

# EGONIGHT: TOWARDS EGOCENTRIC VISION UNDERSTANDING AT NIGHT WITH A CHALLENGING BENCHMARK

Anonymous authors

Paper under double-blind review

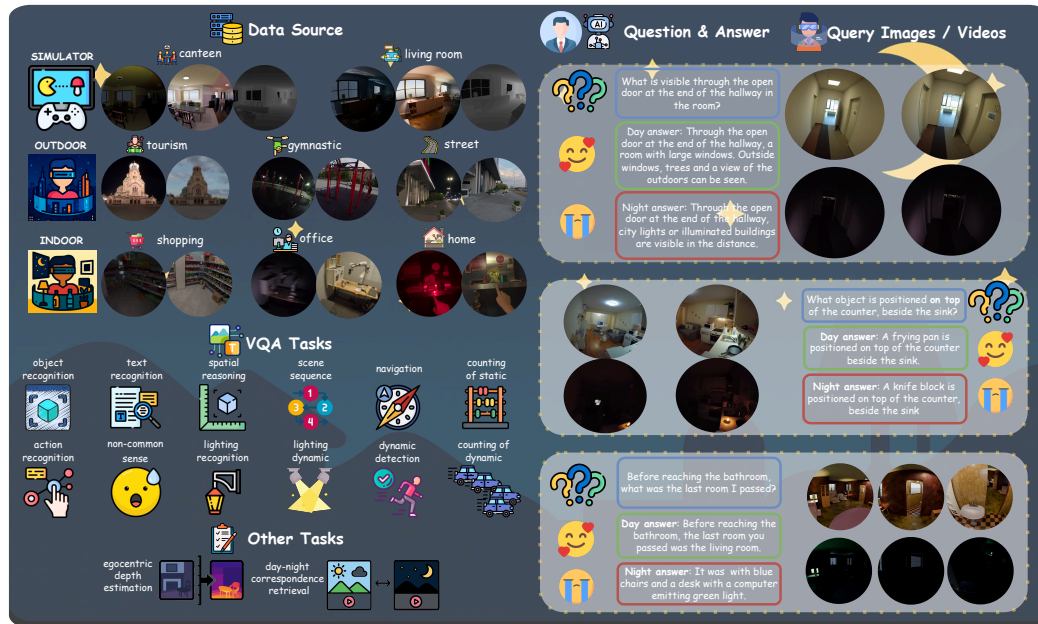


Figure 1: **Overview of the EgoNight.** EgoNight integrates diverse video sources spanning synthetic environments, real-world indoor and outdoor scenes, recorded under both daytime and nighttime conditions, with spatial and temporal alignment. It consists of three benchmarks: (i) *egocentric VQA* as the primary focus, (ii) *day–night correspondence retrieval*, and (iii) *egocentric depth estimation*, all targeting the challenges of low-light egocentric vision. The day–night alignment (illustrated on the right with VQA examples) enables rigorous analysis of illumination gaps in MLLMs.

## ABSTRACT

Most existing benchmarks for egocentric vision understanding focus primarily on daytime scenarios, overlooking the low-light conditions that are inevitable in real-world applications. To investigate this gap, we present **EgoNight**, the first comprehensive benchmark for nighttime egocentric vision, with visual question answering (VQA) as the core task. A key feature of EgoNight is the introduction of day–night aligned videos, which enhance night annotation quality using the daytime data and reveal clear performance gaps between lighting conditions. To achieve this, we collect both synthetic videos rendered by Blender and real-world recordings, ensuring that scenes and actions are visually and temporally aligned. Leveraging these paired videos, we construct **EgoNight-VQA**, supported by a novel day-augmented night auto-labeling engine and refinement through extensive human verification. Each QA pair is double-checked by annotators for reliability. In total, EgoNight-VQA contains 3658 QA pairs across 90 videos, spanning 12 diverse QA types, with more than 300 hours of human work. Evaluations of the state-of-the-art multimodal large language models (MLLMs) reveal substantial performance drops when transferring from day to night, underscoring the challenges of reasoning under low-light conditions. Beyond VQA, EgoNight also introduces two auxiliary tasks, *day–night correspondence retrieval* and *egocentric depth estimation at night*, that further explore the boundaries of existing models. We believe EgoNight-VQA provides a strong foundation for advancing application-driven egocentric vision research and for developing models that generalize across illumination domains. All the data and code will be made available upon acceptance.

# 1 INTRODUCTION

With the rapid development of wearable devices, egocentric vision understanding has become increasingly important. Unlike third-person vision, egocentric perception naturally aligns with the way humans perceive, understand, and interact with the world. A robust egocentric vision system can not only serve as an intelligent assistant in daily activities Yang et al. (2025) but also play a crucial role in embodied AI and robotic learning Li et al. (2025a); Kareer et al. (2025). Beyond these general applications, egocentric vision holds unique potential for assisting specific user groups such as people who are blind or visually impaired Xiao et al. (2025), or physically disabled Zhang et al. (2023a), enabling technologies that enhance navigation, accessibility, and real-time scene understanding.

Significant efforts have been made to advance egocentric vision understanding, including the construction of large-scale ego-centric datasets such as EPIC-KITCHENS Damen et al. (2020), Ego4D Grauman et al. (2022), and Ego-Exo4D Grauman et al. (2024); the design of diverse and challenging benchmarks such as EgoTaskQA Jia et al. (2022), EgoSchema Mangalam et al. (2023), and EgoTempo Plizzari et al. (2025); and the development of egocentric multimodal large language models (MLLMs) such as EgoVLPv2 Pramanick et al. (2023), EgoGPT Yang et al. (2025), and Exo2Ego Zhang et al. (2025a). Despite these advances, almost all prior works focus on daytime scenarios with favorable lighting. In contrast, real-world egocentric systems, for example, intelligent personal assistants for navigation, must operate at night, under low light, uneven illumination, and severely limited visibility. **This motivates us to investigate egocentric vision at night, focusing on complex scene understanding and reasoning tasks.**

**A central challenge in constructing such a benchmark lies in obtaining suitable video sources that capture the characteristics of nighttime environments and developing annotation methods that ensure high labeling quality.** To address this, we place particular emphasis on *day–night aligned videos*, which not only allow us to leverage daytime data to annotate nighttime videos, but also enable rigorous performance comparisons across day and night lighting conditions. However, in practice, collecting perfectly aligned day–night pairs in the real world is highly **non-trivial**. To overcome this, we leverage Blender Iraci (2013), where scene layouts, camera trajectories, and lighting can be precisely controlled, enabling the synthesis of the desired videos. This produces **EgoNight-Synthetic**, a collection of 50 ideally aligned egocentric pairs spanning diverse and complex indoor scenarios with varying illumination levels. To complement synthetic data with real-world evidence, we design a *video-guided recording protocol* to construct **EgoNight-Sofia**, which contains 20 pairs of real-world egocentric videos with spatially and temporally aligned day–night counterparts. These videos cover realistic use cases (e.g., “Where did I put my keys?”, “How much is the item I saw in the grocery shop?”), spanning both indoor and outdoor environments under diverse illumination sources such as streetlights, flashlights, and candles. Finally, we incorporate 20 nighttime videos from the Oxford Day-and-Night dataset Wang et al. (2025b), termed **EgoNight-Oxford**, which serve as an additional testbed despite lacking day–night alignment. Together, these three video sources constitute our **EgoNight** dataset, which is the first egocentric dataset providing day–night aligned correspondences, as mainly summarized in Fig. 1.

The videos in EgoNight pave the way for constructing challenging benchmarks to evaluate the capabilities of existing models. Among many egocentric tasks, we focus on the egocentric video question answering, a flagship task that best reflects high-level understanding in egocentric vision. Specifically, to comprehensively evaluate model abilities, we first propose a diverse set of QA types, spanning well-studied tasks (e.g., object recognition, spatial reasoning, action recognition, counting, text recognition) as well as several underexplored dimensions (e.g., temporal scene sequence understanding, navigation, lighting recognition, and non–common-sense reasoning). These are further organized into paired and unpaired QA types, depending on whether day–night counterparts share the same questions and answers. To construct the benchmark at scale, we then develop a *novel three-stage day-augmented auto-labeling pipeline* that leverages daytime videos to assist in generating question–answer pairs for nighttime clips, followed by extensive human verification to ensure accuracy and reliability. Building EgoNight and annotating VQA required **over 300 hours of human effort**, with each QA pair verified by at least one expert annotator. This process results in the high-quality **EgoNight-VQA** dataset, comprising 3,658 QA pairs. Beyond VQA, we introduce two auxiliary tasks with dedicated testbeds: day–night correspondence retrieval, which evaluates cross-illumination matching, and egocentric depth estimation at night, which is crucial for navigation and interaction in embodied AI. These two tasks further broaden the benchmark and expose new challenges for existing models.

Our extensive experiments across three video sources, three tasks, and 10 state-of-the-art multimodal large language models reveal that nearly all models (including closed-source models such as GPT and Gemini) struggle on this challenging benchmark, with a clear and consistent performance gap between day and night. This highlights the unsolved challenges of egocentric vision at nighttime and calls for more robust models that generalize across illumination conditions. Besides, we highlight that our newly introduced QA types, covering lighting recognition/dynamic, scene sequence reasoning, navigation, and non–common-sense reasoning, are substantially more challenging than well-studied categories, revealing fresh difficulties for MLLMs. **We further prove synthetic data is highly correlated with real data**



and effectively boosts real-world performance through fine-tuning. Our pilot studies further show that fine-tuning on specialized subset of data improves model performance through adapting vision encoder into low-light domain and aligning the language model to the uncertain features during night.

Our main contributions are threefold: i) **EgoNight Dataset**: We present the first egocentric dataset that systematically addresses nighttime conditions, featuring day–night aligned videos from synthetic (EgoNight-Synthetic), real-world (EgoNight-Sofia), and existing (EgoNight-Oxford) sources. ii) **Benchmark Suite**: We build a comprehensive benchmark centered on egocentric VQA with diverse QA types and 3658 fully human-verified QA pairs, complemented by egocentric depth estimation at night and day–night correspondence retrieval tasks. iii) **Empirical Insights**: Extensive evaluations reveal clear day–night performance gaps, underscoring illumination robustness as a key challenge; our newly proposed QA types are also validated to pose practical difficulties for current MLLMs.

## 2 RELATED WORKS

### 2.1 EGOCENTRIC DATASETS AND VQA BENCHMARKS

A series of large-scale egocentric datasets, such as EPIC-KITCHENS Damen et al. (2020), Ego4D Grauman et al. (2022), Ego-Exo4D Grauman et al. (2024), and EgoExoLearn Huang et al. (2024), have laid the foundation for a wide range of tasks, including action recognition Sudhakaran et al. (2019), object detection Ren & Gu (2010), pose estimation Luo et al. (2021), video generation Liu et al. (2021), Ego-Exo correspondence Fu et al. (2025). Among these, we are particularly interested in egocentric visual question answering (VQA) Fan (2019), which provides a natural and human-like framework for comprehensively evaluating model performance through question–answer interactions. In recent years, several egocentric VQA benchmarks have been proposed, including EgoVQA Fan (2019), EgoTaskQA Jia et al. (2022), EgoSchema Mangalam et al. (2023), EgoThink Cheng et al. (2024), EgoTempo Plizzari et al. (2025), EgoCross Li et al. (2025b), EgoBlind Xiao et al. (2025), EgoMemoria Ye et al. (2024), HourVideo Chandrasegaran et al. (2024), EgoLifeQA Yang et al. (2025) with different focuses. However, nearly all of them are confined to daytime or well-lit scenarios, leaving model performance in low-light or nighttime conditions largely unexplored. The Oxford Day-and-Night dataset Wang et al. (2025b) is a partial exception but was not designed for VQA and lacks day–night alignment. This makes EgoNight and EgoNight-VQA fundamentally distinct from prior benchmarks.

### 2.2 MLLMs FOR VIDEO UNDERSTANDING

The rapid development of multimodal large language models (MLLMs) has substantially advanced the frontier of video understanding. Prominent open-source models include Qwen-VL Bai et al. (2023), InternVL Chen et al. (2024b), Video-LLaMA Zhang et al. (2023b), LLaVA-NeXT-Video Li et al. (2024), and GLM-V Hong et al. (2025), while closed-source commercial systems such as GPT-4V Achiam et al. (2023) and Gemini Comanici et al. (2025) demonstrate even stronger capabilities in video captioning, summarization, and open-ended visual question answering. Building on these advances, egocentric MLLMs have emerged to adapt foundation models from exocentric to first-person perspectives. Representative examples include EgoVLPv2 Pramanick et al. (2023) for improved video–language cross-modal fusion, EgoGPT Yang et al. (2025) fine-tuned with egocentric captioning and QA, MM-Ego Ye et al. (2024) with a memory mechanism for long videos, and Exo2Ego Zhang et al. (2025a) leveraging exocentric data for egocentric generalization. These works highlight the potential of MLLMs as egocentric assistants. However, nearly all of them are developed and tested under well-lit daytime conditions, leaving their robustness in low-light or nighttime scenarios unexplored. Nevertheless, almost all existing MLLMs are developed and evaluated under well-lit daytime conditions, with little consideration of low-light or nighttime videos, leaving their robustness in low-light or nighttime scenarios unexplored.

### 2.3 CROSS-DOMAIN GENERALIZATION

Domain generalization Zhou et al. (2022) is a long-standing challenge in computer vision, where models trained on one distribution must adapt to another. Shifts can arise from semantic drift, style changes, or variations in weather and lighting. Many algorithms have been validated across tasks such as image classification Li et al. (2017); Zhou et al. (2021), object detection Fu et al. (2024); Li et al. (2025c), action recognition Pan et al. (2020); Bian et al. (2011), few-shot learning Guo et al. (2020); Fu et al. (2021), and autonomous driving Li et al. (2022; 2023a). In contrast, cross-domain transfer for MLLMs, especially in video understanding, remains underexplored, with only a few recent attempts (e.g., CL-CrossVQA Zhang et al. (2025b), VQA-GEN Unni et al. (2023), Super-CLEVR Li et al. (2023c)). However, none of them are targeted for egocentric video, which is naturally different from exocentric videos in terms of recorded images, camera motion, and contained information. The most relevant effort to us is EgoCross Li et al.

(2025b), an egocentric VQA benchmark that moves beyond daily activities to evaluate model generalization across distinct long-tail, specialized domains such as surgery and industrial settings. In this paper, however, we investigate MLLMs from a different perspective, robustness under nighttime conditions, a common and ubiquitous scenario in daily life, yet previously overlooked dimension of domain generalization in egocentric video understanding.

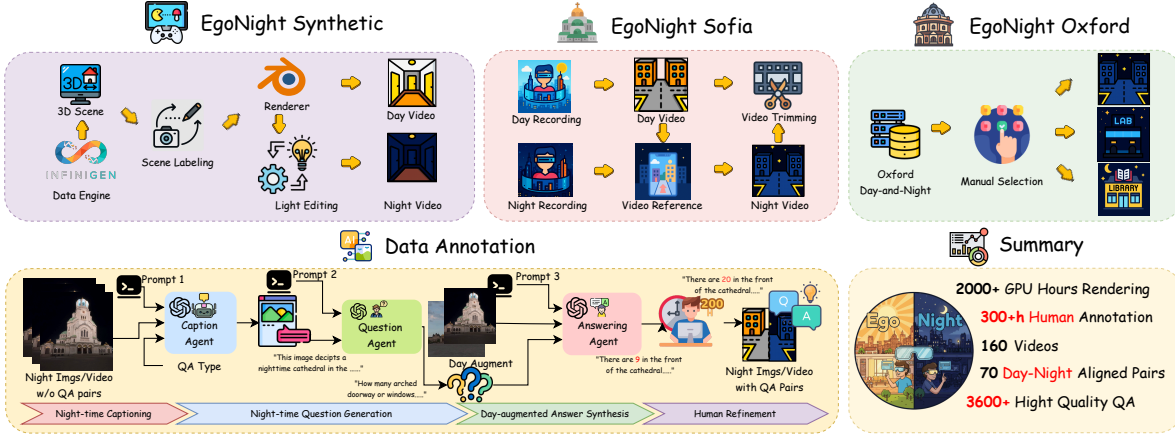


Figure 2: **EgoNight construction and EgoNight-VQA annotation.** EgoNight integrates EgoNight-Synthetic, EgoNight-Sofia, and EgoNight-Oxford sources. Annotation is achieved via a novel three-stage day-augmented Auto QA generation pipeline with 300+ hours of human refinement, resulting in over 3600 high-quality QA pairs.

### 3 EGO NIGHT DATASET & BENCHMARKS

#### 3.1 VIDEO SOURCE COLLECTION

**Overview & Design Principles.** EgoNight is built to systematically evaluate MLLMs under challenging nighttime conditions, which are critical for developing robust intelligent assistants. The collection of video sources follows four principles: ① *Reflect real-world challenges*, such as walking on dimly lit streets or navigating indoors during power outages; ② *Involve natural camera movements* and preferably capture *actions and interactions* with the environment to evaluate both static perception and dynamic understanding; ③ *Ensure diversity of scenarios, illumination, and task difficulties*, spanning indoor, outdoor, office, and grocery settings, lighting from streetlights, flashlights, headlights, and candles, and task levels from easy (relatively clear), through medium (partially visible), to hard (barely visible). ④ *Enable rigorous analysis through day–night paired videos*, where scenes, trajectories, and actions remain consistent across conditions so that differences can be attributed solely to illumination. To meet these requirements, EgoNight integrates three complementary video sources, as illustrated in the upper part of Fig. 2 and detailed below.

**EgoNight-Synthetic.** To obtain perfectly aligned day–night video pairs, we used a simulation environment where every element can be precisely controlled, including the scene layout, camera path, and lighting. This ensures that the day and night videos match exactly at the pixel and frame level, with lighting being the only difference. We first employ Infinigen Raistrick et al. (2023) to generate diverse indoor 3D scenes. Human annotators cleaned and refined these scenes, then simulated walking through the space at a normal speed (1.2 m/s), recording the camera trajectory. We replayed the same trajectory under different lighting conditions. Daytime videos were rendered using Blender Iraci (2013), and we adjusted the lighting to create the corresponding nighttime versions.

In total, EgoNight-Synthetic contains 50 pairs of egocentric videos, covering more than 100 environment assets. These include indoor scenes such as kitchens, bathrooms, and living rooms, populated with over 50 diverse object categories (e.g., windows, tables, beds, chairs, lamps, bookshelves, plates). We design multiple illumination setups, ranging from uniformly lit rooms to sparsely localized lighting, and incorporate three difficulty levels with a different range of motion blur, sensor noise, and illumination level. Besides RGB frames, Blender also allows us to generate ground-truth depth and normals (see Appendix Sec. A.1), making EgoNight-Synthetic richer and more versatile.

**EgoNight-Sofia.** To include realistic human–environment interactions missing from synthetic videos, we also recorded our own day–night paired videos. Capturing perfectly aligned real-world pairs is challenging, so we designed a practical video-guided recording strategy with post-trimming for better alignment.

We first record a daytime video with an ego-wearer exploring an environment while viewing the live camera feed on a phone screen. For the nighttime version, the same person, device, and viewpoint are kept unchanged. Instead of using live preview, the recorded daytime video is played back on the phone, serving as visual guidance to help the ego-wearer match walking speed, viewpoints, and actions. After brief practice, this approach proved more stable and reliable than methods like using landmarks or memorized trajectories. Post-trimming is applied to further refine spatial and temporal consistency.

Our real-world dataset, EgoNight-Sofia, contains 20 day–night paired videos recorded in Sofia, Bulgaria. Despite its modest size, it is a rare resource capturing diverse real-world everyday scenarios, including apartments, offices, grocery stores, streets, tourist spots, and outdoor fitness areas. The recordings include natural actions such as drinking water, locking doors, placing keys, charging devices, or checking price labels—leading to realistic VQA cases (e.g., “Where did I put my keys?”, “How much was the drink?”, “Did I turn left?”). Illumination sources include street lights, lamps, flashlights, and candles.

**EgoNight-Oxford.** Oxford Day–Night Wang et al. (2025b) is a notable exception that also includes egocentric videos captured under both daytime and nighttime conditions. Although it was originally designed for 3D vision tasks such as novel view synthesis, it offers illumination variations across five representative locations in Oxford. However, the day and night videos are not spatially or temporally aligned.

To enrich EgoNight with more realistic nighttime content—particularly for urban outdoor scenes—we manually select 20 nighttime segments to form EgoNight-Oxford, based on two criteria: (i) minimal overlap in trajectories and locations, and (ii) genuinely low-light conditions. These segments serve as a complementary testbed for evaluating model generalization under illumination changes when paired alignment is unavailable.

For both EgoNight-Sofia and EgoNight-Oxford, human annotators categorize each video into easy, medium, or hard levels. Together, these sources provide EgoNight with a balanced combination of precise alignment, natural dynamics, and broad real-world diversity.

### 3.2 EGO NIGHT-VQA BENCHMARK RECONSTRUCTION

**QA Task Taxonomy.** To thoroughly assess models from multiple perspectives, we define a diverse taxonomy of 12 QA tasks. Some of these categories are well-studied and have been explored in previous egocentric VQA benchmarks, such as object/action/text recognition, counting, and spatial reasoning. Others are much less studied or newly proposed in EgoNight-VQA, including scene sequence and navigation (which require not only visual perception but also memory and spatial reasoning), illumination recognition and illumination change (designed to test models’ understanding of lighting concepts), and non–common-sense reasoning (e.g., detecting abnormal cases such as a door inserted into a wall in the synthetic data). More detailed explanations of QA types can be found in the Appendix Sec. A.4.1. We further organize these categories into *paired* and *unpaired* QA types, depending on whether the same questions can be consistently applied across day–night counterparts: 1) **Paired QA Types.** These cover contexts that remain unchanged across day and night, allowing the same QA pairs to be used for both videos and thus providing a clean testbed for measuring performance gaps. Specifically, we include: ① *object recognition*, ② *text recognition*, ③ *spatial reasoning*, ④ *scene sequence*, ⑤ *navigation*, ⑥ *counting of static*, ⑦ *action recognition*, and ⑧ *non–common-sense reasoning*. 2) **Unpaired QA Types.** These include categories that are impractical to pair across day and night, or are only meaningful in the nighttime condition. We consider: ① *lighting recognition*, ② *lighting dynamic*, ③ *dynamic detection*, and ④ *counting of dynamic*. We control QA clip duration by task type. For static or spatial tasks (e.g., object recognition, lighting recognition), we use short clips of 3 seconds to minimize redundancy; For dynamic or temporal tasks (e.g., action recognition, navigation), the entire video is used to capture the complete context. Following recent works Plizari et al. (2025); Xiao et al. (2025), we adopt the *open-ended QA* setting over the closed-form multiple-choice format, as it better reflects natural human–AI interactions.

A detailed summary of each QA type, including whether it is paired or unpaired, clip duration, and example questions, is provided in Fig. 3. This taxonomy makes EgoNight-VQA not only diverse and well-structured but also novel, introducing illumination reasoning and other challenges uniquely tied to nighttime egocentric vision.

**Day-Augmented Auto QA Generation.** Constructing large-scale QA pairs for nighttime videos is particularly challenging due to low visibility, which makes direct annotation both time-consuming and error-prone. To address this, as illustrated in the lower part of Fig. 2, we design a novel three-stage *day-augmented auto QA generation* pipeline that leverages aligned daytime videos as a strong prior for annotating their nighttime counterparts.

Specifically, the pipeline is tailored to each QA type and consists of three stages:

1) **Nighttime captioning.** For each clip, we prompt advanced MLLMs to generate detailed captions with an explicit focus on the target QA type (e.g., highlighting object-related attributes for object recognition or text/logos for text recognition). This ensures that the captions capture the most key information or construct relevant QA pairs.

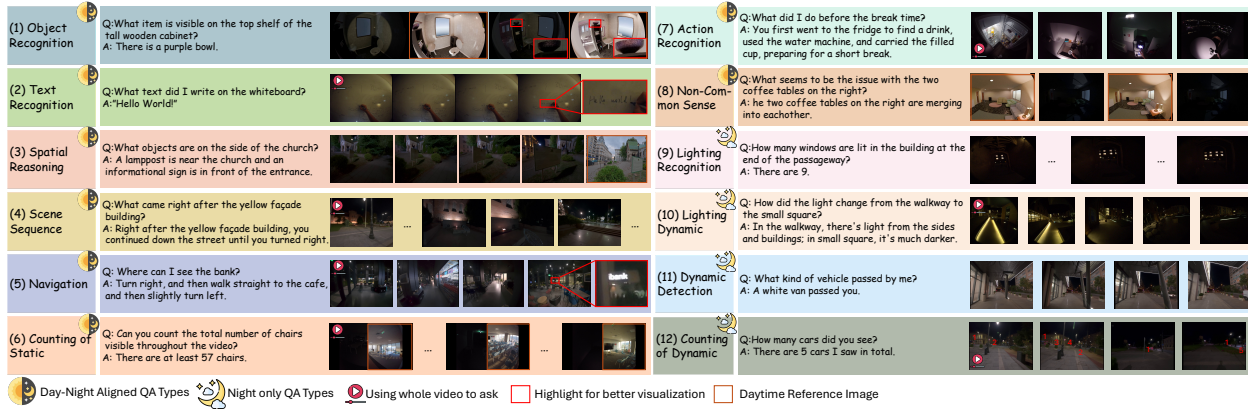


Figure 3: **QA types with examples.** The first eight are *paired* types, where the same question–answer applies to both day and night clips; the last four are *unpaired*, evaluated only at night. QA Types have various durations, with static or spatial tasks (e.g., 1 and 3) using short clips, while dynamic or temporal tasks (e.g., 4 and 5) use full videos.

2) **Nighttime question generation.** The caption, together with the corresponding night clip, is then fed into the same MLLM to produce diverse question candidates centered on the given QA type. This step encourages variety in phrasing and perspective while maintaining fidelity to the visual content.

3) **Day-augmented pseudo answer synthesis.** For paired QA types, pseudo answers are generated by consulting the aligned daytime clip, where content is more visible and less ambiguous. For unpaired QA types or datasets without alignment (e.g., EgoNight-Oxford), answers are instead derived directly from the nighttime clip.

All three stages are powered by GPT-4.1. Empirically, we find that both the QA-type-specific prompting and the inclusion of daytime videos substantially improve the quality and reliability of the generated QA pairs. [More examples of VQA pairs and caption generation can be found in Appendix A.4.2 and A.7.1.](#)

**Human Annotator Refinement.** Finally, human annotators refine QA pairs via three operations: i) **delete**, when QA pairs are meaningless, vague, duplicated, or inconsistent across day–night counterparts (for paired QA types); ii) **modify**, when the question is valid but the answer is wrong (or vice versa), or to resolve ambiguity; iii) **add**, when many pairs are removed or when important, challenging questions, especially about dynamic concepts, are missing.

After the first labeling round, we performed a random double-check to refine low-quality annotations. Thus, although our pipeline combines model generation with human refinement, *every QA pair (3,658 in total) is manually verified at least once.* In total,  $\sim 200$  hours of human effort were invested, ensuring the quality and reliability of EgoNight-VQA.

**Dataset Statistics.** EgoNight-VQA comprises 3,658 high-quality, fully human-verified QA pairs across 12 task types, sourced from EgoNight-Synthetic, EgoNight-Sofia, and EgoNight-Oxford, with an average of 40 pairs per video. Detailed statistics on QA distribution, video durations, task difficulties, scenarios, and illumination are shown in Fig. 4. [The number of videos across the three subsets—Synthetic \(50\), Sofia \(20\), and Oxford \(20\)—is proportionally reflected in the VQA annotations \(2029 : 813 : 816\).](#) This results in an approximately 1:1 balance between synthetic and real (Sofia + Oxford) VQA samples, ensuring that our benchmark is not dominated by synthetic content. [We provide more comparison of Egocentric VQA datasets in Appendix A.3.](#) Overall, EgoNight-VQA provides a diverse and comprehensive benchmark for evaluating egocentric vision models under nighttime conditions.

### 3.3 BENCHMARKS BEYOND EGOCENTRIC VQA

**Day-Night Correspondence Retrieval.** To further assess model capabilities beyond VQA, we introduce *day–night correspondence retrieval*, which evaluates a model’s ability to match visual content across illumination conditions. Specifically, we define two subtasks: i) **Spatial Retrieval (Scene Recognition).** Spatial retrieval, or scene recognition, is a long-standing vision task Arandjelovic et al. (2016); Miao et al. (2024). Here, it is extended: given a query clip and a set of  $N$  candidate clips of equal duration  $s$ , the model must retrieve the one depicting the same scene. This evaluates a model’s ability to capture and relate spatial relations in egocentric videos, e.g., distinguishing a bedroom from a bathroom or another bedroom. We built this benchmark with 1000 randomly generated meta-tasks. Each task samples a query clip, and the candidate set includes its temporally aligned counterpart (with a temporal shift for added difficulty) plus  $N - 1$  negatives from other scenes. Performance is measured by Top-1 accuracy across all tasks. In our setup, we use  $N = 10$ ,  $s = 10$  seconds, and a temporal shift of  $[10, 20]$  frames. Both *Day (query)  $\rightarrow$  Day (database)* and *Day  $\rightarrow$  Night* settings are evaluated.



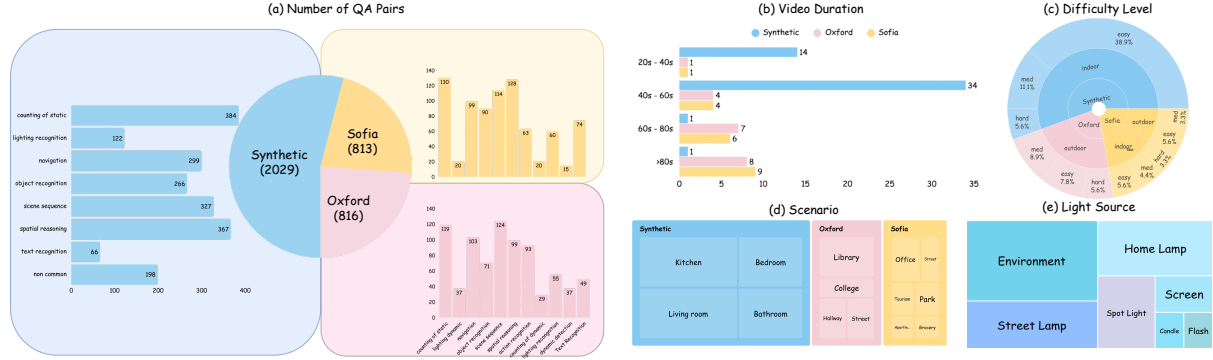


Figure 4: **Statistics of EgoNight-VQA benchmark.** (a) Distribution of QA pairs across QA types and sources. (b) Video duration distribution. (c) Task difficulty levels cross scenarios. (d) Scenario coverage. (e) Illumination coverage.

**ii) Temporal Localization.** We further design a temporal localization task to test whether models can align video segments across dynamics. Given a query clip of duration  $s$ , the model must localize it within the corresponding full video by predicting its start and end timestamps  $(t_i, t_j)$ , directly evaluating temporal reasoning (e.g., grounding “The door is being locked” to 10–20s). We construct 1000 meta-tasks, each generated by randomly sampling one clip from its parent full video that is also randomly selected. Inspired by temporal grounding literature Xin et al. (2024), we adopt mean Intersection-over-Union (mIoU) between the predicted interval  $(t_i, t_j)$  and the ground-truth interval  $(t_i^*, t_j^*)$  as the evaluation metric. Consistent with spatial retrieval, we set  $s = 10$  seconds and evaluate both *Day*  $\rightarrow$  *Day* and *Night*  $\rightarrow$  *Day* settings.

**Egocentric Depth Estimation at Night.** Depth estimation is a fundamental component of computer vision. On the one hand, extensive research Yang et al. (2024a;b); Wang et al. (2025a) has focused on depth estimation in non-egocentric settings (typically not with fisheye cameras), while egocentric depth estimation remains largely underexplored, especially under nighttime conditions. On the other hand, recent works Chen et al. (2024a); Liu et al. (2025) suggest that incorporating depth can enhance models’ spatial reasoning abilities. These two observations motivate us to construct an auxiliary benchmark for *egocentric depth estimation at night*. Specifically, we use EgoNight-Synthetic as the testbed, where ground-truth depth maps are provided by the rendering engine. Thanks to the day–night aligned design, we can quantitatively evaluate models under both controlled daytime and nighttime conditions. For evaluation, we adopt standard depth estimation metrics, including absolute relative error (AbsRel),  $\delta_1(1.25)$ ,  $\delta_2(1.25^2)$ , and  $\delta_3(1.25^3)$ , where  $\delta_k$  measures the percentage of predicted pixels whose relative error is within a threshold of  $1.25^k$ .

## 4 EXPERIMENTS

### 4.1 EVALUATED MLLMS & METRICS

We evaluate a broad set of state-of-the-art MLLMs on the proposed benchmarks. i) For **EgoNight-VQA**, we include two closed-source commercial models, GPT-4.1 Achiam et al. (2023) and Gemini 2.5 Pro Comanici et al. (2025); eight open-source models, Qwen2.5-VL (3B, 7B, 72B) Bai et al. (2023), VideoLLaMA3 (7B) Zhang et al. (2023b), InternVL3 (8B) Chen et al. (2024b), GLM-4.1V (9B-Base) Hong et al. (2025), and LLaVA-NeXT-Video (7B) Li et al. (2024); as well as EgoGPT Yang et al. (2025), one of the few open-source egocentric models tailored for open-ended QA. Following prior work Plizzari et al. (2025); Fan (2019), we adopt an *LLM-as-a-Judge* strategy to assess semantic consistency between predictions and ground truth, and report average accuracy across the test sets. ii) We provide further in-depth analysis on synthetic data quality, and potential model improvements. iii) For **day–night correspondence retrieval**, we benchmark feature-based retrieval methods, DINOv2 Oquab et al. (2023) and Perception Encoder (Percep. Enc.) Bolya et al. (2025), alongside MLLM-based methods, GPT-4.1 and InternVL3 (8B). As described in Sec. 3.3, Top-1 accuracy (Acc-R@1, %) and mIoU (%) are used for evaluating the spatial and temporal subtasks, respectively. iv) For **egocentric depth estimation**, we test a general monocular depth model (Depth Anything Yang et al. (2024a;b)), a 3D reconstruction-based method (VGGTStream Zhuo et al. (2025); Wang et al. (2025a)), and two egocentric fisheye-specific models (DAC Guo et al. (2025) and UniK3D Piccinelli et al. (2025)). For Depth Anything and VGGTStream, input fisheye RGB frames and depth maps are undistorted prior to inference for fair comparison. Additional implementation details (e.g., fps for frame extraction, prompts, and model settings) and discussion about *LLM-as-a-Judge strategy* are provided in the Appendix Sec. A.5.

Models	EgoNight-Synthetic			EgoNight-Sofia			EgoNight-Oxford			Avg. -
	Easy	Medium	Hard	Easy	Medium	Hard	Easy	Medium	Hard	
Closed-Source MLLMs										
GPT-4.1	29.30	<b>26.87</b>	<b>18.87</b>	32.04	<b>29.35</b>	<b>31.69</b>	<b>39.72</b>	<b>37.13</b>	<b>40.72</b>	<b>30.93</b>
Gemini 2.5 Pro	<b>31.05</b>	24.81	16.51	<b>38.24</b>	26.81	28.87	36.75	36.81	27.88	30.60
Open-source MLLMs										
InternVL3-8B	20.21	15.50	16.98	24.03	21.74	20.42	22.90	20.85	16.36	20.06
Qwen2.5-VL-72B	18.39	15.25	12.26	24.03	17.03	20.42	24.81	22.80	16.36	18.99
Qwen2.5-VL-7B	13.01	13.95	13.68	15.44	12.68	12.68	13.74	13.36	12.73	13.44
Qwen2.5-VL-3B	14.69	10.34	7.08	15.50	13.04	12.68	17.18	11.40	12.12	13.41
GLM-4.1V-9B-Base	19.09	13.70	15.57	18.60	18.48	16.20	17.15	22.15	18.79	18.20
VideoLLaMA3-7B	16.85	13.44	14.62	11.11	10.87	9.15	12.26	10.46	9.15	13.64
LLaVA-NeXT-Video-7B	6.36	11.37	1.89	13.95	9.78	14.79	3.05	2.61	3.03	7.28
Egocentric MLLMs										
EgoGPT	15.79	13.55	12.04	12.41	12.13	10.36	12.37	13.58	13.68	14.29

Table 1: **Comparison results on EgoNight-VQA.** Accuracies (%) of OpenQA results across three datasets and three difficulty levels. We compare closed-source models, open-source models, and egocentric-specific models.

## 4.2 RESULTS ON EGO NIGHT-VQA

The main results of all MLLMs are shown in Tab. 1. In addition, we provide per-QA performance comparisons between day (striped bars) and night (solid bars) for paired QA types (Fig. 5(a)) and report nighttime performance across all QA types (Fig. 5(b)), based on averages across all models. Note that non-common case detection is available only in EgoNight-Synthetic, while dynamic events and actions are included only in the real-world data.

From the results in Tab. 1, we observe that almost all MLLMs struggle on our benchmark, with maximum averaged accuracies of 30.93% from the closed-source GPT-4.1, 20.06% from the open-source InternVL3-8B, and 14.29% from the egocentric EgoGPT. **The wide performance spread also confirms that our dataset is sufficiently challenging and effective for distinguishing model capabilities.** Fig. 5(a) further highlights the performance gap, showing declines of 32.8% and 25.0% on EgoNight-Synthetic and EgoNight-Sofia, respectively. Together, these results underscore the substantial challenges posed by our benchmark, exposing the limitations of current MLLMs under nighttime scenarios and highlighting the need for more illumination-robust models. Beyond the overall trends, we note three additional insights from Tab. 1: i) Closed-source models perform best. Within open-source models, Qwen2.5-VL generally improves with scale, yet InternVL outperforms the larger Qwen2.5-VL (72B), suggesting that size alone is insufficient. The relatively low results of EgoGPT further emphasize the need for more robust egocentric models. ii) EgoNight-Oxford achieves the highest scores, but its illumination conditions are more challenging than those in EgoNight-Synthetic and EgoNight-Sofia (Sec. A.4.2, Appendix). This indicates that without paired day videos and our day-augmented auto-labeling strategy, even human annotators face difficulties generating challenging QA pairs, underscoring the practical value of our dataset design; iii) Overall, performance declines across task levels (easy, medium, hard), validating the diversity and difficulty of our benchmark.

From the per-QA results in Fig. 5(a) and Fig. 5(b), we further observe three key trends: i) Models perform better on perception-oriented tasks (e.g., object recognition, text recognition, scene sequence) than reasoning-oriented tasks (e.g., navigation, counting, non-common-sense reasoning cases) under daytime conditions. However, at night, perception tasks suffer larger performance drops, indicating their higher sensitivity to illumination, whereas reasoning tasks, though harder overall, are relatively less affected since they rely more on temporal and contextual cues. ii) MLLMs achieve substantially lower accuracy on our newly proposed tasks, such as lighting recognition, lighting dynamics, scene sequence, dynamic detection, navigation, and non-common-sense reasoning, suggesting that existing MLLMs generalize poorly to novel tasks compared with well-studied ones like object recognition. iii) Each dataset in Fig. 5(b) emphasizes distinct aspects of nighttime challenges, together providing complementary perspectives that ensure EgoNight spans a balanced range of perception-reasoning difficulties under low-light conditions.

## 4.3 MORE IN-DEPTH ANALYSIS

**What is the Quality of EgoNight Synthetic.** More examples of synthetic visualization can be found in Appendix A.4.2 and A.1. To further show the quality of our synthetic dataset, we calculate the Pearson correlation of the average score per-model shown in Appendix A.6 between synthetic and Sofia (0.9359 with p-value  $6.84710^{-5}$ ), synthetic and Oxford (0.8588 with p-value  $1.462 \times 10^{-3}$ ). These strong and statistically significant correlations in-

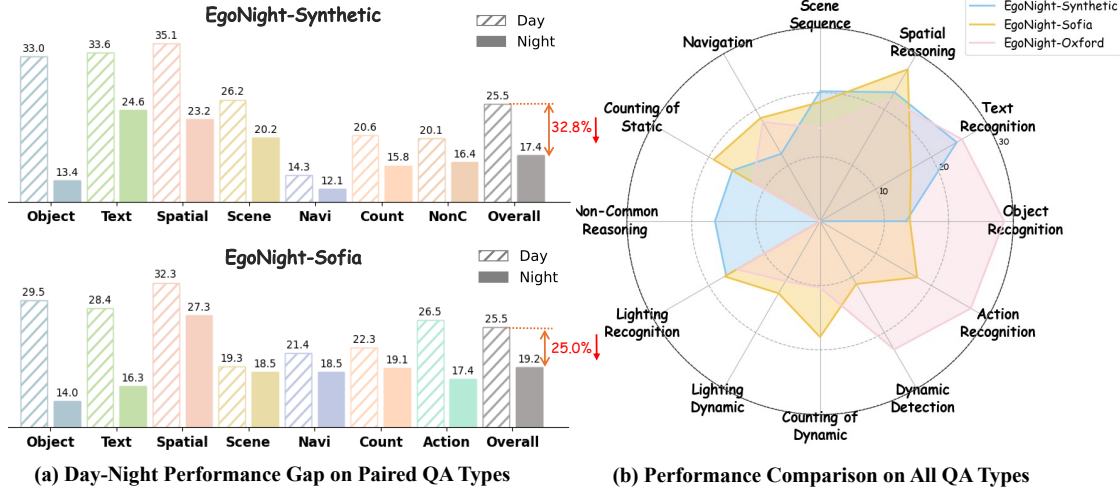


Figure 5: **Performance analysis of MLLMs on EgoNight-VQA.** (a) Day-night performance gap across paired QA types, showing consistent degradation at night. (b) Nighttime performance across all 12 QA types. **NonC** means non-common-sense reasoning.

indicate that performance on synthetic data is highly predictive of performance on real-world data, further validating its representativeness. We also finetune Qwen2.5-VL-7B model with Supervised-Fine-Tuning(SFT) using only the synthetic data and evaluate it on the real dataset. This improves the model accuracy from 14.83% to 20.57%, which demonstrates that our synthetic data can effectively enhance model performance in real-world scenarios.

**Could fine-tuning help to enhance performance?** To explore potential solutions for improving low-light egocentric QA, we proactively conduct pilot studies using the EgoNight benchmark. We split EgoNight into 70% training and 30% testing subsets, and fine-tune Qwen2.5-VL-7B using supervised fine-tuning (SFT) under three configurations: i) Full Fine-Tuning. ii) Fine-tuning vision encoder only. iii) Fine-tuning LLM only. As shown in Tab. 2, our observations are as follows:

- Fine-tuning on EgoNight leads to substantial performance improvements, demonstrating that EgoNight-style nighttime data effectively helps models adapt to low-light egocentric scenarios.
- Both vision-encoder tuning and LLM tuning independently contribute to performance gains. Interestingly, fine-tuning the LLM only yields even larger improvements, suggesting that visual representation is not the only bottleneck. LLM fine-tuning plays a crucial role in aligning uncertain visual features to language space.
- Full fine-tuning consistently outperforms partial fine-tuning, indicating that EgoNight requires both strong visual perception and robust visual-language alignment.

**How perception vs. reasoning-oriented tasks benefit from fine-tuning?** To further dive into the impact of fine-tuning, we compare perception-oriented tasks (Object, Text Recognition) and reasoning-oriented tasks (Navigation, Counting) accuracy in Tab. 3. We can observe that:

- Perception-oriented tasks are significantly easier to enhance through fine-tuning compared to reasoning-oriented tasks, indicating that visual feature learning benefits more from data adaptation than higher-level reasoning.
- Fine-tuning the vision encoder improves perception-oriented tasks by 2–4 times but provides limited gains for reasoning-oriented tasks, revealing the substantial challenges posed by reasoning-centric scenarios.
- Fine-tuning the language model yields improvements for both perception and reasoning-oriented tasks, with larger boosts for perception tasks. This suggests that the main benefit arises from aligning uncertain visual features with language space, while enhancement of true reasoning ability remains limited.

We provide more failure case analysis in Appendix. A.8.1.

#### 4.4 RESULTS ON DAY-NIGHT CORRESPONDENCE RETRIEVAL

Setting	Synthetic	Real
Qwen7B (Base)	23.23	16.40
Enc. FT	29.74	20.92
LLM. FT	35.50	22.26
Full FT	36.25	25.61

Table 2: Fine-tuning performance comparison across datasets. FT means Fine-tuning.

Task	Qwen7B	Enc. FT	LLM FT
Object	8.435	34.718	35.855
Text	18.440	49.890	50.988
Navigation	17.870	19.495	19.918
Counting	16.558	16.945	24.275

Table 3: Fine-tuning performance comparison across tasks.

Models	Spatial Retrieval (Acc - R@1 % $\uparrow$ )				Temporal Localization (mIoU % $\uparrow$ )			
	EgoNight-Synthetic		EgoNight-Sofia		EgoNight-Synthetic		EgoNight-Sofia	
	Day→Day	Night→Day	Day→Day	Night→Day	Day→Day	Night→Day	Day→Day	Night→Day
DINOv2	45.7	28.7	84.5	74.5	-	33.7	-	33.1
Percep. Enc.	65.4	41.6	89.8	80.9	-	32.9	-	33.4
GPT-4.1	75.6	54.1	92.5	84.5	14.7	10.0	21.2	15.5
InternVL3-8B	39.4	27.7	73.9	56.3	10.2	9.9	12.5	13.3

Table 4: Night-to-Day retrieval performance. Each dataset is evaluated on both Day→Day and Night→Day settings.

The results of day-night retrieval are reported in Tab. 4. The gap between Night-Day and Day-Day shows that cross-illumination retrieval remains highly challenging compared with in-domain retrieval. For spatial retrieval, GPT-4.1 consistently outperforms other methods, achieving over 80% accuracy. This suggests that Retrieval-Augmented Generation methods could further improve performance, as Fig. 5(a) already shows that daytime inputs significantly benefit the models. For temporal retrieval, however, GPT-4.1, despite its strong results on egocentric VQA (Tab. 1) and spatial retrieval, shows a substantial drop compared with feature-based methods (DINOv2 and Perception Encoder). A similar degradation is observed for InternVL3-8B. These findings suggest that while MLLMs excel at spatial semantic understanding, they struggle with temporal reasoning, such as timestamp prediction, which is critical for temporal localization. Further results on temporal limitations are provided in Appendix A.6.

#### 4.5 RESULTS ON DEPTH ESTIMATION

Results for depth estimation are reported in Tab. 5. The relatively low scores across all models highlight the difficulty of our EgoNight dataset, which combines egocentric motion, complex geometry, and extreme lighting variations. A clear gap between daytime and nighttime performance again underscores the challenges of low-light conditions. Among the methods, fisheye-based methods (DAC and UniK3D) outperform general depth estimators, suggesting the need for egocentric-specific algorithms. Additional qualitative results are provided in Sec. A.7.3.

Method	Abs Rel $\downarrow$		$\delta_1$ (1.25) $\uparrow$		$\delta_2$ (1.25 <sup>2</sup> ) $\uparrow$		$\delta_3$ (1.25 <sup>3</sup> ) $\uparrow$	
	Day	Night	Day	Night	Day	Night	Day	Night
Depth Anything (U)	0.297	0.302	0.249	0.237	0.463	0.447	0.622	0.60
VGGTStream (U)	0.293	0.298	0.234	0.232	0.447	0.442	0.615	0.609
DAC (F)	0.245	0.292	0.255	0.216	0.495	0.425	0.684	0.602
UniK3D (F)	0.224	0.253	0.280	0.254	0.524	0.481	0.706	0.658

Table 5: Depth estimation results on EgoNight-Synthetic. U: undistorted input; F: fisheye input.

## 5 CONCLUSION

In this work, we introduced EgoNight, the first benchmark suite designed to systematically evaluate egocentric multimodal large language models (MLLMs) under challenging nighttime conditions. EgoNight integrates synthetic and real-world videos with day-night alignment, enabling rigorous analysis of illumination effects. Building upon this data, we proposed EgoNight-VQA, spanning 12 QA types with 3,658 human-verified pairs, alongside two complementary benchmarks: day-night correspondence retrieval and egocentric depth estimation. Experiments reveal that even state-of-the-art MLLMs struggle under low-light conditions, with performance dropping substantially compared to daytime. This highlights that nighttime egocentric vision remains far from being solved, motivating future research into illumination-robust egocentric perception and reasoning. More discussion about contribution and future works is provided in Appendix A.8.2 and A.8.3. We believe EgoNight provides a valuable and timely benchmark that will drive progress toward more reliable egocentric AI assistants.



## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. 2016.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- Wei Bian, Dacheng Tao, and Yong Rui. Cross-domain human action recognition. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 2011.
- Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, Junke Wang, Marco Monteiro, Hu Xu, Shiyu Dong, Nikhila Ravi, Daniel Li, Piotr Dollár, and Christoph Feichtenhofer. Perception encoder: The best visual embeddings are not at the output of the network, 2025. URL <https://arxiv.org/abs/2504.13181>.
- Keshigeyan Chandrasegaran, Agrim Gupta, Lea M Hadzic, Taran Kota, Jimming He, Cristóbal Eyzaguirre, Zane Durante, Manling Li, Jiajun Wu, and Li Fei-Fei. Hourvideo: 1-hour video-language understanding. *NeurIPS*, 2024.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *CVPR*, 2024a.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *CVPR*, 2024b.
- Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. Egothink: Evaluating first-person perspective thinking capability of vision-language models. In *CVPR*, 2024.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blisstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multi-modality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. *TPAMI*, 2020.
- Chenyou Fan. Egovqa-an egocentric video question answering benchmark dataset. In *ICCV Workshop*, 2019.
- Yuqian Fu, Yanwei Fu, and Yu-Gang Jiang. Meta-fdmixup: Cross-domain few-shot learning guided by labeled target data. In *ACM Multimedia*, 2021.
- Yuqian Fu, Yu Wang, Yixuan Pan, Lian Huai, Xingyu Qiu, Zeyu Shangguan, Tong Liu, Yanwei Fu, Luc Van Gool, and Xingqun Jiang. Cross-domain few-shot object detection via enhanced open-set object detector. In *ECCV*, 2024.
- Yuqian Fu, Runze Wang, Yanwei Fu, Danda Pani Paudel, Xuanjing Huang, and Luc Van Gool. Objectrelator: Enabling cross-view object relation understanding in ego-centric and exo-centric videos. *ICCV*, 2025.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *CVPR*, 2022.
- Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *CVPR*, 2024.
- Yuliang Guo, Sparsh Garg, S. Mahdi H. Miangoleh, Xinyu Huang, and Liu Ren. Depth any camera: Zero-shot metric depth estimation from any camera. In *CVPR*, 2025.

- Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *ECCV*, 2020.
- Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, Li-hang Pan, et al. Glm-4.1 v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. pp. *arXiv-2507*, 2025.
- Yifei Huang, Guo Chen, Jilan Xu, Mingfang Zhang, Lijin Yang, Baoqi Pei, Hongjie Zhang, Lu Dong, Yali Wang, Limin Wang, et al. Egoexolearn: A dataset for bridging asynchronous ego-and exo-centric view of procedural activities in real world. In *CVPR*, 2024.
- Bernardo Iraci. *Blender cycles: lighting and rendering cookbook*. Packt Publishing Ltd, 2013.
- Baoxiong Jia, Ting Lei, Song-Chun Zhu, and Siyuan Huang. Egotaskqa: Understanding human tasks in egocentric videos. *NeurIPS*, 2022.
- Simar Kareer, Dhruv Patel, Ryan Punamiya, Pranay Mathur, Shuo Cheng, Chen Wang, Judy Hoffman, and Danfei Xu. Egomimic: Scaling imitation learning via egocentric video. In *ICRA*, 2025.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, 2017.
- Guofa Li, Zefeng Ji, Xingda Qu, Rui Zhou, and Dongpu Cao. Cross-domain object detection for autonomous driving: A stepwise domain adaptative yolo approach. *IEEE Transactions on Intelligent Vehicles*, 2022.
- Jinlong Li, Runsheng Xu, Jin Ma, Qin Zou, Jiaqi Ma, and Hongkai Yu. Domain adaptive object detection for autonomous driving under foggy weather. In *WACV*, 2023a.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023b. URL <https://arxiv.org/abs/2301.12597>.
- Kailing Li, Qi’ao Xu, Tianwen Qian, Yuqian Fu, Yang Jiao, and Xiaoling Wang. Clivis: Unleashing cognitive map through linguistic-visual synergy for embodied visual reasoning. *arXiv preprint arXiv:2506.17629*, 2025a.
- Yanjun Li, Yuqian Fu, Tianwen Qian, Qi’ao Xu, Silong Dai, Danda Pani Paudel, Luc Van Gool, and Xiaoling Wang. Egocross: Benchmarking multimodal large language models for cross-domain egocentric video question answering. *arXiv preprint arXiv:2508.10729*, 2025b.
- Yu Li, Xingyu Qiu, Yuqian Fu, Jie Chen, Tianwen Qian, Xu Zheng, Danda Pani Paudel, Yanwei Fu, Xuanjing Huang, Luc Van Gool, et al. Domain-rag: Retrieval-guided compositional image generation for cross-domain few-shot object detection. *arXiv preprint arXiv:2506.05872*, 2025c.
- Zhuowan Li, Xingrui Wang, Elias Stengel-Eskin, Adam Kortylewski, Wufei Ma, Benjamin Van Durme, and Alan L Yuille. Super-clevr: A virtual benchmark to diagnose domain robustness in visual reasoning. In *CVPR*, 2023c.
- Gaowen Liu, Hao Tang, Hugo M Latapie, Jason J Corso, and Yan Yan. Cross-view exocentric to egocentric video synthesis. In *ACM Multimedia*, 2021.
- Yang Liu, Ming Ma, Xiaomin Yu, Pengxiang Ding, Han Zhao, Mingyang Sun, Siteng Huang, and Donglin Wang. Ssr: Enhancing depth perception in vision-language models via rationale-guided spatial reasoning. *arXiv preprint arXiv:2505.12448*, 2025.
- Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *NeurIPS*, 2021.
- Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *NeurIPS*, 2023.
- Yang Miao, Francis Engelmann, Olga Vysotska, Federico Tombari, Marc Pollefeys, and Dániel Béla Baráth. Scene-graphloc: Cross-modal coarse visual localization on 3d scene graphs. 2024.

- Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.
- Boxiao Pan, Zhangjie Cao, Ehsan Adeli, and Juan Carlos Niebles. Adversarial cross-domain action recognition with co-attention. In *AAAI*, 2020.
- Luigi Piccinelli, Christos Sakaridis, Mattia Segu, Yung-Hsu Yang, Siyuan Li, Wim Abbeloos, and Luc Van Gool. UniK3D: Universal camera monocular 3d estimation. In *CVPR*, 2025.
- Chiara Plizzari, Alessio Tonioni, Yongqin Xian, Achin Kulshrestha, and Federico Tombari. Omnia de egotempo: Benchmarking temporal understanding of multi-modal llms in egocentric videos. In *CVPR*, 2025.
- Shraman Pramanick, Yale Song, Sayan Nag, Kevin Qinghong Lin, Hardik Shah, Mike Zheng Shou, Rama Chellappa, and Pengchuan Zhang. Egovlpv2: Egocentric video-language pre-training with fusion in the backbone. In *ICCV*, 2023.
- Alexander Raistrick, Lahav Lipson, Zeyu Ma, Lingjie Mei, Mingzhe Wang, Yiming Zuo, Karhan Kayan, Hongyu Wen, Beining Han, Yihan Wang, Alejandro Newell, Hei Law, Ankit Goyal, Kaiyu Yang, and Jia Deng. Infinite photorealistic worlds using procedural generation. In *CVPR*, 2023.
- Xiaofeng Ren and Chunhui Gu. Figure-ground segmentation improves handled object recognition in egocentric video. In *CVPR*, 2010.
- Swathikiran Sudhakaran, Sergio Escalera, and Oswald Lanz. Lsta: Long short-term attention for egocentric action recognition. In *CVPR*, 2019.
- Suraj Jyothi Unni, Raha Moraffah, and Huan Liu. Vqa-gen: A visual question answering benchmark for domain generalization. *arXiv preprint arXiv:2311.00807*, 2023.
- Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *CVPR*, 2025a.
- Zirui Wang, Wenjing Bian, Xinghui Li, Yifu Tao, Jianeng Wang, Maurice Fallon, and Victor Adrian Prisacariu. Seeing in the dark: Benchmarking egocentric 3d vision with the oxford day-and-night dataset. *arXiv preprint arXiv:2506.04224*, 2025b.
- Junbin Xiao, Nanxin Huang, Hao Qiu, Zhulin Tao, Xun Yang, Richang Hong, Meng Wang, and Angela Yao. Egoblind: Towards egocentric visual assistance for the blind people. *arXiv preprint arXiv:2503.08221*, 2025.
- Gu Xin, Fan Heng, Huang Yan, Luo Tiejian, and Zhang Libo. Context-guided spatio-temporal video grounding. 2024.
- Jingkang Yang, Shuai Liu, Hongming Guo, Yuhao Dong, Xiamengwei Zhang, Sicheng Zhang, Pengyun Wang, Zitang Zhou, Binzhu Xie, Ziyue Wang, et al. Egolife: Towards egocentric life assistant. In *CVPR*, 2025.
- Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024a.
- Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv:2406.09414*, 2024b.
- Hanrong Ye, Haotian Zhang, Erik Daxberger, Lin Chen, Zongyu Lin, Yanghao Li, Bowen Zhang, Haoxuan You, Dan Xu, Zhe Gan, et al. Mm-ego: Towards building egocentric multimodal llms for video qa. *arXiv preprint arXiv:2410.07177*, 2024.
- Ganlin Zhang, Deheng Zhang, Longteng Duan, and Guo Han. Accessible robot control in mixed reality, 2023a. URL <https://arxiv.org/abs/2306.02393>.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023b.

- Haoyu Zhang, Qiaohui Chu, Meng Liu, Yunxiao Wang, Bin Wen, Fan Yang, Tingting Gao, Di Zhang, Yaowei Wang, and Liqiang Nie. Exo2ego: Exocentric knowledge guided mllm for egocentric video understanding. *arXiv preprint arXiv:2503.09143*, 2025a.
- Yao Zhang, Haokun Chen, Ahmed Frikha, Denis Krompass, Gengyuan Zhang, Jindong Gu, and Volker Tresp. Cl-cross vqa: A continual learning benchmark for cross-domain visual question answering. In *WACV*, 2025b.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL <http://arxiv.org/abs/2403.13372>.
- Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle. *arXiv preprint arXiv:2104.02008*, 2021.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. 2022.
- Dong Zhuo, Wenzhao Zheng, Jiahe Guo, Yuqi Wu, Jie Zhou, and Jiwen Lu. Streaming 4d visual geometry transformer. *arXiv preprint arXiv:2507.11539*, 2025.



## A APPENDIX

### A.1 MORE VIDEO SOURCE CONSTRUCTION DETAILS

**EgoNight-Synthetic Construction.** For EgoNight-Synthetic Construction, we first use the coarse progressive generation method with a fast solver in infinigen Raistrick et al. (2023) to generate 3D scenes in Blender format. Then, a human annotator will edit the scene in the following sequence:

- Explore and edit the scene to remove unreasonable cases and make the indoor scene as natural as possible.
- Add light source in the scene if the generated scene does not include enough illumination to create enough illumination gap between the day and night.
- Record camera trajectory by exploring the whole indoor scene.
- Change the camera model and resolution. Set rendering samples and frames. For all synthetic dataset, we use the Blender build-in Panoramic Fisheye Equisolid camera with Lens 10.5 and field of view 180°.
- Create night scene by modifying the light source, motion blur, and environment map.
- Render the day and night pair using Blender Iraci (2013).

During the dataset construction, we apply home light source during night for 30 scenes and spot light source for 20 scenes to simulate torch light in real life. To create different difficulty levels, we apply different rendering sample size (higher sample size gives lower noise in the final image), spot light size, and motion blur to part of the data as shown in Tab. 6. We also show different modality and difficulty level in Fig. 6

difficulty level	sample size	light condition	motion blur shutter
Easy	4096	105°/ few light on	-
Medium	512	40°-50°spot light / all light off	-
Hard	512	40°-50°spot light / all light off	1-2

Table 6: Difficulty level and corresponding rendering settings.

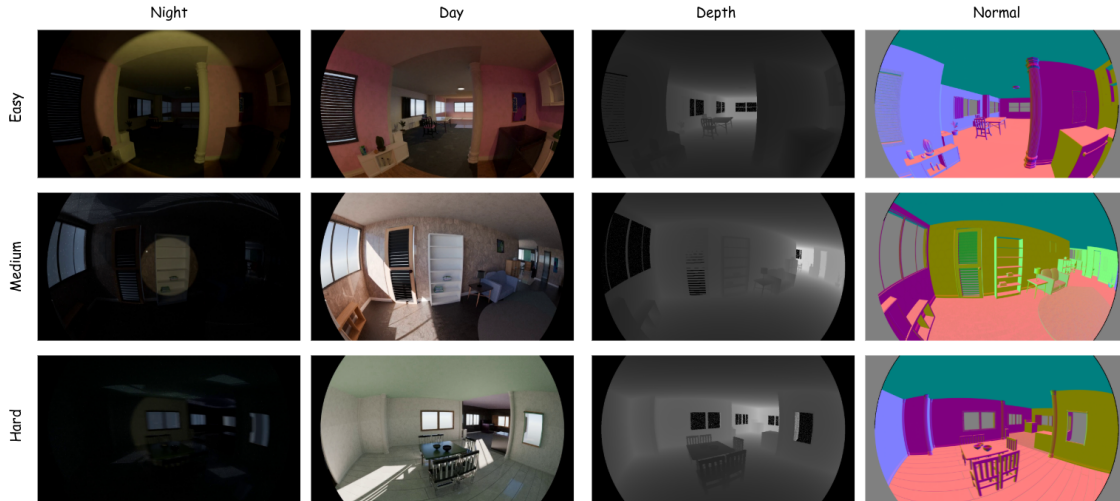


Figure 6: More examples and modalities of synthetic datasets.

**EgoNight-Sofia Construction.** In total, four participants were involved. The recording setup included three different GoPros, a head-mounted rig to fix the camera on the forehead and mimic human-eye perspective, several phones for live preview or daytime video guidance, and diverse lighting sources such as flashlights, spotlights, and candles. The process followed a video-guided recording strategy, as introduced in Sec. 3.1: the ego-wearer first recorded a daytime video while previewing the live feed on a phone, and for the nighttime counterpart replayed the daytime video on the phone as guidance to replicate the same setup, walking speed, viewpoints, and actions. Videos were collected across a wide range of environments, including indoor scenes (apartments, workplaces, grocery shops, building receptions) and outdoor scenarios (fitness areas, tourist landmarks, and street views). Post-trimming was applied to each day-night pair to further ensure alignment. On average, it took around 2-3 hours to produce one paired data.

**EgoNight-Oxford Construction.** We credit the contribution of Seeing in the Dark dataset Wang et al. (2025b), which provides multiple sequences of egocentric videos in the night in a various environment. We built our EgoNight-VQA dataset partially upon this work. Firstly, We enumerated all nighttime clips in Oxford Day–Night and performed a two-stage filtering. (1) Screening for uniqueness of place. We cross-checked scene metadata (route notes/time stamps) to avoid repeated paths within the same landmark. (2) Stratified diversity & quality sampling. Remaining clips were scored on a 1–5 rubric along axes designed for egocentric, low-light evaluation: illumination type (ambient only / mixed artificial / high-contrast point sources), illumination hardness (soft vs. specular/point), exposure stability (auto-gain pumping, blown highlights), scene dynamics (pedestrians/vehicles, occlusions), camera motion pattern (walk, run, head turns), and task context (navigation, road-crossing, object interaction, signage reading). The final set comprises 20 sequences that maximize lighting/task diversity under egocentric night settings while avoiding place overlap and task overlap.

In total, collecting the three video sources required over 100 hours of human effort.

## A.2 MORE BENCHMARK IMPLEMENTATION DETAILS

As described in Sec. 3.2, our auto-labeling pipeline is QA-type specific and involves three customized prompts: one for captioning, one for question generation, and one for answer synthesis. The detailed prompts are shown in Fig. 8.

### A.2.1 EGO NIGHT-VQA HUMAN LABELING

We hired several participants to review and refine the QA pairs generated by our three-stage day-augmented auto-labeling pipeline, compensating them at a rate of €20 per video. Each participant was provided with a detailed labeling instruction document and an onboarding meeting to ensure the guidelines were clearly conveyed. A simplified version of the labeling tutorial is included as in Fig. 7.

#### Annotation Tutorial (Simplified)

**Read Me First:** Please follow the labeling pipeline carefully and complete each step as instructed. On average, annotating one video takes about 2 hours. Easier cases may take less time, but in general, each video should take no less than 1.5 hours to ensure high-quality annotations. (The first video may take longer, as you will need to familiarize yourself with the pipeline.) We will randomly check the labeled data afterward, and annotators will be required to refine their work if the quality does not meet expectations.

**Step 1: Preparation.** You are expected to first download the paired `day.mp4` and `night.mp4` videos (aligned in time, except for unpaired tasks), together with the QA text file (`.txt`), which contains candidate QAs grouped by QA type (e.g., `counting.txt`). Before starting annotation, you should carefully watch both the daytime and nighttime videos to fully understand the scenario and activities.

**Step 2: QA Verification and Refinement.** For each QA pair, you should apply one of three operations:

- **Delete:** You should remove QAs that are meaningless, vague, irrelevant, duplicated, or inconsistent between day–night pairs (for paired QA types).
- **Modify:** If the question is reasonable but the answer is incorrect (e.g., counting errors, wrong action duration), you should correct the answer. You may also rephrase the question to eliminate ambiguity (e.g., clarifying “left/right” as relative to the ego-wearer).
- **Add:** If too many pairs are deleted, or if you notice interesting and challenging cases missing, you need to add new QAs. This is especially important for low-frequency tasks, e.g., dynamic detection or counting of dynamics.

**Step 3: Special Cases.**

- For paired QA types (e.g., object recognition, spatial reasoning), you must ensure the same QAs apply to both day and night videos.
- For unpaired QA types (e.g., lighting recognition, dynamic detection), you only need to ensure correctness on the nighttime video.
- For dynamic events, you are expected to specify temporal spans, e.g., *Q: “Around which time does a red car pass by?” A: “At frames 4–6.”*

**Step 4: Post-processing.** Once QAs were validated, the answer field should be renamed from `''answer''` to `''human_answer''`.

**Appendix: Paired & Non-Paired QA Types.** The same as described in the main file.

Figure 7: Simplified version of annotation tutorial.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

914  
915  
916  
917

914  
915  
916  
917

annotator statistics, we select randomly VQA pairs from 20 videos out of 90, and add four human annotators without GPT annotation. The average pairwise cosine similarity normalized to (0, 1) based on language feature extracted by BLIP-2 Li et al. (2023b) is 0.8458. This shows that human annotators are in general consistent in interpreting scenes. And answers are semantically aligned, even when worded differently.

### A.3 COMPARISONS WITH PRIOR EGOCENTRIC VQA BENCHMARKS.

In Tab. 7, we compare our EgoNight-VQA with prior egocentric VQA benchmarks, including EgoVQA Fan (2019), EgoTaskQA Jia et al. (2022), EgoSchema Mangalam et al. (2023), EgoThink Cheng et al. (2024), EgoTempo Plizzari et al. (2025), EgoCross Li et al. (2025b), EgoMemoria Ye et al. (2024), HourVideo Chandrasegaran et al. (2024), and EgoLifeQA Yang et al. (2025), listing their lighting conditions (mainly daytime or nighttime), video duration length, the number of testing QA examples, number of QA type categories, if temporal-oriented tasks are included or emphasized, and the evaluation metric.

We highlight that EgoNight-VQA is the first to explore nighttime egocentric VQA with aligned day–night video pairs.

Dataset	Lighting	Video Length	# Test	# Categories	Temporal	Metric Type
EgoVQA	☀	(25s, 100s)	250	3	✗	OpenQA
EgoTaskQA	☀	25s	8k	4	✗	OpenQA
EgoSchema	☀	180s	500	-	✗	CloseQA
EgoThink	☀	-	750	12	✗	OpenQA
EgoTempo	☀	45s	500	10	✓	OpenQA
EgoCross	☀	23s	957	15	✓	CloseQA & OpenQA
EgoBlind	☀	(0s, 120s)	5311	6	✓	OpenQA
EgoMemoria	☀	(30s, 1h)	7026	-	✓	CloseQA
HourVideo	☀	(20min, 120min)	12976	-	✓	CloseQA
EgoLifeQA	mostly ☀	44.3 h	6000	5	✓	CloseQA
<b>EgoNight-VQA</b>	Aligned ☀ & ☾	(24s, 214s)	3658	12	✓	OpenQA

Table 7: Comparison between EgoNight-VQA and prior egocentric VQA benchmarks. ☀ means daytime, while ☾ indicates nighttime.

### A.4 MORE EGO NIGHT-VQA EXPLANATIONS AND EXAMPLES

#### A.4.1 QA TYPE DEFINATION

We present the 12 QA types with their detailed definitions in Tab. 8.

Table 8: **Detailed descriptions of the 12 QA types in EgoNight-VQA.** Paired QA types share the same QAs across day–night counterparts, while unpaired QA types are evaluated only on nighttime videos.

QA Type	Attribute	Description
Object Recognition	Paired	Identify and recognize specific objects in the scene (e.g., “What is on the table?”).
Text Recognition	Paired	Read and interpret visible text or logos (e.g., “What does the sign say?”).
Spatial Reasoning	Paired	Understand spatial relations between objects (e.g., “What is left of the chair?”).
Scene Sequence	Paired	Recall the temporal order of visited scenes (e.g., “Which room did I enter after the kitchen?”).
Navigation	Paired	Working as a navigation assistant after watched the whole video (e.g., “How can I reach place B from place A?”).
Counting of Statics	Paired	Count static objects visible in the scene (e.g., “How many chairs are in the room?”).
Action Recognition	Paired	Identify human actions or interactions (e.g., “What action is being performed?”).
Non-Common-Sense Reasoning	Paired	Judge unusual or physically implausible cases, for synthetic videos. (e.g., “Is the door embedded inside the wall?”).
Lighting Recognition	Unpaired	Recognize the illumination source, also include counting. (e.g., “How many light sources are in the room?”).
Lighting Change	Unpaired	Detect changes in lighting conditions (e.g., “Did the light turn off during the clip?”).
Dynamic Detection	Unpaired	Detect dynamic moving objects (e.g., “Is a car/person moving across the scene?”).
Counting of Dynamics	Unpaired	Count the number of dynamic objects or events (e.g., “How many people walked by?”).

#### A.4.2 QA EXAMPLES

We show more QA examples of EgoNight-Synthetic in Fig. 9, EgoNight-Sofia in Fig. 10, and EgoNight-Oxford in Fig. 11. For those paired QA types, we show both day and night frames, while for those unpaired QA types, we demonstrate nighttime frames only. Three frames are shown if the QA is spatial or static related, while more frames are given if the QA is temporal or more dynamic related.



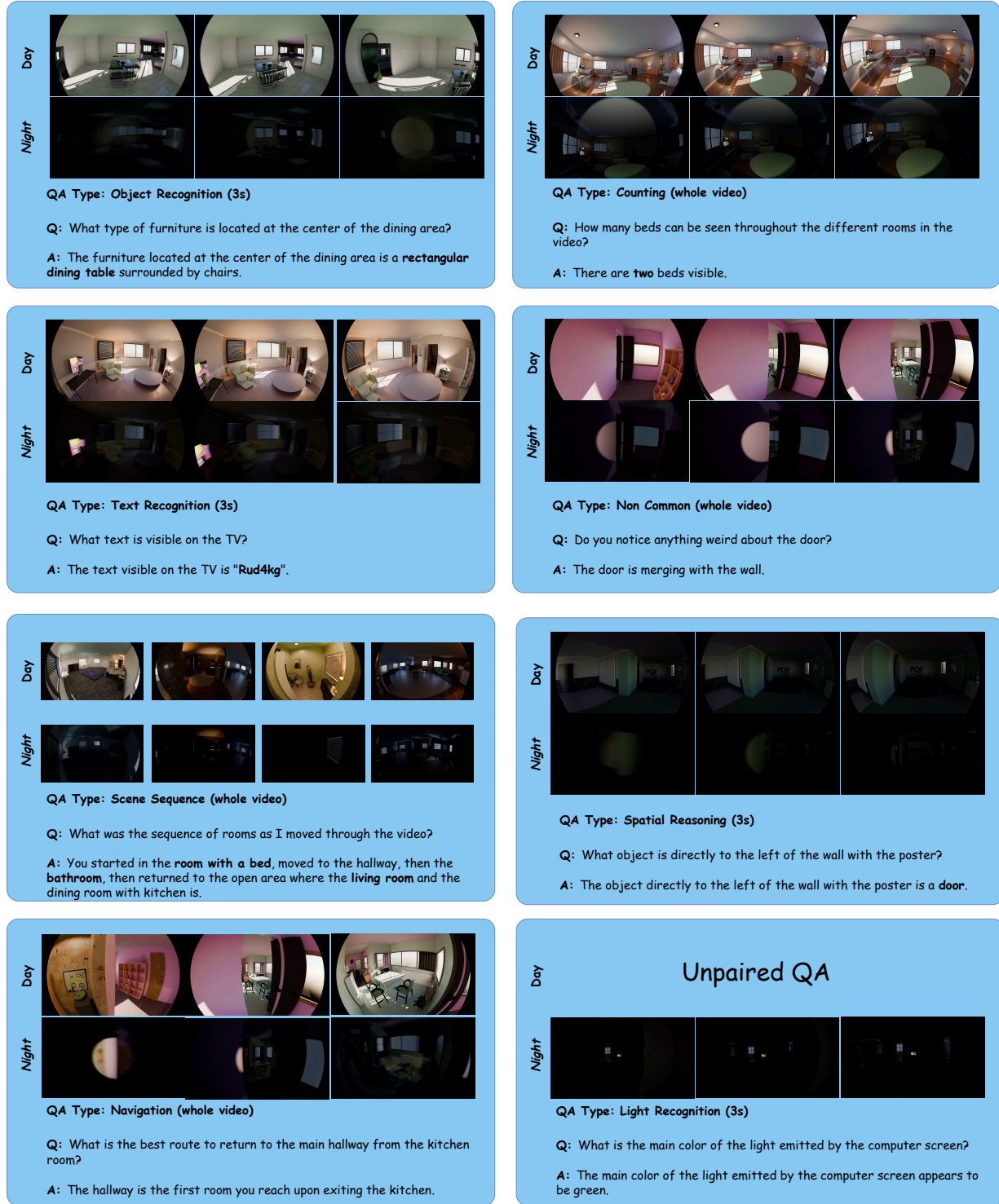


Figure 9: More QA examples from EgoNight-Synthetic dataset.

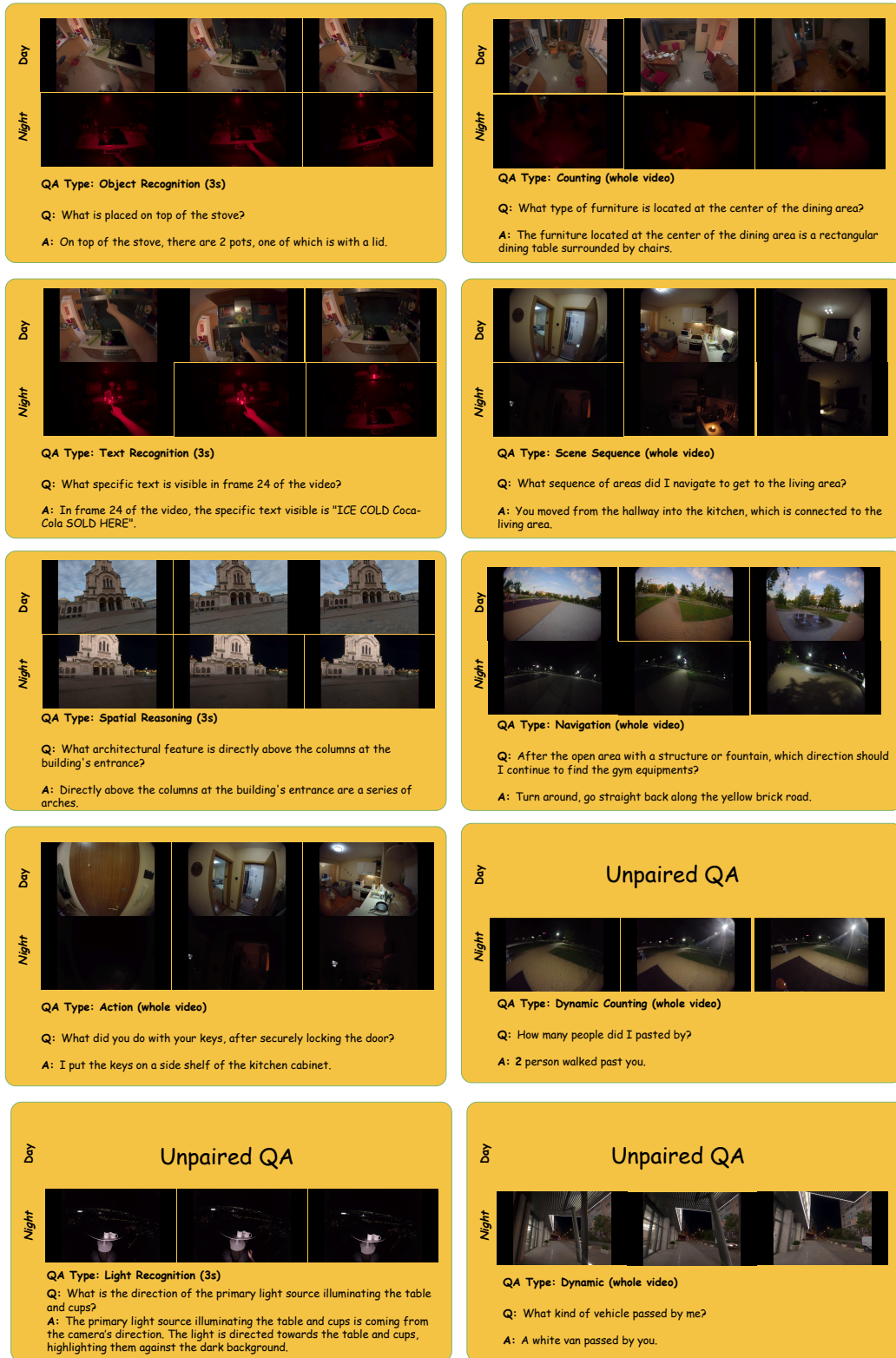


Figure 10: More QA examples from EgoNight-Sofia dataset.

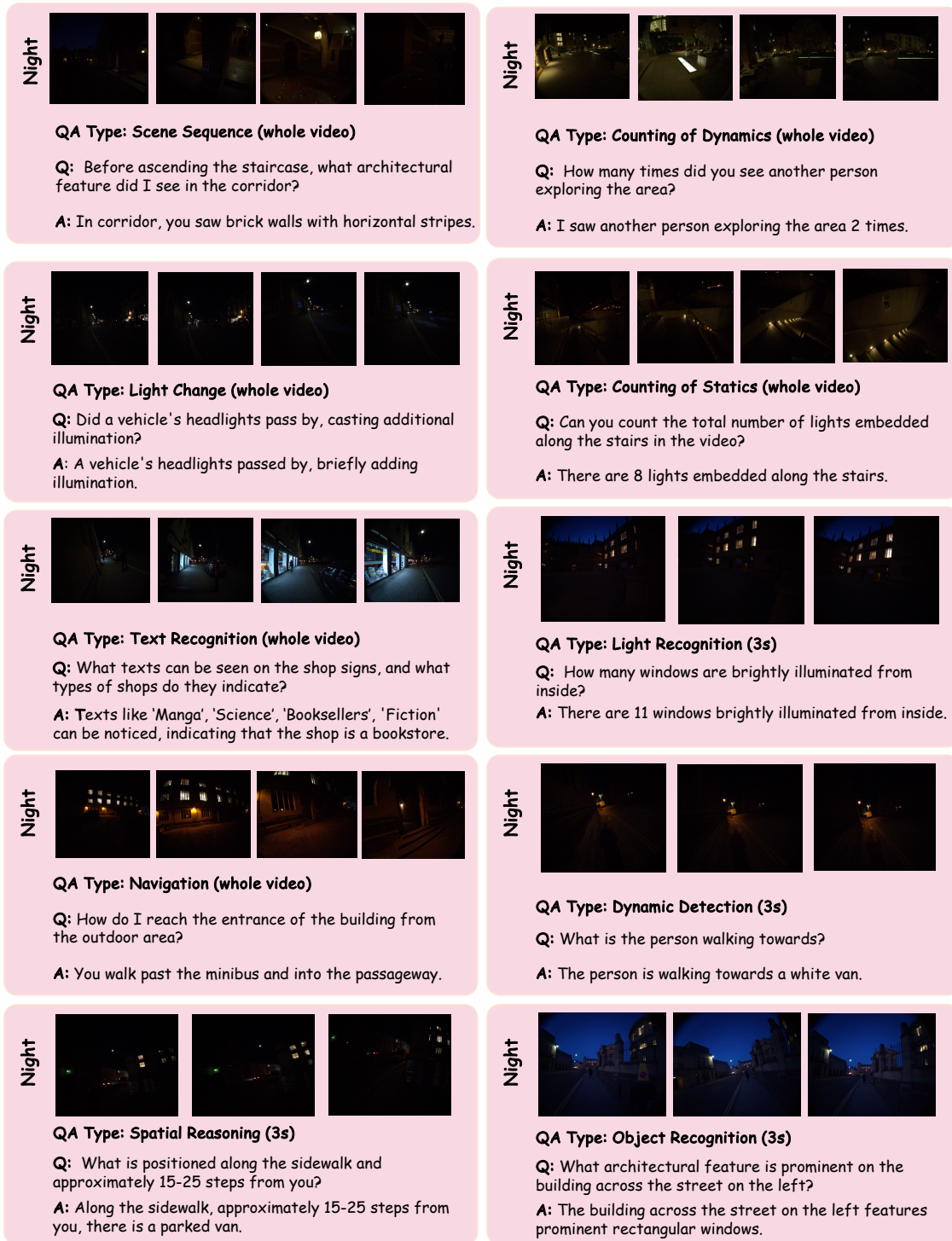


Figure 11: More QA examples from EgoNight-Oxford dataset.

## A.5 MORE EXPERIMENTS SETUPS

### A.5.1 SETUPS FOR EGO NIGHT-VQA EXPERIMENTS.

In this section, we describe the model setup, how GPT is used as the judge, the prompting strategy, the GPU resources, and the approximate runtime for each dataset. For the closed-source model, we directly use the API call.

Model	Inference Speed (min)
GPT-4.1	<5
Gemini 2.5 Pro	<5
InternVL3-8B	<5
Qwen2.5-VL-72B	25
Qwen2.5-VL-7B	<5
Qwen2.5-VL-3B	<5
GLM-4.1V-9B-Base	<5
VideoLLaMA3-7B	<5
LLaVA-NeXT-Video-7B	50
EgoGPT	<5

Table 9: Inference speed for different models (per Video).

For open source models, we use LLama-Factory Zheng et al. (2024) except VideoLLama3 Zhang et al. (2023b) and EgoGPT Yang et al. (2025). We use NVIDIA A6000 GPUs for all the model, except for Qwen2.5-VL-72B, we use 2 NVIDIA H200 GPUs for larger GPU memory. The inference speed for each model is shown in Tab. 9. Video frames are sampled at 2 fps for EgoNight-Synthetic, and 1 fps for EgoNight-Sofia and EgoNight-Oxford, without imposing a maximum frame limit. To further ensure fairness and consistency, the exact prompts used for each task are provided below.

**Model Evaluation Prompt.** For evaluating the language model, we use the following prompt:

Please carefully read the question, use the visual cues in the {video} to answer the question: {question}.

The original FPS of the video is {original\_video\_fps}. This image set is obtained by sampling at {sampling} fps.

Do not include any other content. You need to answer the question in any case and not demand additional context information.

Note: All the actions mentioned refer to the person who recorded the video.

**Evaluation Protocol.** Since the questions and answers are open-ended, we utilize GPT-4.1 Achiam et al. (2023) as a judge. Here is the prompt for evaluating the score given the model prediction, the ground truth answer, and the corresponding question:

role: system,  
content: You are an intelligent chatbot designed for evaluating the correctness of AI assistant predictions for question-answer pairs.  
Your task is to compare the predicted answer with the ground-truth answer and determine if the predicted answer is correct or not. Here's how you can accomplish the task:

**INSTRUCTIONS:**

1. Focus on the correctness and accuracy of the predicted answer with the ground-truth.
2. Consider uncertain predictions, such as 'it is impossible to answer the question from the video', as incorrect, unless the ground truth answer also says that.

role: user,  
content: Please evaluate the following video-based question-answer pair:  
Question: {question}  
Ground truth correct Answer: {answer}  
Predicted Answer: {predicted\_answer}

Provide your evaluation as a correct/incorrect prediction along with the score where the score is an integer value between 0 (fully wrong) and 5 (fully correct). The middle score provides the percentage of correctness. For question that counting the number of objects, if the predicted answer falls in the range of the ground truth answer, it should be considered as correct. Please generate the response in the form of a Python dictionary string with keys 'pred', 'score' and 'reason', where value of 'pred' is a string of 'correct' or 'incorrect', value of 'score' is in INTEGER, not STRING and value of 'reason' should provide the reason behind the decision."

To further validate the LLM-as-a-Judge strategy, we divide all annotations into two groups: (a) answers verified but not modified by humans (preserving the GPT style) and (b) answers modified or created by humans. The ratio is approximately 4 : 6, indicating a high human modification rate. When evaluating the accuracy of GPT-4.1 separately in the two subsets, we obtain 26.73% for group (a) and 27.87% for group (b). The similar scores show that GPT-as-Judge does not prefer GPT-generated answers over human-authored ones. Further more, We randomly sampled 300 QA pairs

(questions, ground truth, model answers, and the corresponding LLM-assigned scores) and asked human evaluators to judge whether each score from the LLM was correct. This yielded an agreement rate of 95.67%, indicating strong alignment between human judgment and the LLM-as-a-Judge decisions, and thus demonstrating the reliability of the LLM-based evaluation.

#### A.5.2 SETUPS FOR DAY-NIGHT CORRESPONDENCE RETRIEVAL.

In this section, we describe the setup of the model, method, e.g. how feature-based retrieval for vision encoders, how prompt VLMs, metric, result, GPU, cost time, and other details.

**i) Spatial Retrieval (Place Recognition).** For feature-based methods Oquab et al. (2023); Bolya et al. (2025), we calculate the CLS tokens  $f^i$  of each frame within the video clip with the vision encoder, with  $i$  the frame index in the clip. Then, the "best matching" strategy is implemented to calculate the similarity between the query clip  $v_q$  and the database clip  $v_d$ . The best cosine similarity between the features of the query clip  $f_q^i$  and the database clip  $f_d^j$ ,

$$\sigma(v_q, v_d) = \max_{i \in [0, s-1], j \in [0, s-1]} \cos(f_q^i, f_d^j). \quad (1)$$

The database video clips are ordered based on the similarity  $\sigma$  and then the most similar clip is retrieved.

Similarly, for MLLM-based methods Achiam et al. (2023); Chen et al. (2024b), we ask the MLLM to assess the "pairwise" similarity between each query-database pair and order the database clips by similarity. The prompt to the MLLM is as follows:

*You are given two video clips from different scenes.  
Your task is to evaluate how similar these two scenes are based on their spatial layout, furniture, objects, architectural features, and overall room structure.*

**CLIP STRUCTURE:**  
 – Images 1– $(s-1)$  from Query Scene  
 – Images  $s$ – $\{2s-1\}$  from Database Scene

**TASK:**  
*Please carefully analyze and compare the spatial layout, furniture placement, objects, architectural features, and overall room structure between these two video clips.*

**IMPORTANT:** *Please respond with ONLY a single numerical similarity score between 0.0 and 1.0, where:*  
 – 0.0 = Completely different scenes (different rooms/locations)  
 – 1.0 = Identical or nearly identical scenes (same room/location)  
 – Values in between represent varying degrees of similarity

*Example responses: "0.85", "0.23", "0.67"*  
*1.0 should be used when the two scenes are identical, so don't use 1.0 if the two scenes are not 100% identical.*  
*Please provide only the numerical score without any additional text or explanation.*

It is noticeable that existing MLLMs have difficulty in processing long-horizon and multi-scene videos. We also conduct the "all-in-one-prompt" experiments by inputting all the images of the query clip and the database clips in one prompt and asking the MLLM to output the ordered database clips. The "all-in-one-prompt" strategy leads to largely degraded performance, as shown in Tab. 10.

Prompt Strategy	Spatial Retrieval R@1 - Synthetic	
	Day $\rightarrow$ Day	Night $\rightarrow$ Day
Pairwise	75.6	54.1
All-in-one	10.5	28.5

Table 10: Ablation on prompting for night-to-day spatial retrieval task.

**ii) Temporal Localization.** The mIoU metric is defined as:

$$\text{mIoU} = \frac{1}{M} \sum_{m=1}^M \frac{|[t_i, t_j] \cap [t_i^*, t_j^*]|}{|[t_i, t_j] \cup [t_i^*, t_j^*]|}, \quad (2)$$



where  $M$  denotes the total number of meta-tasks (1000 in our setup). For feature-based temporal localization, we apply the "best-match" strategy similar as spatial localization, localizing the query clip to the frame stamp with the best clip-to-clip similarity:

$$i = \arg \max_i \sigma(v_q, v_d), v_d = v_D[i : i + s], \quad (3)$$

$v_D$  is the parent full video and the end frame will be  $i + s - 1$ . For MLLM-based method, we input the query clip and the parent full video in the prompt and ask the MLLM to output the start and end frame of the query within the full video. The prompt is as follows:

*You are given a query video clip and a complete video sequence from the same scene.  
Your task is to find the exact temporal position where the query clip appears in the complete video sequence.*

**IMPORTANT CONTEXT:**

- The query clip shows  $s$  consecutive frames from a video sequence
- The complete video sequence shows ALL frames from the same scene in chronological order
- The query clip appears as a consecutive subsequence somewhere within the complete video sequence
- You need to find the exact start and end frame numbers where this subsequence appears

**IMAGE STRUCTURE:**

- Images 1– $s$ : Query video clip (consecutive frames to find)
- Images  $\{s+1\}$ – $\{s+1+\text{video\_len}\}$ : Complete video sequence (all frames in chronological order)

**TOTAL IMAGES:**  $\{\text{query\_count} + \text{database\_count}\}$  images

**TASK:**

1. Look at the query clip to understand what sequence you're looking for
2. Search through the complete video sequence to find where this exact sequence appears
3. The query sequence should appear as consecutive frames in the complete video sequence
4. Pay attention to camera movements, object positions, and scene changes to identify the matching sequence

**FRAME NUMBERING:**

- The complete video sequence frames are numbered from  $\{\min(\text{database\_frame\_numbers})\}$  to  $\{\max(\text{database\_frame\_numbers})\}$
- You need to return the actual frame numbers from this range

**RESPONSE FORMAT:**

*Respond with ONLY two numbers separated by a comma: "start\_frame,end\_frame"*

- *start\_frame*: The frame number where the query clip begins in the complete video sequence
- *end\_frame*: The frame number where the query clip ends in the complete video sequence

*Example: If the query clip appears at frames 15–19 in the complete sequence, respond: "15,19"*

*Valid frame range:  $\{\min(\text{database\_frame\_numbers})\}$  to  $\{\max(\text{database\_frame\_numbers})\}$*

### A.5.3 SETUPS FOR EGOCENTRIC DEPTH ESTIMATION AT NIGHT.

We evaluate four off-the-shelf monocular depth systems without night-specific fine-tuning. For each, we highlight features pertinent to our setting. (F) denotes support for fisheye egocentric images; (U) denotes undistorted/pinhole images.

1. **Depth Anything V2 (metric).** (U) Foundation MDE model (DPT head with DINOv2 backbone) trained on large-scale synthetic labels plus pseudo-labeled real images. We use the *official metric* checkpoints: *Indoor* (Hypersim-tuned) for indoor frames and *Outdoor* (VKITTI2-tuned) for outdoor frames. Outputs metric depth in meters and is known for strong zero-shot generalization.
2. **StreamVGGT.** (U) A causal/streaming transformer for video geometry that processes frames sequentially with state caching to improve temporal consistency and enable real-time inference. We run it in streaming mode to obtain per-frame depth on egocentric sequences.
3. **Depth Any Camera (DAC).** (F) Zero-shot *metric* depth across diverse camera models via a unified ERP (equirectangular) representation with pitch-aware image-to-ERP conversion and FoV alignment. We use the official release with default settings on our pinhole inputs.
4. **UniK3D.** (F) Universal-camera monocular 3D estimation with a spherical 3D formulation and a learned "pencil-of-rays" camera module, enabling accurate metric depth across pinhole, fisheye, and panoramic views. We run the official model in eval mode; when available, we provide intrinsics for pinhole frames.

Model	Object Rec.	Text Rec.	Spatial	Scene Seq.	Nav.	Light Rec.	Counting	Non-Common	Overall
<i>Closed-Source MLLMs</i>									
Gemini	25.94	39.39	32.43	35.47	30.77	31.97	21.88	15.15	28.34
GPT-4.1	25.19	54.55	35.42	28.44	27.09	35.25	20.83	16.67	27.75
<i>Open-Source MLLMs</i>									
InternVL3-8B	17.29	10.61	28.34	20.80	10.37	18.03	16.93	21.21	18.97
Qwen2.5-VL-72B	15.41	16.67	28.88	21.41	7.02	12.30	10.94	21.21	17.15
Qwen2.5-VL-7B	6.77	13.64	17.98	11.93	9.03	15.57	14.06	18.69	13.26
Qwen2.5-VL-3B	7.89	22.73	17.17	14.37	11.37	8.20	13.02	12.63	13.06
GLM-4.1V-9B-Base	13.16	36.36	23.71	21.71	7.02	15.57	19.27	14.14	17.69
LLaVA-NeXT-Video-7B	5.26	10.61	10.08	4.59	3.01	16.39	4.69	9.60	6.85
VideoLLaMA3-7B	10.90	21.21	19.07	23.55	7.02	7.38	18.49	16.67	15.97
<i>Egocentric MLLMs</i>									
EgoGPT	6.02	19.70	18.53	19.88	8.36	8.20	17.71	17.68	14.79
<i>Average across all models</i>									
Average	13.38	24.55	23.16	20.21	12.11	16.89	15.78	16.36	17.38

Table 11: Night-time VQA accuracy (%) per model across all QA categories for EgoNight-Synthetic.

## A.6 MORE EXPERIMENTAL RESULTS

Model	Object Rec.	Text Rec.	Spatial	Scene Seq.	Action	Nav.	Light Rec.	Counting	Dyn. Light	Dynamic	Dyn. Count	Avg.
<i>Closed-Source MLLMs</i>												
GPT-4.1	24.44	33.78	41.32	30.09	38.10	27.27	38.33	24.62	30.00	13.33	25.00	31.06
Gemini	32.22	47.30	35.54	24.78	34.92	29.29	43.33	27.69	15.00	26.67	40.00	32.67
<i>Open-Source MLLMs</i>												
InternVL3-8B	16.67	17.57	33.88	24.78	22.22	21.21	20.00	22.31	25.00	6.67	15.00	22.61
Qwen2.5-VL-72B	14.44	21.62	36.36	18.58	19.05	25.25	18.33	13.85	20.00	6.67	20.00	20.99
Qwen2.5-VL-7B	7.78	10.81	21.49	18.58	11.11	18.18	1.52	15.38	9.52	6.25	15.00	14.02
Qwen2.5-VL-3B	11.11	8.11	23.97	13.27	11.11	16.16	10.00	15.38	5.00	13.33	10.00	14.16
GLM-4.1V-9B-Base	12.22	12.16	30.58	15.04	14.29	17.17	11.67	25.38	15.00	6.67	10.00	18.14
VideoLLaMA3-7B	3.33	5.41	13.22	11.50	6.35	9.09	8.33	18.46	10.00	20.00	15.00	10.68
LLaVA-NeXT-Video-7B	8.89	5.41	18.18	12.39	12.70	15.15	11.67	13.85	0.00	13.33	20.00	12.67
<i>Egocentric MLLMs</i>												
EgoGPT	9.18	3.41	19.26	16.54	6.67	7.96	10.00	14.97	0.00	0.00	10.00	11.47
<i>Average across all models</i>												
Average	13.99	16.31	27.29	18.53	17.45	18.53	17.16	19.13	12.94	11.26	18.00	18.76

Table 12: Night-time VQA accuracy (%) per model across all QA categories for EgoNight-Sofia.

Models	Object Rec.	Text Rec.	Spatial	Scene Seq.	Action	Nav.	Light Rec.	Counting	Dyn. Light	Dynamic	Dyn. Count	Avg.
<i>Closed-Source MLLMs</i>												
GPT-4.1	64.52	34.88	41.35	34.13	65.59	43.75	30.43	18.49	18.52	37.84	18.52	38.95
Gemini 2.5 Pro	56.14	46.51	38.30	27.27	59.55	32.65	31.11	18.97	23.33	37.14	3.70	34.83
<i>Open-source MLLMs</i>												
InternVL3-8B	40.00	27.91	18.68	14.66	36.78	20.83	6.98	9.91	6.67	37.14	7.41	20.57
Qwen2.5-VL-72B	38.18	39.53	29.67	13.79	22.99	23.96	16.28	17.12	16.67	11.43	11.11	22.07
Qwen2.5-VL-7B	9.09	23.26	20.88	11.21	10.34	15.62	9.30	11.71	3.33	11.43	18.52	13.35
Qwen2.5-VL-3B	14.55	13.95	9.89	13.79	19.54	20.83	6.67	10.81	3.70	17.14	6.98	13.62
GLM-4V	22.81	30.23	21.43	15.52	28.74	9.38	19.57	17.12	6.67	37.14	14.81	19.57
VideoLLaMA3-7B	20.00	11.63	15.38	9.57	10.34	7.29	9.52	9.01	0.00	17.65	7.41	10.81
LLaVA-NeXT-Video-7B	9.09	4.65	10.99	0.00	0.00	0.00	6.98	0.00	0.00	2.86	0.00	2.86
<i>Egocentric MLLMs</i>												
EgoGPT	21.13	27.08	14.14	3.42	14.61	6.25	11.11	6.25	21.62	18.92	17.86	12.44
<i>Average across all models</i>												
Avg.	29.81	25.98	22.32	14.55	27.16	18.09	14.96	12.01	10.75	22.95	10.33	19.04

Table 13: Night-time VQA accuracy (%) per model across all QA categories for EgoNight-Oxford.

In this section, we show models per QA accuracy on each dataset for EgoNight-Synthetic in Tab. 11, EgoNight-Sofia in Tab. 12, and EgoNight-Oxford in Tab. 13.



## A.7 MORE VISUALIZATION RESULTS

### A.7.1 EGO-NIGHT-VQA

We further investigate failure reasons by systematically altering illumination, motion blur, and camera noise (sampling quality) in synthetic videos, while keeping scene, annotations, and trajectory fixed in the synthetic dataset. The procedure of generating such difficulty level is described in Appendix. A.1, we calculate the averaged accuracy from Tab. 1 as 18.47% for Easy (moderately dark), 15.88% for Medium (very dark with camera noise), 12.95% for Hard (as dark as medium, with motion blur). From the results, we highlight that: ii) Lower illumination together with camera sensor noise leads to the most significant drop due to loss of contrast and missing fine details. ii) Motion blur further harms performance by causing temporal ambiguity and object shape distortion. iii) These results confirm that night conditions VQA is challenging, highlighting the need for specialized night-egocentric benchmarks. We also visualize more failure cases, and give more analysis in Appendix. A.7.1. Here we provide more failure cases, which compare the day and night VQA output in Fig. 12. This clearly shows the gap between the day and night video understanding. Also, we provide examples of question generated and the corresponding caption to show caption reasoning ability in Fig. 13. Here we provide a more detailed analysis of the failure reason:

- **Extreme Illumination (Counting – Static):** Due to strong red lighting, two doors in the scene become nearly invisible, causing the model to undercount objects.
- **Small Object Disappearance (Text Recognition):** The price tag becomes too small and poorly illuminated at night, making it unreadable and leading to text recognition failure.
- **Spatial Confusion from Limited View (Navigation):** Restricted field of view at night hides key spatial cues (e.g., landmarks, corridor orientation), causing incorrect navigation decisions.
- **Motion Blur (Action Recognition):** Fast hand movement introduces motion blur, making the model misinterpret the content displayed on the screen.

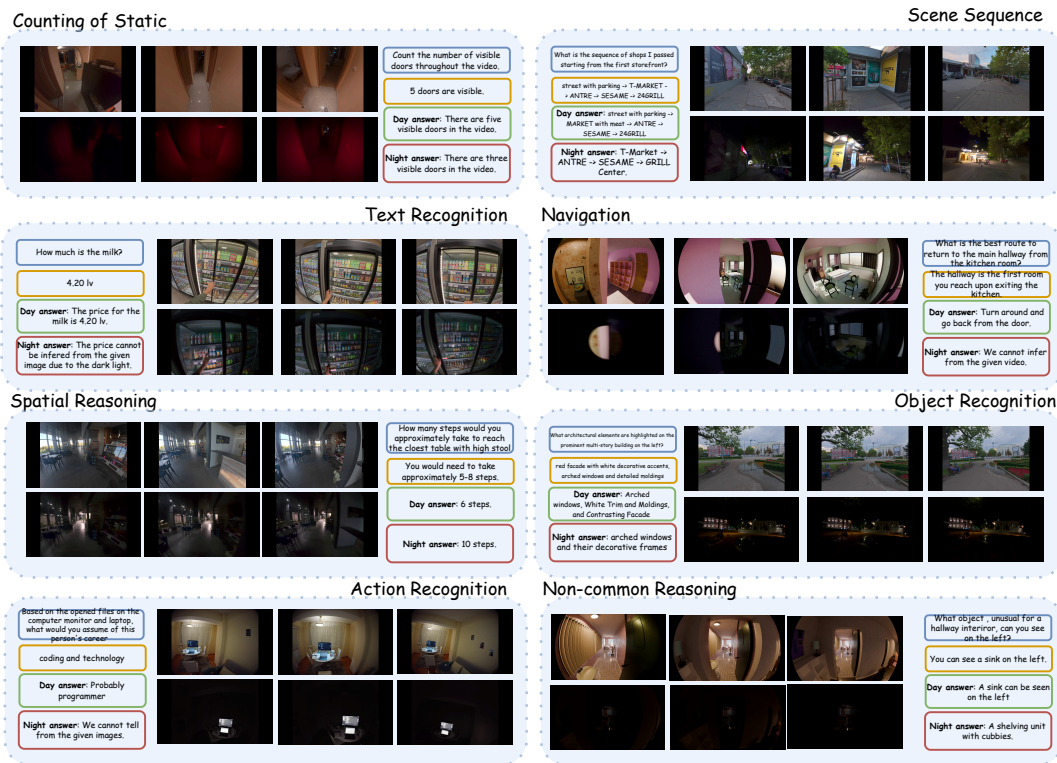


Figure 12: More QA examples with day and night answer produced by the same model.

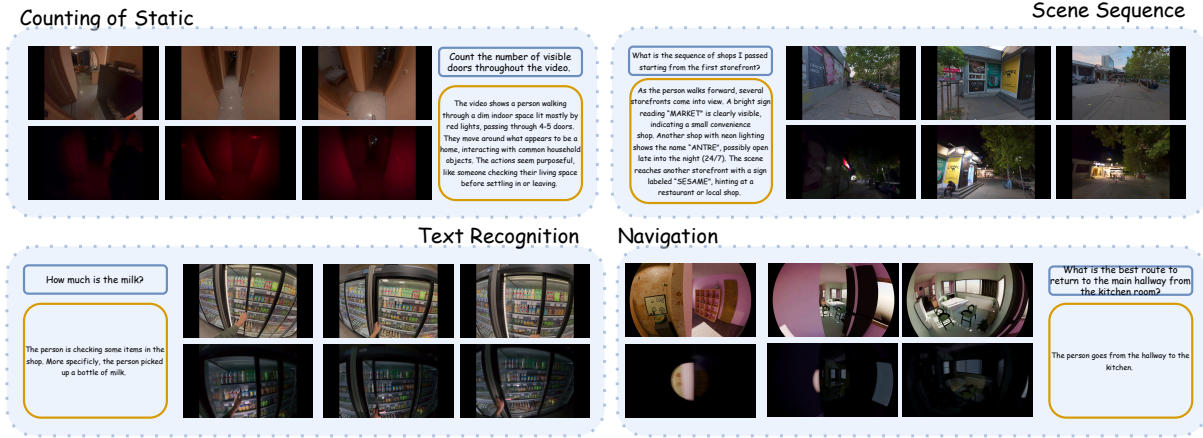


Figure 13: More examples shows the caption together with generated questions.

### A.7.2 DAY-NIGHT CORRESPONDENCE RETRIEVAL

We visualize the qualitative result on one meta sample of Night-to-Day spatial retrieval to better demonstrate the experiment setup and the performance of the benchmarked methods. As shown in Fig. 14, the light condition of the query video clip is drastically different from that of the database clips. Such a difference imposes a great challenge for existing methods in distinguishing the target scenes from the other candidate databases' clips, showing the value of the dataset in the place recognition task.



Methods	DB clip 1	DB 2	DB 3	DB 4	DB 5	DB 6	DB 7	DB 8	DB 9	DB 10
DINOv2	0.67	0.66	0.58	0.54	0.67	<b>0.70</b>	0.61	0.66	0.59	0.67
Percep. Enc.	0.83	0.82	0.78	0.81	0.81	0.80	0.79	0.81	0.78	<b>0.85</b>
GPT-4.1	<b>0.92</b>	0.72	0.62	0.18	0.18	0.18	0.18	0.12	0.12	0.12
InternVL 8B	<b>0.55</b>	0.32	0.32	0.30	0.30	0.25	0.25	0.21	0.15	0.15

Figure 14: Qualitative Result on one meta sample of spatial retrieval. The query video clip and the database video clips are visualized in the image. The table below the figure shows the similarity score between the query and the database clips calculated with different methods. The most similar one is in **bold**, and correct retrieval is in **green**, and the incorrect one is in **red**.

### A.7.3 EGOCENTRIC DEPTH ESTIMATION AT NIGHT

We provide additional qualitative results across day–night conditions (Figs. 15, 16, 17, 18). Consistent with the main paper, nighttime is substantially more challenging: low SNR, head-motion blur, extreme dynamic range, color/white-balance shifts, and auto-exposure fluctuations amplify scale ambiguity and erode edge fidelity, leading to over-smoothed surfaces, depth collapse in dark regions, halos around bright point sources, and temporal instability. UniK3D remains the strongest overall in preserving scene structure under these conditions, though performance still degrades under extreme darkness and sparse texture. By contrast, *StreamVGGT* and *DAC* are notably brittle at night, frequently washing out structure, misinterpreting specular highlights, and producing flattened or unstable depth in large low-illumination areas. The effect is most pronounced outdoors in EgoNight-Sofia and EgoNight-Oxford, where wide dynamic range, sparse texture, and point-light saturation further depress accuracy across methods.

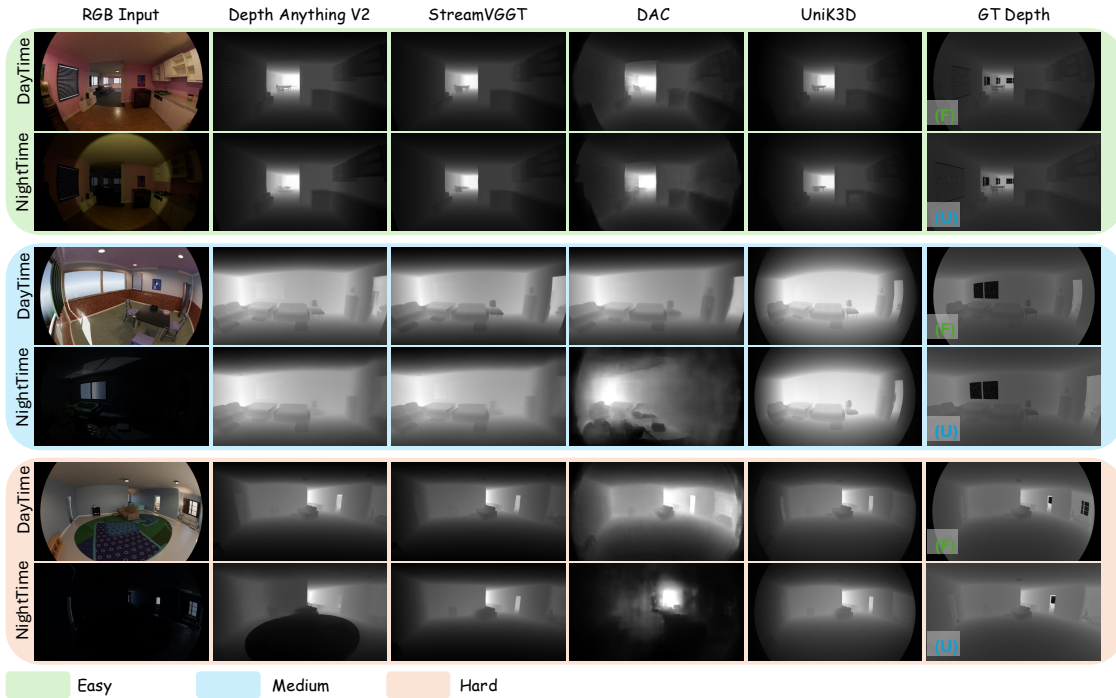


Figure 15: Qualitative results of monodepth estimation in day and night on EgoNight-Synthetic dataset according to different difficulty levels.

## A.8 MORE ANALYSIS

### A.8.1 FAILURECASE ANALYSIS

We further investigate failure reasons by systematically altering illumination, motion blur, and camera noise (sampling quality) in synthetic videos, while keeping scene, annotations, and trajectory fixed in the synthetic dataset. The procedure of generating such difficulty level is described in Appendix. A.1, we calculate the averaged accuracy from Tab. 1 as 18.47% for Easy (moderately dark), 15.88% for Medium (very dark with camera noise), 12.95% for Hard (as dark as medium, with motion blur). From the results, we highlight that: ii) Lower illumination together with camera sensor noise leads to the most significant drop due to loss of contrast and missing fine details. ii) Motion blur further harms performance by causing temporal ambiguity and object shape distortion. iii) These results confirm that night conditions VQA is challenging, highlighting the need for specialized night-egocentric benchmarks. We also visualize more failure cases, and give more analysis in Appendix. A.7.1.

### A.8.2 LIMITATIONS AND FUTURE WORKS

We acknowledge limitations of EgoNight and insights for future research directions in night video understanding. (1) The dataset scale remains modest compared to large-scale vision–language corpora. However, as a testbed, we argue that the current scale of 3,600+ human-verified QA pairs is already sufficient for benchmarking. In future work, we

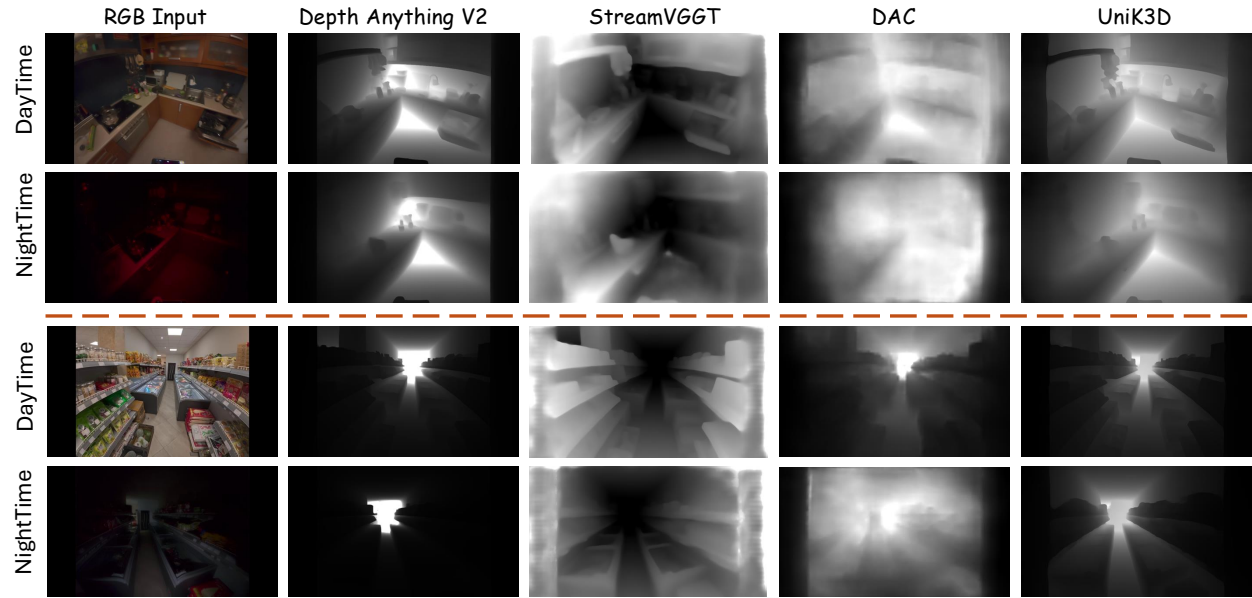


Figure 16: Qualitative results of monodepth estimation in day and night on EgoNight-Sofia dataset indoor part.

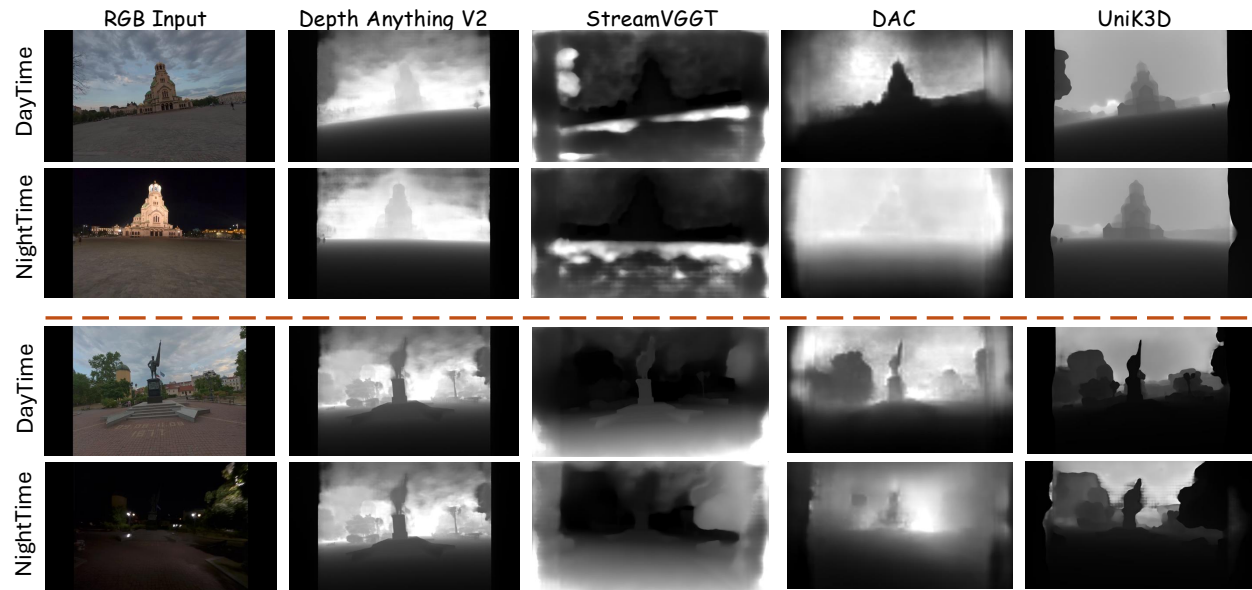


Figure 17: Qualitative results of monodepth estimation in day and night on EgoNight-Sofia dataset outdoor part.

plan to further scale up nighttime videos by synthesizing more data and recording additional real-world footage, which will enable not only benchmarking but also pretraining and fine-tuning to improve MLLM performance.

(2) We show in the main content that fine-tuning on synthetic data improves real world performance. Therefore, we can scale-up synthetic data to build a training set, that can be used to fine-tune the model, thus generalize to real world scenario.

(3) We encourage the community to explore broader research avenues, such as:

- Leveraging unlabeled nighttime or partially annotated video corpora, even if not strictly egocentric;
- Integrating low-level vision techniques for illumination enhancement or robust feature extraction;



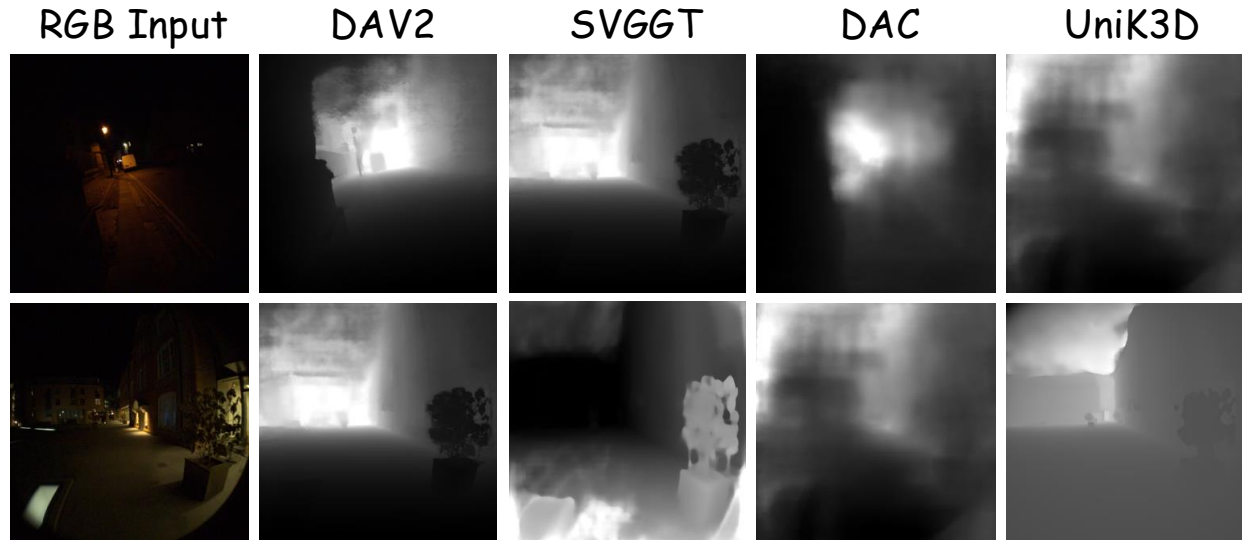


Figure 18: Qualitative results of monodepth estimation in day and night on EgoNight-Oxford dataset, note that DAV2 and SVGGT are shortened for Depth Anything V2 and StreamVGGT respectively.

- Exploring multimodal signals, particularly depth from EgoNight-Synthetic, to improve low-light understanding;
- Developing training-data-free or lightweight adaptation approaches, which are more generalizable across MLLMs.

(4) EgoNight primarily focuses on day–night illumination shifts, while other real-world challenges such as weather variations (rain, fog) and extreme camera motion are not covered. We view these as promising directions for future extensions of EgoNight.

### A.8.3 CONTRIBUTION TO THE COMMUNITY

We believe EgoNight will serve as a valuable resource for the research community in several ways. First, it provides the *first benchmark suite* dedicated to egocentric nighttime vision, a long-overlooked but practically critical setting for robust AI assistants. Second, the dataset’s unique day–night alignment enables rigorous analysis of illumination effects, offering insights that cannot be obtained from prior egocentric benchmarks. Third, by covering multiple tasks, VQA, day–night correspondence retrieval, and depth estimation, EgoNight provides a comprehensive testbed that can catalyze progress across both perception and reasoning. Finally, with all data, annotations, and evaluation code to be released publicly, EgoNight is designed to be easily accessible, extensible, and reproducible, supporting future research on egocentric vision understanding learning.

### A.8.4 USAGE OF LARGE LANGUAGE MODELS (LLMs)

Our annotation pipeline and benchmark evaluation both leverage large language models (LLMs). For data construction, advanced multimodal LLMs are used to generate initial captions, questions, and pseudo answers, which are then refined by human annotators. This hybrid model–human approach substantially reduces annotation cost while ensuring quality. For evaluation, we adopt the *LLM-as-a-Judge* paradigm to assess the semantic correctness of model outputs against ground-truth answers, following recent practice in egocentric VQA. Beyond annotation and evaluation, we also used LLMs to support paper preparation, such as generating icons for illustration figures and assisting with proof-reading. Importantly, while LLMs serve as practical tools throughout our workflow, all core ideas, dataset design, experiments, and analyses are conceived and conducted independently by the authors.

### A.8.5 ETHIC STATEMENT

All indoor egocentric recordings in EgoNight-Sofia were collected with explicit informed consent, and outdoor data is fully anonymized by blurring faces, license plates, house numbers, and other identifiable details, with audio removed,

in compliance with GDPR and privacy standards. Before release, all videos will be verified to contain no personally identifiable information. For EgoNight-Oxford, the subset is derived from the publicly available Oxford Day-and-Night dataset under the BSD-3-Clause license, which permits redistribution and modification with proper attribution. We will retain all required license notices and appropriately acknowledge the original dataset.