# 🦙 ComboBench: Can LLMs Manipulate Physical Devices to Play Virtual Reality Games?

**Shuqing Li**[1]    **Jiayi Yan**[1*]    **Chenyu Niu**[1*]    **Jen-tse Huang**[1]
**Yun Peng**[1]    **Wenxuan Wang**[1†]    **Yepang Liu**[2]    **Michael R. Lyu**[1]
[1]Chinese University of Hong Kong    [2]Southern University of Science and Technology

## Abstract

Virtual Reality (VR) games require players to translate high-level semantic actions into precise device manipulations using controllers and head-mounted displays (HMDs). While humans intuitively perform this translation based on common sense and embodied understanding, whether Large Language Models (LLMs) can effectively replicate this ability remains underexplored. This paper introduces a benchmark, ComboBench, evaluating LLMs' capability to translate semantic actions into VR device manipulation sequences across 262 scenarios from four popular VR games: Half-Life: Alyx, Into the Radius, Moss: Book II, and Vivecraft. We evaluate seven LLMs, including GPT-3.5, GPT-4, GPT-4o, Gemini-1.5-Pro, LLaMA-3-8B, Mixtral-8x7B, and GLM-4-Flash, compared against annotated ground truth and human performance. Our results reveal that while top-performing models like Gemini-1.5-Pro demonstrate strong task decomposition capabilities, they still struggle with procedural reasoning and spatial understanding compared to humans. Performance varies significantly across games, suggesting sensitivity to interaction complexity. Few-shot examples substantially improve performance, indicating potential for targeted enhancement of LLMs' VR manipulation capabilities. We release all materials at `https://sites.google.com/view/combobench`.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable proficiency in general-purpose task solving (Qin et al., 2023), conquering complex domains such as code (Lee et al., 2024; Lam et al., 2025) or math (Lu et al., 2024) problems. While they exhibit increasingly more human-like characteristics (Huang et al., 2024; Liang et al., 2023), an essential attribute of human intelligence is still underexplored: the ability to rapidly learn and apply unfamiliar concepts by leveraging common sense, prior experiences, and a repertoire of cognitive skills.

This is particularly evident in novel interactive environments like video games, where players quickly master device manipulations (atomic actions) and combine them to achieve complex semantic goals. Virtual Reality (VR) games elevate this challenge. They demand not only the execution of atomic actions via physical devices (*e.g.*, Head-Mounted Displays (HMDs) and controllers) but also the inference of complex, often uninstructed, semantic actions. For instance, in *Half-Life: Alyx* (Valve, 2020), when asked to "surrender," players might instinctively raise their controller-held hands even if not explicitly taught.

Such translation of high-level intent into a sequence of physical device manipulations engages a suite of cognitive abilities: (1) *Task decomposition*: Breaking down a high-level semantic action (*e.g.*, "tame the horse" and "plant wheat") into a coherent series of intermediate steps. (2) *Procedural reasoning*: Understanding the logical and temporal order of these steps, including prerequisite conditions or concurrent actions (*e.g.*, the need to till soil before planting seeds). (3) *Spatial reasoning & contextual awareness*: Interpreting instructions within a 3D spatial context (*e.g.*, "move HMD towards the Creeper" and "crouch through the gap") and understanding environmental cues or object

---

[*]Equal contributions.
[†]Corresponding author.

states (*e.g.*, recognizing a door is open/closed and acting accordingly). (4) *Object interaction & tool use understanding*: Correctly mapping intended sub-actions to specific VR device manipulations (*e.g.*, knowing which button to press to "use" an item, and how to manipulate a controller to simulate "swinging" a tool like a pickaxe). This involves understanding the affordances of virtual objects and tools. (5) *Motor action mapping & VR procedural transfer*: Translating abstract actions (*e.g.*, "press," "move," and "trigger") into specific, executable VR controller commands, potentially by adapting from provided examples or general knowledge of VR interaction paradigms. This touches upon a form of simulated embodied reasoning. (6) *Judgment of termination/continuation conditions*: Recognizing when a sub-task or a looped action is complete (*e.g.*, "mine until the block breaks" and "water until the plant grows"). Therefore, playing VR games serves as a rich testbed for evaluating if LLMs can bridge this gap between abstract understanding and grounded, physical interaction.

To systematically evaluate LLMs' ability to perform this crucial translation, we introduce ComboBench, which stands for Cognitive-Oriented Manipulation Benchmark for game combos using physical VR devices. It comprises 262 scenarios derived from four popular VR games: *Vivecraft* (Vivecraft, 2013) (Minecraft in VR), *Half-Life: Alyx* (Valve, 2020), *Moss: Book II* (Polyarc, 2022), and *Into the Radius* (CMGames, 2019). Each scenario presents a high-level semantic action, and the ground truth consists of a fine-grained sequence of VR device manipulations required to achieve it. These sequences are annotated by experienced VR players, allowing us to analyze LLM-generated outputs at the step-level and map their successes and failures to the aforementioned cognitive abilities. For example, failing to "press the X button" after "moving the HMD towards the Creeper" might indicate a lapse in procedural reasoning or object interaction understanding for that specific step.

We evaluate seven LLMs, including GPT-3.5 (OpenAI, 2022), GPT-4 (OpenAI, 2023), GPT-4o (Hurst et al., 2024), Gemini-1.5-Pro (Team et al., 2024), LLaMA-3-8B (Grattafiori et al., 2024), Mixtral-8x7B (Jiang et al., 2023), and GLM-4-Flash (GLM et al., 2024). We design a multi-dimensional scoring approach that assesses: (1) high-level semantic action understanding, (2) procedural step correctness, and (3) device-specific manipulation accuracy, allowing for fine-grained analysis of where each model succeeds or struggles in the translation process. Our findings reveal significant variation in model performance across cognitive capabilities. All models demonstrate strong task decomposition abilities but show pronounced weaknesses in motor action mapping and procedural reasoning. Gemini-1.5-Pro exhibits the most balanced performance across capabilities, while even advanced models like GPT-4 struggle with spatial reasoning compared to human performance. Few-shot examples substantially improve outcomes, particularly for procedural understanding, with diminishing returns beyond three examples. Performance also varies considerably across games, with models generally performing better in environments with more consistent interaction patterns (Vivecraft) than those requiring nuanced controller manipulations (Half-Life: Alyx). These results highlight specific cognitive gaps in current LLMs' ability to perform simulated embodied reasoning for VR interactions and identify targeted areas for improvement toward more capable virtual agents. Our contributions are:

- We introduce ComboBench, the first benchmark designed to evaluate LLMs' fine-grained cognitive abilities in translating high-level semantic actions into VR device manipulations, comprising 262 human-annotated scenarios from four diverse VR games.
- We define a set of key cognitive abilities crucial for VR interaction and design ComboBench to enable step-level analysis of LLM performance against these dimensions.
- We conduct a comprehensive evaluation of six state-of-the-art LLMs, providing a nuanced analysis of their strengths and weaknesses across these cognitive abilities and offering insights into the current frontiers of LLM-driven VR interaction.

## 2 COMBOBENCH: DESIGN AND CURATION

The ComboBench dataset is meticulously curated to evaluate the capability of LLMs in translating high-level semantic actions into sequences of physical VR device manipulations. This section details the game selection process, scenario definition, and the annotation methodology.

## 2.1 Cognitive Capability Taxonomy Development

To establish a theoretically grounded framework for evaluating LLMs in VR contexts, we conducted structured interviews with three domain experts specializing in cognitive science and educational psychology. The experts were selected based on their research backgrounds in spatial cognition, procedural learning, and embodied interaction, which areas that are highly relevant to VR interaction.

**Expert Interview.** Each expert participated in a 90-minute semi-structured interview focused on identifying and categorizing the cognitive abilities required for translating semantic goals into physical actions in virtual environments. The interviews followed a three-phase structure: (1) open-ended discussion about cognitive processes in VR interaction, (2) systematic review of preliminary capability categories identified from literature, and (3) expert suggestions for refinement and expansion of these categories.

**Taxonomy Refinement.** Following the interviews, we synthesized the experts' insights through thematic analysis. Areas of consensus were directly incorporated into our taxonomy, while divergent perspectives were reconciled through follow-up consultations. This iterative process resulted in the identification of six core capability dimensions that comprehensively capture the cognitive demands of VR interaction: (1) Task decomposition: The ability to break down high-level goals into sequentially ordered sub-tasks. (2) Procedural reasoning: Understanding causal relationships between actions and their temporal dependencies. (3) Spatial reasoning & contextual Awareness: Processing spatial relationships and interpreting environmental cues for action selection. (4) Object interaction & tool use understanding: Comprehending affordances and functional properties of virtual objects. (5) Motor action mapping & VR procedural transfer: Translating conceptual actions into specific physical device manipulations. (6) Judgment of termination/continuation conditions: Recognizing completion states or conditions requiring repeated action.

## 2.2 Game Selection Criteria and Process

To ensure a diverse and relevant set of VR interaction paradigms, we selected games based on a systematic process. First, we queried the Steam store (web, 2023) filtering for titles tagged as "VR Only" and available in "English," sorting the results by user review scores in descending order. We then iteratively examined games from this ranked list, focusing on their primary genre as categorized by Steam. To ensure genre diversity, we prioritized games from genres not yet represented in our collection. A crucial selection criterion was the availability of comprehensive textual walkthroughs. For each candidate game, we searched for detailed guides using keywords such as "walkthrough," "guide," or "tutorial." A walkthrough was deemed sufficiently detailed if it provided unambiguous, step-by-step instructions enabling the completion of core game objectives or specific complex tasks. Following this methodology, we selected four popular and critically acclaimed VR games representing distinct genres and interaction styles for ComboBench: (1) *Vivecraft* (Vivecraft, 2013) (Open-world sandbox, crafting) (2) *Half-Life: Alyx* (Valve, 2020) (First-person shooter, puzzle-solving, physics-based interaction) (3) *Moss: Book II* (Polyarc, 2022) (Third-person action-adventure, puzzle-platformer) (4) *Into the Radius* (CMGames, 2019) (First-person survival shooter, exploration) Such selection provides a rich variety of control schemes and task complexities for evaluating LLMs.

## 2.3 Scenario Definition: Semantic Action Identification

For all selected games, eight data annotators, comprising undergraduate and postgraduate computer science students with at least two years of programming experience and sufficient knowledge about VR games, manually identified salient semantic actions from the collected textual walkthroughs. Semantic actions were defined as high-level, goal-oriented tasks described in the walkthroughs (e.g., "tame the horse," "kill the creeper," "solve the gravity glove puzzle") that necessitate a sequence of fine-grained VR device manipulations to accomplish. We focused on scenarios that: (1) involve complex interactions not always explicitly detailed in in-game tutorials, (2) often constitute essential steps or objectives required for game progression. This process resulted in the identification of 262 distinct scenarios across the four games.

## 2.4 Annotation of VR Device Manipulations

Experienced VR users from our annotation team then played through each identified semantic action in the respective games using Oculus Quest 2 VR hardware. The objective was to record the precise sequence of device manipulations required to complete each semantic action. The annotation process captured the following details for each step within a manipulation sequence:

- **Device used:** Specification of whether the HMD or a controller was used.
- **Controller specificity:** If a controller was used, and the action was hand-specific (e.g., primary hand for a tool), the annotation indicated whether the left or right controller was required. If either controller could perform the action, this was noted as "left or right controller."
- **Operation type and parameters:**
  - *Movement:* For actions involving device movement (HMD or controller), the direction (e.g., "towards the Creeper," "upwards") or target position was recorded.
  - *Button presses:* The specific button involved and the action (e.g., "press X button," "release trigger") were noted.
  - *Joystick/thumbstick manipulation:* The direction of joystick push (e.g., "push left thumbstick forward") was recorded.
- **Sequential composition:** For complex semantic actions composed of multiple, distinct sub-actions that might have been annotated individually, the sequence and composition of these simpler actions were explicitly recorded.

## 2.5 Cognitive Capability Labeling Using LLMs

A critical aspect of ComboBench is the annotation of each manipulation step with the specific cognitive capabilities it engages. This fine-grained labeling enables precise analysis of where LLMs succeed or fail in the VR interaction translation process.

**Initial Human Annotation.** To begin, our annotators manually labeled a subset of 50 manipulation sequences (approximately 20% of the dataset), assigning relevant capability categories to each step based on the taxonomy described in Section 2.1. For example, in the sequence required to "tame a horse" in Vivecraft, the step "equip the saddle by pressing the Y button while looking at the inventory slot containing the saddle" was labeled with "Object Interaction & Tool Use Understanding" and "Motor Action Mapping."

**LLM-Assisted Annotation Pipeline.** We then developed an LLM-assisted annotation pipeline to scale this process to the entire dataset. Specifically: ① We used the human-annotated examples as few-shot demonstrations for GPT-4o. ② For each unlabeled manipulation step, we provided the LLM with: [2.a] The semantic action context (e.g., "taming a horse in Vivecraft"). [2.b] The specific manipulation step to label. [2.c] The preceding and following steps (when available). [2.d] Detailed descriptions of each capability category. [2.e] Three few-shot examples with explanations of why each capability was assigned. ③ The LLM generated capability labels along with justifications for each assignment. ④ Human annotators reviewed the LLM-generated labels, making corrections when necessary. The review process revealed an 89.7% agreement rate between LLM-assigned labels and human judgments.

**Multi-label Distribution.** Most manipulation steps engaged multiple cognitive capabilities simultaneously. On average, each step was associated with 2.3 capability categories ($\sigma = 0.8$). The most frequently co-occurring capabilities were "Motor Action Mapping" and "Object Interaction & Tool Use Understanding" (present together in 68% of steps), reflecting the inherent coupling between understanding virtual object affordances and translating this understanding into physical manipulations.

## 2.6 Contextualization and Verification

To further contextualize the annotated actions and aid in verification, we sourced or recorded gameplay videos corresponding to the textual walkthroughs for each game. For each annotated semantic action and its constituent manipulation steps, we recorded the corresponding timestamps in these videos. This allows for visual verification of the annotated sequences and provides richer context for

understanding the actions. If suitable public gameplay videos matching the exact walkthrough steps were unavailable, our annotators recorded their own gameplay sessions while performing the actions.

## 3 EXPERIMENTS

### 3.1 MODEL SELECTION

We evaluate six state-of-the-art LLMs: GPT-4o, GPT-4-turbo, GPT-3.5-turbo, Gemini-1.5-Pro, LLaMA-3-8B, and Mixtral-8x7B. We also perform human evaluation to validate the average human capabilities for comparison, when humans are given exactly the same input as LLMs. For all experiments, we used the official APIs for proprietary models and Hugging Face implementations for open-source models. Temperature was set to 0 across all models to minimize non-deterministic outputs. For embedding calculations, we utilized OpenAI's text-embedding-3-large model via their API.

### 3.2 EVALUATION METRICS

To comprehensively evaluate the capability of LLMs in translating semantic actions into VR device manipulations, we propose a multi-dimensional evaluation framework with four distinct metrics. These metrics collectively capture different aspects of model performance in ComboBench, ranging from strict matching to more flexible semantic alignment.

**Strict Step-by-Step Matching (SSM).** Our first metric evaluates the exact matching between model-generated and ground truth steps, enforcing both sequence length equivalence and semantic alignment: $\text{SSM} = \frac{\text{Number of correctly predicted sequences}}{\text{Total number of sequences}}$. A sequence is considered correctly predicted only when: the number of steps in the generated sequence equals that of the ground truth, and every step in the generated sequence has a cosine similarity above a threshold of 0.8387 with its corresponding step in the ground truth This strict metric serves as a measure of precision in reproducing exact device manipulation sequences and rewards models that can generate complete, step-accurate instructions.

**Common Subsequence Evaluation.** We further introduce two complementary metrics based on common subsequence alignment to assess partial correctness: (1) **Normalized Step Alignment Score (NSAS)** This metric quantifies the alignment between the model-generated sequence and ground truth while accounting for missing and additional steps: $\text{NSAS} = \frac{(|C|-|M|-|A|)-\min_{\text{all\_samples}}}{|G|\cdot(\max_{\text{all\_samples}}-\min_{\text{all\_samples}})}$, where: $|C|$ represents the count of correctly matched steps in the common subsequence, $|M|$ represents missing steps from the ground truth, $|A|$ represents additional steps generated by the model, $|G|$ represents the total number of steps in the ground truth, $\min_{\text{all\_samples}}$ and $\max_{\text{all\_samples}}$ represent the minimum and maximum raw scores across all evaluations, enabling consistent normalization This score is normalized across the entire dataset to ensure fair comparison across different models and scenarios. (2) **Sequential Order Preservation (SOP)** The SOP metric specifically assesses the model's ability to maintain the correct procedural ordering of steps: $\text{SOP} = \frac{|\text{Steps correctly ordered and matched}|}{|G|}$. This metric evaluates whether the steps in the matched subsequence maintain their ordinal positions (e.g., step 1 followed by step 2, etc.) in both the ground truth and model output, capturing the model's procedural reasoning capabilities.

**Semantic Step Coverage (SSC).** Our final metric adopts a more flexible matching approach to evaluate semantic coverage of critical actions: $\text{SSC} = \frac{|\text{MR steps matched to any GT step}|}{|\text{MR}|}$, where a model result (MR) step is considered matched if it has a cosine similarity above the threshold (0.8387) with any step in the ground truth (GT). This metric computes the proportion of generated steps that semantically align with at least one ground truth step, regardless of position.

### 3.3 EXPERIMENTAL RESULTS

We analyze and answer the following Research Questions (RQs): **(RQ1)** How do state-of-the-art LLMs perform in translating semantic actions into VR device manipulations across different VR games? **(RQ2)** How does the number of few-shot examples affect LLMs' ability to execute this

Table 1: Overall performance comparison of LLMs across VR games (5-shot setting). Best model performance per metric is **bolded**, second best is <u>underlined</u>.

| Model | Half-Life: Alyx | | | | Into the Radius | | | | Moss: Book II | | | | Vivecraft | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NSAS↑ | SOP↑ | F1$_{SOP}$↑ | SSC↑ | NSAS↑ | SOP↑ | F1$_{SOP}$↑ | SSC↑ | NSAS↑ | SOP↑ | F1$_{SOP}$↑ | SSC↑ | NSAS↑ | SOP↑ | F1$_{SOP}$↑ | SSC↑ |
| GPT-3.5 | **0.858** | 0.123 | **0.287** | 0.143 | 0.662 | 0.169 | 0.226 | 0.137 | 0.782 | 0.169 | 0.207 | 0.186 | 0.922 | 0.043 | 0.098 | 0.067 |
| GPT-4 | <u>0.853</u> | <u>0.125</u> | 0.258 | **0.172** | <u>0.693</u> | 0.189 | <u>0.328</u> | 0.177 | **0.824** | 0.218 | 0.336 | <u>0.220</u> | 0.927 | <u>0.137</u> | 0.437 | 0.081 |
| GPT-4o | 0.804 | 0.022 | 0.075 | <u>0.167</u> | **0.698** | **0.291** | **0.414** | **0.190** | **0.824** | **0.300** | <u>0.342</u> | **0.222** | <u>0.931</u> | **0.190** | **0.489** | **0.096** |
| Mixtral | 0.839 | **0.126** | 0.246 | 0.147 | 0.666 | 0.123 | 0.228 | 0.097 | 0.756 | 0.117 | 0.191 | 0.121 | 0.926 | 0.060 | 0.239 | 0.070 |
| LLaMA-3 | 0.848 | **0.126** | <u>0.279</u> | 0.162 | 0.644 | <u>0.242</u> | 0.317 | 0.168 | <u>0.823</u> | <u>0.283</u> | **0.349** | 0.200 | 0.929 | 0.039 | 0.122 | 0.042 |
| GLM-4 | 0.836 | 0.076 | 0.183 | 0.149 | 0.618 | 0.096 | 0.186 | 0.149 | 0.749 | 0.087 | 0.174 | 0.165 | 0.909 | 0.000 | 0.045 | 0.061 |
| Human | 0.845 | 0.090 | 0.240 | 0.110 | 0.684 | 0.148 | 0.257 | <u>0.181</u> | 0.817 | 0.112 | 0.328 | 0.174 | **0.935** | 0.122 | <u>0.482</u> | <u>0.084</u> |

Table 2: Overall performance across VR games and settings. We report the average scores for our four evaluation metrics: Strict Step-by-Step Matching (SSM), Normalized Step Alignment Score (NSAS), Sequential Order Preservation (SOP), and Semantic Step Coverage (SSC). Higher is better for all metrics. Bold indicates best model performance, underline indicates second best.

| Model | Average Across Settings | | | | Zero-Shot | | | | 5-Shot | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SSM (%) | NSAS | SOP | SSC | SSM (%) | NSAS | SOP | SSC | SSM (%) | NSAS | SOP | SSC |
| GPT-3.5 | 1.4 | 0.781 | 0.063 | 0.066 | 0.8 | 0.771 | 0.003 | 0.046 | 2.1 | 0.791 | 0.128 | 0.095 |
| GPT-4 | 3.7 | 0.806 | 0.107 | 0.124 | 1.0 | 0.788 | 0.015 | 0.107 | 8.8 | 0.825 | 0.184 | 0.140 |
| GPT-4o | <u>5.3</u> | 0.797 | 0.138 | <u>0.141</u> | 0.6 | 0.785 | 0.015 | 0.108 | <u>10.9</u> | 0.806 | 0.228 | <u>0.161</u> |
| Gemini-1.5 | **5.8** | **0.813** | **0.146** | **0.142** | **2.1** | **0.795** | 0.010 | **0.124** | **11.7** | **0.832** | **0.236** | **0.162** |
| Mixtral | 1.1 | 0.784 | 0.068 | 0.079 | 0.0 | 0.777 | 0.002 | 0.040 | 2.2 | 0.796 | 0.105 | 0.107 |
| LLaMA-3 | 1.2 | 0.787 | 0.088 | 0.111 | 0.1 | 0.783 | 0.011 | 0.088 | 1.8 | 0.794 | 0.163 | 0.132 |
| GLM-4 | 0.0 | 0.761 | 0.038 | 0.077 | 0.0 | 0.762 | 0.006 | 0.052 | 0.0 | 0.765 | 0.071 | 0.120 |
| Human | 1.2 | 0.833 | 0.122 | 0.159 | – | – | – | – | – | – | – | – |

translation? **(RQ3)** Do LLM and human performance exhibit significant variations across the four different VR games, potentially indicating sensitivity to game mechanics and interaction complexity? **(RQ4)** Which cognitive capabilities do current LLMs excel at, and where do they struggle? **(RQ5)** How do LLMs compare to human performance in VR device manipulation tasks?

## 3.4 RQ1 & RQ3: LLM PERFORMANCE ACROSS VR GAMES

Tables 1, 2, and 3 present comprehensive performance metrics for all evaluated LLMs across the four VR games. Our analysis reveals substantial variations in model capabilities and game-specific challenges. Gemini-1.5-Pro emerges as the strongest performer overall, achieving the highest NSAS scores in three of the four games (Half-Life: Alyx: 0.863, Moss: Book II: 0.848, Vivecraft: 0.938), while maintaining competitive performance in Into the Radius (0.682). GPT-4o demonstrates particular strength in Into the Radius with the highest SOP score (0.291)

Table 3: Cross-game performance variation (standard deviation across games) w/ 5-shot examples.

| Model | NSAS $\sigma$↓ | SOP $\sigma$↓ | F1$_{SOP}$ $\sigma$↓ | Game Gap↓ |
|---|---|---|---|---|
| GPT-3.5 | 0.110 | 0.061 | 0.084 | 0.085 |
| GPT-4 | 0.059 | 0.051 | 0.081 | 0.074 |
| GPT-4o | 0.068 | 0.137 | 0.184 | 0.127 |
| Gemini-1.5 | **0.099** | **0.093** | **0.127** | **0.095** |
| Mixtral | 0.114 | 0.031 | 0.065 | 0.070 |
| LLaMA-3 | 0.112 | 0.103 | 0.120 | 0.113 |
| GLM-4 | 0.135 | 0.049 | 0.069 | 0.084 |
| Human | 0.105 | 0.029 | 0.117 | 0.084 |

and F1$_{SOP}$ (0.414), suggesting superior procedural reasoning capabilities in this specific game context. GPT-4-turbo maintains consistently strong performance across all games, positioning itself as a reliable general-purpose model for VR interaction translation.

A striking pattern emerges in the SOP metrics, which vary dramatically across both models and games (0.00-0.30 range). While NSAS scores remain relatively high (mostly >0.75), indicating models can identify relevant steps, the low SOP values reveal fundamental difficulties in maintaining correct temporal ordering. This discrepancy is particularly pronounced in Vivecraft, where models achieve high NSAS scores (0.909-0.938) but struggle with step ordering (SOP: 0.00-0.19), suggesting that simpler interaction patterns may paradoxically lead to overconfidence in step sequencing.

Analysis of performance variations (Table 3) reveals significant game-dependent effects. Vivecraft exhibits the highest average performance across models (NSAS: 0.922-0.938), likely due to its consistent block-based interaction paradigm inherited from Minecraft. In contrast, Into the Radius presents the greatest challenge, with notably lower NSAS scores (0.618-0.698) and high performance variance. This pattern suggests that games featuring realistic physics simulations and complex inventory management pose particular difficulties for current LLMs.

Interestingly, different models exhibit distinct strengths across game types. GPT-4o shows remarkable adaptability in Into the Radius (SOP: 0.291) compared to other models, while struggling in Half-Life: Alyx (SOP: 0.022). Gemini-1.5-Pro maintains the most balanced performance profile across games (Game Gap: 0.095), suggesting more robust generalization capabilities. Smaller models like Mixtral-8x7B and GLM-4-flash show disproportionate performance degradation in complex environments, with GLM-4-flash achieving zero SOP in Vivecraft despite reasonable NSAS scores. The substantial performance variations across games highlight the impact of interaction design on LLM capabilities. Games with discrete, well-defined actions (Vivecraft) enable higher model performance, while those requiring nuanced controller manipulation and spatial reasoning (Half-Life: Alyx, Into the Radius) expose current limitations. The correlation between game complexity and performance degradation is non-linear, moderate complexity (Moss: Book II) sometimes yields better results than simpler environments, suggesting that models may benefit from richer contextual cues in certain scenarios.

These findings collectively demonstrate that while state-of-the-art LLMs have made significant progress in understanding VR interactions, their performance remains highly sensitive to specific game mechanics and interaction paradigms. The gap between high NSAS scores and low SOP values across all games indicates that current models can identify relevant actions but struggle with the procedural reasoning required to sequence them correctly, which is an important capability for successful VR interaction.

## 3.5 RQ2: IMPACT OF FEW-SHOT EXAMPLES

Table 4 demonstrates that few-shot examples substantially improve LLM performance in VR device manipulation tasks, with the most dramatic gains observed in Sequential Order Preservation (SOP), where scores increase by 10–20x from near-zero baselines. All models benefit from in-context examples, though with diminishing returns, the improvement

Table 4: Effect of few-shot examples on model performance (average across all games)

| Model | Zero-shot | | | 3-shot | | | 5-shot | | |
|---|---|---|---|---|---|---|---|---|---|
| | NSAS↑ | SOP↑ | SSM↑ | NSAS↑ | SOP↑ | SSM↑ | NSAS↑ | SOP↑ | SSM↑ |
| GPT-3.5 | 0.781 | 0.003 | 0.006 | 0.783 | 0.113 | 0.013 | 0.806 | 0.128 | 0.022 |
| GPT-4 | 0.799 | 0.015 | 0.012 | 0.807 | 0.106 | 0.039 | 0.824 | 0.167 | 0.066 |
| GPT-4o | 0.785 | 0.015 | 0.012 | 0.802 | 0.166 | 0.059 | 0.815 | 0.159 | **0.077** |
| Gemini-1.5 | 0.797 | 0.010 | 0.016 | 0.827 | 0.174 | 0.048 | **0.833** | **0.207** | 0.056 |
| Mixtral | 0.784 | 0.002 | 0.000 | 0.766 | 0.103 | 0.010 | 0.797 | 0.106 | 0.021 |
| LLaMA-3 | 0.786 | 0.011 | 0.001 | 0.787 | 0.139 | 0.014 | 0.811 | 0.172 | 0.020 |
| GLM-4 | 0.747 | 0.005 | 0.000 | 0.756 | 0.036 | 0.000 | 0.778 | 0.065 | 0.000 |

from zero-shot to 3-shot (average NSAS gain: 2.1%, SOP: 10-fold increase) significantly exceeds that from 3-shot to 5-shot (NSAS: 1.4%, SOP: 20-50% relative gain). Gemini-1.5-Pro exhibits the strongest adaptability, achieving the highest 5-shot performance (NSAS: 0.833, SOP: 0.207), while maintaining consistent improvements across all metrics. The differential impact across metrics reveals that few-shot examples primarily address procedural sequencing challenges (massive SOP improvements) more effectively than exact step matching (modest SSM gains), suggesting that demonstrations help models understand temporal dependencies in VR interactions but do not fully resolve the complexity of translating semantic actions into precise device manipulations.

## 3.6 RQ4: COGNITIVE CAPABILITIES ANALYSIS

We analyzed model performance across six cognitive capabilities required for effective VR interaction (Figure 1). By mapping evaluation metrics to capability scores (0-10 scale), we identified specific strengths and limitations in how LLMs approach spatial-mechanical reasoning tasks.

**Areas of Strength:** All evaluated LLMs demonstrate strong task decomposition capabilities (7.8-8.5), with minimal performance gap compared to humans (8.2). Gemini-1.5-Pro leads with a score of 8.5, while even smaller models like Mixtral-8x7B (8.0) and GLM-4-flash (7.8) perform admirably. This suggests that segmenting high-level actions into component steps aligns well with the sequential reasoning abilities developed during language model pre-training.

**Areas of Weakness:** Motor action mapping emerges as the most significant challenge (0.5-4.5), with all models struggling to precisely translate abstract actions into specific VR control manipulations. GPT-4o performs best in this dimension (4.5), but still falls short of robust capability. Procedural reasoning also shows substantial variation (2.3-7.0), with only Gemini-1.5-Pro approaching adequate performance. Judgment of termination conditions represents another challenge area, with most models scoring below 5.0 (except Gemini-1.5-Pro at 6.0), compared to human performance (6.5).

**Model Comparison:** Gemini-1.5-Pro demonstrates the most balanced performance profile, consistently outperforming other models in procedural reasoning (7.0), spatial reasoning (7.5), and termination judgment (6.0). GPT-4 variants show strong task decomposition and object interaction (5.3-5.7) but lag in procedural sequencing. LLaMA-3-8B shows surprisingly competitive performance in procedural reasoning (5.7), outperforming larger models like GPT-3.5-Turbo (4.3), suggesting architecture differences may be as important as scale.

## 3.7 RQ5: Comparison with Human Performance

To contextualize our findings, we compare LLM performance against human baselines across our evaluation metrics. As shown in Tables 1 and 2, state-of-the-art LLMs demonstrate competitive performance with humans on several key dimensions. Our results reveal a nuanced performance landscape. While top-performing models (Gemini-1.5-Pro, GPT-4o) achieve comparable or superior NSAS scores relative to humans in certain games (e.g., Vivecraft: 0.931-0.938 vs. 0.935 for humans), human participants maintain a decisive advantage in Sequential Order Preservation, particularly for games requiring complex interaction sequences. In Half-Life: Alyx, humans achieve only 0.090 SOP compared to the best model performance of 0.209 (Gemini-1.5-Pro), yet this reflects the challenging nature of the task rather than superior model performance, both humans and models struggle with the intricate procedural requirements of this game.
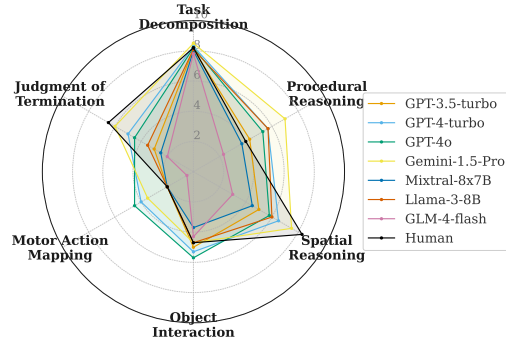


Figure 1: Cognitive capabilities of LLMs and humans in translating semantic actions to VR device manipulations. Higher scores (0-10 scale) indicate stronger abilities.

Analysis of performance variance across games (Table 3) reveals striking similarities between human and high-performing model behavior. The standard deviation of human performance (0.084) closely aligns with that of GPT-4 (0.074) and Mixtral-8x7B (0.070), suggesting that both humans and advanced LLMs exhibit similar sensitivity patterns to game-specific interaction complexities. This convergence is particularly evident in structured environments like Vivecraft, where the consistency gap between humans and LLMs has substantially narrowed. Figure 1 illustrates the capability-wise performance comparison, revealing critical gaps in embodied reasoning. Humans maintain superior performance in spatial reasoning (8.3 vs. 7.5 for Gemini-1.5-Pro) and judgment of termination conditions (6.5 vs. 6.0). These differences are statistically significant ($p < 0.05$, Wilcoxon signed-rank test) and persist across all evaluated models. This performance gap suggests that while LLMs have achieved remarkable progress in understanding VR interaction semantics, they lack the grounded physical intuition that humans naturally apply when reasoning about three-dimensional manipulations and determining action completion states.

The convergence of human and LLM performance on certain metrics, coupled with persistent gaps in spatial and termination reasoning, indicates that current language models can effectively decompose VR tasks but struggle with aspects requiring embodied experience. This finding has important implications for the development of future VR-capable AI systems, suggesting the need for training paradigms that better incorporate spatial and physical reasoning capabilities.

# 4 RELATED WORK

Recent work has explored the use of large language models (LLMs) as generalist agents for embodied and interactive reasoning tasks. In robotics, *SayCan* (Ahn et al., 2022) and *PaLM-E* (Driess et al., 2023) combine LLMs with affordance-based skill models or multimodal inputs to plan and execute robot actions in real-world settings. These methods demonstrate that LLMs can decompose high-level goals into actionable steps when grounded in sensory input and executable primitives. Similar capabilities have been investigated in virtual domains such as Minecraft through agents like *Voyager* (Wang et al., 2023) and simulation platforms like *MineDojo* (Fan et al., 2022), which showcase in-context learning and autonomous skill acquisition by prompting LLMs to generate and refine code or sub-goals based on environmental feedback. However, these systems are typically tuned for code-level or symbolic outputs and do not focus on physical device manipulation or spatially grounded motor control as required in VR environments.

Task decomposition and procedural reasoning have been studied extensively via prompting strategies such as Chain-of-Thought (Wei et al., 2022) and ReAct (Yao et al., 2022), which interleave reasoning with action selection to improve coherence in multi-step planning. LLMs have also been used to generate structured action sequences or API calls from natural instructions in domains like household tasks (Shridhar et al., 2020) and scientific procedures (Wang et al., 2022). Code-as-policy paradigms (Liang et al., 2022) show that LLMs can output executable policy code that integrates logical control flow, enabling conditional and iterative actions. However, these approaches often abstract away the complexity of physical or spatial execution, making them less suited for evaluating embodied skills involving real-time device input, object affordances, or 3D spatial constraints.

To assess grounded reasoning, several benchmarks have been proposed in interactive settings. Animal-AI (Mecattaf et al., 2024) evaluates embodied cognition through physics-based tasks adapted from animal intelligence experiments, highlighting LLMs' partial competence in navigation, tool use, and physical causality. Similarly, platforms like ALFWorld (Shridhar et al., 2021) and Science-World (Wang et al., 2022) test instruction-following via text or symbolic interfaces, while MacGyver-style tasks (Tian et al., 2024) probe object-use innovation in constrained settings. These works underscore known limitations in spatial reasoning, persistence, and tool-use generalization among LLMs.

Concurrently, capability-oriented embodied evaluations for multimodal/embodied LLMs have emerged. EmbodiedBench (Yang et al., 2025) unifies high- and low-level tasks with fine-grained error taxonomies, while VLABench (Zhang et al., 2024a) targets long-horizon, language-conditioned manipulation for VLA policies; the Embodied Agent Interface (EAI) (Li et al., 2025) standardizes modules and metrics for step-level diagnostics. In parallel, GUI/OS/mobile control benchmarks, including OSWorld (Xie et al., 2024), SPA-Bench (Zhang et al., 2024b), WebArena (Zhou et al., 2023), Mind2Web (Deng et al., 2023), AndroidEnv (Toyama et al., 2021), and AppAgent/AppAgent v2 (Zhang et al., 2023; Li et al., 2024b), evaluate precise, procedure-level device interactions (click/tap/drag/typing) with programmatic success checks, offering a complementary view of grounded action competence. On the robotics side, VLA policies such as RT-1/RT-2 (Brohan et al., 2022; 2023) and OpenVLA (Kim et al., 2024) map visual observations and language to action tokens, improving semantic generalization in manipulation; large-scale 3D suites like Habitat 2.0/HAB (Savva et al., 2021), BEHAVIOR-1K/OmniGibson (Li et al., 2024a), and CALVIN (Mees et al., 2021) stress long-horizon rearrangement and physics.

In contrast, our benchmark, ComboBench, targets the translation of semantic goals into fine-grained, physically grounded VR device manipulations, enabling a more precise step-level analysis of embodied cognitive abilities critical for real-world interaction.

# 5 CONCLUSION

We present ComboBench, a comprehensive benchmark evaluating seven state-of-the-art LLMs on their ability to translate semantic actions into VR device manipulations across 262 scenarios from four popular VR games. Our evaluation reveals that while advanced models like Gemini-1.5-Pro demonstrate strong task decomposition capabilities (NSAS > 0.8), they exhibit significant weaknesses in motor action mapping and procedural reasoning, with Sequential Order Preservation scores often below 0.3 even in the best cases. Few-shot examples dramatically improve procedural understanding

(10-20x increase in SOP scores) but provide limited benefit for exact step matching, suggesting that in-context learning helps models understand action relationships but cannot fully bridge the gap in physical manipulation reasoning. The pronounced performance variations across games and cognitive capabilities indicate that current text-trained LLMs lack the embodied understanding necessary for reliable VR interaction, pointing to the need for multimodal training approaches that incorporate spatial, visual, and haptic information. These findings highlight fundamental limitations in how language models represent physical interactions and suggest that achieving human-level VR manipulation capabilities will require architectural innovations beyond scaling current approaches, with important implications for the development of embodied AI systems in virtual and augmented reality applications.

## REFERENCES

VR Content on Steam App Store. `https://store.steampowered.com/search/?vrsupport=401`, 2023.

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2022.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023.

CMGames. Into the Radius, 2019. URL `https://www.into-the-radius.com/`.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023.

Danny Driess, Fei Xia, Arjun Srinivas, Wenlong Huang, Julius Müller, Roberto Martín-Martín, Tobias Bücheler, Yevgen Chebotar Du, Karol Hausman, Saran Tunyasuvunakool, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *arXiv preprint arXiv:2206.08853*, 2022.

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, and Michael Lyu. On the humanity of conversational ai: Evaluating the psychological portrayal of llms. In *The Twelfth International Conference on Learning Representations*, 2024.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model, 2024.

Man Ho Lam, Chaozheng Wang, Jen-tse Huang, and Michael R Lyu. Codecrash: Stress testing llm reasoning under structural and semantic perturbations. *arXiv preprint arXiv:2504.14119*, 2025.

Cheryl Lee, Chunqiu Steven Xia, Jen-tse Huang, Zhouruixin Zhu, Lingming Zhang, and Michael R Lyu. A unified debugging approach via llm-based multi-agent synergy. *arXiv preprint arXiv:2404.17153*, 2024.

Chengshu Li, Josiah Wong, Michael Lingelbach, Roberto Martín-Martín, Jim Fan, et al. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation, 2024a.

Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu Zhou, Sanjana Srivastava, Cem Gokmen, Tony Lee, Li Erran Li, Ruohan Zhang, Weiyu Liu, Percy Liang, Li Fei-Fei, Jiayuan Mao, and Jiajun Wu. Embodied agent interface: Benchmarking llms for embodied decision making, 2025.

Yanda Li, Chi Zhang, Wanqi Yang, Bin Fu, Pei Cheng, Xin Chen, Ling Chen, and Yunchao Wei. Appagent v2: Advanced agent for flexible mobile interactions, 2024b.

Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*, 2022.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024.

Matteo G. Mecattaf, Ben Slater, Marko Tešić, Jonathan Prunty, Konstantinos Voudouris, and Lucy G. Cheke. A little less conversation, a little more action, please: Investigating the physical common-sense of llms in a 3d embodied environment. *arXiv preprint arXiv:2410.23242*, 2024.

Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks, 2021.

OpenAI. Introducing chatgpt. *OpenAI Blog Nov 30 2022*, 2022. URL https://openai.com/index/chatgpt/.

OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Polyarc. Moss: Book II, 2022. URL https://www.polyarcgames.com/games/moss-book-ii.

Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. Is chatgpt a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*, 2023.

Manolis Savva, Abhishek Kadian, Erik Wijmans, Shengyi Qian, Angel Chang, et al. Habitat 2.0: Training home assistants to rearrange their habitat, 2021.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10740–10749, 2020.

Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. URL https://arxiv.org/abs/2010.03768.

Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.

Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ronan Le Bras, Rami Marjieh, Nanyun Peng, Yejin Choi, Thomas L Griffiths, and Faeze Brahman. Macgyver: Are large language models creative problem solvers? In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5303–5324, 2024.

Daniel Toyama, Philippe Hamel, Anita Gergely, Gheorghe Comanici, Amelia Glaese, Zafarali Ahmed, Tyler Jackson, Shibl Mourad, and Doina Precup. Androidenv: A reinforcement learning platform for android, 2021.

Valve. Half-Life: Alyx, 2020. URL https://www.half-life.com/en/alyx/.

Vivecraft. Vivecraft – Virtual Reality Minecraft for SteamVR, 2013. URL https://www.vivecraft.org/.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.

Yujia Wang, Tushar Khot, Ashish Sabharwal, and Peter Clark. Scienceworld: Is your agent smarter than a 5th grader? *arXiv preprint arXiv:2203.07540*, 2022.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.

Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments, 2024.

Rui Yang, Hanyang Chen, Junyu Zhang, Mark Zhao, Cheng Qian, Kangrui Wang, Qineng Wang, Teja Venkat Koripella, Marziyeh Movahedi, Manling Li, Heng Ji, Huan Zhang, and Tong Zhang. Embodiedbench: Comprehensive benchmarking multi-modal large language models for vision-driven embodied agents, 2025.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users, 2023.

Shiduo Zhang, Zhe Xu, Peiju Liu, Xiaopeng Yu, Yuan Li, Qinghui Gao, Zhaoye Fei, Zhangyue Yin, Zuxuan Wu, Yu-Gang Jiang, and Xipeng Qiu. Vlabench: A large-scale benchmark for language-conditioned robotics manipulation with long-horizon reasoning tasks, 2024a.

Zhaofeng Zhang, Yiyan Qi, Jinjie Ni, Jiayi Yuan, Fangkai Yang, et al. Spa-bench: A comprehensive benchmark for smartphone agent evaluation, 2024b.

Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents, 2023.

# A Preliminaries on Virtual Reality

Virtual Reality (VR) represents a fundamentally distinct paradigm of human-computer interaction that transcends traditional interface boundaries. Unlike conventional computing systems that rely on indirect manipulation through keyboards, mice, and two-dimensional displays, VR creates immersive digital environments where users experience presence and embodiment. This paradigm shift necessitates a comprehensive understanding of both the technological infrastructure and the cognitive demands placed on users who must translate abstract intentions into concrete physical manipulations within virtual spaces.

The evolution of VR technology has progressed through several generations, from early tethered systems requiring substantial computational infrastructure to modern standalone devices that integrate processing, display, and tracking capabilities within compact form factors. Contemporary VR systems can be broadly categorized into three architectural approaches. PC-tethered headsets leverage external computational resources to deliver high-fidelity experiences with complex graphics and physics simulations. Standalone headsets, exemplified by devices like the Meta Quest series, incorporate integrated processors that balance performance with portability. Mobile-phone-based solutions represent an accessible entry point, utilizing smartphones as both display and processor, though with inherent limitations in tracking precision and computational capability.

The core hardware components enabling VR interaction form an integrated ecosystem of sensory input and output devices. Head-Mounted Displays (HMDs) serve as the primary visual interface, providing stereoscopic rendering that creates depth perception while simultaneously tracking head orientation and position through integrated sensors. This tracking enables natural viewing behaviors where users can examine virtual objects by physically moving their heads, mirroring real-world visual exploration patterns. Motion controllers, typically deployed in pairs to represent both hands, enable direct manipulation of virtual objects through a combination of positional tracking, button inputs, trigger mechanisms, and thumbstick controls. These devices must balance ergonomic considerations with functional complexity, providing sufficient input channels while maintaining intuitive operation. Spatial tracking systems, whether implemented through external sensors (outside-in tracking) or integrated cameras (inside-out tracking), monitor user movements with six degrees of freedom, capturing both translational and rotational motion to enable natural locomotion and interaction within virtual environments.

The ongoing evolution of VR hardware continues to introduce novel interaction modalities. Haptic gloves promise to deliver tactile feedback through actuators that simulate texture, resistance, and temperature. Full-body tracking systems capture skeletal motion to enable more nuanced avatar control and gesture recognition. Specialized peripherals, from steering wheels for racing simulations to weapon replicas for combat games, demonstrate the trend toward application-specific controllers that enhance immersion through physical affordances that match virtual interactions.

## A.1 Interaction Paradigms and Design Principles

The design of VR interaction paradigms represents a delicate balance between leveraging users' existing motor skills and introducing novel control schemes that exploit the unique capabilities of virtual environments. Direct manipulation forms the foundation of most VR interactions, where users employ hand controllers to simulate natural actions like grasping, throwing, and pushing. This approach capitalizes on users' lifetime of experience with physical object manipulation but requires careful calibration of virtual physics to match expectations. The mapping between controller inputs and virtual hand movements must account for the absence of tactile feedback, often employing visual or auditory cues to confirm successful interactions.

Ray-casting emerged as an elegant solution to the fundamental challenge of interacting with objects beyond physical reach. By projecting virtual rays from controllers, users can select, manipulate, and activate distant objects without locomotion. This technique exemplifies how VR interaction design often augments natural human capabilities rather than strictly simulating physical constraints. Advanced ray-casting implementations incorporate features like ray curvature for improved ergonomics, variable ray length based on context, and visual feedback mechanisms that indicate interaction possibilities.

Gesture recognition systems interpret temporal patterns of controller or hand movement as discrete commands, enabling a rich vocabulary of interactions without relying on button combinations. These systems must balance recognition accuracy with user comfort, avoiding gestures that cause fatigue or require precise movements difficult to perform consistently. Machine learning approaches have enhanced gesture recognition capabilities, allowing for more natural and varied input patterns while maintaining reliable detection rates.

Symbolic input mechanisms address scenarios where direct physical analogues are impractical or inefficient. Virtual keyboards present unique challenges in VR, as users lack tactile feedback and must rely on visual confirmation of key presses. Solutions range from laser-pointer selection of virtual keys to gesture-based text entry systems that map hand movements to characters. Voice commands offer an alternative input modality that bypasses manual interaction entirely, though they introduce considerations around recognition accuracy, latency, and social acceptability in shared spaces.

## A.2 Development Platforms and Technical Considerations

The creation of VR applications relies on sophisticated development ecosystems that abstract hardware complexity while providing fine-grained control over interaction mechanics. Unity and Unreal Engine have emerged as dominant platforms, offering comprehensive toolsets that handle rendering pipelines, physics simulation, spatial audio, and cross-platform deployment. These engines provide specialized VR interaction frameworks that standardize common patterns like object grabbing, teleportation, and menu systems, significantly reducing development complexity.

Hardware software development kits (SDKs) serve as the bridge between high-level application logic and device-specific capabilities. Meta's OpenXR initiative represents an industry effort to standardize VR/AR interfaces, enabling applications to target multiple hardware platforms without extensive modifications. Platform-specific SDKs like SteamVR and Oculus SDK continue to play important roles, offering access to proprietary features and optimizations that enhance performance on particular hardware.

Technical constraints fundamentally shape VR interaction design decisions. Maintaining consistent frame rates above 72Hz (and preferably 90Hz or higher) prevents motion sickness and ensures responsive interactions. This performance requirement influences every aspect of application design, from polygon counts and texture resolution to the complexity of physics simulations. Tracking precision varies across hardware platforms and environmental conditions, necessitating interaction designs that accommodate occasional tracking losses or reduced accuracy. Developers must also consider the diverse computational capabilities across the VR ecosystem, implementing scalable solutions that provide acceptable experiences on entry-level hardware while leveraging the capabilities of high-end systems.

## A.3 Challenges in VR Interaction

Despite remarkable technological progress, VR interaction continues to face fundamental challenges that impact user experience and limit application domains. The locomotion problem exemplifies the tension between physical and virtual spaces. While users may explore vast virtual environments, they remain constrained by finite physical play areas. Teleportation offers a practical solution but breaks immersion and can cause spatial disorientation. Artificial locomotion through thumbstick control risks motion sickness in susceptible users. More exotic solutions like omnidirectional treadmills or redirected walking techniques remain impractical for consumer applications.

The absence of comprehensive haptic feedback represents perhaps the most significant limitation in current VR systems. While controllers provide basic vibration feedback, they cannot simulate the rich tactile experiences of real-world interaction: the weight of objects, surface textures, temperature variations, or resistance to movement. This sensory gap creates a fundamental disconnect between visual expectations and physical sensations, requiring users to adapt their interaction strategies and often leading to reduced precision in manipulation tasks.

Interaction discoverability poses ongoing challenges as VR applications lack standardized interface conventions comparable to desktop or mobile platforms. Users encountering new VR experiences must often learn application-specific control schemes, gesture sets, and interaction patterns. The absence of persistent visual UI elements (to maintain immersion) exacerbates this challenge, as users

cannot easily reference control schemes during gameplay. This lack of standardization increases cognitive load and creates barriers to entry for new users.

Precision manipulation tasks highlight the limitations of current tracking systems and input devices. Tasks requiring fine motor control, such as threading a virtual needle or manipulating small components, prove challenging due to tracking jitter, lack of physical surfaces for hand stabilization, and absence of tactile confirmation. These limitations restrict the types of applications suitable for VR and influence interaction design toward larger, more forgiving target sizes and simplified manipulation schemes.
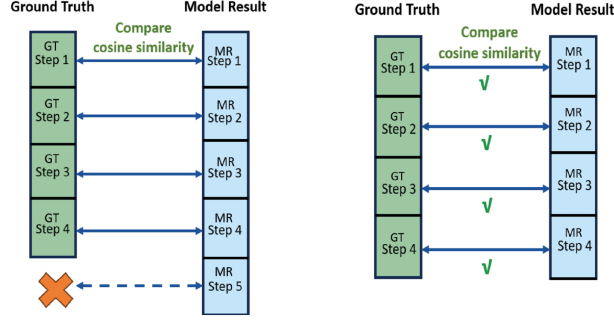


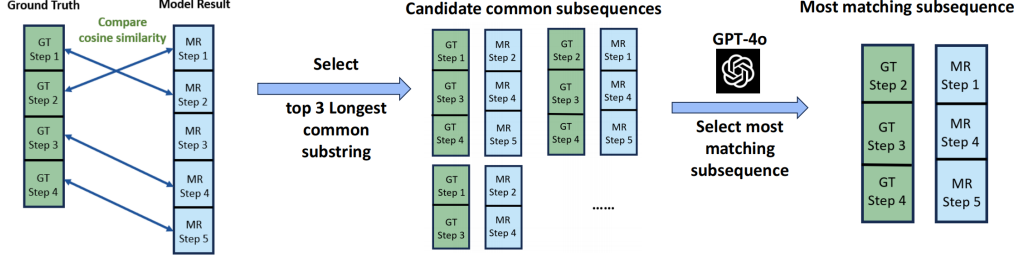Figure 2: Overview of Strict Step-by-Step Matching (SSM) Calculation



Figure 3: Overview of Common Subsequence Evaluation

## B  EXPLANATION OF EVALUATION METRICS

### B.1  STRICT STEP-BY-STEP MATCHING (SSM)

Figure 2 illustrates the Strict Step-by-Step Matching (SSM) calculation process. SSM represents our most stringent evaluation metric, requiring exact correspondence between model-generated sequences and ground truth annotations. The calculation process operates as follows:

In the left panel, we observe a scenario where the ground truth contains 4 steps while the model result contains 5 steps. For SSM to register a match, two conditions must be satisfied: (1) the number of steps must be identical between ground truth and model output, and (2) each step must have a cosine similarity score above our threshold of 0.8387 with its corresponding ground truth step. In this example, the length mismatch alone disqualifies the sequence from being counted as correct, resulting in an SSM score of 0. The orange X symbol on the fifth model step visually indicates this length mismatch failure.

The right panel demonstrates a successful SSM match where both sequences contain 4 steps. Each model step is compared with its corresponding ground truth step using cosine similarity of their text embeddings. The green checkmarks indicate that all four step pairs exceed the similarity threshold, resulting in a successful match and contributing 1 to the SSM score. This metric's strictness explains

why even high-performing models achieve relatively low SSM scores—any deviation in sequence length or individual step similarity results in complete failure for that sequence.

### B.2 Common Subsequence Evaluation

Figure 3 details our Common Subsequence Evaluation approach, which underlies the Normalized Step Alignment Score (NSAS) and Sequential Order Preservation (SOP) metrics. This evaluation method provides more nuanced assessment than SSM by identifying partial matches and preserved ordering within sequences.

The process begins with comparing each step in the ground truth and model result sequences using cosine similarity, as shown by the crossing blue lines in the leftmost panel. Unlike SSM's strict position-based matching, this approach allows steps to match regardless of their positions in the sequences. The algorithm then identifies the top 3 longest common subsequences where matched steps maintain their relative ordering.

In the example shown, multiple candidate subsequences are generated, each representing different ways steps from both sequences can be aligned while preserving order. The model (shown as GPT-4o) then selects the most matching subsequence based on the highest cumulative similarity scores. The final selected subsequence shows GT Steps 2, 3, and 4 matching with MR Steps 1, 4, and 5 respectively. This flexible matching approach allows the metrics to capture semantic correctness even when models include additional steps or present steps in slightly different positions.

The NSAS metric is calculated by considering the correctly matched steps ($|C|$), missing steps from ground truth ($|M|$), and additional steps in the model output ($|A|$), normalized by the total ground truth steps and scaled across the dataset. The SOP metric specifically evaluates whether matched steps maintain their sequential order, providing insight into the model's procedural reasoning capabilities.

## C Detailed Experiment Results

This section provides comprehensive analysis of our experimental results, including detailed performance breakdowns across models, games, and experimental conditions. We present both aggregated metrics and fine-grained analyses that illuminate specific strengths and weaknesses in current LLMs' ability to reason about VR device manipulations.

### C.1 Overall Performance Analysis

The table 5 below presents a holistic view of model performance across all experimental conditions. The results reveal a clear performance hierarchy, with Gemini-1.5-Pro achieving the highest average Normalized Step Alignment Score (NSAS) of 0.845, followed closely by GPT-4o (0.832) and GPT-4 (0.824). Notably, even the best-performing models achieve relatively modest Strict Step-by-Step Matching (SSM) scores, with Gemini-1.5-Pro reaching only 8.7% exact sequence matches. This discrepancy between NSAS and SSM scores indicates that while models can identify appropriate actions, they struggle with precise sequencing and complete reproduction of manipulation sequences.

The Sequential Order Preservation (SOP) scores reveal perhaps the most significant challenge facing current LLMs. Even top-performing models achieve SOP scores below 0.3, indicating difficulty in maintaining correct procedural ordering of steps. This limitation is particularly pronounced in zero-shot settings, where SOP scores approach zero for most models, suggesting that procedural reasoning for VR interactions requires exposure to examples rather than emerging from general language understanding.

Human performance provides an important baseline for contextualizing model achievements. While humans achieve comparable NSAS scores (0.817) to top LLMs, they show notably lower SOP scores (0.124) than leading models. This counterintuitive result reflects the challenging nature of the tasks even for experienced VR users and suggests that perfect procedural recall may be less important than adaptive problem-solving in real-world VR interaction.

16

Table 5: Performance of LLMs across VR Games (Best Few-Shot Setting)

| Model | NSAS | SOP | SSC | SSM | Best FS |
|---|---|---|---|---|---|
| Gemini-1.5-Pro | 0.845 | 0.251 | 0.151 | 0.087 | 5 |
| GPT-4o | 0.832 | 0.291 | 0.190 | 0.135 | 5 |
| GPT-4 | 0.824 | 0.218 | 0.177 | 0.095 | 5 |
| LLaMA-3-8B | 0.823 | 0.283 | 0.200 | 0.040 | 5 |
| Human | 0.817 | 0.124 | 0.174 | 0.021 | - |
| Mixtral-8x7B | 0.790 | 0.123 | 0.142 | 0.039 | 5 |
| GPT-3.5 | 0.778 | 0.169 | 0.137 | 0.037 | 5 |
| GLM-4-Flash | 0.749 | 0.096 | 0.165 | 0.000 | 5 |

## C.2 GAME-SPECIFIC PERFORMANCE PATTERNS

The table 6 below reveals substantial variations in model performance across different VR games, highlighting how game design and interaction complexity influence LLM reasoning capabilities. Vivecraft consistently yields the highest performance across all models, with NSAS scores ranging from 0.909 to 0.938. This strong performance likely reflects the game's discrete, block-based interaction paradigm inherited from Minecraft, which provides clear action-object mappings that align well with linguistic descriptions.

In contrast, Into the Radius proves most challenging, with NSAS scores dropping to 0.618-0.698 across models. This game's emphasis on realistic physics simulation, complex inventory management, and weapon manipulation requires understanding of nuanced spatial relationships and multi-step procedures that current LLMs struggle to capture. The high standard deviation in performance (0.135 for GLM-4-flash) indicates inconsistent model behavior when confronting complex interaction scenarios.

Half-Life: Alyx and Moss: Book II occupy intermediate positions in the difficulty spectrum. Half-Life: Alyx's physics-based puzzles and combat scenarios require precise timing and spatial reasoning, reflected in extremely low SOP scores (0.022 for GPT-4o). Moss: Book II's third-person perspective and puzzle-platforming elements introduce unique challenges in translating camera-relative directions into controller movements, though models show more consistent performance than in Half-Life: Alyx.

Table 6: Performance comparison across different VR games (5-shot setting). We report NSAS scores (primary metric) and SOP scores (in parentheses).

| Model | Half-Life: Alyx | Radius | Moss | Vivecraft |
|---|---|---|---|---|
| GPT-3.5-turbo | 0.858 (0.123) | 0.662 (0.169) | 0.782 (0.169) | 0.922 (0.043) |
| GPT-4-turbo | 0.852 (0.125) | 0.693 (0.189) | 0.824 (0.218) | 0.927 (0.137) |
| GPT-4o | 0.804 (0.022) | 0.698 (0.291) | 0.824 (0.300) | 0.931 (0.190) |
| Gemini-1.5-Pro | 0.863 (0.209) | 0.682 (0.102) | 0.848 (0.265) | 0.938 (0.250) |
| Mixtral-8x7B | 0.839 (0.126) | 0.666 (0.123) | 0.756 (0.117) | 0.926 (0.060) |
| LLaMA-3-8B | 0.848 (0.126) | 0.644 (0.242) | 0.823 (0.283) | 0.929 (0.039) |
| GLM-4-flash | 0.836 (0.076) | 0.618 (0.096) | 0.749 (0.087) | 0.909 (0.000) |
| Human | 0.845 (0.090) | 0.684 (0.148) | 0.817 (0.112) | 0.935 (0.122) |

## C.3 IMPACT OF FEW-SHOT LEARNING

The table 7 below demonstrates the transformative effect of few-shot examples on model performance. The most dramatic improvements occur in SOP scores, which increase by factors of 10-20x from zero-shot to 5-shot settings. GPT-3.5-turbo exemplifies this pattern, improving from 0.036 to 0.226 in SOP F1 score, representing a 527.8% relative gain. This massive improvement suggests that examples primarily help models understand the expected format and level of detail for procedural instructions rather than teaching fundamental VR interaction principles.

The diminishing returns pattern is consistent across models, with the largest gains occurring between zero-shot and 1-shot conditions. The jump from 3-shot to 5-shot provides minimal additional benefit,

Table 7: Impact of few-shot examples on model performance. We report F1 scores for Sequential Order Preservation (SOP) across different number of examples. Higher is better.

| Model | Zero-shot | 1-shot | 3-shot | 5-shot | Relative Gain (%) |
|---|---|---|---|---|---|
| GPT-3.5-turbo | 0.036 | 0.096 | 0.190 | 0.226 | 527.8 |
| GPT-4-turbo | 0.102 | 0.171 | 0.254 | 0.301 | 195.1 |
| GPT-4o | 0.112 | 0.224 | 0.257 | 0.287 | 156.3 |
| Gemini-1.5-Pro | 0.085 | 0.187 | 0.260 | **0.330** | **288.2** |
| Mixtral-8x7B | 0.031 | 0.110 | 0.204 | 0.201 | 548.4 |
| LLaMA-3-8B | 0.090 | 0.165 | 0.254 | 0.299 | 232.2 |
| GLM-4-flash | 0.033 | 0.069 | 0.121 | 0.146 | 342.4 |

indicating that models quickly extract relevant patterns from limited examples. Gemini-1.5-Pro shows the most efficient few-shot learning, achieving top performance with fewer examples than competing models, suggesting superior in-context learning capabilities for procedural tasks.

Interestingly, few-shot examples have differential effects across game types. Complex games like Into the Radius show continued improvement with additional examples, while simpler environments like Vivecraft plateau quickly. This pattern indicates that few-shot learning is most beneficial when dealing with diverse interaction patterns and complex procedural sequences.

## C.4 Cognitive Capability Analysis

The figure 1 shows model performance across six cognitive dimensions, revealing distinct capability profiles. All models demonstrate strong task decomposition abilities (7.8-8.5), indicating that breaking down high-level goals into subtasks aligns well with LLMs' training on hierarchical text structures. Gemini-1.5-Pro leads in this dimension with a score of 8.5, though even smaller models like Mixtral-8x7B achieve respectable scores of 8.0.

Motor action mapping emerges as the most challenging capability across all models (0.5-4.5), highlighting the difficulty of translating abstract action concepts into specific button presses and controller movements. This limitation likely stems from the absence of embodied experience in text-based training data. GPT-4o performs best in this dimension but still falls far short of human-level capability, suggesting a fundamental gap in current architectures.

Procedural reasoning shows high variance across models (2.3-7.0), with Gemini-1.5-Pro again leading. The correlation between procedural reasoning scores and few-shot learning gains suggests that this capability can be partially addressed through examples, though the ceiling remains well below human performance. Spatial reasoning capabilities (4.8-7.5) reveal another significant gap, particularly evident in games requiring 3D navigation and object manipulation.

## C.5 Statistical Significance and Variance Analysis

The tables 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, and 19 below provide detailed statistical analyses of model performance, revealing important patterns in consistency and reliability. And the figures 4, 5, 6 The standard deviation measurements across different games and shot settings illuminate which models maintain stable performance versus those exhibiting high variability. For instance, in Vivecraft, GPT-3.5-turbo shows remarkably consistent NSAS scores in zero-shot settings (std = 0.0248), but this consistency deteriorates with few-shot examples (std = 0.0734 at 3-shot), suggesting that additional examples introduce uncertainty in the model's approach to task completion.

The variance patterns differ significantly between metrics. NSAS scores generally show lower standard deviations (0.02-0.21 range) compared to SOP scores (0.00-0.34 range), indicating that models more consistently identify relevant steps than maintain proper ordering. This pattern is particularly pronounced in complex games like Into the Radius, where SOP standard deviations exceed 0.3 for several models in few-shot settings. Such high variance suggests that models employ different strategies across different runs, sometimes achieving correct ordering by chance rather than through systematic understanding.

Comparison with human variance provides crucial context for interpreting model stability. Human annotators show standard deviations comparable to mid-tier models (0.084 in cross-game performance), suggesting that some degree of variance is inherent to the task rather than a model limitation. However, humans maintain more consistent SOP performance (std = 0.029) compared to all models except Mixtral-8x7B, indicating more reliable procedural reasoning despite overall lower scores.

Table 8: Average and standard deviation of Normalized Step Alignment Score (NSAS) scores comparison of LLMs on *Vivecraft* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-shot | 0.9258 | 0.0248 | 0.9255 | 0.0238 | 0.9191 | 0.0306 | 0.9209 | 0.0334 | 0.9312 | 0.0207 | 0.9244 | 0.0329 |
| 1-shot | 0.921 | 0.0309 | 0.9349 | 0.0506 | 0.9358 | 0.0735 | 0.9362 | 0.0553 | 0.9219 | 0.0636 | 0.9101 | 0.0765 |
| 3-shot | 0.9284 | 0.0734 | 0.914 | 0.1167 | 0.9212 | 0.1115 | 0.9381 | 0.0781 | 0.9005 | 0.1125 | 0.9022 | 0.1051 |
| 5-shot | 0.9218 | 0.0385 | 0.9274 | 0.0674 | 0.9305 | 0.0689 | 0.9378 | 0.0708 | 0.9256 | 0.0477 | 0.9289 | 0.0364 |

Table 9: Average and standard deviation of Sequential Order Preservation (SOP) scores comparison of LLMs on *Vivecraft* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0029 | 0.0312 | 0.0007 | 0.0078 | 0.0012 | 0.0125 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0059 | 0.0624 |
| 1-shot | 0.015 | 0.0734 | 0.1203 | 0.2157 | 0.1568 | 0.2337 | 0.1794 | 0.291 | 0.0812 | 0.164 | 0.0351 | 0.1215 |
| 3-shot | 0.1302 | 0.2352 | 0.1143 | 0.2136 | 0.2826 | 0.3417 | 0.2335 | 0.3417 | 0.0986 | 0.2024 | 0.1124 | 0.2026 |
| 5-shot | 0.0395 | 0.1226 | 0.1366 | 0.2388 | 0.1837 | 0.278 | 0.2495 | 0.3358 | 0.0553 | 0.158 | 0.0374 | 0.1371 |

Table 10: Average and standard deviation of Semantic Step Coverage (SSC) scores comparison of LLMs on *Vivecraft* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.049 | 0.11 | 0.1301 | 0.1443 | 0.1272 | 0.1511 | 0.2221 | 0.2151 | 0.024 | 0.0999 | 0.1088 | 0.1549 |
| 1-shot | 0.1274 | 0.1988 | 0.544 | 0.3672 | 0.6598 | 0.3322 | 0.5747 | 0.3359 | 0.4914 | 0.3617 | 0.3165 | 0.3415 |
| 3-shot | 0.4755 | 0.3526 | 0.6486 | 0.3373 | 0.6817 | 0.3204 | 0.6538 | 0.3373 | 0.5414 | 0.37 | 0.5299 | 0.3785 |
| 5-shot | 0.18 | 0.2337 | 0.5035 | 0.3772 | 0.6183 | 0.3416 | 0.608 | 0.3546 | 0.3579 | 0.3555 | 0.1606 | 0.2605 |

Table 11: Average and standard deviation of Normalized Step Alignment Score (NSAS) scores comparison of LLMs on *Half-Life: Alyx* under different shot settings

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.838 | 0.0413 | 0.8456 | 0.0366 | 0.8376 | 0.0424 | 0.8447 | 0.032 | 0.8376 | 0.0331 | 0.848 | 0.0317 |
| 1-shot | 0.8354 | 0.0582 | 0.8427 | 0.0489 | 0.8472 | 0.0629 | 0.8627 | 0.0482 | 0.807 | 0.1099 | 0.8131 | 0.1289 |
| 3-shot | 0.8452 | 0.0551 | 0.845 | 0.0467 | 0.838 | 0.0757 | 0.8701 | 0.0603 | 0.8255 | 0.0819 | 0.8449 | 0.0707 |
| 5-shot | 0.8577 | 0.0773 | 0.8523 | 0.0613 | 0.8039 | 0.0694 | 0.8625 | 0.0691 | 0.8394 | 0.0834 | 0.848 | 0.0976 |

Table 12: Average and standard deviation of Sequential Order Preservation (SOP) scores comparison of LLMs on *Half Life: Alyx* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0098 | 0.0802 | 0.0252 | 0.1265 | 0.0396 | 0.1745 | 0.0082 | 0.0669 | 0.0019 | 0.0158 | 0.0123 | 0.0704 |
| 1-shot | 0.0447 | 0.0764 | 0.0402 | 0.1224 | 0.024 | 0.0733 | 0.0198 | 0.1263 | 0.0425 | 0.0816 | 0.0447 | 0.0967 |
| 3-shot | 0.0725 | 0.1159 | 0.0312 | 0.0725 | 0.0701 | 0.1261 | 0.1349 | 0.2187 | 0.0703 | 0.1094 | 0.087 | 0.1687 |
| 5-shot | 0.123 | 0.1834 | 0.1248 | 0.2382 | 0.0216 | 0.0809 | 0.2089 | 0.2938 | 0.1257 | 0.2409 | 0.1259 | 0.2385 |

Table 13: Average and standard deviation of Semantic Step Coverage (SSC) scores comparison of LLMs on *Half Life: Alyx* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0785 | 0.1843 | 0.2231 | 0.2089 | 0.2424 | 0.2111 | 0.1989 | 0.1982 | 0.0716 | 0.1187 | 0.1662 | 0.172 |
| 1-shot | 0.2562 | 0.235 | 0.3485 | 0.2336 | 0.4184 | 0.2413 | 0.3859 | 0.2654 | 0.3256 | 0.1934 | 0.3872 | 0.2058 |
| 3-shot | 0.3072 | 0.2444 | 0.3648 | 0.2414 | 0.5611 | 0.229 | 0.5494 | 0.2887 | 0.3544 | 0.2202 | 0.4599 | 0.2371 |
| 5-shot | 0.425 | 0.2814 | 0.6127 | 0.2856 | 0.6934 | 0.2359 | 0.6299 | 0.315 | 0.4642 | 0.2957 | 0.5152 | 0.2708 |

Table 14: Normalized Step Alignment Score (NSAS) scores comparison of LLMs on *Moss: Book II* under different shot settings

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.7819 | 0.0403 | 0.8055 | 0.0717 | 0.7871 | 0.0596 | 0.7994 | 0.0548 | 0.7913 | 0.0595 | 0.7916 | 0.0572 |
| 1-shot | 0.776 | 0.0616 | 0.7993 | 0.0771 | 0.803 | 0.0924 | 0.8139 | 0.0778 | 0.7663 | 0.0793 | 0.7938 | 0.0763 |
| 3-shot | 0.7776 | 0.0889 | 0.818 | 0.0925 | 0.8016 | 0.1242 | 0.8302 | 0.0935 | 0.7613 | 0.1341 | 0.7895 | 0.1371 |
| 5-shot | 0.782 | 0.0952 | 0.8243 | 0.102 | 0.8237 | 0.1092 | 0.8478 | 0.1017 | 0.756 | 0.1469 | 0.8232 | 0.105 |

Table 15: Average and standard deviation of Sequential Order Preservation (SOP) scores comparison of LLMs on *Moss: Book II* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0091 | 0.0581 | 0.0263 | 0.1533 | 0.0113 | 0.0494 | 0.0197 | 0.1029 | 0.0052 | 0.033 | 0.0 | 0.0 |
| 1-shot | 0.034 | 0.1568 | 0.0663 | 0.2044 | 0.1084 | 0.2486 | 0.1145 | 0.2495 | 0.0739 | 0.1721 | 0.0578 | 0.1486 |
| 3-shot | 0.1581 | 0.242 | 0.1678 | 0.2522 | 0.2324 | 0.3185 | 0.2272 | 0.3324 | 0.1351 | 0.252 | 0.2584 | 0.3089 |
| 5-shot | 0.1686 | 0.244 | 0.2182 | 0.2801 | 0.2998 | 0.3062 | 0.2652 | 0.3596 | 0.1169 | 0.247 | 0.2831 | 0.3097 |

Table 16: Average and standard deviation of Semantic Step Coverage (SSC) scores comparison of LLMs on *Moss: Book II* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0715 | 0.1792 | 0.2491 | 0.2844 | 0.2313 | 0.2567 | 0.1763 | 0.221 | 0.0407 | 0.0991 | 0.1208 | 0.1779 |
| 1-shot | 0.0748 | 0.1719 | 0.259 | 0.2771 | 0.3682 | 0.3018 | 0.3449 | 0.3396 | 0.1749 | 0.2393 | 0.2319 | 0.293 |
| 3-shot | 0.3349 | 0.3069 | 0.4593 | 0.3309 | 0.5001 | 0.3444 | 0.5238 | 0.3738 | 0.3207 | 0.3105 | 0.4689 | 0.3399 |
| 5-shot | 0.3737 | 0.3213 | 0.4951 | 0.3373 | 0.5562 | 0.3319 | 0.6091 | 0.3476 | 0.2974 | 0.3385 | 0.4567 | 0.3173 |

Table 17: Average and standard deviation of Normalized Step Alignment Score (NSAS) scores comparison of LLMs on VR game *Into the Radius* under different shot settings

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.6165 | 0.0755 | 0.6408 | 0.0955 | 0.5939 | 0.1306 | 0.6492 | 0.1018 | 0.6644 | 0.0718 | 0.6447 | 0.0882 |
| 1-shot | 0.641 | 0.1177 | 0.6519 | 0.1421 | 0.6282 | 0.1687 | 0.6875 | 0.1159 | 0.6285 | 0.1684 | 0.6285 | 0.1346 |
| 3-shot | 0.6305 | 0.128 | 0.6802 | 0.1645 | 0.6491 | 0.2057 | 0.6634 | 0.1638 | 0.618 | 0.1633 | 0.6479 | 0.1606 |
| 5-shot | 0.6621 | 0.1291 | 0.6927 | 0.1721 | 0.6984 | 0.2136 | 0.6818 | 0.1191 | 0.666 | 0.1495 | 0.6443 | 0.211 |

Table 18: Average and standard deviation of Sequential Order Preservation (SOP) scores comparison of LLMs on *Into the Radius* under different shot settings.

| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0091 | 0.0581 | 0.0263 | 0.1533 | 0.0113 | 0.0494 | 0.0197 | 0.1029 | 0.0052 | 0.033 | 0.0 | 0.0 |
| 1-shot | 0.034 | 0.1568 | 0.0663 | 0.2044 | 0.1084 | 0.2486 | 0.1145 | 0.2495 | 0.0739 | 0.1721 | 0.0578 | 0.1486 |
| 3-shot | 0.1581 | 0.242 | 0.1678 | 0.2522 | 0.2324 | 0.3185 | 0.2272 | 0.3324 | 0.1351 | 0.252 | 0.2584 | 0.3089 |
| 5-shot | 0.1686 | 0.244 | 0.2182 | 0.2801 | 0.2998 | 0.3062 | 0.2652 | 0.3596 | 0.1169 | 0.247 | 0.2831 | 0.3097 |

Table 19: Average and standard deviation of Semantic Step Coverage (SSC) scores comparison of LLMs on *Into the Radius* under different shot settings.

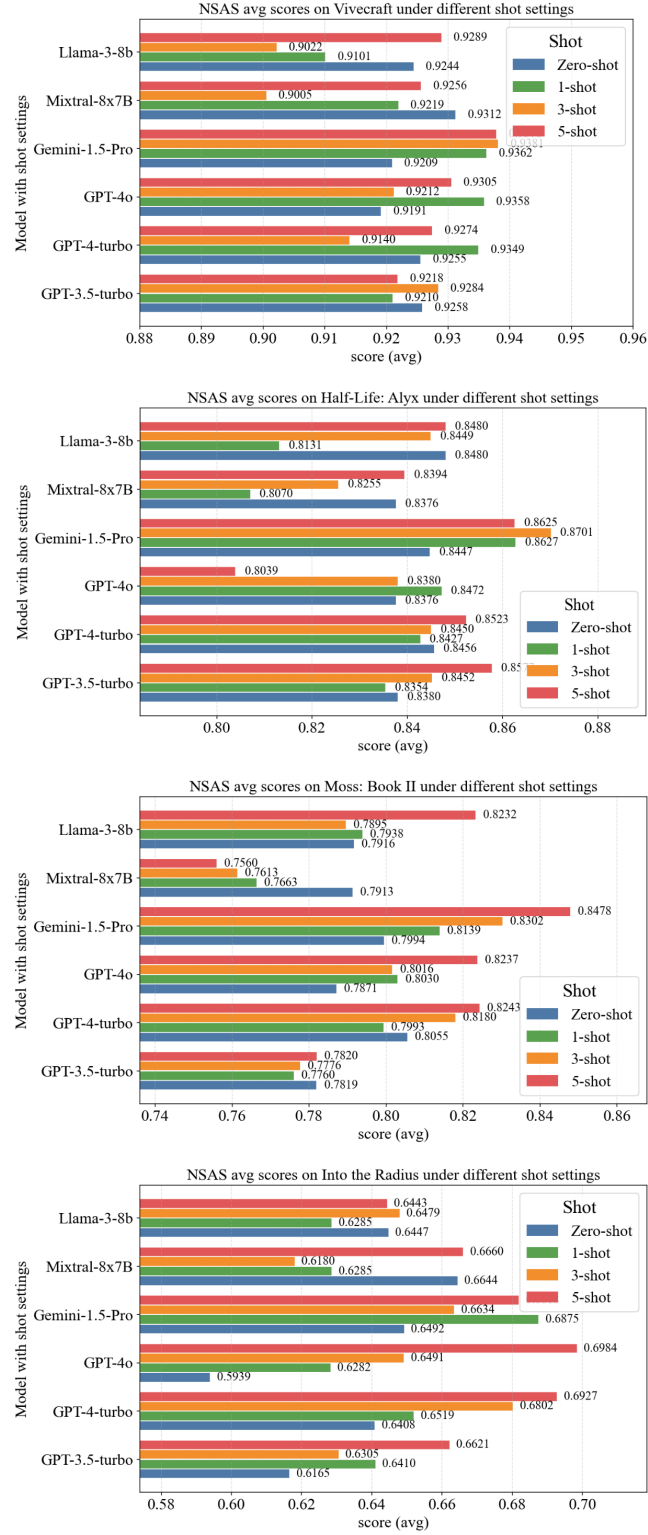| Model | GPT-3.5-turbo | | GPT-4-turbo | | GPT-4o | | Gemini-1.5-Pro | | Mixtral-8x7B | | LLaMA-3-8b | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | avg | std | avg | std | avg | std | avg | std | avg | std | avg | std |
| Zero-Shot | 0.0354 | 0.0982 | 0.1528 | 0.1774 | 0.2199 | 0.1745 | 0.1783 | 0.2107 | 0.0406 | 0.0899 | 0.1243 | 0.1655 |
| 1-shot | 0.1511 | 0.2246 | 0.304 | 0.2983 | 0.4102 | 0.2585 | 0.2823 | 0.3277 | 0.2593 | 0.3249 | 0.3171 | 0.269 |
| 3-shot | 0.2321 | 0.2713 | 0.4463 | 0.3099 | 0.5623 | 0.2976 | 0.3402 | 0.3379 | 0.3544 | 0.2621 | 0.5319 | 0.2382 |
| 5-shot | 0.3302 | 0.3115 | 0.5082 | 0.3063 | 0.6194 | 0.2877 | 0.2971 | 0.3187 | 0.285 | 0.3384 | 0.5314 | 0.2886 |

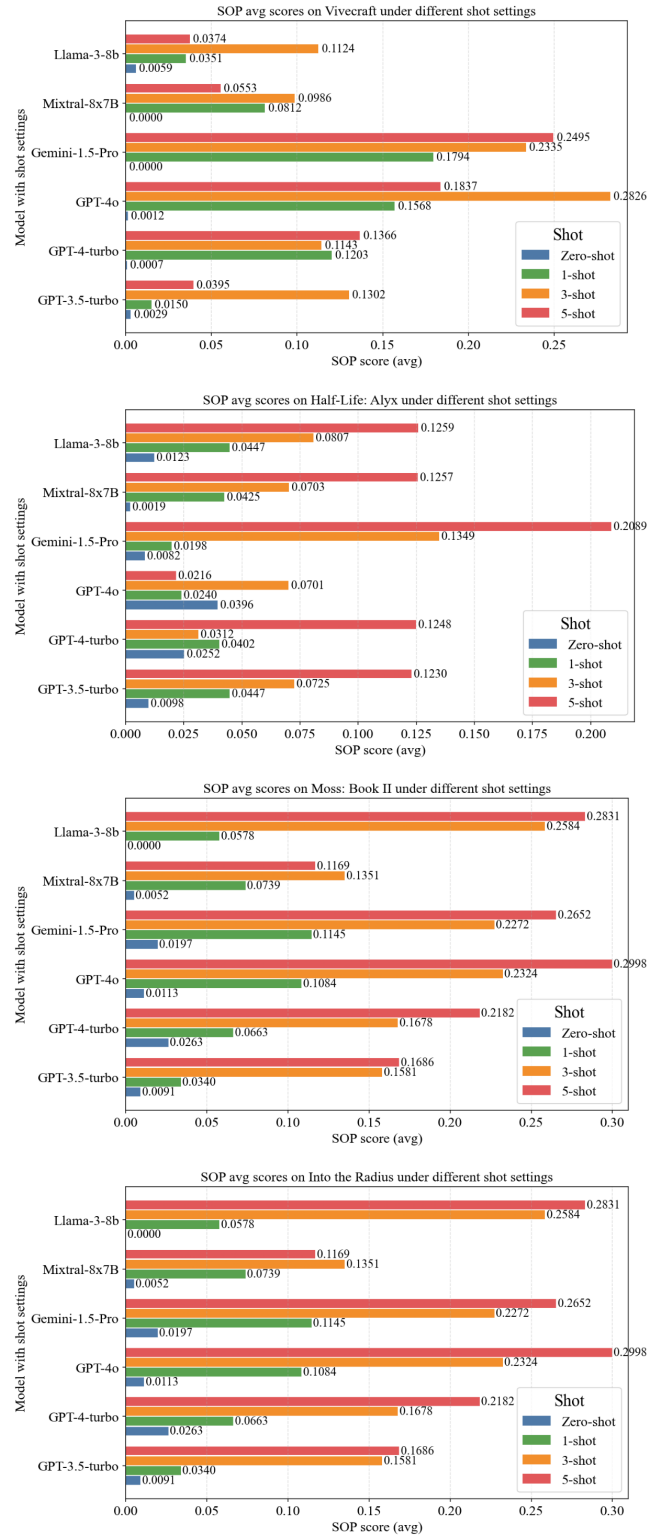Figure 4: LLMs NSAS (avg) by Different Shot Setting Across Four VR Games

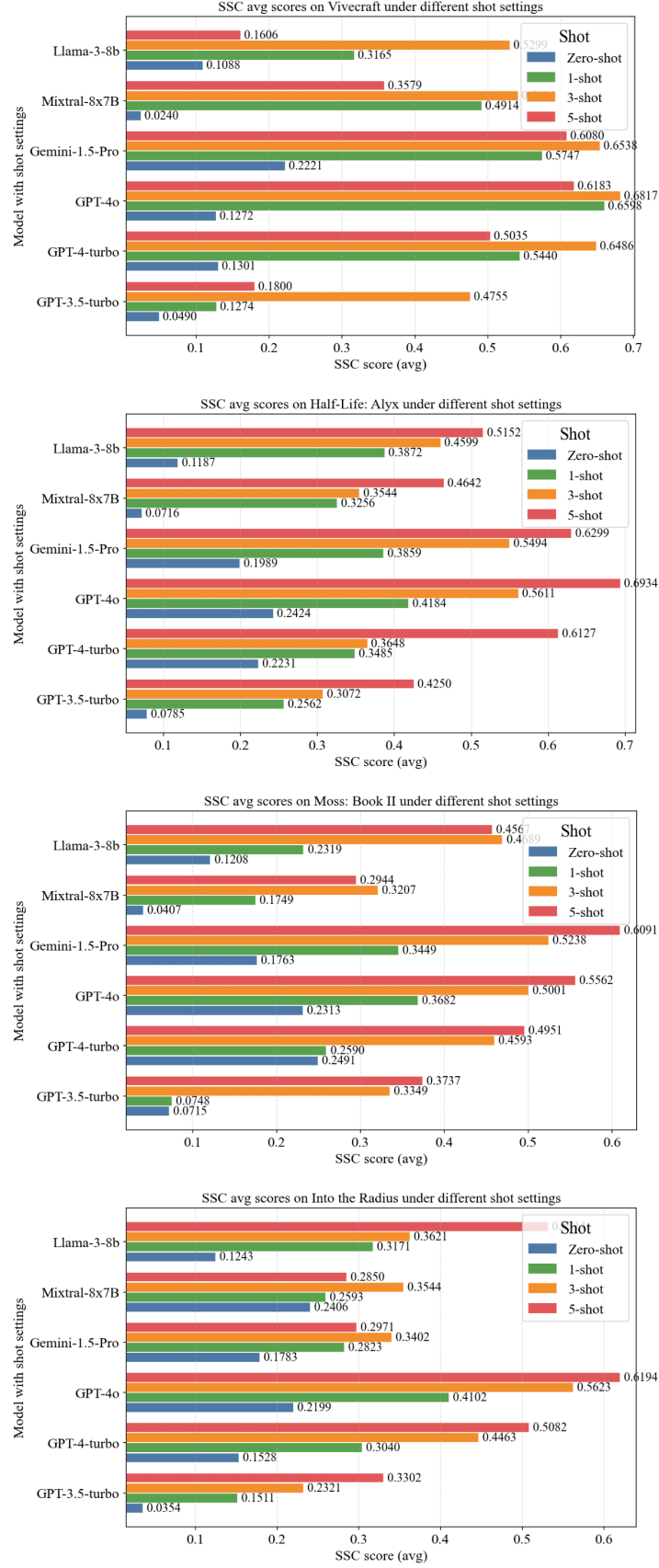Figure 5: LLMs SOP (avg) by Different Shot Setting Across Four VR Games

Figure 6: LLMs SSC (avg) by Different Shot Setting Across Four VR Games

### C.6 CROSS-GAME GENERALIZATION PATTERNS

The cross-game performance analysis reveals important insights about model generalization capabilities. Models that perform well on one game do not necessarily maintain their advantage across others. For example, while GPT-4o achieves the highest SOP score in Into the Radius (0.291), it performs poorly in Half-Life: Alyx (0.022). This game-specific variation suggests that models may overfit to particular interaction patterns rather than developing general VR manipulation capabilities.

The "Game Gap" metric in the table 3 quantifies this generalization challenge. Lower values indicate more consistent cross-game performance. Mixtral-8x7B achieves the lowest Game Gap (0.070), despite not leading in any individual game. This consistency might make it more suitable for applications requiring reliable performance across diverse VR experiences. In contrast, GPT-4o's high Game Gap (0.127) reflects its specialized strengths and weaknesses across different interaction paradigms.

Analysis of confusion patterns reveals that models struggle most when transitioning between games with different control schemes. The shift from Vivecraft's discrete block interactions to Half-Life: Alyx's continuous physics manipulation represents a fundamental change in how actions map to controller inputs. Models trained primarily on text lack the embodied experience to navigate these transitions smoothly, often applying inappropriate interaction patterns learned from one context to another.

### C.7 TEMPORAL DYNAMICS IN SEQUENTIAL TASKS

Detailed examination of step-by-step performance reveals how models handle temporal dependencies in VR interactions. Early steps in sequences generally show higher accuracy (NSAS > 0.9) across all models, with performance degrading for later steps. This degradation is particularly severe for steps that depend on the successful completion of previous actions. For instance, in a sequence like "pick up object, aim at target, throw object," models may correctly identify all three actions but fail to recognize that aiming requires successfully completing the pickup action first.

The SOP metric specifically captures these temporal dependencies, and the low scores across all models highlight a fundamental limitation in current architectures. Even with few-shot examples that demonstrate correct ordering, models struggle to internalize the causal relationships between steps. This suggests that improved performance may require architectural innovations that better capture temporal and causal reasoning, rather than simply scaling existing approaches.

Error analysis reveals common patterns in temporal mistakes. Models frequently suggest parallel actions that must be performed sequentially (e.g., "press trigger while reaching for object" when the trigger can only be meaningfully pressed after grasping). They also struggle with iterative processes, often omitting loop conditions or termination criteria. These patterns indicate that models lack an understanding of the physical constraints that govern VR interactions.

### C.8 DETAILED PERFORMANCE TABLES AND VISUALIZATIONS

The table 2 provides granular data for researchers seeking to understand specific model behaviors. These tables reveal several noteworthy patterns. First, the relationship between different metrics is non-linear. High NSAS scores do not guarantee good SOP performance, and models with similar average scores may achieve them through different strengths. This multidimensional performance landscape suggests that selecting models for specific applications requires careful consideration of which capabilities are most critical.

The table 2 illustrates the strict matching process, highlighting why SSM scores remain low even for generally capable models. The requirement for exact sequence length and step-by-step correspondence proves extremely demanding. Even minor variations in phrasing or step granularity result in match failures. This visualization helps explain why SSM may be overly strict for practical applications, where functional equivalence matters more than exact replication.

The table 2 demonstrates the more nuanced evaluation approach that underlies our NSAS and SOP metrics. By identifying the longest common subsequences with semantic matching, these metrics better capture functional understanding while still penalizing significant deviations from ground truth.

The visualization shows how models might achieve reasonable NSAS scores by identifying most relevant actions while still failing SOP evaluation due to ordering errors.

The heat maps of model performance across game-task combinations reveal clustering of difficulty. Certain task types (e.g., combat sequences in Half-Life: Alyx, inventory management in Into the Radius) consistently challenge all models, while others (e.g., block placement in Vivecraft) show near-ceiling performance. These patterns suggest that targeted improvements for specific interaction types might yield better results than general capability enhancement.

### C.9    IMPLICATIONS FOR FUTURE RESEARCH

The detailed experimental results paint a complex picture of current LLM capabilities and limitations in VR interaction reasoning. While models demonstrate competence in identifying relevant actions and decomposing high-level goals, they consistently struggle with the procedural and embodied aspects of VR interaction. The strong effect of few-shot examples suggests that current models possess latent capabilities that can be activated through appropriate prompting, but fundamental architectural limitations prevent them from achieving human-like understanding of physical manipulation sequences.

The high variance in performance across games and tasks indicates that robustness remains a significant challenge. Models that excel in one context may fail dramatically in another, limiting their practical applicability. This brittleness likely stems from the discrete nature of text-based training, which lacks the continuous, embodied experience that humans leverage when learning new physical tasks.

Moving forward, these results suggest several promising research directions. Multimodal models that incorporate visual and proprioceptive information alongside text may better capture the embodied nature of VR interactions. Explicit modeling of temporal and causal relationships could address the procedural reasoning gaps identified in our experiments. Finally, training on synthetic VR interaction data or through simulated embodiment might provide models with the experiential knowledge currently lacking in text-only approaches.

The detailed results also highlight the importance of comprehensive evaluation frameworks that assess multiple dimensions of capability. Single metrics fail to capture the complexity of VR interaction reasoning, and future benchmarks should continue to embrace multidimensional evaluation approaches that can identify specific strengths and weaknesses in model capabilities.

## D    DISCUSSION, LIMITATIONS & BROADER IMPACTS

Our investigation into LLMs' ability to translate semantic actions into VR device manipulations reveals both promising capabilities and fundamental limitations that reflect broader challenges in bridging linguistic understanding and embodied interaction. The relatively low Sequential Order Preservation (SOP) scores across all evaluated models indicate that current LLMs struggle with the temporal reasoning required for complex procedural tasks. This limitation suggests that while LLMs can identify relevant actions and understand their purposes, they lack the embodied experience necessary to accurately sequence physical manipulations.

The substantial performance variations across different VR games highlight how interaction complexity and consistency impact model performance. Games with standardized, discrete actions (like Vivecraft's block-based interactions) prove more amenable to LLM reasoning than those requiring nuanced controller movements or complex spatial reasoning (like Half-Life: Alyx). This pattern suggests that current language models may benefit from more structured representations of physical actions and explicit training on procedural sequences.

The significant improvement from few-shot examples demonstrates that LLMs possess latent capabilities for VR interaction reasoning that can be activated through appropriate prompting. However, the fact that performance plateaus with additional examples indicates fundamental architectural limitations rather than simple lack of exposure to relevant examples. This finding suggests that advances in VR-capable AI may require new training paradigms that incorporate spatial and temporal reasoning more directly.

From a broader perspective, this work carries important implications for the future of human-computer interaction and AI development. On the positive side, LLMs that can effectively reason about VR interactions could dramatically improve accessibility for users with motor impairments, enable more intuitive natural language interfaces for VR applications, and accelerate the development of intelligent tutoring systems for VR training scenarios. The potential transfer of these capabilities to robotic systems could enable more sophisticated human-robot collaboration in both virtual and physical environments.

However, we must also consider potential negative implications. As LLMs gain greater agency in controlling virtual (and potentially physical) systems, questions of safety, security, and user autonomy become paramount. The ability to translate high-level commands into detailed manipulation sequences could be exploited for unauthorized system control or social engineering attacks. Additionally, the computational resources required for training and deploying such models raise environmental concerns that must be balanced against their benefits.

The digital divide may be exacerbated as advanced VR-AI systems require substantial hardware investments and technical expertise. Ensuring equitable access to these technologies will require conscious effort from researchers, developers, and policymakers. Privacy concerns also emerge as these systems necessarily monitor and analyze detailed user movement patterns and interaction behaviors.

Moving forward, the field must pursue responsible development practices that prioritize user safety, privacy, and autonomy while advancing the technical capabilities of VR-AI systems. This includes developing robust evaluation frameworks that assess not only task performance but also failure modes, implementing transparent systems that users can understand and control, and ensuring that advances in VR interaction AI serve to augment rather than replace human agency in virtual environments.

## E  LARGE LANGUAGE MODELS USAGE STATEMENT

This work incorporated LLMs to aid in editorial refinement and linguistic improvement of the manuscript. The models provided assistance with stylistic enhancements and clarity optimization, including tasks such as rephrasing sentences and correcting grammatical errors.

We explicitly note that LLMs played no role in the conceptualization, theoretical development, or experimental design aspects of this research. The authors retain full responsibility for the entirety of the manuscript's content, including sections improved with LLM support. All LLM-assisted text has been carefully reviewed to ensure adherence to academic standards and ethical research practices.