



ELSEVIER

Speech Communication 32 (2000) 61–77

SPEECH
COMMUNICATION

www.elsevier.nl/locate/specom

Experiments in spoken document retrieval using phoneme n -grams

Corinna Ng ^{a,*}, Ross Wilkinson ^b, Justin Zobel ^a

^a *Department of Computer Science, RMIT, GPO Box 2476V, Melbourne 3001, Australia*

^b *Division of Mathematical and Information Science, CSIRO, 723 Swanston Street, Carlton, Vic. 3053, Australia*

Abstract

In spoken document retrieval (SDR), speech recognition is applied to a collection to obtain either words or subword units, such as phonemes, that can be matched against queries. We have explored retrieval based on phoneme n -grams. The use of phonemes addresses the out-of-vocabulary (OOV) problem, while use of n -grams allows approximate matching on inaccurate phoneme transcriptions. Our experiments explored the utility of word boundary information, stopword elimination, query expansion, varying the length of phoneme sequences to be matched and various combinations of n -grams of different lengths. Given word-based recognition (WBR), we can match queries to speech using a phoneme representation of the words, permitting us to test whether it was the recognition or the matching process that was most crucial to retrieval performance. Our experiments show that there is some deterioration in effectiveness, but the particular form of matching is less vital if the sequence of phonemes was correct. When phone sequences are recognised directly, with higher error rates than for words, it was more important to select a good matching approach. Varying gram length trades precision against recall; combination of n -grams of different lengths, in particular 3-grams and 4-grams, can improve retrieval. Overall, phoneme-based retrieval is not as effective as word-based retrieval, but is sufficient for situations in which word-based retrieval is either impractical or undesirable. © 2000 Elsevier Science B.V. All rights reserved.

Zusammenfassung

Document Retrieval mittels Sprache kann durch Spracherkennung umgesetzt werden, welche auf eine Datensammlung angewendet wird, um Worte oder Teilworte wie zum Beispiel Phoneme zu generieren, die wiederum in Anfragen Verwendung finden können. Wir haben Datenretrieval basierend auf n -Gram-Phoneme untersucht. Der Gebrauch von Phonemen adressiert das “out-of-vocabulary” (OOV) Problem, während n -Gram-Phoneme den approximierenden Abgleich von ungenauen phonetisch Beschreibungen erlaubt. In unsere Experimenten haben wir den Nutzen von verschiedene Techniken untersucht: die Bedeutung von Wortgrenzen, Stopworteliminierung, Anfrageexpansion, Längenvariierung der Phonemsequenzen zum Abgleich und verschiedene Längenkombinationen der n -Gramme. Mittels wortbasierende Spracherkennung ist es möglich, Sprache in eine Phonemrepräsentation abzubilden, um eine Anfrage zu formulieren. Dies erlaubt es zu ermitteln, ob die Spracherkennung oder der Abbildungsprozess hinsichtlich der Retrieval Leistung eine größere Rolle spielt. Unsere Experimente zeigen eine Verschlechterung der Effektivität, aber die genaue Form des Abgleichs ist weniger wichtig, wenn die Phonemsequenz genau gegeben ist. Wenn die genaue Lautsequenz mit einer höheren Fehlerrate als für Worte erkannt wird, ist es umso wichtiger, dass eine gute Abbildungsmethode gewählt wird. Wir beobachten einen Tradeoff zwischen Precision and Recall wenn die Länge der

* Corresponding author.

E-mail addresses: chienn@cs.rmit.edu.au (C. Ng), ross.wilkinson@cmis.csiro.au (R. Wilkinson), jz@cs.rmit.edu.au (J. Zobel).

Grame variiert. Im besonderen verzeichnen wir eine Verbesserung der Anfrageleistung, wenn verschiedene n -Gram-Längen etwa 3-Grame und 4-Grame miteinander kombiniert werden. Abschließend läßt sich sagen, daß phonem-basierende Anfragesysteme nicht so effektiv sind wie wortbasierende, aber es genügt in Situationen, in welchen wortbasierende Anfragen nicht praktikabel oder unerwünscht sind. © 2000 Elsevier Science B.V. All rights reserved.

Résumé

Dans les systemes recherche d'information vocale, la reconnaissance de la parole est appliquee a une collection pour obtenir certaines information qui correspond aux contraintes des requetes. Ces information peuvent etre des mots ou des composants de mots. On a explore la recherche basee sur phoneme de " n -grams". L'utilisation de ces phonemes concerne le probleme du "vocabulaire exterieure", et l'utilisation du concept de " n -gram" permet le matching approximatif dans un environnement imprecis de transcription. Les experiences (c'est-a-dire, les tests) on permit d'explorer beaucoup d'aspects, comme par exemple, "word boundary information", elimination des points d'arret, expansion des requetes, variation de la longueur des sequence de phonemes, et la combinaison des n -grams. Pour un mot donne, on peut associer des requetes a la parole en utilisant la representation phonetique des mots. Cela a permis de decider les aspects important dans la recherche, comme par exemple, la reconnaissance ou le processus de matching. Nos resultat experimentals ont montre qu'il y a une deterioration dans l'efficacite de la reche, mais pour des cas particulier de matching, cela n'était pas important parce que la sequencee des phonemes est correcte. Dans les cas ou les suite de phoneme sont directement reconnues, il etait important de selectionner une bonne approche de matching. La combinaison de n -grams de different longueurs (3-grams and 4-grams) on permit d'ameliorer l'efficacite de la methode de recherche. Dans le cas general, la recherche basee sur le phoneme n'est pas efficace quand celle ci est comparee a la methode qui est basee sur les mots. Cette methode (qui est basee sur le phoneme) est interessante dans les situations ou les methodes basees sur les mots ne some pas utilisable. © 2000 Elsevier Science B.V. All rights reserved.

1. Introduction

Information retrieval techniques are widely used in text databases to identify documents that are likely to be relevant to free-text queries (Salton and McGill, 1983). The aim of spoken document retrieval is to provide similar functionality for databases of spoken documents. Such documents, in the form of audio signals, are recorded from many different sources, such as news broadcasts on radio and television. It would be valuable to be able to retrieve such documents in response to textual or spoken queries.

In order to query such media interactively, speech signals are converted into words, phonemes, or other subword units (Lee, 1989), using a speech recognition system. This work focuses on the phoneme-based approach, where spoken documents are recognised as phoneme sequences and retrieval is based on matching the n -grams (sequences of n symbols) of these transcriptions.

These experiments were conducted as part of the Text REtrieval Conference (TREC), sponsored

by NIST and DARPA to encourage research in information retrieval (Voorhees and Harman, 1998). In the TREC spoken document retrieval (SDR) experiments, word-based approaches have consistently outperformed phoneme approaches. However, there are several reasons for using phonemes.

In word-based recognition (WBR), three assumptions can be made about the recognition process. First, the recogniser uses a large vocabulary (Woodland et al., 1995) containing most of the words to be recognised. In TREC-7 SDR, a collection of 100 h of spoken documents contains about 23,000 unique words. Therefore, obtaining good recognition performance requires a vocabulary of at least that size. Second, the spoken documents must be consistent with the language model used. The model determines the recognition language structure, embodying certain assumptions about exactly what word can precede another word. If the model can predict the language well, higher recognition performance can be achieved. Finally, the computational resources required to perform WBR must be available.

We consider the situation where the above assumptions cannot be made. If the recognition process is on a small hand-held device, where the computational resources are limited, a large vocabulary may not be possible. Furthermore, the resources required to build a language model may not be available and, even with a small language model, the spoken documents may not be consistent with it, thus limiting its usefulness. Therefore, we consider a basic recognition system, such as a phoneme recogniser, which can run on a small processor. Due to its simplicity, it suffers from poor recognition performance compared to the WBR approach. Second, we consider a spoken document collection containing a significantly high proportion of out-of-vocabulary (OOV) words which can adversely affect retrieval when misrecognised. An example is a spoken collection of names and places. From a textual retrieval perspective, these words, due to their rarity in the document collection, usually assist in identification of relevant documents. The WBR system, not having seen sufficient examples of these words in the training process, may well incorrectly recognise them as some other words.

In our work, a basic phoneme recogniser is treated as a black box. Techniques are investigated here to improve retrieval and tested on standardised corpora such as the TREC-7 and TREC-6 SDR collections. These collections do not exhibit any of the properties outlined previously and hence they do not allow us to test our assumptions fully. Therefore, we do not expect our phoneme n -gram approach to be more effective than the word-based approach to retrieval. But using these available corpora allows us to test some of the retrieval properties of phoneme n -grams.

The technique of using n -grams in retrieval (Cavnar, 1994) has shown reasonable performance in text-based collections. Previous experiments by Wechsler et al. (1998), Wechsler and Schäuble (1995), Ng (1998), Ng and Zue (1997, 1998) and Smeaton et al. (1997) on n -gram retrieval from phoneme transcriptions obtained directly from a phoneme recogniser showed that phoneme n -gram retrieval can be effective in practice. Other experiments have shown that phoneme retrieval can be used to complement word retrieval (Witbrock and

Hauptmann, 1997) when word recognition has failed, especially in situations, where names and unknown words are misrecognised.

In this paper, we refine methods for using phoneme n -grams. We used two kinds of phoneme recognition: direct recognition as a series of phonemes and transformation to phonemes from recognised words. In a typical word-based approach to retrieval, word boundary information and stopping (removal of common words) in the documents and queries usually aids retrieval performance (Fuller et al., 1997). We investigate whether the same properties hold for the phoneme n -gram approach to retrieval. The first set of experiments compare retrieval effectiveness for the two kinds of phoneme n -grams, either from direct phoneme recognition or transformation from words. Queries are varied to test the effects of stopping and word boundary information. We combined phoneme n -grams of varying lengths to investigate whether different evidence was being provided by various n -gram sizes. In addition, we examined an approximate technique of query expansion, adding neighbour terms to the query set. In another approach to query expansion, a list of the most frequent misrecognised query terms was obtained from a training collection and used to augment the queries.

Phoneme n -gram retrieval using phonemes transformed from recognised words showed better retrieval performance than from direct phoneme recognition but was not as good as word-based retrieval. Combination of phoneme n -grams showed some improvements in retrieval, in particular combination of 3-grams and 4-grams. Overall, the phoneme-based approach to retrieval is less effective than the word-based approach, but is nonetheless effective enough to be used in practice. This paper expands and consolidates work previously reported at TREC-7 (Fuller et al., 1998) and TREC-6 (Fuller et al., 1997).

2. Document and query collections

The experiments described in this paper were based on two sets of spoken document collections provided by the linguistic data consortium (LDC).

Most of the spoken documents consist of speech segments on a variety of topics. At times, there is simultaneous speech spoken by different speakers as well as background music at varying levels. Other resources used are the pronunciation dictionary used to translate the words to phonemes and phoneme recognisers.

2.1. SDR collection from TREC-7

The TREC-7 SDR track test collection consists of 100 h of news broadcast (Linguistic Data Consortium, 1997). The full TREC-6 SDR track collection was used as the training corpus. There were approximately 2870 documents, containing about 23,000 unique words, with an average length of 275 words. For this collection, another 2580 words were added to the pronunciation dictionary to transform all words, including query words, to phoneme sequences. In practice, this is done via an automated process (Wechsler et al., 1998). Here, four versions of the transcriptions were used for retrieval:

- Reference (REF-97), human transcribed reference transcripts, assumed to have little or minimum recognition error.
- Baseline 1 (WORD₁-97), recognised transcript from CMU (33.8% word error rate).
- Baseline 2 (WORD₂-97), recognised transcript from CMU (46.6% word error rate).
- Phoneme (PHN-97), recognised phoneme-based transcript of our own, built using the HTK system, described later (45% phoneme error rate).

The figures for the word error rate were obtained from the overview given at TREC-7 (Garofolo et al., 1998). The first three sets of transcripts are word-based documents while the last set of transcripts is phoneme-based. The retrieval task here is the traditional ad hoc relevance task, where topics and assessments were generated by human assessors. In this task, there was an average of about 17 relevant documents per query.

The query set consists of 23 queries, with an average length of 16 words. In these experiments, the original set of queries is denoted as (full); a stopped derivation is denoted as (stopped). A stoplist of approximately 370 words was used to

derive the stopped queries. The average length of stopped queries was 9 words.

2.2. SDR collection from TREC-6

The TREC-6 SDR collection was taken from the LDC 1996 Broadcast News corpus, as used in the DARPA 1996 “Hub-4” experiments (Linguistic Data Consortium, 1996a; Linguistic Data Consortium, 1996b). There were approximately 1450 documents in 50 h of broadcast news. The document collection consisted of 18,000 unique words, 2000 of which were not in the pronunciation dictionary.

Three versions of the transcriptions were used for experimentation:

- Reference (REF-96), human-transcribed reference transcripts.
- Baseline (WORD-96), recognised transcript from IBM, with a 50% word error rate (Garofolo et al., 1998).
- Speech (PHN-96), recognised phoneme-based transcript from ETH (Swiss Institute of Technology), with a 45% phoneme error rate.

Both the reference (REF-96) and baseline (WORD-96) collections are word-based while the speech (PHN-96) collection is phoneme-based. The retrieval task here is a known-item search, where it is assumed that there is only one relevant document per query. For this TREC-6 collection, a set of 49 queries were used. The stoplist of about 370 words used in TREC-7 was applied here to stop the query set. The average length of the full queries was 12 words, while stopped queries averaged 6 words.

2.3. Phoneme recognisers

Each of the TREC-7 and TREC-6 collections had a phoneme transcript, but from different phoneme recognisers. The recognisers were developed using the HTK toolkit (Young et al., 1995). Both used a speaker-independent phoneme recogniser based on hidden Markov models (Rabiner and Juang, 1993).

For TREC-6, the phoneme recognition system used was developed at the Swiss Institute of Technology. Models used included acoustic models for 40 monophones, trained using the TIMIT

speech corpus and context-dependent biphone models, trained on the TREC-6 SDR training collection. The recognition process used a stochastic phone-bigram language model to eliminate generation of less probable phone sequences. These phone sequences were then post-processed by clustering acoustically similar monophones into 30 broader classes called phonemes (Mateev et al., 1997). A recognition accuracy of about 45% was obtained when evaluated on 7.5 h of the TREC-6 SDR training collection.

For TREC-7, the models used were 39 monophones and around 1500 right-context biphones. These models were trained on about 19 h of clean TREC-6 SDR data. The right-context models were chosen instead of the left-context models because informal experiments showed that the right-context models achieved better recognition performance. A back-off bigram language model was also used to increase recognition performance. Recognition accuracy of about 45% was obtained when evaluated for about 5 h of the TREC-7 SDR training collection.

2.4. Translation of words to phonemes using a pronunciation dictionary

The CMU pronunciation dictionary¹ was used to translate words to phoneme sequences for all the documents and queries. This approach required a large pronunciation dictionary, which must include all the words in the spoken document collection. For TREC-6, the original pronunciation dictionary had about 110,000 entries and 2000 words had to be added from the test collection. A year later, for TREC-7, an additional 2580 words from the test set had to be added. The kind of words that had to be added were mainly compounds including names as well as inflections of words already in the dictionary. For compound words, pronunciation entries were concatenated; inflections of known words were added by modifying the

phoneme sequences of the suffix to those of the required inflections. Unknown words such as names were added by first using a rule-based word-to-phoneme algorithm before converting to the same phoneme set as that used in the dictionary. These OOV words did not have a significant impact on the retrieval process because they were not used in the queries. None of the queries contained OOV words.

Experiments using documents translated to phonemes are shown with the extension (-phn). Queries, both (full) and (stopped), were translated to phonemes using the same dictionary.

2.5. Creation of *n*-grams from phonemes

To allow matching of phoneme strings we represented them as *n*-grams. Any sequence of symbols can be transformed into a sequence of *n*-grams; e.g., transforming the sequence “ABC-DEF” into a sequence of 3-grams yields “ABC BCD CDE DEF”. We varied *n*-gram size between 2 and 5 inclusive; documents and queries in phoneme *n*-gram form are shown with an extension of (-phn-*N*) where *N* denotes the size of the *n*-gram.

Two types of *n*-gram sequence were formed for each type of query. The first form had *n*-grams created across word boundaries, so that the information as to where each word started or ended was removed. The second form disallowed creation of *n*-grams across word boundaries. Therefore, *n*-grams were only formed within each word’s phoneme sequence. These are denoted as (bound). Bounded *n*-grams of stopped queries (stopped, bound) were similarly created.

3. Retrieval system

In a retrieval environment, queries and documents can be represented by vectors in a high-dimensional vector space. The similarity between a query *q* and document *d* can then be estimated by the cosine measure. One formulation that is effective for word-based retrieval (Salton and Buckley, 1988) is

¹ Carnegie Mellon University Pronouncing Dictionary, 1995. CmuDict.0.4. Available at <http://www.speech.cs.cmu.edu/cgi-bin/cmuDict>.

$$\text{sim}(q, d) = \frac{\sum_{t \in q \cap d} (w_{q,t} \cdot w_{d,t})}{\sqrt{\left(\sum_{t \in q} (w_{q,t})^2 \cdot \sum_{t \in d} (w_{d,t})^2 \right)}}, \quad (1)$$

where $w_{x,t}$ is defined as the weight of term t in either document d or query q . The weight of an index term t in document d is

$$w_{d,t} = \log_2(f_{d,t} + 1), \quad (2)$$

where the term frequency $f_{d,t}$ is the number of occurrences of term t in document d . Similarly, $w_{q,t}$ is defined as

$$w_{q,t} = \log_2 \left(\frac{N}{f_t} + 1 \right), \quad (3)$$

where N is the total number of documents in the collection and f_t is the number of documents containing term t . This measure was used for the TREC-6 SDR experiments (Fuller et al., 1997). The equivalent of this cosine measure in the SMART system is approximated to lxc.btx (Salton and Buckley, 1988; Singhal, 1997).

Another similarity measure is the Okapi formulation (Walker et al., 1997), defined as

$$\begin{aligned} \text{sim}(q, d) &= \sum_{t \in q \cap d} \frac{(k_1 + 1) \cdot f_{d,t}}{k_1 \cdot [(1 - b) + b \cdot (W_d / (\text{avr}_W))]] + f_{d,t}} \\ &\quad \cdot \frac{(k_3 + 1) \cdot f_{q,t}}{k_3 + f_{q,t}} \cdot \log \frac{N - f_t + 0.5}{f_t + 0.5}, \end{aligned} \quad (4)$$

where k_1 , k_3 and b are constants, set in our experiments to 1.2, 1000 and 0.75, respectively (as suggested by the City University group (Walker et al., 1997)). The value W_d is the length of document d in bytes and avr_W is the average document length (in bytes) across the collection. The value N is the total number of documents in the collection, f_t is the number of documents containing term t and $f_{x,t}$ is the frequency of term t in either document d or query q . This Okapi measure was used for the TREC-7 SDR experiments (Fuller et al., 1998). The MG (Witten et al., 1994) text engine was used for all retrieval experiments. The more complex Okapi formulation (4) was preferred to the simpler similarity formula (1) be-

cause previous experiments found that it led to improved retrieval performance.

The parameters used in the Okapi formulation were based on ad hoc text collections, several times the size of the collections used here. The parameters were tailored for TREC ad hoc collections and had not been tested on the TREC SDR collections. The parameters used in the current Okapi formula (4) may not be suitable on the SDR collections, especially in the phoneme n -gram forms and in practice may need to be modified.

4. Experimental questions

The primary objective of the experiments described in this paper were to investigate whether techniques that usually improve retrieval performance for text-based collections can improve retrieval performance for SDR retrieval using phoneme n -grams.

A technique used for word-based retrieval is to remove stopwords, usually consisting of high-frequency function words such as conjunctions, prepositions and pronouns. These words do not contribute much to the overall weight of the document or query. If a general stopword list is used in a non-typical document collection, then it is possible that aggressive stopword removal may degrade retrieval effectiveness. Tested on the TREC ad hoc collection (Fuller et al., 1997), we found that stopping improved retrieval effectiveness by about 12% (after casefolding). From a phoneme-based SDR perspective, it is not possible to apply stopping to the documents, but it can be applied to the text-based queries. In the experiments conducted here, a stopword list of about 370 words was used to stop the queries. The query sets, both stopped (stopped) and non-stopped (full), were translated to phonemes before being converted to n -grams.

Documents recognised directly as phonemes do not have word boundary information but the queries, which are in textual word form, do. The effect of word boundaries in queries can be tested by not permitting phoneme n -grams to be formed across query words. It is possible that phoneme n -grams formed across word boundaries are

causing too many false matches. An example of how word boundary information should aid retrieval is the query phrase “olympic torch”, whose phoneme equivalent is “aWlmpik tOrJ”. Assuming no recognition error for the phrase and retrieval using 3-grams, the gram “ikt” could possibly match the word “predictably”, whose phoneme transcription is “pridiktabliw”.

The lengths of phoneme n -grams of the documents and queries were varied from 2 to 5. The effects of different sizes of n -grams on retrieval effectiveness were investigated and the combination of phoneme n -grams of different lengths was also tested. For each document, n -grams of differing sizes were extracted and the resulting sequences of n -grams were concatenated together to form one large document containing n -grams of different lengths. Various combinations were tried, varying from combination of two n -gram lengths to combining all n -gram lengths. By combining n -grams of different length, we can investigate whether different n -gram sizes provide different forms of evidence for retrieval (Jones et al., 1996). Combination permitted us to investigate whether boundary information and stopping affected retrieval on larger collections.

When phoneme recognition was performed on the training collection of TREC-6 SDR, we found that, for each recognised word, similar recognition errors were occurring throughout the training collection. Therefore, it is likely that similar recognition errors should occur in the test collection. On the smaller TREC-6 test collection, we experimented with augmentation of the original query terms by a list of their incorrectly recognised forms, obtained from the training collection, selected based on their frequency of occurrence in the training collection. From the training documents in words, a list of documents containing the query words was found. By manually going through the same documents recognised as phonemes, we were able to detect the incorrectly recognised query sequences. For example, the query word “olympic”, whose phoneme equivalent is (in one representation of the international phonetic alphabet) “aWlmpik”, was incorrectly recognised as “ilimpik” and “awTawbik”. This set of augmented queries is labelled as *err. n-Grams* of the

augmented queries were created. The rationale for this is that the incorrect transcription may be able to match relevant documents, which may also contain the incorrect transcription. This method is similar in concept of using a confusion matrix-based approach by Wechsler et al. (1998) on the training collection, which can be used to determine which recognised phoneme is most likely to be recognised incorrectly as another. This technique, though not 100% accurate, is the only feasible approach for a larger collection.

Another method of query expansion, which in contrast does not make use of the training collection, is to find all n -grams, or neighbour terms, that differs in at most one character from the query n -gram. This approach used the technique of string edit distance based on phoneme distances to find neighbour terms. It had been shown to be useful for name matching on text (Zobel and Dart, 1996). Although this technique is likely to increase the false alarm rate by exhaustively finding all possible neighbour n -grams, it allowed us to preliminarily test the feasibility of retrieval using approximate string matching techniques without prior information of the recognition process. Idiosyncrasy in the recognition process will not unduly affect retrieval, in contrast to the previous expansion technique. These neighbour terms are added to the original set of terms. This new set of queries is labelled as *nbr*. This approach does not require recognition-dependent information such as confusion matrices. This tentative experiment was conducted only on the smaller TREC-6 SDR collection.

5. Results

Our experiments use the document sets described in Section 2 with the query variations based on stopping (full or stopped), word boundaries (bound or otherwise), n -grams of different sizes, n -gram combination of different sizes, and expansion (*err* or *nbr*). Table 1 summarizes the types of documents used in the experiments. To be consistent with TREC (Voorhees and Harman, 1997) evaluation methods, we retrieve up to 1000 documents per query. The retrieved documents are

Table 1

Types of documents used in experiments

REF-97	Human transcribed reference transcripts (TREC-7)
WORD ₁ -97	Recognised transcript of 33.8% word error rate (TREC-7)
WORD ₂ -97	Recognised transcript of 46.6% word error rate (TREC-7)
PHN-97	Recognised phoneme-based transcript of 45% phoneme error rate (TREC-7)
REF-96	Human transcribed reference transcripts (TREC-6)
WORD-96	Recognised transcript of 50% word error rate (TREC-6)
PHN-96	Recognised phoneme-based transcript of 45% phoneme error rate (TREC-6)

ranked according to their likely relevance based on the similarity measures described in Section 3.

Retrieval effectiveness is most often compared in terms of precision (the proportion of retrieved documents that are relevant) at different and fixed levels of recall (the proportion of relevant documents that have been retrieved). We use average precision (AP) as one measure of effectiveness, which is precision across all the queries. This is appropriate in the TREC-7 environment. Given that the task in TREC-6 SDR was a known-item search, the queries were assumed to have only one relevant document; effectiveness can therefore be determined as the reciprocal rank of the relevant document. To compare the different parameters of each retrieval experiment, the mean reciprocal rank (MRR) is computed. Another performance measure is to calculate the total number of relevant documents retrieved in the top 5 or top 10 returned documents across all queries.

5.1. Results on TREC-7 data

Results for word-based retrieval using both full and stopped queries are shown in Table 2. Documents were neither stemmed nor stopped by default. When documents were stemmed, queries were similarly stemmed prior to retrieval, to prevent mismatches between words in the documents and queries. The individual and combined effects of stopping and stemming on retrieval were investigated. Among the different transcripts, similar trends in retrieval performance were observed when documents were stopped, stemmed, or both. From Table 2, it can be seen that retrieval degraded as the word error rate increased, for both full and stopped queries. When the documents were not stopped and stemmed, a small improvement was observed using stopped queries. It was also better to use stopped queries when documents were stopped and stemmed.

Table 2

Average precision for word-level matching with different query types^a

Document set	Query set	
	Full	Stopped
REF-97	0.389	0.395
REF-97 + stopping	0.396	0.396
REF-97 + stemming	0.446	0.454
REF-97 + stopping + stemming	0.412	0.455
WORD ₁ -97	0.310	0.319
WORD ₁ -97 + stopping	0.318	0.318
WORD ₁ -97 + stemming	0.405	0.417
WORD ₁ -97 + stopping + stemming	0.364	0.420
WORD ₂ -97	0.245	0.253
WORD ₂ -97 + stopping	0.255	0.255
WORD ₂ -97 + stemming	0.340	0.326
WORD ₂ -97 + stopping + stemming	0.274	0.326

^a When documents are stemmed, the queries are similarly stemmed, to prevent mismatches of words. TREC-7 data.

AP results using the phoneme n -gram representations of the documents and queries are shown in Table 3. Phoneme n -grams of the word-based documents were allowed to cross word boundaries. The documents were not stopped or stemmed prior to translation to phonemes. n -Grams of queries, either stopped or full, were formed with and without crossing word boundaries. Comparing the average precision figures of Table 2 with Table 3, we can see that phoneme n -gram retrieval is much less effective than word-based retrieval. For direct phoneme recognised documents (PHN-97), retrieval was significantly lower than for word-based recognised documents. Phoneme n -gram matching of word-based documents increased the number of false matches, because a word sequence in phonemes had been broken into several n -grams and there had been partial matching. Experiments using phoneme 2-grams have consistently shown them to be much poorer and we do not comment on 2-grams further.

In Table 3, the effects of using stopped query terms in phoneme n -gram retrieval were different

to those observed for word-based retrieval. In most of the cases, stopped queries did not perform as well as the full queries. Additional experiments investigated how stopping word-based documents prior to converting to phoneme n -grams would affect retrieval; AP results are shown in Table 4. Retrieval effectiveness, comparing different query strategies, was slightly improved from that shown in Table 3 for phoneme 3-grams for all types of transcripts and for most cases with 4-grams. As in the word-based retrieval case, the effect of using stopped queries on retrieval effectiveness is relatively small compared to other factors, such as stopping documents, varying n -gram sizes, or using boundary information. Stopping documents is more effective, but circumstances do not always permit it. Stopping of queries on documents recognised directly as phonemes had little effect.

Table 3 shows results with phoneme n -grams of documents created across word boundaries. The effect of word boundary information in queries was investigated using both the full (full, bound) and stopped (stopped, bound) queries. Retrieval

Table 3

Average precision for different speech recognition processes, gram lengths and query types^a

Document set (unbounded phoneme n -grams)	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 2-grams</i>				
REF-97	0.104	0.156	0.117	0.148
WORD ₁ -97	0.071	0.134	0.106	0.127
WORD ₂ -97	0.061	0.108	0.095	0.104
PHN-97	0.017	0.028	0.020	0.027
<i>Experiments with phoneme 3-grams</i>				
REF-97	0.307	0.303	0.305	0.302
WORD ₁ -97	0.260	0.276	0.256	0.261
WORD ₂ -97	0.205	0.202	0.213	0.207
PHN-97	0.042	0.056	0.050	0.059
<i>Experiments with phoneme 4-grams</i>				
REF-97	0.323	0.301	0.300	0.299
WORD ₁ -97	0.294	0.253	0.256	0.251
WORD ₂ -97	0.227	0.187	0.198	0.187
PHN-97	0.101	0.098	0.085	0.098
<i>Experiments with phoneme 5-grams</i>				
REF-97	0.293	0.258	0.265	0.257
WORD ₁ -97	0.251	0.213	0.217	0.213
WORD ₂ -97	0.209	0.153	0.168	0.153
PHN-97	0.076	0.075	0.070	0.076

^a Documents were not stopped nor stemmed prior to conversion to phoneme n -grams. TREC-7 data.

Table 4

Average precision, for word-based documents which were stopped prior to conversion to phoneme n -grams, varying gram lengths and query types; TREC-7 data

Document set + stopped (unbounded phoneme n -grams)	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 2-grams</i>				
REF-97	0.094	0.132	0.163	0.169
WORD ₁ -97	0.074	0.098	0.132	0.143
WORD ₂ -97	0.070	0.082	0.118	0.121
<i>Experiments with phoneme 3-grams</i>				
REF-97	0.308	0.310	0.319	0.308
WORD ₁ -97	0.271	0.280	0.272	0.275
WORD ₂ -97	0.218	0.214	0.226	0.219
<i>Experiments with phoneme 4-grams</i>				
REF-97	0.316	0.310	0.313	0.307
WORD ₁ -97	0.266	0.258	0.255	0.257
WORD ₂ -97	0.211	0.196	0.207	0.195
<i>Experiments with phoneme 5-grams</i>				
REF-97	0.274	0.256	0.262	0.253
WORD ₁ -97	0.215	0.212	0.207	0.212
WORD ₂ -97	0.174	0.154	0.170	0.154

degraded using these queries except, marginally, in some cases of phoneme 3-grams, on the first baseline (WORD₁-97) and on the direct phoneme recognised (PHN-97) documents. Additional experiments were used to investigate the effect of boundary information in documents for phoneme n -gram retrieval. Results are shown in Table 5. In almost all cases of using bounded queries, both (full, bound) and (stopped, bound) retrieval was more effective on bounded documents. Again, we see that boundary information in the documents for phoneme n -gram retrieval has a greater influence than in the queries. For documents recognised as phonemes (PHN-97), retrieval was improved using bounded queries when phoneme 4-grams were used. However, word boundary information is not available in documents recognised directly as phonemes.

Experiments on the combined effect of word boundary and stopping in documents yielded the AP figures shown in Table 6. By comparing the results of Table 6 with Tables 3–5, we can observe that word boundaries had a stronger effect than stopping in the documents.

The effect of varying phoneme n -gram sizes was also observed. No particular n -gram size per-

formed well across the different types of queries on all versions of the document collections. For word-based documents translated to phoneme n -grams, the effects of word boundary information and stopping on different n -gram sizes were difficult to evaluate. Depending on the retrieval strategies, whether the queries were stopped or bounded, retrieval performance varied between 3-grams and 4-grams. Overall, 3-grams and 4-grams gave better retrieval results than 2-grams and 5-grams.

With the TREC-7 SDR collection, we combined evidence from phoneme n -grams of different lengths (Fuller et al., 1998). For each document, phoneme n -grams of different sizes were created, then all n -grams were concatenated to give one large document, which was then indexed in the usual way. Phoneme n -grams of queries were combined in a similar manner. Extensive experiments were used to test the combination of phoneme n -grams ranging from 2-grams to 5-grams. We varied the combination by combining two types of n -gram sizes to four types of n -gram sizes (i.e., adding 2-, 3-, 4- and 5-grams together). The best combination we found in our experiments was the combination of 3-grams and 4-grams, reported in Table 7. For the case of combination, bounded

Table 5

Average precision, for word-based documents converted to bounded phoneme n -grams, where n -grams were not formed across word boundaries, varying gram lengths and query types; TREC-7 data

Document set (bounded phoneme n -grams)	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 2-grams</i>				
REF-97	0.107	0.177	0.130	0.173
WORD ₁ -97	0.091	0.146	0.102	0.138
WORD ₂ -97	0.076	0.104	0.086	0.098
<i>Experiments with phoneme 3-grams</i>				
REF-97	0.293	0.327	0.302	0.320
WORD ₁ -97	0.253	0.284	0.267	0.282
WORD ₂ -97	0.198	0.226	0.208	0.223
<i>Experiments with phoneme 4-grams</i>				
REF-97	0.298	0.335	0.304	0.322
WORD ₁ -97	0.260	0.276	0.259	0.270
WORD ₂ -97	0.195	0.210	0.194	0.206
<i>Experiments with phoneme 5-grams</i>				
REF-97	0.259	0.301	0.254	0.295
WORD ₁ -97	0.213	0.247	0.210	0.237
WORD ₂ -97	0.157	0.193	0.157	0.185

Table 6

Average precision, for word-based documents which were stopped prior to conversion to bounded phoneme n -grams, where n -grams were not formed across word boundaries, varying gram lengths and query types; TREC-7 data

Document set + stopping (bounded phoneme n -grams)	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 2-grams</i>				
REF-97	0.105	0.159	0.147	0.181
WORD ₁ -97	0.089	0.117	0.121	0.155
WORD ₂ -97	0.078	0.101	0.097	0.125
<i>Experiments with phoneme 3-grams</i>				
REF-97	0.305	0.329	0.306	0.329
WORD ₁ -97	0.261	0.297	0.274	0.292
WORD ₂ -97	0.206	0.229	0.218	0.230
<i>Experiments with phoneme 4-grams</i>				
REF-97	0.299	0.328	0.300	0.328
WORD ₁ -97	0.259	0.274	0.262	0.272
WORD ₂ -97	0.195	0.209	0.192	0.207
<i>Experiments with phoneme 5-grams</i>				
REF-97	0.260	0.296	0.253	0.296
WORD ₁ -97	0.210	0.237	0.207	0.235
WORD ₂ -97	0.158	0.184	0.158	0.186

queries do not in general improve retrieval though there is a slight improvement for phoneme-recognised documents. The effects of stopping and word boundary information in documents had a

greater effect than in the queries. Retrieval was more effective on bounded documents using bounded queries and stopped queries were more effective on stopped documents.

Table 7

Average precision using combination of 3-grams and 4-grams, for different speech recognition processes and query types; TREC-7 data

	Document set	Query set			
		Full	Full, bound	Stopped	Stopped, bound
Unbounded	REF-97	0.341	0.329	0.316	0.320
	WORD ₁ -97	0.307	0.289	0.273	0.275
	WORD ₂ -97	0.231	0.205	0.211	0.208
	PHN-97	0.082	0.105	0.085	0.101
Unbounded + stopping	REF-97	0.322	0.326	0.329	0.329
	WORD ₁ -97	0.284	0.285	0.269	0.280
	WORD ₂ -97	0.222	0.216	0.225	0.219
Bounded	REF-97	0.315	0.338	0.312	0.332
	WORD ₁ -97	0.273	0.293	0.269	0.285
	WORD ₂ -97	0.209	0.227	0.207	0.220
Bounded + stopping	REF-97	0.314	0.334	0.312	0.335
	WORD ₁ -97	0.275	0.297	0.270	0.291
	WORD ₂ -97	0.217	0.225	0.211	0.225

In Section 3, we discussed better modelling of document length by changing the b parameter in the Okapi measure. Tentative experiments were conducted, varying the parameter between 0.5 and 1.0. We were unable to identify a suitable value, but values between 0.75 and 0.83 were given comparable retrieval performance to those reported here.

The experiments on the TREC-7 SDR data showed that retrieval was ineffective for direct phoneme recognised documents, compared to word-based documents. The effects of stopping and word boundary information in word-based documents were more significant than those in the queries. Bounding phoneme n -grams in the query set slightly reduced retrieval effectiveness. In terms of n -gram sizes, we found that phoneme n -grams of 3- and 4-grams retrieved better than 2- and 5-grams. TREC-7 SDR consists of only about 2800 documents and 23 queries.

The collection is too small to draw any firm conclusions from the experiments, but the results are indicative. We used the Wilcoxon significance test to explore the results further. These showed that stemming was consistently significantly superior to the combination of stemming and stopping, but most of the other differences were not statistically significant.

5.2. Results on TREC-6 data

Experiments using TREC-6 SDR data are shown in Tables 8–12. Documents represented as words retrieved using both full and stopped queries for both automatic (i.e., WORD-96) and manual (i.e., REF-96) recognised spoken documents are shown in Table 8. MRR results using phoneme n -gram representations to query the REF-96 and WORD-96 versions and the phoneme-recognised documents (PHN-96), are shown in Table 9, while the number of documents retrieved in the top 5 and 10 ranks are shown in Table 10.

Comparing the MRR of Tables 8 and 9, once again word-based retrieval was shown to be more effective than phoneme n -gram retrieval. An interesting result was observed by comparing effectiveness of word-based and phoneme n -gram retrieval on the automatically word-recognised version (i.e., WORD-96). Comparing the results for WORD-96 in Tables 8 and 9, the effectiveness of phoneme 3-grams was greater than that of word-based retrieval. Similar improvements were also observed using phoneme 3-grams and 4-grams based on stopped queries, as shown in Table 9. Although the number of times the relevant documents retrieved in the top 10 was reduced, the use of phoneme n -grams did improve the ranks of

Table 8

Retrieval effectiveness using documents represented as sequences of words for different speech recognition processes and query types; TREC-6 data

Document set	Query set	Mean reciprocal rank	No. relevant in the top	
			5	10
REF-96	Full, word	0.702	41	44
REF-96	Stopped, word	0.700	38	43
WORD-96	Full, word	0.558	36	40
WORD-96	Stopped, word	0.521	33	38

Table 9

Mean reciprocal rank for different speech recognition processes, gram lengths and query types; TREC-6 data

Document set	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 3-grams</i>				
REF-96	0.698	0.616	0.629	0.564
WORD-96	0.582	0.579	0.548	0.531
PHN-96	0.204	0.230	0.198	0.222
<i>Experiments with phoneme 4-grams</i>				
REF-96	0.687	0.569	0.581	0.552
WORD-96	0.538	0.486	0.536	0.488
PHN-96	0.236	0.205	0.178	0.187
<i>Experiments with phoneme 5-grams</i>				
REF-96	0.661	0.505	0.598	0.486
WORD-96	0.542	0.453	0.504	0.452
PHN-96	0.172	0.167	0.183	0.172

Table 10

Number of relevant documents retrieved in the top 5 and 10 for different speech recognition processes, gram lengths and query types^a

Document set	Query set			
	Full	Full, bound	Stopped	Stopped, bound
<i>Experiments with phoneme 3-grams</i>				
REF-96	(38,43)	(35,42)	(38,42)	(34,40)
WORD-96	(35,38)	(35,37)	(32,33)	(33,35)
PHN-96	(13,19)	(13,17)	(14,18)	(15,18)
<i>Experiments with phoneme 4-grams</i>				
REF-96	(37,43)	(32,38)	(36,40)	(31,37)
WORD-96	(32,38)	(29,35)	(30,32)	(31,32)
PHN-96	(15,18)	(12,17)	(13,17)	(12,15)
<i>Experiments with phoneme 5-grams</i>				
REF-96	(35,41)	(30,33)	(32,38)	(29,32)
WORD-96	(31,35)	(27,31)	(28,34)	(26,31)
PHN-96	(12,16)	(8,15)	(8,15)	(9,16)

^a Figures are in the form of (top 5, top 10). TREC-6 data.

Table 11

Retrieval effectiveness for phoneme n -grams (-phn- n) without crossing word boundaries (bound) and with query expansion using incorrect phoneme sequences of query terms (err) for different speech recognition processes and gram lengths; TREC-6 data

Document set	Query set	Mean reciprocal rank	No. relevant in the top	
			5	10
REF-96-phn-3	Stopped, phn-3, bound, err	0.423	29	36
REF-96-phn-4	Stopped, phn-4, bound, err	0.509	31	33
REF-96-phn-5	Stopped, phn-5, bound, err	0.456	28	32
WORD-96-phn-3	Stopped, phn-3, bound, err	0.463	28	31
WORD-96-phn-4	Stopped, phn-4, bound, err	0.479	28	29
WORD-96-phn-5	Stopped, phn-5, bound, err	0.442	25	29
PHN-96-phn-3	Stopped, phn-3, bound, err	0.153	10	13
PHN-96-phn-4	Stopped, phn-4, bound, err	0.197	11	13
PHN-96-phn-5	Stopped, phn-5, bound, err	0.148	7	10

Table 12

Retrieval effectiveness for phoneme n -grams (-phn- n) without crossing word boundaries (bound) and with query expansion using string edit distance of one to obtain neighbour terms (nbr), for different speech recognition processes and gram lengths; TREC-6 data

Document set	Query set	Mean reciprocal rank	No. relevant in the top	
			5	10
REF-96-phn-3	Stopped, phn-3, bound, nbr	0.014	1	3
REF-96-phn-4	Stopped, phn-4, bound, nbr	0.136	10	13
WORD-96-phn-3	Stopped, phn-3, bound, nbr	0.011	0	3
WORD-96-phn-4	Stopped, phn-4, bound, nbr	0.117	9	12
PHN-96-phn-3	Stopped, phn-3, bound, nbr	0.014	1	1
PHN-96-phn-4	Stopped, phn-4, bound, nbr	0.056	3	4

relevant documents retrieved. Analysis showed that by, translating words to phonemes, incorrectly recognised words were translated to phoneme sequences that were similar to the phoneme sequence of the correct word. This result supports those of Witbrock and Hauptmann (1997) in the Informedia project, where they combined word-based transcript and phoneme transcripts to improve retrieval performance. That is, phoneme n -gram retrieval has the capacity to correct some of the word-based recognised documents. Therefore, by taking the n -grams of the phonemes, relevant documents were retrieved at higher ranks.

Stopping did not lead to retrieval of more relevant documents or improve the ranks of relevant documents. We found that some relevant documents were being retrieved using phrasal linkage between stopwords and keywords.

Given that phoneme recognition has a high error rate, a large proportion of incorrectly recognised phonemes is inherent to SDR retrieval based on a phoneme recogniser. As discussed before, following our belief that transcription errors of the same word in the training collection would also occur in the test collection, the original query was augmented by a set of potential match terms obtained from the training collection. These terms were chosen based on their frequency of occurrence. Table 11 shows the retrieval performance using a combination of the initial query with the mistranscribed terms. Overall effectiveness significantly degrades with the use of these additional terms. This technique of query expansion is similar in concept to using confusion matrices and calculating the probability of recognition error between phonemes to determine the phoneme sequences to

use. Experiments by Wechsler et al. (1998) and Ng (1998) using dynamic programming techniques have shown that the confusion matrix approach can be effective.

Experiments with using neighbour terms to expand the query set, shown in Table 12, showed that the technique was a complete failure. It did not use the phoneme similarity (or otherwise) of the substituted characters. By using all the terms that differed in one character from the original query term, many incorrect or irrelevant terms were included. A more considered approach to selecting additional query terms – that takes into account the similarity of the phonemes involved – appears to be required.

Although these results consistently show word-based retrieval to be more effective than phoneme-based retrieval for SDR, limitations of the document and query collections may have affected these results. These collections are not an environment in which phoneme-based retrieval is likely to be favoured. In comparison to SDR in practice, recognition processes will be less reliable. Furthermore, it is possible that the collection is too small, so that the results may not generalise for larger collections. Translating words to phonemes using a pronunciation dictionary required unknown words from the collections to be manually added. This was a time-consuming task and would not occur in a practical system. A rule-based algorithm for translating words to phonemes would have been less costly and does not require additional work for unknown words.

6. Conclusions and future work

We have explored methods of SDR based on phoneme n -grams. We investigated the effects of word boundary information, stopping of queries, varying n -gram sizes, combination of phoneme n -grams and query expansion using incorrectly transcribed phoneme sequences and similar query terms. Two sets of spoken documents were used to compare retrieval performance. For each set, representations of the documents were generated from manual transcription, automatic WBR and automatic phoneme-based recognition; the word-

based documents were translated to phonemes using a pronunciation dictionary. Queries were similarly translated, using the pronunciation dictionary. These experiments confirm that phoneme-based retrieval is less effective than word-based retrieval, but is nonetheless reliable enough to be used in practice.

When translating queries into phoneme n -grams, we tested the effect of using or ignoring word boundaries. We found that word boundary information in the queries did not have much impact. Stopping the queries did not improve retrieval, in contrast to retrieval experiments on text-based collections. Word boundary and stopping the word-based documents prior to conversion to phoneme n -grams had a greater impact on retrieval effectiveness. Such processing cannot, of course, be applied to documents recognised directly as phonemes. Phrasal retrieval may be useful for retrieval when a phrase is formed by a stop-word and keyword and this was observed only in the TREC-6 SDR collection.

Varying phoneme n -gram of 3-grams and 4-grams led to some improvement in effectiveness in comparison to other n -gram sizes. The combination of 3-grams and 4-grams was also more effective than other combinations in the TREC-7 collections. Other combinations did not perform as well, apparently because shorter n -grams were retrieving more irrelevant documents and the longer n -grams were not finding some relevant documents.

It would be expected that transforming word sequences into phoneme sequences would lose information and thus result in deterioration of retrieval. This expectation was confirmed in experiments with the TREC-7 data, but there was no difference for the TREC-6 data. One explanation is that the task of finding a single item is much easier than that of finding all relevant documents, so that the TREC-7 experiment is more able to separate techniques based on different hypotheses. The TREC-6 experiments indicated that occasionally word errors are remedied if the word and its inaccurate replacement both map to a similar phoneme sequence.

Query expansion was tested using two different approaches. The first used information about

recognition errors derived from a training collection – the most frequent incorrect phoneme sequences generated for the query words was used to augment the original queries. The second approach was to use neighbour n -grams obtained using string edit distance. Each query term was augmented by a list of potential substitutes from a list of unique terms from the test collection. Both methods of query expansion degraded effectiveness, with the – admittedly simplistic – neighbour method in particular giving very poor results. Although the first approach performed better, in some individual cases improving effectiveness, it required information from the training collection that would not always be available.

The high quality recognition environment characterized by the TREC experiments is not always available. Degraded recognition due to noise and uncertainty as well as the lack of key words in the dictionary can lead to poor recognition. Proper names are often the most important terms in a query and are generally not available in pronunciation dictionaries. Each of these conditions needs further experimentation and particularly in experimental environments that reflect the difficult conditions that so often characterise the speech retrieval environment. The experiments here on phoneme n -grams investigated some of the retrieval techniques which might be useful in such conditions.

Acknowledgements

Thanks to Peter Schäuble, Eugene Munteanu and Martin Wechsler of the Swiss Institute of Technology (ETH) for providing phoneme transcriptions of the spoken document collections. Thanks also to Michael Fuller and Victor Poznanski. We are particularly grateful to the anonymous referees for their detailed, constructive comments. This work was supported by the Australian Research Council.

References

- Canvar, W.B., 1994. Using an n -gram-based document representation with a vector processing retrieval model. In: Proceedings of the Third Text REtrieval Conference (TREC-3). pp. 269–277.
- Fuller, M., Kaszkiel, M., Ng, C., Vines, P., Wilkinson, R., Zobel, J., 1997. MDS TREC-6 report. In: Voorhees, E.M., Harman, D.K. (Eds.), Proceedings of the Sixth Text REtrieval Conference (TREC-6). pp. 241–258.
- Fuller, M., Kaszkiel, M., Kim, D., Ng, C., Robertson, J., Wilkinson, R., Wu, M., Zobel, J., 1998. TREC 7 Ad Hoc, speech and interactive tracks at MDS/CSIRO. In: Voorhees, E.M., Harman, D.K. (Eds.), Proceedings of the Seventh Text REtrieval Conference (TREC-7). pp. 465–474.
- Garofolo, J., Voorhees, E., Auzanne, C., Stanford, V., Lund, B., 1998. 1998 TREC-7 spoken document retrieval track overview and results. In: Voorhees, E.M., Harman, D.K. (Eds.), Proceedings of the Seventh Text REtrieval Conference (TREC-7). pp. 79–90.
- Jones, G.J.F., Foote, J.T., Sparck Jones, K., Young, S.J., 1996. Retrieving spoken documents by combining multiple index sources. In: Proceedings of the 19th ACM-SIGIR Conference on Research and Development in Information Retrieval. pp. 30–38.
- Lee, K.F., 1989. Automatic Speech Recognition, The Development of the SPHINX System. Kluwer Academic Publishers, Dordrecht.
- Linguistic Data Consortium, 1996a. Continuous speech recognition corpus-V: 1996 broadcast news speech (CSR-V hub-4), CD-ROM. ldc@ldc.upenn.edu, Philadelphia, PA 19104-2608, USA.
- Linguistic Data Consortium, 1996b. DARPA continuous speech recognition corpus-IV: Radio broadcast news (CSRIV hub-4), CD-ROM. ldc@ldc.upenn.edu, Philadelphia, PA 19104-2608, USA.
- Linguistic Data Consortium, 1997. Continuous speech recognition corpus-VI: 1997 broadcast news speech (CSR-VI hub-4), CD-ROM. ldc@ldc.upenn.edu, Philadelphia, PA 19104-2608, USA.
- Mateev, B., Munteanu, E., Sheridan, P., Wechsler, M., Schäuble, P., 1997. ETH TREC-6: Routing, chinese, cross-language and spoken document retrieval. In: Voorhees, E.M., Harman, D.K. (Eds.), Proceedings of the Sixth Text REtrieval Conference (TREC-6). pp. 623–636.
- Ng, K., 1998. Towards robust methods for spoken document retrieval. In: Proceedings of International Conference on Spoken Language Processing. Vol. 3, pp. 939–942.
- Ng, K., Zue, V.W., 1997. Subword unit representations for spoken document retrieval. In: Proceedings of the European Conference on Speech Communications and Technology, EUROSPEECH, Rhodes, Greece, pp. 1607–1610.
- Ng, K., Zue, V., 1998. Phonetic recognition for spoken document retrieval. In: Proceedings of International Conference on Acoustic Speech and Signal Processing ICASSP. pp. 325–328.
- Rabiner, L., Juang, B.H., 1993. Fundamentals of Speech Recognition. Prentice-Hall, Englewood Cliffs, NJ.
- Salton, G., Buckley, C., 1988. Term-weighting approaches in automatic text retrieval. Information Processing and Management 24 (5), 513–523.

- Salton, G., McGill, M.J., 1983. *Introduction to Modern Information Retrieval*. McGraw-Hill, New York.
- Singhal, A., 1997. AT&T at TREC-6. In: Voorhees, E.M., Harman, D.K. (Eds.), *Proceedings of the Sixth Text REtrieval Conference (TREC-6)*. NIST Special Publication 500-240, pp. 215–226.
- Smeaton, A.F., Morony, M., Quinn, G., Scaife, R., 1998. Taiscéalái: information retrieval from an archive of spoken radio news. In: *Proceedings of the Second European Conference on Research and Advanced Technology for Digital Libraries (EuroDL)*.
- Voorhees, E., Harman, D., 1997. Overview of the Sixth Text REtrieval Conference. In: Voorhees, E.M., Harman, D.K. (Eds.), *Proceedings of the Sixth Text REtrieval Conference (TREC-6)*. pp. 1–24.
- Voorhees, E., Harman, D., 1998. Overview of the Seventh Text REtrieval Conference. In: Voorhees, E.M., Harman, D.K. (Eds.), *Proceedings of the Seventh Text REtrieval Conference (TREC-7)*. pp. 1–24.
- Walker, S., Robertson, S.E., Boughanem, M., Jones, G.J.F., Sparck Jones, K., 1997. Okapi at TREC-6 Automatic ad hoc, VLC, routing, filtering and QSDR. In: Voorhees, E.M., Harman, D.K. (Eds.), *Proceedings of the Sixth Text REtrieval Conference (TREC-6)*. pp. 125–136.
- Wechsler, M., Schäuble, P., 1995. Speech retrieval based on automatic indexing. In: *Workshop in Computing Science-MIRO*. Springer, Berlin.
- Wechsler, M., Munteanu, E., Schäuble, P., 1998. New techniques for open-vocabulary spoken document retrieval. In: *Proceedings of the 21st ACM-SIGIR International Conference on Research and Development in Information Retrieval*. pp. 20–27.
- Witbrock, M.J., Hauptmann, A.G., 1997. Using words and phonetic strings for efficient information retrieval from imperfectly transcribed spoken documents. In: *Proceedings of the Digital Library Conference*. Philadelphia, PA, USA, pp. 30–35.
- Witten, I., Moffat, A., Bell, T., 1994. *Managing Gigabytes: Compressing and Indexing Documents and Images*. Van Nostrand Reinhold, New York.
- Woodland, P.C., Leggetter, C.J., Odell, J.J., Valtchev, V., Young, S.J., 1995. The 1994 HTK large vocabulary speech recognition system. In: *Proceedings of International Conference on Acoustic Speech and Signal Processing ICASSP*.
- Young, S., Jansen, J., Odell, J., Ollason, D., Woodland, P., 1995. *The HTK Book*. Entropic Cambridge Research Laboratory.
- Zobel, J., Dart, P., 1996. Phonetics string matching: lessons from information retrieval. In: *Proceedings of the 19th ACM-SIGIR Conference on Research and Development in Information Retrieval*. pp. 166–172.