

WalkVLM : Aid Visually Impaired People Walking by Vision Language Model

Anonymous CVPR submission

Paper ID 12248

Abstract

Approximately 200 million individuals around the world suffer from varying degrees of visual impairment, making it crucial to leverage AI technology to offer walking assistance for these people. With the recent progress of vision-language models (VLMs), employing VLMs to improve this field has emerged as a popular research topic. However, most existing methods are studied on self-built question-answering datasets, lacking a unified training and testing benchmark for walk guidance. Moreover, in blind walking task, it is necessary to perform real-time streaming video parsing and generate concise yet informative reminders, which poses a great challenge for VLMs that suffer from redundant responses and low inference efficiency. In this paper, we firstly release a diverse, extensive, and unbiased walking awareness dataset, containing 12k video-manual annotation pairs from Europe and Asia to provide a fair training and testing benchmark for blind walking task. Furthermore, a WalkVLM model is proposed, which employs chain of thought for hierarchical planning to generate concise but informative reminders and utilizes temporal-aware adaptive prediction to reduce the temporal redundancy of reminders. Finally, we have established a solid benchmark for blind walking task and verified the advantages of WalkVLM in stream video processing for this task compared to other VLMs. Our dataset and code will be released at anonymous link <https://walkvml2024.github.io>.

1. Introduction

Approximately 200 million people worldwide suffer from varying degrees of visual impairment, with 36 million completely blind [1, 2]. These visually impaired people (VIPs) are facing severe challenges in daily activities such as walking, which may be alleviated by contemporary artificial intelligence technologies [3, 4].

The current walking assistance works primarily concentrate on electronic assistive devices, sensory substitution devices, and computer vision-based assistive systems [5–7]. Among them, vision-based assistive systems can be roughly divided into detection-based methods and

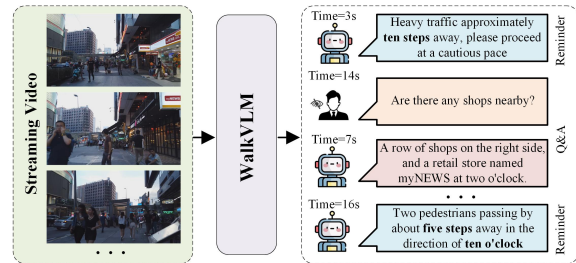


Figure 1. WalkVLM provides opportune, concise, informative walking reminders and answers for visually impaired people based on hierarchical planning and temporal-aware adaptive prediction.

semantic-based methods [8–10]. Detection-based methods have been studied for a long time, aiming to detect potential obstacles in the field of view, so as to let VIPs avoid them [11, 12]. Semantic-based methods utilize vision-language models (VLMs) to analyze images, thereby generating responses to VIPs’ questions [13, 14]. In recent days, with the development of VLMs [15, 16], semantic-based methods have gained significant attention. Some studies have tested VLMs in a zero-shot manner to analyze the effectiveness of these models in blind walking [11, 14]. Moreover, some studies have fine-tuned VLMs using traditional visual question-answer (QA) datasets in this field or a small quantity of self-built datasets, so that the model can better answer user questions [3, 17]. These studies have empowered blind walking tasks with VLMs and already achieved attractive application results.

Although some VLM-based models for blind walking have been developed, these models still face challenges before they can be applied in practice. Firstly, most current research relies on a small number of self-collected image-text pairs and lacks a consistent and extensive benchmark [17, 18]. Moreover, the images and text in traditional datasets are predominantly in a question-and-answer paradigm, which makes it challenging for VLMs to proactively generate guided responses rather than specific answers to questions [13, 19]. Secondly, in blind walking task, it is necessary to perform real-time streaming video parsing and generate concise yet informative reminders, which poses a great challenge for VLMs that suffer from redun-

dant responses and low inference efficiency[20, 21].

In this paper, we propose a WalkVLM for the blind walking task and establish a new benchmark to promote the development of this field. Specifically, we first introduce a diverse, extensive, and unbiased Walking Awareness Dataset (WAD), which contains 12k video-manual annotation pairs from Europe and Asia to provide a fair training and testing baseline. After that, as shown in Figure 1, we introduce the WalkVLM model to interpret video streaming, which employs a chain of thought to hierarchically direct VLM in generating concise yet informative reminders, and achieves opportune reminders by the proposed temporal-aware adaptive prediction. Comprehensive experiments show that, compared to other VLM models, WalkVLM can generate more concise reminders and has better temporal adaptability when handling video streaming in blind walking task. The main contributions of our work are as follows:

- We construct a diverse, extensive, and unbiased walking awareness dataset, providing extensive data support for blind walking task.
- A WalkVLM model for streaming video parsing has been proposed to adaptively provide concise yet informative walking reminder for visually impaired people.
- To the best of our knowledge, this is the first work to utilize VLM to provide opportune walking guidance for visually impaired individuals, laying a solid foundation for the practical application of VLM in this field.

2. Related Work

Vision Datasets for Blind Walking. Existing datasets for blind walking can be roughly divided into two types: detection-based [8, 22–24] and semantic-based [9, 13]. Detection-based datasets have been extensively studied in the blind walking, where researchers utilize these datasets to train the obstacle detection model, thereby reducing the accident rate of VIPs in this task. For example, Zhang *et al.* [22] recently developed a TP-Dataset for detecting visual tactile paving surfaces and offered guidance for the visually impaired through provide walking routes. Islam *et al.* [23] introduced a dataset for improving real-time object recognition systems to aid VIPs in navigation tasks, which contains 90 object annotations from 31 video clips. Compared with detection-based datasets, semantic-based datasets are relatively rare, which contain question-answering properties and provide an enhanced human-computer interaction experience. Gurari *et al.* [9] constructed a VQA dataset for VIPs, which contains 31k visual questions, each with 10 crowdsourced answers. In addition, some researchers have constructed several self-built question-answer datasets with specific attributes during their studies [3, 13], however, these self-built datasets are not open-sourced and are relatively small in scale, making them unsuitable for large-scale and unified benchmarking.

Vision-based Methods for Blind Walking. Similar to the division of datasets, the vision-based methods that help VIPs walking can also be divided into detection-based methods [11, 12] and semantic-based methods [13]. Detection-based methods typically use detectors to obtain obstacles during walking, thereby providing users with specific object locations. Liu *et al.* [12] proposed an open scene understanding system, which improves detection performance by using SAM [25] to generate pixel-level dense segmentation masks. Tian *et al.* [26] proposed a system for understanding dynamic crosswalk scenes, including crosswalks, vehicles, and pedestrians, thereby providing VIPs with indications of when and where to cross the road. The semantic-based approach provides VIPs with the scene understanding in the form of question-answer. Merchant *et al.* [17] verified that vision-language models can generate correct and useful instructions for VIPs, and studied methods to provide users with context-related guidance. Yang *et al.* [3] explored how to utilize VLMs to provide reliable visual question answers for VIPs, and they fine-tuned the VLMs by LoRA on a small amount of self-built dataset to generate detailed and practical suggestions. Moreover, a few applications such as Be My AI ¹ have also adopted semantic-based methods to enable VIPs to take photos for answering questions. However, these applications also only support the question-and-answer paradigm and struggle to provide concise and opportune reminders during walking.

Vision-language Models. With the popularity of large language models (LLM), vision-language models have also begun to receive significant attention [18, 27, 28]. Liu *et al.* [29] proposed the LLaVa, which employ the ViT visual encoder to encode images, follow by mapping them through an MLP to the LLM, yields favorable outcomes in benchmark tests when answering pertinent questions. Subsequently, a plethora of studies emerged based on LLaVa, which greatly impacted various fields [30–33]. Furthermore, multimodal models like Qwen, Gemini, and MiniCPM-V [34–36] have progressively adopted support for multi-frame image inputs and have undergone optimizations for scenarios such as edge devices, significantly enhancing the usability of VLMs in a wide range of applications. Despite the existing studies validating the viability of multimodal large-scale models [37], there remains a dearth of related applications within specific vertical sectors. For instance, only a limited number of studies [3, 13, 17] have focused on the applicability of VLMs in the blind walking task, with a notable absence of unified and systematic modeling approaches.

3. Walking Awareness Dataset

In this section, we have constructed a walking awareness dataset to provide open data support for blind walking task.

¹<https://www.bemyeyes.com>



Figure 2. The data annotation pipeline for constructing the walking awareness dataset. Appendix A.5 provides more random sampling examples to observe the diversity and complexity of WAD dataset.

171 3.1. Data Collection

172 The WAD dataset has a wide range of geographical sources,
173 which originate from 10 different locations in Europe and
174 Asia. 20% of the original data in the WAD dataset comes
175 from the annotators’ recordings, and the rest comes from
176 YouTube². During the recording, six recorders positioned
177 the camera at a height corresponding to chest level, employ-
178 ing focal lengths of 13mm, 20mm, and 26mm, as well as
179 resolutions ranging from 1080p to 4k at 60fps, to enhance
180 the variability of the data. Lastly, we have amassed approx-
181 imately 13 hours of walking video, and see Appendix A for
182 the duration of data gathered from various regions.

183 3.2. Annotation Strategy

184 Figure 2 shows the overall annotation pipeline of walking
185 awareness dataset. Next, we will elaborate from two aspects:
186 scene annotation and response annotation.

187 **Scene annotation.** Scene annotation aims to label the in-
188 herent attributes of the current scene. We requested nine
189 annotators to label the video scene in terms of weather con-
190 ditions, location type, traffic flow rating, danger level, and
191 scene description. When outdoors, weather conditions are
192 divided into six categories such as sunny and rainy, while
193 the status is empty when indoors. The location type is di-
194 vided into eight categories, such as corridors and pedestrian
195 walkway. The traffic flow rating is divided into three levels,
196 which are defined based on the person number in the video
197 stream. The danger level is defined as the walking hazard
198 in the current scene, which is qualitatively divided by the
199 traffic flow rating and road smoothness. The scene descrip-
200 tion is an overview of the current environment, including an
201 expansion on factors such as pedestrian flow, vehicle traffic,
202 road conditions, and the surrounding environment. Subse-
203 quently, we employed the open-world detection model [39]
204 for the preliminary detection of targets, and carried out a
205 corresponding human review to uphold the result accuracy.

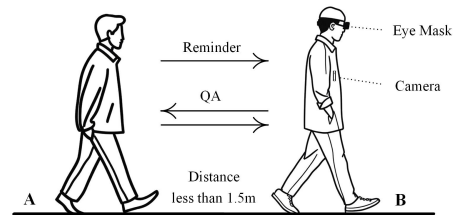


Figure 3. Blind test experiment for analyzing the most critical in-
formation needed by users in blind walking. We required two in-
dividuals to collaborate as a team, where the participant at the rear
provided directions to enable the individual at the front to arrive at
a specific location safely in the absence of any visual information.

206 **Response annotation.** Response denotes the concise re-
207 minders that the model is required to generate, as well as the
208 answer that reply to user’s question in blind walking task.
209 In order to analyze the most critical information needed by
210 users in blind walking, we conducted a blind test experi-
211 ment as shown in Figure 3. In the experiment, we requested
212 two people to collaborate in pairs, with the person A be-
213 hind giving directions, so that the person B in front with
214 eye mask could reach a certain destination without any col-
215 lisions. In such a scenario, the instructions received by per-
216 son B during walking come entirely from person A, and
217 the route priors possessed by real blind people are avoided,
218 which can help us analyze what types of information are
219 necessary for the blind walking task. In a large number of
220 such experiments, we have verified that such guidance can
221 guide visually impaired people to walk safely, indicating
222 that the information provided by person A is sufficiently ef-
223 fective for person B. We recorded the video and audio that
224 occurred during this process, analyzed the information in-
225 teraction between the subjects, and thus provided the fol-
226 lowing valid information types that need to be marked for
227 subsequent *reminder* and *QA* annotations:

- 228 • *Reminder type.* Based on the blind test experiment, as
229 shown in Figure 4, we divided the reminders during walk-
230 ing into six types. (a) Obstacle reminder: Trigger a re-

²<https://www.youtube.com/@poptravelorg>

Dataset	Type	#Sample	Modality	Bounding Box	Weather	Danger level	Scene Summary	QA	Reminder	Open
Obstacle Dataset (2023)[24]	\mathcal{T}	8k	Image	✓	✗	✗	✗	✗	✗	✓
WOTR (2023)[8]	\mathcal{T}	13k	Image	✓	✗	✗	✗	✗	✗	✓
ISLAM <i>et al.</i> (2024)[38]	\mathcal{T}	31	Image & Video	✓	✗	✗	✗	✗	✗	✓
Wang <i>et al.</i> (2024)[11]	\mathcal{T}	50	Video	✓	✓	✗	✗	✗	✗	✗
VizWiz (2018)[9]	\mathcal{S}	31k	Image	✗	✗	✗	✗	✓	✗	✓
Zain <i>et al.</i> (2024)[17]	\mathcal{S}	48	Image	✗	✗	✗	✗	✗	✓	✗
WAD (Ours)	\mathcal{TS}	12k / 120k	Video / Image	✓	✓	✓	✓	✓	✓	✓

Table 1. Static information comparison of different datasets in blind walking task. For dataset types, \mathcal{T} and \mathcal{S} denote the target-based and semantic-based dataset, respectively. WAD dataset holds a significant advantage in terms of sample numbers, categories, and modalities.

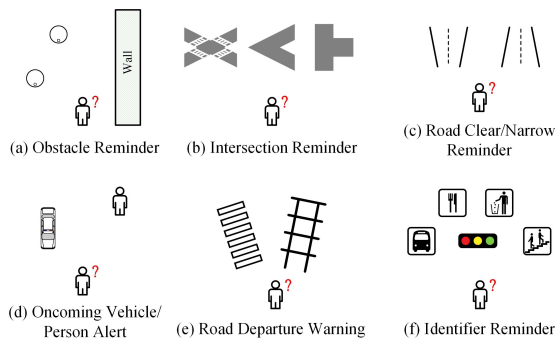


Figure 4. Visualization of six scenarios that require reminders, which were summarized through multiple blind experiments.

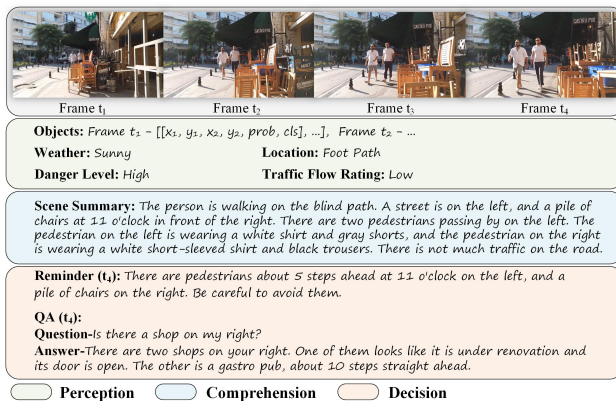


Figure 5. Visualization of the walking awareness dataset. Each sample contains a video clip and multiple labels, with the label hierarchy divided into perception, comprehension, and decision.

conducted manual verification. 260

3.3. Dataset Analysis 261

Figure 5 shows a sample of the WAD dataset, and we divide the annotations into three parts following lower to higher levels: perception, comprehension, and decision. The perception label reflects the basic attributes of the video, such as obstacle location, weather conditions, *etc.*, while the comprehension label reflects the model’s understanding of the entire scene. The decision label contains reminder and QA, reflecting the model’s decision on the user’s walking based on its understanding of the current scenario. 262 263 264 265 266 267 268 269 270

Table 1 illustrates the comparison between the WAD dataset and other prevalent datasets utilized in blind walking tasks, with \mathcal{T} representing the detection-based dataset and \mathcal{S} indicating the semantic-based dataset. Compared to other different types of datasets, WAD has a larger data size while containing more static attributes of the environment, scene summaries, QA, and reminder, thus providing more supervision to train the model. It is worth emphasizing that the samples we furnish are exclusively video clips, which possess a greater volume of information in comparison to the images supplied by other datasets. Moreover, for each video clip, we have extracted 10 keyframes to streamline researchers’ use. The walking awareness dataset contains 3.47 million instances, with categories and the respective proportions shown in Figure 6(a). The category-related dis- 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285

- 231 reminder when there is a non-moving obstacle on the walk-
- 232 ing route. (b) Intersection reminder: Trigger a reminder
- 233 when the current road has intersections, turns, *etc.* (c)
- 234 Road clear/narrow reminder: Provide reminders about
- 235 the width and pass ability of the road. (d) Oncoming
- 236 vehicle/person reminder: When there are moving obsta-
- 237 cles on the walking route, trigger a reminder for potential
- 238 dangers. (e) Road departure warning: Issue a warning
- 239 when there is an angular offset between the walking route
- 240 and the current road. (f) Identifier reminder: Provide re-
- 241 minders for prominent landmarks in the scene, such as
- 242 road signs and traffic lights.
- 243 • *QA type.* For QA type, we proceed from three aspects:
- 244 scene perception, road inquiry, and detailed consultation.
- 245 (a) Scene perception: Macro-level insights such as the
- 246 understanding of the scene. (b) Road inquiry: Route plan-
- 247 ning to reach a certain location within visible range. (c)
- 248 Detailed consultation: Knowledge QA on local details,
- 249 such as road sign content, shop names, *etc.*

250 When marking reminder and QA, we require annotators
251 to indicate the specific location of obstacles in the video. In
252 the annotations, the distances are represented by steps on a
253 scale of 5, the directions are indicated by clock positions,
254 so as to reduce the offset caused by the camera perspective.
255 We require nine annotators to annotate the above content,
256 and the relevant annotation interface is shown in Appendix
257 A.3. After the annotation is completed, in order to further
258 standardize the annotation content to remove potential bias,
259 we used GPT [40] to rephrase the annotated content and

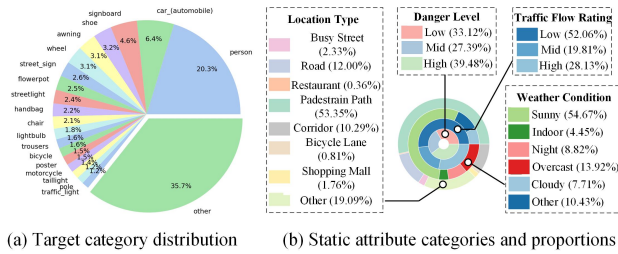


Figure 6. Visualization of the proportion of targets and categories in our walking awareness dataset.

tribution in the WAD dataset is shown in Figure 6(b). We have selected 1.5k samples as a test set based on different static tag types, different reminder types, and different QA types to ensure the diversity and completeness in evaluation.

3.4. Possible Sources of Bias

Although the WAD dataset is collected from a wide range of geographical sources, we are aware of a few biases in our dataset. The regions are still limited, which is still a long way from complete coverage of the globe. The position of the camera and the divergence of focal length are also concerns for us, which need to obtain more general data to compensate for this. In addition, the linguistic preferences of the annotators can introduce specific biases into the generated reminder, which implies that during the walking process, the model might provide information that are more appropriate for the area where the annotation was made.

4. WalkVLM

This section proposes WalkVLM, attempting to empower the blind walking task using a vision-language model based on the WAD dataset. The overall architecture of WalkVLM is shown in Figure 7. We will start with problem formulation and proceed with hierarchical planning and temporal-aware adaptive prediction to generate concise and opportune walking reminders.

4.1. Problem Formulation

We aim to steer a VLM to process video streams, enabling it to provide walking reminders that include temporal attributes, and to enable the model to answer specific questions in human-machine interactions. Specifically, at time t_0 , given the newly appeared frames $[I_{t-N}, \dots, I_{t-1}, I_{t_0}]$, category and obstacle position in the image $[O_{t-N}, \dots, O_{t-1}, O_{t_0}]$, VLM is hoped to generate a concise and informative reminder $T_{t_0}^R$ based on visual information. During walking, VIPs can also raise a question Q_{t_0} to communicate with the VLM at any time, so as to inquire about information such as the current scene and route. Additionally, since generating reminders at every frame may lead to a poor walking guidance experience and impose significant real-time processing pressure on hardware,

WalkVLM needs to be able to predict the current VLM trigger state s_{t_0} based on historical states $[s_{t-N}, \dots, s_{t-1}, s_{t_0}]$ and the previous N frames, so as to choose specific moments to output reminders.

4.2. CoT-Based Hierarchical Planning

We attempt to make VLM conduct step-by-step derivation by a Chain of Thought (CoT) [41], enabling it to summarize from comprehensive information such as the static attributes and the summary of the scene, thereby refining out concise and informative reminders. The model architecture integrates a vision transformer encoder and a large language model (LLM). The vision encoder generates image tokens, while an attention-based extractor aligns these tokens with the LLM, enabling comprehensive understanding and information processing. WalkVLM combines multi-frame information to make reminders, ensuring that the model has a comprehensive perception of the environment.

We divide the process of reminder generation into three levels: perception, comprehension, and decision. At the **perception level**, the model extracts static visual attributes from the current frame, such as location type, weather conditions, and traffic flow rating. To enhance the VLM model’s focus on significant elements and improve visual perception accuracy, we incorporate a priori-object location module (POLM). The POLM initially uses a generic object detector [39] to identify and locate objects in the scene, then filters them based on size and confidence scores to highlight crucial items that reflect road conditions and potential danger. The filtered information and basic environmental attributes provide the necessary input for the model to perceive the external world. At the **comprehension level**, the model integrates all outputs from the perception layer, merging local detection results and fragmented scene information into a comprehensive global summary. Relying on the capabilities of the VLM and the detailed attributes from the perception stage, this stage ensures that the model has a clear understanding of the current environment. At the **decision level**, we focus on training the WalkVLM model to achieve visual QA and reminder. At this stage, the model already possesses an understanding of the static attributes and overall situation of the environment. Therefore, with appropriate guidance, the model is expected to briefly analyze potential hazards in the scene.

During training, we adopted a CoT approach to gradually feed information from three levels into the VLM, and during testing, we let the model predict the aforementioned attributes and generate the corresponding responses.

4.3. Temporal-Aware Adaptive Prediction

Although VLMs are capable of scene parsing across multiple frames and generating the required output, directly applying them to video streaming will lead to unavoidable issues. For instance, when utilizing VLM to generate walking

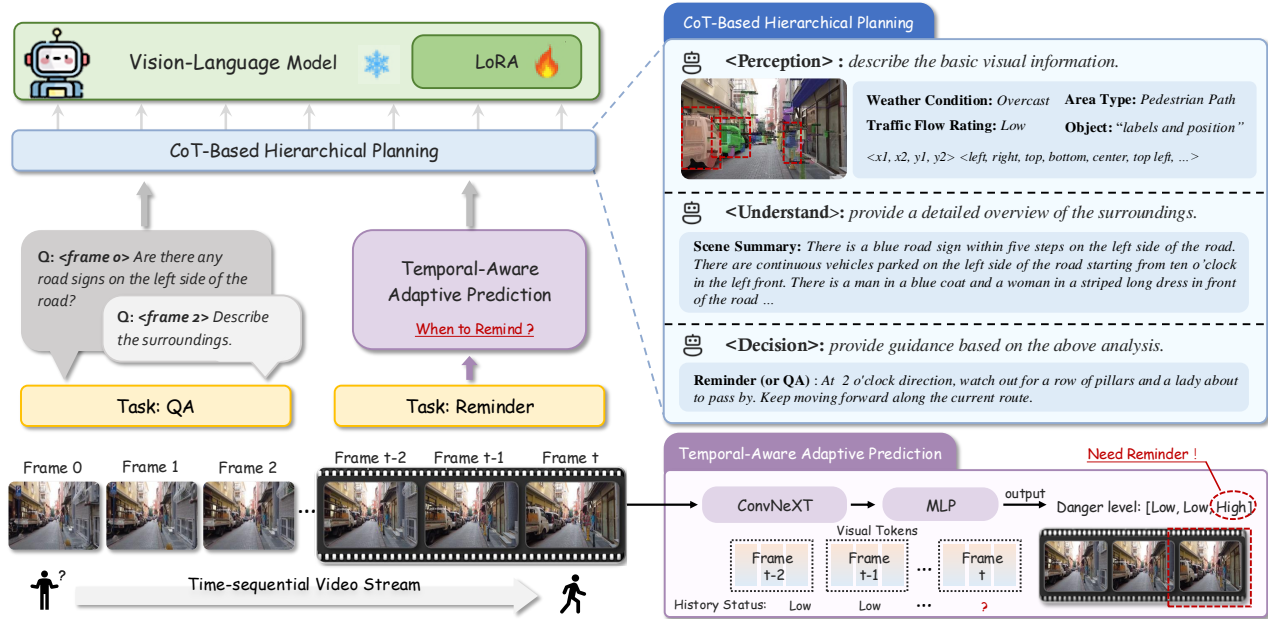


Figure 7. An overview of the proposed WalkVLM framework. WalkVLM employs CoT-based hierarchical planning to summarize the static attributes and understanding of scenes, thereby facilitating the subsequent reminder and QA tasks. Furthermore, temporal-aware adaptive prediction has been proposed to calculate the trigger state of VLM, thereby reducing the temporal redundancy of outputs.

reminders frame by frame or at regular intervals, it will produce a substantial amount of temporal redundancy for the user, resulting in a suboptimal user experience. Secondly, continuous VLM inference also brings computational pressure to hardware devices. Identifying and implementing solutions to this challenge is a key component in the effective utilization of VLM for video streaming processing.

To address the aforementioned issues, we come up with a temporal-aware adaptive prediction (TAP) module that incorporates historical information to pre-calculate whether to trigger the VLM currently, thereby reducing the inference pressure on hardware. Specifically, as shown in the right of Figure 7, we utilize a lightweight model to analyze historical N frames and determine whether to trigger the VLM at the current moment based on the historical output states. Given the frames $[I_{t-N}, \dots, I_{t-1}, I_{t_0}]$, we utilize a 3D convolutional model to extract the features f_v from the sequence. Simultaneously, the predicted trigger states from the previous N moments are independently embedded, concatenated, and then passed through multiple layers of perceptrons to generate the state feature f_s . Furthermore, f_v and f_s are integrated by a multi-layer MLP to generate the current trigger probability \mathcal{P}_t . Three levels of triggers are defined, which correspond to the degrees of danger in the WAD dataset.

The TAP model is used to trigger the reminder of VLM, and subsequent experiments have verified that this module can effectively reduce the temporal redundancy when generating walking guidance.

5. Experiments

5.1. Settings

Models & Details. WalkVLM is implemented with the MiniCPM-V2.6 model [36], which is an 8B multimodal model built upon Qwen2-7B [46]. We add LoRA to all the linear layers of MiniCPM-V2.6 with a rank of 64, while maintaining the video stream sampling rate of 2 FPS. The number of historical frames N is set to 3, and the visual extraction backbone in the TAP module is ConvNext3D [47]. We compared WalkVLM with multiple popular multimodal models, including GPT-4o [44], Qwen2-VL(7B) [45], MiniCPM-V2.6(8B) [36], DeepSeek(1.3B&7B) [42], Yi-VL(6B) [43]. All the prompts of the large models used in this paper can be found in Appendix B.

Metrics. We use the following metrics to evaluate the models: (a) **ROUGE**. This metric measures the similarity between the generated text and the reference text by comparing overlapping words or phrases, including ROUGE-1, ROUGE-2, and ROUGE-L [48]. (b) **TF-IDF Similarity (TF-IDF)**. Combine term frequency and inverse document frequency to evaluate the weight of words, represent the text as a TF-IDF vector, and then measure the semantic similarity between texts [49]. (c) **GPT Score**. GPT4 is used to evaluate the superiority ratio between the generation results of different multimodal models and the ground truth (GT) [50, 51]. (d) **Temporal Redundancy F1-Score (TRF)**. Given the historical model state and historical frames, let the model predict the danger level of the current moment,

Model	Reminder Task					QA Task				
	TF-IDF	ROUGE-1	ROUGE-2	ROUGE-L	GPT Score	TF-IDF	ROUGE-1	ROUGE-2	ROUGE-L	GPT Score
DeepSeek (1.3B) [42]	0.073	0.098	<u>0.015</u>	0.090	0.060	0.182	0.103	0.020	0.095	0.042
DeepSeek (7B) [42]	<u>0.132</u>	0.073	0.009	0.068	0.006	0.189	0.088	0.021	0.081	<u>0.125</u>
Yi-VL (6B) [43]	0.112	0.093	0.009	0.085	0.054	0.113	0.091	0.012	0.082	0.021
MiniCPM-V2.6 (8B) [36]	0.111	0.071	0.007	0.064	0.010	0.192	0.139	0.025	0.120	0.104
GPT-4o [44]	0.116	0.078	0.008	0.072	<u>0.405</u>	0.242	0.163	0.034	0.145	<u>0.125</u>
Qwen2-VL (7B) [45]	0.106	<u>0.107</u>	0.010	<u>0.097</u>	0.018	<u>0.232</u>	<u>0.182</u>	<u>0.037</u>	<u>0.162</u>	0.063
WalkVLM	0.166	0.191	0.062	0.173	0.447	0.189	0.202	0.051	0.174	0.521

Table 2. Quantitative comparison of different methods on reminder and QA tasks. WalkVLM leads in almost all the TF-IDF, ROUGE, and GPT Score metrics. The higher the metric, the better the result. **Bold** and underline indicate the best and the second-best, respectively.

434 and calculate the F1-Score between the prediction and the
 435 GT. TF-IDF and ROUGE evaluate similarity from seman-
 436 tic similarity and word granularity, respectively, while the
 437 GPT Score determines the optimal result by comparing re-
 438 sults with GT. TRF measures the temporal redundancy of
 439 the model’s output; the higher it is, the less temporal redun-
 440 dancy is generated.

441 5.2. Quantitative Results

442 Table 2 presents the quantitative metrics of different mod-
 443 els on the reminder and QA task. On the ROUGE metric,
 444 WalkVLM has achieved the best results in both tasks, veri-
 445 fying that the model’s output is closest to the GT at the word
 446 granularity. On the TF-IDF metric for measuring semantic
 447 similarity, WalkVLM performs the best in reminder tasks,
 448 indicating that the model can generate more concise and ac-
 449 curate results like GT. While in QA tasks, WalkVLM’s per-
 450 formance on TF-IDF scores does not stand out significantly.
 451 This could be attributed to the fact that during training, the
 452 model is encouraged to generate concise answers, which
 453 may inadvertently diminish its capacity to offer elaborate
 454 explanations of the questions. The GPT score represents
 455 the overall evaluation of the LLM on the generated results
 456 and the GT. WalkVLM outperforms other models such as
 457 GPT-4o in terms of GPT scores for reminder and QA tasks,
 458 validating that the model’s output has the most consistent
 459 distribution with the GT.

Model	Yi-VL	MiniCPM-V2.6	GPT-4o	Qwen2-VL	WalkVLM
TRF	0.341	0.396	0.430	<u>0.449</u>	0.505

Table 3. Temporal redundancy assessment of the reminder task, our method achieved the highest TRF score.

460 We use TRF to evaluate the temporal redundancy of the
 461 output from various VLMs. Specifically, we utilize multi-
 462 ple frames of images along with historical dangerous states
 463 as inputs, letting the model to generate a dangerous level
 464 discrimination identifier, thereby determining whether a re-
 465 minder should be triggered currently. As shown in Table
 466 3, compared to other models, WalkVLM has achieved the
 467 highest TRF indicator, which indicates that this model can
 468 better reduce the redundancy of reminders in temporal.

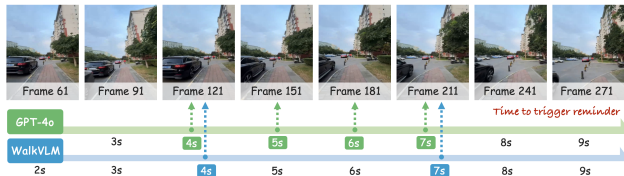


Figure 8. Visualization of triggering moments of GPT-4o and WalkVLM. WalkVLM triggers with less redundancy, providing information to users in a more timely manner.

Model	Reminder Task		QA Task	
	Concise.	Semantic.	Concise.	Semantic.
DeepSeek(1.3B)	0.026	0.080	0.091	0.114
DeepSeek(7B)	0.002	0.197	0.061	0.114
Yi-VL	0.085	0.023	<u>0.121</u>	0.022
MiniCPM-V2.6	0.026	0.122	0.061	0.205
GPT-4o	0.056	<u>0.195</u>	0.030	0.205
Qwen2-VL	<u>0.121</u>	0.168	0.061	<u>0.170</u>
WalkVLM	0.683	0.216	0.576	<u>0.170</u>

Table 4. User study results on conciseness and semantic similarity across different tasks. Higher score indicates better performance.

5.3. Qualitative Results

469 Figure 9 presents the visual comparison in reminder task
 470 between different VLM models. Compared to other meth-
 471 ods such as GPT-4o, WalkVLM can generate more concise
 472 and informative responses, thus providing a better experi-
 473 ence for users. In the left case, whereas other models of-
 474 fer highly detailed responses, WalkVLM simply provides a
 475 concise prompt to the user, effectively highlighting the cru-
 476 cial aspect. As in the right case, WalkVLM perceives the car
 477 coming from the one o’clock direction and conveys the fo-
 478 cus to the user, which other models have not accomplished.
 479

480 Figure 8 shows a qualitative comparison of GPT-4o and
 481 WalkVLM in terms of temporal redundancy. Our model
 482 triggers VLM with lower temporal redundancy and can pro-
 483 vide information to users in a more timely manner. Ap-
 484 pendix C presents more qualitative results, including the
 485 comparison with other VLMs on actual video streams.

5.4. Subjective Results

486 As illustrated in Table 4, we requested nine annotators to
 487 perform a subjective evaluation of various VLM models
 488 with respect to language conciseness and semantic similar-
 489 ity to the GT. Participants are required to rank the results
 490 individually, and we use the top-1 superiority ratio to eval-
 491

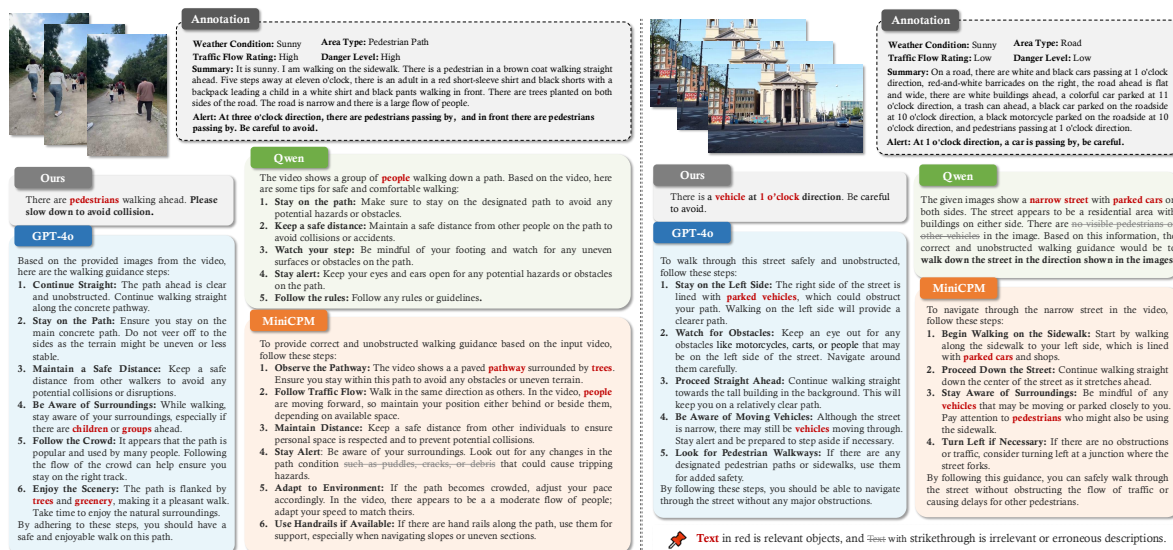


Figure 9. Visualization comparison of different VLM models. Compared to other models, WalkVLM is able to generate concise and informative answers, providing users with a good experience in blind walking. Refer to Appendix C for more visualization comparisons.

492 uate the performance of the method, where a higher score
 493 indicates better performance. Compared with other meth-
 494 ods, WalkVLM has far surpassed other models in terms of
 495 conciseness, both in reminder and QA tasks. In the se-
 496 mantic similarity evaluation against the GT, WalkVLM per-
 497 forms marginally better than GPT-4o in the reminder task
 498 but slightly worse in the QA task. The suboptimal perfor-
 499 mance of WalkVLM in the semantic evaluation of QA tasks,
 500 can be attributed to the conciseness of its output, which
 501 means that a small amount of output information is difficult
 502 to cover all the semantics of the GT.

5.5. Ablative Study

503 The ablation study of WalkVLM is shown in Table 5 to
 504 verify the effectiveness of CoT-based hierarchical planning
 505 (CHP) and POLM prior. We conducted three sets of abla-
 506 tion experiments: (a) **w/o CHP**. Remove the CHP mecha-
 507 nism and generate reminder directly based on the input vi-
 508 sual information. (b) **w/o Pos Prior**. Remove the approx-
 509 imate position of significant obstacles in POLM. (c) **w/o**
 510 **POLM Prior**. Remove the input filtered target exact lo-
 511 cation and category. In these experiments, when the CHP
 512 mechanism was removed, the model’s degradation was sig-
 513 nificant, which may be due to the model’s inability to fully
 514

Configuration	TF-IDF	ROUGE-1	ROUGE-2	ROUGE-L
w/o CHP	0.094	0.073	0.007	0.066
w/o Pos Prior	0.151	0.189	0.062	0.171
w/o POLM Prior	0.152	0.178	0.056	0.164
Full	0.166	0.191	0.062	0.173

Table 5. Ablation study on reminder task. CHP stands for CoT-based hierarchical planning, Pos Prior stands for the general area where obstacles are located in POLM, and POLM Prior stands for the pixel point where the filtered target is exactly located.

515 perceive the scene, resulting in the inconsistency between
 516 the distribution of generated reminder and the GT distribu-
 517 tion. While CHP, enables the model to conduct more de-
 518 tailed analysis from static attributes and scene summaries,
 519 thereby obtaining more concise results. For the case of
 520 lacking POLM prior, the model’s ROUGE performance is
 521 worse compared to lacking position prior, indicating that
 522 the model relies more on the visual details.

6. Conclusion

523 To fulfill the mission of technology for good, this pa-
 524 per presents WalkVLM, a vision-language model for blind
 525 walking task, which employs chain of thought for hierarchi-
 526 cal planning to generate concise but focused reminders, and
 527 utilizes temporal-aware adaptive prediction to reduce the re-
 528 dundancy of reminders in the time series. Additionally, we
 529 have constructed a diverse, extensive, and unbiased walking
 530 awareness dataset, aimed at providing a more robust data
 531 foundation for this field. Comprehensive experiments show
 532 that, compared to other VLM models, WalkVLM can gener-
 533 ate more concise reminder and better temporal adaptability
 534 when handling video streaming in blind walking task.

7. Limitations

535 This paper proposes a WAD dataset and systemati-
 536 cally establishes the blind walking task based on the
 537 vision-language model, thereby setting up an extensive
 538 benchmark and offering valuable data support to this
 539 field. Although the WAD dataset covers dozens of cities,
 540 its generalization capability is still relatively limited in
 541 practical applications, making the collection of additional
 542 data an essential endeavor. Moreover, we devised the
 543 WalkVLM to make the reminders concise and opportune,
 544 but still leave considerable room in inference efficiency.

References

- 548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
- [1] Erin Brady, Meredith Ringel Morris, Yu Zhong, Samuel White, and Jeffrey P Bigham. Visual challenges in the everyday lives of blind people. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 2117–2126, 2013. 1
- [2] Santiago Real and Alvaro Araujo. Navigation systems for the blind and visually impaired: Past work, challenges, and open problems. *Sensors*, 19(15):3404, 2019. 1
- [3] Bufang Yang, Lixing He, Kaiwei Liu, and Zhenyu Yan. Viasist: Adapting multi-modal large language models for users with visual impairments. *arXiv preprint arXiv:2404.02508*, 2024. 1, 2
- [4] Askat Kuzdeuov, Olzhas Mukayev, Shakhizat Nurgaliyev, Alisher Kunbolsyn, and Huseyin Atakan Varol. Chatgpt for visually impaired and blind. In *2024 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC)*, pages 722–727. IEEE, 2024. 1
- [5] Nicholas A Giudice, Benjamin A Guenther, Toni M Kaplan, Shane M Anderson, Robert J Knuesel, and Joseph F Cioffi. Use of an indoor navigation system by sighted and blind travelers: Performance similarities across visual status and age. *ACM Transactions on Accessible Computing (TACCESS)*, 13(3):1–27, 2020. 1
- [6] Gaurav Jain, Yuanyang Teng, Dong Heon Cho, Yunhao Xing, Maryam Aziz, and Brian A Smith. "i want to figure things out": Supporting exploration in navigation for people with visual impairments. *Proceedings of the ACM on Human-Computer Interaction*, 7(CSCW1):1–28, 2023.
- [7] Marion A Hersh and Michael A Johnson. *Assistive technology for visually impaired and blind people*, volume 1. Springer, 2008. 1
- [8] Haiying Xia, Cong Yao, Yumei Tan, and Shuxiang Song. A dataset for the visually impaired walk on the road. *Displays*, 79:102486, 2023. 1, 2, 4
- [9] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617, 2018. 2, 4
- [10] Zongming Yang, Liang Yang, Liren Kong, Ailin Wei, Jesse Leaman, Johnell Brooks, and Bing Li. Seeway: Vision-language assistive navigation for the visually impaired. In *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 52–58. IEEE, 2022. 1
- [11] Hao Wang, Jiayou Qin, Ashish Bastola, Xiwen Chen, John Suchanek, Zihao Gong, and Abolfazl Razi. Visiongpt: Llm-assisted real-time anomaly detection for safe visual navigation. *arXiv preprint arXiv:2403.12415*, 2024. 1, 2, 4
- [12] Ruiping Liu, Jiaming Zhang, Kunyu Peng, Junwei Zheng, Ke Cao, Yufan Chen, Kailun Yang, and Rainer Stiefelwagen. Open scene understanding: Grounded situation recognition meets segment anything for helping people with visual impairments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1857–1867, 2023. 1, 2
- [13] Yi Zhao, Yilin Zhang, Rong Xiang, Jing Li, and Hillming Li. Vialm: A survey and benchmark of visually impaired assistance with large models. *arXiv preprint arXiv:2402.01735*, 2024. 1, 2 605
606
607
608
- [14] Jingyi Xie, Rui Yu, He Zhang, Sooyeon Lee, Syed Masum Billah, and John M Carroll. Emerging practices for large multimodal model (Imm) assistance for people with visual impairments: Implications for design. *arXiv preprint arXiv:2407.08882*, 2024. 1 609
610
611
612
613
- [15] Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C Li, Adrien Bardes, Suzanne Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, et al. An introduction to vision-language modeling. *arXiv preprint arXiv:2405.17247*, 2024. 1 614
615
616
617
618
- [16] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*, 2022. 1 619
620
621
- [17] Zain Merchant, Abrar Anwar, Emily Wang, Souti Chattopadhyay, and Jesse Thomason. Generating contextually-relevant navigation instructions for blind and low vision people. *arXiv preprint arXiv:2407.08219*, 2024. 1, 2, 4 622
623
624
625
- [18] Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1, 2 626
627
628
629
- [19] Maria De Marsico, Chiara Giacanelli, Clizia Giorgia Mangano, Alessio Palma, and Davide Santoro. Vqask: a multimodal android gpt-based application to help blind users visualize pictures. In *Proceedings of the 2024 International Conference on Advanced Visual Interfaces*, pages 1–5, 2024. 1 630
631
632
633
634
635
- [20] Md Adnan Arefeen, Biplob Debnath, Md Yusuf Sarwar Uddin, and Srimat Chakradhar. Vita: An efficient video-to-text algorithm using vlm for rag-based video analysis system. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2266–2274, 2024. 2 636
637
638
639
640
- [21] Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, et al. Sparsevlm: Visual token sparsification for efficient vision-language model inference. *arXiv preprint arXiv:2410.04417*, 2024. 2 641
642
643
644
645
- [22] Xingli Zhang, Lei Liang, Shenglu Zhao, and Zhihui Wang. Grfb-UNET: A new multi-scale attention network with group receptive field block for tactile paving segmentation. *Expert Systems with Applications*, 238:122109, 2024. 2 646
647
648
649
- [23] Md Touhidul Islam, Imran Kabir, Elena Ariel Pearce, Md Alimoor Reza, and Syed Masum Billah. A dataset for crucial object recognition in blind and low-vision individuals' navigation. *arXiv preprint arXiv:2407.16777*, 2024. 2 650
651
652
653
- [24] Wu Tang, De-er Liu, Xiaoli Zhao, Zenghui Chen, and Chen Zhao. A dataset for the recognition of obstacles on blind sidewalk. *Universal Access in the Information Society*, 22(1):69–82, 2023. 2, 4 654
655
656
657
- [25] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 2 658
659
660
661
662

- 663 [26] Shishun Tian, Minghuo Zheng, Wenbin Zou, Xia Li, and
664 Lu Zhang. Dynamic crosswalk scene understanding for the
665 visually impaired. *IEEE transactions on neural systems and*
666 *rehabilitation engineering*, 29:1478–1486, 2021. 2
- 667 [27] Jialu Xing, Jianping Liu, Jian Wang, Lulu Sun, Xi Chen,
668 Xunxun Gu, and Yingfei Wang. A survey of efficient fine-
669 tuning methods for vision-language models—prompt and
670 adapter. *Computers & Graphics*, 119:103885, 2024. 2
- 671 [28] Yizhang Jin, Jian Li, Yexin Liu, Tianjun Gu, Kai Wu,
672 Zhengkai Jiang, Muyang He, Bo Zhao, Xin Tan, Zhenye
673 Gan, et al. Efficient multimodal large language models: A
674 survey. *arXiv preprint arXiv:2405.10739*, 2024. 2
- 675 [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
676 Visual instruction tuning, 2023. 2
- 677 [30] Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng
678 Mou, and Jian Tang. Llava- ϕ : Efficient multi-modal
679 assistant with small language model. *arXiv preprint*
680 *arXiv:2401.02330*, 2024. 2
- 681 [31] Shaoxiang Chen, Zequn Jie, and Lin Ma. Llava-mole:
682 Sparse mixture of lora experts for mitigating data con-
683 flicts in instruction finetuning mllms. *arXiv preprint*
684 *arXiv:2401.16160*, 2024.
- 685 [32] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li,
686 Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave:
687 Tackling multi-image, video, and 3d in large multimodal
688 models. *arXiv preprint arXiv:2407.07895*, 2024.
- 689 [33] Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang,
690 See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. Math-
691 llava: Bootstrapping mathematical reasoning for multimodal
692 large language models. *arXiv preprint arXiv:2406.17294*,
693 2024. 2
- 694 [34] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan
695 Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren
696 Zhou. Qwen-vl: A versatile vision-language model for un-
697 derstanding, localization, text reading, and beyond. *arXiv*
698 *preprint arXiv:2308.12966*, 2023. 2
- 699 [35] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry
700 Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu
701 Soriccut, Angeliki Lazaridou, Orhan Firat, Julian Schrit-
702 twieser, et al. Gemini 1.5: Unlocking multimodal under-
703 standing across millions of tokens of context. *arXiv preprint*
704 *arXiv:2403.05530*, 2024.
- 705 [36] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui,
706 Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He,
707 et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv*
708 *preprint arXiv:2408.01800*, 2024. 2, 6, 7
- 709 [37] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang,
710 Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and
711 Hang Zhao. Drivevlm: The convergence of autonomous
712 driving and large vision-language models. *arXiv preprint*
713 *arXiv:2402.12289*, 2024. 2
- 714 [38] Md Touhidul Islam, Imran Kabir, Elena Ariel Pearce, Md Al-
715 imoor Reza, and Syed Masum Billah. Identifying crucial ob-
716 jects in blind and low-vision individuals’ navigation. *arXiv*
717 *preprint arXiv:2408.13175*, 2024. 4
- 718 [39] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp
719 Krähenbühl, and Ishan Misra. Detecting twenty-thousand
720 classes using image-level supervision. In *ECCV*, 2022. 3, 5
- [40] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Ab-
hishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil
Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The
llama 3 herd of models. *arXiv preprint arXiv:2407.21783*,
2024. 4
- [41] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al.
Chain-of-thought prompting elicits reasoning in large lan-
guage models. *Advances in neural information processing*
systems, 35:24824–24837, 2022. 5
- [42] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai
Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li,
Hao Yang, et al. Deepseek-vl: towards real-world vision-
language understanding. *arXiv preprint arXiv:2403.05525*,
2024. 6, 7
- [43] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang,
Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen,
Jing Chang, et al. Yi: Open foundation models by 01. ai.
arXiv preprint arXiv:2403.04652, 2024. 6, 7
- [44] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perel-
man, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda,
Alan Hayes, Alec Radford, et al. Gpt-4o system card.
arXiv preprint arXiv:2410.21276, 2024. 6, 7
- [45] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan,
Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin
Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui
Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Jun-
yang Lin. Qwen2-vl: Enhancing vision-language model’s
perception of the world at any resolution. *arXiv preprint*
arXiv:2409.12191, 2024. 6, 7
- [46] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen
Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng
Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
arXiv:2407.10671, 2024. 6
- [47] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feicht-
enhofer, Trevor Darrell, and Saining Xie. A convnet for the
2020s. In *Proceedings of the IEEE/CVF conference on com-
puter vision and pattern recognition*, pages 11976–11986,
2022. 6
- [48] Chin-Yew Lin. Rouge: A package for automatic evaluation
of summaries. In *Text summarization branches out*, pages
74–81, 2004. 6
- [49] Juan Ramos et al. Using tf-idf to determine word relevance
in document queries. In *Proceedings of the first instructional
conference on machine learning*, volume 242, pages 29–48.
Citeseer, 2003. 6
- [50] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan
Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with
mt-bench and chatbot arena. *Advances in Neural Information
Processing Systems*, 36:46595–46623, 2023. 6
- [51] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ah-
mad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida,
Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al.
Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*,
2023. 6