



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Speech Communication 45 (2005) 455–470

SPEECH
COMMUNICATION

www.elsevier.com/locate/specom

Confidence measures for speech recognition: A survey

Hui Jiang *

Department of Computer Science, York University, 4700 Keele Street, Toronto, Ont., Canada M3J 1P3

Received 3 August 2004; received in revised form 26 November 2004; accepted 27 December 2004

Abstract

In speech recognition, confidence measures (CM) are used to evaluate reliability of recognition results. A good confidence measure can largely benefit speech recognition systems in many practical applications. In this survey, I summarize most research works related to confidence measures which have been done during the past 10–12 years. I will present all these approaches as three major categories, namely CM as a combination of predictor features, CM as a posterior probability, and CM as utterance verification. Then, I also introduce some recent advances in the area. Moreover, I will discuss capabilities and limitations of the current CM techniques and generally comment on today's CM approaches. Based on the discussion, I will conclude the paper with some clues for future works.

© 2005 Elsevier B.V. All rights reserved.

Keywords: Automatic speech recognition (ASR); Confidence measures (CM); Word posterior probability; Utterance verification; Likelihood ratio testing (LRT); Bayes factors

1. Introduction

Automatic speech recognition (ASR) has achieved some substantial successes in past few decades mostly attributing to two prevalent technologies in the field, namely hidden Markov modeling (HMM) of speech signals and efficient dynamic programming search (also known as *decoding*) techniques for very-large-scale networks. Today, in many aspects, it has become a standard

routine to build a state-of-the-art speech recognition system for any particular task if sufficient training data is provided for the target domain. However, when we migrate speech recognition systems from laboratory demonstrations to real-world applications, even the best ASR systems available today still encounter some serious difficulties. First of all, system performance usually dramatically degrades in the real fields because of ambient noises, speaker variations, channel distortions and many other mismatches. How to maintain and/or improve ASR performance in real-field conditions has been extensively studied in speech community under the topic of *robust*

* Tel.: +1 416 736 2100x33346; fax: +1 416 736 5872.

E-mail address: hj@cs.yorku.ca

speech recognition. Many good tutorial and overview papers, such as Juang (1991), Gong (1995), Lee (1998b) and many others, can be easily found in the literature with regard to this topic. Secondly, since every speech recognizer inevitably will make some mistakes during recognition, outputs from any ASR system are always fraught with a variety of errors. Thus, in any real-world application, it is extremely important to be able to make an appropriate and reliable judgement based on the error-prone ASR results. This requires the ASR systems to automatically assess reliability or probability of correctness for every decision made by the systems. Nowadays, to certain degree, the capability to evaluate reliability of speech recognition results has been regarded as a crucial technique to increase usefulness and “intelligence” of an ASR system in many practical applications. In this area, researchers have proposed to compute a score (preferably between 0 and 1), called *confidence measure* (CM), to indicate reliability of any recognition decision made by ASR systems. For example, a CM can be computed for every recognized word to indicate how likely it is correctly recognized or for an utterance to indicate how much we can trust the results for the utterance as a whole. Despite a large amount of research efforts in the past, we still believe that robust speech recognition and confidence measure will remain as two most active and influential research topics in speech community for a foreseeable future. Due to importance of CM in ASR systems, it has attracted considerable research attention from most major speech research groups all over the world and an excessive amount of research works have been reported in the past decade. But, unlike robust speech recognition, so far we have not seen too many overview papers in the literature to survey this important and active topic. This largely motivates me to write a comprehensive survey to summarize the CM-related research works reported mostly in the past 10–12 years. In the survey, I will mainly highlight the major progresses we have achieved in the CM area during the past decade. And I will stress some promising CM computation approaches which are theoretically sound and experimentally superior, and also discuss their capabilities and limitations. Finally, I will present

some comparative discussions with respect to all reported CM computation methods and conclude the paper with some clues for possible future works from my personal perspective. Throughout the paper, I will attempt to present the CM techniques from a fairly high level and avoid technical and experimental details as much as possible, for which readers may wish to refer to the original papers. At the end of this paper, I also compose a comprehensive list of reference papers for the convenience of readers, which includes most of published works relevant to confidence measures in ASR. To my best knowledge, Lee (2001) seems to be the only CM-related overview paper which gives some good tutorials on statistical nature of confidence measure problems and also enumerates many potential CM applications for ASR.

First of all, we can backtrack some early research works on confidence measure (CM) to non-keyword rejection in word-spotting systems which were proposed to handle unconstrained speech inputs, such as Wilpon et al. (1990), Mathan and Miclet (1991), Chigier (1992), Rose (1992), Sukkar and Wilpon (1993), etc. In these works, they first adopted the so-called *garbage* or *sink* models to explicitly model non-keywords, extraneous speech and background noises in unconstrained input utterances, with which keyword spotting systems first recognize speech inputs to detect all embedded keywords as well as other speech segments corresponding to non-keywords or noises. Besides all of these, they all noticed a need to build additional rejection module to effectively distinguish non-keywords from the detected keywords in order to reduce false alarms in non-keyword rejection. Apparently, the rejection module can be viewed as a stage to investigate reliability or confidence measures for the decisions made by word-spotters. Secondly, other early CM-related works lie in automatic detection of new words (out of the current lexicon) in large vocabulary speech recognition, such as Asadi et al. (1990), Young and Ward (1993) and Young (1994), etc. In addition to modeling out-of-vocabulary (OOV) words with a (or a set of) generic hidden Markov model(s), Young and Ward (1993) proposed to use word score normalization to detect misrecognition and out-of-vocabulary words

for continuous speech recognition. Young (1994) first elucidated how to use the posterior probability as a confidence measure for speech recognition, where she employed the acoustic score normalization based on a separate all-phone recognition to approximate such a posterior probability. Young (1994) also tried to combine the normalized acoustic scores with other high-level knowledge sources, e.g. semantic, pragmatic and discourse analysis, to improve quality of confidence measures. Thirdly, Sukkar and Wilpon (1993), Sukkar (1994) and others realized that confidence measures for ASR become extremely important when speech recognition technology is applied to any practical applications or services for end-users. Some intensive research efforts at the former AT&T Bell Labs resulted in lots of fruitful works related to confidence measure for ASR, but under a different name, namely *utterance verification*. Rose et al. (1995a) first formally cast the confidence measure problem in speech recognition as a statistical hypothesis testing problem in classic statistics and proposed to use likelihood ratio testing (LRT) to solve the problem. Then, the LRT-based formulation has become the basic theoretical foundation for the follow-up works (e.g. Sukkar and Lee, 1996; Rahim et al., 1997; Rahim and Lee, 1997a,b; Sukkar et al., 1997, etc). In these works, they discovered that some discriminative training techniques, such as minimum classification error (MCE) and minimum verification error (MVE) methods, can significantly improve performance of modeling both the *null* and *alternative* hypotheses in utterance verification. Finally, a tremendous amount of research activities have been carried out more recently in this area to seek for some reliable confidence measures for ASR, mainly driven by an increasing number of dialogue applications. Based on confidence measures, spoken language systems will be able to handle error-prone ASR outputs more intelligently in those post-recognition modules, such as language understanding and dialogue management (e.g., see Hazen et al., 2000). Some representative works include Eide et al. (1995), Cox and Rose (1996), Chase (1997), Gillick et al. (1997), Neti et al. (1997), Kemp and Schaaf (1997), Schaaf and Kemp (1997), Weintraub et al. (1997), Rueber (1997), Jiang and Huang

(1998), Willett et al. (1998), Siu and Gish (1999), Benitez et al. (2000), Wessel et al. (2000), Kamppari and Hazen (2000), Kamppari and Hazen (2000), Jiang et al. (2001), Jiang and Lee (2002), Jiang and Lee (2003), San-Segundo et al. (2001), Zhang and Rudnicky (2001) and many others.

Generally speaking, all methods proposed for computing confidence measures (CMs) in speech recognition can be roughly classified into three major categories. Firstly, a large portion of works aim to compute confidence measures based on a combination of the so-called *predictor features*, which are collected during decoding procedure and may include acoustic as well as language information about recognition decisions. Then all predictor features are combined in a certain way to generate a single score to indicate correctness of the recognition decision. We will briefly summarize these methods in Section 2. Secondly, it is well known that the *posterior* probability in the standard maximum a posterior (MAP) decision rule is a good candidate for CM in speech recognition since it is an absolute measure of how well the decision is. However, it is very hard to estimate the posterior probability in a precise manner due to its normalization term in the denominator. In practice, many different approaches have been proposed to approximate it, ranging from simple *filler*-based methods to complex word-graph-based approaches. We will introduce these methods in Section 3. Next, as already mentioned above, under the name of *utterance verification* (UV), lots of works have been conducted to verify the claimed content of a spoken utterance. The content can be hypothesized by a speech recognizer or keyword detector or human transcriber. Under the framework of utterance verification (UV), the CM problem can be formulated as a statistical hypothesis testing. In Section 4, we will briefly present all proposed methods in this category, ranging from the LRT-based non-Bayesian approach (based on Neyman–Pearson Lemma) to the Bayes-factors-based Bayesian approach. We also introduce how to use some discriminative training methods to improve modeling in UV. In the remainder of the paper, in Section 5, I will first mention several samples of very recent research advances regarding CM computation, including

an in-search data selection for improving verification models in UV, and a novel idea to compute confidence measures based on neighborhood information in model space, and how to annotate confidence measures based on semantic information measured by latent semantic analysis (LSA). In Section 6, I will discuss about performance comparison issues among all different CM methods and focus on capabilities and limitations of the current CM techniques in a variety of potential applications. Finally, I will make some general comments on CM methods in ASR and conclude this paper with some clues for future works.

2. CM as combination of predictor features

In the literature, a very large portion of CM-related works aim to search for a predictor feature (or a set of features) which is informative to distinguish correctly recognized results from other possible recognition errors. Any feature can be called a predictor if its probabilistic distribution (e.g., p.d.f.) of correctly recognized words is clearly distinct from that of misrecognized words. Usually, the predictor features may have to be collected within the recognition process at levels of acoustics, language model, syntax, and semantics. Some common predictor features reported in the literature may include:

- *Pure normalized likelihood score* related: acoustic score per frame.
- *N-best* related: count in the N-best list, N-best homogeneity score (the weighted ratio of all paths passing through the hypothesized word in N-best list), top N recognition scores, top N – 1 difference in adjacently ranked recognition scores, etc.
- *Acoustic stability*: a number of alternative hypotheses are generated based on different language model weights in decoding and acoustic stability of any given word is defined as the number of times the word occurs in the list divided by the number of alternatives in the list.
- *Hypothesis density*: the number of alternative arcs spanning the time segment of the recognized word in word graph.
- *Duration* related: HMM state duration, phoneme duration, word duration.
- *Language model (LM)* related: LM score, LM back-off behavior, etc.
- *Parsing* related: whether or not a word is parsed by grammar in robust parsing, position of each parsed word within the semantic slot (either edge or middle), the language model back-off mode of the whole parsed slot, etc.
- *Posterior probability* related: see Section 3 for details.
- *Log-likelihood-ratio* related: see Section 4 for details.

For more predictor features, please refer Cox and Rose (1996), Schaaf and Kemp (1997), Chase (1997), Benitez et al. (2000), San-Segundo et al. (2001), etc. An ideal predictor feature should provide strong information to separate the correctly recognized words from other misrecognitions and the distribution overlap between the two classes should be minor. However, none of the above predictor features is ideal in this sense. As reported in many papers, the overlap is actually quite large even for the best predictor feature. Therefore, some people attempt to combine several different predictor features for a better performance. Many different combinational models have been reported in the literature, including linear discriminant function (Sukkar, 1994; Sukkar and Lee, 1996), generalized linear model (Gillick et al., 1997; Siu and Gish, 1999), single or mixture Gaussian classifier (Chigier, 1992), neural networks (Mathan and Miclet, 1991; Weintraub et al., 1997; San-Segundo et al., 2001), decision tree (Eide et al., 1995; Neti et al., 1997), support vector machine (Zhang and Rudnicky, 2001), boosting (Moreno et al., 2001) and others. In most cases, parameters of combinational models are estimated from some discriminative training procedures based on some criteria such as cross-entropy, classification error rate (see Weintraub et al., 1997 for more details about this).

A combination approach can improve the overall performance only when all individual components are statistically independent. Obviously, this is not the case for the above predictor features. It has been observed in many experiments that all

these predictor features are highly correlated,¹ refer to Kemp and Schaaf (1997), Schaaf and Kemp (1997), Jiang and Huang (1998) and etc. Usually the combination methods cannot significantly improve over the best predictor feature. We believe that an interesting combination strategy lies in combining the best acoustic features, such as posterior probability and N-best related ones, with other pure language features, such as parsing-and/or semantic-related ones, as in Zhang and Rudnicky (2001), Guo et al. (2004) etc. However, so far we have not seen any compelling results by combining various predictor features in confidence measure estimation.

3. CM as posterior probability

It is well known that the conventional ASR algorithm is usually formulated as a pattern classification problem using the *maximum a posterior* (MAP) decision rule to find the most likely sequence of words \hat{W} which achieves the maximum posterior probability $p(W|X)$ given any acoustic observation X , i.e.,

$$\begin{aligned} \hat{W} &= \arg \max_{W \in \Sigma} p(W | X) \\ &= \arg \max_{W \in \Sigma} \frac{p(X | W) \cdot p(W)}{p(X)} \\ &= \arg \max_{W \in \Sigma} p(X | W) \cdot p(W) \end{aligned} \quad (1)$$

where Σ denotes the set of all permissible sentences, $p(W)$ is the probability of W evaluated with a language model, $p(X)$ is the probability of observing X , and $p(X|W)$ is the probability of observing X by assuming that W is the underlying word sequence for X . In theory, the posterior probability $p(W|X)$ is a good confidence measure for the recognition decision that X is recognized as W . However, as shown in Eq. (1), most practical ASR systems simply ignore the term $p(X)$ in decision-making because it is constant across different words W . This explains why the raw ASR scores

are inadequate as confidence measures to judge recognition reliability. However, after being normalized by $p(X)$, the posterior probability $p(W|X)$ can serve as a good confidence measure since it represents the absolute quantitative measure of the match between X and W . In theory, we should compute $p(X)$ as follows:

$$p(X) = \sum_H p(X, H) = \sum_H p(H) \cdot p(X | H) \quad (2)$$

where H denotes any a hypothesis for X , and the above summation must be done over all possible hypotheses for X , including all combinations of words, phonemes, noises and other events. Obviously, without any further constraint, it is impossible to enumerate and model all these hypotheses so that it is extremely difficult to estimate $p(X)$ in a precise manner. In practice, we have to either impose certain assumptions or adopt some approximate methods when estimating $p(X)$ for the posterior probability.

In the first category, it includes the so-called *filler*-based methods which try to calculate $p(X)$ from a set of general filler or background models, i.e., all-phone recognition (Young, 1994), catch-all model (Kamppari and Hazen, 2000), the highest score in recognizing the word from decoder (Cox and Rose, 1996), etc. These approaches are very straightforward and usually can achieve an reasonable performance in many cases. In another category, there are the so-called *lattice*-based methods which attempt to calculate $p(X)$, then the posterior probability $p(W|X)$ in turn, from a word lattice or graph based on the forward–backward algorithm, such as Kemp and Schaaf (1997) and Wessel et al. (1998, 1999, 2000, 2001). Usually, one word lattice or graph is generated by the ASR decoder for every utterance. Then the posterior probability of each recognized word or the entire hypothesized sentence can be calculated based on the word-graph from an additional post-processing stage. Since word graph is a compact and fairly accurate representation of all alternative competing hypotheses of the recognition result which usually dominate the summation when computing $p(X)$ over a variety of hypotheses in Eq. (2), the posterior probability calculated from a word graph can approximate the true $p(W|X)$ pretty well. Therefore, the

¹ Refer to final remarks in Section 7 for a possible explanation for this.

resultant confidence measures generally achieve better performance than all other CMs mentioned in the above. However, generating word graphs and scoring word-graphs for posterior probabilities are relatively complicated and quite demanding in computation, especially in large vocabulary ASR systems. Thus, for the sake of simplicity, an N-best list can also be used in place of word graph for this purpose, such as Rueber (1997), Wessel et al. (1999), etc. Due to its superior performance as a CM for ASR, in the following, I will review some details about how to compute posterior probabilities from a word graph as originally reported by Wessel et al. (2001).

3.1. Word graph notations

Usually the ASR decoder generates a word graph \mathcal{X} for each utterance X . Here, the word graph is represented as a directed, acyclic, weighted graph. All its nodes represent discrete points in time. Each arc is labeled with three variables, i.e. $[w]_s^e$, where w is the hypothesized word attached to the arc, and s and e denote the starting and ending time instances of the arc respectively. Also, each arc is associated with a weight, $\mathcal{B}(w)_s^e$, which is actually acoustic score of generating acoustic feature vectors from time s to e based on the HMM of word w . In every word graph, there are two special nodes: one is called START node which corresponds to the beginning of the utterance and one END node for the end of the utterance. Any path from START node to END node is called a complete path which represents a sentence (a sequence of words) hypothesis for the underlying utterance. Let us assume a complete path in word graph \mathcal{X} of an utterance X , which consists of n different arcs as $C = \{[w_1]_{s_1}^{e_1}, [w_2]_{s_2}^{e_2}, \dots, [w_n]_{s_n}^{e_n}\}$. Obviously, it is straightforward to compute the probability of this complete path given the word graph \mathcal{X} as follows:

$$p(C | \mathcal{X}) = \prod_{i=1}^n \mathcal{B}(w_i)_{s_i}^{e_i} \cdot p(w_i | h_i) \quad (3)$$

where h_i denotes the history of word w_i along the path and $p(w_i|h_i)$ is the language model score computed with n -gram language models.

3.2. Posterior probability of an arc

Based on the above notations, it is easy to compute the posterior probability of any arc $a = [w]_s^e$ given the word graph \mathcal{X} , namely $p(a | \mathcal{X})$.² Normally, $p(a | \mathcal{X})$ is calculated as a ratio between the total probability of all complete paths passing through the arc a to that of all complete paths in \mathcal{X} , i.e.,

$$p(a | \mathcal{X}) = \frac{\sum_{C \in \mathcal{X}, a \supset C} p(C | \mathcal{X})}{\sum_{C \in \mathcal{X}} p(C | \mathcal{X})} \quad (4)$$

where $C \in \mathcal{X}$ denotes C is a complete path in word graph \mathcal{X} and $a \supset C$ denotes that the complete path C passes through the arc a . The posterior probability $p(a | \mathcal{X})$ can be efficiently computed based on a forward–backward algorithm. A forward probability $\alpha_s(a)$ is recursively computed from the starting node of the arc a backward until the START node of the word graph as:

$$\alpha_s(a) = \mathcal{B}(w)_s^e \cdot \sum_{a'} \alpha_{s'}(a') \cdot p(w | h') \quad (5)$$

where the summation is conducted for all arcs a' (s' is the starting time of a' merging into the starting node of a and h' is word history of w in language model computation. Analogously, a backward probability $\beta_e(a)$ is computed from the ending node of a forward until END node of the word graph as:

$$\beta_e(a) = \sum_{a''} \beta_{e''}(a'') \cdot \mathcal{B}(w'')_{s''}^{e''} \cdot p(w'' | h'') \quad (6)$$

where the summation is conducted over all arcs a'' (e'' is the ending time of a'' and w'' is word id in a'') leaving the ending node of a and h'' is word history of w'' when calculating language model score. Obviously, the numerator in Eq. (4) can be computed as the product of $\alpha_s(a)$ and $\beta_e(a)$. And the denominator in Eq. (4) can be recursively computed as forward probability $\alpha_s(a)$ in Eq. (5) starting from START node until END node of the word graph or backward probability $\beta_e(a)$ in Eq. (6) from END node backward to START node of word graph.

² Note that the posterior probability of an arc given the word graph \mathcal{X} differs from the original posterior probability given the utterance X .

3.3. Posterior probability of a recognized word

We can directly use the posterior probability, $p(a | \mathcal{X})$, of the arc, $a = [w]_s^e$, as confidence measure for the recognized word w . But it has been shown that it does not perform very well as a confidence measure for w . We know that except the arc $a = [w]_s^e$, there are usually lots of other arcs in word graph that have the same word id w but slightly different starting time s and ending time e . It will underestimate confidence measure of w if we only count $p(a | \mathcal{X})$ for w . Thus, it is very important to take into account other arcs which have the same word id w but slightly different s and e . Wessel et al. (2001) proposes three different ways to solve this problem. In the first method, called C_{sec} , when calculating confidence measure for the word w in an arc $a = [w]_s^e$, we sum over all arcs in word graph which have the same word id w and intersect with the current arc $a = [w]_s^e$ in time domain. In the second method, called C_{med} , we only accumulate posterior probability for all arcs with the same word id which intersect the median time frame of the arc under consideration. In the third approach, called C_{max} , we determine a best-case probability for word w in an arc $a = [w]_s^e$. We accumulate posterior probability for all arcs (with the same word id) which not only intersect the median time frame but also all other time frames between s and e , and then choose the maximum one from these sums as the confidence measure for the word w in the underlying arc. Based on Wessel et al. (2001), the third method, namely C_{max} , yields the best performance.

There are many other implementation issues to consider when computing posterior probability in word graph, e.g., scaling of probabilities in summation, elimination of redundant silence edges, etc. Readers are referred to Wessel et al. (2001) for more details.

4. CM as utterance verification

Mainly motivated by speaker verification problem, Rose et al. (1995a), Sukkar and Lee (1996), Rahim et al. (1997) have proposed to tackle confidence measure problems from a different pers-

pective. Under the framework of utterance verification (UV), the confidence measure problem in ASR is formulated as a statistical hypothesis testing problem. For a given speech segment X , assume that an ASR system recognizes it as word W which is represented by an HMM λ_W . Utterance verification is a post-processing stage to examine the reliability of the hypothesized recognition result. Under the framework of UV, we first propose two complementary hypotheses, namely the *null* hypothesis H_0 and the alternative hypothesis H_1 as follows:

$$\begin{aligned} H_0 : & \quad X \text{ is correctly recognized and truly} \\ & \quad \text{comes from model } \lambda_W \\ H_1 : & \quad X \text{ is wrongly classified and is} \\ & \quad \text{NOT from model } \lambda_W \end{aligned} \quad (7)$$

Then we test H_0 against H_1 to determine whether we should accept the recognition result or reject it. According to Neyman–Pearson Lemma, under some conditions, the optimal solution to the above testing is based on a likelihood ratio testing (LRT), i.e.,

$$\text{LRT} = \frac{p(X | H_0)}{p(X | H_1)} \underset{H_1}{\overset{H_0}{\geq}} \tau \quad (8)$$

where τ is the critical decision threshold. The LRT-based utterance verification provides a good theoretical formulation to address confidence measure problems in ASR. As pointed out by Lee (2001), the above LRT score can be transformed to a confidence measure based on a monotonic one-to-one mapping function. The major difficulty with LRT is how to model the alternative hypothesis which usually represents a very complex and composite event, where the true distribution of data is unknown. In practice, as in (Rose et al., 1995a; Sukkar and Lee, 1996; Rahim et al., 1997), the same HMM structure is adopted to model the alternative hypothesis, which can be a general background model, or hypothesis-specific *anti-model*, or a set of competing models, or a combination of all the above. In these works, a variety of training methods have been used to estimate HMMs for the alternative hypothesis. It is generally agreed that a discriminative training procedure plays a crucial role in improving modeling

performance for the alternative hypothesis. In (Sukkar and Lee, 1996), a GPD (generalized probabilistic descent) based discriminative training procedure is used to estimate parameters of a linear discriminant function based on a criterion to minimize sub-word level verification error counts represented by a sigmoid function. In (Rahim et al., 1997), it is found that the minimum classification error (MCE) training, which is originally proposed to reduce recognition errors, can contribute to improving performance of UV. In (Rahim and Lee, 1997a; Sukkar et al., 1997), a GPD-based training algorithm is proposed to achieve minimum verification error (MVE) estimation for utterance verification with respect to optimizing verification HMM parameters. In MVE, the string-level verification errors are approximated by using a sigmoid function embedded with a mis-verification function, which actually is negative log-likelihood ratio used in verification. Then the total empirical verification errors can be minimized over all training data by optimizing the verification HMM parameters corresponding to the both null and alternative hypotheses. The optimization can be iteratively achieved by using the GPD algorithm. Experiments clearly show all these discriminative training methods can largely improve performance of the LRT-based utterance verification.

Alternatively, if we consider the above UV problem from a Bayesian viewpoint, the final solution ends up with calculating and evaluating the so-called *Bayes factors* as in (Jiang and Deng, 2001). Bayes factors has its solid foundation from Bayesian theory. Given the speech data X along with the above two hypotheses H_0 and H_1 , Bayes factors is computed as:

$$\begin{aligned} \text{BF} &= \frac{\hat{p}(X | H_0)}{\hat{p}(X | H_1)} \\ &= \frac{\int f(X | \lambda_0, H_0) \cdot p(\lambda_0 | H_0) d\lambda_0}{\int f(X | \lambda_1, H_1) \cdot p(\lambda_1 | H_1) d\lambda_1} \end{aligned} \quad (9)$$

where, for $k = 0, 1$, λ_k is the model parameter under H_k , $p(\lambda_k | H_k)$ is its prior density, and $f(X | \lambda_k, H_k)$ is the likelihood function of λ_k under H_k .

Bayes factors offers a way to evaluate evidence in favor of the *null* hypothesis H_0 because Bayes

factors is the ratio of the posterior odds of H_0 to its prior odds, regardless of the value of the prior odds.³ Therefore, Bayes factors can be used to compare with a threshold, just like the likelihood ratio in Neyman–Pearson lemma, to make a decision with regard to H_0 . In other words, if $\text{BF} > \tau$, where τ is a pre-set critical threshold, then we accept H_0 , otherwise reject it. Like LLR, the BF value can also be transformed or formulated as a confidence measure for ASR.

As shown in (Jiang and Deng, 2001), Bayes factors is a powerful statistical tool to model composite hypotheses and can be used to solve many different verification problems. The same formulation proposed for speaker verification in (Jiang and Deng, 2001) is also equally applicable to the above UV problem though no research work has been reported about this. The key issues are what role the prior distributions $p(\lambda_0 | H_0)$ and $p(\lambda_1 | H_1)$ will play in utterance verification and how to use them as a flexible tool to incorporate a variety of information sources useful for UV.

5. Some recent efforts

Confidence measure or utterance verification aims to verify reliability of speech recognition outputs, which significantly differs from other typical verification problems in statistics, such as test for goodness-of-fit, and outlier detection in statistical data analysis. We believe that it is beneficial not to isolate confidence measure (or utterance verification) from its prior recognition stage. In acoustic level, it is very important to know the distribution properties of competing sources in recognition phase in order to optimize performance of CM or UV. In the following, I will first present two pieces of recent research works along this direction. Besides, I will also briefly summarize some other research works to integrate some high-level knowledge (beyond acoustic information) for CM or UV.

³ Any probability can be converted to the odds scale, i.e., odds = probability/(1 – probability). Thus, $\frac{\text{Pr}(H_0|y)}{\text{Pr}(H_1|y)}$ is called the posterior odds in favor of H_0 , and $\frac{\text{Pr}(H_0)}{\text{Pr}(H_1)}$ is prior odds in favor of H_0 .

5.1. In-search data selection for accurate competing models

Under the UV framework, it is not an easy job to model the alternative hypothesis. Rose et al. (1995a), Sukkar and Lee (1996), Rahim et al. (1997) propose to use the so-called *anti-models* for this purpose. However, it is still unclear what data should be used to estimate these *anti-models*. In their works, some heuristic methods are adopted, such as performing forced-alignment against a wrong or random transcript to generate training data for each anti-model. More recently, Jiang et al. (2001) propose a well-defined in-search data selection procedure to collect the most representative competing tokens for each HMM in the system. Then the selected tokens can be used to estimate highly accurate competing models for the utterance verification purpose.

In (Jiang et al., 2001), we first define competing tokens (CT) of any a given HMM as data segments which are misrecognized to this model during recognition. A dynamic in-search data selection method is proposed to collect competing tokens for every HMM automatically from training data set. In the method, every utterance in training set is recognized with the Viterbi beam search algorithm just as in regular recognition phase. During the Viterbi search, all potential segments located in all active partial paths within the search beam width are compared with the reference segmentation generated from a forced-alignment procedure to determine whether each segment should be a competing token or true token of the model. The procedure is carried out for all training data to collect two token sets, namely the competing token set $\mathcal{S}_C(a)$ and the true token set $\mathcal{S}_T(a)$ for every HMM a in the system. The competing information collected in this way is very valuable for utterance verification. Given that a speech observation X is recognized as W by the decoder, the original hypotheses in Eq. (7) can be re-phrased as follows:

$$\begin{aligned} H_0: X \text{ belongs to } W\text{'s true token set } \mathcal{S}_T(W), \\ \text{i.e., } X \in \mathcal{S}_T(W) \\ H_1: X \text{ belongs to } W\text{'s competing token set } \mathcal{S}_C(W), \\ \text{i.e., } X \in \mathcal{S}_C(W) \end{aligned} \quad (10)$$

Comparing with the original hypotheses, both the null hypothesis H_0 and the alternative hypothesis H_1 in the above are well-defined from available data, which in turn make our modeling problem easier. The simplest way to model them is to estimate two different models A_T and A_C for $\mathcal{S}_T(W)$ and $\mathcal{S}_C(W)$ respectively, based on all tokens collected from training data. Then the LRT-based utterance verification is operated as follows:

$$\begin{aligned} \eta &= \frac{p(X | H_0)}{p(X | H_1)} = \frac{\Pr(X \in \mathcal{S}_T(W))}{\Pr(X \in \mathcal{S}_C(W))} \\ &= \frac{p(X | A_T)}{p(X | A_C)} \underset{H_1}{\overset{H_0}{\geq}} \tau \end{aligned} \quad (11)$$

where τ is the critical decision threshold. The above models A_T and A_C can be estimated based on different criteria, such as maximum likelihood (ML), or minimum verification error (MVE), etc. Jiang et al. (2001) shows the ML-trained models already significantly surpass the conventional UV methods, such as in (Sukkar and Lee, 1996; Sukkar et al., 1997; Rahim et al., 1997; Rahim and Lee, 1997a).

5.2. UV based on neighborhood information in model space

In (Jiang and Lee, 2002, 2003), a novel approach is proposed for utterance verification based on competing information in model space. First of all, let us look at the model space \mathcal{T} of HMM. Each HMM λ in the system can be viewed as a point in the model space \mathcal{T} . Intuitively, we can imagine two nested neighborhoods surrounding the underlying model λ , namely a small neighborhood A_1 and a medium neighborhood A_2 . The small neighborhood A_1 is a tiny neighborhood which tightly surrounds the underlying model λ . As indicated in (Jiang et al., 1999), a neighborhood with a relatively small size contains all variants of the original model due to estimation errors and possible mismatches in testing. It serves as a robust representation of the original model. On the other hand, the medium neighborhood A_2 is significantly larger than A_1 . As the neighborhood size increases, it starts to cover all of its competing models in the model space, which by definition should be close

to the original one in some sense. Based on the concept, we can translate the original hypotheses in Eq. (7) in another way.

Once again, assume a speech observation X is recognized as W which is represented by the model λ_W . We are interested in verifying the reliability of the decision. Given the decision that X is recognized as model λ_W , if X is not from the model λ_W (as stated in the alternative hypothesis), it is reasonable to consider that X probably comes from some competing model of λ_W . Therefore, we can translate the original hypotheses in Eq. (7) as:

$$\begin{aligned} H_0 : & \text{The true model of } X \\ & \text{locates in the small neighborhood } A_1 \\ H_1 : & \text{The true model of } X \\ & \text{locates in the region } A_2 - A_1 \end{aligned} \quad (12)$$

where $A_2 - A_1$ denotes the holed region inside the medium neighborhood by excluding the small neighborhood, as shown in Fig. 1.

In (Jiang and Lee, 2002, 2003), an approach based on Bayes factors is proposed to solve the above hypothesis testing problem.

$$\eta = \frac{\int_{A_1} f(X|\lambda) \cdot p_0(\lambda) d\lambda}{\int_{A_2 - A_1} f(X|\lambda) \cdot p_1(\lambda) d\lambda} \underset{H_1}{\overset{H_0}{\geq}} \tau \quad (13)$$

where $f(X|\lambda)$ is likelihood function, and $p_0(\lambda)$ and $p_1(\lambda)$ represent the prior distribution of model parameters under the hypothesis H_0 and H_1 respectively. Furthermore, Jiang and Lee (2002) propose two simple methods, i.e. a parametric definition and a non-parametric one, to quantitatively define the neighborhoods as well as the prior distributions for the above formulation. Some preliminary experiments show some promising results for utterance verification based on the above framework. Obviously, much more research efforts are needed to define the neighborhoods in a more precise and controllable manner.

5.3. Incorporation of high-level information for CM

So far we have concentrated on confidence measures which solely rely on acoustic information. However, other syntactical or semantic information is also reported to provide certain clues for the purpose of confidence measure, such as Young (1994), Pao et al. (1998), Zhang and Rudnicky (2001), etc. More recently, Cox and Dasmahapatra

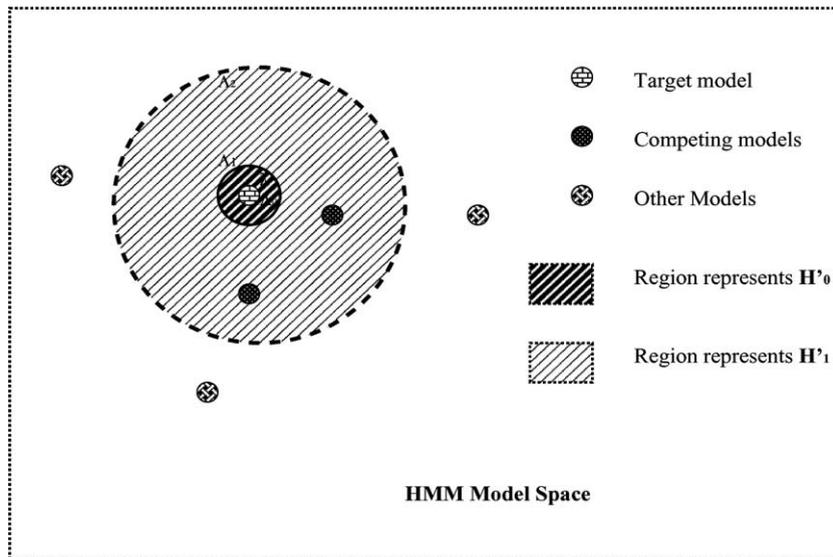


Fig. 1. Illustration of the hypothesis testing based on the neighborhood information in model space.

(2002) report that human can clearly identify a certain portion of recognition errors in recognizer outputs on purely semantic grounds. They also propose to use latent semantic analysis (LSA) to annotate confidence scores for recognized words. Latent semantic analysis (LSA) is a technique for associating words that are “semantically coherent”. The semantic coherence between any two words is computed as the cosine of the angle between the two vectors corresponding to these two words in a reduced subspace. Thus, confidence measure of a recognized word is calculated as an average of coherence of this word with all other recognized words in a close context. Although CMs based on this kind of semantic information is generally not as good as the best CMs in the acoustic level, a combination probably will yield a better performance due to their clear independence. More recently, Guo et al. (2004) conduct some comparative experiments to combine high-level confidence measures, which are based on LSA and/or inter-word mutual information, with the word posterior probability based CMs calculated from word graphs. Some moderate gains have been shown in two large vocabulary speech recognition tasks.

6. Performance and applications of CM: Capabilities and limitations

It is well known that good confidence measures will largely benefit a variety of ASR applications, e.g., to smartly reject non-speech noises, detect/reject out-of-vocabulary words, detect/correct some potential recognition mistakes, clean up human transcription errors in large training corpus, guide the system to perform un-supervised learning, provide side information to assist high-level speech understanding and dialogue management, and so on and so forth. However, confidence measures for ASR is an extremely difficult problem. Even today’s best available CMs are generally not good enough to effectively support most of the above-mentioned applications. In this section, we first briefly talk about the assessment problem of confidence measures in ASR. Then, based on my personal understanding, I will discuss on per-

formance issues of various CMs. At last, I will point out several promising applications for the current CMs in ASR even though many other applications are apparently beyond the capability of today’s techniques.

When evaluating confidence measure annotation, we usually encounter two types of errors, namely false alarm errors and false rejection errors. Obviously, receiver operating characteristic (ROC) curve gives a full picture of verification performance at all operating points. In many cases, it is convenient to use a single-number metric for CM assessment. Some widely used metrics include equal error rate (EER), confidence error rate, normalized cross entropy, etc. Refer to Kemp and Schaaf (1997), Siu and Gish (1999), Maison and Gopinath (2001) and Wessel et al. (2001) for details. Another important issue in CM evaluation is to take recognition boundaries into account. For example, a correctly recognized word may have a very low confidence measure because its boundary is wrong (though its identity is correct). Thus, it is helpful to use the concept of “*word-correctness*” proposed by Weintraub et al. (1997) in evaluating CMs.

As far as CMs performance issues are concerned, it has been widely reported that N-best related feature predictors perform much better than other predictors introduced in Section 2 (see Chase, 1997; Rueber, 1997; Williams and Renals, 1999, etc). Moreover, Wessel et al. (1998, 1999, 2001) clearly shows that posterior probabilities calculated from word graphs significantly outperform N-best-related confidence measures. On the other hand, along a totally different line, Sukkar and Lee (1996, 1997) and Rahim et al. (1997,a) demonstrate that MVE-based discriminative training significantly improve performance of utterance verification. Furthermore, Jiang et al. (2001) shows the performance of utterance verification is largely improved over the previous UV approaches by using an in-search data selection method to train some highly accurate competing models. In (Garcia et al., 1999), the conventional LRT-based utterance verification is compared with posterior probabilities in word graph albeit their implementation is an approximate one. The results show word-graph-based posterior

probabilities outperform the LRT-based utterance verification methods. However, it still remains unclear how the approach in (Jiang et al., 2001) compares with word-graph-based posterior probabilities in (Wessel et al., 2001). Moreover, it will be more informative if all CMs are evaluated in a common corpus for several well-designed verification tasks. Generally speaking, the CMs based on posterior probabilities derived from word graphs are advantageous since language model scores can be naturally incorporated in CM computation in addition to acoustic information. But once being strictly implemented as in (Wessel et al., 2001), the performance of CMs cannot be easily improved within the same paradigm because word graphs with various sizes usually generate CMs with similar performance. On the other hand, performance of utterance verification can be progressively improved by estimating better and better verification models. And the hypothesis testing paradigm, as formulated in LRT- or Bayes-factors-based testing, provides a flexible framework to incorporate a variety of knowledge sources which may be useful for CM computation.

As already mentioned above, the overall performance of CMs (even the best ones) remains fairly poor, which largely limits their applications. Some early research on CM aimed to detect out-of-vocabulary (OOV) words in large-vocabulary ASR system. However, even by today, an effective detection of OOV words in continuous speech recognition remains as an open question. It seems only feasible to use CMs to reject OOV words in some constrained small vocabulary applications, such as isolated voice command controlling, etc. Besides, a large amount of works have been conducted to improve ASR performance with assistance of various CMs (e.g., Neti et al., 1997; Jitsuhiro et al., 1998; Vergyri, 2000; Wessel et al., 2000; Tan et al., 2000; Lleida and Rose, 2000; Koo et al., 2001 and others). A consistent and significant error reduction over the state-of-the-art performance is still not an easy goal to achieve⁴ unless the performance of CMs is enhanced further. Moreover,

CMs have been included in many spoken dialogue systems to provide certain level of support for language understanding and dialogue management. But the CMs themselves are found not robust and reliable enough to be a solid basis for decision-making in many cases. In spite of these, the current CM techniques still have a chance to shine if they are applied to a proper place. Although it is hard to detect or correct errors made by ASR systems by using CMs, it seems much easier to use CMs to detect human-made errors. Thus, it is promising to use CMs to clean-up or verify transcription in a large corpus, such as some preliminary studies in (Arslan and Hansen, 1999; Li et al., 2002). In addition, there are two other successful stories to apply CMs to verify some decisions not hypothesized by ASR systems. One is verbal information verification (VIV) in (Li et al., 2000) where the authors verify each input utterance against pre-stored information, such as birthday, address, phone number, etc, based on its verbal content by using the LLR-based utterance verification technique. Another one is Liu et al. (2001) and Afify et al. (in press) where an LLR-based CM is computed for a look-ahead speech segment at each time instant to discard some unlikely phonemes from further consideration in full search stage. They demonstrate effectiveness of using the LLR-based CMs in such a fast-match stage for search space pruning prior to recognition stage. Moreover, another interesting area to apply CMs is un-supervised adaptation where CMs are used to select more reliable speech segments from recognizer's outputs to self-improve recognition models (e.g. Wallhoff et al., 2000; Goronzy, 2002 and many others). One important issue here is that the operating point in verification stage should be set up to guarantee a low false acceptance rate.

7. Final remarks

Although there are various types of CMs reported for ASR, almost all CMs in acoustic level fundamentally rely almost entirely on a single information source, namely how much the underlying decision can overtake other possible competitors. The larger the difference is the more confident we will believe the decision to be. This

⁴ If counting the correct recognition results which are mistakenly rejected.

explains why most research works to combine a variety of CMs usually do not yield better results. The various CM or UV methods mentioned in this paper attempt to explore this discrepancy in different ways (direct or indirect). For example, in the posterior probability method based on a word graph, if the recognition result significantly surpasses other competing choices in the word graph, the contribution of the recognized path will dominate the total posterior probability computed based on the forward–backward algorithm. In this case, the derived CM will be large (close to 1). If other competing paths in the word graph come very close to the recognized results, the contribution of the recognized path will be relatively small when computing the posterior probability. Thus, the derived CM will be small (close to 0). Similarly in UV, if the recognized result largely surpasses other competitors, the likelihood under the null hypothesis will be significantly larger than that of the alternative hypothesis. As a result, the likelihood ratio will be large. On the other hand, the likelihood ratio will be small if the competing sources from the alternative hypothesis gives comparable results with the recognized one in the null hypothesis. This also explains why it is very important to model distribution properties of competing hypotheses when deriving CMs for ASR. Apparently, it is a real challenge to compute any effective CMs beyond this sole source. Besides, one major drawback of almost all CM or UV methods is that we only verify segment identities but never question the correctness of segmentation hypothesized by ASR systems. It is common that most recognition errors accompany with segmentation mistakes in continuous speech recognition. A preliminary study on boundary adjustment for UV can be found in (Matsui et al., 2001). Three heuristic methods to calculate word posterior probability CM in word graph by Wessel et al. (2001), i.e., C_{sec} , C_{med} and C_{max} as reviewed in Section 3.3, also tackle the problem in an ad hoc way. We believe it is critical to improve performance of CMs by taking this segmentation issue into account. How to consider it effectively in any formal way in CM estimation still remains unclear. Finally, despite a large number of research activities in the past, confidence measure estimation for ASR still

remains unsolved in so many aspects. Due to its importance in practice and its difficulty in theory, we expect much more research efforts will be devoted into this topic in coming years.

References

- Afify, M., Liu, F., Jiang, H., Siohan, O., 2004. A new verification-based fast-match for large vocabulary continuous speech recognition. *IEEE Trans. Speech Audio Processing*, in press.
- Arslan, L., Hansen, J., 1999. Selective training for hidden Markov model with applications to speech classification. *IEEE Trans. Speech Audio Process.* 7 (1), 46–54.
- Asadi, A., Schwartz, R., Makhoul, J., 1990. Automatic detection of new words in a large vocabulary continuous speech recognition system. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 125–128.
- Benitez, M.C., Rubio, A., Torre, A., 2000. Different confidence measures for word verification in speech recognition. *Speech Commun.* 32, 79–94.
- Chase, L., 1997. Word and acoustic confidence annotation for large vocabulary speech recognition. *Proc. of European Conference on Speech Communication Technology*, pp. 815–818.
- Chigier, B., 1992. Rejection and keyword spotting algorithms for a directory assistance city name recognition application. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. II-93–II-96.
- Cox, S., Rose, R., 1996. Confidence measures for the Switchboard database. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 511–514.
- Cox, S., Dasmahapatra, S., 2002. High-level approaches to confidence estimation in speech recognition. *IEEE Trans. Speech Audio Process.* 10 (7), 460–471.
- Eide, E., Gish, H., Jeanrenaud, P., Mielke, A., 1995. Understanding and improving speech recognition performance through the use of diagnostic tools. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 221–224.
- Garcia, M.C., Reichl, W., Ortmanns, S., 1999. On combining confidence measures in HMM-based speech recognizers. *Proc. of Automatic Speech Recognition and Understanding Workshop*.
- Gillick, L., Ito, Y., Young, J., 1997. A probabilistic approach to confidence estimation and evaluation. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 879–882.
- Gong, Y.-F., 1995. Speech recognition in noisy environments: A survey. *Speech Commun.* 16, 261–291.
- Goronzy, S., 2002. Confidence measures. In: *Robust Adaptation to Non-native Accents in Automatic Speech Recognition*. Springer-Verlag, Berlin, pp. 57–78, Chapter 7.

- Guo, G., Huang, C., Jiang, H., Wang, R.-H., 2004. A comparative study on various confidence measures in large vocabulary speech recognition. Proc. of the fourth International Symposium on Chinese Spoken Language Processing, Hong Kong, China.
- Hazen, T.J., Burianek, T., Polifroni, J., Seneff, S., 2000. Recognition confidence scoring for use in speech understanding systems. Proc. of ISCA ITRW Workshop on ASR, Paris, France.
- Jiang, H., Deng, L., 2001. A Bayesian approach to the verification problem: Applications to speaker verification. *IEEE Trans. Speech Audio Process.* 9 (8), 874–884.
- Jiang, L., Huang, X.-D., 1998. Vocabulary-independent word confidence measure using subword features. Proc. of International Conference on Spoken Language Processing.
- Jiang, H., Lee, C.-H., 2002. Utterance verification based on neighborhood information and Bayes Factors. Proc. of International Conference on Spoken Language Processing.
- Jiang, H., Lee, C.-H., 2003. A new approach to utterance verification based on neighborhood information in model space. *IEEE Trans. Speech Audio Process.* 11 (5), 425–434.
- Jiang, H., Hirose, K., Huo, Q., 1999. Robust speech recognition based on Bayesian prediction approach. *IEEE Trans. Speech Audio Process.* 7 (4), 426–440.
- Jiang, H., Soong, F.K., Lee, C.-H., 2001. A data selection strategy for utterance verification in continuous speech recognition. Proc. of European Conference on Speech Communication and Technology, pp. 2573–2576.
- Jitsuhiro, T., Takahashi, S., Aikawa, K., 1998. Rejection of out-of-vocabulary words using phoneme confidence likelihood. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 217–220.
- Juang, B.-H., 1991. Speech recognition in adverse environments. *Computer Speech Language* 5, 275–294.
- Kamppari, S.O., Hazen, T.J., 2000. Word and phone level acoustic confidence scoring. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 1799–1802.
- Kemp, T., Schaaf, T., 1997. Estimating confidence using word lattices. Proc. of European Conference on Speech Communication Technology, pp. 827–830.
- Koo, M.-W., Lee, C.-H., Juang, B.-H., 2001. Speech recognition and utterance verification based on a generalized confidence score. *IEEE Trans. Speech Audio Process.* 9 (8), 821–832.
- Lee, C.-H., 1998b. On stochastic feature and model compensation approaches to robust speech recognition. *Speech Commun.* 25, 29–47.
- Lee, C.-H., 2001. Statistical confidence measures and their applications. Proc. of ICSP, August.
- Li, Q., Juang, B.-H., Zhou, Q., Lee, C.-H., 2000. Automatic verbal information verification for user authentication. *IEEE Trans. Speech Audio Process.* 8 (5), 585–596.
- Li, Q., Jiang, H., Zhou, Q., Zheng, J., 2002. Automatic enrollment for speaker authentication. Proc. of International Conference on Spoken Language Processing.
- Liu, F., Afify, M., Jiang, H., Siohan, O., 2001. A new verification-based fast match approach to large vocabulary speech recognition. Proc. of European Conference on Speech Communication and Technology, pp. 851–854.
- Lleida, E., Rose, R.C., 2000. Utterance verification in continuous speech recognition: Decoding and training procedures. *IEEE Trans. Speech Audio Process.* 6, 558–568.
- Maison, B., Gopinath, R., 2001. Robust confidence annotation and rejection for continuous speech recognition. Proc. of International Conference on Acoustics, Speech and Signal Processing.
- Matsui, T., Soong, F.K., Juang, B.-H., 2001. Verification of multi-class recognition decision using classification approach. Proc. of Automatic Speech Recognition and Understanding Workshop.
- Mathan, L., Miclet, L., 1991. Rejection of extraneous input in speech recognition applications, using multi-layer perceptrons and the trace of HMMs. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 93–96.
- Moreno, P.J., Logan, B., Raj, B., 2001. A boosting approach for confidence scoring. Proc. of European Conference on Speech Communication Technology.
- Neti, C.V., Roukos, S., Eide, E., 1997. Word-based confidence measures as a guide for stack search in speech recognition. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 883–886.
- Rahim, M.G., Lee, C.-H., 1997a. String-based minimum verification error (SB-MVE) training for speech recognition. *Computer Speech Language* 11, 147–160.
- Rahim, M.G., Lee, C.-H., 1997b. A study on robust utterance verification for connected digits recognition. *J. Acoust. Soc. Am.* 101 (5), 2892–2902.
- Rahim, M.G., Lee, C.-H., Juang, B.-H., 1997. Discriminative utterance verification for connected digits recognition. *IEEE Trans. on Speech and Audio Processing* 5 (3), 266–277.
- Pao, C., Schmid, P., Glass, J., 1998. Confidence scoring for speech understanding systems. Proc. of International Conference on Spoken Language Processing.
- Rose, R.C., 1992. Discriminant wordspotting techniques for rejecting non-vocabulary utterances in unconstrained speech. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. II-105–II-108.
- Rose, R.C., Juang, B.H., Lee, C.H., 1995a. A training procedure for verifying string hypothesis in continuous speech recognition. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 281–284.
- Rueber, B., 1997. Obtaining confidence measures from sentence probabilities. Proc. of European Conference on Speech Communication Technology.
- San-Segundo, R., Pellom, B., Hacioglu, K., Ward, W., 2001. Confidence measures for spoken dialogue systems. Proc. of International Conference on Acoustics, Speech and Signal Processing.
- Schaaf, T., Kemp, T., 1997. Confidence measures for spontaneous speech recognition. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. 875–878.

- Siu, M., Gish, H., 1999. Evaluation of word confidence for speech recognition systems. *Computer Speech Language* 13, 299–319.
- Sukkar, R.A., 1994. Rejection for connected digit recognition based on GPD segmental discrimination. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. I-393–I-396.
- Sukkar, R.A., Lee, C.-H., 1996. Vocabulary independent discriminative utterance verification for nonkeyword rejection in subword based speech recognition. *IEEE Trans. Speech Audio Process.* 4 (6), 420–429.
- Sukkar, R.A., Wilpon J.G., 1993. A two pass classification for utterance rejection in keyword spotting. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. II-451–II-454.
- Sukkar, R.A., Setlur, A.R., Lee, C.-H., Jacob, J., 1997. Verifying and correcting recognition string hypotheses using discriminative utterance verification. *Speech Commun.* 22, 333–342.
- Tan, B.T., Gu, Y., Thomas, T., 2000. Utterance verification based speech recognition system. *Proc. of International Conference on Spoken Language Processing*.
- Vergyri, D., 2000. Use of word level side information to improve speech recognition. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 1823–1826.
- Wallhoff, F., Willett, D., Rigoll, G., 2000. Frame-discriminative and confidence-driven adaptation for LVCSR. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 1835–1838.
- Weintraub, M., Beaufays, F., Rivlin, Z., Konig, Y., Stolcke, A., 1997. Neural-network based measures of confidence for word recognition. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 887–890.
- Wessel, F., Macherey, K., Schluter R., 1998. Using word probabilities as confidence measures. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 225–228.
- Wessel, F., Macherey, K., Ney, H., 1999. A comparison of word graph and N-best list based confidence measures. *Proc. of European Conference on Speech Communication Technology*, pp. 315–318.
- Wessel, F., Schluter, R., Ney, H., 2000. Using posterior word probabilities for improved speech recognition. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 1587–1590.
- Wessel, F., Schluter, R., Macherey, K., Ney, H., 2001. Confidence measures for large vocabulary continuous speech recognition. *IEEE Trans. Speech Audio Process.* 9 (3), 288–298.
- Willett, D., Worm, A., Neukirchen, C., Rigoll, G., 1998. Confidence measures for HMM-based speech recognition. *Proc. of International Conference on Spoken Language Processing*.
- Williams, G., Renals, S., 1999. Confidence measures from local posterior probability estimates. *Computer Speech Language* 13, 395–411.
- Wilpon, J.G., Rabiner, L.R., Lee, C.-H., Goldman, R., 1990. Automatic recognition of keywords in unconstrained speech using hidden Markov models. *IEEE Trans. Acoustics Speech Signal Process.* 38 (11), 1870–1878.
- Young, S.R., Ward, W., 1993. Learning new words from spontaneous speech. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. II-590–II-591.
- Young, S.R., 1994. Detecting misrecognitions and out-of-vocabulary words. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. II-21–II-24.
- Zhang, R., Rudnicky, A.I., 2001. Word level confidence annotation using combinations of features. *Proc. of European Conference on Speech Communication Technology*.

Further reading

- Bourlard, H., D'hoore, B., Boite, J.-M., 1994. Optimizing recognition and rejection performance in word-spotting systems. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. I-373–I-376.
- Charlet, D., Mercier, G., Juvet, D., 2001. On combining confidence measures for improved rejection of incorrect data. *Proc. of European Conference on Speech Communication Technology*.
- Dolfing, J.G.A., Wendemuth, A., 1998. Combination of confidence measures in isolated word recognition. *Proc. of International Conference on Spoken Language Processing*.
- Gunawardana, A., Hon, H.-W., Jiang, L., 1998. Word-based acoustic confidence measures for large-vocabulary speech recognition. *Proc. of International Conference on Spoken Language Processing*, pp. 791–794.
- Gupta, S.K., Soong, F.K., 1998. Improved utterance rejection using length dependent thresholds. *Proc. of International Conference on Spoken Language Processing*, pp. 795–798.
- Hacioglu, K., Ward, W., 2002. A concept graph based confidence measure. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. I-225–I-228.
- Hernandez, A.G., Marino, J.B., 2000. Contextual confidence measures for continuous speech recognition. *Proc. of International Conference on Acoustics, Speech and Signal Processing*, pp. 1803–1806.
- Lee, C.-H., 1998. A tutorial on speaker and utterance verification. *Proc. of NORSIG*, pp. 9–16.
- Lin, Q., Das, S., Lubensky, D., Picheny, M., 1998. A new confidence measure based on rank-ordering subphone scores. *Proc. of International Conference on Spoken Language Processing*.
- Lin, Q., Lubensky, D., Roukos, S., 1999. Use of recursive mumble models for confidence measuring. *Proc. of European Conference on Speech Communication Technology*.
- Mengusoglu, E., Ris, C., 2001. Use of acoustic prior information for confidence measure in ASR applications. *Proc. of European Conference on Speech Communication Technology*.

- Modi, P., Rahim, M., 1997. Discriminative utterance verification using multiple confidence measures. Proc. of European Conference on Speech Communication Technology, pp. 103–106.
- Moreau, N., Juvet, D., 1999. Use of a confidence measure based on frame level likelihood ratios for the rejection of incorrect data. Proc. of European Conference on Speech Communication Technology.
- Palmer, D.D., Ostendorf, M., 2001. Improved word confidence estimation using long range features. Proc. of European Conference on Speech Communication Technology.
- Rohlicek, J.R., Jeanrenaud, P., Ng, K., Gish, H., Musicus, B., Siu, M., 1993. Phonetic training and language modeling for word spotting. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. II459–II462.
- Rose, R.C., 1995b. Keyword detection in conversational speech utterances using hidden Markov model based continuous speech recognition. *Computer Speech Language* 9, 309–333.
- Rose, R.C., Riccardi, G., 1999. Automatic speech recognition using acoustic confidence conditioned language models. Proc. of European Conference on Speech Communication Technology.
- Rose, R.C., Yao, H., Riccardi, G., Wright, J., 2001. Integration of utterance verification with statistical language modeling and spoken language understanding. *Speech Commun.* 34, 321–331.
- Sankar, A., Wu, S.-L., 2003. Utterance verification based on statistics of phone-level confidence scores. Proc. of International Conference on Acoustics, Speech and Signal Processing, pp. I-584–I-587.
- Tan, B.T., Gu, Y., Thomas, T., 2001. Word level confidence measures using N-best sub-hypotheses likelihood ratio. Proc. of European Conference on Speech Communication Technology.
- Uhrik, C., Ward, W., 1997. Confidence metrics based on N-gram language model backoff behaviors. Proc. of European Conference on Speech Communication Technology.