

Evaluation of Textual Feedback for Incorrect Usage of an Augmented Reality Application

JULIAN LÜKEN, MEHMED MUSTAFA, JAN SCHNEIDER,
STEFFEN TUNKEL, CHRIS WARIN

March 16, 2020

Georg-August University, Göttingen

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Introduction

In recent years, the number of available mobile Augmented Reality (AR) applications for smartphones has increased [4]. Although mobile AR technology is no novelty, it was not widely available for society until a variety of smartphones were capable of running AR applications. Since most users are newly discovering AR technology, they don't know how to use its features yet, thus using unsupported gestures as input, which leads to frustration. Our main motivation for conducting this research is the lack of good feedback in case of incorrect usage, which could help a lot of users to avoid repeating common mistakes while using AR applications for smartphones. Providing good feedback could also help and teach users to control such AR applications without the need of a separate introduction.

Our research goal is to provide appropriate feedback to unsupported gestures and decide what type of textual feedback is the best. For validation of our approach we conducted a

usability case study.

The organization of the paper is as follows. The section Foundations gives information about usability in general and the AR environment we provide the feedback for. In the section Related Work we discuss similar work on AR and feedback done in the past. The Approach section gives a detailed illustration of the feedback we present to users. In the section Case Study, we talk about the setup case study we conducted in order to validate our approach, followed by our results and the discussion thereof. Lastly, we summarize and conclude in the Summary and Outlook.

Foundations

Since the main goal of our work presented in this paper is to find suitable user feedback for augmented reality environments, we will further clarify what AR means and which technologies we used. AR itself describes a variety of technologies that place 3D virtual objects into a 3D real environment in real time, for example via a handheld device like a mobile phone (mobile AR) or using wearable goggles like the Microsoft HoloLens. In this paper we focus on AR for mobile phones, where recorded imagery of the phone is being augmented and shown to the user to display a mix of virtual objects and the real world at runtime.

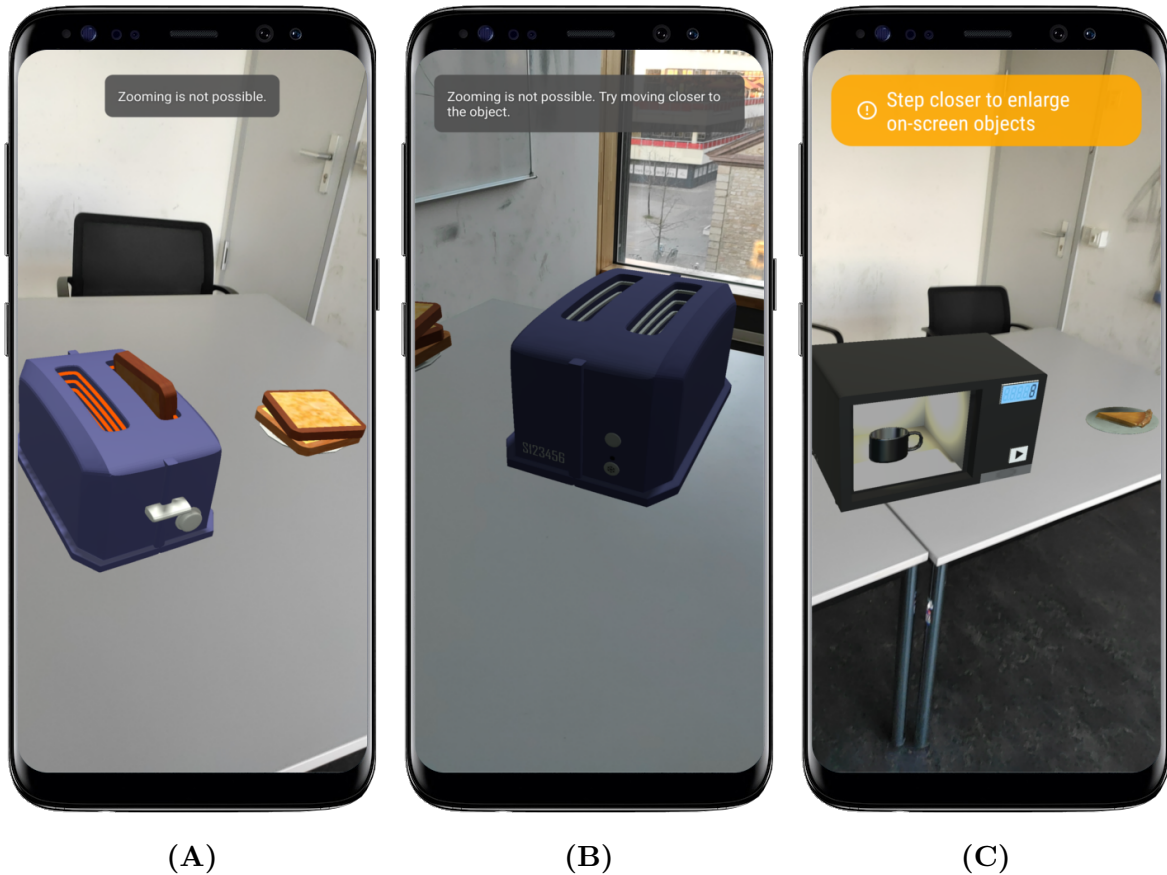


Figure 1: Three screenshots of an app made with the Vivian framework

A virtual prototype is a "[...] *computer simulation of a physical product that can be presented, analyzed, and tested from concerned product life-cycle aspects such as design/engineering, manufacturing, service, and recycling as if on a real physical model.*" [Wang2002]. The Vivian framework is an extension for Unity3D, that can be used to create mobile AR implementations of prototypes using a 3D model and one or more finite-state machines. (sauce?)

In a mobile app made with the Vivian framework, all user input relies on one finger gestures on the touch screen. In the example in Figure 1 (A and B) we can observe a 3D model of a toaster projected onto a live video of a real world scenario using Vivian. The user can interact with the different 3D models on the screen. Buttons can be pressed by simply tapping the screen where the respective button is displayed. Knobs and levers can be moved by putting a finger on the screen where the knob or lever is displayed, moving the finger on the screen until finding the desired position, and removing the finger from the screen afterwards. We refer to this gesture as *drag-and-drop*. The Vivian framework makes a clear distinction between movable and non-movable objects in a scene. The slice of bread in the aforementioned example is movable, whereas the toaster itself cannot be moved. Movable objects can be moved in the same fashion as a lever. To cover greater distances, one can move the device itself while dragging the object. Rotation similarly requires the user to hold their finger on the object and rotate the device until they find the desired orientation. The limit of one-finger-interactions exists due to the peculiarity of AR environments where two-finger-interactions (like a pinch zoom or a two-finger-rotation) would be contradicting with the intended design of such an environment. If one wants to take a closer look at the previously mentioned toaster for example the user would need to step closer to the toaster.

To evaluate different types of feedback, we decided to use usability testing as a tool to quantify the corresponding results. According to Nielsen, "*usability is a quality attribute that assesses how easy user interfaces are to use. The word "usability" also refers to methods for improving ease-of-use during the design process.*" [Nielsen2012]. While usability testing refers to a "*observational methodology to uncover problems and opportunities in designs*" with the goal to identify problems in the design of the product or service, uncover opportunities to improve and learn about the target user's behavior and preferences [Moran2019]. Therefore we chose to have the users fill out a System Usability Scale (SUS) questionnaire which we tailored to fit our needs. The basic SUS consists of ten questions regarding a systems usability. We added a few questions which target the feedback more as well as an open feedback item, where the participants were able to share their thoughts on the system with us.

Related Work

The topic of error handling in AR-applications has not been broadly studied yet. There are several guidelines on how to handle wrong user input. Microsoft states in its design guidelines, that error messages should alert users of an already occurred problem. Another point they make is that the message should be suppressed if it does not make the users change their behavior or perform an action [Microsoft]. In addition to that, Nielsen's *Error Message Guidelines* demand

error messages to be phrased politely. This is to avoid the implication that the user did something stupid or inherently wrong. It is also important that the messages are precise descriptions of exact problems, offering constructive advice on how to fix the problem[Nielsen2001]. A 2002 IEEE article discusses two approaches of how user help might be implemented in AR environments. What they did was to offer an icon which the users might activate themselves when they seek advice for a given task. The other feature was an automated help message that pops up when the user moves an object associated to a certain task closer to his face and tilts it. They state, that "[t]his approach is more suitable for AR interfaces than traditional desktop help systems, which either distract users with a constant barrage of help messages or interrupt their work by making them search explicitly for help." [IEEE2002].

Approach

In order to find feedback in response to incorrect usage of a mobile AR application, we have to define what kinds of gestures we consider incorrect in the first place. The user interface provided by the Vivian framework is based on one finger inputs (see Foundations), therefore we should take into consideration each input that is not done with a single finger. Prominent examples for multiple finger gestures are the *pinch zoom* and the *two finger rotation*. The pinch zoom is a gesture in which two fingers are placed on the touch screen and moved apart in opposite directions. The two finger rotation is a gesture in which two fingers are placed on the touch screen also and are together moved clockwise or counterclockwise about the center of the positions of the fingers. Common misconceptions by inexperienced mobile AR users might include using said pinch zoom to either zoom in on objects or to bring them closer, or using two finger rotation to rotate objects.

To make users quickly and easily understand the controls of Vivian-based apps, we extended the Vivian framework to provide three feedback message implementations that differ in content, colors and size. The content of the first feedback implementation is a *critique*. If the user inputs a gesture that we consider incorrect, the app outputs a message that said input is not possible. For example, if a user tried pinch zooming, they would receive a message saying: "Pinch zooming is not possible". The content of the second feedback implementation is *critique* and *support*. We therefore refer to this feedback type as *combined* in the course of this paper. The user is provided with a message stating that above mentioned input is not possible and hinting them towards what they should try instead. In the scope of the previous example, such a message would say: "Pinch zooming is not possible. Try moving around the object." In the first two implementations, the messages have the same size and color scheme (see Figure 1, A and B). The third implementation is more concise *support*. For our previous example the exact message is saying: "Step closer to enlarge on-screen objects." It also differs in color scheme and size (see Figure 1, C). Each feedback message in every implementation is displayed for 5 seconds.

	Description	Gesture (input)	Response (output)
#1	critique	two finger rotation	Rotating the object is not possible.
		pinch zoom	Zooming is not possible.
#2	combined	two finger rotation	Rotating the object is not possible. Try moving around the object.
		pinch zoom	Zooming is not possible. Try moving closer to the object.
#3	support	two finger rotation (movable object)	Hold the phone and move the object to rotate
		two finger rotation (elsewhere)	Move around the objects
		pinch zoom	Step closer to enlarge on-screen objects

Table 1: The responses for each input in each implementation.

The contents of each message can be found in Table 1. The description column holds the names we assigned to the three different feedback implementations. In the gesture and response columns you can find the corresponding inputs and outputs respectively. If for example a user tried a two finger rotation in the critique implementation, the app would display a message saying "Rotating the object is not possible". In the support implementation we made the distinction between two finger rotations on movable objects and elsewhere. If the aforementioned center point of the two finger rotation is placed on a movable object (see Foundations), the app displays the message that corresponds to "two finger rotation (movable object)". If the two finger rotation is executed elsewhere, the "two finger rotation (elsewhere)" message is displayed instead.

Case Study

Setup of the Case Study

To evaluate the quality of our feedback implementations we conducted a case study in the domain of usability engineering. We used two different prototypes supported by the Vivian Framework (see Foundations). Both of them are quite simple kitchen devices: a toaster prototype and a microwave prototype (as seen in Figure 1).

The microwave's functionality is limited to heating an object inside it with constant power. It has one button to add 10 seconds to the heating duration and one to open the door. The door can be closed by moving it. The status of the microwave is visually indicated by a small display showing the remaining heating time and a light inside the device, which turns on when it is open or heating. The toaster can toast one or two pieces of bread at a time. The toasting can be started by pulling down a handle and stopped by pushing it up again or by pressing the stop button, which is on the backside. Otherwise, the toasting stops automatically after a time which is defined by the position of a rotatable knob. The time mode is divided into 'low', 'medium' and 'high'. Further functionality is provided by the unfreezing mode, which can

be activated by the snowflake button on the backside. The activation of the unfreezing mode results in a longer toasting duration and is indicated by a light above the button. The toasting process itself is displayed by the glowing of the heating elements. Both prototypes have a serial number on the backside. Both scenes contain different additional objects, which the user can interact with. These are pieces of bread for the toaster and for the microwave a cup which is already inside the microwave in the beginning and a piece of cake next to it.

#	Microwave	Toaster
1	Read serial number off the microwave	Read serial number off the toaster
2	Heat up the cup	Toast the bread
3	Remove the cup, put the pie in, set the timer to 20 seconds and remove the plate at 5 seconds	Toast the bread on high heat and put the toaster in unfreezing mode

Table 2: The tasks for the different prototypes.

Based on these functionalities we designed 3 tasks per prototype, which are shown in Table 2. Important about this tasks is that they are increasing in complexity. An example of this increase is that task 2 of the microwave asks to heat the cup, which is already in the microwave, for any time. Therefore just the start button has to be pressed once to fulfill the task. Afterwards, task 3 is to replace the cup by the pie and heat for a specific time. This contains a lot more necessary steps: open the door, move the cup out, move the pie in, close the door, press the start button multiple times and finally also press the stop button after the desired time is over.

For the case study, we divided the participants into 4 groups. Each of these groups tested another version of the feedback implementation. We gave no feedback at all to the participants in group 0. Participants in group 1 got our first implementation with *critique* like feedback, in group 2 they got our second implementation with a *combined* feedback version and in group 3 they got our third implementation with a *supportive* feedback (see Approach). Each participant was asked to fulfill the tasks for each prototype in the order provided in Table 2. One half of the participants were asked to do the tasks of the toaster prototype first, while the other half did the tasks for the microwave prototype first. We were able to evaluate subgroups with a certain implementation level and a certain order of prototypes they tested. In the following, the subgroups are named *TM*, which translates to "first toaster, then microwave" and *MT*, which translates to "first microwave then toaster", followed by the number of their feedback group. For example, a participant who was assigned the toaster first and then the microwave without any feedback was in the subgroup TM0.

Before the actual start of the test, we gave every participant a short introduction. This introduction contains information about the AR technology and reasons behind usability testing to prepare the participants and create a similar level of expectation. It deliberately does not contain any information about how to use the application or the scope of the case study. This is important so the participants are not biased about wrong usage and the feedback they would get. During the test, we used screen recordings to measure each participants use of the application.

These enabled us to acquire completion times for the different tasks and count the number of wrong usages. We used a notepad to register further observations.

Following the test of the application the participants were asked to fill in a questionnaire. This questionnaire consists of 3 parts. First we asked participants if they have used AR applications before. The second part is an SUS (see Foundations) about the usability of the AR application. The third part is an open question directly aimed to the feedback provided by the application, namely: *"If you have any suggestions on what kind of feedback from the app you would find better, please let us know"*.

We conducted the case study with 10 participants per group of feedback implementation. This leads to 40 people participating in the case study in total. The smallest sub groups, defined by implementation and task order, were still filled with 5 participants.

Results of the Case Study

In our approach we consider especially two gestures as incorrect usage, the pinch zoom and the two finger rotation (see Approach). With the setup of our case study the two finger rotation appeared more often than the pinch zoom gesture. In total over all participants and all tasks, two finger rotations appeared 225 times, while pinch zoom just appeared 51 times. Since 82% of the incorrect usages are two finger rotations, the results of our case study concentrate on those. We especially evaluate the 2nd task for the toaster since it was most demanding for the users. As shown in diagram 2, 91% of the overall two finger rotations occurred in that task. Also the majority of pinch zoom gestures occurred there.

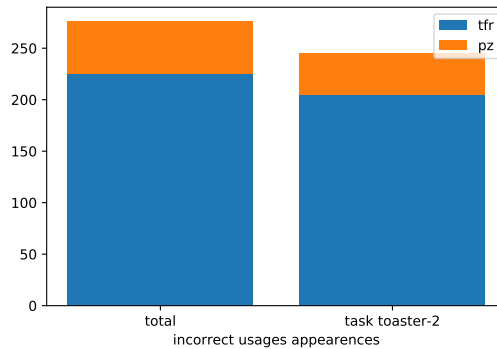


Figure 2: DUMMY: Appearances of two finger rotations and pinch zooms in total and for the 2nd task for the toaster prototype.

To evaluate the effectiveness of the different implementations we compare the users behavior for the 2nd task of the toaster, with different measures. Diagram 3 shows the two most important measures. Diagram 3-A gives the mean time needed for the task with the different implementations. The black lines on the bar plot give additional information with the standard deviation for each implementations time. Diagram 3-B is similar. This diagram shows the mean values of two finger rotation tries per user for the different implementations. Also here the black lines on the bars give the standard deviation for each of this values. Important for this graph

is that users, who had zero two finger rotation attempts, are excluded from the statistic. The reason behind it is, that these users didn't receive any feedback message independent of their given implementation. To trigger any feedback a incorrect usage must appear at least once. In total 4 of the 40 participants didn't tried any two finger rotations on this task. The maximal number for one specific implementation is 2. That means for each implementation at least 8 participants data was used to take the statistical values.

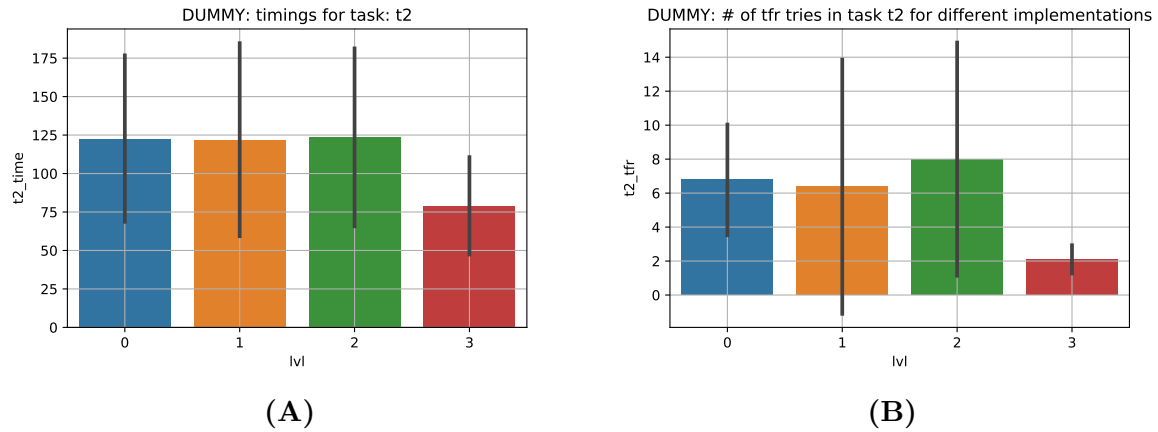


Figure 3: DUMMY: Evaluation of the time and two finger rotation tries needed for the 2nd task for the toaster prototype

To check, whether the result is statistically significant we perform statistical tests on both, the time needed for the task and the number of two finger rotation attempts. The main test is *Welch t-test* to compare the different implementations with the zero group, which didn't receive any kind of feedback messages. *Welch t-test* has the assumption that both populations are normally distributed. We used the *Shapiro-Wilk* test on the residuals between the given implementations group and the zero feedback group. Generally we consider p-values below 0.05% as statistically significant. Table 4 shows the p-values of the tests for the implementation with *concise* feedback (used by group 3). We focus on this results because the mean values of this implementations metrics implicate the biggest difference to the zero feedback group. Since all of the p-values of the *Shapiro-Wilk* test are above the threshold, the normality assumption for *Welch t-test* holds. The result of the *Welch t-test* for the number of two finger rotation tries in task 2 of the toaster implies a significant difference to the zero group, which can not be confirmed by the result for the timings for the task.

```
normality tfr diff: 0.9708724021911621
normality time diff: 0.5073676109313965
Welch tfr: 0.010198448041868305
Welch time: 0.05448460045107597
```

Figure 4: DUMMY: p-values for Shapiro-Wilk test and Welch t-test with both metrics

In addition we evaluate the relationship between our two metrics for the task. The correlation coefficient between the both is 0.43, which indicates a moderate positive relationship between

the metrics. Diagram 5 visualizes the relationship between the time needed for the task and the number of two finger rotation tries. The color of a data point is set by the implementation used by this participant. It is visible that the participants with implementation 3 are all clustered with a low time and a low amount of two finger rotations, while all other groups are more randomly spread.

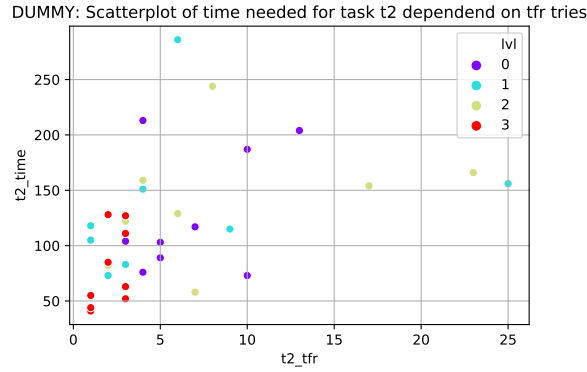


Figure 5: DUMMY: Time needed to fulfill the 2nd task for the toaster prototype depended on the number of tfr tries.

For the evaluation of the questionnaire we calculate the SUS scores for every participant. Diagram 6 shows the scores mean for each implementation. The horizontal red line is at a score of 68, which means average usability.

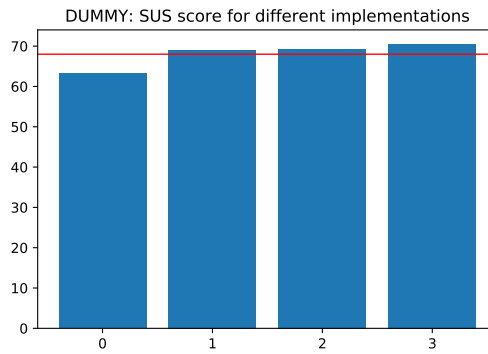


Figure 6: DUMMY: SUS score for different implementations.

24 of the 40 participants answered the open question (*"If you have any suggestions on what kind of feedback from the app you would find better, please let us know"*). We categorize the answers with 4 different labels, divided by the topic the answer is about. These are *app feedback*, *app controls*, *prototype design* and *other*. The appearances of the different types of answers are listed in diagram 7. In case an answer fit to multiple labels, we counted it for all the fitting categories. In addition to the type of answer we also take the implementation the participant tested and the specific text into account to gain knowledge from them.

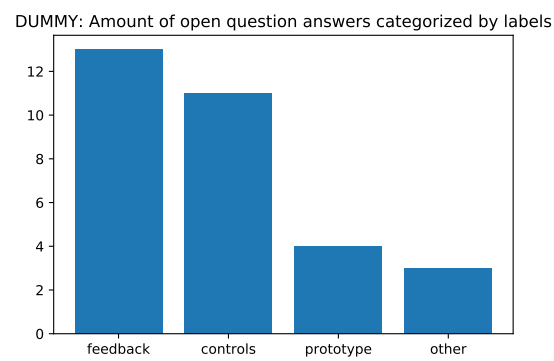


Figure 7: DUMMY: Amount of open question answers categorized by labels.

Discussion of the Case Study Results

Summary and Outlook

[1] [2] [3] AR AR

References

- [1] Arindam Dey et al. “A Systematic Review of Usability Studies in Augmented Reality between 2005 and 2014”. In: *2016 IEEE International Symposium on Mixed and Augmented Reality (ISMAR-Adjunct)*. IEEE, Sept. 2016. DOI: 10.1109/ismar-adjunct.2016.0036.
- [2] Susanna Nilsson and Björn Johansson. “Fun and usable”. In: *Proceedings of the 2007 conference of the computer-human interaction special interest group (CHISIG) of Australia on Computer-human interaction: design: activities, artifacts and environments - OZCHI '07*. ACM Press, 2007. DOI: 10.1145/1324892.1324915.
- [3] I. Poupyrev et al. “Developing a generic augmented-reality interface”. In: *Computer* 35.3 (Mar. 2002), pp. 44–50. DOI: 10.1109/2.989929.
- [4] Tractica. *Mobile Augmented Reality Market to Reach 1.9 Billion Unique Monthly Active Users by 2022*. Apr. 2017. URL: <https://tractica.ondia.com/newsroom/press-releases/mobile-augmented-reality-market-to-reach-1-9-billion-unique-monthly-active-users-by-2022-according-to-tractica/>.