



VOXReality

VOICE DRIVEN INTERACTION IN XR SPACES

Definition and Analysis of VOXReality Use Cases V2

D2.2

31-07-2024



Funded by
the European Union

Version	2.0
WP	WP2
Dissemination level	Public
Deliverable lead	NWO-I
Authors	Moonisa Ahsan, Pablo Cesar, Irene Viola, Sueyoon Lee (NWO-I), Olga Chatzifoti (MAG); Eleni Oikonomou, Ioannis Radin (AEF); Stavroula Bourou, Dimitris Kontopoulos (SYN); Manuel Toledo (VRDays); Leesa Joyce, Clayton Gordy (HOLO); Konstantinos Koutoulakis (ADAPT); Petros Drakoulis, Georgios Papadopoulos, Sotiris Karavarsamis, Athanasios Ntovas, Alexandros Doumanoglou, Konstantoudakis Konstantinos, Mpiliouisis Stefanos, Papadopoulos Georgios, Zarpalas Dimitrios (CERTH); Yusuf Can Semerci (UM).
Reviewers	Spiros Borotis (MAG); Ana Luísa Alves (F6S);
Abstract	This deliverable provides an overview of the activities of WP2, focusing on the description of the use cases, existing and new user requirements for VOXReality. It details the user-centred methodology used to gather requirements, including weekly calls, focus groups, in-situ visits, revisions on requirements, pilot and a user study. The core contributions include updated requirements based on lessons learned from first phase of pilots, with a concluding section summarizing key findings and recommendations for future deliverables.
Keywords	<i>Use Cases, New User Requirements, Technical Requirements, User-Centred Methodology, Focus Groups, Pilots</i>
License	 <p>This work is licensed under a Creative Commons Attribution-No Derivatives 4.0 International License (CC BY-ND 4.0). See: https://creativecommons.org/licenses/by-nd/4.0/</p>

Dissemination Level

PU **Public**

PP Restricted to other programme participants (Including the Commission Services)

RE Restricted to a group specified by the consortium (Including the Commission Services)

CO Confidential, only for members of the consortium (Including the Commission Services)

Nature	
PR	Prototype
RE	Report
SP	Specification
TO	Tool
OT	Other

Version History (Extended)

Version	Date	Owner	Author(s) / Editor(s)	Changes to previous version
0.1	2022-12-22	NWO-I	Pablo Cesar	ToC and initial contents
0.2	2023-01-10	NWO-I	Sueyoon Lee	Section 4.2 - Initial list of user requirements before the design sessions
0.3	2023-03-03	CERTH	Drakoulis Petros, Konstantoudakis Konstantinos, Mpiliouis Stefanos, Papadopoulos Georgios, Zarpalas Dimitrios	Section 6.2
0.4	2023-03-03	UM	Yusuf Can Semerci	Added Chapter 3
0.5	2023-03-09	NWO-I	Sueyoon Lee	Section 2, Section 4.1
0.6	2023-03-09	VRDays	Manuel Toledo	Section 5.2
0.7	2023-03-10	Athens Festival	Elena Oikonomou	Section 5.1
0.8	2023-03-10	HOLO	Clayton Gordy	Section 5.3
0.9	2023-03-11	NWO-I	Moonisa Ahsan	Section 4.2
0.10	2023-03-12	NWO-I	Pablo Cesar and Irene Viola	Sections 1 and 6; consolidation of the text
0.11	2023-03-12	NWO-I	Moonisa Ahsan	Section 5.2
0.12	2023-03-16	NWO-I	Sueyoon Lee	Section 7
0.13	2023-03-20	NWO-I	Moonisa Ahsan, Sueyoon Lee	Section 7 (Table 9, 10, 11)
0.14	2023-03-22	NWO-I	Moonisa Ahsan	Section 7.1, Overall Revisions
0.15	2023-03-22	NWO-I	Pablo Cesar, Irene Viola, Sueyoon Lee, Moonisa Ahsan	Full Draft
0.16	2023-03-28	SYN	Stavroula Bourou	Peer revision
0.17	2023-03-28	F6S	Ana Luisa Alves	Peer revision
0.18	2023-03-30	MAG	Spiros Borotis	Peer revision
1.0	2023-04-04	NWO-I	Moonisa Ahsan Irene Viola Pablo Cesar Sueyoon Lee	Implementation of feedback from partners and overall revisions
Version 2				
1.1	2024-07-11	NWO-I	Moonisa Ahsan Pablo Cesar	Extension of the TOC and document structure
1.2	2024-07-16	NWO-I	Moonisa Ahsan	Section 2, 8,

1.3	2024-07-25	HOLO NWO-I ADAPT	Leesa Joyce Moonisa Ahsan Konstantinos Koutoulakis	Section 9 (Training Requirements)
1.4	2024-07-29	MAG AEF NWO-I	Olga Chatzifoti Eleni Oikonomou Ioannis Radin Moonisa Ahsan	Section 9 (Theatre Requirements)
1.5	2024-07-29	SYN VRDays NWO-I	Stavroula Bourou Dimitris Kontopoulos Manuel Toledo Moonisa Ahsan	Section 9 (Conference Requirements)
1.6	2024-08-05	NWO-I	Moonisa Ahsan	Section 1, 9, 10
1.7	2024-08-08	NWO-I	Moonisa Ahsan Pablo Cesar	Full Draft
1.8	2024-08-26	MAG	Spiros Borotis,	Internal Peer Review
2.0	2024-08-30	NWO-I	Moonisa Ahsan Pablo Cesar	Final Revision and Submit

Table of Contents

Version History (Extended)	4
Table of Contents	6
List of Abbreviations & Acronyms	8
List of Figures	9
List of Tables	11
Executive Summary	12
1 Introduction	13
1.1 VOXReality concept and approach	13
1.2 Purpose of the deliverable	13
1.3 Intended Audience	14
1.4 Structure	14
2 Methodology	15
2.1 Year 1 Overview	15
2.2 Year 2 Overview	17
2.3 Conclusion	22
3 Details of Weekly calls (Year 1)	23
3.1 Weekly calls <i>modus operandi</i>	23
3.2 Results	24
3.3 Coordination with technical partners	24
4 Use cases initial requirements (Year 1)	25
4.1 VR Conference	25
4.2 Augmented Theatre	27
4.3 Training Assistant	29
5 In-Situ Visits (Year 1)	31
5.1 Methodology	31
5.2 Data Analysis and Results	42
6 Use Cases descriptions (Year 1)	53
6.1 VR Conference	53
6.2 Augmented Theatre	57
6.3 Training Assistant	59
7 Requirements (Year 1)	63
7.1 User Requirements	63
7.2 Technical Requirements	79
8 Summary of Pilots (Year 2)	87
8.1 VR Conference	88
8.2 Augmented Theatre	91

8.3	Training Assistant.....	94
9	New Requirements (Year 2)	97
9.1	VR Conference Requirements.....	97
9.2	Augmented Theatre.....	102
9.3	Training Assistant.....	107
10	Conclusion	112
	References	113

PENDING APPROVAL

List of Abbreviations & Acronyms

3D	:	Three-dimensional
AI	:	Artificial Intelligence
AR	:	Augmented Reality
AR3S	:	Augmented Reality Engineering Space
ASR	:	Automatic Speech Recognition
CAD	:	Computer-Aided Design
CV	:	Computer Vision
DA	:	Digital Agent
FAQ	:	Frequently Asked Questions
HMD	:	Head-Mounted Display
ISAR	:	Interactive Streaming for Augmented Reality
ML	:	Machine Learning
NLG	:	Natural Language Generation
NLP	:	Natural Language Processing
NLU	:	Natural Language Understanding
NMT	:	Neural Machine Translation
OC/OC+	:	Open Calls
PC	:	Personal Computer
PDF	:	Portable Document Format
Q&A	:	Questions and Answers
SDK	:	Software Development Kit
UI	:	User Interface
VCE	:	Virtual Conference Environment
VFX	:	Visual Effects
VL	:	Vision-Language
VR	:	Virtual Reality
XR	:	eXtended Reality

List of Figures

Figure 1 Flow of deliverables leading to Definition and Analysis of VOXReality Use cases.....	12
Figure 2. User Requirements Methodology Timeline	15
Figure 3. Participants discussing their ideas during the VRDays, AEF and HOLO workshop.....	16
Figure 4. Results of the activity sheets from the AEF workshop.....	16
Figure 5 Screenshot of revised user-requirements for each of the use-cases.	17
Figure 6 Technology Maturity Diagram	18
Figure 7 (Left) VR Conference App and (Right) AR Theatre App Interaction Flow (Source D5.1)	19
Figure 8 Pilot Implementation Plan	19
Figure 9 Pilot Demonstrations (A) VR Conference; (B) Augmented Theatre; (C) Training Assistant	20
Figure 10 (Left) Static Captions; (Right) Dynamic Captions from a participant's perspective	21
Figure 11 Flow of deliverables leading to Definition and Analysis of VOXReality Use cases.....	22
Figure 12. Participants working on activities during the VR Conference (focus-group) workshop	31
Figure 13. Virtual agent (left) and language (right) translation examples.....	32
Figure 14. Activity 2-1, 2-2 sheet for the VR Conference use case	33
Figure 15. Scenario case and supplementary image on navigating in VR conference	33
Figure 16. Activity 3-1: Brainstorming sheet for drawing ideal VR conference navigation scenario ..	34
Figure 17. (Left) participants sharing their brainstormed ideas, (Right) participants voting for the two best ideas for each scenario	34
Figure 18. Different forms of virtual assistant in VR conference	35
Figure 19. The moderator introducing the workshop to the participants in focus group workshop	36
Figure 20. Activity sheet 1-2: My visit to the theatre (drawing timeline of a visit to the theatre)	37
Figure 21. Scenario case and supplementary image on watching a play in a foreign language	38
Figure 22. Activity sheet 3-1: Brainstorming sheet for drawing AR glass UI and scenarios for the augmented theatre use case	38
Figure 23. The moderator introducing the workshop to the participants	39
Figure 24. A user persona, James, with a beginner level machine assembly experience	40
Figure 25. Scenario case and supplementary image on machine assembly experience with AR glasses	41
Figure 26. Activity sheet 3-1: Brainstorming sheet for ideal AR UI and scenario during machine assembly training.....	41
Figure 27. Selected user-activities scans from the data worksheets filled by participants	43
Figure 28. Insights Canvas for VR Conference (VRDays) Use-case.....	47
Figure 29. Insights Canvas for Augmented Theatre (AEF - Athens Festival) Use-case.....	48
Figure 30. Insights Canvas for Training Assistant (HOLO) Use-case.....	49
Figure 31. Social Area - Immersive Tech Week 2022 (©SnapBoys.nl)	54
Figure 32. Tradeshow - Immersive Tech Week 2022 (© JaynoBrk)	54
Figure 33. Main conference - Immersive Tech Week 2022 (© JaynoBrk).....	55
Figure 34. Virtual Conference space – Conceptual floorplan (©VRDays)	56
Figure 35. A potential stage at the Athens Epidaurus Festival's, (Left) Distance between the stage and the audience, (Right) Chairs for the audience	59
Figure 36 Technology Maturity Diagram	87
Figure 37 Live Presentation demonstration within VR Conference Application	88
Figure 38 VR conference application screenshots	89
Figure 39 Participants in VR Conference Experience	90
Figure 40 Left: Menu items from "Main Menu" scene. Right: VFX spatial matching to theatrical stage in XR client.	91
Figure 41 AR Theatre Pilot 1 performance from different perspectives	92
Figure 42 Assembly Demonstration in (AR) Training Assistant Use Case.....	94

Figure 43 Participants - Training Assistant Pilot 1	95
Figure 44 Status of User Requirements	97
Figure 45 User Requirements Methodology Timeline	112

PENDING APPROVAL

List of Tables

Table 1 VOXReality Components used in all three use cases.	13
Table 2 VOXReality Components used in all three use cases.	21
Table 3. Description of the use case definition template	23
Table 4. Initial requirements for VR Conference use case	25
Table 5. Initial requirements for Augmented Theatre use case	27
Table 6. Initial requirements for Training Assistant use case	29
Table 7. Participants in VRDays focus group	44
Table 8. Participants in AEF focus group	44
Table 9. Participants in HOLO focus group	44
Table 10. Total number of insights for each of the category/sub-categories	50
Table 11. Final Requirements for VR Conference use case	65
Table 12. Final Requirements for Augmented Theatre Use Case	71
Table 13. Final Requirements for Training Assistant use case	75
Table 14 Lessons Learnt from Pilot 1 (VR Conference)	91
Table 15 Lessons Learnt from Pilot 1 (Augmented Theatre)	93
Table 16 Lessons Learnt from Pilot 1 (Training Assistant)	96
Table 17 VR Conference: Completed Requirements in Pilot 1	98
Table 18 VR Conference: Out-of-scope Requirements from Pilot 1	100
Table 19 VR Conference: Open Requirements from Pilot 1	101
Table 20 VR Conference: Extended Requirements from Pilot 1	102
Table 21 Augmented Theatre: Completed Requirements from Pilot 1	103
Table 22 Augmented Theatre: Out-of-scope Requirements from Pilot 1	104
Table 23 Augmented Theater: Open Requirements from Pilot 1	105
Table 24 Extended AR Theatre user requirements	106
Table 25 Training Assistant: Completed Requirements in Pilot 1	107
Table 26 Training Assistant: Out-of-scope Requirements from Pilot 1	109
Table 27 Training Assistant: Open Requirements from Pilot 1	109
Table 28 Extended AR Training user requirements	111

Executive Summary

This deliverable provides a comprehensive overview of the series of activities, project's progress, methodologies, and outcomes over its first two years within WP2 (and correlated activities from other work packages), focusing on the description of the use cases and an extended version of user requirements. The report begins with an account of the user-centred methodology employed to gather requirements and capture the essence of the project's use cases. The sections 3, 4, 5, 6 and 7 are already coming from the previous deliverable D2.1 (Definition and Analysis of VOXReality Use Cases V1). The latest methodology, detailed in [Section 2](#), includes activities for year 1 and 2; such as weekly calls, focus groups, in-situ visits, requirements revision using key performance indicators (KPIs), interaction design, technology maturity planning, and pre-pilot and pilot studies and brief summary of data analysis of the pilots (see Figure 1). [Section 3](#) focuses on the initial weekly calls between project partners—spanning technical teams, user-cantered groups, and use case experts—establishing a shared understanding and gathering an initial set of requirements. Building on this foundation, [Section 4](#) outlines the initial set of requirements that emerged from this collaborative process. [Section 5](#) contains details of the in-situ visits conducted with use case owners. These visits, organized as interactive focus group workshops, were instrumental in refining the requirements and deepening the user-cantered partners' understanding of the specific contexts within which the use cases operate. Following this, [Section 6](#) provides a detailed definition of the three core use cases, offering concrete and specific descriptions that guide the project's development. [Section 7](#) enlists the user and technical requirements identified earlier, serving as a foundational reference for the development of the pilots. [Section 8](#) then summarizes the first round of pilots, offering insights into the initial execution, implementation, and data analysis, which collectively provide an overview of first round of the pilots. The core contributions of this deliverable are presented in [Section 9](#), which introduces the Extended Requirements (Version 2) with sub-divided tables of Achieved, Open/Pending and Opt-out requirements. This section also reflects the key-findings for lessons learned¹ and feedback received during the initial phases, providing an updated list of requirements for each use case that will inform the project's future direction. The report concludes with a summary of key findings and recommendations that will directly support the future deliverables D2.5 (Organisational preparation for VOX pilot scenarios and PRESS analysis V2) and D5.2 (Pilot planning and validation V2).

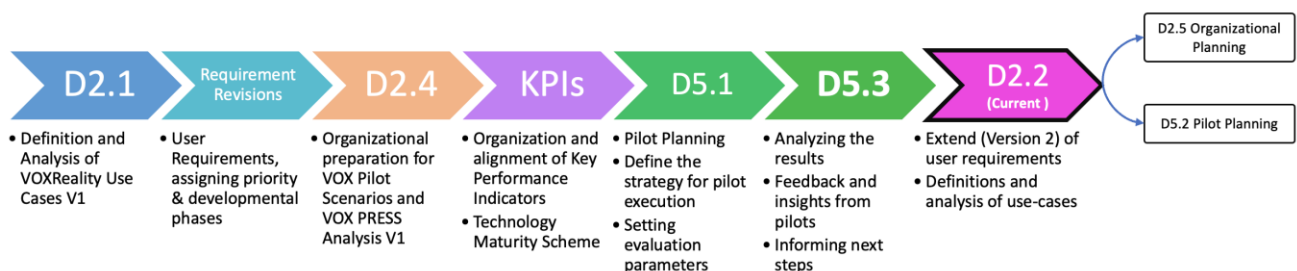


Figure 1 Flow of deliverables leading to Definition and Analysis of VOXReality Use cases

¹ Further detailed account of Lessons Learnt would be covered in WP6 and its related deliverables.

1 Introduction

1.1 VOXReality concept and approach

VOXReality is an ambitious project with the goal of facilitating and exploiting the convergence of two important technologies, Natural Language Processing (NLP) and Computer Vision (CV). Both technologies are experiencing a huge performance increase due to the emergence of data-driven methods, specifically ML and AI. CV/ML are driving the XR revolution beyond what was possible up to now, and speech-based interfaces and text-based content understanding are revolutionizing human-machine and human-human interaction. VOXReality employs an economical approach to combine these two. This NLP focused project pursues the integration of language-based and vision-based AI models with either unidirectional or bidirectional exchanges between the two modalities. Vision systems drive both AR and VR, while language understanding adds a natural way for humans to interact with the back-ends of eXtended-Reality (XR) systems or create multimodal XR experiences combining vision and sound.

The results of the Project are intended to be twofold: [a] a set of pretrained next-generation XR models combining, in various levels and ways, language and vision AI and enabling richer, more natural immersive experiences that are expected to boost XR adoption, and [b] a set of example applications to showcase the capabilities of these models in order to demonstrate innovations in various sectors. The above technologies are being validated through three use cases: VR Conferences, Augmented Theatres and Training Assistants. Following table shows VOXReality components used in each of the use-case scenarios.

Table 1 VOXReality Components used in all three use cases.

Use case Components	ASR	Neural Machine Translation	Vision language Models	Dialogue System
VR Conference	X	X	X	X
Augmented Theatres	X	X	X	
Training Assistant	X			X

1.2 Purpose of the deliverable

This deliverable stems from the first version of **D2.1 Definition and Analysis of VOXReality Use Cases V1**, which earlier provided a base analysis of the VOXReality's use cases, a list of technical and user requirements, initial information flow diagram towards the architecture of the system and the development/deployment plan. This current version provides an **updated version** of the detailed analysis of the VOXReality's use cases, with an extended list of user and technical requirements.

It also elaborates on the user-centric methodology explained in section, with Year 1 and Year 2 overview, detailing the activities, strategies, and outcomes of the activities in chronological order (*Weekly Calls, In-Situ Visits, Requirements Revision with KPIs, Interaction Design, Technology Maturity Plan, Pre-Pilots, Pilots, and Data Analysis of Pilots*) to gather the current refinement of the **Definition and Analysis of VOXReality Use Cases Version 2**.

1.3 Intended Audience

This deliverable is public in its nature and it is addressed, first of all, to the consortium, with an emphasis on the technical research team and the end-users. It is also addressed to the interested public and OpenCall participants with an interest to build over and extend VOXReality Technologies.

1.4 Structure

The document is divided into 10 sections, providing a list of new requirements and a detailed description of the use cases. First, [Section 2](#) details of the year 1 and 2 summary of activities, with the overall user-centered methodology used for gathering the requirements and capturing the essence of the use cases (*Weekly Calls, Focus Groups In-Situ Visits, Requirements Revision with KPIs, Interaction Design, Technology Maturity Plan, Pre-Pilots, Pilots, and Data Analysis of Pilots*). [Section 3](#) details the initial weekly calls between the different partners (technical, user-centered and use cases) which allowed gathering of an initial set of requirements and created a shared understanding. [Section 4](#) provides the initial set of requirements, resulting from this first phase. Later, [Section 5](#) describes the in-situ visits to the use case owners, as interactive focus group workshops, for further detailing the requirements and use cases; and providing the user-centered partners a better understanding of the context of the use cases. Based on the methodology, [Section 6](#) provides a concrete and specific definition of the three use cases. [Section 7](#) lists the previously identified users and technical requirements. [Section 8](#) provides the Summary of Pilots (First Round), summarizing the pilot results for each use case, providing insights into the first phase of pilot execution, implementation and data analysis overview. [Section 9](#) contains the core contributions of the deliverable with an Extended Requirements (Version 2) section that provides an updated list of requirements for each use case, reflecting the lessons learned and feedback received during the initial phases. The report concludes with a Conclusion section, which summarizes the entire report and providing final thoughts and recommendations next deliverables.

1.4.1 Extended Version History

The Version History (Extended) includes both the version history from the parent deliverable (D2.1) of the same origin and builds on top of it the version history for the current deliverable, documenting the complete timeline of changes and updates made to the report over time, ensuring transparency and traceability of the deliverable from Version 1.0 to Version 2.0.

1.4.2 Chapter Revisions

Since this deliverable extends from the previous deliverable D2.1 (Definition and Analysis of VOX Use Cases V1), many sections overlap and are reused to maintain the chronological readability of the work done. For the sake of transparency and ease, the following is a list of chapters and their content status (revised, new, same as previous):

- 1) [Introduction](#): New
- 2) [Methodology](#): Heavily Revised
- 3) [Details of Weekly calls \(Year 1\)](#): Same as D2.1
- 4) [Use cases initial requirements \(Year 1\)](#): Same as D2.1
- 5) [In-Situ Visits \(Year 1\)](#): Same as D2.1
- 6) [Use Cases descriptions \(Year 1\)](#): Same as D2.1
- 7) [Requirements \(Year 1\)](#): Same as D2.1
- 8) [Summary of Pilots \(Year 2\)](#): New
- 9) [Extended Requirements \(Version 2\)](#): New
- 10) [Conclusion](#): New

2 Methodology

In this Project, the consortium has agreed to adopt a user-centric approach in the requirement-gathering procedure which has now extended greatly from the existing two phases that were 1) Weekly calls and 2) In-Situ visits. The requirements were then further refined through discussions, highlighting key performance indicators. Organizational preparation took place next, defining high, medium, and low priority requirements. Interaction design and evaluation were performed, leading to the planning for pilot stages. Pilots were then executed, and feedback was collected for analysis and preparation for final evaluations. Finally, the requirements were extended and refined based on the feedback received as shown in the detailed Figure 2 below.

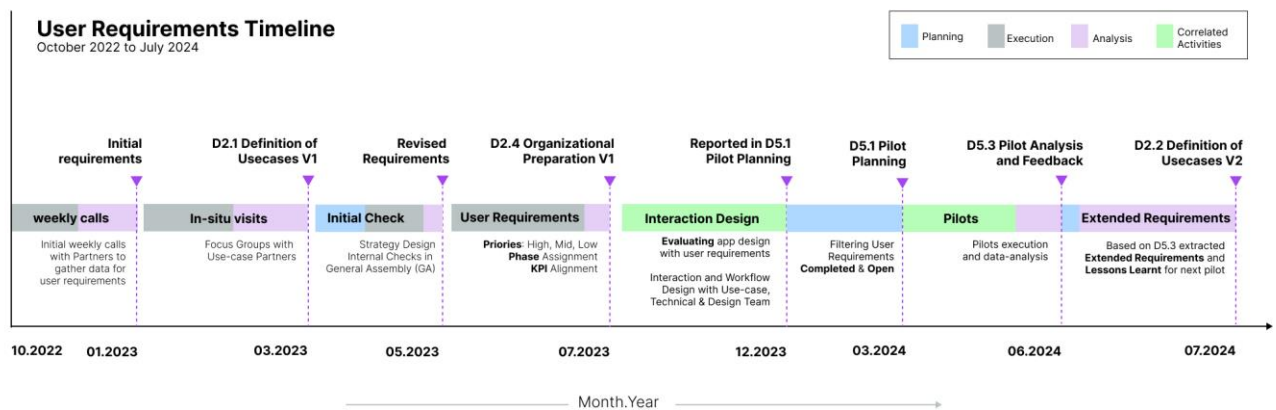


Figure 2. User Requirements Methodology Timeline

2.1 Year 1 Overview

In the first year of the VOXReality project, we laid a solid foundation for the subsequent phases by intensive data collection for requirements. Weekly calls and in-situ visits were organized, which involved dedicated focus groups comprising stakeholders from technical, use case, and design teams. These focus groups played a critical role in facilitating direct interactions and obtaining detailed insights from all involved parties. This data was analysed to identify and formulate user requirements. The evaluation process involved analysing the collected information to understand the needs and expectations of end-users, thereby ensuring that the project's development would be user-centric. The comprehensive set of user requirements derived from this evaluation laid the groundwork for the next technical development ensuring that our approach was well-informed and strategically aligned with user needs at all stages.

2.1.1 Weekly calls

Weekly calls had taken place between the use case partners and the consortium (technical and user-centered team) to scope the use cases and align viewpoints. Seven online meetings took place from October to December 2022 via Microsoft Teams platform, each of which lasted around 2 hours. At least one representative from each partner participated in the calls, and minutes were recorded, documented, and distributed to the consortium after each call to align understanding between partners.

The outcome of the meetings includes three use case documents filled by the use case partners (VRDays, AEF, HOLO), which provided an initial version of the use case with details such as scenario,

required technology, needed equipment, potential assessment protocol, and its target population. The minutes and three use case documents were analysed and turned into an initial list of user requirements, which then were used to scope and prepare the in-situ visits (see section 2.2).

2.1.2 In-Situ Visits

Based on the initial user requirements gathered during the weekly calls, extended user research was undertaken via in-situ visits to each use case partner. Each meeting was prepared as an ‘interactive focus group workshop’ where the goal was to understand the latent user needs and gather deeper insights through playful and participatory activities.

The workshops were conducted on January 27th (Amsterdam, VRDays), January 30th (Athens, AEF), and February 28th (Ismaning, HOLO) 2023. Each session lasted up to three hours, depending on the number and characteristics of participants (Figure 3).



Figure 3. Participants discussing their ideas during the VRDays, AEF and HOLO workshop

The workshops were conducted by a moderator accompanied by an assistant. Tailored presentation slides and activities sheet materials were prepared to form a creative and structured workshop (Figure 4). These materials were based on the initial list of requirements and description of the use cases. The detailed experiment setup can be found in Sections 4.1 VR Conference, 4.2 Augmented Theatre and 4.3 Training Assistant.

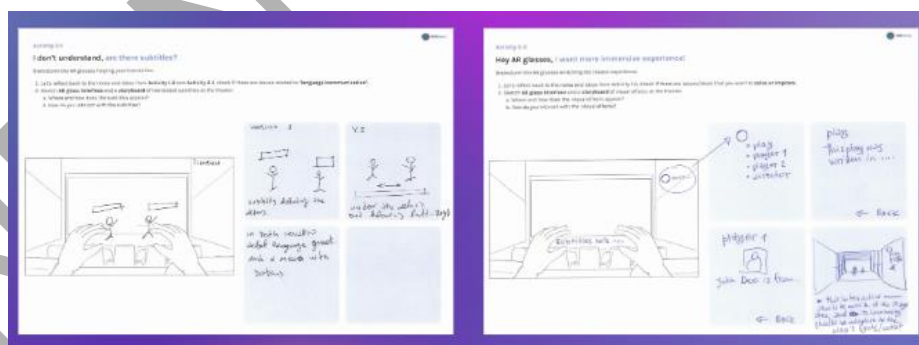


Figure 4. Results of the activity sheets from the AEF workshop

After the workshop, all the activity sheets were collected, documented, and the voice recordings of the participants were transcribed. The different forms of workshop outcomes were analysed by the user-centred research team and turned into a final list of user requirements, which can be found in Sections 6.1 VR Conference, 6.2 Augmented Theatre and 6.3 Training Assistant. These visits helped the use case partners to further reflect over their use cases, resulting in a more detailed description of them, which can be found in Section 5.

2.1.3 Revised User Requirements

We revised user requirements for each of the three use cases. The initial step in this process was to prioritize these requirements (mentioned in D2.1) based on their criticality for the successful implementation of Pilot 1. Each requirement was assigned a priority level—**high**, **medium**, or **low**—based on its importance and urgency. These terms and their context for VOXReality use cases is defined below. The full process for revision is documented in the dedicated chapter User Requirements Revision in the deliverable D2.4 Organisational preparation for VOX pilot scenarios and PRESS analysis V1. This prioritization ensured that the most crucial aspects were addressed promptly and effectively. A snapshot of the user requirements list is given in figure below whereas the complete lists of revised requirements are available in the section [Extended Requirements \(Version 2\)](#).

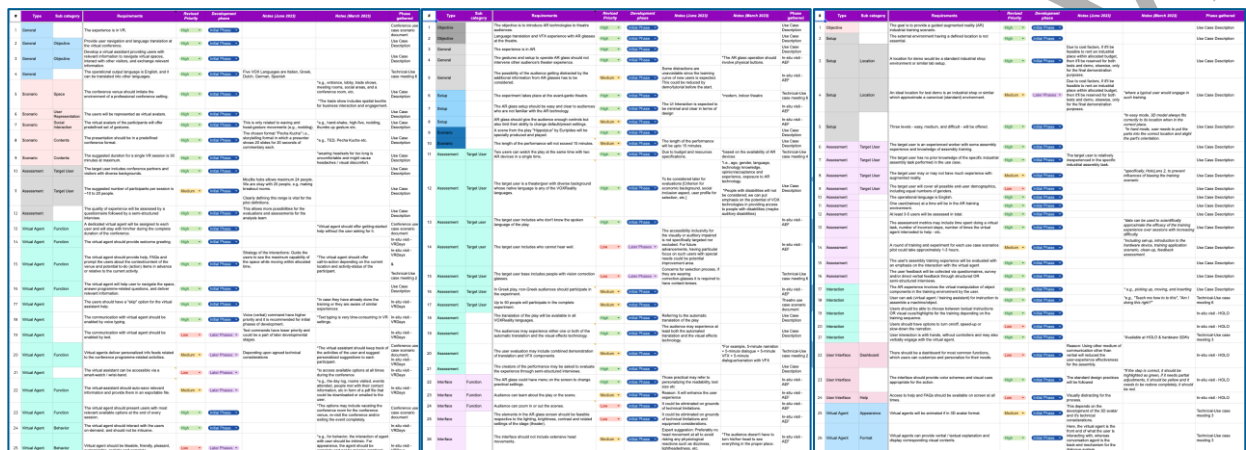


Figure 5 Screenshot of revised user-requirements for each of the use-cases.

- **High-priority** requirements, which were considered critical for the success of the project, received immediate attention and resources, ensuring that the core functionalities and essential features were implemented effectively.
- **Medium-priority** requirements, while not as critical as high-priority ones, held significant importance and contributed to the overall user experience. Allocating resources to medium-priority requirements ensured a balanced approach that catered to both core functionalities and additional enhancements.
- **Low-priority** requirements, while less critical, still provided valuable features or improvements that can be addressed in later stages of development or future iterations.

2.2 Year 2 Overview

For the second year the focus shifted towards refining and operationalizing the user requirements identified in the first year. The primary objective was to align these prioritized user requirements with the project's key performance indicators (KPIs) and the necessary conditions for the first phase of pilots. This alignment was essential for ensuring that the pilot phases were not only feasible but also capable of meeting the project's strategic goals. Detailed planning and analysis were conducted, as documented in deliverables D5.1 (Pilot Planning and Validation V1) and D5.3 (Pilot Analysis and Feedback V1), to establish a coherent framework for the pilot implementations. We'd like to refer the reader for to these two deliverables D5.1 and D5.3 for a detailed information regarding the pilots and the analysis outcomes respectively. Whereas, here we have provided a brief summary of the pilots to provide better context for subsequent sections.

Throughout the second year, significant efforts were dedicated to preparing for the execution of Pilot 1 across all three use cases—[VR Conference](#), [Augmented Theatre](#), and [Training Assistant](#). This involved rigorous planning to ensure that all technical, logistical, and operational aspects were thoroughly covered.

The technology maturity diagram presented (see Figure 6) stems from initial app design to leading towards the series of pilot phases and following by a user-study in parallel. Pilot 0 marked the initiation of internal testing, reserved for developers and consortium members, to gain preliminary insights. The subsequent phases, Pilot 1 and Pilot 2, were based on an ascending scale – starting with a smaller user base and moving to larger-scale experiments involving a more extensive user community. The user studies involving research prototypes ran in parallel to this journey, navigating the realms of both functional and non-functional technologies. This comprehensive diagram was also furnished in year 2.

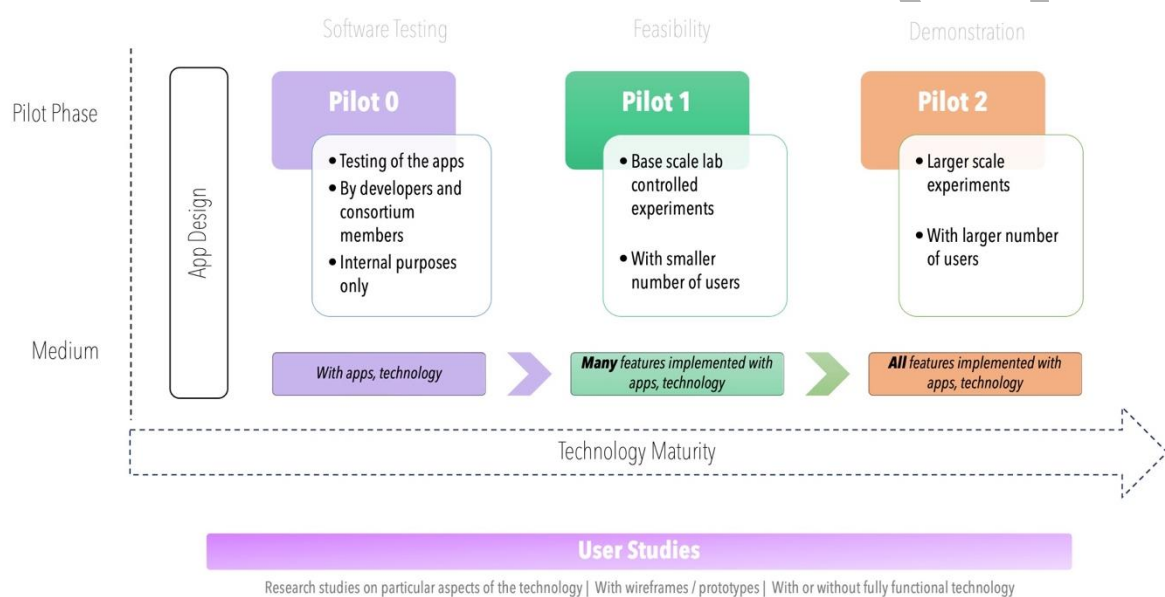


Figure 6 Technology Maturity Diagram

It explained the stage increments of technological development, offering a nuanced view of how VOXReality Technologies evolved and transformed into a fully realized and user-tested solution in XR Spaces over the course from App Design → Pilot 0 → Pilot 1, to test the initial functionalities of the models within the use case applications.

2.2.1 Interaction Design and Workflow (App Design)

In a collaborative effort between NWO-I (CWI), SYN, MAG and related Use-Case Partners conducted a series of meetings with aim to evaluate whether the application interfaces and behaviours were well-aligned with the existing user requirements. For this purpose, the teams created software flows, visual representations, and curated user interaction flows (see one example in Figure 7) on the Miro Boards to enhance the overall user experience and usability of the applications. This is documented in detail in the Chapter 3 Piloting Strategy and Plan of the internal deliverable D5.1 (Pilot Planning and Validation V1).



Figure 7 (Left) VR Conference App and (Right) AR Theatre App Interaction Flow (Source D5.1)

This series of creative meetings and discussion sessions among design and development teams of the applications, along with use-case stakeholders on board, to understand and agree on several aspects of the technologies. The focus was to mutually agree on best practices for user interactions, gather feedback, and validate design decisions depending on the user requirements. These extensive meetings proved to be helpful in making informed iterative refinement cycles based on user requirements and KPIs, ensuring the applications meet the needs and expectations of their intended users. This phase set the stage for a comprehensive design and development collaboration to optimize the usability and effectiveness of the applications prior to their wider deployment.

2.2.2 Pilot 0 Overview (Preparing for Pilot 1)

The objective of this initial piloting activity was to test the developed models and prototype applications internally as a means of preparation for testing with external users in Pilot 1. Pilot 0 was conducted with low fidelity prototypes targeting only the core features of each use case. Pilot 0 phase was intended for the internal and initial phase of software testing, where rigorous testing of the applications was conducted by developers and consortium members. This phase was primarily focused on internal purposes, emphasizing the validation and refinement of our project's software components and subsequent phases of development.

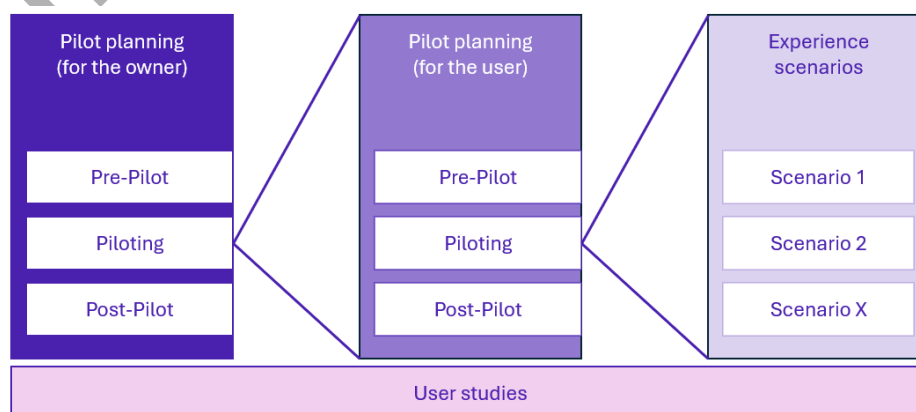


Figure 8 Pilot Implementation Plan

Pilot 0 served a number of beneficial outcomes for the project:

- Provided proof of concept to validate the AI models and XR applications in well-described settings, with adequate quality of prototypes, and with measured and verifiable correspondence to the targeted user requirements for each of the use cases.
- Actively progressed the mapping and understanding of the involved challenges, not only in terms of technological feasibility and performance but also in terms of usability and user acceptance.
- Provided user feedback from a diverse and multidisciplinary testing group—feedback which was interpreted into design and development guidelines, thus advising future work.
- It allowed the consortium to evaluate the adequacy of the risk mitigation plans and informed the design of system extensions and/or modifications to help ensure robustness and stability in the delivery of Pilot 1.
- Produced verifiable evidence material for the promotion of the VOXReality project efficiency and value, thus indirectly contributing to the attraction of interest for the Open Calls (OC+) and Exploitation plans.

2.2.3 Pilot 1 Overview



Figure 9 Pilot Demonstrations (A) VR Conference; (B) Augmented Theatre; (C) Training Assistant

The first round of pilots for the VOXReality project are based on three use cases and one experimental study (see Figure 8). The pilots—VR Conference, Augmented Theatre, and Training Assistant—along with the complementary User Study, were well-executed within realistic environments that closely simulate real-world scenarios, to ensure the outcomes are relevant and applicable for future applications. We direct the readers to the deliverable **D5.3 Pilot Analysis and Feedback V1** for detailed and well-documented results. Here we present a brief overview of the pilot results for each use-case and user study.

The **VR Conference** pilot aimed to assess the impact of AI-assisted VR technology for navigation and speech translation on conference-goers, focusing on (A) navigation technology response and (B) machine translation response. The goal was to enhance virtual conferencing experiences with the support of a virtual agent facilitating navigation in the VR space. The **Augmented Theatre pilot**, based on the Greek play Hippius, evaluated the value of AI-generated translation, AR-displayed subtitles, and VFX as perceived by theatre-goers. The **Training Assistant pilot** focused on the impact of AI-assisted augmented reality (AR) technology on training factory workers in industrial assembly tasks. This section is elaborated further in chapter [Summary of Pilots \(First Round\) \[NEW\]](#)

Table 2 VOXReality Components used in all three use cases.

Use case Components	ASR	Neural Machine Translation	Vision language Models	Dialogue System
VR Conference	X	X	X	X
Augmented Theatres	X	X	X	
Training Assistant	X			X

Each pilot was designed taking into account the specific use case and the features of VOXReality models and applications (see

Table 2) that best fit for those use cases. Subsequently, the most suitable environment to execute each pilot was studied, including defining the number of participants, setup, resources, evaluation parameters and analysis methods for gathering results in this current deliverable. This definition of a pilot framework assisted in selecting and shortlisting the most appropriate metrics and testing guidelines for each of the use-cases.

2.2.4 User Study

Additionally, an independent **User Study** compared static (fixed to gaze) and dynamic (fixed to objects) subtitle placements in a VR theatre environment. This study involved the virtual reality extension of the AR Theatre use case to test usability of static (image) and dynamic (image) captions in various viewing scenarios. After careful analysis of the existing literature and state-of-the art in the related domain aligned with the scientific interests, we choose subtitles placement and design attributes in theatre as the core of our user study in a virtually lab-controlled environment. This user study aimed to explore the optimal subtitle position integration, particularly static and dynamic in the context of a VR theatre environment.



Figure 10 (Left) Static Captions; (Right) Dynamic Captions from a participant's perspective

We used the same Greek play “Hippolytus” mention in page 57 (section 6.2) of the Augmented Theatre. The study examined subtitle placement and design attributes for VR theatre play and VFX, providing user-centric insights into subtitle design for virtual reality. These insights are useful for informing the future design and development decisions. This user study was assistive in informing many internal design-decisions for the AR Theatre application considering it is the most artistic and aesthetically particular use-case.

2.3 Conclusion

The journey from the initial identification of user requirements in Deliverable D2.1 to their refinement and implementation has been marked by systematic revisions and strategic planning (see Figure 11). Initially, a comprehensive set of user requirements was outlined. These requirements were subsequently revised and prioritized, with developmental phases assigned in Deliverable D2.4. The alignment with Key Performance Indicators (KPIs) was then detailed in Deliverable D5.1, setting the foundation for effective pilot planning. As detailed in Deliverable D5.3, the pilot analysis provided critical insights, leading to the identification of new extended requirements and the extraction of valuable lessons learned.

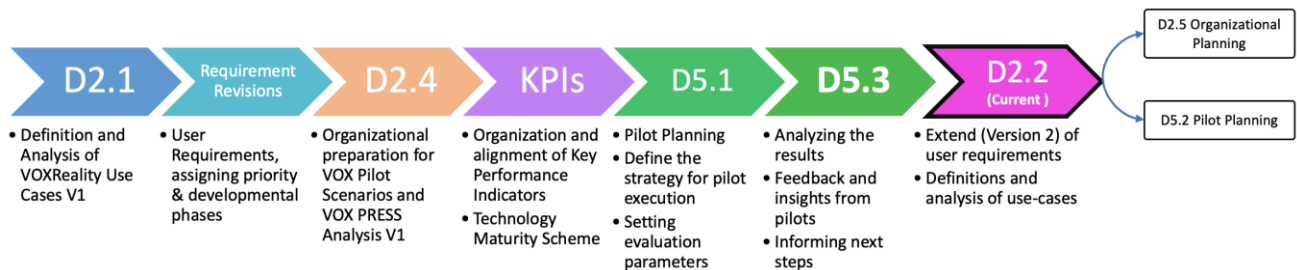


Figure 11 Flow of deliverables leading to Definition and Analysis of VOXReality Use cases

The inputs gathered from the quantitative and quantitative analysis of the pilots for each use case informed the [Extended Requirements \(Version 2\)](#). Overall, the benefits and drawbacks gained from these pilots offered a benchmark for future research and development, guiding open call projects and the second pilot phase. The pilot findings are critical for the upcoming deliverables (see Figure 11) D2.5 (Organisational preparation for VOX pilot scenarios and PRESS analysis V2) and D5.2 (Pilot planning and validation V2), eventually supporting the planning of the second phase of pilots and Open Calls for the project.

3 Details of Weekly calls (Year 1)

The VOXReality consortium comprises teams specialised in user-centred research, technical research and end-users. As expected, each team possesses unique terminologies, perspectives on use cases, and levels of familiarity with the expertise of other teams. In the beginning of the Project the consortium scheduled weekly online calls led by the Scientific Coordinator to synchronize the teams' viewpoints, enhance their comprehension of each other's terminologies, facilitate familiarity with the diverse disciplines represented within the consortium and drive the discussions towards user requirements, technical requirements and formal use case definitions.

3.1 Weekly calls *modus operandi*

The consortium conducted the 2-hour alignment meetings via Microsoft Teams platform, utilizing the project-specific channel. Partners were instructed to ensure at least one representative from each partner attended and participated in the discussions. A total of seven online meetings took place on October 21st, October 31st, November 7th, November 14th, November 21st, December 5th, and December 12th, 2022. Simultaneously with each meeting, the minutes were recorded and collected and the minutes from the previous meeting were distributed to the consortium before the next meeting.

The first two meetings focused on comprehending the partners' perspectives regarding the use cases. During the initial meeting held on October 21st, the discussions revolved around the questions of the use case designers (VRDAYS, AEF and HOLO) related to the utilization of technological components in each use case. On October 31st, the second meeting was devoted to the technical partners' point of view, where the questions related to the execution of each use case experiment were addressed by the use case designers.

The structure of the remaining meetings consisted of three equally separated sessions dedicated to each use case individually and a final shorter session to discuss technical or use case related topics that would affect all use cases and technical components. Furthermore, these meetings were used to take initial decisions on certain technical and use case related aspects by taking advantage of the participation of all members of VOXReality consortium.

The outcome of these meetings was intended to drive the discussions towards user, functional and technical requirements. To facilitate these discussions in the next steps (described in chapter 4), a use case definition template has been prepared to produce complete documents where all details could be collaboratively recorded regarding each use case. Table 3 presents the sections of the mentioned template document with the description of each section. The meetings held in December were used to complete the contribution of each partner to these use case documents and the meetings contributed to the clarification of any misunderstandings and missing details.

Table 3. Description of the use case definition template

Section	Explanation
Title	Title of the use case
Description	High-level description and goals of the overall use case
Scenario	Description of the specific scenario. (How long is it? Description of a user's experience in the scenario. Maybe a step-by-step explanation of one walkthrough.)
Technologies	The software components that will be used. (How and where they should be utilized? What should they provide considering their and the proposal's limitations?)

Equipment	Hardware configurations of the scenario such as laptop, headset, mobile, controller, microphone, speaker etc.
External Context	The surrounding sentences are internal context. What could be the external ones in the scenario? For example, slides, a dictionary of words, description of the scene, agents' knowledge base e.g., the venue, the manual of a training.
Assessment Protocol	The criteria (what will be assessed) and tools (questionnaires etc.) that will be used in the scenario. Where will the scenario take place? The number of participants in each session. Total number of participants in the scenario to be reached.
Target Population	The target users/audience of the scenario. Demographic background to be considered e.g., age, gender, language, technology acceptance, technology experience. Languages to be used. If possible/needed, exclusion criteria (maybe eyeglasses for headsets?).

3.2 Results

The alignment meetings were held until the start of the focus group meetings, resulting in four significant outcomes:

- The partners involved in the Project were able to familiarize themselves with the various disciplines represented within the consortium.
- All members were informed about the terminologies that will be utilized throughout the Project.
- There was a consolidation of viewpoints regarding the design, implementation, and execution of each use case.
- The use cases were defined and documented using a common structure.

The outcomes of these meetings enabled the consortium to easily transition from internal meetings to focus groups involving other stakeholders, such as end-users and application designers of each use case; including conference organizers for VR conference use case, trainers for Training-Assistant use case and directors and actors for Augmented-Theatre use case. Furthermore, the common language, mutually agreed terminologies and the formal representation of use case definitions used among the VOXReality researchers provided a common voice to reach these stakeholders easily. Finally, the use case definitions provided the foundation of the functional and technical requirements as well as enabled the researchers to initiate the implementation of independent technological components of VOXReality.

3.3 Coordination with technical partners

Following the initial definition of the use cases, the technical partners met in person on November 30th, in a general meeting held during Immersive Tech Week in Rotterdam, The Netherlands, to refine and consolidate the technology requirements based on the input from the use-case owners. During the meeting, the partners identified a common architecture to be used in all use cases, and decided on its modularity and on the timeline of its implementation. For each use case, they identified which modules should be implemented and/or activated, what input would be provided (both as direct input and as context), and what is the envisioned output. The outcome of the meeting was a roadmap for technical partners to proceed with their development, which informed the subsequent weekly meetings.

4 Use cases initial requirements (Year 1)

This section lists the initial requirements for each VOXReality's use cases (VR Conference, Augmented Theatre and Training Assistant), based on the analysis of the weekly calls MoM and on the initial description of the use cases provided by the use case owners.

4.1 VR Conference

The goal of the VR conference use case is to provide user navigation and language translation at a virtual conference. The conference will only be held in VR, not hybrid, but a stream of video can happen in the virtual environment. The virtual venue consists of entrances, lobby areas, trade shows, meeting rooms, social areas and a conference/plenary room. The language transition will mainly support the 1 to many conferences setting.

User interaction

Users will be represented by virtual avatars in the conference venue. User avatars will be programmed to deliver visual feedback complementing available real-time multilingual translation and captioning. Users can ask questions to the virtual agents during the navigation, and also allowed to engage with speakers and guests during Q&A sessions.

Virtual agents

Virtual agents are virtual only setting, and help navigate the user by informing about the context and contents of the rooms. Everything happening on the way to the rooms could potentially be considered a part of the virtual agents (digital navigation assistants). Virtual agents can answer questions from the users and deliver personalized info related to the conference programme-related activities. A dedicated virtual agent will engage with users when they log into the virtual conference venue.

Table 4. Initial requirements for VR Conference use case

Type	Requirements	Notes	Phase gathered
Environment	The experience is in VR.		conference use case scenario document
Objective	User navigation and language translation at the virtual conference.		
Interaction	VR controllers will provide navigation, menu selections and haptic interactions.		conference use case scenario document
Setup	The conference is only in virtual, not hybrid.		technical-use case meeting 1
Setup	A stream of video can happen in the virtual environment.		technical-use case meeting 1
Input	Contextual visual information will be used for path generation (navigation)		technical-use case meeting 6
Input	Voice communication will be used for path generation (navigation) will use.		technical-use case meeting 6
Input	Instruction generation (instruction) will be based on the dialog (text) and visual information.		technical-use case meeting 6
Source contents	The length of the speech should not be longer than 30'.		conference use case scenario document

Source contents	The speech is based on a predefined conference format.	*ideally, ex: TED, Pecha Kucha	conference use case scenario document
Output	Language output is in English, but can be translated.		technical-use case meeting 6
Output	All outputs from the models will be text.		technical-use case meeting 3
Venue design	The virtual venue consists of entrances, lobby areas, trade shows, meeting rooms, social areas and a conference/plenary room.		conference use case scenario document
User interaction	Users will be represented by virtual avatars in the venue.		conference use case scenario document
User interaction	User avatars will be programmed to deliver visual feedback complementing available real-time multilingual translation and captioning.		conference use case scenario document
User interaction	Users are allowed to engage with speakers and guests during Q&A session.		conference use case scenario document
User interaction	In the conference room, users can choose between two in-room view options.	*The options include a full screen and conference room setting.	conference use case scenario document
User interaction	Both in-room views (full screen and conference room) include on-demand virtual assistance and real-time multilingual translation and captioning.		conference use case scenario document
User interaction	At the end of every conference session, users will be presented with all available options.	*The options include vacating the conference room for the conference venue, re-visit the conference and/or exiting the event completely.	conference use case scenario document
User interaction	Users can ask the agent questions during the navigation.	*ex) "Please describe the scene for me.", "What is in this room?", "Is there an empty chair in the room?" etc.	technical-use case meeting 3
Translation	Language translation assists the 1 to many conferences setting.	*1: many is the main focus use case scenario	1:1 meeting - CWI
Virtual agents	Virtual agents help navigate the user.		technical-use case meeting 2
Virtual agents	Virtual agents are virtual only settings.	*no physical agents	technical-use case meeting 1
Virtual agents	Virtual agents inform the users about the context/content of the room.		technical-use case meeting 2
Virtual agents	Everything happening on the way to the rooms could potentially be considered a part of the digital navigation assistant.		technical-use case meeting 3
Virtual agents	Virtual agent can integrate scene description from the Visual Language models.		technical-use case meeting 3
Virtual agents	Virtual agents answer questions from the users.		conference use case scenario document

Virtual agents	A dedicated virtual agent will engage with users when they log into the virtual conference venue.		conference use case scenario document
Virtual agents	Virtual agents deliver personalized info feeds related to the conference programme-related activities.		conference use case scenario document
Virtual agents	The agent will act as the dialogue generator and the component communicating with the user.		technical-use case meeting 6

4.2 Augmented Theatre

The goal of the AR Theatre use-case is to provide language translation and VFX experience with AR glasses during the play. The target audience is Greek or Greek play international audiences with different age, gender, technology acceptance and using different language. The experiment will take place at the avant-garde theatre, modern, indoor, instead of the outdoor ancient theatre.

Play design

The play, which will be used for the AR experience, is an ancient tragedy with the length of 10-12 minutes. The play is in English and not dialogue heavy. Maximum 2-3 actors are on the stage at the same time to avoid the technical and design complexity. Additionally, one narrator, who is not on the stage, could be included.

VFX

VFX will be triggered by certain words or phrases from the play. A narrator could be the person who triggers the VFX. The style of the visual effects should be artistically relevant to the opera, and delivered in a high quality rather than having ambitious challenges. The location of the visual effects should not exceed the range of the stage.

Evaluation

Two users will watch the play at the same time due to the availability of the AR devices, and up to 50 people will participate in the experiment. The evaluation may include combined demonstration – e.g., 5-minute narration + 5-minute dialogue, 5-minute VFX + 5-minute dialogue/narration with VFX. While user experience is the main evaluation of the Project, the creator's experience of how their work is presented through the technological medium (VFX) could also be evaluated.

Table 5. Initial requirements for Augmented Theatre use case

Type	Requirements	Notes	Phase gathered
Environment	The experience is in AR.	Environment	
Objective	Language translation and VFX experience with AR glasses at the theatre.	Objective	
Setup	The experiment takes place at the avant-garde theatre	<i>Modern, indoor theatre</i>	technical-use case meeting 6
Source contents	Contextual information (summary, description etc) will be provided instead of the entire script.		technical-use case meeting 1
Device	The simulator sickness limits for the selected AR glasses should be checked.		technical-use case meeting 2
Target user	The audience is Greek or Greek play international audience		technical-use case meeting 1

Target user	The audience has different technology acceptance or experience		technical-use case meeting 1
Target user	The audience is different age, gender and use different language	<i>*People with disabilities will not be considered; we can put emphasis on the potential of VOX technologies in providing access to people with disabilities (maybe auditory disabilities)</i>	technical-use case meeting 1
Play design	The play is in English.		technical-use case meeting 1
Play design	The play is not dialogue heavy		technical-use case meeting 2
Play design	Maximum 2-3 actors are on stage at the same time.		technical-use case meeting 2
Play design	A narrator, who is not on stage, could be included.		technical-use case meeting 2
Play design	The play includes a short dialogue on stage.		technical-use case meeting 2
Play design	The length of the play will be 10-12 minutes.		technical-use case meeting 4
Play design	Play is an ancient tragedy	<i>*planned to be</i>	technical-use case meeting 6
VFX	Certain words or phrases will trigger VFX.		technical-use case meeting 1
VFX	A narrator could be the person who triggers the VFX		technical-use case meeting 2
VFX	Style of the visual effects should be artistically relevant to the opera	<i>Not too cartoonish</i>	technical-use case meeting 6
VFX	Location of the visual effects should not exceed the range of the stage		technical-use case meeting 6
VFX	The visual effect should be delivered in a high quality rather than having ambitious challenges.		technical-use case meeting 6
Evaluation	Up to 50 people will participate in the experiment		theatre use case scenario document
Evaluation	At least two users watch the play at the same time with AR devices	<i>*based on the availability of AR devices)</i>	technical-use case meeting 4
Evaluation	The evaluation should focus on the experience of users with the technological features of VOX.		technical-use case meeting 2
Evaluation	Users should evaluate VFX.		technical-use case meeting 2
Evaluation	The evaluation may include combined demonstration.	<i>For example, 5-minute narration + 5-minute dialogue + 5-minute VFX + 5-minute dialogue/narration with VFX</i>	technical-use case meeting 2
Evaluation	Creator experience could also be evaluated based on rehearsal performance including VFX.		technical-use case meeting 2

Evaluation	Creators could evaluate how their work is presented through the technological medium regarding VFX.		technical-use case meeting 2
Evaluation	The testing of the Project should not influence the normal production of the theatre play.		technical-use case meeting 6
Evaluation	People without glasses should not be neglected from the experience.		technical-use case meeting 6
Evaluation	VFX and the subtitles should not have delays on AR screen.		technical-use case meeting 6

4.3 Training Assistant

The goal of the training use case is to provide machine assembly training using AR glasses. Users will assemble machines with virtual tools, screws, and objects. User interaction is set with hands, without extra controllers. The target user is a beginner trainee with some knowledge about the scenario.

Training design

Users can select different difficulty modes during the assembly experience. In easy mode, 3D model will always fit correctly to its location when in the correct place; in hard mode, user needs to assemble the parts into the correct location and orientation.

User interaction

Users communicate with the agents through voice and hand input. Users can ask for instruction to assemble a machine/object.

Virtual agents

Virtual agents use the visual scenes as context. Virtual agents can provide verbal/textual explanation and display corresponding visual contents. They can show and demonstrate the movements and also direct the user towards objects when necessary. If the virtual agents are in 3D avatar format, they will be animated.

Evaluation

Evaluation may include the speed of the assembly with and without the virtual agent.

Table 6. Initial requirements for Training Assistant use case

Type	Requirements	Notes	Phase gathered
Environment	The experience is in AR.		technical-use case meeting 4
Objective	Machine assembly training using AR glasses.		technical-use case meeting 4
Interaction	User interaction is set with hands, without controllers.	*Available at HOLO & hardware SDKs	technical-use case meeting 3
Setup	The place to assemble machine is a physical grey base.		technical-use case meeting 6
Setup	Tools, screws, objects are virtual.		technical-use case meeting 6
Setup	The visual instructions, step-by-step guide can be pre-set.		technical-use case meeting 3

Input	Hand motions and audio input from the user will be used.	*HOLO collect the hand movement videos for additional contents	technical-use case meeting 4
Target user	The target user is a beginner trainee with some knowledge about the scenario.		technical-use case meeting 3
Training design	User can select different difficulty mode.		technical-use case meeting 6
Training design	In easy mode, 3D model always fits correctly to its location when in the correct place		technical-use case meeting 6
Training design	In hard mode, user needs to assemble the parts into the correct location and align the part's orientation.		technical-use case meeting 6
User interaction	Users communicate with the agents through voice and controller input.		technical-use case meeting 1
User interaction	User can ask for instruction to assemble a machine/object.	*e.g., "Teach me how to this", "Am I doing this right?"	technical-use case meeting 6
Virtual agents	Virtual agents will be animated if in 3D avatar format.		technical-use case meeting 3
Virtual agents	Virtual agents direct the user towards objects that are necessary for the application.	*e.g., certain stance for assembly of a machine part	technical-use case meeting 3
Virtual agents	Virtual agents can show and demonstrate movements that are necessary for the application.		technical-use case meeting 3
Virtual agents	Virtual agents can provide verbal/textual explanation and display corresponding visual contents.		technical-use case meeting 3
Virtual agents	Virtual agents use the visual scene as context.		technical-use case meeting 4
Evaluation	Evaluation may include speed of the assembly with and without the agent.	*Planned to be	technical-use case meeting 3

5 In-Situ Visits (Year 1)

The focus group workshops were conducted with use case owners to understand the latent user needs and gather deeper insights and requirements. Each workshop was conducted at the lead partner location, specifically on January 27th in Amsterdam (VRDays), January 30th in Athens (AEF), and February 28th in Ismaning (HOLO). The following sections describe the detailed experiment setup and procedure (5.1 Methodology) and provide the results of the sessions (5.2 Data Analysis and Results) respectively.

5.1 Methodology

The followed methodology allowed interactive focus group workshops [3][4] of 2-3 hours, where the partners and users were invited to share and brainstorm ideas for a new experience with VR/AR experience. Every workshop was conducted in English by one moderator and one assistant. Activity sheets and presentation slides were accompanied during the workshop as supplementary materials. These materials were developed based on the initial list of requirements. The meetings were recorded for later analysis. No preparatory tasks were required.

5.1.1 VR Conference

The goal of the focus group workshop was to:

- 1) Understand the users' needs and organizers' wants on virtual agents and language translation at the VR conference; and
- 2) brainstorm ideas for the role and design of virtual agents and language translation at the VR conference.

The workshop was conducted on January 27th 2023 from 14:00 to 17:00 at Spaces Herengracht², in Amsterdam. In total 6 participants attended the workshop, including VRDays conference organizers and end-users who had previous experience attending the conference.

The structure of the workshop was as follows:

- Introduction;
- Part A: My current conference experience;
- Part B: Bringing virtual agents and language translation;
- Part C: Designing a future VR conference;
- Conclusion.



Figure 12. Participants working on activities during the VR Conference (focus-group) workshop

² <https://www.spacesworks.com/nl/amsterdam-nl/herengracht/>

Introduction

In the beginning of the workshop, the moderator shortly introduced the study and provided an activity workbook sheet, a pen, and stickers. Participants, moderator and assistant introduced themselves with a warm up activity, information including name, position, expertise in VR, etc.

Part A: My current VR conference experience

Part A activities were designed to understand the current VR conference experiences and to find a design space for adopting language translation and virtual agent solutions.

Activity 1-1: My experience with a VR conference

The first activity was to ask participants to recall the time they attended a VR conference. Using the activity sheet, participants reflected and wrote down their experience of participating the VR conference on:

- When and where was it?
- What was the purpose or goal of attending the conference?
- What did you enjoy the most?
- What is the one thing that could have been improved?

Activity 1-2: My conference experience

Following activity 1-1, participants were asked to write down or visualize the activities they did on the timeline when they attended the VR conference, from entering the VR venue to leaving it. After the drawing, they marked positive and negative moments with green and red stickers, marking at least 3 moments that they thought could be improved/assisted/enriched in some way, assuming they have a superpower, with blue stickers. After individuals completed the timeline drawing, participants shared and discussed the similarities and differences of the timelines as a pair.

Activity 2-1: My experience with virtual agents

The moderator first introduced the concept of virtual agents to help understand the possibilities with visual and audio examples about existing virtual agents (see Figure 13).

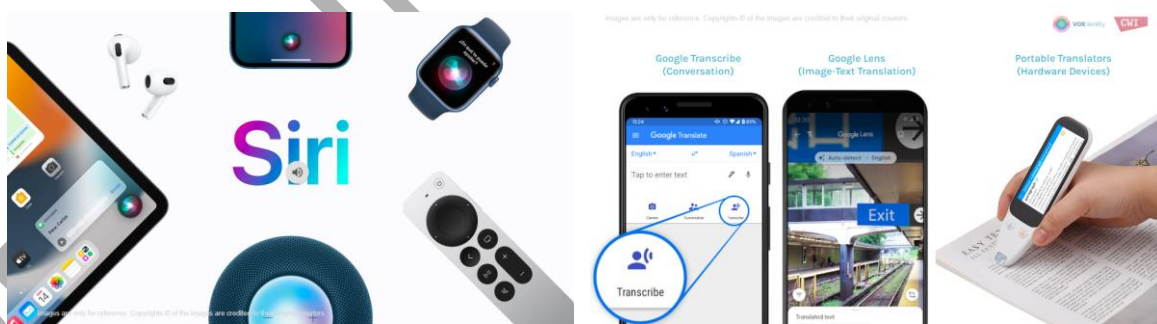


Figure 13. Virtual agent (left) and language (right) translation examples

The participants were later asked to reflect on their experience with virtual agents. They selected one situation and briefly sketched and explained the situation, with supplementary questions (Figure 14):

- Why did you use/had to interact with the virtual agent at that time?
- How did it help solve your problem?
- What was the pain point?

Activity 2-1

My experience with virtual agents

Reflect on your experience with virtual agents

1. Select one situation and briefly sketch / describe it.

a. Why did you use/had to interact with the virtual agent at that time?

b. How did it **help** solve your problem?

c. What was the **pain point**?

Activity 2-2

My experience with language translation

Reflect on your experience with language translation models

1. Select one situation and briefly sketch / describe it.

a. Why did you use/had to interact with the language translation at that time?

b. How did it **help** solve your problem?

c. What was the **pain point**?

Figure 14. Activity 2-1, 2-2 sheet for the VR Conference use case

Activity 2-2: My experience with language translation

The moderator introduced the concept of language translation models to help understand the possibilities with visual examples.

Participants were asked to reflect on their experience with language translation models. They selected one situation and briefly sketched and explained the situation, with supplementary questions:

Why did you use/had to interact with the language translation model at that time?

How did it help solve your problem?

What was the pain point?

Part B: Bringing superpowers: virtual agents and language translation

Part B activities were designed to help the participants brainstorm ideas on how to apply language translation and virtual agents to solve issues raised in Part A.

Activity 3-1: [Virtual assistant], please help me navigate!

For activity 3, the moderator provided a specific scenario of attending a VR conference and asked the participants to imagine they were faced with the situation. In the first scenario, one had just entered the VR conference venue and wants to attend the 'XR haptics' session in Room A-2, without knowing the room location (Figure 15).

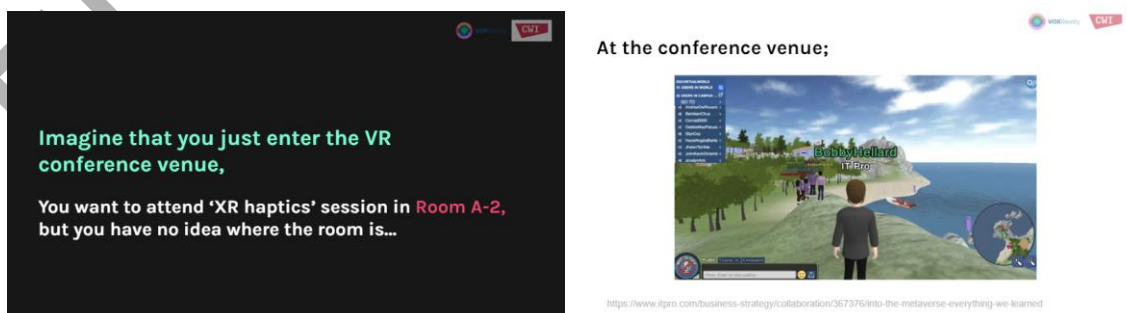


Figure 15. Scenario case and supplementary image on navigating in VR conference

The participants were asked to brainstorm on virtual agents helping their navigation, considering:

- How does the virtual agent look like?
- How do you interact with the virtual agent?

The brainstorming sheet included eight blank boxes – where to sketch the imaginary scenario or screen interface. After finishing the drawing, each participant shared and explained their brainstormed idea.

Activity 3-1

[Virtual assistant], please help me navigate!

Brainstorm the virtual agents helping your navigation.

1. Let's reflect back to the notes and ideas from Activity 1-2 and Activity 2-1; check if there are issues/ideas related to 'navigation'.
 2. Sketch **two storyboards (or interface)** of a virtual assistant helping you how to navigate at the VR conference.
 a. How does the virtual agent **look** like?
 b. How do you **interact** with the virtual agent?

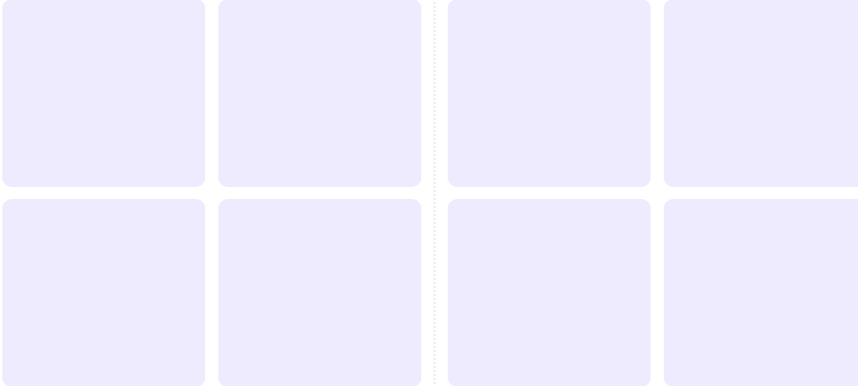


Figure 16. Activity 3-1: Brainstorming sheet for drawing ideal VR conference navigation scenario

Activity 3-2: I do not understand, are there subtitles?

The second scenario was about the subtitles. Participants arrived at Room A-2 and the speaker started to give a presentation. However, the speaker started to talk in French, a language they do not understand, and they seek the language translation. The participants were asked to brainstorm on virtual agents helping the translation, considering:

- Where and how does the subtitles appear?
- How do you interact with the subtitles?

The brainstorming sheet with the same format as activity 3-1 was provided. After finishing the drawing, each participant shared and explained their idea (Figure 17, left). The moderator and assistant placed the activity sheets with the brainstormed ideas on the whiteboard. Participants were provided with four stickers and had to vote for the two best ideas for each scenario, namely navigation and language translation (Figure 17, right).



Figure 17. (Left) participants sharing their brainstormed ideas, (Right) participants voting for the two best ideas for each scenario

Part C: Designing a future VR conference

Part C activities were designed to come together and design the ideal future VR conference as a group based on the accumulated ideas from the previous parts.

Activity 4-1: Four scenarios

The moderator introduced four different forms of virtual assistants that could possibly be integrated into VR conference platforms (Figure 18). Participants were asked to write the order that they would like to use to start exploring the VR space, by selecting a minimum of 1 and a maximum of 5, including their own idea as the fifth option.

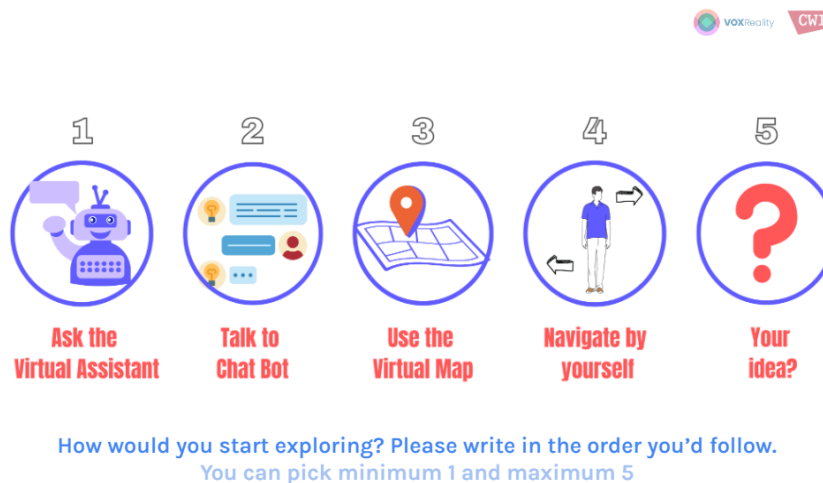


Figure 18. Different forms of virtual assistant in VR conference

Activity 4-2: Brain-drawing – designing an ideal VR conference

Participants were divided into two groups to work on two scenarios: navigation and language translation. They were provided a blank sheet to start the brain-drawing activity, with the steps:

- As a group, sketch an ideal interaction on one topic:
 - How does it look like?
 - How to interact?
- Write the list of components it should have and not have;
- Swap the sketch between two groups and continue on drawing. Repeat once more.
- Each group pitch the generated idea and together reflect on the list of requirements.

Conclusion

The moderator concluded the workshop by appreciating the participants for their active participation and input. The moderator and assistant collected the activity sheets filled with responses and drawing from participants. After the workshop, the researchers analysed activity sheets, along with the voice recordings, and presented the process (5.2 Data Analysis and Results).

5.1.2 Augmented Theatre

The goal of the focus group workshop with AEF was to:

- 1) Understand the needs and desires of both users' and organizers for subtitles (Language translation) and VFX while watching a theatre play;
- 2) And brainstorm ideas for the UI and interaction method between the audience and AR glasses.

The workshop (see Figure 19) was conducted on January 30th 2023 from 10:00 to 13:00 at the Theatre Peiraos 260³, in Athens, Greece. In total 6 people attended the workshop, with both AEF organizers and theatre enthusiasts with diverse profile.



Figure 19. The moderator introducing the workshop to the participants in focus group workshop

The structure of the workshop included the following parts:

- Introduction;
- Part A: My current theatre experience;
- Part B: Bringing superpowers: subtitles and visual effects with AR;
- Conclusion, wrap-up.

Introduction

At the start of the experiment, the moderator gave the participants a short introduction to the study and provided an activity workbook sheet, pen, and stickers. Participants, moderator and assistant shortly introduced themselves with a warm up activity, sharing their name, position, expertise in theatre, expertise in AR, favourite play, etc.

Part A: My current theatre experience

Part A activities were designed to understand the current theatre experiences of the audiences and to find a design space for integrating AR solutions. There were three activities.

Activity 1-1: My experience with a theatre

The first activity was to ask participants to recall the time they watched a theatre play. Using the activity sheet, participants reflected and wrote down their experience of visiting a theatre on:

- When and where was it

³ <https://aefestival.gr/venues/peiraos-260/?lang=en>

- Which language that you watched?
- What did you enjoy the most?
- What is the one thing that could have been improved?

Activity 1-2: My visit to the theatre

On an extension to the first activity, participants were asked to write down or visualize the activities they did on the timeline when they visited a theatre, from entering the theatre to leaving it (see Figure 20). After the drawing, they marked positive and negative moments with green and red stickers. They marked at least 3 moments that they thought could be improved/assisted/richer in some way, assuming they have a superpower, with brown stickers. After individuals completed the timeline drawing, participants shared and discussed the similarities and differences of the timelines as a pair.

Activity 1-2 - 10 min

My visit to the theater

Recall the time you went to watch a theater play

1. Write down / visualize the activities you did on the timeline, from entering the theater leaving the theater.
2. Use **green** stickers to mark the positive moments.
3. Use **red** stickers to mark the negative moments.
4. Use **brown** stickers to mark **at least 3 moments** you think can be improved/assisted/richer in some way, assuming you have a superpower.



Figure 20. Activity sheet 1-2: My visit to the theatre (drawing timeline of a visit to the theatre)

Activity 2-1: My experience with language/communication

The next activity was to reflect on **language/understanding/communication issues** during the theatre play. Participants were asked to select one situation and briefly sketch and explain the situation, with supplementary questions:

- Why was it a problem to you?
- Did you take any action to solve the problem? Why? Or why not?
- Will there be a better way to solve the problem? How?

Part B: Bringing virtual agents and language translation

Part B activities were designed to help participants learn about AR and brainstorm ideas on how to apply AR technology to solve issues raised in Part 1.

The moderator first introduced the concept of AR and played the video on “The Future of Augmented Reality: 10 Awesome Use cases⁴” to easy the understanding of AR possibilities with visual examples.

⁴ <https://youtu.be/WxzcD04rwc8>

Activity 3-1: I do not understand, are there subtitles?

For activity 3, the moderator provided a specific scenario of watching a theatre and asked the participants to imagine they were faced with the situation. In the first scenario they were seated in the theatre wearing AR glasses, and two actors on the stage started to speak Korean – a language that the participant could not understand (Figure 21).

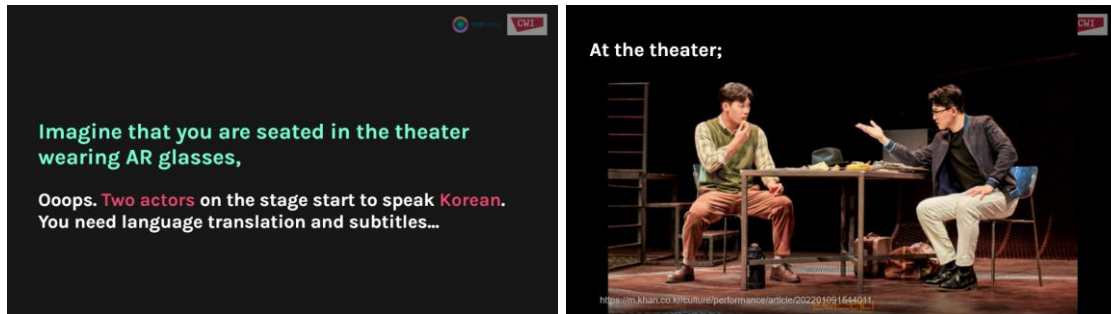


Figure 21. Scenario case and supplementary image on watching a play in a foreign language

The participants were asked to brainstorm on use of AR glasses to assist in the translation, considering:

- Where and how do the subtitles appear?
- How do you interact with the subtitles?

The brainstorming sheet consisted of an illustration of theatre – where to draw the UI of AR glasses – and four boxes – where to sketch the scenario (Figure 22). After finishing the drawing, the moderator took a photo of each idea and shared it on the screen. Each participant shared and explained their idea.



Figure 22. Activity sheet 3-1: Brainstorming sheet for drawing AR glass UI and scenarios for the augmented theatre use case

Activity 3-2: Hey AR glasses, I want more immersive experience!

The second scenario was that they were at the theatre wearing AR glasses, and any text or visual effects could assist them or enrich their experiences. The participants were asked to brainstorm the AR glasses enriching the theatre experience, considering:

- Where and how does the visual effects appear?

- How do you interact with the visual effects?

The brainstorming sheet with the same format as activity 3-1 was provided. After finishing the drawing, the moderator took a photo of each idea and shared it on the screen. Each participant shared and explained their idea.

Conclusion

The moderator concluded the workshop by appreciating the participants for their active participation and input. The activity sheets filled in by the participants were collected for researchers' analysis, together with the meeting recordings. The process can be found later in this document, in section 5.2 Data Analysis and Results.

5.1.3 Training Assistant

The goal of the focus group workshop with HOLO was to:

- 1) Understand the requirements for providing machine assembly training via virtual agents using AR glasses;
- 2) And brainstorm ideas for the role and design of virtual agents using AR glasses.

The workshop (Figure 23) was conducted on February 28th 2023 from 15:00 to 18:00 at Holo-Light offices⁵ in Ismaning, Germany. Two people – a chief project manager and a scientific researcher – participated in the workshop as company representatives.



Figure 23. The moderator introducing the workshop to the participants

The structure of the workshop is as follows:

- Introduction;
- Part A: Current machine assembly experience;
- Part B: Brining superpower: voice enabled virtual agent;
- Conclusion.

⁵ <https://holo-light.com/contact/>

Introduction

At the start of the experiment, the moderator gave the participants a short introduction to the study and provided an activity workbook sheet, a pen, and stickers. Participants, moderator and assistant introduced themselves with a warm up activity, information including name, position, expertise in AR, etc.

Part A: Current machine assembly experience

Part A activities were designed to understand the current machine assembly experiences of the users and to find a design space for integrating AR solutions. However, as the participant themselves were not the target user of the product and service, the moderator provided an imaginary persona to help the ideation process (Figure 24). Participants were asked to imagine him/herself as James when working on the activities.

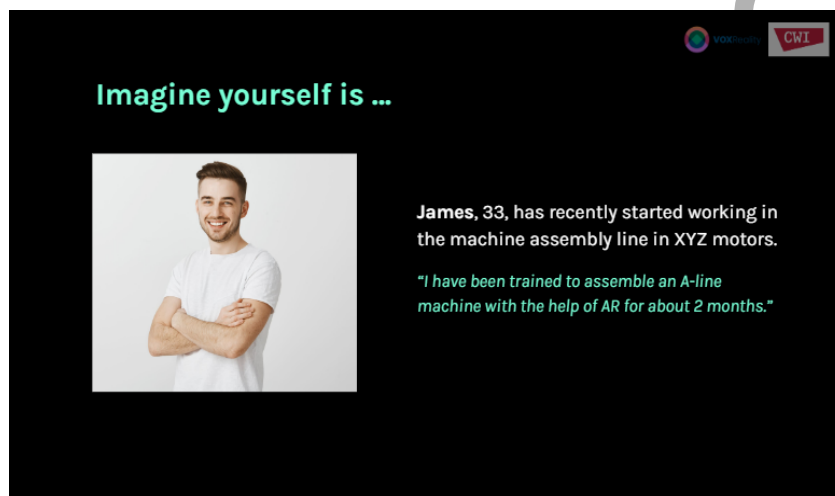


Figure 24. A user persona, James, with a beginner level machine assembly experience

Activity 1-1: My machine assembly experience

For the first activity, the participants were asked to recall the time they assembled machines with AR glasses. Then they had to write down or visualize the activities they did on the timeline, from starting to completing the machine assembly. After the drawing, they marked at least 3 moments that they thought could be improved/assisted/richer in some way, assuming they have a superpower, with brown stickers. After individuals completed the timeline drawing, participants shared and discussed the similarities and differences of the timelines as a pair.

Activity 2-1: I think there should be a better way...?!

The next activity was to select one situation from marked moments from activity 1-1; and briefly sketch and explain the situation, with supplementary questions:

- Why was it a problem to you?
- Did you take any action to solve the problem? Why? Or why not?
- Will there be a better way to solve the problem? How?

Part B: Brining superpower: voice enabled virtual agent

Part B activities were designed to help the participants learn about voice enabled virtual agent and brainstorm ideas on how to apply the technology to solve issues raised in Part A.

The moderator first introduced the concept of voiced enabled virtual agents to help understand the possibilities with visual and audio examples. After, the moderator provided a specific scenario and asked participants to imagine that they were faced with the situation, again, assuming they were James (Figure 25).

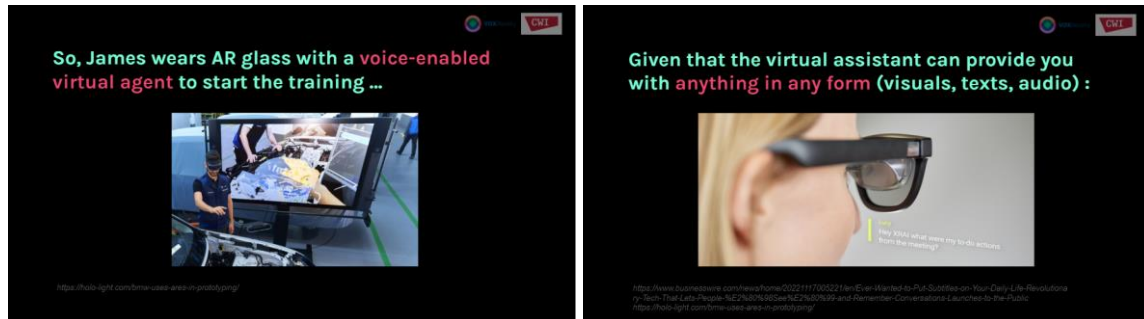


Figure 25. Scenario case and supplementary image on machine assembly experience with AR glasses

Activity 3-1: [Before] So, how do I start the assembly training?

The first scenario was that they wore the AR glasses to start the training but were unclear on how to initiate or operate the training with the AR glasses. The participants were asked to brainstorm starting experience with a virtual agent, considering:

- How would you start the virtual assistant helping your experience?
- What feedback do you expect to see/hear/get from the virtual assistant?

The brainstorming sheet consisted of an illustration of a car with the bonnet opened – where to draw the UI of the AR glasses – and four boxes – where to sketch the scenario (Figure 26). After finishing the drawing, each participant shared and explained their brainstormed idea.

Activity 3-1
So, how do I start the assembly training?
 Brainstorm the starting experience with a virtual agent.

1. Sketch **AR glass interface** and a **storyboard** of starting experience.
 - a. How would you **start** the virtual assistant helping your experience?
 - b. **What feedback** do you expect to see/hear/get from the virtual assistant?

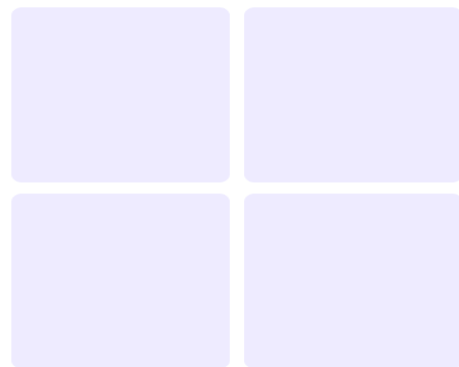
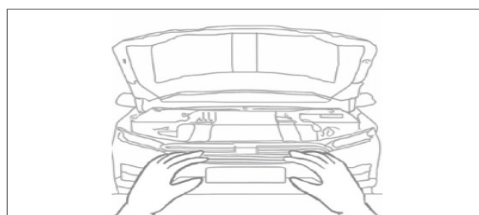


Figure 26. Activity sheet 3-1: Braining storming sheet for ideal AR UI and scenario during machine assembly training

Activity 3-2: [During] Please guide me to the next step!

The second scenario was that they had followed until step 5 of the machine assembly procedure, but now they were stuck and unsure what is the next step. The participants were asked to brainstorm the guidance experience with a virtual agent, considering:

- What action will you take?
- What feedback do you expect to see/hear/get from the virtual assistant?

The brainstorming sheet with the same format as activity 3-1 was provided. After finishing the drawing, each participant shared and explained their brainstormed idea.

Activity 3-3: [Feedback] Hey agent, how was my performance today?

The third scenario was that they had assembled a complete machine, however, they wanted feedback on the entire process so that they can reflect on today's performance. The participants were asked to brainstorm the ending experience with feedback, considering:

- What kind of feedback do you want to receive?
- What form should the feedback be?

The brainstorming sheet with the same format as activity 3-1 was provided. After finishing the drawing, each participant shared and explained their brainstormed idea.

Conclusion

The moderator concluded the workshop by appreciating the participants for their active participation and input. The moderator and assistant collected the activity sheets filled with responses and drawing from participants. After the workshop, the researchers analysed activity sheets, along with the voice recordings, and presented the process (5.2 Data Analysis and Results).

5.2 Data Analysis and Results

This section presents a comprehensive analysis of the qualitative data collected from the user-focused group workshops, highlighting the key themes and insights that emerged from the discussions. The focused groups successfully provided deeper insights from various workshop activities including some in (See Figure 27) about the responses of participants on user-needs and wants on AR/VR agents, translation mechanisms and subtitles in the AR/VR use-cases. The aim for each of the focused group's workshop were:

VR Conference (VRDays):

- 1) Understand the users' needs and organizers' wants on virtual agents and language translation at the VR conference;
- 2) And brainstorm ideas for the role and design of virtual agents and language translation at the VR conference.

Augmented Theatre (AEF):

- 1) Understand the users' needs and organizers' wants on subtitles (Language translation) and VFX while watching a theatre play;
- 2) And brainstorm ideas for the UI and interaction method between the audience and AR glasses.

Training Assistant (HOLO):

- 1) Understand the requirements for providing machine assembly training via virtual agents using AR glasses;
- 2) And brainstorm ideas for the role and design of virtual agents using AR glasses.

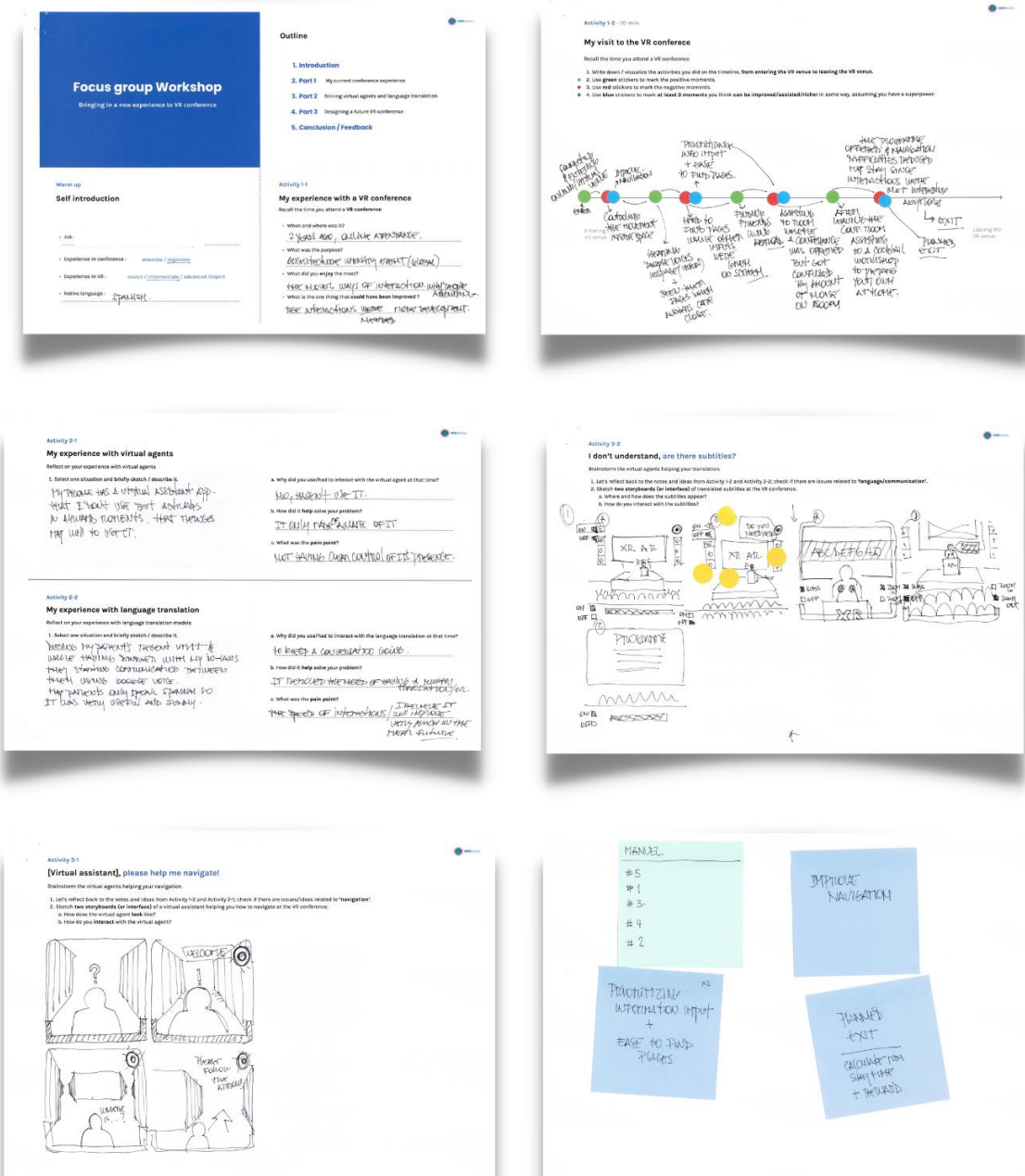


Figure 27. Selected user-activities scans from the data worksheets filled by participants

5.2.1 Participants

The success of virtual conferences largely depends on the needs and preferences of the users who attend them. In order to gather the maximum number of inputs, our focus group brought together a diverse group of participants who had different backgrounds, interests, and experiences with AR/VR technology. By bringing these participants together, we aimed to gather insights and perspectives that could improve the design and implementation of our use-case within the VOXReality project. This section describes the participants in the three focused group workshops, including their demographics, job type, experience in conference/theatre, and prior experience with AR/VR technology. Details of the participants can be found in following Table 5, Table 6 and Table 7.

Table 7. Participants in VRDays focus group

User #	Job	Experience in conference	Experience in VR	Native Language
1	Marketing Manager	Attendee, Organizer	Intermediate	Dutch
2	Director	Organizer	Expert	Dutch
3	Head of Production	Organizer	Intermediate	Spanish
4	Digital Marketing Specialist	Organizer	Intermediate	Hindi
5	Head of Partnerships and Business Development	Organizer	Advanced	English / French
6	Marketing Coordinator	Organizer	Advanced	Hindi

Table 8. Participants in AEF focus group

User #	Job Title	Experience in theatre	Experience in AR	Native Language
1	Retired School Teacher	Intermediate	Novice	Greek
2	Project Manager	Intermediate	Novice	Greek
3	Creative Labourer	Expert	Novice	Italian
4	Actress, Director	Expert	Novice	Greek
5	Composer, Soundtrack Artist, Curator of Mixed Media Projects	Expert	Intermediate	Greek
6	Musician, Guitar Teacher	Intermediate	Novice	Greek

Table 9. Participants in HOLO focus group

User #	Job	Experience in VR	Native Language
1	Chief Project Manager	Expert	German
2	Researcher, Writer	Novice	English

5.2.2 Type of Data

From each of the sessions, a significant amount of data was collected, including:

Text:

- Worksheets with questions and answers referring to the user's introduction, and level of expertise on the VR/AR and related domain(s);
- Their past experiences and anecdotes regarding virtual conferences and events;
- Descriptions of their preferences and interaction details.

Images/drawings:

- Timeline drawings visualizing their experience from start to finish as participants in the given scenarios / activity roles;
- Visual markers (stickers) and sticky notes indicating various user-choices during the activities.

Audio

- Recordings of the interactive discussion sessions, which were transcribed for easier analysis.

5.2.3 Requirements Analysis Methodology

We used Thematic Analysis [6] to analyse the collected qualitative data from the focus groups. It is a method of analysing qualitative data that involves mining of patterns or themes within the data and organizing those themes into categories. It involves reviewing the data multiple times to get a sense of its content, generating labels or tags called codes [7] that reflect relevant features within the data, grouping those codes into potential themes, and then refining and naming those themes in a way that accurately reflects their meaning and relevance to the scientific goal. Due to the wide applications [8] of thematic analysis in contexts of social sciences research and content analysis of media or literature, we used this method to process the acquired data. We identified and assessed the patterns and themes within our data from the workshops and their results are compiled into tabular categories in [Chapter 7](#). We conclude that the accumulated results strongly support the aims and objectives of the focused groups of VR Conference, Augmented Theatre and Training Assistant.

5.2.4 Data Analysis and Classification

This section enlists the comprehensive classification of the qualitative data collected from VOXReality's user-focused workshops, highlighting the major categories and themes that emerged from the discussions. After collecting the data from the participants, it was critical to organize and interpret the data into meaningful knowledge. For this purpose, we classified each dataset into themes and categories defined below, which allowed to identify patterns and relationships in the data, leading to insights and recommendations for the design and development partners in the VOXReality team.

User Journey (Tags)

- **Positive Moments:** The useful or memorable instances in the AR/VR experiences for the user.
- **Negative Moments:** The disturbing or poor instances in the AR/VR experiences for the user.
- **Could be improved (moments):** The instances that were not significantly bad, but could be improved or fixed for a better user-experience.

User Requirements (Tags)

- **New function(s) / feature(s):** Suggestions for new functionalities in the respective application.
- **Behaviour / action / pattern:** A set of actionable practices agreed to be useful by the users.
- **Requirement(s):** Technical or functional requirements for the system(s).
- **Added value:** Features or experiences that provided additional value to the overall experience.
- **Raised issue / concern:** Considerations for the issues or problems to be addressed.
- **Interface (suggestions):** UI ideas for the applications / platforms.
- **Interaction:** User interaction practices for intuitive immersive exploration.



Figure 28. Insights Canvas for VR Conference (VRDays) Use-case



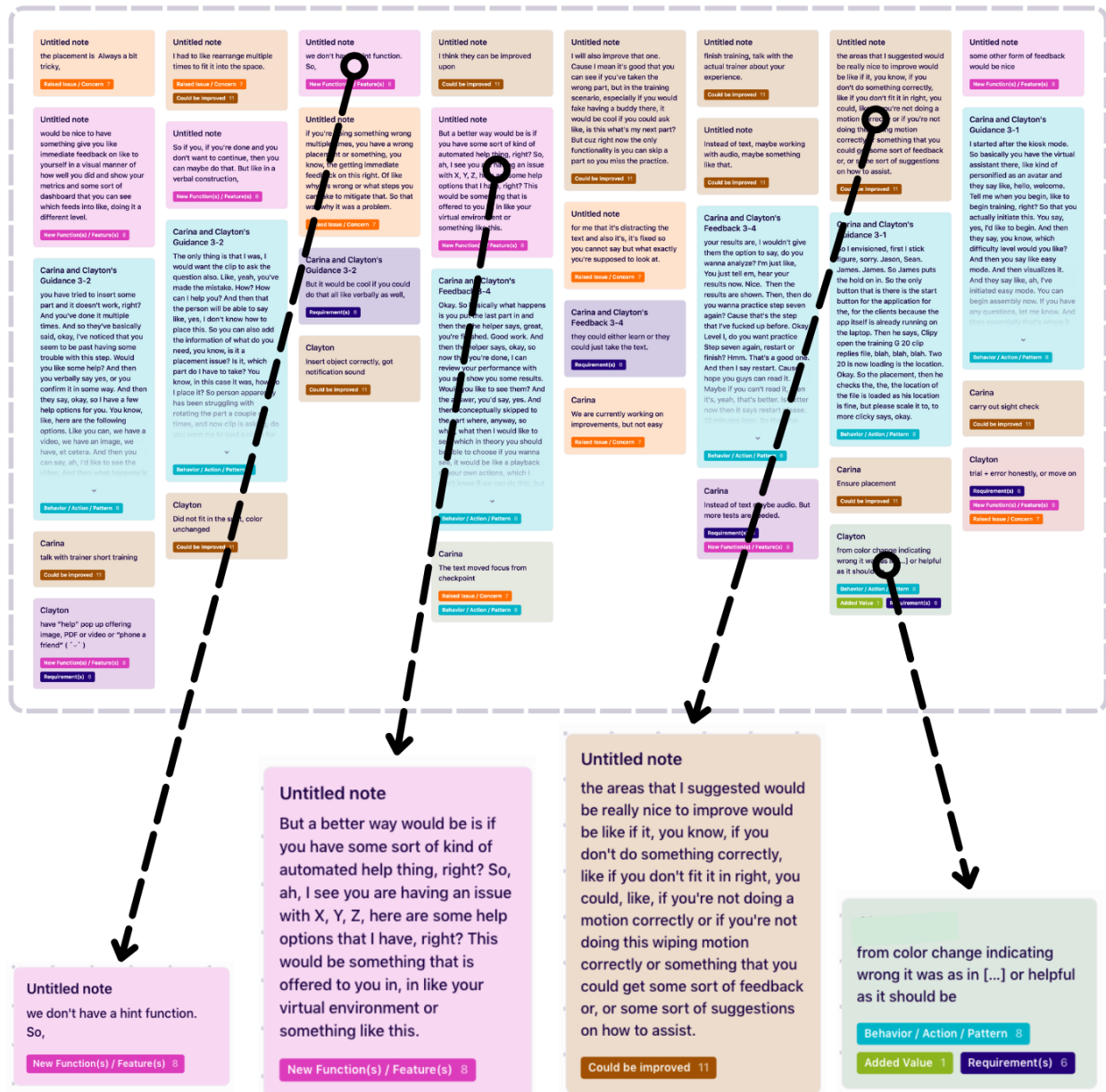


Figure 30. Insights Canvas for Training Assistant (HOLO) Use-case

Table 10 summarizes the number of insights for each category and theme that emerged from the analysis of the data. Each category represents the key considerations discussed by the participants in response to dedicated tasks, questions and group activities.

Table 10. Total number of insights for each of the category/sub-categories

#	Categories	VR Conference (VRDays)		Augmented Theatre (Athens Festival)		Training Assistant (HOLO)
1	Positive Moments	20		21		N/A
2	Negative Moments	25		12		N/A
3	Could be improved	25		22		11
Sub-themes		Virtual Assistant	Subtitles	VFX	Subtitles	Training (AR)
4	New function(s) / feature(s)	40	27	16	15	08
5	Behaviour / action / pattern	17	5	4	8	08
6	Requirement(s)	27	28	18	19	06
7	Added value	21	-	2	-	01
8	Raised issue / concern	10	15	11	1	07
9	Interface (suggestions)	08	29	-	-	-
10	Interaction	05	-	-	-	-

5.2.5 Results

By analysing and synthesising the data, we revealed important themes and patterns, as well as gained a better understanding of the attitudes and behaviours of the participants in augmented and VR applications of all the project use-cases.

The next paragraphs provide insights about the user-experience:

- **Intuitive and easy-to-use interface:** The AR/VR application should have an intuitive and user-friendly interface that allows attendees to navigate and interact with the environment and other attendees easily. It should be easy to access the platforms for virtual conference program, theatre experience or virtual training session, while supporting a quick learning curve for all three use-cases.
- **Comfortable and immersive experience:** VR and AR experiences can be overwhelming and cause motion sickness for some users. Therefore, the application(s) should be designed to create a comfortable and immersive experience, with well-balanced audio and visuals, and smooth movements. Also, timings of the content presented should suit the considerations of wearing headsets for long durations.
- **Interactive and engaging content:** The application should offer a range of interactive and engaging content to keep attendees/viewers/trainees interested and involved throughout the conference/theatre/training sessions.

- **Personalization and customization:** The virtual agent should be able to personalise its responses and recommendations based on attendee preferences and behaviour. This may include recommending sessions or events based on an attendee's past interests (if applicable), or offering personalised assistance based on an attendee's needs.
- **Clear and concise communication:** In all use-cases, (virtual conference room, theatre and training session), very clear and concise communication is crucial. The virtual environment should support high-quality voice and video communication, with minimal latency and interruptions. Also, ability to turn on/off main audio with translation audio and noises from surroundings should be included to make experience more user friendly.
- **Intuitive navigation and controls:** The UI should be intuitive and easy to use, with simple navigation controls. Users should be able to move around, adjust their position (if needed), and interact with other virtual objects and controls with ease.
- **Customizable avatars:** Attendees should be able to customize their avatars to represent their personalities and preferences. This can help foster a sense of community and create a more engaging experience. Many users reported lack of engagement in their experience due to fixed or stereotypical representations of avatars.
- **Realistic environment:** The conference room should be designed to create a realistic environment that immerses attendees in the virtual experience. With precise care on size of the modalities, including screen size, height and placement that gives naturalistic dimensions to users in virtual and augmented environments. Since some users reported problems with too high or low viewing angles of virtual objects and scenes was quite uncomfortable.
- **Multimodal communication:** The virtual agents should be able to communicate with attendees using multiple modes of communication, including voice, text, and visual cues. This will help ensure that attendees can interact with the agent in a way that is most comfortable and convenient for them. Many users suggested that virtual agents should be helpful but not intrusive.
- **Assistance and guidance:** The virtual agent should be able to provide attendees with intuitive assistance and guidance throughout the conference/theatre/training without being asked several times. This may include proactive answering to questions about the conference/theatre/training schedule, providing directions to specific locations or sessions, or helping users connect with other attendees/viewers/trainers (if required).

The next paragraphs provide insights about the expected interfaces for all use cases:

- **Platform Interface:** The AR/VR application(s) interface should display the most frequently used features e.g., conference program and schedule, controls for the theatrical play, or training demos etc to be easily visible so that users can easily view and interact with the systems.
- **Navigation and movement:** The interface should include intuitive navigation controls that allow attendees to move around the environment with ease, e.g., walking, fast run, flying etc.

- **Avatar customization:** The interface should allow attendees to customize their avatars, such as changing their appearance, selecting different outfits, or adding personal touches. Some users suggested various avatar representations such as a wrist/hand band (like smartwatch), sun or cloud character, blobby (cartoony character for virtual assistance) etc.
- **Non-intrusive Communication:** The interface should include tools for communication, such as voice chat, text chat, or virtual hand gestures, that allow attendees to interact with each other in real-time but without being forced or intrusive in the overall experience.
- **Virtual objects:** The interface should include virtual objects that attendees can interact with, such as whiteboards, presentation screens, or interactive models. These objects can be used for presentations, display screen, discussions, or interactive activities.
- **Personalization options:** The interface should allow attendees to personalize their experience, such as adjusting the audio and visual settings, selecting their preferred language or time-zone, or customizing their profile.
- **Help and support:** The interface should include tools for help and support, such as a quick demo, short tutorials, knowledge base, FAQs, etc.

The above recommendations are derived after analysis of the available data and suggestions by the participants of the workshops; however, their feasibility or conformance with the project's goals has not been validated. Therefore, their implementation is subjective to the feasibility in accordance with technical requirements. We report them here for completeness, and to provide guidance for third party developers in furthering the development of the VOXReality software.

6 Use Cases descriptions (Year 1)

In this section, we describe the scenario, technologies, equipment, external contexts, assessment protocol and target population for each of the use case of VR Conference, Augmented Theatre and Training Assistant respectively.

6.1 VR Conference

This scenario objective is to provide virtual conferencing environments with real-time multilingual translation and captioning, and to develop a virtual assistant providing users with relevant information to navigate virtual spaces, interact with other visitors, and receive relevant information.

6.1.1 Scenario

The ideal scenario should emulate the most recognizable features of a professional conference setting: Conference rooms, social spaces (cafes and lounges), business dedicated areas (trade show stands and business meeting areas). In addition, all information delivered to users within the virtual venue should be structured, prioritizing venue navigation, programme information and business-related interactions.

The length of the scenario should be 30 minutes max, ideally based on a predefined conference format, ex: TED, Pecha Kucha, etc. Minimum scenario times are defined by the full range of actions needed to validate the scenario successfully.

The virtual conferencing user case is a virtual door-to-door, assisted experience, including real-time multilingual translation and captioning assisting in conferences, virtual trade shows and social areas. Once users log into the virtual conference venue, a dedicated virtual assistant will engage with them. This virtual assistant is programmed to facilitate navigation, trained to answer questions and conditioned to deliver info feeds related to the conference programme-related activities.

All users of the VCE scenario will benefit from real-time translation and captioning service during conferences (1:X – one-to-many interactions), social interactions and business exchanges (1:1 one-to-one interactions). Once inside the virtual venue, users will be represented by virtual avatars. These avatars will be programmed to deliver visual feedback complementing available real-time multilingual translation and captioning.

Entering the venue

Once users log in to the virtual conference space/environment, they will face recognizable spatial features common to conference/event spaces: Entrances, lobby areas, trade shows, meeting rooms, social areas and ultimately, a conference/plenary room. From the user perspective, the virtual conference space will provide a door-to-door experience. Once at the virtual conference spaces, users can find their way around it thanks to a virtual assistant, who will answer questions, provide wayfinding guidance and inform visitors about relevant conference programme information.

Social interactions

Social interactions within virtual conference spaces are defined as 1:1 interaction. While an informal setting may produce unstructured conversation patterns, business interactions may follow a more structured communication and predefined vocabulary.



Figure 31. Social Area - Immersive Tech Week 2022 (©SnapBoys.nl)

In addition to real-time translation and captioning, virtual conference users will use avatars to provide a recognizable feature in the virtual space.

Trade shows

As part of the virtual conference space, a section dedicated to business interactions will be included: a Trade show where exhibitors and their spatial representations (stand/booth) will be allowed to engage with users/visitors via direct marketing on a visual level, or messages mediated by the personal assistant and the user's preferences.



Figure 32. Tradeshow - Immersive Tech Week 2022 (© JaynoBrk)

Into the conference rooms

Once users access the conference room area, they may choose between two in-room view options: Full screen and conference room setting. Both in-room views include on-demand virtual assistance and real-time multilingual translation and captioning.

Since this is a 1:X (one to many) kinds of interaction, users will be at the receiving end most of the time. However, exceptions may be allowed when Q&A opportunities are available, allowing the audience to engage with speakers/guests. Third parties/agents will enable these interactions.

At the end of every conference session, users will be presented with all available options, including vacating the conference room for the conference venue, re-visit the conference and/or exiting the event completely.



Figure 33. Main conference - Immersive Tech Week 2022 (© JaynoBrk)

A step-by-step walkthrough of the use case includes:

- 1) The user logs into the VCE with the assistance of VR headsets.
- 2) Once connected to VCE, users will be assigned a virtual assistant. This virtual assistant will help users to navigate the space, answer programme-related questions, and deliver relevant information feeds from social and business activities.
- 3) Each virtual assistant will present users with information on available spaces and interactions, such as social and business-related options. Among the options available in the VCE social, exhibitor stands and business meeting areas.
- 4) Users may choose between Full screen and Conference room settings once they have entered any conference rooms. Interactions with speakers and/or presenters may be available.
- 5) Once conferences finish, users can opt to vacate the conference room for the conference venue, where further interactions are possible, re-visit the conference (pre-recorded) and/or exit the event completely.

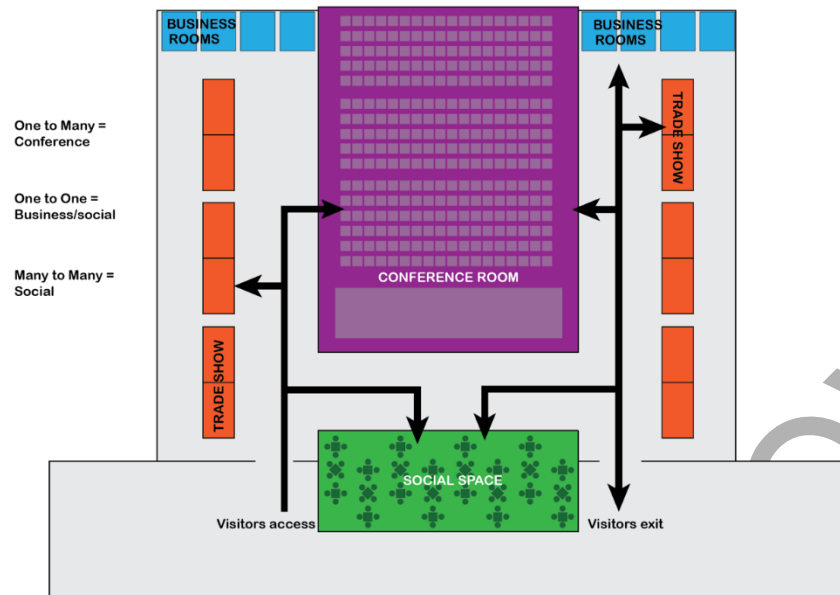


Figure 34. Virtual Conference space – Conceptual floorplan (©VRDays)

6.1.2 Technologies

The Virtual Conference case will utilize the ASR component developed by VOXReality to generate textual counterparts of the speech of users interacting with each other, users interacting with a digital assistant and speakers in the keynote speech.

The textual information acquired from the ASR will be utilized in the NMT component for automatic translations in the interaction cases (keynote speech and 1-1 interaction). NMT component will also make use of contextual information such as terminology used in the conference and abstract of the keynote speech.

Furthermore, the contextual information will also be included from the outputs of the visual language component that provides descriptions of the scene the user sees and the description of the events in a room that the user is in. Apart from the scene description, the VL component will be also capable of answering simple questions regarding various properties or spatial relations between the visible entities via dialog, using the DA. The textual information acquired from ASR will also be utilized in the DA component where a dialogue system helps the users to navigate through the virtual venue. DA component will facilitate a dialogue between the user and the corresponding navigation instructions and descriptions generated by VL component using English as a mediation language, hence, if the user speaks in a different language than English, the NMT model will provide the English translation for DA to generate its responses. The user will be able to ask questions and keep a continues dialogue going with the agent while navigating through the venue. The response generated from DA will be translated to the language of the user from English and off-the-shelf text to speech tools are envisioned to be utilized to generate an audio response in the desired language.

For navigation purposes the DA will make use of predefined map of venue.

6.1.3 Equipment

This user case supports its functioning on VR technology (VR Headsets), allowing users interactions via an in-headset microphone. In addition, VR controllers will provide navigation, menu selections and haptic interactions. The VR application will be accessible via a web-based application using a VR-ready computer (windows-based PC), enabling a live conferencing component.

6.1.4 External Context

External contexts considered for the VR Conferences use case are conference presentation slides, exhibitors' promotional content, TED data sets, conference presentation recordings, social interactions, technical word dictionaries, industry-related papers, industry-related press releases, among others that may be considered relevant.

6.1.5 Assessment Protocol

A single session will invite less than 20 conference goers with different backgrounds. The questionnaires will be used to assess the quality of the experience, followed by a semi-structured interview to get qualitative insights. Dutch and English (Dutch translated to English) will be used.

6.1.6 Target Population

The target population includes conference partners and visitors with different backgrounds.

6.2 Augmented Theatre

This use case aims to introduce AR technologies to theatre audiences. A special performance will be produced for the purpose of VOXReality. A small audience will be invited to experience AR applications in theatre. Specifically, this use case will create an AR experience for the audience via providing AI-generated and by incorporating visual effects in the performance.

6.2.1 Scenario

Up to 50 theatre goers, split in small groups, will arrive at the venue in assigned timeslots. Each group, upon arrival, will be informed about the Project, the technologies and the performance by AEF personnel and informational materials. They will sign the relevant consent forms before being led to the theatre where they will be given a short presentation of the use of the AR equipment prior to watching the performance. Technical staff will assist the audience in wearing and familiarising with the equipment.

A scene of the play "Hippolytus" by Euripides will be performed. The duration of the performance will not exceed 15 minutes. Afterwards the audience will be asked to fill out a questionnaire regarding the performance and the use of AR technologies. The duration of each groups experience, including pre- and post-performance activities, should be no more than 30 minutes.

It is yet to decide if all members of the audience will experience both technologies used, the automatic translation and the visual effects, or if there will be members of the audience without AR equipment. The creators of the performance may also be asked to evaluate the experience, through semi structured interviews [5].

6.2.2 Technologies

The use case will utilize the ASR component developed by the VOXReality consortium to generate textual counterparts of the speech of actors on stage and if the play includes a narrator, the speech of the narrator. The textual information acquired from the ASR will be utilized in the NMT component for automatic translations. NMT component will also make use of contextual information such as director notes, summary of the play, scenography and description of the characters and the theme of the play. Furthermore, the contextual information will also be included from the outputs of the visual language component. The VL component will use as input video from the viewpoint of the audience and periodically output the description of the scene in textual form.

6.2.3 Equipment

The AR equipment to be used as well as other specialised equipment is yet to be determined. The purchase costs of the equipment will be analysed against budget allowance and project needs, without compromising the expected results.

6.2.4 External Context

The external context that will be provided is the following:

- The summary of the ancient Greek play “Hippolytus” by Euripides;
- Mythological context regarding the characters of the tragedy;
- The summary of the scene, and of the events preceding and following;
- Description of the play’s themes, historical context, character relationships;
- Information regarding scenography and kinesiology: information on the set, costumes and movement of actors on the stage;
- Director notes on the scene;
- Music used, if any.

6.2.5 Assessment Protocol

The use case will take place at the Athens Epidaurus Festival’s Peireos 260 venue located at Peireos Str. Athens, Greece. The total participants, 50 theatre goers, will be split into smaller groups for each performance. The size of each session’s group will be determined based on the AR equipment that is to be used. The audience will be selected so as to represent demographic diversity, considering age, gender, language, technological knowledge, opinion/acceptance and experience. Specific demographic targets will be set later on during the project implementation.

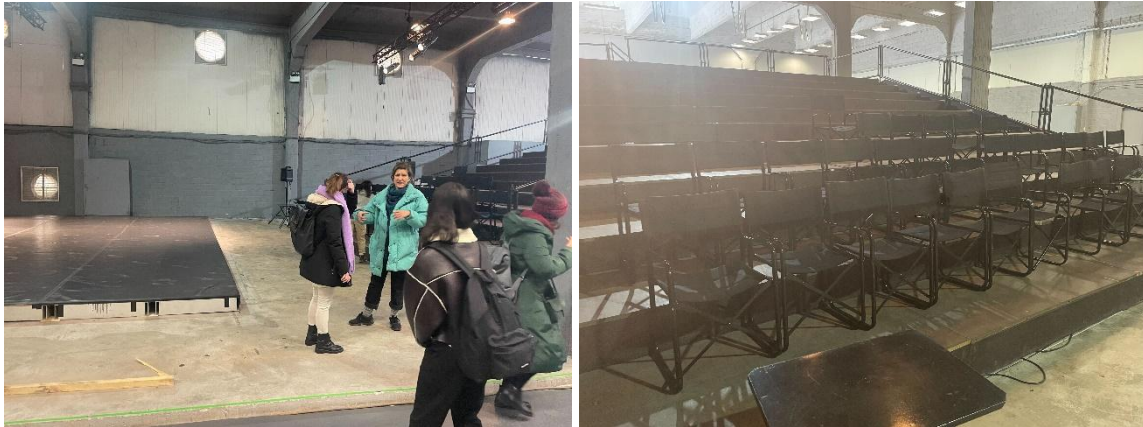


Figure 35. A potential stage at the Athens Epidaurus Festival's, (Left) Distance between the stage and the audience, (Right) Chairs for the audience

The language of the performance will be either Greek or English and the translation will be available in all the languages under the scope of the VOXReality project. Questionnaires will be developed to assess the audience experience. Similarly, semi structure interviews with selected personnel from the performance's creator team will be constructed.

6.2.6 Target Population

Theatre goers, native speakers of any of the VOXReality languages or English. The selected group will fulfil diversity criteria regarding age, gender, knowledge of and exposure to AR technologies and opinion on AR and AI technology.

6.3 Training Assistant

This use case aims to provide a guided AR industrial training scenario. In this use case, a trainee will work on an assembly task through interaction with 3D virtual content superimposed on the real world. During this AR assembly task, users will be guided and supported by a virtual assistant. Users will verbally engage with the virtual assistant who will be able to provide support when prompted and/or offer support when elicited based on trainee performance in the task. The AR assembly task here will involve the visualization of a holographic computer-aided-design (CAD) file. This holographic CAD file object will be interactable and composed of multiple parts which need to be assembled together by the trainee. Physical manipulation (e.g., picking up, moving, and inserting) of specific object components will be used to assemble these constituent parts into the object's frame. Different levels of training will be offered, ranging from easy to difficult, and assembly levels can be performed in repetition. The virtual assistant will provide a unique language-centric interface with the AR training environment and aid end-user trainees in industrial assembly instruction.

6.3.1 Scenario

This training use case scenario will focus on the experience of a single user participating in the AR training environment. In this scenario, a user will be provided with a HoloLens 2 (and/or an android) AR smart glass device. The user will engage with the training experience as detailed below in the user journey walkthrough, but in general are predicted to consist of running through the individual difficulty levels at least one time each. The use case is currently envisioned to take approximately around 3 hours (for one person), which includes time for set-up, introduction to the hardware device, training application scenarios, clean-up, and feedback assessment periods.

The tentative predicted user journey walkthrough will in general occur with the following assumptions: 1) Only one user will be in the AR training session at a time. 2) Trainees will be inexperienced in using the AR glasses, and 3) Trainees will be relatively inexperienced in the industrial assembly task (i.e., unfamiliar with the CAD file to be assembled).

A user will first be introduced to the HoloLens 2 and will be provided with an explanation of what will be happening during their AR training experience. They will then be provided with relevant consent forms to sign, and after signing they will be fitted with the AR Glasses and the training scenario will begin. After launching the application, either by a VOXReality researcher or the participant verbally with the virtual assistant, the user will be greeted by the virtual assistant and prompted to confirm they would like to begin training. After confirming, the user will then be asked by the virtual assistant which difficulty level they would like to be trained on. Following a verbal decision by the user, the training will begin. Users will then be able to see a holographic 3D industrial CAD file object which is to be assembled. This object will not be superimposed over a physical twin object in real life, but rather existing over the local surrounding environment. This object will be composed of a certain number of parts which are loaded initially external to a core frame. The constituent parts are then to be picked up and inserted in the correct spatial position into the object's frame. Depending on the difficulty level, the parts may be visually ordered in a sequential fashion or have a corresponding color-coded deposit location on the object frame. A user would then begin moving the holographic parts in space and attempting to position them into the frame. Throughout this process, the virtual assistant will be readily available to respond to help inquiries and offer support. Support is predicted to come in the form of e.g., documentation such as images or PDFs. As users pick and move the individual assembly parts into the frame, relevant UI elements are predicted to provide information which may aid in the assembly process or provide feedback on how they are performing, which could scale depending on the difficulty level selected. The virtual assistant will also monitor a users' process and intervene when the user appears to need help. For example, if a user is attempting to insert a part in a wrong location or takes a long period of time to move forward with the assembly, the virtual assistant could engage verbally with the user and offer support.

This process is intended to continue iteratively until all constituent parts of the CAD object are assembled. The user will have learned how to assemble this object and have an understanding of the assembly process. The virtual assistant would then be able to offer the user the opportunity to restart the training again along with the option to choose a different, likely more challenging, difficulty level. The user would then be expected to try the different difficulty levels in sequence. This serial assembly practice and continued help/assistance with the virtual assistant will, in the end, have provided a unique training environment for the user. Following the completion of each difficulty level, it would be beneficial if the user could be provided with metrics associated with their performance (e.g., time spent doing the entire task, per each step, etc.) in some format. After the final difficulty level is reached, the user will end the training session by instruction to the virtual assistant. The user would then remove the AR Glasses and the session would be closed. Following, users will be provided with paper assessment forms to get feedback on their experience as well as discuss with VOXReality experimenters on site.

6.3.2 Technologies

The specific technologies intended to be utilized in this use case consist of hardware devices, which include the HoloLens 2 (and/or other android AR Glasses) , and likely a PC laptop where the AR assembly application will run. This application will be based in Unity and will incorporate Holo-Light's Interactive Streaming for Augmented Reality (ISAR) SDK. ISAR-enabled applications provide a client-

server environment which offloads the computationally heavy rendering process from the AR device to a dedicated server environment (e.g., laptop) with a high-quality graphics processing unit and streams the application. ISAR allows the streaming of this entire application to the HoloLens 2, which itself has an ISAR-client application installed, which permits high resolution AR visualizations with low latency. This application will further make use of Holo-Light's Augmented Reality Engineering Space (AR3S), which will provide the basis for visualization and interaction of 3D CAD files for the assembly training tasks. The application developed for this use case, which combines both the ISAR streaming functionality and the AR training and assembly features of AR 3S will be bolstered by the virtual assistant which will be described below.

The use case will utilize the ASR component developed by the consortium to generate textual counterparts of the speech of users interacting with a digital assistant. The textual information acquired from ASR will be utilized in the DA component where a dialogue system helps the users to assemble a machine. DA component will facilitate a dialogue between the user and the corresponding assembly instructions and descriptions generated by VL component using English as a mediation language. The response generated from DA will be in English and off-the-shelf text to speech tools are envisioned to be utilized to generate an audio response.

Moreover, the agent will be composed of a NLU model that will understand what the user is talking about as well as a NLG model, that generates meaningful responses. The agent will make use of manual as an unstructured textual document that guides the agent in interacting with users. The agent will be able to display pdf, corresponding videos, or specific files, if it is asked by the user.

6.3.3 Equipment

As mentioned above, this use case will make use of the HoloLens 2 as the AR device of choice that the end user (trainee) will utilize to interact in the AR training environment. In addition to allowing holographic visualizations of the 3D CAD data, the HoloLens 2 possesses a microphone and speaker, enabling a full audio experience and thus leverages the ASR component, NLU, and NLG models in order to have a dynamic language experience with the virtual assistant during the training sessions. As mentioned above, a PC laptop will likely be used to host the server-side AR application which will be streamed directly to the HoloLens 2. No additional equipment is envisioned to be necessary for this use case as all interactions within the training environment occur with virtual holographic content and audio information.

6.3.4 External Context

Beyond the holographic CAD file objects that users will interact with in this use case, external context will include support materials which the virtual assistant will provide/offer to the user when engaging with the training session. Such content will likely consist of different images of the CAD object and its constituent parts that the user could view. Additionally, instruction manuals may also be offered in the form of PDFs which the user could view and read at their discretion. While not guaranteed to be present, a further external factor which would be beneficial to have would be videos which can show an individual part being assembled in AR space. Such factors are envisioned to be incorporated into the AR application stream which would allow the user to visualize them during the training session. Given that the entirety of the session will be occurring in AR, the external environment having a defined location is not essential for this use case. However, an ideal location would be an industrial shop floor or similar location which would approximate a canonical environment where a typical user would engage in such training.

6.3.5 Assessment Protocol

Following the training session, users would be asked to participate in a feedback session designed to collect information on their experiences with assembly training with emphasis on interaction with the virtual assistant. Participants will be asked to complete a questionnaire/survey and/or provide direct verbal feedback to VOXReality experimenters about the quality of their experience either verbally or by questionnaire. The results from user profiling will provide insights into what went well, what could be better, and overall impressions of the use case and its technology provided. In addition, assessment protocols are envisioned to be application metrics which are profiled during the training sessions such as e.g., time spent doing a virtual task, number of incorrect steps taken, number of times the virtual assistant had to intercede and offer help due to the latter, etc. This data can be used to scientifically approximate the efficacy of the training experience over sessions with increasing difficulty. This data will be collected from the application itself.

This assessment protocol is intended to be acquired from 3-5 users; more would be ideal if possible. Despite such a presumptive low sample size however, a robust coverage of all possible end user demographics is a goal of this use case. In particular, equal numbers of all genders will be aimed for. Language of operation will likely be in English. Users are envisioned to be readily open to technology, with little experience in XR technologies and some experience in industrial tasks.

6.3.6 Target Population

Target populations in this use case will be users who have a desire to work with virtual assistants for task performance during industrial assembly training. Ideal participants would be those who have no prior knowledge of the task to be performed in this use case, and those who have not had much experience with AR, so they can be instructed without prior influence biasing the training scenario. Users will likely be factory workers with some assembly experience and knowledge of assembly training in general.

7 Requirements (Year 1)

This section lists all the requirements gathered in the first six months of VOXReality implementation, based on the weekly calls and in-situ visits. They are divided into user (Section 7.1) and technical (Section 7.2) requirements.

7.1 User Requirements

This section provides the assessed requirements of each use case – VR Conference, Augmented Theatre and Training Assistant. The list is the result of the cumulated user requirement gathering process, from the initial weekly calls (Section 3), in-situ visits (Section 0), and the final use case description (Section 6).

The initial user requirement (Section 4) was derived from the weekly calls where all the consortium partners worked on aligning the scenarios with a mutual understanding of its technologies and possibilities. Based on the initial user requirements, the user-centred partner designed activity materials and conducted interactive focus group workshops with each use case partner through In-situ visits (Section 0). The focus of the workshops was to elicit the latent needs of each use case owners and the actual users so that can help develop the initial user requirements to in details.

The workshops benefited the user-centred partner and use case owners bidirectionally. On the one hand, the analysis of the activity results and discussions during the workshop led to a list of original insights addition to the initial understanding. On the other, the workshop participation of each use case owners helped them to clarify the project setting and define the scenarios into a deeper level. Consequently, the user requirements (Section 7.1) resulted to have modifications from the initial list (Section 4), by combining the outcomes of In-situ visits (Section 0) and use case description (Section 6).

At this stage it is too early to conclude final requirements from the project, thus these requirements in section 7.1 are not binding to development, instead it provided an overarching perspective collected from users. The implementation of features is subjective to technical feasibility mutually agreed among requirements and technical partners. The following sub-sections present updated user requirement lists with tables and summaries per use case scenario.

7.1.1 VR Conference

Following section describes the overview, scenario, assessment, virtual-agent(s), navigation, subtitles and translation, interface and user-interaction themes derived after carefully analysing the data by thematic analysis, whereas the individual requirements within each theme is enlisted in Table 11.

Overview

- The VR experience should primarily focus on user navigation and language translation features along with support for all the standard functionality required for a VR conference platform. A virtual agent should also be developed to provide users with relevant information to navigate virtual spaces, interact with other visitors, and exchange information. The operational output language will be English but can be translated into other languages.

Scenario

- The virtual conference venue should simulate a professional conference setting, including an entrance, lobby, trade shows, meeting rooms, social areas, and a conference room. The trade show should consist of spatial booths for business interaction and engagement. Participants

will be represented as virtual avatars, complete with predefined gestures and features such as hand-shakes, high-fives, nods, and thumbs up gestures. Presentations should follow a specific conference format, such as TED talks or Pecha Kucha, with a maximum duration of 30 minutes, in order to keep time limit feasible for users, and to prevent discomfort caused by prolonged headset use.

Assessment

- The VR conference targets a diverse group of conference partners and visitors, with a recommended limit of 20 participants per session (room-space). The experience will be assessed using a questionnaire and a semi-structured interview to ensure high quality.

Virtual agents

- A virtual agent will be assigned to each user for the duration of the conference, providing help and navigation assistance. The agent should offer getting-started help, a welcome greeting, FAQs, and personalized info feeds related to the conference program. The agent should also suggest call-to-action depending on the current location and activity status of the user, and present all available options at the end of each session. Communication with the agent should be enabled by voice typing and/or text, and users should have the option to skip the virtual assistant help. The agent should be likeable, friendly, pleasant, customizable, non-intrusive, and realistic. Users should be able to access the agent via a smartwatch/wristband, and the agent should auto-save relevant information for export in a downloadable or email-able PDF.

Navigation

- The navigation experience should be clear with consistent navigation cues and visual landmarks that enabled attendees to easily find their way around the virtual venue. To enhance the user experience, several features related to navigation should be available during the virtual conference. Users should have the option to take a quick tour of the venue, which should last between 30-120 seconds and can be replayed or skipped. A virtual map should be accessible to help users navigate the conference venue, allowing them to identify their current location. To facilitate quick and easy navigation between different locations, users should be provided with an option of a virtual taxi or cab. In case the user takes the wrong direction, they should receive a notification. Arrows should guide users towards their target location. Additionally, a flyover view of the venue should be available with zoom in/out options.

Subtitles and Translation

- The users should be able to communicate and translate both text and voice. They should also have the option to mute the entire room or individual speakers and control the volume. Users should be able to turn on/off the program, subtitles, auto-translation, and voice using interactive buttons. For multi-speakers, the subtitles should highlight the active speaker among the group, with priority given to the speaker that the user is currently looking at.
- The subtitles should be displayed relative to the context using various suggestions like speech bubbles or text strips, with the default placement being at the bottom of the screen. To avoid eye-strain, subtitles should be placed in level to viewing angle and also ensure to not mask/block important parts of the display or screen. For HMDs, subtitles should be located over or in the middle of the screen. The subtitles and translation experience should promote inclusivity and accessibility, enabling attendees from around the world to fully participate in the conference and engage with conference content in a meaningful way.

Interface

- There should be a customizable dashboard that is always visible, allowing easy access to frequently used features. A help button should also be present on the corner of the screen for quick assistance. Additionally, two in-room view options should be available: full screen and conference room. To facilitate audience engagement, questions and feedback should appear near the presenter's screen, ensuring that they are noticed and answered promptly. When an attendee asks a question, it should be sent visually or textually to the speaker to avoid it going unanswered.

User Interaction

- The system should provide options for users to interact with speakers and guests during Q&A sessions without interrupting the presentation. Users should be able to ask questions to an agent during navigation, using voice commands such as "Please describe the scene for me," "What is in this room?" or "Is there an empty chair in the room?" Additionally, users should be able to use hand gestures or interactive buttons on the screen for navigation.

Table 11. Final Requirements for VR Conference use case

Type	Sub category	Requirements	Priority	Notes	Phase gathered
General		The experience is in VR.	High		Conference use case scenario document
General	Objective	Provide user navigation and language translation at the virtual conference.	High		Use Case Description
General	Objective	Develop a virtual assistant providing users with relevant information to navigate virtual spaces, interact with other visitors, and exchange relevant information.	High		Use Case Description
General		The operational output language is English, and it can be translated into other languages.	High		Technical-Use case meeting 6
Scenario	Space	The conference venue should imitate the environment of a professional conference setting.	High	*e.g., entrance, lobby, trade shows, meeting rooms, social areas, and a conference room, etc. *The trade show includes spatial booths for business interaction and engagement.	Use Case Description
Scenario	User Representation	The users will be represented as virtual avatars.	High		Use Case Description
Scenario	Social Interaction	The virtual avatars of the participants will offer predefined set of gestures and features.	Low	*e.g., hand-shake, high-five, nodding, thumbs up gesture etc.	Use Case Description

Scenario	Contents	The presentation should be in a predefined conference format.	High	*e.g., TED, Pecha Kucha etc.	Use Case Description
Scenario	Contents	The suggested duration for a single talk is 30 minutes at maximum.	High	*wearing headsets for too long is uncomfortable and might cause headaches / visual discomfort.	Use Case Description
Assessment	Target User	The target user includes conference partners and visitors with diverse backgrounds.	High		Use Case Description
Assessment	Target User	The suggested number of participants per suggestion is 20 people.	Low		Use Case Description
Assessment		The quality of experience will be assessed by a questionnaire followed by a semi-structured interview.	High		Use Case Description
Virtual Agent	Function	A dedicated virtual agent will be assigned to each user and will stay with him/her during the complete duration of the conference.	High	*Virtual agent should offer getting-started help without the user asking for it.	Conference use case scenario document
Virtual Agent	Function	The virtual agent should provide welcome greeting.	High		In-situ visit - VRDays
Virtual Agent	Function	The virtual agent should provide help and FAQs.	High		In-situ visit - VRDays
Virtual Agent	Function	The virtual agent will help user to navigate the space, answer programme-related questions, and deliver relevant information.	High		Use Case Description
Virtual Agent	Function	Virtual agents inform the users about the context/content of the venue and potential to-do items in advance or relative to the current activity.	Low	*The virtual agent should offer call-to-action depending on the current location and activity-status of the participant.	Technical-Use case meeting 2
Virtual Agent		The users should have a skip option for the virtual assistant help.	High	*In case they have already done the training or they are aware of similar experiences	In-situ visit - VRDays
Virtual Agent		The communication with virtual agent should be enabled by voice typing and/or text.	High	*Text typing is very time-consuming in VR settings.	In-situ visit - VRDays
Virtual Agent	Function	Virtual agents deliver personalized info feeds related to the conference programme-related activities.	Low	*The virtual assistant should keep track of the activities of the user and suggest personalized suggestions to each participant.	Conference use case scenario document, In-situ visit - VRDays

Virtual Agent		The virtual assistant can be accessible via a smart-watch / wrist-band.	Low	*to access available options at all times during the conference	In-situ visit - VRDays
Virtual Agent	Function	The virtual-assistant should auto-save relevant information and provide them in an exportable file.	Medium	*e.g., the day log, rooms visited, events attended, people met with their contact information, etc in form of a pdf file that could be downloaded or emailed to the user.	In-situ visit - VRDays
Virtual Agent	Function	The virtual agent should present users with all available options at the end of every session.	High	*The options may include vacating the conference room for the conference venue, re-visit the conference and/or exiting the event completely.	Conference use case scenario document
Virtual Agent	Behavior	The virtual agent should interact with the users on-demand, and should not be intrusive.	High		In-situ visit - VRDays
Virtual Agent	Behavior	Virtual agent should be likeable, friendly, pleasant, customizable, realistic and complete.	Low	*e.g., for behavior, the interaction of agent with user should be intrinsic. For appearance, the agent should be complete and not be missing graphical details like missing limbs, parts or other natural features.	In-situ visit - VRDays
Virtual Agent	Appearance	Virtual agent could look like a cartoony avatar.	Medium	e.g., cloud, taxi, humanic personification, a robot, an egg (blobby), magic wand etc.	In-situ visit - VRDays
Navigation		Quick tour option should be available.	High	*Suggested duration is 30-120 seconds. *It should be skipped and/or replayed.	In-situ visit - VRDays
Navigation	Interface	Virtual map should be available to navigate the conference venue.	Medium	*Users should be able to identify their current location.	In-situ visit - VRDays
Navigation		Users should be assisted for quick and easy navigation between different places in the venue.	Low	*e.g., through a virtual taxi/cab option to go to the desired location in the conference. * Users should be notified in case s/he goes in wrong directions.	In-situ visit - VRDays
Navigation	Interface	Navigation to the target location should be guided by arrows.	High		In-situ visit - VRDays
Navigation	Interface	Flyover view of the venue should be available with zoom in/out options.	High		In-situ visit - VRDays

Subtitles	Function	Translation should be available in both text and voice formats.	Medium		In-situ visit - VRDays
Subtitles	Function	Option to mute the complete room or individual speaker(s) should be available.	High		In-situ visit - VRDays
Subtitles	Function	Users should be able to control the volume +/-	High		In-situ visit - VRDays
Subtitles	Function	Users should be able to turn on/off the program, subtitles, auto-translation, and voice using interactive buttons.	Medium		In-situ visit - VRDays
Subtitles	Interface	For multi-speakers, the subtitles should highlight and different the active speaker among group of speakers.	Medium	<p>*e.g., semi-transparent speech bubble; highlighting the active speaker</p> <p>*When the users look at a particular speaker among many, that speaker's subtitles should have priority of view/voice over other presenters/speakers.</p>	In-situ visit - VRDays
Subtitles	Interface	The subtitles should be displayed relative to the context.	High	<p>*Subtitles display suggestions:</p> <ul style="list-style-type: none"> - strip on the top, - moving text strip on top, - speaker-attached speech bubble, - strip of text on bottom, - text box on side of the speaker. 	In-situ visit - VRDays
Subtitles	Interface	The default placement of subtitles should be on the bottom of the screen.	High	<p>*To match the standard/classic subtitles style in videos and movies.</p> <p>*Also, to prevent eye-strain, the subtitles should not be too high or too low and relative to the screen.</p>	In-situ visit - VRDays
Subtitles	Interface	The subtitles should be visible without masking other important parts in the display/screen.	High		In-situ visit - VRDays
Subtitles	Interface	For HMD, the subtitles should be over or in middle of the screen, not on the bottom of the screen.	High		In-situ visit - VRDays
Interface	Dashboard	There should be a customizable dashboard visible at all times.	Medium	* The dashboard should allow participants to access most frequently used features.	In-situ visit - VRDays
Interface	Help	A help button should be visible on the corner of the screen.	High		In-situ visit - VRDays

Interface		In the conference room, users can choose between two in-room view options i.e., full screen and conference room.	Medium		Conference use case scenario document
Interface	Feedback	Questions/feedback from the audience should appear textually/visually near the presenter's screen.	Medium	*When a user asks a question, his/her question should be sent visually / textually to the speaker, to avoid being unanswered.	In-situ visit - VRDays
User Interaction		Users should be able to engage with speakers and guests during Q&A session(s).	High	*The interaction between the speaker/presenter and the attendee/user should be available without interrupting the presentation/talk.	Conference use case scenario document
User Interaction		Users can interact with and ask the agent questions during the navigation.	High	*e.g., "Please describe the scene for me.", "What is in this room?", "Is there an empty chair in the room?" etc.	Technical-Use case meeting 3
User Interaction		Users should have options to use hand-gestures or interactive-buttons on the screen.	Low		In-situ visit - VRDays
Extra		Chatbot should be available on demand.	Low		In-situ visit - VRDays
Extra	Payment	Users should be able to make digital and crypto payments in the conference items shop.	Low	*This is low-priority / extra suggestion	In-situ visit - VRDays

7.1.2 Augmented Theatre

This section describes the most prominent themes of overview, setup, scenario (play-design), assessment, interface, subtitles, VFX, and additionally recommended requirements that have been acquired after a thorough evaluation of the qualitative data through thematic analysis, whereas the Table 12 enlists details of each individual requirement within relative theme.

Overview

- The main objectives are to understand the users' needs and the organizers' wants for subtitles and VFX while watching a theatre play, and to brainstorm ideas for the UI and interaction method between the audience and AR glasses. It's important to ensure that the gestures and setup required to operate the AR glasses do not interfere with other audience members' theatre experience, and that the possibility of distracting the audience with additional information from the glasses is taken into consideration.

Setup

- The AR glass setup should be user-friendly for those who are not familiar with AR technology. Also, it should provide enough control to the audience while limiting their ability to alter the default settings.

Scenario: Play design

- The play of "Hippolytus", story of a young man who has sworn a sacred oath of chastity and devotion to the goddess Artemis, written by Euripides will be specially produce for a performance of not more than 15 minutes.

Assessment

- The scenario involves two users watching the play simultaneously using AR devices, subject to availability. The target user base includes diverse theatregoers whose native language is any of the VOXReality languages. It may also include some participants with minor hearing or vision impairments. Up to 50 participants will be involved in the complete experiment, including non-Greek speaking audiences in the Greek play. Automatic translation and visual effects technology will be available for the audience, and the user evaluation may include a combined demonstration of narration, dialogue, and VFX. Creators of the performance may also be asked to provide feedback through semi-structured interviews.

Interface

- The AR glass will provide a menu on the screen to adjust practical settings, allowing the audience to learn more about the play or scene and zoom in or out. The elements on the AR screen should adapt to the luminosity of the stage and theatre, and the interface should not require extensive head movements from the audience. Additionally, the UI of the AR screen should consider accessibility, such as font size for those with minor visual impairments.

Subtitles

- The AR Theatre should provide various options for subtitles, such as real-time translations in different languages, adjustable sizes, and the ability to turn them on and off. The audience should be able to move the subtitles up or down on the screen and avoid overlapping with the stage setup, allowing them to watch the actors' faces while reading the subtitles text. Additionally, the glasses should offer different audio language options, with the default being the original spoken language of the play. However, the audience should have the option to switch to a different language or to listen to both the original and local languages.

VFX

- During the early production stage, the implementation of VFX should be discussed with the script and plot, and the director should specify the users' capability to change AR Glass settings and subtitles. VFX can accompany or reflect the narration of the scene, and certain words or phrases may trigger the VFX. The VFX could be more appropriate for substituting supplementary features. However, the location of the visual effects should not exceed the stage's range, and they should not affect the actor's actions. The VFX should help the audience immerse in the performance, and low-quality VFX should be avoided as they may decrease the user experience. Additionally, the style of the visual effects should be artistically relevant to the opera and not too cartoonish.

Additional

- The stage background can be interactive, allowing the audience to learn about the play or read fun facts while waiting for the performance. The audience can also learn more about the player's or director's biographies, as well as other information about the play. Additionally, the audience can watch behind-the-scenes footage or learn about the social/historical background of the play before and after watching it.

Table 12. Final Requirements for Augmented Theatre Use Case

Type	Sub category	Requirements	Priority	notes	Phase gathered
Objective		The objective is to introduce AR technologies to theatre audiences.	High		Use Case Description
Objective		Language translation and VFX experience with AR glasses at the theatre.	High		
General		The experience is in AR.	High		
General		The gestures and setup to operate AR glass should not intervene other audience's theater experience.	High	*The AR glass operation should involve physical buttons.	In-situ visit - AEF
General		The possibility of the audience getting distracted by the additional information from AR glasses has to be considered.	Medium		In-situ visit - AEF
Setup		The experiment takes place at the avant-garde theatre.	High	*modern, indoor theatre	Technical-Use case meeting 6
Setup		The AR glass setup should be easy and clear to audiences who are not familiar with the AR technology.	High		In-situ visit - AEF
Setup		AR glass should give the audience enough controls but also limit their ability to change default/preset settings.	Medium		In-situ visit - AEF
Scenario		A scene from the play "Hippolytus" by Euripides will be specially produced and played.	High		Use Case Description
Scenario		The length of the performance will not exceed 15 minutes.	Medium		Use Case Description
Assessment	Target User	Two users can watch the play at the same time with two AR devices in a single time.	High	*based on the availability of AR devices	Technical-Use case meeting 4
Assessment	Target User	The target user is a theatergoer with diverse background whose native language is any of the VOXReality languages.	High	*i.e., age, gender, language, technology knowledge, opinion/acceptance and experience, exposure to AR technology.	Use Case Description

				*People with disabilities will not be considered; we can put emphasis on the potential of VOX technologies in providing access to people with disabilities (maybe auditory disabilities)	
Assessment	Target user	The target user includes who don't know the language or cannot hear well.	Low		In-situ visit - AEF
Assessment	Target User	The target user base includes people with vision correction glasses.	Low		Technical-Use case meeting 6
Assessment	Target User	In Greek play, non-Greek audiences should participate in the experiment.	Medium		In-situ visit - AEF
Assessment	Target User	Up to 50 people will participate in the complete experiment.	Medium		Theatre use case scenario document
Assessment		The translation of the play will be available in all VOXReality languages.	High		Use Case Description
Assessment		The audiences may experience either one or both of the automatic translation and the VFX technology.	High		Use Case Description
Assessment		The user evaluation may include combined demonstration.	Medium	*For example, 5-minute narration + 5-minute dialogue + 5-minute VFX + 5-minute dialogue/narration with VFX	Technical-Use case meeting 2
Assessment		The creators of the performance may be asked to evaluate the experience through semi-structured interviews.	High		Use Case Description
Interface	Function	The AR glass could have menu on the screen to change practical settings.	High		In-situ visit - AEF
Interface	Function	Audience can learn about the play or the scene.	Low		In-situ visit - AEF
Interface	Function	Audience can zoom in or out the scenes.	Low		In-situ visit - AEF
Interface		The elements in the AR glass screen should adapt to the luminosity of the stage and theater.	Low		In-situ visit - AEF

Interface		The interface should not include extensive head movements.	Low	*The audience doesn't have to turn his/her head to see everything in the proper place.	In-situ visit - AEF
Interface		The UI of the AR screen should consider audience accessibility.	Low	e.g., small font size	In-situ visit - AEF
Subtitles	Default	The default subtitle starts with the local language.	High		In-situ visit - AEF
Subtitles	Setting	The subtitles should be real-time, provide different languages, sizes, and turn on and off functions.	Medium		In-situ visit - AEF
Subtitles	Location	The audience should be able to change the location of the subtitles, either up or down.	Medium		In-situ visit - AEF
Subtitles	Location	The subtitles should not overlap with the stage setup.	High	*The audience should be able to watch the actors' face while reading the subtitles.	In-situ visit - AEF
Subtitles	Location	The subtitles can be placed above or below the screen and should follow each actor.	Medium	*Greek audience is used to subtitles underneath the screen.	In-situ visit - AEF
Subtitles	Default	The default audio starts with the original spoken language of the play.	High		In-situ visit - AEF
Subtitles	Audio	The audience should have an option to change the language audio.	Low	*Some people don't like to hear the specific language.	In-situ visit - AEF
Subtitles	Audio	The audience should have option to listen to the original and/or local language.	Low		In-situ visit - AEF
VFX		The VFX implementation should be discussed earlier with script and the plot.	Low	*early production stage	In-situ visit - AEF
VFX		The director specifies the capability of the users to change AR Glass settings and subtitles.	Low		In-situ visit - AEF
VFX	Function	VFX can reflect or accompanies the narration of the scene.	Low		In-situ visit - AEF
VFX	Function	Certain words or phrases will trigger VFX.	Low		Technical-Use case meeting 1
VFX	Function	A narrator could be the person who triggers the VFX.	Low		Technical-Use case meeting 2

VFX	Scope	The VFX may be more suitable for substituting the supplementary features.	Low		In-situ visit - AEF
VFX	Scope	Location of the visual effects should not exceed the range of the stage.	Low		Technical-Use case meeting 6
VFX	Scope	The implemented VFX should not affect the actor's actions on the stage.	Low	*The added VFX should not confuse the actors in how they perform.	In-situ visit - AEF
VFX		The VFX should help the audience be immersed in the performance.	Low	*Low quality VFX should be avoided since they are reported to decrease the user experience.	In-situ visit - AEF
VFX		Style of the visual effects should be artistically relevant to the opera.	Low	*not too cartoonish	Technical-Use case meeting 6
Extra	Function	Background of the stage can be interactive.	Low		In-situ visit - AEF
Extra	Function	Audience can read about the stage or fun facts while waiting for the performance.	Medium		In-situ visit - AEF
Extra	Function	Audience can learn more about the biography of the player, director, or more information about the play.	Medium		In-situ visit - AEF
Extra	Function	Audience can watch "behind the scene" or social/historical background of the play before and after watching the play.	Medium		In-situ visit - AEF

7.1.3 Training Assistant

In this section, we define the most prominent themes of assessment, user-interaction, interface, virtual-agent, and feedback requirements derived after the careful thematic analysis of the data regarding personal assistant and its related features, while the individual requirements within each theme is enlisted in Table 13.

Assessment

- **Target User:** The *target user* is a factory worker with some assembly experience and limited AR knowledge. They have no prior knowledge of the task and are relatively inexperienced in industrial assembly tasks.
- **Training:** The *training* will be conducted in English, with 3-5 users assessed using metrics like time spent, incorrect steps, and virtual assistant intercession. User feedback will be collected via questionnaire/survey or verbal feedback on assembly training and interaction with the virtual assistant. Each use case scenario takes around 3 hours.

User interaction

- The AR experience involves physical manipulation of object components, where the user can ask for instructions to assemble a machine/object. Users can choose between textual or visual instructions, with options to control the narration speed. User interaction is hands-on without controllers, and verbal engagement with the virtual agent is also possible.

Interface

- The interface should have a dashboard for commonly used functions that users can customize and personalize to their needs. It should also have appropriate colour schemes and visual cues for each action. Access to help and FAQs should be available on the screen at all times.

Virtual agents

- Virtual Agents provide verbal/textual explanations with corresponding visual content. It would use voice commands, visual cues, and other interactive features to help trainees identify the correct tools and parts, navigate the assembly process, and troubleshoot any issues or errors that may arise. They direct users to necessary objects and offer on-demand help, including hints and quick tips in text and/or audio format. The intention is to understand user needs and enhance efficiency of the machine assembly training in real-time, providing trainees with the support and guidance they need to master complex assembly procedures. The virtual assistant monitors user actions, intervenes after a certain number of incorrect steps, highlights errors, and intuitively guides users towards solutions. Support from the virtual assistants may include documentation such as images or PDFs.

Feedback

- Users should receive immediate and self-explanatory feedback on their actions, with instructions clearly stating any placement errors and pointing towards areas that need to be fixed/reassembled. Visual and audible notifications should indicate the correctness of tasks, and suggestions and feedback should appear automatically in case of incorrect assembly. The training and assembly should end with results and feedback, including time taken, number of correctly and incorrectly assembled parts, and other related information collected during the interactive session.

Table 13. Final Requirements for Training Assistant use case

Type	Sub category	Requirements	Priority	Notes	Phase gathered
Objective		The goal is to provide a guided AR industrial training scenario.	High		Use Case Description
Setup		The external environment having a defined location is not essential.	High		Use Case Description
Setup		An ideal location is an industrial shop or similar which approximate a canonical environment.	High	<i>*where a typical user would engage in such training.</i>	Use Case Description
Setup		Three levels - easy, medium, and difficult - will be offered.	High	<i>*In easy mode, 3D model always fits correctly to its location when in the correct place. *In hard mode, user needs to put the parts into the</i>	Use Case Description

				<i>correct location and align the part's orientation.</i>	
Assessment	Target User	The target user is a factory worker with some assembly experience and knowledge of assembly training.	Low		Use Case Description
Assessment	Target User	The target user has no prior knowledge of the task performed in the use case.	High		Use Case Description
Assessment	Target User	The target user has not much experience with AR.	High	<i>*specifically, HoloLens 2, to prevent influences of biasing the training scenario</i>	Use Case Description
Assessment	Target User	The target user is relatively inexperienced in the industrial assembly task.	High		Use Case Description
Assessment	Target User	The target user will cover all possible end-user demographics, including equal numbers of genders.	Low		Use Case Description
Assessment		The operational language is English.	High		Use Case Description
Assessment		One user at a time will be in the AR training session.	High		Use Case Description
Assessment		At least 3-5 users will be assessed.	High		Use Case Description
Assessment		The assessment metrics include - time spent doing a virtual task, number of incorrect steps, number of times the virtual assistant interceded to help - etc.	High	<i>*data can be used to scientifically approximate the efficacy of the training experience over sessions with increasing difficulty.</i>	Use Case Description
Assessment		A round of use case scenarios takes approximately 3 hours.	Low	<i>*including set-up, introduction to the hardware device, training application scenario, clean-up, feedback assessment.</i>	Use Case Description
Assessment		The user's assembly training experience with emphasis on interaction with the virtual assistant will be asked.	High		Use Case Description
Assessment		The user feedback will be collected via questionnaire/survey and/or direct verbal feedback.	High		Use Case Description
Interaction		The AR experience involves the user's physical manipulation of object components.	Low	<i>* e.g., picking up, moving, and inserting</i>	Use Case Description
Interaction		User can ask for instruction to assemble a machine/object.	High	<i>*e.g., "Teach me how to do this", "Am I doing this right?"</i>	Technical-Use case meeting 6
Interaction		Users should be able to choose between textual instructions or visual demo for the training.	High		In-situ visit - HOLO

Interaction		Users should have options to turn on/off, speed-up or slow-down the narration.	Low		In-situ visit - HOLO
Interaction		User interaction is with hands, without controllers and may also verbally engage with the virtual agent.	High	<i>*Available at HOLO & hardware SDKs</i>	Technical-Use case meeting 3
User Interface	Dashboard	There should be a dashboard for most common functions, which users can customize and personalize for their needs.	Medium		In-situ visit - HOLO
User Interface		The interface should provide color schemes and visual cues appropriate for the action.	High	<i>*If the step is correct, it should be highlighted as green, if it needs partial adjustments, it should be yellow and it needs to be redone completely, it should be red.</i>	In-situ visit - HOLO
User Interface	Help	Access to help and FAQs should be available on screen at all times.	Low		In-situ visit - HOLO
Virtual Agent	Appearance	Virtual agents will be animated if in 3D avatar format.	Medium		Technical-Use case meeting 3
Virtual Agent	Format	Virtual agents can provide verbal/textual explanation and display corresponding visual contents.	High		Technical-Use case meeting 3
Virtual Agent		Virtual agents direct the user towards objects that are necessary for the application.	High	<i>*e.g., certain stance for assembly of a machine part</i>	Technical-Use case meeting 3
Virtual Agent		There should be on-demand help during the interaction.	High	<i>*e.g., How-to-do suggestion' in a display/dialogue box/ The virtual assistant will provide support when prompted or/and elicited based on trainee performance in the task.</i>	
Virtual Agent		The help should include hints and quick tips to guide users.	High		In-situ visit - HOLO
Virtual Agent	Format	The tutorial should be in both text and/or audio.	Medium		In-situ visit - HOLO
Virtual Agent		The virtual assistant should monitor the user action and offer instructions to help demonstrate the user on how to assemble a part correctly.	High	<i>*the form could be in format of video tutorials, demo, visual, textual and audible instructions.</i>	In-situ visit - HOLO
Virtual Agent		The virtual assistant should intervene and offer automatic help after a certain number of incorrect steps by users.	Medium		In-situ visit - HOLO
Virtual Agent		The virtual agent should highlight the errors and intuitively guide users towards the solution.	Medium		In-situ visit - HOLO

Virtual Agent	Format	Support from a virtual assistant may come in the form of documentation, such as images or PDFs.	High		Use Case Description
Feedback		Trainee should be provided with immediate and self-explanatory feedback on actions.	Low	<i>*If there is any error in the placement, the instructions should clearly state the issue and point towards the area to be fixed/reassembled</i>	In-situ visit - HOLO
Feedback		Users should have both visual and audible notifications regarding the correctness of the tasks at hand.	High		In-situ visit - HOLO
Feedback		The suggestions and feedback should appear automatically in case of incorrect assembly.	High		In-situ visit - HOLO
Feedback		The training and assembly should end with results and feedback.	High	<i>*feedback includes the time taken, number of correctly assembled part, number of incorrectly assembled parts and other related information collected during the interactive session.</i>	In-situ visit - HOLO

7.2 Technical Requirements

In the three following subsections, one for each use case, we categorize the various technical requirements into Functional and Non-Functional ones, in compliance with the common standard ISO/IEC 25010. The prioritization of the requirements is developed under the MoSCoW model.

Besides the use case-specific requirements, there are universal ones, applicable to all scenarios. These are:

Title	Components' connectivity
Description	All system components must be connected using the most convenient interface and topology, regarding both performance and usability.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	UNI-NFR-01

Title	Context-aware reasoning
Description	The system must consider live audio-visual environmental information.
Typology	Non-Functional
Priority	MUST
Requirement ID	UNI-NFR-02

Title	Real-time performance
Description	The system must operate and allow interactions at near-real time.
Typology	Non-Functional
Priority	MUST
Requirement ID	UNI-NFR-03

7.2.1 VR Conference

Use case-specific technical requirements, tailored for the VR Conferences scenario are:

Title	VR HMD
Description	The system must support the use of VR glasses since the experience will be fully virtual.
Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-01

Title	VR HMD audio capture
Description	The HMD must support audio capturing to input what the user says into the system.
Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-02

Title	VR HMD sound reproduction
Description	The HMD should have built-in speakers to playback sounds to the user.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	VC-NFR-03

Title	Communication with agent
Description	The system must allow users to communicate with virtual agents.
Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-04

Title	Communication with agent
Description	Users must be able to input their feedback through voice and hand-held controllers.
Typology	Functional
Priority	MUST
Requirement ID	VC-FR-01

Title	Virtual conference space
Description	The system must have a virtual conference 3D space.

Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-05

Title	Virtual conference space
Description	The 3D space must comprise entrances/exits, lobby areas, trade shows, meeting rooms, social areas and a grand conference/plenary room.
Typology	Functional
Priority	MUST
Requirement ID	VC-FR-02

Title	Virtual avatars
Description	The system must allow users to be represented by virtual avatars in the conference venue.
Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-06

Title	Virtual avatars
Description	These avatars could be programmed to deliver visual feedback complementing available real-time multilingual translation and captioning (sensible avatars).
Typology	Functional
Priority	COULD
Requirement ID	VC-NFR-03

Title	Virtual agents
Description	The users must be able to communicate with virtual agents that help users explore the virtual conference venue.
Typology	Non-Functional
Priority	MUST
Requirement ID	VC-NFR-07

Title	Virtual agents
Description	The agents should be able to provide information about the context and contents of the rooms and answer questions from the users.
Typology	Functional
Priority	SHOULD
Requirement ID	VC-FR-04

Title	Virtual agents
Description	The agents should be able to provide navigation information to the users.
Typology	Functional
Priority	SHOULD
Requirement ID	VC-FR-05

Title	Virtual agents
Description	The agents could have the form of a humanized avatar.
Typology	Functional
Priority	COULD
Requirement ID	VC-FR-06

Title	One-to-one communication
Description	The platform should allow user avatars to talk between them.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	VC-NFR-08

Title	Talk session one-to-many
Description	The platform should allow a user to give a speech.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	VC-NFR-09

Title	Talk session many-to-many
Description	The system could allow users to engage in Q&A many-to-many sessions.
Typology	Non-Functional
Priority	COULD
Requirement ID	VC-NFR-10

Title	Input video stream
Description	The system could allow users to up-stream video in order to make presentations or give a speech.
Typology	Non-Functional
Priority	COULD

Requirement ID	VC-NFR-11
----------------	-----------

Title	Language translation
Description	The system should be able to provide translation from and into the consortium-selected languages.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	VC-NFR-12

Title	Language translation
Description	The system should be able to provide translation in the form of live captions.
Typology	Functional
Priority	SHOULD
Requirement ID	VC-FR-07

7.2.2 Augmented Theatre

Use case-specific technical requirements, tailored for the Augmented Theatre scenario are:

Title	AR HMD
Description	The system must support the use of AR glasses to provide transparent overlays with captions and effects.
Typology	Non-Functional
Priority	MUST
Requirement ID	TR-NFR-01

Title	AR HMD video capture
Description	The HMD must support video capturing to input what the user sees into the system.
Typology	Non-Functional
Priority	MUST
Requirement ID	TR-NFR-02

Title	Scene audio capture
Description	An array of microphones may be placed on-stage and record the action clearly.
Typology	Non-Functional
Priority	MUST
Requirement ID	TR-NFR-03

Title	Language translation
Description	The system should be able to provide translation from and into the consortium-selected languages.
Typology	Non-Functional
Priority	SHOULD
Requirement ID	TR-NFR-04

Title	Language Translation
Description	The system should be able to provide translation in the form of live captions.
Typology	Functional
Priority	SHOULD
Requirement ID	TR-FR-01

Title	Reactive VFX overlays
Description	The system must be able to provide unscripted visual effects in the form of aligned AR overlays.
Typology	Non-Functional
Priority	MUST
Requirement ID	TR-NFR-05

Title	Actor tracking
Description	The system must be able to track actors on-stage to ensure proper visual effects' alignment.
Typology	Non-Functional
Priority	MUST
Requirement ID	TR-NFR-06

7.2.3 Training Assistant

Use case-specific technical requirements, tailored for the Training Assistant(s) scenario are:

Title	AR HMD
Description	The system must support the use of AR glasses to provide transparent overlays with instructions.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-01

Title	AR HMD video capture
Description	The HMD must support video capturing to input what the user sees into the system.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-02

Title	AR HMD audio capture
Description	The HMD must support audio capturing to input what the user says into the system.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-03

Title	AR HMD sound reproduction
Description	The HMD must have built-in speakers to playback sounds to the user.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-04

Title	AR HMD hand-tracking support
Description	The HMD must support hand-tracking to enable virtual tools' manipulation.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-05

Title	Virtual components library
Description	The system must provide a library of virtual tools and components that can be manipulated by the user in the training scenario.

Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-06

Title	Difficulty-level system
Description	The system could provide users with the option to select a difficulty mode during the training scenario.
Typology	Non-Functional
Priority	COULD
Requirement ID	PA-NFR-07

Title	Difficulty-level system
Description	In easy mode, the 3D model always fits correctly to its location when in the correct place; in hard mode, the user needs to put the parts into the correct location and orientation.
Typology	Functional
Priority	COULD
Requirement ID	PA-FR-01

Title	Communication with agent
Description	The system must allow users to communicate with virtual agents.
Typology	Non-Functional
Priority	MUST
Requirement ID	PA-NFR-08

Title	Communication with agent
Description	Users must be able to input their feedback through voice and hand-tracking channels.
Typology	Functional
Priority	MUST
Requirement ID	PA-FR-02

8 Summary of Pilots (Year 2)

This section provides a brief summary of Pilot 1 for the VR Conference, Augmented Theatre, and Training Assistant scenarios. The focus is on presenting the core aspects of each application, participant information, and key conclusions along with key-findings for lessons learned⁶ from the pilot phase. These pilot findings are critical for the future deliverables D2.5 (Organisational preparation for VOX pilot scenarios and PRESS analysis V2) and D5.2 (Pilot planning and validation V2), eventually supporting the planning of the second phase of pilots and Open Calls for the project. For (internal readers only) a more detailed analysis and comprehensive feedback, please refer to the private deliverable D5.3 (Pilot Analysis and Feedback). This summary aims to offer an overview while directing readers to the complete report for in-depth insights and evaluations.

The technology maturity diagram presented (see Figure 6) provide a complete overview of the pilot's pipeline. Starting from initial app design to leading towards the series of pilot phases and following by a user-study in parallel. Pilot 0 marked the initiation of internal testing, reserved for developers and consortium members, to gain preliminary insights. The subsequent phases, Pilot 1 and Pilot 2, were based on an ascending scale – starting with a smaller user base and moving to larger-scale experiments involving a more extensive user community. The user studies involving research prototypes ran in parallel to this journey, navigating the realms of both functional and non-functional technologies. This comprehensive diagram was also furnished in year 2.

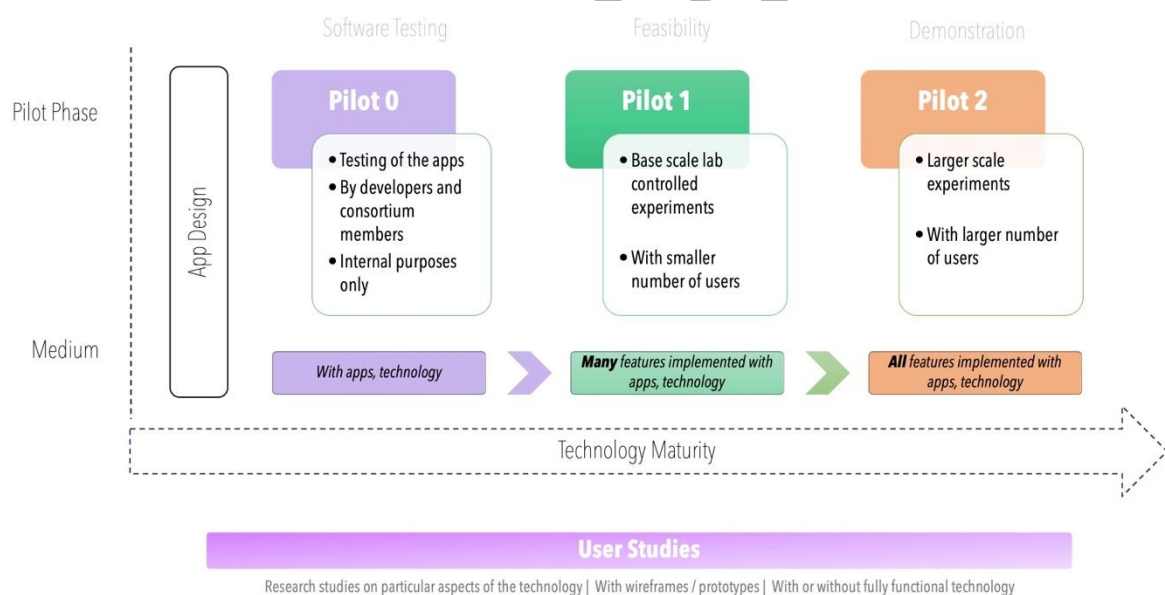


Figure 36 Technology Maturity Diagram

It explained the stage increments of technological development, offering a nuanced view of how VOXReality Technologies evolved and transformed into a fully realized and user-tested solution in XR Spaces over the course from App Design → Pilot 0 → Pilot 1, to test the initial functionalities of the models within the use case applications. This approach is important in order to adequately test and evaluate the user requirements implemented within each scenarios.

⁶ Detailed account of Lessons Learnt would be covered in WP6 and its related deliverables.

8.1 VR Conference

VOXReality VR Conference use case offers a virtual, assisted experience, emulating the features of a professional conference: venue lobby, tradeshow, social spaces, business meeting areas, and conference room (see Figure 37). The VR Conference case includes a virtual personal assistant providing instant translation in multiple languages, venue navigation assistance, and conference program advice, enabling seamless navigation and multilingual interactions.



Figure 37 Live Presentation demonstration within VR Conference Application

During Pilot 1, users interacted with the Virtual Agent, receiving a short tutorial on the Virtual Agent, venue navigation, and user interactions. The pilot aim was to gather and analyse quantitative and qualitative results to evaluate the overall VR conference experience. The pilot intended to assess the impact of AI-assisted VR technology for navigation and speech translation on conference-goers, focusing on (A) navigation technology response and (B) machine translation response. We follow a mixed-methods approach, incorporating quantitative (questionnaires, implicit metrics) and qualitative (post-VR experience interviews) measures.

8.1.1 VR Conference Application

The VR Conference application (Figure 38) emulates a professional conference setting in VR, equipped with real-time translation of a conference presentation and a virtual agent that offers navigation instructions. The application is developed based on Mozilla Hubs, while the VR spaces and 3D objects are created in Blender.



Figure 38 VR conference application screenshots

The features and functionalities that were available in Pilot 1 are:

- 1) The VR space consists of five rooms: Lobby, trade show with 4 booths, conference space, social area, and business area. Three of them (Lobby, trade show, conference space) are already accessible; for the other rooms, doors indicate the location of these spaces.
- 2) The user selects their desired language before entering the application, and then the entire VR application is translated into the user's selected language to enhance the user experience.
- 3) Lobby room with an interactive tutorial. It's a task-based tutorial with tasks that the user should complete, along with instructions on how to navigate the VR space and utilize various functionalities. This way, the user familiarizes themselves with movement and enabling/disabling features.
- 4) User panel that provides options to toggle all features of the VR application. The user panel includes buttons for the virtual agent, mute function, floor map, help, and translation service. Users can access the user panel by looking upwards.
- 5) Tracking of the user's position in the VR space.
- 6) Virtual agent, named VOXY, is a fully interactive, non-intrusive virtual assistant. Users can move it around the space and place it wherever is more convenient for them. The virtual agent provides greetings, navigation guidance, and instructions within the virtual environment. The agent's responses are displayed at separate boxes with arrow navigation for the next part. The user can turn on and off the agent's responses. Moreover, there are two different symbols indicating that VOXY is listening to the user and that VOXY is processing the user's response.
- 7) Navigation guidance and instructions are provided both verbally, with step-by-step guidance to the destination, as well as visually, with green lines on the floor.
- 8) Floor map provides a top view of the space and information about the user's current position and direction.
- 9) Mirrors and other objects are placed in the lobby, the trade show, and the conference rooms to allowing participants to become familiar with their virtual size.

- 10) Help button informs users on how to use the app and navigate within the VR space. Users have access to this information whenever they want via the user panel.
- 11) Translation system that supports 6 languages (IT, ES, DU, DE, GR, EN). In Pilot 1 the 1-to-many translation system is used, which is ideal for a conference room. The translated text appears above the presenter. The users can select if they want to view the translated text or not.



Figure 39 Participants in VR Conference Experience

The pilot study involved 18 participants (11 females, 7 males) with an average age of 37 (range: 24-53). Participants were selected to provide a diverse background professional and linguistic and representation, including English (17), Greek (5), Dutch (13), German (8), Spanish (6), and Italian (5) speakers. Participants (Figure 39) reported their experiences through mixed methods using both industry standard questionnaires and semi-structured interviews, which allowed for in-depth exploration of specific areas of interest while maintaining the flexibility to uncover individual perspectives. The interviews were analysed using Thematic Analysis by Braun & Clark [9].

8.1.2 Conclusion

The findings from the VR conference Pilot Phase 1 highlighted both the strengths and limitations of the Pilot 1 implementation. The results showed several positive patterns, with participants acknowledging the potential of machine translation and showing interest in its use despite technical glitches. There was a notable willingness to attend future VR events, with many participants expressing a preference for virtual conferences due to their time and resource-saving benefits. Additionally, the value of agent assistance was clearly established, as participants appreciated the support from virtual agents in navigating the new VR environment.

However, the pilot also revealed several areas of improvement, including improving translation inaccuracies, issues with acronyms, names, and context, out-of-sync audio and translations. Some navigation issues were also noted due to disconnections and glitches, indicating areas needing refinement to enhance the overall user experience in future VR conferences. While some limitations were inherent to the technological flexibility of the application, others were intentional design and technology-dependent choices made to focus on specific aspects of the navigation and translation modules. These limitations were crucial in identifying lessons learnt (see Table 14) and New (extended) Requirements (see Table 20 in next section [New Requirements \(Year 2\)](#)) for moving forward to Pilot 2 scheduled for next year.

Table 14 Lessons Learnt from Pilot 1 (VR Conference)

#	Category	Sub-category	Lesson Learnt	Phase Gathered
1	User Interaction	Connectivity	Provide stable and reliable connections. Ensure robust and reliable connections. Minimize disconnections and glitches during the conference.	Pilot 1
2	User Interaction	Usability	Incorporate it in designing the tutorial (provide further concise instructions in tutorial session)	Pilot 1
3	Navigation	Replay/Skip Option	Highlight it in the pre-training and tutorial that users can use (existing) replay or skip instructions as needed	Pilot 1

8.2 Augmented Theatre

The Augmented Theatre use case aims to assess the value of AI-generated and AR-displayed subtitles and VFX as perceived by theatre goers. To this end, an excerpt of the play “Hippolytus” was chosen to be performed for the pilot. The play “Hippolytus” was chosen as an indicative sample of a larger category of ancient Greek theatrical plays that pose accessibility difficulties to an international audience and incommensurately disseminate their cultural significance. Unlike in the case of subtitles for 2D monitors used in cinemas or theatres, the novel case of AR captions does not enjoy the benefit of industry standards. In lack of documented best practices and with a goal to prioritize accessibility and usability, extensive user customization options and alternative displaying methods were provided for evaluation by the users.

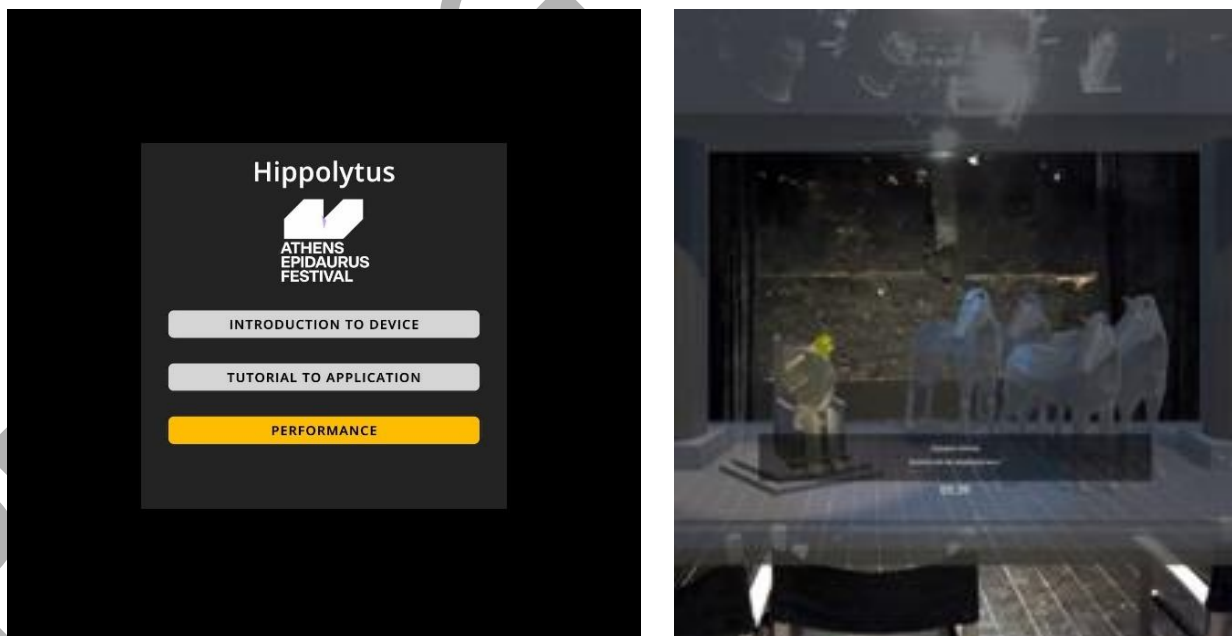


Figure 40 Left: Menu items from "Main Menu" scene. **Right:** VFX spatial matching to theatrical stage in XR client.

The XR client application is comprised of a “Main menu” scene, a “Introduction to the XR device” scene, a “Tutorial to the application” scene, a “Performance” scene and an “Ending credits” scene.

On start-up, the application loads the “Main menu” scene which allows the user to set the language of the application by choosing one of the six (6) available VOXReality languages, and the application mode by choosing one of the three (3) available modes: one with captions positioned in relation to the user (termed “2D” for short), one with captions positioned in relation to the scene (termed “3D” for short) and one with any of the available subtitle positioning scheme, as chosen by the user, accompanied by VFX (termed “VFX”). Choosing an application mode affects the contents of the “Tutorial to the application” and the “Performance” scene, but not the rest of the scenes.

From the “main menu” scene (Figure 41, left) scene the user can navigate to the any of the other scenes and it is suggested to do so in the following order as indicated by the hierarchy of the respective buttons and stressed by automated visual highlights: 1) “Introduction to the XR device”, 2) “Tutorial to the application” and 3) “Performance”. Both the “Introduction to the XR device” (~1’ duration) and the “Tutorial to the application” (~3’ duration) are interactive, step-by-step processes that familiarize the user with the device and the application. The user can opt to perform the introduction and tutorial steps as many times as they need but cannot access the performance scene unless they have completed both at least one time.

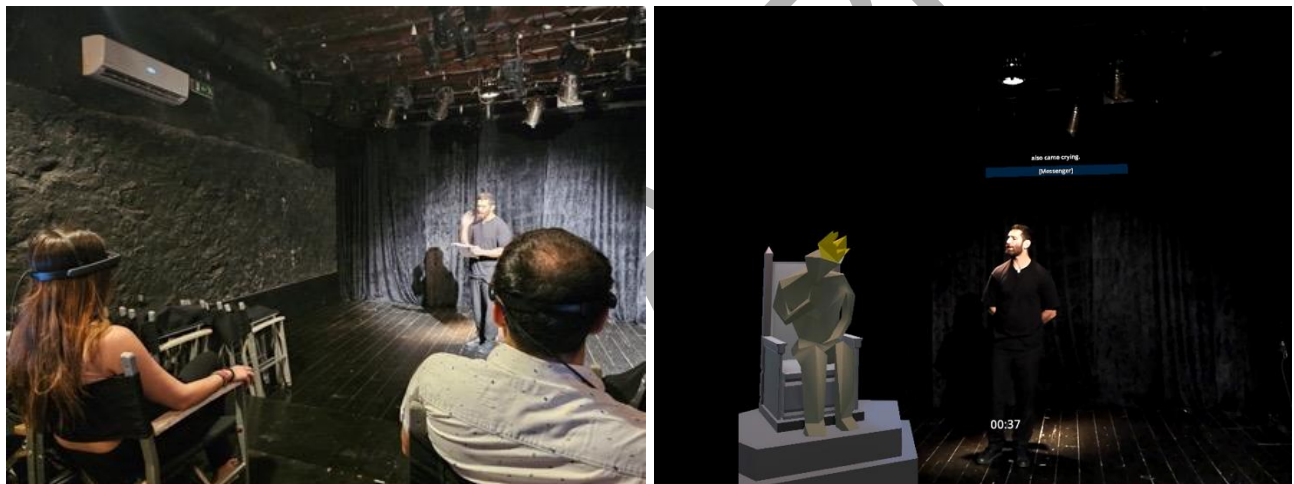


Figure 41 AR Theatre Pilot 1 performance from different perspectives

In the “Introduction to the XR device” the user is instructed to execute all the available input methods at least once. The input methods were limited to single button presses with no button combinations or modifiers like double-taps or press-holds for improved accessibility. In the “Tutorial to the application” the user is instructed to interact with all the available user controls at least once and is advised to customize their preferences, if any, using the provided controls. The user customization settings performed in the “Tutorial to the application” are retained for the “Performance” scene. For the “2D” and “3D” application modes, the settings determine the subtitle font size, subtitle reading distance, text placement options (sub- or sur-titles), text background contrast (opacity of a black background) and overall dimness of the display. The user can also change their language setting at any point through this interface. For the “VFX” mode the user can perform any of the above in addition to toggling the captions on/off, switching between subtitle mode “2D” and “3D” in real-time and toggling the VFX on/off.

The “Performance” scene is an application state which receives, interprets and displays content to the user, as streamed by the server. In this state, the user can finetune their preferences using the same interface as in the “Tutorial to the application” scene if needed. Otherwise, the user is expected to

experience the performance with no other action necessary. Upon completion of the performance, a server message automatically triggers a scene transition to the “Ending credits” scene. In the “Ending credits” the user can see a scrolling list of the consortium team, after which the application automatically closes.

The pilot study involved 12 participants (8 females, 4 males). 25% of the participants belonged in the 18-24 age range, the same percentage was 25-34, while a plurality of 40% belonged in the 35-44 age range, lastly 8% had an age 45-54. Participants represented a variety of native languages: Italian was the most represented with 5 participants, Greek was the second most represented language with 3 speakers, Spanish followed at third place with 2 speaker and French and English were last with 1 participant respectively. All participants reported having at least a university degree, with 6 also having a postgraduate degree and 1 PhD. The most common occupation answer was “Pupil, student, trainee, unpaid intern” with 4 answers, 3 were employees in the public sector, 2 in the private sector, 2 participants were self-employed without employees and lastly 1 participant was self-employed with employees. Furthermore, 50% reported currently working in the cultural sector, while only 2 out of 12 worked in the technology sector.

8.2.1 Conclusion

Our Pilot 1 findings demonstrated both the innovation and the potential of the AR Theatre use case. On the other hand, they revealed limitations, issues and weaknesses that we need to address in the implementation of Pilot 2. The Pilot 1 analysis was crucial in identifying new requirements, derived from both quantitative metrics and qualitative participant feedback. The combination of the quantitative and qualitative analysis including the automated application data analysis has instructed the design modifications of the user requirements for Pilot 2. We are now aiming to redesign the AR Theatre use case accordingly before the 2nd pilot next year (see Augmented Theatre > Requirements). Also, the lessons learnt are summarize below (Table 15).

Table 15 Lessons Learnt from Pilot 1 (Augmented Theatre)

#	Type	Sub-category	Lesson Learnt	Phase Gathered
1	Interface	UI Design	We had anticipated more willingness to explore the UI and personalisation settings however the pilot 1 participants were sometimes hesitant to interact with the UI during the performance fearing they will miss the performance.	Pilot 1
2	Application Tutorial	-	A more intuitive and step-by-step tutorial for the pre-training	Pilot 1
3	Equipment Fitting	-	There was no time limitation to adjust the AR but it was proposed an optional step for comfort. Some users avoided it to save time, but this cost comfort in the long term. We should make this a mandatory step and ideally, find a way to check externally if the fit is good.	Pilot 1
4	Equipment Performance	Avoid overheating of the AR glasses	Performances should not be as condensed as in Pilot 1. Device heating is a limitation as well. That is why we'll try our best to reduce the effects and make participants aware in pre-training. But we cannot completely eliminate it.	Pilot 1

5	Setup / Performance	Connectivity	Provide stable and reliable connections, minimize disconnections and glitches during the performance.	Pilot 1
6	Setup / Performance	Technical Performance	Improve the technical performance to reduce latency, ensure proper pacing, and avoid line delivery issues and other glitches.	Pilot 1
7	Subtitles	Customization	There are no one fits all solutions, responses varied with users preferences, however we are adding some features i.e., black and white frames for contrast with captions	Pilot 1
8	Subtitles	Translation Quality	Ensure translations are accurate and contextually appropriate to maintain the context during the performance.	Pilot 1
9	VFX	Spatial Placement	Ensure VFX are appropriately placed, considering the narrow field of vision of the AR glasses.	Pilot 1
10	Extra	Organisation	Develop reserve lists of participants, ensure back up participants are available for the smooth implementation of Pilot 2	Pilot 1

8.3 Training Assistant

The AR industrial assembly training application focuses on the enhancement of training experience by integrating VOXReality's ASR model and an advanced dialogue agent. The application enables trainees to visualize and manipulate 3D CAD files in an AR setting. It showcases an interactive virtual training assistant with real-time monitoring of performance and a dynamic dialogue system driven by NLP and speech-to-text capabilities. The training initiation involves loading a 3D CAD file, which assists the trainee in accurately assembling parts within the CAD object's structure. It facilitates interaction with pre-designed asset bundles containing essential scene details, interaction scripts, menu functions, and pertinent algorithms. Real-time tracking of performance metrics, such as time allocation and incorrect placements, is done against predetermined benchmarks. Moreover, the training assistant utilizes Hologlight Stream to remotely transmit the application from a high-performance GPU laptop to AR smart glasses, thus overcoming device rendering constraints. The AR training application resides on a laptop (server) and is streamed to the HoloLens 2 (client), amalgamating remote rendering and streaming capabilities to deliver a seamless training encounter.

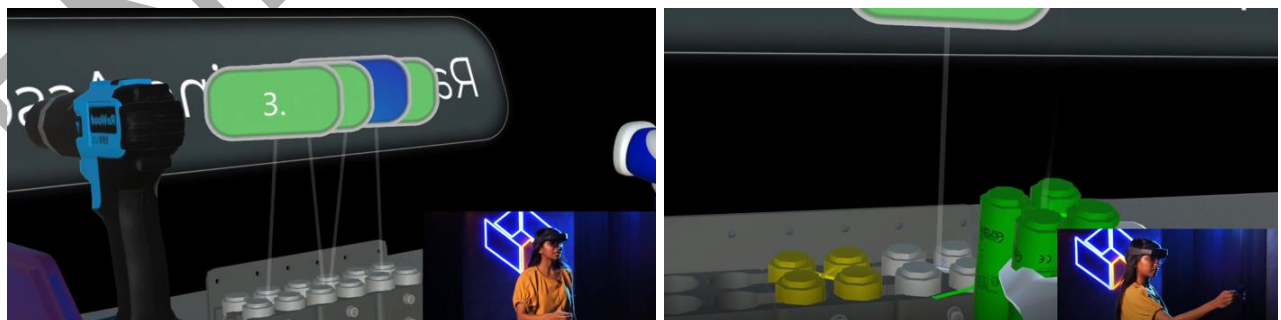


Figure 42 Assembly Demonstration in (AR) Training Assistant Use Case

The training assistant at its completed stage is designed to support the user to communicate with the application without using the hand menu to allow a complete immersion in the assembly process. The functionalities outlined include starting the training, skipping steps, providing documentations / images on demand, assisting proactively when no actions are taken by the user for a fixed period or when the user makes crosses the threshold for repeated wrong moves. The training assistant will also have audio visual cues for feedback.

The Training Assistant at the Pilot 1 stage included the following features and functionalities:

1. Difficulty modes with easy, medium, and hard levels. Easy mode includes assembly preview, visual cues, guiding lines, and object locking. Hard mode lacks visual aids and object locking.
2. Various accuracy modes with simple and advanced snapping. Simple snapping mode facilitates easier attachment, suitable for beginners. Advanced snapping mode requires precise alignment for realistic object placement.
3. Performance indicators with color-coded highlights for correct (green) and incorrect (red) selections.
4. Visual cues and color changes at sequence flags to indicate completion status.
5. Audio cues during the assembly process.
6. Pre-assembly feature allows pre-assembly of parts on the table or shelf before starting the main task.
7. VOXY Avatar is introduced to the assembly environment.
8. VOXY launches and welcomes the user through a text panel display.
9. "VOXY" is used as the keyword for audio trigger to process the audio input.
10. Visual feedback: VOXY turns green and displays the text "listening".
11. VOXY follows the user's visual field to help the user orient. VOXY has additional functionalities:
 - a) start the training; b) open/close text panels; and, c) spatial positioning on command.
12. VOXY answers the user in texts and actions are performed directly without an audio response



Figure 43 Participants - Training Assistant Pilot 1

Ten participants (see Figure 43) volunteered from HOLO's customer base from two locations of MXR Tactics GMBH. Five participants from each location took part in the pilot. They all had a background in software development and had a moderate expertise in manufacturing. The participants had previous experience in AR technology.

8.3.1.1 Conclusion

The pilot 1 evaluation confirms the future prospects and user interest in a voice-driven, AI-powered assistant for AR industry training. In terms of application design and development, more effort is required to provide a functional user experience. Assuming a robust prototype, pilot 2 could benefit from a larger and more complete evaluation protocol to allow a more thorough investigation. Also, the lessons learnt from first phase of Training Assistant Pilot are mentioned in Table 16 below:

Table 16 Lessons Learnt from Pilot 1 (Training Assistant)

#	Type	Sub-category	Lesson Learnt	Phase Gathered
1	Performance	Technical Integration	Improve the microphone input.	Pilot 1
2	Performance	Technical Integration	Improve the dialogue agent integration in the training assistant. The TA (Training Agent) should be reliable, immediate and consistent in their responses and should always return user feedback, even in case of errors.	Pilot 1
3	Performance	User Testing	A table of how many times the service was called, when, how long it took to respond, what command was triggered, etc would help a lot to contextualize the user feedback. Also, performing internal tests and generating such data to compare would also help understand if something went wrong that particular day or if the service integration quality is low overall.	Pilot 1
4	User Interface	Improvement	Improve the UI to avoid overlapping information with the virtual training scene	Pilot 1
5	User Interface	Feedback	Provide clear visual or audio responses before the virtual agent performs any actionable tasks. General note: It would be interesting to see which instructions were given to the users in advance to communicate with Training Assistant (Voxy). E.g., if it was in the form of a demonstration with "example sentences" or a video or verbal instructions etc.	Pilot 1
6	User Interface	Feedback	Provide visual/audible indication when the virtual agent is listening	Pilot 1
7	Virtual Agent	Agent Location	Place virtual agent at a stable, defined location in the AR environment	Pilot 1
8	Virtual Agent	Behavior	Virtual Agent initiates help that is beneficial to the user	Pilot 1

9 New Requirements (Year 2)

This chapter provides a comprehensive overview of the user requirements for the VR Conference, Augmented Theatre, and Training Assistant use cases. It is structured to clearly convey the status of each set of requirements into two broader categories (A) Closed and (B) Open (see Figure 44). **Closed** requirements refer to the ones that are fully completed or have been opted-out due to budgetary limitations, developmental constraints, requiring disproportionate effort relative to their utility in demonstrating model performance or being out of the project's scope. These requirements are no longer active and are considered finalized. Whereas, the **Open** category consists of requirements that are still active and need to be addressed. This category is further divided into three subcategories: Pending requirements from previous phases, Implementation confirmed for Pilot 2, which includes requirements that have been approved for execution in the next development phase, and New (extended) requirements from Pilot 1, which are additional requirements identified during the first pilot that need to be integrated into the project moving forward.

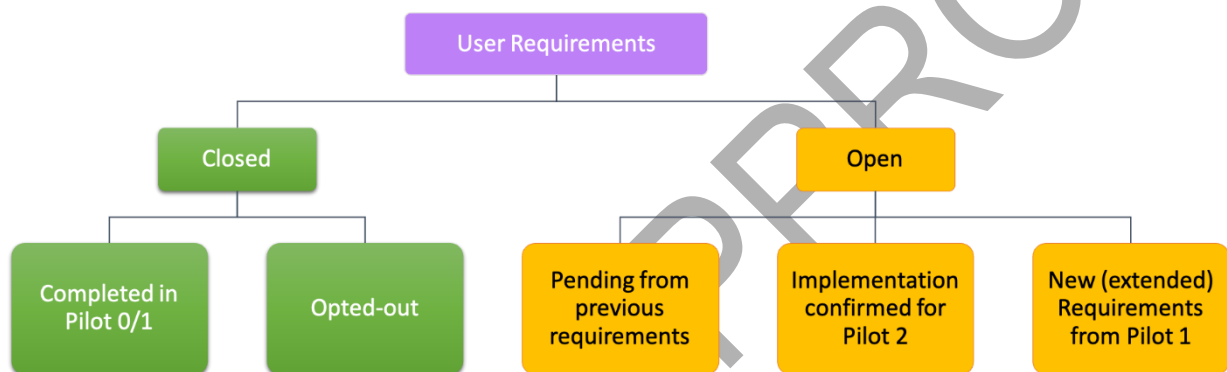


Figure 44 Status of User Requirements

The coming sections includes detailed tables for each of the categories that outline the relevant requirements, offering insights into what has been accomplished, what remains to be addressed, and what has been deliberately excluded or identified for future consideration. This organization aims to showcase what has been achieved so far until pilot 1, set clear expectations and provide a structured summary of progress and planning for each use case for next pilot.

9.1 VR Conference Requirements

This section addresses the user requirements for the VR Conference use case, categorizing them into completed, open/pending, out-of-scope, and extended requirements. The objective is to provide a clear distinction between what has been achieved in Pilot 1, identify functionalities still under consideration for Pilot 2, and outline the requirements excluded from the scope. Additionally, it highlights newly discovered extended requirements. Although the out-of-scope requirements will not be implemented in the current VOX project, they are included in this review due to their relevance and potential for informing future work or Open-Calls Projects. This approach ensures a comprehensive understanding of the project's progress and future directions.

9.1.1 Completed Requirements

In Phase 1 of the Pilot 1 for the VR Conference use case, a number of key requirements were successfully fulfilled (see Table 17). The implementation included the creation of a fully immersive VR experience, with features such as user navigation and language translation seamlessly integrated into the virtual conference environment. A virtual assistant was developed to assist users with navigation, interaction, and information exchange, and the system supported multiple languages, with English as the primary operational language.

The virtual conference venue effectively mirrored a professional setting, with users represented as virtual avatars and presentations adhering to a predefined conference format. Key aspects such as session duration, participant demographics, and assessment methods were all addressed. Additionally, dedicated virtual agents were provided to each user, offering support throughout the conference, including navigation assistance and on-demand interaction.

Furthermore, users benefited from features such as quick-start tutorials, virtual maps, visual navigation cues, and comprehensive subtitle options. The interface and user interactions were designed to enhance usability, with accessible help options and control features. Overall, the successful implementation of these requirements has established a robust foundation for the VR Conference application, setting the stage for future enhancements and refinements.

Table 17 VR Conference: Completed Requirements in Pilot 1

#	Type	Completed Requirements	Revised Priority	Achieved in Pilot 1
1	General	The experience is in VR.	High	YES
2	General	Provide user navigation and language translation at the virtual conference.	High	YES
3	General	Develop a virtual assistant providing users with relevant information to navigate virtual spaces, interact with other visitors, and exchange relevant information.	High	YES
4	General	The operational output language is English, and it can be translated into other languages.	High	YES
5	Scenario	The conference venue should imitate the environment of a professional conference setting.	High	Partially YES
6	Scenario	The users will be represented as virtual avatars.	High	YES
7	Scenario	The presentation should be in a predefined conference format.	High	YES
8	Scenario	The suggested duration for a single VR session is 30 minutes at maximum.	High	YES
9	Assessment	The target user includes conference partners and visitors with diverse backgrounds.	High	YES
10	Assessment	The suggested number of participants per session is ~10 to 20 people.	Medium	YES

11	Assessment	The quality of experience will be assessed by a questionnaire followed by a semi-structured interview.	High	YES
12	Virtual Agent	A dedicated virtual agent will be assigned to each user and will stay with him/her during the complete duration of the conference.	High	YES
13	Virtual Agent	The virtual agent should provide welcome greeting.	High	YES
14	Virtual Agent	The virtual agent will help user to navigate the space, answer programme-related questions, and deliver relevant information.	High	YES
15	Virtual Agent	The users should have a "skip" option for the virtual assistant help.	High	YES
16	Virtual Agent	The communication with virtual agent should be enabled by voice typing.	High	YES
17	Virtual Agent	The virtual agent should interact with the users on-demand, and should not be intrusive.	High	YES
18	Virtual Agent	The virtual agent could look like a cartoony avatar.	Medium	YES
19	Navigation	Quick get started tutorial option should be available.	High	YES
20	Navigation	Virtual map should be available to navigate the conference venue.	High	YES
21	Navigation	Navigation to the target location should be guided by visual cues, arrows, marks/lines on floor etc depending on the context.	High	YES
22	Navigation	Users should be assisted for quick and easy navigation between different places in the venue.	High	YES
23	Subtitles	Translation should be available in textual format.	High	YES
24	Subtitles	Option to mute the complete room or individual speaker(s) should be available.	High	YES
25	Subtitles	Users should be able to control the volume +/-	High	YES
26	Subtitles	Users should be able to turn on/off the program, subtitles, auto-translation and voice using interactive buttons.	Medium	YES
27	Subtitles	For multi-speakers, the subtitles should highlight and differentiate the active speaker among group of speakers.	Medium	YES
28	Subtitles	For head-mounted display, the subtitles should be visible and placed without masking other important parts in the display/screen.	High	YES

29	Interface	A help button should be visible on the corner of the screen.	High	YES
30	User Interaction	Users can interact with the agent and ask the agent questions during the navigation.	High	YES
31	User Interaction	Users should have options to use interactive-buttons on the screen to interact with the system.	Medium	YES
32	Extra	Dialogue Agent should be available on demand.	High	YES
Total Completed Requirements				32/49

9.1.2 Opted-out Requirements

Several user requirements (see Table 18) were intentionally excluded from the VR Conference use case due to being out of scope for the project. The decision to drop these requirements was based on several factors including their misalignment with the project's development cost, their limited impact on the overall user experience, and their lack of suitability for the application. The excluded requirements encompass various aspects such as predefined gestures for virtual avatars, personalized info feeds from virtual agents, accessibility via smart-watches, and the auto-saving of information. Additional dropped features include flyover views with zoom options, voice-based translation, customizable dashboards, and options for hand-gesture interactions. While these features were considered initially, their development costs were deemed disproportionate to the benefits they would provide, and some were found to be inconsistent with the core objectives and realistic vision of the VR Conference. Although these out-of-scope requirements will not be implemented in the current phase, they were documented for their potential value in future projects or research initiatives.

Table 18 VR Conference: Out-of-scope Requirements from Pilot 1

#	Type	Out of Scope Requirements	Revised Priority	Achieved in Pilot 1
1	Scenario	The virtual avatars of the participants will offer predefined set of gestures.	High	NO
3	Virtual Agent	Virtual agents deliver personalized info feeds related to the conference programme-related activities.	Medium	NO
4	Virtual Agent	The virtual assistant can be accessible via a smart-watch / wrist-band.	Low	NO
5	Virtual Agent	The virtual-assistant should auto-save relevant information and provide them in an exportable file.	Medium	NO
6	Virtual Agent	Virtual agent should be likeable, friendly, pleasant, customizable, realistic and complete.	Low	NO
7	Navigation	Flyover view of the venue should be available with zoom in/out options.	Low	NO
8	Subtitles	Translation should be available in voice (audible) format.	Low	NO
9	Subtitles	The standard/default subtitles should be displayed on the bottom of the screen and in particular cases the subtitles should be placed relative to the virtual environment context.	High	NO

10	Interface	There should be a customizable dashboard visible at all times.	Medium	NO
11	Interface	In the conference room, users can choose the in-room view options i.e., full screen and/or conference room view.	Medium	NO
12	User Interaction	Users should have options to use hand-gestures to interact with the system.	Low	NO
13	Extra	Users should be able to make digital and crypto payments in the conference items shop.	Low	NO
Total Out-of-scope Requirements				12/49

9.1.3 Open/Pending Requirements

For the VR Conference use case, several requirements remain open or pending as the project progresses to Pilot Phase 2. These include enhancements to the virtual agent's functionality, such as providing contextual help, FAQs, and actionable prompts relative to user activities. Additionally, the capability for text-based communication with the virtual agent and the presentation of relevant options at the end of sessions are yet to be fully realized. Interface improvements for displaying audience questions near the presenter and enabling user engagement with speakers during Q&A sessions are also being addressed. These open requirements will be the focus for further development in the upcoming phase to enhance the overall user experience.

Table 19 VR Conference: Open Requirements from Pilot 1

#	Type	Open Requirements	Revised Priority	Achieved in Pilot 1	Status for Pilot 2
1	Virtual Agent	The virtual agent should provide help, FAQs and prompt the users about the context/content of the venue and potential to-do (action) items in advance or relative to the current activity.	High	NO	Partially yes
2	Virtual Agent	The communication with virtual agent should be enabled by text.	Low	NO	Open
3	Virtual Agent	The virtual agent should present users with most relevant available options at the end of every session.	High	NO	Open
4	Interface	Questions/feedback from the audience should appear textually/visually near the presenter's screen.	Medium	NO	yes
5	User Interaction	Users should be able to engage with speakers and guests during Q&A session(s).	High	NO	yes
Total Open Requirements				05/ 49.	

9.1.4 New (Extended) Requirements from Pilot 1

Following the data analysis from Pilot 1, several extended user requirements (see Table 20) have been identified for the VR Conference use case. These requirements were discovered to further enhance the user experience and address areas of improvement. The visual design of the VR

conference environment should be refined by incorporating warm colors, ambient lighting, and additional conference elements such as posters and banners to improve aesthetic appeal.

Stable and reliable connectivity must be ensured to provide a seamless experience. Additionally, the text-box responses should be adjusted for better readability and clarity. Improvements are also needed in the accuracy and contextual relevance of translations for subtitles. Furthermore, while the addition of audible instructions by virtual agents is recognized as a valuable enhancement, it remains open for further development. These extended requirements will be considered for implementation in subsequent phases to optimize the overall functionality and user experience of the VR Conference.

Table 20 VR Conference: Extended Requirements from Pilot 1

#	Type	Sub-category	Extended Requirement	Priority	Status	Phase Gathered
1	Interface	Visual Design	Improve the aesthetic appeal of the VR conference environment by using warm colors, ambient lighting and adding conference items (posters, banners etc.).	Medium	Yes	Pilot 1
2	User Interaction	Connectivity	Provide stable and reliable connections	High	Yes	Pilot 1
3	Navigation	Instructions	Adjust the response text in the text-box appropriately for better readability.	High	Yes	Pilot 1
4	Subtitles	Translation	Enhance translation accuracy and contextual relevance	High	Yes	Pilot 1
5	Virtual Agent	Function	Enable agents to provide audible instructions in addition to already existing text-based instructions	Medium	Open	Pilot 1
Total New Extended Requirements						5

9.2 Augmented Theatre

This section provides a comprehensive overview of the user requirements for the Augmented Theatre use case, categorizing them into completed, open/pending, out-of-scope, and extended requirements. The purpose is to clearly outline what has been achieved in Pilot 1, identify functionalities that are still under consideration for Pilot 2, and clarify which requirements were excluded due to being out-of-scope. Additionally, newly discovered extended requirements are highlighted. Although the out-of-scope requirements will not be implemented directly in the current VOXReality project, they are included in this review due to their relevance and potential value for future initiatives or OpenCalls Projects.

9.2.1 Completed Requirements

In Phase 1 of the pilot for the Augmented Theatre use case, several key requirements were successfully met (see Table 21). The successful completion of these requirements includes the introduction of AR technologies to theatre audiences, effective language translation, and VFX integration. Additionally, the experience was designed to ensure minimal disruption to the audience's theatre experience and provided clear and accessible AR glass setup. With these foundational

requirements in place, the theater application skeleton is now functional and ready to support more advanced features in subsequent development phases for pilot 2.

Table 21 Augmented Theatre: Completed Requirements from Pilot 1

#	Type	Completed Requirements	Revised Priority	Achieved in Pilot 1
1	Objective	The objective is to introduce AR technologies to theatre audiences.	High	YES
2	Objective	Language translation and VFX experience with AR glasses at the theatre.	High	YES
3	General	The experience is in AR.	High	YES
4	General	The gestures and setup to operate AR glass should not intervene other audience's theater experience.	High	YES
5	General	The possibility of the audience getting distracted by the additional information from AR glasses has to be considered.	Medium	YES
6	Setup	The experiment takes place at the avant-garde theatre.	High	YES
7	Setup	The AR glass setup should be easy and clear to audiences who are not familiar with the AR technology.	High	YES
8	Setup	AR glass should give the audience enough controls but also limit their ability to change default/preset settings.	Medium	YES
9	Scenario	A scene from the play "Hippolytus" by Euripides will be specially produced and played.	High	YES
10	Scenario	The length of the performance will not exceed 15 minutes.	Medium	YES
11	Assessment	Two users can watch the play at the same time with two AR devices in a single time.	High	YES
12	Assessment	The target user is a theatergoer with diverse background whose native language is any of the VOXReality languages.	High	YES
13	Assessment	The target user includes who don't know the spoken language of the play.	High	YES
14	Assessment	In Greek play, non-Greek audiences should participate in the experiment.	Medium	YES
15	Assessment	Up to 50 people will participate in the complete experiment.	Medium	YES
16	Assessment	The translation of the play will be available in all VOXReality languages.	High	YES
17	Assessment	The audiences may experience either one or both of the automatic translation and the visual effects technology.	High	YES
18	Assessment	The user evaluation may include combined demonstration of translation and VFX components	Medium	YES
19	Assessment	The creators of the performance may be asked to evaluate the experience through semi-structured interviews.	High	YES
20	Interface	The AR glass could have menu on the screen to change practical settings.	High	YES
21	Interface	The interface should not include extensive head movements.	Medium	YES
22	Subtitles	The subtitles should be real-time and available in any of the VOXReality languages	High	YES
23	Subtitles	The subtitles should provide different caption sizes and toggle options wherever feasible.	Medium	YES
24	Subtitles	The audience should be able to fine-tune the placement of the subtitles.	Low	YES

25	Subtitles	The subtitles should not overlap with the stage setup.	High	YES
26	Subtitles	There should be various caption styles (standard, speech bubble etc.).	High	YES
27	VFX	VFX can reflect or accompanies the narration of the scene.	High	YES
28	VFX	Certain words or phrases will trigger VFX.	High	YES
Total Completed Requirements				28/48

9.2.2 Opted-out Requirements

In the AR Training Assistant use case, several requirements were intentionally excluded (see Table 22) from the scope due to their limited impact on overall effectiveness, high development costs, and the significant budget and man-month efforts required for their implementation. The decision to drop these requirements was influenced not only by the substantial resources needed but also by the fact that some features did not align well with the artistic vision of the Theater Application. For instance, interactive features like zooming were deemed counterproductive to the intended immersive experience. Consequently, these features were deemed unsuitable for pilot 2. Although these out-of-scope requirements will not be implemented with VOXReality, they are enlisted here because they hold relevance and potential value for future initiatives or OpenCalls Projects.

Table 22 Augmented Theatre: Out-of-scope Requirements from Pilot 1

#	Type	Out-of-Scope Requirements	Revised Priority	Achieved in Pilot 1
1	Assessment	The target user includes who cannot hear well.	Low	NO
2	Assessment	The target user base includes people with vision correction glasses.	Low	NO
3	Interface	Audience can zoom in or out the scenes.	Low	NO
4	Interface	The UI of the AR screen should consider audience accessibility.	Low	NO
5	Subtitles	The default subtitle starts with the local language.	Low	NO
6	VFX	The director specifies the capability of the users to change AR Glass settings and subtitles.	Low	NO
Total Out-of-Scope Requirements				06 / 48

9.2.3 Open/Pending Requirements

In preparation for Pilot Phase 2 of the Augmented Theatre use case, several requirements remain open (see Table 23) or pending from Phase 1. These include enhancements to the interface, such as providing information about the play or scene and ensuring the AR glass elements are compatible with stage lighting and settings. Subtitle functionalities that are still to be addressed involve contextual actor-specific subtitles, options for language changes, and maintaining the original language audio. Additionally, there are pending aspects related to VFX, including the need for pre-discussion of VFX implementation with the script, ensuring that visual effects do not exceed the stage boundaries or interfere with actor performance, and enhancing immersion and artistic relevance.

Table 23 Augmented Theater: Open Requirements from Pilot 1

#	Type	Open Requirements	Revised Priority	Achieved in Pilot 1	Status for Pilot 2
1	Interface	Audience can learn about the play or the scene.	Medium	NO	yes
2	Interface	The elements in the AR glass screen should be feasible respective to the lighting, brightness, contrast and related settings of the stage (theatre).	Low	NO	yes
3	Subtitles	The subtitles could follow each actor depending on the context.	Low	NO	Open
4	Subtitles	The default audio starts with the original spoken language of the play.	Low	NO	Open
5	Subtitles	The audience should have an option to change the language audio.	Low	NO	Open
6	Subtitles	The audience should have option to listen to the original language.	Low	NO	Open
7	VFX	The VFX implementation should be discussed earlier with script and the plot.	High	NO	yes
8	VFX	A narrator could be the person who triggers the VFX.	Low	NO	Open
9	VFX	Location of the visual effects should not exceed the range of the stage.	Low	NO	Open
10	VFX	The implemented VFX should not affect the actor's actions on the stage.	Medium	NO	Open
11	VFX	The VFX should help the audience be immersed in the performance.	Low	NO	Open
12	VFX	Style of the visual effects should be artistically relevant to the scene.	Medium	NO	Open
13	Extra	Audience can read about the play and characters or fun facts during the session.	Medium	NO	Open
14	Extra	Audiences can watch "behind the scene" (BTS) of the play before and after watching the play.	Low	NO	Open
Total Open Requirements					14 / 48

9.2.4 New (Extended) Requirements from Pilot 1

Pilot 1 results supported the identification of four (4) new user requirements, as detailed in Table 24. The extended user requirements were categorized under the types of the initial user requirements with two (2) new user requirements falling under the type "Setup", one (1) under the type "Subtitles" and one (1) under the type "VFX".

Table 24 Extended AR Theatre user requirements

#	Type	Sub-category	Extended Requirement	Priority	Status	Phase Gathered
1	Pre-training	-	The pre-training phase should include a step-by-step introduction guide (using verbal and visual format) for the complete AR theatre experience	High	Yes	Pilot 1
2	Application Tutorial	-	A simple step-by-step application tutorial to guide the users for the AR Theatre experience	High	Yes	Pilot 1
3	Subtitles	Customization	Offer extensive readability options for displaying the captions to improve the user experience	Medium	Yes	Pilot 1
4	VFX	Aesthetics, immersion	Develop more elaborate VFX to increase enjoyment and immersion	High	Yes	Pilot 1
Total New Extended Requirements						4

Sub-categories were used to provide further overview. Under “Setup”, one requirement was tagged as “Pre-training” and one as “Application tutorial”, while under “Subtitles” the requirement was tagged as “Customization” and under “VFX” the requirement was tagged as both “Aesthetics” and “Immersion”.

The “Setup/Pre-training” user requirement was phrased as “The pre-training phase should include a detailed presentation (using verbal and visual format) for the complete AR theatre experience “. This pre-training phase is proposed to be added at the beginning of the pilot, right after the Consent Forms have been signed. The presentation should include a brief overview of the concept of the AR theatre, the experiment's goal and protocol, and general information about the XR device's capabilities. The presentation should be delivered using traditional presentation media (verbal, handouts, demonstration). A demonstration by a human facilitator of how to wear and use the XR device could also be included. This requirement was added to ease the anxiety and uncertainty of the participants who were unfamiliar with the equipment, the technology and the overall concept of AR Theatre, and thus restrained in their engagement.

The “Setup/Application Tutorial” user requirement was phrased as “A step-by-step application tutorial should be provided to guide users with no AR familiarity.” The requirement to have a tutorial in a use case does not need much elaboration. The important keyword in this requirement though is the phrase “step-by-step”. Whereas in Pilot 1, the interactive “Introduction to the XR device” was a step-by-step process, the “Application tutorial” was structured as a sandbox, with a prompt to users’ to experiment with the settings as much as they wanted until they were satisfied with the result. This was designed with the reasoning that users would very likely have never used the particular XR device before, therefore a linear approach for the XR device would be suitable, but users could have varying levels of interaction skills, therefore the application-specific tutorial should be more freeform to accommodate various user backgrounds. However, this configuration led to conflicting user feedback, such as complaints of overexplaining simple concepts, but not adequately informing of core concepts. Also, although no time limit was imposed, users reported “needing more time”, while others completely missed the prompt to experiment and the fact that interactions were even possible to begin with. For those reasons, it was decided that all tutorial content should be structured as a step-by-step process, at the risk of imposing on more XR-savvy users.

The “Subtitles/Customization” user requirement was phrased as “Offer extensive readability options for displaying the captions”. This was in acknowledgement of the fact that users appreciated the existing customization options, even if using them required higher interactions skills. In fact, users actively suggested including further customization options, such as adding a colour scheme inversion option (black letters on white background or white letters on black background), or supporting multiple font weights, or the ability to choose a dyslexic-friendly font, etc., most of which were with accessibility and readability in mind. Despite the risk of introducing yet more mental load to an audience struggling with the technology being used, the requested features are aligned with an inclusive design approach and were included as an extended user requirement.

Finally, the “VFX/Aesthetics, immersion” requirement was phrased as “Develop elaborate VFX to increase enjoyment and immersion.” Although the VFX in Pilot 1 were low-level proxies, they generated overwhelmingly positive user feedback which highlighted the VFX value as an element increasing the users’ sense of presence and as an aesthetic artifact with its own enjoyment potential. The users’ commented on the fact that VFX extending from the field of view of the AR glasses seemed to reduce their enjoyment and increase their fear of missing out, and that cases of VFX overlapping with subtitle elements or physical elements seemed to break the plausibility illusion, thus reducing their sense of presence.

9.3 Training Assistant

This section provides a comprehensive overview of the user requirements for the Training Assistant use case, categorizing them into completed (sec 9.3.1), open/pending (sec 9.3.3), out-of-scope (sec 9.3.2), and extended requirements (sec 9.3.4). The purpose is to clarify what has been achieved in Pilot 1, outline the functionalities planned for Pilot 2, identify requirements that were intentionally excluded due to being out-of-scope, and address newly discovered extended requirements. Although the out-of-scope requirements will not be implemented in the current VOXReality project, they are included here to acknowledge their origin from a high-quality, user-centric pipeline and to highlight their potential value for future developments or Open Calls (OC+) Funded Projects.

9.3.1 Completed Requirements

In Phase 1 of the pilot, a number of critical requirements for the AR Training Assistant were successfully fulfilled (Table 25), as detailed in D5.3. This phase established the fundamental objective of delivering a guided augmented reality industrial training scenario. The interface has been designed to provide appropriate color schemes and visual cues, and feedback mechanisms include both visual and audible notifications to indicate task correctness. These completed requirements formulate a base infrastructure for subsequent phases of the pilot 2.

Table 25 Training Assistant: Completed Requirements in Pilot 1

#	Type	Requirements	Revised Priority	Achieved in Pilot 0
1	Objective	The goal is to provide a guided augmented reality (AR) industrial training scenario.	High	YES
2	Setup	The external environment having a defined location is not essential.	High	YES
3	Setup	Three levels - easy, medium, and difficult - will be offered.	High	YES
4	Assessment	The operational language is English.	High	YES

5	Assessment	One user(trainee) at a time will be in the AR training environment.	High	YES
6	Assessment	At least 3-5 users will be assessed in total.	High	YES
7	Interaction	The AR experience involves the virtual manipulation of object components in the training environment by the user.	High	YES
8	Interaction	User interaction is with hands, without controllers and may also verbally engage with the virtual agent.	High	YES
9	User Interface	The interface should provide color schemes and visual cues appropriate for the action.	High	YES
10	Feedback	Users should have both visual / audible notifications regarding the correctness of the tasks at hand.	High	YES
11	Feedback	The feedback should appear automatically in case of incorrect assembly.	High	YES
12	Assessment	The target user has no prior knowledge of the specific industrial assembly task performed in the use case.	High	YES
13	Assessment	A round of training and experiment for the use case scenarios pilot could take approximately 1-3 hours.	Medium	Yes
14	Assessment	The user's assembly training experience will be evaluated with an emphasis on the interaction with the virtual agent	High	Yes
15	Assessment	The user feedback will be collected via questionnaires, survey and/or direct verbal feedback through structured OR semi-structured interviews.	High	Yes
Total Completed Requirements				15/38

9.3.2 Opted-out Requirements

In developing the AR Training Assistant, three user requirements (Table 26) were opted out due to their out-of-scope nature and high development costs relative to their effective output. The first requirement was a customizable dashboard for common functions, which would have significantly affected development time and costs without a proportional increase in user effectiveness or satisfaction. The second requirement was the option for users to choose between textual instructions or visual cues/highlights. Despite its high priority, this feature was complex and counter-productive for the application usability, there it was decided to focus on core functionality. The third requirement was providing users with options to control narration, such as turning it on/off or adjusting its speed. This added considerable complexity and cost, and given its lower impact on overall training effectiveness, it was deemed out of scope for the initial phase. These decisions were made to prioritize essential features and deliver a functional and effective pilot version of the AR Training Assistant in Pilot 2.

Table 26 Training Assistant: Out-of-scope Requirements from Pilot 1

#	Type	Out-of-Scope Requirements	Revised Priority	Achieved in Pilot 0
1	User Interface	There should be a dashboard for most common functions, which users can customize and personalize for their needs.	Low	NO
2	Interaction	Users should be able to choose between textual instructions OR visual cues/highlights for the training depending on the training sequence.	High	NO
3	Interaction	Users should have options to turn on/off, speed-up or slow-down the narration.	Low	NO
Total Out-of-Scope Requirements				03 / 38.

While the out-of-scope requirements will not be implemented directly in the current VOX project, they are included in this deliverable due to their technical relevance and potential value for future initiatives or Open Calls Projects.

9.3.3 Open Requirements

Several existing requirements (Table 27) identified in Pilot 1 remain open or pending for implementation in Pilot 2 of the AR Training Assistant. These include enhancements across various areas such as virtual agent functionality, assessment metrics, user interface, interaction capabilities, feedback mechanisms, and setup conditions. These pending requirements are critical for refining the training tool and ensuring it meets the comprehensive needs of its users.

Table 27 Training Assistant: Open Requirements from Pilot 1

#	Type	Open Requirements	Revised Priority	Achieved in Pilot 0	Status for Pilot 2
1	Virtual Agent	The training session should be in both text and/or audio/video format.	Medium	NO	Open
2	Assessment	The assessment metrics may include time spent doing a virtual task, number of incorrect steps, number of times the virtual agent interceded to help - etc.	High	NO	yes
3	User Interface	Access to help and FAQs should be available on screen at all times.	Low	NO	Open
4	Interaction	User can ask (virtual agent / training assistant) for instruction to assemble a machine/object.	High	NO	Yes
5	Feedback	Feedback should be as close to immediate as possible and self-explanatory when requested or needed.	High	NO	Yes
6	Virtual agent	Virtual agents can provide verbal / textual explanation and display corresponding visual contents.	High	NO	Yes
7	Virtual agent	Virtual agents will be animated in 3D avatar format.	Medium	NO	Yes

8	Virtual agent	If/when assistance is required, virtual agents may direct the user towards objects that are necessary for the application or assembly task.	High	NO	Yes
9	Virtual agent	There should be on-demand help during the interaction.	High	NO	Yes
10	Virtual agent	The help should include hints and quick tips to guide the users	High	NO	Open
11	Virtual agent	The virtual agent should monitor the user action and aid help demonstrate the user on how to assemble a part correctly.	High	NO	Open
12	Virtual agent	In case of errors, the virtual agent should intuitively guide users towards the solution.	Low	NO	Yes
13	Virtual agent	Support from a virtual agent may come in the form of documentation, such as images or PDFs, videos.	Medium	NO	Yes
14	Setup	A location for demo would be a standard industrial shop environment or similar lab setup.	High	NO	Open
15	Setup	An ideal location for test demo is an industrial shop or similar which approximate a canonical (standard) environment.	Medium	NO	Open
16	Assessment	The target user is an experienced worker with some assembly experience and knowledge of assembly training.	High	NO	Yes
17	Assessment	The target user may or may not have much experience with augmented reality.	Medium	NO	Yes
18	Assessment	The target user will cover all possible end-user demographics, including equal numbers of genders.	Low	NO	Open
19	Virtual Agent	The virtual agent should intervene and offer automatic help after a certain number of incorrect steps by users.	High	NO	Yes
20	Feedback	The training and assembly should end with results and feedback.	High	NO	Yes
Total Open Requirements					20/38

9.3.4 New (Extended) Requirements from Pilot 1

The data analysis from Pilot 1 of the AR Training Assistant has identified several extended user requirements that will be implemented in the next phase of the pilot to enhance the system's performance, user interface, and virtual agent's behavior. Firstly, in the performance category, the system will enable microphone input transmission from the HL2 client to the laptop application,

facilitating better interaction between the hardware components. Additionally, the available interactions between the XR application and the dialogue agent will be extended, allowing for more dynamic and responsive user experiences. Furthermore, internal technical testing will be performed to gather metrics in laboratory conditions for comparison with the metrics from Pilot 2.

Table 28 Extended AR Training user requirements

#	Type	Sub-category	Extended Requirement	Priority	Status	Phase Gathered
1	Performance	Technical Integration	Enable microphone input transmission from HL2 client to laptop application	High	Yes	Pilot 1
2	Performance	Technical Integration	Extend the available interactions between the XR application and the dialogue agent	High	Yes	Pilot 1
3	Performance	User Testing	Perform internal technical testing to gather metrics for comparison in laboratory conditions to compare with the metrics gathered in pilot 2	High	OPEN	Pilot 1
4	User Interface	Improvement	Improve the UI to avoid overlapping information with the virtual training scene	High	Yes	Pilot 1
5	User Interface	Feedback	User's should be provided with either visual or audible feedback before agent initiates a task	High	Yes	Pilot 1
6	User Interface	Feedback	Provide visual/audible indication when the virtual agent is listening	Medium	Yes	Pilot 1
7	Virtual Agent	Agent Location	Place virtual agent at a stable, defined location in the AR environment	Medium	Yes	Pilot 1
8	Virtual Agent	Behaviour	Virtual Agent initiates help that is beneficial to the user	Medium	OPEN	Pilot 1
Total New Extended Requirements						8

In the user interface category, several improvements are planned. To avoid overlapping information with the virtual training scene, the UI will be enhanced. Users will be provided with visual or audible feedback before the agent initiates a task, ensuring clarity and preparedness. Additionally, there will be visual or audible indications when the virtual agent is listening, enhancing user interaction and engagement. For the virtual agent, medium-priority requirements will also be addressed. The virtual agent will be placed at a stable, defined location in the AR environment, ensuring consistency and reliability in its presence. Efforts will also be made to ensure that the virtual agent initiates help that is beneficial to the user. These extended requirements, derived from Pilot 1, will be utilized in the next phase to refine the AR Training Assistant, ensuring a more robust, user-friendly, and effective training tool.

10 Conclusion

This deliverable is an extended version of D2.1 (Definition and Analysis of VOXReality Use Cases V1), which earlier provided a base analysis of the VOXReality's use cases, a list of technical and user requirements, initial information flow diagram towards the architecture of the system and the development/deployment plan. The initial results for previous deliverable were based in a two-phases user centric process, where end-users are placed in the centre of the methodology. Initial requirements were gathered in the first phase, consisting of weekly calls between the different stakeholders of the project and an initial description of the use cases. Based on the initial requirements, targeted interactive focus groups were organised in three in-situ visits (VRDays, AEF, HOLO). These focus groups provided sufficient results for the initial identification of the final version of the requirements and the final description of the use cases.

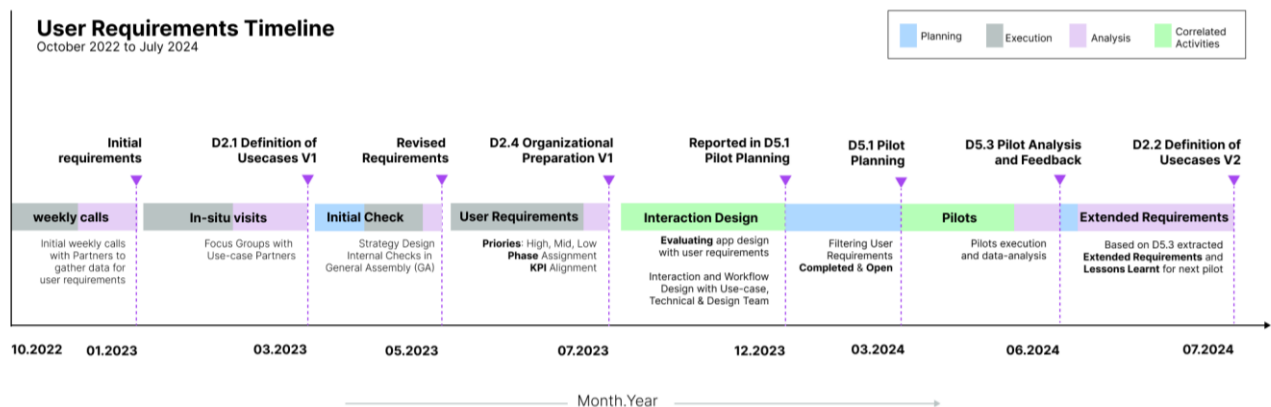


Figure 45 User Requirements Methodology Timeline

This current version provides an updated version of the detailed analysis of the VOXReality's use cases, with an extended list of user and technical requirements. Mainly elaborating on the user-centric methodology adopted in a series of activities (see Figure 45) i.e., Weekly Calls, In-Situ Visits, Requirements Revision with KPIs, Interaction Design, Technology Maturity Plan, Pre-Pilots, Pilots, and Data Analysis of Pilots in Year 1 and Year 2 of the project. Detailed planning and analysis were conducted, as documented in previous deliverables D5.1 (Pilot Planning and Validation V1) and D5.3 (Pilot Analysis and Feedback V1), to establish a coherent framework for the pilot implementations. We'd like to refer the reader for to these two deliverables D5.1 and D5.3 for a detailed information regarding the pilots and the analysis outcomes respectively. Whereas, here we have provided a brief summary of the pilots here to provide better context for our methodology. The outcomes of the activities were used to gather the current refinement of the Definition and Analysis of VOXReality Use Cases Version 2.

The core contributions of this deliverable are presented in [Section 9](#), which introduces the Extended Requirements (Version 2) with sub-divided tables of Achieved, Open/Pending and Opt-out requirements. This section also reflects the lessons learned and feedback received during the initial phases, providing an updated list of requirements for each use case that will inform the project's future direction. Overall, the benefits and drawbacks gained from these pilots offered a benchmark for future research and development, guiding open call projects and the second pilot phase. The pilot findings are critical for the upcoming deliverables (see Figure 10) D2.5 (Organisational preparation for VOX pilot scenarios and PRESS analysis V2) and D5.2 (Pilot planning and validation V2). This deliverable report concludes with a summary of key findings and recommendations eventually supporting the planning of the second phase of pilots and Open Calls for the project.

References

- [1] M. Grieves and J. Vickers, "Digital twin: Mitigating unpredictable, undesirable emergent behaviour in complex systems," *Transdisciplinary Perspectives on Complex Systems: New Findings and Approaches*, pp. 85–113, Jan. 2016, Doi: 10.1007/978-3-319-38756-7_4/FIGURES/7
- [2] K. M. Alam and A. el Saddik, "C2PS: A digital twin architecture reference model for the cloud-based cyber-physical systems," *IEEE Access*, vol. 5, pp. 2050–2062, 2017, Doi: 10.1109/ACCESS.2017.2657006.
- [3] van Boeijen, A., Daalhuizen, J., & Zijlstra, Delft Design Guide: Perspectives, models, approaches, methods. (2nd ed.) BIS Publishers, 2020.
- [4] Kaner, Sam. Facilitator's guide to participatory decision-making. John Wiley & Sons, 2014.
- [5] KALLIO, Hanna, et al. Systematic methodological review: developing a framework for a qualitative semi-structured interview guide. *Journal of advanced nursing*, 2016, 72.12: 2954-2965.
- [6] G. Guest, K. M. MacQueen, and E. E. Namey, *Applied thematic analysis*. Sage Publications, 2011.
- [7] L. S. Nowell, J. M. Norris, D. E. White, and N. J. Moules, "Thematic analysis: Striving to meet the trustworthiness criteria," *International journal of qualitative methods*, vol. 16, no. 1, p. 1 609 406 917 733 847, 2017.
- [8] J. Fereday and E. Muir-Cochrane, "Demonstrating rigor using thematic analysis: A hybrid approach of inductive and deductive coding and theme development", *International journal of qualitative methods*, vol. 5, no. 1, pp. 80–92, 2006.
- [9] V. Braun and V. Clarke, "Using thematic analysis in psychology," *Qual. Res. Psychol.*, Jan. 2006, Accessed: Jun. 10, 2024. [Online]. Available: <https://www.tandfonline.com/doi/abs/10.1191/1478088706qp063oa>



VOXReality

Voice driven
interaction in XR spaces



**Funded by
the European Union**

Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or Directorate-General for Communications Networks, Content and Technology (DG CNECT). Neither the European Union nor the granting authority can be held responsible for them.