

SpeechSkimmer: A System for Interactively Skimming Recorded Speech

BARRY ARONS

MIT Media Lab

Listening to a speech recording is much more difficult than visually scanning a document because of the transient and temporal nature of audio. Audio recordings capture the richness of speech, yet it is difficult to directly browse the stored information. This article describes techniques for structuring, filtering, and presenting recorded speech, allowing a user to navigate and interactively find information in the audio domain. This article describes the SpeechSkimmer system for interactively skimming speech recordings. SpeechSkimmer uses speech-processing techniques to allow a user to hear recorded sounds quickly, and at several levels of detail. User interaction, through a manual input device, provides continuous real-time control of the speed and detail level of the audio presentation. SpeechSkimmer reduces the time needed to listen by incorporating time-compressed speech, pause shortening, automatic emphasis detection, and nonspeech audio feedback. This article also presents a multilevel structural approach to auditory skimming and user interface techniques for interacting with recorded speech. An observational usability test of SpeechSkimmer is discussed, as well as a redesign and reimplementaion of the user interface based on the results of this usability test.

Categories and Subject Descriptors: D.2.2 [**Software Engineering**]: Tools and Techniques—*user interfaces*; H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval; H.5.1 [**Information Interfaces and Presentation**]: Multimedia Information Systems—*audio input/output*; H.5.2 [**Information Interfaces and Presentation**]: User Interfaces—*evaluation/methodology; input devices and strategies; interaction styles*

General Terms: Design, Experimentation, Human Factors

Additional Key Words and Phrases: Audio browsing, interactive listening, nonspeech audio, speech as data, speech skimming, speech user interfaces, time compression

1. INTRODUCTION

Speech is a powerful communications medium that is rich and expressive. Speech is natural, portable, and can be used while doing other things. It is faster to speak than it is to write or type [Gould 1982]; however, it is slower

This work was sponsored by Apple Computer Corp. and Interval Research Corporation.

Author's address: Speech Interaction Research, P. O. Box 14, Cambridge, MA 02142; email: barons@media.mit.edu.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1997 ACM 1073-0516/97/0300-0003 \$03.50

to listen than it is to read. Therefore, recording speech is efficient for the talker, but hearing recorded speech is usually a burden on the listener. Skimming, browsing, and searching are traditionally considered visual tasks that one readily performs while reading a newspaper, window shopping, or driving a car. However, there is no natural way for humans to skim speech information because of the transient nature of audio—the ear cannot skim in the temporal domain the way the eyes can browse in the spatial domain.

This article describes SpeechSkimmer, a system for skimming speech recordings that attempts to overcome these problems of slowness and the inability to browse audio. SpeechSkimmer uses simple speech-processing techniques to allow a user to hear recorded sounds quickly, and at several levels of detail. User interaction through a manual input device provides continuous real-time control over the speed and detail level of the audio presentation.

SpeechSkimmer explores a new paradigm for interactively skimming and retrieving information in speech interfaces. This research takes advantage of knowledge of the speech communication process by exploiting structure, features, and redundancies inherent in spontaneous speech. Talkers embed lexical, syntactic, semantic, and turn-taking information into their speech as they have conversations and articulate their ideas [Levelt 1989]. These cues are realized in the speech signal, often as hesitations or changes in pitch and energy.

Speech also contains redundant information; high-level syntactic and semantic constraints of English allow listeners to understand speech when it is severely degraded by noise, or even if entire words or phrases are removed. Within words there are other redundancies that allow partial or entire phonemes to be removed while still retaining intelligibility.

This research attempts to exploit acoustic cues to segment recorded speech into semantically meaningful chunks. The recordings are then time-compressed to further remove redundant speech information. While there are practical limits to time compression, there are compelling reasons to be able to quickly skim a large speech document. For skimming, redundant as well as nonredundant segments of speech must be removed. Ideally, as the skimming speed increases, the segments with the least information content are eliminated first.

When searching for information visually, we tend to refine our search over time, looking successively at more detail. For example, we may glance at a shelf of books to select an appropriate title, flip through the pages to find a relevant chapter, skim headings to find the right section, then alternately skim and read the text until we find the desired information. To skim and browse recorded speech in an analogous manner the listener must have interactive control over the level of detail, rate of playback, and style of presentation. SpeechSkimmer allows a user to control the auditory presentation through a simple interaction mechanism that changes the granularity, time scale, and style of presentation of the recording.

1.1 Speech as Sound

Along with the meaning of our spoken words, our emotions and important syntactic and semantic information are captured by the pitch, timing, and volume of our speech. At times, more significance can be transmitted with silence than by the use of words. Such information is difficult to convey in a textual or graphical form and is best captured in the sounds themselves. Transcripts can be useful for browsing visually or for electronic keyword searches. However, transcripts are expensive, and automated transcriptions of spontaneous speech, meetings, or conversations are not practical in the foreseeable future [Roe and Wilpon 1993]. A waveform, spectrogram, or other graphical representation can be displayed (see also Section 3.2), yet this does not indicate what was spoken, or how something was said. Speech needs to be heard.

A graphical user interface may make some speech-searching and skimming tasks easier, but there are two reasons for exploring interfaces without a visual display. First, there are a variety of situations where a graphical interface cannot be used, such as while walking, driving, or if the user is visually impaired. Second, an important issue addressed in this research is structuring and extracting information from the speech signal and then presenting it in an auditory form. Once techniques are developed to process and present speech information that take advantage of the audio channel, they can be applied to visual interfaces.

Early versions of SpeechSkimmer therefore explored moving through speech recordings without a visual display. In the usability test of the system, users requested some graphical feedback to help them navigate through a speech recording. The revised SpeechSkimmer user interface incorporates a small amount of visual feedback, but it still can be used without looking at it.

2. BASE TECHNOLOGIES AND SEGMENTATION

This section introduces the core speech technologies used in the SpeechSkimmer system including time compression of speech, adaptive speech detection, emphasis detection, and segmenting recordings based on pauses and pitch. Readers interested in the user interface design and testing of SpeechSkimmer should skip to Section 3.

2.1 Time-Compressing Speech

The length of time needed to listen to an audio recording can be reduced through a variety of time compression methods (reviewed in Arons [1992a]). These techniques allow recorded speech to be sped up (or slowed down) while maintaining intelligibility and voice quality. Time compression can be used in many application environments including voice mail, recordings for the blind, and human-computer interfaces.

2.1.1 *Time Compression Techniques.* A recording can simply be played back with a faster clock rate than it was recorded at, but this produces an

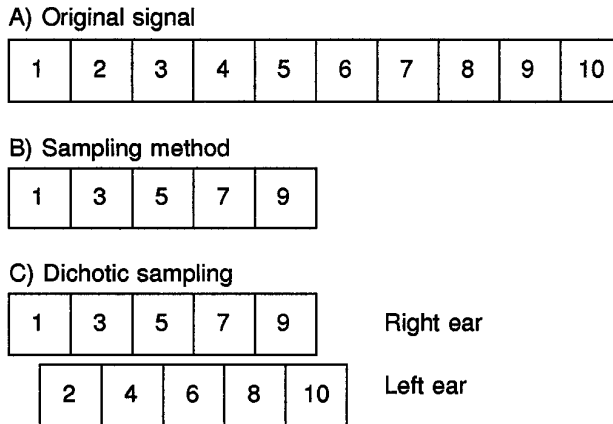


Fig. 1. For a $2\times$ speed increase using the sampling method (b), every other chunk of speech from the original signal is discarded (50ms chunks are used). The same technique is used for dichotic presentation, but different segments are played to each ear (c).

increase in pitch causing the speaker to sound like Mickey Mouse. This frequency shift results in an undesirable decrease of intelligibility. The most practical time compression techniques work in the time domain and are based on removing redundant information from the speech signal. In the *sampling* method [Fairbanks et al. 1954], short segments¹ are dropped from the speech signal at regular intervals (Figure 1). Cross fading, or smoothing, between adjacent segments improves the resulting sound quality.

The *synchronized overlap add method* (SOLA) is a variant of the sampling method that is becoming prevalent in computer-based systems [Hejna 1990; Roucos and Wilgus 1985]. Conceptually, the SOLA method consists of shifting the beginning of a new speech segment over the end of the preceding segment (Figure 2) to find the point of highest cross correlation (i.e., maximum similarity). The overlapping frames are averaged, or smoothed together, as in the sampling method. SOLA can be considered a type of selective sampling that effectively removes entire pitch periods. SOLA produces the best-quality speech for a computationally efficient time domain technique.

Sampling with dichotic presentation² is a variant of the sampling method that takes advantage of the auditory system's ability to integrate information from both ears. It improves on the sampling method by playing the standard sampled signal to one ear and the "discarded" material to the other ear [Scott 1967] (Figure 1(c)). Intelligibility and comprehension increase under this dichotic presentation condition when compared with standard presentation techniques [Gerber and Wulfeck 1977].

¹The segments are typically 30–50 microseconds—longer than a pitch period, but shorter than a phoneme.

²A different signal is played to each ear through headphones.

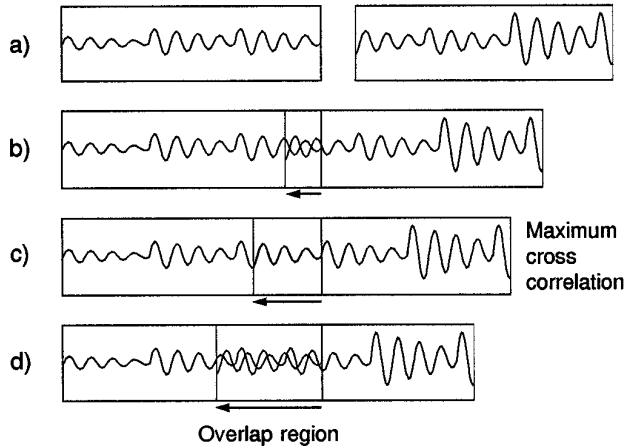


Fig. 2. SOLA: shifting the speech segments (as in Figure 1) to find the maximum cross correlation. The maximum similarity occurs in case (c), eliminating the whole pitch period.

SpeechSkimmer incorporates several time compression techniques for experimentation and evaluation purposes. All of these speech-processing algorithms run in real time on the main processor of a laptop computer (an Apple Macintosh PowerBook 170 was used) and do not require special signal-processing hardware.

2.1.2 Perception of Time-Compressed Speech. *Intelligibility* usually refers to the ability to identify isolated words. *Comprehension* refers to understanding the content of the material (obtained by asking questions about a recorded passage). Early studies showed that isolated words that are carefully selected and trained can remain intelligible up to 10 times normal speed, while continuous speech remains comprehensible up to about twice ($2\times$) normal speed. Time compression decreases comprehension because of a degradation of the speech signal and a processing overload of short-term memory. A $2\times$ increase in speed removes virtually all redundant information [Heiman et al. 1986]; with greater compression, critical nonredundant information is also lost.

Both intelligibility and comprehension improve with exposure to time-compressed speech. Beasley and Maki [1976] informally reported that, following a 30-minute exposure to time-compressed speech, listeners became uncomfortable if they were forced to return to the normal rate of presentation. They also found that subjects' listening rate preference shifted to faster rates after exposure to compressed speech. Perception of time-compressed speech is reviewed in more detail in Arons [1992a], Beasley and Maki [1976], and Foulke [1971].

2.2 Pauses in Speech

Removing (or shortening) pauses from a recording can be used as a form of time compression. The resulting speech is "natural, but many people find it exhausting to listen to because the speaker never pauses for breath"

[Neuburg 1978]. In the perception of normal speech, it has been found that pauses exert a considerable effect on the speed and accuracy with which sentences were recalled, particularly under conditions of cognitive complexity—“Just as pauses are critical for the speaker in facilitating fluent and complex speech, so are they crucial for the listener in enabling him [or her] to understand and keep pace with the utterance” [Reich 1980]. Pauses, however, are only useful when they occur between clauses within sentences—pauses within clauses are disrupting. Pauses suggest the boundaries of material to be analyzed and provide vital cognitive processing time.

Hesitation pauses are not under the conscious control of the talker, and they average 200–250 microseconds. *Juncture* pauses are under talker control, usually occur at major syntactic boundaries, and average 500–1000 microseconds [Minifie 1974]. Recent work, however, suggests that such categorical distinctions of pauses based solely on length cannot be made [O’Shaughnessy 1992]. Juncture pauses are important for comprehension and cannot be eliminated or reduced without interfering with comprehension [Lass and Leeper 1977].

2.2.1 Adaptive Speech Detection. Speech is a time-varying signal; silence (actually background noise) is also time-varying. Background noise may consist of mechanical noises such as fans, that can be defined temporally and spectrally, but can also consist of conversations, movements, and door slams that are difficult to characterize. Speech detection involves classifying these two types of signals. Due to the variability of the speech and background noise patterns, it is desirable to use an adaptive solution that does not rely on arbitrary fixed thresholds [de Souza 1983; Savoji 1989]. The most common error made by speech detectors is the misclassification of unvoiced consonants, or weak voiced segments, as background noise.

An adaptive speech detector (based on Lamel et al. [1981]) was developed for shortening and removing pauses and to provide data for segmentation. Digitized speech files are analyzed in several passes. The first pass gathers energy³ and zero crossing rate⁴ (ZCR) statistics for 10ms frames of audio. The background noise level is determined by smoothing a histogram of the energy measurements and finding the peak of the histogram. The peak corresponds to an energy value that is part of the background noise. A value several decibels above this peak is selected as the dividing line between speech and background noise. The noise level and ZCR metrics provide an initial classification of each frame as speech or background noise.

Additional passes through the sound data are made to refine this estimation based on heuristics of spontaneous speech. This processing fills

³Average magnitude is used as a measure of energy [Rabiner and Sambur 1975].

⁴A high zero crossing rate indicates low-energy fricative sounds such as “s” and “f.” For example, ZCR greater than 2500 crossings/second indicates the presence of a fricative [O’Shaughnessy 1987]. Note that the background noise in most office environments does not contain significant energy in this range.

in short gaps between speech segments [Gruber 1982], removes isolated islands initially classified as speech, and extends the boundaries of speech segments so that they are not inadvertently clipped [Gruber and Lee 1983]. For example, two or three frames initially classified as background noise amid many high-energy frames identified as speech are treated as part of that speech, rather than as a short silence.

This speech detector is fast and works well under a variety of microphone and noise conditions. Audio files recorded in an office environment with computer fan noise and in a lecture hall with over 40 students have been successfully segmented into speech and background noise. See Arons [1994a] for a review of other speech detection techniques and details of the algorithm.

2.3 Acoustically Based Segmentation

Speech recordings need to be segmented into manageable pieces before presentation. Salient audio segments can be automatically selected from a recording by exploiting properties of spontaneous speech. Segmenting audio and finding its inherent structure are essential for the success of future recording-based systems. “Finding the structure” means locating important or emphasized portions of a recording, and selecting the equivalent of paragraph or new-topic boundaries, for the purpose of creating audio overviews or outlines. Ideally, a hierarchy of segments can be created that roughly corresponds to the spoken equivalents of sections, subsections, paragraphs, and sentences of a written document.

Several acoustic cues were explored for segmenting speech:

- Pauses* can suggest the beginning of a new sentence, thought, or topic. Studies have shown that pause lengths are correlated with the type of pause and its importance (see Section 2.2).
- Pitch* is similarly correlated with a talker’s emphasis and new-topic introductions.
- Speaker identification* for separating talkers in a conversation.

None of these techniques are 100% accurate at finding the important boundaries in speech recordings—they all miss some of the desired boundaries and incorrectly locate others. While it is important to minimize these errors, it is perhaps more important to be able to handle errors when they occur, as no such recognition technology will ever be perfect. SpeechSkimmer addresses using such error-prone cues by providing the user with an interface to navigate in a recording and control what segments get played and how they are presented, allowing the user to listen to exactly what he or she wants to hear.

2.3.1 *Speech Detection for Segmentation.* The adaptive speech detector developed for finding and shortening pauses (Section 2.2.1) can also be used for segmentation. Since long pauses typically correspond with juncture pauses that occur at important boundaries (Section 2.2), the lengths of the pauses in a recording can be used to segment the speech. For example,

segmenting a recording with a pause threshold of 0.02 would select the segments of speech that occur just after the longest 2% of the pauses in the recording. Note that a relative, rather than absolute, pause length is used to adapt to the pausing characteristics of the talker.

2.3.2 Pitch-Based Emphasis Detection for Segmentation. Pitch⁵ provides important information for human comprehension and understanding of speech and can also be exploited in machine-mediated systems. For example, there tends to be an increase in pitch range when a talker introduces a new topic [Hirschberg and Grosz 1992; Hirschberg and Pierrehumbert 1986; Silverman 1987], which is an important cue for listeners.

Chen and Withgott [1992] trained a Hidden Markov Model (HMM) [Rabiner 1989] to summarize recordings based on training data hand-marked for emphasis, combined with the pitch and energy content of conversations. They successfully created summaries of the recordings by selecting emphasized portions that were in close temporal proximity. This prosodic approach is promising for extracting high-level information from speech signals. An alternative technique was developed for SpeechSkimmer to detect salient segments and summarize a recording without using statistical models that require large amounts of training data.

A variety of simple pitch metrics were generated and manually correlated with a hand-marked transcript of a 15-minute recording. The metrics were gathered over one-second windows of pitch data (100 frames of 10ms). The number of frames above a threshold and the standard deviation were most strongly correlated with new-topic introductions and emphasized portions of the transcript. Note that these two metrics essentially measure the same thing: significant range and variability in F0. The metric “number of frames above a threshold” was used in the subsequent development of the algorithm.

Since the range and baseline pitch vary considerably between talkers, it is necessary to adaptively determine the pitch threshold for a given speaker. A histogram of the pitch data is used to normalize talker variability, and a threshold is chosen to select the top 1% of the pitch frames. The number of frames in each one-second window that are above the threshold is counted as a measure of emphasis. The scores of nearby windows are then combined for phrase- or sentence-sized segments of the speech recording.

This pitch-based segmentation technique has been successfully used to provide a high-level summary of speech recordings for a variety of talkers. High-scoring salient segments are selected and used by SpeechSkimmer to enable efficient skimming. For further information on the emphasis detector see Arons [1994b].

2.3.3 Evaluating the Emphasis Detection. Stifelman [1995] compared the segmentation of the emphasis detection algorithm with a hierarchical

⁵“Pitch” in this context means the fundamental frequency of voiced speech and is usually denoted as F0.

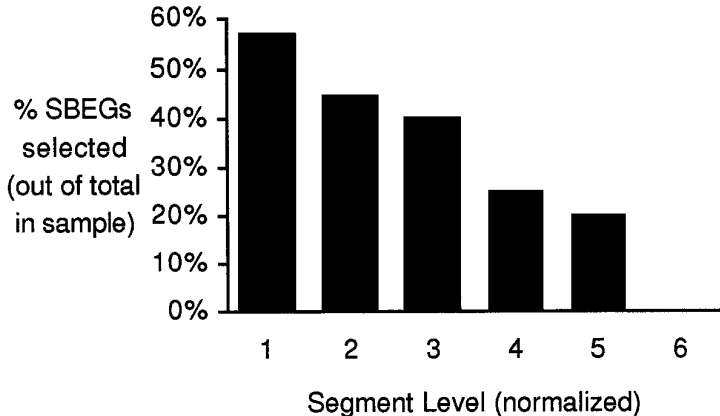


Fig. 3. Percent of segment beginnings selected for each level in the discourse hierarchy (after Stifelman [1995]).

segmentation based on Grosz and Sidner's [1986] theory of discourse structure. Stifelman found that the emphasis algorithm has a relatively high precision (82%), but a low recall (25%), for selecting discourse boundaries in the speech sample tested. This means that the majority of segments selected by the algorithm were good segments in the discourse structure, but that the algorithm did not find all the desired segments.

The discourse structure of a monologue can be thought of as an outline. To extract high-level ideas from a recording the major points in the outline are of most interest, rather than those that are deeply embedded. The outermost segments in the discourse structure need to be found for high-level skimming or summarization. Figure 3 shows the percentage of segments the emphasis detector selected compared to those manually selected at each level in the discourse hierarchy. A greater proportion of the major points in the discourse structure were found, rather than embedded ones.

While there is room for improvement, these results appear promising. The emphasis detection algorithm did select a number of high-level points from the discourse hierarchy without too many false alarms (see also Section 4.2.4). Unfortunately, the algorithm did select some segments that were deeply embedded in the discourse, and the recall rate could be improved. Using this emphasis algorithm as a starting point, it may be possible to improve these scores by tuning the algorithm, or combining it with other acoustic features such as pauses.

2.3.4 Segmentation by Speaker Identification. Acoustically based speaker identification [Kimber et al. 1995; Reynolds and Rose 1995] can provide a powerful cue for segmentation and information retrieval in speech systems. For example, when searching for a piece of information within a recording, the search space can be greatly reduced if individual talkers can be identified (e.g., "play only things Marc said"). See Section 3.2.1 to see how speaker identification data were used in SpeechSkimmer.

1. Unprocessed (normal)
2. Time-compression
 - Pause removal
 - Pause shortening
 - Sampling
 - SOLA
 - Dichotic sampling or SOLA
 - Combined time compression techniques (e.g., SOLA with pause removal)
 - Backward sampling (for intelligible rewind)
3. Skimming
 - Isochronous skimming (equal time intervals)
 - Speech synchronous skimming. Segmentation based on:
 - Pauses
 - Pitch
 - Energy
 - Speaker identification
 - Word spotting
 - User selected segments
 - Combined segmentation techniques (e.g., pauses, pitch, and energy)
 - Backward skimming

Fig. 4. Techniques of time compression and skimming.

3. SPEECHSKIMMER PROTOTYPE

This section integrates the technologies described in the previous section into a coherent system for interactive listening. A framework is described for presenting a continuum of time compression and skimming techniques. This allows a user to quickly skim a speech recording to find portions of interest, then use time compression and pause shortening for efficient browsing, and then slow down further to listen to detailed information. A multilevel approach to auditory skimming, along with user interface techniques for interacting with the audio and providing feedback, is presented.

3.1 Time Compression and Skimming

SpeechSkimmer incorporates ideas and techniques from conventional time compression algorithms and attempts to go beyond the $2\times$ perceptual barrier typically associated with time-scaling speech. These new skimming techniques are intimately tied to user interaction to provide a range of audio presentation speeds. Backward variants of the techniques are also developed to allow audio recordings to be played and skimmed backward as well as forward. Some of the possible time compression and skimming technologies that can be used are shown in Figure 4. Corresponding ranges of speed increases for the different classes of techniques are shown in Figure 5.

Time compression can be considered as “content lossless,” since the goal is to present all the nonredundant speech information in the signal. The skimming techniques are designed to be “content lossy,” as large parts of the speech signal are explicitly removed. This classification is not based on the traditional engineering concept of lossy versus lossless, but is based on the intent of the processing. For example, isochronous skimming selects

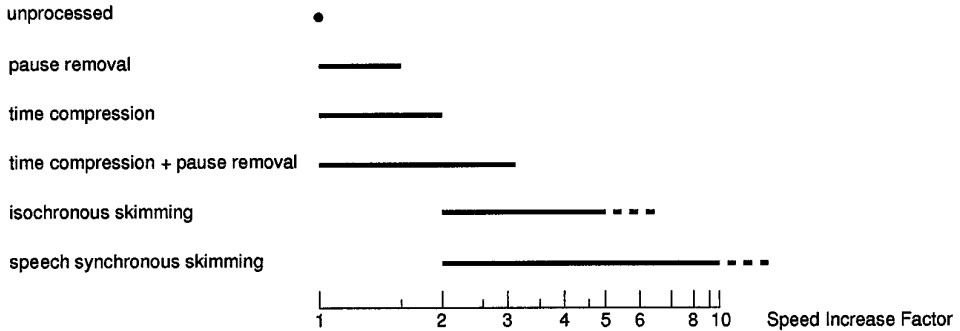


Fig. 5. Schematic representation of the range of speed increases for different time compression and skimming methods.

and presents speech segments based on equal time intervals. If only the first five seconds of each minute of speech are played, this can be considered coarse and lossy sampling. In contrast, a speech-synchronous technique that selects important words and phrases using the natural boundaries in the speech will provide more information content to the listener.

3.2 Skimming Levels

There have been a variety of attempts to present hierarchical or “fisheye” views of visual information [Furnas 1986; Mackinlay et al. 1991]. These approaches are powerful, but inherently they rely on a spatial organization. Temporal video information has been displayed in a similar form [Davis 1995; Elliott 1993; Mills et al. 1992], yet this primarily consists of mapping time-varying spatial information into the spatial domain. Graphical techniques can be used for a waveform or similar display of an audio signal, but such a representation is inappropriate—*sounds need to be heard, not viewed*. This research attempts to present a hierarchical (or “fish ear”) representation of audio information that *only* exists temporally.

A continuum of time compression and skimming techniques has been designed, allowing a user to efficiently skim a speech recording to find portions of interest, then listen to it time-compressed to allow quick browsing, and then slow down further to listen to detailed information. Figure 6 presents one possible “fish ear” view of this continuum. For example, what may take 60 seconds to listen to at normal speed may take 30 seconds when time-compressed and only five or ten seconds at successively higher levels of skimming. If the speech segments are chosen appropriately, it is hypothesized that this mechanism provides a summarizing view of a speech recording.

Four distinct skimming levels have been implemented (Figure 7). Within each level the speech signal can also be time compressed. The lowest skimming level (level 1) consists of the original speech recording without any processing and thus maintains the pace and timing of the original signal. In level 2, the pauses are selectively shortened or removed. Pauses less than 500ms are removed, and the remaining pauses are shortened to

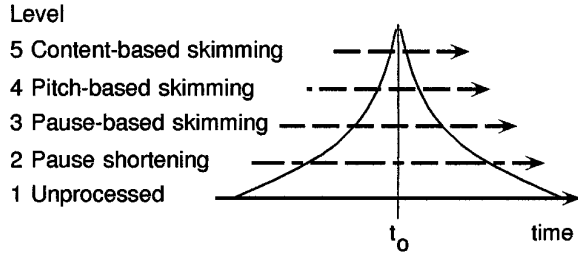


Fig. 6. A hierarchical “fish ear” time scale continuum. Each level in the diagram represents successively larger portions of the levels below it. The curves represent iso-content lines, i.e., an equivalent time mapping from one level to the next. The current location in the sound file is represented by t_0 ; the speed and direction of movement of this point depends upon the skimming level.

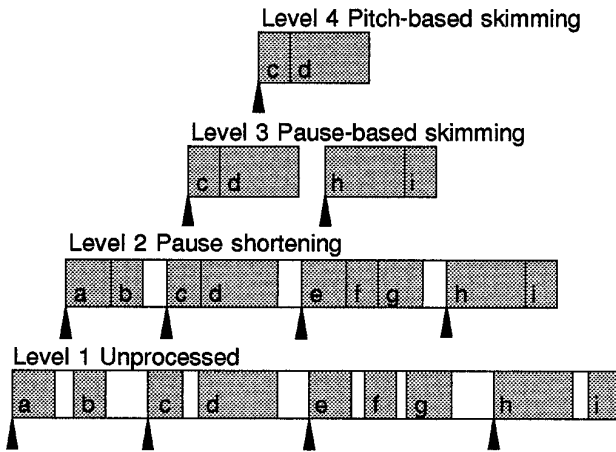


Fig. 7. Speech and silence segments played at each skimming level. The gray boxes represent speech; white boxes represent background noise. The pointers indicate valid segments to go to when jumping or playing backward.

500ms.⁶ This technique speeds up listening yet provides the listener with cognitive processing time and cues to the important juncture pauses (Section 2.2).

Level 3 is based on the premise that long juncture pauses tend to indicate either a new topic, some content words, or a new talker. For example, filled pauses (i.e., “uhh” or “um”) usually indicate that the talker does not want to be interrupted, while long unfilled pauses (i.e., silences) signify that the talker is finished and that someone else may begin speaking [Levelt 1989; O’Shaughnessy 1992]. Thus level 3 skimming attempts to play salient segments based on a simple heuristic: only the speech that occurs just after a significant pause in the original recording is played. For example, after a pause over 750ms is detected, the subsequent five seconds of speech are played (with pauses shortened). Note again that this segmentation process

⁶All thresholds are determined adaptively based on the content of the speech recording.

is error prone, but these errors can be overcome by giving the user interactive control of the presentation.

Level 4 is similar to level 3 in that it attempts to present segments of speech that are highlights of the recording. Salient segments for level 4 are chosen using the emphasis detector (Section 2.3.2) to summarize the recording. In practice, either level 3 or level 4 is used as the top skimming level.

It is somewhat difficult to listen to level 3 or level 4 skimmed speech, as relatively short unconnected segments are played in rapid succession. It has been informally found that playing the segments at normal speed (i.e., not time compressed), or even slowing down the speech, is useful when skimming unfamiliar material. At the highest skimming levels, a short (e.g., 600ms) pure silence is inserted between each of the speech segments to separate them perceptually. An early version of SpeechSkimmer played recorded ambient noise between the selected segments, but this fit in so naturally with the speech that it was difficult to distinguish between segments.

3.2.1 Alternative Skimming Levels Using Speaker Identification. The SpeechSkimmer system has also been used with speaker identification-based segmentation. A two-person conversation was analyzed with speaker identification software [Reynolds and Rose 1995] that determined when each talker was active. These data were translated into SpeechSkimmer format such that level 1 represented the entire conversation; jumping took the listener to the next turn change in the conversation. Level 2 played only the speech from one talker, while level 3 played the speech from the other. Jumping within these levels brought the listener to the start of that talker's next conversational turn.

3.3 Skimming Backward

Besides skimming forward through a recording, it is desirable to play intelligible speech while interactively searching or “rewinding” through a digital audio file [Arons 1991a; Elliott 1993]. Analog tape systems provide little useful information about the signal when it is played completely backward. This is analogous to taking the text “going to the store” and presenting it as the unintelligible “erots eht ot gniog.” Digital systems allow word- or phrase-sized chunks of speech to be played forward individually, with the segments themselves presented in reverse order (resulting in “store, to the, going”). While the general sense of the recording is reversed and jumbled, each segment is identifiable and intelligible. It can thus become practical to browse backward through a recording to find a particular word or phrase. This method is particularly effective if the segment boundaries are at natural pauses in the speech. Note that this technique can also be combined with time-compressed playback, allowing both backward and forward movement at high speeds.

In addition to the forward skimming levels, the speech recordings can also be skimmed backward. Small segments of sound are each played

normally, but are presented in reverse order. When level 1 and level 2 sounds are played backward (i.e., considered level -1 and level -2), short segments are selected based upon speech detection and are played in inverse order. In Figure 7 level -1 would play segments in this order: h-i, e-f-g, c-d, a-b. Level -2 is similar, but without the pauses.

3.4 Jumping

Along with controlling the skimming and time compression, it is desirable to be able to interactively jump between segments within each skimming level. If the user decides that the segment being played is not of interest, it is possible to go on to the next segment without being forced to listen to each entire segment [Arons 1991b; Resnick and Virzi 1992a]. For example, in Figure 7 at level 3, segments c and d would be played, then a short silence, then segments h and i. At any time while the user is listening to segment c or d, a jump forward command would immediately interrupt the audio output and start playing segment h. While listening to segment h or i, the user could jump backward, causing segment c to be played. Valid segments to jump to are indicated with pointers in Figure 7.

The skimming user interface includes a control that jumps backward one segment and drops into normal play mode (level 1, no time compression). The intent of this control is to encourage high-speed browsing of time-compressed level 3 or level 4 speech. When the user hears something of interest, it is easy to use this control to back up a bit, rehear the piece of interest, and then continue listening at normal speed.

3.5 Interaction Mappings

Finding an appropriate mapping between an input device and controlling the skimmed speech is subtle, as there are many independent variables that can be controlled. For the SpeechSkimmer prototype, the primary variables of interest are time compression and skimming level, with all others (e.g., pause-shortening parameters and skimming timing parameters) held constant.

Several mappings of user input to time compression and skimming level were tried. A two-dimensional controller, such as a mouse, allows two variables to be changed independently. For example, the y-axis can be used to control the amount of time compression while the x-axis controls the skimming level (Figure 8). Movement toward the top increases the time compression; movement toward the right increases the skimming level. The right half is used for skimming forward, the left half for skimming backward. Moving to the upper right thus presents skimmed speech at high speed.

The two primary variables can also be controlled by a one-dimensional input device. For example, as the controller is moved forward, the sound playback speed is increased using time compression. As it is pushed forward further, time compression increases until crossing a boundary into the next skimming level. Pushing forward within each skimming level

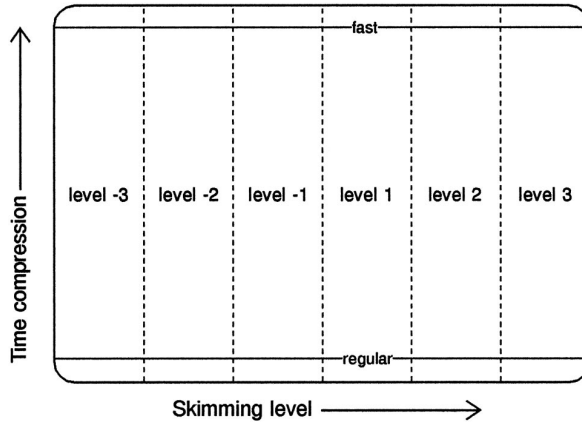


Fig. 8. Schematic representation of two-dimensional control regions. Vertical movement changes the time compression; horizontal movement changes the skimming level.

similarly increases the time compression (Figure 9). Pulling backward has an analogous but reverse effect.

One consideration in all these schemes is the continuity of speeds when transitioning from one skimming level to the next. In Figure 9, for example, when moving from fast level 2 skimmed speech to level 3 speech there is a sudden change in speed at the border between the two skimming levels. Depending upon the implementation, fast level 2 speech may be effectively faster or slower than regular level 3 speech. This problem also exists with a 2D control scheme—to monotonically increase the effective playback speed may require a zigzag motion through skimming and time compression levels.

3.6 Interaction Devices

The speech-skimming system has been used with a mouse, small trackball, touchpad, and a joystick in both the one- and two-dimensional control configurations (two independent controls, one for speed and one for skimming level, were not tried). A mouse provides accurate control, but as a relative pointing device [Card et al. 1991] it is difficult to use without a display. A small hand-held trackball (e.g., controlled with the thumb) eliminates the desk space required by the mouse, but is still a relative device and is therefore also inappropriate for a nonvisual task.

A joystick can be used as an absolute position device. However, if it is spring-loaded (i.e., automatic return to center), it requires constant physical force to hold it in position. If the springs are disabled, a particular position (i.e., time compression and skimming level) can be automatically maintained when the hand is removed (see Lipscomb and Pique [1993] for a discussion of such physical considerations). The home (center) position, for example, can be configured to play forward (level 1) at normal speed. Touching or looking at the joystick's position provides feedback to the current settings. However, in either configuration, an off-the-shelf joystick

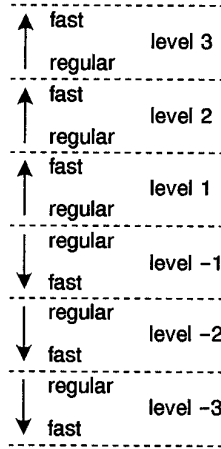


Fig. 9. Schematic representation of one-dimensional control regions.

does not provide any physical feedback when the user is changing from one discrete skimming level to another, and it is difficult to jump to an absolute location.

A small touchpad can act as an absolute pointing device and does not require any effort to maintain the last position selected. A touchpad can be easily modified to provide a physical indication of the boundaries between skimming levels. Unfortunately, a touchpad does not provide any physical indication of the current location once the finger is removed from the surface.

3.7 Touchpad Configuration

The SpeechSkimmer prototype uses a small (7×11 cm) touchpad [Micro-touch 1992] with a two-dimensional control scheme. Small strips of paper were added to the touch-sensitive surface as tactile guides to indicate the boundaries between skimming regions (Figure 10). In addition to the six regions representing skimming levels, two additional regions were added to jump directly to the beginning and end of the sound recording. Four buttons provide jumping and pausing capabilities (Figure 11). Note that the template used in the touchpad only contains static information; it is not necessary to look at it to use the system.

The time compression control (vertical motion) is not continuous, but provides a “finger-sized” region around the “regular” mark that plays at normal speed (Figure 12). To enable fine-grained control of the time compression [Stifelman 1994], a larger region is allocated for speeding the speech up than for slowing it down. The areas between the tactile guides form virtual sliders that control the time compression within a skimming level (note that only one slider is active at a time).

3.8 Nonspeech Audio Feedback

SpeechSkimmer uses recorded sound effects to provide feedback when navigating [Buxton et al. 1991; Gaver 1989]. Nonspeech audio was selected

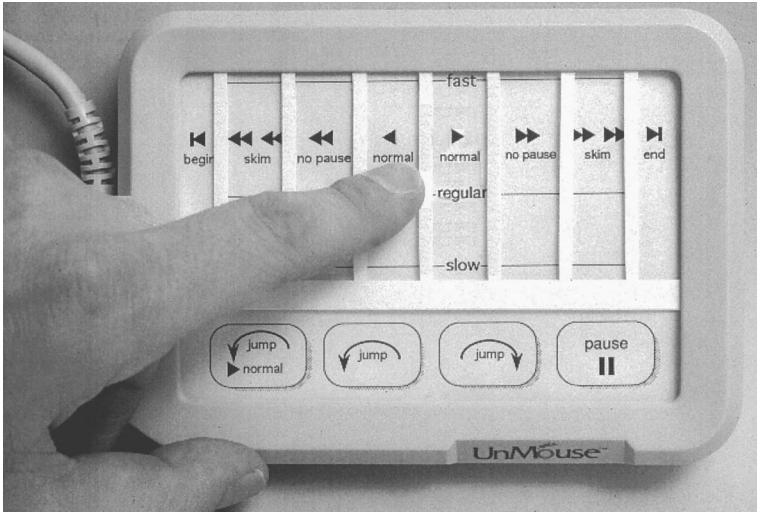


Fig. 10. The touchpad with paper guides for tactile feedback.

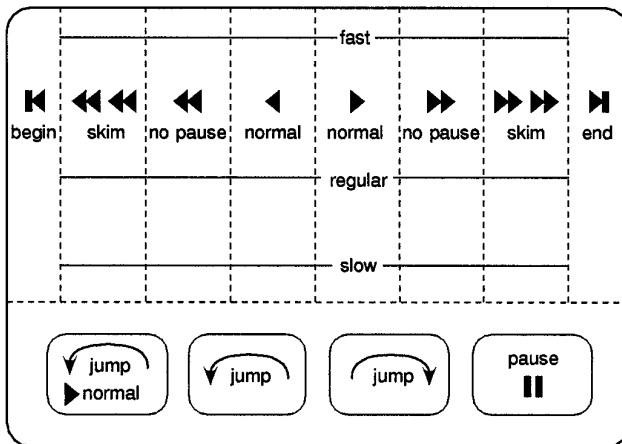


Fig. 11. Template used in the touchpad (a printed version of this fits behind the touch-sensitive surface of the pad). The dashed lines indicate the location of the paper guide strips.

to provide terse, yet unobtrusive navigational cues [Stifelman et al. 1993]. For example, if the user attempts to play past the end of a sound, a cartoon “boing” is played. No explicit feedback is provided for changes in time compression. The speed changes occur with low latency and are readily apparent in the speech signal itself.

When the user transitions to a new skimming level, a short tone is played. The frequency of the tone increases with the skimming level (Figure 13). A double beep is played when the user changes to normal (level 1). This acts as an audio landmark, clearly distinguishing it from the other tones and skimming levels.

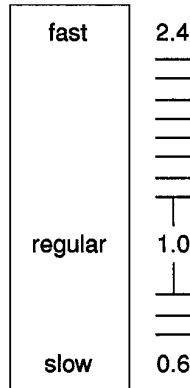


Fig. 12. Mapping of the touchpad control to the time compression range ($0.6\times$ to $2.4\times$).

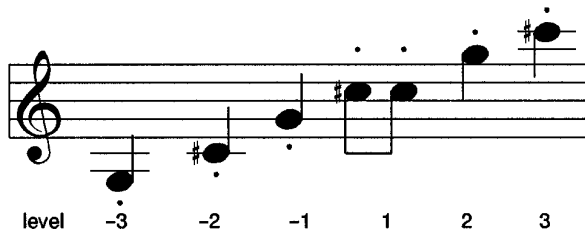


Fig. 13. A musical representation of the tones played at different skimming levels. Notice the double beep “landmark” for normal (level 1) playing. The small dots indicate short and crisp (*staccato*) notes.

A different sound is played when each of the buttons is touched. An attempt was made to create sounds that could be intuitively linked with the function of the button. The feedback played when pausing and unpausing are reminiscent of a piece of machinery stopping and starting. Jumping forward is associated with a rising pitch while jumping backward is associated with a falling pitch.

3.9 Software Architecture

Each recording is postprocessed with the speech detection and emphasis detection algorithms. A single file is created that contains all the segmentation data used for skimming, jumping, and pause shortening.

The run-time application consists of three primary modules: a main event loop, a segment player, and a sound library (Figure 14). The skimming user interface is separated from the underlying mechanism that presents the skimmed and time-compressed speech. This modularization allows for the rapid prototyping of new interfaces using different interaction devices. SpeechSkimmer is implemented in a subset of C++ and runs on Apple Macintosh computers.

The main event loop gathers raw data from the user and maps it to the appropriate time compression and skimming ranges for each input device.

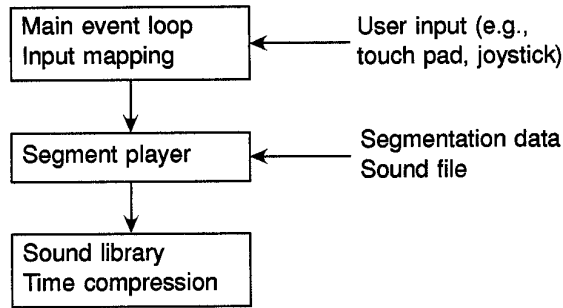


Fig. 14. Software architecture of the skimming system.

The event loop sends requests to the segment player to start and stop playback, jump between segments, and set the time compression and skimming levels.

The segment player combines user input with the segmentation data to select the appropriate portion of the sound to play. When the end of a segment is reached, the next segment is selected based on the current skimming level and data in the segmentation file. Audio data are read from the sound file and passed to the sound library. The size of the audio data buffers is kept to a minimum to reduce the latency between user input and the corresponding sound output.

The sound library provides a high-level interface to the audio playback hardware (based on the functional interface described in Arons [1992b]). Several different time compression algorithms are built into the sound library.

4. USABILITY TESTING

The goal of this test was to find usability problems and successes in the SpeechSkimmer user interface. The usability test was primarily an observational thinking-out-loud study [Ericsson and Simon 1984] that is intended to quickly find major problems in the user interface to an interactive system [Nielsen 1993a].

4.1 Method

4.1.1 Subjects. Twelve volunteer subjects between the ages of 21 and 40 were selected from the Media Laboratory environment. Six of the subjects were administrative staff, and six were graduate students; eight were female, and four were male. None of the subjects were familiar with SpeechSkimmer, but all had experience using computers. Test subjects were not paid, but were offered snacks and beverages to compensate them for their time.

4.1.2 *Procedure.* The tests were performed in an acoustically isolated room with a subject, an interviewer, and an observer.⁷ The sessions were videotaped and later analyzed by both the interviewer and observer. A testing session took approximately 60 minutes and consisted of five parts:

- (1) *A background interview to collect demographic information and to determine what experience subjects had with recorded speech and audio.* Subsequent questions were tailored based on the subject's experiences. For example, someone who regularly recorded lectures would be asked in detail about his or her use of the recordings, how they located specific pieces of information in the recordings, etc.
- (2) *A first look at the touchpad.* Subjects were given the touchpad (Figure 10) and asked to describe their first intuitions about the device. This was done without the interviewer revealing anything about the system or its intended use, other than "it is used for skimming speech recordings." Everything in the test was exploratory; subjects were not given any instructions or guidance.⁸ The subjects were asked what they thought the different regions of the device did, how they expected the system to behave, what they thought backward did, etc.
- (3) *Listening to a trial speech recording with the SpeechSkimmer system.* The subjects were encouraged to explore and "play" with the device to confirm, or discover, how the system operated. While investigating the device, the interviewer encouraged the subjects to "think aloud," to describe what they were doing, and to say if the device was behaving as they expected.
- (4) *A skimming comparison and exercise.* This portion of the test compared two different skimming techniques. A recording of a 40-minute lecture⁹ was divided into two 20-minute parts (half of the subjects had attended the lecture when it was originally presented). Each subject listened to both halves of the recording; one part was segmented using the pitch-based emphasis detector (Section 2.3.2); the other was segmented isochronously (i.e., at equal time intervals). All SpeechSkimmer controls were active under both conditions; users could change speed, skimming level, jump, and so on, the only difference was in the top-level segmentation. The test was counterbalanced for effects of presentation order and portion of the recording (Figure 15).

When skimming was used, both of the segmentation techniques provided a 12:1 compression for this recording (i.e., on average five seconds out of each minute were presented). Note that these figures are

⁷Lisa Stifelman conducted the test; the system designer (Arons) observed.

⁸However, if a subject said something like "I wish it did X," and the system did perform that function, the feature was revealed to them by the interviewer through directed questions (e.g., "Do you think this device can do that? If so, how do you think you could get it to do it? What do you think that button does?").

⁹Nicholas Negroponte's talk titled *Confusion: Media in the Next Millennium* presented October 19, 1993.

# of subjects	first presentation	second presentation
3	pitch-based part 1	isochronous part 2
3	isochronous part 1	pitch-based part 2
3	isochronous part 2	pitch-based part 1
3	pitch-based part 2	isochronous part 1

Fig. 15. Counterbalancing of experimental conditions.

for normal speed (1.0×); by using time compression the subjects could achieve over a 25:1 time savings.

The subjects first skimmed the entire recording at whatever speed they felt most comfortable. The subjects were asked to judge (on a seven-point scale) how well they thought the skimming technique did at providing an overview and selecting indexes into major points in the recording. The subjects were then given a printed list of three questions that could be answered by listening to the recording. The subjects were asked to locate the answer to any of the questions in the recording and to describe their auditory search strategy. This process was repeated for the second presentation condition.

- (5) The test concluded with followup questions regarding the subject's overall experience with the interaction device and the SpeechSkimmer system, including what features they liked or disliked, what they thought was missing from the user interface, etc.

4.2 Results and Discussion

This section summarizes the features of SpeechSkimmer that were frequently used or most liked by the subjects, as well as areas for improvement in the user interface design.

4.2.1 Background Interviews. All subjects had some experience in searching for recorded audio information on compact discs, audio cassettes, or video tape. Subjects' experience included transcribing lectures and interviews, taking personal notes on a microcassette recorder, searching for favorite songs on tape or compact disc, editing video documentaries, and receiving up to 25 voice mail messages per day. Almost all the subjects referred to the process of searching in audio recordings as time consuming; one subject added that it takes "more time than you want to spend."

4.2.2 First Intuitions. Most subjects found the interface intuitive and easy to use and were able to use the device without any training. This ability to quickly understand how the device works is partially because the touchpad controls are labeled similarly to consumer devices such as compact disc players and video cassette recorders. While this familiarity allowed the subjects to initially feel comfortable with the device, and enabled rapid acclimatization to the interface, it also caused some confusion, since a few of the SpeechSkimmer functions behave differently than on the consumer devices.

Level 2 on the skimming template is labeled “no pause,” but most subjects did not have any initial intuitions about what it meant. The label baffled most subjects, since current consumer devices do not have pause removal or similar functionality. Some subjects thought that once they started playing in “no pause” they would not be able to stop or pause the playback. Similarly, the function of “jump and play normal button” was not obvious. The backward play levels were sometimes intuitively equated with traditional (unintelligible) rewind.

4.2.3 Warmup Task. The recording used in the trial task consisted of a loose free-form discussion, and most subjects had trouble following the conversation. Most said that they would have been able to learn the device in less time if the trial recording was more coherent or if they were already familiar with the recording. However, subjects still felt the device was easy to learn quickly.

Subjects were not sure how far the jumps took them. Several subjects thought that the system jumped to the next utterance of the male talker when exploring the interface in the trial task (the first few segments selected for jumping in this recording did occur at a change of talker because the pause-based segmentation algorithm was used).

4.2.4 Skimming. Most subjects found that the pitch-based skimming was effective at extracting interesting points to listen to and for finding information. One user who does video editing described it as “grabbing sound-bite material.” When comparing pitch-based skimming to isochronous skimming a subject said “it is like using a rifle versus a shotgun” (i.e., high accuracy instead of dispersed coverage). Other subjects said that the pitch-based segments “felt like the beginning of a phrase . . . [and were] more summary oriented” and that there was “a lot more content or keyword searching going on” than in the isochronous segmentation.

A few subjects requested that longer segments be played (perhaps until the next pause) or that the length of the segments could be controllable. One subject said “I felt like I was missing a lot of his main ideas, since it would start to say one, and then jump.”

Subjects were asked to rank the skimming performance under the different segmentation conditions. A score of 7 indicates the best possible summary of high-level ideas; a score of 1 indicates very poorly selected segments. The mean score for the pitch-based segmentation was $M = 4.5$ ($SD = 1.7$, $N = 12$); the mean score for the isochronous segmentation was $M = 2.7$ ($SD = 1.4$, $N = 12$). The pitch-based skimming was rated better than isochronous skimming with a statistical significance of $p < 0.01$ using a t test for paired samples. No statistically significant difference was found on how subjects rated the first versus the second part of the talk or on how subjects rated the first versus second sound presented.

Most subjects, including the two that did not think the pitch-based skimming gave a good summary, used the top skimming level to navigate through the recording. When asked to find the answers to questions on the printed list, most started off by saying something like “I’ll go to the

beginning and skim till I get to the right topic area in the recording,” or in some cases “I think its near the end, so I’ll jump to the end and skim backward.”

4.2.5 *No Pause.* While there was some initial confusion regarding the level “no pause,” if a subject discovered its function, it often became a preferred way to quickly listen and search for information. One subject that does video editing said “that’s nice . . . I like the no-pause function . . . it kills dead time between people talking . . . this would be really nice for interviews [since you normally have to] remember when he said [the point of interest]; then you can’t find where it was, and must do a binary search of the audio track . . . For interviews it is all audio—you want to get the sound bite.”

4.2.6 *Jumping.* The function of the button “jump and play normal” was not always obvious. However, subjects that did not understand the button found ways to navigate and perform the same function using the basic controls. This button is a short-cut: a combination of jumping backward and then playing level 1 speech at regular speed. One subject had a moment of inspiration while skimming along at a high speed and tried the button after passing the point of interest. After using this button the subject said in a confirming tone “I liked that, OK.” The subject proceeded to use the button several more times and said “now that I figured out how to do that jump-normal thing . . . that’s very cool. I like that.” It is important to note that after discovering the button “jump and play normal” this subject felt more comfortable skimming at faster speeds. Another subject said “that’s the most important button if I want to find information.”

While most of the subjects used, and liked, the jump buttons, the size or granularity of jumps was not obvious. Subjects assumed that jumping always brought them to the next sentence or topic (in the SpeechSkimmer prototype the granularity of a jump depends on the current skimming level). While using the jump button and the level “backward no pause,” one subject said “Oh, I see the difference . . . I can relisten using the jump key.”

4.2.7 *Backward.* Most subjects figured out the backward controls during the warmup trial, but avoided using them. This is partially attributable to the subject’s initial mental models that associate backward with the conventional rewind of a tape player. Some subjects, however, did find the backward levels useful in locating particular words or phrases that had just been heard.

While listening to the recording played backward, one subject noted “It’s taking units of conversation—and goes backwards.” Another subject said that “It’s interesting that it is so seamless” for playing intelligible segments and that “Compared to a tape where you’re constantly shuffling back and forth, going backward and finding something was much easier, since [while] playing backwards you can still hear the words.” One subject

suggested providing feedback to indicate when the recording was being played backward, to make it more easily distinguishable from forward.

4.2.8 Time Compression. Some subjects thought there were only three discrete speeds and did not initially realize that there was a continuum of playback speeds. A few subjects did not realize that the ability to change speeds extended across all the skimming levels. These problems can be attributed to the three speeds marked on the template (slow, regular, and fast; Figure 11). One subject noted that the tactile strips on the surface break the continuity of the horizontal “speed” lines and made it less clear that the speeds work at all skimming levels. Two subjects suggested using colors to denote the continuum of playback speeds and that the speed labels should extend across all the skimming levels.

Several subjects thought there was a major improvement when listening over headphones. One subject was “really amazed” at how much better the dichotic time-compressed speech was for comprehension than the monotic speech presented over the loudspeaker. Another subject commenting on the dichotic speech said “It’s really interesting—you can hear it a lot better.”

4.2.9 Buttons. The buttons were generally intuitive, but there were some problems of interpretation and accidental use. The “begin” and “end” regions were initially added next to the level 3 and -3 skimming regions on the template to provide a continuum of playback granularity (i.e., normal, no pause, skim, jump to end). Several subjects thought that the begin button should seek to the beginning of the recording and start playing (the prototype seeks to the beginning and waits for user input). One subject additionally thought the speed of playback could be changed by touching at the top or bottom of the begin button.

One subject wanted to skim backward to rehear the last segment played, but accidentally hit the adjacent begin button instead. This frustrated the subject, since the system jumped to the beginning of the recording and hence lost the location of interest. Note also that along with these conceptual and mechanical problems, the words “begin” and “start” are overloaded and could mean “begin playing” as well as “seek to the beginning of the recording.” By far the biggest problem encountered during the usability test was “bounce” on the jump and pause buttons.¹⁰ This was particularly aggravating when it occurred with the pause button, as the subject would want to stop the playback, but the system would temporarily pause, then moments later un-pause. The bounce problem was partially exacerbated by the subjects’ use of their thumbs to touch the buttons. While the touchpad and template were designed to be operated with a single finger for maximum accuracy (as in Figure 10), most of the subjects held the touchpad by the right and left sides and touched the surface with their thumbs during the test. This is partially attributable to the arrangement of

¹⁰Button “bounce” is usually associated with mechanical switches that made several temporary contact closures before settling to a stable state. The difficulties here are associated with the way the touchpad driver software was configured.

the subject and the experimenters during the test. Subjects had to hold the device, as there was no table for placing the touchpad.

4.2.10 *Nonspeech Feedback.* The nonspeech audio was successful at unobtrusively providing feedback. One subject, commenting on the effectiveness and subtlety of the sounds said “After using it for a while, it would be annoying to get a lot of feedback.” Another subject said that the nonspeech audio “helps because there is no visual feedback.” None of the subjects noted that the frequency of the feedback tone changes with skimming level; most did not even notice the existence of the tones. However, when subsequently asked about the device many noted that the tones were useful feedback to what was going on. The cartoon “boings” at the beginning and end were good indicators of the end points (one subject said “it sounds like you hit the edge”), and the other sounds were useful in conveying that something was going on. The “boing” sounds were noticed most often, probably because the speech playback stops when the sound effect is played.

4.2.11 *Search Strategies.* Several different navigation and search strategies were used when trying to find the point in the recording that answered a question on the printed list. Most subjects skimmed (level 3) the recording to find the general topic area of interest, then changed to level 1 (playing) or level 2 (pauses removed), usually with time compression. One subject started searching by playing normally (no time compression) from the beginning of the recording to “get a flavor” for the talk before attempting to skim or play it at a faster rate. One subject used a combination of skimming and jumping to quickly navigate through the recording and efficiently find the answers to the list of questions.

4.2.12 *Followup Questions.* Most subjects thought that the system was easy to use, since they made effective use of the skimming system without any training or instructions. Subjects rated the ease of use of the system on a seven-point scale where 1 is difficult to use; 4 is neutral; and 7 is very easy to use. The mean score for ease of use was $M = 5.4$ ($SD = 0.97$, $N = 10$).

Most subjects liked the ability to quickly skim between major points in a presentation and to jump on demand within a recording. Subjects liked the time compression range, particularly the interactive control of the playback speed. A few subjects were enamored with other specific features of the system including the “fast-forward no pause” level, the “jump and play normal” button, and the dichotic presentation.

One subject commented “I really like the way it is laid out. It’s easier to use than a mouse.” Another subject experimented with turning the touchpad 90 degrees so that moving a finger horizontally, rather than vertically, changed the playback speed.

Most subjects said they could envision using the device while doing other things, such as walking around, but few thought they would want to use it while driving an automobile. Most of the subjects said they would like to

use such a device, and many of them were enthusiastic about the SpeechSkimmer system.

4.2.13 *Desired Functionality.* In the followup portion of the test, the subjects were asked what other features might be helpful for the speech-skimming system. For the most part these items were obtained through probing the subjects and were not mentioned spontaneously.

Some subjects were interested in marking points in the recording that were of interest to them, so they could go back and easily access those points later. A few of the subjects called these “bookmarks.”

Some subjects wanted to be able to jump to a particular place in a recording or have a graphical indicator of their current location. There is a desire, for example, to access a thought discussed “about three-quarters the way through the lecture” by using a “time line” for locating a specific time point.

4.3 Comments on Usability Testing

Informal heuristic evaluation of the interface [Jeffries et al. 1991; Nielsen 1991; Nielsen and Molich 1990] was performed throughout the system design. In addition, the test described in Section 4.1 was very helpful in finding usability problems. The test was performed relatively late in the SpeechSkimmer design cycle, and in retrospect, a preliminary test should have been performed much earlier. Most of the problems in the template layout could have been easily uncovered with only a few subjects. This could have led to a more intuitive interface, while focusing on the features most desired by users.

Note that while 12 subjects were tested here, only a few are needed to get helpful results. Nielsen has shown that the maximum cost-to-benefit ratio for a usability project occurs with around three to four test subjects and that even running a single test subject is beneficial [Nielsen 1993b].

5. REVISING THE SKIMMING INTERFACE

Note again that the usability test was performed without any instruction or coaching of the subjects. It may be easy to fix most of the usability problems by modifying the touchpad template or through a small amount of instruction.

5.1 Thoughts for Redesign Based on the Usability Test

After establishing the basic system functionality, the touchpad template evolved quickly. Figure 16 shows three early templates as well as the one used in the usability test. The “sketch” in Figure 17 shows a revised design that addresses many of the usability problems encountered, and it incorporates the new features requested. The labels and icons are modified to be more consistent and familiar. Notably, “play” has replaced “normal,” and “pauses removed” has replaced the confusing “no pause.”

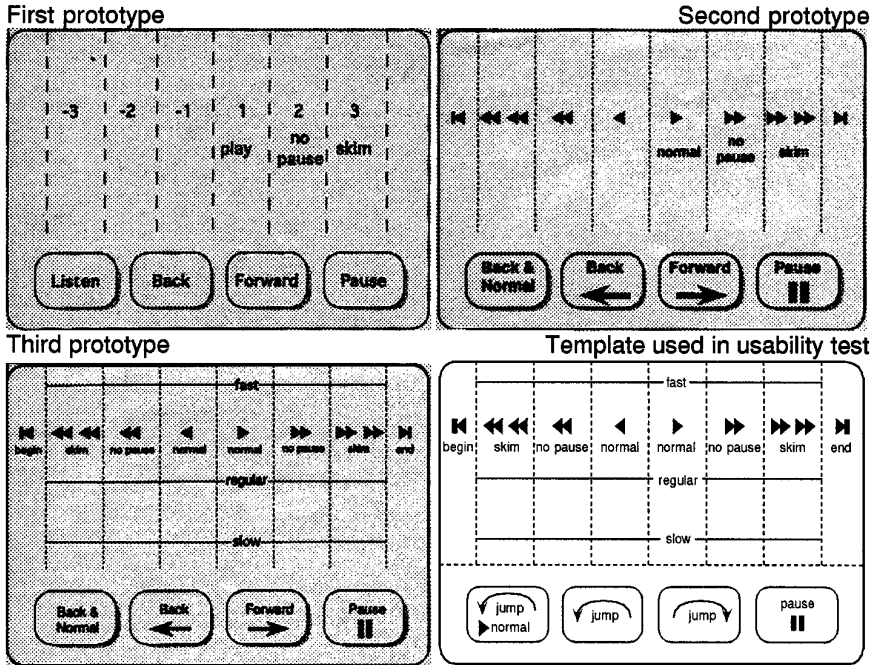


Fig. 16. Early evolution of SpeechSkimmer templates.

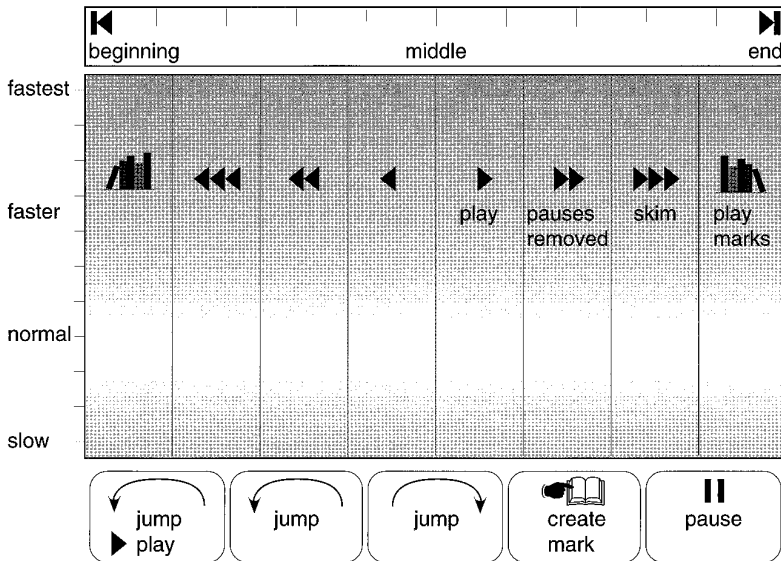


Fig. 17. Sketch of a revised template based on the usability results.

The speed labels are moved, renamed, and accompanied by tick marks to indicate a continuum of playback rates. The shaded background is an additional cue that the speeds extend across all levels. Colors, however, may be more effective than shading. For example, the slow-to-normal range

could fade from blue to white, while the normal-to-fastest range could go from white to red, suggesting a cool-to-hot transition.

Bookmarks, as requested by the subjects, can be implemented in a variety of ways, but are perhaps best thought of as yet another level of skimming. In this case, however, the user interactively creates the list of speech segments to be played. In this design a button “create mark” is added along with new regions for playing the user-defined segments.

A time line is added to directly access time points within a recording. It is located at the top of the template where subjects pointed when talking about this feature. The time line also naturally incorporates the begin and end buttons, removing them from the main portion of the template and out of the way of accidental activation.

The layout and graphic design of this template is somewhat cluttered, and the button “jump and play normal” remains problematic. However, the intuitiveness of this design, or alternative designs, could be quickly tested by asking a few subjects for their initial impressions.

One of the subjects commented that a physical control (such as real buttons and sliders) would be easier to use than the touchpad. Another approach to changing the physical interface is to use a jog and shuttle control, as is often found in video editing systems. Alternatively, a foot pedal could be used in situations where the hands are busy, such as when transcribing or taking notes.

5.2 A New Interface to SpeechSkimmer

A new user interface based on the results of the usability test, and the design sketched in Section 5.1, was implemented using an Apple Newton MessagePad 100 as an input and output device. The MessagePad has a digitizing surface and a graphics display, so it can be used both as an input device and for presenting status information. The touch-sensitive surface works with a stylus or a fingernail (Figure 18) rather than the tip of a finger as with the original touchpad. The MessagePad is rotated 90 degrees from its normal orientation into a landscape configuration to provide more screen real estate for the skimming controls.

Why use a touchpad with a display instead of a traditional screen and mouse? A touchpad was originally selected as an input device so that the system could be used without looking at it, or while doing other things. While the MessagePad interface does display a small amount of status information, it can still be used without looking at it (especially when tactile strips are added to the surface). The MessagePad provides an input and display mechanism in a small portable package that is designed to be handheld.

5.2.1 Time Line for Input and Output. A time line (Figure 19, top) is used both for displaying the current location within a recording and for going to a particular time point. The current position in the speech recording is shown using a small vertical bar in the time line. This bar stays synchronized with the recording and moves slowly as the audio plays,

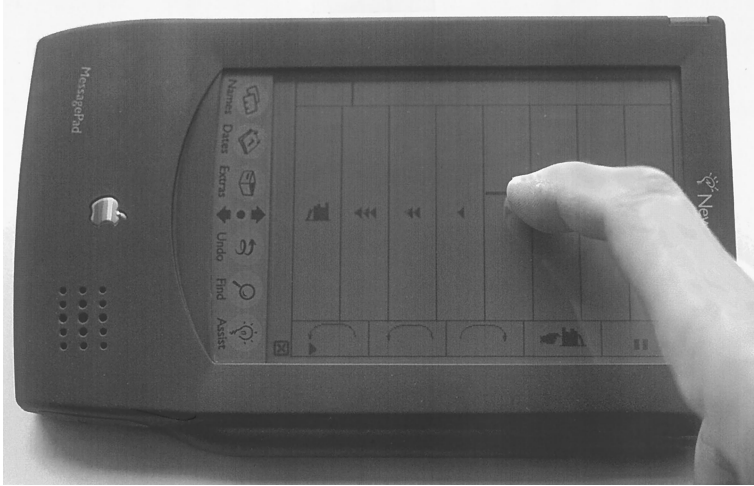


Fig. 18. An early version of the MessagePad interface.

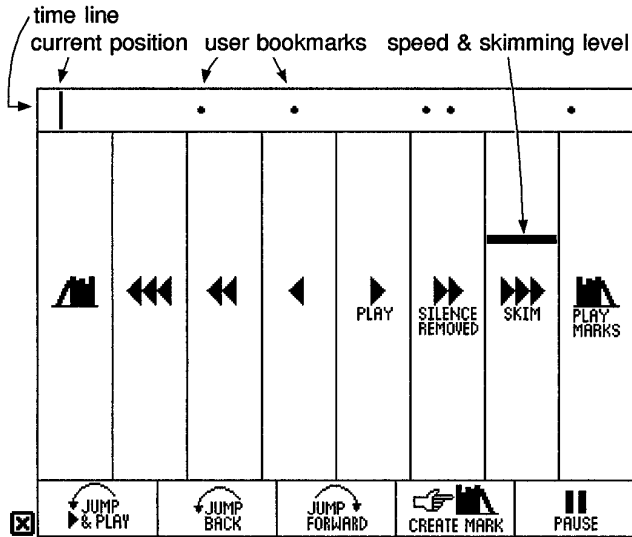


Fig. 19. Screen image from the MessagePad interface.

acting as a percent-done indicator [Myers 1985]. The time line can also be touched to jump to a specific point in the recording.

5.2.2 Bookmarks. A listener can set a personalized bookmark at any point in a recording by touching the button “create mark” (Figure 19, bottom). This causes two events to happen. First, a visual indication for a bookmark (a small circle) is added to the time line, allowing the listener to get a sense of his or her location within the recording. Second, a new speech segment is added to SpeechSkimmer’s internal representation of audio to be played at the highest skimming level.

The bookmarks can be accessed manually by touching the circles in the time line or through the new skimming level “play marks.” This level plays only the segments selected by the user. Thus, this top skimming level represents a user-defined summary of the recording.

5.2.3 Display of Skimming Level and Speed. The playback speed is set by sliding a finger up or down in one of the vertical regions of the MessagePad. The current skimming level and speed are visually indicated by a horizontal bar in one of the slider regions (Figure 19). Note that as with the original touchpad interface only one of these virtual sliders can be selected at a time (i.e., only one speed and skimming level is active).

5.2.4 System Architecture. The MessagePad is connected to an Apple Macintosh computer that performs all speech processing and audio playback. A small amount of processing is done on the MessagePad to update the display and translate the raw coordinates from the digitizing tablet into higher-level events that are sent to the Macintosh (possibly over a wireless infrared link). Ideally, the entire speech-skimming system could be implemented on a portable device such as a MessagePad. However, current-generation handheld computers and PDAs (personal digital assistants) have limitations in the areas of audio input and output, sound data storage, and software tools for managing and manipulating audio.

6. RELATED WORK

6.1 Speech Interfaces

SpeechSkimmer draws many of its ideas from earlier speech systems. While SpeechSkimmer automatically structures recordings from information in the speech signal, many of these predecessor systems structure audio through interaction with the user, placing some burden on the creator of the speech data.

Phone Slave [Schmandt and Arons 1984] and the Conversational Desktop [Schmandt and Arons 1987; Schmandt et al. 1985] explored interactive message gathering and speech interfaces to simple databases of voice messages. VoiceNotes [Stifelman et al. 1993] investigated the creation and management of a self-authored database of short speech recordings. VoiceNotes investigated many of the user interface issues addressed by SpeechSkimmer in the context of a handheld computer.

Hyperspeech is a speech-only hypermedia system that explores issues of speech user interfaces, browsing, and the use of speech as data in an environment without a visual display [Arons 1991b]. The system uses speech recognition input and synthetic speech feedback to aid in navigating through a database of recorded speech. Resnick designed several voice bulletin board systems accessible through a touch tone interface [Resnick 1992; Resnick and Virzi 1992]. These systems addressed issues of navigation and provided shortcuts to “skip and scan” among speech recordings.

The “human memory prosthesis” was envisioned to run on a wireless notepad-style computer to help people remember things such as names and reconstruct past events [Lamming 1991]. Information gathered through a variety of sources, such as notetaking, permits jumping to timestamped points in the audio (or video) stream. Filochat indexed an audio recording with pen strokes on an LCD tablet, allowing the notes to be used to access the recording [Whittaker et al. 1994]. Stifelman’s Audio Notebook also synchronizes handwritten notes with an audio recording [Stifelman 1996]. However, rather than writing on a computer screen, the notes are taken with an ink pen in a paper notebook providing a familiar interface. Both handwriting and page turns are used as indices into the audio. Moran et al. [1996] captured meetings and indexed them through several notetaking tools.

6.2 Gisting and Skimming

Word-spotting and gisting (obtaining the essence of a message) systems are appealing for summarizing and accessing messages, but have limited domains of applicability; the skimming techniques presented here do not use any domain-specific knowledge and will work across all topics.

Several systems have attempted to obtain the gist of a recording using keyword spotting [Wilcox and Bush 1992; Wilpon et al. 1990] in conjunction with syntactic and/or timing constraints in an attempt to broadly classify a message [Houle et al. 1988; Maksymowicz 1990]. Rose’s system takes speech messages and extracts the message category according to a pre-defined notion of topic [Rose 1991]. Similar work has been reported in the areas of retrieving speech documents [Brown et al. 1996; Glavitsch and Schäuble 1992] and editing applications [Wilcox et al. 1992].

Maxemchuk [1980] suggested three techniques for skimming speech messages: using text descriptors for selecting playback points, jumping forward or backward in the message, and increasing the playback rate. Stevens and Edwards [1994] and Raman [1994] developed systems for reading and browsing math equations and structured documents with a text-to-speech synthesizer. These systems addressed issues of navigating in auditory documents and methods of presenting “auditory glances.”

6.3 Segmentation and Display

Kato and Hosoya [1992; 1993] investigated several techniques to enable fast telephone-based message searching by breaking up messages on hesitation boundaries and presented either the initial portion of each phrase or high-energy segments. Kimber et al. [1995] used speaker identification to segment audio recordings that could be browsed with a graphical interface. See Pfeiffer et al. [1996] and Hawley [1993] for attempts at automatically analyzing and structuring audio recordings.

Wolf and Rhyne [1992] present a method for reviewing meetings based on characteristics captured by a pen-based meeting support tool. They found that turn categories of most interest for browsing meetings were preceded

by longer gaps in writing than the other turn types. Several techniques for capturing and structuring office discussions, telephone conversations, and lengthy recordings are described in Hindus et al. [1993], which emphasized graphical representations of recordings in a workstation environment.

7. FUTURE WORK

In addition to the small amount of status information displayed on the MessagePad interface, the skimming system could also take advantage of a full graphical user interface for displaying information. Along with mapping the fundamental SpeechSkimmer controls to a mouse, it is possible to add a variety of visual cues, such as displaying a real-time version of the segmentation information (Figure 7), to aid in the skimming process.

Video editing and display systems can also be used with a speech-skimming interface. For example, when we quickly browse through a set of video images, only the high-level segments of speech could be played, rather than random snippets of audio associated with the displayed frames. Similarly, a SpeechSkimmer-like interface can be used to skim through the audio track while the related video images are synchronously displayed.

The automatic structuring of spontaneous speech is an important area for future work. Integrating multiple acoustic cues (e.g., pitch, energy, pause, speaker identification) will ultimately produce the most successful segmentation techniques. Word spotting can also be used to provide text tags or summaries for flexible information retrieval. Summarizing or gisting systems will advance as speech recognition technology evolves, but may be most useful when combined with the skimming ideas presented here.

8. CONCLUSION

Speech is naturally slow to listen to and difficult to skim. This research attempts to overcome these limitations by making it easier and more efficient to consume recorded speech. By combining techniques that extract structure from spontaneous speech, with a hierarchical representation and an interactive listener control, it is possible to overcome the time bottleneck in speech-based systems. When asked if the system was useful, one test subject commented “Yes, definitely. It’s quite nice. I would use it to listen to talks or lectures that I missed . . . It would be super. I would do it all the time. I don’t do it now, since it would require me to sit through the duration of the two-hour [presentations] . . .”

This article presents a framework for thinking about and designing speech-skimming systems. SpeechSkimmer allows “intelligent” filtering of recorded speech; the intelligence is provided by the interactive control of the human, in combination with the speech segmentation techniques. The fundamental mechanisms presented here allow other types of segmentation or new interface techniques to be easily plugged in. SpeechSkimmer is intended to be a technology that is incorporated into any interface that uses recorded speech, as well as a standalone application. Skimming techniques enable speech to be readily accessed in a range of applications and devices,

empowering a new generation of user interfaces that use speech. When discussing the SpeechSkimmer system, one of the usability test subjects put it succinctly: “it is a concept, not a box.”

This research provides insight into making one’s ears as usable as one’s eyes as a means for accessing stored information. Tufte said “Unlike speech, visual displays are simultaneously a wideband and a perceiver-controllable channel” [Tufte 1990, p. 31]. This work attempts to overcome these conventional notions, increasing the information bandwidth of the auditory channel and allowing the perceiver to interactively access recorded information. Speech is a powerful medium, and its use in computer-based systems will expand in unforeseen ways when users can interactively skim, and efficiently listen to, recorded speech.

ACKNOWLEDGMENTS

Lisa Stifelman provided valuable input in user interface design, assisted in designing and conducting the many hours of the usability test, and helped edit this document. Chris Schmandt provided helpful feedback on the SpeechSkimmer system. Doug Reynolds ran his speaker identification software on my recording. Dorée Seligmann suggested using the MessagePad. Two of the anonymous reviewers helped me focus this article. Many others have contributed to this work and have been thanked elsewhere.

REFERENCES

- ARONS, B. 1991a. Authoring and transcription tools for speech-based hypermedia systems. In *Proceedings of the 1991 Conference of the American Voice I/O Society*. American Voice I/O Society, 15–20.
- ARONS, B. 1991b. Hyperspeech: Navigating in speech-only hypermedia. In *Proceedings of Hypertext*. ACM, New York, 133–146.
- ARONS, B. 1992a. Techniques, perception, and applications of time-compressed speech. In *Proceedings of the 1992 Conference of the American Voice I/O Society*. American Voice I/O Society, 169–177.
- ARONS, B. 1992b. Tools for building asynchronous servers to support speech and audio applications. In *Proceedings of the ACM Symposium on User Interface Software and Technology (UIST)*. ACM Press, New York, 71–78.
- ARONS, B. 1994a. Interactively skimming recorded speech. Ph.D. dissertation, MIT, Cambridge, Mass.
- ARONS, B. 1994b. Pitch-based emphasis detection for segmenting speech recordings. In *Proceedings of the International Conference on Spoken Language Processing*. Vol. 4. Acoustical Society of Japan, Tokyo, 1931–1934.
- BEASLEY, D. S. AND MAKI, J. E. 1976. Time- and frequency-altered speech. In *Contemporary Issues in Experimental Phonetics*, N. J. Lass, Ed. Academic Press, New York, 419–458.
- BROWN, M. G., FOOTE, J. T., JONES, G. J. F., JONES, K. S., AND YOUNG, S. J. 1996. Open vocabulary speech indexing for voice and video mail retrieval. In *Proceedings of ACM Multimedia 96*. ACM, New York, 307–316.
- BUXTON, W., GAVER, B., AND BLY, S. 1991. The use of non-speech audio at the interface. Tutorial Notes. In *Proceedings of ACM SIGCHI*. ACM, New York.
- CARD, S. K., MACKINLAY, J. D., AND ROBERTSON, G. G. 1991. A morphological analysis of the design space of input devices. *ACM Trans. Inf. Sys.* 9, 2, 99–122.
- CHEN, F. R. AND WITHGOTT, M. 1992. The use of emphasis to automatically summarize spoken discourse. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*. IEEE, New York, 229–233.

- DAVIS, M. 1995. Media streams: An iconic visual language for video representation. In *Readings in Human-Computer Interaction: Toward the Year 2000*, R. M. Baecker, J. Grudin, W. A. S. Buxton, and S. Greenberg, Eds. 2nd ed. Morgan Kaufmann, San Francisco, Calif., 854–866.
- DE SOUZA, P. 1983. A statistical approach to the design of an adaptive self-normalizing silence detector. *IEEE Trans. Acoustics Speech Sig. Process.* ASSP-31, 3, 678–684.
- ELLIOTT, E. L. 1993. Watch-Grab-Arrange-See: Thinking with motion images via streams and collages. Master's thesis, Media Arts and Sciences Section, MIT, Cambridge, Mass.
- ERICSSON, K. A. AND SIMON, H. A. 1984. *Protocol Analysis: Verbal Reports as Data*. MIT Press, Cambridge, Mass.
- FAIRBANKS, G., EVERITT, W. L., AND JAEGER, R. P. 1954. Method for time or frequency compression-expansion of speech. *Trans. Inst. Radio Eng. Prof. Group Audio AU-2*, 7–12. Reprinted in G. Fairbanks, *Experimental Phonetics: Selected Articles*. University of Illinois Press, 1966.
- FOULKE, E. 1971. The perception of time compressed speech. In *Perception of Language*, P. M. Kjeldergaard, D. L. Horton, and J. J. Jenkins, Ed. Merrill, Columbus, Ohio, 79–107.
- FURNAS, G. W. 1986. Generalized fisheye views. In *Proceedings of CHI*. ACM, New York, 16–23.
- GAVER, W. W. 1989. Auditory icons: Using sound in computer interfaces. *Hum. Comput. Interact.* 2, 167–177.
- GERBER, S. E. AND WULFECK, B. H. 1977. The limiting effect of discard interval on time-compressed speech. *Lang. Speech* 20, 2, 108–115.
- GLAVITSCH, U. AND SCHAUBLE, P. A. 1992. A system for retrieving speech documents. In the *15th Annual International SIGIR '92*. ACM, New York, 168–176.
- GOULD, J. 1982. Writing and speaking letters and messages. *Int. J. Man Mach. Stud.* 16, 147–171.
- GROSZ, B. J. AND SIDNER, C. L. 1986. Attention, intentions, and the structure of discourse. *Comput. Ling.* 12, 3, 175–204.
- GRUBER, J. G. 1982. A comparison of measured and calculated speech temporal parameters relevant to speech activity detection. *IEEE Trans. Commun.* COM-30, 4, 728–738.
- GRUBER, J. G. AND LE, N. H. 1983. Performance requirements for integrated voice/data networks. *IEEE J. Sel. Areas Commun.* SAC-1, 6, 981–1005.
- HAWLEY, M. 1993. Structure out of sound. Ph.D. dissertation, MIT, Cambridge, Mass.
- HEIMAN, G. W., LEO, R. J., LEIGHBODY, G., AND BOWLER, K. 1986. Word intelligibility decrements and the comprehension of time-compressed speech. *Percep. Psychophys.* 40, 6, 407–411.
- HEJNA, D. J., JR. 1990. Real-time time-scale modification of speech via the synchronized overlap-add algorithm. Master's thesis, Dept. of Electrical Engineering and Computer Science, MIT, Cambridge, Mass.
- HINDUS, D., SCHMANDT, C., AND HORNER, C. 1993. Capturing, structuring, and representing ubiquitous audio. *ACM Trans. Inf. Syst.* 11, 4, 376–400.
- HIRSCHBERG, J. AND GROSZ, B. 1992. Intonational features of local and global discourse. In *Proceedings of the Speech and Natural Language Workshop*. Morgan Kaufmann, San Mateo, Calif., 441–446.
- HIRSCHBERG, J. AND PIERREHUMBERT, J. 1986. The intonational structuring of discourse. In *Proceedings of the Association for Computational Linguistics*. ACL, 136–144.
- HOULE, G. R., MAKSYMOWICZ, A. T., AND PENAFIEL, H. M. 1988. Back-end processing for automatic gisting systems. In *Proceedings of the 1988 Conference of the American Voice I/O Society*. American Voice I/O Society.
- JEFFRIES, R., MILLER, J. R., WHARTON, C., AND UYEDA, K. M. 1991. User interface evaluation in the real world: A comparison of four techniques. In *Proceedings of CHI*. ACM, New York, 119–124.
- KATO, Y. AND HOSOYA, K. 1992. Fast message searching method for voice mail service and voice bulletin board service. In *Proceedings of the 1992 Conference of the American Voice I/O Society*. American Voice I/O Society, 215–222.

- KATO, Y. AND HOSOYA, K. 1993. Message browsing facility for voice bulletin board service. In *Human Factors in Telecommunications '93*. 321–330.
- KIMBER, D., WILCOX, L., CHEN, F., AND MORAN, T. 1995. Speaker segmentation for browsing recorded audio. In *CHI '94 Conference Companion*. ACM, New York, 212–213.
- LAMEL, L. F., RABINER, L. R., ROSENBERG, A. E., AND WILPON, J. G. 1981. An improved endpoint detector for isolated word recognition. *IEEE Trans. Acoustics Speech Sig. Process. ASSP-29*, 4, 777–785.
- LAMMING, M. G. 1991. Towards a human memory prosthesis. Tech. Rep. EPC-91-116, Xerox EuroPARC, Cambridge, U.K.
- LASS, N. J. AND LEEPER, H. A. 1977. Listening rate preference: Comparison of two time alteration techniques. *Percep. Motor Skills* 44, 1163–1168.
- LEVELT, W. J. M. 1989. *Speaking: From Intention to Articulation*. MIT Press, Cambridge, Mass.
- LIPSCOMB, J. S. AND PIQUE, M. E. 1993. Analog input device physical characteristics. *SIGCHI Bull.* 25, 3, 40–45.
- MACKINLAY, J. D., ROBERTSON, G. G., AND CARD, S. K. 1991. The perspective wall: Detail and context smoothly integrated. In *Proceedings of CHI*. ACM, New York, 173–179.
- MAKSYMOWICZ, A. T. 1990. Automatic gisting for voice communications. In *IEEE Aerospace Applications Conference*. IEEE, New York, 103–115.
- MAXEMCHUK, N. 1980. An experimental speech storage and editing facility. *Bell Syst. Tech. J.* 59, 8, 1383–1395.
- MICROTOUCH. 1992. *UnMouse User's Manual*. Microtouch Systems, Wilmington, Mass.
- MILLS, M., COHEN, J., AND WONG, Y. Y. 1992. A magnifier tool for video data. In *Proceedings of CHI*. ACM, New York, 93–98.
- MINIFIE, F. D. 1974. Durational aspects of connected speech samples. In *Time-Compressed Speech*, S. Duker, Ed. Scarecrow, Metuchen, N.J., 709–715.
- MORAN, T., CHIU, P., HARRISON, S., KURTENBACH, G., MINNEMAN, S., AND VAN MELLE, W. 1996. Evolutionary engagement in an ongoing collaborative work process: A case study. In *Proceedings of the ACM 1996 Conference on Computer Supported Cooperative Work*. ACM, New York, 150–159.
- MYERS, B. A. 1985. The importance of percent-done progress indicators for computer-human interfaces. In *Proceedings of the ACM CHI '85 Conference on Human Factors in Computing Systems*. ACM, New York, 11–17.
- NEUBURG, E. P. 1978. Simple pitch-dependent algorithm for high quality speech rate changing. *J. Acoustic Soc. Am.* 63, 2, 624–625.
- NIELSEN, J. 1991. Finding usability problems through heuristic evaluation. In *Proceedings of CHI*. ACM, New York, 373–380.
- NIELSEN, J. 1993a. *Usability Engineering*. Academic Press, San Diego, Calif.
- NIELSEN, J. 1993b. Is usability engineering really worth it? *IEEE Softw.* 10, 6, 90–92.
- NIELSEN, J. AND MOLICH, R. 1990. Heuristic evaluation of user interfaces. In *Proceedings of CHI*. ACM, New York.
- O'SHAUGHNESSY, D. 1987. *Speech Communication: Human and Machine*. Addison-Wesley, Reading, Mass.
- O'SHAUGHNESSY, D. 1992. Recognition of hesitations in spontaneous speech. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*. IEEE, New York, 1521–1524.
- PFEIFFER, S., FISCHER, S., AND EFFELSBERG, W. 1996. Automatic audio content analysis. In *Proceedings of ACM Multimedia 96*. ACM, New York, 21–30.
- RABINER, L. R. AND SAMBUR, M. R. 1975. An algorithm for determining the endpoints of isolated utterances. *Bell Syst. Tech. J.* 54, 2, 297–315.
- RABINER, L. R. 1989. A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. IEEE* 77, 2, 257–286.
- RAMAN, T. V. 1994. AsTeR: Audio system for technical readings. Ph.D. dissertation, Cornell Univ., Ithaca, N.Y.
- REICH, S. S. 1980. Significance of pauses for speech perception. *J. Psycholing. Res.* 9, 4, 379–389.

- RESNICK, P. 1992. HyperVoice: Groupware by telephone. Ph.D. dissertation, MIT, Cambridge, Mass.
- RESNICK, P. AND VIRZI, R. A. 1992. Skip and scan: Cleaning up telephone interfaces. In *Proceedings of CHI*. ACM, New York, 419–426.
- REYNOLDS, D. A. AND ROSE, R. C. 1995. Robust text-independent speaker identification using Gaussian mixture speaker models. *IEEE Trans. Speech Audio Process.* 3, 1, 72–83.
- ROE, D. B. AND WILPON, J. G. 1993. Whither speech recognition: The next 25 years. *IEEE Commun. Mag.* 31, 11, 54–62.
- ROSE, R. C. 1991. Techniques for information retrieval from speech messages. *Lincoln Lab. J.* 4, 1, 45–60.
- ROUCOS, S. AND WILGUS, A. M. 1985. High quality time-scale modification for speech. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*. IEEE, New York, 493–496.
- SAVOJI, M. H. 1989. A robust algorithm for accurate endpointing of speech signals. *Speech Commun.* 8, 45–60.
- SCHMANDT, C. AND ARONS, B. 1984. A conversational telephone messaging system. *IEEE Trans. Consumer Electron. CE-30*, 3, xxi–xxiv.
- SCHMANDT, C. AND ARONS, B. 1987. Conversational desktop. *ACM SIGGRAPH Video Rev.* 27. Videotape.
- SCHMANDT, C., ARONS, B., AND SIMMONS, C. 1985. Voice interaction in an integrated office and telecommunications environment. In *Proceedings of the 1985 Conference of the American Voice I/O Society*. American Voice I/O Society.
- SCOTT, R. J. 1967. Time adjustment in speech synthesis. *J. Acoustic Soc. Am.* 41, 1, 60–65.
- SILVERMAN, K. E. A. 1987. The structure and processing of fundamental frequency contours. Ph.D. dissertation, Univ. of Cambridge, Cambridge, U.K.
- STEVENS, R. AND EDWARDS, A. 1994. Mathtalk: The design of an interface for reading algebra using speech. In *Computers for Handicapped Persons: Proceedings of ICCHP '94*. Lecture Notes in Computer Science, vol. 860. Springer-Verlag, Berlin, 313–320.
- STIFELMAN, L. 1994. A study of rate discrimination of time-compressed speech. *J. Am. Voice I/O Soc.* 16, 69–81.
- STIFELMAN, L. 1996. Augmenting real-world objects: A paper-based audio notebook. In *CHI '96 Conference Companion*. ACM, New York, 199–200.
- STIFELMAN, L. J. 1995. A discourse analysis approach to structured speech. In *Empirical Methods in Discourse Interpretation and Generation*. AAAI, Menlo Park, Calif., 162–167.
- STIFELMAN, L. J., ARONS, B., SCHMANDT, C., AND HULTEEN, E. A. 1993. VoiceNotes: A speech interface for a hand-held voice notetaker. In *Proceedings of INTERCHI*. ACM, New York, 179–186.
- TUFTE, E. 1990. *Envisioning Information*. Graphics Press, Cheshire, Conn.
- WHITTAKER, S., HYLAND, P., AND WILEY, M. 1994. Filochat: Handwritten notes provide access to recorded conversations. In *Proceedings of CHI*. ACM, New York, 271–277.
- WILCOX, L. AND BUSH, M. 1992. Training and search algorithms for an interactive wordspotting system. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*. IEEE, New York.
- WILCOX, L., SMITH, I., AND BUSH, M. 1992. Wordspotting for voice editing and audio indexing. In *Proceedings of CHI*. ACM, New York, 655–656.
- WILPON, J. G., RABINER, L. R., LEE, C., AND GOLDMAN, E. R. 1990. Automatic recognition of keywords in unconstrained speech using hidden Markov models. *IEEE Trans. Acoustics Speech Sig. Process.* 38, 11, 1870–1878.
- WOLF, C. G. AND RHYNE, J. R. 1992. Facilitating review of meeting information using temporal histories. Working Paper 9/17, IBM T. J. Watson Research Center, Yorktown Heights, N.Y.

Received June 1996; revised December 1996; accepted December 1996