# GesPrompt: Leveraging Co-Speech Gestures to Augment LLM-Based Interaction in Virtual Reality

Xiyun Hu*
School of Mechanical Engineering
Purdue University
West Lafayette, USA
hu690@purdue.edu

Dizhi Ma*
Elmore Family School of Electrical
and Computer Engineering
Purdue University
West Lafayette, USA
ma742@purdue.edu

Fengming He
Elmore Family School of Electrical
and Computer Engineering
Purdue University
West Lafayette, IN, USA
he418@purdue.edu

Zhengzhe Zhu
Elmore Family School of Electrical
and Computer Engineering
Purdue University
West Lafayette, USA
zhu714@purdue.edu

Shao-Kang Hsia
School of Mechanical Engineering
Purdue University
West Lafayette, USA
shsia@purdue.edu

Chenfei Zhu
School of Mechanical Engineering
Purdue University
West Lafayette, USA
zhu1237@purdue.edu

Ziyi Liu
School of Mechanical Engineering
Purdue University
West Lafayette, USA
liu1362@purdue.edu

Karthik Ramani
School of Mechanical Engineering
Purdue University
West Lafayette, USA
ramani@purdue.edu

## ABSTRACT

Large Language Model (LLM)-based copilots have shown great potential in Extended Reality (XR) applications. However, the user faces challenges when describing the 3D environments to the copilots due to the complexity of conveying spatial-temporal information through text or speech alone. To address this, we introduce GesPrompt, a multimodal XR interface that combines co-speech gestures with speech, allowing end-users to communicate more naturally and accurately with LLM-based copilots in XR environments. By incorporating gestures, GesPrompt extracts spatial-temporal reference from co-speech gestures, reducing the need for precise textual prompts and minimizing cognitive load for end-users. Our contributions include (1) a workflow to integrate gesture and speech input in the XR environment, (2) a prototype VR system that implements the workflow, and (3) a user study demonstrating its effectiveness in improving user communication in VR environments.

## CCS CONCEPTS

• **Human-centered computing** → **Gestural input**; **Mixed / augmented reality**; **Virtual reality**; **Natural language interfaces**.

*Both authors contributed equally to this research.

## KEYWORDS

Co-speech Gesture, Multimodal Interaction, Large Language Models, Extended Reality, Virtual Reality

## 1 INTRODUCTION

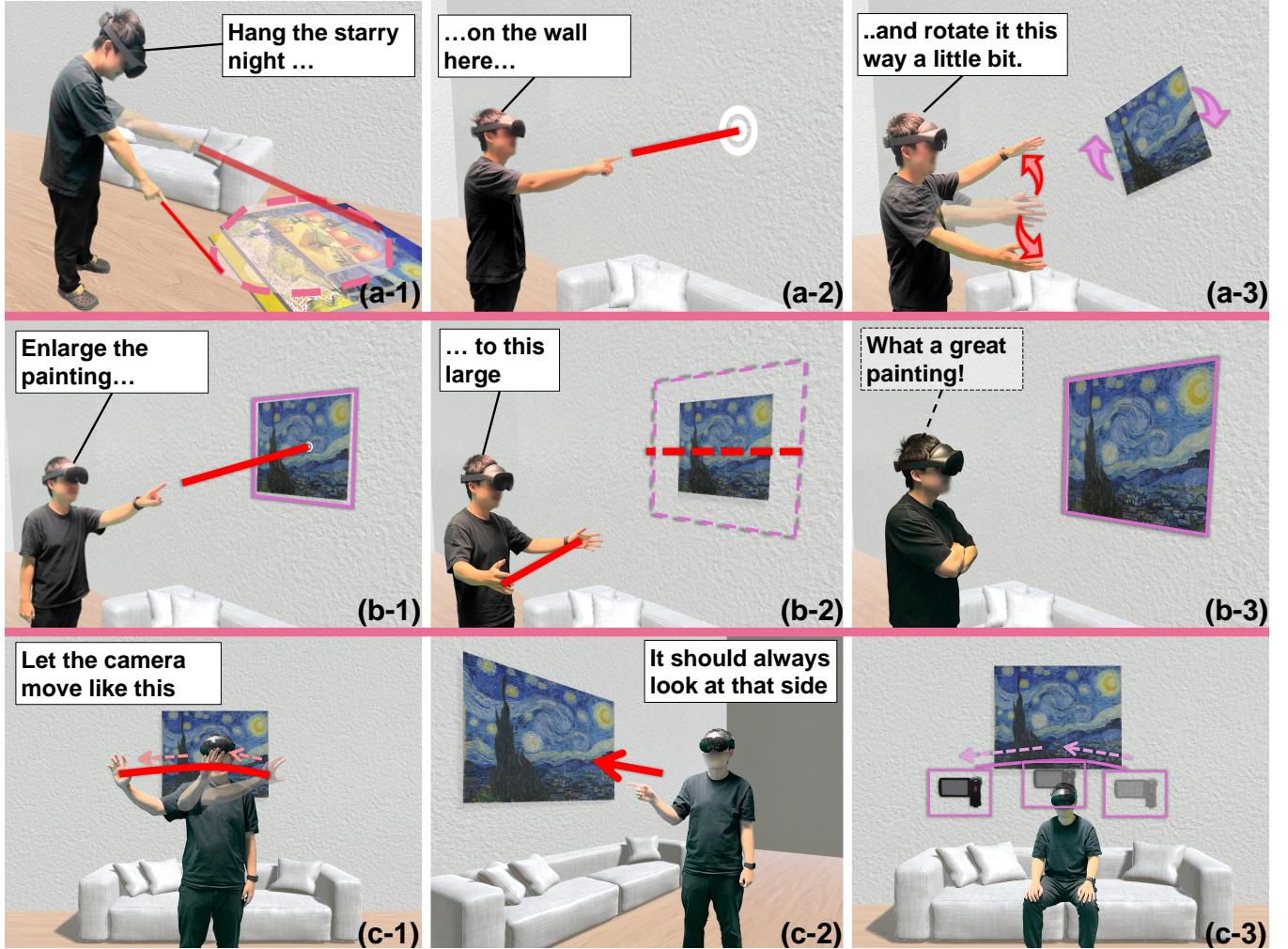With the development of Large Language Models (LLMs), LLM-based copilots on desktops and smartphones that assist end-users with auto-generated suggestions received a wide range of recognition in various domains, such as coding [37], writing [76], and day planning [5]. Typically, these applications leverage the textual context understanding ability of the LLM to understand the end-user's intention and offer assistance to fulfill the user's need. Recently, researchers have explored integrating LLM-based agents into extended reality (XR) environments for various applications, such as social games [99], education [25, 61], and spatial design [116]. Although the above-mentioned LLM-based copilots show great potential in assisting users in XR, interacting with 3D environments using text-only or speech-only prompts remains challenging for end-users to express complex spatial–temporal intents.

First, end-users must craft carefully engineered prompts or precise descriptions to ensure accurate LLM interpretation, often resulting in trial-and-error interactions. Second, tasks in 3D environments involve intricate spatial-temporal information, and such information is difficult to describe precisely and efficiently in natural language. Most copilots that focus on text/speech input struggle

**Figure 1: An example usage of *GesPrompt*: (a) An end-user wants to hang a virtual painting on the wall using *GesPrompt*: (a-1) The user asks the system to select the "Starry Night" painting from a pile of paintings on the floor by gesturing the rough location of the painting; (a-2) in the same utterance, the user also specifies the location of the painting using gesture; (a-3) realizing that the painting is not level, the user instructs the system to adjust its alignment while gesturing the appropriate amount of rotation; (b) The user stands back and sees that the painting is too small, then asks the system to enlarge the painting:(b-1) (b-2) the user sets the size of the painting by gesturing; (b-3) the user is satisfied with the painting, now properly placed and resized to fit the room's aesthetics. (c) Then the user is excited and wants to take a virtual video with the painting: (c-1) The user asks the system to set up the path of the camera by gesturing; (c-2) they also want the camera to look at a certain direction, which is pointed out by the user; (c-3) the camera moves along the path and records the video as the user desires.**

to enable end-users to accurately convey complex spatial-temporal details, as these text/speech-only inputs are prone to pronoun ambiguity [9, 56, 116]. For example, an end-user wants to refer to an unfamiliar object in a 3D environment to a text-only copilot, then the user has no choice but to describe the physical characteristics of the object in detail. This multi-step prompting increases cognitive load and disrupts workflow [17, 79]. Consequently, enabling effortless interaction for complex 3D tasks in XR remains an open challenge.

Some prior works address the problem of end-users' difficulty in describing spatial-temporal information by incorporating scene understanding. In these efforts, researchers guide end-users to describe the spatial-temporal information in an XR environment from high-level to low-level details [28, 109]. While these methods help and enable users to craft accurate prompts that are understandable by LLMs, the sheer number of users' prompts remains demanding. This results in a lengthy process that is cumbersome for end-users to complete. On the other hand, multimodal input—combining LLMs with human modalities such as gesture and gaze—offers a promising path forward. Gestures, as a natural interaction modality, have long been adopted for interacting with 3D environments in HCI research [39, 86, 90, 112]. Specifically, end-users can use

various gestures such as grabbing [111, 114], touching [13, 83], and air-taping [44, 67, 75] to manipulate objects in XR scenarios. Furthermore, compared with speech-only communication, gestures can convey essential non-verbal information that adds emphasis and clarity [11, 12, 15, 20, 94]. For example, end-users may want to use gestures to illustrate the desired size of a virtual object when describing it to the copilot. Recent XR systems (e.g., GazePointAR [56], VRCopilot [116]) leverage pointing alongside speech but remain limited to simple deictic input. Prior HCI work such as Grasp & Speech [80], Mixed Speech–Gesture Interfaces [33], and Multimodal Command Tools [24] typically follow a two-stage pattern: users first perform a manipulative gesture (e.g., pinch, drag) to select or highlight a virtual object, then issue a spoken command (e.g., "rotate," "delete") to trigger an action. While this sequential coupling supports straightforward object manipulation, it treats speech and gesture as largely independent phases and fails to exploit the rich semantic alignment that naturally occurs in co-speech gestures. Thus, we harness *co-speech gestures*—gestures naturally produced alongside spoken utterances—to fill gaps in verbal references and reduce prompting overhead.

We propose GesPrompt, a multimodal XR interface powered by an LLM that allows users to express their intentions through integrated speech and co-speech gestures, without requiring engineered commands or the mental translation of spatial-temporal information into numerical values. The GesPrompt workflow consists of two interconnected main components: (1) an LLM system that possesses both object information and the XR system function information. The LLM system takes the user's speech and converts it into a sequence of function calls that can be executed by the back-end XR system. Additionally, the LLM system identifies ambiguous phrases in the speech that cannot be understood through speech alone and sends them to the gesture processor. (2) a gesture processor that processes gesture inputs. The processor converts raw skeleton gesture data into meaningful parameter values (position, object, direction, rotation, size, and path) using cues, which are discovered by the LLM system, from speech. By jointly reasoning over speech and gesture, GesPrompt enables intuitive, low-effort interaction with 3D scenes.

In summary, we highlight our contributions as:

- A novel co-speech gesture integration method for LLM-powered XR copilots, extending beyond simple deictic pointing.
- A prototype VR copilot that implements the workflow for object manipulation using voice and gesture.
- A user study to evaluate the feasibility of the workflow and investigate user behavior when interacting with the system.

## 2 RELATED WORKS

### 2.1 LLM-based Copilots for 3D Scenes

As LLMs have emerged recently and gained widespread attention [7, 18, 95, 96], LLM-based copilots are becoming a convenient way for end-users without coding expertise to interact and communicate in the HCI area. With the assistance of copilots, end-users can easily create customized digital contents, including images [18, 93, 109], UI designs [48, 78], and 3D models/scenes [23, 43, 45, 60, 82, 104], using simple prompts.

Notably, significant advances in generative models for 3D scenes have captured HCI researchers' attention and motivated them to explore the integration of LLM-based copilots into interactive 3D environments. LumiMood [74] proposed a system that leverages LLMs to automatically adjust lighting and post-processing in 3D scenes, allowing the mood of scenes to dynamically adapt accordingly. Bozkir et al. [16] and Zhu et al. [119] forecast that integrating LLM into XR as narrative agents or visual avatars with different strategies will enhance users' engagement in XR. Aikawa et al. [9] propose an LLM-based system that assists end-users in brainstorming within the AR space. Meanwhile, MagicItem [54] and DreamCodeVR [38] enable end-users without coding expertise to prototype VR content's behavior. By utilizing LLM-based copilots to generate code within VR space, these two works empower users to program basic behavior of 3D contents, such as "move", "jump", and "rotate". Yet, the programmed behaviors are limited in scope as end-users are incapable of modifying more detailed aspects of the virtual content, such as altering the distance of movement. Besides LLM-based copilots, most state-of-the-art copilots such as GPT-4 [7] apply multimodal models, which accept different sources of inputs including texts and images. VR-GPT [52] provides end-users a natural way to complete object collection tasks by interacting with a Vision Language Model (VLM)-based copilot, which is capable of processing both images and language prompts within VR environments.

While the aforementioned works have received recognition for applying robust LLM-based copilots in XR, most of the works focus on entry-level design and prototyping for end-users. Consequently, end-users remain incapable of performing complicated 3D tasks involving detailed spatial-temporal information or dynamic interactions, such as precise movements of virtual content or responsive adjustments according to users' status. To let end-users accurately generate prompts in interactive XR environments, recent work such as LLMR [28] that incorporates scene understanding breaks down complex 3D tasks into several subtasks and employs an ensemble of multiple language models, each addressing a specific task. While this approach enhances the robustness of language models in XR environments, the sheer number of subtasks and associated prompt numbers makes the overall progress overwhelming and time-consuming for end-users. Motivated by these challenges, we strive to develop a system that streamlines the process of prompting LLM-based copilots for complex 3D tasks and interactions, so that end-users can effortlessly utilize these copilots in XR environments.

### 2.2 Gestures in XR Multimodal Interaction

Gestures, being a natural and intuitive human input, have long been recognized in the HCI area for facilitating interactions with 3D environments in daily activities [13, 19, 62, 101, 114, 118]. Through gestures, end-users can effortlessly manipulate 3D objects via body or hand movements. For example, end-users can easily use simple gestures, such as pointing [58], pinching [44, 67], and grabbing [114, 118] to select or move 3D content in XR environments. Beyond basic interactions, gestures also enable end-users to modify the complex behavior of virtual content. MagicHands [13] empowers end users to perform customized gestures to author particle

animations. GesturAR [101] enables end-users without coding expertise to prototype AR applications and modify AR content behaviors through customized freehand gestures. ProGesAR [113] integrates gestures with IoT devices, allowing end-users to trigger IoT functions through customized gestures. Besides providing intuitive interactions, gestures also allow end-users to effectively convey spatial-temporal information in XR environments [22, 73]. SketchingWithHands [49] allows users to use gestures to describe hand-held objects' shape or contour. KinÊtre [22] introduces a system that allows virtual content to change poses by mimicking end-users' postures in the 3D environment.

Combining gestures with speech offers a powerful multimodal approach. Early interfaces (e.g., Grasp & Speech [80], Freehand–Voice Manipulation [57], Mixed Speech–Gesture Interfaces [33], and Multimodal Interaction with Virtual Agents [24]) follow a two-phase pattern: users first perform a gesture to designate or highlight an object, then issue a spoken command (e.g., "rotate," "delete") to specify the action. While effective for discrete tasks, this sequential coupling treats speech and gesture as independent stages and relies on engineered command vocabularies.

Recently, researchers attempted to leverage the intuitiveness and expressiveness of gestures to assist end-users in overcoming the challenges of using LLM-based copilots for 3D tasks in XR environments. When end-users communicate with copilots in 3D environments, one major issue is the pronoun ambiguity for describing spatial-temporal information [8, 9, 56, 116]. Aghel et al. [8] present a Wizard-of-Oz elicitation study and highlight the importance of using gestures for reference when end-users generate prompts in 3D environments. Further, Wu et al. [108] investigate when and how full-body gestures can be used together with LLM-based copilots. Targeting the problem of pronoun ambiguity, GazePointAR [56] enables end-users to leverage eye gaze as well as pointing gestures to describe specific locations in 3D scenes when talking with copilots. Similarly, VRCopilot [116] combines speech and pointing gestures to allow end-users to specify their creation needs in immersive environments. However, these approaches support only simple location cues and do not leverage the full expressiveness of co-speech gestures, which naturally accompany spoken utterances and can encode rotation, size, and motion trajectories.

To address this gap, we leverage *co-speech gestures*—gestures produced concurrently with speech—to enrich LLM understanding of spatial-temporal intents. By jointly interpreting speech and gesture in our multi-modal XR interface, GesPrompt reduces the need for carefully engineered prompts or mental translation of spatial data into numerical values, enabling fluid and intuitive interactions with 3D scenes.

## 2.3 Synergy Between Speech and Gesture: Co-speech Gestures

In human communication, co-speech gestures can provide additional context, complementing speech to communicate intention more clearly [26, 47], especially regarding spatial and physical aspects [79, 110]. Following the definition from a pioneering work - Hand and Mind [94], co-speech gestures can be categorized into four types: deictic gestures, iconic gestures, metaphoric gestures,

and beats. This classification is widely adopted in the design and analysis of co-speech gestures [15, 29, 50, 66, 107].

**Deictic gestures**, often known as the pointing gesture [29], are used to refer to entities in the environment. The pioneer work Put-that-there [15] in the HCI area can serve as an example: the user says "put that there" while pointing to an object on screen and pointing to another location, adding missing information from the speech.

**Iconic gestures** are used to depict the spatial or action aspect of the elements in the speech [50, 66, 107]. For example, when someone says "you cannot believe the huge bird I saw today", they may spread both arms wide apart to visually indicate a large size.

**Metaphoric gestures** are similar to iconic gestures. They both describe the physical aspect of a concept, but the difference is that iconic gestures illustrate concrete contents while metaphoric gestures is used to present abstract concepts [92]. For example, the spreading arm gesture can also demonstrate the level of welcoming, when someone says "I welcome you with my whole heart."

**Beats** are often biphasic gestures to emphasize some words or phrases in a sentence [66]. For example, when someone says "you have to be there <u>now</u>" with a downbeat gesture, they mean to emphasize the timing.

Among the above gestures, we mainly focus on deictic gestures and iconic gestures due to their expressiveness in representing the spatial-temporal information.

In order to extract the spatial-temporal information from the deictic gestures and iconic gestures, the detection of such gestures is the first step. Deictic gestures, which convey only the meaning of referring to something, are usually static and can be detected with current gesture recognition models [53, 64]. In commercial systems such as HoloLens 2 [1] and Meta Quest Pro [3], a ray is constantly cast from the center of the user's hand, and the user can perform a pinch gesture to confirm the selection or reference of objects or entities. The iconic gesture, on the other hand, can represent the action aspect of an object, making the gesture sometimes dynamic. Moreover, the wide variety of shapes and actions that the iconic can represent make it even harder for the traditional gesture detection model to recognize [91]. Recent work by Ghaleb et al. [35] explores a framework that leverages speech information in detecting gestures. Zhang et al. [117] find it promising to use LLM in aid of synthesizing co-speech gestures. Building on these insights, and the fact that co-speech gesture will occur when the user says the word that is relevant to it [42, 98], our system utilizes speech to detect the appearance of the co-speech gesture, and analyzes the co-speech gesture based on the context given in speech.

Although deictic and iconic gestures are effective in representing certain actions or spatial relationships, they naturally have limitations in conveying more complex or abstract meanings. Speech, which evolved in part from gestures, complements these limitations by covering the meanings that gesture alone cannot easily express. This synergy between speech and gesture motivates our design of a system that integrates both inputs, allowing speech to guide the interpretation of gestures, and gesture to complete the missing information from speech.

# 3 GESPROMPT

In this work, we aim to let end-users communicate with LLM-based copilots while conveying spatial-temporal information in XR environments. To achieve this, we design a workflow that processes natural speech and gesture input from users simultaneously. When a user starts a speech and provides gesture command to GesPrompt, the system follows the following steps: (1) extracts users' inputs, including speech and gestures, (2) interprets user input by referencing appropriate virtual content and behaviors, and (3) determines the backend **functions** (F) responsible for controlling targeted behaviors.

Here, we formally define the process as follows: Let $F$ represent a set of available system *functions* $F_1, F_2, ..., F_n$, and let $S$ denote the current state of the XR environment. Calling these functions will alter the environment's state $S$. Each function $F_i$ takes a set of **function parameters** $P_i = \{p_1, p_2, ..., p_k\}$, where $p_i$ represents properties of the XR environment such as objects and spatial positions.

Given users inputs of (1) **text input** $T$: verbal descriptions of the desired state of the XR environment, and (2) **gestural input** $G$: gestural cues that provide additional context to the *text input*, such as hand movements for outlining the shape of an object in the environment, our goal is to:

- Infer the correct sequence of function calls
  $C = \{(F_{i1}, P_{i1}), (F_{i2}, P_{i2}), ..., (F_{im}, P_{im})\}$ from the *text input*.
- Extract the appropriate values for the *function parameters* $P_{im}$ from both text and gesture input.
- Apply the function calls $C$ to transition the XR environment state: $S_{current} \xrightarrow{C} S_{intended}$.

In the following sections, we describe how GesPrompt achieves the goal and its detailed design.

## 3.1 System Walk-through

We introduce GesPrompt, a workflow designed for users to interact with LLM-based copilots through speech and gesture. As illustrated in Figure 2, this workflow is composed of three main components: the **LLM System**, the **Gesture Processor**, and the **XR system**. The workflow of GesPrompt is demonstrated with an example scenario in Figure 2.

An end-user wants to hang a painting on the wall but struggles to verbally describe the exact position of the painting. With our system, which continuously listens to the user's command, the user simply says "Hang the Starry Night painting" while pointing to the paintings, then says "...on the wall here" while pointing to the wall. GesPrompt handles the speech and gesture input with the following steps: first, the *LLM System* identifies the appropriate backend function (F), for the task from the speech, i.e., selecting $F_{move}(object, position)$ for this scenario to place the painting on the wall. Then, the LLM System maps the **parameter tokens** from the speech to *function parameters*. In this scenario, "Starry Night painting" (object name) is mapped to the *object parameter* (a virtual object in the scene), and "here" is mapped to the *position parameter*. With the scene information (e.g., scene objects' names, positions, rotations, and scales) fed into the LLM System, the target painting

can be identified through the object parameter. However, the position parameter "here" is identified as an **ambiguous parameter** ($P_i^{amb}$), as its value cannot be determined.

To resolve this ambiguity, GesPrompt forwards the ambiguous parameter and the corresponding *timestamps* of its parameter token to the *Gesture Processor* for further analysis. In this case, a position parameter and the timestamps associated with the user's utterance of "here" are sent to the Gesture Processor. The gesture that occurred during the interval is then segmented. Since the parameter type is "position", the segmented gesture is identified as a pointing gesture, and the position where the user points is further used as the exact position parameter.

Finally, GesPrompt checks for any invalid or incomplete functions and parameter values in the *XR System*. If all functions and values are valid, the XR System then executes the *function calls* (C), transitioning the XR environment state to the user's intended state. In this case, the Starry Night painting is moved to the position where the user points. When the system detects missing or invalid parameters for the target functions, a textbox will pop up prompting the user to provide further clarification. In the following sections, we describe the design of each component of GesPrompt.
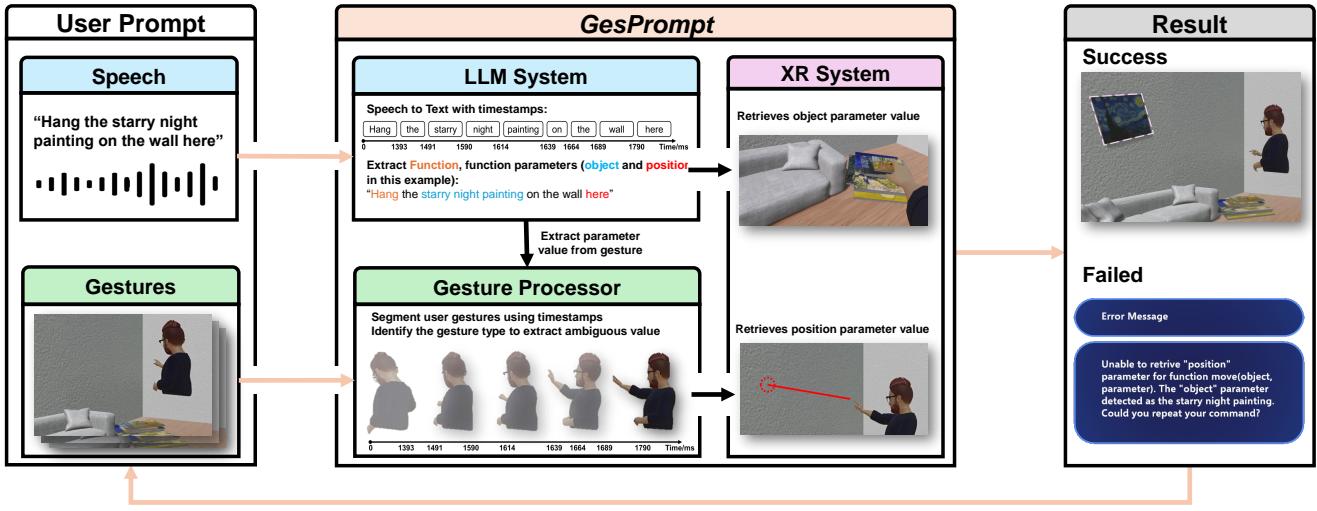
## 3.2 Spatial-Temporal Parameter and Co-Speech Gestures

Although it's possible to articulate spatial-temporal parameters through speech alone, integrating co-speech gestures often provides a more efficient and intuitive means of conveying these parameters. In this section, we concentrate on scenarios where users employ co-speech gestures to depict spatial-temporal parameters. As discussed in Section 2.3, co-speech gestures are context-dependent. A single gesture may convey different meanings based on the accompanying speech. We summarize the characteristics of co-speech gestures and their representations of spatial-temporal parameters in the context of an extended reality environment.
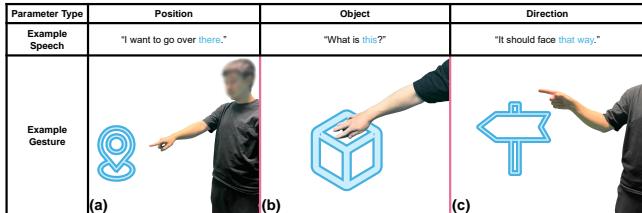
*3.2.1 Deictic Gestures.* This type of gesture indicates a location, entity, or direction within the environment. In a natural manner, index-finger-pointing is commonly adopted for entity selection, while whole-hand-pointing can be used in situations that require less precision [27]. At the same time, the pointing gesture also varies in different cultures. To accommodate variations in pointing deictic gestures, we visualize the user's pointing direction by constantly casting a ray from their hand. This visualization, which serves as guidance for users to indicate desired entities, is a common practice in virtual and augmented reality applications [69, 70].

As shown in Figure 3, the same deictic gesture can represent three types of spatial parameters under different speech contexts: (a) position, (b) object, and (c) direction.

**Position** refers to a 3D world coordinate. The user may point to indicate a specific location (Figure 3(a)). The intersection point between the pointed location and the ray cast from the user's hand provides the value for the position parameter. If the intended location is obscured by other objects, our system uses the user's speech context to disregard any mentioned obstacles. For instance, when organizing a virtual room, the user wants to move a chair behind a table in front of the user. If the user points behind the table and says

**Figure 2: System Workflow: The system modifies the XR environment using speech and gesture prompts, processed by two core components: the LLM system and the Gesture Processor. The LLM system converts speech to text with timestamps, identifies operation functions and parameters, and extracts unambiguous values from the XR scene. The Gesture Processor resolves ambiguous parameters by analyzing gestures segmented using speech timestamps and extracting values from the XR scene. Once all parameters are identified, the system updates the XR environment. If the system fails to update the environment, an error message will be sent to the user.**



**Figure 3: Example showcase of deictic gestures and the spatial parameters. For the Position parameter, the gesture could be pointing toward a specific location. For the Object parameter, the gesture could be pointing directly at an object. For the Direction parameter, the gesture could be a directional pointing motion to indicate the desired orientation**

"The chair should go there, behind the table.", the system interprets the table's presence and places the chair in the correct position.

**Object** refers to an entity in the environment, such as virtual contents, humans, or part of an object. The user may point to or directly touch an entity to refer to it (Figure 3(b)). We apply a similar method to process the object parameter as we do for the position parameter, with some variations. Instead of returning coordinates, the object that intersects with the ray is the output. Selecting an object requires less precision compared to setting the position parameter because an object's region typically encompasses a larger area than a precise location. Therefore, the touch gesture depicted in Figure 3(b) can be interpreted as an up-close pointing gesture intended to select the object being touched by the user.

**Direction** refers to a spatial orientation or a directional vector in the environment. By pointing within the environment (Figure 3(c)), the user specifies a direction. The directional vector of the ray cast from the user's hand provides the value for this parameter.



**Figure 4: Examples of iconic gestures and the corresponding spatial parameters. For the Rotation parameter, gestures could include the rotation of one (a) or both (b) hands to represent rotational movements. For the Size parameter, gestures could indicate the relative size of an object using the distance between fingers (c), hands (d), or between a hand and a reference entity (e). For the Path parameter, gestures could involve circular hand motions (f) outlining the shape of a sphere or demonstrating the trajectory of an object.**

*3.2.2 Iconic Gestures.* Instead of referencing external entities, the user can also refer to the hand gestures themselves through iconic gestures to express spatial parameters. We summarized three categories of spatial parameters that can be represented by iconic
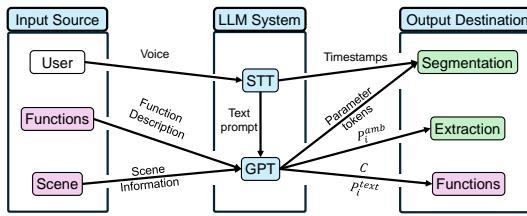
gestures, as illustrated in the examples in Figure 4: (a) Rotation, (b) Size, (c) Path.

**Rotation** also refers to spatial orientation, but in this case, the gesture mimics the full rotation process. For this parameter, there are two cases: the user imagines holding and rotating the object with one hand (Figure 4(a)) or two (Figure 4(b)). By monitoring changes in the position of the user's hands during the gesture, we can determine whether one or both hands are involved. If only one hand is used, the palm's rotation change during the gesture will be the output. If both hands are involved, we calculate the rotation change based on the line formed between the hands.

**Size** refers to the dimension or scale of an object. Deriving from it, length, extent, and degree parameters share the same characteristics as the size parameter, as they all relate to measurable aspects of entities. The size parameter can be represented using three types of gesture: using the distance between fingertips (i.e. thumb and index) (Figure 4(c)); using the distances between hands (Figure 4(d)); using the distance between the hand and another entity (e.g. hand hovering over the tabletop surface to indicate height) (Figure 4(e)). We calculate the distance accordingly using the gesture data, and the resulting distance becomes the value for this parameter.

**Path** refers to a series of positions, rotations, and timestamps. By demonstrating a path, the user can sketch a shape in mid-air using both hands, as shown in Figure 4(f). Alternatively, the user can generate a trajectory for a camera movement, as depicted in Figure 1(c). The sequence of hand joint positions and rotations will be the value for the path parameter.

## 3.3 LLM System



**Figure 5: LLM System consists of a speech-to-text (STT) module and a GPT module. It processes user voice input and data from the XR environment. It then outputs timestamps, parameter values, and tokens to the gesture processor and the XR environment for further processing.**

The LLM system in the workflow is composed of two key modules, as shown in Figure 5: the speech-to-text (STT) module and the GPT module.

The STT module serves as a pre-processor for users' speech. Specifically, it converts users' speech into text and records the timestamp for each word. After the processing, the converted text is forwarded to the GPT module for further analysis, and the timestamps of words are forwarded to the *Gesture Processor* for receiving gesture inputs.

Given the converted text from the STT module, along with scene information and function descriptions-including each function's

purposes and *function parameters* (Figure 5) -the GPT module interprets the converted text to determine appropriate functions and corresponding parameter values so that the XR environment can be changed to users' intended state. In detail, the GPT module consists of three steps: (1) Convert the user prompt into a sequence of function calls. (2) For each function call, attempt to identify the values of the parameters in the prompt. (3) Forwards *ambiguous parameters* $p_i^{amb}$ to *Gesture Processor* for further calculation. The system prompt of the GPT module is included in Appendix A.1.

Now we illustrate how the GPT module accomplishes the above steps. Based on the user's prompt, the GPT module first determines which function calls are necessary to control the XR environment's state. For complex prompts that require multiple function calls, the GPT module breaks the prompt into sub-prompts, allowing each to be resolved by an individual function call. Next, the GPT module traverses users' prompts to extract the appropriate values for the function parameters $P_i^{text}$. The types of $P_i^{text}$ extend beyond the types of spatial parameters. Any parameter value that can be determined from the user's speech can be considered as $P_i^{text}$, such as colors, numeric values, and basic shapes.

If a parameter's value cannot be identified due to ambiguity in the textual description, it is marked as *ambiguous parameters* $P_i^{amb}$, and the parameters have the following constraints:

$$P_i = P_i^{text} \cup P_i^{amb} \tag{1}$$

and

$$P_i^{text} \cap P_i^{amb} = \emptyset \tag{2}$$

It means that the $P_i$ is composed of $P_i^{text}$ and $P_i^{amb}$, while there is no overlap between $P_i^{text}$ and $P_i^{amb}$.

Then, the $P_i^{amb}$ and corresponding *parameter tokens* will be passed to the *Gesture Processor* for further processing.

## 3.4 *Gesture Processor*



**Figure 6: Gesture Processor comprised of a segmentation module and an extraction module. It handles input from the user's gestures, parameter data from the LLM system, and scene information from the XR environment. It then generates gesture parameters for integration within the XR environment.**

We designed a *Gesture Processor* to retrieve the spatial-temporal parameters. Specifically, given $P_i^{amb}$, *parameter tokens*, and timestamps of each word from the LLM system, the *Gesture Processor* analyzes gestural input and calculates the values for the *ambiguous parameters* $P_i^{amb}$.

The *Gesture Processor* consists of two components: segmentation and extraction. As shown in Figure 6, in the segmentation step,

the *Gesture Processor* first retrieves user gestures using timestamps of *parameter tokens*. For example, in Figure 1 (b), the user says 'enlarge the painting to this large'. *Gesture Processor* segments users' gestures when user says *parameter tokens* 'the painting' and 'this large'. In the subsequent extraction step, based on the parameter type, the *Gesture Processor* uses the method described in Section 3.2 to calculate parameter values ($P_i^{Gestural}$) from the segmented gesture data.

*3.4.1  Segmentation.* People tend to semi-automatically synchronize key phrases (*parameter tokens*) with the co-speech gesture [42, 98]. During the segmented co-speech gesture, the gesture encapsulates the majority of the meanings integral to the speech [36, 51, 87]. Based on these findings, GesPrompt utilizes the timing of when *parameter tokens* are articulated to extract the corresponding second stage of the co-speech gesture, subsequently processing the gesture according to the parameter type.



**Figure 7: (a) Hand joints and hand ray used in GesPrompt. (b) Process of creating a cone to select multiple object: (b-1) a set of original hand rays from the segmented gesture; (b-2) the average direction of the hand ray and the circle fit to the hand ray origins, the vertex is 30 cm from the circle; (b-3) Cone volume generated based on the vertex and circle, the height of the cone is set to 10 m.**

*3.4.2  Extraction.* For the gesture data, the positions and rotations of the thumb, index, middle tip joints, index metacarpal joints, and the palm joints for both hands are recorded (Figure 7(a)). The difference in palm positions between each frame and the initial frame is computed. A threshold is then applied to this difference to determine if there is hand movement within the segment. If the segmented gesture corresponds to spatial parameters that can be expressed by the deictic gestures, the hand ray (Figure 7(a)) will be used to determine the parameter value. The hand ray originates from the index metacarpal joint and points to the index fingertip joint. The length of the ray is set to 10 m. For the position parameter, we calculate the average position of intersection points.

For the object parameter, a cone-shaped bounding volume is generated to encompass all intersecting objects. This is accomplished using the method illustrated in Figure 7(b). From a segment of gesture data, the hand ray is computed for each frame, as shown in Figure 7(b-1). The cone's directional vector is derived by averaging these hand-ray vectors. Assuming the index tip positions describe a circle, this circle is fit using the least-squares method on a plane, with the plane's normal corresponding to the directional vector and the average palm positions lying on it. Empirically, the vertex of the cone is positioned 30 cm behind this plane, as indicated by the blue point in Figure 7(b-2). The cone's height is defined as 10 m.

Using the directional vector, circle, vertex, and height parameters, the cone-shaped volume is formed, as depicted in Figure 7(b-3).
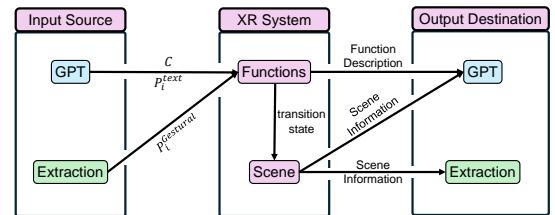
For the direction parameter, the average hand ray is calculated to be the vector for this parameter.

For the rotational parameter, if a single hand is involved, the rotational difference of the palm between the initial and the final frame is utilized. For gestures involving two hands, the line connecting the palm joints is considered, and the rotational difference of this line between the first and the last frame is used as the parameter's value.

For the size parameter in one-hand gestures, the average distance between the index fingertip joint and the thumb tip joint is utilized. If the average distance between the index tip and middle tip joints falls below a certain threshold, a ray is projected from the palm joints. Unlike the hand ray, the direction of this ray follows the palm joint's forward direction (the direction the palm faces) and terminates at the closest surface. The average length of these rays will determine this parameter value. In the case of two-hand scenarios, the mean distance between the palm joints of each hand is considered.

For the path parameter, the positions and rotations of the six joints for each hand are utilized as the parameter value.

## 3.5  XR System



**Figure 8: XR system includes available functions and scene information. It provides and obtains data to the LLM system and gesture processor, and updates the scene state correspondingly.**

The *XR System* manages available functions and scene information as shown in Figure 8. It receives function calls $C$ and associated parameters $P_i^{text}$ from the LLM system, as well as $P_i^{Gestural}$ from the *Gesture Processor*. If any parameters remain unresolved, the user will be prompted for further clarification through speech or gestures. GesPrompt presents a text box requesting the user to provide clarification, as illustrated in Figure 2 failed cases. The error message details the function that GesPrompt intends to execute, along with the parameters that are missing. When the user attempts to call functions that are not available in the XR system, a general error message will be shown to the user ("Sorry, the system is unable to do that, the system is able to do [list of functions]."). Once all parameters are confirmed, the *XR system* executes the function calls to transition the scene to the user's intended state.

In the prototype system, a set of XR functions that assist fundamental XR manipulation was implemented to showcase the feasibility of GesPrompt. Researchers in the future may develop specialized functions to facilitate additional interactions or opt for a system

such as LLMR [28], which employs LLMs to generate scripts derived from the user's prompt. When an XR function is added to the XR system, it is necessary to revise the metaprompt for the GPT module within the LLM system to incorporate the function's description and parameters.

Here we summarize the XR functions that are implemented in the prototype system:

- select(list of objects, color): selects objects of a certain color from a list of objects.
- move(object, position): move an object to a new location.
- rotate_dir(object, direction): set the forward direction of the object to the target direction.
- rotate(object, rotation): applies the rotation to the object.
- resize(object, size): change the size of the object to the target size
- move_path(object, path): let an object move along a trajectory back and forth.
- draw_path(path, shape_type): sketch a predefined shape given the path. The shape_type includes straight lines, circles, and sine waves. We employ a linearized least-squares approach for circle fitting, ordinary least-squares regression for straight lines, and nonlinear least-squares optimization for sine waves.
- set_color(object, color): set the color of an object to a given color.

The scene information sent to the GPT module is formatted in JSON and contains the name, position (meters), rotation (degrees), and scale (meters) of every object present in the scene. Each parameter's numerical value is rounded to three decimal places. An example can be found in Appendix A.3.

## 3.6 Implementation

We developed a prototype VR system that implements the proposed framework in Sec 3. The application is developed on the Unity3D platform (2022.3.20f1) [6] and runs on the Meta Quest Pro headset [3]. The user's voice input is converted into text using the Azure AI speech service [4]. For the LLM system, we used GPT-4o model [77] to process the given text. Skeletal gesture data are acquired using native sensors from the Meta Quest Pro headset and the Meta XR SDK [69].
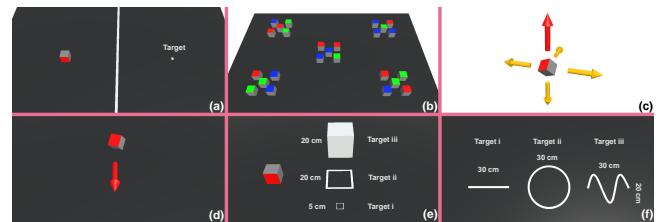
## 4 USER STUDY

We conducted a three-session IRB-approved user study to evaluate: (1) the accuracy of combining natural language and gesture input using GesPrompt , (2) the difference between GesPrompt, gesture-only system, and voice-only system, and (3) the usage of GesPrompt in interior design scenario. The entire study lasted approximately 1.5 hours, and each participant received a $20 e-gift card. The study was conducted in a 5m by 5m indoor area and was screen-recorded for post-analysis. Upon arrival, participants were given an introduction to the system, including a detailed explanation of the workflow and user interface. After three sessions, each participant completed a 5-point Likert-scale questionnaire and the standard System Usability Scale (SUS). We concluded the study with a conversation-style interview to gather subjective feedback on the system.

## 4.1 Participants

We recruited 12 participants (ten males and two females), ranging in age from 18 to 26. Of these, 1 is an experienced developer of head-mounted AR/VR applications, 8 had some experience with head-mounted AR/VR applications, while the remaining three had only heard of AR and VR. Four participants use ChatGPT often, and eight participants use ChatGPT almost every day. Additionally, two participants had no experience with Prompt Engineering, five participants had heard of this concept, two participants had tried it before, and three participants almost always uses prompt engineering in communicating with ChatGPT. None of the participants had used our system prior to the study.

## 4.2 Session 1



**Figure 9: User study session 1 setup: (a) position task; (b) object task; (c) direction task; (d) rotation task; (e) size task; (f) path task.**

*4.2.1 Procedure.* In this session, we evaluated the system's performance in analyzing both speech and gesture data. In alignment with earlier studies on fundamental object manipulation within virtual environments [21, 68], we designed six tasks to examine each spatial parameter independently:

(1) **Position**: move a cube on the left side of the table to the target (Figure 9(a)). The initial position of the cube is randomized in the left region.

(2) **Object**: select cubes of a certain color from piles of cubes (Figure 9(b)). A total of 25 colored cubes, 9 red, 8 green, and 8 blue on the table. The target color and the position of the cubes are randomized for each trial. We count the number of correct cubes selected (TP), incorrect cubes selected (FP), not selected correct cubes (TN), and not selected incorrect cubes negative selection (FN).

(3) **Direction**: set the front face (marked in red) of the cube to one of the five directions (Figure 9(c)). The initial rotation of the cube is randomized in 360° on the x, y, z axis. The target direction is highlighted in red. The order of the target direction is randomized for each participant.

(4) **Rotation**: rotate a cube so that the front face is facing the participant (Figure 9(d)). The initial rotation of the cube is randomized in 180° on the x, y, z axis, and the red side is always visible to the participant. Half of the participants are asked to use one hand, and the other half to use two hands to rotate the object.

(5) **Size**: resize the cube to the size of the target (Figure 9(e)). The initial length of the cube is set to 10 cm. One-third of
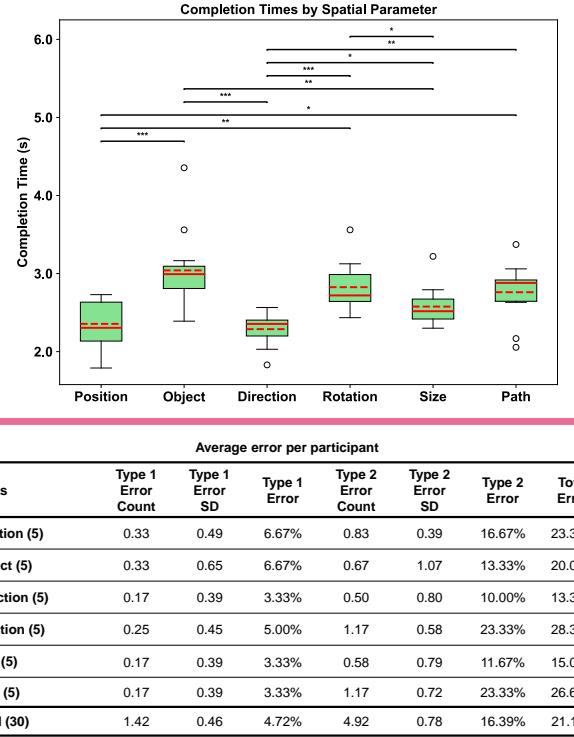
the participants each are asked to use (i) one hand (target length 5 cm), (ii) two hands (target length 20 cm), and (iii) use the table as a reference (target length 20 cm) to resize the object.

(6) **Path**: create a path base on the reference path (Figure 9(f)). The reference path is shown on the table. The participants were asked to draw the path mid-air using the index finger of their dominant hand while giving the speech command. One-third of the participants each are asked to define (i) a straight line (target length 30 cm), (ii) a circle (target diameter 30 cm), and (iii) a sine wave (target length 30 cm, amplitude 20 cm, period of 1.5).

In order to evaluate one spatial parameter per task, the cube object will be automatically chosen for tasks aside from the object task. Additionally, following prior work [28, 101], the following evaluation metrics for failure cases are also used: (1) whether the LLM system could identify $P_i^{amb}$ and the *parameter tokens* (type 1 error) and (2) whether the gesture data were correctly analyzed based on the associated parameters and phrases (type 2 error). A failed attempt is characterized by the authors being able to comprehend the participant's command through speech and gesture, recognizing it as appropriate for the task, yet GesPrompt does not generate an accurate outcome. Participants were instructed to execute each task five times. The successful trials were utilized to assess the parameter metrics, while unsuccessful attempts were included for error analysis. Completion time refers to the interval between the moment the participant initiates the command and when the system finalizes the execution of the function calls. The Mann-Whitney U test is used to compare the results.

Participants were asked to treat the target as a mental reference, and they should not explicitly mention it in their commands (i.e., they should not say "move the cube to the target" without any gesture). The LLM system does not have information about the target. We instructed the participants to interact with the system as if they were communicating with a human and to request the system to adjust the cube to the target state. The participant is informed of the task type before each trial. trail No additional text or voice instructions were provided to minimize influence on their commands.

*4.2.2 Results and Discussion.* All users completed the tasks with an average completion time per command of 3.28 seconds (SD = 0.85) and an average total error rate of 21.11%. The quantitative comparison results of completion time for each task are shown in Figure 10 top. We observe that the average completion time for commands involving parameters "Position" (AVG = 2.71, SD = 0.60), "Direction" (AVG = 2.58, SD = 0.44), "Size" (AVG = 3.15, SD = 0.50) are lower than "Object" (AVG = 4.08, SD = 1.03), "Rotation" (AVG = 3.65, SD = 0.63), "Path" (AVG = 3.52, SD = 0.72). Of the system errors, shown in Figure 10 bottom, the type 1 error rate out of all attempts is 4.72% and the type 2 error rate is 16.39%. The highest total error rate occurred in the rotation task (28.33%), followed by the path task (26.67%), position task (23.33%), object task (20.00%), size task (15.00%), and direction task (13.33%). The gestures used in the object, rotation, and path task are dynamic, implying that even a minor alteration in gesture can lead to significant changes in the outcome. Thus, tasks with dynamic gestures require more time. At



| Tasks | Type 1 Error Count | Type 1 Error SD | Type 1 Error | Type 2 Error Count | Type 2 Error SD | Type 2 Error | Total Error |
|---|---|---|---|---|---|---|---|
| Position (5) | 0.33 | 0.49 | 6.67% | 0.83 | 0.39 | 16.67% | 23.33% |
| Object (5) | 0.33 | 0.65 | 6.67% | 0.67 | 1.07 | 13.33% | 20.00% |
| Direction (5) | 0.17 | 0.39 | 3.33% | 0.50 | 0.80 | 10.00% | 13.33% |
| Rotation (5) | 0.25 | 0.45 | 5.00% | 1.17 | 0.58 | 23.33% | 28.33% |
| Size (5) | 0.17 | 0.39 | 3.33% | 0.58 | 0.79 | 11.67% | 15.00% |
| Path (5) | 0.17 | 0.39 | 3.33% | 1.17 | 0.72 | 23.33% | 26.67% |
| Total (30) | 1.42 | 0.46 | 4.72% | 4.92 | 0.78 | 16.39% | 21.11% |

**Figure 10: User study session 1 result. Top: completion times by spatial parameter, "*" indicates p-value <0.05, "**" indicates p-value < 0.01, "***" indicates p-value < 0.001. The red solid line represents the median, and the red dashed line represents the mean. Bottom: average error per participant.**

the same time, accurately segmenting the portion of the gesture when it occurs is crucial in analyzing the dynamic gesture data, as it contributes to an elevated error rate in these tasks.

For successful trials, we assess the accuracy of the results for individual tasks (Table 1). The position metric is determined by comparing the object's final position with its target (AVG = 3.87 cm, SD = 1.61 cm). The system exhibits a good understanding of the intents involved in the position task. For object metrics, we evaluate precision ($\frac{TP}{TP+FP}$) and recall ($\frac{TP}{TP+FN}$). The average precision achieved is 100%, indicating that all selected cubes are of the correct colors, demonstrating that GesPrompt can comprehend object attributes from speech. The average recall is 96.56%, suggesting that the selection of objects by cone misses some objects outside of it, likely due to an incomplete gesture segment. For direction and rotation, the same degree difference metric is employed. We calculate the least amount of rotation needed for the object to rotate to the target rotation. Compared with the direction task (AVG = 2.37°, SD = 2.30°), the rotational difference is significantly higher for the rotation task (AVG = 12.52°, SD = 8.93°, $p = 0.0025$). This discrepancy could be attributed to the fact that the rotation task requires dynamic gestures, whereas the meaningful segment of the gesture in the direction task is more static. Furthermore, since the cube remains stationary during the command, it may leave the

**Table 1: Summary of quantitative measurements for different spatial parameters. Position represents the distance between the object and the target's placement. Direction captures the angular difference between the object and the target. Rotation indicates the least rotation required to align the object with the target. Size measures the percentage difference in size between the object and the target. Path evaluates the similarity between the user's path and the intended path using DTW-based similarity.**

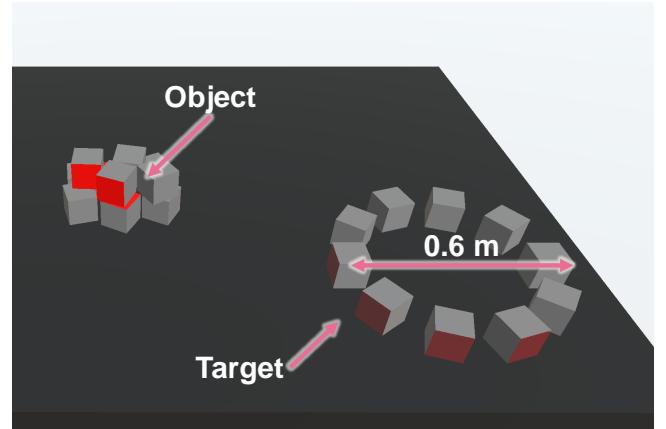| Parameter | Mean | STD |
|---|---|---|
| Position (m) | 0.0387 | 0.0161 |
| Direction (°) | 2.37 | 2.30 |
| Rotation (°) | 12.52 | 8.93 |
| Size (%) | 16.84 | 12.36 |
| Path (%) | 91.47 | 8.23 |
| **Object Selection** | | |
| Precision (%) | 100.00 | |
| Recall (%) | 96.56 | |

**Table 2: Quantitative measurements for rotation, size, and path tasks under different conditions. Rotation (i): one hand; Rotation (ii): two-hand; Size (1): one hand; Size (ii): two-hand; Size (iii): refer to the table; Path (i): straight line; Path (ii): circle; Path (iii): sine wave**

| | Mean | STD |
|---|---|---|
| Rotation (i) | 17.89° | 8.17° |
| Rotation (ii) | 10.82° | 8.22° |
| Size (i) | 23.41% | 11.39% |
| Size (ii) | 12.54% | 11.03% |
| Size (iii) | 18.18% | 13.28% |
| Path (i) | 95.25% | 2.34% |
| Path (ii) | 90.67% | 3.45% |
| Path (iii) | 88.50% | 7.37% |

user uncertain about how much hand rotation is necessary. For the size parameter, we calculated the percentage size difference to assess its accuracy, i.e., $\frac{||ObjectLength-TargetLength||}{TargetLength}$. For the path parameter, dynamic time wrapping-based similarity was used as a metric. The average similarity is 91.47% and the standard deviation is 8.23%.

Table 2 provides a detailed examination of system accuracy for rotation, size, and path tasks under various conditions. Performance with two hands (AVG = 10.82°, SD = 8.22°) in the rotation task is more precise than with one hand (AVG = 17.89°, SD = 8.17°), though not significantly. For the size task, two-hand usage results in a smaller size percentage difference (AVG = 23.41%, SD = 11.39%) compared to one-hand (AVG = 12.54%, SD = 11.03%, $p = 0.044$) and the reference (AVG = 18.18%, SD = 13.28%, $p = 0.035$). In the path task, drawing a line leads to higher similarity (AVG = 95.25%, SD = 2.34%) than drawing circles (AVG = 90.67%, SD = 3.45%, $p = 0.0031$) or sine waves (AVG = 88.50%, SD = 7.37%, $p = 0.004$).

## 4.3 Session 2



**Figure 11: User study session 2 setup**

*4.3.1 Procedure.* In this session, we evaluated the difference between GesPrompt and two baseline systems. Guided by [28, 55], we designed an A/B comparison study to compare user satisfaction when using GesPrompt, gesture-only system, and voice-only system to complete the same spatial task, which is to manipulate the 10 red cubes on the table so that they match the target state. The setup for session 2 can be found in Figure 9. The task is structured to encompass all six spatial parameters and is executable with all three systems. This task incorporates all six spatial parameters as previously indicated: (1) the user is required to select an **object**, (2) generate a circular **path**, (3) define the **position** of this path, (4) determine the **size** (diameter) of the circle, (5) as well as establish the **rotation** or **direction** of each cube.

For the gesture-only system, only the manipulative gestures are enabled. The user is restricted from using predefined Meta XR SDK hand gestures to grab and release an object. Once the object is grabbed, it will follow the position and rotation of the user's hand. While with GesPrompt, none of the predefined gestures are enabled. GesPrompt only takes speech and raw gesture data as input. For the voice-only system, we remove the gesture processor from GesPrompt. The metaprompt, shown in Appendix A.2, is modified to instruct the GPT agent to find the values of all parameters solely on the basis of the prompt and scene information given. We asked the user to treat the target as a mental reference and gave them the same instructions as in session 1. The time limit for using each system was restricted to 10 minutes. We recorded the total interaction time with the system. For the gesture-only system, the interaction time is defined as the time between the first manipulation of the object and the last manipulation. For the other two systems, the interaction time is defined as the time between the first command and the last command. The order of use of GesPrompt and the baseline systems was randomized for each user. At the start of the task, the user was informed of which system they were using. The user was directed to arrange the cubes in a configuration as similar as possible to the target. They were allowed to stop whenever they believed they had achieved their best effort in organizing the cubes.
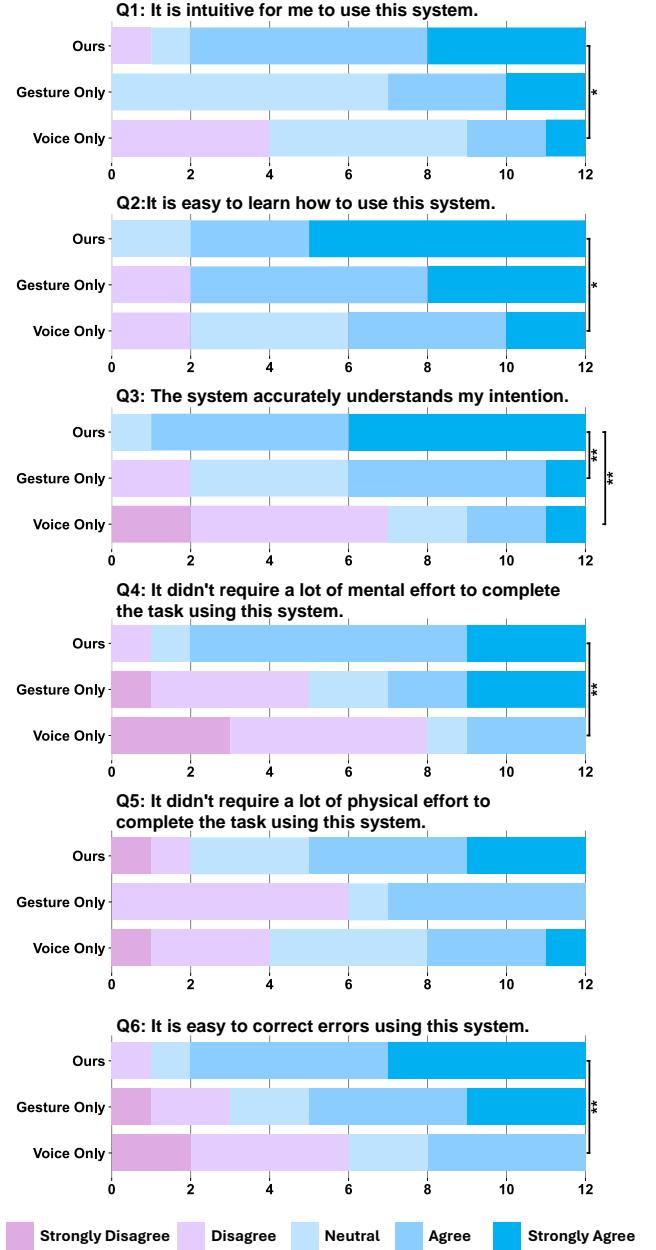
We applied Mann-Whitney U tests to compare the results from the questionnaire.

| | | GesPrompt | Gesture-only | Voice-only |
|---|---|---|---|---|
| | Input | Voice & free-form co-speech gesture | Predefined single-hand grabbing gesture | Voice commands |
| | Example Action | Say "Put the cubes in a circle" while gesturing the **circle** shape | Grab and place each cube to form a **circle** | Say "Put the cubes in a circle" without gesturing |
| | Average Interaction Time(s) | 43.65 (SD = 15.55) | 91.11 (SD = 29.46) | 108.58 (SD = 29.46) |
| Example Participant Results | | | | |

**Figure 12: Session 2 conditions and example participant results.**

*4.3.2 Results and Discussion.* The 12 users completed the task using the three systems within the time limit. The Likert-type results collected from this session are shown in Figure 13. The average interaction time with each system is: gesture-only system 91.11 seconds (SD = 23.00); voice-only system 108.58 seconds (SD = 29.46); and our system 43.65 seconds (SD = 15.55). The SUS results show that GesPrompt has good usability, with a mean score of 86 out of 100. Moreover, the gesture-only system obtained a SUS score of 74, while the voice-only system received a score of 59.

Overall, users found it intuitive (Q1: AVG = 4.08, SD = 0.25) and easy to learn how to use our system (Q2: AVG = 4.42, SD = 0.21). "I would use gesture plus speech to communicate with others in my daily life. (P5)". In comparison to the voice-only system (Q1: AVG = 3.00, SD = 0.26; Q2: AVG = 3.5, SD = 0.28), GesPrompt is more accessible as users are not required to construct their prompts. Our system (Q3: AVG = 4.14, SD = 0.86) demonstrates intent understanding ability higher than the gesture-only system (Q3: AVG = 3.42, SD = 0.25) and voice-only system (Q3: AVG = 2.58, SD = 0.0.34). *"I did not expect that your system can actually understand what I mean by saying 'that'. (P2)"*. In the gesture-only system, some users (P4, P7) with less HMD AR/VR experience report that the hand gesture for grabbing and releasing objects feels unnatural, causing frequent failures to grab their intended object. In the voice-only system, some users felt frustrated when the cubes did not line up in a circle, even after they specifically told the system to create 10 positions in a circular layout and move the cubes there. It did not require much mental effort to use our system (Q4: AVG = 4.00, SD = 0.23) compared to the voice-only system (Q4: AVG = 2.33, SD



**Figure 13: Session 2 Likert-type questionnaire results; "*" indicates p-value <0.05, "**" indicates p-value < 0.01.**

= 0.32). *"I don't need to calculate or refer to the size of the table to change the diameter of the circle (with our system). I can just do a gesture. While in the voice-only system, I have to think about them and explicitly tell the GPT. (P9)"*. Regarding the physical effort, the users find GesPrompt (Q5 : AVG = 3.58, SD = 0.34) demands less than the other two systems (Gesture-only (Q5 : AVG = 2.92, SD = 0.27); Voice-only (Q5 : AVG = 3.00, SD = 0.31)). *"It is frustrating to manually organize the 10 cubes in a circle (with gesture-only system).*

*When I add a new cube to the circle, I also have to adjust the positions of other cubes to make it look like a circle. (P10)"*. Most users agree that it is easier to correct errors (Q6 : AVG = 4.17, SD = 0.25) using our system in comparison to the voice-only system (Q6 : AVG = 2.66, SD = 0.32). *"I just don't know what I should say for the voice-only system to understand me. (P4)"* Yet, one user said that *"I would like to know that the system understands which object I am talking about before I can give the following command. (P8)"*. We will further discuss this concern in Section 5.1.
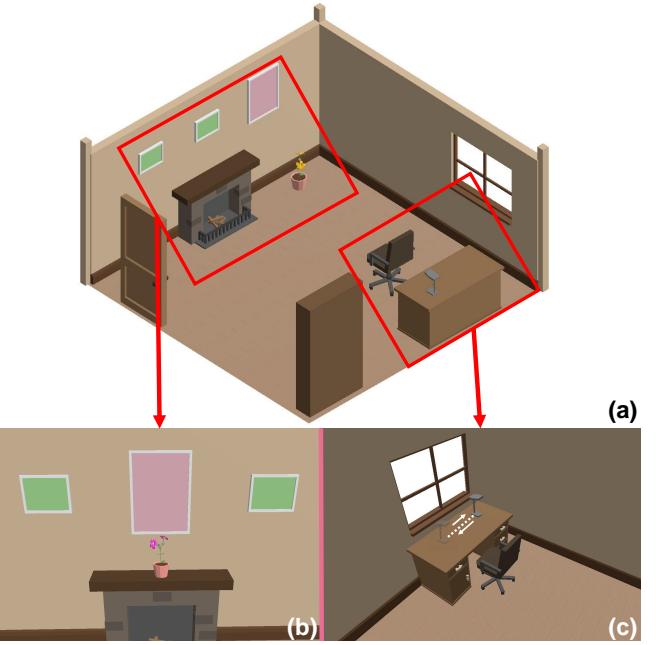
Although the interaction time for the gesture-only system is close to that of the voice-only system, the resulting cube configuration using the gesture-only system did not match the target cubes. After placing the cubes on the edge of a rough circle, most of the users requested to skip the task and said that they were satisfied with the result. The longer interaction time of the gesture-only system compared to our system is due to that the user has to try multiple times to rotate the cube to a suitable orientation and space the cubes evenly along the circle. For the voice-only system, we observed that users gave simple prompts like *"put the cubes in a circle"*, expecting the LLM agent to put the cubes in a spatially reasonable position (i.e., on the table). This observation agrees with Manesh et al. [8]. Although the LLM agent has all scene information, the outcome of executing this command is that the cubes are on the floor rather than the table, and the form resembles an ellipse rather than a circular shape (Figure 12). As a result, the user needed to experiment and incrementally add spatial constraints to the prompt, increasing the interaction time with the voice-only system compared to GesPrompt. For GesPrompt, the users also tried simple prompts at first. When that fails, our system will let them know which parameter value was missing. The users quickly understand the system, and in the next command, add extra co-speech gestures that represent the value. Thus, the interaction time is shorter than that of the two baseline systems.

## 4.4 Session 3

*4.4.1 Procedure.* In this session, we aim to evaluate the overall usability of GesPrompt through a complex interior design task. Participants were instructed to rearrange furniture within a VR room using GesPrompt. The VR room contains manipulatable objects, specifically, three pictures, a table, a lamp, a chair, and a flower. There are also non-manipulatable elements: four walls, a floor, a door, a window, a fireplace, and a bookshelf (Figure 14(a)). To reflect the complexity of manipulating VR objects, we designed five sub-tasks that involve multiple spatial parameters:

(1) Swap the positions of two pictures (Figure 14(b)).
(2) Move the table under the window while keeping the lamp's relative position on it (Figure 14(c)).
(3) Adjust the chair's rotation and size to fit beneath the table (Figure 14(c)).
(4) Reposition the flower to the fireplace with its color matching that of the central picture (Figure 14(b)).
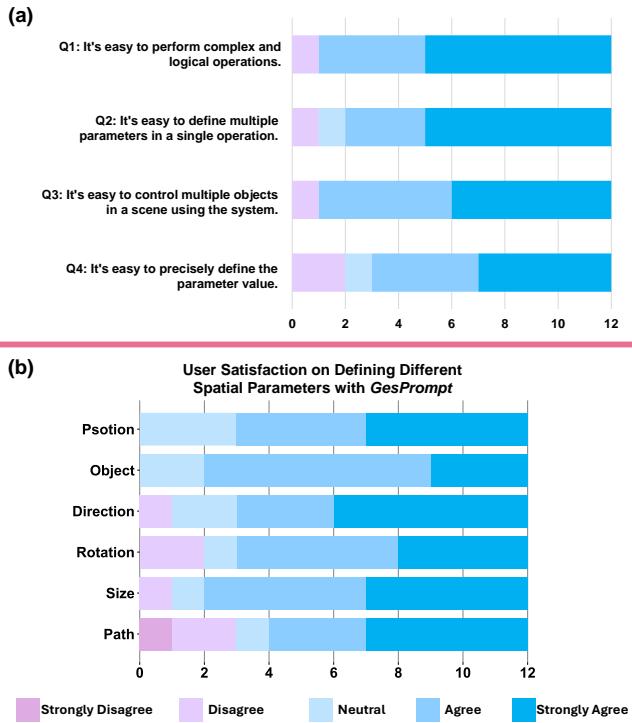(5) Let the lamp move left-and-right on table (Figure 14(c)).

Participants received instructions for these tasks via a picture shown in Figure 14(b) and (c). The participant was instructed to note the distinctions among the pictures, flowers, tables, lamp, and chair, and subsequently complete the assigned task.



**(a)**



**(b)**          **(c)**

**Figure 14: User study session 3 setup: (a) The initial room setup; (b) The target configuration of three pictures and the flower; (c) The target configuration of the table, lamp, and flower.**

*4.4.2 Results and Discussion.* All 12 users completed the sub-tasks and organized the room to the target state. The overall Likert-type results collected from this session are shown in Figure 15(a). Overall, the users find it easy to perform complex and logical manipulations using GesPrompt (Q1: AVG = 4.33, SD = 0.25). *"If I want to swap the locations of two pictures, I don't want to move the pictures directly with my hands. It is much simpler to instruct your system to handle that for me. (P6)"*. In instances where several parameters require definition, users expressed satisfaction with our system (Q2: AVG = 4.33, SD = 0.27). *"When adjusting the chair, it was easy to change my hand gesture from indicating size to indicating rotation, which saved a lot of time. (P11)"*. Users feel it is straightforward to manage multiple objects simultaneously with our system (Q3: AVG = 4.42, SD = 0.25). *"I felt natural to refer to the table and the lamp as 'them' while pointing to them, and the system gets what I meant by 'them'. (P7)"*. The user also finds it easy to precisely define the parameter values (Q4: AVG = 4.00, SD = 0.31). *"I could not describe the exact size of the chair with only my speech. It is easy to define it by gesturing. (P3)"*.

Additionally, the user was asked to rank their satisfaction with defining each spatial parameter using GesPrompt. The results are shown in Figure 15(b). With a high satisfaction rating, GesPrompt was well-received with average scores of 4.17 (SD = 0.23) for position and 4.08 (SD = 0.18) for object parameters. Participants were confident that the specified object and position parameters aligned with their expectations. However, one participant remarked *"When I command the system to move the table, I expect that the rotation of the table will also be changed to facing inside the room. (P11)"*. At the

**(a)**

Q1: It's easy to perform complex and logical operations.

Q2: It's easy to define multiple parameters in a single operation.

Q3: It's easy to control multiple objects in a scene using the system.

Q4: It's easy to precisely define the parameter value.

0  2  4  6  8  10  12

**(b)**

**User Satisfaction on Defining Different Spatial Parameters with *GesPrompt***

Psotion

Object

Direction

Rotation

Size

Path

0  2  4  6  8  10  12

■ Strongly Disagree    ■ Disagree    ■ Neutral    ■ Agree    ■ Strongly Agree

**Figure 15: Session 3 Likert-type questionnaire results. (a) Qualitative Results; (b) User satisfaction for defining spatial parameters using GesPrompt, all p-values are larger than 0.05.**

same time, another participant said *"I want the system to also move the lamp with the table, even if I only selected the table. (P12)"*. These expectations highlight the need to incorporate the spatial context understanding ability in future versions of the system. Potential approaches are: (i) pre-evaluating whether a function call would lead to a spatially sensible environmental state, and (ii) defining the spatial relationship between objects based on the names and positions of the objects with the help of GPT or allowing manual user definition. Participants expressed considerable satisfaction with both direction (AVG = 4.17, SD = 0.28) and rotation (AVG = 3.92, SD = 0.25) parameters. A majority of the 9 participants set the direction of the table, with the others opting for the two-hand rotation. For chair rotation, four opted for one-hand rotation, three chose two-hand rotation, and the remainder used deictic gestures to set the direction. Additionally, 4 participants attempted to directly grab the chair while manipulative gestures are not enabled in the system. The initial position, rotation, and size of the object are possible factors of which gesture is chosen to perform the rotation task. Further study is required to evaluate those factors. For the size parameter (AVG = 4.17, SD = 0.26), all participants used the two-hand size gesture while referencing the length of the hole under the table to resize the chair. Possible factors for choosing the size gesture include the value of the target size, the position, and the size of the reference. For the path parameter (AVG = 3.75, SD = 0.39), although most users are satisfied with the result, it ranks

the lowest among the six spatial parameters. When defining the movement of the lamp on the table, the majority of participants (10) demonstrated the left-to-right and right-to-left motion more than 2 times in a single command. Lacking path editing functions, the prototype system requires participants to redefine the path parameter for any edits. The participant's behaviors suggested that further research is needed to investigate how users would define the animation of objects with speech and gesture.

## 5 LIMITATIONS AND FUTURE WORK

### 5.1 Concurrent Feedback

Although the users are mostly satisfied with our system, some users mention *"The system should show me that it understands my command when I am giving the command (P4)."* This emphasizes the importance of integrating concurrent feedback, meaning the system will interpret words as they are spoken, indicating the system's status when necessary. The feedback could be highlighting the objects the user is currently talking about, changing the objects' size according to the user's hand motion, etc. Currently, the system waits until the users complete a full sentence or command before processing the request and giving feedback by changing the state of the environment objects or asking the user for more information. GesPrompt's workflow ensures that the command is analyzed in full context from both speech and gesture input, further eliminating ambiguities in text. However, this workflow will create a slight delay between interaction and feedback, especially when the user tries to complete multiple tasks in a single command. The delay will lead to uncertainties and frustrations on the user side. With timely feedback, users could correct the system or clarify themselves as soon as feedback arrives, making the whole process more efficient. Future improvement in GesPrompt could be adding real-time processing capabilities similar to GitHub Copilot's [37] code auto-complete feature.

### 5.2 Time Disparity between Gesture and Speech

The gesture segmentation and parameter extraction method was proven to be accurate in most of the commands in the user study. It was theorized on the basis that the co-speech gesture will appear approximately when *parameter token* is spoken [42, 98]. GesPrompt also adds padding before and after the when *parameter token* are spoken to ensure that the entire co-speech gesture is segmented. However, two specific edge cases were observed. In the interview, some users said that *"sometimes I need to think about what kind of gesture I should do after giving the command to the system (P3)"*. This situation results in the time disparity between *parameter tokens* and the segmentation window of the co-speech gesture. Thus, the gesture data segmentation method is not always reliable, especially with a novice user. At the same time, some users mention that *"(when manipulating the cube)...for me, I would just say '[point to the cube] go there [point to a position]' or 'You [point to the cube] move [point to a position]'."* With a shortened command, i.e., an elliptical sentence (ellipsis), the basic elements of a complete sentence could be missing; there would not be any *parameter tokens* to segment and process some of the gesture data. GesPrompt prototype system handles these two edge cases by asking the user to clarify the parameters that the system cannot determine. While this approach

ensures that the parameter value is what the user intended, the extra interaction steps may frustrate some users.

In future work, incorporating a tailored gesture recognition model [10, 72] could further improve the robustness of the system. The gesture recognition model can serve as additional gesture segmentation cues. Together with the *parameter tokens*, they mutually determine the meaningful segments of the gesture data that are relevant to the parameters. Moreover, a calibration module can be integrated into GesPrompt to understand users' habits of merging speech and gestures. A neural network model could be designed to capture how users typically integrate their speech with gestures.

## 5.3 Richer Multimodal Inputs: Complex Gestures, Gaze, & More

While GesPrompt successfully handles both static (position, object, direction, size) and dynamic (path, rotation) co-speech gestures in our proof-of-concept prototype, the space of more complex, composite gestures remains underexplored. In our current implementation, the path parameter is rendered via a simple draw_path function that uses only index-finger positions—omitting richer cues such as hand orientation, multi-finger trajectories, or multi-phase motion. Extending the XR function library—for example, to jointly utilize simultaneous translation and rotation, staged gesture sequences, or parameterized curved trajectories—would enable far more nuanced interactions. Future work could test and investigate specialized XR functions that can use these parts of the path parameter. Potential applications include: Robotics teleoperation [84], where operators could guide both the path and orientation of robotic arms through continuous path gestures. Sports training systems [63], where complex gesture patterns could represent specific motion trajectories (e.g., coaching swings, throws, or strokes). Medical training simulations [106], where fine-grained gesture combinations are needed to mimic surgical tool manipulations or multi-step procedures. AR 3D modeling systems [31], where the co-speech gesture can be used to draw a quick sketch, denote surface normals, and specify solid feature extruding directions. These extensions present both a challenging research agenda and a promising avenue for richer, more expressive multimodal XR interaction.

GesPrompt focus on hand gestures due to their expressiveness and popularity in the XR domain [39, 86]. Besides the hand, other types of human input, such as gaze [81], brain activity [34], and more [115], can also enrich interaction. For example, users may find it more intuitive and convenient to use gaze to indicate the object they are talking about when they are holding something heavy with their hands. Thus, adding other input modalities would increase the usability of the system. The statement is consistent with [65].

Other input modalities can be added to GesPrompt in addition to the gesture processor. The processors for multi-input could run in parallel and make their own guesses of the values for $P_i^{amb}$ based on the corresponding inputs. Then, a Bayesian network model could be applied to determine the most probable value of $P_i^{amb}$.

## 5.4 Adapting GesPrompt to Other Platforms

Although the GesPrompt prototype is implemented for virtual reality evaluation, the framework is adaptable to augmented reality, desktops.

HMD AR and VR have numerous similarities. Based on the reality–virtuality continuum [71], AR is distinct because it incorporates real-world entities. The real-world entities can be added to the XR system by scanning to produce a digital twin of the real-world environment. Once the digital twin is formed with object labels, meshes, position, and rotation, GesPrompt can be adapted to HMD AR. Prior work GazePointAR [56] shows that object parameters can be determined using eye gaze, hand pointing, and 2D vision models. With the help of GesPrompt framework, additional spatial parameters and interactions can be facilitated.

On the desktop platform, the gesture data can be acquired by adding peripheral sensors like Leap Motion [2]. The scene information should include the position of the computer screen and the contents on the screen. With modified metaprompt and XR system functions, GesPrompt has the potential to enable users to interact with speech and gesture for software such as 3D modeling and graphic design.

## 5.5 Context Awareness

In our implementation, we assign a meaningful name to each object in the Unity scene (e.g., "Starry Night" painting), consistent with common scene-building practice. These names supply the textual context that the system needs to identify target objects. While assigning names to objects is optional and varies by individual habit, this step provides an ad hoc solution for the LLM to identify each object. Without proper naming, the scene information may lack critical metadata, potentially causing system failures. Future works could explore the use of visual language models (VLM) [59, 100, 103] to automatically assign descriptive names for objects during the scene construction or in real-time.

While the scene information that contains the name, position, rotation, and size is sent to the LLM system, the implicit spatial relationships are not taken into account. Examples of the spatial relationship include "the lamp should always be on the table", "the front of the table should have space for people", "the painting on the wall should be perpendicular to the ground", etc. User study observations indicate that users expect the system to comprehend these spatial relationships, whereas the GPT agent fails to deduce them from the objects' positions and names. In order to incorporate the spatial relationships of objects in the future iterations of GesPrompt, a validation subsystem can be added to detect the spatial relationships. If the function call produces a spatially invalid state of the environment (i.e., spatial relationship constraint not met), the validation subsystem will notify the user for clarification. Alternatively, the subsystem can be configured to determine the nearest spatial parameters that yield a valid state based on the user's input.

The user's perspective is another element of spatial context. Orientation terms such as "left", "right", "front", and "back" can represent various world directions based on the user's and object's positioning and rotation. Currently, GesPrompt needs the user to

clarify orientation through gestures. To enable the system to interpret these terms, relative direction vectors can be computed, allowing the system to use these vectors when necessary. Furthermore, users might want distant objects to be resized from their perspective, which means that the value of the size parameter acquired from the gesture processor should be proportionally scaled based on the position of the user and the object. To implement perspective resizing, an additional XR system function can be introduced, applying a threshold on distance to choose the appropriate resize function.

## 5.6 Future Applications

By adding other existing XR applications to the function component of GesPrompt XR system, the speech + co-speech gesture interaction can be merged into those applications with proper modifications.

*5.6.1 Interactive XR Content Creation.* The need for XR content creation has been demonstrated in previous works with the content being object [97], UI [32, 40, 41], tutorial [89], and scene [14, 85]. With GesPrompt, the user can easily define spatial parameters without mentally measuring and translating, which means that they can simply use gestures to indicate the value of the parameters, thus reducing the cognitive load. For example, the user can say *"I want to create an object with a shape like this."* while tracing the outline of the imaginary object with their hands; say *"Turn this part red"* while touching part of the object; say *"It can be opened like this."* while demonstrating how it should be opened using gestures.

*5.6.2 Communicating with XR Agent.* User describing their needs to the XR agents has been a challenge when referring to the spatial entities [30]. Our system allows the user to demonstrate spatial and physical concepts by enabling the use of co-speech gestures. For example, if the user is having trouble assembling a virtual mechanical component, they can not think of a way to describe the direction of the component. With the help of our system, they can now ask the agent *"should I insert the component this way or that way?"* while presenting two different directions with a gesture.

*5.6.3 Interacting with XR Environment.* Previous work [46, 88, 101, 102, 105] has explored creating personalized interactions in XR environments. In addition to embodied demonstrations, these approaches often require visual programming to manually connect triggers and actions. GesPrompt allows users to communicate their abstract programming intentions through speech while still demonstrating with gestures, thus reducing the physical effort of "connecting lines" between triggers and actions. In a situation where the users try to program an interaction in which the user uses gesture and voice to trigger an action of an object, the lengthy authoring process could be shortened with our system. The user can simply say, *"When I say 'wingardium leviosa' and do this gesture, the feather will flow with my hand like this."* to create the desired interaction.

## 6 CONCLUSION

In this work, we present GesPrompt, an LLM-powered multimodal XR interface that allows users to express their intention with both speech and gesture. With our system, the user can present spatial-temporal information to the system via co-speech gesture, alleviating the burden of creating detailed prompts that describe the spatial-temporal information. We first discuss the need and method of incorporating the co-speech gestures into the communication with the copilot. Then we present the overall system workflow and a walk-through of the copilot system. GesPrompt consists of two main components: the *LLM system* and the *Gesture Processor*. The *LLM system* directly processes the user's speech, outputs function calls and part of the parameter values to the XR system, and cues for analyzing gesture to the *gesture processor*. The *gesture processor* uses the cues to segment and extract spatial-temporal information from the co-speech gesture, sending the remaining parameter values to the XR system. Finally, the XR system executes the function call, changing the environment to the user-intended configuration. Next, we illustrate the spatial-temporal parameters and their relations to the co-speech gesture. The details of each component in GesPrompt are then explained. We implemented this workflow to create a prototype VR system that is capable of helping the user manipulate objects in the virtual environment. In addition, to explore the feasibility of the proposed workflow and user behavior with GesPrompt , we conducted a three-session user study. The results indicate that GesPrompt significantly improves user satisfaction, decreases cognitive load, and improves the overall experience. Through the user study, we also identified some of the limitations of the current implementation and possible solutions for future studies. In summary, we believe that GesPrompt presents a novel approach to combining gesture and voice input in the HCI area, opening up new possibilities to integrate co-speech gestures into LLM-based systems and inspiring the next generation of intelligent, user-centered XR environments.

## ACKNOWLEDGMENT

## REFERENCES

[1] 2024. Hololens 2. https://www.microsoft.com/en-us/hololens.
[2] 2024. Leap Motion. https://www.ultraleap.com/.
[3] 2024. Meta Quest Pro. https://www.meta.com/quest/quest-pro/.
[4] 2024. Microsoft Azure Speech to Text. https://azure.microsoft.com/en-us/products/ai-services/ai-speech.
[5] 2024. Microsoft Planner. https://tasks.office.com/. Accessed: 2024-09-11.
[6] 2024. Unity Engine. https://unity.com/.
[7] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
[8] Setareh Aghel Manesh, Tianyi Zhang, Yuki Onishi, Kotaro Hara, Scott Bateman, Jiannan Li, and Anthony Tang. 2024. How People Prompt Generative AI to Create Interactive VR Scenes. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*. 2319–2340.
[9] Yuya Aikawa, Ryoya Tamura, Chunchen Xu, Xiao Ge, and Daigo Misaki. 2023. Introducing Augmented Post-it: An AR Prototype for Engaging Body Movements in Online GPT-Supported Brainstorming. In *Adjunct Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–3.

[10] Muneer Al-Hammadi, Ghulam Muhammad, Wadood Abdul, Mansour Alsulaiman, Mohammed A Bencherif, Tareq S Alrayes, Hassan Mathkour, and Mohamed Amine Mekhtiche. 2020. Deep learning-based approach for sign language gesture recognition with efficient hand gesture representation. *Ieee Access* 8 (2020), 192527–192542.

[11] Martha W Alibali. 2005. Gesture in spatial cognition: Expressing, communicating, and thinking about spatial information. *Spatial cognition and computation* 5, 4 (2005), 307–331.

[12] Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. 2022. Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM Transactions on Graphics (TOG)* 41, 6 (2022), 1–19.

[13] Rahul Arora, Rubaiat Habib Kazi, Danny M Kaufman, Wilmot Li, and Karan Singh. 2019. Magicalhands: Mid-air hand gestures for animating in vr. In *Proceedings of the 32nd annual ACM symposium on user interface software and technology*. 463–477.

[14] Mark Billinghurst, Sisinio Baldis, Lydia Matheson, and Mark Philips. 1997. 3D palette: a virtual reality content creation tool. In *Proceedings of the ACM symposium on Virtual reality software and technology*. 155–156.

[15] Richard A Bolt. 1980. "Put-that-there" Voice and gesture at the graphics interface. In *Proceedings of the 7th annual conference on Computer graphics and interactive techniques*. 262–270.

[16] Efe Bozkir, Süleyman Özdel, Ka Hei Carrie Lau, Mengdi Wang, Hong Gao, and Enkelejda Kasneci. 2024. Embedding large language models into extended reality: Opportunities and challenges for inclusion, engagement, and privacy. In *Proceedings of the 6th ACM Conference on Conversational User Interfaces*. 1–7.

[17] Sara C Broaders and Susan Goldin-Meadow. 2010. Truth is at hand: How gesture adds information during investigative interviews. *Psychological Science* 21, 5 (2010), 623–628.

[18] Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).

[19] Volkert Buchmann, Stephen Violich, Mark Billinghurst, and Andy Cockburn. 2004. FingARtips: gesture based direct manipulation in Augmented Reality. In *Proceedings of the 2nd international conference on Computer graphics and interactive techniques in Australasia and South East Asia*. 212–221.

[20] Judee K Burgoon, Thomas Birk, and Michael Pfau. 1990. Nonverbal behaviors, persuasion, and credibility. *Human communication research* 17, 1 (1990), 140–169.

[21] Han Joo Chae, Jeong-in Hwang, and Jinwook Seo. 2018. Wall-based space manipulation technique for efficient placement of distant objects in augmented reality. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*. 45–52.

[22] Jiawen Chen, Shahram Izadi, and Andrew Fitzgibbon. 2012. KinÊtre: animating the world with the human body. In *Proceedings of the 25th annual ACM symposium on User interface software and technology*. 435–444.

[23] Yang Chen, Yingwei Pan, Haibo Yang, Ting Yao, and Tao Mei. 2024. Vp3d: Unleashing 2d visual prompt for text-to-3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4896–4905.

[24] Zhaorui Chen, Jinzhu Li, Yifan Hua, Rui Shen, and Anup Basu. 2017. Multimodal interaction in augmented reality. In *2017 IEEE international conference on systems, man, and cybernetics (SMC)*. IEEE, 206–209.

[25] Vuthea Chheang, Shayla Sharmin, Rommy Márquez-Hernández, Megha Patel, Danush Rajasekaran, Gavin Caulfield, Behdokht Kiafar, Jicheng Li, Pinar Kullu, and Roghayeh Leila Barmaki. 2024. Towards anatomy education with generative AI-based virtual assistants in immersive virtual reality environments. In *2024 IEEE International Conference on Artificial Intelligence and eXtended and Virtual Reality (AIxVR)*. IEEE, 21–30.

[26] Sharice Clough and Melissa C Duff. 2020. The role of gesture in communication and cognition: Implications for understanding and treating neurogenic communication disorders. *Frontiers in Human Neuroscience* 14 (2020), 323.

[27] Hélène Cochet and Jacques Vauclair. 2014. Deictic gestures and symbolic gestures produced by adults in an experimental context: Hand shapes and hand preferences. *Laterality: Asymmetries of Body, Brain and Cognition* 19, 3 (2014), 278–301.

[28] Fernanda De La Torre, Cathy Mengying Fang, Han Huang, Andrzej Banburski-Fahey, Judith Amores Fernandez, and Jaron Lanier. 2024. Llmr: Real-time prompting of interactive worlds using large language models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–22.

[29] Holger Diessel and Kenny R Coventry. 2020. Demonstratives in spatial language and social interaction: An interdisciplinary review. *Frontiers in Psychology* 11 (2020), 555265.

[30] Mustafa Doga Dogan, Eric J. Gonzalez, Karan Ahuja, Ruofei Du, Andrea Colaço, Johnny Lee, Mar Gonzalez-Franco, and David Kim. 2024. Augmented Object Intelligence with XR-Objects. https://doi.org/10.1145/3654777.3676379 arXiv:2404.13274 [cs.HC]

[31] Runlin Duan, Xiyun Hu, Min Liu, Jingyu Shi, and Karthik Ramani. 2025. pARametric: Empowering In Situ Parametric Modeling in Augment Reality for Personal Fabrication. *Journal of Computing and Information Science in Engineering* 25, 4 (2025), 041001.

[32] João Marcelo Evangelista Belo, Mathias N Lystbæk, Anna Maria Feit, Ken Pfeuffer, Peter Kán, Antti Oulasvirta, and Kaj Grønbæk. 2022. Auit–the adaptive user interfaces toolkit for designing xr applications. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*. 1–16.

[33] Markus Friedrich, Stefan Langer, and Fabian Frey. 2021. Combining Gesture and Voice Control for Mid-air Manipulation of CAD Models in VR Environments.. In *VISIGRAPP (2: HUCAPP)*. 119–127.

[34] Adélaïde Genay, Anatole Lécuyer, and Martin Hachet. 2021. Being an avatar "for real": a survey on virtual embodiment in augmented reality. *IEEE Transactions on Visualization and Computer Graphics* 28, 12 (2021), 5071–5090.

[35] Esam Ghaleb, Ilya Burenko, Marlou Rasenberg, Wim Pouw, Ivan Toni, Peter Uhrig, Anna Wilson, Judith Holler, Aslı Özyürek, and Raquel Fernández. 2024. Leveraging Speech for Gesture Detection in Multimodal Communication. *arXiv preprint arXiv:2404.14952* (2024).

[36] Esam Ghaleb, Ilya Burenko, Marlou Rasenberg, Wim Pouw, Peter Uhrig, Judith Holler, Ivan Toni, Aslı Özyürek, and Raquel Fernández. 2024. Co-Speech Gesture Detection through Multi-Phase Sequence Labeling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 4007–4015.

[37] GitHub. 2023. GitHub Copilot. https://github.com/features/copilot Accessed: 2024-09-11.

[38] Daniele Giunchi, Nels Numan, Elia Gatti, and Anthony Steed. 2024. DreamCodeVR: Towards Democratizing Behavior Design in Virtual Reality with Speech-Driven Programming. In *2024 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*. IEEE, 579–589.

[39] Lin Guo, Zongxing Lu, and Ligang Yao. 2021. Human-machine interaction sensing technology based on hand gesture recognition: A review. *IEEE Transactions on Human-Machine Systems* 51, 4 (2021), 300–309.

[40] Fengming He, Xiyun Hu, Xun Qian, Zhengzhe Zhu, and Karthik Ramani. 2024. AdapTUI: Adaptation of Geometric-Feature-Based Tangible User Interfaces in Augmented Reality. *Proceedings of the ACM on Human-Computer Interaction* 8, ISS (2024), 44–69.

[41] Fengming He, Xiyun Hu, Jingyu Shi, Xun Qian, Tianyi Wang, and Karthik Ramani. 2023. Ubi Edge: Authoring Edge-Based Opportunistic Tangible User Interfaces in Augmented Reality. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–14.

[42] Henning Holle, Thomas C Gunter, Shirley-Ann Rüschemeyer, Andreas Hennenlotter, and Marco Iacoboni. 2008. Neural correlates of the processing of co-speech gestures. *NeuroImage* 39, 4 (2008), 2010–2024.

[43] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 2023. 3d-llm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems* 36 (2023), 20482–20494.

[44] Zhanpeng Huang, Weikai Li, and Pan Hui. 2015. Ubii: Towards seamless interaction between digital and physical worlds. In *Proceedings of the 23rd ACM international conference on Multimedia*. 341–350.

[45] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 867–876.

[46] Rahul Jain, Jingyu Shi, Runlin Duan, Zhengzhe Zhu, Xun Qian, and Karthik Ramani. 2023. Ubi-TOUCH: Ubiquitous Tangible Object Utilization through Consistent Hand-object interaction in Augmented Reality. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco, CA, USA) *(UIST '23)*. Association for Computing Machinery, New York, NY, USA, Article 12, 18 pages. https://doi.org/10.1145/3586183.3606793

[47] Spencer D Kelly and Quang-Anh Ngo Tran. 2023. Exploring the Emotional Functions of Co-Speech Hand Gesture in Language and Communication. *Topics in Cognitive Science* (2023).

[48] Tae Soo Kim, DaEun Choi, Yoonseo Choi, and Juho Kim. 2022. Stylette: Styling the web with natural language. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–17.

[49] Yongkwan Kim and Seok-Hyung Bae. 2016. SketchingWithHands: 3D sketching handheld products with first-person hand posture. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*. 797–808.

[50] Sotaro Kita and Asli Özyürek. 2003. What does cross-linguistic variation in semantic coordination of speech and gesture reveal?: Evidence for an interface representation of spatial thinking and speaking. *Journal of Memory and language* 48, 1 (2003), 16–32.

[51] Sotaro Kita, Ingeborg Van Gijn, and Harry Van der Hulst. 1998. Movement phases in signs and co-speech gestures, and their transcription by human coders. In *Gesture and Sign Language in Human-Computer Interaction: International Gesture Workshop Bielefeld, Germany, September 17–19, 1997 Proceedings*. Springer, 23–35.

[52] Mikhail Konenkov, Artem Lykov, Daria Trinitatova, and Dzmitry Tsetserukou. 2024. Vr-gpt: Visual language model for intelligent virtual reality applications.

arXiv preprint arXiv:2405.11537 (2024).

[53] Okan Köpüklü, Ahmet Gunduz, Neslihan Kose, and Gerhard Rigoll. 2019. Real-time hand gesture detection and classification using convolutional neural networks. In *2019 14th IEEE international conference on automatic face & gesture recognition (FG 2019)*. IEEE, 1–8.

[54] Ryutaro Kurai, Takefumi Hiraki, Yuichi Hiroi, Yutaro Hirao, Monica Perusquia-Hernandez, Hideaki Uchiyama, and Kiyoshi Kiyokawa. 2024. MagicItem: Dynamic Behavior Design of Virtual Objects with Large Language Models in a Consumer Metaverse Platform. *arXiv preprint arXiv:2406.13242* (2024).

[55] David Ledo, Steven Houben, Jo Vermeulen, Nicolai Marquardt, Lora Oehlberg, and Saul Greenberg. 2018. Evaluation strategies for HCI toolkit research. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–17.

[56] Jaewook Lee, Jun Wang, Elizabeth Brown, Liam Chu, Sebastian S Rodriguez, and Jon E Froehlich. 2024. GazePointAR: A Context-Aware Multimodal Voice Assistant for Pronoun Disambiguation in Wearable Augmented Reality. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–20.

[57] Minkyung Lee, Mark Billinghurst, Woonhyuk Baek, Richard Green, and Woontack Woo. 2013. A usability study of multimodal input in an augmented reality environment. *Virtual Reality* 17 (2013), 293–305.

[58] Minkyung Lee, Richard Green, and Mark Billinghurst. 2008. 3D natural hand interaction for AR applications. In *2008 23rd International Conference Image and Vision Computing New Zealand*. IEEE, 1–6.

[59] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. 2024. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326* (2024).

[60] Fangfu Liu, Hanyang Wang, Weiliang Chen, Haowen Sun, and Yueqi Duan. 2024. Make-Your-3D: Fast and Consistent Subject-Driven 3D Content Generation. *arXiv preprint arXiv:2403.09625* (2024).

[61] Ziyi Liu, Zhengzhe Zhu, Lijun Zhu, Enze Jiang, Xiyun Hu, Kylie A Peppler, and Karthik Ramani. 2024. ClassMeta: Designing Interactive Virtual Classmate to Promote VR Classroom Participation. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–17.

[62] Hao Lü and Yang Li. 2012. Gesture coder: a tool for programming multi-touch gestures by demonstration. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2875–2884.

[63] Dizhi Ma, Xiyun Hu, Jingyu Shi, Mayank Patel, Rahul Jain, Ziyi Liu, Zhengzhe Zhu, and Karthik Ramani. 2024. avaTTAR: Table Tennis Stroke Training with Embodied and Detached Visualization in Augmented Reality. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. 1–16.

[64] Dr M Madhiarasan, Prof Roy, and Partha Pratim. 2022. A comprehensive review of sign language recognition: Different types, modalities, and datasets. *arXiv preprint arXiv:2204.03328* (2022).

[65] Daniel Martin, Sandra Malpica, Diego Gutierrez, Belen Masia, and Ana Serrano. 2022. Multimodality in VR: A survey. *ACM Computing Surveys (CSUR)* 54, 10s (2022), 1–36.

[66] Ingrid Masson-Carro, Martijn Goudbeek, and Emiel Krahmer. 2017. How what we see and what we know influence iconic gesture production. *Journal of nonverbal behavior* 41 (2017), 367–394.

[67] Daniel Mendes, Fernando Fonseca, Bruno Araujo, Alfredo Ferreira, and Joaquim Jorge. 2014. Mid-air interactions above stereoscopic interactive tables. In *2014 IEEE Symposium on 3D User Interfaces (3DUI)*. IEEE, 3–10.

[68] Daniel Mendes, Filipe Relvas, Alfredo Ferreira, and Joaquim Jorge. 2016. The benefits of dof separation in mid-air 3d object manipulation. In *Proceedings of the 22nd ACM conference on virtual reality software and technology*. 261–268.

[69] Meta XR SDK 2024. https://developer.oculus.com/downloads/package/meta-xr-sdk-all-in-one-upm/.

[70] Microsoft. 2024. Mixed Reality Toolkit (MRTK). https://github.com/microsoft/MixedRealityToolkit-Unity. Accessed: 2024-12-11.

[71] Paul Milgram, Haruo Takemura, Akira Utsumi, and Fumio Kishino. 1995. Augmented reality: A class of displays on the reality-virtuality continuum. In *Telemanipulator and telepresence technologies*, Vol. 2351. Spie, 282–292.

[72] Abdullah Mujahid, Mazhar Javed Awan, Awais Yasin, Mazin Abed Mohammed, Robertas Damaševičius, Rytis Maskeliūnas, and Karrar Hameed Abdulkareem. 2021. Real-time hand gesture recognition based on deep learning YOLOv3 model. *Applied Sciences* 11, 9 (2021), 4164.

[73] Hiroaki Nishino, Kouichi Utsumiya, and Kazuyoshi Korida. 1998. 3d object modeling using spatial and pictographic gestures. In *Proceedings of the ACM symposium on Virtual reality software and technology*. 51–58.

[74] Jeongseok Oh, Seungju Kim, and Seungjun Kim. 2024. LumiMood: A Creativity Support Tool for Designing the Mood of a 3D Scene. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–21.

[75] SK Ong and ZB Wang. 2011. Augmented assembly technologies based on 3D bare-hand interaction. *CIRP annals* 60, 1 (2011), 1–4.

[76] OpenAI. 2023. ChatGPT: GPT-4, Large Language Model. https://chat.openai.com/ Accessed: 2024-09-11.

[77] OpenAI GPT-4o 2024. =https://openai.com/index/hello-gpt-4o/.

[78] Xiaohan Peng, Janin Koch, and Wendy E Mackay. 2024. DesignPrompt: Using Multimodal Interaction for Design Exploration with Generative AI. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*. 804–818.

[79] Raedy Ping and Susan Goldin-Meadow. 2010. Gesturing saves cognitive resources when talking about nonpresent objects. *Cognitive Science* 34, 4 (2010), 602–619.

[80] Thammathip Piumsomboon, David Altimira, Hyungon Kim, Adrian Clark, Gun Lee, and Mark Billinghurst. 2014. Grasp-Shell vs gesture-speech: A comparison of direct and indirect natural interaction techniques in augmented reality. In *2014 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 73–82.

[81] Thammathip Piumsomboon, Gun Lee, Robert W Lindeman, and Mark Billinghurst. 2017. Exploring natural eye-gaze-based interaction for immersive virtual reality. In *2017 IEEE symposium on 3D user interfaces (3DUI)*. IEEE, 36–39.

[82] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).

[83] Sara Price, Carey Jewitt, and Nikoleta Yiannoutsou. 2021. Conceptualising touch in VR. *Virtual Reality* 25, 3 (2021), 863–877.

[84] Wen Qi, Salih Ertug Ovur, Zhijun Li, Aldo Marzullo, and Rong Song. 2021. Multi-sensor guided hand gesture recognition for a teleoperated robot using a recurrent neural network. *IEEE Robotics and Automation Letters* 6, 3 (2021), 6039–6045.

[85] Xun Qian, Fengming He, Xiyun Hu, Tianyi Wang, Ananya Ipsita, and Karthik Ramani. 2022. Scalar: Authoring semantically adaptive augmented reality experiences in virtual reality. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–18.

[86] Siddharth S Rautaray and Anupam Agrawal. 2015. Vision based hand gesture recognition for human computer interaction: a survey. *Artificial intelligence review* 43 (2015), 1–54.

[87] Raúl A Sánchez-Ancajima, Sarajane Marques Peres, Javier A López-Céspedes, José L Saly-Rosas-solano, Ronald M Hernández, and Miguel A Saavedra-López. [n. d.]. Gesture Phase Segmentation Dataset: An Extension for Development of Gesture Analysis Models. ([n. d.]).

[88] Anika Sayara, Emily Lynn Chen, Cuong Nguyen, Robert Xiao, and Dongwook Yoon. 2023. Gesturecanvas: A programming by demonstration system for prototyping compound freehand interaction in vr. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–17.

[89] Jingyu Shi, Rahul Jain, Seungguen Chi, Hyungjun Doh, Hyunggun Chi, Alexander J Quinn, and Karthik Ramani. 2025. CARING-AI: Towards Authoring Context-aware Augmented Reality INstruction through Generative Artificial Intelligence. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. https://doi.org/10.1145/3706598.3713348

[90] Jingyu Shi, Rahul Jain, Hyungjun Doh, Ryo Suzuki, and Karthik Ramani. 2023. An HCI-centric survey and taxonomy of human-generative-AI interactions. *arXiv preprint arXiv:2310.07127* (2023).

[91] Timo Sowa and Ipke Wachsmuth. 2001. Interpretation of shape-related iconic gestures in virtual environments. In *International Gesture Workshop*. Springer, 21–33.

[92] Benjamin Straube, Antonia Green, Bianca Bromberger, and Tilo Kircher. 2011. The differentiation of iconic and metaphoric gestures: Common and unique integration processes. *Human brain mapping* 32, 4 (2011), 520–533.

[93] Hendrik Strobelt, Albert Webson, Victor Sanh, Benjamin Hoover, Johanna Beyer, Hanspeter Pfister, and Alexander M Rush. 2022. Interactive and visual prompt engineering for ad-hoc task adaptation with large language models. *IEEE transactions on visualization and computer graphics* 29, 1 (2022), 1146–1156.

[94] Michael Studdert-Kennedy. 1994. Hand and Mind: What Gestures Reveal About Thought. *Language and Speech* 37, 2 (1994), 203–209.

[95] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).

[96] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[97] Fernando Vera and J Alfredo Sánchez. 2016. A model for in-situ augmented reality content creation based on storytelling and gamification. In *Proceedings of the 6th Mexican Conference on Human-Computer Interaction*. 39–42.

[98] Petra Wagner, Zofia Malisz, and Stefan Kopp. 2014. Gesture and speech in interaction: An overview. , 209–232 pages.

[99] Hongyu Wan, Jinda Zhang, Abdulaziz Arif Suria, Bingsheng Yao, Dakuo Wang, Yvonne Coady, and Mirjana Prpa. 2024. Building LLM-based AI Agents in Social Virtual Reality. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–7.

[100] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191* (2024).

[101] Tianyi Wang, Xun Qian, Fengming He, Xiyun Hu, Yuanzhi Cao, and Karthik Ramani. 2021. Gesturar: An authoring system for creating freehand interactive augmented reality applications. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. 552–567.

[102] Tianyi Wang, Xun Qian, Fengming He, Xiyun Hu, Ke Huo, Yuanzhi Cao, and Karthik Ramani. 2020. CAPturAR: An augmented reality tool for authoring human-involved context-aware applications. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. 328–341.

[103] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 2023. CogVLM: Visual Expert for Pretrained Language Models. (2023). arXiv:2311.03079 [cs.CV]

[104] Zehan Wang, Haifeng Huang, Yang Zhao, Ziang Zhang, and Zhou Zhao. 2023. Chat-3d: Data-efficiently tuning large language model for universal dialogue of 3d scenes. *arXiv preprint arXiv:2308.08769* (2023).

[105] Zeyu Wang, Cuong Nguyen, Paul Asente, and Julie Dorsey. 2021. Distanciar: Authoring site-specific augmented reality experiences for remote environments. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–12.

[106] Rong Wen, Wei-Liang Tay, Binh P Nguyen, Chin-Boon Chng, and Chee-Kong Chui. 2014. Hand gesture guided robot-assisted surgery based on a direct augmented reality interface. *Computer methods and programs in biomedicine* 116, 2 (2014), 68–80.

[107] Roel M Willems, Aslı Özyürek, and Peter Hagoort. 2007. When language meets action: The neural integration of gesture and speech. *Cerebral Cortex* 17, 10 (2007), 2322–2333.

[108] Liwei Wu, Ben Lafreniere, Tovi Grossman, Thomas White, and Stephanie Santosa. 2024. Body Language for VUIs: Exploring Gestures to Enhance Interactions with Voice User Interfaces. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*. 133–150.

[109] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. In *Proceedings of the 2022 CHI conference on human factors in computing systems*. 1–22.

[110] Ying Choon Wu and Seana Coulson. 2014. Co-speech iconic gestures and visuo-spatial working memory. *Acta psychologica* 153 (2014), 39–50.

[111] Yukang Yan, Chun Yu, Xiaojuan Ma, Xin Yi, Ke Sun, and Yuanchun Shi. 2018. Virtualgrasp: Leveraging experience of interacting with physical objects to facilitate digital object retrieval. In *Proceedings of the 2018 Chi conference on human factors in computing systems*. 1–13.

[112] LI Yang, Jin Huang, TIAN Feng, WANG Hong-An, and DAI Guo-Zhong. 2019. Gesture interaction in virtual reality. *Virtual Reality & Intelligent Hardware* 1, 1 (2019), 84–112.

[113] Hui Ye and Hongbo Fu. 2022. ProGesAR: Mobile AR Prototyping for Proxemic and Gestural Interactions with Real-world IoT Enhanced Spaces. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–14.

[114] Shahrouz Yousefi, Mhretab Kidane, Yeray Delgado, Julio Chana, and Nico Reski. 2016. 3D gesture-based interaction for immersive experience in mobile VR. In *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2121–2126.

[115] Difeng Yu, Qiushi Zhou, Tilman Dingler, Eduardo Velloso, and Jorge Goncalves. 2022. Blending on-body and mid-air interaction in virtual reality. In *2022 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 637–646.

[116] Lei Zhang, Jin Pan, Jacob Gettig, Steve Oney, and Anhong Guo. 2024. VRCopilot: Authoring 3D Layouts with Generative AI Models in VR. *arXiv preprint arXiv:2408.09382* (2024).

[117] Zeyi Zhang, Tenglong Ao, Yuyao Zhang, Qingzhe Gao, Chuan Lin, Baoquan Chen, and Libin Liu. 2024. Semantic Gesticulator: Semantics-Aware Co-Speech Gesture Synthesis. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–17.

[118] Qian Zhou, Sarah Sykes, Sidney Fels, and Kenrick Kin. 2020. Gripmarks: Using hand grips to transform in-hand objects into mixed reality input. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–11.

[119] Jiarui Zhu, Radha Kumaran, Chengyuan Xu, and Tobias Höllerer. 2023. Free-form Conversation with Human and Symbolic Avatars in Mixed Reality. In *2023 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 751–760.

# A   APPENDIX

## A.1   GesPrompt GPT Module Metaprompt

```
### Objective
Use user commands to generate precise JSON outputs that
    dictate functions for a Unity environment. This
    involves identifying manipulated objects, handling
    missing parameters, and outputting functions in a
    structured JSON format. Different functions require
    different parameters.

### Function Types
1. **Select**: Selects objects of certain color from a
    list of objects. Require "objects" and "color".
2. **Move**: Moves a GameObject to a specified position.
    Require "object" and "position".
3. **Rotate_Dir**: Rotates a GameObject to a specified
    orientation base on the direction. Require "object"
    and "direction".
4. **Rotate**: Rotates a GameObject to a specified
    orientation. Require "object" and "rotation".
5. **Resize**: Changes the size of a GameObject. Require
    "object" and "size".
6. **Move_Path**: Defines an animation path for a
    GameObject through multiple positions. Require "
    object" and "path".
7.  **Draw_Path**: Draws a path in the scene. Require "
    path" and "shape_type".
8.  **Set_Color**: Alters the color of a GameObject.
    Require "object" and "color".

### Parameter Types
1. **Position**: A Vector3 value denoting the position.
2. **Object**: A GameObject.
3. **List_Object**: A list of gameobjects.
4. **Direction**: A Vector3 value denoting the direction.
5. **Rotation**: A Quaternion value denoting the rotation
    .
6. **Size**: A Vector3 value denoting the size.
7. **Path**: A list of positions and rotation.
8. **Color**: A Vector 3 value denoteing the RGB of the
    color.
9. **shape_type**: An int value to define shape type, 0
    for line, 1 for circle, 2 for sine wave.
10. **Type**: An int value to define function type.

### Missing Parameters
If one spatial parameter is not clearly defined in the
    command, like "move this object to that place", "
    move the blue object to the right side of the flat
    surface", where "this object", "that place", "the
    blue object" and "right side of the flat surface"
    are missing "object" and "position". In this case,
    you need to define the parameters with descriptions,
    like "this", "that", "blue" and "right side of the
    flat surface".
If the description itself cannot represent the parameter
    alone you need to supplement the description with
    the reference, like "move the object on the left
    side of the table to the right side", the "
    description" of the "parameter" position should be "
    the right side of the table".

### Spatial Parameters
In case the description includes demonstrative pronouns
    like "this", "that", "there", etc., an the word
    index of the first word in the phrase related to the
    missing spatial parameter is needed for the missing
    parameter. Give each object an "id" so if it occurs
    in multiple functions you can use "id" to define it
    . This also works for "position". Its idea is
    independent from "object". Both indices should start
    from 0. If a missing parameter is defined with "id
    ", then you use only "id" to represent it.

**Note:**
```

```
The word is separated by space, so the index of the first
    word related to the missing parameter is needed for
    the missing parameter. The index starts from 1.
    Make sure the index does not exceed the number of
    words in the command. The description of the
    parameter is the phrase related to the missing
    parameter. The description must be exactly the same
    as the phrase in the original command related to the
    missing parameter.

### Process:
The scene information accompanied should also be put into
    consideration.
1. **Action Identification and Mapping:** Parse the
    command to extract functions and match them with the
    given functions.
2. **Multiple functions:** If the command contains
    multiple functions or the command needs to be broken
    down into smaller functions, you need to extract
    them separately.
3. **Parameter Extraction:**
   - Extract necessary parameters for functions from the
        command and the scene information.
   - You need to consider the scene information to
        determine the parameter values.
   - If you cannot find the parameter value in the scene
        information and the user command, you need to
        mark it as missing.
   - Identify missing parameters and use indices to note
        the first words related to missing details.
4. **JSON Output Construction:** Create a structured JSON
    string that details all functions, including
    parameters, missing parameters, and pronoun indices.
5. **JSON Output Format:**
   - "functions": List of functions with "id", "
        functionType", and "parameters".
   - "id": Unique identifier for the function.
   - "functionType": Type of function.
   - "parameters": List of parameters with "id", "
        parameterType", and additional parameters
        specific to the action type.
   - "parameterType": Type of parameter.


### Output Specification
- A JSON formatted string encapsulating the functions and
    parameters, focusing strictly on providing
    actionable data without additional narrative or
    explanations.

- Example JSON Output for Command "Take that object from
    the table, make it twice as big, and place it on the
    right side of the blue box ":
```json
{
    "functions": [
        {
            "id": 0,
            "functionType": 5,
            "parameters": [
                {
                    "id": 0,
                    "description": "that object from the
                        table",
                    "index": 2,
                    "parameterType": 2
                },
                {
                    "id": 1,
                    "parameterType": 6,
                    "size": {
                        "x": 2.0,
                        "y": 2.0,
                        "z": 2.0
                    }
                }
            ]
```

```
        },
        {
            "id": 1,
            "functionType": 2,
            "parameters": [
                {
                    "id": 0,
                    "parameterType": 2
                },
                {
                    "id": 2,
                    "description": "right side of the
                        blue box",
                    "index": 17,
                    "parameterType": 1
                },
            ]
        }
    ]
}
```
**Note:**
"parameters":{[{"id": 0}]} represents the "that object
    from the table" which is defined with "id": 0, so
    you just need to give "id": 0 for the "parameterType
    :2"  parameter in the "move" function.
```

## A.2   Voice-only System GPT Module Metaprompt

```
### Objective
Use user commands to generate precise JSON outputs that
    dictate functions for a Unity environment. This
    involves identifying manipulated objects, handling
    missing parameters, and outputting functions in a
    structured JSON format. Different functions require
    different parameters.

### Function Types
1. **Select**: Selects objects of certain color from a
    list of objects. Require "objects" and "color".
2. **Move**: Moves a GameObject to a specified position.
    Require "object" and "position".
3. **Rotate_Dir**: Rotates a GameObject to a specified
    orientation base on the direction. Require "object"
    and "direction".
4. **Rotate**: Rotates a GameObject to a specified
    orientation. Require "object" and "rotation".
5. **Resize**: Changes the size of a GameObject. Require
    "object" and "size".
6. **Move_Path**: Defines an animation path for a
    GameObject through multiple positions. Require "
    object" and "path".
7.  **Draw_Path**: Draws a path in the scene. Require "
    path" and "shape_type".
8.  **Set_Color**: Alters the color of a GameObject.
    Require "object" and "color".

### Parameter Types
1. **Position**: A Vector3 value denoting the position.
2. **Object**: A GameObject.
3. **List_Object**: A list of gameobjects.
4. **Direction**: A Vector3 value denoting the direction.
5. **Rotation**: A Quaternion value denoting the rotation
    .
6. **Size**: A Vector3 value denoting the size.
7. **Path**: A list of positions and rotation.
8. **Color**: A Vector 3 value denoteing the RGB of the
    color.
9. **shape_type**: An int value to define shape type, 0
    for line, 1 for circle, 2 for sine wave.
10. **Type**: An int value to define function type.

### Missing Parameters
```

If one spatial parameter is not clearly defined in the
    command, like "move this object to that place", "
    move the blue object to the right side of the flat
    surface", where "this object", "that place", "the
    blue object" and "right side of the flat surface"
    are missing "object" and "position". In this case,
    you need to define the parameters with descriptions,
    like "this", "that", "blue" and "right side of the
    flat surface".
If the description itself cannot represent the parameter
    alone you need to supplement the description with
    the reference, like "move the object on the left
    side of the table to the right side", the "
    description" of the "parameter" position should be "
    the right side of the table".

### Spatial Parameters
In case the description includes demonstrative pronouns
    like "this", "that", "there", etc., an the word
    index of the first word in the phrase related to the
    missing spatial parameter is needed for the missing
    parameter. Give each object an "id" so if it occurs
    in multiple functions you can use "id" to define it
    . This also works for "position". Its idea is
    independent from "object". Both indices should start
    from 0. If a missing parameter is defined with "id
    ", then you use only "id" to represent it.

**Note:**
The word is separated by space, so the index of the first
    word related to the missing parameter is needed for
    the missing parameter. The index starts from 1.
    Make sure the index does not exceed the number of
    words in the command. The description of the
    parameter is the phrase related to the missing
    parameter. The description must be exactly the same
    as the phrase in the original command related to the
    missing parameter.

### For the missing parameters use the scene information
    and the user command to deduct the value to your
    best ability.

### Process:
The scene information accompanied should also be put into
    consideration.
1. **Action Identification and Mapping:** Parse the
    command to extract functions and match them with the
    given functions.
2. **Multiple functions:** If the command contains
    multiple functions or the command needs to be broken
    down into smaller functions, you need to extract
    them separately.
3. **Parameter Extraction:**
   - Extract necessary parameters for functions from the
       command and the scene information.
   - You need to consider the scene information to
       determine the parameter values.
   - If you cannot find the parameter value in the scene
       information and the user command, do not mark it
       as missing, use the scene information and the
       user command to deduct the value to your best
       ability.
4. **JSON Output Construction:** Create a structured JSON
    string that details all functions, including
    parameters, and pronoun indices.
5. **JSON Output Format:**
   - "functions": List of functions with "id", "
       functionType", and "parameters".
   - "id": Unique identifier for the function.
   - "functionType": Type of function.
   - "parameters": List of parameters with "id", "
       parameterType", and additional parameters
       specific to the action type.
   - "parameterType": Type of parameter.

### Output Specification
- A JSON formatted string encapsulating the functions and
    parameters, focusing strictly on providing
    actionable data without additional narrative or
    explanations.

- Example JSON Output for Command "Take that object from
    the table, make it twice as big, and place it on the
    right side of the blue box ":
```json
{
    "functions": [
        {
            "id": 0,
            "functionType": 5,
            "parameters": [
                {
                    "id": 0,
                    "description": "that object from the
                        table",
                    "index": 2,
                    "parameterType": 2
                },
                {
                    "id": 1,
                    "parameterType": 6,
                    "size": {
                        "x": 2.0,
                        "y": 2.0,
                        "z": 2.0
                    }
                }
            ]
        },
        {
            "id": 1,
            "functionType": 2,
            "parameters": [
                {
                    "id": 0,
                    "parameterType": 2
                },
                {
                    "id": 2,
                    "description": "right side of the
                        blue box",
                    "index": 17,
                    "parameterType": 1
                },
            ]
        }
    ]
}
```
**Note:**
"parameters":{[{"id": 0}]} represents the "that object
    from the table" which is defined with "id": 0, so
    you just need to give "id": 0 for the "parameterType
    :2"  parameter in the "move" function.

## A.3   Scene Information

```
{
  "EnvObjects": [
    {"name": "table", "position": {"x": 0.0, "y": 0.055,
        "z": 1.9}, "rotation": {"x": 0.0, "y": 90.0, "z
        ": 0.0}, "size": {"x": 3.0, "y": 0.1, "z": 2.0}}
  ],
  "TargetObjects": [
    {"name": "Red Cube 00", "position": {"x": 0.005, "y":
        0.734, "z": 1.85}, "rotation": {"x": 90.0, "y":
        269.9, "z": 0.0}, "size": {"x": 0.1, "y": 0.1,
        "z": 0.1}},
```

```
{"name": "Blue Cube 01", "position": {"x": 0.007, "y
    ": 0.734, "z": 1.859}, "rotation": {"x": 0.0, "y
    ": 165.4, "z": 270.0}, "size": {"x": 0.1, "y":
    0.1, "z": 0.1}},
{"name": "Red Cube 02", "position": {"x": -0.207, "y
    ": 0.734, "z": 1.833}, "rotation": {"x": 0.0, "y
    ": 3.2, "z": 0.0}, "size": {"x": 0.1, "y": 0.1,
    "z": 0.1}},
{"name": "Blue Cube 03", "position": {"x": 0.253, "y
    ": 0.734, "z": 1.777}, "rotation": {"x": 0.0, "y
    ": 173.3, "z": 0.0}, "size": {"x": 0.1, "y":
    0.1, "z": 0.1}},
```

```
{"name": "Red Cube 04", "position": {"x": -0.149, "y
    ": 0.734, "z": 1.956}, "rotation": {"x": 0.0, "y
    ": 263.4, "z": 0.0}, "size": {"x": 0.1, "y":
    0.1, "z": 0.1}}
]
}
```