

International Journal of Speech Technology

Predicting Search Term Reliability For Spoken Term Detection Systems

--Manuscript Draft--

Manuscript Number:	IJST-D-13-00001R1
Full Title:	Predicting Search Term Reliability For Spoken Term Detection Systems
Article Type:	Manuscript
Keywords:	spoken term detection, voice keyword search, information retrieval
Corresponding Author:	Joseph Picone, Ph.D. Temple University Elkins Park, Pennsylvania UNITED STATES
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Temple University
Corresponding Author's Secondary Institution:	
First Author:	Joseph Picone, Ph.D.
First Author Secondary Information:	
Order of Authors:	Joseph Picone, Ph.D. Amir Hossein Harati Nejad Torbati, MS
Order of Authors Secondary Information:	
Abstract:	Spoken term detection is an extension of text-based searching that allows users to type keywords and search audio files containing recordings of spoken language. Performance is dependent on many external factors such as the acoustic channel, language, pronunciation variations and acoustic confusability of the search term. Unlike text-based searches, the likelihoods of false alarms and misses for specific search terms, which we refer to as reliability, play a significant role in the overall perception of the usability of the system. In this paper, we present a system that predicts the reliability of a search term based on its inherent confusability. Our approach integrates predictors of the reliability that are based on both acoustic and phonetic features. These predictors are trained using an analysis of recognition errors produced from a state of the art spoken term detection system operating on the Fisher Corpus. This work represents the first large-scale attempt to predict the success of a keyword search term from only its spelling. We explore the complex relationship between phonetic and acoustic properties of search terms. We show that a 76% correlation between the predicted error rate and the actual measured error rate can be achieved, and that the remaining confusability is due to other acoustic modeling issues that cannot be derived from a search term's spelling.

General Comments:

This paper is the first comprehensive analysis of search term performance on a problem of scale. Keyword search is an increasingly important problem. The genesis for this work was an evaluation conducted by NIST in 2006. The technology analyzed in this paper has been in operational use, particularly in the intelligence community, for some time, on extremely large amounts of operational data. The paper is an attempt to bring some basic science to problems that have been consistently observed by operational systems. The authors can assure you that the observations in this paper are informed by extensive operational use. Unfortunately, this information is classified, and no such comparable resources exist in the public domain on the unclassified side. This work is intended to begin a movement in this direction, and to influence an upcoming keyword search evaluation.

The reviewers' comments about the lack of basic science in the paper must be tempered by the observation that the human language technology field typically focuses on applications of technology. Algorithms are most often borrowed from other fields, such as statistics, and rarely invented with the field. Most journal papers consist of adaptations of these algorithms to specific problems of interest to the community. This paper falls in the latter category. We make no claims to the development of novel algorithms. It is study of the correlates between search terms and keyword search performance. It is the first such study on a critical mass of data. Keyword search technology, similar to speech recognition technology, needs vast amounts of data if one is to attempt to draw conclusions about the trends in performance for a wide variety of words.

The importance of this work is in the exploration of what features correlate with keyword search term performance. There are a vast number of linguistic features that could be considered. We have, in fact, considered over 150 combinations of various features. It would make sense to explore some of the specific combinations in greater detail. For example, reviewer no. 1 commented about the lack of discussion about BPC. However, in reality, most of these linguistic features carry about the same amount of information and individually are not great predictors of performance. Hence, we pursued an aggregate approach.

A significant contribution is to reduce this vast space to something more manageable. For example, it was well known before this work that duration was a primary correlate. However, we have done a good job of exploring this relationship in more detail. The goals of this work are very clearly stated in the introduction to the paper. This topic has not been previously explored in the literature, though error analysis has been performed on a much more limited and anecdotal basis previously. This work will form the basis for a more extended evaluation being conducted this year. We hope as data from these new studies become available we can extend the models and gain new insight into the problem.

Below we address the specific concerns of the reviewers:

Reviewer #1:

- 1. Authors have not given any insights about the explored features and models.**
- 2. The paper is not technically strong due to lack of analysis of results**
- 3. The paper is not scientific, in the sense that there is no logic or justification in using various features and models.**
- 6. In the present draft, there is no technical contribution. Authors have used existing features and models blindly, and got some numbers. There is no interpretation regarding**

the selection of features and models. There is no analysis and interpretation about the obtained results.

These four concerns are essentially restatements of the same issue.

The revised paper contains slightly more discussion of the explored features and models. It must be pointed out that the entire section on acoustic analysis resulted from an analysis of the poor performance obtained from the phonetic distance approach. We originally approached the work only from a phonetic basis. Unfortunately, those features had limited value. That triggered a far deeper analysis of these factors, which in turn gave rise to a major piece of this paper. So the claim that there was no detailed analysis provided seems to be due to a lack of understanding of the relationship between these two approaches and the motivations for the acoustic approaches.

The underlying, or root causes, of errors is discussed in the context of the degree to which acoustic matching is influenced by phonetic constraints. We have added some additional comments to better underscore the analysis presented in this paper. However, it must be emphasized that the key contribution of the paper is an exploration of which features are most highly correlated. That in itself represents an analysis of search term performance. The underlying causes of why particular features are good or bad is discussed broadly in terms of some acoustic and phonetic issues, but it is difficult to draw such conclusions for particular features in a statistically meaningful manner. Anecdotal evidence only goes so far in this type of analysis.

Some analysis is provided throughout the paper whenever possible. For example, on page 10 we discuss the role the parameter “count” plays in the results. We also discuss duration on page 5. We discuss and analyze the construction of the data sets in Table 2 and 3.

We have added more detailed explanations of the experiments throughout Section IV. Hopefully this will not put us over the page limits.

4.Details about existing literature are missing.

5.There is no comparison with the state of the art in this area.

As mentioned previously, this is the first comprehensive study of this type of problem that we are aware of. That is what makes the paper unique. We have accurately cited all the previous work related to the specific systems and data presented in this study. In fact, the authors worked closely with the providers of these data and systems. This paper was not meant to be an extensive or exhaustive tutorial on keyword search.

Further, all prediction algorithms used are referenced, even though the details of these algorithms is outside the scope of this work.

If the reviewer feels we left out particular references, we would appreciate more specific guidance, as long as these are relevant to the study presented.

7. In Table-1 sounds are categorized into various groups based on their manner of articulation. Later there is no discussion on, how these groups effect the recognition performance.

Space limitations, of course, impact the extent to which every experiment can be discussed. There are over 1,000 experiments using over 150 features that were conducted to bring the study to this

point, and it is difficult to give justice to each one of these. The entries labeled SFS in the various tables represent combinations of various features. If necessary, we can go into greater detail of the specific subset of features selected. We decided not to do this since the relative differences between these combinations are fairly small. We focused on the explicit features that had the greatest impact on the problem.

Reviewer #2:

First, this paper is really screaming for some sort of user evaluation. I'm not talking about a controlled user study, I am talking about some preliminary qualitative work as to whether speech-naïve users actually improve their search performance using the tool.

There is a web interface, but it is mentioned only in passing. I would like to see more discussion of this. If the authors can provide, my evaluation of the paper will go up considerably.

User evaluations are difficult to conduct. We have provided a web-based demo and have observed some naïve and expert users use of this demo. However, the demo is connected to a relatively limited database of indexed audio provided by Microsoft. There is no doubt that the feedback provided by our prediction algorithm is useful. However, stating this in a scientifically meaningful manner is difficult because user experiences with keyword search are complex. Despite significant advances in the technology in the past decade, user perceptions of the technology, particularly in the intelligence community, vary greatly.

At our university lab, we do not have the resources to conduct such studies. It should be noted that the 2006 evaluation, upon which this work was based, also did not conduct user studies. Few HLT papers present such results in a meaningful manner, and traditionally focus more on objective measures of performance.

One of the authors has extensive experience with the use of this technology in intelligence applications. As mentioned in the paper, dealing with users' frustrations with keyword search formed the basis for this work. Unfortunately, this information is classified and cannot be discussed in the open literature. However, let me emphasize that keyword search technology has been extensively evaluated by the authors in very real operational conditions. The author has worked directly with analysts on strategies for selecting keywords, and those experiences have informed the research in this paper.

For example, it was long conjectured that search terms with long durations were better. Unfortunately, many search terms used in operational settings tend to be short. It has been a complicated process to train users on the shortcomings of these approaches. There are strong interactions between search terms, the domain of interest, and the acoustic confusability of these terms. Anecdotal evidence of a few users searching the limited audio archives available on the Internet (such as the site our web-based demo links to) have only limited scientific value.

Again, our lab does not have the resources to conduct such evaluations on the unclassified side, and that really goes beyond the scope of this paper. This paper is not meant to be a comprehensive evaluation of search term technology. It is meant to explore a small piece of this overall puzzle – what phonetic and acoustic features correlate with search term performance.

Second, the approach considers only the term being searched, not the corpus or the task. Thus it somewhat confounds standard text-based search techniques, which users are probably already using. For example, based on Figure 2, searching for "shopping" may be a *bad* idea, if the term is found frequently, as it will give too many results and these will have to be searched tediously. So there is also some tension here. Perhaps the work would be best thought of as part of a total search strategy, that might also consider frequency of terms in the corpus as well as similar sounds terms in the corpus which might be candidate hits (see, e.g. the work of Vemuri).

We acknowledge that data is always a problem in these kinds of studies. When we began this study, we attempted to acquire evaluation data from all the major systems presented at the 2006 NIST evaluation. We were able to capture most of it. However, the database used in that study turned out to be too small and suffered from some of the problems the reviewer alluded to. Therefore, we worked with several of the system providers to get more data, and that is why we used Fisher. Even Fisher is moderate for this kind of analysis. If a search term appears less than 50 times in a corpus, it is very hard to draw conclusions about its behavior. On the other hand, analysts want to be able to search for terms that are not commonly occurring words.

This is why, for example, we explored underlying linguistic features such as N-grams of phonemes. Unfortunately, even N-grams of phonemes occur somewhat sparsely, even in Fisher. This motivated us to look at features such as N-grams of BPC. Such features occur in multiple words and hence we get more stable estimates of their overall performance. This is a detail that might not be obvious to the average reader.

We believe we have been very careful to emphasize the differences between text-based searches and audio-based searches. We agree that for a given false alarm rate, if a term occurs frequently, the number of actual hits returned will be greater than for a less frequently occurring term. Unfortunately, for the less frequently occurring term, the recognizer will also have fewer examples to train on. Hence, there is a complex relationship between acoustic decoding and what amounts to language model probabilities.

Intelligence analysts, or people who use keyword search technology frequently in their day-to-day jobs, do understand these tradeoffs at an intuitive level. They tend not to use generic terms that occur frequently. Since they are searching vast amounts of data, they must use terms that are generally occurring infrequently and of high value. However, within this range of terms, there is significant latitude in which specific terms they can pick. A good example might be choosing the difference between a word like "vote" and a word like "democracy" or "government". The technology provided in this paper gives them some insight into the difference between such terms.

Summary Comments:

All in all I think the paper is solid and a small contribution. I do question the overall utility of the approach in a real search task and would like to see some evidence that it does have value there. But it is also a short paper and I would be happy to see it published if other reviewers are positive.

We would like to thank the reviewers for their thoughtful and insightful feedback. We have done our best to address each of their concerns. We hope our responses have adequately addressed your concerns and we can proceed with publication.

We want to emphasize that this is the first study of its type to be published in a major journal as far as we know. The timing is important since there are significant evaluations of keyword search planned for the coming years. It makes sense to establish this as an area of science that can support and motivate additional analyses. As more data becomes available, the scientific quality of these studies can be refined.

PREDICTING SEARCH TERM RELIABILITY FOR SPOKEN TERM DETECTION SYSTEMS

Amir Hossein Harati Nejad Torbati and Joseph Picone

Abstract— Spoken term detection is an extension of text-based searching that allows users to type keywords and search audio files containing recordings of spoken language. Performance is dependent on many external factors such as the acoustic channel, language, pronunciation variations and acoustic confusability of the search term. Unlike text-based searches, the likelihoods of false alarms and misses for specific search terms, which we refer to as reliability, play a significant role in the overall perception of the usability of the system. In this paper, we present a system that predicts the reliability of a search term based on its inherent confusability. Our approach integrates predictors of the reliability that are based on both acoustic and phonetic features. These predictors are trained using an analysis of recognition errors produced from a state of the art spoken term detection system operating on the Fisher Corpus. This work represents the first large-scale attempt to predict the success of a keyword search term from only its spelling. We explore the complex relationship between phonetic and acoustic properties of search terms. We show that a 76% correlation between the predicted error rate and the actual measured error rate can be achieved, and that the remaining confusability is due to other acoustic modeling issues that cannot be derived from a search term's spelling.

Keywords— spoken term detection, voice keyword search, information retrieval

Manuscript submitted December 30, 2012. Manuscript revised and resubmitted on May 8, 2013.

A. Harati and J. Picone are with the Department of Electrical and Computer Engineering at Temple University, 1947 North 12th Street, Philadelphia, Pennsylvania 19027 USA (phone: 215-204-4841; fax: 215-204-5960; email: joseph.picone@isip.piconepress.com).

I. INTRODUCTION

The goal of a Spoken Term Detection (STD) system is “to rapidly detect the presence of a word or phrase in a large audio corpus of heterogeneous speech material” (Fiscus et al., 2007). As shown in Figure 1, STD systems typically index the audio data as a preprocessing step, allowing users to rapidly search the index files using common information retrieval approaches. Indexing can be done using a speech to text (STT) system (Miller et al., 2007), or simpler engines based on phoneme recognition (Nexidia, 2008). Like most detection tasks, STD can be characterized in terms of two kinds of errors: false alarms and missed detections (Martin et al., 1997). The overall error can be defined as a linear combination of these two errors. In this paper, we give equal weights to both types of errors.

Search engines have been used extensively to retrieve information from text files. Regular expressions (Duford, 1993) and statistically-based information retrieval algorithms (Manning et al., 2008) have been the foundations of such searches for many years. Text-based search algorithms use simple character recognition and character matching algorithms in which the identity of a character is known with probability 1 (no ambiguity). Unlike searching text data, searching through audio data requires handling ambiguity at the acoustic level. Determining the presence of a particular phone or word is not an exact science and must be observed through probabilities. A similarity measure used in such searches is typically based on some kind of score computed from a machine learning system. For text-based search systems, the performance of the system is independent of the term being searched (at least for a language like English where words are explicitly separated using spaces). For audio-based searches, however, the performance of the system depends on many external factors including the acoustic channel, speech rate, accent, language, vocabulary size and the inherent confusability of the search terms. Here we address only the latter problem – predicting the reliability of a search term based on its inherent confusability.

The motivation for this work grew out of observations of typical users interacting with both word-based (Miller et al., 2007) and phone-based (Nexidia, 2008) voice keyword search systems over the past seven years. While it is well known that some aspects of search term performance, such as the duration of the word, correlate with search term performance (Doddington et al., 1999; Harati & Picone,

1
2
3
4 2013), selecting robust and accurate search terms can be as much art as science. Users can quickly
5
6 become frustrated because the nuances of the underlying speech processing engine don't always align
7
8 with users' expectations based on their experiences with text-based searches. Therefore, our goal in this
9
10 work was to develop a technology similar to password strength checking which displays the predicted
11
12 strength of a keyword as a user types a search term.
13

14
15 A demonstration of the system is available at
16
17 http://www.isip.piconepress.com/projects/ks_prediction/demo/current/. A screenshot of the user interface
18
19 is shown in Figure 2. The output of the tool is a visual feedback to the user in the form of a numeric score
20
21 in the range [0,100%] that indicates the quality of the search term (e.g., 100% means the search term is
22
23 strong and less likely to result in inaccurate hits). If a search term is likely to cause inaccurate results, that
24
25 results in users having to sift through many utterances to find content of interest. The tool is an attempt to
26
27 provide users with an interactive indication of the quality of a proposed term before they execute the
28
29 search. Our experience with users is that, without this type of feedback, they often gravitate towards short
30
31 search terms that are highly confusable. The tool makes it very easy for users to understand the value of
32
33 selecting alternate search terms. Though not currently included in this tool, an obvious extension is to
34
35 provide users with a list of alternate terms that are semantically similar yet have better reliability. Though
36
37 we have not conducted extensive user evaluations with this tool, anecdotal results suggest that the
38
39 feedback is very useful to casual users, and that users quickly understand the importance of selecting
40
41 good search terms.
42
43
44
45

46
47 Our general approach in this work was to analyze error patterns produced by existing keyword search
48
49 systems and to develop a predictive model of these errors. To build predictors of errors, we investigated
50
51 both the acoustic phonetic distance between words and similarity measures of the underlying phone
52
53 sequences. The use of acoustic measures resulted from a detailed analysis of the limited predictive power
54
55 of phonetic or linguistic information. Our hypothesis for the acoustic phonetic approach was that
56
57 acoustically similar words should have the same average error rate for a given speech recognizer. The
58
59 similarity measure-based approach calculates an edit distance between the underlying phone sequences
60
61
62
63
64
65

1
2
3
4 (Picone et al., 1990). These two approaches provided simple but useful baseline performance. A third
5
6 approach, which is a major focus of this work, is based on extracting a variety of features from the
7
8 spelling of a word and uses machine learning algorithms to estimate the error rate for that word.
9

10 A block diagram of our general approach is demonstrated in **Error! Reference source not found.**
11
12 The input, a keyword search term that can consist of a word or phrase, is first transformed into features.
13
14 These features result from the conversion of a word into several linguistic representations (e.g., phones,
15
16 syllables). The preprocessor forms an augmented feature vector from an analysis of these linguistic
17
18 representations (e.g., N-grams of phones or broad phonetic class). The machine learning block estimates
19
20 one or more reliability scores, and passes these to the postprocessor for aggregation and normalization.
21
22 For the machine learning task, we have implemented several statistical models based linear regression
23
24 (Bishop, 2011), feed-forward neural networks (Bishop, 2011) and random forests (Breiman, 2001). The
25
26 feature extraction process is central to this work since we have investigated what underlying linguistic
27
28 properties of a word are the strongest predictors of search error rates. Since different approaches predict
29
30 the error rate in different ways, we also explored combining predictors using a simple linear averaging
31
32 that employs particle swarm optimization (PSO) to find the optimal weights (Kennedy & Eberhart, 1995).
33
34
35
36
37

38 The problem of predicting search term reliability is a relatively new problem and for the first time is
39
40 addressed comprehensively in this paper. Researchers have often performed error analysis on speech
41
42 recognition or keyword search experiments, but these have often been focused on system optimization
43
44 and have been very specific to the data under consideration. The goal of the approaches explored in this
45
46 paper is to develop a predictive tool that generalizes across corpora and can be used for vast audio
47
48 archives found in YouTube and through search engines such as Google ad Bing. Hence, it is important
49
50 that the methodology mix both linguistic and acoustic knowledge. In this paper, we present an extensive
51
52 analysis of the predictive power of various types of features derived from this type of information.
53
54
55
56
57
58
59
60
61
62
63
64
65

II. FEATURE GENERATION

In this section we explore several approaches to generating features that can be used to measure the similarity between words. Our goal is to determine feature combinations that have the highest correlation with measured error rates. Since this type of analysis is relatively new, there are no widely accepted set of baseline features for this problem. Our approach in this paper is to hypothesize a wide range of linguistic and acoustic features, and then to employ feature selection methods, discussed in Section III, to select the most relevant ones.

A. Linguistically-derived Features

Our original approach, motivated by the need to develop application-independent metrics, was based on a phonetic distance measure. Each token was converted into a phonetic representation using a dictionary or letter to sound rules (Elovitz et al., 1976). An edit distance (Wagner and Fischer, 1974) was computed using a standard dynamic programming approach. This approach was an attempt to model the underlying phonetic similarity between words, particularly compound words or words that shared morphemic representations.

Next we introduced a family of algorithms based on features extracted from the linguistic properties of words. These features included duration, length (number of letters), number of syllables, number of syllables/length, number of consonants/length, number of vowels/length, a ratio of the number of vowels to the number of consonants, number of occurrences in the language model (count), monophone frequency, broad phonetic class (BPC) frequency, consonant-vowel-consonant (CVC) frequency, biphone frequency, 2-grams of the BPC and CVC frequencies, and 3-grams of the CVC frequencies. We have used a simple phoneme-based duration model (Harati and Picone, 2013) to estimate the duration. The total number of linguistic features is 150, which includes a variety of N-grams of the above features.

The correlation between duration and the average error rate is shown in Figure 4. The average error rate decreases as the duration increases. This correlates with our general experiences with users of these systems. On the surface, it would appear that the more syllables contained in a search term, the lesser its

1
2
3
4 likelihood of being confused. However, as we will see shortly, the variance of this predictor is too high to
5
6 be useful in practical applications, due to some issues related to acoustic training in speech recognition.
7

8
9 The number of syllables was determined using a dictionary or syllabification software (Fisher, 1997)
10
11 for terms not in the dictionary. Mapping phones to consonant and vowel classes was easily accomplished
12
13 using a table lookup. The frequency of occurrence of a word, which we refer to as count, was measured
14
15 on the Fisher Corpus. A summary of the BPC classes used in our study is shown in Table 1. The
16
17 frequency measures used with these features consisted of the fraction of times each symbol appears in a
18
19 word.
20

21 22 *B. Acoustic-Based Features*

23
24 Based on our observation that linguistically-derived units had limited predictive power (to be explored
25
26 more fully in Section IV), we hypothesized that words with similar acoustic properties will result in
27
28 similar error rates. One possibility to exploit this behavior is to cluster words with similar acoustic
29
30 properties and average their associated error rates. We explored two ways to do this based on their
31
32 acoustic and phonetic properties. For an acoustic-based distance algorithm, the criterion used was a
33
34 Euclidian distance in the acoustic space. The acoustic space is constructed from features vectors based on
35
36 a concatenation of standard MFCC features (with derivatives and acceleration components) and duration
37
38 (Young et al., 2006; Davis & Mermelstein, 1980),
39
40

41
42 The acoustic data was, of course, extracted from a different, non-overlapping corpus:
43
44 SWITCHBOARD (SWB) (Godfrey et al., 1992). A list of words was extracted from our target database,
45
46 the Fisher Corpus (Cieri et al., 2004). All instances of these words were located in SWB using the
47
48 provided time alignments (Deshmukh et al., 1998). Durations of the corresponding tokens were
49
50 normalized using a variation of an averaging approach developed by Karsmakers et al. (2007). Feature
51
52 vectors were constructed using three different approaches.
53
54

55
56 In the first approach, each token was divided into three sections by taking its total duration in frames
57
58 and splitting that duration into three sections with durations arranged in 3-4-3 proportions (e.g., a token of
59
60 20 frames was split into three sections of lengths 6, 8 and 6 frames respectively). The average of the
61
62
63
64
65

1
2
3
4 corresponding feature vectors in each segment was computed, and the three resulting feature vectors were
5
6 concatenated into one composite vector. The final feature vector was obtained by adding the duration of
7
8 the token to the three 39-dimensional MFCC feature vectors, bringing the total dimension of the feature
9
10 vector to $3*39+1=118$.

11
12 We then created an alternate segmentation following the procedure described above that was based on
13
14 a 10-24-32-24-10 proportion. This resulted in a feature vector of dimension $5*39+1=196$ elements. In our
15
16 third approach, we divided the utterance into 10 equal-sized segments, which resulted in a feature vector
17
18 of dimension $39*10+1=391$.

19
20 Since there are so many word tokens, we used a combination of *K*-MEANS clustering and *k*-nearest
21
22 neighbor classification (kNN) to produce an estimate of a test token's error rate. All feature vectors for a
23
24 given word were clustered into *K* representative feature vectors, or cluster centroids, using *K*-MEANS
25
26 clustering. We then used kNN classification to locate the *k* nearest clusters for a test token. The overall
27
28 error rate for a word was computed as the weighted average of the *k* clusters, with the weighting based on
29
30 an acoustic distance:
31
32

$$33 \quad err(w_i) = A \sum_{j \in D_k} \frac{1}{dist_{\text{Euclidean}}(w_i, w_j) + \varepsilon} err(w_j), \quad (1)$$

$$34 \quad A = \sum_{j \in D_k} \frac{1}{dist_{\text{Euclidean}}(w_i, w_j) + \varepsilon}, \quad (2)$$

35
36 where w_i is the word in question, D_k is the set of *k* nearest neighbors, and ε is a small positive constant
37
38 that guarantees the denominator will be non-zero.
39
40

41 42 43 44 45 46 47 **III. MACHINE LEARNING**

48
49 We evaluated three types of machine learning algorithms to map features to error rates. These algorithms
50
51 were chosen because they are representative of the types of learning algorithms available, provide a good
52
53 estimate of what type of performance is achievable, and also give us insight into the underlying
54
55 dependencies between features. Some have historical significance (e.g., linear regression) as a baseline
56
57
58
59
60
61
62
63
64
65

1
2
3
4 algorithm while others are known to provide state of the art performance (e.g., random forests). The
5
6 models used in this paper can be regarded as a baseline for future research on this topic.

7
8
9 Linear regression (LR) (Bishop, 2011) is among the simplest methods that can be used to explore
10 dependencies amongst features. We assume that the predictive variable (e.g. error rate) can be expressed
11 as linear combination of the features:
12

$$15 \quad y = X\beta + \varepsilon, \quad (3)$$

$$17 \quad \hat{\beta} = (X'X)^{-1} X'y. \quad (4)$$

18
19
20 where X represents the input feature vector for a word, y represents the predicted error rate, ε is the
21 prediction error and β represents the weights to be learned from the training data.
22
23

24
25
26 Feed-forward neural networks (NN) (Bishop, 2011) are among the most efficient ways to model a
27 nonlinear relationship and have demonstrated robust performance across a wide range of tasks. As before,
28 we assume a simple predictive relationship between X and y :
29
30

$$31 \quad y = f(X) + \varepsilon. \quad (5)$$

32
33 In our implementation, $f()$, the function to be estimated, is approximated as a weighed sum of sigmoid
34 functions. We have used a network with one hidden layer. The output node is chosen to be linear.
35
36 Training was implemented the back-propagation algorithm.
37
38

39
40
41 A random forest (RF) (Breiman, 2001) gives performance that is competitive with the best algorithms
42 and yet does not require significant parameter tuning. The merits of the RF approach include speed,
43 scalability and, most importantly, robustness to overfitting. A common approach for implementing a
44 random forest is to grow many regression trees, each referred to as a base learner, using a probabilistic
45 scheme. The training process for each base learner seeks the best predictor feature at each node from
46 among a random subset of all features. A random subset of the training data is used that is constructed by
47 sampling with replacement so that the size of the dataset is held constant. This randomization helps
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 ensure the independence of the base learners. Each tree is grown to the largest extent possible without any
5
6 pruning.

7
8
9 RFs can also be used for feature selection using a bagging process that is implemented as follows. For
10 one-third of trees in the forest, we generate the training subset using a special scheme: for the k^{th} tree we
11 first put aside one-third of the data from the bootstrap process (sampling with replacement), and label this
12 data out-of-bag (OOB) data. We apply the OOB data to each tree and compute the mean square error
13 (MSE). Next, we randomly permute the value of a specific feature, rerun the OOB data, and compute the
14
15 difference between old and new MSE. The value of this difference, averaged across all trees, shows the
16
17 degree of sensitivity to this feature, and can be interpreted as the importance of that variable.
18
19
20
21
22
23

24 **IV. BASELINE EXPERIMENTS**

25
26 The data used in this project was provided by BBN Technologies (BBN) and consisted of recognition
27 output for the Fisher 2300-hour training set (Cieri et al., 2004). The speech recognizer was trained on
28 370 hours of SWB. The decoder used was configured to run 10 times faster than real time and was similar
29 to a decoder used for keyword search (Miller et al., 2007). Recognition output consisted of word lattices,
30 which we used to generate 1-best hypotheses and average duration information.
31
32
33
34
35
36

37
38 Though it is preferable to have disjoint training and evaluation sets, because the data available is
39 limited, we used a cross-validation approach. We divided the data into 10 subsets and at each step use one
40 of these subsets as the evaluation set and other 9 subsets as training data. At each step we trained models
41 from a flat-start state using the corresponding training data. After rotating through all 10 subsets, we
42 concatenated the results to obtain the overall estimate of performance. Statistics on both the training and
43 evaluation sets are reported in terms of MSE, correlation and R values.
44
45
46
47
48
49
50

51 We have used two feature selection algorithms to explore which features are most important:
52 sequential feature selection (the function `sequentialfs` in MATLAB) (Aha & Bankert, 1996) and random
53 forests (the function `TreeBagger` in MATLAB) (Breiman, 2001). We began with a set of 150 features. We
54 generated 7 subsets of these features as shown in Table 2. Set 1 was generated using sequential feature
55 selection (SFS) and linear regression with correlation as the criterion function. Set 2 was similar to set 1
56
57
58
59
60
61
62
63
64
65

1
2
3
4 except it used MSE as the criterion. Sets 3 and 4 used sequential feature selection with a neural network,
5
6 with correlation and MSE as criteria. Sets 5 and 6 used a regression tree (built using the MATLAB
7
8 function `RegressionTree.template`), with correlation and MSE as criteria respectively. Set 7 used the RF
9
10 approach previously described. We see in Table 2 that approximately 50 features seems to be optimal but
11
12 as few as 7 features gives reasonable performance. SFS selected features such as duration, length and
13
14 count as the most relevant, particularly for the case of 7 features. It also appears the training data is large
15
16 enough to support these kinds of investigations as the results are well-behaved as a function of the
17
18 number of features.
19
20

21
22 A plot of feature importance as determined by the RF algorithm is shown in Figure 5. Count, which
23
24 represents the frequency of occurrence of a word, is recognized as the most important feature (its removal
25
26 causes the highest increase in error.) Note that this does not mean that count is the most relevant feature in
27
28 predicting the error rate. It simply means that other features are highly correlated with each other, so
29
30 removing any one of these does not appreciably reduce the information content in the feature vector.
31
32 Figure 5 demonstrates that no individual feature stands out as having a large predictive power. For
33
34 example, N-grams of phonemes individually occur so infrequently that it is very hard for any one N-gram
35
36 to influence the error rate. On the other hand, duration, length and other such aggregate features are
37
38 correlated to each other and hence in combination don't provide a significant amount of new information.
39
40 Therefore, we must explore more sophisticated combinations of these features.
41
42
43

44
45 In Table 3, we present the correlation of the predicted error rates for the acoustic-based features using
46
47 the *K*-MEANS/kNN approach previously described. Performance is optimal for $K=2$ and $k=\text{inf}$, which
48
49 simply means the feature vectors were clustered into 2 clusters, and every element of each cluster was
50
51 used in the kNN computation. However, overall performance is not extremely sensitive to the parameter
52
53 settings, and the correlation of performance between the training and evaluation sets is good.
54
55

56
57 In Table 4, we show similar results as a function of the number of nearest neighbors for the
58
59 phonetic-based distance metric. Though the MSEs are comparable for both methods, the R values are
60
61 higher for the acoustic-based metric, indicating a better prediction of the error rates. This seems to
62
63
64
65

1
2
3
4 indicate that acoustic modeling in speech recognition plays a more dominant role than the linguistic
5 structure of a search term. Optimal performance is obtained with $k=30$, which is on the order of the
6 number of phonemes in our phoneme inventory, indicating that an excessive number of degrees of
7 freedom are not needed in these feature sets.
8
9

10
11
12
13 In Table 5, we compare three different classification algorithms as a function of the feature sets. The
14 acoustic-based metric resulted in an R value of 0.6 on the evaluation set, while the phonetic-based
15 methods resulted in an R value of 0.5, and the feature-based methods resulted in an R of 0.7. The RF and
16 NN classification methods resulted in similar R values. Approximately 80% of the R value in these cases
17 was due to duration. The remaining features accounted for a very small increase in the R value. There is
18 no strong preference for features such as BPC and CVC since they were roughly comparable in their
19 contribution to the overall R value.
20
21
22
23
24
25
26
27

28
29 The result of this section shows that some of the features like duration, count, bigram frequencies and
30 acoustic distance have a relatively good correlation with the expected word error rate. A combination of
31 these features can explain about 50% of the variance in the prediction results. Our intuition indicates that
32 duration reduces the acoustic ambiguity while bigram frequencies reflect both the occurrence of the word
33 in the training database and the acoustic confusability of certain phoneme sequences.
34
35
36
37
38
39

40 V. SYSTEM COMBINATION

41
42 In order to investigate whether we can build a better predictor by combining different machines, we
43 examined the correlation between predictors. As shown in Table 6, the acoustic-based distance is least
44 correlated with the phonetic-based approach, indicating there could be a benefit to combining these
45 predictors. We have explored combining systems using a weighted average of systems, where optimum
46 weights are learned using particle swarm optimization (PSO) (Kennedy and Eberhart, 1995). The training
47 process for PSO followed the same procedure described previously: the data, in this case word error rates
48 for individual words, is divided into 10 equal subsets. One subset is used for evaluation, the remaining 9
49 subsets are used for training, and the process is repeated by selecting each of the 10 subsets as the
50 evaluation set. The 9 subsets are used to train 75 different classifiers representing a variety of systems
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 selected across the three approaches (acoustic, phonetic and feature-based). PSO is applied to the
5
6 predicted error rates produced by these 75 models on the held-out training data (referred to as
7
8 development data). The result of this process is a vector representing the optimum weight of each
9
10 machine. This process is repeated for each of the 10 partitions. The 10 vectors that result are then
11
12 averaged together to produce the overall optimum weights. These weights are used to combine all 75
13
14 machines into a single model. The error rate predictions of this model are then evaluated against the
15
16 reference error rates measured from the speech recognition output.
17
18

19
20 In this work we have a linearly constrained problem in which we want to find optimum weights for
21
22 our classifiers under the constraint that these weights sum to one. We have used Paquet and
23
24 Engelbrecht (2003) for this constrained optimization problem. In Table 7, we show the results obtained by
25
26 combining all 75 machines using PSO. These 75 machines are composed of 27 machines that use the
27
28 acoustic-based approach, 8 machines using the phonetic-based approach and 40 machines using the
29
30 feature-based approach. We also investigated removing the 8 linear regression machines, reducing the
31
32 number of systems from 75 to 67. This is shown in the second row of Table 7. The last three columns
33
34 show the percent that each machine contributes to the overall score.
35
36

37
38 Acoustic-based and feature-based machines contribute equally to the overall score, and both
39
40 contribute significantly more than the phonetic-based approaches. In fact, when all 75 machines are
41
42 pooled, 43 of these machines (57%) have weights that are zero, implying they add no additional
43
44 information. The 43 machines included 12 from the acoustic-based machines (out of 27), 6 from the
45
46 phonetic-based machines (out of 8), and 25 from the feature-based machines (out of 40). By manually
47
48 excluding the 8 linear regression machines performance increases slightly. Prior to using PSO, our best
49
50 performance was an R value of 0.708. Our best R value with PSO and system combination was 0.761,
51
52 which is an improvement of 7.5%. Figure 6 shows the predicted error rate versus the reference error rate
53
54 for the system representing the second row of Table 7, demonstrating that there is good correlation
55
56 between the two.
57
58
59
60
61
62
63
64
65

VI. SUMMARY

We have demonstrated an approach to predicting the quality of a search term in a spoken term detection system that is based on modeling the underlying acoustic phonetic structure of the word. Several similarity measures were explored (acoustic, phonetic and feature-based), as were several machine learning algorithms (regression, neural networks and random forests). The acoustic-based and feature-based representations gave relatively good performance, achieving a maximum R value of 0.7. By combining these systems using a weighted averaging process based on particle swarm optimization, the R value was increased to 0.761.

To further improve these results, we need to find better features. One of the more promising approaches to feature generation involves an algorithm that predicts the underlying phonetic confusability of a word based on inherent phone-to-phone confusions (Picone et al., 1990). We also, of course, need more data, particularly data from a variety of keyword search engines. It is hoped that such data will become available with the upcoming Spoken Term Detection evaluation to be conducted by NIST in 2013.

VII. ACKNOWLEDGMENTS

The authors would like to thank Owen Kimball and his colleagues at BBN for providing the data necessary to perform this study. This research was supported in part by the National Science Foundation through Major Research Instrumentation Grant No. CNS-09-58854.

VIII. REFERENCES

- Aha, D. W., & Bankert, R. L. (1996). A comparative evaluation of sequential feature selection algorithms. In D. Fisher & H.-J. Lenz (Eds.), *Learning from Data: Artificial Intelligence and Statistics V* (1st ed., pp. 199–206). New York City, New York, USA: Springer.
- Bishop, C. (2011). *Pattern Recognition and Machine Learning* (2nd ed., p. 738). New York, New York, USA: Springer.
- Breiman, L. (2001). Random Forest. *Machine Learning*, 45(1), 5–32.
- Cieri, C., Miller, D., & Walker, K. (2004). The Fisher Corpus: a Resource for the Next Generations of Speech-to-Text. *Proceedings of the International Conference on Language Resources and Evaluation* (pp. 69–71). Lisbon, Portugal.

- 1
2
3
4 Davis, S., & Mermelstein, P. (1980). Comparison of Parametric Representations for Monosyllabic Word
5 Recognition in Continuously Spoken Sentences. *IEEE Transactions on Acoustics, Speech, and*
6 *Signal Processing*, 28(4), 357–366.
7
- 8 Deshmukh, N., Ganapathiraju, A., Gleeson, A., Hamaker, J., & Picone, J. (1998). Resegmentation of
9 Switchboard. *Proceedings of the International Conference on Spoken Language Processing* (pp.
10 1543–1546). Sydney, Australia.
11
- 12 Doddington, G., Ganapathiraju, A., Picone, J., & Wu, Y. (1999). Adding Word Duration Information to
13 Bigram Language Models. Presented at the IEEE Automatic Speech Recognition and Understanding
14 Workshop. Keystone, Colorado, USA.
15
- 16 Duford, D. (1993). *crep: a regular expression-matching textual corpus tool* (p. 84). Technical Report
17 No. CUCS-005-93. Department of Computer Science, Columbia University, New York, New York,
18 USA. doi: <http://hdl.handle.net/10022/AC:P:12304>.
19
- 20 Elovitz, H., Johnson, R., McHugh, A., & Shore, J. (1976). *Automatic Translation of English Text to*
21 *Phonetics by Means of Letter-to-Sound Rules* (NRL Report No. 7948) (p. 102). Washington, D.C.,
22 USA. doi: <http://www.dtic.mil/dtic/tr/fulltext/u2/a021929.pdf>.
23
- 24 Fiscus, J., Ajot, J., Garofolo, J., & Doddington, G. (2007). Results of the 2006 Spoken Term Detection
25 Evaluation. *Proceedings of the SIGIR 2007 Workshop: Searching Spontaneous Conversational*
26 *Speech* (pp. 45–50). Amsterdam, Netherlands.
27
- 28 Fisher, W. (1997). Tsylb syllabification package. url: <ftp://jaguar.ncsl.nist.gov/pub//tsylb2-1.1.tar.Z>. Last
29 accessed on December 24, 2012.
30
- 31 Godfrey, J., Holliman, E., & McDaniel, J. (1992). SWITCHBOARD: Telephone speech corpus for
32 research and development. *Proceedings of the IEEE International Conference on Acoustics, Speech*
33 *and Signal Processing* (pp. 517–520). San Francisco, California, USA.
34
- 35 Harati, A., & Picone, J. (2013). Assessing Search Term Strength in Spoken Term Detection. To be
36 presented at the IEEE International Multi-Disciplinary Conference on Cognitive Methods in
37 Situation Awareness and Decision Support. San Diego, California, USA.
38
- 39 Karsmakers, P., Pelckmans, K., Suykens, J., & Van hamme, H. (2007). Fixed-Size Kernel Logistic
40 Regression for Phoneme Classification. *Proceedings of INTERSPEECH* (pp. 78–81). Antwerp,
41 Belgium.
42
- 43 Kennedy, J., & Eberhart, R. (1995). Particle swarm optimization. *Proceedings of the IEEE International*
44 *Conference on Neural Networks* (pp. 1942–1948). Washington, D.C., USA.
45
- 46 Manning, C., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval* (p. 496).
47 Cambridge, UK: Cambridge University Press.
48
- 49 Martin, A., Doddington, G., Kamm, T., Ordowski, M., & Przybocki, M. (1997). The DET curve in
50 assessment of detection task performance. *Proceedings of Eurospeech* (pp. 1895–1898). Rhodes,
51 Greece.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 Miller, D., Kleber, M., Kao, C.-L., Kimball, O., Colthurst, T., Lowe, S., & Schwartz, R. (2007). Rapid
5 and Accurate Spoken Term Detection. *Proceedings of INTERSPEECH* (pp. 314–317). Antwerp,
6 Belgium.
7
8
9 Nexidia, Inc. (2008). *Phonetic Search Technology* (p. 17). Atlanta, Georgia, USA. Retrieved from
10 [http://www.nexidia.com/government/files/Static Page Files/White Paper Phonetic Search](http://www.nexidia.com/government/files/Static_Page_Files/White_Paper_Phonetic_Search_Tech%2Epdf)
11 [Tech%2Epdf](http://www.nexidia.com/government/files/Static_Page_Files/White_Paper_Phonetic_Search_Tech%2Epdf).
12
13 Paquet, U., & Engelbrecht, A. P. (2003). A new particle swarm optimiser for linearly constrained
14 optimisation. *Proceedings of the IEEE Congress on Evolutionary Computation* (pp. 227–233).
15 Canberra, Australia.
16
17 Picone, J., Doddington, G., & Pallett, D. (1990). Phone-mediated word alignment for speech recognition
18 evaluation. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 38(3), 559–562.
19
20
21 Wagner, R., & Fischer, M. J. (1974). The String-to-String correction problem. *Journal of the ACM*, 21(1),
22 168–173.
23
24 Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., Moore, G., et al. (2006). *The HTK*
25 *Book* (p. 384). Cambridge, U.K. (v3.4.1, url: <http://htk.eng.cam.ac.uk/docs/docs.shtml>).
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

IX. LIST OF FIGURES

Figure 1. Spoken term detection can be partitioned into two tasks: indexing and search. One common approach to indexing is to use a speech to text system (after Fiscus et al., 2007).

Figure 2. A prototype of a web-based application that predicts voice keyword search term reliability is shown. The search term reliability is automatically updated as the user types a search term. A demonstration is available at http://www.isip.piconepress.com/projects/ks_prediction/demo/current/.

Figure 3. In our approach to predicting search term reliability, we decompose terms into features, such as N-grams of phonemes and the number of phonemes, and apply these features to a variety of machine-learning algorithms.

Figure 4. The relationship between duration and error rate shows that longer words generally result in better performance, but the overall variance of this measure is high.

Figure 5. Feature importance based on the RF algorithm is shown. The feature "count," which represents the frequency of occurrence of a word, is by far the singlemost valuable feature since it is not correlated with any of the other features.

Figure 6. The predicted error rate is plotted against the reference error rate, demonstrating good correlation between the two.

X. FIGURES

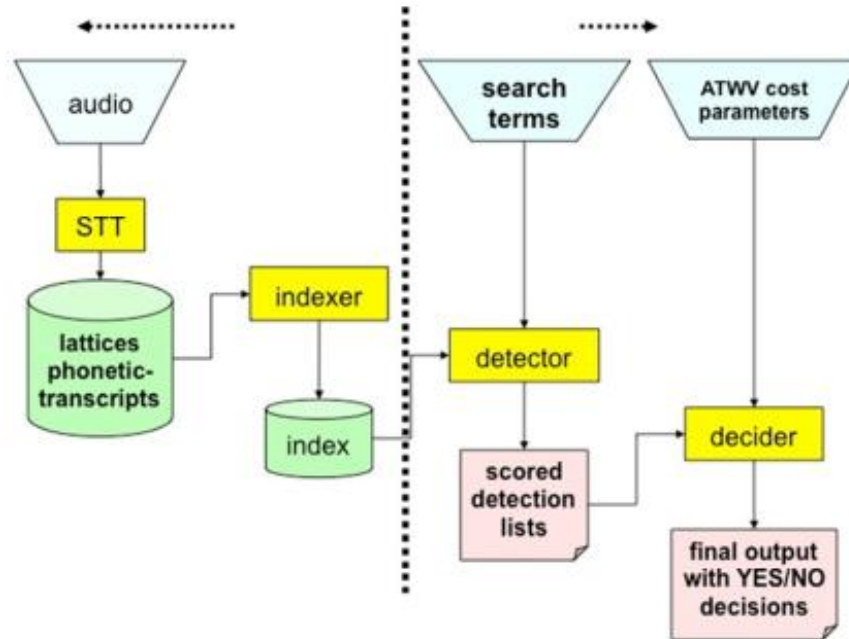


Figure 1. Spoken term detection can be partitioned into two tasks: indexing and search. One common approach to indexing is to use a speech to text system (after Fiscus et al., 2007).

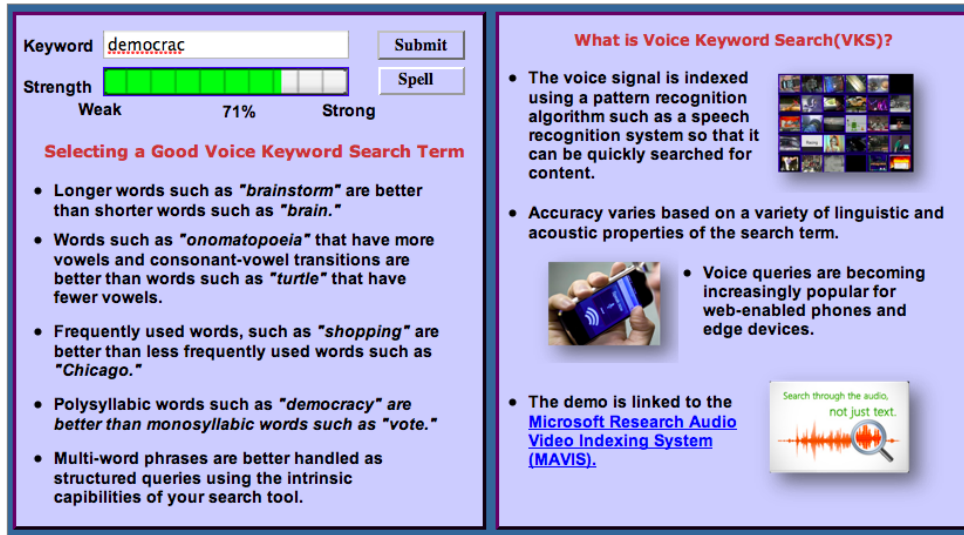


Figure 2. A prototype of a web-based application that predicts voice keyword search term reliability is shown. The search term reliability is automatically updated as the user types a search term. A demonstration is available at http://www.isip.piconepress.com/projects/ks_prediction/demo/current/.

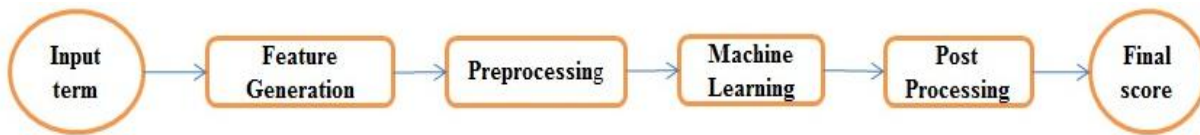


Figure 3. In our approach to predicting search term reliability, we decompose terms into features, such as N-grams of phonemes and the number of phonemes, and apply these features to a variety of machine-learning algorithms.

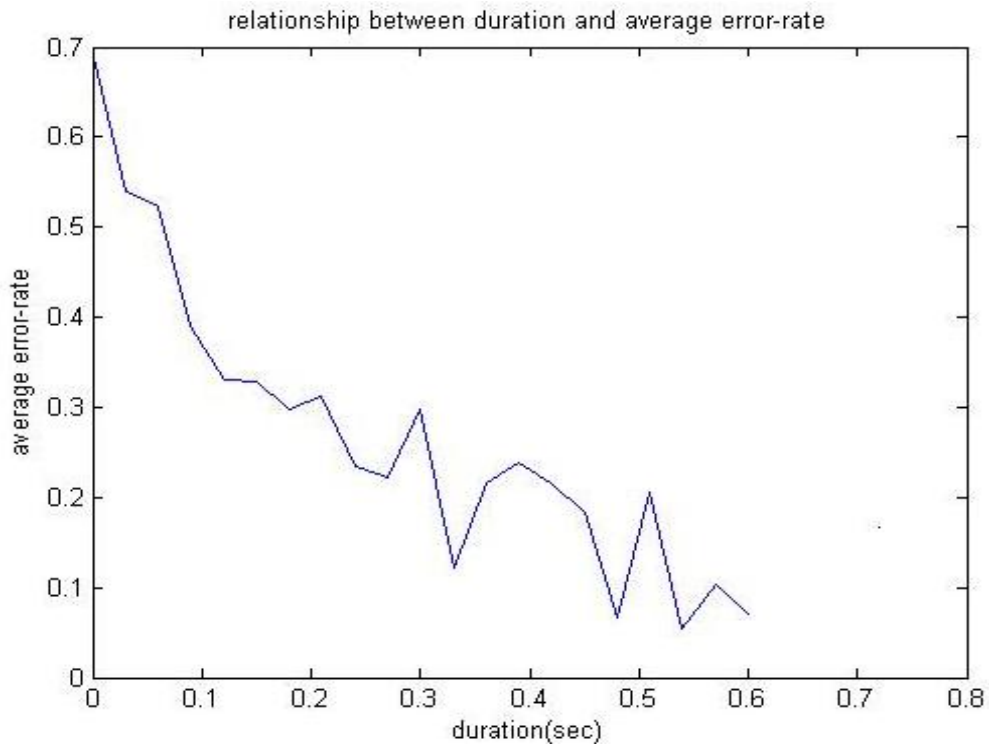


Figure 4. The relationship between duration and error rate shows that longer words generally result in better performance, but the overall variance of this measure is high.

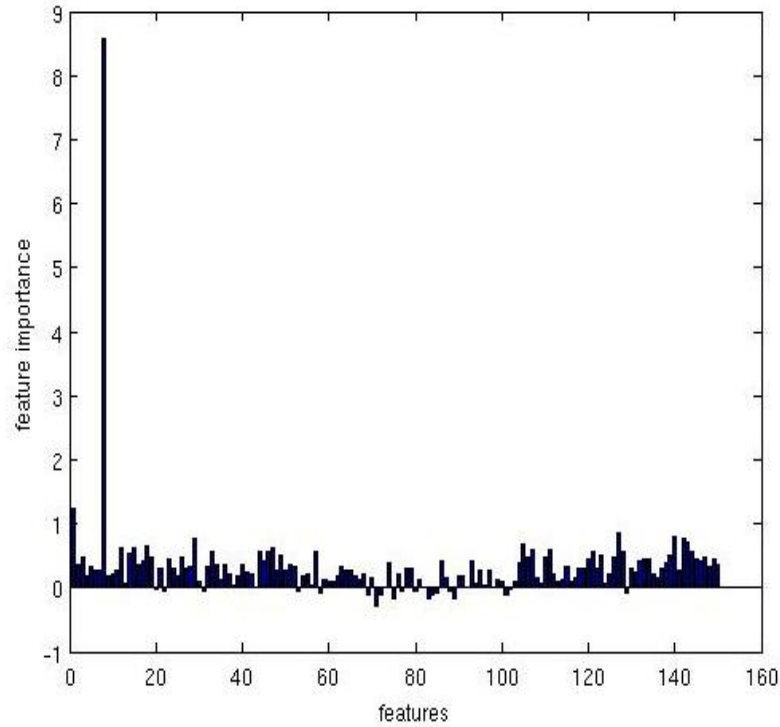


Figure 5. Feature importance based on the RF algorithm is shown. The feature "count," which represents the frequency of occurrence of a word, is by far the singlemost valuable feature since it is not correlated with any of the other features.

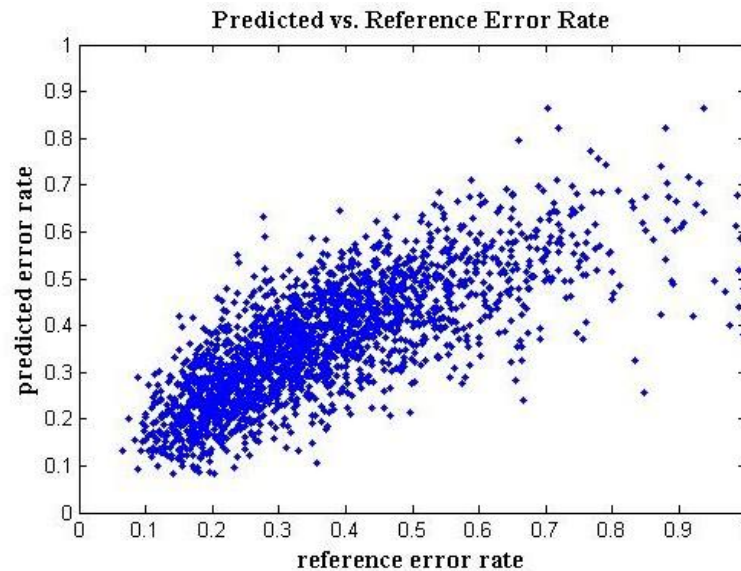


Figure 6. The predicted error rate is plotted against the reference error rate, demonstrating good correlation between the two.

XI. LIST OF TABLES

Table 1. A mapping of phones to broad phonetic classes is shown.

Table 2. The number of features is shown for different feature selection methods as a function of the mean square error (MSE) on both the training and test sets. Performance for the correlation and MSE criteria was comparable.

Table 3. The correlation of predicted error rates with actual error rates is shown for our acoustic distance measure. Performance on the eval set is comparable for sets 1 and 2 for a broad range of parameter settings. The correlation between open set and closed set performance is also good.

Table 4. Results are shown for the phonetic distance algorithm as a function of the number of nearest neighbors used in kNN.

Table 5. A comparison of the different classification algorithms as a function of the feature sets is shown. R values are shown (the MSE results follow the same trend). Random forests (RF) give very stable results across a wide range of conditions.

Table 6. The correlation between various classifiers is shown. The acoustic-based distance is least correlated with the phonetic-based approach, indicating there could be a benefit to combining these predictors.

Table 7. Performance improves slightly by combining many predictors using PSO. The acoustic and feature-based metrics contribute equally to the overall result.

XII. TABLES

Class	Phonemes
Silence	sp sil
Stops	b p d t g k
Fricatives	jh ch sh s z zh f th v dh hh
Nasals	m n ng en
Liquids	l el r w y
Vowels	iy ih eh ey ae aa aw ay ah ao ax oy ow uh iw er

Table 2. A mapping of phones to broad phonetic classes is shown.

Method	No. Feats	MSE (Train)	MSE (Eval)
All Features / LR / Corr	150	0.015	0.018
SFS / LR / Corr	55	0.016	0.017
SFS / LR / MSE	54	0.016	0.017
SFS / NN / Corr	12	0.015	0.015
SFS / NN / MSE	14	0.015	0.015
SFS / Tree / Corr	7	0.015	0.020
SFS / Tree / MSE	7	0.016	0.019
RF	56	0.006	0.014

Table 1. The number of features is shown for different feature selection methods as a function of the mean square error (MSE) on both the training and test sets. Performance for the correlation and MSE criteria was comparable.

Set	K	k	Train		Eval	
			MSE	R	MSE	R
1	1	1	0.027	0.227	0.027	0.270
1	1	3	0.025	0.340	0.025	0.370
1	1	5	0.024	0.394	0.023	0.425
1	1	30	0.021	0.528	0.020	0.543
1	1	inf	0.023	0.456	0.022	0.471
1	2	1	0.026	0.293	0.025	0.330
1	2	3	0.024	0.414	0.023	0.444
1	2	5	0.022	0.461	0.022	0.473
1	2	30	0.019	0.569	0.019	0.583
1	2	inf	0.018	0.601	0.018	0.615
1	3	5	0.022	0.475	0.022	0.497
1	3	30	0.019	0.565	0.019	0.579
1	3	inf	0.018	0.600	0.018	0.614
1	4	5	0.022	0.477	0.021	0.499
1	4	30	0.020	0.542	0.020	0.559
1	4	inf	0.019	0.578	0.018	0.595
1	12	5	0.024	0.397	0.023	0.432
1	12	30	0.021	0.503	0.021	0.520
1	12	inf	0.021	0.519	0.020	0.542
2	2	5	0.024	0.387	0.024	0.407
2	4	inf	0.020	0.550	0.019	0.568
2	15	inf	0.021	0.526	0.020	0.551
2	17	inf	0.021	0.526	0.020	0.551

Table 3. The correlation of predicted error rates with actual error rates is shown for our acoustic distance measure. Performance on the eval set is comparable for sets 1 and 2 for a broad range of parameter settings. The correlation between open set and closed set performance is also good.

k	Train		Eval	
	MSE	R	MSE	R
1	0.026	0.296	0.026	0.322
3	0.024	0.405	0.024	0.421
5	0.023	0.434	0.023	0.451
30	0.021	0.502	0.021	0.519
50	0.021	0.503	0.021	0.519
100	0.021	0.499	0.021	0.515
300	0.022	0.483	0.022	0.498
inf	0.023	0.459	0.022	0.478

Table 4. Results are shown for the phonetic distance algorithm as a function of the number of nearest neighbors used in kNN.

Classifier Method	No. Feats	LR		NN		RF	
		Train	Eval	Train	Eval	Train	Eval
All Features / LR / Corr	150	0.683	0.618	0.724	0.624	0.895	0.708
SFS / LR / Corr	55	0.654	0.629	0.753	0.692	0.875	0.701
SFS / LR / MSE	54	0.654	0.629	0.735	0.686	0.857	0.697
SFS / NN / Corr	12	0.571	0.573	0.697	0.691	0.776	0.676
SFS / NN / MSE	14	0.573	0.574	0.697	0.689	0.799	0.679
SFS / Tree / Corr	7	0.561	0.564	0.674	0.669	0.761	0.659
SFS / Tree / MSE	7	0.561	0.564	0.674	0.669	0.761	0.659
RF	56	0.635	0.604	0.734	0.675	0.882	0.703

Table 5. A comparison of the different classification algorithms as a function of the feature sets is shown. R values are shown (the MSE results follow the same trend). Random forests (RF) give very stable results across a wide range of conditions.

	Acoustic	Phonetic	Feature
Acoustic	1	0.4	0.6
Phonetic	0.4	1	0.7
Feature	0.6	0.7	1

Table 7. The correlation between various classifiers is shown. The acoustic-based distance is least correlated with the phonetic-based approach, indicating there could be a benefit to combining these predictors.

Machines	Train		Eval		Relative Contribution		
	MSE	R	MSE	R	Acoustic	Phonetic	Feature
All	0.00092	0.913	0.012	0.760	41.1%	10.5%	48.3%
NN+RF	0.00084	0.918	0.012	0.762	44.7%	15.7%	39.5%

Table 6. Performance improves slightly by combining many predictors using PSO. The acoustic and feature-based metrics contribute equally to the overall result.

I. LIST OF FIGURES

Figure 1. Spoken term detection can be partitioned into two tasks: indexing and search. One common approach to indexing is to use a speech to text system (after Fiscus et al., 2007).

Figure 2. A prototype of a web-based application that predicts voice keyword search term reliability is shown. The search term reliability is automatically updated as the user types a search term. A demonstration is available at http://www.isip.piconepress.com/projects/ks_prediction/demo/current/.

Figure 3. In our approach to predicting search term reliability, we decompose terms into features, such as N-grams of phonemes and the number of phonemes, and apply these features to a variety of machine-learning algorithms.

Figure 4. The relationship between duration and error rate shows that longer words generally result in better performance, but the overall variance of this measure is high.

Figure 5. Feature importance based on the RF algorithm is shown. The feature "count," which represents the frequency of occurrence of a word, is by far the singlemost valuable feature since it is not correlated with any of the other features.

Figure 6. The predicted error rate is plotted against the reference error rate, demonstrating good correlation between the two.

II. FIGURES

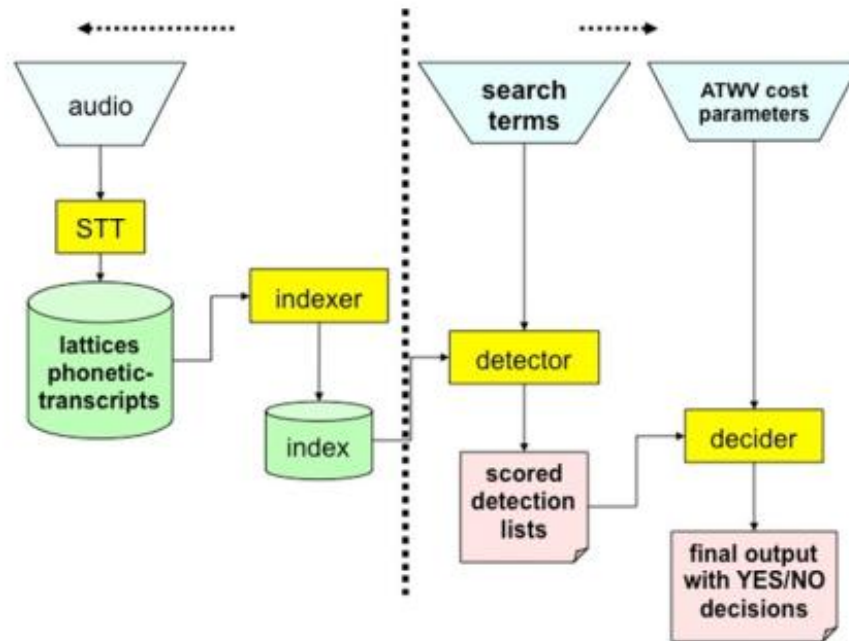


Figure 1. Spoken term detection can be partitioned into two tasks: indexing and search. One common approach to indexing is to use a speech to text system (after Fiscus et al., 2007).

Keyword

Strength

Weak 71% Strong

Selecting a Good Voice Keyword Search Term

- Longer words such as *"brainstorm"* are better than shorter words such as *"brain."*
- Words such as *"onomatopoeia"* that have more vowels and consonant-vowel transitions are better than words such as *"turtle"* that have fewer vowels.
- Frequently used words, such as *"shopping"* are better than less frequently used words such as *"Chicago."*
- Polysyllabic words such as *"democracy"* are better than monosyllabic words such as *"vote."*
- Multi-word phrases are better handled as structured queries using the intrinsic capabilities of your search tool.

What is Voice Keyword Search(VKS)?

- The voice signal is indexed using a pattern recognition algorithm such as a speech recognition system so that it can be quickly searched for content.
- Accuracy varies based on a variety of linguistic and acoustic properties of the search term.
- Voice queries are becoming increasingly popular for web-enabled phones and edge devices.

The demo is linked to the [Microsoft Research Audio Video Indexing System \(MAVIS\)](#).

Figure 2. A prototype of a web-based application that predicts voice keyword search term reliability is shown. The search term reliability is automatically updated as the user types a search term. A demonstration is available at http://www.isip.piconepress.com/projects/ks_prediction/demo/current/.

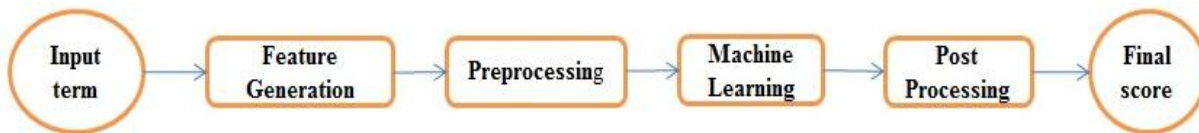


Figure 3. In our approach to predicting search term reliability, we decompose terms into features, such as N-grams of phonemes and the number of phonemes, and apply these features to a variety of machine-learning algorithms.

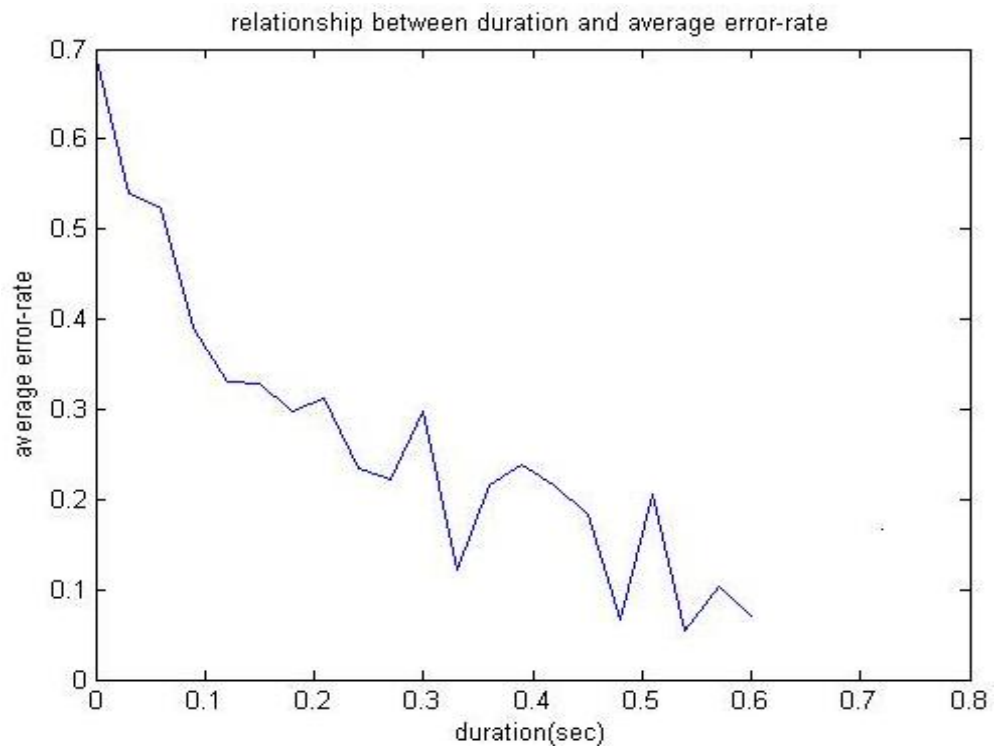


Figure 4. The relationship between duration and error rate shows that longer words generally result in better performance, but the overall variance of this measure is high.

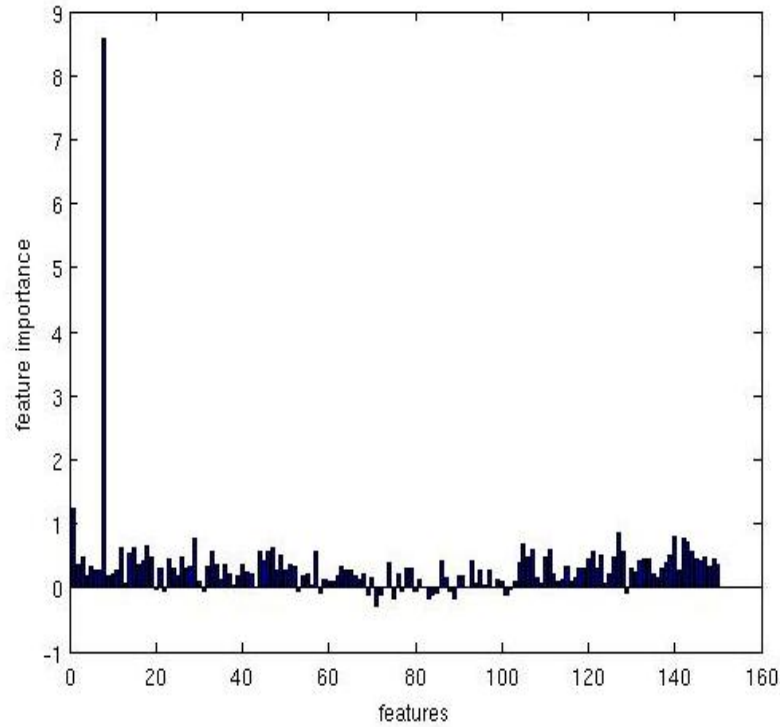


Figure 5. Feature importance based on the RF algorithm is shown. The feature "count," which represents the frequency of occurrence of a word, is by far the singlemost valuable feature since it is not correlated with any of the other features.

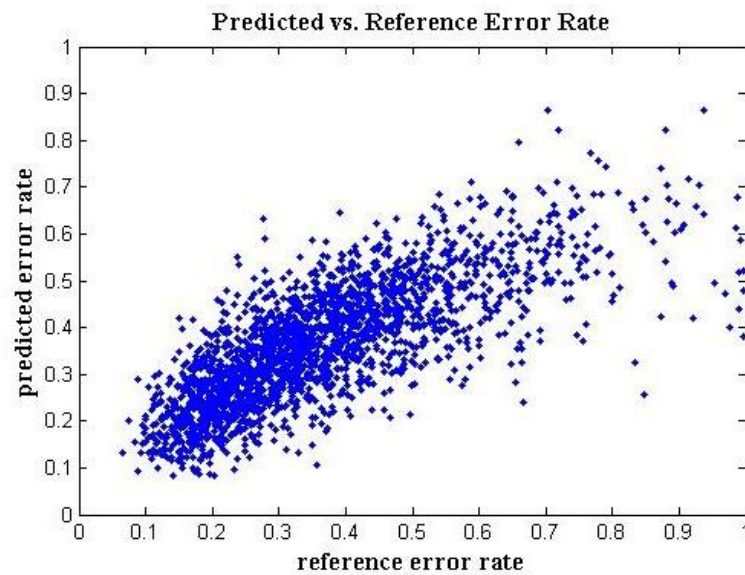


Figure 6. The predicted error rate is plotted against the reference error rate, demonstrating good correlation between the two.

I. LIST OF TABLES

Table 1. A mapping of phones to broad phonetic classes is shown.

Table 2. The number of features is shown for different feature selection methods as a function of the mean square error (MSE) on both the training and test sets. Performance for the correlation and MSE criteria was comparable.

Table 3. The correlation of predicted error rates with actual error rates is shown for our acoustic distance measure. Performance on the eval set is comparable for sets 1 and 2 for a broad range of parameter settings. The correlation between open set and closed set performance is also good.

Table 4. Results are shown for the phonetic distance algorithm as a function of the number of nearest neighbors used in kNN.

Table 5. A comparison of the different classification algorithms as a function of the feature sets is shown. R values are shown (the MSE results follow the same trend). Random forests (RF) give very stable results across a wide range of conditions.

Table 6. The correlation between various classifiers is shown. The acoustic-based distance is least correlated with the phonetic-based approach, indicating there could be a benefit to combining these predictors.

Table 7. Performance improves slightly by combining many predictors using PSO. The acoustic and feature-based metrics contribute equally to the overall result.

II. TABLES

Class	Phonemes
Silence	sp sil
Stops	b p d t g k
Fricatives	jh ch sh s z zh f th v dh hh
Nasals	m n ng en
Liquids	l el r w y
Vowels	iy ih eh ey ae aa aw ay ah ao ax oy ow uh iw er

Table 2. A mapping of phones to broad phonetic classes is shown.

Method	No. Feats	MSE (Train)	MSE (Eval)
All Features / LR / Corr	150	0.015	0.018
SFS / LR / Corr	55	0.016	0.017
SFS / LR / MSE	54	0.016	0.017
SFS / NN / Corr	12	0.015	0.015
SFS / NN / MSE	14	0.015	0.015
SFS / Tree / Corr	7	0.015	0.020
SFS / Tree / MSE	7	0.016	0.019
RF	56	0.006	0.014

Table 1. The number of features is shown for different feature selection methods as a function of the mean square error (MSE) on both the training and test sets. Performance for the correlation and MSE criteria was comparable.

Set	K	k	Train		Eval	
			MSE	R	MSE	R
1	1	1	0.027	0.227	0.027	0.270
1	1	3	0.025	0.340	0.025	0.370
1	1	5	0.024	0.394	0.023	0.425
1	1	30	0.021	0.528	0.020	0.543
1	1	inf	0.023	0.456	0.022	0.471
1	2	1	0.026	0.293	0.025	0.330
1	2	3	0.024	0.414	0.023	0.444
1	2	5	0.022	0.461	0.022	0.473
1	2	30	0.019	0.569	0.019	0.583
1	2	inf	0.018	0.601	0.018	0.615
1	3	5	0.022	0.475	0.022	0.497
1	3	30	0.019	0.565	0.019	0.579
1	3	inf	0.018	0.600	0.018	0.614
1	4	5	0.022	0.477	0.021	0.499
1	4	30	0.020	0.542	0.020	0.559
1	4	inf	0.019	0.578	0.018	0.595
1	12	5	0.024	0.397	0.023	0.432
1	12	30	0.021	0.503	0.021	0.520
1	12	inf	0.021	0.519	0.020	0.542
2	2	5	0.024	0.387	0.024	0.407
2	4	inf	0.020	0.550	0.019	0.568
2	15	inf	0.021	0.526	0.020	0.551
2	17	inf	0.021	0.526	0.020	0.551

Table 3. The correlation of predicted error rates with actual error rates is shown for our acoustic distance measure. Performance on the eval set is comparable for sets 1 and 2 for a broad range of parameter settings. The correlation between open set and closed set performance is also good.

k	Train		Eval	
	MSE	R	MSE	R
1	0.026	0.296	0.026	0.322
3	0.024	0.405	0.024	0.421
5	0.023	0.434	0.023	0.451
30	0.021	0.502	0.021	0.519
50	0.021	0.503	0.021	0.519
100	0.021	0.499	0.021	0.515
300	0.022	0.483	0.022	0.498
inf	0.023	0.459	0.022	0.478

Table 4. Results are shown for the phonetic distance algorithm as a function of the number of nearest neighbors used in kNN.

Classifier Method	No. Feats	LR		NN		RF	
		Train	Eval	Train	Eval	Train	Eval
All Features / LR / Corr	150	0.683	0.618	0.724	0.624	0.895	0.708
SFS / LR / Corr	55	0.654	0.629	0.753	0.692	0.875	0.701
SFS / LR / MSE	54	0.654	0.629	0.735	0.686	0.857	0.697
SFS / NN / Corr	12	0.571	0.573	0.697	0.691	0.776	0.676
SFS / NN / MSE	14	0.573	0.574	0.697	0.689	0.799	0.679
SFS / Tree / Corr	7	0.561	0.564	0.674	0.669	0.761	0.659
SFS / Tree / MSE	7	0.561	0.564	0.674	0.669	0.761	0.659
RF	56	0.635	0.604	0.734	0.675	0.882	0.703

Table 5. A comparison of the different classification algorithms as a function of the feature sets is shown. R values are shown (the MSE results follow the same trend). Random forests (RF) give very stable results across a wide range of conditions.

	Acoustic	Phonetic	Feature
Acoustic	1	0.4	0.6
Phonetic	0.4	1	0.7
Feature	0.6	0.7	1

Table 7. The correlation between various classifiers is shown. The acoustic-based distance is least correlated with the phonetic-based approach, indicating there could be a benefit to combining these predictors.

Machines	Train		Eval		Relative Contribution		
	MSE	R	MSE	R	Acoustic	Phonetic	Feature
All	0.00092	0.913	0.012	0.760	41.1%	10.5%	48.3%
NN+RF	0.00084	0.918	0.012	0.762	44.7%	15.7%	39.5%

Table 6. Performance improves slightly by combining many predictors using PSO. The acoustic and feature-based metrics contribute equally to the overall result.