

An overview of audio information retrieval

Jonathan Foote*

Institute of Systems Science, National University of Singapore, Heng Mui Keng Terrace, Singapore 119597

Abstract. The problem of audio information retrieval is familiar to anyone who has returned from vacation to find an answering machine full of messages. While there is not yet an “AltaVista” for the audio data type, many workers are finding ways to automatically locate, index, and browse audio using recent advances in speech recognition and machine listening. This paper reviews the state of the art in audio information retrieval, and presents recent advances in automatic speech recognition, word spotting, speaker and music identification, and audio similarity with a view towards making audio less “opaque”. A special section addresses intelligent interfaces for navigating and browsing audio and multimedia documents, using automatically derived information to go beyond the tape recorder metaphor.

1 Introduction

The problem of audio information retrieval is familiar to anyone who has returned from vacation to find an answering machine full of messages. If you are not fortunate, you may have to listen to the entire tape to find the urgent one from your boss. Determining there is no such message is guaranteed to be time-consuming.

While there is not yet an “AltaVista” for the audio data type, many workers are finding ways to automatically locate, manipulate, skim, browse, and index audio using recent advances in speech recognition and machine listening. Such methods will be indispensable to cope with the burgeoning amounts of audio available both on the Internet and elsewhere. (One German site¹ advertises more than 3 days’ worth of available audio).

This paper will attempt to link a variety of research efforts across many fields, at a level of detail suitable for the non-specialist. A full list of references and URL pointers is given for the interested reader. The focus here will be on systems that can automatically extract information from an

audio signal, rather than approaches that depend on human-generated annotations or decisions. (See, for example, the music recommendation service of Firefly² as an example of the latter). Because of the wide variety of disciplines involved, from speech recognition, information retrieval, audio analysis, signal processing, psychoacoustics, and machine learning, it is difficult to include all pertinent work, for which I must apologize in advance³.

1.1 Technologies

A variety of technological advances can be used to make audio less “opaque”, that is, provide some insight into the content of an audio file, and perhaps ways of using it as other than a monolithic block of digital data. The available methods can be roughly divided into those that assume some speech content in the audio, and those that don’t. The general outline of this paper will follow this division; Sect. 2 will consider approaches based on automatic speech recognition (ASR), while Sect. 3 will consider more general audio analysis suitable for a wider range of audio such as music or sound effects (while not excluding speech). This parallels Smoliar *et al.*’s classification of multimedia retrieval into “expression” and “semantic” approaches [1]. In the first, objects are retrieved by some *physical* description, or an example from a similar medium. Thus, most content-based image retrieval systems require an image or sketch as a search query. The second approach requires some *semantic* analysis or knowledge, for example, annotations or picture captions. While semantic retrieval is perhaps the ideal, as it would hopefully retrieve what an intelligent human would, it requires either intensive human-assisted annotations or impractically sophisticated automatic content analysis (imagine the difficulty of automatically finding a photo of, say, Martin Luther King in a large database of uncaptioned images). Recent advances in ASR, however, have yielded audio retrieval systems that may approach the semantic ideal, as they can actually recognize (if not understand) the words uttered in an

² <http://www.firefly.com>

³ In particular, though there is much interesting work in the fields of psychoacoustics and auditory scene analysis, I must reluctantly consider them as beyond the scope of this overview.

* Present address: FX Palo Alto Laboratory, Inc., 3400 Hillview Ave, Building 4, Palo Alto, CA 94304, USA

¹ <http://www.icf.de/RIS/>

audio stream. Thus, it is entirely possible to automatically locate audio of someone *talking* about Martin Luther King, or indeed one of his speeches.

1.1.1 Information retrieval

Conventional information retrieval (IR) research has been mainly based on (computer-readable) text [2, 3], and is familiar to many through the popular web search engines such as Lycos or AltaVista. The classic IR problem is to locate desired text documents using a search query consisting of a number of keywords. Typically, matching documents are found by locating query keywords within them. If a document has a high number of query terms, it is regarded as being more “relevant” to the query than other documents with fewer or no query terms. Documents can then be ranked by relevance and presented to the user for further exploration, as the web search engines do. Though powerful IR algorithms are available for text, it is clear that for audio, or multimedia in general, common term-matching approaches are useless due to the simple lack of identifiable words (or comparable entities) in audio documents. The problem becomes even more open-ended when one considers audio, such as music, which may have no speech.

Even when a desired audio document can be located in a large archive, another problem to overcome is the linearity of audio files. To ensure that nothing important is missed, the entire audio file must be auditioned from start to finish, which takes significant time. In contrast, the transcription of a minute-long message is typically a paragraph of text, which may be scanned by eye in a matter of seconds. Even if there is a “fast-forward” button it is generally a hit-or-miss operation to find a desired section in a lengthy file. A typical interface treats audio as an undifferentiated stream: the “tape recorder” metaphor (with stop, play, rewind and fast-forward buttons) is ubiquitous. In contrast, most text-processing software has a “find” command that does simple string-match word searching to locate desired information in large files. Besides a scroll bar, many word processors can scroll by paragraph, page, function, or chapter. As above, an audio interface of similar flexibility is impossible unless suitable indexing entities, analogous to “words” or “pages”, can be located in the audio. Section 4 discusses approaches to precisely this problem.

1.1.2 Automatic speech recognition

ASR is a technology rapidly coming out of the research laboratories into everyday use. While there have been decades of hard effort on the task⁴, recent advances both in search algorithms and commonly available computing power are rapidly making ASR practical. A perfect ASR system that could quickly transcribe spoken audio documents would be an ideal solution to most audio indexing and retrieval tasks (at least for speech). Such a system would essentially reduce the audio retrieval problem to the straightforward text retrieval problem described above.

⁴ Aleksandr Solzhenitsyn’s novel *The First Circle* (1968) describes speech recognition research in a Stalin-era Soviet labor camp.

2 Speech recognition

Practically all ASR systems in use are based on hidden Markov models (HMMs) [4]. A hidden Markov model is a statistical representation of a speech event like a word; model parameters are typically trained on a large corpus of labeled speech data. Given a trained set of HMMs, there exists an efficient algorithm for finding the most likely model sequence (the recognized words), given unknown speech data. This approach has proved successful not only for large-vocabulary recognition systems, but for “keyword-spotting” systems where the location of only a few words or phrases is desired. Typically, this is done by training HMMs for both the desired keywords and a “filler” model that attempts to match everything not a keyword [5–7]. Such systems can be both accurate and computationally far less expensive than a large-vocabulary recognition system, while being flexible enough to handle unconstrained real-world speech [8].

Large-vocabulary recognition systems, in contrast, typically use a “sub-word” approach: rather than building an explicit HMM for every one of the tens of thousands of words in the vocabulary, a few hundred sub-word models are used, typically phonetically based. Given a phonetic dictionary, the appropriate sub-word models can be concatenated to form a word model. For example, the word “right” could be constructed by concatenating the three sub-word models for the phones “R AY T”. In addition, a large-vocabulary system requires a statistical “language model” that defines likely word combinations. For example, in English, the word bigram “of the” would be far more likely than the bigram “oaf the”. The language model can thus constrain the recognizer to word combinations that are more likely, and hence more likely correct. To be useful, language models must be trained on example text (typically millions of words) from a similar domain, which may be more practical for certain domains such as news, where large corpora of newspaper text is available, than others such as conversational speech.

A particular advantage of using ASR for audio information retrieval is that – unlike dictation or voice-command tasks – most or all of the desired audio is already present, and thus ASR can be performed off-line rather than in real time. A disadvantage is that there can be orders of magnitude more data to actually recognize, and an ASR system sufficiently fast for dictation may be far too slow to use on hours of audio at search time.

The primary drawback of ASR systems is their limited accuracy. Though the best continuous speech recognition systems can achieve better than 90% word accuracy on carefully recorded, limited-domain tasks such as the Wall Street Journal corpus [9], similar systems achieve little better than 50% or 60% word accuracy on real-world tasks such as telephone conversations or news broadcasts [10–13]. Even though word accuracy rates might appear unusably bad, the results of ASR transcription can still be surprisingly helpful for information retrieval. The reason for this is as follows: even with a 50% word error rate, the chance that a recognizer will miss every word in a three-word query (and thus a desired document) is $(\frac{1}{2})^3 = 0.125$, and will be even lower assuming more relevant documents have more keyword occurrences. A similar effect has been termed “semantic co-occurrence filtering” [14]. Other approaches based on key-

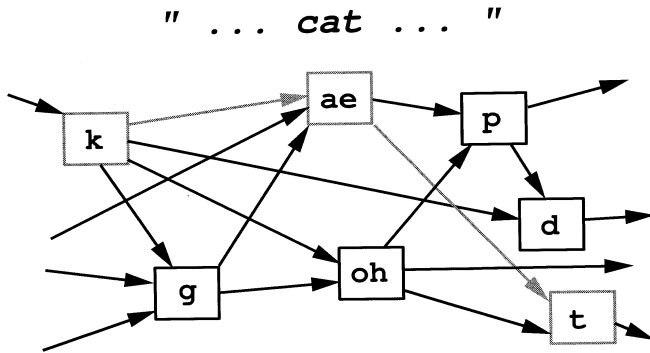


Fig. 1. Example phone lattice for utterance “cat”

word spotting or sub-word units suffer similar problems. A particular drawback of sub-word- or phonetic-based systems is the extreme difficulty in correctly recognizing sub-word units such as phones. Information retrieval based on such systems must be robust to recognition errors.

2.1 Keyword spotting

Automatically detecting words or phrases in unconstrained speech is usually termed “word spotting”; this technology is the foundation of several audio-indexing efforts from a number of groups. A popular test is the SWITCHBOARD corpus⁵, which contains recordings of spontaneous telephone conversations. Because recording subjects were asked to converse about certain topics (e.g., pets, weather, gun control), several research groups attempted to automatically determine the topic of each conversation. Workers at BBN have used both large-vocabulary ASR and keyword spotting to investigate topic identification [15]. Their approach allowed for a variable number of keywords, determined automatically; they report results from using from 100–2500 keywords. Their system results in 88.3% topic identification accuracy (from among ten topics), given both sides of a SWITCHBOARD conversation are used.

A group at Enigma used statistical models of keyword co-occurrence to discriminate between radio news topics, such as sports or weather, with very good accuracy [16, 17]. The Video Mail Retrieval (VMR) group⁶ at Cambridge University (which included the author) has investigated using 35 pre-selected keywords for audio information retrieval [18]. They report retrieval accuracy near 90% of that obtainable from a perfect wordspotter. Kate Knill, while at the Cambridge University Engineering Department, has been investigating fast keyword spotting for hand-held PDAs [19].

2.2 Sub-word indexing

Large-vocabulary ASR for audio indexing suffers, as we have seen, from several drawbacks: if a word is not present in the phonetic dictionary, it will not be recognized. Also, a language model must be used, and finding sufficient example text may not be possible. Thirdly, large-vocabulary ASR

S: [um] IN DOCUMENTS, WHICH MIGHT BE video mail messages say, [pause]
 R: AND IN DOCUMENT WHICH MIGHT BE THE DEAL MINISTRY IS A

S: [um] AND THAT will tell you WHICH ONES are relevant. You COULD compute
 R: HOME AND THAT WOULD REVIEW WHICH ONES OF REVENUE COULD INCLUDE

S: SOME relevance score depending ON HOW CERTAIN you are of FINDING THEM
 R: SOME OF ITS BOARD RULING ON HOW CERTAIN OF FINDINGS ON THEM

S: and so on. [pause] [um] [pause] ANOTHER APPROACH WOULD BE RATHER THAN
 R: ANOTHER APPROACH WOULD BE RATHER THAN

S: have A SET OF [pause] distinct DOCUMENTS you HAVE might have ONE
 R: A SET OF THIS THING DOCUMENT AND HAVE MUCH OF ONE

S: ongoing SOURCE OF [pause] SOUND SUCH AS [um] [pause] say a
 R: ON GROWING SOURCE OF SOUND SUCH AS A HOME SO

S: news bulletin, [pause] BBC nine o'clock news or something on the
 R: NEWSPRINT IN THE RECENT RESULTS AND

S: RADIO. [loud_breath] [um] [pause] AND THEN and YOU COULD [pause]
 R: RADIO FIRM AND THEN YOU COULD

S: listen to THE WHOLE OF IT [um] [pause] AND if you suddenly spotted
 R: LOSE INTO THE WHOLE OF ITS OWN AND A FEW SOUTHERN IS OFF TWO

S: [pause] [um] THE KEY WORD IN THAT [pause] [um] [pause] then then you
 R: AND THE KEY WORD IN THAT FOUND THAT THE END TO

S: DECIDE that THE SURROUNDING AREA OF text, [pause] [um] I mean MAYBE
 R: DECIDE OF THE SURROUNDING AREA OF TAX ON LINE MAY BE

Fig. 2. Automatic transcription of spontaneous speech. S: spoken words (lower case are misrecognized words) R: recognized speech

may be very expensive in terms of computation and storage (though recent advances in search algorithms are making this much less of a concern). Though this may be acceptable for typical speech recognition applications such as dictation, it is clearly unacceptable to incur several hours of computation when searching an audio corpus of similar length. To avoid some of these drawbacks, several alternatives to large-vocabulary ASR have been pursued. A common feature is the use of sub-word-indexing units, such as phones or phone clusters. Typically, these are smaller than words, hence there are fewer possible units, dramatically reducing the search space. An unfortunate drawback is that as units get smaller the recognition accuracy will typically decrease as well. Saving the cost of building a detailed language model will unfortunately impact recognition accuracy.

A group at ETH Zürich⁷, has reported speech retrieval using a large set of automatically chosen vowel-consonant-vowel sub-word units as indexes [20, 21]. A recent paper from MIT’s Spoken Language Systems Group⁸ investigates a number of possible sub-word units, from sequences of phones or broader phone classes, automatically discovered phone “multigrams”, to multiple-phone syllable units [22]. Retrieval performance is again related to recognition accuracy; results show no clear benefit of any particular unit, though some might be better than others for certain domains.

Another promising approach is “lattice-based” word spotting. A lattice is a compact representation of multiple best hypothesis generated by a phone or word recognition system. If the phone lattice is generated before need, it can then be searched extremely rapidly to find phone strings corresponding to desired query words. James reports the lattice scanner approach working about 1000 times real time; in other words an hour of audio may be searched in 3.6 s [23]. Figure 1 shows an example lattice for the word “cat”. Keep-

⁵ <http://morph ldc.upenn.edu/ldc/news/newsletter/v1.2/Switch.html>

⁶ <http://svr-www.eng.cam.ac.uk/Research/Projects/vmr/>

⁷ <http://www-ir.inf.ethz.ch/>

⁸ <http://sls-www.lcs.mit.edu/>

ing multiple hypotheses makes the system much more robust to recognition errors. For example, in the figure, even though the phone “T” was not the first choice, it is still in the lattice, thus, the phone string K AE T can still be found. If the lattice contains too many hypotheses, however, recognition accuracy will suffer from too many false alarms, as words that were not uttered can be found in a deep enough lattice. For example, the word “goat” (G OH T) can be found in the lattice of Fig. 1, even though the uttered word was presumably “cat”. The Video Mail Retrieval group at Cambridge University has successfully used phone lattices for open-vocabulary voice message retrieval [24, 25]; other workers have investigated lattice-based keyword spotting on the TIMIT corpus [26].

2.3 Large-vocabulary ASR

Several research groups have used large-vocabulary recognition for audio characterization. Dragon Systems⁹ were perhaps the first to use large-vocabulary recognition for speech document characterization; in 1993, they investigated topic identification on the Switchboard corpus [27]. A large effort is the Informedia project¹⁰ at Carnegie Mellon University. There, a combination of verbatim text transcriptions and large-vocabulary recognizer output are used for search indexes of video data such as television broadcasts [28]. Section 4.5 describes the the Informedia project in more detail. The Video Mail Retrieval group at Cambridge University has investigated combining small vocabulary keyword spotting (35 keywords) with large-vocabulary ASR for voice mail retrieval. Their findings suggest that a combination of the two approaches can be superior to either alone, as the keyword spotting allows detection of words not in the large-vocabulary lexicon [10]. Figure 2 shows the results of using a 5 000-word large-vocabulary recognizer on a spontaneously-spoken video mail message. Note that the language model is clearly inappropriate for the domain. Section 4.4 gives more information about the VMR project.

2.4 Speaker identification

A somewhat easier proposition than speech recognition is to simply identify differences between voices, rather than determine what the voices are saying. Technology to do this, termed *speaker identification* or *speaker ID*, can be very accurate in the right circumstances [29]. The applications to audio indexing are immediate: work at Xerox PARC allowed recorded meetings to be segmented and analyzed by speaker. A timeline display showed when particular speakers were talking during the meeting, as well as random access to play back a desired portion of the recording [30, 31]. Figure 3 shows how various speakers are displayed versus time [32]. Note how non-speech audio events like silence and applause may be located as well, which may be important cues for automatic segmentation.

A novel application of speaker ID was to align multi-hour recordings of the US House of Representatives with

the text transcription published as the Congressional Record. Because the published text has been corrected and amended, it is often of interest to determine what a given lawmaker actually said. By using a dynamic-programming approach to align the audio record with the (differing) printed version, different speaker turns can be accurately located, adding much to the utility of the audio transcript [33].

Segmenting multimedia streams is a promising area for speaker identification. If the ID technique works well enough on a sub-second time scale, it could be used to detect speaker changes in the soundtrack of a video or multimedia source, allowing it to be indexed and some sort of structure (such as dialogs) to be determined. Wyse and Smoliar report a “novelty measure” based on the cepstral difference between a short (0.75 s) and long (3 s) analysis window [34]. Differences are only compared in similar regions of the feature space to prevent intra-speaker variation (such as different vowels) from generating a high novelty score. When the difference exceeds a threshold, this can signal a new speaker or a significant change in the audio stream. These could be the audio equivalents of scene or camera changes, cuts, fades and wipes. It should be possible to fuse data intelligently extracted from both the video and audio streams, yielding more complete and robust information (about key frames, for example) than is available from either mode alone. A group at the University of Mannheim has been looking at the automatic analysis of film and video soundtracks. They have presented a system for automatic film genre classification based on low-level video and long-term audio frequency and amplitude characteristics [35]. In more recent work, they have attempted to automate violence detection in movie soundtracks by recognizing shots, cries and explosions. This was done by matching characteristics such as amplitude, frequency, and pitch [36].

3 Music and audio analysis

While ASR can give valuable clues to the content of speech, the universe of possible audio is, of course, much wider than speech alone. Music is a large, and extremely variable (and hence challenging) audio class¹¹. Considering the range of sounds that people might want to archive or classify, from the gamut of musical genres through sound effects to animal cries to synthesizer samples, and it is clear that speech-based methods alone are woefully inadequate for general audio discovery. Complicating the task is that any of the above can and will occur in combination; e.g., a narrator speaks over music or natural sounds, a translated version of a speech may be mixed over the original, and widely different sounds may occur sequentially in the same stream.

3.1 Music discrimination

A general problem in audio analysis is to simply discriminate speech from non-vocal music or other sounds. This has immediate applications for speech recognition: in general,

⁹ <http://www.dragonsys.com/>

¹⁰ <http://informedia.cs.cmu.edu/>

¹¹ Few can even agree on a satisfactory definition of “music”; consider John Cage’s controversial composition *4’ 33”* consisting of 4 min and 33 s of silent performance.

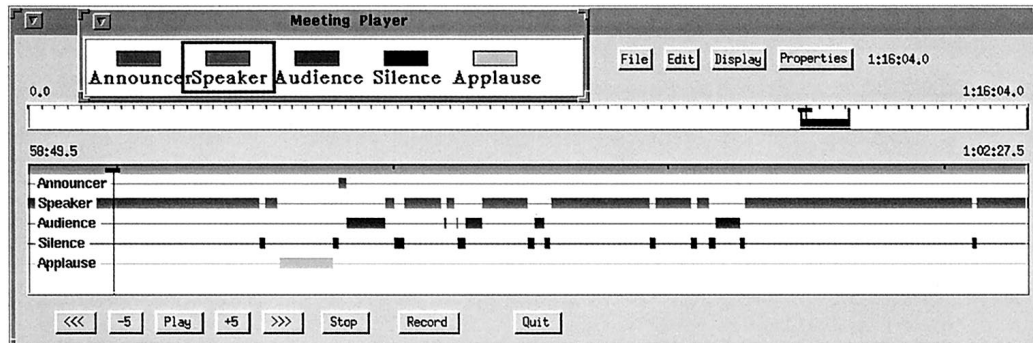


Fig. 3. PARC audio browser with speaker segmentation

there is no guarantee that a given multimedia audio source contains speech, and it is important not to waste valuable resources attempting to perform speech recognition on music, silence, or other non-speech audio. A straightforward approach to discriminating music from speech was presented by John Saunders [37] based on the statistics of the energy contour and the zero-crossing rate. Saunders reports a 98% classification accuracy on commercial radio broadcasts. Eric Scheirer of the MIT Media Lab and Malcolm Slaney of Interval Research report a speech/music discriminator based on various combinations of 13 features, such as 4-Hz modulation energy, “spectral centroid”, and zero-crossing rate [38]. Several classification strategies, including Gaussian mixture models and K-nearest-neighbor classifiers were evaluated. They report a 1.4% error rate on a large and diverse collection of FM radio broadcasts, when looking at relatively long-term (2.4 s) windows. Michelle Spina and Victor Zue of MIT’s Spoken Language Systems Group report experiments on 2 h of NPR radio news [39]. Using a maximum *a posteriori* approach on mel-frequency cepstral coefficients, the authors report an 80.9% classification accuracy into seven classes, including clean, telephone, and noisy speech, silence, music, and speech plus music. They include a plot of classification performance versus window size, and report that a window size near 0.5 s resulted in optimal classification. Finally, on a much less extensive data set, consisting of short samples (0.5–5 s) of music or speech, the author reports excellent discrimination accuracy using a classifier based on the statistics derived from an MMI-based vector quantizer [40], albeit on a much smaller sample of audio clips.

4 Advanced audio interfaces

Arguably more useful than classifying audio would be a system that would let users find audio items of interest, from either a large database of recordings or within a long recording itself. This section presents approaches to audio information retrieval¹². Many use the “ranked-list” user interface familiar from Web search engines. Even when a potentially relevant audio item is found, there are better ways of browsing the audio than the “tape recorder” metaphor, whose controls are typically limited to stop, play, and fast-forward and rewind.

¹² Though an interesting topic, we will not consider mere playback interfaces that do not extract information from the audio.

4.1 SpeechSkimmer

Barry Arons’s SpeechSkimmer is an excellent example of pushing audio interaction beyond the tape recorder metaphor [41]. It has long been known that humans can understand speech much quicker than the rate it is typically spoken. Using time-compression processing that alters the audio playback rate without changing the pitch, SpeechSkimmer users can audition spoken documents at several times real time, as well as backwards. Pauses can be identified (and removed) by detecting speech, using an adaptive algorithm said to be robust to background noise. In addition, the audio may be segmented by analyzing the speaker’s pitch to find cues associated with new topics, or by using speaker identification to locate conversational turns. SpeechSkimmer’s hierarchy of summarization techniques enable, in Aron’s words, a “fish-ear” view of an audio document.

4.2 Audio retrieval by content

Given the proliferation of audio databases on the Internet and elsewhere (some commercial sound effects libraries contain as many as 100 CDs), there is interest in doing for sound what Web search engines do for text. This requires some measure of audio similarity, which is a complicated and subjective matter. Measures of text similarity can be simple as counting the number of words in common. Most approaches to general audio retrieval take a perceptual approach, using measures derived from the audio that reflect perceptual characteristics such as brightness or loudness. A group at Technische Universität Berlin has used a neural net to map a sound clip to a text description, which could be inverted to find sounds by description [42]. An obvious drawback here is the subjective nature of audio descriptions. Sounds that a particular listener describes as “sharp” may be quite different from another’s. A later paper used a self-organizing map (SOM) on perceptually derived spectral features[43]. The net effect was to organize a set of 100 sample-synthesizer sounds into a 2D matrix such that similar sounds were closer and more disparate sounds were found further away on the grid.

Work by a group at Muscle Fish LLC¹³ has resulted in a compelling audio retrieval-by-similarity demonstration¹⁴

¹³ <http://www.musclefish.com>

¹⁴ <http://www.musclefish.com/cbrdemo.html>

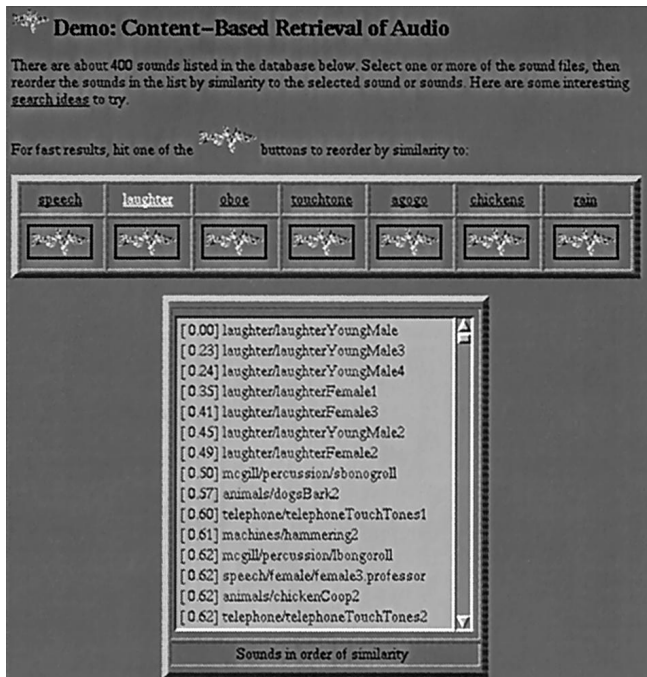


Fig. 4. Example Muscle Fish audio retrieval results, searching for audio similar to “laughter” (from <http://www.musclefish.com>)

(Fig. 4). Muscle Fish’s approach is to analyze sound files for a specific set of psychoacoustic features. This results in a vector of attributes that include loudness, pitch, bandwidth and harmonicity [44]. Given enough training samples, a Gaussian classifier can be constructed, or for retrieval, a covariance-weighted Euclidean (Mahalanobis) distance is used as a measure of similarity. For retrieval, the distance is computed between a given sound example and all other sound examples (about 400 in the demonstration). Sounds are ranked by distance, with the closer ones being more similar.

Recent work by the author, using an entirely different approach, has resulted in a similar retrieval framework. Here, distance measures are computed between histograms derived from a discriminatively trained vector quantizer. Audio is first parametrized into a spectral representation (mel-frequency cepstral coefficients). A learning algorithm then constructs a quantization tree that attempts to put samples from different training classes into different bins. A histogram of an audio file can be made by looking at the relative frequencies of samples in each quantization bin. If histograms are considered vectors, then simple Euclidean or cosine measures can be used to determine similarity between them, and hence their source audio [40]. This approach has been used for speaker identification [45, 46], as well as music and audio retrieval [47].

The Muscle Fish and author’s retrieval performance have been compared on the same audio corpus of 409 short sound files [47]. When retrieving simple, mono-component sounds like isolated instrument samples, the Muscle Fish retrieved sounds of similar timbre but varying pitch, while the author’s approach retrieved sounds of similar pitch from instruments of varying timbre. This demonstrates the subjective nature of audio similarity: it is not clear which criterion is more

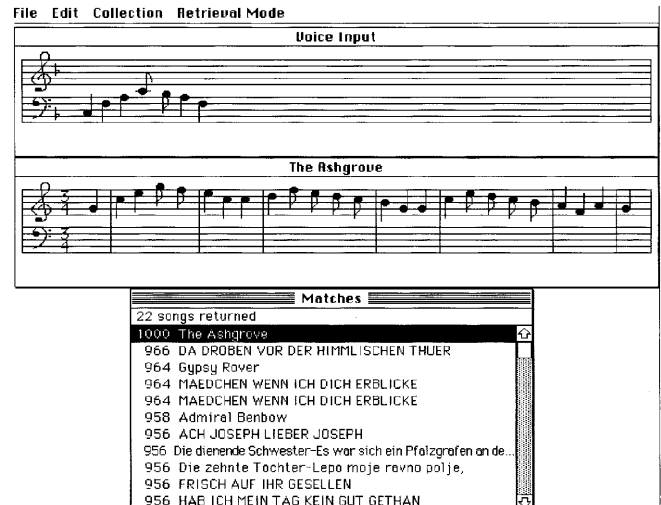


Fig. 5. Interface of Waikato tune retrieval application

important – the appropriate choice is probably application-dependent. (Similarly for image retrieval, the relative importance of shape vs. color is not clear.)

4.3 Music and MIDI retrieval

While information retrieval for text relies on simple text queries, the structure of a query for sound or music is not so obvious. Though textual descriptions can be assigned to sounds, they are not always obvious or indeed well-defined. The content-based retrieval applications of the previous section have avoided the problem somewhat by using audio examples as the query, in effect saying “look for things that sound like *this*”. Some recent work improves on this by allowing the user to sing or whistle a desired tune. Pitch extraction algorithms convert this to a note-like representation, which can be used to query a database of music. Unfortunately, extracting score-like attributes for anything but the simplest pieces has proved extremely difficult. Researchers in the area have finessed this problem by using archives of MIDI (Musical Instrument Digital Interface) files, which are score-like representations of music intended for musical synthesizers or sequencers. Given a melodic query, then, the MIDI files can be searched for similar melodies. Researchers at Cornell report surprisingly effective retrieval using query melodies that have been quantized to three levels, depending on whether each note was higher, lower, or similar pitch as the previous one [48]. Besides simplifying the pitch extraction, this allows for less-than-expert singing ability on the part of the user! A similar, if more advanced, application has been developed at the University of Waikato in New Zealand [49]. The Waikato system uses flexible string-matching algorithms to locate similar melodies located anywhere in a piece. Figure 5 (after McNab et al. [50]) shows the recognized tune sung by the user on the top staff, and the retrieved musical tunes¹⁵.

¹⁵ The tune retrieval application (for Macintosh computers) may be downloaded at <ftp://ftp.cs.waikato.ac.nz/pub/mac/MRv2.0.1.sea.hqx>

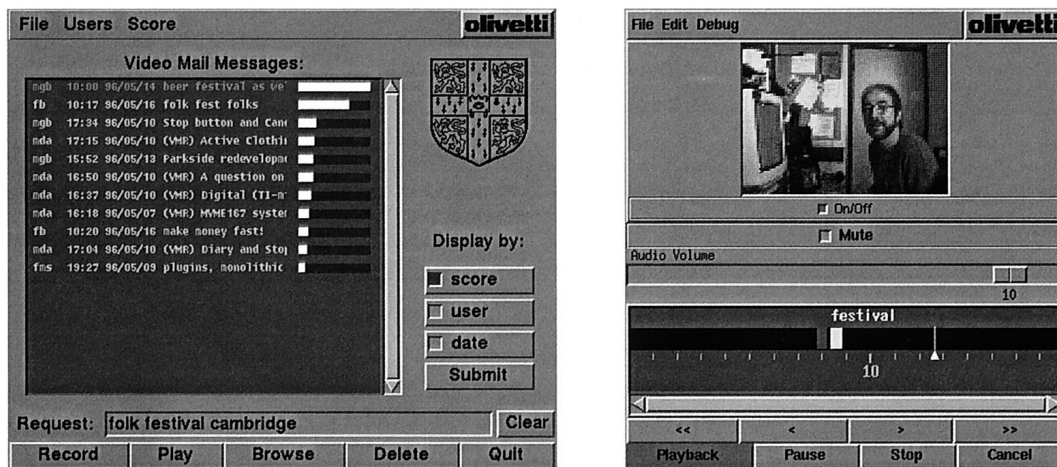


Fig. 6. Video mail retrieval user interface. Search engine (left) and mail browser (right)

4.4 Video mail retrieval

As mentioned in Sect. 2.1, the VMR group at Cambridge University has used open-vocabulary word spotting (based on phone lattices) to retrieve messages from an archive of video mail. In operation, the user types in a text search query, exactly as in a Web search engine. The lattice-based wordspotter finds instances of each query word spoken in each message (if any). A matching score for each message is computed by the retrieval engine, and the interface then displays a list of messages ranked by score. Scores are represented by bar graphs, as in the left image of Fig. 6 (after Brown et al. [51]).

After the ranked list of messages is displayed, the user must still investigate the listed messages to find the relevant one(s). The video mail browser (the right image for Fig. 6) is an attempt to represent a dynamic time-varying process (the audio/video stream) by a static image that can be taken in at a glance. A message is represented as a horizontal timeline, and events are displayed graphically along it. Putative keyword hits are displayed along the timeline, as in Fig. 6. The timeline is the black bar; the scale indicates time in seconds. When pointed at with the mouse, keyword names are highlighted in white (so it may be read in the presence of many other keyword hits). Clicking on the desired time in the time bar starts message playback at that time; this lets the user selectively play regions of interest, rather than the entire message.

4.5 Informedia

The Informedia project¹⁶ at Carnegie Mellon University is an impressive combination of video and audio analysis and text-based information retrieval techniques [12, 28, 52]. Given a video broadcast and a text transcription (from production notes; closed captions are insufficiently accurate), an HMM-based approach can accurately time-align the spoken words with the transcript. The Informedia project uses this to provide an abstract of the video by extracting both key frames from the video and important words from the text. Word

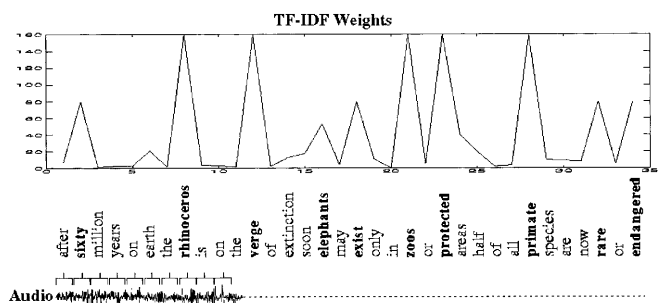


Fig. 7. Informedia term weighting approach

importance is determined from *tf/idf weights*, which stands for “term frequency/inverse document frequency”. Basically, terms (words) which appear frequently in a single document have a high “term frequency” (tf) and are thus more important. Similarly, terms which appear across many documents in the collection are less important, and have a lower “document frequency” (df). The product of tf and the inverse of df gives a measure of word importance, as shown in Fig. 7. Extracting words with high *tf/idf* from the audio stream and concatenating them yields a digest that (hopefully) contains the gist of the audio, and which can be auditioned extremely rapidly. Also conventional text-based retrieval techniques can be used to locate desired broadcasts or relevant portions within them.

5 Future directions

This paper has described some novel and powerful ways of extracting information from audio data. Clearly there is a lot of room for further research. One promising direction is the inclusion of Radio Broadcast News into the DARPA speech recognition evaluation effort (CSRIV Hub-4). This will encourage speech recognition researchers to tackle the difficult problem of general audio, rather than just clean speech, as has been the case with most research to date. Audio information retrieval is also now a sub-task of the TREC (Text Retrieval Conference), encouraging those in the text retrieval community to consider audio as well. Some interesting work,

¹⁶ <http://http://www.informedia.cs.cmu.edu/>

somewhat beyond the scope of this review, is being done in automatic translation of audio and text, as well as cross-language information retrieval. A particularly exciting area is how to combine information from various modes, such as audio and video. Efforts like that at the Informedia project hint at the power that fusing low-level information from different media can bring to the general problems of multimedia recognition, segmentation, and retrieval. Research at Ryukoku University in Japan, which deserves a wider audience, has pioneered combining video analysis and word spotting to segment articles from TV news [53]. Finally, issues of scale will need to be addressed: how well do various methods cope with vast numbers of large documents? As multimedia archives proliferate on the WWW and elsewhere, technology like that presented here will become indispensable to locate, retrieve, and browse audio and multimedia information.

References

- Smoliar SW, Baker JD, Nakayama T, Wilcox L (1996) Multimedia search: An authoring perspective. In: Proceedings of the First International Workshop on Image Databases and Multimedia Search, IAPR, August 1996, pp 1–8
- van Rijsbergen CJ (1979) Information Retrieval, 2nd edition. Butterworths, London
- Frakes W, Baeza-Yates R (1992) Information Retrieval: data structures and algorithms. Prentice Hall, Englewood Cliffs, N.J.
- Rabiner LR (1989) A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. IEEE* 77(2):257–286
- Rose RC, Paul DB (1990) A hidden Markov model based keyword recognition system. In: *Proc. ICASSP 90*, IEEE CS Press, Piscataway, N.J., pp 129–132
- Wilcox LD, Bush MA (1992) Training and search algorithms for an interactive wordspotting system. In: *Proc. ICASSP 92*, vol. 2, IEEE CS Press, Piscataway, N.J., pp 97–100
- Jeanrenaud P, Ng K, Siu M, Rohlicek J, Gish H (1993) Phonetic-based word spotter: Various configurations and application to event spotting. In: *Proc. Eurospeech 93*, Berlin, Germany, ESCA, pp 2145–2148
- Rose RC (1991) Techniques for information retrieval from speech messages. *Lincoln Lab J* 4(1):45–60
- Pallett D et al. (1995) Benchmark tests for the ARPA spoken language program. In: *Proc. ARPA SLS Technology Workshop*, January 1995
- Jones GJF, Foote JT, Spärck Jones K, Young SJ (1996) Robust talker-independent audio document retrieval. In: *Proc. ICASSP 96*, volume 1, April 1996, Atlanta, Ga. IEEE CS Press, Piscataway, N.J., pp 311–314
- Jeanrenaud P, Eide E, Chaudhari U, McDonough J, Ng K, Siu M, Gish H (1995) Reducing word error rate on conversational speech from the Switchboard corpus. In: *Proc. ICASSP 95*, May 1995, Detroit, Mich. IEEE CS Press, Piscataway, N.J., pp 53–56
- Hauptmann A, Witbrock M (1997) Informedia: News-on-demand – multimedia information acquisition and retrieval. In: Maybury MT (ed) *Intelligent Multimedia Information Retrieval*, chapter 10. MIT Press, Cambridge, Mass., pp 215–240 (available on the internet: <http://www.cs.cmu.edu/afs/cs/user/alex/www/>)
- Kubala F, Jin H, Matsoukas S, Nguyen L, Schwartz R (1997) Broadcast news transcription. In: *Proc. ICASSP 97*, volume 1, April 1997, IEEE CS Press, Piscataway, N.J., pp 203–206
- Kupiec J, Kimber D, Balasubramanian V (1994) Speech-based retrieval using semantic co-occurrence filtering. In: *Proc. ARPA Human Language Technology Workshop*, March 1994, Plainsboro, N.J. ARPA
- McDonough J, Ng K, Jeanrenaud P, Gish H, Rohlicek JR (1994) Approaches to topic identification on the Switchboard corpus. In: *Proc. ICASSP 94*, volume 1, Adelaide, Australia. IEEE CS Press, Piscataway, N.J., pp 385–388
- Wright JH, Carey MJ, Parris ES (1995) Improved topic spotting through statistical modelling of keyword dependencies. In: *Proc. ICASSP 95*, May 1995, Detroit, Mich. IEEE CS Press, Piscataway, N.J., pp 313–316
- Wright JH, Carey MJ, Parris ES (1995) Topic discrimination using higher-order statistical models of spotted keywords. *Comput Speech Lang* 9(4):381–405
- Jones GJF, Foote JT, Spärck Jones K, Young SJ (1995) Video Mail Retrieval: the effect of word spotting accuracy on precision. In: *Proc. ICASSP 95*, volume 1. IEEE CS Press, Piscataway, N.J., pp 309–312
- Knill KM, Young SJ (1994) Speaker dependent keyword spotting for hand-held devices. Technical Report 193. Engineering Department, Cambridge University, Cambridge, UK
- Schäuble P, Wechsler M (1995) First experiences with a system for content based retrieval of information from speech recordings. In: *IJCAI Workshop: Intelligent Multimedia Information Retrieval*, August 1995 (available on the Internet: <ftp://ftp.inf.ethz.ch/pub/publications/papers/is/ir/ijcai95.ps.gz>)
- Wechsler M, Schäuble P (1995) Speech retrieval based on automatic indexing. In: Rijsbergen CJ van (ed) *Proceedings of the MIRO Workshop*, September 1995, University of Glasgow, Glasgow, UK
- Ng K, Zue V (1997) Subword unit representations for spoken document retrieval. In: *Proc. Eurospeech 97*. ESCA (available on the Internet: <http://www.sls.lcs.mit.edu/~kng/papers/sir-eurospeech97.ps>)
- James DA, Young SJ (1994) A fast lattice-based approach to vocabulary independent wordspotting. In: *Proc. ICASSP 94*, volume 1, Adelaide, Australia. IEEE CS Press, Piscataway, N.J., pp 377–380
- Foote JT, Jones GJF, Spärck Jones K, Young SJ (1997) Unconstrained keyword spotting using phone lattices. *Comput Speech Lang* (in press)
- Young SJ, Brown MG, Foote JT, Jones GJF, Spärck Jones K (1997) Acoustic indexing for multimedia retrieval and browsing. In: *Proc. ICASSP 97*, volume 1, April 1997, Munich, Germany. IEEE CS Press, Piscataway, N.J., pp 199–202
- Gelin P, Wellekens C (1996) Keyword spotting for video soundtrack indexing. In: *Proc. ICASSP 96*, volume 1, April 1996, Atlanta, Ga. IEEE CS Press, Piscataway, N.J., pp 299–302
- Gillick L, Baker J, Bridle J, et al. (1993) Application of large vocabulary continuous speech recognition to topic and speaker identification using telephone speech. In: *Proc. ICASSP 93*, volume 11, May 1993, San Francisco, Calif. IEEE CS Press, Piscataway, N.J., pp 471–474
- Smith MA, Christel MG (1995) Automating the creation of a digital video library. In: *Proc. ACM Multimedia 95*, November 1995, San Francisco, Calif. ACM Press, New York, pp 357–358
- Furui S (1994) An overview of speaker recognition technology. In: *ESCA Workshop on Automatic Speaker Recognition, Identification, and Verification*, April 1994, ESCA, pp 1–9
- Wilcox L, Chen F, Balasubramanian V (1994) Segmentation of speech using speaker identification. In: *Proc. ICASSP 94*, volume S1, April 1994, IEEE CS Press, Piscataway, N.J., pp 161–164
- Kimber D, Wilcox L (1996) Acoustic segmentation for audio browsers. In: *Proc. Interface Conference*, July 1996, Sydney, Australia (available on the Internet: <http://www.fxpal.xerox.com/abstracts/kim96.htm>)
- Chen F, Hearst M, Kimber D, Kupiec J, Pedersen J, Wilcox L (1997) Managing Multimedia Data: Using Metadata to Integrate and Apply Digital Data (Chapter Metadata for Mixed Media Access) McGraw-Hill, New York
- Roy D, Malamud C (1997) Speaker identification based text to audio alignment for an audio retrieval system. In: *Proc. ICASSP 97*, April 1997, Munich, Germany. IEEE CS Press, Piscataway, N.J., pp 1099–1102
- Wyse L, Smoliar S (1998) Toward content-based audio indexing and retrieval and a new speaker discrimination technique. In: Rosenthal DF, Okuno HG (eds) *Readings In Computational Auditory Scene Analysis*. Lawrence Erlbaum, New York
- Fischer S, Effelsberg W (1995) Automatic film genre classification. In: *Proc. ACM Multimedia '95*, November 1995, San Francisco, Calif. ACM Press, New York, pp 295–304
- Pfeiffer S, Fischer S, Effelsberg W (1996) Automatic audio content analysis. Technical Report TR-96-008, University of Mannheim, Mannheim, Germany. (available on the Internet: <ftp://pi4.informatik.uni-mannheim.de/pub/techreports/1996/TR-96-008.ps.gz>)
- Saunders J (1996) Real-time discrimination of broadcast speech/music. In: *Proc. ICASSP 96*, volume 11, May 1996, Atlanta, Ga. IEEE CS

- Press, Piscataway, N.J., pp 993–996
38. Scheirer E, Slaney M (1997) Construction and evaluation of a robust multifeature music/speech discriminator. In: Proc. ICASSP 97, volume 11, April 1997, IEEE CS Press, Piscataway, N.J., pp 1331–1334
 39. Spina M, Zue V (1996) Automatic transcription of general audio data: Preliminary analyses. In: Proc. International Conference on Spoken Language Processing, October 1996, Philadelphia, Pa., pp 594–597
 40. Foote JT (1997) A similarity measure for automatic audio classification. In: Proc. AAAI 1997 Spring Symposium on Intelligent Integration and Use of Text, Image, Video, and Audio Corpora, March 1997, Stanford, Palo Alto, Calif
 41. Arons B (1997) SpeechSkimmer: A system for interactively skimming recorded speech. *ACM Trans Comput Hum Interaction* 4(1):3–38 (available on the Internet: <http://barons.www.media.mit.edu/people/barons/papers/ToCHI97.ps>)
 42. Feiten B, Ungvary T (1991) Organizing sounds with neural nets. In: Proc. 1991 Int. Computer Music Conf., San Francisco, Calif. International Computer Music Association
 43. Feiten B, Günzel S (19) Automatic indexing of a sound database using self-organizing neural nets. *Comput Music J* 18(3):53–65
 44. Wold E, Blum T, Keslar D, Wheaton J (1996) Content-based classification, search, and retrieval of audio. *IEEE Multimedia* 3(3):27–36
 45. Foote JT, Silverman HF (1994) A model distance measure for talker clustering and identification. In: Proc. ICASSP 94, volume S1, April 1994, Adelaide, Australia. IEEE CS Press, Piscataway, N.J., pp 317–32
 46. Foote JT (1998) Rapid speaker identification using discrete MMI feature quantisation. *Expert Syst Appl* 13(4):283–289
 47. Foote JT (19) Content-based retrieval of music and audio. In Kuo CCJ, et al. (eds) *Multimedia Storage and Archiving Systems II*, Proc. SPIE, volume 3229, pp 138–147 (available on the Internet: <http://svr-www.eng-cam.ac.uk/~jtf/papers/spie97-abs.html>)
 48. Ghias A, et al. (1995) Query by humming. In: Proc. ACM Multimedia 95, November 1995, San Francisco, Calif. ACM Press, New York, pp 231–236
 49. McNab R, Smith L, Witten I, Henderson C, Cunningham S (1996) Towards the digital music library: Tune retrieval from acoustic input. In: Proc. Digital Libraries 96, pp 11–18 (available on the Internet: <http://www.cs.waikato.ac.nz/~rjmcnab/papers/mt.ps.gz>)
 50. McNab R, Smith L, Witten I, Henderson C (1997) Tune retrieval in the multimedia library. (available on the Internet: <http://www.cs.waikato.ac.nz/~rjmcnab/papers/mmtools.ps.gz>)
 51. Brown MG, Foote JT, Jones GJF, Spärck Jones K, Young SJ (1996) Open-vocabulary speech indexing for voice and video mail retrieval. In: Proc. ACM Multimedia 96, November 1996, Boston, Mass. ACM Press, New York, pp 35–43
 52. Hauptmann A, Witbrock M, Rudnick A, Reed S (1995) Speech for multimedia information retrieval. In: Proc. UIST-95 User Interface Software and Technology, Pittsburgh, Pa. pp 79–80 (available on the Internet: <http://informedia.cs.cmu.edu/research/uist95.ps>)
 53. Arik S (1996) Article extraction and classification of TV news using image and speech processing. In: International Symposium on Cooperative Database Systems for Advanced Applications (CODAS-96), Kyoto, Japan (available on the Internet: <http://banjo.kuis.kyoto-u.ac.jp/~tarumi/J/event/kyoto-camera/L006.ps>)



JONATHAN T. FOOTE was born in 1963 in Hollywood, California. He attended public schools in Santa Monica, California, somehow without learning how to surf. He received a Bachelor of Science (Electrical Engineering) degree in 1985 and a Master of Engineering (Electrical) degree in 1986 from Cornell University. From 1986 to 1988 he worked as a development engineer for Teradyne, Inc., in Boston, Massachusetts. From 1988 to 1993 he was at Brown University, where he earned a Ph.D. in Electrical Engineering. He received an Outstanding Research Award from the Brown University Chapter of Sigma Xi in 1992, and a Brown University Presidential Teaching Award in 1993. From 1993 to 1996 he pursued a postdoctoral fellowship at Cambridge University in England. With colleagues at Cambridge, his work on multimedia indexing using speech recognition received the Best Paper awards at ACM SIGIR 96 and ACM Multimedia 96. In 1996, he was awarded a J. William Fulbright Fellowship at the National University of Singapore, where he spent the following year. In December 1997, he joined FX Palo Alto Laboratory, Inc., as a Senior Research Scientist. He is also a professional musician, and enjoys not surfing in his free time.