Indoor 3.6M: A Benchmark Image Dataset for Global Indoor Geolocation

Anonymous authors

Paper under double-blind review

ABSTRACT

Image geolocation has advanced rapidly for outdoor imagery, driven by large-scale benchmarks and strong visual cues such as landmarks, skylines, and vegetation. In contrast, indoor image geolocation remains underexplored: indoor scenes lack distinctive geographic features, are highly ambiguous, and are not adequately represented in existing datasets. We address this gap by introducing the first large-scale benchmark for indoor geolocation, consisting of **3.6 million** images across **213 countries**, alongside a novel geographic sampling methodology that mitigates dataset bias by incorporating population, land area, and visual diversity metrics. We finetune state-of-the-art CLIP-based models such as Pigeon and GeoCLIP and report performance at country and continent levels using both top-k accuracy as well as distance based accuracy metrics. Results highlight that continent-level geolocation is feasible, but fine grained indoor geolocation e.g street and city level geolocation remains an open challenge. This work defines a new frontier for geolocation research and provides the resources to advance it.

1 Introduction

Image geolocation, which involves determining the geographic origin of a photograph based on its visual content (Hays & Efros, 2008), is a critical vision task with a wide range of applications, including forensic investigations and fraud detection. Current approaches to geolocation typically follow either a retrieval-based or classification-based framework. Retrieval-based methods depend on extensive geotagged image databases, employing similarity metrics to match query images with known locations (Hays & Efros, 2008; Vo et al., 2017). In contrast, classification-based approaches discretize the Earth's surface into geocells, treating geolocation as a multi-class classification task, requiring substantial training data per geocell to achieve high accuracy (Seo et al., 2018; Weyand et al., 2016).

As with other core computer vision tasks—such as object detection, semantic segmentation, scene recognition and image classification—the performance of geolocation models is closely tied to the availability of large, diverse, and high-quality datasets. Datasets such as ImageNet (Krizhevsky et al., 2017), MS COCO (Lin et al., 2014), and Places (Zhou et al., 2017) have been pivotal in driving progress in their respective domains. However, for image geolocation, the need for comprehensive datasets is even more pronounced due to the task's inherent complexity and global scope. The visual appearance of locations can vary dramatically depending on factors such as seasonal changes, time of day, weather conditions, and human-induced modificationsPramanick et al. (2022). Additionally,



Figure 1: Images (a) and (b) show outdoor landmarks near the Arc de Triomphe, while (c)–(f) depict diverse indoor scenes from a nearby hotel. This highlights the geolocation challenge posed by visually similar indoor environments compared to distinctive outdoor environments.

the global nature of the task requires representation across a wide variety of geographic regions, each possessing unique and sometimes subtle visual characteristics. Fine-grained geolocation further necessitates high-density, geotagged imagery to achieve precise localization.

To support the training and evaluation of geolocation models, high-quality geotagged images annotated with precise geographic location information are essential. The coverage, diversity, and geographic balance of these datasets directly influence the generalizability and accuracy of geolocation models in various contexts. Despite the growing availability of geotagged imagery from social media, curating datasets that are comprehensive, balanced, and representative of global geographic diversity remains challenging. Urban areas and popular tourist destinations are often overrepresented, while rural or less-frequented regions suffer from data scarcity, leading to models that are less effective in underrepresented areas.

While significant advances have been made in outdoor and mixed-environment (or hybrid) image geolocation, indoor geolocation remains under-explored and presents a unique set of challenges. Unlike outdoor environments, where landmarks, street signs, skylines, and natural features offer rich contextual cues, indoor spaces are more constrained and visually repetitive. The interiors of buildings, rooms, and enclosed areas typically lack the expansive contextual markers found in outdoor settings. Moreover, variations in design, layout, and lighting across indoor spaces introduce additional layers of complexity. These factors highlight the need for datasets specifically tailored to indoor environments.

An indoor image typically depicts a scene from an enclosed or semi-enclosed space, such as a home, office, or public building, and is defined by elements like furniture, walls, artificial lighting, and interior structural elements. These spaces can range from small rooms to vast halls, each with distinct characteristics. The line between indoor and outdoor environments can also blur in transitional spaces like covered patios or parking garages, where structural openness is combined with indoor elements such as artificial lighting and furniture, producing environments that straddle the boundary between the enclosed and the open.

The feasibility of indoor image geolocation lies in the distinctive visual markers inherent in the design, utilization and layout of interior spaces. Regional, cultural, religious, economic, and political factors shape architectural styles, materials, decor, and spatial layouts, resulting in distinct visual characteristics that vary geographically. Furniture, decor, artwork, religious symbols, and fixtures like electrical outlets provide valuable locational clues. Additionally, the layout of indoor spaces is often tailored to human needs and influenced by local aesthetics, making their visual structure identifiable and learnable. Despite lacking the prominent landmarks typical of outdoor settings, indoor environments offer a rich array of details that can support effective geolocation.

Given the current emphasis on outdoor geolocation and the limited focus on indoor environments, it is evident that a geographically diverse dataset dedicated to indoor image geolocation is crucial for advancing this field. Such a dataset would capture the unique characteristics of indoor spaces across a broad range of geographic locations and functional areas. Its development represents a critical step toward addressing the existing gap in indoor geolocation research and enables the creation of models capable of fine-grained localization of complex, enclosed environments. To empower research into indoor image geolocation we make the following contributions:

- We introduce a global dataset of geotagged indoor images featuring diverse scnes including living spaces, functional areas, leisure and public facilities. This extensive collection, enriched with comprehensive metadata, will empower indoor-specific geolocation research, addressing a critical gap in the current literature.
- We introduce a sampling framework that generates geographically representative image subsets by jointly accounting for visual diversity, land area, and population distribution, mitigating geographic bias in large-scale datasets used to train geolocation models.
- We present a geographically representative benchmark test set designed to evaluate the performance of both indoor-specific and hybrid geolocation models on diverse indoor scenes. This benchmark provides a standardized evaluation framework for fairly evaluating and comparing advancements in indoor geolocation research.
- Finally, we finetune GeoCLIP-yielding a specialized GeoCLIP model that outperforms the Geo-CLIP baseline across all levels of geographic granularity, establishing a benchmark for indoor image geolocation and paving the way for future innovations in this field.

The dataset along with evaluation scripts are available at: https://github.com/anonymous-for-doubleblind-review.

RELATED WORK

108

109

110 111

112 113 114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

In recent years, image geolocation has seen remarkable advancements, driven by a convergence of cutting-edge computer vision techniques, deep learning architectures, and the availability of largescale geotagged image datasets. These innovations have significantly improved the ability of models to accurately predict the geographic origin of images. The evolution of geolocation techniques has largely been defined by two primary paradigms: retrieval-based approaches and classification-based approaches. While retrieval-based methods rely on matching query images with similar images in a large geotagged database, classification-based methods divide the Earth's surface into discrete regions or geocells (Weyand et al., 2016), treating geolocation as a multi-class classification problem. More recently, hybrid approaches have emerged, combining the strengths of both paradigms to enhance geolocation accuracy.

State-of-the-art systems like PIGEON/PIGEOTTO (Haas et al., 2023) and Geoclip (Vivanco Cepeda et al., 2024) exemplify this advancement. These models utilize CLIP Vision Transformers (ViTs) (Dosovitskiy et al., 2020; Radford et al., 2021) and leverage large-scale geotagged image datasets to infer geographic locations based on visual content. The success of these systems highlights the effectiveness of modern neural architectures in capturing complex visual features tied to specific geographic locations, and the importance of combining such architectures with comprehensive, high-quality datasets.

Despite the remarkable advancements in image geolocation, global-scale geolocation remains a significant challenge, pushing researchers to focus on a more limited scope of the problem by directing attention towards closed-domain geolocation tasks. This shift arises due to the difficulties of tackling geolocation on a global scale, which necessitates access to an extensive, diverse, and truly global dataset---an asset that remains elusive. As a result, researchers have concentrated on more constrained tasks such as geolocating images of skylines (Ramalingam et al., 2010), beaches (Cao et al., 2012), deserts (Tzeng et al., 2013), the Alps (Saurer et al., 2016), hotel rooms (Stylianou et al., 2019), or specific urban areas like San Francisco (Berton et al., 2022), or even individual countries like USA Suresh et al. (2018), by leveraging tailored datasets. While these focused efforts have yielded impressive results and enhanced our understanding of geolocation techniques, they leave

139 140 141

142

143

144 145

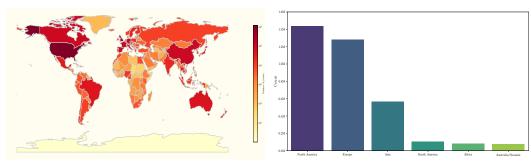
Table 1: Comparison of geolocation datasets. The "Benchmark" column indicates whether the dataset provides a dedicated test or evaluation set specifically designed to benchmark the performance of geolocation models.

Dataset	Year	Size	Scene Type	Scale	Type	Bench -mark
Im2GPS (Hays & Efros, 2008)	2008	6.5M	Mixed	Global	Geotagged	Х
San Francisco Landmarks (Chen et al., 2011)	2011	1.1M	Outdoor	City	Geotagged	Х
YFCC100M (Thomee et al., 2016)	2016	100M	Mixed	Global	Geotagged, Multimodal	х
MP-16 (Larson et al., 2017)	2017	5M	Mixed	Global	Geotagged	Х
PlaNet (Weyand et al., 2016)	2016	126M	Outdoor	Global	Geotagged	Х
Im2GPS3k (Vo et al., 2017)	2017	3k	Mixed	Global	Geotagged	√
YFCC4k (Vo et al., 2017)	2017	4K	Mixed	Global	Geotagged	√
YFCC26k (Muller-Budack et al., 2018)	2018	26K	Mixed	Global	Geotagged	√
Hotels50K (Stylianou et al., 2019)	2019	1M	Indoor (Hotel rooms)	Global	Geotagged	х
GWS15K (Clark et al., 2023)	2023	15K	Outdoor	Global	Geotagged	√
INDOOR-3.6M	2024	3.6M	Indoor (Scene agnostic)	Global	Geotagged, Multimodal	х
INDOOR-40K	2024	40K	Indoor (Scene	Global	Geotagged, Multimodal	✓

159

156 157

agnostic)



- (a) Global distribution of INDOOR-3.6M data
- (b) Distribution of Images Across Continents

Figure 2

an important gap in the field—specifically, the geolocation of scene-agnostic indoor imagery on a global scale.

Indoor image geolocation presents a unique challenge with valuable applications in fields such as digital forensics, law enforcement, and augmented reality. However, it requires geotagged indoor imagery with global diversity, which current datasets lack. For instance, indoor datasets like NYU Depth V2 (Silberman et al., 2012), SUN RGB-D (Song et al., 2015), and Places365 (Zhou et al., 2017) are designed for tasks such as object detection and scene recognition but do not provide the geographic metadata necessary for geolocation. Similarly, mixed-environment datasets such as MediaEval Placing Task (MP-16) (Larson et al., 2017) and YFCC100M (Thomee et al., 2016), which encompass both indoor and outdoor environments, also fall short for indoor geolocation as they tend to prioritize outdoor scenes. While these datasets have led to the development of powerful geolocation algorithms, models trained on them often perform poorly on indoor imagery due to the substantial differences in visual characteristics between indoor and outdoor environments.

Existing image geolocation benchmark datasets, such as IM2GPS (Hays & Efros, 2008) and its successor IM2GPS3k (Vo et al., 2017), along with subsets of YFCC100M like YFCC4K (Vo et al., 2017) and YFCC26K (Muller-Budack et al., 2018), have been instrumental in evaluating geolocation systems. However, these datasets predominantly comprise outdoor imagery, rendering them inadequate for assessing indoor-specific geolocation tasks. Indoor environments present unique challenges, necessitating the interpretation of more complex and nuanced visual features, including variations in room layout, furniture arrangements, lighting conditions, and decorative elements. Consequently, to facilitate accurate indoor geolocation, it is imperative to develop specialized indoor-specific datasets for both training and benchmarking purposes.

3 DATASET OVERVIEW

To achieve accurate and reliable indoor image geolocation, a large and diverse dataset covering various indoor environments is essential. The INDOOR-3.6M dataset addresses this need by being agnostic to specific indoor scenes, enabling generalization across a wide range of locations, including residential, office, shopping, leisure, and public spaces. Geolocation data, provided either as GPS coordinates or text-based location labels, is included alongside textual information such as descriptions and metadata as supplementary features. This multimodal approach enhances the dataset's versatility, particularly for tasks that benefit from both visual and textual data. It is important to note that the dataset does not explicitly identify specific locations in the manner typical of place recognition tasks. However, the accompanying text, and descriptions may contain useful information that could inform place recognition applications.

Data sources and collection methods: The INDOOR-3.6Mdataset was constructed using three primary sources: Flickr (flickr.com, 2024), a popular photo-sharing platform where users upload and tag images with metadata; Wikidata (wikidata.org, 2024), a free, collaborative knowledge base that provides structured data to support Wikipedia and other Wikimedia projects; and Booking.com (booking.com, 2024), a popular hotel booking website. For the image repositories (Flickr and Wiki-

data), we formulated search terms based on indoor scene categories from the Places365 dataset¹, and appended "indoor" to categories typically associated with outdoor environments. To ensure usability and proper attribution, we restricted our search to images with Creative Commons licenses and included only those with latitude and longitude coordinates. However, the initial search terms yielded few results because of these constraints. To address this, we manually refined the search terms, generalizing specific categories like "ski resort" to broader terms such as "resort" and expanding our vocabulary with synonyms and colloquial terms. Additionally, we introduced new categories that seemed relevant but were absent from the original list. Productive search terms included "living room", "indoor", "villa", "cottage", "diner", "office space", and "beach house". For the web scraping component, we employed country labels to initialize a crawler that retrieved images from search results for each country.

This collection process yielded approximately 10 million candidate images combined. In addition to visual data, we also collected associated textual metadata from these platforms, such as usergenerated tags, descriptions, and captions. The textual data varies significantly in length, language, and detail, ranging from brief labels to detailed narratives or contextual information.

To ensure the dataset's focus remained on indoor scenes, we filtered the candidate images using the Places365 ResNet indoor/outdoor image classifier (Zhou et al., 2017). Recognizing that the distinction between indoor and outdoor scenes can sometimes be blurred, we used the classifier to quantify the "indoorness" of each image. We retained only those images with a probability of being indoor, $P(indoor) \geq 0.5$. Additionally, we recorded this likelihood score for each image in the metadata, placing images on a continuum between relatively indoor spaces (P(indoor) = 0.5) and purely indoor spaces (P(indoor) = 1.0).

Scale and Distribution: The INDOOR-3.6Mdataset comprises 3.6 million images spanning a wide variety of scenes from 223 countries worldwide, uploaded between 1978 and 2024. While the dataset aims to be representative of indoor environments, it is not entirely geographically representative due to its reliance on online sources (See Figure 2a). This dependence introduces inherent biases in geographic distribution, resulting in over-representation of regions with a strong digital footprint and larger populations such as United States (which represent 30% of the data), and under-representation of areas with less online activity or smaller populations. Figure 5 illustrates the dataset's distribution according to the MIT indoor scenes label set. A significant portion of the images are labeled as "tv studio", which predominantly corresponds to spaces where a TV is present—-commonly living rooms.

Metadata enrichment: The dataset incorporates metadata enrichment encompassing geospatial information, scene classification, and object segmentation. Using the GPS data, we use the Nominatim API(Nominatim, 2024) to perform reverse geocoding, yielding detailed location information including building names, street addresses, suburbs, and cities. This granular metadata facilitates fine-grained, location-based classification tasks. In addition, for each image, we include top 10 scene category labels obtained from Places365 and a ViT trained on MIT indoor Scenes dataset, as well as segmentation masks extracted using Segment Anything Model (SAM)(Kirillov et al., 2023) for pixel-level segmentation and YOLOv8(Jocher et al., 2022) for object detection and labeling. Scene labels, segmentation masks, and object detection results enhance the dataset by providing additional cues for geolocation. These annotations help models identify important features like furniture, signage, or cultural artifacts, which are critical for pinpointing locations. Such features also align with real-world practices, like Europol's 'Trace an Object'tra initiative, where visual clues in scenes are used to infer locations. By including these annotations, the dataset supports more advanced and accurate geolocation methods.

4 Indoor Image Geolocation Benchmark Dataset

Our analysis reveals that current benchmark datasets for image geolocation predominantly consist of outdoor scenery. Figure A.1 illustrates the percentage of indoor images identified at various likelihood thresholds across existing mixed-environment image geolocation benchmark datasets. Furthermore, current geolocation models trained predominantly on outdoor imagery struggle signif-

¹https://github.com/CSAILVision/places365/blob/master/categories_ places365.txt

icantly with indoor environments. The performance gap between outdoor and indoor settings can exceed hundreds of kilometers in positioning error, highlighting the need for specialized approaches to indoor geolocation.

To address the limitations of current benchmark datasets, which predominantly focus on outdoor environments, we introduce a new benchmark dataset specifically for indoor geolocation: INDOOR-40K. This dataset is curated to minimize the visual biases of existing benchmarks by providing a diverse collection of 15,000 images from various indoor environments across 193 countries. To ensure the dataset is distinct from those used to train existing geolocation models, we carefully selected images captured after 2017—following the release of YFCC100M—and exclusively sourced from booking.com, ensuring each image contains GPS metadata. This curation process resulted in a initial pool of approximately 800,000 images, from which we sampled the final benchmark set according to the methodology outlined in the next section.

5 Sampling Strategy

To address the inherent geographic bias in large-scale image collections, we propose a sampling strategy that incorporates geographic and visual factors into subset construction. The method proceeds in three steps. First, we compute visual diversity scores from image embeddings obtained using the CLIP ViT-L/14 model. Diversity is quantified by the average pairwise cosine distance between embeddings, stratified across semantic scene categories within each country. This ensures that the sampled images capture a range of scene types rather than over representing a small set of dominant visual contexts. Second, we incorporate external country-level statistics—land area and population size—into a regression framework that estimates their relative contributions to observed visual diversity. The resulting weights are expressed as: $w_c = \alpha \cdot \text{Population}_c + \beta \cdot \text{LandArea}_c$, where w_c is the sampling weight for country c. Several regression models were evaluated; Random Forest regression achieved the highest performance ($R^2 = 0.861$) in predicting country-level diversity, substantially outperforming linear and polynomial baselines. Feature importance analysis further indicated that both factors contributed nearly equally (population: 0.493, land area: 0.507).

Finally, samples are drawn proportionally to these weights across scene categories, yielding a geographically representative dataset (Figure 3). We applied this sampling strategy to both our benchmark test set and training subsets from INDOOR-3.6M, ensuring geographic and scene-level diversity throughout the evaluation pipeline.

To evaluate the impact of this sampling strategy, we perform ablation studies comparing model performance on datasets sampled with and without our method. Specifically, we fine-tune Geo-CLIP to establish a baseline for indoor geolocation using a subset of the INDOOR-3.6M dataset, following the sampling strategy described. GeoCLIP was selected for its state-of-the-art performance in environment-agnostic geolocation. We retained most of the training parameters from Vivanco Cepeda et al. (2024), including a constant learning rate of 1e-6 and a batch size of 256. The model converged after 10 epochs and outperformed the original GeoCLIP on our test set. Table 2 highlights the improved performance across all levels of granularity.

We also assessed the zero-shot classification performance of CLIP on a location classification task. For this, using the INDOOR-3.6M dataset, we divided the Earth into semantic geocells based on the approach in Haas et al., ensuring each geocell contained between 1,000 and 2,000 images. This resulted in approximately 1,300 geocells. We utilized the image encoder from the clip-vit-large-patch14 Radford et al. (2021) architecture to perform zero-shot classification of geocells. The encoder extracted visual embeddings, which were then used to predict geocells without additional training. For GPS prediction, the latitude and longitude of an image were approximated by averaging the GPS coordinates of all images within the predicted geocell. The results of these experiments are presented in Table 2.

The study underscores the potential of domain-specific training in enhancing geolocation models, particularly for indoor environments. Our experiments with GeoCLIP on the INDOOR-3.6M dataset reveal critical insights into model performance across various geographic scales, with the fine-tuned GeoCLIP consistently outperforming its counterparts. The reduction in mean distance error from 4089.11 km for the baseline GeoCLIP to 3598.02 km for the fine-tuned version is especially remarkable given the inherent complexity of indoor geolocation. The most striking

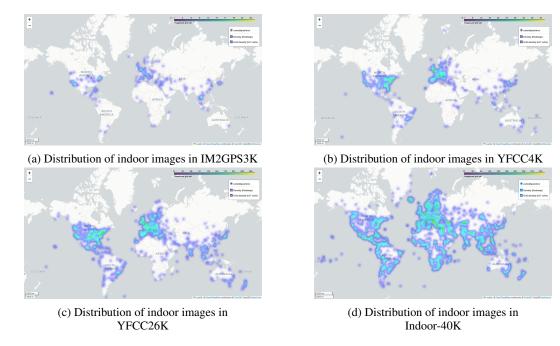


Figure 3: Indoor image distributions ($p_{indoor} \ge 0.5$) across geolocation benchmarks. Existing datasets (IM2GPS3K, YFCC4K, YFCC26K) are geographically skewed, while Indoor40K (ours) achieves more representative global coverage.

observations emerge at broader scales, where fine-tuned GeoCLIP demonstrates pronounced gains, such as improving continent-level accuracy from 53% to 61% and country-level accuracy from 25% to 35%. These results highlight the ability of the model to leverage the diversity and richness of INDOOR-3.6M to capture geographically meaningful features. While the gains at finer scales, such as street and city levels, are more modest, the consistent improvements across all levels reinforce the importance of domain-specific datasets in overcoming the unique challenges of indoor geolocation.

To evaluate the impact of the proposed sampling strategy, we conducted ablation studies using datasets prepared with random sampling and the strategic sampling methodology described in the Appendix. The model finetuned on the dataset created using our sampling strategy yields better performance on geolocating both over represented classes and underrepresented classes.

Table 2: Comparison of GeoCLIP, finetuned GeoCLIP (GeoCLIP*), and Pigeon on indoor images in current benchmark dataset and Indoor-40K.

	Model	Continent	Country	Region	City	Street	Mean Dist.
		(2500 km)	(750 km)	(200 km)	(25 km)	(1 km)	Error (km)
im2gps3k (n = 530)	GeoCLIP	0.5981	0.3940	0.2064	0.1660	0.0774	2991.76
	Pigeon*	0.456	0.218	0.085	0.060	0.013	5220.12
	GeoCLIP*	0.7132	0.4396	0.2604	0.1717	0.0717	2682.21
yfcc4k (n = 1769)	GeoCLIP	0.6337	0.3825	0.1724	0.0888	0.0407	3253.73
	Pigeon*	0.382	0.160	0.057	0.028	0.008	5035.26
	GeoCLIP*	0.6733	0.4036	0.1860	0.0955	0.0367	2982.24
yfcc26k (n = 9144)	GeoCLIP	0.6215	0.3757	0.1742	0.0988	0.0441	3343.35
	Pigeon*	0.580	0.320	0.140	0.075	0.025	3850.0
	GeoCLIP*	0.6666	0.4193	0.1976	0.1012	0.0418	2965.46
indoor-40k (n=40000)	GeoCLIP	0.5175	0.2399	0.1004	0.0479	0.0204	4316.22
	Pigeon*	0.519	0.290	0.194	0.158	0.134	4468.52
	GeoCLIP*	0.5822	0.3097	0.1418	0.0632	0.0191	3897.20

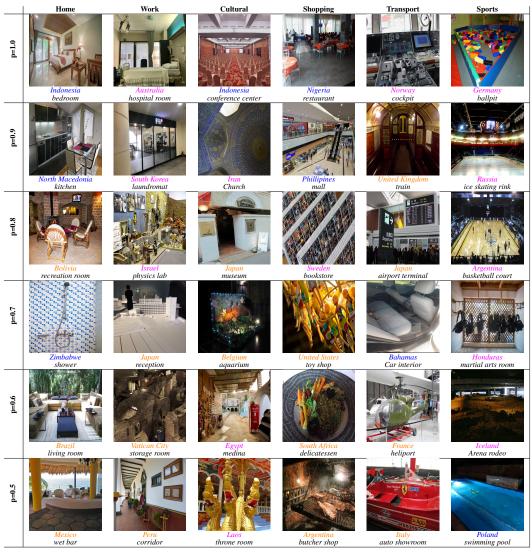


Figure 4: Samples of images from the dataset representing different parts of the world. The rows correspond to the indoor likelihood score P(indoor), while the columns categorize the scene types according to Places365 indoor scene categories at Level 2. Country names in blue, magenta, and yellow are sourced from Booking.com, Wikidata, and Flickr, respectively.

6 CHALLENGES

Large-scale indoor datasets face two critical challenges. First, INDOOR-3.6M exhibits significant **geographic and demographic biases** due to over-representation of regions with higher internet penetration and tourism, hindering model performance in underrepresented areas. Second, **GPS data validation** remains problematic—user-uploaded geotags from photo-sharing platforms are of-

Table 3: Comparison of Zero Shot CLIP, CLIP Linear probing, Pigeon* and GeoCLIP* on indoor-40K.

Model	Continent Accuracy	Top-1 Country Accuracy	Top-3 Country Accuracy	Top-5 Country Accuracy
CLIP zero-shot	0.594	0.152	0.281	0.359
CLIP Linear Probing	0.568	0.196	0.326	0.402
Pigeon*	0.659	0.237	0.416	0.491
GeoCLIP*	0.661	0.256	0.428	0.516

ten unreliable due to device limitations, poor satellite coverage, or manual tagging errors. While hotel booking platforms provide more accurate GPS data, they are limited to residential scenes.

The emergence of large vision models like Vision Transformers (ViTs) (Dosovitskiy et al., 2020) and CLIP (Radford et al., 2021) introduces challenges related to potential **data leakage**. These models are pretrained on vast datasets scraped from the internet, including Flickr-sourced collections like YFCC100M. Consequently, there is no guarantee that new publicly sourced datasets, such as INDOOR-3.6M, do not introduce a data leak when fine-tuning such models. This overlap could artificially inflate performance metrics during model evaluation. To mitigate this risk for our benchmark test set, we deliberately selected images captured after 2017, the publication year of YFCC100M, reducing the likelihood of overlap with this widely used dataset.

Indoor geolocation datasets introduce additional difficulties for geolocation systems due to **intraclass variation**. Unlike outdoor environments, where variations are often limited to views in the four cardinal directions (North, South, East, and West), indoor spaces exhibit far more complexity. In settings like hotels, different floors and rooms have distinct layouts, styles, and views, making it harder for models to establish consistent visual cues. This issue is exacerbated by the absence of clear landmarks, necessitating more nuanced feature extraction. Moreover, indoor environments are more subject to **temporal dynamics**. Frequent renovations, redecorations, and repurposing result in visual instability, which can quickly render models obsolete. Continuous updating or adaptive learning is required to ensure that models remain effective over time. To truly advance the field of indoor geolocation, it is crucial for future work to actively confront these issues, ensuring that models are both reliable and adaptable across diverse and evolving environments.

7 ETHICS

The INDOOR-3.6M dataset has been developed with careful attention to ethical considerations. The dataset contains geotagged indoor images sourced from public platforms, without the intention of identifying specific individuals or private spaces. License and owner information are included in the metadata to allow proper attribution. Geographic bias is acknowledged, particularly the overrepresentation of urban areas, and researchers are encouraged to apply sampling strategies and imbalance mitigation techniques to achieve fairer regional representation in model training. The dataset is *strictly* for research purposes, and misuse for purposes such as unauthorized surveillance or invasive applications is strongly discouraged. Researchers are urged to handle the data responsibly, especially during algorithm development and when implementing public-facing technologies.

There are concerns about the harmful applications of this dataset for geolocation technology, including privacy violations and unauthorized surveillance. We encourage researchers to remain mindful of the societal impact of their work, implementing safeguards to prevent abuse and adhering to privacy laws and ethical standards. It is essential that the research community stays actively engaged in discussions about the ethical development and use of indoor geolocation technologies, ensuring that advancements prioritize individual privacy and security. Misuse for invasive purposes is explicitly discouraged.

8 CONCLUSION

We introduce a new specialised dataset for indoor image geolocation (INDOOR-3.6M) as well as a benchmark dataset–INDOOR-40K. These contributions represent a significant step toward addressing the unique challenges of indoor image geolocation, where traditional outdoor models often struggle. Our dataset offers global coverage of diverse indoor spaces, enabling geolocation models to learn fine-grained features that are critical for accurately predicting the locations of indoor scenes. Our results demonstrate the utility of this dataset in improving the performance of geolocation models on indoor environments. Fine-tuning CLIP based geolocation models e.g Pigeon and GeoCLIP with INDOOR-3.6M yielded measurable improvements across various levels of geographic granularity. However, indoor geolocation remains a challenging problem, with mean distance errors on the INDOOR-40K test set still exceeding 3,500 km. Despite these challenges, INDOOR-3.6M lays a strong foundation for advancing indoor geolocation.

REFERENCES

- Stop Child Abuse Trace an Object Europol europol.europa.eu. https://www.europol.europa.eu/stopchildabuse. [Accessed 07-09-2024].
- Gabriele Berton, Carlo Masone, and Barbara Caputo. Rethinking visual geo-localization for large-scale applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4878–4888, 2022.
- booking.com. Booking.com. https://www.booking.com, 2024. Online travel agency.
- Liangliang Cao, John R Smith, Zhen Wen, Zhijun Yin, Xin Jin, and Jiawei Han. Bluefinder: estimate where a beach photo was taken. In *Proceedings of the 21st International Conference on World Wide Web*, pp. 469–470, 2012.
 - David M Chen, Georges Baatz, Kevin Koser, Sam S Tsai, Ramakrishna Vedantham, Timo Pylvanainen, Kimmo Roimela, Xin Chen, Jeff Bach, Marc Pollefeys, et al. City-scale landmark identification on mobile devices. In *CVPR* 2011, pp. 737–744. IEEE, 2011.
 - Brandon Clark, Alec Kerrigan, Parth Parag Kulkarni, Vicente Vivanco Cepeda, and Mubarak Shah. Where we are and what we're looking at: Query based worldwide image geo-localization using hierarchies and scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23182–23190, 2023.
 - Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
 - flickr.com. Flickr. https://www.flickr.com, 2024. Photo-sharing platform.
 - Lukas Haas, Silas Alberti, and Michal Skreta. Pigeon: Predicting image geolocations. *arXiv* preprint *arXiv*:2307.05845, 2023.
 - James Hays and Alexei A Efros. Im2gps: estimating geographic information from a single image. In 2008 ieee conference on computer vision and pattern recognition, pp. 1–8. IEEE, 2008.
 - Glenn Jocher, Ayush Chaurasia, Alex Stoken, Jirka Borovec, Yonghye Kwon, Kalen Michael, Jiacong Fang, Colin Wong, Zeng Yifu, Diego Montes, et al. ultralytics/yolov5: v6. 2-yolov5 classification models, apple m1, reproducibility, clearml and deci. ai integrations. *Zenodo*, 2022.
 - Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
 - Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
 - Martha Larson, Mohammad Soleymani, Guillaume Gravier, Bogdan Ionescu, and Gareth JF Jones. The benchmarking initiative for multimedia evaluation: Mediaeval 2016. *IEEE MultiMedia*, 24 (1):93–96, 2017.
 - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
 - Eric Muller-Budack, Kader Pustu-Iren, and Ralph Ewerth. Geolocation estimation of photos using a hierarchical model and scene classification. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 563–579, 2018.
 - Nominatim. Nominatim. https://www.nominatim.org, 2024. Geocoding API.

- Shraman Pramanick, Ewa M Nowara, Joshua Gleason, Carlos D Castillo, and Rama Chellappa.
 Where in the world is this image? transformer-based geo-localization in the wild. In *European Conference on Computer Vision*, pp. 196–215. Springer, 2022.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
 - Srikumar Ramalingam, Sofien Bouaziz, Peter Sturm, and Matthew Brand. Skyline2gps: Localization in urban canyons using omni-skylines. In 2010 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 3816–3823. IEEE, 2010.
 - Olivier Saurer, Georges Baatz, Kevin Koser, Lubor Ladicky, and Marc Pollefeys. Image based geo-localization in the alps. *International Journal of Computer Vision*, 116:213–225, 2016.
 - Paul Hongsuck Seo, Tobias Weyand, Jack Sim, and Bohyung Han. Cplanet: Enhancing image geolocalization by combinatorial partitioning of maps. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 536–551, 2018.
 - Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V 12*, pp. 746–760. Springer, 2012.
 - Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 567–576, 2015.
 - Abby Stylianou, Hong Xuan, Maya Shende, Jonathan Brandt, Richard Souvenir, and Robert Pless. Hotels-50k: A global hotel recognition dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 726–733, 2019.
 - Sudharshan Suresh, Nathaniel Chodosh, and Montiel Abello. Deepgeo: Photo localization with deep neural network. *arXiv preprint arXiv:1810.03077*, 2018.
 - Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.
 - Eric Tzeng, Andrew Zhai, Matthew Clements, Raphael Townshend, and Avideh Zakhor. User-driven geolocation of untagged desert imagery using digital elevation models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 237–244, 2013.
 - Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired alignment between locations and images for effective worldwide geo-localization. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Nam Vo, Nathan Jacobs, and James Hays. Revisiting im2gps in the deep learning era. In *Proceedings* of the IEEE international conference on computer vision, pp. 2621–2630, 2017.
 - Tobias Weyand, Ilya Kostrikov, and James Philbin. Planet-photo geolocation with convolutional neural networks. In *European Conference on Computer Vision*, pp. 37–55. Springer, 2016.
 - wikidata.org. Wikidata. https://www.wikidata.org, 2024. Free and open knowledge base.
 - Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3920–3928, 2017.

A SUPPLEMENTARY MATERIALS

A.1 INDOOR-40K DISTRIBUTIONS

Figure 5 shows the scene and continent distribution in INDOOR-40K. The dataset exhibits diverse coverage across indoor scene categories, with the most represented scenes being bathrooms, bedrooms, restaurants, and various commercial spaces. The geographic distribution shows reasonable global coverage with Europe (33.4%), Asia (25.9%), and Africa (17.6%) being the most represented continents.

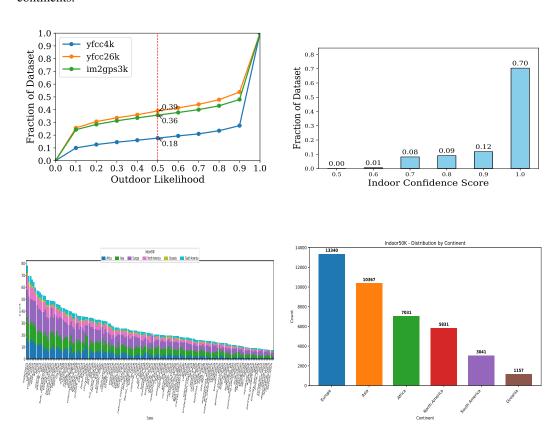


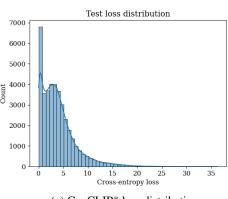
Figure 5: Scene (left) and continent (right) distribution in INDOOR-40K.

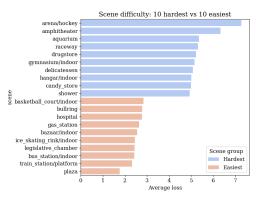
A.2 GEOCLIP* ERROR ANALYSIS

Figure 6 presents an analysis of indoor geolocation difficulty on INDOOR-40K. The loss distribution reveals a long-tailed pattern where many indoor images can be localized with moderate error, but a substantial fraction are extremely challenging. Scene-level analysis shows that average geolocation error varies substantially across categories—highly distinctive environments (e.g., religious sites, laboratories) are easier to localize, while visually generic spaces (e.g., corridors, gyms) yield much higher errors.

A.3 CLIP vs GeoCLIP* Embedding Analysis

We analyze clustering behavior across three embedding spaces: (i) Pretrained CLIP (ViT-L/14) embeddings, trained for generic vision-language alignment; (ii) Finetuned GeoCLIP image embeddings (GeoCLIP*), which incorporate geographic supervision into indoor image features; and (iii) Geo-CLIP* location embeddings, which embed GPS coordinates into the same latent space as the indoor images.





(a) GeoCLIP* loss distribution

(b) Scene-level difficulty

Figure 6: Analysis of indoor geolocation difficulty on INDOOR-40K. (a) The long-tailed loss distribution shows that while many indoor images can be localized with moderate error, a substantial fraction are extremely challenging. (b) Average geolocation error varies substantially across scene categories—highly distinctive environments are easier to localize, while visually generic spaces yield much higher errors.

We apply hierarchical agglomerative clustering to the test set embeddings for each embedding type. Figures 7 and 8 show sampled images from resulting clusters of bedroom images with country and continent labels.

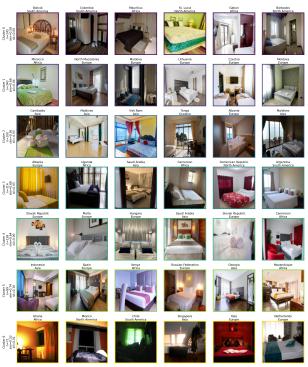
To examine global structure, we also project all images onto a world map and color them by cluster assignment. This allows direct inspection of whether clusters correspond to coherent geographic regions.

We observe that **Pretrained CLIP embeddings** produce clusters that reflect semantic similarity but are geographically incoherent, while **GeoCLIP image embeddings** show improved regional consistency, with clusters more often localized to specific continents. **GeoCLIP location embeddings**, which are derived from GPS coordinates rather than images, exhibit the strongest geographic alignment, with clusters corresponding to contiguous world regions. This suggests that the joint embedding space effectively bridges image and location representations, enabling geographically coherent structure to emerge.



(b) CLIP cluster world map

Pretrained CLIP (ViT-L/14) embeddings: Clusters are semantically coherent but geographically incoherent. Images with similar layouts or styles appear grouped together, yet cluster assignments scatter across multiple



(c) GeoCLIP image bedroom clusters



(d) GeoCLIP image cluster world map

GeoCLIP image embeddings: Compared to CLIP, clusters show improved continent and regional consistency. Images from the same continent more frequently appear in the same cluster, and map projections (right) reveal clearer geographic grouping.

Figure 7: Clustering comparison between CLIP and GeoCLIP image embeddings. Left: sampled clusters of bedroom images annotated with country and continent labels. Right: world maps with images colored by cluster assignment.





(a) GeoCLIP location bedroom clusters

(b) GeoCLIP location cluster world map

GeoCLIP location embeddings: Despite being derived from GPS coordinates rather than images, these embeddings exhibit the strongest geographic alignment. Clusters correspond to contiguous world regions, showing the effectiveness of the joint image-location embedding space.

Figure 8: Clustering results for GeoCLIP location embeddings. Left: sampled clusters of bedroom images from Indoor-40K annotated with country and continent labels. Right: world maps with images colored by cluster assignment.