

User Awareness for Collaborative Multi-touch Interaction

Markus Schlattmann¹, Yuelong Yu², Nils Gruendl³, Manfred Bogen², Alexander Kulik³,
David d'Angelo², Bernd Froehlich³ and Reinhard Klein⁴

¹AGT Group (R&D) GmbH, Darmstadt, Germany

²Fraunhofer IAIS, Sankt Augustin, Germany

³Bauhaus-Universität, Weimar, Germany

⁴University of Bonn, Bonn, Germany

Keywords: Multi-touch, User Aware, Applications, Interaction Metaphors.

Abstract: Multi-touch enables direct manipulation of graphical computer interfaces. This intuitive interaction paradigm rapidly became popular, particularly in the domain of individual mobile devices. Collaborative work on large scale multi-touch devices, instead, suffers from mutual interferences if many simultaneous touch events cannot be attributed to individual users. On the example of a large, adjustable, high-resolution (4K) multi-touch device, we describe a lightweight and robust method to solve this issue. An additional depth camera above the tabletop device tracks the users around the table and their respective hands. This environment tracking and the multi-touch sensor are automatically calibrated to the same coordinate system. We discuss relevant implementation details to help practitioners building similar systems. Besides improved multi-user coordination we reveal general usability benefits including the reduction of false positives. Exploiting the developed system, we implemented several novel test applications to analyze the capabilities of such a system with regard to different interaction metaphors. Finally, we combined several of the analyzed metaphors to a novel application serving as an intuitive multi-touch application and environment for seismic interpretation.

1 INTRODUCTION

Many professional computer applications could benefit from immediate collaboration of colleagues and require technologies that support this. In the particular realm of seismic data interpretation in the oil and gas industry, we observed several issues originating from isolated desktop workplaces including the loss of information and the misinterpretation of results.

Tabletop computers with multi-touch (MT) input offer a promising alternative. They can be easily operated by non-computer scientists and facilitate smooth communication and cooperation among different experts.

Unfortunately, most existing multi-touch systems suffer from missing user awareness. If a multi-touch system detects two touch points on the screen, it generally cannot distinguish whether the touch points belong to one hand, two hands or even different users. Therefore, multiple users and

multiple hands can only work in the same context, which often results in interference.

To solve these problems and allow for more natural collaboration, some previous research systems already included additional sensor technologies. However, these systems all suffered from severe limitations, restricting the quality of the tabletop display or its surrounding either, or even the movements/locations of the users themselves. Using a depth camera (e.g. Microsoft Kinect), instead, provides additional information that enables reliable and robust user tracking. This in turn enables the robust association of detected touch-points to individual users and their hands, which provides many new possibilities to improve the expressiveness of multi-touch gestures and realize software-supported multi-user multi-touch input coordination. In particular we identify the following applications:

- Individual users can associate different tools to their input; thus different functions may even be operated simultaneously by cooperating users.

- Software-controlled access management eliminates involuntary interference (e.g. during manipulation an object access is blocked for other users)
- Automatic partitioning of the screen and input space with respect to user positions.

The main objective of our work was the development of a functional tabletop prototype that provides all the features and the quality required for co-located collaboration on the interpretation of seismic data. For the moment we were focusing only on the acceptance of the system by expert users. In future work we are aiming to gain further insights into the usability of the system and its integration in the professional workflow.

2 RELATED WORK

Many researchers implemented multi-user coordination policies using the commercially available DiamondTouch system (Dietz and Leigh, 2001); (Morris et al., 2004); (Morris et al., 2006); (Morris et al., 2004). Unfortunately, the system limits the choice of display component to front projection. It is furthermore limited to four users that must maintain physical contact to the corresponding receiver unit while avoiding to touch each other.

Marquard et al., (2010) used a glove that tagged features of the hand with unambiguous optical markers for robust hand and finger identification. Roth et al. propose attaching small IR-emitting devices to the users' hands for cryptographically sound user identification (Roth et al., 2010). This is obviously the most secure implementation of user tracking for tabletop interaction, but requiring the users to wear an electronic device which breaks with the walk-up and use paradigm of most tabletop applications.

Dang et al., instead, proposed a system that does not require the user to wear a glove. They suggested a heuristic to identify the hands belonging to a touch point based on the orientation of the tracked ellipse (Dang et al., 2009). Besides the limitation that users must always touch the surface with the finger pad instead of the tip, the association is immediately lost once the fingers are released from the screen. More recently, Ewerling et al. proposed multi-touch sensing based on the maximally stable extremal regions algorithm that implicitly organizes touch points in hierarchies corresponding to individual hands (Ewerling et al., 2012). However, this approach only works with additional depth information above the tabletop and cannot distinguish individual users.

Towards the association of touch points to users, Walther-Franks et al. proposed to use additional proximity sensors in the frame of the tabletop device (Walther-Franks et al., 2008). The system provides rough user tracking for many purposes, but due to the dislocation of the touch sensing on the screen surface and the user tracking around its housing, a robust correlation of touch-points with a user's hand cannot be ensured. Touch points detected in close proximity to the user's body position, tracked at the edge of the tabletop device, may also belong to somebody else reaching into her proximity.

Recently, Annett et al. presented an improved version of such a system, using a much larger number of proximity sensors to derive a higher tracking density (Annett et al., 2011). In particular, their system incorporates sensors in upward orientation, tracking the user's arm above the display frame. This adaptation enables a more accurate association of tracked touch points to individual users. While this system is a significant improvement, the general limitations related to the missing overlap of both tracking systems remain. In a similar spirit, Richter et al. (Richter et al., 2012) suggested to capture the users' shoes for identification. The association with touch points cannot be realized robustly, but individual settings can automatically be applied based on the coming and going of users - if they keep wearing the shoes that are linked to them in the database.

Dohse et al. used an additional camera mounted above the tabletop display for the association of touch-points with users. The hands are tracked above the display using color segmentation or shadow tracking respectively (Dohse et al., 2008). The authors suggest cancelling the light from the screen with polarization filters to avoid interference with the color of displayed items. As an alternative, they propose tracking the dark silhouettes of the hands above the illuminated screen.

Another approach by Andy Wilson suggests using a depth sensing camera attached to the ceiling as an all-in-one sensor both for touch detection with defined interaction surfaces and context awareness (Wilson, 2010). An advanced follow up on this approach has been presented by Klompaker et al., (2012). Their framework provides the means for tracking tangible input devices and the users' hands for gesture recognition and touch input on arbitrary surfaces. In both cases, the same sensor that is exploited for touch sensing covers the whole surrounding area in depth, touch data can directly be associated with tracked users. As a downside of this concept, the relative low resolutions of the depth

cameras impede accurate multi-touch input. Martinez et al., (2011) consequently proposed to combine the context tracking of the depth camera with accurate touch sensing at the surface of the tabletop device. Unfortunately, they did not provide much detail about their implementation.

Jung et al. (Jung et al., 2011) proposed the combination of the Kinect depth sensor for body tracking and RFID for the automatic authentication of users on a multi-touch tabletop device. They realized a prototype demonstrating the benefits of the concept for applications with critical security issues. The overall usability of the system and the technical details necessary to achieve this were not discussed in the paper.

3 MULTI-TOUCH CONTEXT TRACKING

In this section, the method for computing the association of touch-points with respective users and hands is described.

3.1 Hardware Setup

Our high-resolution multi-touch table (3840x2160 pixels) measures 56" display diagonal. With the Kinect device mounted about 3 meters above the floor (see Figure 1) we can capture the entire screen area and about 50 cm of its surrounding in each direction.

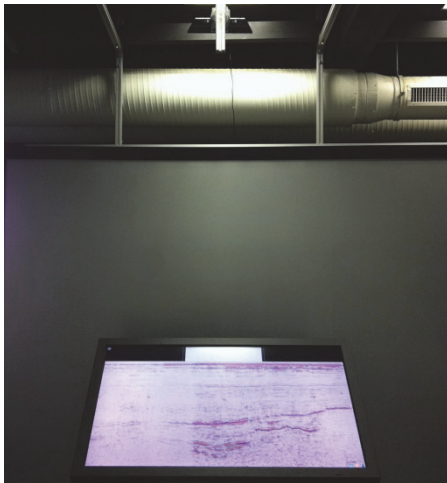


Figure 1: The arrangement of MT-table (Bottom) with Microsoft Kinect (Top).

A standard desktop PC is used to drive the

application on the multi-touch table. A second machine runs the environment tracking. Both workstations share the recorded user input data via network using the TUIO protocol (Kaltenbrunner et al., 2005).

3.2 Automatic Calibration and Calibration Detection

The internal calibration of the different sensors embedded in the Kinect (color, infrared, depth) must be performed only once as their relation does not change over time. We use the method published by Nicolas Burrus for this purpose (Burrus, 2011).

The extrinsic parameters defining the relation between the display area and one of the Kinect sensors, instead, must frequently be re-calibrated. These parameters have to be re-computed every time the display or the Kinect has been moved. Due to the adaptability of our assembly this can happen quite often. For ergonomic usage in different situations our touch table allows to adjust the height (75 - 125 cm) and the orientation (0 - 70 degrees) of the interactive display. Even the pivot of the display can be changed between landscape and portrait orientation.

We propose a fully automatic calibration routine without the necessity of any additional calibration object such as a printed chessboard pattern. The screen itself is employed to display a calibration pattern (chessboard) that can be recorded by the color sensor of the Kinect and processed for calibration with the method available in OpenCV (Bradski, 2000). Note that the calibration pattern on the display is not visible for the infrared or depth sensor of the Kinect. Additionally, the need for re-calibration induced by height or inclination changes is detected automatically using motion sensors attached to the multi-touch table.

3.3 Combined Segmentation

The sensors of the Microsoft Kinect enable different ways of segmenting the foreground (user body parts) from the background (e.g. floor and display). However, every sensor and according method has its benefits and drawbacks. We achieved the best results with a combination of the IR-intensity image and the depth information. Depth segmentation facilitates the background subtraction in the uncontrolled surrounding of the tabletop device. However, the imprecision and quantization of depth values obtained by the Kinect impede precise depth segmentation close to the display surface.

Furthermore, the effective image resolution of the depth image is low (approximately 320x240), wherefore small body parts (e.g. fingers or small hands) tend to vanish in the segmentation (see Figure 2(Left)).

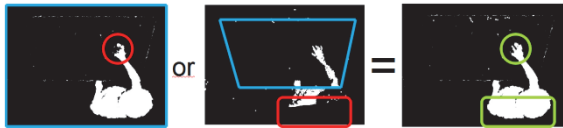


Figure 2: A combination of depth segmentation (left image) and IR-intensity segmentation (middle image) leads to better segmentation results (right image) and solves the respective problems (see regions in green circles). The red circles indicate the problematic regions in depth-based segmentation (left) and IR-based segmentation (middle), respectively.

The background subtraction performed on infrared (IR) intensity image from the Kinect provides a higher resolution (640x480) and also works close to the display surface as the display does not emit infrared light. Using a logical or-operation, we combine the depth-based segmentation of the entire image with the infrared segmentation, but only inside the image region corresponding to the display surface. The result is shown in Figure 2.

3.4 User Separation and Tracking

After the segmentation, only noise and those regions belonging to the users remain in the image. To filter high frequency noise we first perform an erosion step followed by a dilatation step. Then, we search for the largest connected components. A simple threshold on the minimum component size filters further artifacts of noise. The remaining connected components (CC) can be assumed to correspond to individual users. They are assigned a user ID and tracked over time (Figure 3).

A drawback of this approach is that as soon as two users touch or occlude each other, their separate CCs will merge. Thus, we apply a second processing step to separate the user regions also if a CC comprises multiple users. We exploit the depth information provided by the Kinect camera to identify the upper body of each user.

Our algorithm searches for height peaks within each connected component that are at least 40 cm above the known height of the tabletop surface. If only one of such peaks exists, the entire component is interpreted as a single user region. Otherwise, the region is separated into individual areas surrounding each peak based on frame coherence (i.e. pixels of

the CC are divided based on the nearest pixel assignment of the previous frame). Apparently, as this solution is based on frame coherence, it only works if at a certain time in the past the users did not occlude each other. Furthermore, it sometimes leads to erroneous separation in extreme cases, e.g. if the occlusion continues for a long time during complex movement of the users.

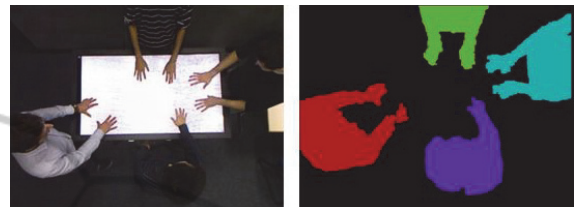


Figure 3: Example segmentation of four users.

3.5 Associating Touch-Points to Users and Hands

For each touch-point the corresponding user and hand have to be estimated. To this end, we project the touch-points to the image containing the user regions using the transformation matrix derived from the calibration (see Figure 4).

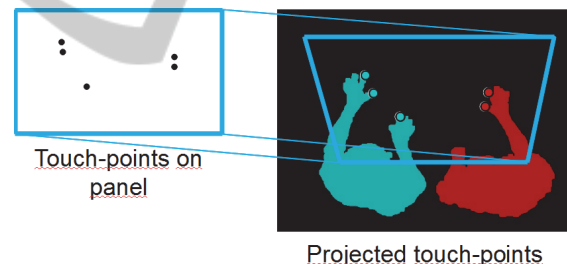


Figure 4: Projection and assignment of the touch-points (black dots on the left) to the image containing the user regions (right). The blue rectangle indicates the area of the display surface. The projected touch-points (colored dots with black margin on the right) are assigned to the closest user region.

If a touch-point is located inside a user region, the corresponding user ID is transcribed directly. Otherwise the closest user region is computed using a breadth-first search. We further distinguish both hands of a user based on the geodesic distances of touch-points in the graph representing the user component. Our procedure builds on the Dijkstra algorithm to find the shortest path between the touch-points assigned to one user. The procedure stops either if all the other touch-points of the respective user have been found or if the Dijkstra radius exceeds a certain threshold D (30 cm).

The resulting clusters are interpreted as individual hands. Unlike Euclidean distances, geodesic distances support robust clustering even if both hands are in close proximity. However, this approach is still limited to cases where the hand regions do not merge to one region in the camera images. If the hands merge, we can only transcribe the hand state of touch points from a previous state (if existing).

3.6 Identification of Involuntary Touches

If, for one respective user, more than two clusters are found we ignore those with the lowest frame coherence as the respective touch events most likely occurred involuntarily. We also classify a touch-point to be involuntary if the area of the touch footprint on the multi-touch panel is larger than a certain threshold (5cm^2). This way, touches too large to originate from a finger or soft-touch stylus are ignored. This simple thresholding allows users to rest their hands on the tabletop surface while operating with their fingers or a stylus.

A one-time invalid touch-point stays invalid. We also tried other mechanisms allowing touch-points to become valid again such as (adaptive) hysteresis thresholding. However, we observed this simple mechanism to work best. We observed that voluntary touches were hardly ever classified as involuntary.

4 APPLICATIONS AND FINDINGS

We first developed several test applications to explore the novel possibilities of our user aware tabletop system. These applications were informally, but continuously tested by visitors of our lab in order to identify which functionality best suites our objective. Note that these demonstrators were realized with an earlier multi-touch table prototype based on a back-projection display. After these tests we could identify the interaction techniques most suitable for the seismic interpretation application.

We later implemented a novel system based on a large high resolution LCD display. It was equipped with a multi-touch sensor frame that provides high precision and tracking robustness combined with low latency. The final system including the user tracking was set up with custom software for the exploration and interpretation of seismic data.

4.1 Test Applications

Territoriality has been shown to be of particular relevance for the coordination of multi-user cooperation (Scott et al., 2004). The separation between personal and group territories facilitate the coordination of individual and cooperative subtasks.

In the context of a game based on the exchange of virtual tokens we implemented a technique that assigns a particular area of the shared tabletop surface to each involved user (see Figure 5). The circular area is automatically placed in closest possible position to the user it is assigned to and serves as a storage area for managing the virtual tokens. Other users could not access the collected items.

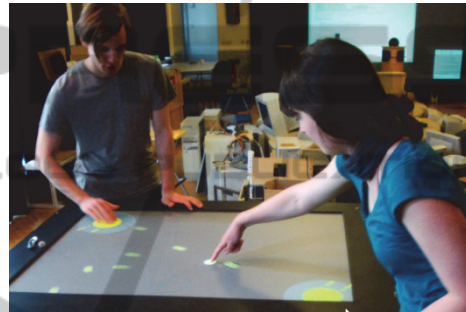


Figure 5: The yellow circles with gray surrounding define the users' personal storage space.

We observed that users immediately understood the concept of the personal storage areas that were following them. They appreciated that the areas were always situated relative to their body, thus facilitating access even while walking around the table. We also frequently observed the annoyance of users about the disappearing of their personal storage area when they stepped away too far from the table for a moment. We introduced some latency to the automatic area management which alleviated this issue a little. Only a robust user identification could solve this problem (e.g. (Jung et al., 2011); (Richter et al., 2012); (Roth et al., 2010)).

The **orientation of GUI elements** is a crucial issue in the realm of tabletop interfaces. Users generally surround the devices from all accessible sides, while some elements might only be legible from one direction. Shen et al. described a metaphoric "magnet" feature allowing the reorientation of all GUI-elements to improve legibility. They observed user conflicts with this global functionality (Shen et al., 2004). Turning all items to face a particular user necessarily turns the same items away from the others.

We reasoned that users would generally not be interested in orienting all available GUI elements towards themselves, but only try to improve the legibility of a few items they are focusing on. Thus we implemented an automatism that turned a selected item automatically towards the user who selected it. This concept works well if users only seek for their own reading comfort. In test applications including a photo viewer and a document reader, however, we observed that users also want to show items to others. The automatic alignment was interfering in these situations. Therefore we adapted the algorithm such that items still turned toward the user who selected them, but if they are moved (e.g. in the direction of another user) they turned toward the user being closest in the movement direction. This combination was well accepted by test users, but eventually we realized that the items can just always be oriented in the direction of a movement. This simplification works similarly well and does not even require user awareness.

Associating particular tools to the fingers of individual hands is probably the most relevant augmentation for multi-user multi-touch interaction. In a drawing application we observed that users just expected to keep a selected color for drawing until they had chosen a new one (see Figure 6). In fact, most of our visitors were puzzled when we explained this to be a novel feature. They could not imagine this to be any different. Consequently, they were surprised that the chosen association was lost when they stepped further away from the screen such that they left the tracking area for a moment.

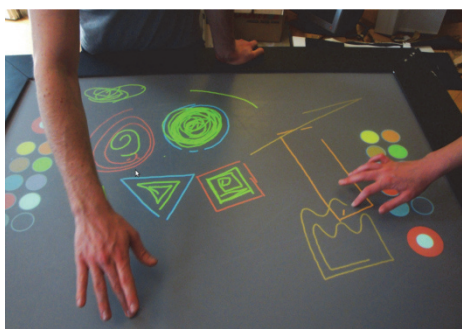


Figure 6: User awareness allows assigning individual tools to each touch event. Here, each user draws with an individually selected color.

Automatic **input coordination** can reduce mutual interference of multi-user input enormously. In the most basic configuration, GUI elements cannot be acquired by users while they are already manipulated by someone else. This ensures that each

user can complete a desired manipulation without disturbance. We experimented with this feature in a photo viewer application.

We found that more subtle adaptations to the users' input are often preferable compared to simple locking. In fact users often want to stop others from moving an element. To this end they touch the moving element in an attempt to stop it. Without any context awareness such conflicting input from multiple users generally results in scaling the respective graphics element. We implemented a coordination policy that only locks the type of possible manipulations to those already operated by the user who first acquired the element. Thus, if one is moving an item, others may still interfere to hold it, but this interference will not cause any other type of transformation, e.g. scaling.

In combination with individual tool selection this approach also allows multiple users to apply different operations simultaneously on the same element. One may for example continue to draw lines on an object while it is moved around by somebody else.

4.2 Seismic Interpretation Application

In the context of an industrial research project we realized an application for the collaborative analysis of seismic data on a multi-touch tabletop device. The application builds on two of the above mentioned techniques for user awareness. Furthermore, the display quality was improved by using a high-resolution LCD panel that offers multi-touch input. Our earlier prototype did not satisfy the visual quality required for the interpretation of the fine grained data (see Figure 7).

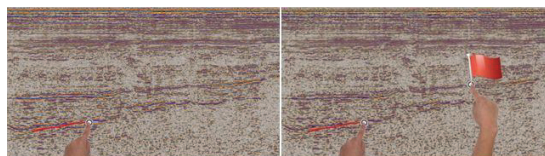


Figure 7: Illustration of individual tool selection. The left hand belongs to a user having enabled the line annotation tool while the right hand belongs to a different user having enabled the point annotation tool.

The application displays cutting planes derived from a volumetric dataset. The collaborative task is the interpretation of seismic features, searching for places with a high probability of hydrocarbon accumulation. The application primarily supports navigation along a horizontal cutting plane as well as selection and annotation of so called seismic

lines. In Figure 8 the collaborative work on such seismic lines is shown for two users including the depth image.



Figure 8: Left: Two users working on the same seismic line. Right: According depth image.

From the above mentioned functionality we implemented *input coordination* and *individual tool selection* (see Figure 7). From our earlier experiences with the test applications we expected that technologically unaware users would not even realize that they are tracked to ensure fluid multi-user multi-touch interaction.

Experts from geology were using our context-aware multi-touch system and we received very positive feedback. Also the display size, its impressive image quality with ultra-high resolution and the accuracy of the multi-touch input were highly appreciated. With the expectation of becoming frequent users of the system our visitors were delighted about the adaptation features of the assembly, which easily allow changing the height and inclination of the display. As expected, most of these experts did not realize the implicit input coordination to be an obvious feature. The system just worked as expected.

5 CONCLUSIONS

We analyzed state-of-the-art methods to achieve user awareness of multi-touch tabletop displays and derived an improved method from our experiments. Based on the sensor data from a depth camera we achieved robust context tracking. We described an implementation of this method including automatic re-calibration, sensor-fused segmentation, separation of users, robust hand identification based on geodesic distances and a detection method for involuntary touches. Based on the resulting user awareness, we suggested interaction techniques for fluent co-located collaboration. While the general idea of context tracking with an overhead camera has been proposed earlier, we contribute a detailed description of a timely method that is robust and easy to implement. Finally, we described a high-

fidelity system prototype and a collaborative application for the exploration and interpretation of seismic data.

We are looking forward to gain further insights on the usability of the system in long-term studies with expert users.

REFERENCES

- Annett, M., Grossman, T., Wigdor, D., & Fitzmaurice, G., (2011). Medusa: a proximity-aware multi-touch tabletop. *Proc. UIST 2011* (pp. 337-346). New York, NY, USA: ACM-Press.
- Bradski, G., (2000). The OpenCV Library. *Dr. Dobb's Journal of Software Tools*.
- Burrus, N. (2011, June). *Kinect RGB Demo v0.5.0*. Retrieved from <http://nicolas.burrus.name/index.php/Research/KinectRgbDemoV5>
- Dang, C. T., Straub, M., & André, E., (2009). Hand distinction for multi-touch tabletop interaction. *Proc. ITS 2009* (pp. 101-108). New York, NY, USA: ACM-Press.
- Dietz, P., & Leigh, D. (2001). DiamondTouch: a multi-user touch technology. *Proc. UIST 2001* (pp. 219-226). New York, NY, USA: ACM-Press.
- Dohse, K. C., Dohse, T., Still, J. D., & Parkhurst, D. J., (2008). Enhancing Multi-user Interaction with Multi-touch Tabletop Displays Using Hand Tracking. *Proc. Advances in Computer-Human Interaction 2008* (pp. 297-302). Washington, DC, USA: IEEE Computer Society.
- Ewerling, P., Kulik, A., & Froehlich, B., (2012). Finger and hand detection for multi-touch interfaces based on maximally stable extremal regions. *Proc. ITS 2012* (pp. 173-182). New York, NY, USA: ACM-Press.
- Jung, H., Nebe, K., Klompaker, F., & Fischer, H. (2011). Authentifizierte Eingaben auf Multitouch-Tischen. *Mensch & Computer 2011* (pp. 305-308). München: Oldenbourg Wissenschaftsverlag GmbH.
- Kaltenbrunner, M., Bovermann, T., Bencina, R., & Costanza, E., (2005). TUIO - A Protocol for Table-Top Tangible User Interfaces. *Proc. of the 6th International Workshop on Gesture in Human-Computer Interaction and Simulation*.
- Klompaker, F., Nebe, K., & Fast, A., (2012). dSensingNI: a framework for advanced tangible interaction using a depth camera. *Proc. of ACM TEI 2012* (pp. 217-224). New York, NY, USA: ACM-Press.
- Marquardt, N., Kiemer, J., & Greenberg, S., (2010). What caused that touch?: expressive interaction with a surface through fiduciary-tagged gloves. *Proc. ITS 2010* (pp. 139-142). New York, NY, USA: ACM-Press.
- Martinez, R., Collins, A., Kay, J., & Yacef, K., (2011). Who did what? Who said that?: Collaid: an environment for capturing traces of collaborative learning at the tabletop. *Proc. ITS 2011* (pp. 172-181).

- New York, NY, USA: ACM-Press.
- Morris, M. R., Huang, A., Paepcke, A., & Winograd, T., (2006). Cooperative gestures: multi-user gestural interactions for co-located groupware. *Proc. CHI 2006* (pp. 1201-1210). New York, NY, USA: ACM-Press.
- Morris, M. R., Ryall, K., Shen, C., Forlines, C., & Vernier, F., (2004). Beyond "social protocols": multi-user coordination policies for co-located groupware. *Proc. CSCW 2004* (pp. 262-265). New York, NY, USA: ACM-Press.
- Morris, R. M., Ryall, K., Shen, C., Forlines, C., & Vernier, F., (2004). Release, relocate, reorient, resize: fluid techniques for document sharing on multi-user interactive tables. *CHI Extended Abstracts 2004* (pp. 1441-1444). New York, NY, USA: ACM-Press.
- Richter, S., Holz, C., & Baudisch, P. (2012). Bootstrapper: recognizing tabletop users by their shoes. *Proc. CHI 2012* (pp. 1249-1252). New York, NY, USA: ACM-Press.
- Roth, V., Schmidt, P., & Güldenring, B., (2010). The IR ring: authenticating users' touches on a multi-touch display. *Proc. UIST 2010* (pp. 259-262). New York, NY, USA: ACM-Press.
- Scott, S. D., Sheelagh, M., Carpendale, T., & Inkpen, K. M., (2004). Territoriality in collaborative tabletop workspaces. *Proc. CSCW 2004* (pp. 294-303). New York, NY, USA: ACM-Press.
- Shen, C., Vernier, F. D., Forlines, C., & Ringel, M., (2004). DiamondSpin: an extensible toolkit for around-the-table interaction. *Proc. CHI 2004* (pp. 167-174). New York, NY, USA: ACM-Press.
- Walther-Franks, B., Schwarten, L., Krause, M., Teichert, J., & Herrlich, M., (2008). User Detection for a Multi-touch Table via Proximity Sensors. *Proceedings of the IEEE Tabletops and Interactive Surfaces*. IEEE Computer Society.
- Wilson, A. D., (2010). Using a depth camera as a touch sensor. *Proc. ITS 2010* (pp. 69-72). New York, NY, USA: ACM-Press.

