

# GeoChain: Multimodal Chain-of-Thought for Geographic Reasoning

Sahiti Yerramilli<sup>\*1</sup>, Nilay Pande<sup>\*2</sup>, Rynaa Grover<sup>\*1</sup>, Jayant Sravan Tamarapalli<sup>\*1</sup>

<sup>1</sup>Google

<sup>2</sup>Waymo

sahitiy@google.com, nilayp@waymo.com, rynaa@google.com, jayantsravan@google.com

## Abstract

This paper introduces GeoChain, a large-scale benchmark for evaluating step-by-step geographic reasoning in multimodal large language models (MLLMs). Leveraging 1.46 million Mapillary street-level images, GeoChain pairs each image with a 21-step chain-of-thought (CoT) question sequence (over 30 million Q&A pairs). These sequences guide models from coarse attributes to fine-grained localization across four reasoning categories - visual, spatial, cultural, and precise geolocation - annotated by difficulty. Images are also enriched with semantic segmentation (150 classes) and a visual locatability score. Our benchmarking of contemporary MLLMs (GPT-4.1 variants, Claude 3.7, Gemini 2.5 variants) on a diverse 2,088-image subset reveals consistent challenges: models frequently exhibit weaknesses in visual grounding, display erratic reasoning, and struggle to achieve accurate localization, especially as the reasoning complexity escalates. GeoChain offers a robust diagnostic methodology, critical for fostering significant advancements in complex geographic reasoning within MLLMs.

**Code:** <https://github.com/sahitiy/geochain>

**Dataset:** <https://huggingface.co/datasets/sahitiy51/geochain>

## 1 Introduction

As large vision-language models (VLMs) continue to make rapid progress on general visual question answering and captioning tasks [Team *et al.*, 2024; OpenAI *et al.*, 2024; Wang *et al.*, 2024; Dai *et al.*, 2023], their capacity for structured geographic reasoning remains underexplored. The ability to infer a location from visual cues - such as terrain, signage, vehicles, or architecture - considered alongside spatial and cultural knowledge, is crucial for real-world applications like remote sensing, disaster response, and autonomous navigation. More broadly, geographic localization serves as a testbed for grounded intelligence, requiring models to reason over subtle visual features, incorporate world knowledge, and

disambiguate locations that may be visually similar. Despite this, existing benchmarks rarely probe the kind of step-by-step reasoning that such tasks demand.

We introduce GeoChain, a novel multimodal benchmark for evaluating structured geographic reasoning in large language models (MLLMs). As depicted in Figure 1, each GeoChain sample features a street-level image from the Mapillary dataset [Warburg *et al.*, 2020] paired with a 21-step chain-of-thought (CoT) question sequence. These sequences progressively guide models from coarse inferences, such as hemisphere or continent, to fine-grained predictions like city, latitude, and longitude. The complete GeoChain framework comprises 1.46 million images, each with this 21-question CoT structure, yielding over 30 million question-answer pairs. Questions span four core reasoning categories: visual cues, spatial localization, cultural inference, and precise geolocation, all annotated with difficulty levels for granular evaluation. This curriculum-style structure offers vital diagnostic insights into where and why models fail across reasoning stages, moving beyond sole reliance on final predictions.

To facilitate focused evaluations, we curated GeoChain Test-Mini, a diverse and challenging subset. This curation process leverages a locatability score, adapted from GeoReasoner [Li *et al.*, 2024] and computed using features from a pretrained MaskFormer model [Cheng *et al.*, 2021]. This score quantifies the visual identifiability of a location from a single image, allowing us to stratify GeoChain Test-Mini into Easy, Medium, and Hard tiers based on thresholds in the 0.12-0.6 range. The resulting GeoChain Test-Mini contains 2,088 carefully selected images, designed to offer a representative yet manageable scale for robust MLLM assessment.

Our main contributions are as follows:

- **GeoChain Benchmark Framework:** A novel benchmark for evaluating step-by-step MLLM geographic reasoning, derived from 1.46 million Mapillary street-level images and generating over 30 million Q&A pairs through 21-step chain-of-thought questions, all structured across diverse reasoning categories and difficulty levels.
- **Rich Augmentation & Curated Evaluation Set:** A methodology for enhancing images with semantic labels (150 classes) and a human-inspired locatability score for difficulty stratification, culminating in the GeoChain Test-

---

<sup>\*</sup>Equal Contribution.

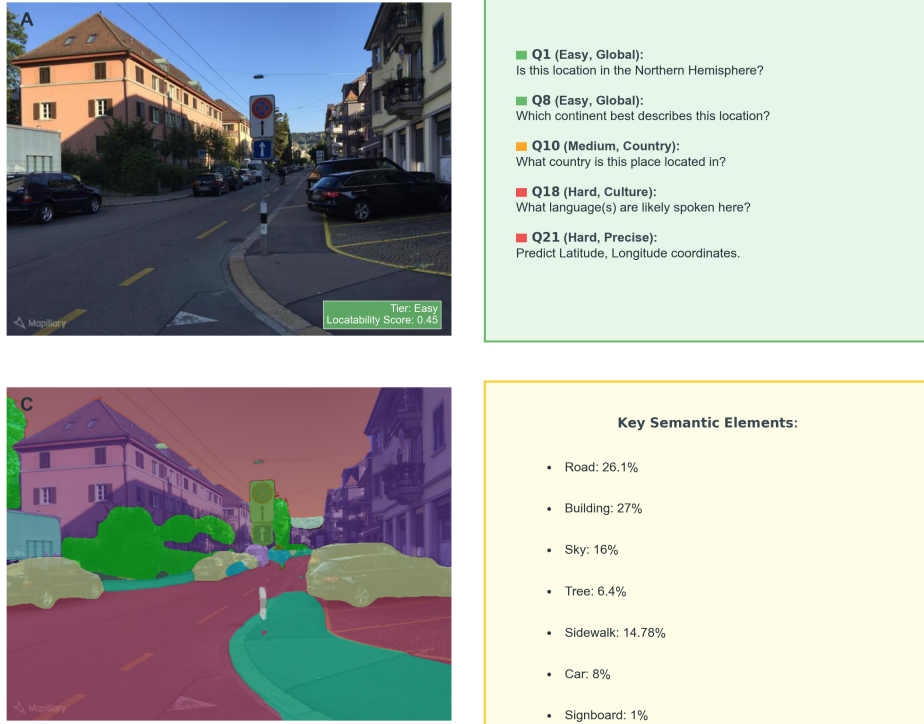


Figure 1: Components of a GeoChain instance: **(Top-Left)** Easy Mapillary Street-Level Sequences (MSLS) image with locatability score of 0.45. **(Top-Right)** Example chain-of-thought questions with difficulty indicators. **(Bottom-Left)** Derived semantic segmentation map. **(Bottom-Right)** Extracted key semantic labels. Together, these elements enable step-by-step diagnostic evaluation of geographic reasoning.

Mini: a quality-controlled 2088-image evaluation set; the resulting rich semantic metadata also offers a valuable resource for broader community research and future investigations.

- **Comprehensive MLLM Benchmarking & Analysis:** Evaluation of leading MLLMs on GeoChain Test-Mini, providing detailed insights into their geographic reasoning capabilities, performance variations, and common failure modes.

## 2 Related Work

### 2.1 Image-Based Geolocation

Early work in visual geolocation predominantly focused on matching query images to large, geotagged image databases, often aiming for direct coordinate prediction. For instance, Im2GPS [Hays and Efros, 2008] pioneered retrieving locations by comparing against a massive photograph dataset. Later, deep learning significantly advanced the field; PlaNet [Weyand *et al.*, 2016] utilized convolutional neural networks (CNNs) for global location prediction, and architectures like NetVLAD [Arandjelovic *et al.*, 2016] learned robust image representations for effective place recognition, improving upon earlier retrieval methods. Other approaches, such as those focusing on urban or cross-view settings [Tian *et al.*, 2017], further specialized these techniques. GeoChain diverges from these paradigms, which primarily target endpoint localization accuracy or image retrieval. Instead, it introduces

a structured multimodal reasoning benchmark where models must articulate a 21-step chain-of-thought (CoT) sequence of answers to geographically relevant questions, thereby enabling finer-grained diagnostic insight into their internal reasoning processes.

### 2.2 Multimodal Geographic Reasoning and Benchmarks

More recent efforts have begun to integrate visual understanding with language-based reasoning for complex geographic tasks. GeoReasoner [Li *et al.*, 2024], for example, introduced a fine-tuning strategy for MLLMs using human gameplay traces, primarily to improve final location prediction by modeling human-like inference. Similarly, other recent studies [Pramanik *et al.*, 2024; Yang *et al.*, 2024] also concentrate on predicting precise latitude and longitude. GeoComp [Song *et al.*, 2025] presents a large-scale dataset of geolocation gameplay data, emphasizing step-wise reasoning rooted in real human gameplay that often involves external metadata, active exploration, and dynamic information gathering. While these approaches offer valuable insights into human-like inference and gameplay dynamics, GeoChain’s contribution is complementary. It does not involve model fine-tuning or rely on gameplay trajectories. Instead, GeoChain employs a fully static, image-grounded evaluation framework: each sample consists of a single image paired with its fixed CoT question sequence, standardized across the entire benchmark. This design facilitates direct and controlled benchmarking of

different models’ inherent reasoning capabilities under uniform conditions, distinct from evaluating exploratory strategies or the ability to process dynamic data.

Other benchmarks, such as GAEA [Campos *et al.*, 2025], generate diverse conversational questions from detailed, place-specific metadata like OpenStreetMap attributes. While this can create rich contextual queries, it introduces challenges related to the temporal stability of dynamic data (e.g., changes in urban landscape) and complicates fair, apples-to-apples model comparisons due to non-uniform question sets. Consequently, disentangling model reasoning failures from idiosyncratic question characteristics becomes difficult. GeoChain mitigates these issues by grounding its standardized questions in more enduring visual semantics such as the presence of characteristic vegetation, architectural styles, or road infrastructure, often identifiable through image segmentation and stable general geographic facts. This focus ensures the evaluation centers on the consistency of the reasoning process itself.

Furthermore, existing geospatial benchmarks like GEO-Bench [Lacoste *et al.*, 2023] primarily target remote sensing applications, offering valuable tools for Earth monitoring with satellite imagery and tasks such as classification or segmentation. In contrast, GeoChain specifically addresses agent-level geographic reasoning from high-resolution, ground-level imagery, emphasizing natural-language understanding of spatial, cultural, and visual cues directly perceivable in such environments.

### 2.3 Mapillary Street-Level Sequences Dataset

GeoChain is built upon the Mapillary Street-Level Sequences (MSLS) dataset [Warburg *et al.*, 2020], a large-scale, crowd-sourced collection of diverse, geo-tagged street-level images. MSLS’s global coverage, with data from numerous cities worldwide reflecting the breadth of the MSLS ecosystem, and its varied capture conditions (diverse cameras, viewpoints, seasons, times of day) make it an ideal foundation for a benchmark aiming to evaluate generalizable geographic reasoning.

## 3 GeoChain Benchmark Construction

The GeoChain benchmark is constructed by augmenting the Mapillary Street-Level Sequences (MSLS) dataset [Warburg *et al.*, 2020]. MSLS provides a diverse collection of geo-tagged street-level imagery (approximately 1.4 million images in its full extent, with a geographical distribution across numerous cities as illustrated in Figure 2), crucial for developing and evaluating geographic localization models. However, to facilitate fine-grained, step-by-step reasoning, we introduce several layers of annotation and metadata. Our contributions enhance the MSLS dataset in three primary ways: semantic class labeling, locatability score computation, and the design of a structured chain-of-thought question battery. These augmentations, followed by a careful test set curation process, collectively enable a more nuanced evaluation of multimodal models’ geographic reasoning capabilities.

### 3.1 Semantic Class Labeling

To ground visual reasoning in explicit semantic content, each image in our benchmark is augmented with a semantic segmentation map. This map provides a detailed understanding of the scene’s composition by identifying various objects and environmental features. We employed MaskFormer [Cheng *et al.*, 2021], a state-of-the-art transformer-based architecture for semantic segmentation. Specifically, we utilized a MaskFormer model pre-trained on the ADE20K dataset [Zhou *et al.*, 2017], which offers a rich label set of 150 distinct classes, encompassing a wide array of objects, environmental elements (e.g., “tree”, “sky”, “road”), and architectural features (e.g., “building”, “window”, “door”).

From the segmentation map, we calculate how much of the image is covered by each category. We do this by working out the percentage of the image’s total area that each specific category takes up. For example, we might find that ‘sky’ covers 30% of an image, and ‘road’ covers 15%. This measurement of what’s in the scene, and how much of it there is, then helps us create the correct answers for many of the visual questions in our benchmark.

### 3.2 Locatability Score Computation

To systematically assess model performance across varying levels of visual ambiguity, we compute a *locatability score* for every image considered for the GeoChain benchmark. This score, ranging from 0 to 1, quantifies how visually identifiable a location is likely to be, with higher scores indicating more distinct and easily locatable scenes. Our methodology for calculating this score is adopted from [Li *et al.*, 2024]. The distribution of these computed locatability scores across the considered images is shown in Figure 3.

The core idea behind this score is to leverage common visual cues that humans, particularly proficient GeoGuessr players [GeoGuessr, 2013], rely on for geolocalization. The process involves several steps:

1. **Identification of Cues:** A set of cues frequently used by GeoGuessr players (e.g., “houses in central Chile are more likely to have terracotta tiled roofs”) is established.
2. **Cue-to-Class Similarity:** The semantic similarity between these cues and the 150 class labels produced by the MaskFormer model (as described in Section 3.1) is computed. This typically involves using text embeddings to represent both the cues and the class labels, followed by a similarity measure (e.g., cosine similarity).
3. **Class Weight Derivation:** The similarities are aggregated across all cues for each class and then subjected to min-max normalization to derive a set of weights  $w_c$  for each class  $c$ . These weights reflect the importance of each visual class for geolocalization.
4. **Weighted Score Aggregation:** The final locatability score for an image is computed as a weighted sum of the percentage areas of the classes present in the image.

This locatability score is then used to stratify the images within the GeoChain test set into three distinct tiers: Easy, Medium, and Hard. This stratification, based purely on visual cues inherent in the imagery, allows for a more granular analysis of model performance and helps identify specific

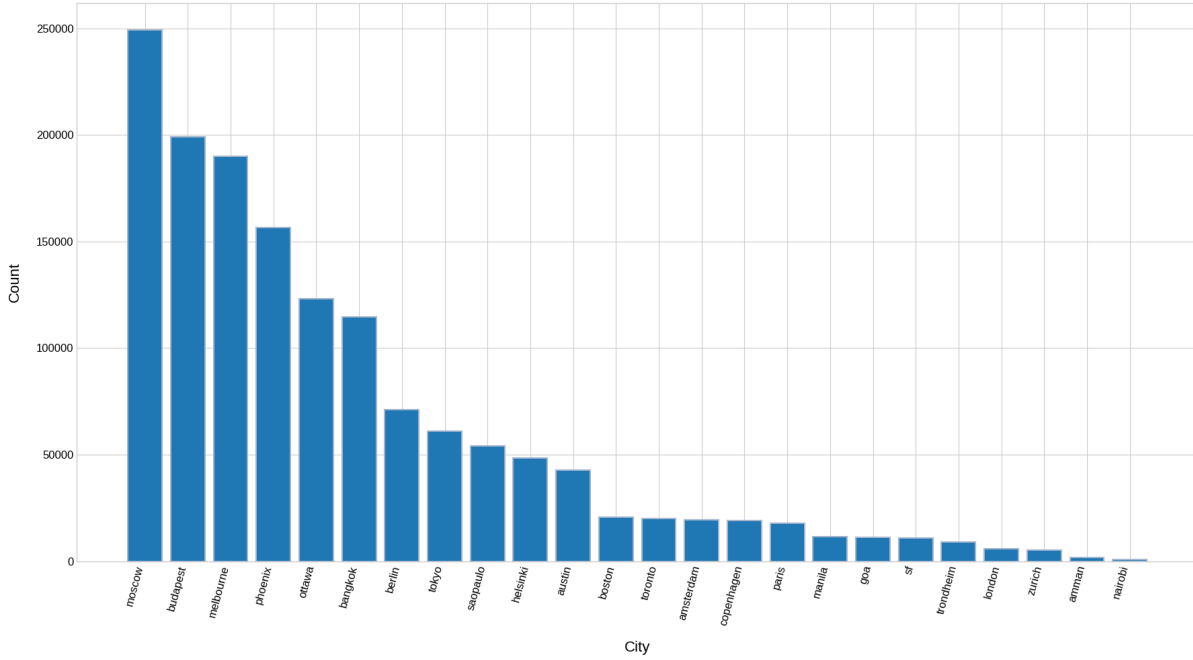


Figure 2: Count of images per city, illustrating the city distribution within the GeoChain dataset.

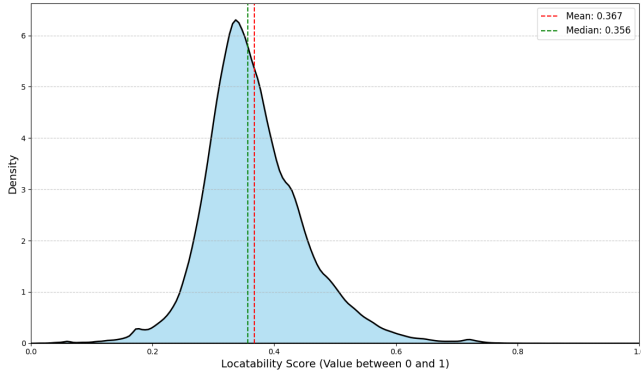


Figure 3: Distribution of Locatability Scores in GeoChain

weaknesses in reasoning about visually challenging environments.

### 3.3 Chain-of-Thought Question Design

A central component of GeoChain is a carefully designed sequence of 21 questions that guide the model through a step-by-step reasoning process, from coarse-grained observations to fine-grained localization. This chain-of-thought (CoT) approach aims to mimic a structured human-like deduction process. The questions are ordered such that earlier questions elicit information or focus attention on attributes that can be instrumental in answering subsequent, more complex questions.

The full list of 21 questions, along with their rank, assigned difficulty (Easy, Medium, Hard), question type (Binary, Multiclass, Free-text), and question category (e.g., Cul-

ture/Infrastructure, Geo Localization, Terrain/Environment), is provided in Appendix A.1. The difficulty annotation (Easy, Medium, Hard) for each question reflects the anticipated challenge of answering that specific question in isolation, based on the type of information required.

The question set is designed to be static across all data points in the benchmark. This uniformity ensures a consistent evaluation framework, allowing for direct, apples-to-apples comparisons of different models’ reasoning capabilities. The questions cover diverse aspects:

- **Visual Object/Attribute Presence:** Some questions directly query the presence of specific objects or attributes identifiable from the image (e.g., “Do you see any boats or ships?”). Ground truth for these questions is primarily derived from the semantic class labels extracted via the MaskFormer model (Section 3.1). For instance, if the class “boat” occupies a non-zero percentage of the image, the answer would be affirmative.
- **Inferential and Contextual Knowledge:** Other questions require more derivative reasoning or contextual knowledge beyond direct object identification (e.g., “Is this place near a coast?”, “What side of the road do vehicles drive on here?”). The MSLS dataset encompasses images from 24 distinct cities globally. For images originating from these locations, we manually curated ground truth answers for city-level attributes or environmental characteristics (e.g., predominant architectural styles, typical climate indicators) that apply broadly to the image’s geographic area.
- **Progressive Localization:** The sequence progresses from general observations (e.g., hemisphere, continent) to specific details (e.g., country, city, precise latitude and longitude coordinates).

The question types include binary (Yes/No), multiclass (selection from a predefined set of options), and free-text (open-ended answers, such as country name or coordinates). This variety tests different aspects of a model’s understanding and generation capabilities.

The semantic segmentation labels generated in Section 3.1 were instrumental in constructing several questions that directly probe the visual understanding capabilities of the models. Beyond their use in the current benchmark, this rich semantic metadata, now part of GeoChain, offers a valuable resource for the community. It can be leveraged to design new questions aimed at further investigating specific aspects of model behavior, such as tendencies towards visual hallucination [Li *et al.*, 2023] [Rohrbach *et al.*, 2018] or the fine-grained ability to identify a wider array of objects. The insights derived from such extended evaluations can subsequently guide targeted improvements in model development.

By analyzing model performance across this structured chain of questions, GeoChain aims to provide deeper insights into the strengths and weaknesses of multimodal geographic reasoning systems.

### 3.4 Test Set Curation and Sampling Strategy

To create the “GeoChain Test-Mini” subset for focused evaluation, we prioritized stratification, visual quality, and diversity. We initially targeted 2100 images, stratified by locatability scores into 700 Easy, 700 Medium, and 700 Hard examples. A tiered, unique-sequence sampling strategy was employed: unique image sequences were randomly sampled first for the Hard tier, then for the Medium tier (from remaining unique sequences), and finally for the Easy tier, ensuring no sequence was reused across tiers. The underlying MSLS dataset exhibits a notable skew in its per-city image distribution (as highlighted by the overall dataset statistics in Figure 2). Consequently, to avoid introducing new biases that could arise from attempting to manually balance city representation or ‘carefully’ over/under-sample from specific locations, our approach was to randomly sample unique image sequences across all available cities within each defined locatability tier. These 2100 candidates underwent manual visual inspection, where 12 images with critical quality issues (e.g., excessive blur, poor exposure) were removed. This rigorous curation yielded a final Test-Mini set of 2088 high-quality, diverse, and appropriately challenging images.

## 4 Analysis

In this section, we evaluate the performance of frontier vision-language models: GPT-4.1, GPT-4.1-mini [OpenAI *et al.*, 2024], Claude 3.7 Sonnet [Sonnet, 2025], Gemini 2.5 Flash [Google, 2025a] and Gemini 2.5 Pro [Google, 2025b] on the GeoChain “Test-Mini” benchmark, focusing on their ability to reason accurately and consistently across a structured 21-step geographic reasoning chain.

### 4.1 Evaluation Metrics

#### Haversine Distance

The final question in each GeoChain sequence (Question 21) requires the model to predict the geographic coordinates (lat-

itude, longitude) of the depicted scene. To evaluate the accuracy of these specific predictions, we compute the *Haversine distance*: the shortest distance over the Earth’s surface between the predicted and ground-truth coordinates, assuming a spherical Earth. A detailed explanation of the Haversine distance calculation is provided in Section A.2.

#### Pass Score

The **Pass Score** is computed as the average fraction of questions correctly answered across the full 21-step reasoning chain for each image. A prediction for any question is considered correct if it matches the ground-truth answer for that specific question, accounting for its type (e.g., exact match for free-text, class match for multiclass, or binary match). Crucially, for the final latitude and longitude prediction (Question 21), a response is deemed correct contributing to the Pass Score if its Haversine distance (as defined in Section 4.1) from the ground truth is less than 50km.

### 4.2 Overall Model Performance

Overall model performance (Table 1) offers nuanced insights into current MLLM geographic reasoning. The leading Gemini models exhibit specialized strengths: Gemini-2.5-pro excels in complex multi-step reasoning (pass score 81.84%), whereas Gemini-2.5-Flash achieves superior localization precision (445.24 km mean error), hinting at differing architectural or training optimizations. This divergence underscores that broad inferential ability and precise geolocalization are distinct skills, likely requiring separate optimization pathways rather than being monolithic capabilities. GPT-4.1 maintains a competitive position; however, the substantial localization inaccuracies of Claude 3.7 Sonnet (1289.04 km error) and GPT-4.1 Mini (1194.77 km error) underscore that robust geospatial grounding is a significant developmental hurdle, indicating a key area for advancement in MLLM capabilities.

The introduction of threshold-based localization accuracies - at City (< 25 km), Region (< 200 km), and Country (< 750 km) levels further refines this performance landscape. Gemini-2.5-pro’s superior performance is reinforced by its top-tier City-level precision (59.38%). Complementing this, Gemini-2.5-Flash excels in broader accuracy, leading at both Region-level (70.02%) and Country-level (90.31%). GPT-4.1 also demonstrates notable strength in City-level performance (57.84%), surpassing Gemini-2.5-Flash in this specific high-precision context. Conversely, Claude 3.7 Sonnet’s previously noted localization challenges are starkly emphasized by its profound difficulties at these finer scales (e.g., 40.34% at City-level), performing below GPT-4.1 Mini (48.61% at City-level) here. These granular metrics effectively highlight that achieving reliable, high-confidence City-level precision is a primary differentiator and a significant challenge across the evaluated MLLMs.

### 4.3 Breakdown by Image Difficulty

Analyzing model performance by image difficulty (Table 2) reveals critical operational characteristics. As expected, ‘Hard’ images significantly challenge all models, leading to substantial increases in mean localization errors often exceeding 1000-2000 km for several models. The Gemini models

Table 1: Overall model-level accuracy and localization metrics.

Model	Pass Score (%)	Mean Dist (km)	< 25 km (%)	< 200 km (%)	< 750 km (%)
Gemini-2.5-pro	<b>81.84</b>	489.51	<b>59.38</b>	69.95	88.51
Gemini-2.5-Flash	79.77	<b>445.24</b>	55.71	<b>70.02</b>	<b>90.31</b>
GPT-4.1	79.25	611.89	57.84	67.36	86.24
Claude 3.7 Sonnet	76.23	1289.04	40.34	47.07	73.31
GPT-4.1 Mini	70.42	1194.77	48.61	52.87	72.77

Table 2: Performance by image difficulty. Accuracy (%) and Haversine distance (km) for each difficulty level.

Model	Diff	Pass Score	M. Dist.
Claude 3.7 Sonnet	Easy	77.2	885.86
	Medium	78.3	989.13
	Hard	73.2	2000.14
GPT-4.1 Mini	Easy	70.8	863.19
	Medium	73.2	827.78
	Hard	67.3	1910.44
GPT-4.1	Easy	79.3	357.36
	Medium	81.6	428.46
	Hard	76.8	1052.13
Gemini-2.5 Flash	Easy	80.5	<b>287.61</b>
	Medium	82.5	<b>188.45</b>
	Hard	76.3	873.78
Gemini-2.5 Pro	Easy	<b>83.3</b>	300.29
	Medium	<b>84.2</b>	304.32
	Hard	<b>78.0</b>	<b>866.62</b>

consistently lead: Gemini-2.5-pro achieves top Pass Scores across all difficulties (e.g., 78.0% on Hard), while Gemini-2.5-Flash generally provides superior localization on 'Easy' and 'Medium' images (e.g., 188.45 km on Medium). Notably, Gemini-2.5-pro performs the best for localization precision on 'Hard' images (866.62 km), possibly where its stronger inferential capacity becomes decisive. An intriguing anomaly is the better localization by some models, like Gemini-2.5-Flash, on 'Medium' versus 'Easy' images, potentially due to bias towards certain cities in pre-training data. Furthermore, Claude 3.7 Sonnet's performance is particularly interesting: despite reasonable Pass Scores (e.g., 73.2% on Hard), its poor localization (2000.14 km on Hard) highlights a profound disconnect between understanding cues and grounding them spatially.

#### 4.4 Breakdown by Question Category

Analyzing Pass Scores by question category (Table 3), informed by the benchmark's diverse question structures (e.g., visual queries versus free-text specific knowledge), reveals distinct performance strata. Foundational "Visual" questions, focusing on direct object presence (e.g., "Do you see any boats?"), yield universally high scores (all models >91.8%), suggesting robust basic visual grounding and low immediate hallucination, with Claude 3.7 Sonnet leading (92.8%).

Similarly, "Terrain" identification is generally strong. In contrast, categories like "Geo Localization" and "Cultural" show mixed results; models likely handle simpler, coarse queries (e.g., continent identification) better than challenging free-text questions requiring specific knowledge (e.g., city/state names, language identification). Unsurprisingly, "Exact Loc" demanding precise latitude/longitude output—is definitively the most challenging category across all models. Within this landscape, Gemini-2.5-pro consistently excels, particularly in the more demanding categories like "Terrain" (87.4%), "Cultural" (77.9%), and "Exact Loc." (63.5%). GPT-4.1 also demonstrates strong performance, notably in "Geo Localization" (76.9%) and "Exact Loc." (61.5%). Claude 3.7 Sonnet's profile, with its excellent "Visual" scores but significantly weaker "Exact Loc." performance (51.0%), starkly illustrates a common theme: a disconnect between initial cue processing and final, precise geospatial grounding, which remains the primary MLLM hurdle.

#### 4.5 Breakdown by City

A city-level view (Fig. 4) shows that performance is far from uniform:

Gemini-2.5-pro is the most stable, topping the leaderboard in 20 / 24 cities and exceeding 85% accuracy in visually distinctive urban centres such as Tokyo, Zurich and Toronto. Gemini-2.5-Flash and GPT-4.1 follow closely, maintaining more than **75%** accuracy in most regions. Performance on Claude 3.7 Sonnet and GPT 4.1 Mini fluctuate sharply: they perform competitively in cue-rich European cities (Paris, Berlin) but collapse in visually ambiguous locales (Nairobi, São Paulo, Amman). Mean Haversine error (Fig. 6) confirms the pattern: Gemini-2.5-pro keeps errors below 300 km in nearly every city, whereas Claude and GPT 4.1 Mini exceed 1000 km in several cases (Helsinki, Melbourne, São Paulo).

These results highlight how regional factors such as vegetation, signage language, traffic orientation and architectural style strongly modulate geolocation accuracy.

## 5 Conclusion

This paper introduced GeoChain, a large-scale, chain-of-thought benchmark designed to dissect multimodal geographic reasoning in MLLMs using street-level imagery and a 21-step diagnostic framework. Our evaluations on the curated GeoChain Test-Mini subset reveal that even leading MLLMs exhibit significant deficiencies in visual grounding, reasoning consistency, and localization accuracy, particularly as task and visual complexity escalate. By enabling a granular, step-by-step analysis, GeoChain moves beyond simple end-task

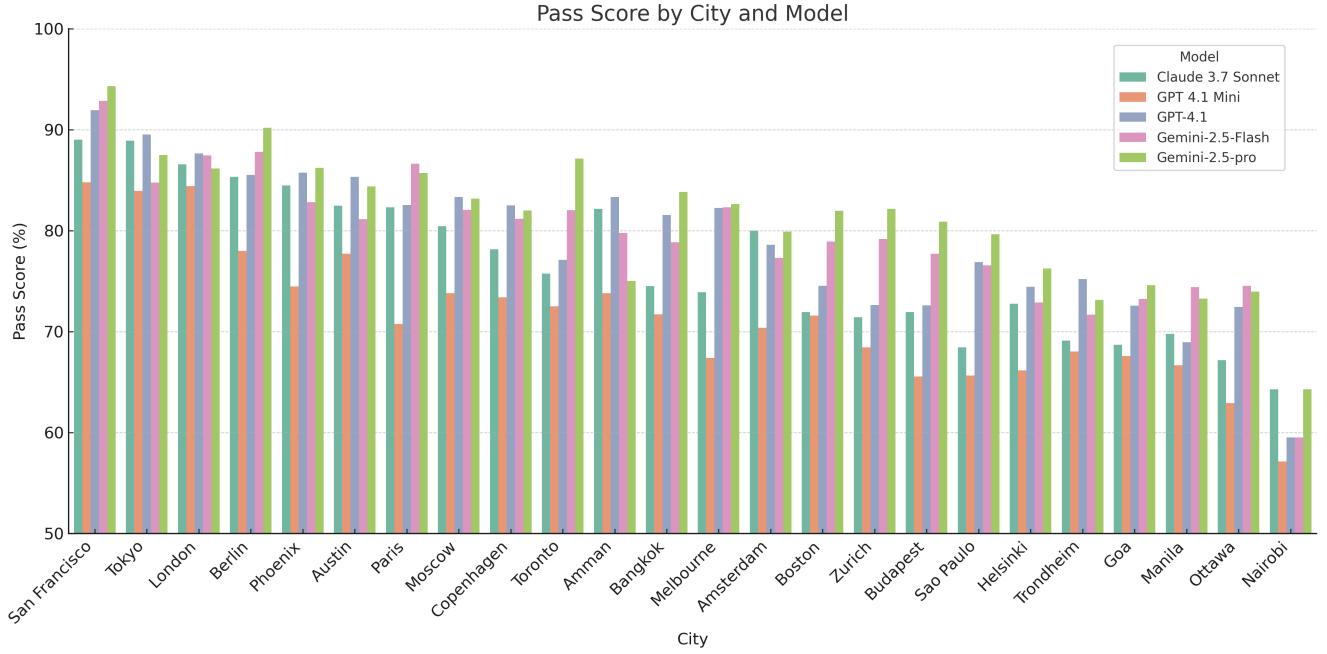


Figure 4: Pass score (%) by city, highlighting the influence of geographical location on model accuracy.

Table 3: Pass score (%) by question category.

Model	Visual	Terrain	Geo Localization	Cultural	Exact Loc.
Claude 3.7 Sonnet	<b>92.8</b>	84.7	69.4	67.4	51.0
GPT-4.1 Mini	92.3	78.7	64.1	56.8	40.7
GPT-4.1	91.8	84.8	<b>76.9</b>	68.3	61.5
Gemini-2.5-Flash	92.4	86.0	73.5	75.3	59.8
Gemini-2.5-pro	92.1	<b>87.4</b>	76.8	<b>77.9</b>	<b>63.5</b>

accuracy to pinpoint these critical failure modes, thereby providing an essential diagnostic resource and methodology. We anticipate that GeoChain will steer future research towards developing more robust, geographically aware, and reliable AI systems capable of nuanced real-world understanding.

## 6 Limitations

While GeoChain offers a novel diagnostic approach, we acknowledge several limitations. GeoChain is built upon the Mapillary Street-Level Sequences training split; consequently, while our chain-of-thought reasoning framework and the overall task are novel, there is a potential that MLLMs have encountered these specific visual scenes or highly similar ones during their extensive pre-training. Evaluating performance on truly "unseen" street-level imagery is an inherent challenge for the field, given the ubiquity of data from sources like Google Street View [Anguelov *et al.*, 2010] and OpenStreetMap [Haklay and Weber, 2008], meaning that performance assessments may partly reflect familiarity with certain visual data rather than solely generalization to entirely new scenes. Additionally, the underlying geographical distribution of the images, though diverse, retains some skew,

potentially affecting the generalizability of the findings in all urban contexts. Furthermore, our locatability score’s precision is contingent upon the accuracy of an upstream semantic segmentation model, which could introduce noise into the difficulty stratification.

## 7 Usage of Generative AI tools

We utilized Generative AI tools to help improve the language, phrasing, and readability of this manuscript.

## References

- [Anguelov *et al.*, 2010] Dragomir Anguelov, Carole Dulong, Daniel Filip, Christian Frueh, Stéphane Lafon, Richard Lyon, Abhijit Ogale, Luc Vincent, and Josh Weaver. Google street view: Capturing the world at street level. *Computer*, 43(6):32–38, 2010.
- [Arandjelovic *et al.*, 2016] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. NetVLAD: CNN architecture for weakly supervised place recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5287–5297, 2016.



- [Campos *et al.*, 2025] Ron Campos, Ashmal Vayani, Parth Parag Kulkarni, Rohit Gupta, Aritra Dutta, and Mubarak Shah. GAEA: A geolocation aware conversational model, 2025.
- [Cheng *et al.*, 2021] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, pages 17864–17875. Curran Associates, Inc., 2021.
- [Dai *et al.*, 2023] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- [GeoGuessr, 2013] GeoGuessr. Geoguessr - let’s explore the world!, 2013. Accessed: [Date you accessed the website].
- [Google, 2025a] Google. Gemini 2.5 flash is now in preview. <https://blog.google/products/gemini/gemini-2-5-flash-preview/>, 2025. Accessed 2025-05-18.
- [Google, 2025b] Google. Gemini 2.5 pro model. <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/gemini-2-5-thinking>, 2025. Accessed through the Gemini API on [Date of access] or Vertex AI.
- [Haklay and Weber, 2008] Muki Haklay and Patrick Weber. Openstreetmap: User-generated street maps. *IEEE Pervasive Computing*, 7(4):12–18, 2008.
- [Hays and Efros, 2008] James Hays and Alexei A. Efros. IM2GPS: estimating geographic information from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8. IEEE, June 2008.
- [Lacoste *et al.*, 2023] Alexandre Lacoste, Nils Lehmann, Pau Rodriguez, Evan David Sherwin, Hannah Kerner, Björn Lütjens, Jeremy Andrew Irvin, David Dao, Hamed Alemohammad, Alexandre Drouin, Mehmet Gunturkun, Gabriel Huang, David Vazquez, Dava Newman, Yoshua Bengio, Stefano Ermon, and Xiao Xiang Zhu. Geo-bench: Toward foundation models for earth monitoring, 2023.
- [Li *et al.*, 2023] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [Li *et al.*, 2024] Ling Li, Yu Ye, Bingchuan Jiang, and Wei Zeng. Georeasoner: Geo-localization with reasoning in street views using a large vision-language model. In *International Conference on Machine Learning (ICML)*, 2024.
- [OpenAI *et al.*, 2024] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rameev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selman, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Mad-



die Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

[Pramanik *et al.*, 2024] Saksham Pramanik, Aayush Mundra, Ashutosh Mittal, Sreyas Mohan, Saim Wani, Shramay S Vernekar, Pranav M Dixit, Shanti Priya, Ankur Beniwal, Ojaswa Sharma, and Senthil Mani. Evaluating precise geolocation inference capabilities of vision language models, 2024.

[Rohrbach *et al.*, 2018] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.

[Song *et al.*, 2025] Zirui Song, Jingpu Yang, Yuan Huang, Jonathan Tonglet, Zeyu Zhang, Tao Cheng, Meng Fang, Iryna Gurevych, and Xiuying Chen. Geolocation with real human gameplay data: A large-scale dataset and human-like reasoning framework, 2025.

[Sonnet, 2025] Claude 3.7 Sonnet. Claude 3.7 sonnet documentation. <https://docs.anthropic.com/en/docs/overview>, 2025. Anthropic AI Model.

[Team *et al.*, 2024] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal, Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry Lepikhin, Michela Paganini,

Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vodrahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gwoon Chen, Chris Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek, Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen, Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes Dotiwalla, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Bat-saikhan, Mateo Wirth, James Qin, Ivo Danihelka, Tulse Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, Madhu Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz, Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, Quan Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang, Johan Schalkwyk, Christina Butterfield, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Marvin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna, Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, HyunJeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt Harvey, Sean Sechrist, Meire Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang, James Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys, Kate

Baumli, Connie Tao, Rajkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen, Nemanja Rakićević, Dominik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayanan Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert Webson, Hyo Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su, Jake Walker, Jane Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit Bagri, Arnar Mar Hrafnkelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou, Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, Ruibo Liu, Tara Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeynep Cankara, Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihla, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey, Anastasija Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty, Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rustemi, Pen Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause, Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur, Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, Duc Dung Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech, Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma,

Clara Huiyi Hu, Aurko Roy, Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng, Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams Yu, Sarah Hodgkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Bloniarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sherman, Jay Pavagadhi, Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qiuqia Li, Anton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, Alberto Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb Arnold, Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin Polacek, Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao, Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Rubenstein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aishwarya Kamath, Ted Klimenko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis Daun, Khushen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Felix de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng, Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwanicki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srinivasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Garrette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki Iturrate, Michael B. Chang, Jackie Xi-ang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty, Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed,

Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong, Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirschall, Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhschskaya, Ashwin Sreevatsa, Shuang Song, Luis C. Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal, Christel Ngani, Kati Goshvadi, Rebecca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen, Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang Luong, Boyu Wang, Sujeewan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jacyln Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao, Adnan Ozturk, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Milad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Vilella, Haibin Zhang, Harry Richardson, James Martens, Matko Bosnjak, Shreyas Ramamohan Belle, Jeff Seibert, Mahmoud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills, Joseph Pagadora, Tsendsuren Munkhdalai, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon Qureshi,

Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnappalli, Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Bartek Perz, Wooyeol Kim, Nandita Dukkupati, Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder, Chester Kwak, Matt Miecznikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer, Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, Austin Matthews, CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, Shereen Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadisy, Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang, Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh Keshava, Shubham Agrawal, Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Pettrini, DongHyun Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Marcus Wain-

wright, Mario Cortes, Frederick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa, François-Xavier Aubet, Anton Algymr, Dan Banica, Agoston Weisz, Andras Orban, Alexandre Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al Meray, Martin Baeuml, Trevor Strohmman, Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Koray Kavukcuoglu, Jeff Dean, and Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024.

[Tian *et al.*, 2017] Yicong Tian, Chen Chen, and Mubarak Shah. Cross-view image matching for geo-localization in urban environments. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1998–2006, 2017.

[Wang *et al.*, 2024] 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 Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution, 2024.

[Warburg *et al.*, 2020] Frederik Warburg, Søren Hauberg, Gregory D. D. Funke, and Yoko Yuki. Mapillary street-level sequences: A dataset for lifelong place recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 931–940. IEEE, June 2020.

[Weyand *et al.*, 2016] Tobias Weyand, Ilya Kostrikov, and James Philbin. PlaNet - photo geolocation with convolutional neural networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016, Proceedings, Part II*, volume 9906 of *Lecture Notes in Computer Science*, pages 37–55. Springer, 2016.

[Yang *et al.*, 2024] Ruixiang Yang, Cheng Zhang, Lingxi Meng, He Wang, Xiaoyan Li, Yuke Li, Shuo Wang, Hao-ran Wei, Yiyang Li, Wentao Qu, Pengchuan Zhang, Jiazheng Xu, Bihan Wen, Diyi Yang, Kangkang Lu, Saurabh Gupta, Guanzhong Wang, Zhiqiang Shen, Baining Guo, Ruoming Jin, Song-Chun Zhu, and Hongjing Lu. VLMs as GeoGuessr masters—exceptional performance, hidden biases, and privacy risks, 2024.

[Zhou *et al.*, 2017] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ADE20K dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 633–641. IEEE, July 2017.

## A Appendix

### A.1 Implementation Details

#### Questions

This section details the complete 21-question sequence (Table 4) that forms the core of the GeoChain benchmark, designed to evaluate the step-by-step geographic reasoning capabilities of Multimodal Large Language Models (MLLMs). Each

question in the sequence is characterized by its rank, designated difficulty level (Easy, Medium, or Hard), expected response format (Binary, Multiclass, or Free-text), and its primary Question Category (Visual Cues, Geographical localization, Culture/Infrastructure, Terrain/Environment, or Exact Location). This comprehensive listing provides a transparent foundation for understanding the specific tasks underpinning the performance evaluations discussed throughout this paper.

#### System Prompt

To guide the Multimodal Large Language Models (MLLMs) and standardize their responses for the GeoChain benchmark tasks, the following system prompt was consistently employed:

#### System Prompt

You are an accurate geolocation model. Given the image, answer the following questions in order. Please provide your best guess. Each question is also provided with question type. For Binary questions, answer Yes/No only. For Multiclass questions, answer as one of the provided options in brackets. Final question type is a free text question, answer it as a free string text. If you are not sure about the answer, give your best guess. Answer format should be a json dict with question indices as keys (0 indexed) and values as Answer: <answer>, Reasoning: <reasoning>.

#### Tools and Infrastructure

The execution of model inference was managed by Prompt-foo<sup>1</sup>, a platform that ensures reproducibility in benchmarking by offering versatile prompt configuration and effective API linkage. We used the transformers library in Hugging Face<sup>2</sup>; to run the MaskFormer model for computing segmentation masks. These calculations were performed on an NVIDIA GeForce RTX 3060 graphics processing unit.

### A.2 Haversine Distance

Haversine Distance the shortest distance over the Earth’s surface between the predicted and ground-truth coordinates, assuming a spherical Earth.

The Haversine formula is given by:

$$\Delta\phi = \phi_2 - \phi_1$$

$$\Delta\lambda = \lambda_2 - \lambda_1$$

$$a = \sin^2\left(\frac{\Delta\phi}{2}\right) + \cos(\phi_1)\cos(\phi_2)\sin^2\left(\frac{\Delta\lambda}{2}\right)$$

$$d = 2R \cdot \arcsin(\sqrt{a})$$

Here,  $d$  is the Haversine distance between two points  $(\phi_1, \lambda_1)$  and  $(\phi_2, \lambda_2)$ . This metric provides an interpretable and robust way to measure geographic prediction error.

Table 4: The GeoChain 21-Step Benchmark Question Set.

Rank	Difficulty	Question	Question Type	Question Category
1	Easy	Do you see any boats or ships?	Binary	Visual Cues
2	Easy	Do you see one or more of the following vehicles: Bus, Truck, Car, Van, Motorbike, Minibike, Bicycle?	Binary	Visual Cues
3	Easy	Can you see any traffic lights?	Binary	Visual Cues
4	Easy	Can you see any flag?	Binary	Visual Cues
5	Easy	Would you say this location is near the Equator?	Binary	Geographical localization
6	Easy	Does this location seem to be close to the Poles?	Binary	Geographical localization
7	Easy	Is this place located in the Northern Hemisphere?	Binary	Geographical localization
8	Easy	Which continent best describes where this location is? (7 continents: North America/South America/Europe/Africa/Asia/Oceania/Antarctica)	Multiclass	Geographical localization
9	Medium	What side of the road do vehicles drive on here? (Left/Right)	Multiclass	Culture/Infrastructure
10	Medium	What country is this place located in?	Free-text	Geographical localization
11	Medium	Is this place near coast?	Binary	Terrain/Environment
12	Medium	Does this location appear to be an island?	Binary	Terrain/Environment
13	Easy	Is this place located in a desert region?	Binary	Terrain/Environment
14	Easy	Does this location seem to be in a mountainous or hilly region?	Binary	Terrain/Environment
15	Medium	What is the most likely climate type for this location? (5 main climate types: Tropical/Dry/Temperate/Continental/Polar)	Multiclass	Terrain/Environment
16	Easy	Does this place look like a big city?	Binary	Culture/Infrastructure
17	Medium	Would you classify this place as a small town?	Binary	Culture/Infrastructure
18	Hard	What language(s) are most likely spoken at this place?	Free-text	Culture/Infrastructure
19	Hard	Can you name the state or province this place belongs to?	Free-text	Geographical localization
20	Hard	What is the name of the city, town, or village seen here?	Free-text	Geographical localization
21	Hard	Based on everything observed, what are the latitude and longitude coordinates of this place? Please give a tuple of float coordinates (lat, lon)	Free-text	Exact Location

Table 5: Pass score (%) across question difficulty and image difficulty. Each row shows performance on a given question difficulty across images of increasing ambiguity.

Model	Question Difficulty	Easy Images	Medium Images	Hard Images
Claude 3.7 Sonnet	Easy	89.3	89.1	87.7
	Medium	76.0	75.7	72.0
	Hard	45.8	52.7	34.8
GPT 4.1 Mini	Easy	86.7	86.7	84.2
	Medium	66.1	67.3	62.1
	Hard	37.4	44.6	27.9
GPT-4.1	Easy	87.9	87.5	84.3
	Medium	75.5	76.5	71.8
	Hard	51.8	60.0	43.3
Gemini-2.5-Flash	Easy	90.7	91.0	89.4
	Medium	76.7	77.8	74.2
	Hard	47.3	54.9	38.4
Gemini-2.5-pro	Easy	<b>91.6</b>	<b>91.3</b>	<b>89.8</b>
	Medium	<b>78.2</b>	<b>79.9</b>	<b>75.7</b>
	Hard	<b>52.4</b>	<b>61.6</b>	<b>45.9</b>

### A.3 Additional Analysis

#### Image Difficulty vs Question Difficulty Interaction

To analyze how visual and reasoning difficulty interact, we compute a two-dimensional pass rate matrix over **question difficulty** (Easy, Medium, Hard) and **image difficulty** (Easy, Medium, Hard). Table 5 presents this breakdown for each model.

We observe a consistent trend across all models: accuracy declines with both increasing *image* difficulty and *question* difficulty. Importantly, hard questions on hard images represent the most challenging setting, with pass rates often below 40%—even for state-of-the-art models.

Gemini-2.5-pro shows the strongest resilience across the board, maintaining high scores even on hard questions in ambiguous scenes. In contrast, Claude 3.7 Sonnet and GPT 4.1 Mini exhibit large drops in performance under compounding difficulty, confirming their brittleness in multi-factor reasoning.

This matrix allows us to quantify model *sensitivity to visual ambiguity* and pinpoint failure modes. For example, a model that performs well on hard questions from easy images but poorly on the same questions from hard images may lack robustness in interpreting noisy visual cues. Conversely, a model that fails uniformly on hard questions indicates weaknesses in logical chaining or symbolic inference. Together, this analysis emphasizes the need for benchmarks that probe cross-modal interactions, rather than evaluating visual or linguistic difficulty in isolation.

#### Breakdown by Question Difficulty

To better understand how models handle increasing reasoning complexity, we group questions by their annotated difficulty levels: **Easy**, **Medium**, and **Hard**. These difficulty tags were

assigned manually based on the subtlety, required external knowledge, and ambiguity of each question.

Table 6: Pass score (%) by question difficulty.

Model	Easy	Medium	Hard
Claude 3.7 Sonnet	88.7	74.6	44.5
GPT 4.1 Mini	85.9	65.2	33.4
GPT-4.1	87.3	75.8	54.7
Gemini-2.5-Flash	90.8	76.2	51.3
Gemini-2.5-pro	<b>91.1</b>	<b>78.4</b>	<b>55.1</b>

Across all models, accuracy decreases consistently with question difficulty. Gemini-2.5-pro achieves the highest pass rates at all levels, followed closely by Gemini-2.5-Flash and GPT-4.1. Interestingly, Claude 3.7 Sonnet and GPT 4.1 Mini both exhibit sharp drops on hard questions, with performance falling below 45% and 35%, respectively.

These findings suggest that while many models can answer surface-level geographic questions accurately, their reasoning falters as complexity increases especially when fine-grained localization or symbolic inference is required. The relatively better performance of Gemini-2.5-pro on hard questions indicates more stable multi-hop reasoning or greater robustness to subtle visual signals.

#### Accuracy vs. Reasoning Depth

Figure 5 reveals a typical degradation pattern: All models perform well in the initial questions (1–9), which ask about visual or global cues such as vehicles, hemisphere, or continent. These are relatively easy to infer on the basis of surface-level features.

As the questions become more complex and semantically demanding, the accuracy drops sharply, especially at questions 10 and 17. These questions requiring nuanced interpre-

<sup>1</sup><https://www.promptfoo.dev>

<sup>2</sup><https://huggingface.co/docs/transformers/en/index>

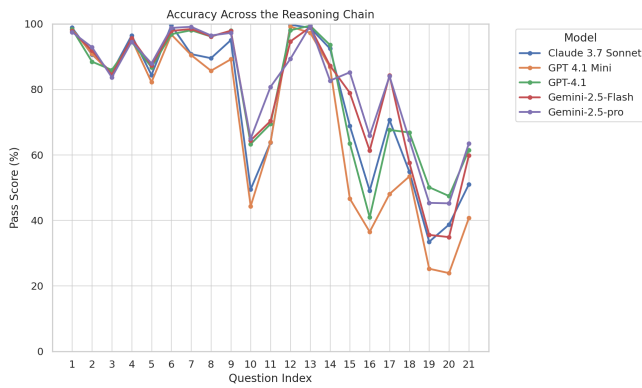


Figure 5: Average pass score across the 21-step Geochain reasoning chain. Accuracy decreases as questions progress from coarse global inference to fine-grained localization.

tation of environmental and infrastructure signals.

In particular, we observe a performance bump around questions 12–14. Despite appearing later in the sequence, these questions ask about relatively easy visual features (e.g., desert, hills, or city size). This reinforces the value of structuring questions not just by logical sequence but also by measured difficulty, allowing finer-grained diagnostics of model capability.

The final steps of the chain (questions 18–21) see the steepest drop in performance, as models are asked to predict language, administrative region, city name, and exact coordinates - tasks that require multi-modal reasoning, robust world knowledge, and low-level visual grounding.

This progressive breakdown highlights GeoChain’s utility as a diagnostic benchmark. By tracking model accuracy at each reasoning step, researchers can isolate failure modes (e.g. visual hallucination vs. failure to capture cultural cues) and understand how performance degrades under deeper spatial inference chains.

### Breakdown by Question Type

To assess how models handle varying degrees of response constraint, we analyzed Pass Scores across three fundamental question types: Binary, Multiclass, and Free-text, with results presented in Table 7 and Figure 7. This breakdown reveals a distinct performance hierarchy directly correlated with the open-endedness of the required answer.

Across all evaluated MLLMs, a clear difficulty gradient was observed: Binary questions yielded the highest success rates, followed by Multiclass questions, with Free-text questions proving to be the most challenging by a substantial margin. For instance, Gemini-2.5-pro achieved 88.9% on Binary and an exceptional 92.9% on Multiclass questions, but its score dropped to 56.7% for Free-text tasks. This pattern of significantly lower performance on Free-text questions was universal, underscoring the inherent difficulty in precise, open-ended generation and factual recall compared to selecting from constrained options.

In the structured formats, Gemini-2.5-pro consistently led, achieving the top scores for both Binary (88.9%) and Multiclass (92.9%) questions, with Gemini-2.5-Flash also per-

forming strongly. Notably, for the more demanding Free-text questions, GPT-4.1 emerged as the top performer with a Pass Score of 57.8%, slightly ahead of Gemini-2.5-pro (56.7%). This suggests a particular strength in GPT-4.1’s generative capabilities for unconstrained answers. Claude 3.7 Sonnet demonstrated robust performance on Binary (86.2%) and Multiclass (84.5%) questions, often comparable to GPT-4.1, but its accuracy significantly declined on Free-text questions (45.5%), reaffirming its challenges with precise, unprompted generation. As anticipated, GPT-4.1 Mini generally recorded the lowest scores across all types. This analysis by question type effectively highlights that while current MLLMs are largely proficient with constrained-choice tasks, open-ended free-text responses remain a key area for improvement.

Table 7: Pass score (%) by question type.

Model	Binary	Multiclass	Free-text
Claude 3.7 Sonnet	86.2	84.5	45.5
GPT-4.1 Mini	82.3	73.8	37.5
GPT-4.1	85.9	85.8	57.8
Gemini-2.5-Flash	88.5	90.9	50.5
Gemini-2.5-pro	<b>88.9</b>	<b>92.9</b>	<b>56.7</b>



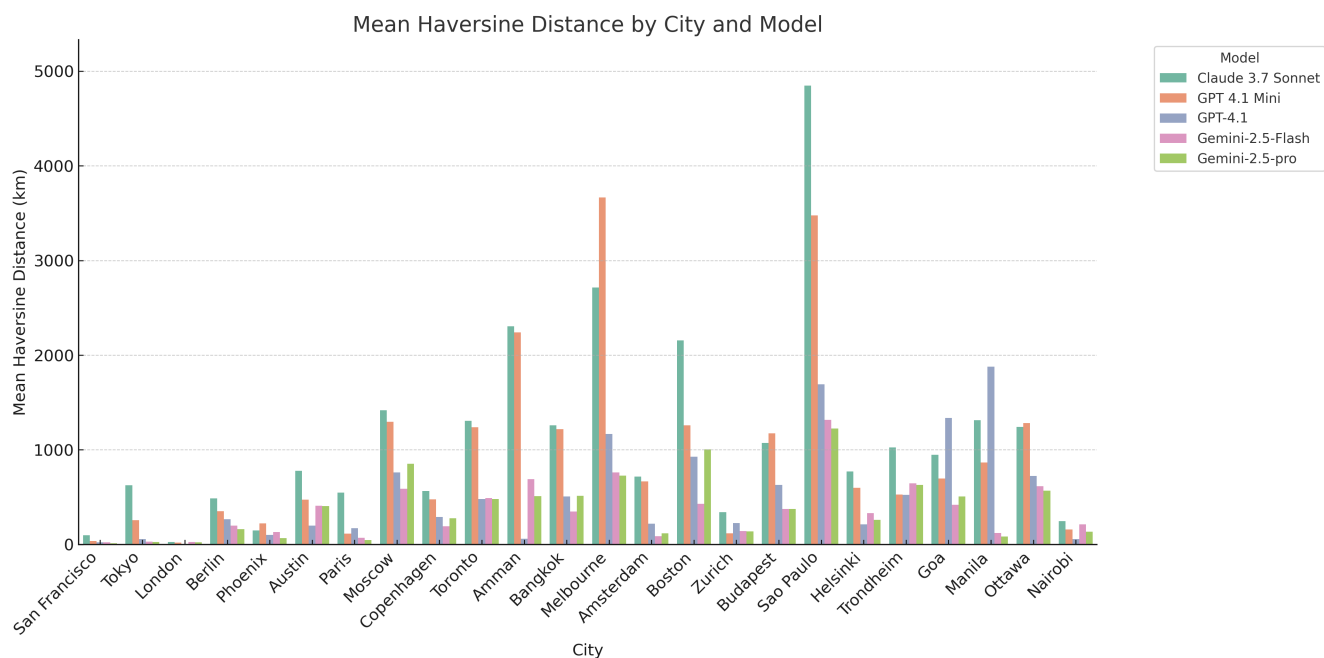


Figure 6: Mean Haversine distance (km) by city and model. Larger values indicate poor localization precision.

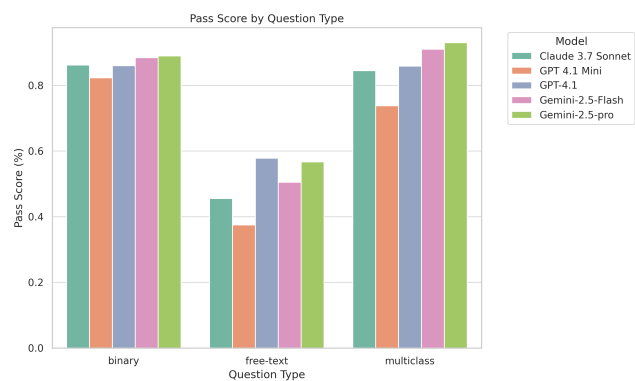


Figure 7: Model vs Question Type