

Autonomous Workflow for Multimodal Fine-Grained Training Assistants Towards Mixed Reality

Anonymous ACL submission

Abstract

Autonomous artificial intelligence (AI) agents have emerged as promising protocols for automatically understanding the language-based environment, particularly with the exponential development of large language models (LLMs). However, a fine-grained, comprehensive understanding of multimodal environments remains under-explored. This work designs an autonomous workflow tailored for integrating AI agents seamlessly into extended reality (XR) applications for fine-grained training. We present a demonstration of a multimodal fine-grained training assistant for LEGO brick assembly in a pilot XR environment. Specifically, we design a cerebral language agent that integrates LLM with memory, planning, and interaction with XR tools and a vision-language agent, enabling agents to decide their actions based on past experiences. Furthermore, we introduce LEGO-MRTA, a multimodal fine-grained assembly dialogue dataset synthesized automatically in the workflow served by a commercial LLM. This dataset comprises multimodal instruction manuals, conversations, XR responses, and vision question answering. Last, we present several prevailing open-resource LLMs as benchmarks, assessing their performance with and without fine-tuning on the proposed dataset. We anticipate that the broader impact of this workflow will advance the development of smarter assistants for seamless user interaction in XR environments, fostering research in both AI and HCI communities.

1 Introduction

The advent of “Industry 4.0”, centered on the concept of smart manufacturing, presents a landscape with both opportunities and challenges for enhancing production efficiency (Goel and Gupta, 2020; Bécue et al., 2021; Jan et al., 2023). Training assistance for automating and accelerating industrial



(a) Industrial Car Assembly.



(b) LEGO Brick Assembly. We illustrate several use cases in the demo of BrickDream.¹

Figure 1: Examples of fine-grained assembly in XR systems.

assembly is in huge demand across various manufacturing applications, such as furniture manufacturing (You et al., 2022), industrial product assembly (Funk et al., 2017), and car assembly (Belalouna et al., 2020).

Mixed reality (MR), encompassing both virtual reality (VR) and augmented reality (AR), spans a spectrum from fully real environments to “matrix-like” virtual environments, showing promise for industrial manufacturing assembly tasks (Gavish et al., 2015; Stender et al., 2021; Butaslac et al., 2022). These multimodal, interactive, user-centric environments provide a solution for trainees who experience significant cognitive workload for training (Hou and Wang, 2013; Botto et al., 2020; Dalim et al., 2020). However, the assistance of a senior person as a trainer is required, either in person or remotely (Fidalgo et al., 2023).

To advance intelligent virtual assistants, tradi-

¹<https://www.youtube.com/watch?v=KkZKL3aKMJs>

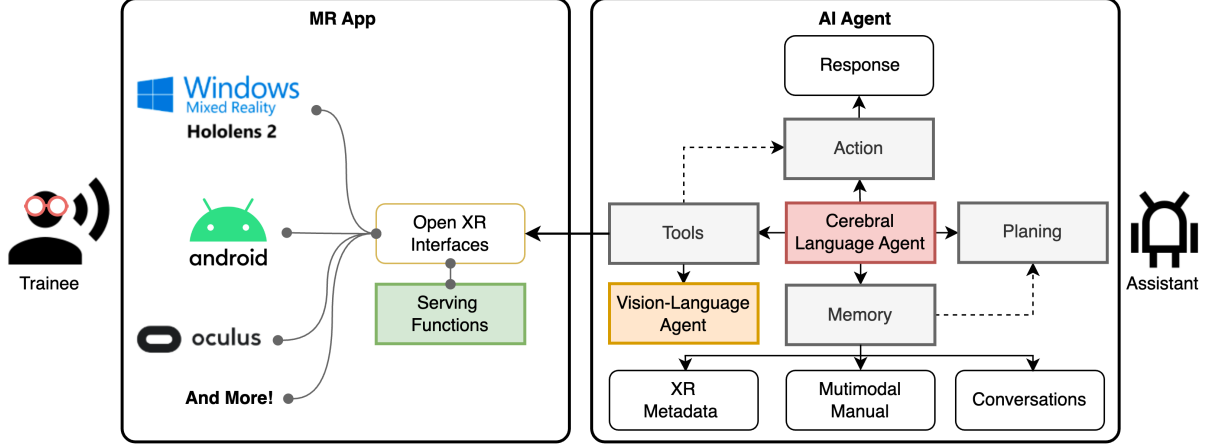


Figure 2: The proposed autonomous workflow, involving an AI agent interacting with a MR application. The AI agent comprises a core cerebral language agent, which interacts with a vision-language agent to interpret multimodal context into metadata, which can be utilized by the cerebral language agent iteratively. The MR application seamlessly interacts with AI agents by serving functions as external tools.

tional work leverages natural language processing (NLP) techniques (Li and Yang, 2021; Li et al., 2021, 2022; Colabianchi et al., 2023) and reinforcement learning (Sloan et al., 2022) to promote human-machine interactions. LLMs, as the new era of prevalent NLP techniques, have been observed to elicit diverse interaction patterns across tasks, demonstrating its versatility and feasibility (Mahmood et al., 2023). However, (i) tailoring assistant services by grounding interactions; and (ii) understanding users’ situated multimodal contexts remain challenging and under-explored (Dong et al., 2023).

To this end, we introduce an autonomous workflow (see Figure 2) tailored for seamlessly integrating AI agents into XR applications for fine-grained training. We present a demonstration of a multimodal fine-grained training assistant within a toy XR application for LEGO brick assembly. Specifically, we design a cerebral language agent that integrates LLM with memory, planning, and interaction with XR serving functional tools and a vision-language agent, enabling agents to decide their actions based on experiences. Then, we introduce LEGO-MRTA, a multimodal fine-grained assembly dataset synthesized automatically by a commercial LLM. This dataset comprises 65 multimodal instruction manuals, 1,423 conversations with vision question answering, serving usages of 18 functional tools in an XR environment. Additionally, several prevailing open-resource LLMs are presented as benchmarks, assessing their per-

formance with and without fine-tuning on the proposed dataset. Furthermore, we anticipate that the broader impact of this workflow will advance the development of smarter assistants for seamless user interaction in XR environments, fostering research in both AI and HCI communities.

We summarize our contributions as follows:

- We design a workflow, which seamlessly integrates autonomous AI agents for fine-grained assembly assistance in an XR demonstration.
- We create a multimodal manual-grounded fine-grained assembly conversation dataset in the XR context.
- We assess several open-resource LLMs as benchmarks, evaluating their performance with and without fine-tuning on the proposed dataset.

2 Related Work

For thoroughness, we provide preliminaries concerning multimodal datasets and virtual dialogue assistants within the realm of mixed reality (MR).

2.1 Multimodal Datasets towards MR

Traditional multimodal datasets focus on the interactions with sensor data (Patrik et al., 2018) between human-human or human-robot, and only a few of them provide small-scale task-oriented dialogues, such as OFAI-MMTD (Schreiter and Krenn, 2016) and (Kontogiorgos et al., 2018), and Chinese Whispers (Dimosthenis et al., 2020). ScanScribe (Zhu et al., 2023b) releases a 3D scene-text pairs dataset for 3D vision and text alignment learn-

	Domain	#Conv.	#Utt.	#Token	#AvgUtt.	#AvgToken
MDC	Minecraft building	509	15,926	113,116	30.7	7.9 (Architect) / 2.9 (Builder)
CerealBar	Instruction following	1,202	23,979	3,641	19.9	14.0 (Instructor) / 8.5 (Follower)
CVDN	Navigation	2,050	12,361	2,223	6.0	33.5 (Navigators) / 48.1 (Oracles)
TEACH	Household	3,215	45,000	3,429	13.7	5.7 (Commander) / 3.8 (Follower)
LEGO-ARTA	LEGO Assembly	1,423	35,131	7,173	24.8	26.6 (Trainer) / 12.7 (Trainee)

Table 1: Comparison of dialogue datasets towards MR.

ing. HoloAssistant (Wang et al., 2023c) provides a dataset containing 350 unique instructor-performer pairs with AR metadata to perceive, reason, and interact in the physical world. However, the conversations are not publicly available.

Recent studies have concentrated on multimodal datasets with conversations. MDC (Narayan-Chen et al., 2019) presents a collection of 509 human-human conversations in the Minecraft VR games. CerealBar (Suhr et al., 2019) creates 1,202 human-to-human conversations that map user instructions to system actions in a situated VR game environment. CVDN (Thomason et al., 2020) collects 2,050 human-robot conversations on Amazon Mechanical Turk for improving parsing and perception for natural language commands. Teach (Padmakumar et al., 2022) builds over 3,000 human-human, interactive dialogues to complete household tasks in the simulation.

Different from those aforementioned datasets, LEGO-MRTA gathers 1,423 human-human natural conversations between trainers and trainees. Unlike robotic commands, the length of utterances is relatively longer. Furthermore, these conversations are generated by grounding both on an instruction manual and responses from an XR, ensuring that the simulated conversations closely resemble natural human language. We compare the statistics of the above datasets in Table 1.

2.2 Virtual Dialogue Assistants for MR

Conventional efforts focus on creating virtual assistants for human-machine interactions using NLP techniques (Li and Yang, 2021; Li et al., 2021, 2022; Colabianchi et al., 2023) and reinforcement learning (Sloan et al., 2022). LLMs, representing the forefront of contemporary NLP techniques, hold tremendous promise for advancing towards the next generation of intelligent assistants (Naveed et al., 2023). The recent remarkable achievements of LLMs have spurred a growing interest in utilizing them to address a variety of complex tasks (Zhang et al., 2023), with particular attention being drawn to LLM-augmented autonomous

agents (Yao et al., 2022; Huang et al., 2022; Shinn et al., 2023; Madaan et al., 2023).

Autonomous agents expand the capabilities of LLMs into sequential action execution, demonstrating their proficiency in interacting with environments and addressing complex tasks through data collection (Wang et al., 2023b; Liu et al., 2023). A crucial aspect of this advancement relies on the capacity of LLMs to generate and interpret images, enabling them to access visual content and provide inputs, thereby integrating with mixed reality (MR) environments (Oyanagi et al., 2023; Wei et al., 2024). Regarding skill training, autonomous agents and LLMs can create immersive learning experiences that blend virtual and physical environments. For instance, students can utilize them to explore workflows and concepts in a more interactive and engaging manner (Gong et al., 2023; Li et al., 2024). In the context of MR serving as a sandbox (Li et al., 2023b) for LLMs and autonomous agents, the relationship is mutually beneficial. MR offers a secure (Naihin et al., 2023), adaptable, and regulated setting for training models. Together, LLMs, autonomous agents, and MR hold the potential to revolutionize our interaction with the digital world (Xu et al., 2023).

The convergence of LLMs, autonomous agents, and MR presents both excitement and challenges. As MR training experiences become more realistic and personalized, they demand greater amounts of data, encompassing detailed information about trainees’ behaviors, preferences, and interactions. Ensuring the availability and reusability of this data poses a significant challenge. Overall, our workflow aims to enhance MR training experiences by facilitating more natural language interactions, generating precise 3D models of real-world objects (Li et al., 2023a), and fostering dynamic and interactive experiences. While challenges remain (Xi et al., 2023; Ayache et al., 2023), the potential of this powerful technological fusion offers numerous exciting possibilities that could revolutionize personalization in virtual experiences. This entails the development of dedicated workflows and datasets.

3 Fine-grained Training Workflow

In this section, we describe the proposed workflow (See Figure 2) that advances AI agents towards XR guided fine-grained training.

3.1 Definition of Fine-Grained Training

In the context of fine-grained training, we anticipate the ability to (i) accurately follow professional training instructions documented in an instruction manual and; (ii) be sensitive to detailed visual information, ultimately for complex industrial assembly tasks, as illustrated in Figure 1 (a).

We define the following two roles during a training session:

- **User:** A human trainee who aims to acquire expertise and will work on fine-grained assembly tasks through interaction with the XR environment.
- **Assistant:** A virtual AI agent who will be able to assist the trainees in training and respond to their inquiries. It offers support with (i) a conversation agent that replies to trainees' requests and provides guidance grounded in the instruction manual; (ii) an interface for users to interact with XR environment; and (iii) a vision-language agent that understands and transmits users' visual context to language.

3.2 Autonomous AI Agent

We design the autonomous AI agent with a chain of two agents, namely (i) a cerebral language agent that serves to reply to trainees' requests, provide guidance, interact with XR and the vision-language agent; and (ii) a vision-language agent that understands and transmits users' visual context to language, which is then utilized by the cerebral language agent for planning.

3.2.1 Cerebral Language Agent

Inspired by the concept of LLM-powered autonomous agents (Wang et al., 2023b), we develop a cerebral language agent that incorporates an LLM with *memory*, *planning*, and functional *tools* that can interact with XR application, thereby enabling agents to make decisions regarding their *actions* based on past experiences. It can handle multimodal inputs, such as instruction manuals, historical conversations, and metadata within XR environments, and subsequently generate actions (i.e., responses or API calls for the XR application). The scope of responsibility of the agent is defined in

a system prompt (See P2, Table 5, Appendix A). Notably, it is able to alleviate the challenges (See §1): (i) it tailors assistant services by seamlessly interaction with XR applications to discover the business needs gradually; (ii) it interacts with a vision-language agent (See §3.2.2), which facilitates the capability of understanding the multimodal context in XR environments.

3.2.2 Vision-Language Agent

The vision-language agent's mission is to bridge the gap between understanding visual context and language, enabling effective utilization by the cerebral language agent (See §3.2.1) to conduct comprehensive planning for global optimization. Its core is the vision-language model (VLM) which is a task-driven large model that transmits vision input into language output needed by specific tasks.

In the context of LEGO assembly training, we observe two distinct patterns in LEGO instruction manuals (See an example in Figure 3) and define the following two tasks: **(T1) Object detection.** Given an image or a sequence of images as input, the objective is to predict the position of an object requested in a query and generate output in the format of "<Object> <Xleft> <Ytop> <Xright> <Ybottom>". For example, during assembly step 2, the AI trainer might direct the trainee, "Please gather the earth blue pair of legs and the silver metallic upper part of the body." In response, the trainee may ask, "Is this the one?" The vision-language agent is tasked with recognizing the object the trainee is referring to. **(T2) Assembly state detection.** Given an image or a sequence of images as input, the objective is to identify if the current assembly state matches the reference state provided in the instruction manual. For example, during assembly step 3, the vision-language agent is responsible for assisting the user's request, such as "Am I assembling them correctly?"

3.3 Pilot XR Application Design

We design an XR application as a pilot to show intuitive demonstration. First, we utilized a commercial LLM to generate candidate user requirements using the prompt (See P1, Table 5, Appendix A) as input. Then, we brainstormed and discussed the generated user requirements within a group of researchers and developers and finalized 7 user requirements (See Table 6, Appendix A) and 18 serving functional tools (See Table 2). We develop standard application programming interfaces (APIs) to en-

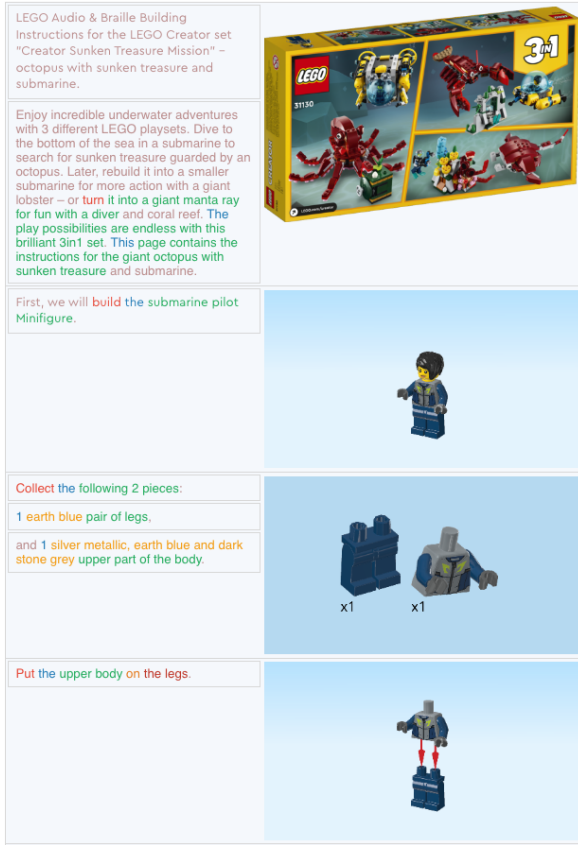


Figure 3: An example of LEGO instruction manual. It consists of a summary section at the beginning followed by three sequential instruction steps. Each step includes textual instructions paired with corresponding images to guide the assembly process.

able seamless interactions between functions in XR application and AI agent.

4 Dataset Creation

In this section, we introduce how to use the proposed workflow (§3) to create a multimodal dialogue dataset in XR environment.

4.1 Instruction Manual Crawling

We crawled 65 multimodal instruction manuals for fine-grained training from LEGO official website². A manual provides illustrated images and textural instructions on how to use, operate, assemble, and install a LEGO brick set. The key sections of an instruction manual include: (i) a summary that describes the general information, such as topics and candidate parts for assembly. It is followed by (ii) a sequence of multimodal step instructions. Each step contains a set of textual instructions and

²<https://legoaudioinstructions.com/instructions>

an illustration by image. Key functional phrases such as theme entities are highlighted in textual instructions. Here we show an example of a LEGO instruction manual in Figure 3.

4.2 Tool Response Generation

First, we use the crawled instruction manuals and the well-designed prompt template to produce prompts as an LLM input to generate user functional requirements and decide the serving functional tools. Then, we randomly choose up to 6 tools for each conversation session and record the simulated responses generated from templates.

4.3 VLM-based QA Construction

First, we use the step instruction to construct a query, containing a special token (“[detection]”) for the object detection task and a single instruction in a step. Second, we employ a query and the aligned image as inputs for MiniGPT-v2 (Chen et al., 2023; Zhu et al., 2023a), generating inference output as an answer of the query in the format of “<Object> <Xleft> <Ytop> <Xright> <Ybottom>”. Last, we iterate through all instruction steps in a conversation session, repeating the above two steps to construct vision question answering (VQA) pairs.

4.4 Multimodal Context-Aware Conversation Generation

We generate conversations grounded on both the instruction manual and tool responses using a commercial LLM. First, we reconstruct full instruction manuals with a summary and 10 step instructions because the average number of steps per manual is 215.3, which is quite long. This may limit the input tokens of an LLM and potentially distract the LLM with less grounding capability. Second, we instantiate the designed prompt template (See P3, Table 5, Appendix A) with the chunked instruction manuals. Last, we utilize a commercial LLM as the core of the proposed workflow to generate the conversations. Specifically, the system prompt informs the language agent about its responsibilities. The query prompt is used for each round of requests to generate a conversation. The historical rounds of requests are tracked by memory.

4.5 Dataset Statistics

We report the statistics of instruction manuals (See Table 3) and conversations (See Table 1).

We obtain 65 instruction manuals as grounding to lead a commercial LLM to generate 1,423

Tool Name	Description
StartAssemble	Initiate the assembly process.
NextStep	Move to the next assembly step.
FrontStep	Go back to the previous assembly step.
Explode	Trigger an explosion for detailed viewing.
Recover	Restore the initial state of AR objects after explosion.
FinishedVideo	End the assembly process and show a video of the assembled LEGO bricks.
ReShow	Repeat the current assembly step.
Enlarge	Enlarge or zoom out the current object.
Shrink	Shrink or zoom in the current object.
GoToStep	Go to the given assembly step number.
Rotate	Rotate the current object to a direction ("Up", "Down", "Left", "Right", "None").
ShowPieces	Show all candidate LEGO pieces to be assembled.
HighlightCorrectComponents	Highlight correct attachment points and components.
GetCurrentStep	Get the number of the current step.
GetRemainingStep	Get the number of the remaining steps.
CheckStepStatusVR	Check if the current step in Unity is accomplished correctly or not.
APICallObjectRecognitionAR	Call the VLM agent to identify LEGO pieces based on the provided video streaming data from AR glasses and highlight the recognized pieces in the AR environment.
APICallCheckStepStatusAR	Call the VLM agent to determine if the current assembly step is completed correctly or not, using the provided video streaming data from AR glasses as input.

Table 2: Descriptions of serving tools in the pilot XR application.

LEGO-ARTA Instruction Manual	
#Manual	65
#InstructionStep	13,994
#Token	8,676
#Theme entity	2,412
#AvgInstructionStep	215.3
#AvgConversation	28.3
Modalities	Text, Image

Table 3: Statistics of instruction manuals in the LEGO-MRTA dataset.

human-human natural conversations between trainers and trainees. Each instruction manual can make 28.3 on average. Theoretically, the amount of conversations can be enlarged by multiple times of requests. However, we focus on showcasing how to create meaningful datasets automatically. We construct 26,405 context-response pairs from generated conversations and VQA pairs as data samples. The average length is 107 tokens for the context and 145 tokens for the response utterance. We utilize 21.1k samples for fine-tuning open-resource LLMs to enhance the instruction-following capability and evaluate their performance on 5.25k test samples. Compared with existing datasets, LEGO-MRTA ensures that the simulated conversations closely resemble natural human language because of the design of the simulation method.

5 Experimental Setup

5.1 LLM Benchmarks

We consider several prevailing 7B open-source decoder-only LLMs as benchmarks, considering privacy concerns associated with fine-grained training in manufacturing.

- **BLOOM** (Le Scao et al., 2022) is pretrained on the multilingual ROOTS corpus, offering multilingual capabilities for various natural language processing tasks.
- **Falcon-instruct** (Almazrouei et al., 2023) is pre-trained on a large corpus of RefinedWeb data and fine-tuned on mixed chat and instruct datasets.
- **Llama2-Chat** (Touvron et al., 2023) is a pre-trained and fine-tuned generative text model optimized specifically for dialogue tasks, ensuring high-quality conversational responses.
- **Vicuna1.5** (Zheng et al., 2023) is a chat assistant derived by fine-tuning Llama 2 on user-shared conversations collected from ShareGPT.
- **OpenChat3.5** (Wang et al., 2023a) is a chat model fine-tuned with the C-RLFT strategy on mixed-quality data, achieving performance comparable to larger models like ChatGPT.
- **XVERSE**³ is a versatile model supporting 8k context length, ideal for longer multi-round dialogues, knowledge question-answering, and summarization tasks, trained on a diverse dataset of

³https://github.com/xverse-ai/XVERSE-7B/blob/main/README_EN.md

Model	BLEU-4		ROUGE-1		ROUGE-2		ROUGE-L		ToolACC (%)		ThemeACC (%)	
PEFT (LoRA)	/wo	/w	/wo	/w	/wo	/w	/wo	/w	/wo	/w	/wo	/w
BLOOM	2.88	54.07	20.49	61.91	6.50	49.52	3.78	58.63	49.62	77.86	26.30	64.61
Falcon	5.38	10.30	8.68	11.33	4.25	7.41	5.11	10.20	22.79	<u>17.65</u>	12.66	<u>10.81</u>
Llama2-Chat	10.23	30.53	18.59	40.65	7.41	25.48	10.59	32.91	21.37	55.73	47.20	55.51
Vicuna1.5	14.11	54.71	29.30	62.64	14.21	50.47	15.48	59.36	52.67	78.12	69.69	<u>66.79</u>
OpenChat3.5	22.00	<u>6.94</u>	29.70	34.51	15.69	23.69	22.50	11.36	51.97	74.02	58.19	81.90
XVERSE	22.42	53.55	28.45	61.54	14.31	49.77	22.39	58.03	49.62	83.97	57.53	71.10
BlueLM	22.72	55.69	30.40	63.52	14.98	51.58	23.76	60.35	48.15	82.22	47.51	68.08
Qwen	24.82	59.78	31.44	66.95	17.69	55.95	25.66	64.26	45.71	77.14	54.96	71.17
Mistral	25.87	54.17	33.32	62.07	17.99	49.40	26.32	58.62	49.62	78.20	54.80	66.65

Table 4: Benchmarking LLMs on LEGO-MRTA dataset, without (/wo) and with (/w) parameter-efficient fine-tuning (PEFT) using low rank adaption (LoRA). The bold font indicates the highest score in each column. The underline indicates the performance decrease after fine-tuning.

2.6 trillion tokens.

- **BlueLM-Chat**⁴ is a large-scale language model optimized for chat tasks, offering improved context understanding.
- **Qwen-Chat** (Bai et al., 2023) is a chat model that fine-tunes the pretrained Qwen model using human alignment techniques.
- **Mistral-Instruct** (Jiang et al., 2023) is a fine-tuned version of the Mistral-7B-v0.1, specifically tailored for instruction-based tasks using publicly available conversation datasets.

5.2 Evaluation Metrics

We evaluate the performance in terms of both overlap (BLUE-n, ROUGE-n) and informativeness (ToolACC, ThemeACC):

- **BLUE-n** measures precision, which measures the ratio of n-grams in the generated responses that match those in the reference responses. We consider $n = 4$.
- **ROUGE-n** measures recall, which calculates the ratio of n-grams in the reference responses that are captured by the generated responses. Here we consider $n = 1, 2, L$ and L denotes the number of longest common subsequences.
- **ToolACC** is defined as the ratio of correctly mentioned entities by the generated responses, compared to the reference response, from a list of serving tools.
- **ThemeACC** is defined as the ratio of correctly mentioned entities compared to the reference response, from a list of theme entities obtained from the instruction manual.

⁴<https://huggingface.co/vivo-ai/BlueLM-7B-Chat>

5.3 Implementation Details

The implementation of the workflow is based on Langchain.⁵ The “gpt-3.5-turbo-16k-0613” is used as the commercial LLM for generating data, e.g., conversations, user requirements, and serving functions. The MiniGPT4-v2⁶ is used as the VLM to detect the object, followed by simple rules to generate VQAs.

We use LoRA to conduct PEFT 7B open-source LLMs with the proposed dataset based on framework (Hyouga, 2023). Specifically, the maximum sequence length is 1024 and the learning rate is $5e-05$. The model is trained for 3 epochs with a per-device batch size of 4, and accumulated gradients every 4 steps. A cosine learning rate scheduler is employed, with a maximum gradient norm of 1.0. We log results every 5 steps and save model checkpoints every 100 steps. Warm-up steps are set to 0. LoRA is used with a rank of 8 and a dropout rate of 0.1 for regularization. All experiments are run on NVIDIA A100 SXM4 40GB GPUs.

6 Outcomes

6.1 Evaluation on Benchmark LLMs

Table 4 shows the performance on 9 prevailing open-source LLMs, without and with fine-tuning on the LEGO-MRTA dataset.

First, after fine-tuning, the performance of all models dramatically improved in terms of all metrics, except for the results that are underlined. This demonstrates the feasibility and effectiveness of tailoring LLM for fine-grained training in XR environments. In addition, this shows the proposed

⁵<https://python.langchain.com/docs/get-started/introduction>

⁶<https://github.com/Vision-CAIR/MiniGPT-4>

Limitations

The generation of user requirements and the dataset relies solely on the simulation process. This workflow serves as a fast solution to verify the concept of a LLM agent aiding in a specific use case, such as a LEGO assembly assistant. However, we acknowledge that the study user requirements are valuable and needed to build up user-centric AI agents and XR applications. Besides, the demonstration codes do not optimize LLM and VLM simultaneously, potentially leading to suboptimal outcomes. We have only assessed LLMs as benchmarks. However, we have not conducted separate assessments of the influence on the vision-language agent and user experience in MR. We plan to explore these aspects in future work.

Ethics Statement

We realize that there are risks in developing a large language model for users, so it is necessary to pay attention to the ethical issues. Therefore, we use the open-resourced LLMs as benchmarks and consider user-centric points: A user will first be provided an explanation of what will be happening during their XR training experience. Users will then be provided with relevant consent forms to sign, and after signing they will be fitted with the HoloLens 2 and the training scenario will begin. After launching the application, 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.

References

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.

Julia Ayache, Marta Bienkiewicz, Kathleen Richardson, and Benoit Bardy. 2023. extended reality of socio-motor interactions: Current trends and ethical considerations for mixed reality environments design. In *ICMI*, pages 154–158.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Adrien Bécue, Isabel Praça, and João Gama. 2021. Artificial intelligence, cyber-threats and industry 4.0: Challenges and opportunities. *Artificial Intelligence Review*, 54(5):3849–3886.

Fahmi Bellalouna, Mika Luimula, Panagiotis Markopoulos, Evangelos Markopoulos, and Franco Zipperling. 2020. Fiaar: an augmented reality firetruck equipment assembly and configuration assistant technology. In *CogInfoCom*, pages 000237–000244. IEEE.

Carola Botto, Alberto Cannavò, Daniele Cappuccio, Gida Morat, Amir Nematollahi Sarvestani, Paolo Ricci, Valentina Demarchi, and Alessandra Saturnino. 2020. Augmented reality for the manufacturing industry: the case of an assembly assistant. In *VRW*, pages 299–304. IEEE.

Isidro M Butaslac, Yuichiro Fujimoto, Taishi Sawabe, Masayuki Kanbara, and Hirokazu Kato. 2022. Systematic review of augmented reality training systems. *IEEE Transactions on Visualization and Computer Graphics*.

Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechu Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*.

Silvia Colabianchi, Andrea Tedeschi, and Francesco Costantino. 2023. Human-technology integration with industrial conversational agents: A conceptual architecture and a taxonomy for manufacturing. *Journal of Industrial Information Integration*, 35:100510.

Che Samihah Che Dalim, Mohd Shahrizal Sunar, Arindam Dey, and Mark Billingham. 2020. Using augmented reality with speech input for non-native children’s language learning. *International Journal of Human-Computer Studies*, 134:44–64.

Kontogiorgos Dimosthenis, Sibirtseva Elena, and Gustafson Joakim. 2020. Chinese Whispers: A Multimodal Dataset for Embodied Language Grounding. In *LREC*.

Xin Luna Dong, Seungwhan Moon, Yifan Ethan Xu, Kshitiz Malik, and Zhou Yu. 2023. Towards next-generation intelligent assistants leveraging llm techniques. In *SIGKDD*, pages 5792–5793.

Catarina G Fidalgo, Yukang Yan, Hyunsung Cho, Maurício Sousa, David Lindlbauer, and Joaquim Jorge. 2023. A survey on remote assistance and training in mixed reality environments. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2291–2303.

Markus Funk, Andreas Bächler, Liane Bächler, Thomas Kosch, Thomas Heidenreich, and Albrecht Schmidt. 2017. Working with augmented reality? a long-term analysis of in-situ instructions at the assembly workplace. In *PETRA*, pages 222–229.

651	Nirit Gavish, Teresa Gutiérrez, Sabine Webel, Jorge	Chen Li, Andreas Kornmaaler Hansen, Dimitrios	705
652	Rodríguez, Matteo Peveri, Uli Bockholt, and Franco	Chrysostomou, Simon Bøgh, and Ole Madsen. 2022.	706
653	Tecchia. 2015. Evaluating virtual reality and aug-	Bringing a natural language-enabled virtual assistant	707
654	mented reality training for industrial maintenance	to industrial mobile robots for learning, training and	708
655	and assembly tasks. <i>Interactive Learning Environ-</i>	assistance of manufacturing tasks. In <i>SII</i> , pages 238–	709
656	<i>ments</i> , 23(6):778–798.	243. IEEE.	710
657	Ruchi Goel and Pooja Gupta. 2020. Robotics and indus-	Chen Li, Jinha Park, Hahyeon Kim, and Dimitrios	711
658	try 4.0. <i>A Roadmap to Industry 4.0: Smart Produc-</i>	Chrysostomou. 2021. How can i help you? an intel-	712
659	<i>tion, Sharp Business and Sustainable Development</i> ,	ligent virtual assistant for industrial robots. In <i>HRI</i> ,	713
660	pages 157–169.	pages 220–224.	714
661	Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane	Chen Li and Hong Ji Yang. 2021. Bot-x: An ai-based	715
662	Durante, Yusuke Noda, Zilong Zheng, Song-Chun	virtual assistant for intelligent manufacturing. <i>Multi-</i>	716
663	Zhu, Demetri Terzopoulos, Li Fei-Fei, et al. 2023.	<i>agent and Grid Systems</i> , 17(1):1–14.	717
664	Mindagent: Emergent gaming interaction. <i>arXiv</i>	Mingsheng Li, Xin Chen, Chi Zhang, Sijin Chen,	718
665	<i>preprint arXiv:2309.09971</i> .	Hongyuan Zhu, Fukun Yin, Gang Yu, and Tao	719
666	Hiyouga. 2023. Llama factory. https://github.com/	Chen. 2023a. M3dbench: Let’s instruct large mod-	720
667	hiyouga/LLaMA-Factory .	els with multi-modal 3d prompts. <i>arXiv preprint</i>	721
668	Lei Hou and Xiangyu Wang. 2013. A study on the	<i>arXiv:2312.10763</i> .	722
669	benefits of augmented reality in retaining working	Yuan Li, Yixuan Zhang, and Lichao Sun. 2023b. Metaa-	723
670	memory in assembly tasks: A focus on differences in	gents: Simulating interactions of human behav-	724
671	gender. <i>Automation in Construction</i> , 32:38–45.	iors for llm-based task-oriented coordination via	725
672	Wenlong Huang, Pieter Abbeel, Deepak Pathak, and	collaborative generative agents. <i>arXiv preprint</i>	726
673	Igor Mordatch. 2022. Language models as zero-shot	<i>arXiv:2310.06500</i> .	727
674	planners: Extracting actionable knowledge for em-	Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue,	728
675	bodied agents. In <i>ICML</i> , pages 9118–9147. PMLR.	Shelby Heinecke, Rithesh Murthy, Yihao Feng,	729
676	Zohaib Jan, Farhad Ahamed, Wolfgang Mayer, Niki	Zeyuan Chen, Juan Carlos Nieves, Devansh Arpit,	730
677	Patel, Georg Grossmann, Markus Stumptner, and	et al. 2023. Bolaa: Benchmarking and orchestrating	731
678	Ana Kuusk. 2023. Artificial intelligence for industry	llm-augmented autonomous agents. <i>arXiv preprint</i>	732
679	4.0: Systematic review of applications, challenges,	<i>arXiv:2308.05960</i> .	733
680	and opportunities. <i>Expert Systems with Applications</i> ,	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	734
681	216:119456.	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	735
682	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	736
683	sch, Chris Bamford, Devendra Singh Chaplot, Diego	et al. 2023. Self-refine: Iterative refinement with	737
684	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	738
685	laume Lample, Lucile Saulnier, et al. 2023. Mistral	Amama Mahmood, Junxiang Wang, Bingsheng Yao,	739
686	7b. <i>arXiv preprint arXiv:2310.06825</i> .	Dakuo Wang, and Chien-Ming Huang. 2023. Llm-	740
687	Dimosthenis Kontogiorgos, Vanya Avramova, Simon	powered conversational voice assistants: Interaction	741
688	Alexanderson, Patrik Jonell, Catharine Oertel, Jonas	patterns, opportunities, challenges, and design guide-	742
689	Beskow, Gabriel Skantze, and Joakim Gustafson.	lines. <i>arXiv preprint arXiv:2309.13879</i> .	743
690	2018. A multimodal corpus for mutual gaze and	Silen Naihin, David Atkinson, Marc Green, Mer-	744
691	joint attention in multiparty situated interaction. In	wane Hamadi, Craig Swift, Douglas Schonholtz,	745
692	<i>LREC</i> .	Adam Tauman Kalai, and David Bau. 2023. Test-	746
693	Teven Le Scao, Angela Fan, Christopher Akiki, El-	ing language model agents safely in the wild. <i>arXiv</i>	747
694	lie Pavlick, Suzana Ilic, Daniel Hesslow, Roman	<i>preprint arXiv:2311.10538</i> .	748
695	Castagné, Alexandra Sasha Luccioni, François Yvon,	Anjali Narayan-Chen, Prashant Jayannavar, and Ju-	749
696	Matthias Gallé, et al. 2022. Bloom: A 176b-	lia Hockenmaier. 2019. Collaborative dialogue in	750
697	parameter open-access multilingual language model.	minecraft. In <i>ACL</i> , pages 5405–5415.	751
698	<i>arXiv preprint arXiv:2211.05100</i> .	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muham-	752
699	Bai Li, Xinyuan Li, Yaodong Cui, Xuepeng Bian, Siyu	mad Saqib, Saeed Anwar, Muhammad Usman, Nick	753
700	Teng, Siji Ma, Lili Fan, Yonglin Tian, Fei-Yue Wang,	Barnes, and Ajmal Mian. 2023. A comprehensive	754
701	et al. 2024. Integrating large language models and	overview of large language models. <i>arXiv preprint</i>	755
702	metaverse in autonomous racing: An education-	<i>arXiv:2307.06435</i> .	756
703	oriented perspective. <i>IEEE Transactions on Intel-</i>	Akimi Oyanagi, Kazuma Aoyama, Kenichiro Ito, Tomo-	757
704	<i>ligent Vehicles</i> .	hiro Amemiya, and Michitaka Hirose. 2023. Virtual	758
		reality training system using an autonomy agent for	759

760	learning hospitality skills of a retail store. In <i>HCI</i> , pages 483–492. Springer.	815
761		816
762	Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen, Span-	817
763	dana Gella, Robinson Piramuthu, Gokhan Tur, and	818
764	Dilek Hakkani-Tur. 2022. TEACH: Task-driven Em-	819
765	bodied Agents that Chat. In <i>AAAI</i> , volume 36, pages	820
766	2017–2025.	
767		
768	Jonell Patrik, Bystedt Mattias, Fallgren Per, Kontogior-	821
769	gos Dimosthenis, Lopes José, Malisz Zofia, Mas-	822
770	carenhas Samuel, Oertel Catharine, Eran Raveh, and	823
771	Shore Todd. 2018. ARMI: An Architecture for	824
772	Recording Multimodal Interactions. In <i>LREC</i> .	825
773		
774	Stephanie Schreitter and Brigitte Krenn. 2016. The ofai	
775	multi-modal task description corpus. In <i>LREC</i> , pages	
	1408–1414.	
776	Noah Shinn, Federico Cassano, Ashwin Gopinath,	
777	Karthik R Narasimhan, and Shunyu Yao. 2023. Re-	
778	flexion: Language agents with verbal reinforcement	
779	learning. In <i>NeuralIPS</i> .	
780	Hannah Sloan, Richard Zhao, Faisal Aqlan, Hui Yang,	
781	and Rui Zhu. 2022. Adaptive virtual assistant for vir-	
782	tual reality-based remote learning. In <i>ASEE Annual</i>	
783	<i>Conference & Exposition</i> .	
784		
785	Birga Stender, Johannes Paehr, and Thomas N Jambor.	
786	2021. Using ar/vr for technical subjects in vocational	
787	training—of substantial benefit or just another techni-	
	cal gimmick? In <i>EDUCON</i> , pages 557–561. IEEE.	
788		
789	Alane Suhr, Claudia Yan, Jack Schluger, Stanley Yu,	
790	Hadi Khader, Marwa Mouallem, Iris Zhang, and	
791	Yoav Artzi. 2019. Executing instructions in situ-	
792	ated collaborative interactions. In <i>EMNLP-IJCNLP</i> ,	
	pages 2119–2130.	
793	Jesse Thomason, Aishwarya Padmakumar, Jivko	
794	Sinapov, Nick Walker, Yuqian Jiang, Harel Yedid-	
795	sion, Justin Hart, Peter Stone, and Raymond Mooney.	
796	2020. Jointly improving parsing and perception for	
797	natural language commands through human-robot	
798	dialog. <i>Journal of Artificial Intelligence Research</i> ,	
799	67:327–374.	
800	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	
801	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	
802	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	
803	Bhosale, et al. 2023. Llama 2: Open founda-	
804	tion and fine-tuned chat models. <i>arXiv preprint</i>	
805	<i>arXiv:2307.09288</i> .	
806	Guan Wang, Sijie Cheng, Xianyu Zhan, Xiangang Li,	
807	Sen Song, and Yang Liu. 2023a. Openchat: Advanc-	
808	ing open-source language models with mixed-quality	
809	data. <i>arXiv preprint arXiv:2309.11235</i> .	
810	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	
811	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	
812	Xu Chen, Yankai Lin, et al. 2023b. A survey on large	
813	language model based autonomous agents. <i>arXiv</i>	
814	<i>preprint arXiv:2308.11432</i> .	
	Xin Wang, Taein Kwon, Mahdi Rad, Bowen Pan, Is-	815
	hani Chakraborty, Sean Andrist, Dan Bohus, Ashley	816
	Feniello, Bugra Tekin, Felipe Vieira Frujeri, et al.	817
	2023c. Holoassist: an egocentric human interaction	818
	dataset for interactive ai assistants in the real world.	819
	In <i>ICCV</i> , pages 20270–20281.	820
	Yuxi Wei, Zi Wang, Yifan Lu, Chenxin Xu, Changxing	821
	Liu, Hao Zhao, Siheng Chen, and Yanfeng Wang.	822
	2024. Editable scene simulation for autonomous	823
	driving via collaborative llm-agents. <i>arXiv preprint</i>	824
	<i>arXiv:2402.05746</i> .	825
	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	826
	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	827
	Senjie Jin, Enyu Zhou, et al. 2023. The rise and	828
	potential of large language model based agents: A	829
	survey. <i>arXiv preprint arXiv:2309.07864</i> .	830
	Fengli Xu, Jun Zhang, Chen Gao, Jie Feng, and Yong	831
	Li. 2023. Urban generative intelligence (ugi): A	832
	foundational platform for agents in embodied city	833
	environment. <i>arXiv preprint arXiv:2312.11813</i> .	834
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	835
	Shafraan, Karthik Narasimhan, and Yuan Cao. 2022.	836
	React: Synergizing reasoning and acting in language	837
	models. <i>arXiv preprint arXiv:2210.03629</i> .	838
	Yingchao You, Ze Ji, Xintong Yang, and Ying Liu. 2022.	839
	From human-human collaboration to human-robot	840
	collaboration: automated generation of assembly task	841
	knowledge model. In <i>ICAC</i> , pages 1–6. IEEE.	842
	Jianguo Zhang, Kun Qian, Zhiwei Liu, Shelby Hei-	843
	necke, Rui Meng, Ye Liu, Zhou Yu, Silvio Savarese,	844
	and Caiming Xiong. 2023. Dialogstudio: Towards	845
	richest and most diverse unified dataset collection for	846
	conversational ai. <i>arXiv preprint arXiv:2307.10172</i> .	847
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	848
	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	849
	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.	850
	Judging llm-as-a-judge with mt-bench and chatbot	851
	arena. <i>arXiv preprint arXiv:2306.05685</i> .	852
	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and	853
	Mohamed Elhoseiny. 2023a. Minigpt-4: Enhancing	854
	vision-language understanding with advanced large	855
	language models. <i>arXiv preprint arXiv:2304.10592</i> .	856
	Ziyu Zhu, Xiaojuan Ma, Yixin Chen, Zhidong Deng,	857
	Siyuan Huang, and Qing Li. 2023b. 3d-vista: Pre-	858
	trained transformer for 3d vision and text alignment.	859
	In <i>ICCV</i> .	860
	A Appendix	861
	A.1 Prompt templates	862
	A.2 Generated user requirements	863
	A.3 Qualitative analysis of the dataset	864
	The collective results as seen in the figure (6) il-	865
	lustrates that our simulated model has practical	866
	assistance capabilities as summarized below:	867

(P1) Prompt template for user requirement generation

[Task description]

You are an AI agent who acts as a Unity developer for AR applications. Your role is to analyze users' functional needs based on the manuals and then develop the corresponding functions in an AR training system. Note that is not for visually impaired users, but for trainees who are visually healthy and able to wear HoloLen2 AR glasses.

Here are samples of manuals:

[Manuals]

(P2) Prompt template for conversation generation

1. System prompt

[Task description]

Brief version: The task is to generate multiple turns of conversations and called tools between the trainer (assistant) and trainee (user) grounded on the task-specific guidelines and tools in LEGO XR application.

Full version: The trainer aims to teach the trainee how to accomplish the assembly task based on the task-specific guidelines, supported by an XR application. Specifically, the trainee is wearing AR glasses to see both VR environment and real world. The trainee knows nothing about the guidelines before trainer's guidance. For each step, the trainee must ask at least one deep-dive question, or request a troublesome issue if he or she cannot follow the guide, or call tools from XR application and learn how to use those tools; the trainer must answer the question, assist the trainee, show them the responses to the execution of the tools. At the end of a conversation, first, the trainer must ask if the trainee has accomplished the task and the trainee must tell if the trainee can accomplish the task; second, the trainer must ask how is user experiences, and the trainee provide feedback on the user experience. You must add a section title to separate which key point in the guideline in the generated conversation and generate until the final step of the guidelines.

[Tool description as shown in Table 2]

2. Query prompt

[Task description (Brief version)]

[Summary and step instructions in a manual]

Imagine some trainee's utterances have the intent of using the tools with the following responses:

[Tool responses]

LEGO Assembly Assistant prompt (P3)

You are a helpful AI assistant who aims to train the user how to assemble a LEGO car in XR immersive system.

Extended Reality (XR) directs to the assortment of Virtual Reality (VR), Augmented Reality (AR), and Mixed Reality (MR).

Please make sure you complete the objective above with the following rules:

- (1) The user is a trainee who is wearing HoloLen 2 glasses and is able to see XR environments in real-time.
- (2) You are able to call Unity functions in the LEGO AR application.
- (3) You are able to obtain HoloLens 2 Sensor Streaming data.
- (4) Alert if the user asks you something outside of the LEGO assembly task but do not give overconfident answers.

Your task is to answer the user's questions and assist the user in understanding how to complete the LEGO assembly task in XR.

Table 5: Prompt templates used in this work.

User requirement	Description
3D Model Interaction	Create 3D models of the LEGO pieces and the Monster Truck assembly. Trainees can interact with these 3D models using hand gestures and voice commands, making it easier to understand the assembly process.
Step-by-Step Guidance	Display step-by-step instructions directly in the trainees' field of view. This can include both visual instructions and written or spoken guidance.
Real-Time Feedback	Provide real-time feedback to trainees as they assemble the LEGO set. Use AR to highlight the correct attachment points and components, and indicate when they've completed a step correctly.
Object Recognition	Implement object recognition so that HoloLens 2 can identify LEGO pieces and highlight them when trainees look at them. This can help trainees quickly find the right pieces.
Progress Tracking	Keep track of trainees' progress and provide them with an overview of the steps they have completed and those remaining. This can help them stay organized and motivated.
Troubleshooting Assistance	Include a troubleshooting mode that guides trainees through common problems and solutions they might encounter during the assembly.
Data Logging	Collect data on trainees' performance and interaction with the AR training system to analyze their progress and make improvements to the training process.

Table 6: User requirements of the XR training system.

Realistic simulation. LEGO is a well known block building concept. The dataset simulates various real-world scenarios encountered during LEGO assembly tasks. By replicating factors such as piece variability, environmental conditions, and assembly constraints, the dataset provides a realistic training

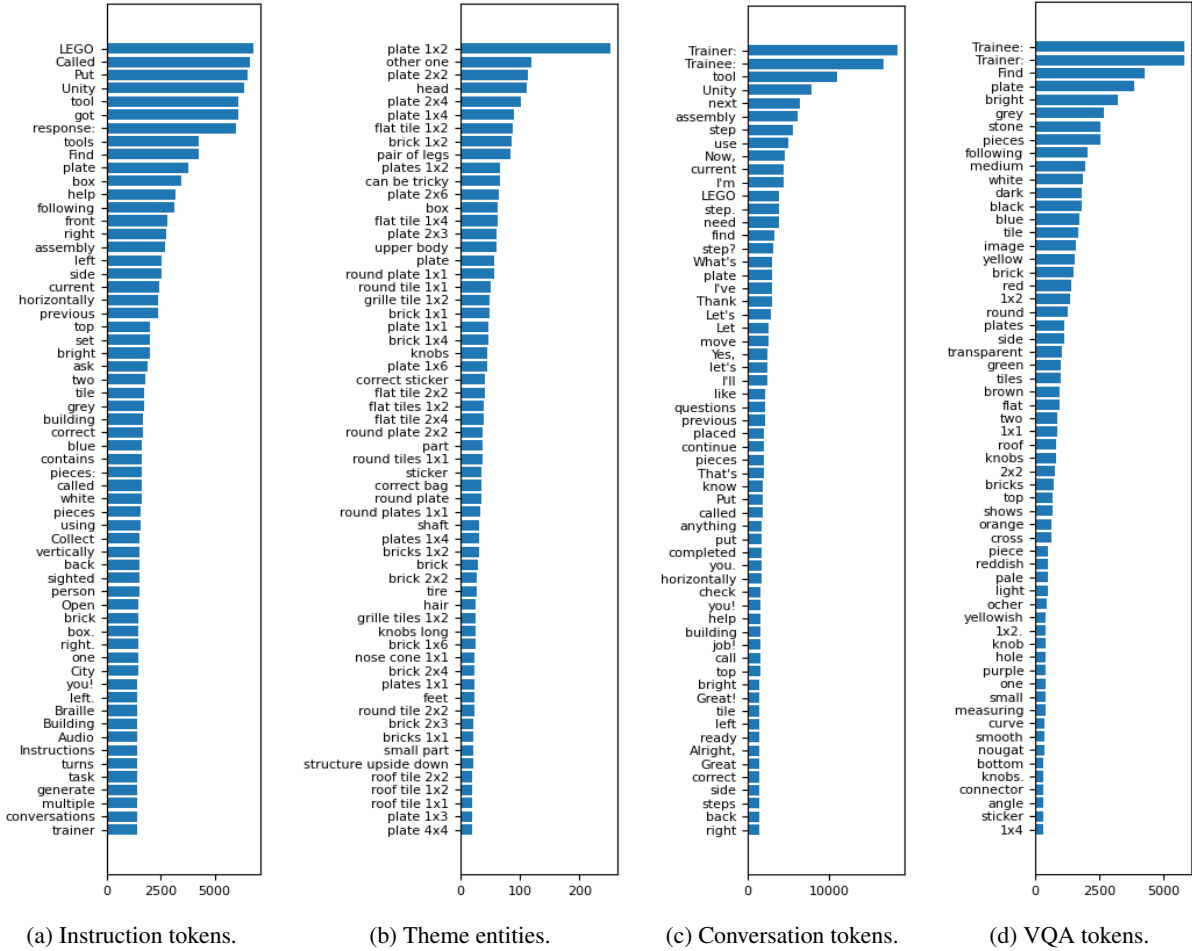


Figure 6: Distribution of top 60 frequent tokens in the above four parts: (a) instructions, (b) entities in the manual, (c) conversations, and (b) VQAs. The x-axis denotes frequency and the y-axis denotes tokens in four parts of the LEGO-MRTA dataset.

environment for machine learning models. This realism enhances the model’s ability to generalize to unseen situations, ensuring reliable performance in diverse assembly settings.

Diversity in task difficulty. From simple structures to intricate designs, the dataset exposes the model to diverse assembly scenarios, enabling it to learn robust representations of LEGO building principles. This diversity fosters adaptability in the model, empowering it to tackle simple to difficult or probably novel assembly tasks with confidence and efficiency.

Transfer learning to other tasks. The dataset is structured to facilitate transfer learning, allowing knowledge and representations learned from one assembly task to be applied to related tasks or domains. By leveraging pre-trained models or features learned from similar assembly tasks, machine learning models can bootstrap their learning

process on new assembly tasks. This transfer learning capability accelerates model adaptation to new environments and tasks, reducing the need for extensive retraining and improving overall training efficiency.

A.3.1 Instruction tokens (Figure 6a)

Analyzing the provided instruction tokens, we can derive several factors that contribute to the usability and effectiveness of our simulated dataset for XR training and assembly training:

- **Clear instructional guidance.** Tokens like “Put,” “Find,” “Collect,” and “Open” provide clear and concise instructions for performing various assembly tasks. These instructions guide users through the assembly process step-by-step, ensuring clarity and direction in the training environment.
- **Spatial orientation and manipulation.** Tokens such as “front,” “right,” “left,” “horizontally,” and

“vertically” offer spatial orientation cues, helping users understand the spatial relationships between LEGO components and how to manipulate them during assembly. This spatial awareness enhances users’ ability to accurately position and align LEGO pieces.

- **Feedback and assistance.** Tokens like “help” and “response” indicate provisions for feedback and assistance within the training environment. Offering assistance and feedback helps users troubleshoot issues, learn from mistakes, and improve their assembly skills over time, enhancing the learning experience.
- **Multimodal learning:** The inclusion of tokens like “Audio Instructions” suggests the incorporation of multimodal learning techniques within the training environment. Integrating audio instructions alongside visual cues enhances usability by catering to different learning styles and preferences, making the training experience more accessible and engaging for users.
- **Adaptive learning.** Tokens such as “current” and “previous” imply a dynamic learning environment where users can track their progress and revisit previous steps if needed. Adaptive learning features enhance usability by allowing users to learn at their own pace, review concepts as needed, and progress through the training material in a structured manner.
- **Interactive learning environment.:** The presence of tokens like “conversations” and “trainer” indicates an interactive learning environment where users can engage in dialogue and receive guidance from trainers or virtual assistants. Interactivity enhances usability by promoting engagement, collaboration, and active participation in the learning process, leading to more effective skill acquisition and retention.

The instruction tokens in our simulated dataset indicate clear guidance, and spatial orientation cues within an interactive learning environment.

A.3.2 Theme entities (Figure 6b)

Based on the theme entities provided, we can analyze the relevance of these tokens to the learning process:

- **Part identification.** Tokens such as “plate 1x2,” “plate 2x4,” and “brick 1x1” provide specific identifiers for different LEGO parts commonly used in assembly tasks. By including a variety of part identifiers, the dataset facilitates part recognition and identification, enabling the model to learn the

characteristics and properties of each component.

- **Spatial orientation and configuration.** Tokens like “head,” “upper body,” and “feet” suggest the inclusion of assembly instructions related to spatial orientation and configuration of LEGO structures. Understanding the spatial arrangement of components is essential for accurate assembly, and these tokens help the model grasp the hierarchical structure of assemblies and the placement of parts within them.
- **Assembly techniques.** Tokens such as “can be tricky” and “structure upside down” hint at the inclusion of assembly techniques and strategies within the dataset. Learning various assembly techniques is crucial for efficiently building complex structures, and these tokens provide guidance on overcoming challenges and optimizing assembly processes.
- **Component variations.** Tokens like “round plate 1x1,” “flat tile 2x4,” and “grille tile 1x2” introduce variations of standard LEGO components, reflecting the diversity of parts encountered in real-world assembly scenarios. By including a range of component variations, the dataset prepares the model to handle different types of parts and adapt to varying assembly requirements.
- **Accessory identification.** Tokens such as “pair of legs,” “tire,” and “hair” denote accessory pieces commonly used in LEGO constructions, adding realism and complexity to assembly tasks. Recognizing and incorporating accessory pieces is essential for creating realistic and detailed models, and these tokens help the model understand the role of accessories in assembly.
- **Quality control and correctness.** Tokens like “correct sticker” and “correct bag” emphasize the importance of quality control and correctness in assembly tasks. Ensuring that the correct parts are used in the right context is essential for achieving accurate and high-quality assemblies, and these tokens highlight the need for attention to detail and accuracy in the assembly process.
- **Structural components.** Tokens such as “shaft,” “structure upside down,” and “roof tile” suggest the inclusion of structural components and building techniques within the dataset. Understanding the role of structural components and mastering advanced building techniques is critical for creating stable and aesthetically pleasing assemblies, and these tokens provide guidance on construct-

ing sturdy and well-balanced structures.

The theme entities included in our simulated dataset provide a realistic representation of the assembly tasks by encompassing part identification, spatial orientation, assembly techniques, component variations, accessory recognition, quality control, and structural components, the dataset equips the model with the knowledge and skills necessary to effectively assemble LEGO structures in virtual environments.

A.3.3 Conversation Tokens (Figure 6c)

We can infer several aspects that contribute to the usability and effectiveness of our simulated dataset for conversations during training and assembly training:

- **Role identification.** The presence of “Trainer” and “Trainee” tokens indicates a clear distinction between the roles of the instructor guiding the training session and the learner receiving instructions. This role identification fosters clarity and structure in the conversation, ensuring effective communication between the trainer and trainee entities.
- **Instructional guidance.** Tokens such as “step,” “plate,” “use,” and “find” suggest the provision of instructional guidance within the conversation. The trainer entity likely provides step-by-step instructions and prompts to the trainee, guiding them through the assembly process and facilitating learning in a structured manner.
- **Interactive dialogue.** The conversation tokens include interactive dialogue cues such as “Let’s,” “Yes, let’s,” and “Thank you!” These cues foster engagement and collaboration between the trainer and trainee entities, creating a supportive and interactive learning atmosphere conducive to effective learning and skill development.
- **Feedback and encouragement.** Tokens like “Great!” and “Alright” suggest the inclusion of positive feedback and encouragement within the conversation. Positive reinforcement enhances motivation and engagement, encouraging active participation and fostering a positive learning experience for the trainee.
- **Error handling and assistance.** The presence of tokens like “check,” “help,” and “completed” indicates provisions for error handling and assistance within the conversation. The trainer entity likely offers guidance and support to the trainee in identifying and correcting errors, ensuring a constructive learning process and facilitating skill

development.

- **Spatial orientation and task management.** Tokens such as “right,” “left,” “back,” and “steps” provide spatial orientation cues and references to assembly tasks. This spatial orientation facilitates effective communication of assembly instructions and task management between the trainer and trainee entities, ensuring accurate placement and alignment of LEGO components during assembly.

The conversation tokens provide instructional guidance, facilitate interactive dialogue, offer feedback and encouragement, handle errors, and provide spatial orientation cues for task management.

A.3.4 VQA Tokens (Figure 6d)

Analyzing the provided tokens from the vision language model, we can identify several factors contributing to its usability and effectiveness:

- **Object recognition.** Tokens such as “plate,” “brick,” “tile,” and “knob” represent common LEGO elements that users encounter during assembly tasks. By including these tokens, the dataset enables the vision language model to recognize and identify various LEGO components accurately, facilitating object recognition and understanding in XR training environments.
- **Color detection.** Tokens like “bright,” “grey,” “white,” and “blue” provide color descriptors for different LEGO pieces. Incorporating color information allows the vision language model to detect and differentiate between LEGO components based on their color, enhancing the model’s ability to interpret and analyze assembly scenes accurately.
- **Shape recognition.** Tokens such as “round,” “flat,” “roof,” and “connector” describe the shapes and configurations of LEGO elements. By including shape descriptors, the dataset enables the vision language model to recognize and classify different types of LEGO pieces based on their shapes, facilitating shape recognition and classification in XR training environments.
- **Size specification.** Tokens like “1x2,” “2x2,” and “1x1” specify the sizes and dimensions of LEGO elements. Incorporating size information allows the vision language model to understand the scale and proportions of LEGO components within assembly scenes, aiding in size estimation and spatial reasoning during XR training tasks.
- **Material and texture.** Tokens such as “smooth,” “nougat,” and “transparent” describe the materi-

als and textures of LEGO elements. Including material and texture descriptors enables the vision language model to identify and distinguish between different surface finishes and textures, enhancing its ability to recognize and characterize LEGO components accurately.

- **Part relationships.** Tokens like “side,” “top,” and “bottom” provide spatial relationship cues between LEGO elements. By including part relationship descriptors, the dataset enables the vision language model to understand the spatial arrangement and orientation of LEGO components within assembly scenes, facilitating the interpretation of complex assembly structures and configurations.
- **Visual context understanding.** Tokens such as “image” and “shows” suggest the inclusion of visual context information within the dataset. Providing visual context cues enables the vision language model to interpret and analyze assembly scenes holistically, incorporating visual information to enhance its understanding of the surrounding environment and improve object recognition accuracy.

Our simulated dataset successfully provides object recognition, color detection, shape recognition, size specification, material and texture characterization, part relationships, and visual context understanding. Altogether, these tokens contribute to the usability and effectiveness of the training environment by providing clear guidance, realistic representation of components and challenges, interactive dialogue, and enhanced vision understanding. These elements collectively enhance the learning experience and skill development in XR assembly tasks.

A.3.5 Called Tools (Figure 7)

As shown in Figure 7, we plot the distribution of the number of tools invoked in the generated conversations. The most frequently called and essential functional tools are those related to process control: “NextStep” (57.02%), “StartAssembly” (4.58%), “CheckStepStatusVR” (4.28%), “GoToStep” (2.44%), “GetRemainingStep” (1.90%), “GetCurrentStep” (1.78%), “1.31%”. This indicates that users prioritize adherence to the assembly procedure during the fine-grained assembly task. Functional tools related to user interactions are also significant, for example, “HighlightCorrectComponents” (4.64%).

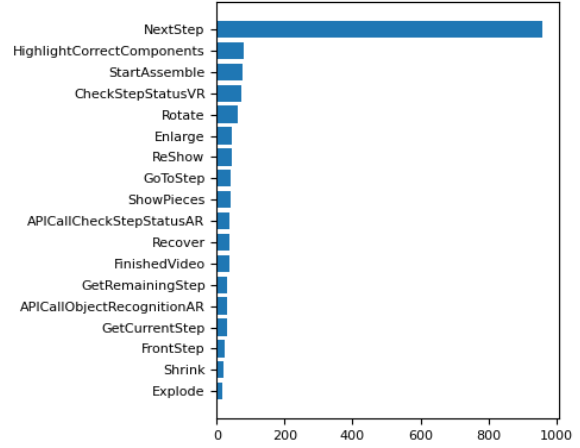


Figure 7: Distribution of called tools in conversations.

A.4 Engineering details in the workflow

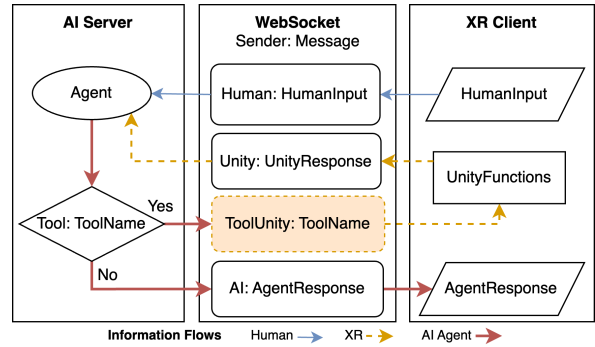


Figure 8: Information flow.