

Interacting with Smart Consumer Cameras: Exploring Gesture, Voice, and AI Control in Video Streaming

Anonymous Author(s)*

ABSTRACT

Livestreaming and video calls have grown in popularity due to the increased connectivity and advancements in mobile devices. Our interactions with these cameras are limited as the cameras are either fixed or manually remote controlled. Here we present a Wizard-of-Oz elicitation study to inform the design of interactions with smart 360° cameras or robotic mobile desk cameras for use in video-conferences and live-streaming situations. There was an overall preference for devices that can minimize distraction as well as preferences for devices that can show they demonstrate an understanding of video-meeting context. We find participants dynamically grow with regards to the complexity of interactions which illustrate the need for deeper event semantics within the Camera AI. Finally, we detail interaction techniques and design insights to inform the next generation of personal video cameras for streaming and collaboration.

ACM Reference Format:

Anonymous Author(s). 2018. Interacting with Smart Consumer Cameras: Exploring Gesture, Voice, and AI Control in Video Streaming. In *TVX '19: ACM International Conference on Interactive Experiences for Television and Online Video, June 05–07, 2019, Manchester, UK*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Personal media streaming has entered a new age [10] with the growth of Internet-enabled cameras and smartphones coupled with various streaming services and video calling applications. While in the past video-streaming was found in enterprise conference calls, it now also enables self-run Internet “TV” shows on YouTube, social media livestreams on Facebook, and personal broadcasts on Twitch. In many of these cases,

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. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

TVX '19, June 05–07, 2019, Manchester, UK

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-9999-9/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

people resort to holding a cameraphone, using a laptop screen mounted webcam, or even having installed setups with green screens, lighting, and picture in picture broadcasts. Consider the solo broadcaster or set of broadcasters trying to queue various cameras. Typically they would rely on a display camera zoomed into a fixed area so they could show a item and then have a second camera to capture themselves or the room. Cameras could be manually triggered or an articulating camera could be remote controlled. Both cases require various levels of interaction from the broadcaster.

Currently, we are faced with two recent advancements. First, cameras now have high resolution captures well beyond 1080p HD. Many consumer cameras shoot video in 4k or 6k and can capture wide or 360° fields of view, allowing one to crop out an HD frame. Second, artificial intelligence (AI) technologies have vastly improved in the areas of object detection, face identification, and voice recognition. Towards this, we ask how would one control such a high-resolution AI-enabled personal camera.

In this article, we present a Wizard-of-Oz elicitation study to understand how people would interact with a smart personal camera that could either (1) pick the person or objects to focus on automatically, (2) could be controlled by visual hand gestures like pointing, or (3) could be voice controlled to focus on a person or object. In the study, we not only surface what kind of visual or audio commands one would use to control such a personal video camera, but also how different cameras’ interaction techniques can lead to the camera becoming a participant in the meeting and can facilitate deeper connections with remote participants. Finally, we illustrate how this research can inform the design of future enterprise or personal video cameras for meetings or livestreaming.

2 RELATED WORK

The broader research work in video conferencing, chat rooms, and livestreams is expansive. Telepresence systems have tried to connect people through extensive camera augmentation and registration [3] as well as simpler side-by-side [19] arrangements. Beyond work contexts, livestreaming on the Internet has grown [10] due to expanding connectivity and network-enabled cameras (and cameraphones). In this article, we focus less on improving image fidelity [21] and/or the mechanics of building an AI system for object recognition. Instead we focus on experiences, people, and behaviors.

107 There has been some growth recently in video cameras for
 108 meetings, mostly in the enterprise sector. New devices like the
 109 Meeting Owl¹ and Sony Xperia Hello!² provide tower-like
 110 single lens 360° cameras for video calling. The Meeting Owl
 111 has a top facing camera and is designed for boardroom use.
 112 The Xperia Hello! has a front skewed single lens camera for
 113 meetings or personal calling. Other devices, like the Amazon
 114 Echo Show and Spot, have a similar camera to that on a laptop
 115 of tablet device. Finally, the Google Snap camera is a device
 116 designed to be placed in vantage points where one might need
 117 photos or video marketed for non-enterprise environments.
 118

119 User Behaviors in Live Streams & Video Conferences

120 One common scenario for live streaming includes a conversation
 121 between individuals or a group of people who are sitting
 122 in a room or fixed location. However, as streaming becomes
 123 easier and more mobile, people may wish to view multiple
 124 streams from different perspectives, such as at a large public
 125 event [4]. Audience members may also interact with stream-
 126 ers in other ways such as through chat or voting mechanisms
 127 to influence what happens [6]. However, enabling presenters
 128 to easily share their perspectives and activities with viewers
 129 continues to be an important area of focus as well.
 130

131 One issue that has emerged in single, or fixed-camera video
 132 streaming settings has to do with how people go beyond being
 133 “talking heads” to showing physical items or highlighting
 134 movements, as in a performance [14]. A study of how people
 135 share physical artifacts (whiteboards, notebooks, mobile
 136 phone screens) with remote meeting participants over video
 137 conferencing revealed that this behavior is difficult and can
 138 be disruptive to the flow of the meeting [11].

139 Other work [20] has looked at how multiple fixed cameras
 140 might be used to automatically switch between different
 141 views of people engaging in a distributed task. Initial ex-
 142 plorations in this space suggest that automatically switching
 143 camera viewpoints could be a scalable solution when people
 144 are distributed across multiple sites; however, such a setup
 145 still requires that rooms be outfitted with several pieces of
 146 hardware that are not easily portable.
 147

Automatic Camera Operators

148 Another area of research is on the automatic camera operator.
 149 This includes proposing a set of rules for an automatic camera-
 150 person [18] to learning where to focus a pan/tilt/zoom (PTZ)
 151 camera based on previous user interactions [8] to automatic
 152 positioning of multiple PTZ cameras to maximize the cap-
 153 tured video information for passive viewers—providing the
 154 best close-up view with limited cameras [9]. Typically, when
 155 multiple cameras are in operation, camera selection is done
 156

¹<https://www.owlabs.com/>

²<http://www.sony.jp/xperia-smart-products/products/G1209/>

160 via algorithm, such as a constraint satisfaction problem [5],
 161 or done via community cooperation [15]. More sophisticated
 162 controls create camera-sensing networks [13] or use a set
 163 of distributed smart cameras and with gesture sensing [7] to
 164 control the queuing.
 165

Interactions

166 Hand gestures have been a common interaction technique for
 167 video and for remote control [2]. This is of particular interest
 168 given the current advancements in AI and the availability of
 169 large scale, crowdsourced datasets for gesture recognition
 170 such as the 20BN-Jester Dataset V1. One question we wish
 171 to investigate is how people would interact through gestures
 172 for a personal desk robot. Bearing the most similarity to our
 173 work is the recent Wizard-of-Oz Human-Drone Interaction [1]
 174 experiment. While hand gestures have been explored before
 175 for remote controlling quadcopter drones [17], Cauchard et
 176 al.’s elicitation study showed how people would control a
 177 personal quadcopter through gesture and voice.
 178

179 Also similar to our work is a recent observational lab study
 180 in which participants critiqued pairs of physical prototypes
 181 (prosthetic hands) for a face-to-face or remote collabora-
 182 tor. [12]. Here, details in physical prototypes were shared
 183 in two different contexts. The authors found that traditional
 184 media sharing experiences are optimized for an upwards
 185 gaze whereas tabletops focus downward; they suggest head-
 186 mounted displays to keep media spaces in peripheral vision
 187 while downward attention is focused on an artifact or object.
 188 They also suggest prioritizing the camera preview window as
 189 theirs was rendered too small in their experiment.
 190

3 METHOD & EXPERIMENT

191 We present a Wizard-of-Oz elicitation study to understand
 192 the kind of interactions people would have with smart video
 193 cameras in a livestream or video-conference context. We used
 194 an office space with a table and a whiteboard and a set of
 195 printed material to simulate a conversation across two lo-
 196 cations where two collocated participants shared materials
 197 and information with a remote attendee using four different
 198 camera setups and interaction modalities: (1) A stationary 4k
 199 360° camera that automatically identifies the salient person or
 200 object and streams it. (2) A small moving camera-robot that
 201 pans on the table to automatically identify the salient person
 202 or object and streams it. (3) A small moving camera-robot
 203 that pans on the table and responds to visual gestures to turn
 204 to the salient person or object and stream it. (4) A small mov-
 205 ing camera-robot that pans on the table and responds to voice
 206 commands to turn to the salient person or object and stream
 207 it. The 4k 360° camera was a Ricoh Theta V and the small
 208 moving camera-robot was an Anki Cozmo. Both cameras are
 209 consumer-level devices retailing for \$400 and \$180 respec-
 210 tively. An experimenter remote controlling the Cozmo over
 211

213 WiFi was in an adjacent office with a confederate who was
214 connected to the participants via a speakerphone to ask ques-
215 tions and comment. This setup mirrors a remote participant
216 on a video call with no camera of their own (thus not drawing
217 eye contact from the participants) which is also analogous to a
218 typical livestream. The difference between the two scenarios
219 is that the confederate can give audio feedback whereas a
220 livestream would rely on text comments as feedback.

221 For each case, minimal instruction was given to the parti-
222 cipants. For each condition this was:

- 223 • Condition 1: “The camera will use AI to automatically
224 find and stream people and objects in the room relevant
225 to the conversation.”
- 226 • Condition 2: Same as Condition 1.
- 227 • Condition 3: “The robot-camera will respond to hand
228 gestures to move.”
- 229 • Condition 4: “The robot-camera will respond to voice
230 commands if you first address it as ‘hey robot’.” This
231 condition mimics modern voice appliances such as
232 Amazon Alexa, Apple Siri.

233 We did not refer to either camera by brand or name to avoid
234 anthropomorphism of the Cozmo robot and while the Cozmo
235 has various expressive face/eye animations, we kept its ex-
236 pression in a neutral repose.

237 Participants

238 For this study, we recruited 10 volunteers (2 females, 8 males)
239 from our institution. All of the participants were college
240 educated and were somewhat familiar with 360° or robot-
241 articulating cameras. Additionally, all of the participants re-
242 ported they generally felt comfortable around cameras in the
243 room. There was also an even distribution of participants who
244 said they were easily distracted or not easily distracted.

245 Task

246 The basic task among each participant pair was to discuss var-
247 ious electronics or accessories to purchase for a company. We
248 separated this into four tasks, one per category type: (a) lap-
249 top, (b) tablet, (c) phone, and (d) backpack/messenger bag.
250 Two items from each task type were selected and printed out
251 as technical specifications and prices from their associated
252 web pages. The whiteboard contained a handwritten (but in-
253 complete) table of regions, employee ranks, and budget prices.
254 Conditions and tasks were matched randomly and counterbal-
255 anced for each pair (four in total). Each participant had one
256 item assigned to them to explain and share with the remote
257 participant. For example, for Condition 3 and task *phone*, one
258 participant received printouts for an iPhone X and the other
259 received printouts for a Google Pixel 2. Both participants
260 were asked to discuss each device as it relates to a price point,

261 region, or employee rank on the white board. The camera-
262 robot was the only broadcasting feed and responded to visual
263 gestures (Condition 3) to move.

264 Procedure

265 An office equipped with a table, chairs, and a whiteboard
266 was outfitted with several printouts. A small (GoPro) camera
267 was mounted to the far corner of the whiteboard, opposite
268 side of the experiment related whiteboard content, to record
269 the experiment and participants were instructed to ignore
270 it. Participants were first given 5 minutes to review all the
271 printouts. Conditions and tasks were randomly assigned then,
272 before each condition, an experimenter brought a camera
273 device into the room, briefly explained its capabilities and
274 discussed the task. The participants were told that an executive
275 from corporate headquarters (the confederate) was connected
276 via speakerphone.

277 The experimenter then left and the confederate prompted
278 the session to begin. The confederate minimally prompted
279 for information or visuals about the whiteboard and printouts
280 if these were not spontaneously provided by the participants.
281 As the Cozmo’s resolution is rather low, the confederate had
282 a copy of all the handouts but did not disclose this to the
283 participants. This allowed the confederate to “play along”
284 and overcome the current limitations with the camera-robot.
285 Figure 2a shows an example of the camera’s fidelity.

286 Upon completion of the condition-task pair, the experi-
287 menter returned, collected the device and handed the partici-
288 pants a short survey to fill out. Once completed, the surveys
289 were collected, a device was returned to the room and the next
290 condition-task pair started. There were four condition-task
291 pairs in total and the full experimental session lasted around
292 30 minutes overall. At the end of the last condition-task pair,
293 a final survey was completed by the participants, and they
294 participated in a short interview/discussion with the two ex-
295 perimenters in order to collect additional comments about
296 their reactions to the different conditions.

297 Surveys and Exit Interview

298 After each condition-task pair, participants were handed a
299 short survey. After the final pair, they were handed an ad-
300 ditional overall survey. Once completed, the experimenters
301 asked the participants if there was any additional feedback,
302 comments, or questions. Each of the condition-task surveys
303 asked questions (on a 5 point Likert scale from strongly dis-
304 agree to strongly agree) related to the camera condition: was
305 the camera useful, did they know what the camera was looking
306 at during the meeting, did they want to control the camera/did
307 they feel they could control the camera. Additionally, there
308 were two open-ended questions asking participants what they
309 liked and disliked about the technology they used in that
310 condition.

319 The final survey gathered demographics (age, gender, education)
 320 experience with robotic or 360° cameras, how comfortable
 321 they are with being recorded, if the camera was distracting,
 322 and a self-reported assessment of whether they were
 323 easily distracted in general. We also asked if they wanted to
 324 be able to toggle the camera's recording state, and we asked
 325 them to rank the four conditions in order of preference. Finally
 326 we asked if there was any other way they would have
 327 liked to control the camera and if there was anything else
 328 aside from people, papers, and whiteboards that they would
 329 want the camera to know about.
 330

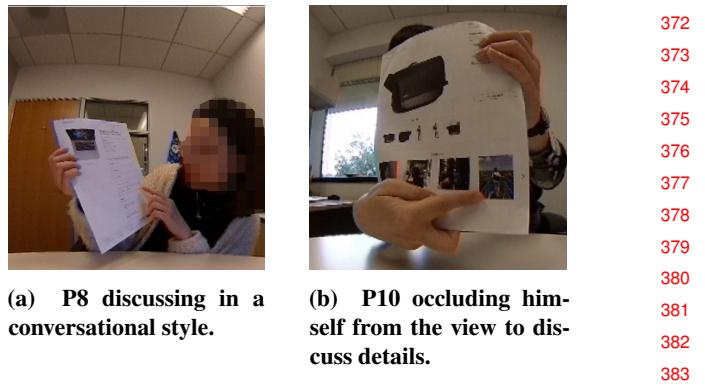
331 4 RESULTS

332 When asked to rate the four conditions in order of preference,
 333 the 360° degree camera was the favorite of 5/10 participants.
 334 The next most favored condition was the AI robot (3/10 participants), followed by voice control (2/10 participants). The
 335 gesture control condition was selected as the second favorite
 336 option by 4 out of 10 participants. To further unpack the variation
 337 in preferences between the different interaction methods,
 338 we now discuss behaviors observed during the experiment as
 339 well as open-ended comments provided by participants for the
 340 four different conditions. These provide some insights into the
 341 ways in which the different technological affordances of the
 342 360° and robot cameras influenced the behaviors and attitudes
 343 of participants with regards to sharing their environments with
 344 the meeting viewer.
 345

346 360°

347 During the 360° session, most participants were fairly comfortable
 348 with the technology even though not all of them had used such a device before. For this condition, 3 participants
 349 selected "Strongly Agree" with the statement that the camera
 350 was useful, 5 selected "Agree" and 2 were neutral. However,
 351 the 360-degree camera was rated significantly lower than the
 352 three robotic conditions in terms of agreement with the question
 353 "I knew what the camera was looking at." Participants
 354 brought up the point that the 360° camera required the least
 355 amount of input or interaction from the two presenters in the
 356 room (because the remote person could see everything and
 357 control what they focused on). The benefit of this setup was
 358 that it gave greater agency to the remote viewer: "The camera
 359 was for you [the remote participant], not for me . . . it's
 360 not my job." (P6) However, while this benefited the remote
 361 participant, the two co-located meeting participants lacked
 362 awareness of the remote party's focus. This was sometimes
 363 challenging because
 364

365 "it's always visible and the remote user can watch
 366 me . . . though I feel being monitored I also wonder
 367 whether the remote users have time and effort
 368 to monitor all the stuff in the 360° camera" (P9)



384 Figure 1: Two participants on the 360° camera condition showing
 385 details on the printout. P8 (1a) discussed the paper in a
 386 conversational manner, showing herself with the document. P10
 387 (1b) totally occluded himself from view and put the document
 388 forward as the sole artifact.

390 In terms of presenting the documents, two participants (P3
 391 and P4) left their documents flat on the table when discussing
 392 them. They had a general assumption that "It sees everything,
 393 I don't have to control it." (P3) despite the fact that there is a
 394 loss of legibility at a sharp angle—a case not considered by
 395 the participants. The rest of the participants would hold up
 396 the document to the camera to aid the remote in seeing the
 397 visuals (See Figure 1).

398 AI Camera Robot

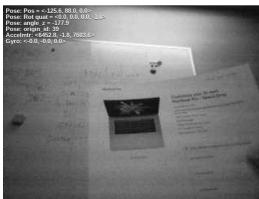
400 The AI Camera Robot began by doing a 360° sweep of the
 401 room. This took roughly 3 seconds and allowed the participants
 402 to feel as if the AI (actually WOZ) scanned the room so it could track the video session automatically. One participant
 403 was neutral to the usefulness of this condition, the rest Agreed
 404 (7) or Strongly Agreed (2) to the usefulness. Here we saw the
 405 AI plus Robot combination making a connection between the
 406 participant and remote confederate.

407 "I felt the meeting person from the remote site listen to me when the robot camera turned towards
 408 me when I was talking." (P4)

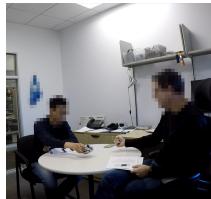
409 Others felt the robot was a participant following along, "It's great that the robot could understand the conversation and get
 410 to know where he's supposed to turn towards." (P7) No other
 411 condition gathered such feedback about making a stronger
 412 connection with the (invisible) remote participant. One participant contrasted the AI robot condition with the 360° camera
 413 condition by saying:

414 "It's fun and having a robot moving feels like
 415 you're interacting with a real user, real people . . . [the]
 416 360° camera is just like a device, which is cold
 417 and doesn't feel very human" (P5).

418
 419
 420
 421
 422
 423



(a) The robot camera has been directed via gesture to look at the whiteboard.



(b) P5, unhappy with the response time of the robot, picks it up and points it at P6.

Figure 2: a A participant (P2) wanted to show a handout which he just put in the camera’s field of vision instead of directing the robot camera back to him. **(b)** P5 and P6 gesture to the robot mostly out of field. Eventually, P5 just picks up the robot and points it to P6.

Gesture Driven Camera Robot

One participant (P8) disagreed that the Gesture-Driven Camera Robot was useful, while the rest agreed or strongly agreed. However, many participants cited there was a bad recognition lag time on the robot or that it was unclear to them what gestures they could actually use to make it work.

“It seems gesture sensor is not so sensitive. And sometimes, feel confused how to make the right gesture to the robot to work well.” (P7)

Also, “(I was) not sure how/what gestures can direct the robot.” (P4).

Of the participants who made some gestures in the robots field of view, they still cited latency issues (likely partially due to the WiFi streaming rate of the robot). However, it was the case that most participants would gesture outside of the field of view of the camera and then slowly drop the gesture into the robots vision. Once the wizard could actually see the gesture from the robot’s camera, it would be enacted. Two participants continually gestured outside of the robot’s vision. Upon giving up, one participant just picked up the robot to turn it manually to the other person (Figure 2b). Four participants physically moved the robot in this condition (P7, P8, P9, and P10). One other participant asked if he was allowed to move the robot, but only did so after the sessions were completed. In this case, he thought “why am I using my hands [to gesture] when I can just do this? (*mimes picking up robot*)” (P8).

However, there were also some positive points made about gesture control being natural “It’s like interacting with real people, natural.” (P5) and not distracting “I don’t have to stop talking to give it instructions, nice.” (P7).

Voice controlled Camera Robot

This condition had unanimous usefulness ratings: 4 with strongly agree and 6 with agree. We mimicked the interaction style of modern day speech interfaces by having the “hey robot” trigger. During the experiment, the robot (wizard) was rather flexible as many participants said “hello robot” or “hi robot.” Many people enjoyed this interaction. “It’s interesting to talk to the robot, it’s like having a secretary.” (P7), “(It’s) easy to change views, flexible in conversation” (P5), and “I didn’t need to learn a new control interface” (P8). Others cited difficulty: “Takes effort to describe orientation direction and how far to rotate.” (P3) and “If within arms reach, easier to manually position or let remote person control it.” (P6).

There were a few different interaction dynamics in this condition. A few participants (P3, P4, P5) gave relative-rotational instructions like *Turn {n}° to the {d}*. Others (P7, P8) would say the degrees of rotation with a logical object (like themselves) instead of clockwise or counterclockwise directive: e.g. “Turn 90° towards me”. Some participants (P9, P10) would just say the trigger “hey robot” as the implicit *Turn to me* command. Finally some would give a logical object like “turn to the whiteboard” (P9) and others, notably P4, stared with hard coordinate directions “Turn 70° to the left” but grew to say “Robot get the whiteboard, the whiteboard” by the end of the 10 minute condition session.

Issues with the Camera Robot

There were some common issues with the Cozmo camera robot. First, many participants mentioned it was noisy (P3, P4, P5, P8). While not overwhelmingly loud, its motor driven tracks and head do make a grinding noise that is uncommon in most video sessions (different than the hum of a CPU or high-pitch sound of some electric devices). This is more so recognized next to the silent 360° camera. In the context of recording, it introduced some difficulties: “Made obvious noise when moving so i had to talk over it.” (P8)

There was some small frustration with not knowing the vocabulary of commands, be them visual or audio, to give the robot (P3, P4, P6, P8) but this was the part of the exercise of this elicitation study. Also a few participants (P6, P8) asked what was the point of having the co-located presenters control the robot, as opposed to letting remote person just have control of the robot in the first place?

5 DISCUSSION

We investigated how people would interact with small smart camera devices and how would they change the experience. Three advancements motivate this work: more specialized devices are entering the consumer market, modern deep AI is driving advancements in gesture and voice interfaces, and finally, the growth of video communication (in teleconferences,

531 video calls, and livestreams). This study illustrates not the
532 commands and manipulations people use with smart cameras.
533 It points to an understanding of the context of use, not just the
534 visual content in the scene. While the 360° camera was most
535 often chosen as the most preferred when ranking the four
536 conditions, most of the conditions were thought to be useful.
537 We did find several social and environmental observations
538 which should be considered for the future design of personal
539 smart cameras.

540 First, we note that most of the participants asked for feed-
541 back from the device. No feedback was given from the Cozmo
542 robot despite it having colored lights, text-to-speech, and can
543 nod to say yes or shake in disagreement. This was to keep it
544 close to the 360° camera which had no feedback capabilities.
545 Further, many participants asked for a preview window of
546 the devices and previous work has shown its value in media
547 spaces [12] and livestream chats [16]. This was not in the
548 study design as we wanted to focus on the interaction with
549 the device and not self positioning in the preview frame.

550 In the cases where the robot had a gesture or voice con-
551 trol, many participants initially paused wondering what to
552 gesture or say. Typically, when one participant would invent
553 a command, the other would repeat a similar command. This
554 dissipated over the course of each session as the participants
555 started inventing newer commands, but was something to
556 note. While most all the participants would generally control
557 the camera for themselves, a few participants would give a
558 command for the other participant. In one case, P5 gestured to
559 move the camera to P6 (Figure 2). In another, P10 generally
560 just moved the camera by hand for P9.

562 Awareness

563 By withdrawing feedback, our experiment highlighted the
564 importance of feedback delivery to the participants, especially
565 with gesture control. Adding audio (tones or voice) output
566 or visual displays showing what the robot is capturing could
567 help alleviate these issues but also run the risk of distracting
568 the video with sound or another screen in the frame. A screen
569 would need to be on the camera-device as not to have someone
570 look one direction to fix the camera in the other direction.

571 When it comes to voice commands, some participants di-
572 rectly addressed the robot by naming objects (e.g. “get the
573 whiteboard”, “back to me”). The robot needs to develop an
574 awareness of its surroundings, remember its surroundings, or
575 slowly build and remember the context to be fully operational
576 as soon as possible. In particular, the AI will need to recog-
577 nize common objects and their location given the context and
578 task (meeting room, cooking show, family video call, etc.).

579 Overall, having a camera that knows where to look was
580 seen as having value if one knows where it is looking. And
581 while the 360° camera’s silent omniscience was preferred, the
582 robot having a focus and attention brought it into the meeting

583 as more of a participant. Conversely, the design of how a
584 360° camera could indicate its focus point merits exploration.
585 In either case, there is a trade-off between awareness and
586 interruption that must be taken into consideration.

587 Camera Position

588 Most participants held papers up in front of the camera, both
589 with the robot and the 360° camera. In some cases, they were
590 holding the paper and still pointed at figures or paragraphs. If
591 the 360° camera had instead been mounted above the table
592 looking down, the users might have felt no need to hold
593 their paper sheets. However, the aforementioned benefit of
594 the robot on the desk makes it feel like more of a tool for
595 engagement. This prompted several participants to pick it up
596 and move it. A very portable and light-weight robot certainly
597 encourages that kind of interaction. The physical design of
598 a desk robot should have several view points. While many
599 media spaces have conversations that are gaze-forward (to
600 another) or gaze-down (to an object) [12], the design of such
601 a desk camera robot should support both and possibly correct
602 perspective. This is something that is lacking in common
603 cameras on laptops and tablets and even more modern devices
604 like the Meeting Owl.

605 Laptops and Other Screens

606 While this experiment used papers and a whiteboard, many
607 meetings and video sessions involve sharing screens from
608 laptops, tablets, and smart phones—which are all part of the
609 context of a video meeting. Indeed, half of our participants
610 (P4, P7, P8, P9, P10) stated they would like the Cozmo to look
611 at a laptop or device screen. This could be treated as any other
612 media object or document where the participants turn their
613 laptops toward the camera, like they did holding paper sheets.
614 However, a small portable focused device could sit to the side
615 of a user at the table’s edge or even be positioned between a
616 user and the device. Here, a device like a Cozmo could track
617 a screen or hand-held device and, as our study elicited, keep
618 users informed on what the camera is capturing. By contrast,
619 placing a 360° camera between a device and a person could
620 lead to some potentially unflattering vantage points for the
621 participants. Even if the device has the AI smarts not to show
622 faces from under ones nose, for example, users still might
623 wonder if that is what it is capturing.

624 6 CONCLUSIONS

625 We describe insights from a Wizard-of-Oz elicitation study to
626 inform the design of AI-powered interactions for smart cam-
627 eras for video calling and livestreaming. While the silent, still
628 360° camera was preferred by the participants, it carried the
629 understanding that it was auto-panning and cropping by itself;
630 with that, participants still wondered where it was looking.
631 Cameras that can “look” in a direction, such as the Cozmo,

provide such an indication of where it is pointed but introduce some distractions. Coupled with an auto-AI engine, participants considered the robot more of a meeting participant than the 360° camera—however this requires the AI modeling the event semantics (in this experiment, a teleconference) to facilitate meaningful interactions. However, most general purpose consumer devices are not well designed for these interactions and we have illustrated design and interaction techniques required for automated assisted cameras for streaming.

REFERENCES

- [1] Jessica R. Cauchard, Jane L. E, Kevin Y. Zhai, and James A. Landay. 2015. Drone & Me: An Exploration into Natural Human-drone Interaction. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp '15)*. ACM, New York, NY, USA, 361–365. <https://doi.org/10.1145/2750858.2805823>
- [2] Niloofer Dezfuli, Mohammadreza Khalilbeigi, Jochen Huber, Florian Müller, and Max Mühlhäuser. 2012. PalmRC: Imaginary Palm-based Remote Control for Eyes-free Television Interaction. In *Proceedings of the 10th European Conference on Interactive TV and Video (EuroITV '12)*. ACM, New York, NY, USA, 27–34. <https://doi.org/10.1145/2325616.2325623>
- [3] Tony Dunnigan, John Doherty, Daniel Avrahami, Jacob Biehl, Patrick Chiu, Chelhwon Kim, Qiong Liu, Henry Tang, and Lynn Wilcox. 2015. Evolution of a Tabletop Telepresence System Through Art and Technology. In *Proceedings of the 23rd ACM International Conference on Multimedia (MM '15)*. ACM, New York, NY, USA, 775–776. <https://doi.org/10.1145/2733373.2807400>
- [4] William A Hamilton, John Tang, Gina Venolia, Kori Inkpen, Jakob Zillner, and Derek Huang. 2016. Rivulet: Exploring participation in live events through multi-stream experiences. In *Proceedings of the ACM International Conference on Interactive Experiences for TV and Online Video*. ACM, 31–42.
- [5] Michael Janzen, Michael Horsch, and Eric Neufeld. 2008. Camera Selection Using SCSPs. In *Proceedings of the 2008 Conference on Future Play: Research, Play, Share (Future Play '08)*. ACM, New York, NY, USA, 252–253. <https://doi.org/10.1145/1496984.1497038>
- [6] Pascal Lessel, Michael Mauderer, Christian Wolff, and Antonio Krüger. 2017. Let's Play My Way: Investigating Audience Influence in User-Generated Gaming Live-Streams. In *Proceedings of the 2017 ACM International Conference on Interactive Experiences for TV and Online Video*. ACM, 51–63.
- [7] Chang Hong Lin, Marilyn Wolf, Xenefon Koutsoukos, Sandeep Neema, and Janos Sztipanovits. 2010. System and Software Architectures of Distributed Smart Cameras. *ACM Trans. Embed. Comput. Syst.* 9, 4, Article 38 (April 2010), 30 pages. <https://doi.org/10.1145/1721695.1721704>
- [8] Qiong Liu and Don Kimber. 2003. Learning automatic video capture from human's camera operations. In *Image Processing, 2003. ICIP 2003. Proceedings. 2003 International Conference on*, Vol. 2. IEEE, II–543.
- [9] Qiong Liu, Xiaojin Shi, Don Kimber, Frank Zhao, and Frank Raab. 2005. An online video composition system. In *Multimedia and Expo, 2005. ICME 2005. IEEE International Conference on*. IEEE, 490–493.
- [10] Danielle Lottridge, Frank Bentley, Matt Wheeler, Jason Lee, Janet Cheung, Katherine Ong, and Cristy Rowley. 2017. Third-wave Livestreaming: Teens' Long Form Selfie. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services (MobileHCI '17)*. ACM, New York, NY, USA, Article 20, 12 pages. <https://doi.org/10.1145/3098279.3098540>
- [11] Jennifer Marlow, Scott Carter, Nathaniel Good, and Jung-Wei Chen. 2016. Beyond Talking Heads: Multimedia Artifact Creation, Use, and Sharing in Distributed Meetings. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing (CSCW '16)*. ACM, New York, NY, USA, 1703–1715. <https://doi.org/10.1145/2818048.2819958>
- [12] Terrance Mok and Lora Oehlberg. 2017. Critiquing Physical Prototypes for a Remote Audience. In *Proceedings of the 2017 Conference on Designing Interactive Systems (DIS '17)*. ACM, New York, NY, USA, 1295–1307. <https://doi.org/10.1145/3064663.3064722>
- [13] Faisal Z. Qureshi. 2010. Collaborative Sensing via Local Negotiations in Ad Hoc Networks of Smart Cameras. In *Proceedings of the Fourth ACM/IEEE International Conference on Distributed Smart Cameras (ICDSC '10)*. ACM, New York, NY, USA, 190–197. <https://doi.org/10.1145/1865987.1866017>
- [14] Stuart Reeves, Christian Greiffenhangen, Martin Flintham, Steve Benford, Matt Adams, Ju Row Farr, and Nicholas Tandavantij. 2015. I'd Hide You: Performing live broadcasting in public. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, 2573–2582.
- [15] Marco Sá, David A. Shamma, and Elizabeth F. Churchill. 2014. Live Mobile Collaboration for Video Production: Design, Guidelines, and Requirements. *Personal Ubiquitous Comput.* 18, 3 (March 2014), 693–707. <https://doi.org/10.1007/s00779-013-0700-0>
- [16] David A. Shamma, Elizabeth F. Churchill, Nikhil Bobb, and Matt Fukuda. 2009. Spinning Online: A Case Study of Internet Broadcasting by DJs. In *Proceedings of the Fourth International Conference on Communities and Technologies (C&T '09)*. ACM, New York, NY, USA, 175–184. <https://doi.org/10.1145/1556460.1556486>
- [17] Adrian Stoica, Federico Salvioli, and Caitlin Flowers. 2014. Remote Control of Quadrotor Teams, Using Hand Gestures. In *Proceedings of the 2014 ACM/IEEE International Conference on Human-robot Interaction (HRI '14)*. ACM, New York, NY, USA, 296–297. <https://doi.org/10.1145/2559636.2559853>
- [18] Hugo Strubbe and Mi Suen Lee. 2001. UI for a Videoconference Camera. In *CHI '01 Extended Abstracts on Human Factors in Computing Systems (CHI EA '01)*. ACM, New York, NY, USA, 333–334. <https://doi.org/10.1145/634067.634264>
- [19] Paul Tanner and Varnali Shah. 2010. Improving Remote Collaboration Through Side-by-side Telepresence. In *CHI '10 Extended Abstracts on Human Factors in Computing Systems (CHI EA '10)*. ACM, New York, NY, USA, 3493–3498. <https://doi.org/10.1145/1753846.1754007>
- [20] Marian F Ursu, Manolis Falakakis, Martin Groen, Rene Kaiser, and Michael Frantzis. 2015. Experimental enquiry into automatically orchestrated live video communication in social settings. In *Proceedings of the ACM International Conference on Interactive Experiences for TV and Online Video*. ACM, 63–72.
- [21] Bolun Wang, Xinyi Zhang, Gang Wang, Haitao Zheng, and Ben Y. Zhao. 2016. Anatomy of a Personalized Livestreaming System. In *Proceedings of the 2016 Internet Measurement Conference (IMC '16)*. ACM, New York, NY, USA, 485–498. <https://doi.org/10.1145/2987443.2987453>

637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742