

# “Does the cafe entrance look accessible? Where is the door?” Towards Geospatial AI Agents for Visual Inquiries

Anonymous ICCV submission

Paper ID \*\*\*\*\*



Figure 1. We introduce our vision for *Geo-Visual Agents*—multimodal AI agents capable of understanding and responding to nuanced visual-spatial inquiries about the world by analyzing large-scale repositories of geospatial images combined with traditional GIS data sources. For example, *StreetViewAI* [14] (above) makes street view accessible to blind users by combining geographic context, user information, and dynamic street view images into an MLLM, accessed via an AI chat interface and accessible screen reader controls.

## Abstract

001 *Interactive digital maps have revolutionized how people*  
 002 *travel and learn about the world; however, they rely on pre-*  
 003 *existing structured data in GIS databases (e.g., road net-*  
 004 *works, POI indices), limiting their ability to address geo-*  
 005 *visual questions related to what the world looks like. We*  
 006 *introduce our vision for Geo-Visual Agents—multimodal*  
 007 *AI agents capable of understanding and responding to nu-*  
 008 *anced visual-spatial inquiries about the world by analyz-*  
 009 *ing large-scale repositories of geospatial images, including*  
 010 *streetscapes (e.g., Google Street View), place-based photos*  
 011 *(e.g., TripAdvisor, Yelp), and aerial imagery (e.g., satel-*  
 012 *lite photos) combined with traditional GIS data sources.*  
 013 *We define our vision, describe sensing and interaction ap-*  
 014 *proaches, provide three exemplars, and enumerate key chal-*  
 015 *lenges and opportunities for future work.*

## 1. Introduction

Over the last two decades, precise location sensing, per-  
 vasive internet connectivity, and interactive digital maps  
 have transformed travel planning and in situ navigation,  
 enabling turn-by-turn directions, location-aware search,  
 and dynamic route optimization. Despite these advances,  
 current mapping systems are confined to pre-existing struc-  
 tured geospatial data, leaving a vast repository of visual  
 information—latent within street-level, aerial, and user-  
 contributed imagery—untapped and inaccessible for an-  
 swering what we term *geo-visual questions*. That is,  
 visually-oriented questions about a location. Imagine, for  
 example, a wheelchair user asking “*Are there stairs lead-*  
*ing up to the library on 35th?*” or a blind traveler inquiring  
 “*Where is the door to the cafe and what does it look like?*”

In this workshop paper, we introduce our vision for  
*Geo-Visual Agents*—multimodal AI agents capable of un-  
 derstanding and responding to nuanced visual-spatial in-  
 quiries about the world by analyzing large-scale repositories

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035 of geospatial images (e.g., street-level and aerial imagery)  
 036 combined with traditional GIS databases (e.g., road net-  
 037 works, POI databases, transit schedules). We envision Geo-  
 038 Visual Agents acting as “visual-spatial co-pilots” across a  
 039 spectrum of contexts from *a priori* travel planning to *in situ*  
 040 navigation. Crucially, while we expect many high-value  
 041 user scenarios where a Geo-Visual Agent is actively sens-  
 042 ing and processing visual-spatial data in real-time via AR  
 043 glasses [12, 33, 62] or smartphone cameras [39, 49, 56], an  
 044 equally large set of questions can be answered by analyzing  
 045 existing (and largely untapped) repositories of geo-related  
 046 imagery—either on-demand (e.g., spinning up an AI agent  
 047 to query and analyze sources) or via pre-computation.

048 Our vision moves beyond the current paradigm of  
 049 geospatial artificial intelligence (GeoAI) [17, 30, 35] such  
 050 as *CARTO AI* [7] and *SuperMap* [51], which primarily fo-  
 051 cuses on large-scale data analysis for domain experts. Sim-  
 052 ilarly, our work is related to but distinct from emerging  
 053 paradigms in GIS research such as “Autonomous GIS”—  
 054 AI-based scientific assistants that help “reason, derive, in-  
 055 novate, and advance geospatial solutions to pressing global  
 056 challenges” [38]. Moreover, because our envisioned agents  
 057 work primarily via multimodal conversational AI, we draw  
 058 inspiration from recent work in *Geospatial Visual Question*  
 059 *Answering* (GVQA) such as *MQVQA* [61] and *TAMMI* [6],  
 060 which attempt to imbue multimodal LLMs with domain-  
 061 specific geographic knowledge; however, again these sys-  
 062 tems are aimed at analysts and function primarily on remote  
 063 aerial imagery. While related, our focus is on addressing the  
 064 personal, interactive, and often immediate needs of an indi-  
 065 vidual planning travel or actively navigating a space.

066 Below, we expand on our vision including a breakdown  
 067 of visual-spatial inquiries, modalities of sensing and inter-  
 068 action, and three emerging examples, *StreetViewAI* [14],  
 069 *Access Scout* [27], and *BikeButler*. Throughout, we high-  
 070 light key opportunities and open challenges.

## 071 2. Geo-Visual Queries Across Travel Stages

072 We envision Geo-Visual Agents providing value across the  
 073 full mobility cycle from pre-travel planning to *in-situ* nav-  
 074 igation. Below, we enumerate four travel stages and op-  
 075 portunities for Geo-Visual Agents therein, focusing on ac-  
 076 cessibility but also broader user scenarios such as driving  
 077 and biking. Selecting and fusing data sources will be a  
 078 function of user task and data availability. For example,  
 079 pre-travel planning may rely on streetscape images, user-  
 080 contributed photos, and place-based reviews while in-situ  
 081 navigation might combine these sources with visual content  
 082 from a user’s real-time camera feed (e.g., from AR glasses)  
 083 and context sensing (e.g., travel mode inference, location).

084 **Pre-travel planning.** In this phase, the user is not phys-  
 085 ically present at a location but planning a future visit. The  
 086 agent acts as a remote, interactive guide, enabling detailed

087 investigation and reducing uncertainty before travel. For ex-  
 088 ample: (1) a blind parent planning a trip to a park may ask,  
 089 “What kind of equipment does the playground have, and  
 090 does it seem safe?” (2) A person with a mobility disability  
 091 virtually investigates a route and inquires “Are there curb  
 092 ramps all the way to my doctor’s office?” (3) A potential  
 093 homebuyer may ask neighborhood-related questions such  
 094 as “What do the streets look like?”, “Are there tree-lined  
 095 sidewalks?”, and “How much graffiti is there?”

096 **While navigating.** During travel itself, the user is under  
 097 cognitive and physical load, navigating their environment,  
 098 making route choices, and dynamically avoiding obstacles.  
 099 Here, the agent provides forward-looking information about  
 100 the destination or upcoming maneuvers, enhancing situa-  
 101 tional awareness and facilitating *in situ* travel decisions. For  
 102 example: (1) A driver approaching a destination asks, “You  
 103 said to turn left at the next light. Are there any landmarks?”  
 104 (2) A cyclist nearing a decision point queries, “Is there a  
 105 protected bike lane at the next intersection, and which side  
 106 of the road is it on?” (3) A rail user exiting a train asks,  
 107 “Which exit is closest to the library’s accessible entrance?”

108 **Destination arrival.** When arriving at a destination, the  
 109 user is faced with a litany of “last 10 meters” problems re-  
 110 lated to the appearance of their destination, the path to and  
 111 location of an entrance, and the presence of obstacles or  
 112 safety issues. For example, (1) approaching their destina-  
 113 tion, a delivery driver may inquire “Where is the loading  
 114 zone for this building?”; (2) a person meeting a friend in a  
 115 busy plaza may ask, “I’m looking for the coffee shop; can  
 116 you describe its storefront so I can more easily spot it?”.  
 117 (3) a blind traveler’s ride share arrives for pickup at a busy  
 118 airport and asks, “Can you help me find the silver Toyota  
 119 Camry with license plate KNI667?”.

120 **Indoor exploration.** Finally, upon entering a desti-  
 121 nation, the agent’s role can shift to supporting micro-  
 122 navigation through complex indoor environments like air-  
 123 ports, stores, or office buildings. This stage presents a sig-  
 124 nificant data challenge, as comprehensive visual and map  
 125 datasets for indoor spaces are rare [13]. For example, (1)  
 126 a customer trying to find the location of a specific item in  
 127 a hardware store may ask “Based on the aisle signs, which  
 128 direction do I go to find the plumbing department?” (2)  
 129 A low-vision traveler looking at an airport departure board:  
 130 “Can you tell me which gate Delta Flight 850 is leaving  
 131 from?”; (3) A wheelchair user in a large convention center:  
 132 “Can you guide me to the nearest accessible restroom?”

133 Together, these scenarios illustrate how Geo-Visual  
 134 Agents can transform how we navigate and understand  
 135 places, enhancing accessibility, offering landmark-based  
 136 navigation, improving personal safety, and even leading to  
 137 serendipitous discovery. Below, we describe potential data  
 138 sources and then outline interaction modalities.

139	<b>3. Sensing and Data Sources</b>	192
140	The power of a Geo-Visual Agent lies in its ability to syn-	193
141	thesize heterogeneous data sources, fusing visual evidence	194
142	with structured geospatial data to form a holistic and accu-	195
143	rate understanding of a place or route. We focus below on	196
144	geo-related image sources rather than structured GIS data.	197
145	<b>Streetscape Imagery.</b> Street view imagery (SVI) [25,	198
146	37]—such as <i>Google Street View</i> (GSV), <i>Cyclomedia</i> , and	199
147	<i>Mapillary</i> —provide a rich, large-scale image archive of	200
148	the world. GSV alone has over 220 billion images span-	201
149	ning 10 million miles across 100 countries [19]. Such data	202
150	can be used to analyze road conditions [3], street mark-	203
151	ings (crosswalks [2, 34], bike lanes [45]), sidewalk infras-	204
152	tructure (sidewalk material [23], curb ramps [21, 41]), bus	205
153	stops [32], building facades [31], graffiti [52], trees and veg-	206
154	etation [36], neighborhood health indicators [54, 63], and	207
155	more. Primary limitations include image recency [55], oc-	208
156	clusions due to obstructing objects in front of the SVI cam-	209
157	era ( <i>e.g.</i> , buses) [47], and geographic distribution (images	210
158	are distributed every 10-15 meters along roadways but not	211
159	foot pathways or inside parks or buildings).	212
160	<b>User-Contributed Photos.</b> Place-based platforms like	213
161	<i>Google Places</i> , <i>Yelp</i> , and <i>TripAdvisor</i> contain vast, crowd-	214
162	sourced libraries of photos tied to specific POIs, which pro-	215
163	vide a useful complement to SVI, including building interi-	216
164	ors, curated (business uploaded) shots of storefronts, and	217
165	pictures of menus, food [16], and social activities ( <i>e.g.</i> ,	218
166	[60])—all which are often accompanied by user-contributed	219
167	text ( <i>e.g.</i> , reviews). We found, however, that analysis of	220
168	such multimodal data is less common in the literature. The	221
169	key limitation here is data availability, particularly for un-	222
170	popular or recently opened places, and social biases in <i>who</i>	223
171	uploads and <i>why</i> ( <i>e.g.</i> , see [4, 58]).	224
172	<b>Aerial Imagery.</b> Aerial imagery from satellites, air-	225
173	planes, or drones can provide high-resolution, top-down or	226
174	oblique (45-degree angle) views of spatial structures, in-	227
175	cluding building footprints, parking lots, vegetation, and	228
176	pedestrian infrastructure [24]. While remote sensing and	229
177	photogrammetry research has existed for many decades—	230
178	<i>e.g.</i> , for land use classification, agriculture, disaster re-	231
179	sponse, and military analyses [29, 59]—such techniques	232
180	have not been applied to the Geo-Visual Agent context	233
181	( <i>e.g.</i> , answering queries about parking lot locations, rooftop	234
182	restaurant patios, or unmapped pedestrian shortcuts). Sim-	235
183	ilar to streetscapes, aerial imagery can suffer from occlu-	236
184	sions (from tree cover, clouds), shadows from tall buildings,	237
185	and lack of availability. In the US, high-resolution aerial im-	238
186	agery is often provided by the federal government such as	239
187	USGS [53] and NASA [40].	240
188	<b>Robotic scans.</b> Robots such as autonomous vehicles,	241
189	ground-based delivery robots, and drones [48, 50] infused	
190	with sensor suites (cameras, LiDAR) can generate high-	
191	fidelity scans of the environment, producing not just images	
	but 3D reconstructions with mensuration [26]. While a po-	
	tentially promising future data source, there is currently a	
	lack of open data and APIs.	
	<b>Infrastructure-based Cameras.</b> Infrastructure-based	
	cameras installed for traffic, weather, security, and safety	
	monitoring provide real-time views of cities and uniquely	
	offer dynamic information about pedestrian and car move-	
	ment, human activity, weather conditions, and transient ob-	
	structions [28, 43, 46]; however, while some camera feeds	
	are open— <i>e.g.</i> , DOT traffic cameras—most are not and pri-	
	vacancy is a key consideration. Moreover, there is a lack of	
	density and availability ( <i>e.g.</i> , in rural areas).	
	<b>First-person Camera Streams.</b> Finally, first-person	
	camera streams from AR glasses [12, 33, 62], smartphone	
	cameras [5, 39, 49, 56], and dashcams [42, 57] are criti-	
	cal for in-situ travel stages, offering a real-time, egocen-	
	tric view for navigation, identifying transient obstacles, and	
	reading signs. While primarily used for immediate assis-	
	tance, these streams could also help update or correct exist-	
	ing geospatial datasets in a continuous feedback loop ( <i>e.g.</i> ,	
	[57]). However, key considerations include high computa-	
	tional and power requirements, robust network connectivity,	
	and privacy concerns for both the user and bystanders.	
	<b>4. Processing and Interpreting with AI</b>	
	Our vision relies not just on diverse forms of geospatial im-	
	agery and pre-existing GIS data but also advances in mul-	
	timodal AI ( <i>e.g.</i> , scene understanding [9, 11], object affor-	
	dances [22, 33], and spatial reasoning [8, 10, 15, 44]) to ex-	
	tract semantic information and object relationships. While	
	some analyses could be pre-computed for known high-value	
	entities ( <i>e.g.</i> , presence and location of curb ramps [21, 41]),	
	we expect a long-tail of bespoke queries, which will re-	
	quire a Geo-Visual Agent to seek out, analyze, and synthe-	
	size image-based sources with pre-existing metadata in GIS	
	databases in real-time.	
	<b>5. Delivering the Answers</b>	
	Finally, a crucial aspect of our vision is <i>how</i> the agent deliv-	
	ers information, which is a function of the user’s abilities,	
	their current context, and the complexity and type of data.	
	Regardless of delivery mode, agents need to report uncer-	
	tainty and data provenance to build trust and mitigate error.	
	<b>Audio-First Interfaces:</b> For hands-free and/or eyes-free	
	operation—essential for drivers, cyclists, and blind and low	
	vision users—audio interfaces are critical ( <i>e.g.</i> , using ear-	
	buds or a smart speaker). The challenge, however, is provid-	
	ing well-structured verbal descriptions to convey complex	
	visual information without overwhelming the user.	
	<b>Multimodal Interfaces:</b> Agents should also select and	
	show relevant imagery. For instance, after describing an	
	entrance, the agent could display a photo of the door	



242 (e.g., drawn from SVI or Yelp). The challenge lies in  
243 the AI’s ability to select the most appropriate photo(s)—  
244 appropriately cropped—from large archives.

245 **AI-Generated Abstracted Visualizations:** For highly  
246 complex spatial information, a raw photo or a long verbal  
247 description may be insufficient. An exciting frontier is the  
248 agent’s ability to generate simplified, abstract diagrams on  
249 the fly—akin to a modern *LineDrive* system [1]. Making  
250 these abstractions accessible, perhaps tactilely, is also a crit-  
251 ical area of open research.

## 252 6. Case Study Applications

253 To help showcase and concretize our vision, we highlight  
254 three emerging Geo-Visual Agent prototypes.

255 **StreetViewAI.** Current SVI tools are inaccessible to  
256 blind users. Our group is addressing this problem through  
257 the design of *StreetViewAI* [14] (Figure 1), which uses  
258 context-aware, real-time AI to support virtually explor-  
259 ing routes, inspecting destinations, or even remotely vis-  
260 iting tourist locations such as the Grand Canyon [18].  
261 *StreetViewAI* provides accessible interactive controls for  
262 blind users to pan and move between panoramic images  
263 and dynamically converse with a live, multimodal AI agent  
264 about the scene and local geography. In a lab study, blind  
265 users effectively used *StreetViewAI* to virtually navigate  
266 streetscapes. Key challenges: reconciling users’ mental  
267 models of SVI, a tendency to over-trust AI, and the diffi-  
268 culty of synthesizing rich visual data into concise audio.

269 *AI Agent.* *StreetViewAI* employs three separate AI sub-  
270 systems. Most relevant is the *AI Chat Agent*, which allows  
271 for conversational interactions about the user’s current and  
272 past street views as well as nearby geography. The agent  
273 uses Google’s *Multimodal Live API* [20], which supports  
274 real-time interaction, function calling, and retains memory  
275 of all interactions within a single session. When the user ini-  
276 tiates a chat either via typing or speaking, we transmit each  
277 GSV interaction along with the user’s current view and geo-  
278 graphic context (e.g., nearby places, current heading). Thus,  
279 they can ask about local geography, current and past views,  
280 and object relationships (e.g., “where is the entrance?”).

281 **Accessibility Scout.** Assessing the accessibility of un-  
282 familiar environments is a critical but often laborious job  
283 for people with disabilities. While standardized checklists  
284 exist, they often fail to account for an individual’s unique  
285 and evolving needs. *Accessibility Scout* [27] is an LLM-  
286 based system designed to address this gap by generating  
287 personalized accessibility scans from images—e.g., from  
288 *TripAdvisor*, *Yelp*, and *Airbnb*—to identify potential con-  
289 cerns based on self-reported abilities and interests. In user  
290 studies, we found that *Accessibility Scout*’s personalized  
291 scans were more useful than generic ones and that its col-  
292 laborative Human-AI approach was effective and built trust.

293 *AI Agent.* The *Accessibility Scout* pipeline begins by

294 creating a structured user model in JSON format, initial-  
295 ized from a user’s plain text description of their abilities  
296 and preferences. To assess an environment, the agent mim-  
297 ics how users assess environmental accessibility by first an-  
298 alyzing an image and the user’s intent (e.g., “going on a  
299 date”) to identify potential tasks a user might perform, such  
300 as “dining” or “toileting”. The agent then decomposes  
301 these tasks into primitive motions like “grabbing” that are  
302 required to complete them. For each task, the agent ana-  
303 lyzes the user model, task information, and segmented im-  
304 age to identify and describe environmental concerns. Cru-  
305 cially, the system is designed for Human-AI collaboration;  
306 users can provide feedback on identified concerns which the  
307 agent uses to update the user model.

308 **BikeButler.** Existing mapping tools define optimal bike  
309 routes using objective data like distance and elevation, but  
310 often ignore subjective qualities related to a cyclist’s com-  
311 fort and perceived safety. However, a desirable bike route  
312 depends on factors not found in standard GIS databases,  
313 such as the presence of tree-lined streets, pavement qual-  
314 ity, or bike lane widths. *BikeButler* is an early-stage proto-  
315 type Geo-Visual Agent that generates personalized cycling  
316 routes by fusing structured data from *OpenStreetMap* with  
317 visual analysis of SVI. The system creates routes optimized  
318 for a user’s specific profile (e.g., beginner, expert) and al-  
319 lows them to rate route segments, creating a feedback loop  
320 that refines their preferences for future journeys.

## 321 7. Discussion and Conclusion

322 In this paper, we introduced our vision for Geo-Visual  
323 Agents, dynamic and conversational AI co-pilots that can  
324 see and reason about the world in real-time. Our envisioned  
325 agents answer nuanced visual questions about the visual  
326 world—from a blind user navigating a complex intersec-  
327 tion to a cyclist seeking the safest, most pleasant route. Our  
328 prototypes offer an initial window into this vision, offering  
329 personalized, interactive experiences extending far beyond  
330 current mapping services.

331 Still, significant challenges remain, including: (1) *Dy-*  
332 *namic information synthesis:* creating agents that can intel-  
333 ligently select, fuse, and reason over a heterogeneous set  
334 of real-time and archived data sources; (2) *Trust and trans-*  
335 *parency:* communicating uncertainty and data provenance;  
336 (3) *Speech UIs:* effectively verbalizing complex visual in-  
337 formation concisely via text or speech; (4) *Personalization*  
338 learning from a user’s unique needs and preferences; (5)  
339 *Spatial reasoning* accurately tracking and modeling spatial  
340 relationships between objects; (6) *Generative spatial ab-*  
341 *stractions:* dynamically generating spatial visualizations to  
342 help aid understanding.

343 Addressing these challenges will require a concerted ef-  
344 fort across disciplines from computer vision and HCI to ac-  
345 cessibility and geospatial science. Join us!

346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402

**References**

[1] Maneesh Agrawala and Chris Stolte. Rendering effective route maps: improving usability through generalization. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, page 241–249, New York, NY, USA, 2001. Association for Computing Machinery. 4

[2] Dragan Ahmetovic, Roberto Manduchi, James M. Coughlan, and Sergio Mascetti. Mind your crossings: Mining gis imagery for crosswalk localization. *ACM Trans. Access. Comput.*, 9(4), 2017. 3

[3] Shazab Ali, Meng Xu, and Daehan Kwak. Smart roadway monitoring: Pothole detection and mapping via google street. In *Internet Computing and IoT and Embedded Systems, Cyber-physical Systems, and Applications: 25th International Conference, ICOMP 2024, and 22nd International Conference, ESCS 2024, Held as Part of the World Congress in Computer Science, Computer Engineering and Applied Computing, CSCE 2024, Las Vegas, NV, USA, July 22–25, 2024, Revised Selected Papers*, page 151. Springer Nature, 2025. 3

[4] V. Antoniou and A. Skopeliti. Measures and indicators of vgi quality: An overview. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, II-3/W5: 345–351, 2015. 3

[5] Apple. Detect doors around you using Magnifier on iPhone, 2025. 3

[6] Hichem Boussaid, Lucrezia Tosato, Flora Weissgerber, Camille Kurtz, Laurent Wendling, and Sylvain Lobry. Visual question answering on multiple remote sensing image modalities. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR) Workshops*, pages 2319–2328, 2025. 2

[7] CARTO. Genai — ai-powered spatial insights, 2025. 2

[8] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14455–14465, 2024. 3

[9] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, 2018. 3

[10] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. In *Advances in Neural Information Processing Systems*, pages 135062–135093. Curran Associates, Inc., 2024. 3

[11] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3

[12] Alexander Fiannaca, Ilias Apostolopoulos, and Eelke Folmer. Headlock: a wearable navigation aid that helps blind cane users traverse large open spaces. In *Proceedings of the 16th International ACM SIGACCESS Conference on Computers & Accessibility*, page 19–26, New York, NY, USA, 2014. Association for Computing Machinery. 2, 3

[13] Jon E. Froehlich, Anke M. Brock, Anat Caspi, João Guerreiro, Kotaro Hara, Reuben Kirkham, Johannes Schöning, and Benjamin Tannert. Grand challenges in accessible maps. *Interactions*, 26(2):78–81, 2019. 2

[14] Jon E. Froehlich, Alex Fiannaca, Nimer Jaber, Victor Tsaran, and Shaun Kane. Streetviewai: Making street view accessible using context-aware multimodal ai. In *The 38th Annual ACM Symposium on User Interface Software and Technology*, page 22, New York, NY, USA, 2025. ACM. 1, 2, 4

[15] Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. Scene-llm: Extending language model for 3d visual understanding and reasoning, 2024. 3

[16] Alessandro Gambetti and Qiwei Han. Aigen-foodreview: A multimodal dataset of machine-generated restaurant reviews and images on social media, 2024. 3

[17] Google. Google earth ai: Our state-of-the-art geospatial ai models. <https://blog.google/technology/ai/google-earth-ai/>, 2025. Accessed: August 1, 2025. 2

[18] Google. Treks: Grand canyon, 2025. 4

[19] Google. Celebrate 15 years of exploring your world on Street View, 2025. 3

[20] Google. Vertex ai multimodal live api, 2025. 4

[21] Kotaro Hara, Jin Sun, Robert Moore, David Jacobs, and Jon Froehlich. Tohme: detecting curb ramps in google street view using crowdsourcing, computer vision, and machine learning. In *Proceedings of the 27th Annual ACM Symposium on User Interface Software and Technology*, page 189–204, New York, NY, USA, 2014. Association for Computing Machinery. 3

[22] Mohammed Hassanin, Salman Khan, and Murat Tahtali. Visual affordance and function understanding: A survey. *ACM Comput. Surv.*, 54(3), 2021. 3

[23] Maryam Hosseini, Fabio Miranda, Jianzhe Lin, and Claudio T. Silva. Citysurfaces: City-scale semantic segmentation of sidewalk materials. *Sustainable Cities and Society*, 79: 103630, 2022. 3

[24] Maryam Hosseini, Andres Sevtsuk, Fabio Miranda, Roberto M. Cesar, and Claudio T. Silva. Mapping the walk: A scalable computer vision approach for generating sidewalk network datasets from aerial imagery. *Computers, Environment and Urban Systems*, 101:101950, 2023. 3

[25] Yujun Hou and Filip Biljecki. A comprehensive framework for evaluating the quality of street view imagery. *International Journal of Applied Earth Observation and Geoinformation*, 115:103094, 2022. 3

[26] Dingkun Hu and Jennifer Minner. Uavs and 3d city modeling to aid urban planning and historic preservation: A systematic review. *Remote Sensing*, 15(23), 2023. 3

[27] William Huang, Xia Su, Jon E. Froehlich, and Yang Zhang. Accessibility scout: Personalized accessibility scans of built

460 environments. In *The 38th Annual ACM Symposium on User*  
461 *Interface Software and Technology*, page 18, New York, NY,  
462 USA, 2025. ACM. 2, 4

[28] Gaurav Jain, Basel Hindi, Zihao Zhang, Koushik Srinivasula,  
463 Mingyu Xie, Mahshid Ghasemi, Daniel Weiner, Sophie Ana  
464 Paris, Xin Yi Therese Xu, Michael Malcolm, Mehmet Kerem  
465 Turkan, Javad Ghaderi, Zoran Kostic, Gil Zussman, and  
466 Brian A. Smith. Streetnav: Leveraging street cameras to  
467 support precise outdoor navigation for blind pedestrians. In  
468 *Proceedings of the 37th Annual ACM Symposium on User*  
469 *Interface Software and Technology*, New York, NY, USA,  
470 2024. Association for Computing Machinery. 3

[29] Bhargavi Janga, Gokul Prathin Asamani, Ziheng Sun, and  
471 Nicoleta Cristea. A review of practical ai for remote sensing  
472 in earth sciences. *Remote Sensing*, 15(16), 2023. 3

[30] Krzysztof Janowicz, Song Gao, Grant McKenzie, Yingjie  
473 Hu, and Budhendra Bhaduri. Geoai: spatially explicit arti-  
474 ficial intelligence techniques for geographic knowledge dis-  
475 covery and beyond. *International Journal of Geographical*  
476 *Information Science*, 34(4):625–636, 2020. 2

[31] Hyejin Kim, Seula Park, and Jiyoung Kim. A study on  
477 barrier-free entrance object detection using deep learning in  
478 street view imagery. In *2024 IEEE International Conference*  
479 *on Big Data (BigData)*, pages 8716–8718, 2024. 3

[32] Minchu Kulkarni, Chu Li, Jaye Jungmin Ahn, Katrina  
480 Oi Yau Ma, Zhihan Zhang, Michael Saugstad, Kevin Wu,  
481 Yochai Eisenberg, Valerie Novack, Brent Chamberlain, and  
482 Jon E. Froehlich. Busstopcv: A real-time ai assistant for la-  
483 beling bus stop accessibility features in streetscape imagery.  
484 In *Proceedings of the 25th International ACM SIGACCESS*  
485 *Conference on Computers and Accessibility*, New York, NY,  
486 USA, 2023. Association for Computing Machinery. 3

[33] Jaewook Lee, Andrew D. Tjahjadi, Jiho Kim, Junpu Yu,  
487 Minji Park, Jiawen Zhang, Jon E. Froehlich, Yapeng Tian,  
488 and Yuhang Zhao. Cookar: Affordance augmentations in  
489 wearable ar to support kitchen tool interactions for peo-  
490 ple with low vision. In *Proceedings of the 37th Annual*  
491 *ACM Symposium on User Interface Software and Technol-*  
492 *ogy*, New York, NY, USA, 2024. Association for Computing  
493 Machinery. 2, 3

[34] Meiqing Li, Hao Sheng, Jeremy Irvin, Heejung Chung, An-  
494 drew Ying, Tiger Sun, Andrew Y Ng, and Daniel A Ro-  
495 driguez. Marked crosswalks in us transit-oriented station ar-  
496 eas, 2007–2020: A computer vision approach using street  
497 view imagery. *Environment and Planning B: Urban Analyt-*  
498 *ics and City Science*, 50(2):350–369, 2023. 3

[35] Wenwen Li and Chia-Yu Hsu. Geoai for large-scale image  
499 analysis and machine vision: Recent progress of artificial  
500 intelligence in geography. *ISPRS International Journal of*  
501 *Geo-Information*, 11(7), 2022. 2

[36] Xiaojiang Li, Carlo Ratti, and Ian Seiferling. Quantifying the  
502 shade provision of street trees in urban landscape: A case  
503 study in boston, usa, using google street view. *Landscape*  
504 *and Urban Planning*, 169:81–91, 2018. 3

[37] Yongchang Li, Li Peng, Chengwei Wu, and Jiazhen Zhang.  
505 Street view imagery (svi) in the built environment: A theo-  
506 retical and systematic review. *Buildings*, 12(8), 2022. 3

[38] Zhenlong Li, Huan Ning, Song Gao, Krzysztof Janowicz,  
507 Wenwen Li, Samantha T. Arundel, Chaowei Yang, Budhen-  
508 dra Bhaduri, Shaowen Wang, A-Xing Zhu, Mark Gahegan,  
509 Shashi Shekhar, Xinyue Ye, Grant McKenzie, Guido Cer-  
510 vone, and Michael E. Hodgson. Giscience in the era of arti-  
511 ficial intelligence: A research agenda towards autonomous  
512 gis, 2025. 2

[39] Alice Lo Valvo, Daniele Croce, Domenico Garlisi, Fabrizio  
513 Giuliano, Laura Giarré, and Ilenia Tinnirello. A navigation  
514 and augmented reality system for visually impaired people.  
515 *Sensors*, 21(9), 2021. 2, 3

[40] National Aeronautics and Space Administration and U.S.  
516 Geological Survey. Landsat data access. <https://landsat.gsfc.nasa.gov/data/data-access/>,  
2025. Free access to Landsat satellite imagery archive dating  
back to 1972. Joint NASA-USGS program providing contin-  
uous Earth observation data. 3

[41] John S. O’Meara, Jared Hwang, Zeyu Wang, Michael  
Saugstad, and Jon E. Froehlich. Rampnet: A two-stage  
pipeline for bootstrapping curb ramp detection in streetscape  
images from open government metadata. In *Workshop on Vi-*  
*sion Foundation Models and Generative AI for Accessibility:*  
*Challenges and Opportunities at ICCV 2025*. IEEE, 2025.  
Workshop Paper. 3

[42] Sangkeun Park, Joohyun Kim, Rabeb Mizouni, and Uichin  
Lee. Motives and concerns of dashcam video sharing. In  
*Proceedings of the 2016 CHI Conference on Human Fac-*  
*tors in Computing Systems*, page 4758–4769, New York, NY,  
USA, 2016. Association for Computing Machinery. 3

[43] Yurii Piadyk, Joao Rulff, Ethan Brewer, Maryam Hosseini,  
Kaan Ozbay, Murugan Sankaradas, Srimat Chakradhar, and  
Claudio Silva. Streetaware: A high-resolution synchronized  
multimodal urban scene dataset. *Sensors*, 23(7), 2023. 3

[44] Kanchana Ranasinghe, Satya Narayan Shukla, Omid Pour-  
saeed, Michael S. Ryoo, and Tsung-Yu Lin. Learning to lo-  
calize objects improves spatial reasoning in visual-llms. In  
*Proceedings of the IEEE/CVF Conference on Computer Vi-*  
*sion and Pattern Recognition (CVPR)*, pages 12977–12987,  
2024. 3

[45] Luís Rita, Ricky Nathvani, Miguel Peliteiro, Tudor-Codrin  
Bostan, Emily Muller, Esra Suel, A. Barbara Metzler, Tiago  
Tamagusko, and Adelino Ferreira. Using deep learning and  
google street view imagery to assess and improve cyclist  
safety in london. *Sustainability*, 15(13), 2023. 3

[46] Joao Rulff, Giancarlo Pereira, Maryam Hosseini, Marcos  
Lage, and Claudio Silva. Towards data-informed interven-  
tions: Opportunities and challenges of street-level multi-  
modal sensing, 2024. 3

[47] Manaswi Saha, Michael Saugstad, Hanuma Teja Maddali,  
Aileen Zeng, Ryan Holland, Steven Bower, Aditya Dash,  
Sage Chen, Anthony Li, Kotaro Hara, and Jon Froehlich.  
Project sidewalk: A web-based crowdsourcing tool for col-  
lecting sidewalk accessibility data at scale. In *Proceedings of*  
*the 2019 CHI Conference on Human Factors in Computing*  
*Systems*, page 1–14, New York, NY, USA, 2019. Association  
for Computing Machinery. 3

[48] Hunsoo Song, Joshua Carpenter, Jon E. Froehlich, and Jinha  
Jung. Accessible area mapper for inclusive and sustainable



- 575 urban mobility: A preliminary investigation of airborne point  
576 clouds for pathway analysis. In *1st ACM SIGSPATIAL Work-*  
577 *shop on Sustainable Mobility (SuMob 2023)*, 2023. 3
- [49] Xia Su, Han Zhang, Kaiming Cheng, Jaewook Lee, Qiaochu  
578 Liu, Wyatt Olson, and Jon E. Froehlich. Rassar: Room ac-  
579 cessibility and safety scanning in augmented reality. In *Pro-*  
580 *ceedings of the 2024 CHI Conference on Human Factors in*  
581 *Computing Systems*, New York, NY, USA, 2024. Association  
582 for Computing Machinery. 2, 3
- [50] Xia Su, Ruiqi Chen, Jingwei Ma, Chu Li, and Jon E.  
583 Froehlich. Flymethrough: Human-ai collaborative 3d in-  
584 door mapping with commodity drones. In *The 38th Annual*  
585 *ACM Symposium on User Interface Software and Technol-*  
586 *ogy*, page 14, New York, NY, USA, 2025. ACM. 3
- [51] SuperMap. Ai gis, 2025. 2
- [52] Eric K. Tokuda, Roberto M. Cesar, and Claudio T. Silva.  
590 Quantifying the presence of graffiti in urban environments.  
591 In *2019 IEEE International Conference on Big Data and*  
592 *Smart Computing (BigComp)*, pages 1–4, 2019. 3
- [53] U.S. Geological Survey. Earthexplorer. [https://](https://earthexplorer.usgs.gov/)  
594 [earthexplorer.usgs.gov/](https://earthexplorer.usgs.gov/), 2025. Query and order  
595 satellite images, aerial photographs, and cartographic prod-  
596 ucts. Provides access to over 40 years of Landsat data and  
597 various aerial photography collections. 3
- [54] Zeyu Wang, Koichi Ito, and Filip Biljecki. Assessing the  
599 equity and evolution of urban visual perceptual quality with  
600 time series street view imagery. 145:104704. 3
- [55] Zeyu Wang, Yingchao Jian, Adam Visokay, Don MacKen-  
602 zie, and Jon E. Froehlich. Street view for whom? an initial  
603 examination of google street view’s urban coverage and so-  
604 cioeconomic indicators in the us. Under review, 2025. Sub-  
605 mitted for review. 3
- [56] Chris Yoon, Ryan Louie, Jeremy Ryan, MinhKhang Vu,  
607 Hyegi Bang, William Derksen, and Paul Ruvolo. Leveraging  
608 augmented reality to create apps for people with visual dis-  
609 abilities: A case study in indoor navigation. In *Proceedings*  
610 *of the 21st International ACM SIGACCESS Conference on*  
611 *Computers and Accessibility*, page 210–221, New York, NY,  
612 USA, 2019. Association for Computing Machinery. 2, 3
- [57] Aziza Zhanabatyrova, Clayton Frederick Souza Leite, and  
614 Yu Xiao. Automatic map update using dashcam videos.  
615 *IEEE Internet of Things Journal*, 10(13):11825–11843,  
616 2023. 3
- [58] Guiming Zhang and A-Xing Zhu. The representativeness  
618 and spatial bias of volunteered geographic information: a re-  
619 view. *Annals of GIS*, 24(3):151–162, 2018. 3
- [59] Lefei Zhang and Liangpei Zhang. Artificial intelligence for  
621 remote sensing data analysis: A review of challenges and  
622 opportunities. *IEEE Geoscience and Remote Sensing Maga-*  
623 *zine*, 10(2):270–294, 2022. 3
- [60] Mengxia Zhang and Lan Luo. Can consumer-posted photos  
625 serve as a leading indicator of restaurant survival? evidence  
626 from yelp. *Management Science*, 69(1):25–50, 2023. 3
- [61] Meimei Zhang, Fang Chen, and Bin Li. Multistep question-  
628 driven visual question answering for remote sensing. *IEEE*  
629 *Transactions on Geoscience and Remote Sensing*, 61:1–12,  
630 2023. 2
- [62] Yuhang Zhao, Elizabeth Kupferstein, Brenda Veronica Cas-  
632 tro, Steven Feiner, and Shiri Azenkot. Designing ar visu-  
633 alizations to facilitate stair navigation for people with low  
634 vision. In *Proceedings of the 32nd Annual ACM Symposium*  
635 *on User Interface Software and Technology*, page 387–402,  
636 New York, NY, USA, 2019. Association for Computing Ma-  
637 chinery. 2, 3
- [63] Shengyuan Zou and Le Wang. Detecting individual aban-  
639 doned houses from google street view: A hierarchical deep  
640 learning approach. *ISPRS Journal of Photogrammetry and*  
641 *Remote Sensing*, 175:298–310, 2021. 3
- 642