# Exploring the Effect of Visual Cues on Eye Gaze During AR-Guided Picking and Assembly Tasks

Arne Seeliger*
ETH Zurich

Gerrit Merz
Karlsruhe Institute of Technology

Christian Holz
ETH Zurich

Stefan Feuerriegel
ETH Zurich

## ABSTRACT

In this paper, we present an analysis of eye gaze patterns pertaining to visual cues in augmented reality (AR) for head-mounted displays (HMDs). We conducted an experimental study involving a picking and assembly task, which was guided by different visual cues. We compare these visual cues along multiple dimensions (in-view vs. out-of-view, static vs. dynamic, sequential vs. simultaneous) and analyze quantitative metrics such as gaze distribution, gaze duration, and gaze path distance. Our results indicate that visual cues in AR significantly affect eye gaze patterns. Specifically, we show that the effect varies depending on the type of visual cue. We discuss these empirical results with respect to visual attention theory.

**Index Terms:** Human-centered computing—Ubiquitous and mobile computing—Empirical studies in ubiquitous and mobile computing—

## 1 INTRODUCTION

Eye gaze analysis is frequently employed to better understand how users interact with computer interfaces. It has been extensively used in the context of two-dimensional interfaces like screens. Common applications are, for example, the analysis of eye gaze in web page interactions [10]. Eye gaze analysis has also been used in combination with subtle image modulations to direct visual attention [2].

Typically, evidence is obtained through desktop-based experiments. However, when studying eye gaze using mobile or wearable devices [22], a transfer of findings from desktop-based settings is difficult, especially when body movement and physical surroundings need to be accounted for [16]. Moreover, eye gaze behavior might differ substantially between static and dynamic settings [24].

Considering AR, specifically for HMDs, empirical assessment of eye gaze remains rare. Renner and Pfeiffer [30] compare AR-based attention guiding techniques in assembly tasks. However, they mainly employ evaluation metrics related to task performance (e. g., completion time). Similarly, Burova et al. [5] examine AR-based guidance and safety awareness cues. Despite calculating fixation counts, the work is mainly based on self-assessments (i. e., questionnaires). Moreover, the above studies use virtual reality (VR) devices to simulate AR environments. Yet, it remains open whether such findings are transferable to AR settings.

In this paper, we provide an analysis of eye gaze patterns relating to AR-based guidance cues in picking and assembly tasks. In one user study, 12 participants were guided by visual cues, which were displayed through an HMD. Specifically, eight visual cues were used, which differed along multiple dimensions (in-view vs. out-of-view, static vs. dynamic, sequential vs. simultaneous). We recorded eye gaze data for all participants and inferred quantitative metrics, including gaze distribution, gaze duration, and gaze path distance. Based on this, we discuss the empirical results in light of visual attention theory.

*e-mail: aseeliger@ethz.ch

## 2 RELATED WORK

### 2.1 Visual Attention and Search

At any given moment, the human visual system is exposed to more perceptual information than can be processed at the same time [7], which makes visual attention and search necessary [38]. Visual attention enables the active selection of relevant information and the ignoring of irrelevant information from a complex visual environment [7]. Although attention and eye movement are closely entangled [38], it is possible to fixate one location while attending another [23]. Therefore, it is common to distinguish between overt eye movement and covert deployment of attention [29]. Overt attention can be measured effectively by an eye tracker while the tracking of covert attention poses greater challenges [38].

Visual search characterizes a situation in which a subject looks for a target item among multiple distractor items [37]. Many theories regarding the mechanisms of search and search efficiency have been proposed (e. g., Treisman's feature integration theory (FIT) [35]). In this context, two forms of guidance are often distinguished, namely bottom-up, stimulus-driven attention and top-down, goal-driven attention [38]. In stimulus-driven attention, attention is attracted automatically through a salient stimulus, whereas goal-driven attention is under overt control of the human observer [7].

### 2.2 Attention Guidance in Augmented Reality

Attention guidance towards areas-of-interest (AOIs) can be achieved by directing the user's attention through the use of visual cues that indicate the target location [19]. Visual cues can either be presented at a fixed point on the display (in-view) or affixed to the focus object (in-situ) [30]. In this terminology, a conventional ON-SCREEN ARROW (e. g., [1]) is classified as in-view, whereas an arrow placed within the environment and pointing towards an object of interest is classified as in-situ. In case the object of interest lies outside a user's field-of-view (FOV), the visual cue first needs to guide the user's attention to the off-screen location.

There are various examples of visual cues. HALO is a well-known cue for visualizing off-screen locations [3]. For this, HALO shows circles around off-screen objects, where the circles are sufficiently large to reach into the border region of the display window. However, HALO has limited capacity when displaying locations of multiple objects in the same corner. This is addressed by WEDGE [18], where circles are replaced by isosceles triangles. The ATTENTION FUNNEL [4] uses a tunnel of frames that are drawn from the central FOV towards the location of the target object. Other cues include, for instance, DEADEYE [25]. Note that most of the aforementioned cues were originally developed for two-dimensional interfaces such as screens but they have also been found suitable for HMDs [13].

### 2.3 Eye Gaze Metrics

Common metrics for eye gaze analysis include the following. (1) Eye fixations point to individual locations of user attention, thus yielding gaze distributions. A larger number of fixations can indicate higher importance or noticeability of an AOI [28]. A larger total count of fixations has been associated with inefficient search since this indicates that many irrelevant elements have been sampled [11]. (2) Dwells provide an aggregated metric, often to describe gaze

duration. For this, multiple consecutive fixations on the same AOI are counted as a single dwell [12, 16, 22]. The number of dwells and dwell times thus link to the overall attention per AOI. Dwell times may also be taken as an indicator of task difficulty [33]. (3) The time to first fixation (TTFF) relates to how quickly an object has captured a user's attention. It is defined as the time between stimulus onset and its first fixation. Thus, TTFF can be used to assess an object's property to attract attention [6]. (4) Scanpaths refer to the path defined by successive points-of-regard (PORs) (e. g., [12]). Scanpaths help to understand user search behavior by indicating more or less efficient search [11]. (5) Saccades refer to rapid eye movements between fixations and are quantified typically according to saccadic amplitudes and saccadic velocity [11, 12]. Saccades relate to both cognitive and non-conscious processes and are thus relevant for studying affective information processing, which is outside of our study objective.

## 3 AR GUIDANCE AND EYE TRACKING SYSTEM

To asses gaze patterns in AR, we developed a system consisting of: (1) an HMD for showing visual cues, (2) eye tracking for identifying PORs, and (3) eye gaze analysis through quantitative metrics.

### 3.1 Head-Mounted Display for Guiding Attention in AR

Our system guides user attention by showing different AR-based visual cues through an HMD. Specifically, it is designed to run on HoloLens 2 and, for this, we used the Mixed Reality Toolkit and Unity. For augmented viewing, HoloLens 2 provides a horizontal FOV of 43° and a vertical FOV of 29°. However, no details about the FOV of its eye tracking system have been released so far. To give an approximation, we measured at which angles the eye tracking system of HoloLens 2 was able to record gaze points. We observed a horizontal FOV of approximately 40° in both directions and a vertical FOV of approximately 20° in the upper direction and 40° in the lower direction.

### 3.2 Eye Tracking for Identifying PORs and Fixated AOIs

For eye tracking, the system uses the HMD's built-in eye tracking sensors, which capture eye movement at a maximum sampling rate of 30 Hz. HoloLens 2 provides gaze rays that lie within 1.5° of visual angle [27]. Since manufacturer specifications can be inaccurate in real-world scenarios [9], we validated the hardware by examining accuracy and precision at three distances (0.5 m, 1 m, and 2 m) with 12 participants.[1] At each distance, seven virtual targets were arranged in a cross-format at known spatial coordinates parallel to the y-z-plane. Targets were placed 0°, 10°, and 15° from the center. Targets were presented in random order for three seconds each. Participants were asked to fixate the stimulus and press a button on a PC mouse once they perceived a color change, which occurred after two seconds. This was employed to keep the attention of the participants on the target (cf. [9]). We started recording PORs after one second for one second (∼ 30 PORs), thus stopping when the color change happened. The average accuracy was 1.67° (0.5 m), 0.51° (1 m), and 0.47° (2 m), with corresponding average precision over all targets of 1.21 (0.5 m), 0.30 (1 m), and 0.26 (2 m). Our system estimates PORs by calculating where a user's gaze ray intersects the spatial surroundings. To achieve this, we modeled the spatial surroundings using Unity, so that the virtual model of the real world overlays the physical world. Hence, the dimensions and locations of AOIs correspond to the dimensions and locations of the physical objects of interest. Using this approach, one can determine in which AOI any POR is located by simply assessing to which AOI the coordinates of that POR belong.

---

[1] Accuracy denotes the average angular offset between the calculated gaze ray and an imaginary gaze ray projected from its origin onto the target. Given a target, precision denotes the standard deviation of the calculated PORs.

### 3.3 Implemented Eye Gaze Metrics

To infer quantitative eye gaze metrics, we first identify fixations and saccades from PORs. Here, we use a dispersion-based algorithm designed to detect fixations in 3D, which is described more thoroughly in [36]. This approach makes use of ellipsoidal bounding volumes whose size depends on the distance between the user and the fixation point. For each POR, we identify a fixation by checking whether the point lies within the ellipsoidal bounding volume. Now, given a set of successive PORs that have been classified as a fixation, we attribute that fixation to an AOI if any of the PORs of that set lie within the AOI. In other words, we attribute a fixation consisting of a set of PORs to an AOI if the two intersect. Based on the identified fixations, the system calculates the following gaze metrics:

1. *Number of fixations*: For each AOI, we count the the number of fixations.

2. *Dwell duration*: For each AOI, we define dwell duration as the difference in seconds between the first fixation after entering the AOI and the last fixation before exiting the same AOI.

3. *Inter-POR distance of scanpath*: A series of PORs constitutes a scanpath in three-dimensional space. We define the inter-POR distance of such a scanpath as the average spatial distance between successive PORs.

4. *Angular distance*: We determine the visual angle $\theta$ at time $t$ between two successive gaze rays. Specifically, it is calculated through the dot product of two three-dimensional gaze direction vectors $g_t$ and $g_{t-1}$, i. e., $\theta_t = \arccos\left(\frac{\langle g_{t-1}, g_t \rangle}{\langle |g_{t-1}|, |g_t| \rangle}\right)$.

5. *TTFF*: TTFF on AOIs are calculated as the the time difference between starting the task and fixating the AOI for the first time.

## 4 USER STUDY

We explore gaze patterns under the guidance of different AR-based visual cues through a user study with $N = 12$ participants. Specifically, we assessed $C = 10$ experimental conditions, i. e., 2 baseline conditions and 8 AR-based visual cues.

### 4.1 Experimental Setup and Task

We implemented a simulated assembly task similar industrial ones [17]. Participants assembled parts on a workpiece carrier, specifically washers and nuts that had to be turned on different screws. The screws were located on a central board in front of the participant. The corresponding nuts and washers were distributed across (a) picking bins on a table and (b) picking bins located in shelves at the end of the room (see Fig. 1).



Figure 1: 180° view of the experimental setup. Screws are located on the workpiece carrier, nuts and washers in the picking bins.

We used different assembly parts as follows. Screws varied by size (i. e., 4, 6, 8, and 10 mm) and by type (i. e., "A" and "B"). This gives 24 distinct assembly parts (i. e., 8 screws, and, analogously,

8 nuts and 8 washers). As described above, nuts and washers were placed in picking bins either within the FOV (on the table) or outside the FOV (on the shelves, see Fig. 1). Picking bins and screws were labeled according to the assembly parts by a three-letter code, referring to the size, type, and assembly part. For instance, the picking bin of an 8 mm type "A" nut was labeled by "8-A-N".

In our experimental task, participants were asked to assemble $M$ target screws with the corresponding nuts and washers. During this, participants were supported by different cues as defined by the experiment condition. Depending on condition, the number of target screws to be assembled, $M$, was varied (i.e.; $M = 2$ or $M = 3$; see Sec. 4.2). Participants were allowed to choose the order of assembly. That is, participants were free to attach the nut on the screw first or the washer, if both were highlighted simultaneously in the respective experimental conditions. Further, no time limit was imposed.

## 4.2 Experimental Design

We conducted a within-subject user study with $C = 10$ conditions (2 baselines and 8 AR-based visual cues). The conditions were randomly assigned across $N = 12$ participants. The conditions represent different cues that were used to show the targets in the experimental task. In the baseline conditions, participants received guidance through a piece of paper that stated the labels of the target screws and the labels of the corresponding nuts and washers. For all other conditions, participants were shown one or multiple AR-based visual cues. The complete list of conditions is given in Fig. 2.

The visual cues differed along various dimensions: (a) whether user attention was guided to targets within or outside the FOV; (b) whether user attention was guided by a single visual cue at a time or by multiple visual cues displayed simultaneously; and (c) whether the cue appeared static or dynamic. Depending on the dimension of a visual cue, the experiment task was adapted as follows:

- *Within vs. outside FOV.* In order to guide user attention to targets within the FOV, all target picking bins were located on the table in front of the participant and thus within the user's FOV (see Fig. 1). In contrast, for guidance towards targets outside the FOV, picking bins were located on the two shelves outside of the FOV.
- *Sequential vs. simultaneous.* Sequential cues displayed only a single visual marker at a time for the target picking bins. Simultaneous visual cues showed multiple visual markers at the same time. In the latter case, visual markers were differentiated by coloring them according to the different targets. As described above, participants were free to attend the highlighted picking bins in any order and without time restrictions.
- *Static vs. dynamic.* Static cues had no or only limited movement in the AR environment, while dynamic cues had some form of movement towards the target.

Each visual cue was accompanied by a small semi-transparent highlighting box around the target screw (or screws). This highlighting box was shown in addition to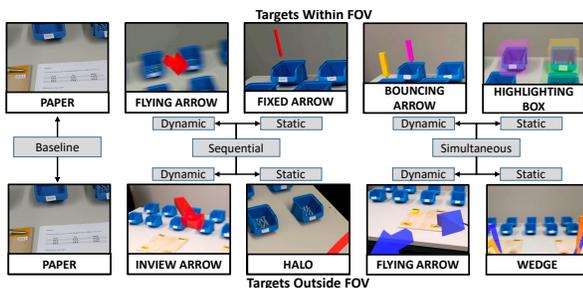 the main visual cue in order to help participants identify the target screw. As a default, participants were asked to assemble $M = 3$ target screws (selected at random). An exception was made for two visual cues (simultaneously displayed cues with targets outside FOV). Here, the number of target screws was set to $M = 2$ in order to reduce visual clutter within the HMD.

With the selection of visual cues (Fig. 2), we cover a range of common visual cues that have been found to be effective regarding task completion time or number of task errors (e.g., see [13, 15, 30]). We discarded guidance cues that were ill-suited for HoloLens, such as the ATTENTION FUNNEL [4], because it is difficult to follow in small FOVs [30] and is limited to a single off-screen object at a time [14]. In implementing the cues, we kept their visual appearance as close as possible to the respective original presentation.

## 4.3 Procedure and Measurements

Upon arrival, each participant received an information sheet summarizing the goals, methods, and compensation of the user study. After having time to ask questions, participants signed a consent form and filled out a demographics questionnaire. Participants were compensated with the equivalent of USD 20.

Participants first calibrated the HMD to their eyes for optimal hologram appearance and stable eye tracking. Before starting the experimental task, participants were introduced to the HMD in a training round in which the setting, task, and HMD controls were explained by an experimenter. Afterwards, each participant performed the experimental task for all $C = 10$ conditions. Prior to conducting the user study, we obtained ethics approval from the Ethics Committee of ETH Zurich.

## 4.4 Participants and Statistical Testing

We recruited 12 participants (6 female, 6 male) aged between 23 and 35 ($M = 26.92$; $SD = 4.14$). No participant had color vision impairment or other binocular vision disorders. All participants had normal or corrected vision. Five participants reported experience with VR technology and one with AR technology.

We used a repeated measures ANOVA. In case the assumption of sphericity was violated, the degrees of freedom were adjusted using Greenhouse–Geisser correction. In case of violated ANOVA test assumptions, we used Friedman tests together with Scheffè's method for multiple comparisons [32]. This is a conservative method for pairwise comparisons and any number of non-pairwise comparisons of group means [26]. We report mean values ($M$) and median values, depending on the underlying statistical test.

## 5 RESULTS

## 5.1 Gaze Distribution

We compare how visual cues guide user attention and thus affect gaze distribution based on the following AOIs:

1. "*On targets*" refers to all *highlighted* targets. It includes all picking bins and assembly parts that are relevant for solving the current task (independent of whether they are within or outside the FOV) and that are marked by the visual cue.

2. "*On potential targets*" refers to all *candidate* targets. It includes all picking bins and assembly parts that have not been assigned to target objects.

Results for gaze distributions (i.e., eye fixations per AOI) are shown in Fig. 3. For the conditions with **targets within FOV**, visual cues had a profound effect on the gaze distribution for "on targets" (Fig. 3a, left): The number of fixations on targets was significantly affected by the presence of a visual cue ($F(4, 44) = 12.5, p < 0.001$). All means were higher than that of the BASELINE. In other words, all AR-based cues captured more eye fixations on targets than the BASELINE (non-pairwise comparison, $p < 0.01$). The largest number of eye fixations was attributed to the HIGHLIGHTING BOXES



Figure 2: Ten experimental conditions, which are grouped by within FOV (top) and outside FOV (bottom).

cue ($M = 43.08, SD = 14.51$), which generated more eye fixations than all other conditions (all pairwise tests with $p < 0.05$, except for the FLYING ARROW). For "on potential targets" (Fig. 3a, right), the number of fixations also differed significantly between conditions ($F(1.61, 17.66) = 18.05, p < 0.001$). More specifically, simultaneous cues ($M = 12.83, SD = 4.30$) attracted fewer eye fixations than sequential cues ($M = 37.38, SD = 20.82, p < 0.001$). Conversely, the FIXED ARROW had the largest number of eye fixations on potential targets ($M = 47.25, SD = 24.10$) with significant differences (pairwise) to all other visual cues ($p < 0.05$).



(a) Conditions: Targets within FOV.
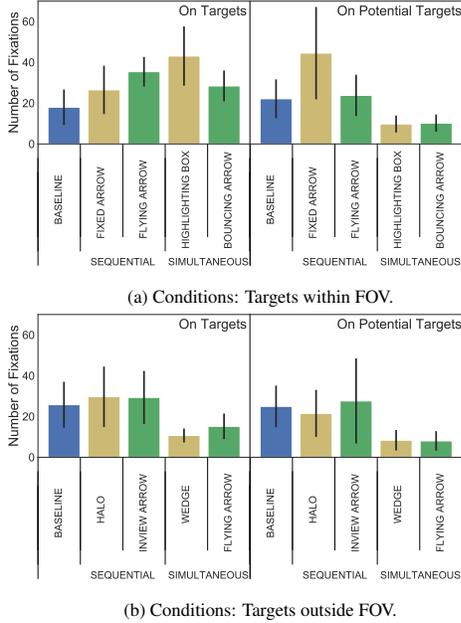


(b) Conditions: Targets outside FOV.

Figure 3: Eye fixations on targets and potential targets across different experimental conditions. Whiskers show standard deviations.

For the condition with **targets outside the FOV** and fixations "on targets" (Fig. 3b, left), group means differed significantly across cues ($F(4, 44) = 11.04; p < 0.001$). Moreover, non-pairwise comparisons of sequential ($M = 29.50; SD = 13.62$) and simultaneous ($M = 12.92; SD = 5.44$) visual cues showed significant differences ($p < 0.001$). In this regard, both baseline and sequential cues captured eye gaze to a similar extent. However, simultaneous cues generated a considerable lower number of eye fixations on targets. Similar patterns are observed for "on potential targets" (Fig. 3b, right). The average number of fixations on potential targets were significantly different among the experimental conditions ($F(4, 44) = 11.15, p < 0.001$). Likewise, comparing sequential ($M = 27.71, SD = 17.20$) and simultaneous ($M = 10.33, SD = 5.83$) visual cues yielded significant differences ($p < 0.001$).

## 5.2 Gaze Duration

We report gaze duration by comparing average dwell times among different conditions (Fig. 4). Here, we draw upon the same AOIs as before, namely "on targets" and "on potential targets".

For conditions with **targets within FOV** (Fig. 4a), the average dwell duration for "on targets" differed significantly ($\chi^2(4) = 12.00; p < 0.05$). However, comparing sequential with simultaneous visual cues and static with dynamic ones did not lead to significant differences. For "on potential targets", the mean dwell duration varied significantly across conditions ($\chi^2(4) = 23.73; p < 0.001$). Likewise, comparing sequential vs. simultaneous cues and static vs. dynamic cues did not lead to significant differences. We further conducted pairwise comparisons: The simultaneous dynamic visual cue

BOUNCING ARROW (median $= 0.18$) led to a significant difference over the baseline condition (median $= 0.39; p < 0.01$), as did the FIXED ARROW (median $= 0.28, p < 0.05$).



(a) Conditions: Targets within FOV.
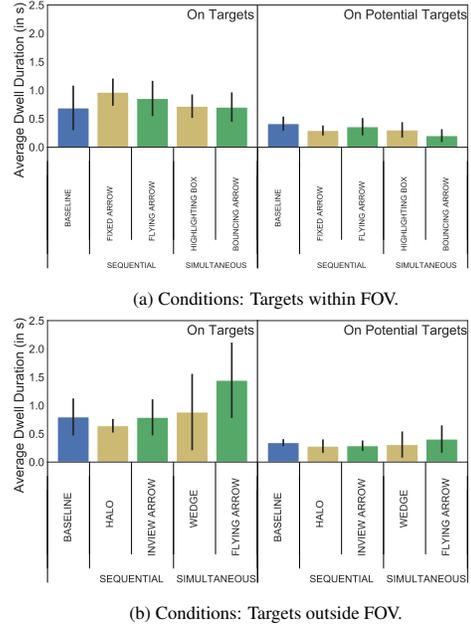


(b) Conditions: Targets outside FOV.

Figure 4: Dwell times on targets and on potential targets across different experimental conditions. Whiskers show standard deviations.

For conditions with **targets outside FOV** (Fig. 4b), average dwell duration for "on targets" varied significantly ($\chi^2(4) = 13.53; p < 0.01$). A non-pairwise comparison of sequential visual cues (median $= 0.72$) and simultaneous ones (median $= 1.16$) showed significant differences ($p < 0.05$). The visual cue FLYING ARROW (simultaneous dynamic) generated the longest average dwell duration (median $= 1.45$). However, a pairwise comparison to all conditions was not significant, except for the simultaneous static visual cue (WEDGE, $p < 0.05$). We did not find any significant deviations in group means for a dwell duration "on potential targets".

## 5.3 Inter-POR Distance of Scanpath

We report the following results for the average distance between consecutive PORs (see Fig. 5). For conditions with **targets within FOV**, group means were significantly different ($F(1.78, 19.54) = 14.41; p < 0.001$). Here, a non-pairwise comparison established that visual cues resulted in a higher average inter-POR distance ($M = 4.53; SD = 2.00$) as compared to the BASELINE ($M = 2.27; SD = 0.86; p < 0.01$). Furthermore, the inter-POR distance was higher for sequential cues ($M = 5.29; SD = 2.30$) as compared to simultaneous cues ($M = 3.77, SD = 1.30, p < 0.05$).

For conditions with **targets outside FOV**, group means significantly differed ($F(4, 44) = 48.36; p < 0.001$). Based on a non-pairwise comparison, we found that the average inter-POR distance was higher for cues ($M = 9.54; SD = 2.04$) as compared to the BASELINE ($M = 4.24; SD = 0.94; p < 0.001$). Furthermore, dynamic cues showed a shorter inter-POR distance ($M = 8.47; SD = 1.66$) as compared to static ones ($M = 10.60; SD = 1.85; p < 0.001$).

## 5.4 Time to First Fixations

We report average TTFFs (in seconds) for the following AOIs: (1) *Cue* represents the AR-based visual cue (i. e., implemented for all conditions except for the two baselines). (2) *Target boxes* and (3) *target screws* refer to the picking bins and screws, respectively.

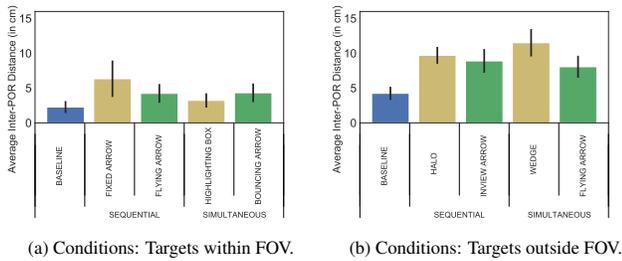(a) Conditions: Targets within FOV.  (b) Conditions: Targets outside FOV.

Figure 5: Inter-POR distances across different experimental conditions. Whiskers show standard deviations.

Both have been highlighted for the participants to solve the current task at hand. (4) *Potential target boxes* and (5) *potential target screws* comprise of any other picking bins and screws, respectively, that are not assigned to the targets and thus not highlighted by any visual marker. Fig. 6 reveals large differences in TTFF among the visual cues and AOIs. We report the highest $p$-value that applies to the mentioned conditions.

For conditions with **targets within FOV**, TTFFs differed significantly among visual cues for each AOI ($p < 0.001$). The only exception are the potential target boxes ($p = 0.16$). For both types of targets (box and screw), TTFF was lower when visual cues were shown as compared to the BASELINE ($p < 0.05$). Moreover, sequential cues had a significantly lower TTFF on potential screw targets as opposed to simultaneous ones ($p < 0.001$).

| | BASELINE | FIXED ARROW | FLYING ARROW | HIGHLIGHTING BOX | BOUNCING ARROW |
|---|---|---|---|---|---|
| Cue | | 6.8 | 5.4 | 0.88 | 8.4 |
| Target Box | 12 | 9.5 | 1.7 | 0.88 | 2.3 |
| Potential Target Box | 8.7 | 2.9 | 4.3 | 5.2 | 3.5 |
| Target Screw | 20 | 3.5 | 6.5 | 9.7 | 13 |
| Potential Target Screw | 14 | 7.8 | 6.1 | 33 | 23 |

(a) Conditions: Targets within FOV.

| | BASELINE | HALO | INVIEW ARROW | WEDGE | FLYING ARROWS |
|---|---|---|---|---|---|
| Cue | | 6.7 | 13 | 4.7 | 18 |
| Target Box | 29 | 12 | 9.7 | 19 | 12 |
| Potential Target Box | 4.3 | 4.3 | 3 | 3.3 | 10 |
| Target Screw | 30 | 6.4 | 6.8 | 28 | 18 |
| Potential Target Screw | 27 | 11 | 11 | 20 | 32 |

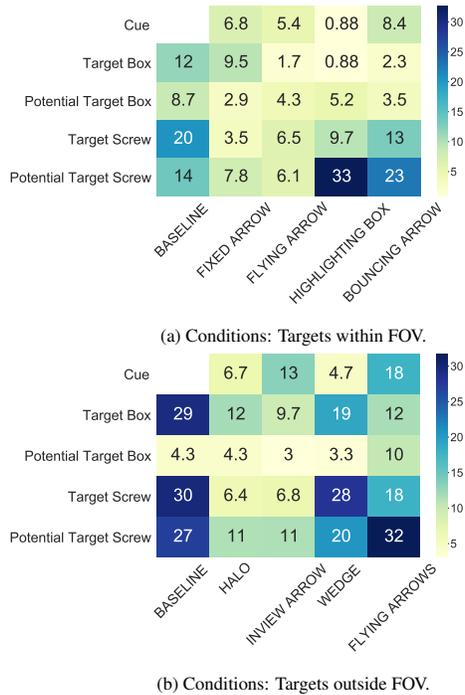(b) Conditions: Targets outside FOV.

Figure 6: TTFF (seconds) across AOIs and experimental conditions.

For the conditions with **targets outside FOV**, we found no significant differences in TTFFs regarding cues and potential targets screws. All other AOIs had significantly different TTFFs across conditions ($p < 0.01$). Looking at target boxes, TTFF was lower when visual cues were shown as compared to the BASELINE ($p < 0.01$). Regarding target screws, sequential visual cues resulted in a lower TTFF than simultaneous ones ($p < 0.05$).

## 6  DISCUSSION

Our results show that visual attention was strongly affected by the presence of cues. This finding is in line with works reviewed above that highlight the potential for AR-based visual cues to shift fixations towards objects of interest (e. g., [5]). Across all visual cues and positions of targets (within FOV and outside FOV), users fixated target objects more quickly when being guided by visual cues, which has been identified as an important criterion for successful attention guidance [30]. Additionally, visual cues increased the average inter-POR distance of participants' scanpaths. Larger distances between successive PORs suggest that visual cues are regarded as meaningful, as the cues guide user attention to the desired targets more directly and with less interim PORs [11]. Similarly, scanpath length has been considered as an indicator for efficiency or productivity of interfaces [11, 28]. Overall, these findings not only corroborate existing studies on the effectiveness of AR-based visual cues for attention guidance (e. g., [21, 30]), but also provide first insights into their underlying eye gaze mechanisms. It is worth mentioning, however, that some effects were more pronounced depending on whether the target object was within or outside the FOV.

### 6.1  The Effect of Simultaneous Cues on Gaze Behavior

Simultaneous cues were successful in shifting fixations away from non-target objects, which can be seen as an indicator for unhindered or efficient search [11]. For cues with targets within the FOV, we observed a substantial number of fixations on the target boxes. For simultaneous cues with targets outside the FOV, however, the opposite was observed (fewer fixations on target boxes). One potential reason might be that participants fixated these cues instead of the target boxes. Irrespective of this, we summarize that simultaneously presented cues were more successful in shifting attention away from non-target objects than sequential ones. This could, possibly, arise from the difference between parallel and serial search, which plays a central role in, for example, FIT [35].

Analyzing other gaze metrics like TTFF, we found some cues to affect gaze patterns very strongly and others to have a small effect. In particular, for cues with targets within the FOV, the simultaneous HIGHLIGHTING BOXES led to many quick fixations on targets. One cause might be that this cue is an exogenous one, attracting bottom-up (stimulus-driven) attention [21]. In contrast, cues like the FLYING ARROW can be seen as endogeneous, since they point towards the real target, thus requiring a process under attentional control [21]. While this is consistent with other studies [34], more empirical evidence is needed to investigate this question. For cues with targets outside the FOV, we found WEDGE to not only lead to the fewest fixations on targets, but also their slowest TTFF. Due to multiple in-view elements, this cue might lead to visual clutter (especially for small FOVs [15]), which is linked to increased difficulty of visual search [31], thereby explaining increased TTFF on targets.

### 6.2  The Effect of Cue Motion on Gaze Behavior

When presented with dynamic cues, users fixated the target boxes more quickly than the visual cue. The opposite pattern occurred across all static conditions, in which users first fixated the visual cue and then the target boxes. This result suggests that users did not need to look at the dynamic cue to extract the required information from it, which is plausible given that visual information is not only gathered at the point of eye fixation, but also in the visual periphery [8].

Our results further show that cues with more motion (e. g., FLYING ARROW) revealed pronounced effects on gaze behavior, as compared to cues with limited motion, such as the rotating INVIEW ARROW. This might be explained by differences in visual saliency among the cues. Visual saliency is is a composite metric of many low-level bottom-up features, including both size and motion, and it is generally linked to attention capture [20].

## 6.3 Conclusion and Future Research

In this paper, we explored the effects of visual guidance cues on eye gaze patterns in AR through a user study with 12 participants. Different visual cues were displayed via an AR HMD while eye movements were tracked simultaneously. Although we found our system to measure with high accuracy and precision, our findings are limited by the maximum sample frequency of 30 Hz. Thus, future studies could strive to utilize low-latency eye tracking systems for AR with high sampling frequencies. This could, for example, be beneficial for assessing additional eye gaze patterns like smooth pursuits. Further, we focused on a variety of conventional gaze-related metrics. However, a multitude of complementing metrics might be informative, such as transition probabilities among AOIs.

## REFERENCES

[1] R. Alghofaili, Y. Sawahata, H. Huang, H. C. Wang, T. Shiratori, and L. F. Yu. Lost in Style: Gaze-Driven Adaptive Aid for VR Navigation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, New York, NY, USA, 5 2019. ACM.

[2] R. Bailey, A. McNamara, N. Sudarsanam, and C. Grimm. Subtle gaze direction. *ACM Transactions on Graphics*, 28(4):1–14, 8 2009.

[3] P. Baudisch and R. Rosenholtz. Halo: A Technique for Visualizing off-Screen Objects. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 481–488, New York, NY, USA, 2003. ACM.

[4] F. Biocca, A. Tang, C. Owen, and F. Xiao. Attention Funnel: Omnidirectional 3D Cursor for Mobile Augmented Reality Platforms. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1115–1122, New York, NY, USA, 2006. ACM.

[5] A. Burova, J. Mäkelä, J. Hakulinen, T. Keskinen, H. Heinonen, S. Siltanen, and M. Turunen. Utilizing VR and Gaze Tracking to Develop AR Solutions for Industrial Maintenance. *Conference on Human Factors in Computing Systems - Proceedings*, pages 1–13, 2020.

[6] M. D. Byrne, J. R. Anderson, S. Douglass, and M. Matessa. Eye tracking the visual search of click-down menus. *Conference on Human Factors in Computing Systems - Proceedings*, pages 402–409, 1999.

[7] M. M. Chun, J. M. Wolfe, and E. B. Goldstein. *Blackwell Handbook of Perception*. Blackwell Pub., Oxford, UK, 2000.

[8] M. P. Eckstein, A. Caspi, B. R. Beutter, and B. T. Pham. The decoupling of attention and eye movements during multiple fixation search. *Journal of Vision*, 4(8):165, 2004.

[9] A. M. Feit, S. Williams, A. Toledo, A. Paradiso, H. Kulkarni, S. Kane, and M. R. Morris. Toward everyday gaze input: Accuracy and precision of eye tracking and implications for design. *Conference on Human Factors in Computing Systems*, 2017-May:1118–1130, 2017.

[10] V. Georges, F. Courtemanche, S. Sénécal, T. Baccino, M. Fredette, and P. M. Léger. UX Heatmaps: Mapping User Experience on Visual Interfaces. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 4850–4860, New York, NY, USA, 2016. ACM.

[11] J. H. Goldberg and X. P. Kotval. Computer interface evaluation using eye movements: Methods and constructs. *International Journal of Industrial Ergonomics*, 24(6):631–645, 1999.

[12] J. H. Goldberg, M. J. Stimson, M. Lewenstein, N. Scott, and A. M. Wichansky. Eye Tracking in Web Search Tasks: Design Implications. In *Proceedings of the 2002 Symposium on Eye Tracking Research & Applications*, pages 51–58, New York, NY, USA, 2002. Association for Computing Machinery.

[13] U. Gruenefeld, A. El Ali, W. Heuten, and S. Boll. Visualizing Out-of-View Objects in Head-Mounted Augmented Reality. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services*, pages 1–7, New York, NY, USA, 9 2017. ACM.

[14] U. Gruenefeld, D. Ennenga, A. E. Ali, W. Heuten, and S. Boll. EyeSee360: Designing a Visualization Technique for out-of-View Objects in Head-Mounted Augmented Reality. In *Proceedings of the 5th Symposium on Spatial User Interaction*, pages 109–118, New York, New York, USA, 10 2017. ACM.

[15] U. Gruenefeld, D. Lange, L. Hammer, S. Boll, and W. Heuten. FlyingARrow: Pointing Towards Out-of-View Objects on Augmented Reality Devices. In *Proceedings of the 7th ACM International Symposium on Pervasive Displays*, volume 18, pages 1–6, New York, NY, USA, 2018. Association for Computing Machinery.

[16] M. Gullberg and K. Holmqvist. Visual Attention towards Gestures in Face-to-Face Interaction vs. on Screen. In *Gesture and Sign Language in Human-Computer Interaction (Lecture Notes in Computer Science)*, volume 2298, pages 206–214. Springer, Berlin, Heidelberg, 2002.

[17] A. Guo, X. Wu, Z. Shen, T. Starner, H. Baumann, and S. Gilliland. Order Picking with Head-Up Displays. *Computer*, 48(6):16–24, 2015.

[18] S. Gustafson, P. Baudisch, C. Gutwin, and P. Irani. Wedge: Clutter-Free Visualization of off-Screen Locations. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 787–796, New York, New York, USA, 2008. ACM Press.

[19] R. Hoffmann, P. Baudisch, and D. S. Weld. Evaluating Visual Cues for Window Switching on Large Screens. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 929–938, New York, New York, USA, 2008. ACM Press.

[20] L. Itti and C. Koch. Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2(3):194–203, 2001.

[21] N. F. S. Jeffri and D. R. A. Rambli. Guidelines for the Interface Design of AR Systems for Manual Assembly. In *Proceedings of the 2020 International Conference on Virtual and Augmented Reality Simulations*, pages 70–77, 2020.

[22] J. Y. Jiang, F. Guo, J. H. Chen, X. H. Tian, and W. Lv. Applying Eye-Tracking Technology to Measure Interactive Experience Toward the Navigation Interface of Mobile Games Considering Different Visual Attention Mechanisms. *Applied Sciences*, 9(16):3242–3254, 2019.

[23] T. A. Kelley, J. T. Serences, B. Giesbrecht, and S. Yantis. Cortical mechanisms for shifting and holding visuospatial attention. *Cerebral Cortex*, 18(1):114–125, 2008.

[24] T. Kinsman, K. Evans, G. Sweeney, T. Keane, and J. Pelz. Ego-motion compensation improves fixation detection in wearable eye tracking. *Eye Tracking Research and Applications Symposium (ETRA)*, pages 221–224, 2012.

[25] A. Krekhov and J. Krüger. Deadeye: A Novel Preattentive Visualization Technique Based on Dichoptic Presentation. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):936–945, 1 2019.

[26] S. Lee and K. Lee. What is the proper way to apply the multiple comparison test? *Korean Journal of Anesthesiology*, 71(5):353–360, 2018.

[27] Microsoft. Eye tracking on HoloLens 2, 2020.

[28] A. Poole and L. J. Ball. Eye Tracking in Human-Computer Interaction and Usability Research: Current Status and Future Prospects. *Encyclopedia of Human Computer Interaction*, pages 211–219, 2005.

[29] M. I. Posner. Orienting of attention. *The Quarterly Journal of Experimental Psychology*, 32(1):3–25, 1980.

[30] P. Renner and T. Pfeiffer. Attention guiding techniques using peripheral vision and eye tracking for feedback in augmented-reality-based assistance systems. In *2017 IEEE Symposium on 3D User Interfaces (3DUI)*, pages 186–194, 4 2017.

[31] R. Rosenholtz, L. Yuanzhen, and L. Nakano. Measuring visual clutter. *Journal of Vision*, 7:1–22, 2007.

[32] H. Scheffe. A Method for Judging All Contrasts in the Analysis of Variance. *Biometrika*, 40(1/2):87–104, 1953.

[33] S. Stork and A. Schubö. Human cognition in manual assembly: Theories and applications. *Advanced Engineering Informatics*, 24(3):320–328, 2010.

[34] J. Theeuwes. Top-down and bottom-up control of visual selection. *Acta Psychologica*, 135(2):77–99, 2010.

[35] Treisman. A and G. Gelade. A feature-integration theory of attention. *Cognitive psychology*, 12(1):97–136, 1980.

[36] S. Weber, R. S. Schubert, S. Vogt, B. M. Velichkovsky, and S. Pannasch. Gaze3DFix: Detecting 3D fixations with an ellipsoidal bounding volume. *Behavior Research Methods*, 50(5):2004–2015, 2018.

[37] J. M. Wolfe. Guided Search 2.0 A revised model of visual search. *Psychonomic Bulletin & Review*, 1(2):202–238, 6 1994.

[38] J. M. Wolfe. Visual Search: How Do We Find What We Are Looking For? *Annual Review of Vision Science*, 6:539–562, 2020.