

An Augmented Reality-based Remote Collaboration Platform for Worker Assistance

Georgios Chantziaras¹, Andreas Triantafyllidis¹, Aristotelis Papaprodromou¹,
Ioannis Chatzikonstantinou¹, Dimitrios Giakoumis¹, Athanasios Tsakiris¹,
Konstantinos Votis¹ and Dimitrios Tzovaras¹

¹ Information Technologies Institute, Centre for Research and Technology Hellas, Thessaloniki,
Greece

{geochan, atriand, arispapapro, ihatz, dgiakoum, atsakir,
kvotis, Dimitrios.Tzovaras}@iti.gr

Abstract. Remote working and collaboration is important towards helping workplaces to become flexible and productive. The significance of remote working has also been highlighted in the COVID-19 pandemic period in which mobility restrictions were enforced. This paper presents the development of an augmented reality platform, aiming to assist workers in remote collaboration and training. The platform consists of two communicating apps intended to be used by a remote supervisor (located e.g., at home) and an on-site worker, and uses intuitive digital annotations that enrich the physical environment of the workplace, thereby facilitating the execution of on-site tasks. The proposed platform was used and evaluated in user trials, demonstrating its usefulness and virtue by assessing its performance, worker satisfaction and task completion time.

Keywords: Augmented Reality, HMD, Remote Collaboration, Spatial Mapping

1 INTRODUCTION

Remote working and collaboration, i.e., the ability to work and collaborate from anywhere-anytime, allows for increased autonomy and flexibility for workers and may enhance their productivity [1]. Limitations in mobility which have been enforced during the recent COVID-19 outbreak, served to underline the importance of remote collaboration digital tools, and have bolstered their application in the workplace [2].

Augmented Reality (AR) is an emerging technology that enhances our perception of the real world by overlaying virtual information on top of it. According to Azuma [3] an AR system must combine real and virtual content, be interactive in real time and be registered in 3D. AR applications are pervasive in our everyday lives and cover various domains such as manufacturing, repairs, maintenance and architecture. The rapid adoption of AR technology can facilitate the development of various AR-based collaboration tools.

Remote collaboration and training can be significantly enhanced through the use of immersive technologies such as AR. According to Regenbrecht et al. [4] collaborative

AR is defined as an AR system where “multiple users share the same augmented environment” locally or remotely and which enables knowledge transfer between different users. AR-enabled collaboration is a relatively young field of research, although the first achievements date back several decades [5]. Nonetheless, the potential of AR for improvement in collaboration efficiency has been reported [6]. Previous research has shown that AR-enhanced remote collaboration can have a major impact in the construction industry [7]. AR-based training platforms allow instructions and annotations to be attached to real world objects without the need of an on-site expert. As the tasks of assembling, operating or maintaining in the construction industry field become more complex the need to reduce costs and training time is essential. The immersion provided by AR-based systems has shown to significantly reduce training time and costs required by employers [8].

The majority of AR platforms developed for workplaces, have been so far application-specific and limited in integrating both remote collaboration and training capabilities. In this direction, we present an AR-based platform aiming to improve collaboration efficiency, productivity, and training. The platform is based on the marker-less augmented reality technology and it can be used on any environment and workplace from any user equipped with a smartphone, a tablet or a Head-Mounted Display (HMD) as the only required equipment. The platform also uniquely takes advantage of augmented reality-enhanced training, through providing the ability to extract keyframe clips with step-by-step instructions, and store them for future reference. The rest of the paper is organized as follows: Section 2 presents related works found in the scientific literature. The proposed system design is described in section 3. Section 4 shows the system development. Preliminary experimental results from user trials are reported in section 5. Concluding remarks are discussed in section 6.

2 AUGMENTED REALITY FOR TRAINING AND COLLABORATION

Piumsomboon et al. [9] report on CoVAR, a remote collaboration Mixed-Reality system that is based on the fusion of Augmented Reality and Augmented Virtuality concepts. A local user’s AR HMD is used to map the environment, which is transmitted and presented to the remote user as a 3D environment. Users may interact through eye gaze, head gaze and hand gestures. The proposed system incorporates several collaboration facilitating features such as 3rd person view, awareness cues and collaborative gaze.

Alem et al. [10] report on HandsOnVideo, an AR-enabled remote collaboration system that is based on the use of natural hand. The remote collaborator uses an overhead fixed camera to capture hand motion and transmit it as video feed to the display of the local collaborator. The local collaborator essentially sees the video feed of the remote collaborator’s hands superimposed over their viewing field and registered with the environment. Thus the remote collaborator is able to guide the local collaborator through hand gestures that are visible in real-time.

Billinghurst et al. [11] propose a face-to-face collaboration system where users manipulate Tangible User Interface (TUI) elements through an AR interface. The elements reported in the paper are materialized in the form of flat markers that are used for identification and tracking. A 3D representation of the corresponding element is superimposed on the tracking marker. Authors present a series of applications of the proposed approach.

Barakonyi et al. [12] report on an AR-augmented videoconferencing system. The system aims to facilitate collaborative control of applications through augmentation of input using marker-based tracking. Application content is superimposed over the markers, and the users may manipulate the markers or place them on their workspace. Thus, the user is able to control an application using marker manipulation, in addition to regular mouse-based input.

Vassigh et al. [13] report on the development and testing of a collaborative learning environment with application to building sciences with the aim of integrating simulation technologies with AR for enhanced decision making in architecture, engineering and construction (AEC). Authors present a system application to the design of an architectural building component, where professionals collaborate through the manipulation of blocks in an AR-enhanced tablet interface.

Weibel et al. [14] propose a platform for multimodal Augmented Reality-based training of maintenance and assembly skills to improve training in the field of maintenance and assembly operations. They report that the skill level of technicians who trained with the developed training platform was higher than the skill level of those who used traditional training methods.

The innovation of our proposed system compared to the related work is its ability to extract keyframes and a summary of the performed steps that can benefit the remote supervisor by reducing her/his cognitive load during a demonstration. It can also operate at any environment to assist any task without the use of markers or a specific hardware setup and it can be deployed both on a mobile device and on an HMD.

3 SYSTEM DESIGN

The platform consists of two applications that communicate with each other, one running on the device of the on-site worker and another running on the remote device used by an expert guide. Fig. 1 shows the schematic overview of the proposed system.

Both apps connect to a backend manager platform. The remote supervisor receives video feed from the AR HMD of the on-site worker, sharing their first-person view of the workspace. The remote expert guide is able to guide the on-site worker by inserting virtual cues and annotations on his workspace view. Annotations become available in the view of the on-site worker.

Using their credentials, users can log in to the local and remote applications. The on-site user can use a mobile device (smartphone) or a Mixed Reality (MR) HMD, such as the Microsoft HoloLens. Through the device's sensors the surrounding environment is scanned. The same application can run on both devices allowing the preferred choice of use depending on the situation and the task at hand.

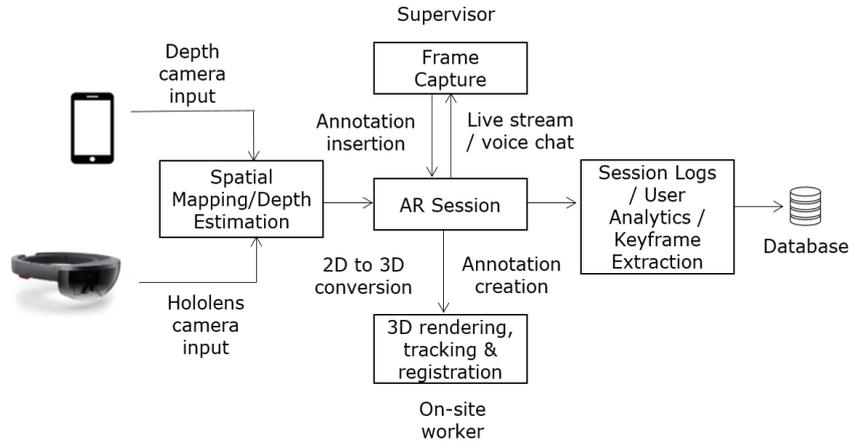


Fig. 1. AR Remote Collaboration tool system architecture

For example an HMD can be more useful for a task that requires both hands to be free while a tablet could be more suitable for the task of simply inspecting a machine and for simple, one-handed interactions. The on-site worker (Hololens) can select an available remote expert guide from a list and call for assistance using hand gestures or touch screen (mobile). The remote expert (mobile) chooses whether to accept or reject the call.

Once the call is setup the remote expert receives live video view from the on-site worker on his mobile device (Fig. 2). Additionally the two users can communicate through real time voice chat. At any time, the remote expert is able to freeze a specific frame from the live view. The expert can zoom and pan on the frame by using pinching and dragging touch gestures in order to focus on a specific part of the worker's view. Subsequently, the expert can insert annotations on the frozen frame, selecting from an array of available symbols (pointing arrows, 3D models), as well as text. Insertion is intuitively performed through touching a point in the viewing field. Through a 2D screen space to 3D world space coordinates transformation, the annotations are sent to the on-site user and rendered as 3D meshes superimposed on his view of the surroundings. Through an on-screen shortcut, the worker and remote expert may clear all annotations with a single interaction.

For each annotation we extract the relevant keyframes and save a clip of the annotation step. During the session the remote expert can access the list of the previous annotation clips and view the corresponding step in real time. Once the call is terminated we generate a session summary containing every annotation step and upload it to a server.

Those sessions can be accessed at any time from any user thus contributing to knowledge sharing and reducing training costs and time. Thus we combine real time collaboration and training with asynchronous AR-enabled step by step instructions.



Fig. 2. Remote supervisor app view with live stream from the Hololens camera (left) On-site worker wearing the Hololens (right)

Furthermore, a log of the whole session containing timestamps for each action performed is uploaded to the server. The saved logs offer valuable insights concerning the time spent on each screen, the total duration of the call and annotations that are used more frequently.

The platform also incorporates a push notification system that informs the remote supervisor about incoming calls. Through the same system and a web based manager back-end platform scheduled calls can be arranged between different users for a specific date and time.

An overview of the remote expert app user interface functionality can be seen in Fig. 3. On the center of the screen the live camera preview of the on-site worker's view is located. On the left panel there are four buttons. The top button is used to capture a frame from the live view. The rest are used for inserting arrow, text and 3D hand annotations respectively. Next to that panel there is the clips preview panel with thumbnails for each recorded action clip. The right panel contains buttons to undo, erase and send annotations or delete the annotations from the on-site user's view. Next to that panel a red button for recording clips is located.

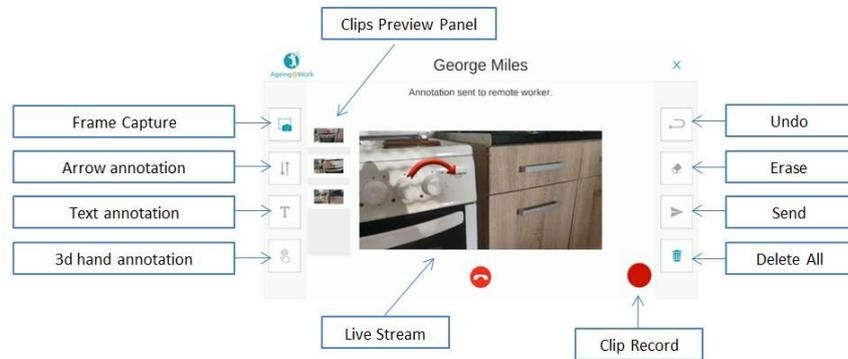


Fig. 3. User interface functionality overview of the guide app

4 IMPLEMENTATION

Both on-site and remote user applications were developed with the Unity3D game engine. Unity3D is ideal for the development of AR applications as it can render 3D meshes on top a device's camera view. The basic component of a Unity application is a scene. The main scene of the on-site app is initially an empty 3D space that consists of a virtual camera that is aligned with the device's physical camera and the camera's live view as a background. The 3D annotations are created as 3D transformations that contain 3D meshes and are continuously tracked as the device moves. The engine offers the ability to build for multiple target platforms so the same core application can be deployed on a mobile device and an HMD. The basic components implemented in the system are presented in Fig. 4 and described below.

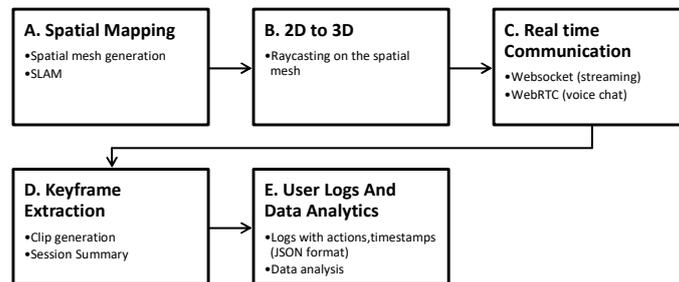


Fig. 4. Impelented modules flow

4.1 Spatial Mapping

The system is based on the markerless augmented reality technology. No previous knowledge of the environment or markers is needed and it can work on any indoor or outdoor space. For the mobile version we implemented the ARCore library that can detect horizontal and vertical planes as well as the ability to reconstruct a spatial mesh of the environment based on the depth camera of the mobile device [15]. Given a 3D point A in the real-world environment and a corresponding 2D point a in the depth image, the value assigned by the Depth API at a is equal to the length of the distance between the camera center C and point A , projected onto the principal axis Z as shown in Fig. 5. This can also be referred as the z -coordinate of A relative to the camera origin C . By assigning a depth value to every point in the camera frame we construct a depth map of the surrounding environment. Through a process called simultaneous localization and mapping, or SLAM, the device understands where it is located relative to the world around it [16]. The app detects and tracks visually distinct feature points in the camera image to compute its orientation and change in location. The visual information is fused with inertial measurements from the device's IMU to estimate the pose (position and orientation) of the camera relative to the world over time. We align the pose of the virtual camera of the 3D scene with the pose of the device's camera in order to render the virtual content from the correct perspective. The rendered virtual image can be overlaid on top of the image obtained from the device's camera, making it appear as if the virtual content is part of the real.

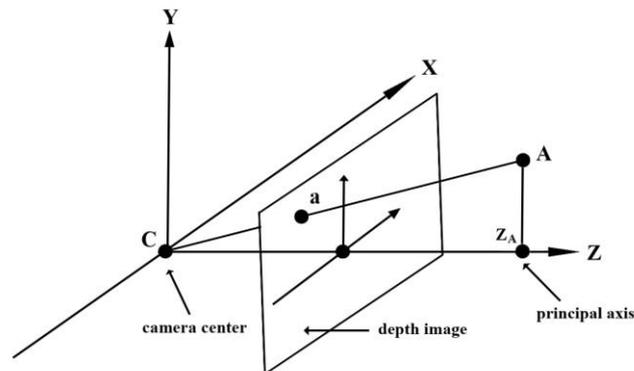


Fig. 5. Projection of a real world point A on the principal axis Z , source: [15]

Similarly on the HoloLens we used Microsoft's spatial mapping API to detect surfaces. The 6 sensors of the HoloLens provide a more detailed spatial map that leads to more precise annotations [17]. The application continuously scans different volumes of the environment in order to receive spatial mapping data. Spatial mapping provides a set of spatial surfaces for the scanned volumes. Those volumes are attached to the HoloLens (they move, but do not rotate, with the HoloLens as it moves through the environment). Each spatial surface is a representation of a real-world surface as a triangle mesh attached to a world-locked spatial coordinate system. During the AR

session new data are continuously gathered from the environment through the sensors and the spatial surface is updated. For each new spatial surface acquired a spatial collider component is calculated that will be later used for the mapping of the 2D coordinates in the 3D space.

4.2 Mapping 2D coordinates in the 3D world

The interaction of the remote expert while inserting annotations is performed on the 2D surface of the tablet. In order to create 3D annotations in the on-site user's view we convert the 2D coordinates of the inserted annotations in the captured frame to 3D world space coordinates. To do so we implement the raycasting method [18]. Raycasting is the process of shooting an invisible ray from a point, in a specified direction to detect whether any colliders lay in the path of the ray. The spatial mesh that is extracted in the spatial mapping process described in the previous section contains a collider component. Each time a frame is captured a virtual camera is stored at the current position and orientation. The virtual camera's projection and world matrix are identical to the device's camera matrices. The ray originates from the virtual camera's position and goes through positions (x,y) pixel coordinates on the screen. The point in 3D world space where the ray intersects with the spatial collider is the origin of the 3d annotation (Fig. 6). Thus the 3D annotations appear in the equivalent positions of their 2D counterparts.

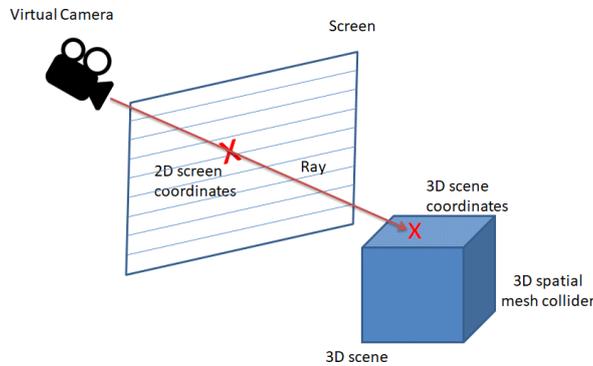


Fig. 6. Ray originating from a virtual camera and intersecting with a 3D collider

4.3 Real time communication

In order to achieve a low latency and fast communication between the two users we implemented the WebSocket protocol [19]. WebSocket is a computer communications protocol, providing full-duplex communication between a client and a remote server over a single TCP connection. Through WebSocket, servers can pass data to a client

without prior client request, allowing for dynamic content updates. We used the web-socket connection to stream frames from the on-site application camera to the supervisor. The video stream resolution and aspect ratio depends on the on-site device's camera resolution so the supervisor application dynamically updates the stream preview frame for each call. Apart from the real time video stream the platform offers real time voice chat to facilitate the communication between the users. The voice chat is based on the WebRTC framework [20]. The WebRTC framework combines two different technologies: media capture devices and peer-to-peer connectivity. Media capture devices includes video cameras and microphones. For cameras and microphones, we used the `mediaDevices` module to capture `MediaStreams`. The `RTCPeerConnection` interface handles the peer-to-peer connectivity. This is the central controller for the connection between two clients in the WebRTC communication.

4.4 Keyframe Extraction

A novel feature of the proposed system is the ability to extract keyframe clips of the generated annotations. A 10 second clip of the on-site user's view can be saved for each created annotation. During the call those clips are available for the supervisor to inspect on a separate panel. Once the call is terminated the clips can be uploaded to a knowledge base and be available for future users to watch for training purposes or to assist them during the performance of a task.

4.5 User Logs And Data Analytics

During each user session the system records a log of the performed actions. Each action is defined by its type and a timestamp. The available actions range from user login and call for assistance to the type of annotation created. Based on those logs valuable information can be extracted such as which users had the most call time or what type of annotations are mostly used. Through acquiring and analyzing such quantitative data collections, we are able to examine the usability and efficiency of the platform, and improve its features.

5 EXPERIMENTAL RESULTS

In order to evaluate our system we performed user trials within a laboratory environment. The basic purpose was to test the usability of the system and the collaborative experience of the users, both remote and on-site.

5.1 Experiment Setup

We chose the use of a 3D printer by an untrained worker as experimental task, because of the value and potential of 3D printers in the modern construction industry. This scenario can be adapted to similar cases in the construction and manufacturing domain as it involves the operation of a machine by a worker, input through buttons

and panels and the use of different devices. The printer used for this task is the Ultimaker 3 3D printer, a desktop 3D printer with a dual extruder (Fig. 7). This printer uses Polyactic Acid (PLA), a thermoplastic polyester, to extrude the plastic on a build platform where it solidifies. The on-site users were equipped with an Android mobile device equipped with a depth sensing camera with a resolution of 2260 x 1080 pixels. The remote supervisors were equipped with an Android tablet and were stationed in a separate room away from the laboratory in which the printer resided. The communication was handled through a high speed Wi-Fi connection.



Fig. 7. The Ultimaker 3 3D printer set up

5.2 Participants and procedure

A total of 8 participants took part in the study, 7 male and 1 female. The mean age of the participants was 25 (± 1.36) years. The participants were asked to perform the task of printing a cube from a usb stick. None of them had used the 3D printer before. The remote guide assisted the on-site user by annotating the usb stick on the table, the usb port on the printer as well as the actions required on the printer input menu (Fig. 8). After the test the users were asked to fill out a questionnaire with questions about the ease of communication during the collaboration session, the usability of the interfaces and the effectiveness of the platform. The questionnaire was based on the System Usability Scale (SUS) [21]. The System Usability Scale (SUS) is an effective tool for assessing the usability of a product such as an application. It consists of a 10-item questionnaire with five response options for respondents; from Strongly agree to Strongly disagree.

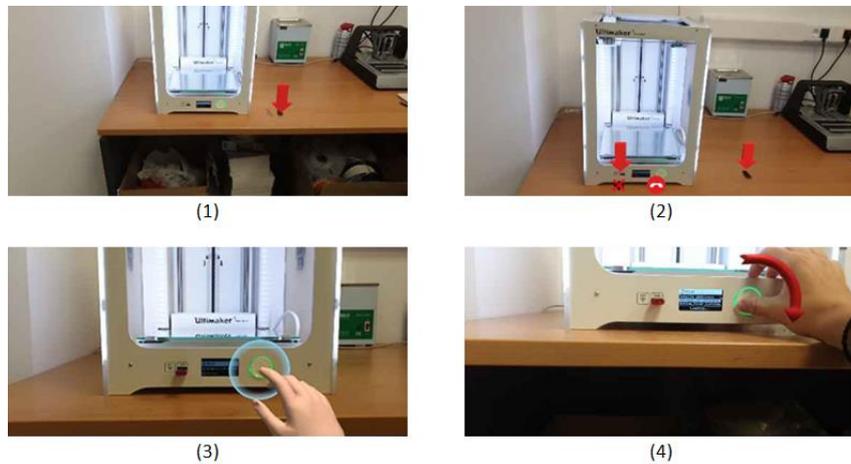


Fig. 8. The basic annotations inserted by the remote guide. (1) Arrow pointing to the usb stick (2) Arrow pointing to the usb port (3) 3D hand animation indicating the button to push (4) rotating arrow indicating the clockwise rotation of the button to select a file

5.3 Results

Every participant was able to complete the task successfully without previous knowledge of using the 3D printer. The average SUS score was 94. According to Bangor et al. [21] any system above 68 can be considered usable. The higher the score the more usable the system is. We can deduce from the score that the AR platform is highly intuitive for the users. The average of positively-worded questions Q1, Q3, Q5, Q7 and Q9 (e.g. questions related to easiness to use and learn) is relatively high while the average of negatively-worded questions Q2, Q4, Q6, Q8 and Q10 (e.g. questions related to complexity and required technical support) is low (Fig. 9). Apart from the questionnaire, quantitative data was collected through the platform's log files. Based on those we were able to measure the mean time required for a user to perform the task compared to the time needed using a traditional manual. The mean execution time was 78 (± 15) secs. The fastest completion time was 60 secs. A separate group of 4 inexperienced users acting as a control group, was asked to perform the same task using only written instructions. Their mean task completion time was 87 (± 26) secs. We notice a 10% improvement in completion times using our proposed system (Fig. 10).

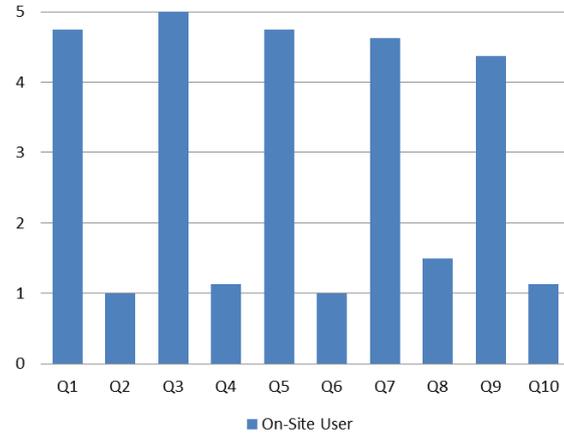


Fig. 9. Average SUS scale rating results from the on-site users (1: strongly disagree, 5: strongly agree) – Q1: I would use this system frequently, Q2: The system is unnecessarily complex, Q3: The system is easy to use, Q4: The support of a technical person is needed, Q5: The functions in this system are well, Q6: Too much inconsistency, Q7: The system is easy to learn, Q8: The system is very cumbersome to use, Q9: I felt very confident using the system, Q10: I needed to learn a lot of things before using

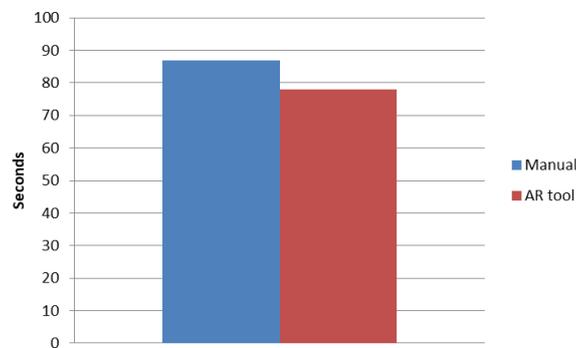


Fig. 10. Mean task completion time using a traditional manual and the AR tool

Regarding the performance of the system, we measured the frame rate of the applications. The on-site app runs at 60 fps with no drops even when multiple annotations are rendered. The average frame rate of the video stream received by the supervisor app is 15 fps over a high speed Wi-Fi connection.

6 CONCLUSION

In this paper we proposed a novel AR remote collaboration platform for workplaces. The platform primarily aims to facilitate the execution of on-site tasks by inexperienced workers, through the guidance by a remote user, using intuitive digital annotations which enrich the physical environment of the on-site worker. In this context, the platform could be particularly significant during the COVID-19 pandemic period, in which several people, including vulnerable groups such as older people or individuals with a chronic disease, were enforced to work from their home .

In addition to real-time collaboration through AR, the platform offers a novel approach to generating intuitive process documentation, through the automatic extraction of keyframe video clips, that enable highlighting the critical parts of the process at hand. To this end, the platform may contribute to the lifelong training paradigm [22].

Preliminary results from our experimental user study, showed that the platform promotes efficiency in task completion, and it is highly usable. As a next step, we aim to conduct longitudinal studies in real-life working environments, to further assess the effectiveness of the proposed system.

Acknowledgment. This work was funded by the European Union’s Horizon 2020 Research and Innovation Programme through Ageing at Work project (under grant agreement no 826299).

References

1. Attaran, M., Attaran, S., & Kirkland, D. (2019). The need for digital workplace: increasing workforce productivity in the information age. *International Journal of Enterprise Information Systems (IJEIS)*, 15(1), 1-23.
2. Waizenegger, L., McKenna, B., Cai, W., & Bendz, T. (2020). An affordance perspective of team collaboration and enforced working from home during COVID-19. *European Journal of Information Systems*, 1-14.
3. Azuma R., A Survey of Augmented Reality, Tele-operators and Virtual Environments, *Presence: Teleoperators and Virtual Environments* 6,4, August 1997, 355-385.
4. Regenbrecht, H. T., Wagner, M., & Baratoff, G. (2002). Magicmeeting: A collaborative tangible augmented reality system. *Virtual Reality*, 6(3), 151–166.
5. Billinghurst, M. and Kato, H., 2002. Collaborative augmented reality. *Communications of the ACM*, 45(7), pp.64-70.
6. Jalo, H., Pirkkalainen, H., Torro, O., Kärkkäinen, H., Puhto, J. and Kankaanpää, T., 2018. How Can Collaborative Augmented Reality Support Operative Work in the Facility Management Industry?. In *KMIS* (pp. 39-49).
7. El Ammari, K., Hammad, A., (2019). ” Remote interactive collaboration in facilities management using BIM-based mixed reality.” *Journal of Automation in Construction*, 107, 102940.

8. Martínez, H., Skournetou, D., Hyppölä, J., Laukkanen, S., & Heikkilä, A. (2014). Drivers and Bottlenecks in the Adoption of Augmented Reality Applications. *Journal of Multimedia Theory and Applications*, 1, 27–44.
9. Piumsomboon T, Day A, Ens B, Lee Y, Lee G, Billingham M. Exploring enhancements for remote mixed reality collaboration. In *SIGGRAPH Asia 2017 Mobile Graphics & Interactive Applications 2017 Nov 27* (pp. 1-5).
10. Alem, L., Tecchia, F. and Huang, W., 2011. HandsOnVideo: Towards a gesture based mobile AR system for remote collaboration. In *Recent trends of mobile collaborative augmented reality systems* (pp. 135-148). Springer, New York, NY.
11. Billingham, M., Kato, H. and Poupyrev, I., 2001, August. Collaboration with tangible augmented reality interfaces. In *HCI international* (Vol. 1, pp. 5-10).
12. Barakonyi, I., Fahmy, T. and Schmalstieg, D., 2004, May. Remote collaboration using augmented reality videoconferencing. In *Graphics Interface* (Vol. 2004, pp. 89-96).
13. Vassigh, S., Newman, W.E., Behzadan, A., Zhu, Y., Chen, S.C. and Graham, S., 2014. Collaborative learning in building sciences enabled by augmented reality. *American Journal of Civil Engineering and Architecture*, 2(2), pp.83-88.
14. Webel, Sabine & Engelke, Timo & Peveri, Matteo & Olbrich, Manuel & Preusche, Carsten. (2011). *Augmented Reality Training for Assembly and Maintenance Skills*. *BIO Web of Conferences*. 1. 10.1051/bioconf/20110100097.
15. <https://developers.google.com/ar/develop/java/depth/overview>
16. C. Cadena et al., "Past, Present, and Future of Simultaneous Localization and Mapping: Toward the Robust-Perception Age," in *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1309-1332, Dec. 2016, doi: 10.1109/TRO.2016.2624754.
17. Evans, G., Miller, J., Pena, M.I., MacAllister, A., & Winer, E. Evaluating the Microsoft HoloLens through an augmented reality assembly application. *Defense + Security* (2017).
18. Roth, S. D.. "Ray casting for modeling solids." *Comput. Graph. Image Process.* 18 (1982): 109-144.
19. V. Wang, F. Salim, and P. Moskovits, *The Definitive Guide to HTML5 Websocket*. New York, NY: Apress, 2013.
20. <http://www.webrtc.org>
21. Bangor, A., P. Kortum, and J. Miller. 2009. "Determining What Individual SUS Scores Mean: Adding an Adjective Rating Scale." *Usability Studies* 4 (3): 114–123. doi:10.5555/2835587.2835589.
22. Lee, M., & Morris, P. (2016). Lifelong learning, income inequality and social mobility in Singapore. *International Journal of Lifelong Education*, 35(3), 286-312