\Lambda')] p(X, q | s, \Lambda') H(X, q, s, \Lambda, \Lambda')$$

This proves that the summand of  $\Delta_D$ ,  $\int_{\mathcal{X}} [\Gamma(\Lambda') + d(s)] p(\chi, q | s, \Lambda') H(\chi, q, s, \Lambda, \Lambda') d\chi$ , is non-negative for any  $s$  (according to (121)), and therefore  $\Delta_D \geq 0$ .

Given (124) and (122),  $D_i$  can be computed according to (59) which is a sufficiently large but bounded value.

## APPENDIX V: UNIFYING DISCRIMINATIVE TRAINING CRITERION USING THE REVERSE JENSEN INEQUALITY

While the proof in Appendix IV above (based on Axelrod et. al.'s work [2]) shows that the GT update formulas for CDHMM are valid given a sufficiently large (but bounded) constant  $D_i$ , no explicit construction of  $D_i$  was given. Therefore, it was an existence proof. In this Appendix, we will outline a constructive proof based on Jebara's work [26][27].

In principle, Jebara's method is applicable to maximizing any rational function, whose numerator and denominator are a mixture of exponential models. Therefore, it is applicable to optimizing our unified discriminative criterion of (26) for all MMI, MCE, and MPE/MWE criteria. In this brief review of Jebara's method, we will introduce the principle of the reverse Jensen inequality and its application to discriminative objective function optimization.

For a rational function in the form of (27) and (28), we desire to maximize the following equivalent function:

$$\log O(\Lambda) = \log G(\Lambda) - \log H(\Lambda) = \log \sum_s p(X, s | \Lambda) C(s) - \log \sum_s p(X, s | \Lambda) \quad (125)$$

The first term of the right-hand side of (125) is a log-sum function similar to log likelihood. Based on the well known Jensen's inequality and with several steps of simplifications, we have

$$\begin{aligned} \log \sum_s p(X, s | \Lambda) C(s) &\geq Q_G(\Lambda; \Lambda') + k \\ &= \sum_s \left( \frac{p(X, s | \Lambda') C(s)}{\sum_{s'} p(X, s' | \Lambda') C(s')} \right) \log p(X, s | \Lambda) C(s) + k, \end{aligned} \quad (126)$$

where  $k$  is a constant irrelevant to  $\Lambda$  (although relevant to  $\Lambda'$ ), and equality holds when  $\Lambda = \Lambda'$ . This is similar to the E-step in the EM algorithm.

The left hand side of (126) is a lower bound of  $\log G(\Lambda)$ , and makes tangential contact with  $\log G(\Lambda)$  at  $\Lambda'$ . Therefore maximizing the auxiliary function  $Q_G(\Lambda; \Lambda')$  guarantees increase of  $\log G(\Lambda)$  iteratively.

However, in order to maximize  $\log O(\Lambda)$ , we need a lower bound for  $\log O(\Lambda)$ , which in turn requires an upper bound of  $\log H(\Lambda)$ . In [26], it was shown (non-trivially) that using reverse Jensen's inequality an auxiliary function  $Q_H(\Lambda; \Lambda')$  can be constructed so that

$$\log \sum_s p(X, s | \Lambda) \leq Q_H(\Lambda; \Lambda') + \tilde{k} = \sum_{Q_H(\Lambda; \Lambda')} (-w_s) \log p(Y_s, s | \Lambda) + \tilde{k} \quad (127)$$

where  $\tilde{k}$  is a  $\Lambda$ -irrelevant constant that makes  $Q_H(\Lambda; \Lambda') + \tilde{k}$  tangential contact with  $\log G(\Lambda)$  at  $\Lambda'$ ,  $w_s$  and  $Y_s$  are positive weights and modified observations, respectively. The reverse Jensen inequality was derived by exploiting the convexity of the cumulant generating function of exponential family in [26] and will not be elaborated further here.

Given (126) and (127), one can construct the auxiliary function for  $\log O(\Lambda)$  as:

$$Q(\Lambda; \Lambda') = Q_G(\Lambda; \Lambda') - Q_H(\Lambda; \Lambda') \quad (128)$$

Then optimizing  $\log O(\Lambda)$  can be achieved by iteratively optimizing  $Q(\Lambda; \Lambda')$ , which takes the same step as the M-step in the conventional EM algorithm for an HMM (i.e., with a closed-form solution in the M-step).

Note for our unifying discriminative objective function (26), the summand of  $G(\Lambda)$  may take negative value for MPE; i.e., for some path  $s$  that has many insertion errors, the corresponding  $\underline{C}(s)$  may be negative. In this case, we can add extra dummy training tokens to the training set, while these dummy tokens can only be recognized as correct references. Appending these dummy tokens to  $s$  can effectively increase its raw accuracy count to be positive. Moreover, since the dummy token will not compete with any other tokens in the training set, it will not affect the training performance.

## APPENDIX VI: BAYES RISK AND MPE

In [54], Schlüter showed that MMI and MCE criteria are an upper bound of the true Bayes sentence- or string-level error rate. Extending that work, we show here that the MPE criterion, after being normalized by the number of training utterance, is an upper bound of the true Bayes risk at the sub-string (e.g., phone) level. In order to compare the MPE criterion with the true Bayes risk, we denote by  $P(s_r | X_r)$  the true posterior probability, and  $P(s_r | X_r, \Lambda)$  the model-based posterior probability. Given the true posterior probability  $P(s | X)$ , The Bayes risk of classifying  $X$  to  $s$  is

$$R(s | X) = \sum_{s'} P(s' | X) \cdot e_{s, s'}$$

where  $e_{s, s'}$  is the loss of classifying  $X$  to class  $s$  while  $X$  actually has reference class  $s'$ . For zero-one loss (e.g., sentence error),  $e_{s, s'} = \delta(s, s')$ . However, for MPE, the error is counted on the phone level. Therefore, the conventional zero-one loss is not suitable. Instead, we need to use a non zero-one loss. For classification tasks that count errors on the sub-string level (e.g., speech recognition for words or phones), we use non zero-one loss for  $e_{s, s'}$  as the raw error counts, including deletion, insertion, and substitution errors of the recognized string  $s$  given the reference string  $s'$ . Specifically, in the MPE framework,  $e_{s, s'} = |s'| - A(s, s')$ , where  $|s'|$  is the number of sub-strings in reference  $s'$ , and  $A(s, s')$  is the raw accuracy as defined for MPE [48]. Therefore, minimum Bayes risk becomes

$$\begin{aligned}
R_{Bayes}(e | X) &= \min_s R(s | X) \\
&= \min_s \left[ \sum_{s'} P(s' | X) \cdot |s'| - \sum_{s'} P(s' | X) A(s, s') \right] \\
&= \sum_{s'} P(s' | X) \cdot |s'| - \max_s \left[ \sum_{s'} P(s' | X) A(s, s') \right]
\end{aligned}$$

Then, the expectation of the true Bayes error is

$$\begin{aligned}
R_{Bayes}(e) &= \lim_{R \rightarrow \infty} \frac{1}{R} \sum_{r=1}^R \left\{ \sum_{s'} P(s' | X) \cdot |s'| - \max_s \sum_{s'} P(s' | X_r) A(s, s') \right\} \\
&= \int dX \cdot p(X) \cdot \left\{ \sum_{s'} P(s' | X) \cdot |s'| - \max_s \sum_{s'} P(s' | X) A(s, s') \right\} \\
&= \int dX \cdot p(X) \cdot \left\{ \sum_{s'} P(s' | X) \cdot |s'| \right\} - \int dX \cdot p(X) \cdot \left\{ \max_s \sum_{s'} P(s' | X) A(s, s') \right\} \\
&\quad \underbrace{1 \ 4 \ 4 \ 4 \ 4 \ 4 \ 4 \ 2 \ 4 \ 4 \ 4 \ 4 \ 4 \ 4 \ 3}_{\text{expectation of number of sub-strings (e.g., phones)}}
\end{aligned}$$

On the other hand, the model-based posterior probability for string  $s$  given the observation sequence  $X$  is denoted by  $P(s | X, \Lambda)$ . In the limit of an infinite amount of training data, the normalized MPE criterion then becomes

$$\begin{aligned}
\mathcal{O}_{MPE}(\Lambda) &= \lim_{R \rightarrow \infty} \frac{1}{R} \sum_{r=1}^R \sum_{s'} P(s' | X_r, \Lambda) A(s', S_r) \\
&= \sum_s \int dX \cdot p(X, s) \cdot \sum_{s'} P(s' | X, \Lambda) A(s', s) \\
&= \int dX \cdot p(X) \sum_s P(s | X) \cdot \sum_{s'} P(s' | X, \Lambda) A(s', s) \\
&= \int dX \cdot p(X) \cdot \sum_{s'} P(s' | X, \Lambda) \sum_s P(s | X) A(s', s) \\
&\leq \int dX \cdot p(X) \cdot \sum_{s'} P(s' | X, \Lambda) \cdot \left[ \max_s \sum_s P(s | X) A(s', s) \right] \\
&\quad \underbrace{1 \ 4 \ 4 \ 2 \ 4 \ 4 \ 3}_{=1} \quad \underbrace{1 \ 4 \ 4 \ 4 \ 4 \ 2 \ 4 \ 4 \ 4 \ 3}_{\text{independent with } s' \text{ after taking max}} \\
&= \int dX \cdot p(X) \cdot \max_{s'} \sum_s P(s | X) A(s', s)
\end{aligned}$$

Since  $\mathcal{O}_{MPE}(\Lambda)$  is actually the expectation of utility instead of loss, the complementary loss incurred by normalized MPE is

$$\begin{aligned}
R_{MPE}(e | \Lambda) &= \int dX \cdot p(X) \cdot \left\{ \sum_{s'} P(s' | X) \cdot |s'| \right\} - \mathcal{O}_{MPE}(\Lambda) \\
&\geq \int dX \cdot p(X) \cdot \left\{ \sum_{s'} P(s' | X) \cdot |s'| \right\} - \int dX \cdot p(X) \cdot \left\{ \max_s \sum_{s'} P(s' | X) A(s, s') \right\} \\
&= R_{Bayes}(e)
\end{aligned}$$

This proves that the MPE criterion provides an upper bound of true (model-independent) minimum Bayes risk on the sub-string level. Moreover, if we can minimize  $R_{MPE}(e | \Lambda)$  to

achieve equality in the above inequality, i.e., if the model distribution  $P(s'|X, \Lambda)$  can be estimated to be

$$P(s|X, \Lambda) = \begin{cases} 1 & \text{iff } s = \max_{s'} \sum_{s''} P(s'|X) A(s'', s') \\ 0 & \text{otherwise} \end{cases}$$

then the minimum Bayes risk can be obtained.

It is important to point out that, for the non zero-one loss as we discussed here, the conventional maximum *a posteriori* probability classifier may not lead to minimum *Bayes* risk at the sub-string level. To address this issue, Minimum Bayes Risk decoding has been developed in [20], to which we would like to refer the interested readers for technical details.

## REFERENCES

- [1] S. Amari. "A theory of adaptive pattern classifiers," IEEE Trans. Electronic Computers, Vol. 16, 1967, pp. 299-307.
- [2] S. Axelrod, V. Goel, R. Gopinath, P. Olsen, and K. Visweswariah. "Discriminative Estimation of Subspace Constrained Gaussian Mixture Models for Speech Recognition," IEEE Trans. Audio, Speech & Language Proc., Vol. 15, 2007, pp. 172-189.
- [3] L. E. Baum and J.A. Eagon, "An inequality with applications to statistical prediction for functions of Markov processes and to a model of ecology", Bull. Amer. Math. Soc. vol. 73, pp. 360-363, 1967.
- [4] L. Baum and G. Sell. "Growth transformations for functions on manifolds", Pacific J. Math., Vol.27, no.2, pp.211-227, 1968.
- [5] R. Battiti, "First- and Second- Order Methods for Learning: Between Steepest Descent and Newton's Method, Neural Computation, 1992, vol. 4, pp. 141-166
- [6] E. Birney. "Hidden Markov models in biological sequence analysis," IBM Journal Research & Dev., Vol. 45, 449-454, 2001.
- [7] P. Brown. "The Acoustic Modeling Problem in Automatic Speech Recognition," Ph.D. thesis, Carnegie Mellon University, 1987.
- [8] W. Chou and L. Li. "A minimum classification error framework for generalized linear classifier in machine learning and its application to text categorization/retrieval", Proc. Int. Conf. Machine Learning & App., pp. 382-390, Dec. 2004.
- [9] W. Chou. "Minimum classification error approach in pattern recognition," in Pattern Recognition in Speech and Language Processing, (W. Chou and B.-H. Juang eds.) CRC Press, Boca Raton, pp. 1-49, 2003.
- [10] M. Collins. "Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms," Proc. Empirical Methods of Natural Language Processing (EMNLP), 2005, pp. 1-8.
- [11] L. Deng and D. O'Shaughnessy. Speech Processing --- A Dynamic and Optimization-Oriented Approach, Marcel Dekker Publishers, New York, NY, 2003.
- [12] L. Deng, D. Yu, and A. Acero. "A generative modeling framework for structured hidden speech dynamics," Neural Information Processing System (NIPS) Workshop, Whistler, BC, Canada, Dec. 2005.
- [13] L. Deng, J. Wu, J. Droppo, and A. Acero. "Analysis and comparison of two feature extraction/compensation algorithms," IEEE Signal Processing Letters, Vol. 12, No. 6, June, 2005, pp. 477-480.

- [14] J. E. Dennis and R. B. Schnabel, "Numerical methods for unconstrained optimization and nonlinear equations", SIAM's Classics in Applied Mathematics, 1996
- [15] J. Droppo and A. Acero. "Joint Discriminative Front End and Back End Training For Improved Speech Recognition Accuracy." Proc. ICASSP 2006.
- [16] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*, Cambridge University Press, 1998.
- [17] S. E. Fahlman, "An empirical study of learning speed in back-propagation networks" Canergie Mellon Univ., Pittsburgh, PA, 1988, Tech. Rep.
- [18] P. Gopalakrishnan, D. Kanevsky, A. Nadas, and D. Nahamoo, "An inequality for rational functions with applications to some statistical estimation problems," *IEEE Trans. Information Theory.*, Vol 37, January. 1991, pp.107 – 113.
- [19] Y. Gao and J. Kuo. "Maximum entropy direct models for speech recognition," *IEEE Trans. Speech & Audio Proc.*, 2006.
- [20] V. Goel and W. Byrne. "Minimum Bayes Risk Methods in Automatic Speech Recognition," in *Pattern Recognition in Speech and Language Processing*, (W. Chou and B.-H. Juang eds.) CRC Press, Boca Raton, 2003.
- [21] A. Gunawardana and W. Byrne, "Discriminative speaker adaptation with conditional maximum likelihood linear regression," Proc. EUROSPEECH, 2001.
- [22] A. Gunawardana, M. Mahajan, A. Acero, and J. Platt. "Hidden conditional random fields for phone classification," Proc. Interspeech, Lisbon, 2005.
- [23] X. He and W. Chou, "Minimum classification error linear regression for acoustic model adaptation of continuous density HMMs," Proc. ICASSP, April 2003.
- [24] X. He, L. Deng, and W. Chou. "A novel learning method for hidden Markov models in speech and audio processing," Proc. *IEEE Workshop on Multimedia Signal Processing*, Victoria, BC, October 2006.
- [25] X. Huang, A. Acero and H. Hon. *Spoken Language Processing*. Prentice Hall, 2001.
- [26] T. Jebara. *Discriminative, Generative and Imitative Learning*, PhD Thesis, Media Laboratory, MIT, December 2001.
- [27] T. Jebara, and A. Pentland. "On Reversing Jensen's Inequality". In *Neural Information Processing Systems Vol. 13*, December 2000.
- [28] B.-H. Juang, and S. Katagiri. "Discriminative learning for minimum error classification, *IEEE Trans. Signal Processing.*, Vol. 40, No. 12, 1992, pp. 3043-3054.
- [29] B.-H. Juang, W. Chou, and C.-H. Lee, "Minimum classification error rate methods for speech recognition," *IEEE Trans. Speech & Audio Proc.*, Vol. 5, May 1997.
- [30] D. Kanevsky, "A generalization of the Baum algorithm to functions on non-linear manifolds", Proc. ICASSP, 1995.
- [31] S. Kapadia, V. Valtchev, and S. J. Young, "MMI training for continuous phoneme recognition on the TIMIT database," in Proc. ICASSP, 1993, vol. 2, pp. 491–494.
- [32] J. Lafferty, A. McCallum, and F. Pereira. "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," Proc. *Intern. Conf. Machine Learning*, 2001, pp. 282-289.
- [33] Y. Li, L. Shapiro, and J. Bilmes. "A generative/discriminative learning algorithm for image classification," *IEEE Int. Conf. Computer Vision*, Beijing, China, pp. 1605-1612, 2005.
- [34] A. McCallum, D. Freitag, and F. Pereira. "Maximum entropy Markov models for information extraction and segmentation." Proc. *Intern. Conf. Machine Learning*, 2000, pp. 591-598.

- [35] W. Macherey, L. Haferkamp, R. Schlüter, and H. Ney, "Investigations on error minimizing training criteria for discriminative training in automatic speech recognition," Proc. Interspeech, 2005, Lisbon, pp. 2133-2136.
- [36] E. McDermott. "Discriminative training for automatic speech recognition using the minimum classification error framework," Neural Information Processing System (NIPS) Workshop, Whistler, BC, Canada, Dec. 2005.
- [37] E. McDermott and T. Hazen, "Minimum Classification Error Training Of Landmark Models For Real-Time Continuous Speech Recognition," Proc. ICASSP 2006
- [38] E. McDermott, T. Hazen, J. Le Roux, A. Nakamura, and S. Katagiri, "Discriminative training for large vocabulary speech recognition using minimum classification error" IEEE Trans. Speech and Audio Processing, Vol. 15. No. 1, 2007, pp. 203-223.
- [39] M. Mahajan, A. Gunawardana and A. Acero, "Training Algorithms for Hidden Conditional Random Fields," Proc. ICASSP 2006
- [40] A. Nádas, "A Decision Theoretic Formulation of a Training Problem in Speech Recognition and a Comparison of Training by Unconditional versus Conditional Maximum Likelihood", IEEE Trans. on Acoustics, Speech and Signal Processing, 1983
- [41] A. Nádas, D. Nahamoo, and M. A. Picheny, "On a Model-Robust Training Method for Speech Recognition," IEEE Trans. on Acoustics, Speech and Signal Processing, 1988
- [42] Y. Normandin, "Hidden Markov Models, Maximum Mutual Information Estimation, and the Speech Recognition Problem," Ph.D. dissertation, McGill University, Montreal, 1991.
- [43] Y. Normandin. "Maximum mutual information estimation of hidden Markov models", in Automatic Speech and Speaker Recognition, (C.-H. Lee, F.K. Soong and K.K. Paliwal eds.) Kluwer Academic Publishers, Norwell, MA, 1996.
- [44] F. Och. "Minimum error rate training in statistical machine translation," Proc. 41th Meeting Assoc. Comp. Linguistics, pp. 160-167, 2003.
- [45] F. Pereira. "Linear models for structure prediction", Proc. Interspeech, Lisbon, 2005, pp. 717-720.
- [46] D. Povey and P.C. Woodland, "Minimum phone error and I-smoothing for improved discriminative training," Proc. ICASSP, 2002.
- [47] D. Povey, M.J.F. Gales, D.Y. Kim, and P.C. Woodland, "MMI-MAP and MPE-MAP for acoustic model adaptation," Proc. Eurospeech, 2003.
- [48] D. Povey, "Discriminative Training for large Vocabulary Speech Recognition," Ph.D. dissertation, Cambridge University, Cambridge, UK, 2004.
- [49] D. Povey, B. Kingsbury, L. Mangu, G. Saon, H. Soltau, and G. Zweig. "fMPE: Discriminatively trained features for speech recognition," Proc. DARPA EARS RT-04 Workshop, Nov. 7-10, 2004, Palisades, NY, Paper No. 35, 5 pages.
- [50] C. Rathinavalu and L. Deng. "Speech trajectory discrimination using the minimum classification error learning," IEEE Trans. Speech and Audio Processing, Vol.6, 1998, pp. 505-515.
- [51] L. Rabiner and B.-H. Juang. Fundamentals of Speech Recognition. Prentice Hall, 1993.
- [52] M. Riedmiller and H. Braun, "A direct adaptive method for faster back propagation learning: The RPROP algorithm," in Proc. IEEE Intl. Conf. Neural Networks, San Francisco, CA, 1993, pp. 586-591
- [53] J. Le Roux and E. McDermott, "Optimization for discriminative training," Proc. INTERSPEECH, Lisbon, 2005.
- [54] R. Schlüter. "Investigations on Discriminative Training Criteria," Ph.D. dissertation, RWTH Aachenm University of Technology, Aachen, Germany, 2000.

- [55] R. Schlüter, W. Macherey, B. Muller, and H. Ney. “Comparison of discriminative training criteria and optimization methods for speech recognition,” *Speech Communication*, Vol. 34, 2001, pp. 287-310.
- [56] S. Sun, C. Greenwood, and R. Neal. “Haplotype inference using a Bayesian hidden Markov model,” *Genetic Epidemiology*, Vol. 31, 2007, pp. 1-12.
- [57] S. Theodoridis and K. Koutroumbas. *Pattern Recognition*. Elsevier Science, Academic Press, 2003.
- [58] V. Vapnik. *Statistical Learning Theory*, Wiley-Interscience, 1998.
- [59] V. Valtchev. *Discriminative Methods for HMM-based Speech Recognition*, Ph.D. Thesis, Cambridge University, England, 1995.
- [60] V. Valtchev, J.J. Odell, P.C. Woodland, S.J. Young. “MMIE training of large vocabulary recognition systems”, *Speech Communication*, Vol. 22 (4), 1997.
- [61] V. Valtchev, J.J., Odell, P.C., Woodland, S.J. Young, “Lattice-based discriminative training for large vocabulary speech recognition”, *Proc. ICASSP*, 1996.
- [62] P. Woodland and D. Povey, “Large scale discriminative training for speech recognition,” *Proc. ITRW ASR*, 2000.
- [63] Y. Xiong, Q. Huo, C. Chan. “A discrete contextual stochastic model for the offline recognition of handwritten Chinese characters” *IEEE Trans. PAMI*, Vol. 23, pp. 774-782, 2001.
- [64] R. Yan, J. Zhang, J. Yang, and A. Hauptmann. “A discriminative learning framework with pairwise constraints for video object classification,” *IEEE Trans. PAMI*, Vol 28, No.4, pp.578-593, 2006.
- [65] J. Yang, Y. Xu, and C.S. Chen. “Hidden Markov model approach to skill learning and its application in telerobotics,” *IEEE Trans. on Robotics & Automation*, Vol. 10, No.5, pp. 621-631, 1994.
- [66] C. Yen, S. Kuo, C.-H. Lee. “Minimum error rate training for PHMM-based text recognition,” *IEEE Transactions on Image Processing*,” Vol. 8, pp. 1120 – 1124, 1999.
- [67] D. Yu, L. Deng, and A. Acero. “A\* lattice search algorithm for a long-contextual-span hidden trajectory model and phonetic recognition,” *Proc. Interspeech*, Lisbon, pp. 553-556, 2005.
- [68] D. Yu, L. Deng, X. He, A. Acero. “Large-margin minimum classification error training for large-scale speech recognition tasks,” *Proceedings of ICASSP*, Honolulu, Hawaii, April 2007.