

# WalkVLM : Aid Visually Impaired People Walking by Vision Language Model

## Supplementary Material

### Contents

<b>1. Walking Awareness Dataset</b>	<b>1</b>
1.1. Data Regional Distribution . . . . .	1
1.2. Dataset Category Definition . . . . .	1
1.3. Annotation Process . . . . .	1
1.4. Detection Model . . . . .	2
1.5. Sample Visualization . . . . .	2
1.6. Data Analysis . . . . .	2
1.7. Benchmark Data Splits . . . . .	2
1.8. Possible Sources of Bias . . . . .	3
<b>2. Model &amp; Details</b>	<b>4</b>
<b>3. Architectural of TAP</b>	<b>4</b>
3.1. All Prompts Used in Paper . . . . .	4
3.2. Evaluation of Temporal Redundancy F1-Score	4
<b>4. Experiment</b>	<b>4</b>
4.1. Visualization of Hierarchical Reasoning . . .	4
4.2. Visual Comparison of Different Models . . .	5
4.3. Reasoning Efficiency . . . . .	5
4.4. Comparison of Video Streaming Inference .	7
<b>5. Discussion</b>	<b>9</b>
<b>6. Societal Impact</b>	<b>9</b>
<b>7. Limitations</b>	<b>10</b>
<b>8. Acknowledgements</b>	<b>10</b>

## 1. Walking Awareness Dataset

### 1.1. Data Regional Distribution

Table 1 shows the data distribution and corresponding duration in the WAD dataset. The WAD dataset covers ten cities and contains a wide range of data sources. Figure 1 illustrates the relevant regional distribution. As illustrated, our dataset is spread across Asia and Europe, showing a relatively balanced distribution between different regions. Furthermore, the sampling across different regions is relatively uniform, with a large number of samples at various locations to avoid bias, which has good generalization characteristics.

City	Country	Hours
Amsterdam	Netherlands	1:21h
Bangkok	Thailand	2:55h
Chiang Mai	Thailand	1:07h
Istanbul	Turkey	1:08h
Kuala Lumpur	Malaysia	1:12h
Singapore	Singapore	1:36h
Stockholm	Sweden	1:06h
Venice	Italy	1:50h
Zurich	Switzerland	1:05h
Beijing	China	2:33h

Table 1. The source region and duration of the WAD dataset. Refer to Fig. 1 for visualization results.

### 1.2. Dataset Category Definition

As shown in Table 2, WAD dataset contains multiple pre-defined data categories. For weather conditions, we have selected the most common types, avoiding scenarios such as rainy or snowy days that make visually impaired people (VIPs) difficult to go outside. For location types, we have selected the types where VIPs are likely to appear, avoiding rare locations. For the traffic flow rating, we instructed annotators to count the number of people in each video segment and used this count as the basis for classification. For scene summarization, during annotation, we required annotators to summarize static attributes such as road conditions, pedestrian flow, and vehicle flow, providing a comprehensive description of the current environment. Currently, the granularity of our dataset is still relatively coarse. In the future, we will continue to refine different fine-grained categories and gradually expand the size of the dataset.

### 1.3. Annotation Process

We use the page shown in Figure 2 to request annotators to make marks. For static tags, we have provided relevant options for the annotators. For scene summary, we require annotators to describe aspects such as the scene, road conditions, pedestrian flow, and vehicle flow. For reminder and QA, we require annotators to expand on different situations, as described in Section 3.2 of the main paper. Since descriptive tags carry a temporal dimension, we have adopted the annotation method in Table 3 for labeling. After the text categorization is completed, we perform a quality inspec-

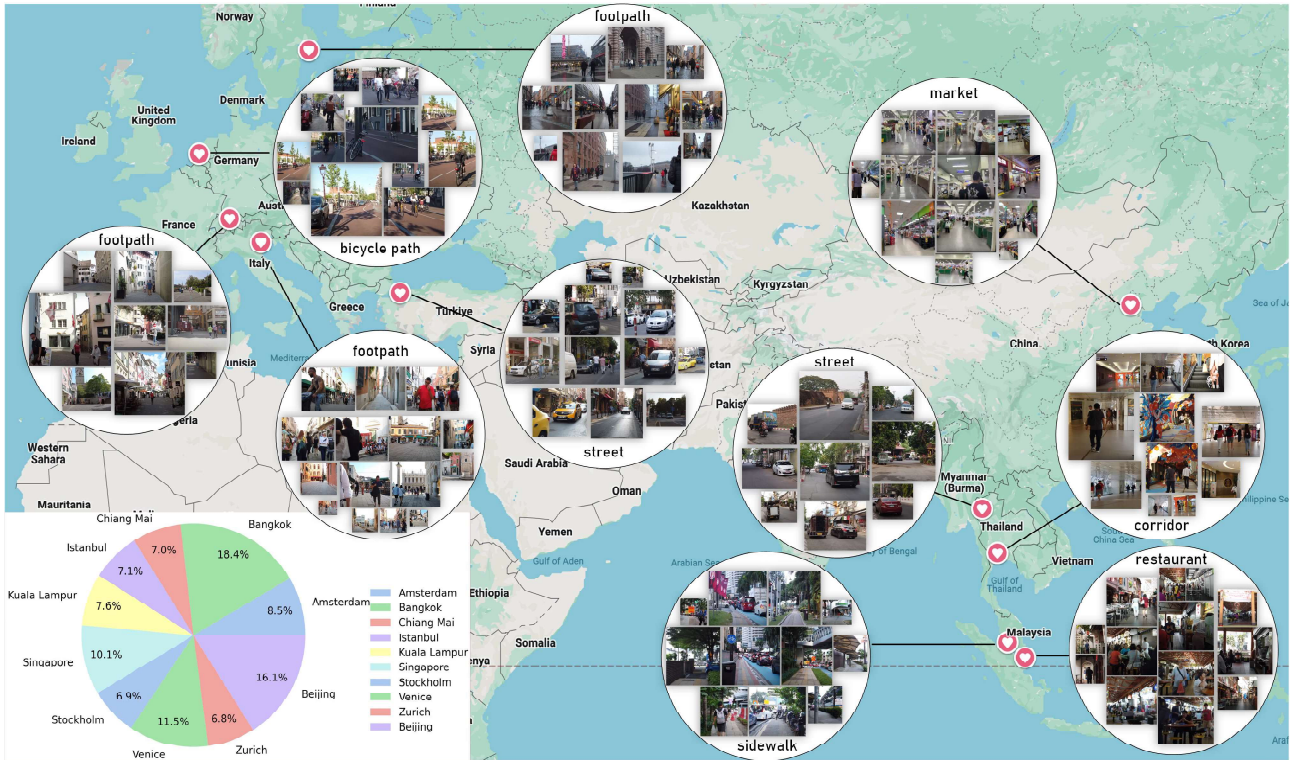


Figure 1. Visualization results of the WAD dataset sorted by region. The WAD dataset has a wide range of sources, and the samples and categories shown are randomly obtained from the dataset. The pie chart in the lower left corner shows the proportion of video length from different regions.

tion on it and use Llama3.1<sup>1</sup> to normalize the samples that pass the inspection to debias.

#### 1.4. Detection Model

The Detic model [5] has achieved excellent results on the LVIS benchmark [2] in open-world detection tasks by training the detector classifier on image classification data. In view of the model’s good generalization ability, we use it to perform preliminary target extraction on the WAD dataset. Figure 3 presents some example of the detection results of the Detic model in the WAD dataset, demonstrating that the model has a strong ability to extract small and complex targets. After using the model for detection, we conducted manual confirmation and deleted some false positive boxes, thus obtaining the final detection results.

#### 1.5. Sample Visualization

Figure 4 and Figure 5 show more sample visualization results in the WAD dataset. Our dataset has wide coverage, diverse types, and possesses ideal reminder attributes to train VLM to have guiding capabilities in blind walking tasks.

<sup>1</sup><https://ai.meta.com/blog/meta-llama-3-1/>

#### 1.6. Data Analysis

Figure 7 shows the distribution of the top 100 categories contained in the WAD dataset, while Table 4 shows all the categories included. Figure 6 presents a word cloud distribution with annotated descriptions, where the most frequently used words include *oclock*, *pedestrain*, *direction*. We have counted the word count distribution in different annotated texts in Figure 8. For reminder and QA scenarios, the data contained in WAD is shorter in length, while for summary scenario descriptions are more detailed.

#### 1.7. Benchmark Data Splits

To ensure the diversity of test data, we adopted a category-based combined clustering method. Through this method, we carefully selected a certain number of samples from the clustering results to form our test set. Ultimately, we selected 1007 reminders and 134 QA pairs as our testset. Furthermore, we conducted a thorough analysis of the distribution of the test set to confirm that they are accurate and that the same type of data is represented in the training set.

Tag Type	Category	Note
Weather Conditions	Sunny	-
	Night	Not make fine-grained distinctions
	Overcast	-
	Cloudy	-
	Indoor	Not make fine-grained distinctions
	Other	Severe weather conditions such as rain and fog for walking
Location Type	Busy Street	Open-air commercial streets
	Road	Roads where vehicles can travel normally
	Restaurant	Food stalls gathered together, inside large canteens
	Pedestrian Path	Walking paths in parks and other places for healthy walking
	Corridor	Indoor walking paths
	Bicycle Lane	Bicycle roads with bicycle signs
	Shopping Mall	Large shopping supermarkets
Other	Niche scenarios	
Traffic Flow Rating	Low	Fewer than 2 people appear in the sliced video
	Mid	Between 2 and 10 people appear in the sliced video
	High	More than 10 people appear in the sliced video
Danger Level	Low	The road is clear, the pedestrian flow is low, and no dangers within 15 steps
	Mid	Other scenarios that do not belong to low or high
	High	Potential collision factors, such as narrow roads, bumpy roads, vehicle warnings
Scene Description	-	Detailed description of the current environment, level of danger, and pedestrian flow
QA	-	The three types of inquiries mentioned in the paper and concise responses
Reminder	-	Brief walking directions to provide to the user based on the current scenario

Table 2. The interpretation of label categories contained in the WAD dataset.

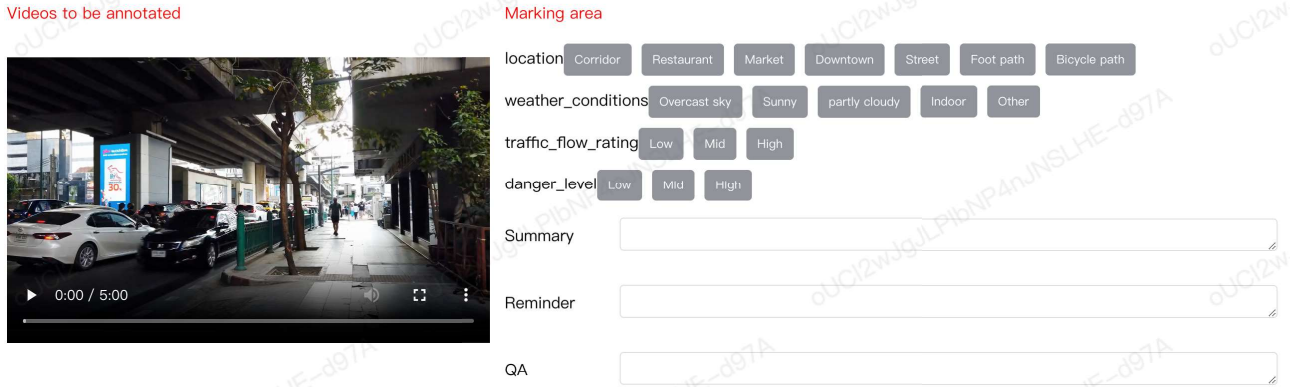


Figure 2. Annotation tool interface. Annotators mark the static attributes of the video in the video, record the time points of reminders and QA, and enter corresponding text descriptions.

### 1.8. 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.





Figure 3. The detection results provided in the WAD dataset, which were pre-detected by the Detic model [5], and then manually reviewed to ensure the correctness of the results. See [here](#) for more detection samples.

---

<b>&lt;time - category1,category2,...&gt;</b>
<b>Annotation</b>
...
<2m30s - A, E>
<i>almost hit the wall, go forward in the 11 o'clock direction to return to the main route.</i>
<2m43s - B>
<i>five steps ahead is the fork in the road, go forward in the 10 - o'clock direction to return to the main route.</i>
...
<3m39s - O>
<i>Q: describe the current scene</i>
<i>A: at a crossroads with many vehicles, keep still to avoid, there are some obstacles ahead, be careful to avoid</i>
...

---

Table 3. Example of reminder and QA result annotation with a temporal dimension. We required annotators to mark the time when events occurred in the video, the question and reminder categories, as well as concise responses.

## 2. Model & Details

### 3. Architectural of TAP

See Fig.9. The historical frames and states are fed into the TAP to predict whether the VLM should be triggered to perform an alter.

#### 3.1. All Prompts Used in Paper

Table 5 displays all the prompts utilized in this paper under various circumstances such as normalizing annotation results, reasoning with VLM, and conducting evaluations. Normalize the annotation results are crucial for ensuring the consistency and uniformity of annotation results, and this

prompt are used in the preprocessing stage to correct bias in the data. For the inference prompt of other models, we input historical multi-frame images and historical states to enable it to generate trigger states and reminders for the user. In the prompt of WalkVLM, we make the model predict different levels of labels step by step and gradually output the results. The evaluation prompt based on GPT4 compares different results with the ground truth to obtain the proportion statistics of the optimal model.

#### 3.2. Evaluation of Temporal Redundancy F1-Score

This section systematically evaluates the redundancy of temporal outputs of different models. Temporal redundancy refers to the excessive frequency of output information in this paper. In order to evaluate the temporal redundancy of different models, we decompose the test video to ensure that each sample contains historical  $N$  frames and  $N$  states, thereby predicting the trigger state under the current situation. We collected 834 such samples as a test set. The predicted labels are divided into three levels, corresponding to the degree of danger. When the degree of danger is high, we regard it as triggering VLM. By comparing the predicted different states with the ground truth, the distribution gap between the two sets of data can be calculated, thereby calculating the F1-score.

## 4. Experiment

### 4.1. Visualization of Hierarchical Reasoning

We have demonstrated the results of hierarchical reasoning using WalkVLM in Figure 10. WalkVLM can effectively extract static attributes from video streams and generate a comprehensive summary of the current scene. After integrating fragmented attributes, the model produces concise and informative walking instructions.





Figure 4. Visual examples of QA samples in WAD dataset. See [here](#) for dynamic samples.

## 4.2. Visual Comparison of Different Models

Figure 11 and 12 presents a comparison of additional visualization results between WalkVLM and other models. Our approach yields more streamlined results, enabling a superior human-machine interaction experience during blind

walking task.

## 4.3. Reasoning Efficiency

The reasoning efficiency of WalkVLM on different devices is shown in Table 6. During the experiment, the API deployed on the A100 was used to provide services. In con-

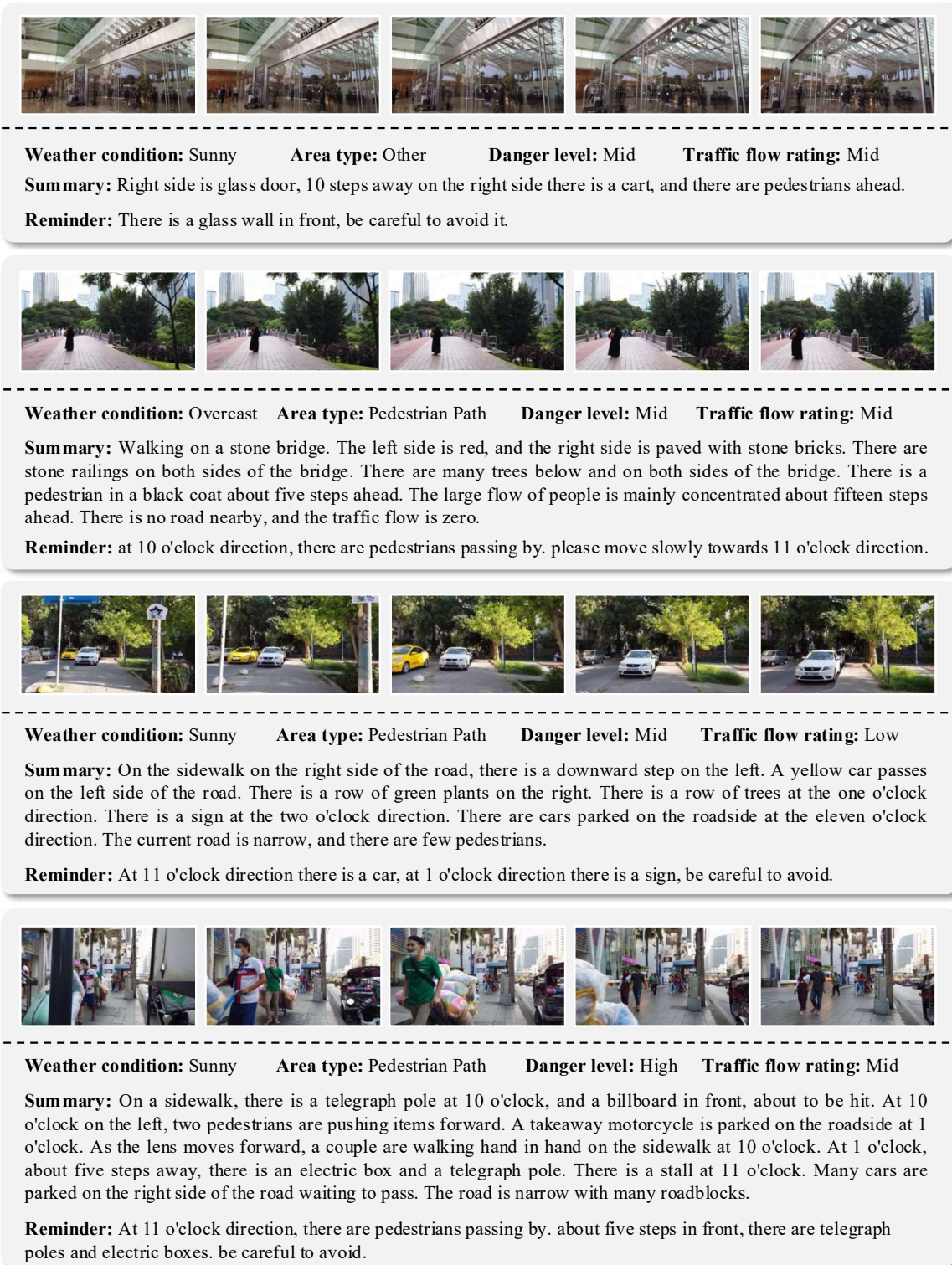


Figure 5. Visual examples of reminder samples in WAD dataset. See [here](#) for dynamic samples.

person, car-(automobile), signboard, shoe, awning, wheel, street-sign, flowerpot, streetlight, handbag, chair, lightbulb, trousers, bicycle, poster, motorcycle, taillight, pole, traffic-light, jean, short-pants, sandal-(type-of-shoe), vent, backpack, flag, license-plate, jersey, headlight, air-conditioner, trash-can, rearview-mirror, umbrella, shirt, dress, strap, jacket, curtain, banner, bench, truck, crossbar, manhole, skirt, cone, telephone-pole, statue-(sculpture), pipe, box, hat, plastic-bag, bus-(vehicle), suitcase, doorknob, boat, dining-table, coat, helmet, bottle, windshield-wiper, watch, bird, lantern, balloon, boot, clock, grill, spotlight, lamppost, baseball-cap, sunglasses, fireplug, beachball, sock, tank-top-(clothing), tag, cellular-telephone, stool, hinge, lamp, shopping-bag, postcard, bolt, billboard, television-set, polo-shirt, cup, reflector, ball, basket, bucket, window-box-(for-plants), antenna, painting, tablecloth, flower-arrangement, bracelet, button, belt, bell, baby-buggy, flagpole, ladder, bowl, spectacles, vase, clock-tower, blouse, book, stop-sign, handle, banana, refrigerator, toy, sunhat, beanie, doughnut, necklace, train-(railroad-vehicle), bottle-cap, fan, tarp, vest, crate, orange-(fruit), magazine, apple, skateboard, parking-meter, postbox-(public), necktie, dog, earring, vending-machine, sweatshirt, barrel, lampshade, chandelier, cowboy-hat, minivan, newsstand, choker, hook, dish-antenna, scarf, camera, pizza, mask, drawer, weathervane, figurine, motor-scooter, magnet, pigeon, speaker-(stereo-equipment), cart, cooler-(for-food), blackboard, roller-skate, hot-air-balloon, flip-flop-(sandal), unicycle, headscarf, cabinet, hatbox, mirror, legging-(clothing), candle, satchel, teddy-bear, lanyard, log, glove, pennant, wall-socket, shower-cap, blinker, canister, pottery, robe, gargoyle, steering-wheel, newspaper, suspenders, dumpster, water-bottle, easel, kite, cushion, apron, horse, wreath, pew-(church-bench), dispenser, tomato, towel, melon, pumpkin, doormat, fire-extinguisher, sombrero, walking-cane, can, telephone-booth, thermostat, wineglass, heart, bandanna, tambourine, cat, jar, peach, carton, ring, frisbee, pot, carrot, watering-can, surfboard, mailbox-(at-home), headband, buoy, coconut, hose, card, sweater, lemon, remote-control, butterfly, grape, plate, knob, gravestone, knocker-(on-a-door), elephant, globe, mast, paper-plate, raincoat, wristlet, projector, watermelon, tote-bag, pirate-flag, mail-slot, tray, bulletproof-vest, brass-plaque, handcart, table, tricycle, towel-rack, laptop-computer, belt-buckle, fire-alarm, bow-(decorative-ribbons), slipper-(footwear), sink, papaya, sawhorse, briefcase, glass-(drink-container), cake, latch, coat-hanger, step-stool, fish-(food), napkin, pastry, motor, shopping-cart, sofa, silo, doll, toilet, tank-(storage-vessel), cookie, crucifix, oven, bamboo, tassel, hairnet, golfcart, fish, bread, cow, monitor-(computer-equipment) computer-monitor, lion, seashell, microwave-oven, earphone, Christmas-tree, water-jug, wagon-wheel, airplane, locker, broom, calendar, pop-(soda), barrette, mammoth, rollerblade, avocado, blazer, scoreboard, hippopotamus, birdbath, shield, rubber-band, paper-towel, music-stool, straw-(for-drinking), poncho, neckerchief, pinwheel, houseboat, crutch, green-bean, birthday-card, sunflower, pickup-truck, grocery-bag, wine-bottle, faucet, halter-top, wine-bucket, sandwich, life-buoy, basketball-backboard, bullhorn, aerosol-can, tapestry, toilet-tissue, bathtub, tripod, goldfish, gourd, fireplace, stepladder, orange-juice, edible-corn, oil-lamp, garden-hose, potato, shower-curtain, water-tower, knife, onion, apricot, tennis-racket, piggy-bank, ashtray, puppet, sculpture, pretzel, fedora, brassiere, milk-can, cantaloup, blimp, blanket, guitar, kiwi-fruit, brake-light, armor, shawl, scissors, table-tennis-table, toothbrush, birdcage, lettuce, cylinder, radiator, turban, kimono, birdhouse, slide, envelope, Dixie-cup, Ferris-wheel, microphone, swimsuit, lime, beer-bottle, shaving-cream, fishbowl, ice-skate, camper-(vehicle), hairpin, pillow, underwear, oar, bonnet, chinaware, cymbal, penguin, sausage, strawberry, costume, dish-towel, gull, sword, bagel, spoon, crown, harmonium, duffel-bag, candle-holder, camcorder, horse-buggy, jumpsuit, clothes-hamper, knee-pad, bathrobe, comic-book, beer-can, giant-panda, map, phonograph-record, bell-pepper, toolbox, solar-array, rhinoceros, booklet, cupcake, shower-head, binoculars, monkey, match-box, hand-towel, deer, pan-(for-cooking), dove, wheelchair, armoire, camel, goose, hair-dryer, dress-hat, tiger, tennis-ball, place-mat, bridal-gown, ottoman, cornice, mug, pear, sail, boxing-glove, passenger-car-(part-of-a-train), cap-(headwear), horse-carriage, urn, wig, wind-chime, thermos-bottle, fume-hood, crock-pot, bubble-gum, cherry, drum-(musical-instrument), wagon, bed, clarinet, eyepatch, tissue-paper, padlock, cigarette, parasol, baseball-bat, teacup, mandarin-orange, aquarium, bun, bowling-ball, telephone, lemonade, dog-collars, windmill, saltshaker, tartan, zucchini, lab-coat, tinsel, radar, pitcher-(vessel-for-liquid), pug-dog, sheep, coffee-maker, folding-chair, pinecone, visor, octopus-(animal), medicine, cassette, yogurt, saddlebag, wardrobe, basketball, persimmon, tape-(sticky-cloth-or-paper), tights-(clothing), baseball-glove, water-hose, cauliflower, cover, garbage-truck, forklift, bath-mat, chopping-board, computer-keyboard, propeller, wristband, gift-wrap, duck, railcar-(part-of-a-train), violin, football-helmet, blueberry, chopstick, piano, starfish, lawn-mower, fork, diaper, frying-pan, shark, wallet, duct-tape, pineapple, elk, toaster, earplug, wall-clock, cab-(taxi), zebra, bow-tie, hog, mallet, boiled-egg, knitting-needle, keycard, condiment, dragonfly, garlic, pepper-mill, drumstick, snowman, thumbtack, gasmask, pouch, teapot, sling-(bandage), barrow, bulldozer, spear, bookmark, mat-(gym-equipment), coffee-table, sleeping-bag, bat-(animal), runner-(carpet), iron-(for-clothing), bath-towel, coatrack, musical-instrument, bulletin-board, pie, tin-foil, overalls-(clothing), bib, pelican, egg, mascot, cistern, bookcase, giraffe, pad, irench-coat, bandage, chalice, flannel, clipboard, dustpan, celery, sweet-potato, headset, bread-bin, howler-hat, walking-stick, saddle-blanket, phonebook, seahorse, clasp, lollipop, desk, broccoli, nailfile, anklet, dress-suit, rag-doll, beanbag, gondola-(boat), bear, mushroom, cider, dishwasher, alcohol, clementine, flap, rifle, ice-cream, ski, snowboard, vacuum-cleaner, automatic-washer, trailer-truck, hamper, television-camera, cigar-box, tobacco-pipe, bouquet, candy-bar, ferry, bead, banjo, ladybug, pacifier, shovel, control, fishing-rod, cruise-ship, wash-basin, whipped-cream, pen, goggles, pan-(metal-container), flipper-(footwear), cucumber, nightshirt, dolphin, water-cooler, cloak, mop, pendulum, canoe, artichoke, heater, hammock, water-gun, almond, paintbrush, shredder-(for-paper), pita-(bread), liquor, eggbeater, scale-(measuring-instrument), dresser, ski-boot, cigarette-case, teakettle, armband, frog, file-cabinet, tow-truck, squid-(food), mouse-(computer-equipment), keg, tongs, deadbolt, quesadilla, hair-curler, koala, asparagus, platter, bobbin, coaster, milk, inhaler, salami, flamingo, life-jacket, coffeepot, urinal, eggplant, business-card, mattress, fig-(fruit), corkboard, raft, cash-register, cabana, suit-(clothing), kitchen-table, corset, gorilla, cocoa-(beverage), yacht, salmon-(fish), spice-rack, parachute, coil, squirrel, ironing-board, projectile-(weapon), coverall, trophy-cup, thread, measuring-stick, dinghy, crowbar, ski-pole, trunk, salad, dartboard, bedpan, award, rabbit, cincture, parka, colander, windsock, home-plate-(baseball), baboon, green-onion, éclair, tooth-paste, saucer, highchair, handkerchief, pajamas, saxophone, potholder, ladle, spatula, first-aid-kit, veil, parakeet, scrubbing-brush, clip, blender, stapler-(stapling-machine), parrot, measuring-cup, owl, ice-maker, sweat-pants, videotape, corkscrew, marker, muffin, tiara, cast, beret, gun, tape-measure, generator, cowbell, sushi, hookah, seabird, row, tachometer, cream-pitcher, battery, Band-Aid, lightning-rod, hamburger, elevator-car, checkbook, hockey-stick, syringe, beeper, gelatin, wrench, water-scooter, hornet, fire-hose, Lego, stove, key, palette, chicken-(animal), deck-chair, chaise-longue, hairbrush, flashlight, smoothie, mitten, flute-glass, crab-(animal), bagpipe, clothespin, soap, lizard, river-boat, boom-microphone, radish, paperweight, fire-engine, candy-cane, bow-(weapon), sponge, wedding-cake, hourglass, ice-pack, tea-bag, cappuccino, eagle, machine-gun, salmon-(food), wet-suit, clutch-bag, cube, brussels-sprouts, wolf, toothpick, kennel, soccer-ball, prawn, hamburger, identity-card, egg-yolk, pegboard, honey, duckling, pencil, ham, saddle-(on-an-animal), gameboard, hot-sauce, amplifier, alarm-clock, tortilla, manatee, brownie, nutcracker, popsicle, funnel, hotplate, trampoline, crib, heron, shampoo, butter, army-tank, date-(fruit), bottle-opener, cornet, camera-lens, jelly-bean, griddle, atomizer, armchair, bass-horn, humming-bird, salsa, baguet, sweatband, arctic-(type-of-shoe), footstool, power-shovel, drone, tractor-(farm-equipment), bunk-bed, food-processor, radio-receiver, cufflink, scarecrow, cock, cougar, chocolate-cake, wok, raspberry, ping-pong-ball, blackberry, dollhouse, space-shuttle, skewer, bobby-pin, school-bus, puffin, car-battery, razorblade, stirrup, drill, truffle-(chocolate), fighter-jet, thermometer, cupboard, screwdriver, sled, eel, pipe-bowl, broach, plume, sofa-bed, ferret, turtle, escargot, crescent-roll, printer, quilt, chocolate-bar, paddle, toaster-oven, motor-vehicle, puffer-(fish), soya-milk, cork-(bottle-plug), cabin-car, walrus, paty-(food), police-cruiser, skullcap, baseball, handsaw, wooden-spoon, pool-table, sewing-machine, pitchfork, cardigan, crayon, manger, kettle, CD-player, barge, flash, rolling-pin, cleansing-agent, dagger, waffle, hardback-book, toast-(food), puppy, egg-roll, chili-(vegetable), kitchen-sink, chocolate-mousse, router-(computer-equipment), pencil-sharpener, pin-(non-jewelry), kayak, sharpener, grater, nut, shoulder-bag, pantyhose, plow-(farm-equipment), mint-candy, crisp-(potato-chip), needle, pea-(food), beef-(food), sherbert, pepper, iPod, bullet-train, polar-bear, headboard, volleyball, bulldog, crape, reamer-(juicer), birdfeeder, table-lamp, pocketknife, jewelry, meatball, pudding, hand-glass, Bible, money, stylus, sugarcane-(plant), cayenne-(spice), shepherd-dog, lip-balm, soup-bowl, cornbread

Table 4. Full list of the target categories present in the walking awareness dataset, sorted by the number of occurrences in the dataset.

junction with the TAP module, the latency was within 1 second. The inference speed on the edge device (iPad Pro M4) is relatively slow, with a delay of 2-3 seconds per call. We will subsequently build a model with lower computational requirements to directly deploy it on edge devices.

#### 4.4. Comparison of Video Streaming Inference

In this section, we deployed Walk VLM and MiniCPM-V2.6 [4] on cloud devices to verify the differences in performance between the two models in real-world scenarios. The visualization results of the two models on the video stream can



Application Scenario	Models	Input Prompt
Normalize the annotation results	Llama 3.1	<p>Please normalize the following manually annotated output to reduce information redundancy and maintain as standardized an output as possible. During the processing, please follow these guidelines:</p> <ol style="list-style-type: none"> <li>Convert all labels to lowercase.</li> <li>Remove any superfluous spaces or special characters.</li> <li>Retain the accurate position of objects in the sentence, such as what time or how many steps away.</li> <li>For similar or repetitive semantic annotations, reduce the redundancy of semantics.</li> <li>Output the result directly.</li> </ol> <p>The annotated text is as follows:</p> <p>...</p>
Instruct VLM to provide guidance on blind walking based on the provided video	DeepSeek (1.3B&7B) Yi-VL (6B) MiniCPM-V2.6 (8B) GPT-4o Qwen2-VL (7B)	<p>"request": "format the sentence below into the format, given in English",</p> <p>"restriction": "answer in json like the format given below without code block",</p> <p>"sentence": "Analyze the following video frames and determine the danger level for a blind person if they were to walk straight ahead. The danger levels are categorized as follows: Low (open areas, with few people or obstacles): Mid (moderate danger, such as areas with some obstacles or moderate activity): High (high danger, such as narrow pathways, crowded areas, or busy roads)arc the danger levels for the first two frames:1:{history_states[0]}2:history_states[1]}on this information and the provided image below, please provide the danger level for Frame 3.that danger level are indicated by single letters only. (A, B, or C)., provide walking instructions based on the provided image to ensure safe navigation.",</p> <p>"format": { "data": { "Frame 3 Danger Level": "string (A, B, or C)", "Walking Guidance": "string" } }</p> <p>...</p>
Training and inference of WalkVLM	WalkVLM	<p>"request": "format the sentence below into the format, given in English",</p> <p>"restriction": "answer in json like the format given below without code block",</p> <p>"sentence": "You are now a guide. I can't see the path and will be walking solely based on your instructions. Each input frame displays the road information ahead. The main objects in each image are { json_str }. Please provide clear and unobstructed walking directions. Describe in order: 1. Location (e.g., corridor, restaurant, market, downtown, street, foot path, bicycle path), 2. Weather conditions (e.g., overcast sky, sunny, partly cloudy, indoor), 3. Traffic flow rating (e.g., low: 0-4 people/minute, medium: 4-10 people/minute, high: 10+ people/minute), 4. Describe the overall scene based on the input images and all the information from the above three points, 5. Please guide me on how to proceed based on the input images and all previous descriptions.",</p> <p>"format": { "data": { "1. Location": "string", "2. Weather conditions": "string", "3. Traffic flow rating": "string", "4. Describe the overall scene in the image": "string", "5. Instructions on how I should proceed": "string" } }</p> <p>...</p>
Use LMM to evaluate the similarity between generated results and ground truth	GPT4	<p>Please act as an impartial judge and evaluate the quality of the responses provided by multiple assistants displayed below. You should choose the assistant that matches the GT answer. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Avoid any positional biases and ensure that the order in which the responses were presented does not influence your decision. Do not favor certain names of the assistants. Be as objective as possible. The answer should be the most closest to the semantics of the GT result and have the most concise answer. After providing your explanation, strictly follow the following format to output your final verdict: if assistant A is better, output "[[A]]", if assistant B is better, output "[[B]]". Request you select a relatively optimal result and directly output the option.</p> <p>{GT}</p> <p>{}</p> <p>The Start of Assistant A's Answer</p> <p>{}</p> <p>The End of Assistant A's Answer</p> <p>The Start of Assistant B's Answer</p> <p>{}</p> <p>The End of Assistant B's Answer</p>

Table 5. All prompts utilized in this paper.

Device	Infrastructure	Precision	Inference Speed
A100	vLLM	FP16	65-70 token/s
A100	vLLM	INT8	92-108 token/s
iPad Pro (M4)	llama.cpp	Q4_K	16-18 token/s

Table 6. Inference speed comparison of WalkVLM under different hardware platforms and quantization methods.

be viewed [here](#), where WalkVLM is capable of generating

less temporal redundancy. As shown in 13, on real-time video streams, for two models with the same size parameters, WalkVLM can generate more concise and accurate walking guidance.

However, the current model still has certain **limitations** in practical applications. **Firstly**, the model has a weak ability to prioritize events, making it difficult to identify the most urgent actions that need reminders in the scene. Facing this issue, our next attempt is to establish an event priority



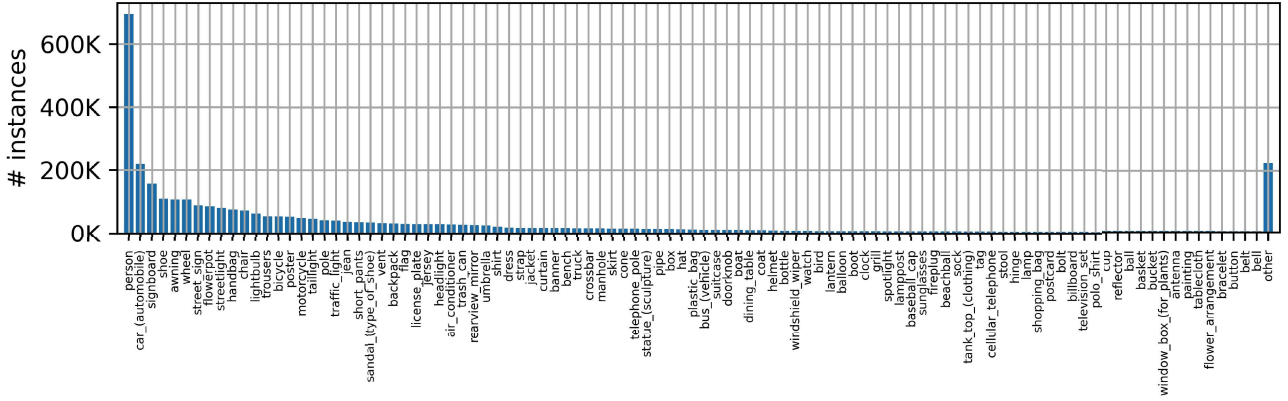


Figure 7. Detect target distribution. For clarity, display the top 100 with the highest frequency of occurrence.

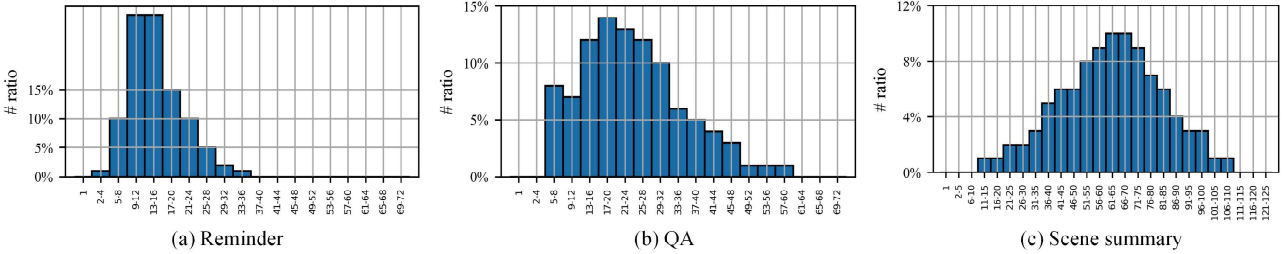


Figure 8. Data length distribution in different text annotation types.

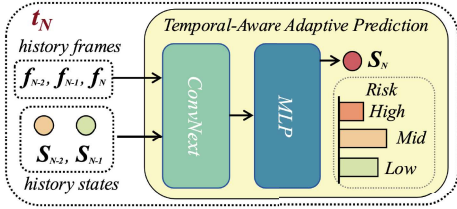


Figure 9. Architectural of TAP when the history window is 3.

blind individuals. This collaborative effort can ultimately result in technologies that are not only more effective but also more widely adopted and accessible.

From an educational standpoint, our work can also play a pivotal role in raising awareness about the challenges faced by the visually impaired community. By showcasing the potential of AI and machine learning in addressing these challenges, we hope to inspire more individuals and organizations to contribute towards creating a more inclusive society. This increased awareness can lead to more supportive policies and initiatives that focus on improving the quality of life for the visually impaired.

Additionally, the WalkVLM model and dataset have the potential to impact various industries beyond assistive technologies. For instance, they can be adapted for use in smart city planning, where understanding pedestrian behavior and

safety is crucial. This broader application can lead to safer and more accessible urban environments for everyone, not just the visually impaired.

In summary, our contribution not only advances the state of the art in AI and machine learning but also has far-reaching societal implications. By providing a robust benchmark and a rich dataset, we are paving the way for innovative solutions that can significantly enhance the lives of blind individuals and promote a more inclusive society.

## 7. Limitations

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

## 8. Acknowledgements

We appreciate the data support from the YouTube platform @poptravelorg, which has greatly enriched our data





Figure 10. Hierarchical inference results visualization of WalkVLM.

sources. We thank all impaired individuals who participated in the research and experiments for this paper, whose real feedback provided valuable insights for our system.

## References

- [1] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023. 9
- [2] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019. 2
- [3] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020. 9
- [4] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024. 7
- [5] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp

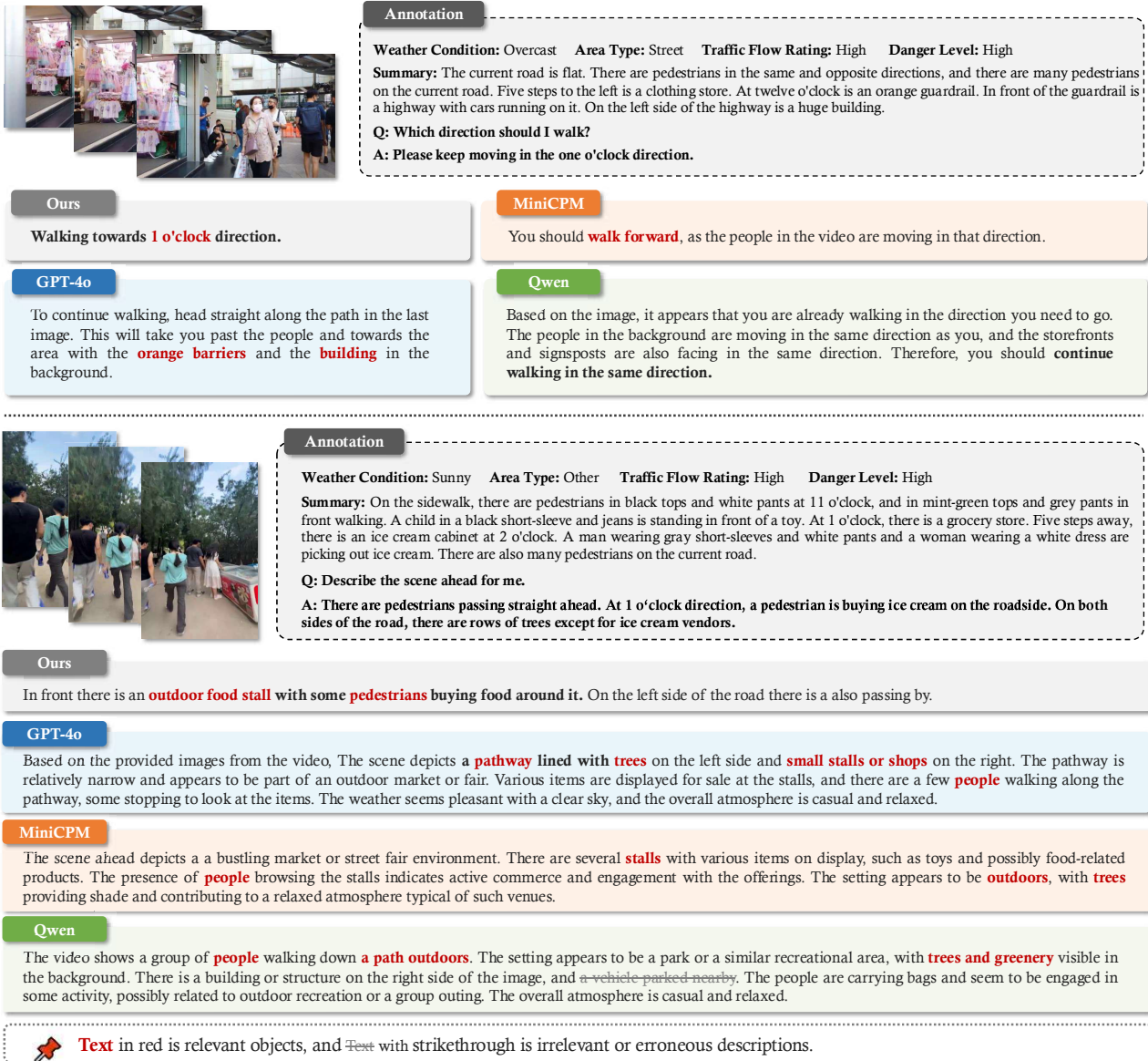


Figure 11. More QA visualization comparisons of different models.

Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022. 2,

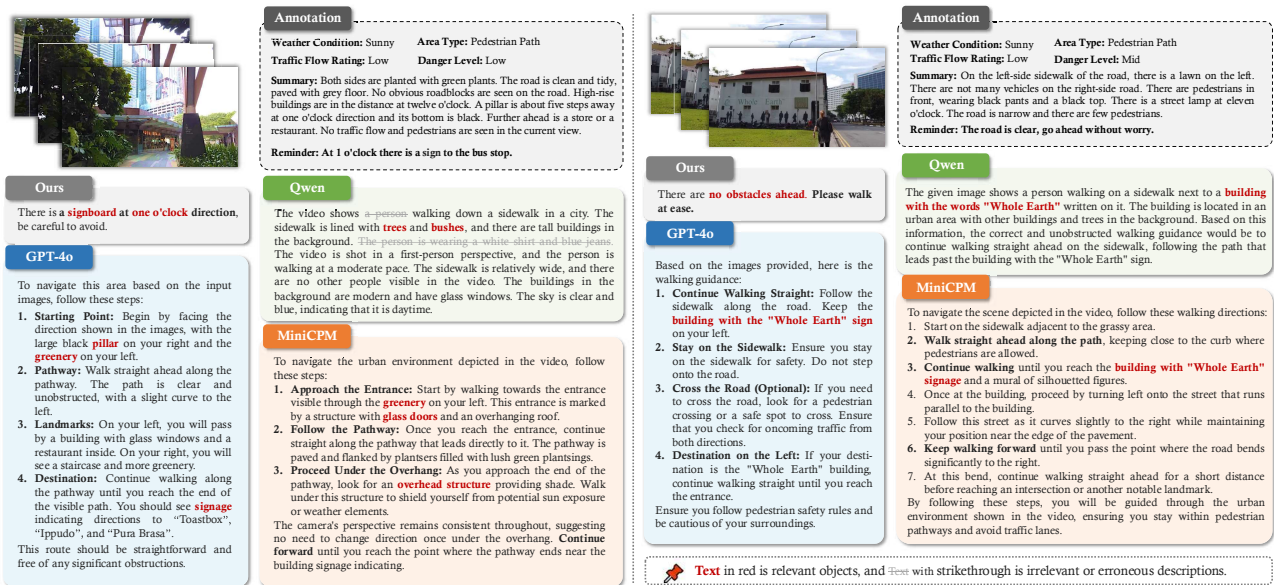


Figure 12. More reminder visualization comparisons of different models.

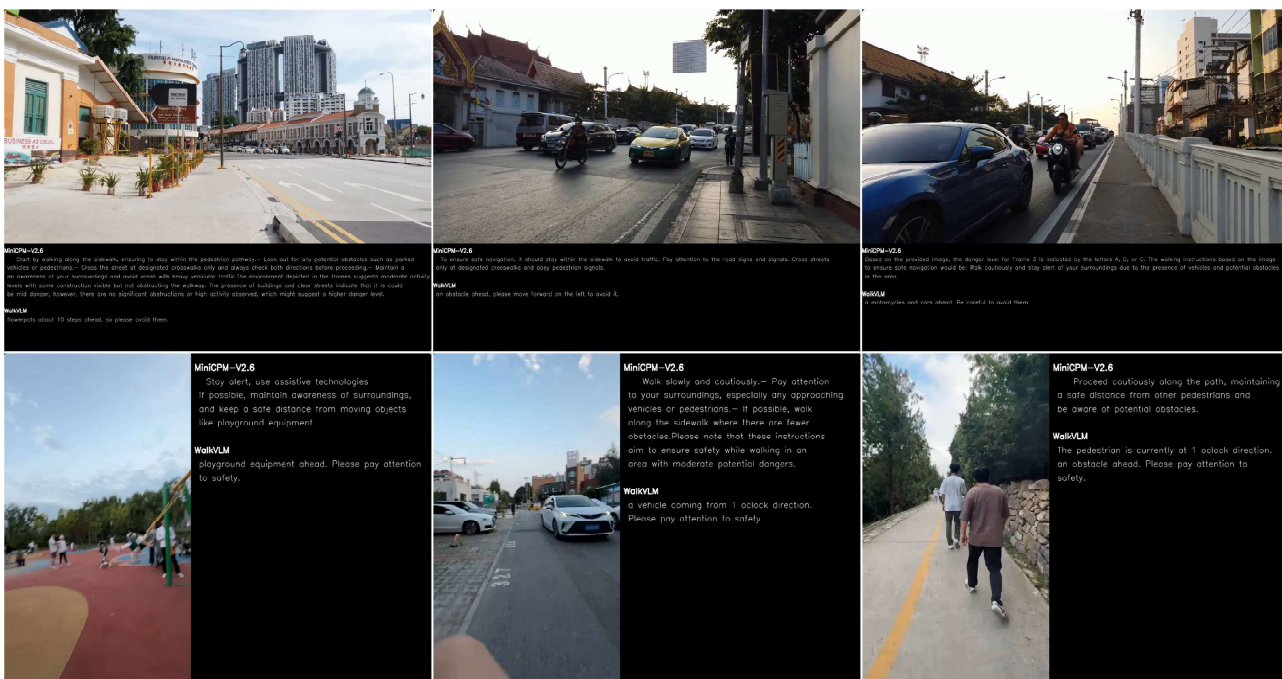


Figure 13. Sampling results of video stream inference in the blind walking task. Zoom in to view the generated results. See here for the video inference results. WalkVLM is capable of generating less temporal redundancy and providing more concise and informative responses.