Benchmarking Video-Language Models for Embodied Motion Cognition in Urban Open-Ended Spaces

Anonymous ACL submission

Abstract

Large multimodal models exhibit remarkable intelligence, yet their embodied cognitive abilities during motion in open-ended urban 3D space remain underexplored. We introduce a benchmark to evaluate whether video-large language models (Video-LLMs) can naturally process continuous first-person visual observations like humans, enabling recall, perception, reasoning, and navigation. We have manually control drones to collect 3D embodied motion video data from real-world cities and simulated environments, resulting in 1.5k video clips. Then we design a pipeline to generate 5.2k multiple-choice questions. Evaluations of 17 widely-used Video-LLMs reveal current limitations in urban embodied cognition. Correlation analysis provides insight into the relationships between different tasks, showing that causal reasoning has a strong correlation with recall, perception, and navigation, while the abilities for counterfactual and associative reasoning exhibit lower correlation with other tasks. We also validate the potential for Sim-to-Real transfer in urban embodiment through fine-tuning. The project is accessible at the following URL (anonymous): https://embodiedagentbenchmark. github.io/CityVideo-Bench/.

1 Introduction

001

004

800

011

012

014

018

023

029

034

039

042

Humans can process continuous first-person visual observations, enabling them to discern direction, judge distance, and navigate in the threedimensional space of the real world (Richardson et al., 2010; Burigat et al., 2017; Grauman et al., 2022; Song et al., 2023). This refers to embodied cognition in motion, which highlights that cognitive processes are deeply rooted in the body's interactions with the world (Shapiro, 2019; Newcombe et al., 2023). Naturally, endowing agents with this embodied cognition capability has been a long-term goal in the field of embodied intelligence (Fan et al., 2022; Singh et al., 2023).

In recent years, large multimodal models 043 (LMMs) (Li, 2023; Wang et al., 2024d) have emerged as a promising approach to achieve this 045 Typically, video-large language models goal. (Video-LLMs) are evaluated on capabilities such 047 as video summarization (Samel et al., 2024; Hua et al., 2024), event question answering (Wang 049 et al., 2024a), and goal localization (Yu et al., 2023). However, the benchmarks used to assess 051 these capabilities are often limited to disembodied third-person video clips, where the agent itself is static (Wu et al., 2024b; Song et al., 2024; Fang 054 et al., 2024). Besides, existing embodied video 055 understanding research (Suglia et al., 2024; Cheng et al., 2024) mainly focuses on robotic arm manip-057 ulation (Nair et al., 2022) or indoor/ground-level movement (Marcu et al., 2024; Yang et al., 2024). However, the embodied cognitive abilities required 060 for three-dimensional motion in urban open-ended 061 spaces have not been well-defined or assessed. The 062 first-person perspective visual continuous observa-063 tions generated in this scenario possess the follow-064 ing characteristics: 065

• Complex Scene and Rich Semantic Information: Urban areas are vast, containing diverse elements like skyscrapers, bridges, and tunnels that provide rich semantic information and pose comprehension and navigation challenges, while dynamic elements like pedestrians and vehicles require real-time adaptation (Yao et al., 2024; Xu et al., 2023).

067

068

069

070

071

072

073

074

075

076

077

081

• Unique Aerial Motion: Aerial navigation involves vertical mobility and a first-person perspective, adding complexity by requiring enhanced embodied cognition for processing diverse motion and observation angles, necessitating advanced spatial awareness and decisionmaking (Gao et al., 2024; Lee et al., 2024).

From these characteristics, we infer that embodied



Figure 1: Example of video-language multiple choice question-answering. This figure presents three video examples corresponding to three data sources: real cities in Guangdong Province, China; the simulator EmbodiedCity constructed based on the real city of Beijing, China; and the simulator AerialVLN built on virtual cities. To ensure logical consistency in the movement trajectories within the videos, all video clips consist of continuous perceptual observations generated from ongoing or completed vision-language navigation tasks. For each video clip, we designed different task types to evaluate the embodied cognitive intelligence of Video-LLMs.

cognition in urban open spaces poses new challenges, and assessing LMMs' embodied cognitive abilities offers insights for future urban applications (Wang et al., 2024c).

We can use drones to capture motion video in urban spaces as they navigate buildings and dynamic elements. Establishing a benchmark presents challenges: 1) Creating a task set to evaluate embodied capabilities in urban spaces. 2) Obtaining video data: Unlike most high-altitude aerial views (Li et al., 2016; Zhu et al., 2021; Wen et al., 2021), our goal is to record drones maneuvering among urban structures with flexible movement and camera angles. Drones face issues like signal loss, limited range, and crashes due to obstructions and interference, making data collection difficult and costly. 3) Designing logical and purposive motion routes to ensure coherent visual observations.

Accordingly, we introduce a benchmark designed for embodied motion cognition from embodied videos in urban airspace. Firstly, we propose a novel task set comprising 16 tasks characterized by urban spatiotemporal features, as shown in Fig-

100

101

102

104

ure 1 and Table 1. Secondly, we manually operate drones to collect embodied video data from 1) the real cities in Guangdong Province, China, 2) a simulator EmbodiedCity (Gao et al., 2024) built on the real city Beijing, China, and 3) a simulator AerialVLN (Liu et al., 2023) built on virtual cities. Using both real devices and simulators helps to rapidly increase the number of videos. The movements in the videos are intentional, directed towards navigating to a specific position within urban space or following a particular route (Wu et al., 2024c). Then, we developed a question-answer generation pipeline with trained human annotators (over 800 hours of effort) and the expertise of LMMs, generating high-quality multiple-choice questions (MCQs). Finally, we quantitatively and qualitatively evaluate widely-used Video-LLMs in zero-shot settings, including both proprietary and open-source models. We additionally attempt supervised fine-tuning (SFT) on two Video-LLMs to validate the effectiveness of our dataset.

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

Overall, the innovation of this research is the establishment of **the first benchmark for embod**-

Table 1: Task set overview. Embodied cognition in motion is divided into four abilities, each of which is manifested in several tasks. Each task is provided with corresponding handcrafted question prototypes.

	Recall
Trajectory Captioning	Summarize the agent's movement path using visual cues / landmarks.
Sequence Recall	What is the agent's next step after [changing direction to the left over the intersection]?
Object Recall	what is [located next to] the [Central ALL-STAR cafe]?
Scene Recall	Describe scene the agent observes [during the descent to a lower height near the destination].
Start/End Position	Where are the starting point and final destination of the agent's movement?
	Perception
Proximity	How does the distance between the agent and the [rooftop with solar panels in the residential area] change after the [agent descends to street level]?
Duration	Which takes longer, [the agent's movement through the skyscraper alley] or [its descent to the balcony on the 9th floor]?
Landmark Position	Given [navigation goal/route] initially, what is agent's current position relative to [landmark]?
Goal Detection	Given [navigation goal] at starting location, is the target currently visible in the field of view, and
	if so, what is its position within the view?
Cognitive Map	Summarize historical movement observations into a cognitive map.
	Reasoning
Causal	Why did the agent [perform a descent after ascending alongside the cylindrical building]?
Counterfactual	Instead of [flying over the elevated highway intersection], if the agent chooses to [fly around the
	cylindrical building], can it complete the task, and how is the alternative route?
Association	Given [navigation goal] at starting location, are there any relevant urban elements or objects in
	sight when the navigation goal is not visible?
	Navigation
Progress Evaluation	(Route-oriented vision-language navigation) Given [navigation route] at starting location, analyze
	which step the navigation is currently perform.
High-level Planning	(Goal-oriented vision-language navigation) Given [navigation goal] at starting location, make
	next plan from the current location.
Action Generation	Given [navigation goal/route] initially, generate the next control action from the current location.

ied cognition specifically tailored for motion in urban open-ended spaces:

- We propose a novel task set comprising 4 categories and 16 tasks to evaluate how Video-LLMs recall, perception, reasoning, and navigation from embodied videos.
- We consequently develop 5.2k multiple-choice questions and 1.5k video clips, derived from real world and simulated environments. The dataset generation pipeline can be extended to other embodied movement videos.
- 17 popular LMMs are evaluated and their shortcomings are analyzed. We also explored the correlation between embodied cognitive abilities and the potential of Sim-to-Real.

2 Related Work

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

Embodied Capabilities of LMMs. Embodied in-144 telligence refers to the concept that cognitive pro-145 cesses are deeply rooted in the body's interactions 146 with the world (Gupta et al., 2021; Shi et al., 2024). 147 148 Large Multimodal Models (LMMs) have demonstrated unprecedented visual understanding capabil-149 ities and are considered the "brains" for developing 150 embodied agents (Tang et al., 2023; Liang et al., 2024; Huang et al., 2023). Unlike past work (Du 152

et al., 2024; Ramakrishnan et al., 2024) that primarily focused on 2D images, static point cloud or language-based spatial understanding, human comprehension of the world is grounded in continuous visual perception (Zhang et al., 2024; Liao et al., 2024; Majumdar et al., 2023), akin to embodied cognition through video streams. Thus, we need relative video benchmarks from diverse sources to comprehensively evaluate the potential of LMMs in various embodied scenarios.

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

Video Benchmarks for LMMs. Traditional video benchmarks cover various tasks like abstract understanding and spatiotemporal analysis (Xu et al., 2017; Wu et al., 2024a; Fu et al., 2024). They mainly concentrate on understanding the video content (Wu et al., 2024b; Song et al., 2024; Fang et al., 2024), lacking exploration of Video-LLMs' embodied cognitive abilities from an embodied, egocentric perspective. While some research has focused on embodied capabilities in indoor or groundlevel scenes (Chandrasegaran et al., 2024; Mangalam et al., 2023; Sima et al., 2024; Marcu et al., 2024; Yang et al., 2024), there has been insufficient exploration of embodied abilities in urban open 3D spaces. Part of above-mentioned video benchmarks are shown in Table 2. Comparatively, we independently record embodied motion video data along with corresponding MCQs to evaluate

Table 2: The proposed and popular benchmarks for video-large language models. Our benchmark's video sources and scenarios are different from others, focusing on evaluating the embodied cognitive abilities of Video-LLMs related to urban 3D aerial motion.

Benchmark	Video Source	Video Theme	Embodied	Environment	Motion	Video Num.	QA Num.
LongVideoBench (Wu et al., 2024b)	Web-Collected	Life, Movie, News	×	1	/	3.8k	6.7k
MovieChat-1K (Song et al., 2024)	Web-Collected	Movie, TV series	×	1	1	1.0k	14.0k
MMBench-Video (Fang et al., 2024)	YouTube	Life	×	1	1	600	2.0k
HourVideo (Chandrasegaran et al., 2024)	Public Dataset	Human Activity	\checkmark	Indoor/Outdoor	Ground-Level	500	13.0k
EgoSchema (Mangalam et al., 2023)	Public Dataset	Human Activity	\checkmark	Indoor/Outdoor	Ground-Level	5.0k	5.0k
Lingoqa (Marcu et al., 2024)	Self-Recorded	Driving	\checkmark	Outdoor	Ground-Level	28.0k	419.0k
VSI-Bench (Yang et al., 2024)	3 Public Datasets	Indoor Motion	\checkmark	Indoor	Ground-Level	288	5.0k
Ours	Self-Recorded	Aerial Agent Motion	~	City	3D Aerial Space	1.5k	5.2k

models' cognitive abilities in complex, dynamic urban 3D spaces.

3 Benchmark Design and Construction

We firstly define 16 embodied tasks that assess the embodied cognitive capabilities of Video-LLMs from different four aspects. Then, we describe the dataset generation process. Finally, we provide the statistical characteristics of the dataset.

3.1 Task Set

181

182

185

186

187

189

190

192

193

194

196

198

199

200

201

204

205

206

210

212

213

215

217

218

219

221

Considering the characteristics of urban openended scenarios, embodied cognition in motion can be divided into four abilities: recall, perception, reasoning, and navigation (Chen and Dolan, 2011; Tang et al., 2021; Chandrasegaran et al., 2024). Each ability is evaluated through several specific tasks, which are outlined in Table 1.

Recall tasks evaluates Video-LLM' ability to cognitively remember key aspects of urban environments. By integrating the city's semantic elements seen in the video, the task includes Trajectory Captioning, where agents summarize their paths using landmarks. Sequence Recall asks agents to sort out the sequence of actions, while Object Recall focuses on identifying objects near landmarks. Scene Recall involves describing details observed during specific actions, ensuring comprehensive memory and understanding in dynamic urban contexts. The "Start/End Position" indicates whether agents are aware of "where I come from" and "where I am."

Perception include static relative spatial relationships (Landmark Position, Goal Detection), dynamic position changes (Proximity), temporal understanding (Duration), and scene-level comprehension (Cognitive Map). The design of these tasks encompasses a range from static to dynamic, spatial to temporal, and micro to macro levels.

Reasoning focuses on analyzing and making sense of its actions and surroundings within the urban environment. This includes understanding the causal relationships behind movements, such as why the agent descends after ascending alongside a building. It also involves counterfactual reasoning, where agents consider alternative routes, like flying around a building instead of over a highway, and assess the viability of these options. Additionally, Video-LLMs are required to recognize associations between visible urban elements and their navigation goals, even when the goals themselves are not in direct view, showcasing urban commonsense and the ability to think critically and adaptively.

222

223

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

251

252

254

255

256

257

259

260

262

Navigation tasks aims to evaluate whether Video-LLMs can plan routes and directly output actions in urban spaces (Wu et al., 2024c). We assessed two types of vision-language navigation (VLN): route-oriented VLN (Zhou et al., 2024), where the agent is provided with a navigation route (e.g., fly forward to the white building, then turn right, and continue to the lakeside to stop) and can assess the progress of the current route (Progress Evaluation) while ultimately outputting control actions (Action Generation); and goal-oriented VLN (Chaplot et al., 2020), where only a navigation goal (e.g., the lakeside) is given, allowing the agent to autonomously perform high-level navigation planning (High-level Planning) and eventually map it into control actions (Action Generation). Three-dimensional aerial navigation in urban environments is one of the most challenging tasks in embodied intelligence.

The proposed tasks reflect real-world challenges that embodied video systems in urban open-ended spaces may encounter, enhancing the practical relevance of the evaluation.

3.2 Dataset Generation Pipeline

As presented in Figure 2a, we will outline the pipeline developed for creating the dataset, which primarily consists of four steps: video curation, multiple-choice question-answering (MCQ) generation, blind filtering, and human refinement.

3.2.1 Video Curation

The primary consideration for this benchmark is obtaining massive high-quality embodied mobility



Figure 2: **Dataset Generation Pipeline and Statistics.** a) The pipeline includes four steps: video curation, MCQ generation, blind filtering, and human refinement. b) Histogram of frame count of videos. c) Histogram of path lengths. d) Histogram of word count of questions. e) Violin plot of word count of different categories of questions. f) Word cloud generated from questions and choices.

294

video data. We collecte drone flight video data in two cities, Shenzhen and Zhaoqing, in Guangdong Province, China. The data collection was conducted using two DJI Mini 4K drones. To further expand the dataset, we chose to acquire data from simulators with two city benchmarks: Embodied-City (Gao et al., 2024) and AerialVLN (Liu et al., 2023). They both have the following advantages: a) Realistic environment modeling; b) Support for aerial agents; c) Existing aerial route reference.

We divided our team into two groups. The first group designs the navigation goal/route and provides text instructions for the endpoint, ensuring that the final destination can be reached based on the instructions and common urban knowledge. The second group, consisting of experienced pilots (with over 1000 hours of real drone flight time), performs flights to perform the VLN tasks, collecting first-person perspective data during the flight. Through this approach, the movement of the agent is purposeful, and the collected video data is more conducive to evaluating the capabilities of Video-LLMs. See Appendix B.1 for details.

3.2.2 MCQ Generation

The objective of this phase is to leverage the foundational video understanding and text generation capabilities of a powerful LMM to quickly produce high-quality MCQs for each task. We have designed a Chain-of-thought (CoT) prompting based on characteristics of embodied movement videos, consisting of the following four steps: a) Narration Generation; b) Structured Extraction of Key Information; c) Role-playing; d) Providing Question/Option Templates and Examples. The detail process and prompts are in Appendix B.2. In order to improve the quality of problem generation, we use Gemini-1.5 Flash (Team et al., 2024) with video understanding capability to assist in generation at this stage.

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

325

3.2.3 Blind Filtering

The purpose of blind filtering is to eliminate MCQs that can be answered correctly based solely on urban common sense, without any video input. This process aims to evaluate the embodied cognition capabilities of Video-LLMs based on mobile video, rather than their world knowledge or common sense. Blind filtering enhances the quality of the dataset. The specific method involves using n Video-LLMs to guess the answer to any given MCQ case without video input. If all n Video-LLMs guess correctly, the MCQ is removed; otherwise, it is retained.

3.2.4 Human Refinement

The generated MCQs may contain invalid questions, ambiguous options, incorrect answers, and various other issues that require further human refinement. These issues stem from two main sources: a) Even the most advanced Video-LLMs lacks the ability to fully understand embodied movement videos. b) The urban aerial agent scenarios are complex, making video comprehension challenging. We approach the human refinement process from four aspects: a) Invalid/ambiguous

questions: For example, when the task specifies 326 for the agent to "navigate to the main entrance," 327 the presence of multiple nearby buildings leads to 328 ambiguity, as the agent is unsure which building's main entrance is intended. The navigation target should be clarified to "navigate to the main entrance 331 of the yellow building on the right." b) Urban ele-332 ment hallucination: This refers to the presence of city elements or objects in the correct option that were never actually present in the video. c) Direc-335 tion: Directional descriptions are often incorrect or 336 imprecise. d) Choices with insufficient differentia-337 tion or errors: One scenario is more than one option 338 is correct. Another is only one is correct, but the ground truth was provided incorrectly. The entire refinement process required over 800 person-hours. 341 See Appendix B.3 for more examples. 342

3.3 Dataset Statistics

The proposed dataset includes 1,547 video clips 345 with resolutions of 1280x720 pixels for real drone frames, 960x720 pixels for EmbodiedCity frames, and 520x520 pixels for AerialVLN frames. As shown in Figure 2b&c, the durations of these video clips range from 10 seconds to 10 minutes, and the trajectories of the UAV in the videos encompasses different directions and movement patterns in 3-351 dimensional spaces, thus suggesting our dataset covers diverse scenarios from brief fine movements, medium-range command execution to long-354 distance navigation. The dataset also possesses over 5.2K multi-choice questions, including low-357 level tasks such as Recall and Perception, and also high-level tasks such as Reasoning and Navigation. Seen from Figure 2d, both low-level and high-level tasks take up to approximately 50% of the total questions, enabling a comprehensive evaluation of the various capabilities required for Video-LLMs to perform embodied cognitive tasks during motion in urban open-ended spaces. For each question, the word count ranges from 50 to 250, varying based on the question category and the complexity of the environment and tasks included in the specific video, as is shown in Figure 2e. We finally generate a word cloud in Figure 2f, demonstrating the richness of urban elements included in our MCQs.

4 Experiments

371

374

We initially evaluated the performance of 16 popular Video-LLMs on various tasks related to embodied cognition in motion. Subsequently, we conducted detailed analyses focusing on the models, tasks, and video data sources. Finally, we summarized and categorized the reasons for failures across different tasks. 375

376

377

378

379

380

381

382

383

384

385

386

387

389

390

391

392

393

394

395

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

4.1 Experimental Setup

Evaluation Metric: Benefiting from the multiplechoice format of each question, we can directly calculate the accuracy for each task type as well as the overall accuracy.

Baselines: The baseline includes both proprietary and open-source Video-LLMs. For proprietary models, we used state-of-the-art models, including GPT-40, GPT-40-mini (OpenAI, 2025), Gemini-1.5 Flash, Gemini-1.5 Pro (Team et al., 2024), Gemini-2.0 Flash (Google, 2025) and Qwen-VL-Max (Cloud, 2025). For the open-source models, we focus on those capable of video/multiple image input, including LLaVA-NeXT-Video-7B (Liu et al., 2024a), Kangaroo (Liu et al., 2024b), Qwen2-VL series (Wang et al., 2024b) and InternVL2 series (Chen et al., 2024). In terms of input frame numbers, the Gemini and Qwen series allow for convenient adjustment of the input parameter fps, while others require input as *frame number*. To ensure a fair comparison while considering local computational resource constraints, the goal is to input as many frames as possible.

Others: All local model inference and finetuning is performed on three NVIDIA RTX A6000. Detail calculation of evaluation metric, baseline settings, and prompts can be found in Appendix C.

4.2 Model Comparison

The quantitative result is shown in Table 3. We can draw the following conclusions:

- Both proprietary models and open-source models exhibit relatively poor embodied cognitive abilities when navigating urban open-ended spaces. The best-performing model, Qwen-VL-Max, achieves an average accuracy of only 46.4%. This underscores the value of this benchmark, highlighting that embodied cognitive abilities in urban three-dimensional spaces have not been adequately addressed.
- Some open-source Video-LLMs outperform part of proprietary models. Specifically, models that have been optimized for video data demonstrate superior performance compared to LMMs that focus on images.

Table 3: Accuracy of different Video-LLMs on embodied tasks. The former section shows existing popular Video-LLMs' results. The latter section demonstrates fine-tuning results for two models, highlighting sim-to-real potential.

	1				Recall				Р	erceptio	on		R	easonii	ıg	N	avigatio	on
Method	Rank	Avg.	Trajectory Captioning	Sequence Recall	Object Recall	Scene Recall	Start/End Position	Proximity	Duration	Landmark Position	Goal Detection	Cognitive Map	Causal	Counterfactual	Association	Progress Evaluation	High-level Planning	Action Generation
Baseline																		
Random	-	19.7	18.5	17.0	20.8	13.5	21.8	37.8	35.6	19.7	18.0	21.9	18.2	25.0	18.3	21.8	15.9	16.4
Proprietary Models (API)																		
Gemini-1.5 Flash[1 fps]	4	40.5	39.7	51.8	61.7	79.3	61.3	47.1	59.8	37.8	28.7	47.9	60.0	42.4	20.0	43.3	32.6	34.4
Gemini-1.5 Pro[1 fps]	3	42.5	58.6	61.6	65.0	72.1	66.2	66.4	63.6	37.4	33.8	46.0	63.6	46.2	23.0	38.8	43.8	31.9
Gemini-2.0-flash[1 fps]	5	38.3	47.9	58.9	63.3	75.7	57.0	66.4	47.7	27.9	27.8	45.3	62.7	24.2	17.8	39.2	48.4	30.5
GPT-4o-mini[32f]	6	36.5	33.0	53.6	48.3	59.5	56.3	69.7	51.5	33.3	31.3	42.4	65.5	47.7	47.7	30.8	57.5	25.4
GPT-40[32f]	2	43.6	47.6	58.9	65.0	67.6	61.3	63.0	47.7	36.8	42.4	52.8	66.4	44.7	44.7	34.2	67.8	33.8
Qwen-VL-Max-latest[32f]	1	45.5	44.9	70.5	64.2	75.7	73.9	78.2	43.9	44.8	44.7	61.1	77.3	49.2	49.2	38.8	70.0	29.6
Open-source Models																		
LLaVA-NeXT-Video-7B-hf[32f]	3	38.6	55.7	39.3	43.3	61.3	40.8	58.8	52.3	49.5	16.7	26.8	44.5	20.5	58.7	36.6	52.3	19.2
Phi-3.5-vision-instruct[32f]	2	38.7	67.0	57.1	57.5	64.9	45.1	48.7	45.5	49.2	17.0	52.1	51.8	34.8	13.9	33.2	59.7	15.6
Kangaroo[64f]	1	39.2	27.0	66.1	60.8	69.4	53.5	75.6	57.6	35.5	37.2	60.0	64.5	42.4	19.1	32.5	41.9	32.4
Qwen2-VL-2B-Instruct[0.5 fps]	5	31.9	29.9	54.5	30.8	57.7	57.6	24.6	47.7	22.0	22.1	64.2	46.4	35.6	28.8	44.2	27.3	27.3
Qwen2-VL-7B-Instruct[0.25 fps]	4	36.2	36.5	50.9	47.5	65.8	65.7	47.2	48.5	25.1	28.4	55.8	55.5	29.5	33.9	59.3	32.7	32.7
InternVL2-2B[32f]	11	27.6	19.2	29.5	37.5	55.9	22.5	57.1	37.9	19.3	24.6	39.2	33.6	45.5	33.5	29.2	37.6	20.9
InternVL2-4B[32f]	10	28.1	19.2	37.5	33.3	62.2	24.6	66.4	42.4	23.2	26.5	32.8	36.4	35.6	29.5	32.2	22.1	28.1
InternVL2-8B[32f]	9	28.1	23.4	23.2	35.0	52.3	22.5	58.0	44.7	23.1	27.4	28.3	33.6	45.5	31.5	35.7	21.4	28.1
InternVL2-26B[32f]	8	28.3	24.1	36.6	35.0	60.4	23.2	53.8	41.7	20.0	28.4	32.8	34.5	45.5	29.0	37.6	23.4	22.8
InternVL2-40B[32f]	7	28.4	22.2	19.6	30.8	54.1	21.1	61.3	50.0	23.2	26.5	34.7	27.3	41.7	32.4	34.9	22.3	28.4
InternVL2-Llama3-76B[32f]	6	28.9	19.5	38.4	37.5	54.1	18.3	65.5	48.5	22.9	28.1	33.6	30.9	43.2	31.3	34.5	23.2	28.9
Fine-Tuning:Training set																		
InternVL2-4B(before)[32f]	4	28.0	17.5	32.9	34.2	61.5	26.9	66.2	41.7	21.0	25.3	37.0	37.3	33.3	25.8	30.8	35.0	21.4
InternVL2-4B(after)[32f]	1	31.4	22.0	34.3	40.8	60.0	24.7	73.0	53.6	21.0	44.1	51.1	31.3	42.9	34.4	36.5	35.0	19.6
InternVL2-8B(before)[32f]	3	29.4	22.4	25.7	35.5	53.8	22.6	59.5	48.8	23.2	30.0	30.4	38.8	40.5	33.3	35.9	34.2	20.8
InternVL2-8B(after)[32f]	2	31.2	21.1	35.7	42.1	61.5	23.7	74.3	51.2	21.4	42.4	52.6	29.9	38.1	34.4	36.3	35.8	19.3
Fine-Tuning:Test set																		
InternVL2-4B(before)[32f]	3	28.3	21.3	45.2	31.8	63.0	20.4	66.7	43.8	27.1	27.9	28.5	34.9	39.6	24.1	27.4	29.7	23.0
InternVL2-4B(after)[32f]	2	31.5	25.5	38.1	34.1	60.9	20.4	66.7	37.5	22.1	38.8	33.1	32.6	50.0	31.4	28.1	39.9	28.9
InternVL2-8B(before)[32f]	4	26.5	24.5	19.0	34.1	50.0	22.4	55.6	37.5	22.8	24.5	26.2	25.6	54.2	22.6	24.3	37.0	22.1
InternVL2-8B(after)[32f]	1	31.7	25.5	35.7	34.1	60.9	18.4	66.7	39.6	23.8	37.4	31.5	34.9	50.0	32.8	27.7	39.1	29.4



Figure 3: Correlation of Cognitive Abilities. A higher value indicates more similar Video-LLMs' performance on the two tasks, implicitly indicating the relationship between the two tasks.

• Smaller parameter models appear to be more unstable. For two models with equivalent average accuracy, the open-source small parameter model tends to have a lower minimum accuracy across all tasks compared to the commercial large parameter model.

4.3 Correlation of Cognitive Abilities

423

424

425

426

427

428

429

430

431

432

We tend to explore the relationships between these tasks and the underlying cognitive abilities, as shown in Figure 3. Specifically, we compute the pairwise correlations between each column (representing individual tasks) in Table 3, excluding the fine-tuning portion. The implicit assumption in this approach is that if multiple models exhibit similar performance across two tasks, it suggests that the embodied cognitive abilities needed for these tasks are similar. We can derive the following conclusions: 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

- The **causal reasoning** task exhibits a high correlation with almost all other tasks. It suggests that the **ability to understand and infer causality is fundamental to a wide range of cognitive processes**. This finding may indicate that causal reasoning is potentially a key factor in the emergence of embodied cognitive in Motion.
- **Recall**-type tasks demonstrate strong intercorrelations among themselves. These tasks all involve memory-related problems, underscoring the ability to recall information is a shared underlying requirement for these tasks, highlighting memory as a central component in cognitive task execution.
- Navigation tasks have a high correlation with



Figure 4: Average performance of 17 Video-LLMs in different video sources.

both Recall and Perception tasks. This observation aligns with prior knowledge that effective action and planning depend on robust memory and perceptual capabilities.

· Counterfactual and Association reasoning -both high-level reasoning tasks-exhibit low correlations with other task types. These tasks rely on distinct cognitive processes that are not shared with the other tasks in our analysis. This suggests that some embodied cognitive abilities may operate independently rather than as components of a general intelligence framework. Therefore, when tasks involve these two high-level abilities, targeted training is necessary.

4.4 Sim-to-Real

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

490

We average the performance of all Video-LLMs and listed the results of four categories of cogni-472 tive abilities across three different data sources, as 473 shown in Figure 4. From the model performance 474 perspective, there is no significant distributional 475 difference among the various video sources over-476 all. Embodied research has long suffered from a 477 lack of real-world data. Here, we used data from 478 EmbodiedCity and AerialVLN as the training set, 479 and real-world data as the test set. We employed 480 LoRA (Hu et al., 2021) to fine-tune the large mod-481 els InternVL-4B and InternVL-8B to explore the 482 potential for Sim2Real transfer (see Appendix D 483 for details on fine-tuning.). The results, as shown 484 in Table 3, indicate that both Goal-Detection and 485 high-level planning on the test set improved by 486 approximately 7% post-fine-tuning. The mean im-487 provements for the two fine-tuned models were 488 2.9% and 4.9%, respectively. 489

Error Analysis 4.5

After examining the reasoning processes of VLMs, 491 three common error types were identified, as shown 492 in Figure 5: 493



Figure 5: Three common errors in Video-LLMs when performing tasks of urban embodied cognition, along with the corresponding examples.

• Urban Elements/Scenes Understanding Error: Complex urban scenes pose challenges to perception-related tasks for Video-LLMs, resulting in insufficient alignment and urban hallucinations. This means that the models guess based solely on textual content and fail to detect urban objects that are entirely absent in the video during analysis.

494

495

496

497

498

499

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

- Motion Understanding Errors: Video-LLMs struggle to distinguish orientation and misinterpret changes in the camera's gimbal angle as vertical movement, indicating limited spatial awareness.
- Egocentric Thinking Error: Video-LLMs fail to perform complex embodied reasoning tasks, such as route planning and extrapolation.

Recall and Perception tasks, reliant on visual abilities, suffer from urban elements/scenes and agent motion understanding error. As for the complex and challenging tasks of reasoning and navigation, various errors are prevalent. See Appendix E for more details.

5 Conclusion

In this work, we propose a benchmark for embodied motion cognition in urban open spaces, comprising 1.5k video clips and 5.2k multiple-choice questions. We evaluated the performance of 17 currently popular Video-LLMs in terms of recall, perception, reasoning, and planning. The experimental results indicate that the best current Video-LLMs achieve only a 45.5% accuracy rate. Our analysis further reveals that causal reasoning is highly correlated with other tasks such as recall, perception, and planning. Fine-tuning large models with simulation data can enhance their performance on real-world embodied video tasks.

530

6

ments.

References

47(6):1029-1039.

Limitations

Our study has the following limitations:

• Dataset Scope and Generalization: Due to the

highly diversity of urban structures and forms, our dataset includes only a subset of urban scenes.

The scale of videos and multiple-choice ques-

tions (MCQs) can be further expanded to cap-

ture the vast complexity of global urban environ-

Static Evaluation Framework: The MCQ for-

mat evaluates Video-LLMs in an offline, non-

interactive setting. This does not fully reflect

real-world embodied navigation, where agents

must dynamically adapt to unforeseen obstacles

• Model Selection Bias: The rapid evolution of

Video-LLMs poses challenges for comprehensive evaluation. While we tested 17 widely-used mod-

els, new architectures are published frequently,

and our analysis may not reflect the capabilities

These limitations highlight opportunities for fu-

ture work in expanding dataset diversity, develop-

ing interactive evaluation protocols, and incorpo-

Stefano Burigat, Luca Chittaro, and Riccardo Sioni.

2017. Mobile three-dimensional maps for wayfind-

ing in large and complex buildings: Empirical com-

parison of first-person versus third-person perspec-

tive. IEEE Transactions on Human-Machine Systems,

Keshigeyan Chandrasegaran, Agrim Gupta, Lea M

Hadzic, Taran Kota, Jimming He, Cristóbal Eyza-

guirre, Zane Durante, Manling Li, Jiajun Wu, and

Li Fei-Fei. 2024. Hourvideo: 1-hour video-language

understanding. arXiv preprint arXiv:2411.04998.

Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi,

Abhinav Gupta, and Russ R Salakhutdinov. 2020.

Object goal navigation using goal-oriented semantic

exploration. Advances in Neural Information Pro-

David Chen and William B Dolan. 2011. Collecting

highly parallel data for paraphrase evaluation. In

Proceedings of the 49th annual meeting of the associ-

ation for computational linguistics: human language

cessing Systems, 33:4247-4258.

technologies, pages 190-200.

of emerging or specialized models.

rating newer models as the field evolves.

or environmental changes.

531

- 532
- 5
- 535
- 537
- 5
- 540
- 541 542
- 54
- 544
- 545
- 546 547
- 548 549
- 55

55

552 553

- 554

555

5

5

50

5

- 564 565
- 566

567 568

569

570 571

- 572 573
- 574

575 576 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.

577

578

579

580

581

583

584

585

586

587

588

589

590

592

594

595

596

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

- Sijie Cheng, Kechen Fang, Yangyang Yu, Sicheng Zhou, Bohao Li, Ye Tian, Tingguang Li, Lei Han, and Yang Liu. 2024. Videgothink: Assessing egocentric video understanding capabilities for embodied ai. *arXiv preprint arXiv:2410.11623*.
- Alibaba Cloud. 2025. Qwen documentation. https: //tongyi.aliyun.com/. Accessed: 2025-01-24.
- Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. 2024. EmbSpatial-bench: Benchmarking spatial understanding for embodied tasks with large vision-language models. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 346–355, Bangkok, Thailand. Association for Computational Linguistics.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. Advances in Neural Information Processing Systems, 35:18343– 18362.
- Xinyu Fang, Kangrui Mao, Haodong Duan, Xiangyu Zhao, Yining Li, Dahua Lin, and Kai Chen. 2024. Mmbench-video: A long-form multi-shot benchmark for holistic video understanding. *arXiv preprint arXiv:2406.14515*.
- Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.
- Chen Gao, Baining Zhao, Weichen Zhang, Jinzhu Mao, Jun Zhang, Zhiheng Zheng, Fanhang Man, Jianjie Fang, Zile Zhou, Jinqiang Cui, et al. 2024. Embodiedcity: A benchmark platform for embodied agent in real-world city environment. *arXiv preprint arXiv:2410.09604*.
- Google. 2025. Gemini api documentation. https: //ai.google.dev/gemini-api/docs. Accessed: 2025-01-24.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. 2022. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012.
- 9

Agrim Gupta, Silvio Savarese, Surya Ganguli, and Li Fei-Fei. 2021. Embodied intelligence via learning and evolution. *Nature communications*, 12(1):5721.

632

633

635

640

642

652

653

654

655

657

661

668

670

671

674

675

677

678

679

681

682

- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Hang Hua, Yunlong Tang, Chenliang Xu, and Jiebo Luo. 2024. V2xum-Ilm: Cross-modal video summarization with temporal prompt instruction tuning. arXiv preprint arXiv:2404.12353.
- Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. 2023. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*.
 - Jungdae Lee, Taiki Miyanishi, Shuhei Kurita, Koya Sakamoto, Daichi Azuma, Yutaka Matsuo, and Nakamasa Inoue. 2024. Citynav: Language-goal aerial navigation dataset with geographic information. arXiv preprint arXiv:2406.14240.
- Chunyuan Li. 2023. Large multimodal models: Notes on cvpr 2023 tutorial. *arXiv preprint arXiv:2306.14895*.
- Jing Li, Dong Hye Ye, Timothy Chung, Mathias Kolsch, Juan Wachs, and Charles Bouman. 2016. Multitarget detection and tracking from a single camera in unmanned aerial vehicles (uavs). In 2016 IEEE/RSJ international conference on intelligent robots and systems (IROS), pages 4992–4997. IEEE.
- Zijing Liang, Yanjie Xu, Yifan Hong, Penghui Shang, Qi Wang, Qiang Fu, and Ke Liu. 2024. A survey of multimodel large language models. In *Proceedings* of the 3rd International Conference on Computer, Artificial Intelligence and Control Engineering, pages 405–409.
- Ruotong Liao, Max Erler, Huiyu Wang, Guangyao Zhai, Gengyuan Zhang, Yunpu Ma, and Volker Tresp. 2024. Videoinsta: Zero-shot long video understanding via informative spatial-temporal reasoning with llms. *arXiv preprint arXiv:2409.20365*.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. Llavanext: Improved reasoning, ocr, and world knowledge.
- Jiajun Liu, Yibing Wang, Hanghang Ma, Xiaoping Wu, Xiaoqi Ma, Xiaoming Wei, Jianbin Jiao, Enhua Wu, and Jie Hu. 2024b. Kangaroo: A powerful videolanguage model supporting long-context video input. *arXiv preprint arXiv:2408.15542*.
- Shubo Liu, Hongsheng Zhang, Yuankai Qi, Peng Wang, Yanning Zhang, and Qi Wu. 2023. Aerialvln: Visionand-language navigation for uavs. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 15384–15394.

Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. 2023. Where are we in the search for an artificial visual cortex for embodied intelligence? *Advances in Neural Information Processing Systems*, 36:655–677. 687

688

690

691

692

693

694

695

696

697

698

699

700

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

718

719

720

721

722

724

725

727

728

729

730

731

733

734

735

736

737

738

739

740

741

- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. 2023. Egoschema: A diagnostic benchmark for very long-form video language understanding. Advances in Neural Information Processing Systems, 36:46212–46244.
- Ana-Maria Marcu, Long Chen, Jan Hünermann, Alice Karnsund, Benoit Hanotte, Prajwal Chidananda, Saurabh Nair, Vijay Badrinarayanan, Alex Kendall, Jamie Shotton, et al. 2024. Lingoqa: Visual question answering for autonomous driving. In *European Conference on Computer Vision*, pages 252– 269. Springer.
- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. 2022. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*.
- Nora S Newcombe, Mary Hegarty, and David Uttal. 2023. Building a cognitive science of human variation: Individual differences in spatial navigation.
- OpenAI. 2025. Openai api documentation. https: //platform.openai.com/docs/overview. Accessed: 2025-01-24.
- Santhosh Kumar Ramakrishnan, Erik Wijmans, Philipp Kraehenbuehl, and Vladlen Koltun. 2024. Does spatial cognition emerge in frontier models? *arXiv preprint arXiv:2410.06468*.
- Michael J Richardson, Kerry L Marsh, and RC Schmidt. 2010. Challenging the egocentric view of coordinated perceiving, acting, and knowing. *The mind in context*, pages 307–333.
- Karan Samel, Apoorva Beedu, Nitish Sontakke, and Irfan Essa. 2024. Exploring efficient foundational multi-modal models for video summarization. *arXiv preprint arXiv:2410.07405*.
- Lawrence Shapiro. 2019. *Embodied cognition*. Routledge.
- Haochen Shi, Zhiyuan Sun, Xingdi Yuan, Marc-Alexandre Côté, and Bang Liu. 2024. OPEx: A component-wise analysis of LLM-centric agents in embodied instruction following. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 622–636, Bangkok, Thailand. Association for Computational Linguistics.
- Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. 2024. Drivelm: Driving with graph visual question answering. In *European Conference on Computer Vision*, pages 256–274. Springer.

- 743 744 745
- 745 746 747
- 740 749 750 751 752
- 753 754
- 755 756 757
- 7 7 7
- 7
- 7
- 765 766
- 767 768
- 7
- 7
- 774 775
- 7
- 778 779 780

78

7

- 7
- 7
- 7
- 790

793 794

795 796

796 797

79 79

798 799 Kunal Pratap Singh, Jordi Salvador, Luca Weihs, and Aniruddha Kembhavi. 2023. Scene graph contrastive learning for embodied navigation. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 10884–10894.

- Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023.
 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009.
- Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. 2024. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232.
 - Alessandro Suglia, Claudio Greco, Katie Baker, Jose L Part, Ioannis Papaioannou, Arash Eshghi, Ioannis Konstas, and Oliver Lemon. 2024. Alanavlm: A multimodal embodied ai foundation model for egocentric video understanding. arXiv preprint arXiv:2406.13807.
- Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. 2023. Video understanding with large language models: A survey. *arXiv preprint arXiv:2312.17432*.
- Zongheng Tang, Yue Liao, Si Liu, Guanbin Li, Xiaojie Jin, Hongxu Jiang, Qian Yu, and Dong Xu. 2021. Human-centric spatio-temporal video grounding with visual transformers. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(12):8238– 8249.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Haibo Wang, Chenghang Lai, Yixuan Sun, and Weifeng Ge. 2024a. Weakly supervised gaussian contrastive grounding with large multimodal models for video question answering. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 5289– 5298.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Shuo Wang, David C Anastasiu, Zheng Tang, Ming-Ching Chang, Yue Yao, Liang Zheng, Mohammed Shaiqur Rahman, Meenakshi S Arya, Anuj Sharma, Pranamesh Chakraborty, et al. 2024c. The 8th ai city challenge. In *Proceedings of the*

IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7261–7272.

800

801

802

803

804

805

806

807

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

- Yuhao Wang, Yusheng Liao, Heyang Liu, Hongcheng Liu, Yanfeng Wang, and Yu Wang. 2024d. MM-SAP: A comprehensive benchmark for assessing selfawareness of multimodal large language models in perception. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9192–9205, Bangkok, Thailand. Association for Computational Linguistics.
- Longyin Wen, Dawei Du, Pengfei Zhu, Qinghua Hu, Qilong Wang, Liefeng Bo, and Siwei Lyu. 2021. Detection, tracking, and counting meets drones in crowds: A benchmark. In *CVPR*.
- Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. 2024a. Star: A benchmark for situated reasoning in real-world videos. *arXiv preprint arXiv:2405.09711*.
- Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. 2024b. Longvideobench: A benchmark for longcontext interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*.
- Wansen Wu, Tao Chang, Xinmeng Li, Quanjun Yin, and Yue Hu. 2024c. Vision-language navigation: a survey and taxonomy. *Neural Computing and Applications*, 36(7):3291–3316.
- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multime-dia*, pages 1645–1653.
- Fengli Xu, Jun Zhang, Chen Gao, Jie Feng, and Yong Li. 2023. Urban generative intelligence (ugi): A foundational platform for agents in embodied city environment. *arXiv preprint arXiv:2312.11813*.
- Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. 2024. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv preprint arXiv:2412.14171*.
- Fanglong Yao, Yuanchang Yue, Youzhi Liu, Xian Sun, and Kun Fu. 2024. Aeroverse: Uav-agent benchmark suite for simulating, pre-training, finetuning, and evaluating aerospace embodied world models. *arXiv preprint arXiv:2408.15511*.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490.
- Jiazhao Zhang, Kunyu Wang, Rongtao Xu, Gengze Zhou, Yicong Hong, Xiaomeng Fang, Qi Wu,

955

Zhizheng Zhang, and He Wang. 2024. Navid: Videobased vlm plans the next step for vision-and-language navigation. *arXiv preprint arXiv:2402.15852*.

- Gengze Zhou, Yicong Hong, and Qi Wu. 2024. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 38, pages 7641–7649.
 - Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. 2020. Vision-language navigation with self-supervised auxiliary reasoning tasks. In *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, pages 10012–10022.
 - Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. 2021. Detection and tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7380–7399.

A Dataset Examples

855

863

864

867

870

871

873

874

875

884

887

896

900

901

902

904

To better illustrate the proposed dataset, we provide MCQ examples and videos from real world, the EmbodiedCity simulator, and the AerialVLN simulator, encompassing all task types. They are presented in Fig. 6, Fig. 7, and Fig. 8.

B Dataset Generation Pipeline

B.1 Details of Video Curation

Selected Simulator Advantages: a) Realistic environment modeling. Both EmbodiedCity and AerialVLN are built on Unreal Engine and include various building types and streets, with more than a hundred categories of micro urban elements, which enriches the semantic information of the obtained embodied mobility video data. b) Support for aerial agents. Both simulators have built-in AirSim plugins, allowing easy control of aerial agents. c) Existing route reference. Research work on visionlanguage navigation (Zhu et al., 2020) has already been conducted in these simulators, allowing us to obtain some route coordinates and instruction data. Although most routes cannot be directly used (for example, most flying paths in AerialVLN contain many meaningless, repetitive flying actions and lack logicality), they provide some reference for our data collection.

Drone Setup: The drone has only one camera equipped with a gimbal, which can tilt from 0 degrees to 90 degrees downward. The second group consists of experienced pilots (with over 1000 hours of real drone flight time) to ensure the rationality of flight operations. Through this approach, the movement of the agent is purposeful,

and the collected video data is more conducive to evaluating the capabilities of Video-LLMs.

B.2 Details of MCQ Generation

In this work, we use LLM to generate the multiple choice questions (MCQ) in the dataset.

a) Narration Generation: Initially, the LMM is prompted to systematically generate a narration of the UAV's trajectory by combining embodied movement videos and destination instructions, based on the videos from Embodied City. Videos from the other two sources are originally combined with trajectory infomation.

b) Structured Extraction of Key Information: The LMM then extracts a list of movements and objects, providing structured text that ensures the subsequent MCQ generation is more aligned with the requirements.

c) Role-playing: In the last MCQ generation prompt, the model is given a specific context and role, enhancing its understanding of the task and adherence to instructions.

d) Providing Question/Option Templates and Examples: For each task, several templates for questions and options are provided, with sections marked for replacement. We also provide detailed task definitions and examples.

In the prompt, we first set up a scenario for the model, one where the model act as a teacher who needs to raise a series of test questions based on the given videos. This role playing trick makes it easier for us to explain the details of the question generation task onward. Then, we break the task into several sequential parts, and give the specific requirements for the sub-tasks. The sub-tasks include video understanding, question generation, answer generation, and finally structured input/output, with requirements, templates and examples. As we have 14 categories of questions focusing on different aspect of the video, it's hard to cover them all in one general instruction, so we adopt respectively written instructions for question and answer generation, each having its own task explanation, template and example output. During tests, the LLM sometimes gives random explanatory text before it raises the question, which could be an obstacle in the processing work, so we added specific instructions in the prompt to prevent the model from doing so. The main part of the prompt we use are shown in Figure 9, though most of the explanations, templates and examples for question generating are left out due to space constraints.



Figure 6: Videos collected from the real drones and MCQs examples.

In question generation, we find that the LLM still has problem understanding complicated city environments in the videos, so we introduced pregenerated structured object list and movement series to assist the model in understanding the videos. Also, this step is done by prompts shown in Figure 10 and Figure 11:

B.3 Details of Human Refinement

957

961

962

Following the automatic question generation phase, 964 we implemented a systematic human refinement process to ensure the quality and reliability of the 966 generated questions. During the refinement phase, 967 we identified and addressed various issues across the generated questions, with statistics shown in 969 970 Table 4. This process focused on four major aspects requiring manual intervention, with detailed 971 examples provided in Table 5. First, ambiguous 972 or unclear questions were identified and refined to ensure precise communication of the intended 974

query. Second, hallucinated urban elements, which were occasionally generated by the LLM but not present in the actual video content, were either corrected or removed to maintain factual accuracy. Third, directional descriptions were standardized by converting absolute directions to relative ones, consistently using the drone as the reference system. This standardization was crucial for maintaining consistency in spatial relationships across the dataset. Finally, multiple-choice options were carefully reviewed and adjusted to ensure appropriate differentiation between choices while maintaining one unambiguously correct answer. 975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

The refinement process significantly enhanced the dataset's quality by eliminating potential sources of confusion and ensuring all questions accurately reflected the video content. Human annotators were specifically trained to identify these issues and apply standardized corrections, resulting in a more robust and reliable benchmark for



Figure 7: Videos collected from the EmbodiedCity simulator and MCQs examples.

evaluating video-based navigation understanding. Table 4: Statistics of Issues in Human Refinement Phase

Issue	Count
Invalid/ambiguous questions	761
Urban element hallucination	184
Direction	446
Choices with insufficient differentiation or errors	2363

C Additional Experimental Settings

C.1 Evaluation Metric Calculation

997

998

999

1000

1001 1002

1004

1006

This study employs a multi-stage workflow to evaluate the accuracy of a model in video-based multiple-choice question answering tasks. This workflow integrates regular expressions and model inference for efficiency, complemented by manual validation to ensure robustness. The combined approach achieves precise and scalable accuracy assessment for the target task.The procedure is structured as follows:

- Answer Extraction and Correction. First, an-1007 swers are extracted from the model-generated text. Regular expressions are applied to iden-1009 tify standardized outputs (e.g., options starting 1010 with "A." or standalone alphabetic characters). 1011 For unstructured or ambiguous text, the GPT-40 1012 model is invoked to infer the most plausible op-1013 tion (A-G) based on the output context. Entries 1014 unresolved by automated methods are flagged 1015 with the special token "l**k". A manual verifi-1016 cation step is subsequently performed to correct 1017 these flagged entries, ensuring the validity of all 1018 answers. 1019
- Data Alignment and Cleaning. The corrected answers are aligned with a reference dataset containing ground-truth answers. Data consistency is ensured through format standardization (e.g., removing ".mp4" suffixes from video IDs and normalizing column names). Valid entries are filtered using an inner join operation, merging 1026



Figure 8: Videos collected from the AerialVLN simulator and MCQs examples.

model predictions with reference data based onvideo IDs and question categories.

• Statistical Analysis and Result Export. Accuracy is computed by category through group-wise comparisons of matched answers, followed by aggregation to derive the overall accuracy. Results are organized into structured tables, detailing percategory and total accuracy values. These outputs facilitate systematic performance evaluation and further analysis.

C.2 Brief Introduction on Baselines

We introduce the participating long-context LMMs as follows:

Gemini-1.5-Flash. Released on February 14, 2024, with an API service, Gemini-1.5-Flash is a model with 150 million parameters. It supports an input token limit of 1 million, an output token limit of 8192, and a maximum video duration of 1 hour. The frame rate is set to 1 fps.

Gemini-1.5-Pro. Released on February 14, 2024, with an API service, Gemini-1.5-Pro is a model with 175B parameters. It supports an input token limit of 2 million, an output token limit of 8192, and a maximum video duration of 2 hours. The frame rate is set to 1 fps.

Gemini-2.0-Flash. Released on December 11, 2024, Gemini-2.0-Flash is the latest lightweight model in the Gemini series, offering improved efficiency and a 12M context length. The input token limit is 1 million, the output token limit is 8192, and the frame rate is 1 fps.

GPT-40-mini. Released on July 18, 2024, GPT-40-mini is a compact version of GPT-40, designed for faster inference with a 64K context length. It has 8B parameters, an input token limit of 128K, an output token limit of 16384, and a frame rate of 32 f.

GPT-40. Released on May 13, 2024, with an API service, GPT-40 is the latest multimodal LMM

Prompt for Question Generation

MAIN INSTRUCTIONS:

You are the teacher of the course "Video comprehension and Spatial reasoning". To test the students, you need to make many test questions from a series of egocentric videos taken by a UAV while it flies towards a specific known destiny in a city scene.

Your goal is to raise multi choice questions about the details and spatial or temporal logic from video that satisfy lateron given requirements.

For each video, you must raise a multi choice question, that must be strictly restricted to and strongly related to the video content.

For each question, you need to raise the question itself, and then give 5 choices, labeled as A, B, C, D and E, including ONLY 1 CORRECT answer and 4 wrong answers. Specially, the "Proximity", "Duration" and "Counterfactual" questions should only contain 3 choices, including 1 correct answer and 2 wrong answers. Specially, the "Viewpoint Invariance" questions should only contain 2 choices, including 1 correct answer and 1 wrong answer. For the choices, you need to make sure that the correct answer is not too obvious, and the wrong answers are not too irrelevant. Finally, you need to give the correct answer to the question, which is one of A, B, C, D and E (or one of A, B and C, or one of A and B). You are allowed to put the correct answer at any position among the 5 choices.

All possible questions falls into 14 categories, that covers different aspects of the student's ability to understand the video content, spatial and temporal relationship and causal logic.

The categories, templates and examples of Question, Choices and Answer are given respectively as follows:

1. Captioning

TASK EXPLANATION: This question requires the student to summarize the UAV's movement route in the video by combining the specific objects along the way. Proper answer should be a concise summary of the UAV's movement route, including its starting location, final destiny and specific buildings, objects along the way.

TEMPLATE Question: Summarize the UAV's [movement route] in the video by combining the [specific objects] along the way.

TEMPLATE Choices: The UAV goes from [specific objects] to the [specific objects] by [specific movements about specific objects].

EXAMPLE OUTPUT:

Question: According to the video, which of the following choices better summarizes the UAV's movement route? Choices:

A. The UAV goes from [between two tall skyscrapers] to the [the balcony on the 31st floor] by flying straight across the opening ground below.

B. The UAV goes from [the center of the square] to the [left of the main entrance] by descending height to the rooftop.

C. The UAV goes from [the main entrance] to the [right of the main entrance] by turning left to face the main entrance.

D. The UAV goes from [the balcony on the 10th floor] to the [the center of the square] by descending to the ground level.

E. The UAV goes from [between two tall skyscrapers] to the [the balcony on the 30th floor] by flying straight across the opening ground below.

Answer: A

-----(The explanation, template and example for other categories) ------

NOTICE: The [category] and the [destiny] is known, given by textual input.

EXAMPLE INPUT:

Please raise a [category] multi choice question based on the video taken along the way towards [destiny]. For your reference, the content of the video is shown in a chronological narration, which is given as follows:

.....(narration text)...... Also, the objects and positions mentioned in the questions and correct choices should come from the video frames, or you can refer to the list below:(object list)......

The frames from the video are as follows:(video frames)......

EXAMPLE OUTPUT:

uestion:	?
hoices:	A
	В
	С
	D
	Е
nswer:	А

Answer:

С C

You output must only contain the Question, Choices and Answer in the form shown above, and any other reply, such as "Certainly, here's a question about ... " is definitely NOT ALLOWED!

The INPUT are as follows:

Please raise a {question_category} multi choice question based on the video taken along the way towards {video destiny}

For your reference, the content of the video is shown in a chronological narration, which is given as follows: {narration text}

Also, the objects and positions mentioned in the questions and correct choices should come from the video frames, or you can refer to the list below:

{object list}

1066

1068

1070

1072

1076

1077

1079

1080

1081

1082

1083

1084

1085

1087

The frames from the video are as follows: {video frames}

Figure 9: The prompt used in MCQ generation. The prompt contains scenario setup, detailed instructions, respectively written template and example output, and finally input template. Role playing and formatted structure help the LLM to better understand the user's intentions.

from OpenAI, featuring a 128K context length. It has 200B parameters, an input token limit of 128K, an output token limit of 4096, and a frame rate of 32 f.

Qwen-VL-Max-latest. Released in April 2024, Qwen-VL-Max-latest is the most advanced model in the Qwen-VL series, supporting multimodal tasks with a 128K context length. It has an input token limit of 128K, an output token limit of 8192, and a frame rate of 32 f.

LLaVA-NeXT-Video-7B-hf. Released in April 2024, LLaVA-NeXT-Video-7B-hf is a videofocused LMM with 7B parameters. In this experiment, we set the frame rate to 32 f and the output token limit to 512.

Phi-3.5-vision-instruct. Released in August 2024, Phi-3.5-vision-instruct is an upgraded version of Phi-3-Vision-Instruct, with 4.2B parameters and a 128K context length. In this experiment, we set the frame rate to 32 f and the output token limit to 512. The temperature coefficient is set to 0, and the maximum possible output is selected.

Kangaroo. Released on July 17, 2024, Kangaroo is a specialized LMM for multi-image and video understanding, with 8B parameters. Since the frame rate cannot be modified in the code, in this experiment, we used the default frame rate of 64 f and set the output token limit to 256. The temperature coefficient is set to 0, and the maximum possible output is selected.

1088

1089

1090

1091

1092

1093

1094

1095

1097

1098

1099

1100

1101

1102

1103

Qwen2-VL-2B-Instruct. Released on August 30, 2025, Qwen2-VL-2B-Instruct is a lightweight instruct-tuned model with 2B parameters and a 64K context length. In this experiment, we set the resolution to 360 * 420 and adopted the model's default settings. The frame rate is 0.5 fps, which is the maximum limit for our GPU's computational capacity. The output token limit is set to 256.

Qwen2-VL-7B-Instruct. Released on August 1104 30, 2025, Qwen2-VL-7B-Instruct is a mid-sized 1105 instruct-tuned model with 7B parameters and a 1106 128K context length. In this experiment, we set the 1107 resolution to 360 * 420 and adopted the model's 1108 default settings. The frame rate is 0.25 fps, which 1109

Prompt for Object List Generation

MAIN INSTRUCTIONS:

In the "Egocentric UAV Video Scene Recall and Object Extract" task, you are working with a chronological series of frames from an egocentric video taken by a UAV while it flies in a city scene.

Your goal is to write a chronological list of the objects and positions the UAV came across along its route shown in the video frames.

You should give the list in a chronological order, and each item should be a detailed phrase describing the object or position the UAV came across in the video frame, so that there'd be no confusion.

Especially, you should put the starting and ending position of the UAV in the list, on top and bottom respectively. A narration of movements of the UAV are given to you in textual form, to which you can refer for generation.

The output should be a list of phrases, following the format below:

EXAMPLE OUTPUT: (.json)

"The building beside the square", "The window to the right of the main gate", "The red car parked on the side of the road", "The traffic light at the closest intersection"

NOTICE: Instead of using " ' " to quote each phrases, you should use " " " to quote each phrases in order to avoid errors in the JSON format. There should be no other characters.

NOTICE: Extra explanation of the output is not allowed. Your output should be a list of phrases ONLY.

The inputs are as follows:

The destiny of the UAV is: {video_destiny} The movements of the UAV are: {narrations} The frames are as follows: {video_frames}

Figure 10: The prompt used in object list generation. The prompt contains scenario setup, detailed instructions and example output. As a simpler task than question generation, input/output template is not needed.

is the maximum limit for our GPU's computational capacity. The output token limit is set to 256.

InternVL2 series (2B, 4B, 8B, 26B, 40B, Llama3-76B). Released on July 4, 2025, the InternVL2 series offers a range of models from 2B to 76B parameters, with context lengths scaling from 32K to 256K, designed for diverse multimodal tasks. In this experiment, we set the resolution to 448 * 448 and adopted the model's default settings. The frame rate is 32 f, and the output token limit is set to 1024.

1121 C.3 Responses of Video-LLMs

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1122The following table presents the format correctness1123rates of six proprietary models through API calls.1124The probabilities indicate each model's capability1125to reliably produce outputs adhering to the specified1126format requirements. These results demonstrate the1127varying levels of format adherence among state-of-1128the-art language models.

C.4 Cost of Proprietary Models

1129

1136

In Table 7, we show the cost of our evaluation of
proprietary models. The cost of input is signifi-
cantly greater than the cost of output, suggesting
that the existing proprietary models still face the
problem of excessive token amount when receiving
video input.11301130
11311131

C.5 Prompt

For each Video-LLM, the input includes an embod-	1137
ied movement video, a single MCQ, and a prompt	1138
that introduces the background and output format.	1139

For models such as Gemini-1.5-flash, the videos1140can be uploaded to cloud storage space, the ques-1141tions are given to the model separately for evalua-1142tion, as is shown in the following prompt:1143

Prompt for Narration Generation

MAIN INSTRUCTIONS:

In the "Egocentric UAV Video Narration" task, you are working with a chronological series of frames from an egocentric video taken by a UAV while it flies towards a specific known destiny in a city scene. Your goal is to write a chronological Narration of the route of the UAV that are shown in the video frames.

For the video Narration, you should first show the [starting position] and the [final destiny] of the UAV. Then you chronologically list the specific [movements] and [positions] according to the video frames.

The [movements] may include but not limited to "goes forward to", "turns left to face", "descend to the height of", etc. The [positions] may include "the center of the square", "the balcony on the 10th floor", "left of main entrance", etc. The [destiny] is known, given by textual input.

The video frame may contain movements taken on the way to the destiny, but it is possible that the destiny is not reached in the video frame.

For each or each several frames, you should only describe the movements and positions of the UAV shown in the video frame.

EXAMPLE OUTPUT: (.json)

{

"In the video, the UAV goes from in front of the building with an antenna on the roof to the main entrance." "Movement 1: The UAV descends to the height of the 1st floor."

"Movement 2: The UAV turns downward to face the main entrance."

•••••

"Movement X: The UAV goes forward to the main entrance."

EXAMPLE EXPLAIN:

The first line summarizes the starting position and the destiny of the UAV. Each of the following lines contains a single, detailed movement of the UAV, Starting with "Movement X:". The movements are in chronological order, and each movement should be a detailed phrase describing the movement or position of the UAV in the video, so that there'd be no confusion.

The inputs are as follows:

The destiny of the UAV is: {video_destiny} The frames are as follows: {video_frames}

Figure 11: The prompt used in movement extraction. The prompt contains scenario setup, detailed instructions, respective explanation of prompt components, and finally input template.

This video presents the perception data of drones moving in the city from a first person perspective. Please answer the following questions: <video input> The template for the answer is: Option: [] (Only output one option from 'A' to 'E' here, do not output redundant content) Reason: [] (Explain why you choose this option)

For models such as GPT-40 and Qwen, the videos are given to the model one by one, together with the all the questions based on this video, and the models are then required to answer all the questions one by one. In this way, we avoid uploading

the same video repeatedly, so as to reduce time consumption and token number. Such prompt are shown below:

1150 1151 1152

1153

This video (captured into multiple frames of images as follows) presents the perception data of drones moving in the city from a first person perspective. Please answer the following questions: <Question 0>, <Question 1>, ... The template for the answer is: QA0: Option: []; Reason: []... QA1: Option: []; Reason: []... The Option only outputs one option from 'A' to 'E' here, do not output redundant content. Reason explains why you choose this option.

1148

1149

Nervice Grade (0.4) Reprint the baseline of the baseline of the second s	
Navigation Goal: [8th floor of the building (with a huge dome)], what is	the
Clerify and complete unclear or direction (spatial relationship) between the current position and the desti-	nation
Invalid/ambiguous questions charing and comprise question storms when the drone reaches the current position?	
Navigation Goal: [Rooftop of a nearby blue building]. Given the [planne	d milestone
(top of the tallest building ahead)], what should be the drone's next action	n?
Correct hallucinated elements Navigation Goal: [A store with a green sign reading "Meiyue Styling" ("	Orange Fresh")]
to accurate elements What should the drone approach next?	
A. Two buildings on sides, highway behind, parking lot below	
B. Highway adjacent, park behind, walkway below	
Remove hallucinated elements C. High-rise ahead, highway behind, no sidewalk below	
D. Roads and trees nearby, many high-rises around, parking lot below.	
E. Three buildings ahead, park with lanes in front	
Answer: D	
The drone is positioned on the right side of the highway	
Use drone as reference system (The highway is positioned on the right side of the drone).	
Direction for directions The target building is positioned on the left side of the drone	
(drone has the target building in its left field of view).	
Convert absolute directions The drone is currently located slightly southwest of the target balcony	
to relative directions (The target building is positioned on the lower left side of the drone).	
Choices:	
A. park walkway -> high-rise	
B. Park -> street -> high-rise level -> building side	
C. Street -> park -> high-rise level -> 15th floor balcony	
D. Park -> street -> high-rise level -> roof edge -> 15th floor balcony	
(Left side of domed building -> right turn-> slight descent -> horizontal	move to building)
Correct incorrect options E. High-rise -> descent -> park	
to be accurate Answer: D	
Choices:	
A. Grass ascent is longer	
Choices with insufficient B. Tree area crossing is longer	
differentiation or errors C. Equal duration	
Answer: A(B)	
Choices:	
A. Ascend	
B. (Backward after) descend	
Modify one option C. Forward (after descent)	
to contain fortual error D. Left turn	
E. Right turn	
F. Rotate camera upward	
G. Rotate camera downward	
Answer: C	

Table 5: Examples of human refinement for generated questions.

Table 6: Format Correctness Rates of Proprietary Models

Proprietary Models (API)	Rate
Gemini-1.5 Flash	0.992762
Gemini-1.5 Pro	0.983810
Gemini-2.0-flash	0.969714
GPT-4o-mini	0.912190
GPT-40	0.961143
Qwen-VL-Max-latest	0.997333

Table 7: Evaluation Cost of Different Models

Model	Input Cost	Output Cost	Total Cost
GPT-40	\$95.49	\$0.27	\$95.77
GPT-4o-mini	\$34.31	\$0.04	\$34.35
Gemini-1.5-flash	\$13.41	\$0.02	\$13.43
Gemini-1.5-pro	\$223.49	\$0.29	\$223.78
Gemini-2.0-flash	\$17.88	\$0.02	\$17.90
Qwen-VL-Max-latest	\$4.31	\$0.24	\$4.55

D Fine-Tuning

1154

In this experiment, we conducted multimodal fine-1155 tuning on the InternVL2-4B and InternVL2-8B 1156 models. During fine-tuning, the visual en-1157 coder (freeze_visual_encoder=True) and the 1158 language model backbone (freeze_llm=True) 1159 were frozen, and only the language model was 1160 lightly adapted using LoRA technology (rank 128, 1161 alpha 256, dropout 0.05). The experimental data 1162 was based on 70% of the multiple-choice ques-1163 tion samples from All_MCQ. jsonl, with the asso-1164

ciated video data totaling 39.56 GB, stored in the video_LLM/video directory. The videos were processed by a pre-trained visual encoder to extract spatiotemporal features, which were then concatenated with text embeddings and fed into the model. The maximum input sequence length was set to 8192 to accommodate the joint modeling of long video segments and complex text.

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

Training employed the AdamW optimizer (learning rate 1e-6, weight decay 0.05) combined with linear warm-up and cosine annealing scheduling. The batch size per GPU was 4, with a gradient accumulation step of 2, resulting in an effective batch size of 8. The maximum training epoch was set to 1. Video data was efficiently loaded using a distributed data loader, with dynamic mixed precision (AmpOptimWrapper) and gradient clipping (maximum norm 1) enabled. Checkpoints were saved every 64 steps to ensure training stability and storage efficiency.

E Error Cases

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205 1206 Table 8 and T/able 9 shows the results and reasoning process of an Video-LLM (Language and Vision Learning Model) answering questions on a dataset. The table has four columns:

Question Category: This column describes the type of question, such as Action Generation, High-level Planning, etc.

Question: This column contains the specific questions and the corresponding options.

Ground Truth: This column shows the correct answers to the questions, which are carefully checked in human refinement.

Video-LLM Answer: This column provides the answer given by the Video-LLM, and also the reasoning behind the Video-LLM model's choices, demonstrating the model's reasoning process.

Through these questions and their options, answers, and explanations, this table highlights the shortcomings of the Video-LLM in the four fundamental abilities: object recall, trajectory capturing, causality, and high-level planning.

Table 8: Error Example Part 1

Question Category	Question	Ground Truth	LVLM Answer
Object Recall	According to the video, what is directly above the balcony on the 24th floor? Choices: A. There is a rooftop garden above the balcony on the 24th floor. B. There is a billboard directly above the balcony on the 24th floor. C. There is another balcony directly above the balcony on the 24th floor. D. There is nothing directly above the balcony on the 24th floor. E. There is a drone stationed above the balcony on the 24th floor.	С	Option D; Reason: There is no mention or indication in the video that there is anything di- rectly above the balcony on the 24th floor.
Object Recall	Question: In the video, what is on the left side of the ""Huaji Stone Rolled Rice Rice Roll"" shop? Choices: A. Large parking lot full of cars. B. Another store with a green storefront and a plus sign that says ""Dashenlin"". C. A shop with a black car parked at the door, and the signboard reads ""NOBEE Nobeni"". D. A small, unpaved area with some overgrown vegetation. E. A bus stop with shelters and benches.	С	Option B; Reason: The video shows that on the left side of the "Hugj Stone Rolled Rice Roll" shop, there is an- other store with a green store- front and a plus sign that says "Dashenlin".
Duration	 Which takes longer, the UAV's descent towards the street level or its maneuvering between the buildings? Choices: A. The UAV's descent towards the street level takes longer than its maneuvering between the buildings. B. The UAV's descent towards the street level takes shorter than its maneuvering between the buildings. C. The UAV's descent towards the street level takes the same time as its maneuvering between the buildings. 	A	Option B; Reason: The UAV's descent towards the street level is a straightfor- ward movement, while maneu- vering between buildings re- quires more complex naviga- tion and adjustments, which would take longer.
Duration	Which takes longer, the UAV's descent to street level near the bus stop or its flight over the center of the street between skyscrapers?Choices:A. The UAV's descent to the street level takes longer than its flight over the center of the street.B. The UAV's descent to the street level takes shorter than its flight over the center of the street.C. The UAV's descent to the street level takes the same time as its flight over the center of the street.	A	Option B; Reason: The video shows that UAV's flight over the center of the street between skyscrapers takes longer than its descent to the street level near the bus stop.
Trajectory Captioning	According to the video, which of the following sentences better summarizes the movement of the drone? Choices: A. The drone ascends from the facade of the high-rise building, moves forward, turns right on the last side and descends, and reaches the top of the seaside fence. B. The drone descends from the city center to street level and crosses the main thoroughfare to the top of a skyscraper. C. The drone flew from an open street next to a building over several rooftops to a park near the ocean. D. The drone travels from a bridge over the river through numerous skyscrapers to the main city park. E. The drone turns left from above the road and follows the path along the water into the forest area.	A	Option D; Reason: The video shows the drone starting from the facade of a tall building, moving forward, turning right on the last side, and descend- ing to reach the top of the sea- side fence. This matches op- tion D.
Start/End Position	Question: According to the video, where are the starting and ending points of the drone's journey? Choice: A. The drone starts at a large intersection on a city street and ends at a store window with a blue canopy. B. The drone starts above the street and ends at a shop window with a white canopy. C. The drone starts from the fire exit in front of the building and ends at the subway entrance. D. The drone starts from a fire escape on the middle floor of the building and ends at a store window with a white canopy. E. The drone starts in the center of the street and ends in the park.	В	Option D; Reason: The video shows the drone starting from a fire escape on the middle floor of the building and end- ing at a store window with a white canopy.
Proximity	Question: After the drone moves to the outside of the park, how does the distance between the drone and the staircase-like building change? Choices: A. The distance between the drone and the staircase-like building decreases. B. The distance between the drone and the staircase-like building increases. C. The distance between the drone and the staircase-like building remains constant.	A	Option B; Reason: As the drone moves to the outside of the park, it appears to be mov- ing away from the staircase- like building, increasing the distance between them.
Scene Recall	Question: What objects did the drone observe when it hovered over the trees in the middle of the intersection? Choices: A. The drone observes road markings and phone booths as it hovers over the woods. B. The drone observes a bridge, a car, and a lamppost as it hovers over the woods. C. The drone observed skyscrapers, rooftop gardens, and flags while hovering over trees. D. The drone observed a playground, fountain, and street vendors as it hovered over the woods. E. The drone observed a church a bell tower and a water tower while hovering over the woods.	A	Option C; Reason: When the drone hovers over the trees in the middle of the intersec- tion, it observes skyscrapers, rooftop gardens, and flags.
Counterfactual	Question: If the drone does not rise and move parallel to the group of skyscrapers, but chooses to fly directly over the roofs of lower buildings, can it still reach the trees in the middle of the zebra crossing at the intersection? Choices: A. If the drone chooses to fly directly over the rooftop, it can complete the mission since it has a direct path, but it will take longer. B. If the drone chooses to fly directly over the roof, it will not be able to complete its mission because it will miss the intersection. C. If the drone chooses to fly directly over the roof, it can complete the mission because the advection.	A	Option C; Reason: If the drone chooses to fly directly over the roofs of lower build- ings, it can still reach the trees in the middle of the zebra crossing at the intersection be- cause this alternative path is more efficient and direct.
Counterfactual	Question: If the drone does not fly along the winding road over the lake, but flies directly over the open fields, can it still reach the narrow bridge across the lake in the park? Choices: A. If the drone chooses to fly directly over an open field, it can accomplish the mission because the field provides a clear path, whereas the alternative would require less time. B. If the drone chooses to fly directly over the open field, it will not be able to complete its mission because it will get lost in the middle. C. If the drone chooses to fly directly over the open field, it can complete the mission because the alternative path is more efficient.	В	Option C; Reason: Flying di- rectly over an open field can help the drone avoid trees, making the task easier and po- tentially faster.
Sequence Recall	Question: What is the next step after the drone crosses the second intersection? Choice: A. The drone ascends to gain a higher vantage point. B. The drone continues its progress over the roof. C. The drone descends close to the roof. D. The drone turns sharply to the right. E. The drone hovers in place.	С	Option: B; Reason: After crossing the second intersec- tion, the drone continues its journey over the roof, as seen in the video.

Table	Q٠	Frror	Exam	nle	Part	2
Table	2.	LIIUI	Елаш	pic	1 art	4

Question Category	Question	Ground Truth	LVLM Answer
Causal	Question: Why do drones keep rising after leaving the beach?	А	Option: C; Reason: The drone
	A. The drone ascends to capture a wider view of the waterfront.		level of the building, as shown
	B. The drone rises, avoiding trees and lamp posts.C. The drone rises to avoid hitting the roof of the building.		in the video frames.
	D. The drone ascends to the upper level of the building.		
Landmark Position	Question: The drone is navigating to (the parking space on the right after entering the community	А	Option: E; Reason: The drone
	entrance). When the drone reaches the current location, what is the direction to the destination? Choice:		is currently located over the community center and needs
	A. The drone is currently located above the parking space on the right and has just entered the		to move towards the parking
	B. The drone is currently located far away from the community, above a wide avenue.		space on the right.
	C. The drone is currently located above the community center and not close enough to the parking space.		
	D. The drone is currently located above the left side of the community, near another high-rise building.		
	E. The drone is currently located above the community center and not close enough to the parking		
Goal Detection	Question: The drone is navigating to the balcony on the 13th floor. Is the destination within sight	A	Option: D; Reason: The video
	at the current location? If yes, where is the destination in the view? Choice:		does not show the balcony on the 13th floor within the
	A. The balcony on the 13th floor is currently located below the center of the drone's field of view.		drone's field of view at the cur-
	of view.		Tent location.
	C. The balcony on the 13th floor is currently located on the right side of the drone's field of view. D. The balcony on the 13th floor is currently located in the lower right corner of the drone's field		
Goal Detection	of view.		Ontion: D: Reason: The video
Goal Delection	at the current location? If yes, where is the destination in the view?	1	does not show the balcony
	Choice: A. The balcony on the 13th floor is currently located below the center of the drone's field of view.		on the 13th floor within the drone's field of view at the cur-
	B. The balcony on the 13th floor is currently located in the upper left corner of the drone's field of view		rent location.
	C. The balcony on the 13th floor is currently located on the right side of the drone's field of view.		
	D. The balcony on the 13th floor is currently located in the lower right corner of the drone's held of view.		
Action Generation	The drone is navigating to the 31st floor balcony. What is the next move for the drone, given the target location for the drone's next move?	В	Option: A; Reason: The drone needs to rise to reach the
	Choices:		31st floor balcony, so the next
	A. Kise. B. Descend.		This aligns with option A.
	C. Fly forward. D. Turn left.		
	E. Turn right.		
	G. Rotate the camera downward.		
Action Generation	The drone is navigating to the 13th floor balcony. What is the next move for the drone, given the target location for the drone's next move?	В	Option: F; Reason: The drone needs to rotate the camera up-
	Choices: A Rise		ward to align with the balcony on the 13th floor
	B. Descend.		
	D. Turn left after descending.		
	E. Fly forward. E. Turn left.		
	G. Turn right.		
	I. Rotate the camera downward.		
High-level Planning	The drone is navigating to a nearby bus stop. Thinking back on past actions and objects seen in the video, what should the drone approach next in order to get from its current location to its	Е	Option: C; Reason: Crossing the street is the next step for
	destination?		the drone to get to the bus
	A. The center of the square.		stop.
	B. The top of the tallest building in the distance. C. Crossing the street.		
	D. The end of the parking lot. E. Street edge.		
High-level Planning	The drone is navigating to the largest intersection nearby. Thinking back on past actions and	В	Option: A; Reason: The top of
	location to its destination?		tion. To navigate to the inter-
	Choices: A. The top of the traffic light at the intersection.		section, approaching the top of the traffic light would pro-
	B. Center of the roundabout.		vide a clear vantage point for identifying the intersection
	D. Street parking.		identifying the intersection.
Progress Evaluation	E. The tallest building visible on the horizon. The drone moves according to a series of movement instructions. What are the drones doing	В	Option: C; Reason: The drone
-	now?		is currently positioned to the right of the trees as it moves
	A. Fly towards the sun.		through the cityscape.
	B. Go to the clear between the buildings. C. Go to the river.		
	D. Go to the brown building. E. Go to the sun.		
Progress Evaluation	The drone moves according to a series of movement instructions. What are the drones doing	В	Option: D; Reason: The drone
	Choices:		ward view of the city, and
	A. Fly towards the sun. B. Turn right to the residential building.		changing the camera gimbal to normal will provide a more
	C. Increase the height of the drone.		standard perspective.
	D. Change the camera gimbal from downward to normal.		