

Models of ATTENTI Communication: From P

Creating computing and communication systems that sense and reason about human attention by fusing together information from multiple streams.

ONE OF THE MAIN RESULTS OF 20TH CENTURY COGNITIVE psychology is that, despite the overall impressive abilities of people to sense, remember, and reason about the world, our cognitive abilities are extremely limited in well-characterized ways. In particular, psychologists have found that people grapple with scarce attentional resources and limited working memory. Such limitations become salient

By Eric Horvitz, Carl Kadie, Tim Paek, and David Hovel

when people are challenged with remembering more than a handful of new ideas or items in the short term, recognizing important targets against a background pattern of items [4], or interleaving multiple tasks [5].

These results indicate that we cannot help but to inspect the world via a limited spotlight of attention. As such, we often generate clues implicitly and explicitly about what we are selectively attending to and how deeply we are focusing. Given constraints on attentional resources, it is no surprise that communication among people relies deeply on attentional signals. Psychologists and linguists studying communication have recognized that signaling and detecting attentional states lies at the heart of the fast-paced and fluid interactions that people have with one another when collaborating or communicating [1, 6]. Attentional cues are central in decisions about when to initiate or to make an effective contribution to a conversation or project. Beyond knowing when to speak or listen in a conversa-

tion, attention is critical in detecting that a conversation is progressing. More generally, detecting or inferring attention is an essential component of the overall process of *grounding*—converging in a shared manner on a mutual understanding of a communication [1].

The findings about our limited attentional resources—and about how we rely on attentional signals in collaborating—have significant implications for how we design computational systems and interfaces. Over the last five years, our team at Microsoft Research has explored, within the *Attentional User Interface* project, opportunities for enhancing computing and communication systems by treating human attention as a central construct and organizing principle. We consider attention a rare commodity—and critical cur-



ON in Computing and Principles to Applications

rency—in reasoning about the information awareness versus disruption of users [11]. We have also pursued the use of attentional cues as an important source of rich signals about goals, intentions, and topics of interest [9]. We seek to build systems that sense, and share with users, natural signals about attention to support conversations and other forms of fluid mixed-initiative collaborations with computers. Moving to considerations of computational efficiency, an assessment of a user's current and future attention can be employed to triage computational resources. Investigations in this realm include selective allocation of resources in rendering graphics via relying on models or on direct observations of visual attention, and in guiding precomputation and prefetching [10] with forecasts of future attention. Finally, although there is a rich history of prior work on attention from cognitive psychology, we have found there is much we do not yet understand. Thus, beyond pooling results from prior psychological studies, we need to continue to perform user studies that adapt or extend prior results on attention and memory from cognitive psychology to real-world computing and communication applications [2, 3].

We shall first describe several principles and methodologies at the heart of research on integrating models of attention into human-computer interaction and communications. Then, we shall review representative efforts illustrating how we can harness these principles in attention-sensitive messaging and mixed-initiative interaction applications.

Models of Attention and Decision Making Under Uncertainty

How might we access and use information about a user's attention? To be sure, subtle clues about attention are often available, and a number of these clues can be taken as direct signals about the attentional status of users. For example, sensing patterns of simple gestures such as the touching and lifting of a device in different settings can relay important information about attention that can be exploited in a number of exciting ways [7]. Moving to higher-precision sensing, several researchers have pursued the use of gaze-tracking systems, and have used signals about the focus of visual attention in a variety of applications. As gaze sensors grow in reliability and decrease in cost, we are seeing the evolution of devices that recognize when and how they are interrogated by the spotlight of visual attention.

Nonetheless, we may often be uncertain about a user's attentional focus and workload in light of observations, and about the value of alternate actions in different contexts. Thus, we turn to models that can be harnessed to reason about a user's attention and about the ideal attention-sensitive actions to take under uncertainty. Such models and reasoning can unleash new functionalities and user experiences.

We have constructed by hand and learned from data Bayesian models viewed as performing the task of an automated "attentional Sherlock Holmes," working to reveal current or future attention *under uncertainty* from an ongoing stream of clues. Bayesian attentional models take

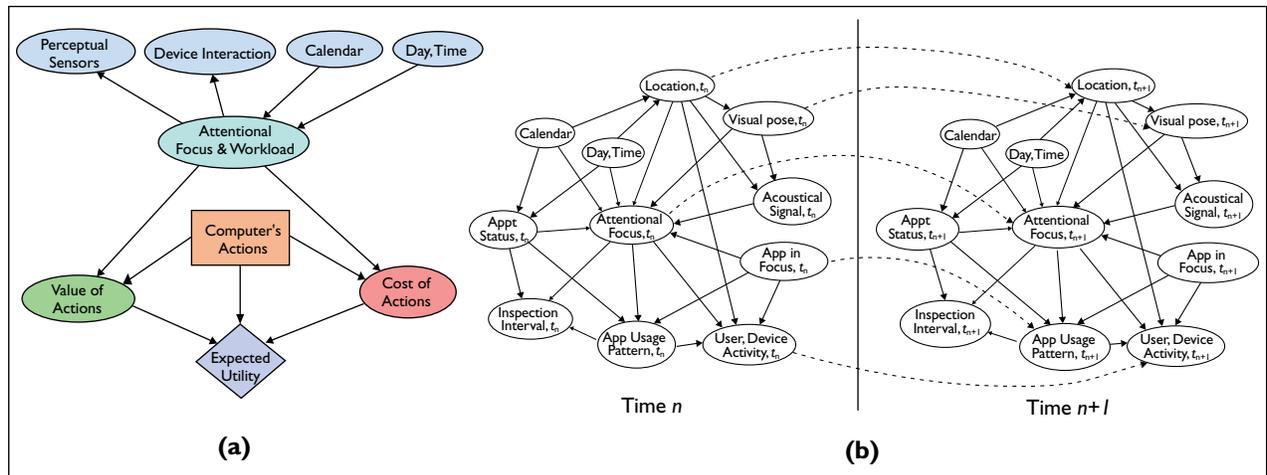


Figure 1. (a) High-level decision model considering a user's attentional focus and workload as a random variable, influenced by the observed states of several sensors. (b) A temporal Bayesian attentional model, highlighting key dependencies (dashed arcs) between variables in adjacent time slices.

as inputs sensors that provide streams of evidence about attention and provide a means for computing probability distributions over a user's attention and intentions.

Perceptual sensors include microphones listening for ambient acoustical information or utterances, cameras supporting visual analysis of a user's gaze or pose, accelerometers that detect patterns of motion of devices, and location sensing via GPS and analysis of wireless signals. However, more traditional sources of events can also offer valuable clues. These sources include a user's online calendar and considerations of the day of week and time of day. Another rich stream of evidence can be harvested by monitoring a user's interactions with software and devices. Finally, background information about the history of a user's interests and prior patterns of activities and attention can provide valuable sources of information about attention.

To build probabilistic attentional models able to fuse evidence from multiple sensors, we leverage the results of accelerated research over the last 15 years on representations for reasoning and decision making under uncertainty. Such work has led to inferential methods and representations including Bayesian networks and influence diagrams—graphical models that extend probabilistic inference to considerations of actions under uncertainty. Algorithms have been developed that enable us to compute probability distributions over outcomes and expected utilities of actions from these graphical representations.

Figure 1a displays a high-level influence diagram representing sensor fusion and decision making in

the context of a user's attention under uncertainty. As portrayed in the figure, a set of variables (oval nodes) representing sensed evidence influence a random variable representing a user's attentional status which, in turn, influences the expected value of alternate actions or configurations. We introduce intermediate cost and benefit variables in the pedagogical model as it can be useful to deliberate about the value and costs associated with different outcomes. Decisions (rectangular node) about ideal computer actions take into consideration the costs and benefits, given uncertainty about a user's attention. In the end, the expected utility (diamond-shaped node) is influenced by the action and the costs and benefits.

We extend such a high-level, pedagogical view by constructing richer models that contain additional intermediate variables and key interdependencies among the variables. Also, as both devices and people are immersed in time, we move beyond pointwise considerations of the states of variables, to build higher-fidelity temporal attentional models that represent changing observations and beliefs with the flow of time. We have employed dynamic Bayesian networks and Hidden Markov Models for representing and reasoning about states of attention and location over time.

Figure 1b displays two adjacent time slices of a temporal attentional model. Such a model provides a probability distribution over a user's workload and task developed for an application that provides selective filtering of messages and communications to users. In this case, the status of attention includes approximately 15 discrete states.

Economic Models of Attention and Information

As we can all attest from personal experiences, computers and communication systems today have little

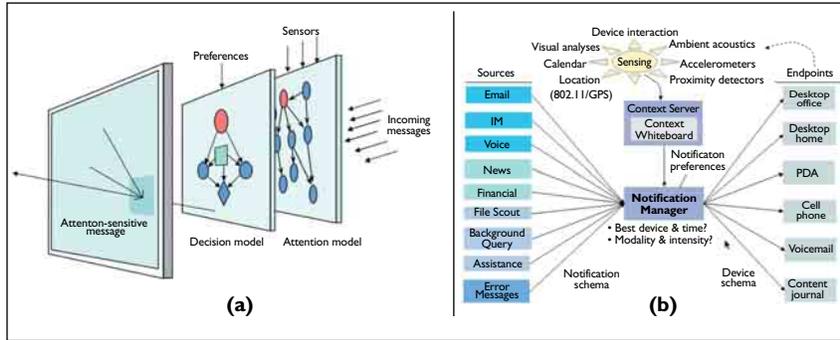


Figure 2. Conceptual overview of the Notification Platform, a cross-device messaging system that balances the costs of disruption with the value of information from multiple message sources. The system employs a probabilistic model of attention and executes ongoing decision analyses about ideal alerting, fidelity, and routing. (b) Constellation of components of Notification Platform, depicting the subscription architecture. Subscribed sources and devices communicate with the Notification Manager via a set of standard interfaces. Sensor findings from multiple devices are considered in deliberations about information value, attention, and the best channel and alerting modality.

awareness of the value and costs of relaying messages, alerts, and calls to users. Research on the *Notification Platform* project has centered on formulating economic principles of attention-sensitive notification—and on implementing a cross-device-alerting system based on these principles. A descendant of the Notification Platform named *Bestcom* applies similar principles to interpersonal communications [12]. We focus here on the Notification Platform.

The Notification Platform system modulates the flow of messages from multiple sources to devices by performing ongoing decision analyses. These analyses balance the expected value of information with the attention-sensitive costs of disruption. As highlighted in Figure 2a, the system serves as an attention-savvy layer between incoming messages and a user, taking as inputs sensors that provide information about a user’s attention, location, and overall situation.

The design of the Notification Platform was informed by several earlier prototypes exploiting context-sensing for identifying a user’s workload, including the *Priorities* system [11, 12]. *Priorities* employs classifiers that predict the urgency of incoming email. The classifiers are trained with sample messages, either obtained via explicit training or by automatically *drafting* data sets by observing a user’s interaction with an email browser. Studies have demonstrated the system performs remarkably well at classifying the urgency of messages (see, for example, the receiver-operator characteristic curve described in [11]). Beyond classifying the urgency of messages, *Priorities* also observes a user’s patterns of presence at a desktop computer based on time of day, and infers

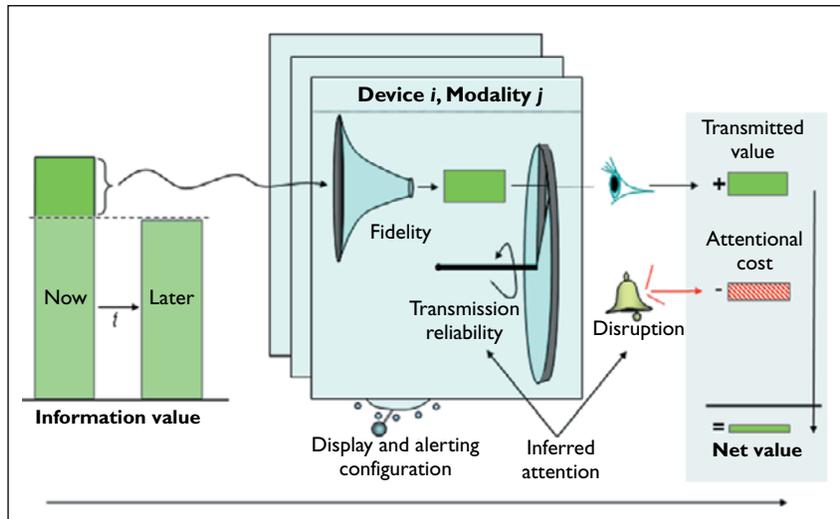
the time until a user will review unread messages. The system computes an expected *cost of delayed review* for each incoming message. This cost is considered, along with a cost of interruption based on activity sensing and calendar information, in automated decisions about if and how to alert and transmit information to a user about email, tasks, and appoint-

ment reminders in mobile and desktop settings.

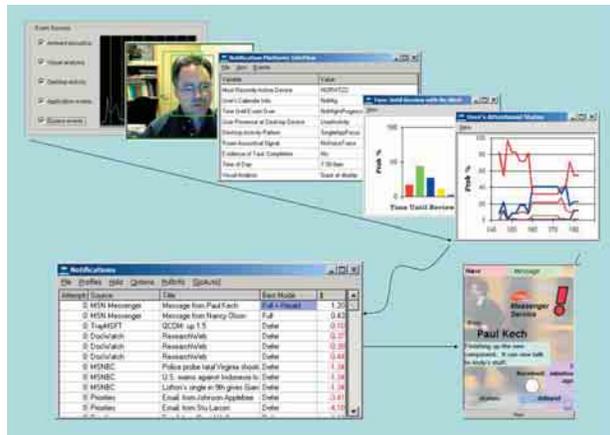
The Notification Platform uses a decision-analytic model for cross-device alerting and relay of information from multiple sources. The analyses consider a user’s attention and location under uncertainty, as well as the fidelity and relevance of potential communication channels. We developed a distributed architecture that executes over multiple devices. Figure 2b displays a schematized view of the architecture of the Notification Platform. Standard interfaces and metadata schemas allow users to subscribe different sources of information and devices to a *Notification Manager*. At the heart of the Notification Manager is a Bayesian attentional model and decision analysis that accesses clues about attention and location from sensors via a module we refer to as a *Context Server*.

The Context Server accesses several states and streams of evidence, including a user’s appointments from Microsoft Outlook, events about device presence and activity, an analysis of ambient acoustics in the room, and a visual analysis of pose using a Bayesian head-tracking system. Key abstractions from the evidence, such as “voice trace detected,” “task completion occurred within five seconds,” “single application focus,” “head-tracked—looking away from display,” and “meeting away from office—ending in 10 minutes,” are posted to a volatile store called the *Context Whiteboard*, which is continually updated by incoming evidence. The Context Whiteboard is contacted for updated information every few seconds by the Bayesian attentional model in the Notification Manager.

The Notification Manager’s decision analysis weighs the expected costs and benefits of alerting a user about messages coming into the system’s *Universal Inbox*. In computing the costs of disruption, the decision model considers the probability distribution over a user’s attentional state and location in several places in its analysis, including the cost of disruption associated with different alerts for each device, the availability of different devices, and the likelihood that the information will reach the user when alerted in a specific manner on a device.



(a)



(b)

The ongoing expected-utility analysis is performed in accordance with a user's preferences, stored in a profile. These include assertions about the cost of disruption for each alert modality, conditioned on the user being in different attentional states. As an example, for the case of a desktop computer, the system makes available a set of display alternatives as the product of different visual displays of the alert (for example, thumbnail, full-display alert) and several auditory cues (for example, no auditory cue, soft chime, louder herald). The placement of the alert with regard to the current focus of visual attention or interaction is also considered.

Figure 3a captures in a graphical manner the deliberation of the Notification Platform about incoming messages. The system computes the expected value of receiving an alert as the difference between the value of alerting the user now and the value that will be obtained when the information is viewed later. Given probability distributions over a user's attention and location inferred from its sensors, Notification Platform iterates over all alerting and display modalities

Figure 3. (a) Graphical depiction of the Notification Manager's analyses. Attention-sensitive costs of disruption and the value of information are considered, along with losses based in decreased fidelity (narrowing funnel) and transmission reliability (spinning slotted disk) associated with the use of each alerting modality of all subscribed devices. (b) View of a portion of the Notification Platform's real-time reasoning. Information from multiple sensors is posted to the Context Whiteboard and fused to infer the user's attentional status and location. Multiple notifications are sorted by net expected value and the best channel and alerting modality with the highest expected utility is selected.

for each device with an expected-utility analysis to decide if, when, and how to alert a user. As represented with the metaphor of a narrowing funnel in Figure 3a, the system considers, for each device and modality, the loss in fidelity of information transmitted. In addition, the system considers the likelihood that an alert will be received, given inferred probability distributions over the attention and location of the user. This reliability of transmission is represented metaphorically in the figure as the chance that a message will make it through a slot in a spinning disk. In the end, the attention-sensitive costs of disruption are subtracted from estimates of the value of alerting, yielding a net value of alerting a user for each channel and alerting modality. The channel and modality with the highest expected value is selected.

Figure 3b displays several aspects of the behind-the-scenes functioning of the Notification Platform. A *context palette* displays current findings drawn from sensor sources. Several views onto components of the decision analysis are displayed, including inference about the time-varying attention of the user. At the current time, the user is inferred to be most likely in a state named "high-focus solo activity," which has competed recently with "low-focus solo activity," "conversation in office," and other less likely states. The Universal Inbox displays messages from several sources, including email, instant messaging, breaking news, and stock prices. Messages have also been received from *DocWatch*, an agent subscribed to by the user that identifies documents of interest for the user. Each message is annotated with the best device and alerting policy, and the associated net expected dollar value of relaying the messages with that channel and mode is indicated. As portrayed in the inbox, it is worthwhile passing two instant messages on to the user. Other alerts are "in the red," as the cost of



Figure 4. Sensing PDA, outfitted with multiple perceptual sensors, including proximity, motion, and touch sensors. In the background, accelerometer signals are displayed showing the motion fingerprint of a user walking while looking at the device.

disruption dominates the net value of information. In this case, the ideal alerting mode and channel for an instant message is determined to be a visual notification in a large format coupled with an audio herald at the user's desktop system.

Ongoing research on the Notification Platform project includes the refinement of preference assessment tools to ease the task of encoding preferences. Currently, users can adjust sliders to change a set of predefined defaults on costs of interruptions. Another key area of work centers on using machine learning for building probabilistic models of attention, location, and cost of disruption from data. Results from ongoing machine-learning efforts by our team have been applied to refine the Notification Platform [12].

As highlighted in Figure 4, we have also been working to make small devices aware of the attentional status and location of users [7]—and either transmitting local sensor information to inform a central Notification Manager, performing entirely local notification management and related services based on the observations, or doing a combination of central and local deliberation about notification. In the latter case, the central Notification Manager makes general decisions about routing, and relies on the endpoint device to perform precision targeting of the timing and alerting modality, based on local sensing and reasoning. As an example, with the use of a method we refer to as “bounded deferral,” a local device commits to relaying a message that it has received before a message-specific deadline is reached; the device does its best to find a good time for interruption within the allotted period. Research on smart endpoints includes the challenges of embedding and leveraging multiple perceptual sensors on small

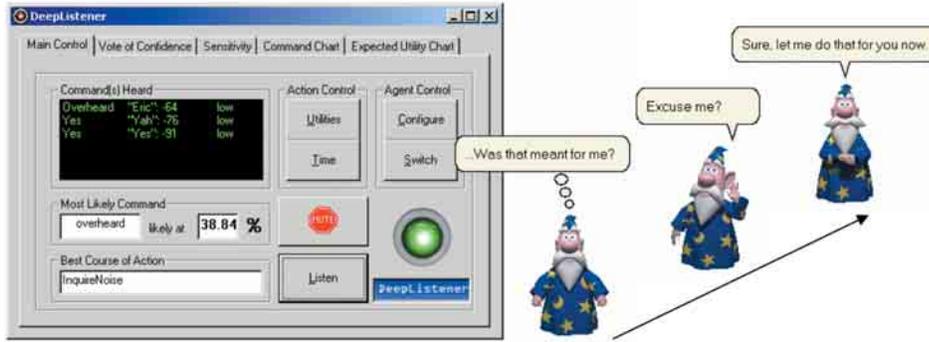
devices, including GPS, 802.11 signal strength, accelerometers, infrared proximity detectors, and touch sensors. Part of this work has explored opportunities for developing devices, such as cell phones that behave with more insight about their disruptiveness by considering the situation at hand, including states derived from coarse models of attention [8].

Additionally, we are continuing to pursue psychological studies of disruption. Formal studies of the costs of disruption began with the early work of Ovsiankina and Zeigarnik nearly 75 years ago. The rich body of work in this realm includes studies on memory, problem solving, and overall task efficiency in the face of disruptions. More recent work includes efforts by our team [2, 3] and other groups to probe the influence of notifications of various types and salencies on the efficiency and satisfaction with performing a variety of computer tasks. The psychological studies and results complement the mathematical models; the economic models provide a principled, flexible foundation that can integrate findings about the costs uncovered by user studies of disruption via the setting of parameters considered in expected-utility decision making.

Attention, Initiative, and Interaction

In another area of investigation within the Attentional User Interface project, we have studied the use models of attention to enhance the robustness and fluidity of human-computer collaboration. Some of this work focuses on the recognition of attentional cues as coordinative signals in *mixed-initiative* interaction with computing devices. In mixed-initiative interaction, both users and computers take turns in contributing to a project or an understanding [9]. The turn taking of conversational dialogue is a prototypical example of mixed-initiative interaction. Psychologists have found that people engaged in conversations rely on attentional cues to signal when a contribution is going to be offered or has been accepted [1]. We have sought to endow computers with an analogous ability to recognize and emit signals to guide the nature and timing of contributions and clarifications in support of mixed-initiative interaction.

DeepListener and Quartet represent efforts in mixed-initiative interaction to incorporate attention in spoken language systems. Both systems tackle what we have referred to as the “speech-target problem:” When a computer with an open microphone and speech recognizer detects an utterance, how is it to recognize that it is being addressed when there are other people or listening devices in a room? DeepListener and Quartet explicitly address this challenge



(a)

Figure 5. (a) DeepListener's deliberation about the target of speech and ideal clarification dialog. The system first makes an expected utility decision to share in a subtle manner its thoughts about the possibility that it is the target of an utterance. Given additional recognitions, it goes ahead to seek clarification, and finally executes an action for the user. (b) Quartet in action. Quartet's partial recognition is displayed at the top of the display. The system's belief about the attentional status of the user, with regards to initiating, maintaining, or breaking out of conversational dialogue, is represented as a dynamically changing probability distribution.



(b)

noisy environment. After analyzing a new utterance a bit later, the system engages the user in a clarification dialogue, and then invokes a desired action.

Quartet operates with a continuous speech recognition system, and incorporates a richer model of attention under uncertainty. It examines keyboard events, an analysis of the content and the coherence

of natural language parsing, and a visual pose analysis to ascertain the attentional status of the user and system with regards to the establishment, maintenance, and disruption of attention between the user and system.

Since Quartet couples speech recognition to a natural language parser, the system can also use the grammatical parse to reason about whether an utterance was misrecognized, or properly recognized but intended for someone else. Figure 5b shows Quartet listening to a user talking *about* the system rather than speaking *to* the system. In this case, Quartet is being used as an assistant to control, via voice commands, the navigation of slides displayed in a presentation. Requests directed *to* Quartet about navigation among slides arise intermittently during the more dominant stream of ongoing utterances associated with the presentation. In this example, the user is talking *about* the computer, and, based on a fusion of the user's language and visual pose, Quartet infers the user is likely speaking to someone else.

DeepListener reasons about a user's attention and intentions to guide clarification dialogue in a spoken command and control setting. The system considers its uncertainty about whether it is the target of speech, what it has heard, and the likelihood of different intentions. DeepListener continues to make expected-utility decisions about taking actions in the world, or about how it should approach users, if necessary, to clarify their intentions before taking such world actions. These decisions take into consideration the utilities of alternate dialogue actions and the stakes of the world actions.

DeepListener shares its attention and availability by gracefully changing the colors and intensities of an *attentional lens* that glows on its control panel, or via gestures and rendered "thoughts" of an animated agent. These affordances provide cues that assist with conversational turn taking.

DeepListener shares its attention and availability by gracefully changing the colors and intensities of an *attentional lens* that glows on its control panel, or via gestures and rendered "thoughts" of an animated agent. These affordances provide cues that assist with conversational turn taking.

Figure 5a displays a situation where DeepListener has detected an utterance first directed elsewhere in a

lenges, including the use of sensed or inferred attention to provide clues about a user's intentions, the content and context at hand, and the nature and ideal timing of the appropriate contributions. This work includes using sensed or inferred attention to inform speech recognition systems about the specific micro-contexts being addressed with utterances. Such narrowing of the spotlight of analysis can be useful for enhancing recognition as it can enable spoken dialogue systems to swap in the appropriate language models and semantics, and adjust the scope of possible actions. Also, robust solutions to the speech-target problem promise to significantly influence the overall sociology of human-computer interaction, by allowing users to interact with multiple devices and people in their proximity with speech and gestures in a manner similar to the way people interact with one another.

In another realm of innovation, computers with an ability to track and to understand attentional patterns among people engaged in conversations can provide new kinds of services and facilities. For example, methods for identifying visual attention among participants in a conversation can be used to automate the control of cinematography, and to capture, organize, and understand a group meeting or videoconference. Thus, beyond enhancing human-computer interaction, sensing and reasoning about attention promises to enhance the way we communicate and collaborate with one another.

Conclusion

We have described efforts to endow computing and communication systems with the ability to sense and reason about human attention. After reviewing some background on the nature and importance of attention in cognition and discourse, we discussed methods for inferring attention from multiple streams of information, and for leveraging these inferences in decision making under uncertainty. Then, we presented illustrative applications of the use of attentional models in cross-device, multichannel messaging and in mixed-initiative interaction. Research on the use of models of attention in computing and communication is still in its youth. We expect that continuing refinement of methods for recognizing, reasoning, and communicating about attention will change in a qualitative manner the way we perceive and work with computing systems and devices. **C**

CONTRIBUTORS ON THE CONSTELLATION OF EFFORTS ON THE ATTENTIONAL USER INTERFACE PROJECT INCLUDE JOHNSON APACIBLE, ED CUTRELL, MARY CZERWINSKI, SUSAN DUMAIS, KEN HINCKLEY, DAVID HOVEL, ANDY JACOBS, CARL KADIE, PAUL KOCH, JOHN KRUMM, NURIA OLIVER, TIM PAEK, JOHN PLATT, DANIEL ROBBINS, CHAITANYA SAREEN, JOE TULLIO, MAARTEN VAN DANTZICH, AND ANDY WILSON.

REFERENCES

1. Clark, H.H. and Schaefer, E.F. Collaborating on contributions to conversations. *Language and Cognitive Processes 2/1*, (1987), 19–41.
2. Cutrell, E., Czerwinski, M., and Horvitz, E. Notification, disruption, and memory: Effects of messaging interruptions on memory and performance. In *Proceedings of Interact 2001*. IFIP. Conference on Human-Computer Interaction, Tokyo, Japan, July 2001.
3. Czerwinski, M., Cutrell, E., and Horvitz, E. Instant messaging: Effects of relevance and time. *People and Computers XIV: Proceedings of HCI 2000*. S. Turner, P. Turner, eds. British Computer Society (2000), 71–76.
4. Eriksen, C.W. and Yeh, Y. Allocation of attention in the visual field. *J. Experimental Psychology: Human Perception and Performance 11* (1985), 583–597.
5. Gillie, T. and Broadbent, D. What makes interruptions disruptive? A study of length, similarity and complexity. *Psychological Research 50* (1989), 243–250.
6. Grosz, B.J. and Sidner, C.L. Plans for discourse. *Intentions in Communication*. MIT Press, Cambridge, MA (1990), 417–444.
7. Hinckley, K., Pierce, J., Sinclair, M., and Horvitz, E. Sensing techniques for mobile interaction. In *Proceedings of the ACM UIST 2000 Symposium on User Interface Software and Technology*. (San Diego, CA, Nov. 2000). ACM Press, NY, 91–100.
8. Hinckley, K. and Horvitz, E. Toward more sensitive mobile phones. In *Proceedings of the ACM UIST 2001 Symposium on User Interface Software and Technology*. ACM Press, NY, 191–192.
9. Horvitz, E. Principles of mixed-initiative user interfaces. In *Proceedings of the ACM SIGCHI 1999*. (Pittsburgh, PA, May 1999). ACM Press, NY, 159–166.
10. Horvitz, E. Principles and applications of continual computation. *Artificial Intelligence J. 126* (2001) Elsevier Science, New York, NY, 159–196.
11. Horvitz, E., Jacobs, A., and Hovel, D. Attention-sensitive alerting. In *Proceedings of the 15th Conference on Uncertainty and Artificial Intelligence*. (Stockholm, Sweden, July 1999). Morgan Kaufmann, San Francisco, CA, 305–313.
12. Horvitz, E., Koch, P., Kadie, C.M., and Jacobs, A. Coordinate: Probabilistic forecasting of presence and availability. In *Proceedings of the 18th Conference on Uncertainty and Artificial Intelligence*. (Edmonton, Canada, July 2002) 224–233.

THE COMPLETE SET OF REFERENCES FOR THIS ARTICLE CAN BE ACCESSED AT RESEARCH.MICROSOFT.COM/~HORVITZ/CACM-ATTENTION.HTM

ERIC HORVITZ (horvitz@microsoft.com) is a senior researcher and group manager of the Adaptive Systems & Interaction Group at Microsoft Research, Redmond, WA.

CARL KADIE (carlk@microsoft.com) is a research software design engineer at Microsoft Research, Redmond, WA.

TIM PAEK (timpack@microsoft.com) is a researcher in the Adaptive Systems + Interaction Group at Microsoft Research, Redmond, WA.

DAVID HOVEL (davidhov@exmsft.com) is a software developer at Microsoft Research, Redmond, WA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.