

Interaction Techniques Using Prosodic Features of Speech and Audio Localization

Alex Olwal^{1,2}

Steven Feiner¹

¹ Department of Computer Science
Columbia University
New York, USA

² Department of Numerical Analysis and Computer Science
Royal Institute of Technology
Stockholm, Sweden

alx@kth.se, feiner@cs.columbia.edu

ABSTRACT

We describe several approaches for using prosodic features of speech and audio localization to control interactive applications. This information can be applied to parameter control, as well as to speech disambiguation. We discuss how characteristics of spoken sentences can be exploited in the user interface; for example, by considering the speed with which a sentence is spoken and the presence of extraneous utterances. We also show how coarse audio localization can be used for low-fidelity gesture tracking, by inferring the speaker's head position.

Categories and Subject Descriptors: H.5.2 [User Interfaces]: Graphical user interfaces, natural language, voice I/O; I.2.7 [Natural Language Processing].

General Terms: Human Factors.

Keywords: Speech, gesture, interaction, voice I/O.

1 INTRODUCTION

We describe a set of nonverbal metrics of speech for use as additional parameters in speech-based interaction. This information allows an application to react explicitly or implicitly to characteristics of the user's speech.

Previously, Igarashi and Hughes [2] showed how duration, pitch and tonguing of nonverbal voice can be used for interactive application control, while Tsukahara and Ward [6] explored the use of prosodic features for appropriate emotional computer response in human-computer dialogues. We extend this work by examining how prosodic features of *verbal voice*, such as speech rate, duration and volume, can be used to control graphical user interfaces, and how audio localization can expand user expressiveness, as well as help resolve ambiguities in speech recognition. We introduce methods for speech-based cursor control and 3D manipulation that use these metrics. After describing our nonverbal metrics, we present our conclusions and plans for future work.

2 SPEECH-BASED CURSOR CONTROL

Speech-based cursor control can make it possible for individuals who are physically disabled, or temporarily unable to use a keyboard or mouse, to interact with a traditional 2D graphical user interface. Problems with previous work on speech-based cursor control include precision and user strain. Karimullah and Sears [3]



Figure 1. A desktop user manipulating a 3D model with speech commands, simple head gesture (tracked with audio) and speech rate. Rotation (shown at right) occurs after speech recognition.

tried to address the lack of precision with a predictive cursor that shows the estimated cursor position for a command after its speech-recognition delay (the time from when the command was spoken to its execution), but they concluded that it did not provide the expected benefits. They found that cursor speed, target size, and speech recognition delays and errors were the factors that were most critical in speech-based cursor control. We address these problems by controlling cursor speed through speech rate (i.e., how fast the user issues a spoken command), in addition to providing visual feedback through a predictive cursor [3].

We are experimenting with nonverbal features in a prototype system in which the cursor speed and direction are controlled by speech commands, speech rate, and the user's position extrapolated from their speech, as shown in Figure 1. In one approach, speech commands provide the direction (right, left, up, and down) and speech rate controls the cursor speed. Mapping speech rate to cursor speed is easy to understand and allows the user to execute slow, high-precision cursor movements by issuing commands at a slower pace, and to move the cursor quickly through fast speech (e.g., "Moouooooove leeeeeeeft!" vs. "Move left!"). The cursor's speed can be changed while it is moving, by reissuing the command at a different pace. Figure 2 illustrates how the cursor speed is changed, after speech processing, when the same command is issued with varying speech rate.

In a second approach, the user provides directional information by leaning to the left or right (Figure 3a), and we use simple audio localization to determine the side to which the user is leaning, as described below. If this second approach is used exclusively (we permit the simultaneous use of both approaches in our testbed), it could make it possible to use a smaller grammar and thus potentially improve speech recognition.

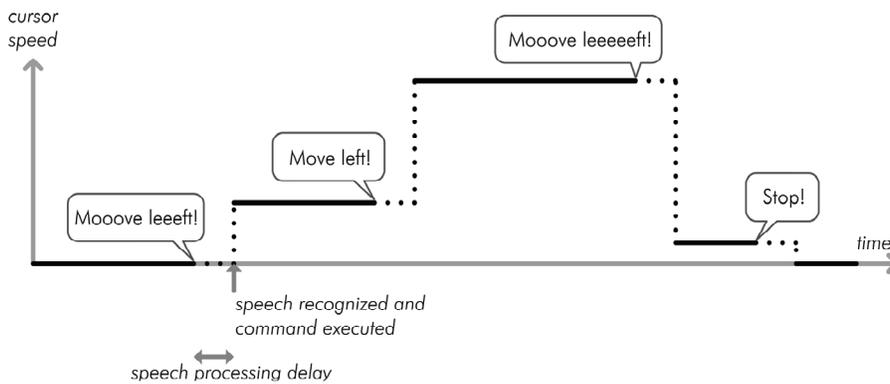


Figure 2. A series of user interactions with our approach to speech-based cursor control. The speed of the continuous cursor movement is mapped to speech rate—the faster the user speaks the command, the faster the cursor moves. Speech processing delays are indicated by dotted linestyle. The user can at any time adjust the cursor speed by reissuing the command at a different pace, as shown. Thus, the user can slow the cursor as it approaches the target and minimize the effects of overshooting due to the delay in speech processing.

3 SPEECH-BASED OBJECT MANIPULATION

We are also experimenting with the use of nonverbal speech features for object manipulation, such as rotation. Consider, for example, rotation around the axis perpendicular to the screen. A head gesture to either side (identified through audio tracking) corresponds to rotation around that axis, and we thus found it appropriate to map the operation to speech and head gesture, as shown in Figures 1 and 3(b). We find this more intuitive than specifying such an operation entirely with speech. The speech command, speech rate and the user’s position are used to determine the operation (rotation), rotational speed, and direction of rotation, respectively.

4 NONVERBAL METRICS

Our prototypes use a set of new interaction techniques for controlling interactive applications, based on nonverbal features in the user’s speech. In contrast to previous work [2], we consider characteristics of the verbal sentence.

4.1 Rate

We approximate the speech rate as the number of spoken syllables per second. This metric is independent of what is said, and indicates how fast the user spoke. We use the speech rate to differentiate sentences that are spoken at varying speed (e.g., “Zoom in!” vs. “Zoououououm iiiiiiiiiiiin!” where higher speech rate could yield a larger zoom step).

4.2 Duration

We define the duration of a sentence as the time it takes the speaker to speak it. In contrast to speech rate, this metric depends on what is said, since it is not normalized. Thus, we consider duration only when different sentences with similar meaning are compared. We use duration to assign different meanings to different sentence formulations (e.g., “Move that there!” vs. “Move that object over there, please!” where the object could move slower if a longer sentence formulation is chosen).

4.3 Volume

We currently calculate two volume metrics, the average and maximum volume level from the start to the end of the sentence. Volume distinguishes sentences that are spoken with different loudness (e.g., “Rotate left!” vs. “ROTATE LEFT!” where louder speech could yield faster rotation).

4.4 Position

Without the need for a separate head tracker, the user’s head position can be approximated by the originating direction of the speech. We use coarse audio localization to distinguish sentences that are spoken from different directions. Position data can be used to make assumptions about user gestures, and have the application react accordingly. For example, in car racing games, players instinctually lean to the left when they want to turn left quickly. The application could in this case let the car turn more when it detects that the user is leaning to either side. The disadvantage is that audio (verbal or nonverbal) is always required in this scenario unless other tracking mechanisms are employed.

We can also use audio localization to improve recognition rate by taking into account redundant information from speech and gesture. Our speech recognizer (and most other available recognizers) provides an ordered n -best list over recognized speech. If such a list contains “Move left” and “Move in,” and the user was leaning to the left during the speech, mutual disambiguation [5] might be used to pick “Move left” over “Move in.” Audio localization can not only help reduce errors in speech recognition, but can also reduce cognitive load (and thus, potentially, user errors), by combining speech and gesture. For example, the user could speak “Move” and lean to the left, to perform a “Move left” action, as shown in Figure 3.

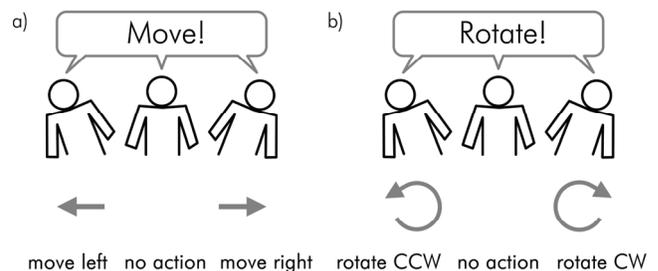


Figure 3. a) The user, as seen from behind, specifies the direction of a translation action by leaning in the corresponding direction. b) Leaning to the left or right can be interpreted as a rotation around a pivot-point on the user’s torso, which we map to the direction of the rotation action: counterclockwise or clockwise, respectively.

5 IMPLEMENTATION

Our prototypes are implemented in Java and communicate locally with IBM ViaVoice 10 (through the Java Speech API), which provides the necessary data for computing speech rate, duration and average volume. For audio localization, we use inexpensive omnidirectional microphones instead of special-purpose hardware (e.g., array microphones that provide localization information). In our experimental setup, we use two microphones, one each on the left and right side of a flat-panel display. Each microphone is connected to a separate computer running the speech recognition software, which computes all metrics except position. The Unit dataflow framework [4] is used to communicate by Ethernet the recognized speech, the computed metrics, and the position of the associated microphone, to an application server that provides audiovisual feedback (Figure 4). Speech position is computed in the application server as the difference between the two average volumes, and if one is significantly higher than the other (i.e., the difference is above a set threshold), we assume that the user is closer to this microphone. The recognition result is taken from this source and its associated position is used to indicate the user's position. (The user is assumed to be in a neutral, centered position if no significant difference is detected.) We used coarse discrete spatial classification because our simple experimental setup is not sufficiently robust to support reliable continuous localization. Despite its lack of sophistication, this approach allowed us to rapidly test this low-fidelity audio tracker and our associated interaction techniques.

6 CONCLUSIONS AND FUTURE WORK

We have introduced a set of new interaction techniques based on spatial audio and prosodic features in speech. We show that even very simple speech analysis can increase interaction bandwidth, and that spatial audio can expand user expressiveness in speech-based applications.

One limitation of using speech features is that they are normally used to convey emotion, rather than for interaction control. Another limitation is that our simple volume metrics do not distinguish between changes in the user's proximity to a microphone and changes in the volume of the speech itself, which could be handled by more sophisticated analysis. It is also possible to use other, more or less intrusive, tracking technologies for gesture tracking, such as cameras or electromagnetic trackers. On the other hand, it might be desirable to avoid additional technology, if sufficient data can be provided through the (already available) speech input.

Since our experimental setup provides us with extremely rudimentary audio localization, we intend to investigate the use of array microphones for accurate 3D audio tracking and better speech recognition.

We are also interested in expanding our set of metrics to include pitch and energy (dB range), which are additional important features for distinguishing emotions [1]. Our experiments naturally expand into the consideration of prosodic speech features for adjectives (e.g., "Faster-faster-faster! Slooooweeer!"). The user could use adjectives for continuous adjustment of the rate of the operation, instead of having to reissue the command differently. Finally, we plan to perform a user study of our speech-based cursor control to investigate its benefits, and further investigate speech-based interaction techniques for 2D and 3D manipulation.

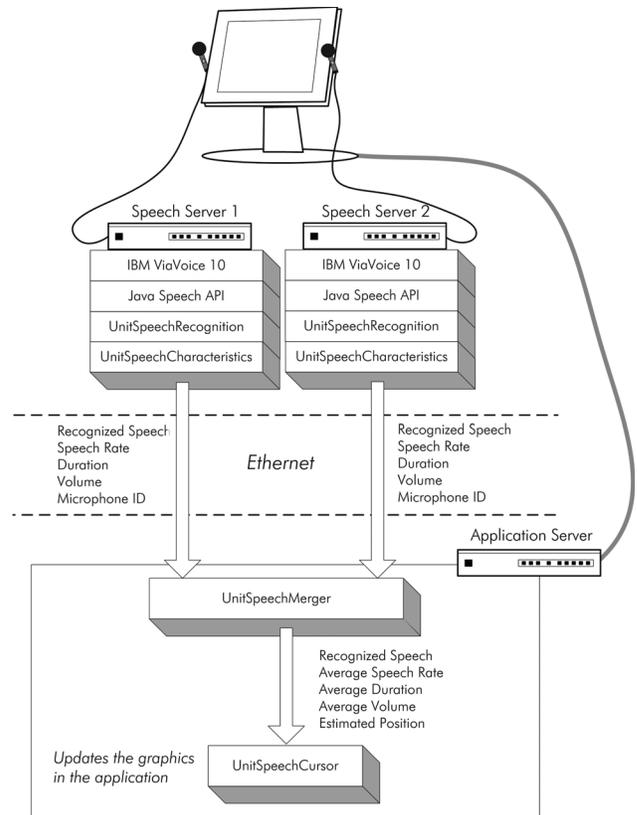


Figure 4. The architecture for our experimental setup. The Unit dataflow framework provides the infrastructure for the application.

ACKNOWLEDGEMENTS

This research was supported in part by Office of Naval Research Contracts N00014-99-1-0249, N00014-99-1-0394, and N00014-04-1-0005; NSF Grants IIS-00-82961 and IIS-01-21239; and gifts from Microsoft Research and Alias Systems.

REFERENCES

1. Hirschberg, J., Liscombe, J., and Venditti, J. Experiments in emotional speech. *Proc. ISCA & IEEE Workshop on Spontaneous Speech Processing and Recognition*, 2003.
2. Igarashi, T. and Hughes, J.F. Voice as sound: Using non-verbal voice input for interactive control. *Proc. UIST 2001*, 2001, 155-156.
3. Karimullah, A.S. and Sears, A. Speech-based cursor control. *Proc. SIGCAPH 2002*, 2002, 178-185.
4. Olwal, A. and Feiner, S. Unit: Modular Development of Distributed Interaction Techniques for Highly Interactive User Interfaces. *Proc. GRAPHITE 2004*, 2004, 131-138.
5. Oviatt, S. Mutual disambiguation of recognition errors in a multimodel architecture. *Proc. CHI '99*, 1999, 576-583.
6. Tsukahara, W. and Ward, N., Responding to subtle, fleeting changes in the user's internal state. *Proc. CHI 2001*, 2001, 77-84.