

ProjecTA: A Semi-Humanoid Robotic Teaching Assistant with In-Situ Projection for Guided Tours

To appear at ACM CHI '26.

HANQING ZHOU*, School of Design, SUSTech, China

YICHUAN ZHANG*, School of Design, SUSTech, China

ZIHAN ZHANG, School of Design, SUSTech, China

WEI ZHANG, School of Psychology, Shenzhen University, China

CHAO WANG, Honda Research Institute Europe, Germany

PENGCHENG AN[†], School of Design, SUSTech, China

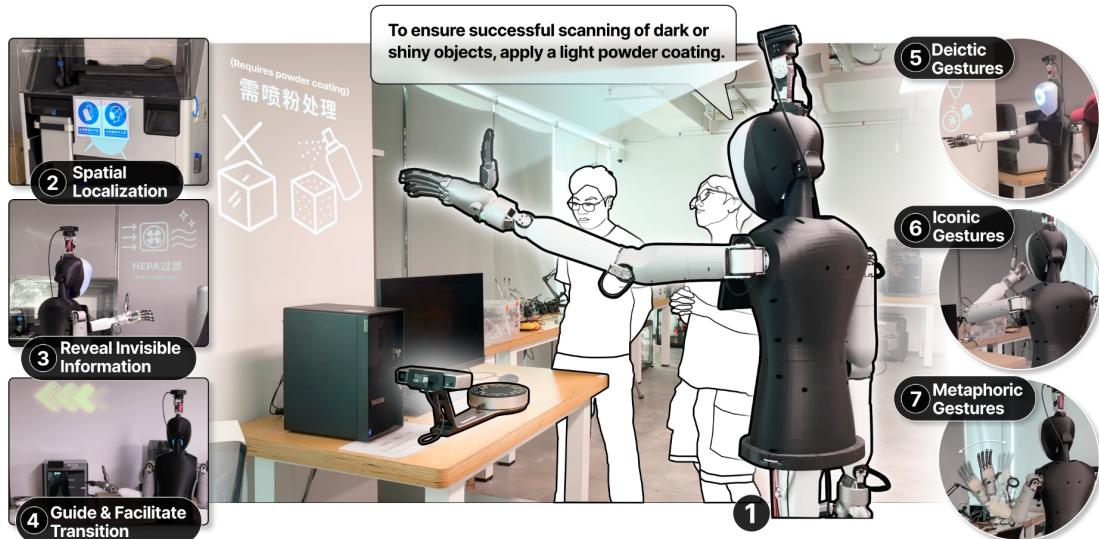


Fig. 1. ① System overview of ProjecTA, offering guided tours in a makerspace. Through in-situ projection, ProjecTA ② provides clear spatial localization, ③ visualizes hidden equipment information, and ④ facilitates transitions. It employs gestures, including ⑤ deictic gestures pointing to concrete objects or locations, ⑥ iconic gestures depicting shapes or operations, and ⑦ metaphoric gestures conveying abstract ideas.

*Both authors contributed equally to this research.

[†]Corresponding Author.

Authors' Contact Information: [Hanqing Zhou](mailto:Hanqing.Zhou@sustech.edu.cn), School of Design, SUSTech, Shenzhen, China, 12331483@mail.sustech.edu.cn; [Yichuan Zhang](mailto:Yichuan.Zhang@sustech.edu.cn), School of Design, SUSTech, Shenzhen, China, 12531639@mail.sustech.edu.cn; [Zihan Zhang](mailto:Zihan.Zhang@sustech.edu.cn), School of Design, SUSTech, Shenzhen, China, zhangzihan654@gmail.com; [Wei Zhang](mailto:Wei.Zhang@szu.edu.cn), School of Psychology, Shenzhen University, Shenzhen, China, zhangwei633@szu.edu.cn; [Chao Wang](mailto:Chao.Wang@honda-ri.de), Honda Research Institute Europe, Offenbach/Main, Germany, chao.wang@honda-ri.de; [Pengcheng An](mailto:Pengcheng.An@sustech.edu.cn), School of Design, SUSTech, Shenzhen, China, anpc@sustech.edu.cn.

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.

Robotic teaching assistants (TAs) often use body-mounted screens to deliver content. In nomadic, walk-and-talk learning, such as tours in makerspaces, these screens can distract learners from real-world objects, increasing extraneous cognitive load. HCI research lacks empirical comparisons of potential alternatives, such as robots with in-situ projection versus screen-based counterparts; little knowledge has been derived for designing such alternatives. We introduce ProjecTA, a semi-humanoid, gesture-capable TA that guides learners while projecting near-object overlays coordinated with speech and gestures. In a mixed-method study (N=24) in a university makerspace, ProjecTA significantly reduced extraneous load and outperformed its screen-based counterpart in perceived usability, usefulness of visual display, and cross-modal complementarity. Qualitative analyses revealed how ProjecTA's coordinated projections, gestures and speech anchored explanations in place and time, enhancing understanding in ways a screen could not. We derive key design implications for future robotic TAs leveraging spatial projection to support mobile learning in physical environments.

CCS Concepts: • Human-centered computing → Empirical studies in HCI; Mixed / augmented reality; • Computer systems organization → External interfaces for robotics.

Additional Key Words and Phrases: Robotic Teaching Assistant, In-situ projection, Guided tours, Physical Learning Space

ACM Reference Format:

Hanqing Zhou, Yichuan Zhang, Zihan Zhang, Wei Zhang, Chao Wang, and Pengcheng An. 2026. ProjecTA: A Semi-Humanoid Robotic Teaching Assistant with In-Situ Projection for Guided Tours. To appear at ACM CHI '26.. In *Proceedings of CHI conference on Human Factors in Computing Systems (Conference CHI '26)*. ACM, New York, NY, USA, 35 pages. <https://doi.org/XXXXXX.XXXXXXX>

1 Introduction

Across formal and informal learning settings, robots serving as teaching assistants (TAs) can meaningfully offload repetitive and standardized tasks [11]. For instance, they can handle routine classroom logistics such as roll call, task reminders, and repeatedly deliver foundational explanations for novices [10]. They also provide step-by-step prompts [113] and safety checks for hands-on procedures [44] in laboratory-like environments. And beyond classrooms and labs, they show potential in supplementing background facts or maintaining tour pace in exhibitions or museums [53, 85].

Many deployed Robotic TAs have upper-torso or fully humanoid forms, ranging from small, desktop units (e.g., NAO [9, 82]) to larger, mobile robots capable of gesturing (e.g., Pepper [103, 113]). A key motivation for upper-torso humanoid platforms is their ability to leverage human-like nonverbal signals, especially gestures, which support learner's comprehension and improve instructional task performance [47, 72]. For instance, pairing robot speech with co-speech gestures or on-screen cues improves word learning in child L2 tutoring [31]; and robot pointing helps learners locate targets faster and reduce misunderstandings [92].

To pair speech with visuals, current robotic TAs typically rely on chest- or head-mounted displays during face-to-face interaction. For example, robotic TAs utilizing the Pepper platform can present instructional content, such as supplementary visuals and structured learning points, on its chest-mounted tablet [97, 103, 113]. Such a screen-based paradigm is especially beneficial for close-range, face-to-face, and robot-centric engagement where the robot and its on-board screen are the learners' primary visual focus during the moment of interaction [9, 82, 103, 113].

However, robotic TAs are also needed in nomadic contexts, where learners and robots, instead of interacting face-to-face, move through space and jointly engage with external targets distributed across different spots. We refer to this arrangement as *nomadic learning*, which is commonly seen in settings such as museums [17, 38, 51], instructional labs [99], or makerspaces. Here, a robot's screen can distract learners' focus from the real-world object of interest, causing

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.
Manuscript submitted to ACM

attention switching between the display and the physical referent. According to Cognitive Load Theory, such mental integration of separate information sources can increase extraneous cognitive load for learners [22].

To support learners in such contexts, the robot, instead of being the center of attention, should be a coordinator, directing focus toward environmental targets and anchoring information around them [88]. Prior work shows that in-situ augmentative displays, enabled via headsets or stationary projectors, can reduce attention switching and lower cognitive load by co-locating information with its physical referent [12, 35, 88]. Recognizing this potential, a few design cases combined (non-humanoid) robotic chassis with on-board projection [53, 85], suggesting benefits of supplementing speech with projected content, while leaving the gesture-projection integration an unaddressed opportunity.

To date, HCI research lacks empirical comparisons to examine whether robotic TAs with in-situ projection would outperform their screen-based counterparts in nomadic learning, how they may shape learners' experiences, or what design implications follow.

Addressing these opportunities, we present *ProjecTA*, a mobile robotic TA offering novice learners on-boarding tours in a makerspace. ProjecTA is equipped with a head-mounted projector, an arm-and-hand actuation module, and a mobile chassis (see Figure 5). As shown in Figure 1 ①, it places near-object overlays in coordination with voice narration and gestures. ProjecTA prepares its tours via a presentation-choreography workflow, which translates teaching goals into executable scripts orchestrating segmented speech, visual display assets, and callable gesture units on a unified timeline. We evaluated ProjecTA via a mixed-methods, within-subject comparative study in a university makerspace, in which 24 participants completed guided tours using ProjecTA and its functionally equivalent counterpart with a chest-screen display, called *Baseline*. Our research question (RQ) is: ***How does a robotic TA with in-situ projection, compared with a screen-based counterpart, affect learners' experiences during makerspace tours?***

Rich empirical data yielded from our mixed-methods evaluation show that quantitatively, ProjecTA significantly lowered learners' extraneous load and resulted in higher perceived usability, usefulness of visual display, and cross-modal complementarity in comparison with Baseline. The qualitative findings further contextualize how the near-object visual overlays projected by ProjecTA reduced referent matching and attention switching in the space, offered a spatially accessible display to nomadic learning, and visualized equipment's critical or hidden details on-site. Concrete examples reveal how the robot's projection, gestures, and speech supplemented one another: for instance, projected visuals disambiguated pointing gestures, and iconic gestures reinforced visual and verbal messages. Relevant design implications are thereby derived to inform future practice.

This work thereby contributes: (1) ProjecTA, a gesture-capable robotic TA with in-situ projection for supporting learners in guided tours; (2) an empirical comparison between such a system and a screen-based counterpart, revealing its positive impacts, and offering rich qualitative accounts of the learners' experiences; (3) a set of design implications for future robotic TAs with in-situ projection to facilitate nomadic learning in physical settings.

2 RELATED WORK

2.1 Robotic Teaching Assistants in Physical Learning Environments

In HCI and educational technology, an increasing body of work has examined robots functioning as teaching assistants in real-world learning environments such as classrooms [9, 10, 31, 113] or supporting place attachment [49]. Prior studies have shown the feasibility of introducing robots in classroom environments, and proven benefits for both educators and learners [87]. Literature views robotic TAs as assistants rather than replacements: they offload repetitive,

standardized routines (e.g., offering roll call, reminders, or repeated knowledge explanations for novices), hence allowing educators to focus on integrative and empathetic work [11, 74].

Many deployed robotic TAs take a fully or upper-torso humanoid form, typically combined with an on-robot display pairing the robot speech with visuals, which is suited for close-range, face-to-face interactions. For example, Pepper pairs speech with a chest-mounted tablet to present instructional content and menu-style prompts [103, 113]; Furhat uses a head-mounted, back-projected face that supports gaze and lip-sync while placing short textual or graphic prompts near the head [2]. QTrobot is a small screen-face cartoon tabletop humanoid used as a peer to nudge children’s handwriting posture [109]. More broadly, desktop humanoids such as NAO are typically positioned within arm’s length and synchronized with a nearby screen for stepwise tasks or quizzes, emphasizing near-field scaffolding [9, 82]. These platforms are thus optimized for deskside tutoring and proximal face-to-face interactions.

However, this on-body screen design presents a key challenge: it forces learners to look at the robot instead of the physical object being discussed. Cognitive Load Theory (CLT) [22] predicts a split-attention effect when verbal explanations and related visuals are spatially separated, increasing extraneous load; co-locating cues on or next to the referent mitigates this cost [22]. While empirically effective for deskside Q&A [89, 90], the on-body screen solution can degrade during multi-object walkthroughs, where learners must shift their gaze between the artifact and the on-body screen [64]. This motivates moving visuals off the robot and onto the target in object-centered, nomadic teaching.

We use ‘walk-and-talk’ to describe instruction delivered while moving through a space and stopping at relevant items. In museums, early guides such as RHINO navigated galleries and stopped at exhibits for explanations, guiding visitors re-orienting to new artifacts along the route [17]. A more recent museum robot proactively approaches visitors for in-situ explanations [38, 51]. In instructional labs, mobile humanoids like Pepper traverse benches and stations, with the group re-forming around different pieces of equipment [99]. These learning settings are thereby object-centered and spatially distributed.

In such settings, a robot is more effective as a mobile coordinator than a screen: it should direct attention and provide information directly in the environment. This can be achieved in two ways: first, by using deictic gestures such as pointing, to establish joint attention on a specific object [92]; second, by presenting information as in-situ visuals, such as explanatory overlays projected onto the object itself. This approach can mitigate the split-attention effect, aligning with both CLT and HRI findings on the benefits of in-place visualization [22, 35, 88]. However, most empirical evidence comes from stationary operations such as assembly and repair tasks in fixed locations, where evaluation mainly focused on task performance [4, 91, 104], rather than on learners’ extraneous cognitive load. A few in-vehicle navigation studies assessed attentional distraction and showed that head-up displays (HUDs) can reduce glances and divided attention compared with a conventional dashboard screen [55], but these settings still involve seated drivers rather than learners moving through space. Thus, HCI still lacks empirical evidence on whether and how in-situ projection, in comparison with conventional screen-based displays, can influence learners’ extraneous cognitive load in nomadic learning with mobile robots.

As instruction shifts from deskside tutoring to walk-and-talk tours, challenges of screen-bound humanoids become clear: first, physical gestures are often not precise enough to identify distant or out-of-view objects, leading to confusion about what is being referenced [50, 92]. Pairing speech with on-body display separates auditory and visual information from the target object, increasing extraneous cognitive load for learners [22]. Third, body-mounted screens are difficult for learners to see from different angles while moving, hindering shared learning experience [64]. These constraints motivate moving visuals off the robot and into the shared environment via public, in-situ augmentation.

2.2 Public Augmentative Displays in Learning Settings

Public augmentative displays present information in the environment, so co-located learners can see and reference it together [3], using it as an anchor for discussion, coordination, and sense-making [8, 81]. Pedagogically, this approach supports joint attention and efficient reference resolution [25, 105]; and with information co-located with its referent, learners avoid mentally integrating separated sources, reducing extraneous cognitive load [22].

In classrooms, such shared interfaces typically appear as wall displays [40, 98, 107], fixed projection [19, 73], and interactive whiteboards [45, 71]: for instance, visually encoding students' progress on a wall display to build shared awareness [107], and projecting learners' ongoing web searches to facilitate timely feedback and whole-class discussion [73]. Beyond classrooms, museum exhibits like DeepTree use large, shared surfaces for collaborative exploration [16]; and an industrial training system projects step-by-step instructions onto a workpiece to guide novices [18]. Across these settings, such public augmentative displays speed learners' shared access to key information and help them coordinate.

Since public augmentative displays are shared by default, any nearby learner can see, point to, and reference without extra individual equipment (e.g., wearables or handhelds). For instance, AAR system uses an actuated projector to help bystanders see the same AR content as HMD users [43]. Radu's meta-review on educational AR argues for the scaffolding effects of shared, in-room overlays in co-located learning settings [81]. By contrast, wearable or handheld AR systems provides individualized, private overlays by default, requiring mirroring or streaming when needed for public access [67]. Such private displays benefit personalized guidance or asymmetrical collaboration [39, 106], such as Lumilo glasses which inform teachers for tailoring instructions to students' real-time needs [46]. However, for maintaining learners' pace and delivering shared content in co-located settings, co-attendant public displays are more convenient, without requiring learners wearing and holding extra devices [43].

However, most public displays are fixed installations (wall-mounted screens, ceiling projectors, etc.). Their utility can be limited in nomadic and spatially distributed learning. In settings like a museum tour or a lab walkthrough, where the focus of instruction moves from one object to another, a fixed display cannot follow the learners. This motivates our exploration of a nomadic, in-situ approach and how it can better support object-centered, walk-and-talk learning.

2.3 Combining In-situ Projection with Robotic Systems

In-situ projection is proven to support fixed-location tasks via information overlays onto the area of operation. For instance, ImproVisAR overlays a 'piano roll' onto a keyboard as intuitive guidance [30], while LuminAR turns a desk into an interactive surface [61]. Evidences from assembly-style tasks show that in-situ projection lowers workload and improves performance relative to paper/tablet/HMD baselines (Pick-by-Projection [5]; LEGO setups [34]). While most of these evidences come from fixed workstations or tabletops, empirical comparisons in educational contexts, especially nomadic, walk-and-talk learning, remain scarce. Leppink's Cognitive Load Scale (CLS) [58] distinguishes learners' Extraneous load (avoidable effort due to presentation) from Intrinsic Load (inherent complexity of learning topics). Since in-situ projection co-locates cues with their referents, it mitigates split-attention: a primary source of extraneous load [23]. Our study thereby uses CLS [58] to test whether in-situ projection, integrated in nomadic settings, would lower learners' extraneous load, extending benefits reported at fixed workstations [5, 34].

In robotics, external, fixed projectors have been used to augment robots and their workspace. For instance, ShapeBots uses a ceiling-projector over a tabletop micro-robot swarm for data physicalization [101]. RobotIST overlays IDE-style feedback (errors/next step/state) on the task surface beside a desktop robotic manipulator, easing procedures more than conventional tools [95].

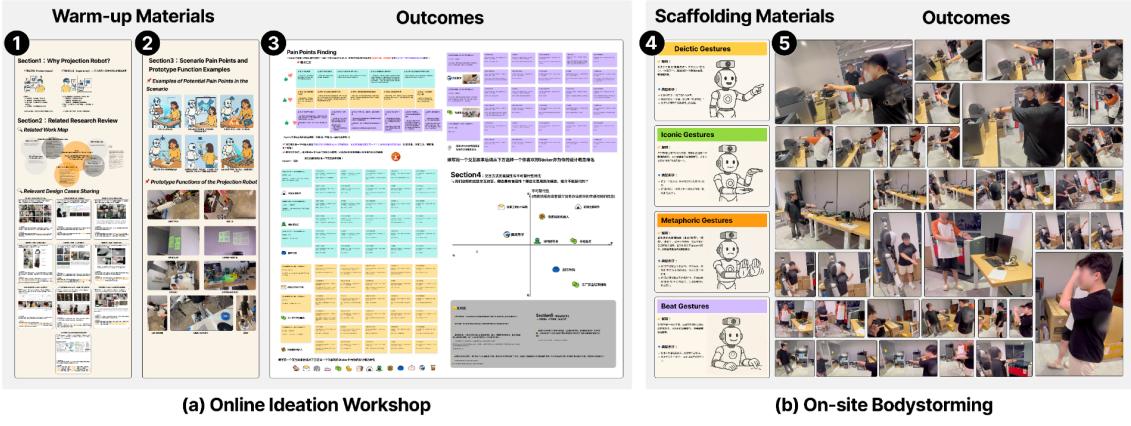


Fig. 2. Formative studies: (a) warm-up materials and outcomes from the online ideation workshop; (b) scaffolding materials and outcomes from the on-site bodystorming.

Beyond fixed setups, on-board projectors have been mounted to actuated gimbals, mobile chassis, or drones to make the display angle-adjustable or fully mobile [100]. PAMI, a stationary meeting-room installation, orients its projected overlays to nearby surfaces for hybrid collaboration [86]. A few systems attach projectors to mobile chassis to deliver near-object overlays in museums [85] and poster exhibitions [53]. Another system navigates indoors and projects an on-demand interactive surface on a chosen spot around exhibits [32]. In children’s play settings, a wheeled pro-cam robot roams the room and projects scene backdrops and prompt cues on the floor/props to guide dramatized activities [1]. Many mobile robots also project arrows, paths, or footprints onto nearby surfaces to convey motion intent and planned trajectory [20, 28, 41, 110]. Related work on mobile projection mapping keeps overlays geometrically aligned while moving [54], and drone-based projection enables pop-up displays where installation is impractical [29, 62, 93]. While these systems effectively make visuals mobile and situated, projection is typically utilized as a single output channel. One exception, Visiobo, times projected visual overlays with LLM-based auditory narration [53], showing the promise of multimodal coordination. Yet prior systems mostly confine co-reference to speech and visuals, with almost no exploration of jointly coordinating gestures, in-situ projection, and narration. In particular, the integration of in-situ projection with gesture-capable robots, such as upper-torso or humanoid platforms, remains unaddressed.

Apart from projection-specific systems, a large body of work have studied multimodal coordination in social robotics. GenComUI synchronizes generated visual aids (map annotations, route cues, animated feedback) with spoken instructions, outperforming speech-only baselines [37]. ELEGNT pairs expressive robot motion with directional lighting/visual signals, making intent more legible and boosting trust [48]. Zhang et al. use gaze and environmental context to ground indirect requests, improving team coordination [115]. Leusmann et al. coordinate questioning with exploratory actions to sustain engagement [60]. Most prior work assumes close-range, face-to-face interaction with information centered on the robot, leaving open how to design the communication channels when the robot and people jointly engage an external object in non-face-to-face, walk-and-talk learning.

Building on the potential benefits of in-situ projection [23, 34], while retaining the instructional value of gestures [26, 27], we explore a gesture-capable robot with an on-board projector that coordinates in-situ projection, arm-and-hand gestures, and speech to support learners in nomadic, walk-and-talk physical learning environments.

3 Formative Study

To explore the value and design opportunities of pairing a teaching assistant (TA) robot with in-situ projection to facilitate physical learning environments, we conducted two co-design activities in our formative study: (1) an online ideation workshop with experts in Robotics and HCI to broadly identify design opportunities; and (2) a bodystorming session in a makerspace to ground our design and collect contextualized requirements.

3.1 Online Ideation Workshop with Robotics and HCI Experts

3.1.1 Participants and Procedure. We recruited six experts (E1–E6). E1 was the Chief Scientist at a Robotics research institute in Europe, and E4 was a university professor in Robotics in Asia, while the remaining experts held master's degrees in Robotics (E3, E5, E6) or HCI (E2). Participants were split into two groups. Each group used Figma¹ to collaborate while having an online meeting facilitated by the researchers, lasting around 2.5 hours.

Each group was offered a Figma template aligned with the co-design agenda to scaffold their collaboration (see Figure 2 (a)). At the start, the experts were briefed about the core objective of the session, i.e., to identify design opportunities for robotic TAs with in-situ projection. As warm-up materials (see Figure 2 ①), relevant design cases from HCI and Robotics on robots presenting information in physical environments were briefly summarized for the experts. These cases were meant to help them understand the background and spark innovations that both build upon and diverge from previous solutions. To ground the discussion, we also shared short demo videos showcasing our robot so experts could understand the functional capabilities of the hardware platform (see Figure 2 ②).

In the next phase, each expert was asked to broadly ideate As-Is scenarios where robots with in-situ projection might be helpful (20 min). Building on these As-Is scenarios, they then created solution-oriented To-Be scenarios (30 min). To facilitate experts' concrete To-Be solutions [56], a structured format prompted entries on how the robot should coordinate in-situ projection, gesture, and voice, and how it should interact with actors and objects in the space.

In a follow-up discussion (30 min), the experts were asked to heuristically evaluate all the created ideas on two aspects: the unique advantages of in-situ projection within each scenario and how these advantages might generalize to other scenarios. Finally, the workshop concluded with a 10-minute collective review, the co-design outcomes were shown in Figure 2 ③.

3.1.2 Major Insights. The experts contributed thirteen detailed To-Be scenarios with rich, vivid commentaries. We analyzed these outcomes using an affinity diagram [65]. The resulting insights, serving as our early design inputs, helped us identify a set of key Design Opportunities (**DO1–DO4** below) that in-situ projection may open up for robotic TAs, as well as typical design scenarios to begin with. We summarize the major insights below:

DO1- Pinpointing key objects and locations through spatial cues: In several To-Be scenarios for in-situ learning or operational tasks, the experts noted that robots' pointing sometimes may be coarse or vague when needing to refer to regions or objects through the space. They envisaged using in-situ projection to anchor clearer references in the scene. **Examples:** E5 envisioned a projection robot highlighting key regions of posters or exhibits, compensating a physical pointing. E3 envisaged a robotic TA overlaying projected positional cues on a workpiece to indicate where to apply the power drill, in a furniture-repairing scenario.

DO2- Presenting learning content alongside its physical referent: The To-Be scenarios related to training or touring highlighted that in-situ projection could place learning materials beside the objects being described simultaneously by a robotic TA. This could avoid learners constantly looking back and forth between an object and its

¹<https://www.figma.com/>

explanation. **Examples:** E5 envisaged that in nursing training, a robotic TA could project task briefings, key steps, etc., next to the mannequin, allowing trainees to glance-and-act. Similarly, E6 envisaged a robotic TA projecting co-located explanations for learners in a makerspace.

DO3- Projecting visuals to complement or summarize robotic TAs' speech: Experts' To-Be scenarios also envisioned that in-situ projection could visualize concise summaries or add-on materials that complement robots' verbal outputs on the spot. **Examples:** E4 envisaged a museum tour, where a robot projected background information and auxiliary materials at places along the way, deepening visitors' understanding. E2 described a similar example: a robot using projection to consolidate its verbalization in a poster exhibition.

DO4- Visualizing hidden or invisible information on-site: A few To-Be scenarios imagined how in-situ projection could help robotic TAs make hidden or occluded content visible to users. **Examples:** In E1's example of organizing items in a domestic setting, a robot projected information onto a nearby wall to remind the user of misplaced items and items stored out of sight, which would otherwise remain invisible to the user.

The To-Be scenarios created by the experts consistently highlighted the most distinctive potential of a robotic TA with in-situ projection: it affords unobtrusive, ad-hoc, on-site information aid for users to both learn about and learn with spatial and physical objects in a particular environment (e.g., museum, makerspace, workshop, exhibition...). We chose one of such typical scenarios: a robotic TA guiding beginners through an introductory tour of a makerspace, which the experts identified as a particularly rich and representative setting even though museum-like environments (e.g., museums[85] and exhibitions[53, 102]) have received more prior attention in related work. In line with the experts' feedback, a makerspace contains various pieces of equipment distributed spatially, each requiring specific instructional materials, suggesting sufficient space to explore all four design opportunities summarized above. For onboarding in a makerspace, each beginner often needs a similar introductory tour, echoing a typical duty of robotic TAs: offloading repetitive, standardized tasks from educators. For these reasons, the beginner onboarding tour in a makerspace has been chosen as our design scenario to proceed with.

3.2 On-site Bodystorming with Makerspace Experts

To collect contextualized requirements and practical inputs for designing a robotic TAs' speech, gestures, and in-situ projection, we conducted another co-design session in a university makerspace(see Figure 2 (b)), adopting a bodystorming method, which has been proven to offer on-the-ground design ideas with situational factors by having participants enact and role-play in the context [78].

3.2.1 Participants and Procedure. We recruited a makerspace educator and manager (M1) and two experienced makers (M2, M3), each with more than 6 years of active engagement in makerspace environments. Under researcher facilitation, they completed a 35-minute in-situ bodystorming session, where they took turns to role-play a robotic TA offering novices an introductory tour, acting out its speech, gestures, and visual presentation using both bodily performance and verbal description. Such enactment was meant to extract the experts' embodied know-how of how to introduce the makerspace, and gather references, examples, and inspirations for creating the robot's behaviors in different communicative channels.

Scaffolding materials were prepared to further support experts' enactment. First, to help experts comprehensively enact potential gestures, we provided them with a set of gesture classification cards (see Figure 2 ④) based on McNeill's classification [70] as widely used in sociology and educational sciences to analyze interpersonal communication (e.g., in classroom teaching or social conversing). The cards provide clear definitions and examples of each gesture type:

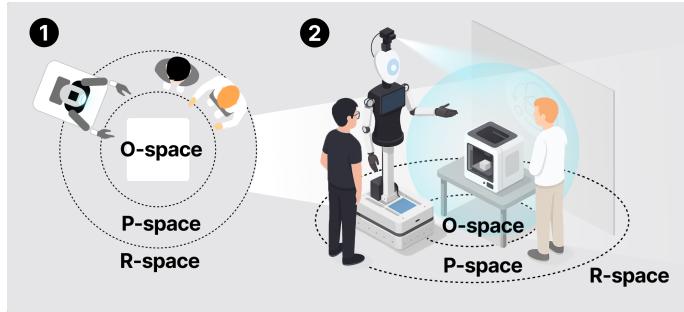


Fig. 3. ① Illustration of F-formation; ② F-formation applied to in-situ projection placement in the makerspace.

(1) Deictic (pointing to entities or locations), (2) Iconic (depicting concrete objects or actions; e.g., holding both hands apart to indicate a part's width), (3) Metaphoric (mapping abstract concepts onto space; e.g., an upward hand sweep to indicate "increase"), and (4) Beat (small rhythmic strokes aligned with speech).

To help experts enact the positional relationships between the robot and the learners, we drew on the theory of F-formation [69], which conceptualizes spatial arrangements in interpersonal interaction (see Figure 3 ①) into three types of spaces: (1) O-space, a shared focus area in the center, where the device or object of interest should be located; (2) R-space, the surrounding interaction circle formed by the robotic TA and learners; (3) P-space, outsiders or bystanders remaining in the outer public space (not considered in our scenario). Related work has similarly leveraged F-formation to improve naturalness for robot–human engagement (e.g., the work by Yousuf et al. [114]). In addition, we offered the experts an on-site live demonstration of the robot's prerecorded gestures and sample visual presentations to help them understand the hardware capabilities.

The experts' enactment involved six makerspace devices (see Figure 2 ⑤). They visited each of the devices sequentially and enacted in turn the explanations they thought a robot should give to beginners. In each episode, they planned the key learning points, selected gestures (from the four types introduced earlier), described the visual presentation and its placement, and demonstrated how these outputs could coordinate with each other. Each enactment episode was followed by a brief discussion and on-the-spot iterations [96]. The whole session was video-recorded for analysis.

3.2.2 Major Insights. The bodystorming session resulted in rich and concrete design inputs. For instance, the video clips of experts' bodily enactments were categorized based on McNeill's framework[70] for later constructing the robot's gestural unit library. The experts' verbal descriptions were subjected to a content analysis[33] to extract key learning points of each device, as well as examples on how the robot's speech and visual display could be designed. To gain structured insights, we summarize all these inputs into six major design requirements (**DR1-DR6**):

DR1- Verbal explanation should cover how it works, how to operate it, and safety information: the experts' verbalization consistently indicated three types of learning points: principles of how the equipment works (e.g., forming principles of FDM and resin printers), basic methods of how to operate the equipment, and safety reminders (e.g., do not touch the soldering iron tip by hand).

DR2- Verbal explanation should show connections between devices and use clear transitions to move between them: the experts' verbal articulation showed the need for coherent narration that integrates both content connections (such as comparing resin-based and FDM 3D printers and contrasting their forming principles) and clear transitions (e.g., "*This equipment has now been introduced; please follow me to the next one*"-M1).

DR3- Leveraging deictic gestures for orienting learners' attention in the makerspace: the experts' bodily demonstrations showed how deictic gestures are essential for orienting learners' attention toward the intended equipment, area, or direction (as shown in Figure 2 ⑤), during a makerspace tour. Deictic gestures were also used to emphasize visual display content. For instance, when introducing the laser cutter, M1 said, “*The operating procedure can be found in the image,*” while pointing toward the projection. Similarly, M3 suggested pointing to the visual display to make learners notice a learning point. In total, the experts' enactments generated 18 deictic gesture exemplars.

DR4- Utilizing iconic and metaphoric gestures to vividly reinforce learning concepts: the experts employed iconic gestures to mimic or illustrate specific objects or operations, thereby making verbal explanations more vivid. For example, using hand spanning or spacing to indicate workpiece size (M3) or 3D scanning distance (M2). Or covering eyes to warn against laser exposure (M1). In addition, metaphoric gestures were enacted to convey and reinforce abstract concepts, such as waving to stress “*don't leave waste*”-M1, or pushing hands forward to indicate prohibition (M2), or finger tapping to highlight cautions (M3). In total, experts generated 15 iconic and 12 metaphoric gesture exemplars on site, which we used to build the robot's gestural unit library.

DR5- Visual displays should use simplistic and clear graphics and concise text to highlight key information: the experts emphasized that the robot's visuals should remain simplistic and clear, avoiding lengthy text or overly complex illustrations. M3 envisaged that, unlike museum or exhibition tours, makerspace tours could benefit from visuals that resemble product manuals—using simple yet clear graphic styles. M2 stressed the importance of highlighting key information points, such as dimensions, distances, or temperatures, by directly labeling numbers on schematic diagrams. M1 further noted the value of visually emphasizing critical parts of the device, such as warning signs or the location of an emergency stop button.

DR6- Visual displays need to be temporally aligned with the robot's speech and gestures: through embodied enactments, the experts commonly emphasized that displayed visuals should appear in synchrony with the relevant narration and gestures, much like slides in a lecture. As M2 noted, when introducing operational steps, the robot should verbally describe each step while simultaneously presenting the corresponding image. M2 further suggested that sometimes, visuals could also follow the pointing gestures, such as highlighting a region with arrows while the hand gesture indicates the same area.

4 The ProjecTA System

In this section, we describe how the inputs from our formative study were translated into the design and implementation of the ProjecTA system.

4.1 Design of ProjecTA

ProjecTA is a prototype system meant to probe how a robotic TA with in-situ projection would support learners in a physical learning environment by offloading certain repetitive, standardized tasks from educational routines. In this study, ProjecTA has been specifically designed to offer novice learners guided tours in a makerspace. ProjecTA could present information through speech, body gestures, as well as in-situ projection onto physical objects or nearby surfaces in the space. As shown in Figure 4, we illustrate its design details through an exemplar usage scenario:

Eager to explore the university makerspace, beginners Jeff and Jerry started a tour led by ProjecTA, a robotic guide. At the entrance, ProjecTA greeted them (see Figure 4 ①), giving them an overview of the tour before leading them to the northeastern corner, where several machines are placed in line.

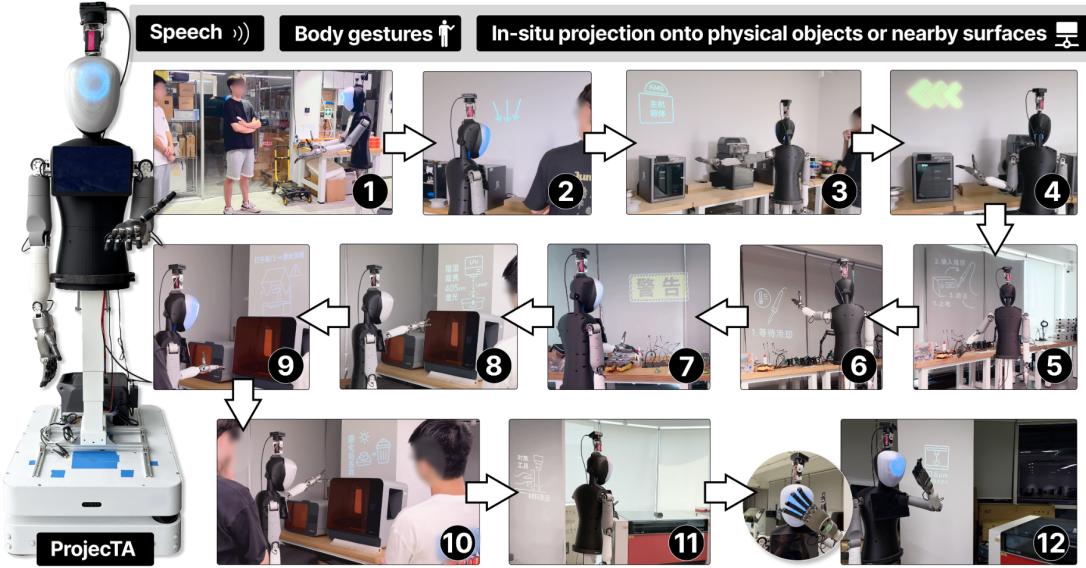


Fig. 4. An exemplar usage scenario of ProjectTA during makerspace guided tours

To make the instructional target clear, ProjecTA projected an arrow above the FDM 3D printer² (**DO1**, Figure 4 ②). “Let’s look at the FDM 3D printer [...] Notice its two main components: the automatic material switching system on top and the chassis below.” As the explanation unfolded, its projection displayed a simplistic and clear diagram (**DR5**, Figure 4 ③): a semicircle and a square, labeled “Automatic Material Switching (AMS) System” and “Chassis,” respectively. This helped the students quickly grasp the printer’s structure. Upon finishing, ProjecTA projected an animation of moving leftward to signal a transition, “Next, we will look at the soldering station, please follow me.” (**DR2**, Figure 4 ④)

At the soldering station, ProjecTA projected the operating steps right next to the soldering iron, so Jeff and Jerry could follow the instructions without looking away (**DO2**, Figure 4 ⑤). When the robot described how to replace the soldering tip, the projected steps updated in real time. First, “Step 1: Wait for cooling,” and then “Step 2: Loosen the nut,” keeping the visual instructions synchronized with the verbal explanation (**DR6**, Figure 4 ⑥). Afterward, the in-situ visuals switched to a warning symbol, “[...] Beyond the basic function and operation,” ProjecTA added, “there are also safety precautions to keep in mind [...].” (**DR1**, Figure 4 ⑦)

A directional projection then guided them to the next machine: the resin-based 3D printer. “Unlike FDM printers, this machine cures resin layer by layer with a laser focused on the tank bottom,” ProjecTA explained. As the curing process is a complex concept and cannot be demonstrated on the spot, the robot projected a diagram to complement and summarize its verbal explanation (**DO3**, Figure 4 ⑧). Moving on, ProjecTA used pointing to direct the learners’ focus to the printer’s lid, explaining that opening this lid would automatically shut off the laser (**DR3**, Figure 4 ⑨). Subsequently, the robot pointed to the newly projected visuals to stress the proper disposal of waste resin (**DR3**, Figure 4 ⑩).

Finally, they stopped at the laser cutter, whose key component, the focusing lens, was not visible to visitors—it was too small and hidden inside the machine. Hence, ProjectTA projected a magnified illustration to show what this

²Fused Deposition Modeling: a 3D printer melts plastic filament and deposits it layer by layer to form an object.

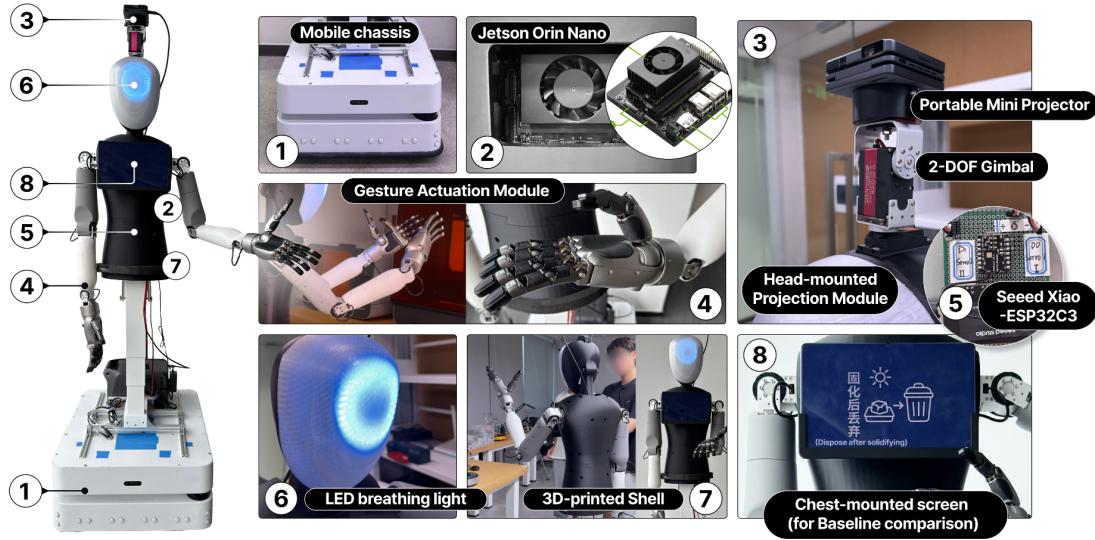


Fig. 5. Hardware components of ProjecTA consist of a mobile chassis, a gesture actuation module, a head-mounted projection module, and a chest-mounted screen (for baseline comparison).

component looked like and how it worked (DO4, Figure 4 (11)). The robot continued its explanation. When it came to safety, the robot raised its hand in a shielding gesture in front of its eyes as a warning and cautioned, “Don’t look directly at the laser.”(DR4, Figure 4 (12))

The tour concluded with ProjecTA stating, “This brings the session to an end. I hope the explanations have been helpful.” Jeff and Jerry had completed their first learning experience in the makerspace.

4.2 System Implementation of ProjecTA

4.2.1 Robotic Hardware. As shown in Figure 5, ProjecTA consists of a Gesture Actuation Module, a head-mounted projection module, and a mobile chassis. The Gesture Actuation Module integrates two 6-DOF³ arms driven by Eyoubot⁴ planetary torque motors and a 20-DOF dexterous hand from Linkerbot⁵ (see Figure 5 (4)), controlled via a Jetson Orin Nano (see Figure 5 (2)). The projection module carries an Aurzen ZIP Tri-Fold Portable Mini Projector⁶ (720p, 100 ANSI lm, USB-C, auto-focus) mounted on a 2-DOF gimbal (see Figure 5 (3)) powered by a Seeed Xiao-ESP32C3 with Wi-Fi pan/tilt control (see Figure 5 (5)). To be noted, this head-mounted solution for projector placement was determined for two key reasons: (1) unlike an arm-mounted solution, it does not interfere with the robot’s ability to perform gestures; (2) it provides a full 360-degree projection range, overcoming the body occlusion that prevents a torso-mounted projector from reaching the robot’s backside.

The mobile chassis is a Wheeltec S300⁷ with dual M10P LiDARs for 360° Simultaneous Localization and Mapping (SLAM), controlled via a Jetson Orin Nano (see Figure 5 (1)). The robot’s shell is 3D-printed in PLA (see Figure 5 (7)), and

³Degrees of freedom: the number of independent axes a mechanism can move or rotate around.

⁴<http://www.eyoubot.com/>

⁵<https://linkerbot.cn/>

⁶<https://aurzen.com/>

⁷<https://wheeltec.net/>

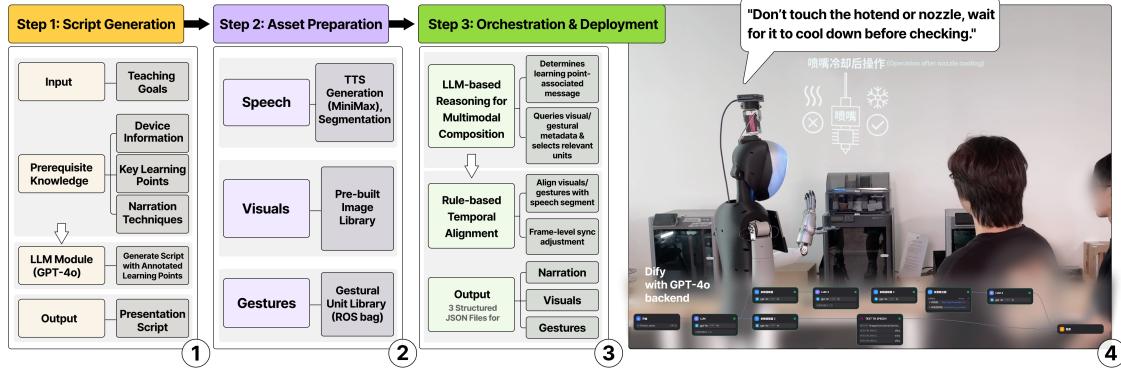


Fig. 6. Presentation choreography workflow, aligning speech, visuals, and gestures on a unified timeline.

its face is made of semi-transparent PETG housing an LED breathing light to indicate speaking status (see Figure 5 (6)). The projector and 2-DOF gimbal communicate with the control PC via MQTT⁸ over Wi-Fi, while the Gesture Actuation Module and mobile base use CAN bus⁹ with Robot Operating System (ROS) topics for control. During guided tours, the robot follows SLAM-based maps with preset device positions, with occasional human intervention for minor adjustments due to navigation limits.

4.2.2 Presentation Choreography Workflow. This workflow enables educators and technicians to edit, preview, and verify ProjecTA’s presentations across modalities. Based on input teaching goals, i.e., which equipment to cover and in what sequence, it compiles an executable choreography script that controls ProjecTA during a tour. The workflow assembles and schedules (1) generated speech segments, (2) existing visual display assets, and (3) predefined gesture units, then coordinates them across channels at runtime. The workflow also makes ProjecTA readily extensible and deployable in future similar learning contexts (other makerspaces or museums, etc.).

Presentation Script Generation: As shown in Figure 6 (1), this LLM-based module generates ProjecTA’s script based on the teaching goals (which equipment to cover and in what sequence) and a set of prerequisite knowledge.

[Prerequisite Knowledge] consists of Device Information, Key Learning Points, and Narration Techniques.

Device Information. The descriptive information about all pieces of equipment in the makerspace was pre-collected from reliable sources.

Key Learning Points. Based on experts’ verbalization in the bodystorming session, we compiled the essential learning points that should be covered when explaining each piece of makerspace equipment to beginners. These points include basic knowledge about how it works, how to operate it, and safety precautions (**DR1**), e.g., “PLA material is recommended to be printed at 210–220 °C” (for more details, please refer to Supplementary Material A). We also referenced resources widely used in maker education to ensure comprehensiveness [6, 63, 68]. These resulting key learning points serve as the central materials: ProjecTA’s presentation scripts are structured around them, and its visual assets and gestural units are also designed to communicate these key learning points effectively to novice learners.

Narration Techniques. Narration techniques are extracted from the bodystorming enactments by the experts, e.g., using everyday analogies to facilitate comprehension, or inserting transitions between, or drawing connections across

⁸Message Queuing Telemetry Transport: a lightweight publish/subscribe messaging protocol widely used in IoT/robots.

⁹Controller Area Network: a robust, low-level communication network originally from automotive systems, used for reliable device control.

related machines (**DR2**). Based on both expert enactments and our own testing, the narration for each device was kept to about 4–5 minutes.

[Script Generation]: We used the *GPT-4o-2024-08-06*¹⁰ model to generate the presentation scripts. The aforementioned *Device Information*, *Key Learning Points*, and *Narration Techniques* were integrated into the instruction prompt (see Supplementary Material A for the prompt). The prompt required the LLM to explicitly mark where each **key learning point** appears in the generated script, enabling easier verification and supporting subsequent multimodal orchestration.

Preparation of Speech, Visual Assets, and Gestural Units: As shown in Figure 6 ②, the preparation of multimodal resources are outlined below:

[Speech Generation and Segmentation]: The generated script was passed to MiniMax’s text-to-speech (TTS) service¹¹. As voice timbre was not a focus of this study, we adopted a standard explanatory voice from the MiniMax platform (see supplementary video). The generated audio was segmented into units according to semantic boundaries, ensuring each segment contained at most one *key learning point* for scheduling. Metadata for each speech segment included its duration and the *key learning point* it was aligned with.

[Visual Display Assets Library]: The visual display assets were organized as a pre-built image library designed to support novices to understand and consolidate each *key learning point*. Note that some *key learning points* require multiple images when they span steps of sub-concepts. As shown in Figure 7 ⑤, each image uses simplistic, clear graphics with minimal labels aligned with the narration (**DR5**).

The metadata includes each visual display asset’s intended projection location in the makerspace (e.g., on the equipment or nearby surfaces), a textual description, and the linked *key learning point* for later orchestration. See Supplementary Material B for details.

[Robot’s Gestural Unit Library]: We constructed a library of 42 robot gesture-control sequences, derived from 45 video-recorded bodily enactments that the experts performed to explain specific *key learning points*. These gestures included deictic (**DR3**), iconic, and metaphoric forms (**DR4**). Highly similar examples were merged, and each remaining sequence was teleoperated and recorded as a ROS bag file, producing a discrete gestural unit that can be triggered via the **Gesture Actuation Module**. Note that a single *key learning point* may map to multiple gestural units to support flexible orchestration.

The metadata of each unit includes the robot’s spatial location and orientation, matching the expert’s original demonstration for the given equipment, a textual description, a contextual note (indicating which equipment and *key learning point* it addresses and the accompanying narration), and its duration. See Supplementary Material C for details.

Orchestration of Speech, Visual Assets, and Gestural Units: As shown in Figure 6 ③, the orchestration process has two steps:

[LLM-based Reasoning for Multimodal Composition]: For each **speech segment**, the system determines whether its learning point-associated message needs further illustration or reinforcement, and then queries visual and gestural library metadata to select and combine the most relevant images, gestural units, or both.

[Rule-based Temporal Alignment]: Selected visuals and gestures are then aligned temporally with the associated speech segment (**DR6**). Selected display assets remain visible for the full segment; gestures are synchronized to playback. When multiple images are needed, each is aligned to the exact narration moment it explains (frame-level

¹⁰<https://platform.openai.com/>

¹¹<https://www.minimaxi.com/>

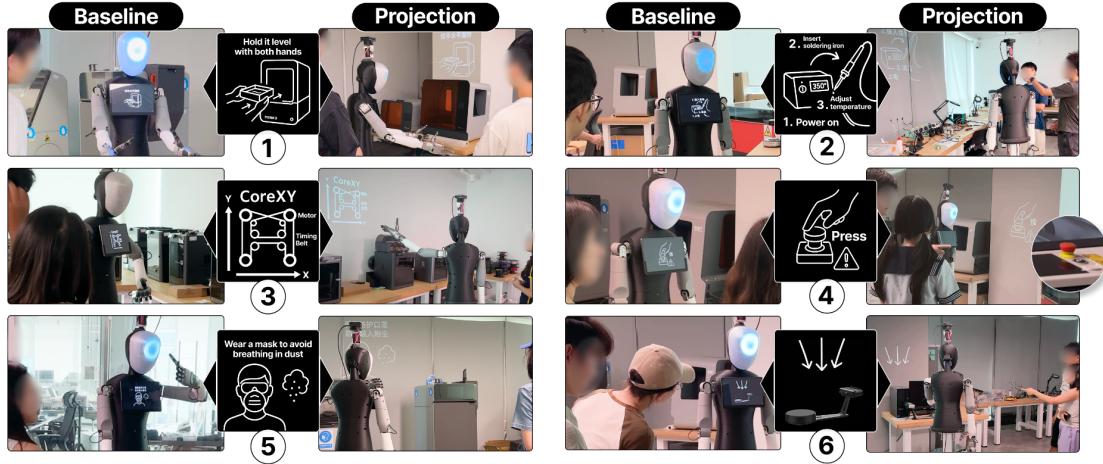


Fig. 7. Comparative examples of Baseline (chest-mounted screen) vs. ProjecTA (in-situ projection) guidance in makerspace tours.

synchronization). If a gesture exceeds the segment length, the inter-segment pause is extended to prevent overlap. Educators and technicians can preview the outcome on the robot and fine-tune timings across the modalities.

In our system, multimodal resources were orchestrated along a unified timeline. The LLM-based agentic pipeline was constructed upon Dify¹², with GPT-4o-2024-08-06 as the base model (refer Figure 6 ④) to achieve multimodal orchestration. All reasoning is performed offline to avoid runtime latency and allow for human preview and verification. The system outputs three structured JSON files encompassing narration, visuals, and gestures, each with execution times; the backend server converts these into synchronized, executable robot behaviors. This choreography system also enables users to pre-generate multiple narration variants and select the preferred one for further edits or final deployment.

5 Methodology

Our study aims to address the follow question: How does a robotic TA with in-situ projection, compared with a screen-based counterpart, affect learners' experiences during makerspace tours? (RQ) To investigate this question, we built and deployed ProjecTA in a real university makerspace (Figure 9). Adopting a mixed-methods, within-subject comparative approach, ProjecTA was contrasted with Baseline, a functionally equivalent counterpart using the chest-screen as its visual display.

5.1 Baseline system for Comparison

For comparison, the Baseline and ProjecTA systems were deployed on identical robotic hardware equipped with both a chest-mounted screen (Figure 5 ⑧) and a head-mounted projector (Figure 5 ③). This configuration extends prior findings of spatial AR and HUD studies in stationary operation (e.g., assembling or driving) to nomadic learning with a mobile robot. In our study, the guided walk-and-talk tour serves as a typical nomadic learning setting, providing a baseline cognitive-load context for comparing the two display modalities. The Baseline condition replicated the

¹²<https://dify.ai/>

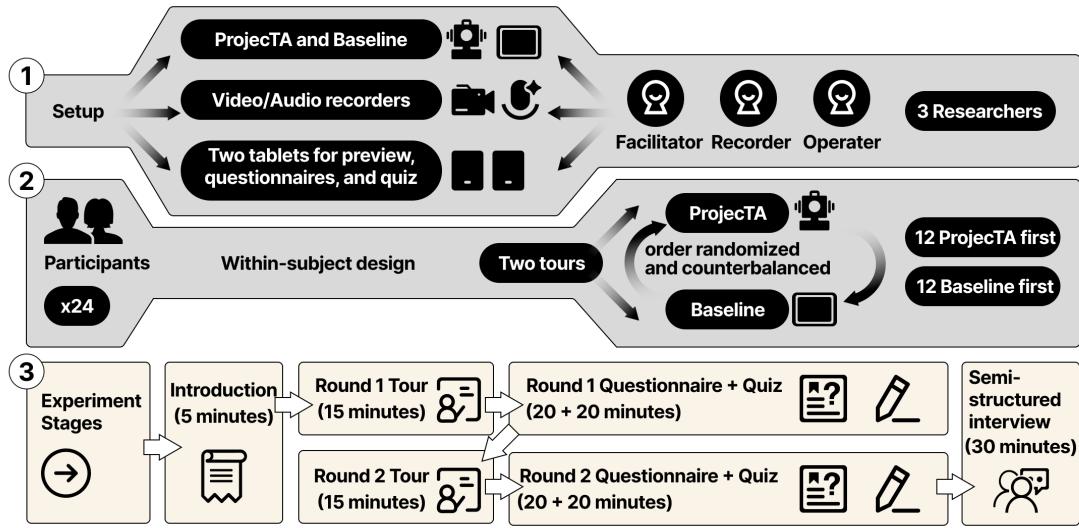


Fig. 8. ① Experimental setup, ② experimental design and ③ experiment stages.

standard screen-based displays of existing robotic platforms like Pepper [80]. Conversely, the ProjecTA condition utilized in-situ projection to examine differences in learner experiences, such as extraneous cognitive load. To ensure a fair comparison, both conditions ran on the same choreography instance generated by the Presentation Choreography Workflow, including the same narration, gestures, visual assets, tour route, and timing. To further ensure fairness for the Baseline, we made two necessary screen-specific optimizations: (1) deictic gestures that, in ProjecTA, pointed to projected overlays were replaced with gestures cueing viewers to look at the screen (Figure 7 ③); (2) spatial markers that ProjecTA placed on or around the physical equipment (arrows, halos, highlights) were mirrored on the screen by showing an image of the same equipment with identical markers (see Figure 7 ⑥). All other assets, including movement animations and learning point illustrations, were unchanged (see Figure 7 ①②④⑤). Therefore, the sole significant difference between the two conditions was the display modality (on-body screen versus in-situ projection), allowing us to assess their effects on learners' cognitive load and experiences during makerspace tours.

5.2 Participants

This study recruited 24 participants (13 male, 11 female; age 20–31, $M = 24.25$, $SD = 2.61$) via social-media posts. We targeted novice learners and screened out registrants who reported with rich makerspace experience before scheduling. We collected demographics via an online questionnaire, and respondents included university students and working professionals, spanning backgrounds in engineering/sciences ($n = 14$), humanities ($n = 3$), social sciences ($n = 4$), and art/design ($n = 3$). Sessions were conducted in dyads (12 pairs): 6 pairs who knew each other and 6 pairs of strangers. We assessed two aspects of participants' past experiences using five-point rating scales (1 = never, 5 = often): hands-on activities in makerspaces (19/24 rated 1) and guided equipment tours (22/24 rated 1). Familiarity with makerspace equipment was also measured on a five-point rating scale (1 = not at all familiar, 5 = very familiar), and participants mostly reported low familiarity. Detailed familiarity ratings are provided in Appendix B.

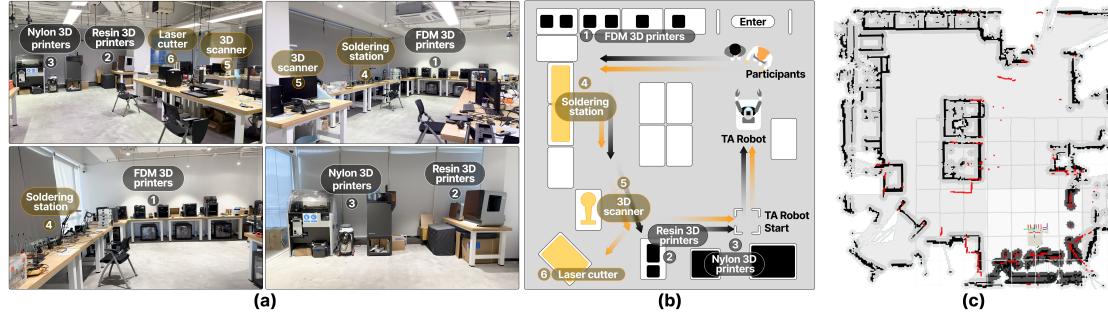


Fig. 9. Experimental environment: (a) photos of the university makerspace, (b) top-down schematic of equipment layout and robot path, and (c) RViz map of the space.

5.3 Setup and Procedure

The experiment was conducted in a university makerspace (Figure 9). As shown in Figure 8 ①, the setup consisted ProjectTA and Baseline; video and audio recorders for data collection; and two additional tablets for participants to preview the procedure and complete questionnaires and post-tour quiz during the session. Two researchers were present in each session: one facilitated the session and one recorded the procedure. A third researcher, located in an adjacent space, operated the backend system.

To compare ProjectTA with the Baseline, we used a within-subject comparative design (see Figure 8 ②). Each participant dyad completed two robot-guided tours, one with ProjectTA and one with the Baseline; modality order was randomized and balanced (12 ProjectTA first, 12 Baseline first). To avoid learning the same content twice, we prepared two non-overlapping tour sets, A and B, each covering three different machines: Set A included the FDM 3D printer (Figure 9 ①), Resin 3D printer (Figure 9 ②), and Nylon 3D printer (Figure 9 ③); Set B included the Soldering station (Figure 9 ④), 3D scanner (Figure 9 ⑤), and Laser cutter (Figure 9 ⑥). The tour sets were treated as a counterbalanced random factor to offset learning effects. Every participant experienced both sets across the two rounds, with set order counterbalanced within each modality order: half of the participants experienced A first and half B first. This scheme isolates the effect of display modality while mitigating cross-condition learning effects [7] and order effects. The experimental stages (Figure 8 ③) were as follows:

Introduction (5 minutes): The researcher introduced the overall procedure, including the structure of the two tour-guided rounds and emphasized safety precautions (e.g., the robot's emergency stop). Participants then signed consent forms and completed a demographic questionnaire, along with questions about prior experience with makerspace equipment and guided tours.

Round 1 Tour (15 minutes): Participants completed the first round of robot-guided learning with one display modality (projector or screen). Two participants jointly followed the explanations.

Round 1 Questionnaire + Quiz (20 + 20 minutes): After the tour, participants first completed questionnaires that included measures of cognitive load, user engagement, and a customized scale, followed by a first-round-specific post-tour quiz aligned with three devices introduced in that round.

Round 2 Tour (15 minutes): Participants experienced the other display modality in the second round. The robot introduced the remaining three devices. The process was identical to Round 1.

| ID | Caption | Details |
|---------|--|--|
| RS1–4 | Intrinsic Load (learning material's inherent complexity) | the Intrinsic Load subscale in Cognitive Load Scale (CLS)[59] |
| RS5–8 | Extraneous Load (avoidable effort caused by presentation) | the Extraneous Load subscale in Cognitive Load Scale (CLS)[59] |
| RS9–11 | Focused Attention (immersiveness and absorption) | the FA subscale in User Engagement Scale Short Form (UES-SF) [79] |
| RS12–14 | Perceived Usability | the PU subscale in User Engagement Scale Short Form (UES-SF) [79] |
| RS15–17 | Aesthetic Appeal | the AE subscale in User Engagement Scale Short Form (UES-SF) [79] |
| RS18–20 | Reward Factor (rewarding experience) | the RW subscale in User Engagement Scale Short Form (UES-SF) [79] |
| RS21 | Usefulness of visual display | The visual display of the robot system was very helpful for my understanding of the presentation content. |
| RS22 | | The visual display of the robot system made it easier for me to focus my attention on the corresponding area of the actual equipment or space. |
| RS23 | Multi-modal complementary (visuals, gestures, and speech) | When the robot system performed an action, I could easily locate the related information on its visual display. |
| RS24 | | The robot system's actions, speech, and visual display content were well-coordinated. |
| RS25 | | The visual display of the robot system effectively complemented its speech and actions. |

Table 1. Item-type distribution in rating scales.

Round 2 Questionnaire + Quiz (20 + 20 minutes): As in Round 1, participants completed a same questionnaire and a second-round-specific post-tour quiz.

Semi-structured interview (30 minutes): After both rounds, participants engaged in a semi-structured interview and participants received 200 CNY as compensation for their time.

5.4 Data Gathering

5.4.1 Questionnaire. The questionnaire comprised 25 rating items (RS1–RS25; Table 1). RS1–RS8 used the streamlined Cognitive Load Scale (CLS) [59] with two subscales: Intrinsic Load (RS1–RS4), capturing the effort caused by the learning materials' inherent complexity and learners prior knowledge, and Extraneous Load (RS5–RS8), capturing the avoidable effort induced by the presentation of the learning materials. Because the equipment sets (learning materials) were counterbalanced within each condition, we did not expect condition differences on Intrinsic Load. Our primary test was whether presentation mode (ProjecTA vs. the Baseline condition) differentially affected Extraneous Load.

RS9–RS20 used the User Engagement Scale Short Form (UES-SF) [79] with four subscales: Focused Attention (RS9–RS11; experienced immersiveness or absorption), Perceived Usability (RS12–RS14), Aesthetic Appeal (RS15–RS17), and Reward Factor (RS18–RS20). We used the UES-SF to assess experiential differences beyond cognitive load.

RS21–RS25 were custom items aligned with our research goals. RS21–RS22 assessed the perceived usefulness of the visual displays for (a) understanding the presented content and (b) maintaining focus on the relevant object or region in the space. RS23–RS25 assessed perceived multimodal complementarity among visual display, gesture, and speech across the two conditions.

5.4.2 Post-tour Quiz. Each tour was followed by a post-tour quiz. Each quiz comprised 30 items covering all three pieces of equipment from that tour. Items targeted each equipment's *Learning Points*: how it works, how to operate it, and safety precautions. Question formats included single-answer multiple-choice and ordering (sequence) items. Each correct response (i.e., a correct option or a correctly placed position in an ordering item) was worth one point; raw scores

were then linearly normalized to a 0–100 scale for analysis. After each tour, participants completed the questionnaire first and then the tour-specific quiz. This was intended to insert a brief, content-related delay, reducing reliance on immediate/short-term memory (e.g., verbatim recall of just-seen labels) and yielding a more valid assessment.

5.4.3 Interview. We conducted semi-structured interviews in dyads to collect in-depth user feedback. Interviews began with questions about participants' experiences in both conditions, probing how they followed the presentations in each round and what difficulties they encountered. We then asked them to describe positive and negative experiences with each condition and to suggest improvements. Next, participants were asked to provide concrete examples from both conditions and state their preference with rationales. Finally, we invited suggestions for other scenarios where the system could be applied. Each interview lasted around 20 minutes and was video-recorded for analysis; the complete question list is provided in [Appendix A](#). We conducted a thematic analysis following Blandford's guidelines [13], two researchers conducted collaborative inductive coding. They initially annotated the transcript to identify relevant quotes, key concepts, and preliminary patterns in the data. These initial insights were further developed through regular discussions among four researchers, leading to a detailed coding scheme aligned with the research objectives. Quotes were then coded and clustered into a hierarchy of emerging themes, continually reviewed, and refined in recurrent meetings, where exemplar quotes were also selected to illustrate each theme and sub-theme. Alongside this, the team reviewed and annotated the session videos, keeping both the research questions and the emerging thematic structure in view. We collected the video segments that served as evidence or exemplars for the thematic analysis results, especially those highlighting behaviors of participants during nomadic guiding tour in makerspace. In addition, we supplemented our analysis with photographic documentation of key interaction moments, spatial arrangements, and notable projection–object–participant configurations observed during the sessions.

6 Findings

6.1 Quantitative Results

The internal consistency of two custom subscales was assessed. Cronbach's α indicated acceptable reliability for Multimodal Complementarity (RS23–RS25) in both conditions (ProjecTA $\alpha = .75$; Baseline $\alpha = .72$; $> .70$). Visual Display Usefulness items (RS21–RS22) were $< .70$ in both conditions and were analyzed separately.

Normality of within-participant difference scores was tested using the Shapiro–Wilk test. Assumptions were met for Intrinsic Load (RS1–RS4), Extraneous Load (RS5–RS8), Focused Attention (RS9–RS11), Perceived Usability (RS12–RS14), and Quiz Performance; these outcomes were compared with paired-samples t-tests. Aesthetic Appeal (RS15–RS17), Reward Factor (RS18–RS20), RS21, RS22, and Multimodal Complementarity (RS23–RS25) violated normality and were analyzed with Wilcoxon signed-rank tests. The results are presented in [Figure 10](#).

Intrinsic Load (RS1–4) captures the effort required by the material's inherent complexity and the learner's prior knowledge. No significant difference was found: $M_{\text{ProjecTA}} = 3.96$, $SD_{\text{ProjecTA}} = 2.31$; $M_{\text{Baseline}} = 4.21$, $SD_{\text{Baseline}} = 2.36$; $M_{\text{diff}} = -0.25$, $SD_{\text{diff}} = 2.50$; $t(23) = -0.489$, $p = 0.629$, Cohen's $d = -0.100$.

Extraneous Load (RS5–8) captures the avoidable effort introduced by how information is presented. ProjecTA was significantly lower than Baseline with a large effect size (Cohen's $|d| > 1$): $M_{\text{ProjecTA}} = 1.90$, $SD_{\text{ProjecTA}} = 1.01$; $M_{\text{Baseline}} = 4.04$, $SD_{\text{Baseline}} = 1.61$; $M_{\text{diff}} = -2.15$, $SD_{\text{diff}} = 1.98$; $t(23) = -5.322$, $p < 0.001$, $d = -1.086$.

Focused Attention (RS9–11) captures immersiveness and absorption in experience. No significant difference found: $M_{\text{ProjecTA}} = 4.00$, $SD_{\text{ProjecTA}} = 0.73$; $M_{\text{Baseline}} = 3.79$, $SD_{\text{Baseline}} = 0.93$; $M_{\text{diff}} = 0.21$, $SD_{\text{diff}} = 0.73$; $t(23) = 1.390$, $p = 0.178$, $d = 0.284$.

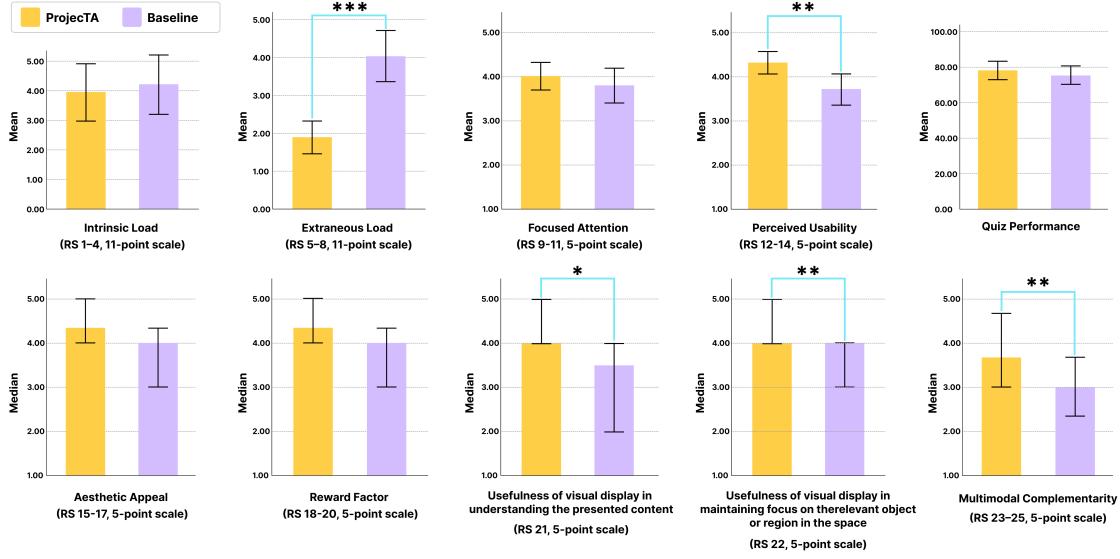


Fig. 10. Quantitative results comparing ProjecTA with Baseline (* $p < .05$, ** $p < .01$, *** $p < .001$; error bars indicate 95% CI).

Perceived Usability (RS12-14) were reverse-worded and were reverse-scored for analysis and presentation (higher scores reflect greater usability). *ProjecTA* was significantly higher than *Baseline*: $M_{\text{ProjecTA}} = 4.31$, $SD_{\text{ProjecTA}} = 0.60$; $M_{\text{Baseline}} = 3.71$, $SD_{\text{Baseline}} = 0.84$; $M_{\text{diff}} = 0.60$, $SD_{\text{diff}} = 0.86$; $t(23) = 3.393$, $p < 0.01$, $d = 0.692$.

Quiz Performance. No significant difference found: $M_{\text{ProjecTA}} = 78.13$, $SD_{\text{ProjecTA}} = 12.11$; $M_{\text{Baseline}} = 75.52$, $SD_{\text{Baseline}} = 12.16$; $M_{\text{diff}} = 2.61$, $SD_{\text{diff}} = 13.77$; $t(23) = 0.927$, $p = 0.364$, $d = 0.189$.

Aesthetic Appeal (RS15-17). No significant difference (Wilcoxon): $Z = -1.371$, $p = 0.171$, $r = 0.280$. Medians (IQR): *ProjecTA* = 3.17 [2.67, 4.00]; *Baseline* = 3.00 [2.42, 3.92]; median difference = 0.00 [-0.33, 0.33].

Reward Factor (RS18-20). No significant difference (Wilcoxon): $Z = -1.703$, $p = 0.089$, $r = 0.348$. Medians (IQR): *ProjecTA* = 4.33 [3.75, 5.00]; *Baseline* = 4.00 [2.75, 4.83]; median difference = 0.17 [0.00, 0.67].

Usefulness of visual display in understanding the presented content (RS21) *ProjecTA* significantly exceeded *Baseline* with a medium-large effect (Wilcoxon): $Z = -2.429$, $p < 0.05$, $r = 0.496$. Medians (IQR): *ProjecTA* = 4.00 [4.00, 5.00]; *Baseline* = 4.00 [3.00, 4.00]; median difference = 0.00 [0.00, 1.00].

Usefulness of visual display in maintaining focus on the relevant object or region in the space (RS22) *ProjecTA* significantly exceeded *Baseline* with a large effect (Wilcoxon): $Z = -3.096$, $p < 0.01$, $r = 0.632$. Medians (IQR): *ProjecTA* = 4.00 [4.00, 5.00]; *Baseline* = 3.50 [2.00, 4.00]; median difference = 1.00 [0.00, 2.00].

Multimodal Complementarity (RS23-25) assesses the perceived coordination and mutual reinforcement among the robot's visual displays, gestures, and speech. *ProjecTA* significantly exceeded *Baseline* with a large effect (Wilcoxon): $Z = -2.958$, $p < 0.01$, $r = 0.604$. Medians (IQR): *ProjecTA* = 3.67 [3.00, 4.67]; *Baseline* = 3.00 [2.33, 3.92]; median difference = 0.50 [0.00, 1.00].

Overall, the most significant finding is that ProjecTA substantially reduced extraneous cognitive load: the avoidable effort demanded by the presentation of learning materials [22]. While both the ProjecTA condition and baseline condition achieved comparable quiz scores, ProjecTA's ability to lessen avoidable cognitive effort suggests strong potential for Manuscript submitted to ACM

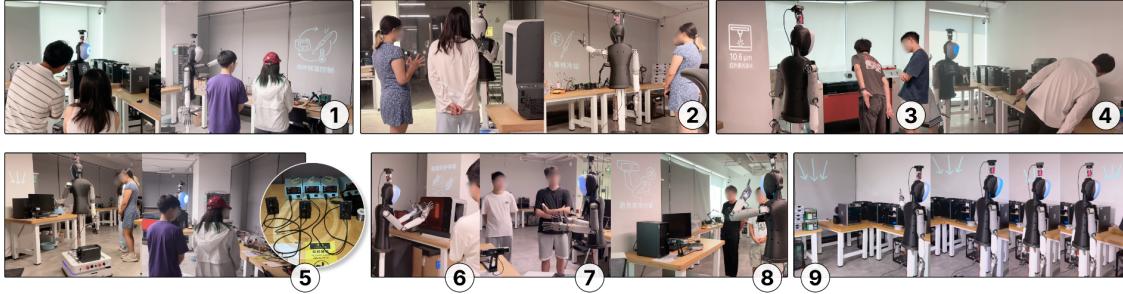


Fig. 11. Illustration of qualitative findings with ProjecTA in real makerspace scenarios. Compared with the Baseline, ProjecTA ① reduced attention switching between physical referents and displayed content, ② enabled shared access to display information for the group, ③ enlarged critical information such as numeric data, ④ highlighted hard-to-see or hidden parts of equipment, ⑤ projected directional arrows and on-object highlights onto hand-indicated targets, ⑥ mimed protective glove use with a projected reminder, ⑦ used gestures to indicate size and range, ⑧ covered its eyes to introduce the "don't look at the laser" safety point, and ⑨ offloaded spatial descriptions by projecting arrows to indicate in-situ 3D printer locations.

future robotic TA implementations. As expected from the counterbalanced design, there was no significant difference in the learning materials' inherent difficulty (intrinsic load). Because the two conditions used identical hardware and the same choreography and differed only in visual modality, there were no significant differences in aesthetic appeal, immersiveness (focused attention), or rewarding experience. By contrast, ProjecTA scored significantly higher on outcomes tied to information delivery and cross-modal coordination: higher perceived usability, more effective visual displays for understanding content and staying oriented to physical referents, and stronger multimodal complementarity, with visuals, gestures, and speech perceived as working together more coherently.

6.2 Qualitative Results

Qualitative results from participants' interviews help us further contextualize and exemplify their experiences:

6.2.1 How ProjecTA's In-Situ Projection Supported Learning in the Physical Space. Our results identified several key patterns in how the in-situ projection of ProjecTA brought unique advantages to participants' learning during the guided tour.

In-situ projection reduced attention switching between physical referents and displayed content. As shown in Figure 11 ①, by anchoring visuals adjacent to the target object, ProjecTA's in-situ projection reduced the learners' effort to switch between the visual display and the physical referents. Participants (23/24) consistently reported that placing cues on or near the objects reduced eye/head movements and shortened the path from locating the device to identifying the exact part. As P6 stated, "*Projection placed the visual content next to the machine, so I can see the content besides the equipment at a glance, without switching back and forth*". Similarly, P13 appreciated the projected content being "*closer to the device*," which "*naturally guides attention to the actual machine*" (P17). Several contrasted this with the chest-mounted screen: "*[with projection], I don't have to turn back to the screen; projection is more direct and makes it easier to engage with the target. With the screen, the constant attention switching makes it harder to focus*" (P1). Additionally, Nearly half of the participants (11/24) mentioned that in-situ projection helped them feel less pressure of missing information. As P5 felt, "*turning from the screen to the machine risks missing an explanation*". As P17 put it: "*The screen feels discontinuous [...] when my attention switched from the screen to the equipment, [the Baseline] was*

still talking and showing visuals, and I felt I suddenly missed something". These accounts underscore in-situ projection's advantage over screen-centric paradigms in reducing misalignment and loss of information caused by gaze shifts.

In-situ projection improved shared access to information in the physical space. In a group setting, the chest-mounted screen created a narrow, private viewing cone, hindering spatial access to the content (Figure 11 ②). Participants reported having to reposition themselves to see the screen (P2-3), crowding and mutual occlusion when standing side-by-side (P7-8), or inconvenience to look at the screen from lateral viewing positions (P4, P7, and P22). Viewed through F-formation [69], Baseline tended to fix the *O-space* (primary interaction space) in front of its screen, forcing the participants to frequently shuffle and adjust head orientation to participate. As P3 noted: "*During the laser cutter explanation, [the Baseline] was angled, so I had to [...] keep adjusting my position to catch the information*". By contrast, in-situ projection created a large, public, and shared visual field that was legible from multiple angles. Participants noted: "*you can see it clearly from all angles, and it's bigger*" (P2), and "*even from the back I can still see it*"(P5).

In-situ projection highlighted key information, and visualized hidden or hard-to-see parts. In the physical learning setting, in-situ projection helped enlarge critical information and key components: P3, P7, P8, and P24 valued the usefulness of enlarging numeric data next to the machines, such as the working temperature of equipment, or the size of the workpiece. As P3 stated, "*the details such as numbers stand out more clearly*" (see Figure 11 ③). P7, P13, P14, and P18 had positive experiences regarding the projection using "zoom-ins" to help locate intricate components. P18 shared an example of this: "*The projected [...] emergency STOP button immediately helped me locate it on the actual machine*" (see Figure 7 ④). Moreover, participants appreciated that certain hard-to-see or hidden parts of equipment were 'brought to the foreground' immediately by the projection. For instance, P13 shared an example of a back-side feeder reel (see Figure 11 ④), which was hidden behind the machine. The projection helped them to notice and learn about this module. Similarly, P7 mentioned an example in which the projection helped them situate internal nozzles without direct sight.

Challenges. Despite its clear advantages, participants identified several practical challenges. Robotic arms' motion sometimes caused image jitter, which might be distracting (P16). Strong ambient lighting could diminish the projection's clarity and color fidelity(P5). Finally, projecting onto uneven or angled surfaces introduced visual distortions. For instance, P18 reported an issue with perspective skew: "*the projection is not a flat plane, there's a bit of oblique perspective*." P15 also reported an issue with occlusion, where a lower part of the projection became less readable when blocked by an object. Despite that the simplistic graphics were generally appreciated for their clarity, some participants (P3, P7, P8, P24) pointed out that when explaining rather complex parts, higher fidelity images such as photos of internal components might be more productive for learning.

6.2.2 Examples of Preferred Complementarity Across ProjecTA's Projection, Gesture, and Speech. Our results offered rich and vivid examples of how ProjecTA's projection, gestures, and speech supplement or amplify one another.

Projection augmented or disambiguated deictic gestures for clear reference. Upon experiencing Baseline, some participants indicated that the robot's deictic gestures (e.g., pointing) alone can be ambiguous sometimes: For instance, as felt by P15, "[Baseline] sometimes pointed not exactly at the content. I could infer its intent, but it didn't pinpoint the target". In contrast, augmenting these deictic gestures with in-situ projections, such as directional arrows and on-object highlights, effectively resolved this ambiguity, enabling a smoother experience, especially when handing off instructions between devices (see Figure 11 ⑤): "*when switching to a new device, ProjecTA first projected an arrow onto the machine*,

followed by the pointing gesture, providing a more natural visual handoff”(P2). In cluttered environments, “*a projected halo on the target surface plus a brief pointing motion helped me identify the object of explanation*” (P5).

Iconic and metaphoric gestures concretized or reinforced spoken and displayed explanations. Over half of participants (15/24) recalled The robot’s use of iconic and metaphoric gestures that demonstrated and reinforced the spoken or displayed content. They reported some memorable examples about its *iconic gestures*, which directly mapped to physical actions or properties. For example, P4 described how ProjecTA “*traced the distance with its hands at the 3D scanner to indicate size and scan range, synchronized with projected numbers*.”(see Figure 11 ⑦) Other telling examples include miming a horizontal insertion (Figure 7 ①) and then a subsequent operation, which aligned with the verbal and projected instruction (P22), or miming putting on protective gloves along with a projected reminder (P19 and P24, refer Figure 11 ⑥), or performing a powder-mixing motion at the nylon printer (P10). The robot’s *Metaphoric gestures* were experienced to symbolically reinforce abstract concepts communicated by speech and visuals at the same time. Memorable examples included conveyance of prohibitions and warnings, such as the robot covering eyes to communicate ‘don’t look at the laser’ (P7, P9, P12, refer Figure 11 ⑧), or using pushing motions or repetitive waving (Figure 1 ⑦) to signal a warning (17/24). These embodied gestures created a strong link between the robot’s physical motions and the knowledge to convey, enhancing the vividness and expressiveness of the spoken and displayed explanation.

Projection offloaded spatial description from speech and clarified verbal content. Apart from disambiguating deictic gestures, projections was also experienced by the participants to offload and clarify the robot’s verbal narration which otherwise might be complex or cumbersome. Namely, with the projected cues illustrating the topic of discussion, the robot could replace verbose descriptions (e.g., ‘the second device from the left’) with simple, direct phrases like ‘this device’ or ‘this part’ (see Figure 11 ⑨). As P1 noted, “*the projection made the part [of the scanner being talked about] explicit, so I knew exactly what part the ProjecTA meant.*” Beyond offloading references to concrete targets from verbalization, projections could also simplify the robot’s verbal reference and description of abstract concepts: such as mentioning ‘this process’ or ‘this phenomenon’ verbally, while projecting an animated diagram or flowchart at the same time. This helped users grasp complex ideas more intuitively and seamlessly. For instance, a P1 experienced, “*when some knowledge was explained [verbally], the visuals of projection helped me grasp how it works.*” Similarly, P14 found the projection offloaded retention, stating, “*I don’t always remember what ProjecTA says, but once I see the projection I remember it immediately.*”

Challenges. Despite the above mentioned benefits of such cross-modal coordination, the participants also pinpointed some pragmatic challenges for future refinement. First, several (P11, P12) reported difficulty switching their focus between the robot’s projections and gestures, noting that the two modalities, occasionally, competed for their attention rather than complementing each other. Second, a few participants (P19, P23) reported that the robot’s pointing sometimes seemed imprecise or skewed, causing them to narrow their focus to the projection alone. Finally, several participants (P12, P17, P18) suggested adding even more verbal pre-announcement of the robot’s intentions for upcoming visual or gestural outputs. This verbal cue would help them mentally prepare for the upcoming action, creating a smoother and more predictable interaction.

6.2.3 Social Expectations and Initiatives from Participants. Participants also articulated social expectations for a robotic TA to support the nomadic learning tour. Several participants (10/24) framed ProjecTA as a social coordinator that mediated group access to content, rather than just delivering information. For instance, P8 noted that “*With the screen [of Baseline] I worried about bumping others; when two of us checked closely together it got crowded around the*

screen [...] whereas in-situ projection felt more naturally shared [...]. Echoing this, P16 noted that ProjecTA was better suited for facilitating “group access and joint attention”, making it “easier for the group to orient to the same part of the equipment.” Moreover, participants perceived the combination of projection and gestures as more socially coherent than the screen-gesture pairing; as P8 put, it “creates a sense of [...] inviting you to engage with the physical machine.” By contrast, the screen-based condition: “*the gesture is the gesture and the screen is the screen [...]*”

More than half of the participants (13/24) also reported that in-situ projections served as a trigger for their self-initiated exploration. P13, for example, recounted noticing a back-side feeder reel (see Figure 11 ④) that had been hidden behind the machine until the projection highlighted it, which prompted him to step in toward the machine, inspect the referenced module, and verify the robot’s explanation. P7 further noted that “*projecting the internal workings felt safer and more efficient for learning.*” These accounts suggest that projected content can create “hooks” that encourage learners to move, look closer, and personally verify robots’ explanation, rather than only passively receiving information.

Many Participants (17/24) preferred the robot to communicate in human-like and emotionally engaging ways. For instance, P6 particularly appreciated the robot’s welcoming poses at the start of the guiding tour and everyday metaphors in its explanations that helped them form intuitive understandings of technical components. P7 wished ProjecTA could adopt more expressive gestures and projected visuals for mobilizing learners’ emotions.

A few participants (3/24) also voiced a desire for more active, two-way interaction: they wanted not only to receive explanations but also to ask the robot questions during the tour. Although this feature was disabled in our study to maintain experimental control, this expectation indicates mixed-initiative dialogue as a natural extension. Finally, as ProjecTA moved between stations, P2 suggested that it could project its walking path and planned moves in the physical space to further increase its intuitiveness to learners in nomadic learning contexts.

7 Discussion

Robotic TAs can handle standardized, repetitive tasks and bring value to educational settings [11, 87]. Yet current deployments typically rely on a screen-based display [103, 113], which makes it hard to move visual information with the narration and place it directly near or right on the object [1]. Exploration on integrating in-situ projection with gesture-capable robotic TAs, as well as on aligning robots’ projected visuals, speech, and gestures, remains scarce. We thereby set out to explore how a robotic TA with in-situ projection, compared with a screen-based counterpart, affects learners’ experiences during makerspace tours (**RQ**).

To address this, we built ProjecTA, a robotic TA that links in-situ projection with speech and gestures to overlay information directly on or near target equipment for supporting learners in guided tours. In a real makerspace, 24 participants experienced ProjecTA and Baseline, its screen-based counterpart, in a controlled within-subject comparison.

Addressing **our research question (RQ)**, we used a mixed-methods approach. The quantitative results show that ProjecTA significantly lowered extraneous load with a large effect size, and enhanced perceived usability, usefulness of the visual display, and multimodal complementarity, indicating its ability to lessen unnecessary cognitive effort to improve learner experiences related to information delivery. We then supplemented these results with qualitative analysis to characterize more vivid experiences: how near-object visual overlays aligned with robotic TAs’ speech and gestures, reducing learners’ referent matching and attention switching in the space, etc., which provided more contextual understandings to **RQ**.

To further extend our findings regarding our **RQ**, we distill design implications from our empirical data and discuss them below to inspire and inform future research.

7.1 Design Implications

7.1.1 Implications 1: Unlocking More Design Possibilities Combining In-situ Projection with Deictic Gestures. Our study shows that in-situ projection can serve as a means of spatial referencing and can be mutually reinforced by robots' deictic gestures. These findings meaningfully extend **DO1** summarized in our formative study, and point to three promising directions:

Hand and projection co-referencing to spatial targets. Our study shows that projecting cues (arrows/halos) directly onto the hand-indicated target clarifies the reference, helping learners locate it quickly, reducing referent matching and gesture ambiguity. Building on this finding, and extending prior systems that implemented verbal-visual co-referencing [53, 75] or projection-only referencing (e.g., supporting nursing [15] and anatomy [36]), we envision hand-projection co-pointing where the robot points to the region with projected overlays marking detailed references such as paths, multiple sub-targets, or no-go zones, furthering current spatial referencing.

Local cues for robot pointing at the projection. To reduce referential ambiguity when robots are pointing at projected content, we suggest adding local cues (e.g., arrows, halos, borders) on referenced targets, similar to ProjecTA pointing to physical targets. Beyond overlaying cues on real objects, we found that when the robot points to projected content, adding local highlights or animations can also be helpful, for reducing visual search and improving comprehension. Participants also explicitly expected the robotic TA to mark the exact region within the projected content it was indicating, so they could immediately see which part of the projection the gesture referred to, extending prior designs, which mainly focused on augmenting physical targets (such as Visiobo [53]).

Projection for hard-to-reach or uncontrollable referents. Many prior systems projected overlays to ease professionals' operational performance [21], while less work focused on helping novice learners grasp hidden, unsafe, or fragile parts of an object without touching it. Building on our finding that in-situ projection revealed otherwise hidden details, we suggest that projection, as a precise, contact-free approach, could extend robot pointing for hard-to-reach or uncontrollable parts during nomadic learning tours (see Figure 6 (4)). For instance, one participant suggested projecting markers on an FDM printer's feed port and hot nozzle to avoid touch while speeding target localization. Beyond unsafe and risky scenarios, in our tours, projections onto hidden or hard-to-see parts prompted novices to move closer, inspect the equipment to check the robot's explanation. These verification and exploration behaviors are consistent with prior educational psychology research showing that actively initiating self-explanation and testing one's own interpretations can deepen conceptual understanding [24]. More designs could draw from this beyond makerspaces, e.g., robots in a chemistry lab could indicate the storage area for concentrated sulfuric acid and project a warning sign on the bottle label, as inspired by AR safety training [52]; In no-touch settings such as exhibitions, robot can indicate key details on the artifact without contact [77].

7.1.2 Implications 2: Thoughtful Placement of Projected Content to Reduce Learners' Attention Switching.

As some participants reported that the Baseline condition made them fear missing content and hesitate to look away from the screen, our qualitative findings in Section 6.2.1 similarly showed that ProjecTA's near-object overlays reduced frequent referent matching and encouraged learners to verify information directly on the equipment.

Our study shows ProjecTA reduced learners' attention switching between presentation content and the physical equipment (**DO2**). As shown in our qualitative results, some participants reported that Baseline made them fear missing content and hesitate to look away from the screen, whereas ProjecTA's near-object overlays reduced frequent referent matching. This extends prior findings in stationary tasks (e.g., assembly, repair, or driving [4, 55, 91, 104]), suggesting

the importance of thoughtful and finer-grained mechanism for the spatial placement of projected content in nomadic learning with mobile robots.

Projectable regions of the artifact as the display Existing systems have extensively explored fixed-position projections, whereas our study surfaces design opportunities and challenges of designing mobile projective systems for learners. In our design, some equipment itself served as the projection surface, for instance, the nylon printer's body (see Figure 1 ②). This way, overlays are placed directly on the suitable surfaces of explained artifacts, which could further reduce referent matching compared with wall projections [84]. To fully leverage this potential, a robotic TA needs robust 3D perception plus on-the-fly geometric and radiometric compensation, as shown in Raskar et al.'s Shader Lamps [83], which calibrates projection to an object's shape and color on-the-spot.

Towards volumetric F-formation: robot–learner–object proxemics in 3D. Projection placement can be further optimized by analyzing the robot–learner–object proxemics. Our work used a pragmatic guideline to define projection placements: an F-formation with projected visuals and referents in O-space, surrounded by the robot and learners in R-space [69]. As shown in Figure 3 ①, the standard F-formation diagram is a preliminary 2D simplification [69]; to extend this, a more sophisticated 3D F-formation model (see Figure 3 ②) can be established to view O- and R-space as volumes and take learners' visual field, displays' viewing cone, and physical occlusions into consideration, as can be supported by spatial reconstruction (e.g., RGB-D mapping/SLAM [76]).

Such a volumetric F-formation modeling could help a robotic TA properly decide when to utilize screen-based display and when to utilize in-situ projection. For example, when the human-robot formation is face-to-face dialogue or when artifacts are not suitable to project on, screens may be favored. Our ProjecTA hardware (see Figure 5) was also built to incorporate both the chest-screen for frontal interaction and the in-situ projection for object-focused guidance. In future real-world implementation, instead of treating projection and screens as competing options, designers can integrate both according to the 3D formation of learners and objects, to maximize the benefits of both.

7.1.3 Implications 3: Extending In-situ Projection as Human-Robot Collaborative Interfaces for Nomadic Learning. Different from whiteboards or fixed projectors, a robotic TA can carry its own projection. Our work illustrated how this could augment nomadic, walk-and-talk learning by visually supplementing the robot's verbalized content (**DO3**) and revealing critical or otherwise unseen information of the referent (**DO4**). These position in-situ projection as a new communicative medium for robotic TAs and learners, pointing its evolution into learner-robot collaborative interfaces.

Visualizing robots' chain of thought for learners' sensemaking. As shown in our results, in addition to equipment explanations, a few participants also wanted projected cues about ProjecTA's movement path and action plan, suggesting that enabling learners to predict a robotic TA's trajectories and status may increase the robot's social competence. Although our current work mainly focused on conveying equipment explanations, in future collaborative learning among robots and learners, in-situ projection could also externalize a robot's intent, plan, and intermediate reasoning to improve transparency and explainability [94]. For instance, Wengefeld et al. used laser projection to signify robot intent [111], and Mirror Eyes displayed mirror reflections on the robot's eyes to convey its focus of interest to humans [57]. As educational research shows exposing AI reasoning can aid learning [14], future work could use spatial projection to depict robots' reasoning about and in the physical environments (e.g., spatial problem-solving) to support learners' sensemaking.

Projected interface for learners' input. While our study primarily used in-situ projection as information display, projection can also function as an input interface, enabling spatial, embodied collaboration among learners and robotic

TAs. Drawing from prior work, for instance, in material- and tool-based skill training, learners could trigger robots' support or feedback via on-projection inputs such as taps, traces, or region selection [42, 112]. While prior studies already showed the value of a separate projector in turning its surrounding areas into interactive surfaces (e.g., [66]), our findings suggest future opportunities of such interactive surfaces being carried around by humanoid robots, serving as pervasive and environment-adaptive user interfaces to manipulate or collaborate with robots.

7.1.4 Implications 4: Fine-Graining Multimodal Orchestration for Projection-enhanced Robotic TAs. Our results show that the visuals, gestures, and speech were perceived to complement one and another significantly better in the projection-based condition than the screen-based counterpart. Qualitative evidences further illustrate examples of such crossmodal complementarity preferred by the learners. Taken together, we outline the following design possibilities.

Pre-action speech cues to better prepare learners for upcoming gestures and visuals. A few participants reported that projecTA's gestures and projections were occasionally unanticipated, primarily during their initial engagement with the system. A brief spoken cue in advance, such as "Please see illustration above," prepares learners for robots' upcoming gestures or visuals. Such pre-action speech cues are especially helpful for collaborative tasks requiring frequent repositioning or viewpoint shifts, because gestures and projections are limited to each learner's field-of-view, whereas speech can broadcast across the space, prompting learners to anticipate upcoming visual communication. Furthermore, as shown in our results, robotic TAs should pair such pre-action cues with brief pauses that not only give learners time to prepare, but also enable them to interrupt the robot, either requesting for more details, verifying what the robot has just explained, or changing the robot's upcoming actions. This helps learners to know when they can act without missing key information, and intervene the robot when needed.

More precise temporal alignment in multimodal orchestration. While our system's temporal alignment technique proved effective, our results point toward the need for a more granular approach. Future work should pursue a finer-grained alignment by breaking down each modality into smaller, timestamped units. Further, individual keywords within speech segments, as well as keyframes within projected animations or videos could all be assigned with timestamps, to enable keyword-to-keyframe level visual-speech coordination. Similarly, gestural units could be decomposed into timestamped micro-steps (e.g., 'raise hand,' then 'cover eyes,' then 'lower hand,' instead of a single 'cover-eyes' action). This aligns with the idea of 'atomic actions' in LaMI [108]. By aligning these multimodal micro-units based on semantic and contextual reasoning, we can achieve smoother and more expressive presentation of robots without sacrificing action transparency.

Additionally, our qualitative findings also revealed a few moments where concurrent modalities competed for learners' attention. This suggests that multimodal orchestration should not only synchronize channels, but also designate at each moment which modality leads and which ones recede into a supporting role. For instance, when an iconic gesture is intended to carry the main semantic load, projected content could briefly dim or simplify, or prompt learners to look back at the robot. Conversely, when detailed projected information should become the primary focus, the robot might momentarily pause unnecessary body movements and point toward the display. Such deliberate role-switching between lead and supporting modalities can reduce unnecessary modality conflicts and leverage the complementary strengths of speech, gesture, and projection.

7.2 Limitations and Future Work

To enable a controlled comparison, both ProjecTA and the Baseline delivered pre-choreographed explanations and did not support in-tour Q&A or dialogue. Future implementation could introduce mixed-initiative interaction for

more natural collaboration. The Presentation Choreography Workflow is promising to evolve into an educator-facing authoring tool for easy creation, customization, and fine-tuning. ProjecTA is an exploratory prototype; occasional image jitter was observed during motion. Future iteration could add a stabilizer or motion-compensated gimbal to the projector mount. Although the substantial reduction in extraneous cognitive load demonstrated ProjecTA's promise for nomadic learning, immediate quiz scores remained comparable across conditions, suggesting the need for more comprehensive assessments of learning outcomes in future research. Due to our focus on university makerspace, our sample consisted of novice makers, mostly from engineering or science backgrounds. Future work can extend this line of inquiry to similarly complex and varied nomadic learning contexts beyond makerspaces (e.g., museums, botanical gardens, and even outdoor sites). It can also engage larger and more varied learner cohorts. These extensions will help better understand and generalize the impacts of similar systems on both learning performance and experience. To focus on non-face-to-face, nomadic learning, the robot did not use facial expressions; in future, projecting facial cues onto the scene may provide social signals without prompting head turns. Finally, this study compared projection and screen presentation separately; combining body-mounted screens with in-situ projection may yield additive effects worth future evaluating.

8 Conclusion

In this study, we empirically examine how a robotic TA with in-situ projection, compared to a screen-based counterpart, affect learners' experiences during makerspace tours. We implemented ProjecTA and Baseline that differed only in display modality. We evaluated them in a real university makerspace with 24 novices across two rounds covering six pieces of equipment. Results show that ProjectA significantly lowered learner's extraneous cognitive load and was rated significantly higher in perceived usability relative to the Baseline. Participants consistently reported that near-object visual overlays reduced attention switching, facilitated easier mapping from visuals to physical targets, and enhanced the multimodal integration of projection, gestures, and speech. Based on these findings, we distill design implications to inform future work in better supporting nomadic learning in physical settings with robotic TAs equipped with in-situ projection.

Acknowledgments

We appreciate all participating domain experts for their contributions to the two co-design activities, with special thanks to Prof. Fang Wan. Our thanks also go to all participants for their time and efforts throughout the process and the reviewers for their invaluable comments that helped us significantly improve this paper. This work is supported by the SUSTech Grant for AI: R01656020.

References

- [1] Jong-gil Ahn, Hyeonsuk Yang, Gerard J Kim, Namgyu Kim, Kyoung Choi, Hyemin Yeon, Eunja Hyun, Miheon Jo, and Jeonghye Han. 2011. Projector Robot for Augmented Children's Play. In *Proceedings of the 6th International Conference on Human-Robot Interaction (HRI '11)* (Lausanne, Switzerland). Association for Computing Machinery, New York, NY, USA, 27–28. [doi:10.1145/1957656.1957666](https://doi.org/10.1145/1957656.1957666)
- [2] Samer Al Moubayed, Jonas Beskow, Gabriel Skantze, and Björn Granström. 2012. Furhat: a back-projected human-like robot head for multiparty human-machine interaction. In *Cognitive behavioural systems: COST 2102 international training school, dresden, Germany, february 21-26, 2011, revised selected papers*. Springer, Berlin, Heidelberg, 114–130.
- [3] Pengcheng An, Kenneth Holstein, Bernice d'Anjou, Berry Eggen, and Saskia Bakker. 2020. The TA Framework: Designing Real-time Teaching Augmentation for K-12 Classrooms. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)* (Honolulu, HI, USA). Association for Computing Machinery, New York, NY, USA, 1–17. [doi:10.1145/3313831.3376277](https://doi.org/10.1145/3313831.3376277)

- [4] Doris Aschenbrenner, Florian Leutert, Argun Çençen, Jouke Verlinden, Klaus Schilling, Marc Latoschik, and Stephan Lukosch. 2019. Comparing human factors for augmented reality supported single-user and collaborative repair operations of industrial robots. *Frontiers in Robotics and AI* 6 (2019), 37.
- [5] Andreas Baechler, Liane Baechler, Sven Autenrieth, Peter Kurtz, Thomas Hoerz, Thomas Heidenreich, and Georg Kruell. 2016. A Comparative Study of an Assistance System for Manual Order Picking—Called Pick-by-Projection—with the Guiding Systems Pick-by-Paper, Pick-by-Light and Pick-by-Display. In *2016 49th Hawaii International Conference on System Sciences (HICSS)* (Koloa, HI, USA). IEEE, Piscataway, NJ, USA, 523–531. doi:[10.1109/HICSS.2016.72](https://doi.org/10.1109/HICSS.2016.72)
- [6] Manuela Barbara, Sarah Pulé, and Lawrence Farrugia. 2024. Using Makerspaces to enrich Design and Technology education. *Techne Serien - Forskning i slöjdpedagogik och slöjdvetenskap* 31, 3 (2024), 33–50. doi:[10.7577/TechneA.5835](https://doi.org/10.7577/TechneA.5835)
- [7] Claudia Bartels, Martin Wegrzyn, Anne Wiedl, Verena Ackermann, and Hannelore Ehrenreich. 2010. Practice effects in healthy adults: a longitudinal study on frequent repetitive cognitive testing. *BMC neuroscience* 11, 1 (2010), 118.
- [8] Rebekka Bärthele, Stefan Sauer, and Jan Wilkenning. 2023. Exploring the potential of an Augmented Reality sandbox for geovisualization. In *Proceedings of the ICA*, Vol. 5. Copernicus Publications, Göttingen, Germany, 1. doi:[10.5194/ica-proc-5-1-2023](https://doi.org/10.5194/ica-proc-5-1-2023)
- [9] Paul Baxter, Emily Ashurst, Robin Read, James Kennedy, and Tony Belpaeme. 2017. Robot education peers in a situated primary school study: Personalisation promotes child learning. *PLoS one* 12, 5 (2017), e0178126.
- [10] Francisco Bellas, Martin Naya-Varela, Alma Mallo, and Alejandro Paz-Lopez. 2024. Education in the AI era: a long-term classroom technology based on intelligent robotics. *Humanities and Social Sciences Communications* 11, 1 (2024), 1–20.
- [11] Tony Belpaeme, James Kennedy, Aditi Ramachandran, Brian Scassellati, and Fumihide Tanaka. 2018. Social robots for education: A review. *Science robotics* 3, 21 (2018), eaat5954.
- [12] Oliver Bimber and Ramesh Raskar. 2005. *Spatial Augmented Reality: Merging Real and Virtual Worlds*. CRC Press, Wellesley, MA, USA.
- [13] Ann Blandford, Dominic Furniss, and Stephann Makri. 2016. *Qualitative HCI research: Going behind the scenes*. Morgan & Claypool Publishers, San Rafael, CA, USA.
- [14] Andrea Blasco and Vicky Charisi. 2024. *AI Chatbots in K-12 Education: An Experimental Study of Socratic vs. Non-Socratic Approaches and the Role of Step-by-Step Reasoning*. Working Paper 5040921. SSRN. doi:[10.2139/ssrn.5040921](https://doi.org/10.2139/ssrn.5040921) Posted Dec 2, 2024; last revised Nov 26, 2025.
- [15] Donna Z Bliss, Adam J Becker, Olga V Gurvich, Cynthia S Bradley, Erica Timko Olson, Mary T Steffes, Carol Flaten, Scott Jameson, and John P Condon. 2022. Projected augmented reality (P-AR) for enhancing nursing education about pressure injury: A pilot evaluation study. *Journal of Wound Ostomy & Continence Nursing* 49, 2 (2022), 128–136.
- [16] Florian Block, Michael S Horn, Brenda Caldwell Phillips, Judy Diamond, E Margaret Evans, and Chia Shen. 2012. The deeptree exhibit: Visualizing the tree of life to facilitate informal learning. *IEEE Transactions on Visualization and Computer Graphics* 18, 12 (2012), 2789–2798.
- [17] Wolfram Burgard, Armin B Cremers, Dieter Fox, Dirk Hähnel, Gerhard Lakemeyer, Dirk Schulz, Walter Steiner, and Sebastian Thrun. 1999. Experiences with an interactive museum tour-guide robot. *Artificial intelligence* 114, 1-2 (1999), 3–55.
- [18] Sebastian Büttner, Michael Prilla, and Carsten Röcker. 2020. Augmented Reality Training for Industrial Assembly Work—Are Projection-based AR Assistive Systems an Appropriate Tool for Assembly Training?. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)* (Honolulu, HI, USA). Association for Computing Machinery, New York, NY, USA, 1–12. doi:[10.1145/3313831.3376720](https://doi.org/10.1145/3313831.3376720)
- [19] Jane E Caldwell. 2007. Clickers in the large classroom: Current research and best-practice tips. *CBE—Life Sciences Education* 6, 1 (2007), 9–20.
- [20] Ravi Teja Chadalavada, Henrik Andreasson, Robert Krug, and Achim J. Lilienthal. 2015. That's on my mind! Robot to Human Intention Communication through On-board Projection on Shared Floor Space. In *2015 European Conference on Mobile Robots (ECMR)* (Lincoln, United Kingdom). IEEE, Piscataway, NJ, USA, 1–6. doi:[10.1109/ECMR.2015.7403771](https://doi.org/10.1109/ECMR.2015.7403771)
- [21] Tathagata Chakraborti, Sarath Sreedharan, Anagha Kulkarni, and Subbarao Kambhampati. 2018. Projection-aware Task Planning and Execution for Human-in-the-Loop Operation of Robots in a Mixed-Reality Workspace. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (Madrid, Spain). IEEE, Piscataway, NJ, USA, 4476–4482. doi:[10.1109/IROS.2018.8593830](https://doi.org/10.1109/IROS.2018.8593830)
- [22] Paul Chandler and John Sweller. 1991. Cognitive Load Theory and the Format of Instruction. *Cognition and Instruction* 8, 4 (1991), 293–332. doi:[10.1207/s1532690xci0804_2](https://doi.org/10.1207/s1532690xci0804_2)
- [23] Paul Chandler and John Sweller. 1992. The split-attention effect as a factor in the design of instruction. *British Journal of Educational Psychology* 62, 2 (1992), 233–246.
- [24] Michelene TH Chi, Nicholas De Leeuw, Mei-Hung Chiu, and Christian LaVancher. 1994. Eliciting self-explanations improves understanding. *Cognitive science* 18, 3 (1994), 439–477.
- [25] Herbert H. Clark and Susan E. Brennan. 1991. Grounding in Communication. In *Perspectives on Socially Shared Cognition*, Lauren B. Resnick, John M. Levine, and Stephanie D. Teasley (Eds.). American Psychological Association, Washington, DC, USA, 127–149. doi:[10.1037/10096-006](https://doi.org/10.1037/10096-006)
- [26] Herbert H Clark and Meredith A Krych. 2004. Speaking while monitoring addressees for understanding. *Journal of memory and language* 50, 1 (2004), 62–81.
- [27] Susan Wagner Cook, Zachary Mitchell, and Susan Goldin-Meadow. 2008. Gesturing makes learning last. *Cognition* 106, 2 (2008), 1047–1058.
- [28] Michael D Coover, Tiffany Lee, Ivan Shindrev, and Yu Sun. 2014. Spatial augmented reality as a method for a mobile robot to communicate intended movement. *Computers in Human Behavior* 34 (2014), 241–248.
- [29] Rajkumar Darbar, Joan Sol Roo, Thibault Lainé, and Martin Hachet. 2019. DroneSAR: Extending Physical Spaces in Spatial Augmented Reality Using Projection on a Drone. In *Proceedings of the 18th International Conference on Mobile and Ubiquitous Multimedia (MUM '19)*.

- Association for Computing Machinery, New York, NY, USA, Article 4, 7 pages. doi:[10.1145/3365610.3365631](https://doi.org/10.1145/3365610.3365631)
- [30] Jordan Aiko Deja, Sandi Štor, Ilonka Pucihar, Mahesha Weerasinghe, Rafael Marco Balbin, Klen Čopić Pucihar, and Matjaž Kljun. 2025. ImproVisAR: designing augmented reality piano roll for teaching improvisation. *Virtual Reality* 29, 3 (2025), 140.
- [31] Ö. Ece Demir-Lira, Junko Kanero, Cansu Oranç, Sümeyye Koşkulu, Idil Franko, Tilbe Göksun, and Aylin C. Küntay. 2020. L2 Vocabulary Teaching by Social Robots: The Role of Gestures and On-Screen Cues as Scaffolds. *Frontiers in Education* 5 (2020), 599636. doi:[10.3389/feduc.2020.599636](https://doi.org/10.3389/feduc.2020.599636)
- [32] Ahmed Elsharkawy, Khawar Naheem, Dongwoo Koo, and Mun Sang Kim. 2021. A UWB-driven self-actuated projector platform for interactive augmented reality applications. *Applied Sciences* 11, 6 (2021), 2871.
- [33] John M Ford. 2004. Content analysis: An introduction to its methodology. *Personnel psychology* 57, 4 (2004), 1110.
- [34] Markus Funk, Thomas Kosch, and Albrecht Schmidt. 2016. Interactive Worker Assistance: Comparing the Effects of In-Situ Projection, Head-Mounted Displays, Tablet, and Paper Instructions. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Heidelberg, Germany) (*UbiComp '16*). Association for Computing Machinery, New York, NY, USA, 934–939. doi:[10.1145/2971648.2971706](https://doi.org/10.1145/2971648.2971706)
- [35] Ramsundar Kalpagam Ganesh, Yash K Rathore, Heather M Ross, and Heni Ben Amor. 2018. Better teaming through visual cues: how projecting imagery in a workspace can improve human-robot collaboration. *IEEE Robotics & Automation Magazine* 25, 2 (2018), 59–71.
- [36] Yuan Gao, Yuyun Zhao, Le Xie, and Guoyan Zheng. 2021. A projector-based augmented reality navigation system for computer-assisted surgery. *Sensors* 21, 9 (2021), 2931.
- [37] Yate Ge, Meiyi Li, Xipeng Huang, Yuanda Hu, Qi Wang, Xiaohua Sun, and Weiwei Guo. 2025. GenComUI: Exploring Generative Visual Aids as Medium to Support Task-Oriented Human-Robot Communication. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (*CHI '25*). Association for Computing Machinery, New York, NY, USA, Article 433, 21 pages. doi:[10.1145/3706598.3714238](https://doi.org/10.1145/3706598.3714238)
- [38] Raphaela Gehle, Karola Pitsch, Timo Dankert, and Sebastian Wrede. 2017. How to Open an Interaction Between Robot and Museum Visitor? Strategies to Establish a Focused Encounter in HRI. In *Proceedings of the 2017 ACM/IEEE International Conference on Human-Robot Interaction* (Vienna, Austria) (*HRI '17*). Association for Computing Machinery, New York, NY, USA, 187–195. doi:[10.1145/2909824.3020219](https://doi.org/10.1145/2909824.3020219)
- [39] Kyungwon Gil, Jimin Rhim, Taejin Ha, Young Yim Doh, and Woontack Woo. 2014. AR Petite Theater: Augmented Reality Storybook for Supporting Children's Empathy Behavior. In *2014 IEEE International Symposium on Mixed and Augmented Reality - Media, Art, Social Science, Humanities and Design (ISMAR-MASH'D)* (Munich, Germany). IEEE, Piscataway, NJ, USA, 13–20. doi:[10.1109/ISMAR-AMH.2014.6935433](https://doi.org/10.1109/ISMAR-AMH.2014.6935433)
- [40] Xiaoshan Zhu Gordy, Ellen M Jones, and Jessica H Bailey. 2018. Technological innovation or educational evolution? A multi-disciplinary qualitative inquiry into active learning classrooms. *Journal of the Scholarship of Teaching and Learning* 18, 2 (2018), 1–23.
- [41] Zhao Han, Jenna Parrillo, Alexander Wilkinson, Holly A. Yanco, and Tom Williams. 2022. Projecting Robot Navigation Paths: Hardware and Software for Projected AR. In *2022 17th ACM/IEEE International Conference on Human-Robot Interaction (HRI)* (Sapporo, Japan). Association for Computing Machinery, New York, NY, USA, 623–628. doi:[10.1109/HRI53351.2022.9889537](https://doi.org/10.1109/HRI53351.2022.9889537)
- [42] Chris Harrison, Hrvoje Benko, and Andrew D. Wilson. 2011. OmniTouch: Wearable Multitouch Interaction Everywhere. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology* (Santa Barbara, CA, USA) (*UIST '11*). Association for Computing Machinery, New York, NY, USA, 441–450. doi:[10.1145/2047196.2047255](https://doi.org/10.1145/2047196.2047255)
- [43] Jeremy Hartmann, Yen-Ting Yeh, and Daniel Vogel. 2020. AAR: Augmenting a Wearable Augmented Reality Display with an Actuated Head-Mounted Projector. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology* (Virtual Event, USA) (*UIST '20*). Association for Computing Machinery, New York, NY, USA, 445–458. doi:[10.1145/3379337.3415849](https://doi.org/10.1145/3379337.3415849)
- [44] Mohammad Nehal Hasnine, Bipin Indurkhyia, and Mahmoud Mohamed Hussien Ahmed. 2024. Socially Assistive Robot as Laboratory Safety Assistant for Science Students. In *2024 IEEE International Conference on Advanced Robotics and Its Social Impacts (ARSO)* (Hong Kong). IEEE, Piscataway, NJ, USA, 37–42. doi:[10.1109/ARSO60199.2024.10557826](https://doi.org/10.1109/ARSO60199.2024.10557826)
- [45] Sara Hennessy. 2011. The role of digital artefacts on the interactive whiteboard in supporting classroom dialogue. *Journal of computer assisted learning* 27, 6 (2011), 463–489.
- [46] Kenneth Holstein, Bruce M. McLaren, and Vincent Aleven. 2018. Student Learning Benefits of a Mixed-Reality Teacher Awareness Tool in AI-Enhanced Classrooms. In *Artificial Intelligence in Education (Lecture Notes in Computer Science, Vol. 10947)*, Carolyn Penstein Rosé, Roberto Martínez-Maldonado, H. Ulrich Hoppe, Rose Luckin, Manolis Mavrikis, Kaska Porayska-Pomsta, Bruce McLaren, and Benedict du Boulay (Eds.). Springer International Publishing, Cham, Switzerland, 154–168. doi:[10.1007/978-3-319-93843-1_12](https://doi.org/10.1007/978-3-319-93843-1_12)
- [47] Autumn B Hostetter. 2011. When do gestures communicate? A meta-analysis. *Psychological bulletin* 137, 2 (2011), 297.
- [48] Yuhui Hu, Peide Huang, Mouli Sivapupurapu, and Jian Zhang. 2025. ELEGNT: Expressive and Functional Movement Design for Non-anthropomorphic Robot. arXiv:2501.12493 [cs.RO] doi:[10.48550/arXiv.2501.12493](https://doi.org/10.48550/arXiv.2501.12493) arXiv preprint, 13 pages.
- [49] Yixin Hu, Anjun Zhu, Catalina L. Toma, and Bilge Mutlu. 2025. Designing Telepresence Robots to Support Place Attachment. In *2025 20th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, Piscataway, NJ, USA, 252–261.
- [50] Annie Huang, Alyson Ranucci, Adam Stogsdill, Grace Clark, Keenan Schott, Mark Higger, Zhao Han, and Tom Williams. 2024. (Gestures Vaguely): The Effects of Robots' Use of Abstract Pointing Gestures in Large-Scale Environments. In *Proceedings of the 2024 ACM/IEEE International Conference on Human-Robot Interaction* (Boulder, CO, USA) (*HRI '24*). Association for Computing Machinery, New York, NY, USA, 293–302. doi:[10.1145/3610977.3634975](https://doi.org/10.1145/3610977.3634975)
- [51] Takamasa Iio, Satoru Satake, Takayuki Kanda, Kotaro Hayashi, Florent Ferreri, and Norihiro Hagita. 2020. Human-like guide robot that proactively explains exhibits. *International Journal of Social Robotics* 12, 2 (2020), 549–566.

- [52] Muhammed Ismael, Roderick McCall, Fintan McGee, Ilyasse Belkacem, Mickaël Stefas, Joan Baixauli, and Didier Arl. 2024. Acceptance of augmented reality for laboratory safety training: methodology and an evaluation study. *Frontiers in Virtual Reality* 5 (2024), 1322543.
- [53] Jiaqi Jiang, Kexin Huang, Hanqing Zhou, Huiying Lu, and Pengcheng An. 2025. Visiobo Demo: Augmenting Static Prints with Projection-based Visual Cueing and Concept Mapping via LLM Reasoning. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–5.
- [54] Misaki Kasetani, Tomonobu Noguchi, Hirotake Yamazoe, and Joo-Ho Lee. 2015. Projection Mapping by Mobile Projector Robot. In *2015 12th International Conference on Ubiquitous Robots and Ambient Intelligence (URAI)*. IEEE, Piscataway, NJ, USA, 13–17.
- [55] SeungJun Kim and Anind K. Dey. 2009. Simulated Augmented Reality Windshield Display as a Cognitive Mapping Aid for Elder Driver Navigation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 133–142.
- [56] Amy Koike, Bengisu Cagiltay, and Bilge Mutlu. 2024. Tangible Scenography as a Holistic Design Method for Human-Robot Interaction. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*. Association for Computing Machinery, New York, NY, USA, 459–475.
- [57] Matti Krüger, Daniel Tanneberg, Chao Wang, Stephan Hasler, and Michael Gienger. 2025. Mirror Eyes: Explainable Human-Robot Interaction at a Glance. arXiv:2506.18466 [cs.RO] doi:10.48550/arXiv.2506.18466 arXiv preprint; accepted to IEEE RO-MAN 2025 (related DOI: 10.1109/ROMAN63969.2025.11217810).
- [58] Jimmie Leppink, Fred Paas, Cees PM Van der Vleuten, Tamara Van Gog, and Jeroen JG Van Merriënboer. 2013. Development of an instrument for measuring different types of cognitive load. *Behavior research methods* 45, 4 (2013), 1058–1072.
- [59] Jimmie Leppink and Angelique Van den Heuvel. 2015. The evolution of cognitive load theory and its application to medical education. *Perspectives on medical education* 4, 3 (2015), 119–127.
- [60] Jan Leusmann, Anna Belardinelli, Luke Haliburton, Stephan Hasler, Albrecht Schmidt, Sven Mayer, Michael Gienger, and Chao Wang. 2025. Investigating LLM-Driven Curiosity in Human-Robot Interaction. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI ’25). Association for Computing Machinery, New York, NY, USA, Article 599, 16 pages. doi:10.1145/3706598.3713923
- [61] Natan Linder and Pattie Maes. 2010. LuminAR: Portable Robotic Augmented Reality Interface Design and Prototype. In *Adjunct Proceedings of the 23rd Annual ACM Symposium on User Interface Software and Technology* (New York, NY, USA) (UIST ’10 Adjunct). Association for Computing Machinery, New York, NY, USA, 395–396. doi:10.1145/1866218.1866237
- [62] Ragavendra Lingamaneni, Thomas Kubitz, and Jürgen Scheible. 2017. DroneCAST: Towards a Programming Toolkit for Airborne Multimedia Display Applications. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services (MobileHCI ’17)*. Association for Computing Machinery, New York, NY, USA, Article 85, 8 pages. doi:10.1145/3098279.3122128
- [63] Tyler S. Love. 2018. Perceptions of Safety in Makerspaces: Examining the Influence of Professional Development. In *Proceedings of the 105th Mississippi Valley Technology Teacher Education Conference* (2018-11-15). Mississippi Valley Technology Teacher Education Conference, Nashville, TN, USA, 1–16. <https://www.mississippivalley.org/wp-content/uploads/2018/11/Perceptions-of-Safety-in-Makerspaces-Love.pdf> Session III: Research, Laboratories, and New Initiatives.
- [64] Adrian Lozada, Uthman Tijani, Villa Keth, Hong Wang, and Zhao Han. 2025. Anywhere Projected AR for Robot Communication: A Mid-Air Fog Screen-Robot System. In *2025 20th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, Piscataway, NJ, USA, 520–528.
- [65] Andrés Lucero. 2015. Using Affinity Diagrams to Evaluate Interactive Prototypes. In *Human-Computer Interaction – INTERACT 2015*. Springer, Cham, Switzerland, 231–248.
- [66] Thomas Ludwig, Michael Döll, and Christoph Kotthaus. 2019. “The Printer is Telling Me about Itself”: Supporting the Appropriation of Hardware by Using Projection Mapping. In *Proceedings of the 2019 on Designing Interactive Systems Conference (DIS ’19)*. Association for Computing Machinery, New York, NY, USA, 331–344.
- [67] Paul Lukowicz, Andreas Poxrucker, Jens Weppner, Benjamin Bischke, Jochen Kuhn, and Michael Hirth. 2015. Glass-Physics: Using Google Glass to Support High School Physics Experiments. In *Proceedings of the 2015 ACM International Symposium on Wearable Computers*. Association for Computing Machinery, New York, NY, USA, 151–154.
- [68] Michael Lundberg and Jay Rasmussen. 2018. Foundational Principles and Practices to Consider in Assessing Maker Education. *Journal of Educational Technology* 14, 4 (2018), 1–12.
- [69] Nicolai Marquardt, Ken Hinckley, and Saul Greenberg. 2012. Cross-Device Interaction via Micro-Mobility and F-Formations. In *Proceedings of the 25th Annual ACM Symposium on User Interface Software and Technology*. Association for Computing Machinery, New York, NY, USA, 13–22.
- [70] David McNeill. 1992. *Hand and Mind: What Gestures Reveal about Thought*. University of Chicago Press, Chicago, IL, USA.
- [71] Neil Mercer, Sara Hennessy, and Paul Warwick. 2010. Using interactive whiteboards to orchestrate classroom dialogue. *Technology, Pedagogy and Education* 19, 2 (2010), 195–209.
- [72] Joseph E. Michaelis and Daniela Di Canio. 2022. Embodied Geometric Reasoning with a Robot: The Impact of Robot Gestures on Student Reasoning about Geometrical Conjectures. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–14.
- [73] Neema Moraveji, Meredith Morris, Daniel Morris, Mary Czerwinski, and Nathalie Henry Riche. 2011. ClassSearch: Facilitating the Development of Web Search Skills through Social Learning. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1797–1806.

- [74] Omar Mubin, Catherine J Stevens, Suleman Shahid, Abdullah Al Mahmud, and Jian-Jie Dong. 2013. A review of the applicability of robots in education. *Journal of Technology in Education and Learning* 1, 209-0015 (2013), 13.
- [75] N Hari Narayanan and Mary Hegarty. 2002. Multimedia design for communication of dynamic information. *International journal of human-computer studies* 57, 4 (2002), 279–315.
- [76] Richard A. Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. 2011. KinectFusion: Real-Time Dense Surface Mapping and Tracking. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, Piscataway, NJ, USA, 127–136.
- [77] Vasiliki Nikolakopoulou, Petros Printezis, Vassilis Maniatis, Dimitris Kontzas, Spyros Vosinakis, Pavlos Chatzigrigoriou, and Panayiotis Koutsabasis. 2022. Conveying intangible cultural heritage in museums with interactive storytelling and projection mapping: the case of the mastic villages. *Heritage* 5, 2 (2022), 1024–1049.
- [78] Antti Oulasvirta, Esko Kurvinen, and Tomi Kankainen. 2003. Understanding contexts by being there: case studies in bodystorming. *Personal and ubiquitous computing* 7, 2 (2003), 125–134.
- [79] Heather L O'Brien, Paul Cairns, and Mark Hall. 2018. A practical approach to measuring user engagement with the refined user engagement scale (UES) and new UES short form. *International Journal of Human-Computer Studies* 112 (2018), 28–39.
- [80] Amit Kumar Pandey and Rodolphe Gelin. 2018. A mass-produced sociable humanoid robot: Pepper: The first machine of its kind. *IEEE Robotics & Automation Magazine* 25, 3 (2018), 40–48.
- [81] Iulian Radu. 2014. Augmented reality in education: a meta-review and cross-media analysis. *Personal and ubiquitous computing* 18, 6 (2014), 1533–1543.
- [82] Aditi Ramachandran, Sarah Strohkorb Sebo, and Brian Scassellati. 2019. Personalized Robot Tutoring Using the Assistive Tutor POMDP (AT-POMDP). In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. AAAI Press, Palo Alto, CA, USA, 8050–8057.
- [83] Ramesh Raskar, Greg Welch, Kok-Lim Low, and Deepak Bandyopadhyay. 2001. Shader Lamps: Animating Real Objects with Image-Based Illumination. In *Rendering Techniques 2001: Proceedings of the Eurographics Workshop on Rendering*. Springer, Vienna, Austria, 89–102.
- [84] Umair Rehman and Shi Cao. 2020. Comparative evaluation of augmented reality-based assistance for procedural tasks: a simulated control room study. *Behaviour & Information Technology* 39, 11 (2020), 1225–1245.
- [85] Hyocheol Ro, Jung-Hyun Byun, Inhwan Kim, Yoon Jung Park, Kyuri Kim, and Tack-Don Han. 2019. Projection-based Augmented Reality Robot Prototype with Human-Awareness. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, Piscataway, NJ, USA, 598–599.
- [86] Hyocheol Ro, Inhwan Kim, JungHyun Byun, Yoonsik Yang, Yoon Jung Park, Seungho Chae, and Tackdon Han. 2018. PAMI: Projection Augmented Meeting Interface for Video Conferencing. In *Proceedings of the 26th ACM International Conference on Multimedia*. Association for Computing Machinery, New York, NY, USA, 1274–1277.
- [87] Violeta Rosanda and Andreja Istenic Starcic. 2019. The Robot in the Classroom: A Review of a Robot Role. In *Emerging Technologies for Education*. Springer, Cham, Switzerland, 347–357.
- [88] Eric Rosen, David Whitney, Elizabeth Phillips, Gary Chien, James Tompkin, George Konidaris, and Stefanie Tellex. 2019. Communicating Robot Arm Motion Intent through Mixed Reality Head-Mounted Displays. In *Robotics Research: The 18th International Symposium ISRR*. Springer, Cham, Switzerland, 301–316.
- [89] Rinat Rosenberg-Kima, Yaakov Koren, Maya Yachini, and Goren Gordon. 2019. Human-Robot-Collaboration (HRC): Social Robots as Teaching Assistants for Training Activities in Small Groups. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, Piscataway, NJ, USA, 522–523.
- [90] Rinat B Rosenberg-Kima, Yaakov Koren, and Goren Gordon. 2020. Robot-supported collaborative learning (RSCL): Social robots as teaching assistants for higher education small group facilitation. *Frontiers in Robotics and AI* 6 (2020), 148.
- [91] Stephanie Rosenthal, Shaun K. Kane, Jacob O. Wobbrock, and Daniel Avrahami. 2010. Augmenting On-Screen Instructions with Micro-Projected Guides: When It Works, and When It Fails. In *Proceedings of the 12th ACM International Conference on Ubiquitous Computing*. Association for Computing Machinery, New York, NY, USA, 203–212.
- [92] Allison Sauppé and Bilge Mutlu. 2014. Robot Deictics: How Gesture and Context Shape Referential Communication. In *Proceedings of the 2014 ACM/IEEE International Conference on Human-Robot Interaction*. Association for Computing Machinery, New York, NY, USA, 342–349.
- [93] Jürgen Scheible, Achim Hoth, Julian Saal, and Haifeng Su. 2013. Displaydrone: A Flying Robot Based Interactive Display. In *Proceedings of the 2nd ACM International Symposium on Pervasive Displays*. Association for Computing Machinery, New York, NY, USA, 49–54.
- [94] Svenja Y Schött, Rifat Mehreen Amin, and Andreas Butz. 2023. A literature survey of how to convey transparency in co-located human–robot interaction. *Multimodal Technologies and Interaction* 7, 3 (2023), 25.
- [95] Yasaman S. Sefidgar, Thomas Weng, Heather Harvey, Sarah Elliott, and Maya Cakmak. 2018. Robotist: Interactive Situated Tangible Robot Programming. In *Proceedings of the 2018 ACM Symposium on Spatial User Interaction*. Association for Computing Machinery, New York, NY, USA, 141–149.
- [96] Elena Segura, Laia Vidal, and Asreen Rostami. 2016. Bodystorming for movement-based interaction design. *Human Technology* 12, 2 (2016), 193–251.
- [97] Thomas Sievers. 2025. A Humanoid Social Robot as a Teaching Assistant in the Classroom. arXiv:2508.05646 [cs.HC] doi:10.48550/arXiv.2508.05646 arXiv preprint, 12 pages.

- [98] Paula AG Soneral and Sara A Wyse. 2017. A SCALE-UP mock-up: Comparison of student learning gains in high-and low-tech active-learning environments. *CBE—Life Sciences Education* 16, 1 (2017), ar12.
- [99] Gavin Sudrey, Adam Jacobson, and Belinda Ward. 2018. Enabling a Pepper Robot to provide Automated and Interactive Tours of a Robotics Laboratory. arXiv:1804.03288 [cs.RO] doi:10.48550/arXiv.1804.03288 arXiv preprint, 8 pages.
- [100] Ryo Suzuki, Adnan Karim, Tian Xia, Hooman Hedayati, and Nicolai Marquardt. 2022. Augmented Reality and Robotics: A Survey and Taxonomy for AR-Enhanced Human-Robot Interaction and Robotic Interfaces. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–33.
- [101] Ryo Suzuki, Clement Zheng, Yasuaki Kakehi, Tom Yeh, Ellen Yi-Luen Do, Mark D. Gross, and Daniel Leithinger. 2019. Shapebots: Shape-Changing Swarm Robots. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*. Association for Computing Machinery, New York, NY, USA, 493–505.
- [102] Aki Tamai, Tetsushi Ikeda, and Satoshi Iwaki. 2019. A Method for Guiding a Person Combining Robot Movement and Projection. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, Piscataway, NJ, USA, 1265–1270.
- [103] Fumihide Tanaka, Kyosuke Isshiki, Fumiki Takahashi, Manabu Uekusa, Rumiko Sei, and Kaname Hayashi. 2015. Pepper Learns Together with Children: Development of an Educational Application. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*. IEEE, Piscataway, NJ, USA, 270–275.
- [104] Arthur Tang, Charles Owen, Frank Biocca, and Weimin Mou. 2003. Comparative Effectiveness of Augmented Reality in Object Assembly. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 73–80.
- [105] Michael Tomasello and Michael Jeffrey Farrar. 1986. Joint Attention and Early Language. *Child Development* 57, 6 (1986), 1454–1463. doi:10.2307/1130423
- [106] Dishita Turakhia, Mark Parent, Tovi Grossman, Michael Glueck, and Ben Lafreniere. 2025. Investigating Augmented Reality for Adaptive Motor-Skill Training. In *Proceedings of the 51st Graphics Interface Conference (GI '25)* (Kelowna, BC, Canada). Association for Computing Machinery, New York, NY, USA, 1–10. https://www.benlafreniere.ca/assets/papers/Turakhia_GI2025_-_AR_Adaptive_Motor-Skill_Training.pdf
- [107] Erik van Alphen and Saskia Bakker. 2015. Lernanto: An Ambient Display to Support Differentiated Instruction. In *Proceedings of the 11th International Conference on Computer Supported Collaborative Learning (CSCL 2015)* (Gothenburg, Sweden). International Society of the Learning Sciences (ISLS), Madison, WI, USA, 759–760. Poster.
- [108] Chao Wang, Stephan Hasler, Daniel Tanneberg, Felix Ocker, Frank Joublin, Antonello Ceravola, Joerg Deigmoeller, and Michael Gienger. 2024. Lami: Large Language Models for Multi-Modal Human-Robot Interaction. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–10.
- [109] Chenyang Wang, Daniel C. Tozadore, Barbara Bruno, and Pierre Dillenbourg. 2024. Co-designing a Child-Robot Relational Norm Intervention to Regulate Children’s Handwriting Posture. In *Proceedings of the 23rd Annual ACM Interaction Design and Children Conference*. Association for Computing Machinery, New York, NY, USA, 934–939.
- [110] Atsushi Watanabe, Tetsushi Ikeda, Yoichi Morales, Kazuhiko Shinozawa, Takahiro Miyashita, and Norihiro Hagita. 2015. Communicating Robotic Navigational Intentions. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, Piscataway, NJ, USA, 5763–5769.
- [111] Tim Wengefeld, Dominik Höchemer, Benjamin Lewandowski, Mona Köhler, Manuel Beer, and Horst-Michael Gross. 2020. A Laser Projection System for Robot Intention Communication and Human Robot Interaction. In *2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*. IEEE, Piscataway, NJ, USA, 259–265.
- [112] Robert Xiao, Chris Harrison, and Scott E. Hudson. 2013. WorldKit: Rapid and Easy Creation of Ad-Hoc Interactive Applications on Everyday Surfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 879–888.
- [113] Kazuyoshi Yoshino and Shanjun Zhang. 2023. Teaching-Assistant Robot Tutoring Students in the Classroom. In *Proceedings of the 15th International Conference on Education Technology and Computers*. Association for Computing Machinery, New York, NY, USA, 113–119.
- [114] Mohammad A. Yousuf, Yoshinori Kobayashi, Yoshinori Kuno, Akiko Yamazaki, and Keiichi Yamazaki. 2013. How to Move towards Visitors: A Model for Museum Guide Robots to Initiate Conversation. In *2013 IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*. IEEE, Piscataway, NJ, USA, 587–592.
- [115] Yan Zhang, Tharaka Sachintha Ratnayake, Cherie Sew, Jarrod Knibbe, Jorge Goncalves, and Wafa Johal. 2025. Can You Pass That Tool?: Implications of Indirect Speech in Physical Human-Robot Collaboration. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–18.

A Semi-structured Interview Questions

- **Noticing and visibility.** During the session, to what extent did you notice information presented via the *projection* and the *screen*? In what situations, if any, was the content difficult to see (e.g., viewing angle, glare, distance)?

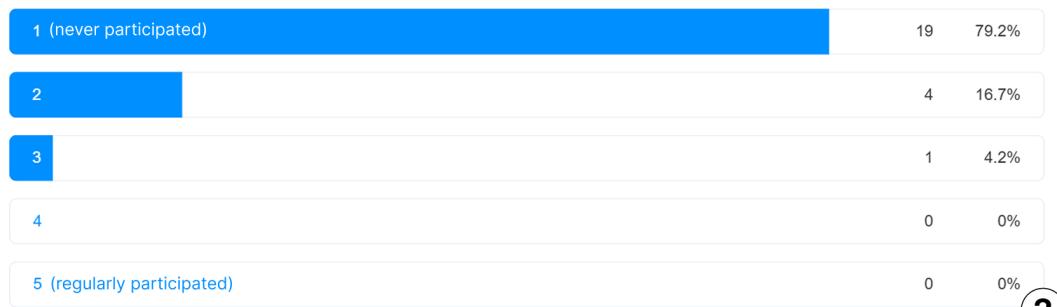
- **Experiences with both display modes.** Looking back on both rounds, how would you characterize your interaction with the robot under each mode? Which aspects worked well or left a strong impression, which did not, and what concrete improvements would you suggest?
- **Comparative evaluation.** Comparing *projection + speech/gestures* with *screen + speech/gestures*, how did the two modes differ in supporting your understanding and attention? Please illustrate with a specific moment (e.g., when gesture–projection coordination worked or failed—what gesture occurred and what was shown), and indicate which mode you would prefer and why.
- **Potential application scenarios.** In what other scenarios do you think this system would be valuable? How do you envision it operating there, and what benefits would it provide?

B Familiarity and Past Engagement Rating Scale Results

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

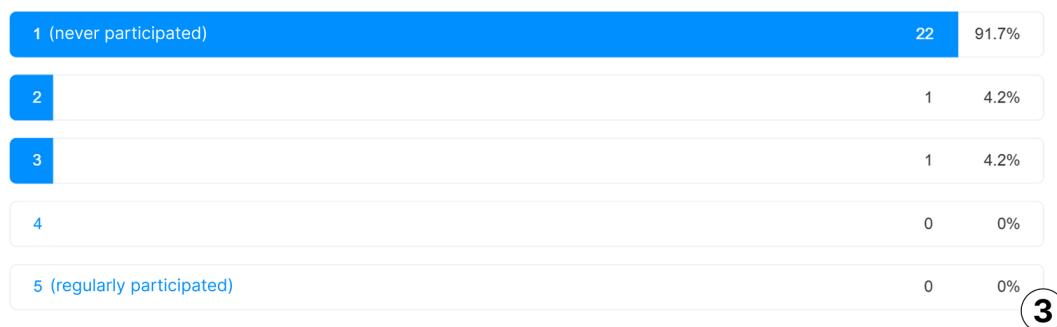
| Makerspace Equipment | 1 (not at all familiar) | 2 | 3 | 4 | 5 (Very familiar) |
|--|-------------------------|-----------|-----------|--------|-------------------|
| FDM 3D Printer (Bambu Lab X1E/P1S, etc.) | 19 (79.2%) | 3 (12.5%) | 2 (8.3%) | 0 (0%) | 0 (0%) |
| Resin 3D Printer (Form3, etc.) | 21 (87.5%) | 3 (12.5%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Nylon SLS 3D Printer (Fuse1/Sift) | 24 (100%) | 0 (0%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Soldering Station (Weller WSD81, etc.) | 20 (83.3%) | 2 (8.3%) | 2 (8.3%) | 0 (0%) | 0 (0%) |
| 3D Scanner (EinScan SP, etc.) | 17 (70.8%) | 7 (29.2%) | 0 (0%) | 0 (0%) | 0 (0%) |
| Laser Cutter (Trotec Speedy 400) | 17 (70.8%) | 4 (16.7%) | 3 (12.5%) | 0 (0%) | 0 (0%) |

1.3 Mean | 1 Median



1

1.1 Mean | 1 Median



2

3

Fig. 12. Demographic questionnaire results on ① participants' familiarity with makerspace equipment, ② past engagement in hands-on activities, and ③ prior experience with equipment-guided tours.