

Multimodal Interfaces

Sharon Oviatt
Center for Human-Computer Communication, Computer Science Dept.
Oregon Graduate Institute of Science & Technology
Beaverton, Oregon, USA

1. What are multimodal systems, and why are we building them?

Multimodal systems process two or more combined user input modes— such as speech, pen, touch, manual gestures, gaze, and head and body movements— in a coordinated manner with multimedia system output. This class of systems represents a new direction for computing, and a paradigm shift away from conventional WIMP interfaces. Since the appearance of Bolt's (1980) "Put That There" demonstration system, which processed speech in parallel with touch-pad pointing, a variety of new multimodal systems has emerged. This new class of interfaces aims to recognize naturally occurring forms of human language and behavior, which incorporate at least one recognition-based technology (e.g., speech, pen, vision). The development of novel multimodal systems has been enabled by the myriad input and output technologies currently becoming available, including new devices and improvements in recognition-based technologies. This chapter will review the main types of multimodal interfaces, their advantages and cognitive science underpinnings, primary features and architectural characteristics, and general research in the field of multimodal interaction and interface design.

The growing interest in multimodal interface design is inspired largely by the goal of supporting more transparent, flexible, efficient, and powerfully expressive means of human-computer interaction. Multimodal interfaces are expected to be easier to learn and use, and are preferred by users for many applications. They have the potential to expand computing to more challenging applications, to be used by a broader spectrum of everyday people, and to accommodate more adverse usage conditions than in the past. Such systems also have the potential to function in a more robust and stable manner than unimodal recognition systems involving a single recognition-based technology, such as speech, pen, or vision.

The advent of multimodal interfaces based on recognition of human speech, gaze, gesture, and other natural behavior represents only the beginning of a progression toward computational interfaces capable of relatively human-like sensory perception. Such interfaces eventually will interpret continuous input from a large number of different visual, auditory, and tactile input modes, which will be recognized as users engage in everyday activities. The same system will track and incorporate information from multiple sensors on the user's interface and surrounding physical environment in order to support intelligent adaptation to the user, task and usage environment. Future adaptive multimodal-multisensor interfaces have the potential to support new functionality, to achieve unparalleled robustness, and to perform flexibly as a multifunctional and personalized mobile system.

2. What types of multimodal interfaces exist, and what is their history and current status?

Multimodal systems have developed rapidly during the past decade, with steady progress toward building more general and robust systems, as well as more transparent human interfaces than ever before (Benoit, Martin, Pelachaud, Schomaker, & Suhm, 2000; Oviatt et al., 2000). Major developments have occurred in the hardware and software needed to support key component technologies incorporated within multimodal systems, as well as in techniques for integrating parallel input streams. Multimodal systems also have diversified to include new modality combinations, including speech and pen input, speech and lip movements, speech and manual gesturing, and gaze tracking and manual input (Benoit & Le Goff, 1998; Cohen, et al., 1997; Stork & Hennecke, 1995; Turk & Robertson, 2000; Zhai, Morimoto & Ihde, 1999). In addition, the array of multimodal applications has expanded rapidly, and presently ranges from map-based and virtual reality systems for simulation and training, to person identification/verification systems for security purposes, to medical and web-based transaction systems that

eventually will transform our daily lives (Neti, Iyengar, Potamianos & Senior, 2000; Oviatt et al., 2000; Pankanti, Bolle & Jain, 2000).

In one of the earliest multimodal concept demonstrations, Bolt had users sit in front of a projection of “Dataland” in “the Media Room” (Negroponte, 1978). Using the “Put That There” interface (Bolt, 1980), they could use speech and pointing on an armrest-mounted touchpad to create and move objects on a 2-D large-screen display. For example, the user could issue a command to “Create a blue square there,” with the intended location of “there” indicated by a 2-D cursor mark on the screen. Semantic processing was based on the user’s spoken input, and the meaning of the deictic “there” was resolved by processing the x,y coordinate indicated by the cursor at the time “there” was uttered. Since Bolt’s early prototype, considerable strides have been made in developing a wide variety of different types of multimodal systems.

Among the earliest and most rudimentary multimodal systems were ones that supported speech input along with a standard keyboard and mouse interface. Conceptually, these multimodal interfaces represented the least departure from traditional graphical user interfaces (GUIs). Their initial focus was on supporting richer natural language processing to support greater expressive power for the user when manipulating complex visuals and engaging in information extraction. As speech recognition technology matured during the late 1980s and 1990s, these systems added spoken input as an alternative to text entry via the keyboard. As such, they represent early involvement of the natural language and speech communities in developing the technologies needed to support new multimodal interfaces. Among the many examples of this type of multimodal interface are CUBRICON, Georal, Galaxy, XTRA, Shoptalk and Miltalk (Cohen et al., 1989; Kobsa, et al., 1986; Neal & Shapiro, 1991; Seneff, Goddeau, Pao & Polifroni, 1996; Siroux, Guyomard, Multon & Remondeau, 1995; Wahlster, 1991).

Several of these early systems were multimodal-multimedia map systems to which a user could speak or type and point with a mouse to extract tourist information or engage in military situation assessment (Cohen et al., 1989; Neal & Shapiro, 1991; Seneff, et al., 1996; Siroux, et al., 1995). For example, using the CUBRICON system a user could point to an object on a map and ask: “*Is this <point> an air base?*” CUBRICON was an expert system with extensive domain knowledge, as well as natural language processing capabilities that included referent identification and dialogue tracking (Neal & Shapiro, 1991). With the Georal system, a user could query a tourist information system to plan travel routes using spoken input and pointing via a touch-sensitive screen (Siroux, et al., 1995). In contrast, the Shoptalk system permitted users to interact with complex graphics representing factory production flow for chip manufacturing (Cohen et al., 1989). Using Shoptalk, a user could point to a specific machine in the production layout and issue the command: “*Show me all the times when this machine was down.*” After the system delivered its answer as a list of time ranges, the user could click on one to ask the follow-up question, “*What chips were waiting in its queue then, and were any of them hot lots?*” Multimedia system feedback was available in the form of a text answer, or the user could click on the machine in question to view an exploded diagram of the machine queue’s contents during that time interval.

More recent multimodal systems have moved away from processing simple mouse or touch-pad pointing, and have begun designing systems based on two parallel input streams that each are capable of conveying rich semantic information. These multimodal systems recognize two natural forms of human language and behavior, for which two recognition-based technologies are incorporated within a more powerful bimodal user interface. To date, systems that combine either speech and pen input (Oviatt & Cohen, 2000) or speech and lip movements (Benoit, et al, 2000; Stork & Hennecke, 1995; Rubin, Vatikiotis-Bateson, & Benoit, 1998) constitute the two most mature areas within the field of multimodal research. In both cases, the keyboard & mouse have been abandoned. For speech & pen systems, spoken language sometimes is processed along with complex pen-based gestural input involving hundreds of different symbolic interpretations beyond pointing¹ (Oviatt et al., 2000). For speech & lip movement systems, spoken language is processed along with corresponding human lip movements during the natural audio-visual experience of spoken interaction. In both of these sub-literatures, considerable work has been directed toward quantitative modeling of the integration and synchronization characteristics of the two rich input modes being processed, and innovative time-sensitive architectures have been developed to process these patterns in a robust manner.

¹ However, other recent pen/voice multimodal systems that emphasize mobile processing, such as MiPad and the Field Medic Information System (Holzman, 1999; Huang, et al., 2000), still limit pen input to pointing.

Multimodal systems that recognize speech and pen-based gestures first were designed and studied in the early 1990s (Oviatt, Cohen, Fong, & Frank, 1992), with the original QuickSet system prototype built in 1994. The QuickSet system is an agent-based collaborative multimodal system that runs on a hand-held PC (Cohen et al., 1997). With QuickSet, for example, a user can issue a multimodal command such as “Airstrips... facing this way <draws arrow>,” and facing this way <draws arrow>,” using combined speech and pen input to place the correct number, length and orientation (e.g., SW, NE) of aircraft landing strips on a map. Other research-level systems of this type were built in the late 1990s. Examples include the Human-centric Word Processor, Portable Voice Assistant, QuickDoc and MVEWS (Bers, Miller, & Makhoul, 1998; Cheyer, 1998; Oviatt et al., 2000; Waibel, Suhm, Vo, & Yang, 1997). These systems represent a variety of different system features, applications, and information fusion and linguistic processing techniques. For illustration purposes, a comparison of five different speech and gesture systems is summarized in Figure 1. In most cases, these multimodal systems jointly interpreted speech and pen input based on a frame-based method of information fusion and a late semantic fusion approach, although QuickSet used a statistically-ranked unification process and a hybrid symbolic/statistical architecture (Wu, Oviatt, & Cohen, 1999). Other very recent systems also have begun to adopt unification-based multimodal fusion and a hybrid architectural approach (Bangalore & Johnston, 2000; Denecke & Yang, 2000; Wahlster, 2001).

Multimodal System Characteristics:	QuickSet	Human-Centric Word Processor	VR Aircraft Maintenance Training	Field Medic Information	Portable Voice Assistant
Recognition of simultaneous or alternative individual modes	Simultaneous & individual modes	Simultaneous & individual modes	Simultaneous & individual modes	Alternative individual modes ¹	Simultaneous & individual modes
Type & size of gesture vocabulary	Pen input, Multiple gestures, Large vocabulary	Pen input, Deictic selection	3D manual input, Multiple gestures, Small vocabulary	Pen input, Deictic selection	Pen input Deictic selection ²
Size of speech vocabulary³ & type of linguistic processing	Moderate vocabulary, Grammar-based	Large vocabulary, Statistical language processing	Small vocabulary, Grammar-based	Moderate vocabulary, Grammar-based	Small vocabulary, Grammar-based
Type of signal fusion	Late semantic fusion, Unification, Hybrid symbolic/statistical MTC framework	Late semantic fusion, Frame-based	Late semantic fusion, Frame-based	No mode fusion	Late semantic fusion, Frame-based
Type of platform & applications	Wireless handheld, Varied map & VR applications	Desktop computer, Word processing	Virtual reality system, Aircraft maintenance training	Wireless handheld, Medical field emergencies	Wireless handheld, Catalogue ordering
Evaluation status	Proactive user-centered design & iterative system evaluations	Proactive user-centered design	Planned for future	Proactive user-centered design & iterative system evaluations	Planned for future

Figure 1. Examples of functionality, architectural features, and general classification of different speech and gesture multimodal applications.

¹ The FMA component recognizes speech only, and the FMC component recognizes gestural selections or speech. The FMC also can transmit digital speech and ink data, and can read data from smart cards and physiological monitors.

² The PVA also performs handwriting recognition.

³ A small speech vocabulary is up to 200 words, moderate 300-1,000 words, and large in excess of 1,000 words. For pen-based gestures, deictic selection is an individual gesture, a small vocabulary is 2-20 gestures, moderate 20-100, and large in excess of 100 gestures.

Although many of the issues discussed for multimodal systems incorporating speech and 2-D pen gestures also are relevant to those involving continuous 3-D manual gesturing, the latter type of system presently is less mature (Sharma, Pavlovic, & Huang, 1998; Pavlovic, Sharma, & Huang, 1997). This primarily is because of the significant challenges associated with segmenting and interpreting continuous manual movements, compared with a stream of x,y ink coordinates. As a

result of this difference, multimodal speech and pen systems have advanced more rapidly in their architectures, and have progressed further toward commercialization of applications. However, a significant cognitive science literature is available for guiding the design of emerging speech and 3-D gesture prototypes (Condon, 1988; Kendon, 1980; McNeill, 1992), which will be discussed further in section 5. Existing systems that process manual pointing or 3-D gestures combined with speech have been developed by Koons and colleagues (Koons, Sparrell, & Thorisson, 1993), Sharma and colleagues (Sharma et al., 1996), Poddar and colleagues (Poddar, Sethi, Ozyildiz, & Sharma, 1998), and by Duncan and colleagues (Duncan, Brown, Esposito, Holmback, & Xue, 1999).

Historically, multimodal speech and lip movement research has been driven by cognitive science interest in intersensory audio-visual perception and the coordination of speech output with lip and facial movements (Benoit & Le Goff, 1998; Bernstein & Benoit, 1996; Cohen & Massaro, 1993; Massaro & Stork, 1998; McGrath & Summerfield, 1987; McGurk & MacDonald, 1976; McLeod & Summerfield, 1987; Robert-Ribes, Schwartz, Lallouache & Escudier, 1998; Sumbly & Pollack, 1954; Summerfield, 1992; Vatikiotis-Bateson, Munhall, Hirayama, Lee, & Terzopoulos, 1996). Among the many contributions of this literature has been a detailed classification of human lip movements (visemes) and the viseme-phoneme mappings that occur during articulated speech. Existing systems that have processed combined speech and lip movements include the classic work by Petajan (1984), Brooke and Petajan (1986), and others (Adjoudani & Benoit, 1995; Bregler & Konig, 1994; Silsbee & Su, 1996; Tomlinson, Russell & Brooke, 1996). Additional examples of speech and lip movement systems, applications, and relevant cognitive science research have been detailed elsewhere (Benoit, et al., 2000). Although few existing multimodal interfaces currently include adaptive processing, researchers in this area have begun to explore adaptive techniques for improving system robustness in noisy environmental contexts (Dupont & Luetin, 2000; Meier, Hürst & Duchnowski, 1996; Rogozan & Deglise, 1998), which is an important future research direction. Although this literature has not emphasized the development of applications, nonetheless its quantitative modeling of synchronized phoneme/viseme patterns has been used to build animated characters that generate text-to-speech output and coordinated lip movements. These new animated characters are being used as an interface design vehicle for facilitating users' multimodal interaction with next-generation conversational interfaces (Cassell, Sullivan, Prevost, & Churchill, 2000; Cohen & Massaro, 1993).

While the main multimodal literatures to date have focused on either speech and pen input or speech and lip movements, recognition of other modes also is maturing and beginning to be integrated into new kinds of multimodal systems. In particular, there is growing interest in designing multimodal interfaces that incorporate vision-based technologies, such as interpretation of gaze, facial expressions, and manual gesturing (Morimoto, Koons, Amir, Flickner & Zhai, 1999; Pavlovic, Berry, & Huang, 1997; Turk & Robertson, 2000; Zhai, et al., 1999). These technologies unobtrusively or *passively* monitor user behavior and need not require explicit user commands to a "computer." This contrasts with *active input modes*, such as speech or pen, which the user deploys intentionally as a command issued to the system. While passive modes may be "attentive" and less obtrusive, active modes generally are more reliable indicators of user intent.

Multimodal interfaces process two or more combined user input modes— such as speech, pen, touch, manual gestures, gaze, and head and body movements— in a coordinated manner with multimedia system output. They are a new class of interfaces that aim to recognize naturally occurring forms of human language and behavior, and which incorporate one or more recognition-based technologies (e.g., speech, pen, vision).

Active input modes are ones that are deployed by the user intentionally as an explicit command to a computer system (e.g., speech).

Passive input modes refer to naturally occurring user behavior or actions that are recognized by a computer (e.g., facial expressions, manual gestures). They involve user input that is unobtrusively and passively monitored, without requiring any explicit command to a computer.

Blended multimodal interfaces are ones that incorporate system recognition of at least one passive and one active input mode. (e.g., speech and lip movement systems).

Temporally-cascaded multimodal interfaces are ones that process two or more user modalities that tend to be sequenced in a particular temporal order (e.g., gaze, gesture, speech), such that partial information supplied by recognition of an earlier mode (e.g., gaze) is available to constrain interpretation of a later mode (e.g., speech). Such interfaces may combine only active input modes, only passive ones, or they may be blended.

Mutual disambiguation involves disambiguation of signal or semantic-level information in one error-prone input mode from partial information supplied by another. Mutual disambiguation can occur in a multimodal architecture with two or more semantically rich recognition-based input modes. It leads to recovery from unimodal recognition errors within a multimodal architecture, with the net effect of suppressing errors experienced by the user.

Visemes refers to the detailed classification of visible lip movements that correspond with consonants and vowels during articulated speech. A *viseme-phoneme mapping* refers to the correspondence between visible lip movements and audible phonemes during continuous speech.

Feature-level fusion is a method for fusing low-level feature information from parallel input signals within a multimodal architecture, which has been applied to processing closely synchronized input such as speech and lip movements.

Semantic-level fusion is a method for integrating semantic information derived from parallel input modes in a multimodal architecture, which has been used for processing speech and gesture input.

Frame-based integration is a pattern matching technique for merging attribute-value data structures to fuse semantic information derived from two input modes into a common meaning representation during multimodal language processing.

Unification-based integration is a logic-based method for integrating partial meaning fragments derived from two input modes into a common meaning representation during multimodal language processing. Compared with frame-based integration, unification derives from logic programming, and has been more precisely analyzed and widely adopted within computational linguistics.

As multimodal interfaces gradually evolve toward supporting more advanced recognition of users' natural activities in context, they will expand beyond rudimentary bimodal systems to ones that incorporate three or more input modes, qualitatively different modes, and more sophisticated models of multimodal interaction. This trend already has been initiated within biometrics research, which has combined recognition of multiple behavioral input modes (e.g., voice, handwriting) with physiological ones (e.g., retinal scans, fingerprints) to achieve reliable person identification and verification in challenging field conditions (Choudhury, Clarkson, Jebara & Pentland, 1999; Pankanti, et al., 2000).

3. What are the goals and advantages of multimodal interface design?

Over the past decade, numerous advantages of multimodal interface design have been documented. Unlike a traditional keyboard and mouse interface or a unimodal recognition-based interface, multimodal interfaces permit flexible use of input modes. This includes the choice of which modality to use for conveying different types of information, to use combined input modes, or to alternate between modes at any time. Since individual input modalities are well suited in some situations, and less ideal or even inappropriate in others, modality choice is an important design issue in a multimodal system. As systems become more complex and multifunctional, a single modality simply does not permit all users to interact effectively across all tasks and environments.

Since there are large individual differences in ability and preference to use different modes of communication, a multimodal interface permits diverse user groups to exercise selection and control over how they interact with the computer (Fell, et al., 1994; Karshmer & Blattner, 1998). In this respect, multimodal interfaces have the potential to accommodate a broader range of users than traditional interfaces— including users of different ages, skill levels, native language status, cognitive styles, sensory impairments, and other temporary illnesses or permanent handicaps. For example, a visually impaired user or one with repetitive stress injury may prefer speech input and text-to-speech output. In contrast, a user with a hearing impairment or accented speech may prefer touch, gesture or pen input. The natural alternation between modes that is permitted by a multimodal interface also can be effective in preventing overuse and physical damage to any single modality, especially during extended periods of computer use (Markinson¹, personal communication, 1993).

Multimodal interfaces also provide the adaptability that is needed to accommodate the continuously changing conditions of mobile use. In particular, systems involving speech, pen or touch input are suitable for mobile tasks and, when combined, users can shift among these modalities from moment to moment as environmental conditions change (Holzman, 1999; Oviatt, 2000b, 2000c). There is a sense in which mobility can induce a state of temporary disability, such that a person is unable to use a particular input mode for some period of time. For example, the user of an in-vehicle application may frequently be unable to use manual or gaze input, although speech is relatively more available. In this respect, a multimodal interface permits the modality choice and switching that is needed during the changing environmental circumstances of actual field and mobile use.

A large body of data documents that multimodal interfaces satisfy higher levels of user preference when interacting with simulated or real computer systems. Users have a strong preference to interact multimodally, rather than unimodally, across a wide variety of different application domains, although this preference is most pronounced in spatial domains (Hauptmann, 1989; Oviatt, 1997). For example, 95% to 100% of users preferred to interact multimodally when they were free to use either speech or pen input in a map-based spatial domain (Oviatt, 1997). During pen/voice multimodal interaction, users preferred speech input for describing objects and events, sets and subsets of objects, out-of-view objects, conjoined information, past and future temporal states, and for issuing commands for actions or iterative actions (Cohen & Oviatt, 1995; Oviatt & Cohen, 1991). However, their preference for pen input increased when conveying digits, symbols, graphic content, and especially when conveying the location and form of spatially-oriented information on a dense graphic display such as a map (Oviatt & Olsen, 1994; Oviatt, 1997; Suhm, 1998). Likewise, 71% of users combined speech and manual gestures multimodally, rather than using one input mode, when manipulating graphic objects on a CRT screen (Hauptmann, 1989).

During the early design of multimodal systems, it was assumed that efficiency gains would be the main advantage of designing an interface multimodally, and that this advantage would derive from the ability to process input modes in parallel. It is true that multimodal interfaces sometimes support improved efficiency, especially when manipulating graphical information. In simulation research comparing speech-only with multimodal pen/voice interaction, empirical work demonstrated that multimodal interaction yielded 10% faster task completion time during visual-spatial tasks, but no significant efficiency advantage in verbal or quantitative task domains (Oviatt, 1997; Oviatt, Cohen, & Wang, 1994). Likewise, users' efficiency improved when they combined speech and gestures multimodally to manipulate 3D objects, compared with unimodal input (Hauptmann, 1989). In another early study, multimodal speech and mouse input improved efficiency in a line-art drawing task (Leatherby & Pausch, 1992). Finally, in a study that compared task completion times for a graphical interface versus a multimodal pen/voice interface, military domain experts averaged four times faster at setting up complex simulation scenarios on a map

¹ R. Markinson, University of California at San Francisco Medical School, 1993.

when they were able to interact multimodally (Cohen, McGee & Clow, 2000). This latter study was based on testing of a fully functional multimodal system, and it included time required to correct recognition errors.

One particularly advantageous feature of multimodal interface design is its superior error handling, both in terms of error avoidance and graceful recovery from errors (Oviatt & van Gent, 1996; Oviatt, Bernard, & Levow, 1999; Oviatt, 1999a; Rudnicky & Hauptmann, 1992; Suhm, 1998; Tomlinson, et al., 1996). There are user-centered and system-centered reasons why multimodal systems facilitate error recovery, when compared with unimodal recognition-based interfaces. For example, in a multimodal speech and pen-based gesture interface users will select the input mode that they judge to be less error prone for particular lexical content, which tends to lead to error avoidance (Oviatt & van Gent, 1996). They may prefer speedy speech input, but will switch to pen input to communicate a foreign surname. Secondly, users' language often is simplified when interacting multimodally, which can substantially reduce the complexity of natural language processing and thereby reduce recognition errors (Oviatt & Kuhn, 1998; see section 5 for discussion). In one study, users' multimodal utterances were documented to be briefer, to contain fewer complex locative descriptions, and 50% fewer spoken disfluencies, when compared with a speech-only interface. Thirdly, users have a strong tendency to switch modes after system recognition errors, which facilitates error recovery. This error resolution occurs because the confusion matrices differ for any given lexical content for the different recognition technologies involved in processing (Oviatt, et al., 1999).

In addition to these user-centered reasons for better error avoidance and resolution, there also are system-centered reasons for superior error handling. A well-designed multimodal architecture with two semantically rich input modes can support *mutual disambiguation* of input signals. For example, if a user says "ditches" but the speech recognizer confirms the singular "ditch" as its best guess, then parallel recognition of several graphic marks can result in recovery of the correct plural interpretation. This recovery can occur in a multimodal architecture even though the speech recognizer initially ranks the plural interpretation "ditches" as a less preferred choice on its n-best list. Mutual disambiguation involves recovery from unimodal recognition errors within a multimodal architecture, because semantic information from each input mode supplies partial disambiguation of the other mode, thereby leading to more stable and robust overall system performance (Oviatt, 1999a, 2000a). Another example of mutual disambiguation is shown in Figure 3. To achieve optimal error handling, a multimodal interface ideally should be designed to include complementary input modes, and so the alternative input modes provide duplicate functionality such that users can accomplish their goals using either mode.

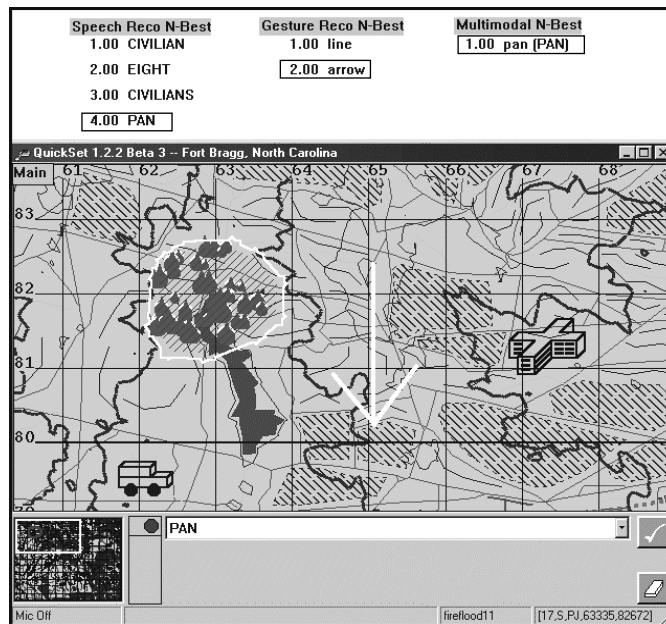


Figure 3. Multimodal command to "pan" the map, which illustrates mutual disambiguation occurring between incoming speech and gesture information, such that lexical hypotheses were pulled up on both n-best lists to produce a correct final multimodal interpretation.

In two recent studies involving over 4,600 multimodal commands, a multimodal architecture was found to support mutual disambiguation and error suppression ranging between 19 and 41% (Oviatt, 1999a, 2000a). Improved robustness also was greater for “challenging” user groups (accented vs. native speakers) and usage contexts (mobile vs. stationary use). These results indicate that a well-designed multimodal system not only can perform more robustly than a unimodal system, but also in a more stable way across varied real-world users and usage contexts. Finally, during audio-visual perception of speech and lip movements, improved speech recognition also has been demonstrated for both human listeners (McLeod & Summerfield, 1987) and multimodal systems (Adjoudani & Benoit, 1995; Tomlinson, et al., 1996).

4. What methods and information have been used to design novel multimodal interfaces?

The design of new multimodal systems has been inspired and organized largely by two things. First, the cognitive science literature on intersensory perception and intermodal coordination during production has provided a foundation of information for user modeling, as well as information on what systems must recognize and how multimodal architectures should be organized. For example, the cognitive science literature has provided knowledge of the natural integration patterns that typify people’s lip and facial movements with speech output (Benoit, Guiard-Marigny, Le Goff & Adjoudani, 1996; Ekman, 1992; Ekman & Friesen, 1978; Fridlund, 1994; Hadar, Steiner, Grant, & Rose, 1983; Massaro & Cohen, 1990; Stork & Hennecke, 1995; Vatikiotis-Bateson et al., 1996), and their coordinated use of manual or pen-based gestures with speech (Kendon, 1980; McNeill, 1992; Oviatt, DeAngeli, & Kuhn, 1997). Given the complex nature of users’ multimodal interaction, cognitive science has and will continue to play an essential role in guiding the design of robust multimodal systems. In this respect, a multidisciplinary perspective will be more central to successful multimodal system design than it has been for traditional GUI design. The cognitive science underpinnings of multimodal system design are described in section 5.

Secondly, high-fidelity automatic simulations also have played a critical role in prototyping new types of multimodal systems (Dahlbäck, Jönsson & Ahrenberg, 1992; Oviatt et al., 1992). When a new multimodal system is in the planning stages, design sketches and low-fidelity mock-ups may initially be used to visualize the new system and plan the sequential flow of human-computer interaction. These tentative design plans then are rapidly transitioned into a higher-fidelity simulation of the multimodal system, which is available for proactive and situated data collection with the intended user population. High-fidelity simulations have been the preferred method for designing and evaluating new multimodal systems, and extensive data collection with such tools preferably is completed before a fully-functional system ever is built.

During high-fidelity simulation testing, a user interacts with what she believes is a fully-functional multimodal system, although the interface is actually a simulated front-end designed to appear and respond as the fully-functional system would. During the interaction, a programmer assistant at a remote location provides the simulated system responses. As the user interacts with the front end, the programmer tracks her multimodal input and provides system responses as quickly and accurately as possible. To support this role, the programmer makes use of automated simulation software that is designed to support interactive speed, realism with respect to the targeted system, and other important characteristics. For example, with these automated tools, the programmer may be able to make a single selection on a workstation field to rapidly send simulated system responses to the user during a data collection session.

High-fidelity simulations have been the preferred method for prototyping multimodal systems for several reasons. Simulations are relatively easy and inexpensive to adapt, compared with building and iterating a complete system. They also permit researchers to alter a planned system’s characteristics in major ways (e.g., input and output modes available), and to study the impact of different interface features in a systematic and scientific manner (e.g., type and base-rate of system errors). In comparison, a particular system with its fixed characteristics is a less flexible and suitable research tool, and the assessment of any single system basically amounts to an individual case study. Using simulation techniques, rapid adaptation and investigation of planned system features permits researchers to gain a broader and more principled perspective on the potential of newly emerging technologies. In a practical sense, simulation research can assist in the evaluation of critical performance tradeoffs and in making decisions about alternative system designs, which designers must do as they strive to create more usable multimodal systems.

To support the further development and commercialization of multimodal systems, additional infrastructure that will be needed in the future includes: (1) simulation tools for rapidly building and reconfiguring multimodal interfaces,

(2) automated tools for collecting and analyzing multimodal corpora, and (3) automated tools for iterating new multimodal systems to improve their performance (see Oviatt et al., 2000, for further discussion).

5. What are the cognitive science underpinnings of multimodal interface design?

The ability to develop multimodal systems depends on knowledge of the natural integration patterns that typify people's combined use of different input modes. In particular, the design of new multimodal systems depends on intimate knowledge of the properties of different modes and the information content they carry, the unique characteristics of multimodal language and its processability, and the integration and synchronization characteristics of users' multimodal interaction. It also relies on accurate prediction of when users are likely to interact multimodally, and how alike different users are in their specific integration patterns. The relevant cognitive science literature on these topics is very extensive, especially when consideration is given to all of the underlying sensory perception and production capabilities involved in different input modes currently being incorporated in new multimodal interfaces. As a result, this section will be limited to introducing the main cognitive science themes and findings that are relevant to the more common types of multimodal system.

5.1 When do users interact multimodally?

During natural interpersonal communication, people are always interacting multimodally. Of course, in this case the number of information sources or modalities that an interlocutor has available to monitor is essentially unlimited. However, all multimodal systems are constrained in the number and type of input modes they can recognize. Also, a user can compose active input during human-computer interaction that either is delivered multimodally or that is delivered entirely using just one mode. That is, although users in general may have a strong preference to interact multimodally rather than unimodally, this is no guarantee that they will issue every command to a system multimodally, given the particular type of multimodal interface available. Therefore, the first nontrivial question that arises during system processing is whether a user is communicating unimodally or multimodally.

In the case of speech and pen-based multimodal systems, users typically intermix unimodal and multimodal expressions. In a recent study involving a visual-spatial domain, users' commands were expressed multimodally 20% of the time, with others just spoken or written (Oviatt et al., 1997). Predicting whether a user will express a command multimodally also depends on the type of action she is performing. In particular, users almost always express commands multimodally when describing spatial information about the location, number, size, orientation, or shape of an object. In one study, users issued multimodal commands 86% of the time when they had to add, move, modify, or calculate the distance between objects on a map in a way that required specifying spatial locations (Oviatt et al., 1997). They also were moderately likely to interact multimodally when selecting an object from a larger array, for example, when deleting a particular object from the map. However, when performing general actions without any spatial component, such as printing a map, users expressed themselves multimodally less than 1% of the time. These data emphasize that future multimodal systems will need to distinguish between instances when users are and are not communicating multimodally, so that accurate decisions can be made about when parallel input streams should be interpreted jointly versus individually. They also suggest that knowledge of the type of actions to be included in an application, such as whether the application entails manipulating spatial information, should influence the basic decision of whether to build a multimodal interface at all.

In a multimodal interface that processes passive or blended input modes, there always is at least one passively-tracked input source providing continuous information (e.g., gaze tracking, head position). In these cases, all user input would by definition be classified as multimodal, and the primary problem would become segmentation and interpretation of each continuous input stream into meaningful actions of significance to the application. In the case of blended multimodal interfaces (e.g., gaze tracking and mouse input), it still may be opportune to distinguish active forms of user input that might be more accurately or expeditiously handled as unimodal events.

5.2 What are the integration and synchronization characteristics of users' multimodal input?

The past literature on multimodal systems has focused largely on simple selection of objects or locations in a display, rather than considering the broader range of multimodal integration patterns. Since the development of Bolt's (1980) "Put That There" system, speak-and-point has been viewed as the prototypical form of multimodal integration. In Bolt's system, semantic processing was based on spoken input, but the meaning of a deictic term such as "that" was resolved by processing the x,y coordinate indicated by pointing at an object. Since that time, other

multimodal systems also have attempted to resolve deictic expressions using a similar approach, for example, using gaze location instead of manual pointing (Koons et al., 1993).

Unfortunately, this concept of multimodal interaction as point-and-speak makes only limited use of new input modes for *selection* of objects— just as the mouse does. In this respect, it represents the persistence of an old mouse-oriented metaphor. In contrast, modes that transmit written input, manual gesturing, and facial expressions are capable of generating symbolic information that is much more richly expressive than simple pointing or selection. In fact, studies of users' integrated pen/voice input indicate that a speak-and-point pattern only comprises 14% of all spontaneous multimodal utterances (Oviatt et al., 1997). Instead, pen input more often is used to create graphics, symbols & signs, gestural marks, digits and lexical content. During interpersonal multimodal communication, linguistic analysis of spontaneous manual gesturing also indicates that simple pointing accounts for less than 20% of all gestures (McNeill, 1992). Together, these cognitive science and user-modeling data highlight the fact that any multimodal system designed exclusively to process speak-and-point will fail to provide users with much useful functionality. For this reason, specialized algorithms for processing deictic-point relations will have only limited practical use in the design of future multimodal systems. It is clear that a broader set of multimodal integration issues needs to be addressed in future work. Future research also should explore typical integration patterns between other promising modality combinations, such as speech and gaze.

It also is commonly assumed that any signals involved in a multimodal construction will co-occur temporally. The presumption is that this temporal overlap then determines which signals to combine during system processing. In the case of speech and manual gestures, successful processing of the deictic term “that square” in Bolt's original system relied on interpretation of pointing when the word “that” was spoken in order to extract the intended referent. However, one empirical study indicated that users often do not speak deictic terms at all, and when they do the deictic frequently is not overlapped in time with their pointing. In fact, it has been estimated that as few as 25% of users' commands actually contain a spoken deictic that overlaps with the pointing needed to disambiguate its meaning (Oviatt et al., 1997).

Beyond the issue of deixis, users' input frequently does not overlap at all during multimodal commands to a computer. During spoken and pen-based input, for example, users' multimodal input is sequentially integrated about half the time, with pen input preceding speech 99% of the time, and a brief lag between signals of one or two seconds (Oviatt et al., 1997). This finding is consistent with linguistics data revealing that both spontaneous gesturing and signed language often precede their spoken lexical analogues during human communication (Kendon, 1980; Naughton, 1996). The degree to which gesturing precedes speech is greater in topic-prominent languages such as Chinese than it is in subject-prominent ones like Spanish or English (McNeill, 1992). Even in the speech and lip movement literature, close but not perfect temporal synchrony is typical, with lip movements occurring a fraction of a second before the corresponding auditory signal (Abry, Lallouache & Cathiard, 1996; Benoit, 2000).

In short, although two input modes may be highly interdependent and synchronized during multimodal interaction, synchrony does not imply simultaneity. The empirical evidence reveals that multimodal signals often do not co-occur temporally at all during human-computer or natural human communication. Therefore, multimodal system designers cannot necessarily count on conveniently overlapped signals in order to achieve successful processing in the multimodal architectures they build. Future research needs to explore the integration patterns and temporal cascading that can occur among three or more input modes, such as gaze, gesture and speech, so that more advanced multimodal systems can be designed and prototyped.

In the design of new multimodal architectures, it is important to note that data on the order of input modes and average time lags between input modes has been used to determine the likelihood that an utterance is multimodal versus unimodal, and to establish temporal thresholds for fusion of input. In the future, weighted likelihoods associated with different utterance segmentations, for example, that an input stream containing speech, writing, speech should be segmented into [S / W S] rather than [S W / S], and with inter-modal time lag distributions, will be used to optimize correct recognition of multimodal user input (Oviatt, 1999b). In the design of future time-critical multimodal architectures, data on users' integration and synchronization patterns will need to be collected for other mode combinations during realistic interactive tasks, so that temporal thresholds can be established for performing multimodal fusion.

5.3 What individual differences exist in multimodal interaction, and what are the implications for designing systems for universal access?

When users interact multimodally, there actually can be large individual differences in integration patterns. For example, previous empirical work on multimodal pen/voice integration has revealed two main types of user—ones who habitually deliver speech and pen signals in an overlapped or simultaneous manner, and others who synchronize signals sequentially with pen input preceding speech by up to 4 seconds (Oviatt, 1999b). These users' dominant integration pattern could be identified when they first began interacting with the system, and then persisted throughout their session (Oviatt, 1999b). That is, each user's integration pattern was established early and remained consistent, although two distinct integration patterns were observed among different users. As mentioned earlier, substantial differences also have been observed in the degree to which manual gestures precede speech for members of different linguistic groups, such as Chinese, Spanish and Americans. All of these findings imply that future multimodal systems capable of adapting temporal thresholds for different user groups potentially could achieve greater recognition accuracy and interactive speed.

Both individual and cultural differences also have been documented between users in modality integration patterns. For example, substantial individual differences also have been reported in the temporal synchrony between speech and lip movements (Kricos, 1996). In addition, lip movements during speech production are less exaggerated among Japanese speakers than Americans (Sekiyama & Tohkura, 1991). In fact, extensive inter-language differences have been observed in the information available from lip movements during audio-visual speech (Fuster-Duran, 1996). These findings have implications for the degree to which disambiguation of speech can be achieved through lip movement information in noisy environments or for different user populations. Finally, non-native speakers, the hearing impaired, and elderly listeners all are more influenced by visual lip movement than auditory cues when processing speech (Fuster-Duran, 1996; Massaro, 1996). These results have implications for the design and expected value of audio-visual multimedia output for different user groups in animated character interfaces.

Finally, gender, age, and other individual differences are common in gaze patterns, as well as speech and gaze integration (Argyle, 1972). As multimodal interfaces incorporating gaze become more mature, further research will need to explore these gender and age specific patterns, and to build appropriately adapted processing strategies. In summary, considerably more research is needed on multimodal integration and synchronization patterns for new mode combinations, as well as for diverse and disabled users for whom multimodal interfaces may be especially suitable for ensuring universal access.

5.4 Is complementarity or redundancy the main organizational theme that guides multimodal integration?

It frequently is claimed that the propositional content conveyed by different modes during multimodal communication contains a high degree of redundancy. However, the dominant theme in users' natural organization of multimodal input actually is complementarity of content, not redundancy. For example, speech and pen input consistently contribute different and complementary semantic information—with the subject, verb, and object of a sentence typically spoken, and locative information written (Oviatt et al., 1997). In fact, a major complementarity between speech and manually-oriented pen input involves visual-spatial semantic content, which is one reason these modes are an opportune combination for visual-spatial applications. Whereas spatial information is uniquely and clearly indicated via pen input, the strong descriptive capabilities of speech are better suited for specifying temporal and other non-spatial information. Even during multimodal correction of system errors, when users are highly motivated to clarify and reinforce their information delivery, speech and pen input express redundant information less than 1% of the time. Finally, during interpersonal communication linguists also have documented that spontaneous speech and manual gesturing involve complementary rather than duplicate information between modes (McNeill, 1992).

Other examples of primary multimodal complementarities during interpersonal and human-computer communication have been described in past research (McGurk & MacDonald, 1976; Oviatt & Olsen, 1994; Wickens, Sandry & Vidulich, 1983). For example, in the literature on multimodal speech and lip movements, natural feature-level complementarities have been identified between visemes and phonemes for vowel articulation, with vowel rounding better conveyed visually, and vowel height and backness better revealed auditorally (Massaro & Stork, 1998; Robert-Ribes, et al., 1998).

In short, actual data highlight the importance of complementarity as a major organizational theme during multimodal communication. The designers of next-generation multimodal systems therefore should not expect to rely on duplicated information when processing multimodal language. In multimodal systems involving both speech and pen-based gestures and speech and lip movements, one explicit goal has been to integrate complementary modalities in a manner that yields a synergistic blend, such that each mode can be capitalized upon and used to overcome weaknesses in the other mode (Cohen, et al., 1989). This approach to system design has promoted the philosophy of using modes and component technologies to their natural advantage, and of combining them in a manner that permits mutual disambiguation. One advantage of achieving such a blend is that the resulting multimodal architecture can function more robustly than an individual recognition-based technology or a multimodal system based on input modes lacking natural complementarities.

Section 5 discusses the growing cognitive science literature that provides the empirical underpinnings needed to design next-generation multimodal interfaces. This cognitive science foundation has played a key role in identifying computational “myths” about multimodal interaction, and replacing these misconceptions with contrary empirical evidence. Ten common myths about multimodal interaction include:

Myth #1: *If you build a multimodal system, users will interact multimodally*

Myth #2: *Speech & pointing is the dominant multimodal integration pattern*

Myth #3: *Multimodal input involves simultaneous signals*

Myth #4: *Speech is the primary input mode in any multimodal system that includes it*

Myth #5: *Multimodal language does not differ linguistically from unimodal language*

Myth #6: *Multimodal integration involves redundancy of content between modes*

Myth #7: *Individual error-prone recognition technologies combine multimodally to produce even greater unreliability*

Myth #8: *All users’ multimodal commands are integrated in a uniform way*

Myth #9: *Different input modes are capable of transmitting comparable content*

Myth #10: *Enhanced efficiency is the main advantage of multimodal systems*

— **Summarized myths taken from “Ten Myths of Multimodal Interaction”
(Oviatt, 1999b)**

Figure 4. Ten myths of multimodal interaction: Separating myth from empirical reality.

5.5 What are the primary features of multimodal language?

Communication channels can be tremendously influential in shaping the language transmitted within them. From past research, there now is cumulative evidence that many linguistic features of multimodal language are qualitatively very different from that of spoken or formal textual language. In fact, it can differ in features as basic as brevity, semantic content, syntactic complexity, word order, disfluency rate, degree of ambiguity, referring expressions, specification of determiners, anaphora, deixis, and linguistic indirectness. In many respects, multimodal language is simpler linguistically than spoken language. In particular, comparisons have revealed that the same user completing the same map-based task communicates fewer words, briefer sentences, and fewer complex spatial descriptions and disfluencies when interacting multimodally, compared with using speech alone (Oviatt, 1997). One implication of these findings is that multimodal interface design has the potential to support more robust future systems than a unimodal design approach. The following is an example of a typical user’s spoken input while attempting to designate an open space using a map system: “Add an open space on the north lake to b-- include the

north lake part of the road and north.” In contrast, the same user accomplished the same task multimodally by encircling a specific area and saying: *“Open space.”*

In previous research, hard-to-process disfluent language has been observed to decrease by 50% during multimodal interaction with a map, compared with a more restricted speech-only interaction (Oviatt, 1997). This drop occurs mainly because people have difficulty speaking spatial information, which precipitates disfluencies. In a flexible multimodal interface, they instead use pen input to convey spatial information, thereby avoiding the need to speak it. Further research is needed to establish whether other forms of flexible multimodal communication generally ease users’ cognitive load, which may be reflected in a reduced rate of disfluencies.

During multimodal pen/voice communication, the linguistic indirection that is typical of spoken language frequently is replaced with more direct commands (Oviatt & Kuhn, 1998). In the following example, a study participant made a disfluent indirect request using speech input while requesting a map-based distance calculation: *“What is the distance between the Victorian Museum and the, uh, the house on the east side of Woodpecker Lane?”* When requesting distance information multimodally, the same user encircled the house and museum while speaking the following brief direct command: *“Show distance between here and here.”* In this research, the briefer and more direct multimodal pen/voice language also contained substantially fewer referring expressions, with a selective reduction in co-referring expressions that instead were transformed into deictic expressions. This latter reduction in coreference would simplify natural language processing by easing the need for anaphoric tracking and resolution in a multimodal interface. Also consistent with fewer referring expressions, explicit specification of definite and indefinite reference is less common in multimodal language (Oviatt & Kuhn, 1998). Current natural language processing algorithms typically rely heavily on the specification of determiners in definite and indefinite references in order to represent and resolve noun phrase reference. One unfortunate by-product of the lack of such specifications is that current language processing algorithms are unprepared for the frequent occurrence of elision and deixis in multimodal human-computer interaction.

In other respects, multimodal language clearly is different than spoken language, although not necessarily simpler. For example, users’ multimodal pen/voice language departs from the canonical English word order of S-V-O-LOC (i.e., Subject-Verb-Object-Locative constituent), which is observed in spoken language and also formal textual language. Instead, users’ multimodal constituents shift to a LOC-S-V-O word order. A recent study reported that 95% of locative constituents were in sentence-initial position during multimodal interaction. However, for the same users completing the same tasks while speaking, 96% of locatives were in sentence-final position (Oviatt, et al., 1997). It is likely that broader analyses of multimodal communication patterns, which could involve gaze and manual gesturing to indicate location rather than pen-based pointing, would reveal a similar reversal in word order.

One implication of these many differences is that new multimodal corpora, statistical language models, and natural language processing algorithms will need to be established before multimodal language can be processed optimally. Future research and corpus collection efforts also will be needed on different types of multimodal communication, and in other application domains, so that the generality of previously identified multimodal language differences can be explored.

6. What are the basic ways in which multimodal interfaces differ from graphical user interfaces?

Multimodal research groups currently are rethinking and redesigning basic user interface architectures, because a whole new range of architectural requirements has been posed. First, graphical user interfaces typically assume that there is a single event stream that controls the underlying event loop, with any processing sequential in nature. For example, most GUIs ignore typed input when a mouse button is depressed. In contrast, multimodal interfaces typically can process continuous and simultaneous input from parallel incoming streams. Secondly, GUIs assume that the basic interface actions, such as selection of an item, are atomic and unambiguous events. In contrast, multimodal systems process input modes using recognition-based technologies, which are designed to handle uncertainty and entail probabilistic methods of processing. Thirdly, GUIs often are built to be separable from the application software that they control, although the interface components usually reside centrally on one machine. In contrast, recognition-based user interfaces typically have larger computational and memory requirements, which often makes it desirable to distribute the interface over a network so that separate machines can handle different recognizers or databases. For example, cell phones and networked PDAs may extract features from speech input, but transmit them to a recognizer that resides on a server. Finally, multimodal interfaces that process two or more

recognition-based input streams require time-stamping of input, and the development of temporal constraints on mode fusion operations. In this regard, they involve uniquely time-sensitive architectures.

7. What basic architectures and processing techniques have been used to design multimodal systems?

Many early multimodal interfaces that handled combined speech and gesture, such as Bolt's "Put That There" system (Bolt, 1980), have been based on a control structure in which multimodal integration occurs during the process of parsing spoken language. As discussed earlier, when the user speaks a deictic expression such as "here" or "this", the system searches for a synchronized gestural act that designates the spoken referent. While such an approach is viable for processing a point-and-speak multimodal integration pattern, as discussed earlier, multimodal systems must be able to process richer input than just pointing, including gestures, symbols, graphic marks, lip movements, meaningful facial expressions, and so forth. To support more broadly functional multimodal systems, general processing architectures have been developed since Bolt's time. Some of these recent architectures handle a variety of multimodal integration patterns, as well as the interpretation of both unimodal and combined multimodal input. This kind of architecture can support the development of multimodal systems in which modalities are processed individually as input alternatives to one another, or those in which two or more modes are processed as combined multimodal input.

For multimodal systems designed to handle joint processing of input signals, there are two main subtypes of multimodal architecture. First, there are ones that integrate signals at the *feature level* (i.e., "early fusion") and others that integrate information at a *semantic level* (i.e., "late fusion"). Examples of systems based on an early feature-fusion processing approach include those developed by Bregler and colleagues (Bregler, Manke, Hild, & Waibel, 1993), Vo and colleagues (Vo et al. 1995), and Pavlovic and colleagues (Pavlovic, et al., 1997; 1998). In a feature-fusion architecture, the signal-level recognition process in one mode influences the course of recognition in the other. Feature fusion is considered more appropriate for closely temporally synchronized input modalities, such as speech and lip movements (Stork & Hennecke, 1995; Rubin, et al., 1998).

In contrast, multimodal systems using the late semantic fusion approach have been applied to processing multimodal speech and pen input or manual gesturing, for which the input modes are less coupled temporally. These input modes provide different but complementary information that typically is integrated at the utterance level. Late semantic integration systems use individual recognizers that can be trained using unimodal data, which are easier to collect and already are publicly available for speech and handwriting. In this respect, systems based on semantic fusion can be scaled up easier in number of input modes or vocabulary size. Examples of systems based on semantic fusion include Put That There (Bolt, 1980), ShopTalk (Cohen, et al., 1989), QuickSet (Cohen et al., 1997), CUBRICON (Neal & Shapiro, 1991), Virtual World (Codella et al., 1992), Finger-Pointer (Fukumoto, Suenaga, & Mase, 1994), VisualMan (Wang, 1995), Human-Centric Word Processor, Portable Voice Assistant (Bers, et al., 1998), the VR Aircraft Maintenance Training System (Duncan et al., 1999) and Jeanie (Vo & Wood, 1996).

As an example of multimodal information processing flow in a late-stage semantic architecture, Figure 5 illustrates two input modes (e.g., speech and manual or pen-based gestures) recognized in parallel and processed by an understanding component. The results involve partial meaning representations that are fused by the multimodal integration component, which also is influenced by the system's dialogue management and interpretation of current context. During the integration process, alternative lexical candidates for the final multimodal interpretation are ranked according to their probability estimates on an n-best list. The best-ranked multimodal interpretation then is sent to the application invocation and control component, which transforms this information into a series of commands to one or more back-end application systems. System feedback typically includes multimedia output, which may incorporate text-to-speech and non-speech audio, graphics and animation, and so forth. For examples of feature-based multimodal processing flow and architectures, especially as applied to multimodal speech and lip movement systems, see Benoit, et al. (2000).

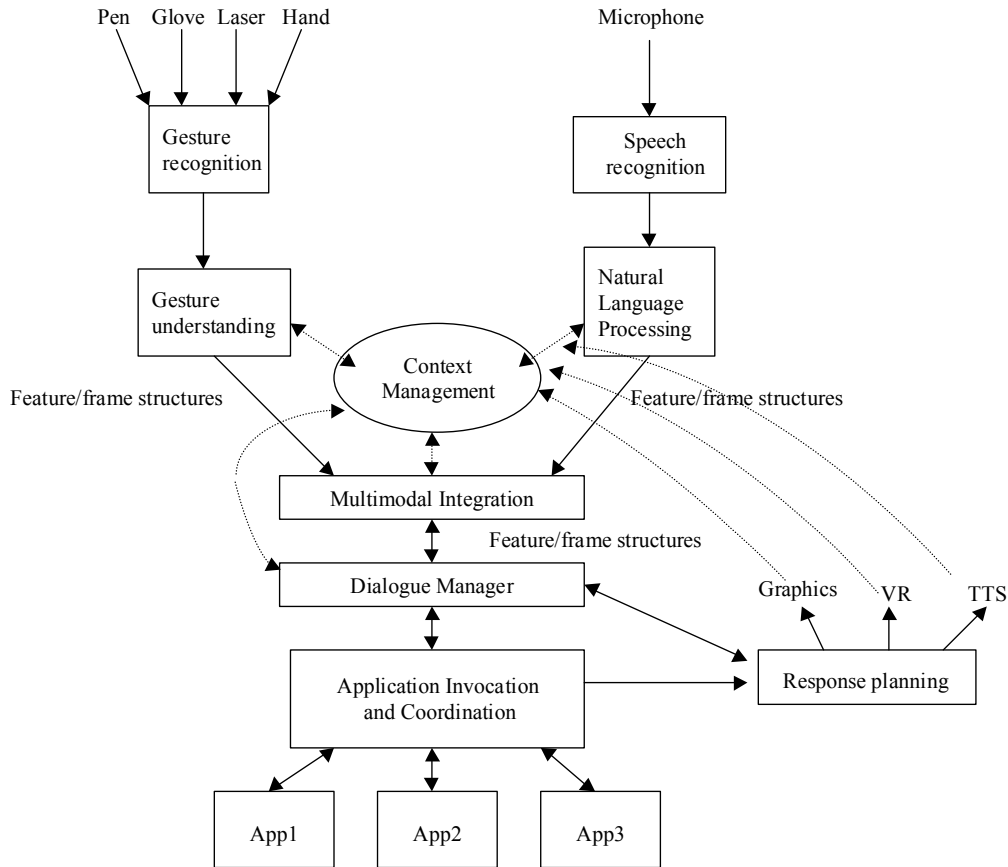


Figure 5. Typical information processing flow in a multimodal architecture designed for speech and gesture.

There are many ways to realize this information processing flow as an architecture. One common infrastructure that has been adopted by the multimodal research community involves *multi-agent architectures*, such as the Open Agent Architecture (Cohen, Cheyer, Wang, & Baeg, 1994; Martin, Cheyer, & Moran, 1999) and Adaptive Agent Architecture (Kumar & Cohen, 2000). In a multi-agent architecture, the many components needed to support the multimodal system (e.g., speech recognition, gesture recognition, natural language processing, multimodal integration) may be written in different programming languages, on different machines, and with different operating systems. Agent communication languages are being developed that can handle asynchronous delivery, triggered responses, multi-casting and other concepts from distributed systems, and that are fault-tolerant (Kumar & Cohen, 2000). Using a multi-agent architecture, for example, speech and gestures can arrive in parallel or asynchronously via individual modality agents, with the results recognized and passed to a facilitator. These results, typically an n-best list of conjectured lexical items and related time-stamp information, then are routed to appropriate agents for further language processing. Next, sets of meaning fragments derived from the speech and pen signals arrive at the multimodal integrator. This agent decides whether and how long to wait for recognition results from other modalities, based on the system's temporal thresholds. It fuses the meaning fragments into a semantically- and temporally-compatible whole interpretation before passing the results back to the facilitator. At this point, the system's final multimodal interpretation is confirmed by the interface, delivered as multimedia feedback to the user, and executed by any relevant applications. In summary, multi-agent architectures provide essential infrastructure for coordinating the many complex modules needed to implement multimodal system processing, and they permit doing so in a distributed manner that is compatible with the trend toward mobile computing.

The core of multimodal systems based on semantic fusion involves algorithms that integrate common meaning representations derived from speech, gesture and other modalities into a combined final interpretation. The semantic fusion operation requires a common meaning representation framework for all modalities, and a well-defined operation for combining partial meanings that arrive from different signals. To fuse information from different modalities, various research groups have independently converged on a strategy of recursively matching and merging attribute/value data structures, although using a variety of different algorithms (Vo & Wood, 1996; Cheyer & Julia, 1995; Pavlovic & Huang, 1998; Shaikh et al., 1997). This approach is considered a *frame-based integration* technique. An alternative logic-based approach derived from computational linguistics (Carpenter, 1990, 1992; Calder, 1987) involves the use of *typed feature structures* and *unification-based integration*, which is a more general and well understood approach. Unification-based integration techniques also have been applied to multimodal system design (Cohen et al., 1997; Johnston et al., 1997; Wu et al., 1999). Feature-structure unification is considered well suited to multimodal integration, because unification can combine complementary or redundant input from both modes, but it rules out contradictory input. Given this foundation for multimodal integration, more research still is needed on the development of canonical meaning representations that are common among different input modes which will need to be represented in new types of multimodal systems.

When statistical processing techniques are combined with a symbolic unification-based approach that merges feature structures, then the multimodal architecture that results is a *hybrid symbolic/statistical* one. Hybrid architectures represent one major new direction for multimodal system development. Multimodal architectures also can be hybrids in the sense of combining Hidden Markov Models (HMMs) and Neural Networks (NNs). New hybrid architectures potentially are capable of achieving very robust functioning, compared with either an early- or late-fusion approach alone. For example, the Members-Teams-Committee (MTC) hierarchical recognition technique, which is a hybrid symbolic/statistical multimodal integration framework trained over a labeled multimodal corpus, recently achieved 95.26% correct recognition performance, or within 1.4% of the theoretical system upper bound (Wu, et al., 1999).

8. What are the main future directions for multimodal interface design?

The computer science community is just beginning to understand how to design innovative, well integrated, and robust multimodal systems. To date, most multimodal systems remain bimodal, and recognition technologies related to several human senses (e.g., haptics, smell, taste) have yet to be well represented or included at all within multimodal interfaces. The design and development of new types of systems that include such modes will not be achievable through intuition. Rather, it will depend on knowledge of the natural integration patterns that typify people's combined use of various input modes. This means that the successful design of multimodal systems will continue to require guidance from cognitive science on the coordinated human perception and production of natural modalities. In this respect, multimodal systems only can flourish through multidisciplinary cooperation, as well as teamwork among those representing expertise in the component technologies.

Most of the systems outlined in this chapter have been built during the past decade, and they are research-level systems. However, in some cases they have developed well beyond the prototype stage, and are being integrated with other software at academic and federal sites, or appearing as newly shipped products. To achieve commercialization and widespread dissemination of multimodal interfaces, more general, robust, and scalable multimodal architectures will be needed, which just now are beginning to emerge. Future multimodal interfaces also must begin incorporating input from multiple heterogeneous information sources, such that when combined they contain the discriminative information needed to support more robust processing than individual recognition technologies. To support increasingly pervasive multimodal interfaces, these information sources also must include data collected from a wide array of input modes and sensors, and from both active and passive forms of user input. Finally, future multimodal interfaces, especially mobile ones, will require active adaptation to the user, task, ongoing dialogue, and environmental context.

In the future, multimodal interfaces could be developed that provide a better balance between system input and output, so they are better matched with one another in expressive power. Multimodal interfaces have the potential to give users more expressive power and flexibility, as well as better tools for controlling sophisticated visualization and multimedia output capabilities. As these interfaces develop, research will be needed on how to design whole multimodal-multimedia systems that are capable of highly robust functioning. To achieve this, a better understanding will be needed of the impact of visual displays, animation, text-to-speech, and audio output on users'

multimodal input and its processability. One fertile research domain for exploring these topics will be multimedia animated character design and its impact on users' multimodal input and interaction with next-generation conversational interfaces.

In conclusion, multimodal interfaces are just beginning to model human-like sensory perception. They are recognizing and identifying actions, language and people that have been seen, heard, or in other ways experienced in the past. They literally reflect and acknowledge the existence of human users, empower them in new ways, and create for them a "voice." They also can be playful and self-reflective interfaces that suggest new forms of human identity as we interact face to face with animated personas representing our own kind. In all of these ways novel multimodal interfaces, as primitive as their early bimodal instantiations may be, represent a new multidisciplinary science, a new art form, and a socio-political statement about our collective desire to humanize the technology we create.

9. Acknowledgments

I'd like to thank the National Science Foundation for their support over the past decade, which has enabled me to pursue basic exploratory research on many aspects of multimodal interaction, interface design and system development. The preparation of this chapter has been supported by NSF Grant No. IRI-9530666 and by NSF Special Extension for Creativity (SEC) Grant No. IIS-9530666. This work also has been supported by contracts DABT63-95-C-007 and N66001-99-D-8503 from DARPA's Information Technology and Information Systems Office, and Grant No. N00014-99-1-0377 from ONR. I'd also like to thank Phil Cohen and others in the Center for Human-Computer Communication for many insightful discussions, and Dana Director and Rachel Coulston for expert assistance with manuscript preparation. I also wish to acknowledge LEA, Inc. for giving permission to reprint figures 1 and 5, and to acknowledge ACM for allowing the reprint of figures 3 and 4.

References

- Abry, C., Lallouache, M.-T., & Cathiard, M.-A. (1996). How can coarticulation models account for speech sensitivity to audio-visual desynchronization? In D.G. Stork & M.E. Hennecke (Eds.), Speechreading by Humans and Machines: Models, Systems and Applications, (pp. 247-255). New York: Springer Verlag.
- Adjoudani, A., & Benoit, C. (1995). Audio-visual speech recognition compared across two architectures. Proc. of the Eurospeech Conference Vol.2 (pp. 1563-1566) Madrid, Spain.
- Argyle, M. (1972). Nonverbal communication in human social interaction. In R. Hinde (Ed.), Nonverbal Communication, (pp. 243-267). Cambridge: Cambridge Univ. Press.
- Bangalore, S., & Johnston, M. (2000). Integrating multimodal language processing with speech recognition. In B. Yuan, T. Huang & X. Tang (Eds.) Proceedings of the International Conference on Spoken Language Processing (ICSLP'2000) Vol.2 (pp. 126-129). Beijing: Chinese Friendship Publishers.
- Benoit, C. (2000). The intrinsic bimodality of speech communication and the synthesis of talking faces. In M. Taylor, F. Neel & D. Bouwhuis (Eds.), The Structure of Multimodal Dialogue II (pp. 485-502). Amsterdam: John Benjamins.
- Benoit, C., Guiard-Marigny, T., Le Goff, B., & Adjoudani, A. (1996). Which components of the face do humans and machines best speechread? In D.G. Stork & M.E. Hennecke, (Eds.), Speechreading by Humans and Machines: Models, Systems, and Applications: Vol. 150 of NATO ASI Series. Series F: Computer and Systems Sciences, (pp. 315-325). Berlin, Germany: Springer-Verlag.
- Benoit, C., & Le Goff, B. (1998). Audio-visual speech synthesis from French text: Eight years of models, designs and evaluation at the ICP. Speech Communication, 26, 117-129.
- Benoit, C., Martin, J.-C., Pelachaud, C., Schomaker, L., & Suhm, B. (2000). Audio-visual and multimodal speech-based systems. In D. Gibbon, I. Mertins & R. Moore (Eds.), Handbook of Multimodal and Spoken Dialogue Systems: Resources, Terminology and Product Evaluation. (pp. 102-203) Kluwer.
- Bernstein, L., & Benoit, C. (1996). For speech perception by humans or machines, three senses are better than one. Proceedings of the International Conference on Spoken Language Processing, (ICSLP 96) Vol. 3, (pp. 1477-1480). New York: IEEE press.
- Bers, J., Miller, S., & Makhoul, J. (1998). Designing conversational interfaces with multimodal interaction. DARPA Workshop on Broadcast News Understanding Systems, 319-321.

- Bolt, R. A. (1980). Put-that-there: Voice and gesture at the graphics interface. Computer Graphics, 14 (3), 262-270.
- Bregler, C., & Konig, Y. (1994). Eigenlips for robust speech recognition. Proc. of the Int'l. Conf. on Acoustics Speech and Signal Processing (IEEE-ICASSP) Vol.2, 669-672. IEEE Press.
- Bregler, C., Manke, S., Hild, H., & Waibel, A. (1993). Improving connected letter recognition by lipreading. Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE-ICASSP) Vol.1 557-560. Minneapolis, MN: IEEE Press.
- Brooke, N.M., & Petajan, E.D. (1986). Seeing speech: Investigations into the synthesis and recognition of visible speech movements using automatic image processing and computer graphics. Proceedings International Conference Speech Input and Output: Techniques and Applications, 258, 104-109.
- Calder, J. (1987). Typed unification for natural language processing. In E. Klein & J. van Benthem (Eds.). Categories, Polymorphisms, and Unification (pp. 65-72). Center for Cognitive Science, University of Edinburgh.
- Carpenter, R. (1990). Typed feature structures: Inheritance, (in)equality, and extensionality. Proceedings of the ITK Workshop: Inheritance in Natural Language Processing, 9-18. Tilburg: Institute for Language Technology and Artificial Intelligence, Tilburg University.
- Carpenter, R. (1992). The logic of typed feature structures. Cambridge, U. K.: Cambridge University Press.
- Cassell, J., Sullivan, J., Prevost, S., & Churchill, E. (Eds.) (2000). Embodied conversational agents. Cambridge, MA: MIT Press.
- Cheyner, A. (1998, January). MVIEW: Multimodal tools for the video analyst. International Conference on Intelligent User Interfaces (IUI'98), 55-62. New York: ACM Press.
- Cheyner, A., & Julia, L. (1995, May). Multimodal maps: An agent-based approach. International Conference on Cooperative Multimodal Communication, (CMC'95), 103-113. Eindhoven, The Netherlands.
- Choudhury, T., Clarkson, B., Jebara, T. & Pentland, S. (1999). Multimodal person recognition using unconstrained audio and video. Proceedings of the 2nd International Conference on Audio-and-Video-based Biometric Person Authentication, (pp. 176-181) Washington, DC.
- Codella, C., Jalili, R., Koved, L., Lewis, J., Ling, D., Lipscomb, J., Rabenhorst, D., Wang, C., Norton, A., Sweeney, P., & Turk C. (1992). Interactive simulation in a multi-person virtual world. Proceedings of the Conference on Human Factors in Computing Systems (CHI'92), 329-334. New York: ACM Press.
- Cohen, M.M., & Massaro, D.W. (1993). Modeling coarticulation in synthetic visual speech. In M. Magnenat-Thalman & D. Thalman (Eds.), Models and Techniques in Computer Animation. Tokyo: Springer-Verlag.
- Cohen, P. R., Cheyer, A., Wang, M., & Baeg, S. C. (1994). An open agent architecture. AAAI '94 Spring Symposium Series on Software Agents, 1-8. AAAI Press. (Reprinted in Huhns and Singh (Eds.). (1997). Readings in Agents (pp. 197-204). San Francisco: Morgan Kaufmann.)
- Cohen, P. R., Dalrymple, M., Moran, D. B., Pereira, F. C. N., Sullivan, J. W., Gargan, R. A., Schlossberg, J. L., & Tyler, S. W. (1989). Synergistic use of direct manipulation and natural language. Proceedings of the Conference on Human Factors in Computing Systems (CHI'89), 227-234. New York: ACM Press. (Reprinted in Maybury & Wahlster (Eds.), (1998). Readings in Intelligent User Interfaces (pp. 29-37). San Francisco: Morgan Kaufmann.)
- Cohen, P. R., Johnston, M., McGee, D., Oviatt, S., Pittman, J., Smith, I., Chen, L., & Clow, J. (1997). Quickset: Multimodal interaction for distributed applications. Proceedings of the Fifth ACM International Multimedia Conference, 31-40. New York: ACM Press.
- Cohen, P. R., McGee, D. R., & Clow, J. (2000). The efficiency of multimodal interaction for a map-based task. Proceedings of the Language Technology Joint Conference (ANLP-NAACL 2000), 331-338. Seattle: Association for Computational Linguistics Press.
- Cohen, P. R., & Oviatt, S. L. (1995). The role of voice input for human-machine communication. Proceedings of the National Academy of Sciences, 92 (22), 9921-9927. Washington, D. C.: National Academy of Sciences Press.
- Condon, W.S. (1988). An analysis of behavioral organization. Sign Language Studies, 58, 55-88.
- Dahlbäck, N., Jönsson, A., & Ahrenberg, L., (1992, January). Wizard of Oz studies – why and how. In Wayne D. Gray, William E. Hefley, & Dianne Murray (Eds.), Proceedings of the International Workshop on Intelligent User Interfaces (pp. 193-200). New York: ACM Press.
- Denecke, M. & Yang, J., (2000). Partial Information in Multimodal Dialogue. Proceedings of the International Conference on Multimodal Interaction, 624-633. Beijing, China.
- Duncan, L., Brown, W., Esposito, C., Holmback, H., & Xue, P. (1999). Enhancing virtual maintenance environments with speech understanding. Boeing M&CT TechNet.
- Dupont, S., & Luettin, J. (2000, September). Audio-visual speech modeling for continuous speech recognition. IEEE Transactions on Multimedia, 2(3) 141-151. Piscataway, NJ: Institute of Electrical and Electronics Engineers.

- Ekman, P. (1992, January). Facial expressions of emotion: New findings, new questions. American Psychological Society, 3(1), 34-38.
- Ekman, P., & Friesen, W. (1978). Facial Action Coding System. Consulting Psychologists Press.
- Fell, H., Delta, H., Peterson, R., Ferrier, L., Mooraj, Z., & Valleau, M. (1994). Using the baby-babble-blanket for infants with motor problems. Proceedings of the Conference on Assistive Technologies (ASSETS'94), 77-84. Marina del Rey, CA.
- Fridlund, A. (1994). Human facial expression: An evolutionary view. New York: Academic Press.
- Fukumoto, M., Suenaga, Y., & Mase, K. (1994). Finger-pointer: Pointing interface by image processing. Computer Graphics, 18(5), 633-642.
- Fuster-Duran, A. (1996). Perception of conflicting audio-visual speech: an examination across Spanish and German, In D.G. Stork & M.E. Hennecke (Eds.), Speechreading by Humans and Machines: Models, Systems and Applications, (pp. 135-143). New York: Springer Verlag.
- Hadar, U., Steiner, T.J., Grant, E.C., & Clifford Rose, F. (1983). Kinematics of head movements accompanying speech during conversation. Human Movement Science, 2, 35-46.
- Hauptmann, A. G. (1989). Speech and gestures for graphic image manipulation. Proceedings of the Conference on Human Factors in Computing Systems (CHI'89), Vol. 1 241-245. New York: ACM Press.
- Holzman, T. G. (1999). Computer-human interface solutions for emergency medical care. Interactions, 6(3), 13-24.
- Huang, X., Acero, A., Chelba, C., Deng, L., Duchene, D., Goodman, J., Hon, H., Jacoby, D., Jiang, L., Loynd, R., Mahajan, M., Mau, P., Meredith, S., Mughal, S., Neto, S., Plumpe, M., Wang, K. & Wang, Y. (2000). MiPad: A next-generation PDA prototype. Proceedings of the International Conference on Spoken Language Processing (ICSLP 2000) Vol.3, (pp. 33-36). Beijing, China: Chinese Military Friendship Publishers.
- Johnston, M., Cohen, P. R., McGee, D., Oviatt, S. L., Pittman, J. A., & Smith, I. (1997). Unification-based multimodal integration. Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics, 281-288. San Francisco: Morgan Kaufmann.
- Karshmer, A. I. & Blattner, M. (organizers). (1998). Proceedings of the 3rd International ACM Proceedings of the Conference on Assistive Technologies (ASSETS'98). Marina del Rey, CA. (URL <http://www.acm.org/sigcaph/assets/assets98/assets98index.html>).
- Kendon, A. (1980). Gesticulation and speech: Two aspects of the process of utterance. In M. Key (Ed.), The Relationship of Verbal and Nonverbal Communication (pp 207-227). The Hague: Mouton.
- Kobsa, A., Allgayer, J., Reddig, C., Reithinger, N., Schmauks, D., Harbusch, K., and Wahlster, W. (1986). Combining Deictic Gestures and Natural Language for Referent Identification. Proc. 11th International Conf. On Computational Linguistics, 356-361. Bonn, Germany.
- Koons, D., Sparrell, C. & Thorisson, K. (1993). Integrating simultaneous input from speech, gaze, and hand gestures. In M. Maybury (Ed.), Intelligent Multimedia Interfaces. (pp. 257-276) Cambridge, MA: MIT Press.
- Kricos, P.B. (1996). Differences in visual intelligibility across talkers, In D.G. Stork & M.E. Hennecke (Eds.), Speechreading by Humans and Machines: Models, Systems and Applications, (pp. 43-53). New York: Springer Verlag.
- Kumar, S., & Cohen, P.R. (2000). Towards a fault-tolerant multi-agent system architecture. Fourth International Conference on Autonomous Agents 2000, 459-466. Barcelona, Spain: ACM Press.
- Leatherby, J. H., & Pausch, R. (1992, July). Voice input as a replacement for keyboard accelerators in a mouse-based graphical editor: An empirical study. Journal of the American Voice Input/Output Society, 11(2).
- Martin, D. L., Cheyer, A. J., & Moran, D. B. (1999). The open agent architecture: A framework for building distributed software systems. Applied Artificial Intelligence, 13, 91-128.
- Massaro, D.W. (1996). Bimodal speech perception: A progress report, In D.G. Stork & M.E. Hennecke (Eds.), Speechreading by Humans and Machines: Models, Systems and Applications, (pp. 79-101). New York: Springer Verlag.
- Massaro, D.W., & Cohen, M.M. (1990). Perception of synthesized audible and visible speech. Psychological Science, 1(1), 55-63, January.
- Massaro, D.W. & Stork, D. G. (1998). Sensory integration and speechreading by humans and machines. American Scientist, 86, 236-244.
- McGrath, M. & Summerfield, Q. (1985). Intermodal timing relations and audio-visual speech recognition by normal-hearing adults. Journal of the Acoustical Society of America, 77 (2), 678-685.
- McGurk, H. & MacDonald, J. (1976). Hearing lips and seeing voices, Nature, 264, 746-748.
- McLeod, A. & Summerfield, Q. (1987). Quantifying the contribution of vision to speech perception in noise. British Journal of Audiology, 21, 131-141.
- McNeill, D. (1992). Hand and mind: What gestures reveal about thought. Chicago, IL: University of Chicago Press.

- Meier, U., Hürst, W. & Duchnowski, P. (1996). Adaptive bimodal sensor fusion for automatic speechreading. Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE-ICASSP), 833-836. IEEE Press.
- Morimoto, C., Koons, D., Amir, A., Flickner, M., & Zhai, S. (1999). Keeping an Eye for HCI. Proceedings of SIBGRAP'99, XII Brazilian Symposium on Computer Graphics and Image Processing, 171-176.
- Naughton, K. (1996). Spontaneous gesture and sign: A study of ASL signs co-occurring with speech. In L. Messing (Ed.), Proceedings of the Workshop on the Integration of Gesture in Language & Speech, (pp. 125-134). Univ. of Delaware.
- Neal, J. G., & Shapiro, S. C. (1991). Intelligent multimedia interface technology. In J. Sullivan & S. Tyler (Eds.), Intelligent User Interfaces (pp.11-43). New York: ACM Press.
- Negroponte, N. (1978, December). The Media Room. Report for ONR and DARPA. Cambridge, MA: MIT, Architecture Machine Group.
- Neti, C., Iyengar, G., Potamianos, G., & Senior, A. (2000). Perceptual interfaces for information interaction: Joint processing of audio and visual information for human-computer interaction. In B. Yuan, T. Huang & X. Tang (Eds.), Proceedings of the International Conference on Spoken Language Processing (ICSLP'2000), Vol. 3, (pp. 11-14). Beijing, China: Chinese Friendship Publishers.
- Oviatt, S. L. (1997). Multimodal interactive maps: Designing for human performance. Human-Computer Interaction [Special issue on Multimodal Interfaces], 12, 93-129.
- Oviatt, S. L. (1999a). Mutual disambiguation of recognition errors in a multimodal architecture. Proceedings of the Conference on Human Factors in Computing Systems (CHI'99), 576-583. New York: ACM Press.
- Oviatt, S.L. (1999b) Ten myths of multimodal interaction. Communications of the ACM, 42(11), 74-81. New York: ACM Press. (Translated into Chinese by Jing Qin and published in the Chinese journal Computer Application.)
- Oviatt, S.L. (2000a). Multimodal system processing in mobile environments. Proceedings of the Thirteenth Annual ACM Symposium on User Interface Software Technology (UIST'2000), 21-30. New York: ACM Press.
- Oviatt, S.L. (2000b). Taming recognition errors with a multimodal architecture. Communications of the ACM, 43 (9), 45-51. New York: ACM Press.
- Oviatt, S. L. (2000c). Multimodal signal processing in naturalistic noisy environments. In B. Yuan, T. Huang & X. Tang (Eds.), Proceedings of the International Conference on Spoken Language Processing (ICSLP'2000), Vol. 2, (pp. 696-699). Beijing, China: Chinese Friendship Publishers.
- Oviatt, S. L., Bernard, J., & Levow, G. (1999). Linguistic adaptation during error resolution with spoken and multimodal systems. Language and Speech, 41(3-4), 415-438 (special issue on "Prosody and Speech").
- Oviatt, S. L., & Cohen, P. R. (1991). Discourse structure and performance efficiency in interactive and noninteractive spoken modalities. Computer Speech and Language, 5(4), 297-326.
- Oviatt, S.L. & Cohen, P.R., (2000, March). Multimodal systems that process what comes naturally. Communications of the ACM, 43(3), 45-53. New York: ACM Press.
- Oviatt, S. L., Cohen, P. R., Fong, M. W., & Frank, M. P. (1992). A rapid semi-automatic simulation technique for investigating interactive speech and handwriting. Proceedings of the International Conference on Spoken Language Processing, 2, 1351-1354. Univ. of Alberta.
- Oviatt, S. L., Cohen, P. R., & Wang, M. Q. (1994). Toward interface design for human language technology: Modality and structure as determinants of linguistic complexity. Speech Communication, 15, 283-300. European Speech Communication Association.
- Oviatt, S.L., Cohen, P.R., Wu, L., Vergo, J., Duncan, L., Suhm, B., Bers, J., Holzman, T., Winograd, T., Landay, J., Larson, J. & Ferro, D. (2000). Designing the user interface for multimodal speech and gesture applications: State-of-the-art systems and research directions. Human Computer Interaction, 15(4), 263-322. (to be reprinted in J. Carroll (Ed.) Human-Computer Interaction in the New Millennium, Addison-Wesley Press: Boston, 2001).
- Oviatt, S. L., DeAngeli, A., & Kuhn, K. (1997). Integration and synchronization of input modes during multimodal human-computer interaction. Proceedings of Conference on Human Factors in Computing Systems (CHI'97), 415-422. New York: ACM Press.
- Oviatt, S. L., & van Gent, R. (1996). Error resolution during multimodal human-computer interaction. Proceedings of the International Conference on Spoken Language Processing, 2, 204-207. University of Delaware Press.
- Oviatt, S. L., & Kuhn, K. (1998). Referential features and linguistic indirection in multimodal language. Proceedings of the International Conference on Spoken Language Processing, 6, 2339-2342. Sydney, Australia: ASSTA, Inc.
- Oviatt, S. L. & Olsen, E. (1994). Integration themes in multimodal human-computer interaction. In Shirai, Furui, & Kakehi (Eds.) Proceedings of the International Conference on Spoken Language Processing, 2, (pp. 551-554). Acoustical Society of Japan.

- Pankanti, S., Bolle, R.M., & Jain, A. (Eds.), (2000). Biometrics: The future of identification. Computer, 33(2), 46-80.
- Pavlovic, V., Berry, G., & Huang, T. S. (1997). Integration of audio/visual information for use in human-computer intelligent interaction. Proceedings of IEEE International Conference on Image Processing, 121-124. IEEE Press.
- Pavlovic, V., & Huang, T. S., (1998). Multimodal prediction and classification on audio-visual features. AAAI'98 Workshop on Representations for Multi-modal Human-Computer Interaction, 55-59. Menlo Park, CA: AAAI Press.
- Pavlovic, V., Sharma, R., & Huang, T. (1997). Visual interpretation of hand gestures for human-computer interaction: A review. IEEE Transactions on Pattern Analysis and Machine Intelligence, 19(7) 677-695.
- Petajan, E.D. (1984). Automatic Lipreading to Enhance Speech Recognition, PhD thesis, University of Illinois at Urbana-Champaign.
- Poddar, I., Sethi, Y., Ozyildiz, E., & Sharma, R. (1998, November). Toward natural gesture/speech HCI: A case study of weather narration. In M. Turk, (Ed.), Proceedings 1998 Workshop on Perceptual User Interfaces (PUI'98) (pp. 1-6). San Francisco, CA.
- Robert-Ribes, J., Schwartz, J-L., Lallouache, T. & Escudier, P. (1998). Complementarity and synergy in bimodal speech: Auditory, visual, and auditory-visual identification of French oral vowels in noise. Journal of the Acoustical Society of America, 103(6) 3677-3689.
- Rogozan, A. & Deglise, P. (1998). Adaptive fusion of acoustic and visual sources for automatic speech recognition. Speech Communication, 26(1-2), 149-161.
- Rubin, P., Vatikiotis-Bateson, E., & Benoit, C. (Eds.). (1998). Audio-visual speech processing [Special issue]. Speech Communication, 26, 1-2.
- Rudnicky, A., & Hauptman, A. (1992). Multimodal interactions in speech systems. In M. Blattner & R. Dannenberg (Eds.), Multimedia Interface Design (pp.147-172). New York: ACM Press.
- Sekiyama, K., & Tohkura, Y. (1991). McGurk effect in non-English listeners: Few visual effects for Japanese subjects hearing Japanese syllables of high auditory intelligibility. Journal of the Acoustical Society of America, 90, 1797-1805.
- Seneff, S., Goddeau, D., Pao, C., & Polifroni, J. (1996). Multimodal discourse modelling in a multi-user multi-domain environment. In T. Bunnell & W. Idsardi (Eds.), Proceedings of the International Conference on Spoken Language Processing, Vol. 1 (pp. 192-195). University of Delaware & A.I. duPont Institute.
- Shaikh, A., Juth, S., Medl, A., Marsic, I., Kulikowski, C., & Flanagan, J. (1997). An architecture for multimodal information fusion. Proceedings of the Workshop on Perceptual User Interfaces (PUI'97), 91-93. Banff, Canada.
- Sharma, R., Huang, T.S., Pavlovic, V.I., Schulten, K., Dalke, A., Phillips, J., Zeller, M., Humphrey, W., Zhao, Y., Lo, Z., & Chu, S. (1996, August). Speech/gesture interface to a visual computing environment for molecular biologists. Proceedings of 13th International Conference on Pattern Recognition (ICPR 96), Vol. 3, 964-968.
- Sharma, R., Pavlovic, V.I., & Huang, T.S. (1998). Toward multimodal human-computer interface. Proceedings IEEE, 86(5) [Special issue on Multimedia Signal Processing], 853-860.
- Silsbee, P.L., & Su, Q (1996). Audiovisual sensory intergration using Hidden Markov Models. In D.G. Stork & M.E. Hennecke (Eds.), Speechreading by Humans and Machines: Models, Systems and Applications, (pp. 489-504). New York: Springer Verlag.
- Siroux, J., Guyomard, M., Multon, F., & Remondeau, C. (1995, May). Modeling and processing of the oral and tactile activities in the Georal tactile system. International Conference on Cooperative Multimodal Communication, Theory & Applications. Eindhoven, Netherlands.
- Stork, D. G., & Hennecke, M. E. (Eds.). (1995). Speechreading by Humans and Machines. New York: Springer Verlag.
- Suhm, B. (1998). Multimodal interactive error recovery for non-conversational speech user interfaces. Ph.D. thesis, Fredericiana University. Germany: Shaker Verlag.
- Summy, W. H. & Pollack, I. (1954). Visual contribution to speech intelligibility in noise. Journal of the Acoustical Society of America, 26, 212-215.
- Summerfield, A.Q. (1992). Lipreading and audio-visual speech perception, Philosophical Transactions of the Royal Society of London, Series B, 335, 71-78.
- Tomlinson, M. J., Russell, M. J. & Brooke, N. M. (1996). Integrating audio and visual information to provide highly robust speech recognition. Proceedings of the International Conference on Acoustics, Speech and Signal Processing (IEEE-ICASSP) Vol. 2, 821- 824. IEEE Press.

- Turk, M. & Robertson, G. (Eds.). (2000). Perceptual user interfaces [Special issue]. Communications of the ACM, 43(3), 32-70.
- Vatikiotis-Bateson, E., Munhall, K.G., Hirayama, M., Lee, Y.V., & Terzopoulos, D. (1996). The dynamics of audiovisual behavior of speech. In D.G. Stork and M.E. Hennecke, (Eds.), Speechreading by Humans and Machines: Models, Systems, and Applications, Vol. 150 of NATO ASI Series. Series F: Computer and Systems Sciences. (pp. 221-232). Berlin, Germany: Springer-Verlag.
- Vo, M. T., & Wood, C. (1996). Building an application framework for speech and pen input integration in multimodal learning interfaces. Proceedings of the International Conference on Acoustics Speech and Signal Processing (IEEE-ICASSP) Vol.6, 3545-3548. IEEE Press.
- Vo, M. T., Houghton, R., Yang, J., Bub, U., Meier, U., Waibel, A., & Duchnowski, P. (1995). Multimodal learning interfaces. Proc. of the DARPA Spoken Language Technology Workshop.
- Wahlster, W. (1991). User and discourse models for multimodal communication. In Joseph W. Sullivan & Sherman W. Tyler (Eds.), Intelligent User Interfaces, chap. 3, (pp. 45-67). New York: ACM Press.
- Wahlster, W. (2001, March). SmartKom: multimodal dialogs with mobile web users. Proceedings of the Cyber Assist International Symposium, 33-34. Tokyo International Forum.
- Waibel, A., Suhm, B., Vo, M.T., & Yang, J. (1997, April). Multimodal interfaces for multimedia information agents. Proceedings of the International Conference on Acoustics Speech and Signal Processing (IEEE-ICASSP) Vol.1, 167-170. IEEE Press.
- Wang, J. (1995). Integration of eye-gaze, voice and manual response in multimodal user interfaces. Proceedings of IEEE International Conference on Systems, Man and Cybernetics, 3938-3942. IEEE Press.
- Wickens, C. D., Sandry, D. L., & Vidulich, M. (1983). Compatibility and resource competition between modalities of input, central processing, and output. Human Factors, 25, 227-248.
- Wu, L., Oviatt, S., & Cohen, P. (1999). Multimodal integration – A statistical view. IEEE Transactions on Multimedia, 1(4), 334-341.
- Zhai, S., Morimoto, C., & Ihde, S. (1999). Manual and gaze input cascaded (MAGIC) pointing. Proceedings of the Conference on Human Factors in Computing Systems (CHI'99), 246-253. New York: ACM Press.