

OpenStreetView-5M: The Many Roads to Global Visual Geolocation

Guillaume Astruc*^{1,2,5} Nicolas Dufour*^{1,6} Ioannis Siglidis*¹ Constantin Aronssohn¹
Nacim Bouia¹ Stephanie Fu^{1,4} Romain Loiseau^{1,2} Van Nguyen Nguyen¹
Charles Raude¹ Elliot Vincent^{1,3} Lintao XU¹ Hongyu Zhou¹ Loic Landrieu¹

¹ LIGM, Ecole des Ponts, CNRS, UGE ² UGE, IGN, ENSG, LASTIG ³ Inria Paris ⁴ UC Berkeley

⁵ CESBIO, Univ de Toulouse, CNES/CNRS/IRD/INRAE/UPS ⁶ LIX, CNRS, Ecole Polytechnique, IP Paris

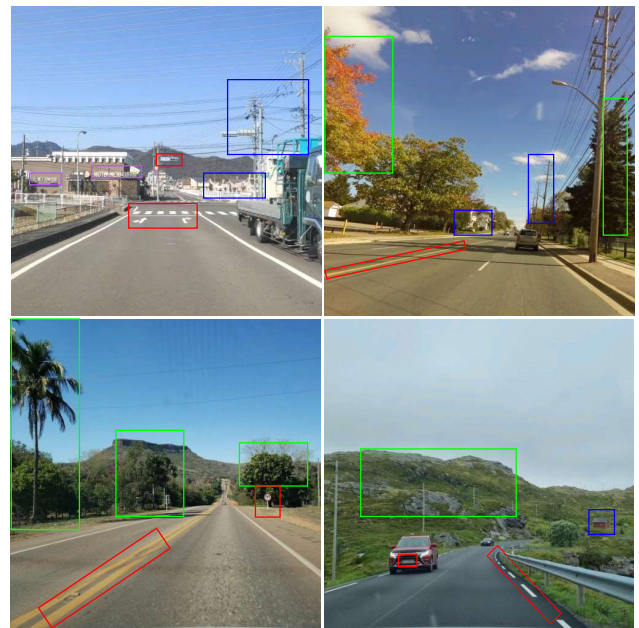
Abstract

Determining the location of an image anywhere on Earth is a complex visual task, which makes it particularly relevant for evaluating computer vision algorithms. Yet, the absence of standard, large-scale, open-access datasets with reliably localizable images has limited its potential. To address this issue, we introduce OpenStreetView-5M, a large-scale, open-access dataset comprising over 5.1 million geo-referenced street view images, covering 225 countries and territories. In contrast to existing benchmarks, we enforce a strict train/test separation, allowing us to evaluate the relevance of learned geographical features beyond mere memorization. To demonstrate the utility of our dataset, we conduct an extensive benchmark of various state-of-the-art image encoders, spatial representations, and training strategies. All associated codes and models can be found at <https://github.com/gastruc/osv5m>.

1. Introduction

While natural image classification is the standard for evaluating computer vision methods [15, 58, 69], global geolocation offers a compelling alternative task. In contrast to classification, where the focus is often a single object, geolocation involves detecting and combining various visual clues, like road signage, architectural patterns, climate, and vegetation. Predicting a single GPS coordinate or location label from these observations necessitates a rich representation of both the Earth’s culture and geography; see Figure 1 for some examples. Furthermore, the abundance of geo-tagged street-view images depicting complex scenes with a clear and consistent point of view makes this task appropriate for training and evaluating modern vision models.

Despite this potential, few supervised approaches are



drivephotograph, and_eng, gciem, bootprint, Mapillary, licensed under CC-BY-SA.

climate/vegetation traffic markers architecture culture/script

Figure 1. **Global Visual Geolocation.** Predicting the location of an image taken anywhere in the world from just pixels requires detecting a combination of clues of various abstraction levels [44]. Can you guess where these images were taken?¹

trained and evaluated for the task of geolocation. We attribute this to the limitations of existing geolocation datasets: (i) Large and open geolocation datasets contain a significant portion of noisy and non-localizable images [26, 32, 70]; (ii) Street view datasets are better suited for the task but are both proprietary and expensive to download [11, 14, 23, 25, 41, 63]. To address these issues, we introduce OpenStreetView-5M (OSV-5M), an open-access

*Denotes equal contributions.

¹ From top left to bottom right: Nagoya, Japan; Ontario, Canada; Mato Grosso, Brazil; Lofoten, Norway.

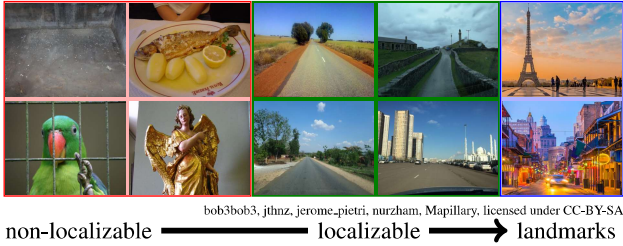


Figure 2. **Localizable vs Non-Localizable.** Images from our dataset (green) occupy the space between weakly localizable images (red) like the ones from the test set of Im2GPS3k [70] and landmark images used to advertise CV conferences (blue).

dataset of 5.1 million high-quality and crowd-sourced street view images. Our ambition is to make both street view images and global geolocation new standards for measuring progress in deep learning.

Automating visual geolocation has significant potential benefits, with direct applications in fields such as journalism, forensics, as well as historical and cultural studies. Learning robust geographical representations may also be valuable for various deep learning challenges, including self-supervised learning and generative modeling, or the development of more interpretable AI systems. Thanks to its size and scope, and its strict train/test split, OSV-5M serves as a robust and reliable benchmark for computer vision models. To demonstrate this, we design an extensive evaluation experiment to measure the impacts of various factors such as pretraining strategies, model scale, spatial representations, fine-tuning approaches, contrastive losses, and auxiliary tasks.

2. Related Work

In this section, we detail the notion of image localizability (Section 2.1), the main existing geolocation datasets (Section 2.2), and geolocation methods (Section 2.3).

2.1. Localizability

As noted by Izbicki et al. [32], images exhibit a range of localizability, an inherently perceptual concept, see Figure 2. Non-localizable images lack information that connects them to a specific location or are of too low quality to properly analyze. Weakly localizable images only contain vague or indirect hints, such as people, animals, and objects in indoor scenes. Localizable images should contain enough information to allow for an informed guess relative to their location. For example, street view images are generally localizable as they typically contain salient features indicative of the local environment such as climate, nature, architecture, or utility and regulatory infrastructure. At the far end of the spectrum, landmark images showcase emblematic monuments or iconic landscapes, making their location instantly identifiable to most viewers.

According to this criteria, a visual inspection suggests that 35% of the images in Im2GPS3k, a dataset commonly used

to benchmark geolocation methods [70], are non-localizable. When used for evaluation, this may lead to unreliable errors or promote methods that have memorized biases of the training distribution. When used for training, non-localizable images can lead to sub-optimal representations or encourage spurious correlations. OSV-5M predominantly comprises localizable street view images whose accurate geolocation requires robust geographical representations.

2.2. Geolocation Datasets

We motivate the need for OSV-5M by reviewing existing geolocation datasets from the two main sources of geotagged images: web-scraped and street view images, see Table 1.

Web-Scraped. Image hosting platforms like Flickr provide a near-endless source of geotagged images, which has been used to create large open datasets, like YFC100M [66]. Most images correspond to personal or amateur photographs representing food, art, and images of pets and friends, and are either weakly or non-localizable. Even strongly localizable images are typically taken in tourist spots, injecting an often Western cultural bias towards recognizable landmarks [36]. The provided location metadata can be occasionally missing or inaccurate, and the online nature of these images implies they can be easily removed, hindering reproducibility². For evaluation purposes, cleaner subsets have been proposed that improve both the image distribution coverage and annotation quality [64, 70], but remain still heavily biased and predominantly non-localizable. Despite their small scope and size, these datasets are currently the primary means of evaluating geolocation models.

Street View. Conversely, street view images tend to be strongly localizable. Captured through panoramic cameras or dash-cams, they depict in high quality a vehicle’s surroundings, which corresponds mostly to outdoor scenes with rich geographical cues. Google famously provides a global street view coverage, which is, however, expensive to acquire for academic purposes (\$1000 for 150k images) and cannot be shared. Existing open datasets from this source either only consist of dense samples from 3 US cities meant for navigation [46, 77], or are inaccessible [14, 23, 41].

Luckily, crowd-sourced platforms such as Mapillary [4] offer a global and diverse source of open-access street view images for various environments, from dense cities and suburbs to remote and inhabited landscapes. These images have been used to construct several benchmarks for multiple tasks other than geolocation, including depth estimation [9], semantic segmentation [49], traffic sign detection and classification [18], place recognition [71] and visual localization [33]. With 5.1M Mappillary images taken across the globe, OSV-5M is the largest open-access street-view image dataset

²60% of the 2014 YFCC-split [47] was deleted by 2020 [32]!

Table 1. **Geolocation Datasets.** OpenStreetView-5M contains strongly localizable street views with access, scope, and size comparable to web-scraped databases.

Image Source	size	open-access	scope
Web-scraped			
Im2GPS [26]	237	✓	biased
Im2GPS3k [70]	2997	✓	biased
YFCC4k [70]	4536	✓	biased
YFCC26k [64]	26k	✓	biased
MP-16 [39]	4.7M	✓	biased
Moussely <i>et al.</i> [47]	14M/6M ²	✓	global
YFCC100M [66]	100M	✓	biased
PlaNet [72]	125M	✗	biased
Street view			
Google-WS-15k [14]	15k	✗	global
GMCP [77]	105k	✗	3 cities
StreetCLIP [23]	1M	✗	unknown
OpenStreetView-5M	5.1M	✓	global

and the only one designed for global geolocation. OSV-5M has a similar order of magnitude to popular YFCC-based geolocation train sets [39, 47], and comes with a clean test set that is 33 times bigger than the current largest street-view image test benchmark [14] (which is not openly accessible).

2.3. Geolocation Methods

Place recognition [78] and visual localization [16, 37, 53, 54, 60] are popular tasks that consist in finding the pose of images in a known scene. In contrast, visual geolocation predicts 2D coordinates or discrete locations (*e.g.*, countries), and aims for lower accuracy and the ability to generalize to unseen areas [27]. Existing geolocation approaches can be categorized by whether they treat geolocation as an image retrieval problem, a classification problem, or both.

Image Retrieval-Based Approaches. A straightforward method for image localization is to find the most similar image in a large image database and predict its location [26]. The first successful approaches involved retrieving the nearest image in a space of handcrafted features such as color histograms [26], gist features [50], or textons [43]. It was later improved with SIFT features and support vector machines [28]. Deep features further boosted the performance of these approaches [70]. While such models typically exhibit high performance given a large and dense enough image database, they do not involve representation learning. Consequently, unless provided with robust features, they may perform poorly in sparsely represented or dynamically changing environments.

Classification-Based Approaches. Geolocation can also be approached as a classification problem by discretizing latitude and longitude coordinates. The choice of partition is critical, ranging from regular [72], adaptive [14], semantic-

driven [65], combinatorial [62], administrative [24, 55], and hierarchical [14, 70] partitions. Classification-based methods must strike a delicate balance between the quantity and size of cells; if the discretization is too coarse, the performance will be limited, while too many small cells may not have enough samples for learning-based methods. Furthermore, a typical classification loss such as cross-entropy does not incorporate the distance between regions: confusing two adjacent cells is equivalent to mistaking the continent.

Hybrid Approaches. Retrieval and classification approaches can be combined to overcome the limitations of discretization. This can be achieved using ranking losses [70] or contrastive objectives [38]. Haas *et al.* [24] follow a classification-then-regression approach based on prototype networks. Finally, Izbicki *et al.* [32] go beyond single-location prediction by estimating probability distributions based on spherical Gaussians.

3. OpenStreetView-5M

OpenStreetView-5M establishes a new open benchmark for geolocation by providing a large, open, and clean dataset. The Appendix details the construction of the dataset. As detailed below, OpenStreetView-5M improves upon several limitations of current geolocation datasets.

Scale. Deep neural networks have historically been selected over other machine learning methods because they benefit from larger amounts of data. OSV-5M consists of 4,894,685 training and 210,122 test images, with a height of 512 pixels and an average width of 792 ± 127 pixels.

Scope. Many geolocation datasets are restricted to a few cities [46, 77] or are significantly biased towards the Western world [36]. In contrast, OpenStreetView-5M images are uniformly sampled on the globe, covering 70k cities and 225 countries and territories, as shown in Figure 3. The distribution of test images across countries has a normalized entropy of 0.78 [73, Eq. 19], suggesting high diversity. Our train set has a normalized entropy of 0.67, which is comparable to the entropy of the distribution of the countries’ area (0.71).

Access. OpenStreetView-5M is based on the crowd-sourced street view images of Mapillary [4] which follow the CC-BY-SA license: free of use with attribution [2].

Quality Evaluation. We estimate through manual inspection of 4500 images that 96.1% ($\pm 0.57\%$) of the images in the OpenStreetView-5M dataset are localizable, with a 95% confidence level [31, Chap. 8]. Among the weakly or non-localizable images, 70% (2.7% total) are low-quality: under- or over-exposed, blurry, or rotated; 30% (1.2% total) are poorly framed, indoor, or in tunnels.

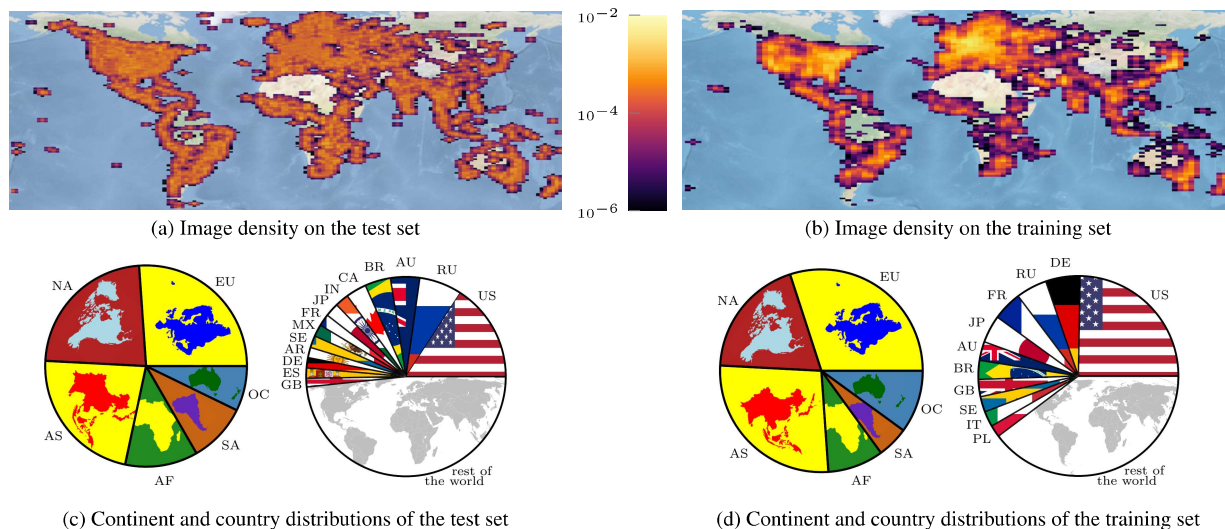


Figure 3. **OpenStreetView-5M**. Image density and proportions per country and continent for the train and test sets. To ensure an unbiased evaluation, we prioritize the uniformity of the test set’s distribution across the globe over the training set distribution.

Spatial Separation. Without carefully enforcing the spatial separation between train and test images, geolocation can reduce to place-recognition. As our goal is to assess the capacity of models to learn robust geographical representations, we ensure that no image in the OSV-5M training set lies within a 1km radius of any image in the test set.

Sequence Separation. Street-view images are typically acquired by a limited number of camera sensors mounted on the top or front of a small fleet of vehicles assigned to a given region. This correlation between location, cars, and sensors can be exploited to simplify the geolocation task. Notoriously, players of the web-based geolocation game GeoGuessr [3] can locate images from Ghana by spotting a piece of duct tape placed on the corner of the roof rack of the Google Street View car [6]. OpenStreetView-5M tries to avoid this pitfall by ensuring that no image sequence (a continuous series of images acquired by the same user) appears in both training and test sets. While this might not prevent images taken with the same vehicle on different days from being in both sets, it limits such occurrences.

Metadata. Rich metadata beyond geographical coordinates can improve the robustness and versatility of geolocation models. Each image in our dataset is associated with four tiers of administrative data: country, region (*e.g.*, state), area (*e.g.*, county), and the nearest city [6]. Note that areas are not defined for one-third of the dataset. We also associate each image with a set of additional information: land cover, climate, soil type, the driving side, and distance to the sea where the image was taken. See the Appendix for more

details on these attributes.

4. Benchmark

We use OSV-5M to benchmark supervised deep learning approaches in the context of visual geolocation. We first present our evaluation metrics (Section 4.1) and framework (Section 4.2). We then explore several design choices, starting with the image encoder backbone (Section 4.3), the prediction objective (Section 4.4), the fine-tuning strategy (Section 4.5), and the choices of contrastive losses (Section 4.6). In each experiment, we select the top-performing designs and integrate them into a *combined model*, which we evaluate and analyze in Section 4.7.

4.1. Evaluation Metrics.

We denote the space of images by \mathcal{I} and the span of longitude and latitude coordinates by $\mathcal{C} = [-180, 180] \times [-90, 90]$. Our objective is to design a model that maps an image from \mathcal{I} to its corresponding location in \mathcal{C} . We measure the accuracy of predicted location across geolocation models with three complementary sets of metrics:

- *Haversine distance* [68] δ , between predicted and ground truth image locations;
- *Geoscore*, based on the famous GeoGuessr game [3], defined as $5000 \exp(-\delta/1492.7)$ [24];
- Accuracy of predicted locations across administrative boundaries: country, region, area, and city.

While the average distance between predictions and ground truth is sensitive to outliers (*i.e.*, a few poor predictions can significantly undermine an otherwise high-performing algorithm), the accuracy metric based on administrative borders can avoid this issue. However, this metric

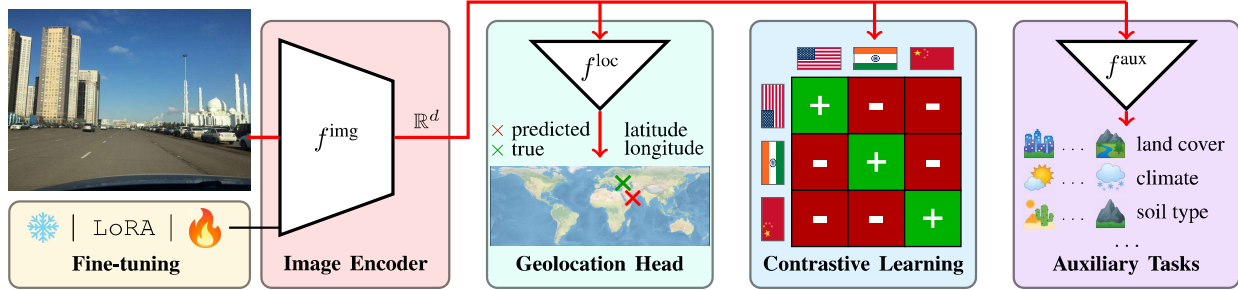


Figure 4. **Visual Geolocation Model.** We propose a simple and versatile framework for visual geolocation and explore the impact of various components of this approach in train-test performance on OpenStreetView-5M. Starting from the left, the input image is converted to a vector representation by an image encoder f^{img} (red). Then a geolocation head f^{loc} maps this vector to a set of geographical predictions (mint). Then a contrastive objective is potentially added (cyan), as well as auxiliary targets to learn better representations for geolocation (lila). We also consider various parameter fine-tuning strategies for training our image encoder, by freezing all or part of f^{img} (yellow).

Table 2. **Impact of Image Encoder.** Several pretrained backbones are evaluated in OpenStreetView-5M. We outline the influence of various architectures, pretraining strategies, and datasets. Best scores are highlighted in **bold**. We denote closed datasets with †.

Architecture	Size ($\times 10^6$)	Pretraining		Train. time (in h)	Geoscore \uparrow	Distance \downarrow	Classification accuracy \uparrow			
		Objective	Dataset				Country	Region	Area	City
1 ViT-B-32	88	CLIP	LAION-2B	22	2052	2992	35.7	7.0	0.5	0.3
2 ResNet50	23	Classification	ImageNet-1k	45	1260	4171	20.8	3.0	0.2	0.1
3 ViT-L-14	300	DINOv2	DINOv2†	316	2530	2233	46.9	10.7	0.7	0.3
4 ViT-L-14	300	CLIP	LAION-2B	206	2474	2358	44.8	10.6	0.8	0.2
5 ViT-L-14	300	CLIP	DATA_COMP	206	2719	1964	50.6	12.8	1.0	0.4
6 ViT-L-14	300	CLIP	Meta-CLIP	206	2724	1888	49.7	12.7	1.1	0.4
7 ViT-L-14	300	CLIP	OpenAI†	206	2888	1688	53.3	14.6	1.2	0.5
8 ViT-L-14	300	StreetCLIP	OpenAI† + GSV†	206	3028	1481	56.5	16.3	1.5	0.7
9 ViT-bigG-14	1800	CLIP	LAION-2B	900	2878	1766	53.4	15.0	1.3	0.5

can be too lenient for large divisions or arbitrarily punitive for small ones. The Geoscore offers a compromise by rewarding precise predictions without being overly sensitive to large but rare errors.

4.2. Framework

The models evaluated in this benchmark follow a consistent architecture, represented in Figure 4. All considered networks contain the two following modules:

- the image encoder $f^{\text{img}} : \mathcal{I} \mapsto \mathbb{R}^d$, which maps an image to a d -dimensional vector;
- the geolocation head $f^{\text{loc}} : \mathbb{R}^d \mapsto \mathcal{C}$, which maps this vector to geographic coordinates.

Implementation details. Unless stated otherwise, f^{img} is always a pretrained and frozen CLIP ViT-B/32 model [57] with $d = 768$ and f^{loc} is a Multilayer Perceptron (MLP) with GroupNorms [75]. This base model directly regresses geographical coordinates and uses the L_1 norm as loss function. The model is trained with a batch size of 512 images for 30 epochs (260k iterations) with a fixed learning rate of 2×10^{-4} . Throughout the paper we will denote in blue the frozen base model, in orange its fine-tuned version, and in green the model combining all top-performing designs.

4.3. Image Encoder

We first benchmark various architectures for the image encoder module f^{img} , with varying backbones, and pretraining strategies and datasets:

- *Architecture.* We test a standard ResNet50 [29], and modern ViTs [17] of multiple sizes (B-32, L-14, and bigG-14).
- *Pretraining.* We consider different types of pretraining objectives, including classification on ImageNet, self-supervised pretraining DINOv2 [52], text supervision CLIP [57], as well as StreetCLIP [23], which is finetuned specifically for geolocation.
- *Dataset.* We consider several pretraining datasets, including LAION-2B [61], DATA_COMP [20], Meta-CLIP [76], and the proprietary datasets of DINOv2, OpenAI, and StreetCLIP [23].

Analysis. Our experimental results are presented in Table 2. Here, we summarize several key takeaways:

- *Model Size.* As shown in Rows 1, 2, 4, and 9 of Table 2, there is a direct correlation between the size of the image encoder and its geolocation performance. The large ViT, bigG-14 model with 1.8 billion parameters (Row 9) improves significantly on the performance of its smaller versions. As

Table 3. **Prediction Modules.** We report the performance of various prediction models and objectives. QuadTrees, hierarchical supervision, and hybrid models all significantly improve on direct regression or classification with administrative borders. We underline the accuracy for divisions that the method is specifically trained to categorize.

		Number classes	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow			
					country	region	area	city
Reg.	Coord.	-	2052	2992	35.7	7.0	0.5	0.3
	Sin/cos	-	1192	4797	13.6	2.1	0.1	0.0
Classification	Country	222	2263	2981	<u>56.3</u>	-	-	-
	Region	2.8k	2683	2858	57.0	<u>30.2</u>	-	-
	Area	9.3k	1935	4454	36.3	19.7	<u>8.8</u>	-
	City	69.8k	2600	3217	52.2	28.5	7.3	<u>4.9</u>
	+ hierarchy	69.8k	2868	2768	<u>58.2</u>	<u>34.3</u>	<u>9.6</u>	<u>6.0</u>
	QuadTree	11.0k	2772	2832	54.8	<u>27.7</u>	5.4	2.8
	+ hierarchy	11.0k	2890	2654	57.4	29.9	5.9	2.9
Hybrid	11.0k	3036	2518	60.8	36.3	9.5	5.7	

the size of models correlates with their training time, we select ViT-L-14 as the best compromise.

- *Pretraining.* As seen in rows 3, 7, and 8, CLIP pretraining leads to better results than DINO or image classification. We thus focus on the latter for further comparisons.

- *Dataset.* Rows 4 to 8 show the significant impact of the choice of pretraining datasets. The geolocation-oriented StreetCLIP (row 8) leads to the best results, followed by OpenAI’s CLIP (row 7). However, both datasets are not open access. We choose DATA.COMP (row 5) as the best open-source dataset for its slightly better country classification rate compared to Meta-CLIP (row 6).

4.4. Prediction Head

We examine three different possible supervision schemes for the geolocation head f^{loc} : regression, classification (including hierarchical classification), and a hybrid approach.

Regression. We start with the most straightforward approach: f^{loc} directly regresses coordinates in \mathcal{C} . We train an MLP supervised with the L_1 loss between true and predicted coordinates. To account for the periodicity of the latitude, we also test an approach where we regress instead the cosine and sine of the longitude and latitude and then recover the real coordinates with trigonometry [42].

Classification. We divide the train set into a set \mathcal{K} of K divisions, such as countries, regions, areas, and cities, which amount to $K = 222, 2.8k, 9.3k,$ and $69.8k$, respectively. As some administrative borders can have vastly different sizes, we also consider an adaptive partition with a QuadTree of depth 10 and maximum leaf size of 1000, corresponding to 11k cells. We then train a classifier $f^{\text{classif}} : \mathbb{R}^d \mapsto \mathbb{R}^K$ which maps an image representation to the probability that the image was taken in each division. Then, to predict the final geographic location, we define f^{lookup} , which associates

Table 4. **Parameter Fine-tuning Strategies.** We compare the performance of different parameter fine-tuning strategies, in terms of performance, number of parameters, and training time.

	Param. (M)	Train. time	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow			
					country	region	area	city
Frozen	0.6	22	2052	2992	35.7	7.0	0.5	0.3
LoRA-32	2.4	44	2101	2760	36.7	6.4	0.4	0.0
Last block	7.7	26	2587	2372	46.7	12.9	1.0	0.5
Fine-tuning	88.0	132	2893	2085	54.9	19.1	1.6	0.8

each division with the average location of its training images: $f^{\text{lookup}} : \mathcal{K} \mapsto \mathcal{C}$. The predicted geolocation can be summarized as: $f^{\text{loc}} = f^{\text{lookup}} \circ \arg \max f^{\text{classif}}$.

In our implementation, f^{classif} is an MLP trained with cross-entropy, while f^{lookup} is a look-up table obtained directly from the training set.

Hierarchical Supervision. We can exploit the nested nature of the administrative divisions and QuadTree cells to supervise all levels simultaneously [48, 70]. More precisely, we predict a probability vector *at the finest resolution* (either city or maximum depth of the QuadTree), which we aggregate recursively to obtain predictions at all levels. We can now supervise with a cross-entropy term for each level.

Hybrid Approach. Inspired by approaches that combine both classification and retrieval [24, 70], we perform regression and classification in a two-step approach. Given the output of our QuadTree classifier $f^{\text{classif}} : \mathbb{R}^d \rightarrow \mathbb{R}^K$, we define $f^{\text{relative}} : \mathbb{R}^d \rightarrow [-1, 1]^{2K}$ that outputs the relative coordinates of the predicted location inside each cell k . We scale these values such that $(0, 0)$ points to the centroid of the training images in the cell and $[-1, 1]^2$ spans the entire bounding box. Using the cell prediction of the classifier f^{classif} and the relative position from f^{relative} , we can predict the location of the image with sub-cell precision.

We train f^{classif} with the cross-entropy, and f^{relative} with the L_2 loss between the predicted and true relative coordinates on the division that contains the true location.

Analysis. We report the performance of different prediction heads in Table 3, and make the following observations:

- *Regression.* Predicting sines and cosines does not improve the regression model’s performance. We hypothesize that this is due to the non-linearity of the trigonometric formula.

- *Classification.* Classification methods generally perform well in Geoscore and starkly improve their respective classification rates, e.g. +23.2% region accuracy for the region classifier compared to the regression model. However, their influence on the average error distance is smaller. Coarse partitions, like countries, are limited by the low precision of f^{lookup} . Inversely, overly refined partitions such as cities lead to a more challenging classification setting where most labels

Table 5. **Contrastive Learning.** We report the impact of adding a contrastive objective to our model, defined by various notions of positive matches between images.

	Pairs	Geoscore \uparrow	Distance \downarrow	Classification accuracy \uparrow			
				country	region	area	city
	no contrastive	2893	2085	54.9	19.1	1.6	0.8
geographic	country	2903	2005	66.8	13.7	0.7	0.3
	region	3028	2131	60.0	33.3	2.9	1.0
	area	2376	2886	43.7	18.9	3.7	1.2
	city	2912	2209	56.3	24.5	3.2	1.2
	cell	2891	2310	55.9	25.4	3.5	1.3
text-based		2812	2171	66.0	13.0	0.7	0.2

have only a few training examples. QuadTree-constructed labels achieve performance close to the administrative division-based classifier across all levels, *e.g.* 54.8% vs. 56.3% for countries and 27.7% vs. 30.2% for regions. This compounds into an overall better performance, which shows that adapting the granularity of the label distribution according to the image density appears to be a successful heuristic.

- *Hierarchical & Hybrid.* Supervising on all levels simultaneously significantly improves the prediction. Hybrid methods bridge the gap between classification and regression, yielding high precision without relying on very fine-grained partitions. These results validate the underlying spatial hierarchical nature of geographical data [67]. We select both hybrid and hierarchical designs for the combined model.

4.5. Parameter Fine-tuning

We evaluate different fine-tuning strategies to quantify the impact of learning dedicated features for geolocation. In all configurations, we learn f^{loc} from random weight, and f^{img} is fine-tuned as follows:

- *Frozen.* f^{img} is initialized with pretrained weights and remains frozen.
- *LoRA-32.* We fine-tune f^{img} using Low Rank Adaption [30] and a rank of 32 (more values in supplementary).
- *Last block.* We unfreeze the last transformer block of f^{img} , responsible for producing the image embedding.
- *Fine-tuning.* We fine-tune all parameters of f^{img} .

Analysis. In Table 4, we report the impact of different fine-tuning strategies. Training only the last transformer block instead of using LoRA leads to a ten times larger Geoscore improvement in only half the training time. This suggests that pretrained models can extract relevant patch embeddings, while image encoding must be significantly adapted for geolocation. Fine-tuning the entire network leads to an even larger improvement but a five-fold increase in training time. However, the resulting performance is comparable to the frozen ViT-bigG-14 shown in Table 2 and trains 9 times faster. We select the fine-tuning configuration as the top-performing approach and denote it in **orange**.

Table 6. **Combined Model.** We report the improvements brought by each top-performing design choice and their combination and compare them with baselines and competing approaches.

	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow			
			country	region	area	city
Base model	2052	2992	35.7	7.0	0.5	0.3
ViT-L-14 DC	+ 667	- 1028	+14.9	+ 5.8	+0.5	+ 0.1
QuadTree	+ 720	- 160	+19.1	+20.7	+5.4	+ 2.5
Hybrid	+ 264	- 314	+ 6.0	+ 8.6	+4.5	+ 2.9
Hierarchical	+ 118	- 178	+ 2.6	+ 0.2	+0.5	+ 0.1
Fine-tuning	+ 841	- 907	+19.2	+12.1	+1.1	+ 0.5
Region contrast.	+ 135	+ 46	- 5.1	+14.2	+2.1	+ 0.2
Combined model	+1309	- 1178	+32.3	+32.4	+9.8	+ 5.6
	3361	1814	68.0	39.4	10.3	5.9
Random	328	8724	20.0	2.0	0.0	0.0
Human Evaluation	1009	6407	48.9	12.2	3.0	0.0
GeoEstimator [48]	3331	2308	66.8	39.4	18.4	4.2
StreetCLIP 0-shot [23]	2273	2854	38.4	20.8	9.9	14.8

4.6. Contrastive Objectives

Contrastive learning builds positive and negative sample pairs from the training set and pushes representations of positive pairs close to each other while contrasting negative ones [12, 13]. Positive pairs can be formed within the same modality, such as different views of an object, or across modalities, such as images and captions. In the geolocation context, we propose two approaches to construct such pairs:

- *Geographic.* We match images if they are within the same administrative division: countries, regions, areas, cities, or QuadTree cells. We modify the dataloader to ensure each image is part of at least one positive pair. Contrary to Haas *et al.* [23], we use the multi-positive MIL-NCE loss [45] as our contrastive objective to account for images in several positive pairs, *e.g.* in the same country.

- *Text-Based.* Similar to Haas *et al.* [23], we pair each image with a textual description of its location formed as the following string: “An image of the city of \$CITY, in the area of \$AREA, in the region of \$REGION, in \$COUNTRY.”.

Analysis. In Table 5, we measure the impact on the **fine-tuned model** of different approaches for constructing contrastive pairs. We observe a consistent improvement in terms of performance when building positive pairs with regions, which may be the division most likely to present unique and homogeneous visual and cultural identities. In contrast, areas appear to hurt the performance when used contrastively. Overall, contrastive learning yields a much higher country and region classification rate compared to the classification-based approaches of Table 3, suggesting that encouraging geographically consistent representations is advantageous for geolocation. We also observe that using text as a proxy when geographically consistent pairs exist is not beneficial.

4.7. Combined Model

Summarizing our previous exploration and analysis, we combine the most impactful design choices for each experiment

Table 7. **Nearest Neighbors.** We report the performance of nearest neighbor retrieval using different encoders.

	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow			
			country	region	area	city
CLIP-ViT-B32-LAION	2511	3455	49.3	29.6	1.9	13.1
DINov2	2994	2542	61.1	37.1	22.9	16.4
CLIP-ViT-L14-DATACOMP	3201	2047	64.5	38.4	23.3	16.6
CLIP-ViT-L14-OpenAI	3545	1458	72.8	44.4	27.5	19.3
StreetCLIP	3597	1386	73.4	45.8	28.4	19.9
Combined model	2734	2608	54.9	24.5	13.6	9.4

into a strong geolocation model, denoted in green: ViT-L-14 backbone pretrained on DATA_COMP, QuadTree partition with hybrid prediction and hierarchical supervision, fully fine-tuned with a region-contrastive loss. As shown in Table 6, this model starkly improves on the base model, with an increase of +1309 in Geoscore, an average distance reduced by 45%, and significantly better accuracy at all levels of administrative divisions.

Analysis. In Table 6, we compare the performance of our combined model to a random baseline (select the location of a random image in the training set) and a human performance obtained by asking 80 annotators to guess the locations of the same 50 images randomly sampled from the test set [44]. Despite the difficulty of the task, the average annotator’s performance is significantly better than chance. Our baseline model, and more substantially our combined model, far surpasses the accuracy of annotators. We also evaluate two state-of-the-art geolocation models: StreetCLIP [23] evaluated in zero-shot using the text string given in Section 4.6, and the GeoEstimator model [48] fine-tuned on our training set. As both models are designed for geolocation, they yield good performance. Owing to its bespoke geocells, GeoEstimator reaches the highest accuracy for area classification, illustrating the benefit of architectures with built-in geographical priors. See the appendix for further experiments, notably on the impact of auxiliary variables.

Nearest Neighbor. We perform retrieval by matching each image from the test set with an image from the train set based on the cosine distance between the features of each image encoder. We perform approximate matching with the FAISS algorithm [35] through the AutoFAISS package [1], without re-ranking [34, 56]. As reported in Table 7, retrieval methods trained through contrastive learning exhibit high performance. However, the supervision of our combined model based on geographic coordinates and cells does not enhance its retrieval performance. In fact, its retrieval score is lower than that of its pretrained image encoder. These findings are consistent with observations that fine-tuning self-supervised models decreases retrieval performance [74].

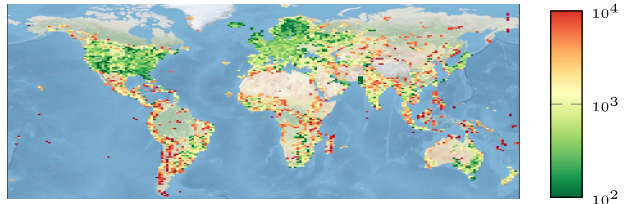


Figure 5. **Spatial Distribution of Errors.** We plot the average prediction error of the combined model in km across the globe.

Error Distribution. We report in Figure 5 a heatmap of the average error distance. Areas sparsely populated with training images, such as South America, have a significantly higher error rate. We report a Pearson correlation coefficient of -0.25 between image density and error, suggesting that image density is not the only factor in the mistakes of our proposed model. See Figure 6 for a visualization of the error distribution. Over half of the combined model’s predictions are within 250km of the true image locations.

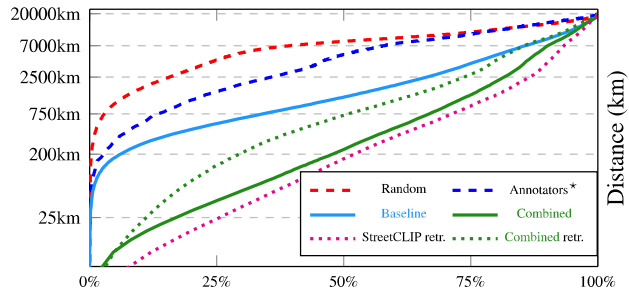


Figure 6. **Error Distribution.** Proportion of predictions within a set distance in the test set. * evaluated on 50 images only.

5. Conclusion

We introduced a new open-access street view dataset of unprecedented size and quality, enabling the consistent training and evaluation of global geolocation models for the first time. Through an extensive experimental framework, we demonstrate that our dataset is a competitive benchmark for developing and evaluating general and bespoke state-of-the-art computer vision approaches for geolocation. Through its scale and quality, we expect OSV-5M to also be useful for self-supervised learning and generative modeling, valuable tasks beyond the scope of visual geolocation.

Acknowledgements. OSV-5M was made possible through the generous support of the Mapillary team, which helped us navigate their vast street view image database. Our work was supported by the ANR project READY3D ANR-19-CE23-0007, and the HPC resources of IDRIS under the allocation AD011014719 made by GENCI. We thank Valérie Gouet for her valuable feedback.

References

- [1] AutoFAISS. <https://github.com/criteo/autofaiss>. Accessed: 2023-10-10. **8**
- [2] CC BY-SA 2.0 DEED: Attribution-ShareAlike 2.0 Generic. <https://creativecommons.org/licenses/by-sa/2.0/deed.en>. Accessed: 2023-10-10. **3, 21**
- [3] GeoGuessr. <https://www.geoguessr.com/>. Accessed: 2023-10-10. **4**
- [4] Mapillary. <https://www.mapillary.com/>. Accessed: 2023-10-10. **2, 3, 12**
- [5] Multimediacommons - yfcc100m core dataset. <https://multimediacommons.wordpress.com/yfcc100m-core-dataset/>. Accessed: 2023-10-10. **13**
- [6] Plonkit guide to Ghana. <https://www.plonkit.net/ghana>. Accessed: 2023-10-10. **4**
- [7] Reverse Geocoder. pypi.org/project/reverse_geocoder/. Accessed: 2023-10-10. **19**
- [8] Yfcc100m. <https://gitlab.com/jfolz/yfcc100m>. Accessed: 2023-10-10. **13**
- [9] Manuel López Antequera, Pau Gargallo, Markus Hofinger, Samuel Rota Bulò, Yubin Kuang, and Peter Kotschieder. Mapillary planet-scale depth dataset. In *ECCV*, 2020. **2**
- [10] Hylke E Beck, Niklaus E Zimmermann, Tim R McVicar, Noemi Vergopolan, Alexis Berg, and Eric F Wood. Present and future Köppen-Geiger climate classification maps at 1-km resolution. *Scientific data*, 2018. **15, 19**
- [11] David M Chen, Georges Baatz, Kevin Köser, Sam S Tsai, Ramakrishna Vedantham, Timo Pylvänäinen, Kimmo Roimela, Xin Chen, Jeff Bach, Marc Pollefeys, et al. City-scale landmark identification on mobile devices. In *CVPR*, 2011. **1**
- [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICLR*, 2020. **7**
- [13] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, 2005. **7**
- [14] 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 *CVPR*, 2023. **1, 2, 3, 14**
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Image Net: A large-scale hierarchical image database. In *CVPR*, 2009. **1**
- [16] Mingyu Ding, Zhe Wang, Jiankai Sun, Jianping Shi, and Ping Luo. CamNet: Coarse-to-fine retrieval for camera re-localization. In *ICCV*, 2019. **3**
- [17] 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. In *ICLR*, 2020. **5**
- [18] Christian Ertler, Jerneja Mislej, Tobias Ollmann, Lorenzo Porzi, Gerhard Neuhold, and Yubin Kuang. The Mapillary traffic sign dataset for detection and classification on a global scale. In *ECCV*, 2020. **2**
- [19] Müller-Budack et al. Geolocation estimation of photos using a hierarchical model and scene classification. In *ECCV*. **14**
- [20] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. DataComp: In search of the next generation of multimodal datasets. *NeurIPS Dataset and Benchmark*, 2023. **5**
- [21] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 2021. **12**
- [22] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *ICLR*, 2018. **12, 20**
- [23] Lukas Haas, Silas Alberti, and Michal Skreta. Learning generalized zero-shot learners for open-domain image geolocalization, 2023. **1, 2, 3, 5, 7, 8, 14**
- [24] Lukas Haas, Silas Alberti, and Michal Skreta. PIGEON: Predicting image geolocations. *arXiv preprint arXiv:2307.05845*, 2023. **3, 4, 6, 14**
- [25] Stephen Hausler, Sourav Garg, Ming Xu, Michael Milford, and Tobias Fischer. Patch-NetVLAD: multi-scale fusion of locally-global descriptors for place recognition. In *CVPR*, 2021. **1**
- [26] James Hays and Alexei A Efros. Im2gps: estimating geographic information from a single image. In *CVPR*, 2008. **1, 3**
- [27] James Hays and Alexei A Efros. Large-scale image geolocalization. *Multimodal location estimation of videos and images*, 2015. **3**
- [28] James Hays and Alexei A Efros. Large-scale image geolocalization. *Multimodal location estimation of videos and images*, 2015. **3**
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. **5**
- [30] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2021. **7, 16**
- [31] Barbara Illowsky and Susan Dean. Introductory statistics. 2018. **3**
- [32] Mike Izbicki, Evangelos E Papalexakis, and Vassilis J Tsotras. Exploiting the Earth's spherical geometry to geolocate images. In *MLKDD*, 2020. **1, 2, 3, 13**
- [33] Ara Jafarzadeh, Manuel López Antequera, Pau Gargallo, Yubin Kuang, Carl Toft, Fredrik Kahl, and Torsten Sattler. CrowdDriven: A new challenging dataset for outdoor visual localization. In *ICCV*. **2**
- [34] Herve Jegou, Hedi Harzallah, and Cordelia Schmid. A contextual dissimilarity measure for accurate and efficient image search. In *CVPR*, 2007. **8**
- [35] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 2019. **8, 18**

- [36] Sebastian Kalkowski, Christian Schulze, Andreas Dengel, and Damian Borth. Real-time analysis and visualization of the YFCC100M dataset. In *Workshop on community-organized multimodal mining: opportunities for novel solutions*, 2015. 2, 3, 13
- [37] Alex Kendall, Matthew Grimes, and Roberto Cipolla. PoseNet: A convolutional network for real-time 6-dof camera relocalization. In *ICCV*, 2015. 3
- [38] Giorgos Kordopatis-Zilos, Panagiotis Galopoulos, Symeon Papadopoulos, and Ioannis Kompatsiaris. Leveraging EfficientNet and contrastive learning for accurate global-scale location estimation. In *International Conference on Multimedia Retrieval*, 2021. 3
- [39] Martha Larson, Mohammad Soleymani, Guillaume Gravier, Bogdan Ionescu, and Gareth JF Jones. The benchmarking initiative for multimedia evaluation: MediaEval 2016. *IEEE MultiMedia*, 2017. 3
- [40] John Latham, Renato Cumani, Ilaria Rosati, and Mario Bloise. Global land cover share (GLC-SHARE) database beta-release version 1.0-2014. *FAO: Rome, Italy*, 2014. 15, 19
- [41] Grace Luo, Giscard Biamby, Trevor Darrell, Daniel Fried, and Anna Rohrbach. g^3 : Geolocation via guidebook grounding. *Findings of EMNLP*, 2022. 1, 2
- [42] Oisín Mac Aodha, Elijah Cole, and Pietro Perona. Presence-only geographical priors for fine-grained image classification. In *CVPR*, 2019. 6
- [43] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, 2001. 3
- [44] Sneha Mehta, Chris North, and Kurt Luther. An exploratory study of human performance in image geolocation tasks. In *HCOMP 2016 GroupSight Workshop on Human Computation for Image and Video Analysis*, volume 308, 2016. 1, 8
- [45] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *CVPR*, 2020. 7, 17
- [46] Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, Denis Teplyashin, Karl Moritz Hermann, Mateusz Malinowski, Matthew Koichi Grimes, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, et al. The StreetLearn environment and dataset. *arXiv preprint arXiv:1903.01292*, 2019. 2, 3
- [47] Hatem Mousselly-Sergieh, Daniel Watzinger, Bastian Huber, Mario Döller, Elöd Eged-Zsigmond, and Harald Kosch. World-wide scale geotagged image dataset for automatic image annotation and reverse geotagging. In *ACM multimedia systems*, 2014. 2, 3, 13
- [48] Eric Müller-Budack, Kader Pustu-Iren, and Ralph Ewerth. Geolocation estimation of photos using a hierarchical model and scene classification. In *ECCV*. 6, 7, 8
- [49] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kotschieder. The Mapillary Vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017. 2
- [50] Aude Oliva and Antonio Torralba. Building the gist of a scene: The role of global image features in recognition. *Progress in brain research*, 2006. 3
- [51] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 17, 18
- [52] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DinoV2: Learning robust visual features without supervision. *TMLR*, 2023. 5
- [53] Nathan Piasco, Désiré Sidibé, Valérie Gouet-Brunet, and Cédric Demonceaux. Improving image description with auxiliary modality for visual localization in challenging conditions. *International Journal of Computer Vision*, 2021. 3
- [54] Noé Pion, Martin Humenberger, Gabriela Csurka, Yohann Cabon, and Torsten Sattler. Benchmarking image retrieval for visual localization. In *3DV*, 2020. 3
- [55] 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 *ECCV*, 2022. 3, 14
- [56] Danfeng Qin, Stephan Gammeter, Lukas Bossard, Till Quack, and Luc Van Gool. Hello neighbor: Accurate object retrieval with k-reciprocal nearest neighbors. In *CVPR*, 2011. 8
- [57] Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 5
- [58] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihí Zelnik-Manor. Imagenet-21k pretraining for the masses. In *NeurIPS Datasets and Benchmarks Track*, 2021. 1
- [59] Pedro A Sanchez, Sonya Ahamed, Florence Carré, Alfred E Hartemink, Jonathan Hempel, Jeroen Huisig, Philippe Lagacherie, Alex B McBratney, Neil J McKenzie, Maria De Lourdes Mendonça-Santos, et al. Digital soil map of the world. *Science*, 2009. 15, 19
- [60] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, et al. Benchmarking 6-DoF outdoor visual localization in changing conditions. In *CVPR*, 2018. 3
- [61] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: An open large-scale dataset for training next generation image-text models. 2022. 5
- [62] Paul Hongsuck Seo, Tobias Weyand, Jack Sim, and Bohyung Han. CPLaNet: Enhancing image geolocation by combinatorial partitioning of maps. In *ECCV*, 2018. 3
- [63] Zachary Seymour, Karan Sikka, Han-Pang Chiu, Supun Samarasekera, and Rakesh Kumar. Semantically-aware attentive neural embeddings for image-based visual localization. *arXiv preprint arXiv:1812.03402*, 2018. 1
- [64] Jonas Theiner, Eric Müller-Budack, and Ralph Ewerth. Interpretable semantic photo geolocation. In *WACV*, 2022. 2, 3
- [65] Jonas Theiner, Eric Müller-Budack, and Ralph Ewerth. Interpretable semantic photo geolocation. In *WACV*, 2022. 3
- [66] 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*, 2016. [2](#), [3](#)
- [67] Waldo R Tobler. A computer movie simulating urban growth in the Detroit region. *Economic geography*, 1970. [7](#)
- [68] Glen Van Brummelen. *Heavenly mathematics: The forgotten art of spherical trigonometry*. Princeton University Press, 2012. [4](#)
- [69] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The iNaturalist species classification and detection dataset. In *CVPR*, 2018. [1](#)
- [70] Nam Vo, Nathan Jacobs, and James Hays. Revisiting IMG2GPS in the deep learning era. In *ICCV*, 2017. [1](#), [2](#), [3](#), [6](#)
- [71] Frederik Warburg, Soren Hauberg, Manuel Lopez-Antequera, Pau Gargallo, Yubin Kuang, and Javier Civera. Mapillary street-level sequences: A dataset for lifelong place recognition. In *CVPR*, 2020. [2](#)
- [72] Tobias Weyand, Ilya Kostrikov, and James Philbin. Planet-photo geolocation with convolutional neural networks. In *ECCV*, 2016. [3](#), [14](#)
- [73] Allen R Wilcox. Indices of qualitative variation. Technical report, Oak Ridge National Lab., Tenn., 1967. [3](#)
- [74] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *CVPR*, 2022. [8](#)
- [75] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018. [5](#), [17](#)
- [76] Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying CLIP data. *ICLR*, 2024. [5](#)
- [77] A.R. Zamir and M. Shah. Image geo-localization based on multiple nearest neighbor feature matching using generalized graphs. 2014. [2](#), [3](#)
- [78] Xiwu Zhang, Lei Wang, and Yan Su. Visual place recognition: A survey from deep learning perspective. *Pattern Recognition*, 2021. [3](#)

Appendix

This supplementary material starts by providing further details on the construction and analysis of our dataset OpenStreetView-5M in Section A, showcasing indicative samples in Figure A. Then, we provide additional experiments in Section B and qualitative results in Figure B. Finally, Section C further implementation details can be found and Section D outlines our Datasheet [21] for OpenStreetView-5M.

A. OpenStreetView-5M Dataset

OpenStreetView-5M is designed to achieve an open, large-scale, balanced, and global geographical coverage. Through the Mapillary API and the support of the Mapillary team, we gained access to the locations of all 1.8B images [4]. To provide a more manageable and better distributed dataset, we design a specific construction approach, presented in this section. The code to reproduce the treatment can be found at <https://github.com/gastruc/osv5m>.

A.1. Construction Approach

Sampling. We start by ensuring that regions with high image density are not disproportionately represented. We define a 100×100 m grid across the entire world and randomly choose one image per cell. Then, both the training and test sets are sampled with a weight proportional to the local image density raised to the power of -0.75 . Such a strategy balances density-based sampling (which tends to be biased towards urban centers) and area-based sampling (which might favor larger countries). We eliminate images from the test set that are either located within a 1km radius of any train image or share a sequence ID.

Handcrafted Filters. We apply a series of handcrafted filters to remove low-quality images

- *Blurriness.* Blurry images indicate low quality and potentially low localizability. We remove images whose average logarithmic magnitude spectrum is below 120dB.

- *Radiometry.* Certain images hosted on Mapillary are too dark to be meaningfully analyzed, while other have a distinct encoding errors giving them a purple tint. To remove those, we first filter out images whose average brightness (average value over pixels and RGB channels) is below 50. To handle purple images, we remove images for which over 50% of pixels meet the following criteria: $R > 60 \ \& \ G > 60 \ \& \ B < 50$.

- *Exposition.* The exposure of Dash-cam images can be badly exposed, for example, when they face the sun. To filter them, we remove images for which 70% of pixels have a brightness over 250 (overexposed) or under 5 (underexposed).

Rotation-Based Filtering. We perform a learning-based filtering based on a pretrained and frozen RotNet network [22]. This model learns self-supervised image representations by training for the pretext task of predicting a random rotation applied to an input image. Although it is used as a pretext task in the original paper, it becomes useful for filtering out images downloaded from Mapillary’s website that are incorrectly rotated. We use the pretrained network to infer the rotation of various images and then use the following filtering strategy depending on RotNet’s prediction:

- 0° (**96% of images**) For normal street view images the cues that signify an absence of rotation are multiple: the sky is up, and cars and pedestrians are upward. We keep these images unchanged.

- 180° (**4%**) Over 90% images predicted to be rotated by 180° are, in fact, actually upside down. We rotate all these images by a half-turn. For the images in the test, we perform an additional visual inspection to remove the small proportion of non-localizable images not removed by the previous filters.

- 90° or 270° (**0.2%**) Images predicted as rotated by a quarter-turn are in the vast majority taken indoors or in tunnels. We remove all such images from both the train and test set.

A.2. Discussion

Why Not Just Subsample YFCC100M? The wide adoption of YFCC100M, with its nearly 50 million geotagged images, might question the need for creating yet another geotagged image dataset. However, several compelling reasons justify creating OpenStreetView-5M instead of subsampling YFCC100M:

- *Data Distribution.* The images shared on Flickr do not aim to capture our world in an objective way, but instead focus on aesthetic and cultural value. For example, recognizable landmarks like the Eiffel Tower or the Louvre, are a cultural symbol of the city of Paris, yet they lack information that is useful in identifying other cities as French or even other streets as Parisian. Additionally, many images are renders or infographics. In contrast, OSV-5M only features dashcam pictures, that offer a consistent front-view perspective, that is more objective as it doesn’t focus on something specific, and thus may be more beneficial for learning visual geographical representations.

- *Localizability.* From a manual inspection of 1000 images we find that fewer than 10% ($\pm 1.3\%$, 95% confidence) of YFCC100M’s images are perceptually localizable. In stark contrast, OSV-5M boosts this perceptual localizability to a rate of 96.1% (± 0.57 , 95% confidence), making it a more



Figure A. **Images from OSV-5M.** The true locations can be found on the next page. The Mapillary users are credited in the subcaptions.

suitable candidate for a standard evaluation benchmark for global geolocation.

- *Geographical Bias.* Images in the YFCC100M dataset exhibit a high cultural bias towards the Western world, with over 35% of images from the US and nearly 70% from North America and Europe [36]. OSV-5M offers a more equitable global representation, as detailed in Figure 2 of the main paper.

- *Selection Challenges.* Subsampling YFCC100M based on metadata alone is ambiguous: 30% of images lack titles, 68% lack descriptions, 30% lack tags, and 50% lack geotagging. The tags “travel” and “nature” cover fewer than 2 million images. Using instead automated selection methods may inadvertently propagate existing biases, such as filtering street views of non-Western countries.

- *Persistence.* As happens with a lot of large research dataset, YFCC comes only as a collection of image URLs that need to be downloaded directly from Flickr. Such a dataset construction approach, even if the only feasible choice for very large datasets, is very volatile and can prevent future reproducibility. For example, 60% of the 2014 YFCC-split [47] was deleted by 2020 [32]. While YFCC100M

used to be hosted on Yahoo’s Webscope, this option is no longer available [5]. Instead users need to create an AWS account, that requires a credit card to acquire API credentials for downloading the data through a designated S3 bucket [8]. Even if no charge is applied, this setting may be prohibitive for academics or residents of certain countries. Also, due to the sensitive nature of the Flickr data, users need to make a formal request to download the dataset, something that isn’t needed for our dataset. Instead, OSV-5M ensures persistence, open and easy access for long-term and broad usage.

To summarize, YFCC100M is a vast and unstructured set of images, a subset of which may be well suited for localization and place recognition. However, the ambiguous localizability, geographical content, metadata, persistence, and access to its images highlight the need for a dedicated dataset like OSV-5M, specifically designed for the task of global visual geolocation.

Visible GeoTags. Due to the diversity in user input data, we found that a small percentage of images (< 5%) have a visible overlaid text on the bottom part that tags their location. This should be taken into consideration when con-

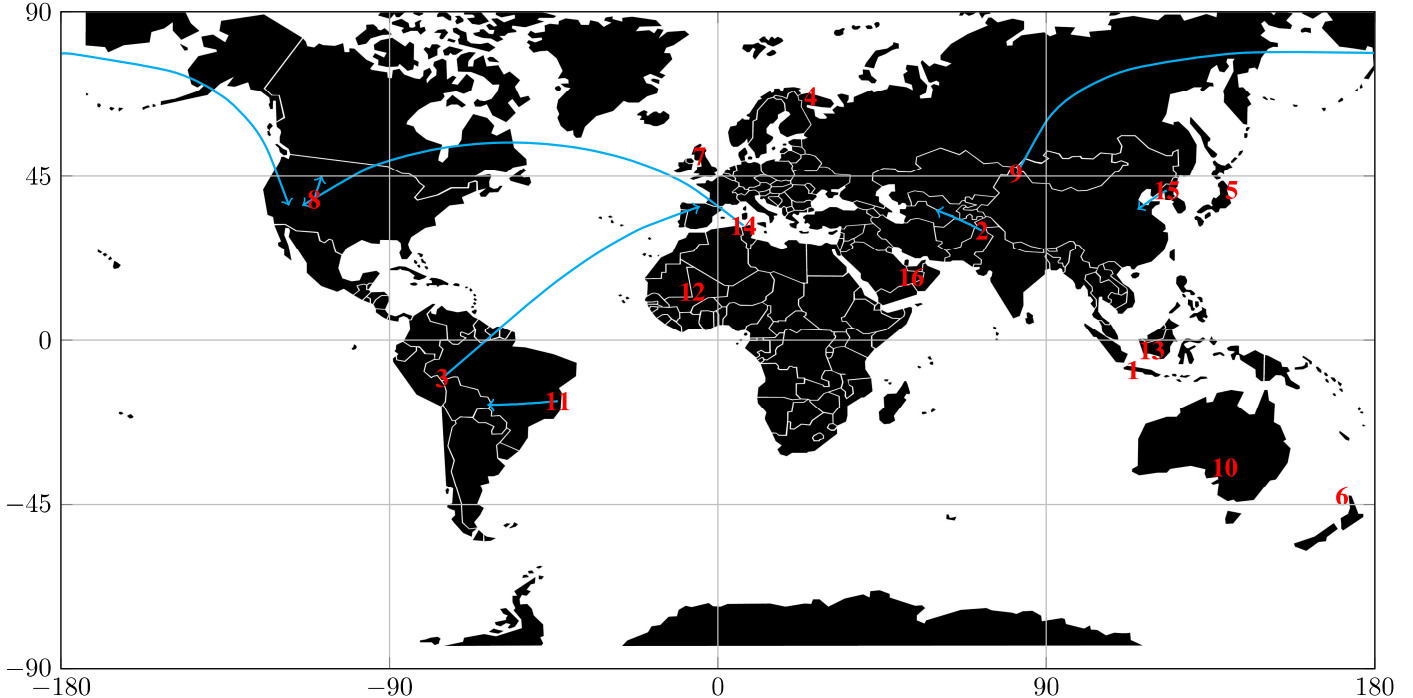


Figure B. **True Locations.** Location of the images of Figure A. With blue we visualize errors of the combined model that are superior to 500 km. Most of the images (9 out of 16) are predicted within 500km of where they were taken. We observe that two difficult images (9 and 14) are erroneously mapped to the US, which could be explained by the geographical bias of the training set.

structuring a benchmark for a future dataset. However, due to the standard ViTs resampling of images to 224×224 , these coordinates become indecipherable, as demonstrated in Figure C. We implement for our data loader the option to add a Gaussian blur with a width of 2 to the bottom 14 rows. When training and/or testing a baseline model with this blur, we observe only small and inconclusive differences in score: training without blurring but testing with it yielded slightly better results than both training and testing without the blur, yet training and testing with the blur produced inferior outcomes. This indicates that (i) the network is not able to read the coordinates, and (ii) the bottom rows do not contain critical geographic information. However, we recommend using the blur for methods that use higher-resolution models to obscure any potential location-specific details in the text.

Limitations. We list three main limitations of our OSV-5M dataset:

(i) *Geographical Bias in Training.* Due to our reliance crowd-sourced from Mapillary users, the distribution of locations is biased towards Western countries. We designed our test set to explicitly balance this distribution, but the training set remains affected by the number of selected images.

(ii) *User separation.* We successfully separated images from the same sequence between training and test sets. How-

ever, we could not separate images uploaded by the same user on different days, as the required metadata was not available at the time of the dataset construction.

(iii) *Resolution.* The dataset provides images with a vertical resolution of 512 pixels. This restricts the ability to zoom in and read distant texts, for example in street signs, potentially obscuring valuable visual cues. However, through our metadata users can access higher-resolution versions of all our images on the mapillary website.

Training SOTA methods on OSV-5M. Many state-of-the-art geolocation methods [14, 23, 24, 72] either rely on private datasets or lack publicly available code, that prevents their evaluation. In our main paper we evaluated the performance of the pretrained StreetCLIP model both for zero-shot retrieval (Tab 6 and Fig 6) and as a pretrained image encoder (Tab 2), yet the implementation required to fine-tune the model is not publicly available. Similarly, the complete training code of Translocator [55] is also not available. We managed to train the publicly available ISN model [19] on OSV-5M, achieving good performance which we attribute to its bespoke geocell module. The aforementioned difficulty in training and evaluating SOTA models show the clear need for open-source datasets and implementations of visual geolocation approaches, that our paper directly addresses.



Figure C. **Visible Geotagging.** A small minority of images ($< 5\%$) have visually overlaid geotags at their bottom left corner (1). For those images as resized by our data loader to 224×224 and as optionally blurred, we empirically measure to not provide any important information that the network can use to improve its performance .

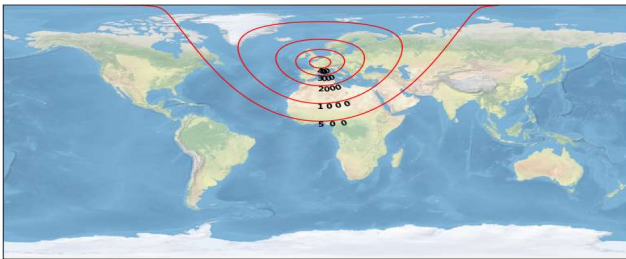


Figure D. **Geoscore.** From a point centered in Paris, red contours highlight level sets of the score along the earth’s spherical geometry.

Geoscore. In our paper geoscore is introduced as a better evaluation method as it strikes a balance between rewarding precision and not being oversensitive to outlier predictions. Let us consider, for example, a model which produces nine accurate predictions but fails on the tenth image, choosing New Zealand instead of Ireland, a 20 000km mistake. Contrast this with another model which consistently mispredicts by 2 000km. Solely examining the mean error might misleadingly favor the latter model, when the first one has a higher geographic proficiency. In terms of geoscore, the model with one major error would achieve an average score close to 4500, while the one that is consistently off would score 1300. In that way, geoscore provides a more intuitive way to compare the performance of models on our dataset. See Figure D for an illustration of Geoscore.

B. Additional Experiments

This section presents further results and analysis of our proposed framework.

Auxiliary Supervision. We start by evaluating the performance gained by learning to predict various auxiliary information. Based on their coordinates, we associate to each image of our dataset the following meta-data, according to its latitude and longitude coordinates:

- *Land Cover.* Relying on the Global Land Cover Share

Table A. **Auxiliary Variables.** We report the impact on geolocation performance of learning to predict various auxiliary variables. We also report the performance on the test set for each variable as the overall accuracy or the average error.

	Num of classes.	Perf. test	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow			
					country	region	area	city
no auxiliary	-	-	2893	2085	54.9	19.1	1.6	0.8
land cover	11	54.8	2821	2102	52.2	16.9	1.4	0.7
climate	31	58.3	2898	2022	53.7	18.8	1.7	0.8
soil type	15	47.7	2826	2111	52.4	17.6	1.5	0.7
driving side	1	94.6	2896	2025	54.5	18.7	1.6	0.7
dist to sea	-	543km	2870	2053	52.5	18.7	1.5	0.7
all	-	-	2910	1987	54.0	19.8	1.6	0.8

Database [40], we classify each image of our dataset into one of 11 land cover types, such as artificial, forest, or crops.

- *Climate.* We use recent Köppen-Geiger climate classification maps [10] to associate each image with a climate type among 31, such as tropical rainforest, arid steppe, or temperate with dry winter.

- *Soil Type.* Thanks to the Digital World Soil Map [59], we characterize the local soil with a 15 class nomenclature, such as acrisols, fluvisols, or ferralsols.

- *Driving Side.* We also add a binary indicator for whether a country uses left or right-hand traffic.

- *Distance to the Sea.* For all locations we compute the distance to their nearest sea.

The maps we used to extract land cover, climate, and soil types come in a resolution of 1 km (or 30 arc-seconds).

We use an MLP f^{aux} to predict the image’s metadata in addition to its coordinates. All categorical variables are supervised with the unweighted sum of cross-entropy terms, while the distance to the sea is supervised with the L1 loss. Adding auxiliary tasks encourages the model to focus on relevant geographical cues. As seen in Table A, we only observe a modest impact, indicating that the large train set of OSV-5M allows our model to already learn good latent variables for geolocation. It should be noted that our model

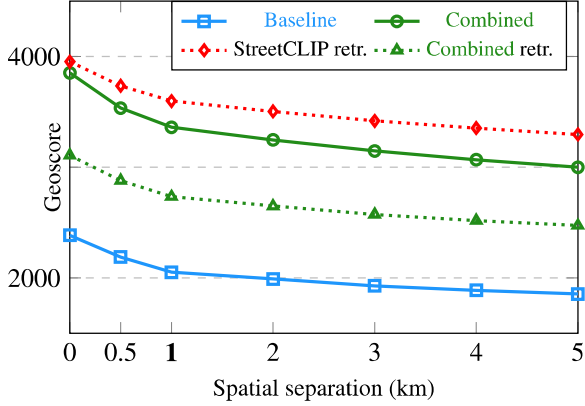


Figure E. **Spatial Separation.** We report the performance of different approaches for test sets defined by various separation radii for the train set.

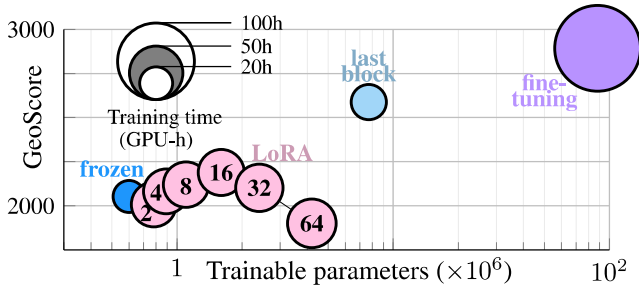


Figure F. **Effect of LoRA Bottleneck Width.** We report the performance of finetuning with LoRA of different bottleneck widths, in comparison to finetuning the last block, or the whole network. For each experiment, the marks' radius are proportional to the training time.

can perform accurate predictions for complex geographic variables in the test set, which may have some useful applications in itself.

Spatial Separation. We study the impact of the radius of spatial separation between the train and the test set. We do this by creating test sets along different radii of separation from the training set: 0m (488k images), 500m (294k), 1km (210k), 2km (166k), 3km (136k), 4km (117k) and 5km (107k). As observed in Figure E, all methods, including retrieval-based approaches, are equally affected by this phenomenon, indicating that, as expected, the problem of global geolocation becomes harder as the separation radius increases. This allows us to define different versions of our test set tiered by difficulty. In particular, if we remove the separation between train and test makes the task becomes significantly easier: 3952 geoscore for StreetCLIP in retrieval mode and 3852 for our best model, corresponding to an average distance error of 1191km.



Figure G. **Erroneous Predictions.** Images that are consistently predicted wrongly despite being sampled from areas with relatively high density of training images.



Figure H. **Attention Maps.** We visualize the self-attention maps of the [CLS] token of the last layer of the image encoder of the combined model. We show the mean across all heads in (2), and manually selected an interesting layer in (3).

LoRA. Fig F shows the results with different widths of the LoRA bottleneck, ranging from 2 to 64. We share similar observations with the LoRA paper [30, 7.2]: higher ranks do not increase or even slightly decrease performance. Un-freezing the last transformer block remains more efficient in terms of training time, and fine-tuning the entire model leads to even better performance.

Erroneous Predictions. In Fig G we illustrate some sources of geolocation errors not related to the density of training images. These include landscapes that are: (i) similar between very distant countries (Fig G (a,b)), or (ii) any key information is far away from the camera (Fig G (b,c)), or are (iii) monotonous and nearly featureless (Fig G (c)).

Humans and Baselines. We compare in Table B our models against two random baselines: selecting randomly a location on the map or the location of a random image from the training set. We also construct an Annotator Ensemble Oracle by selecting the most accurate prediction for each image from all annotators. Our baseline model, and more substantially our combined model, far surpasses the accu-

Table B. **Annotator Performance.** We report the average performance of 80 annotators on a subset of 50 images.

	Geo \uparrow score	Dis \downarrow tance	Classification accuracy \uparrow		
			continent	country	region
Annot. performance	1009	6407	48.9	12.2	3.0
Annot. ensemble oracle	3919	443	98.0	70.0	28.0
Random location	120	10273	16.0	0.0	0.0
Random image	328	8724	20.0	2.0	0.0
Base model	2235	3247	74.0	36.0	8.0
Combined model	3333	1948	86.0	70.0	34.0

racy of individual annotators, but is still outmatched by the Annotator Ensemble Oracle.

Attention Maps. We represent in Figure H the self-attention maps of the [CLS] token of the last layer of the combined model of images from the teaser. We observe that the network focuses on regions of interest containing useful geographical cues, such as the double yellow road line—a specific trait to certain countries—or vegetation and buildings.

C. Implementation Details

In this section, we detail our architecture, loss, metrics, and the retrieval algorithm.

Architecture. All considered networks have a base image encoder $\mathcal{I} \mapsto \mathbb{R}^d$, with a d which depends on each architecture ($d = 768$ for the model ViT-B-32, and $d = 1024$ for all the other encoders). We then add one or several heads to map the image representation to geographical information:

- *Regression f^{loc} .* This network directly predicts the longitude and latitude of an image with a MLP of size $d \mapsto d \mapsto 64 \mapsto 2$ with group norms of 4 groups [75] and without normalizing the last layer.

- *Regression $f^{loc} \sin/\cos$.* For this variation, we predict the cosine and sine of both coordinates with an MLP: $d \mapsto d \mapsto 64 \mapsto 4$ with a normalization that ensures that the squared sum between coordinate 0, 1 and 2, 3 is 1. We then use the **atan2** function to recover the corresponding coordinates.

- *Classification $f^{classif}$.* To predict in which of the K geographic divisions an image was taken, we use an MLP: $d \mapsto d \mapsto 512 \mapsto K$.

- *Hybrid $f^{relative}$.* In the hybrid model, we predict both the division and the position of the image within this cell. The relative position is predicted for all cells with an MLP $\phi^{relative} : d \mapsto d \mapsto 512 \mapsto \mathbb{R}^2 K$ with a specialized normalization for the last layer, explained below. During inference, we select the relative prediction of the cell with the highest prediction score for $f^{classif}$. During training, we only supervise the relative prediction that corresponds to the true cell.

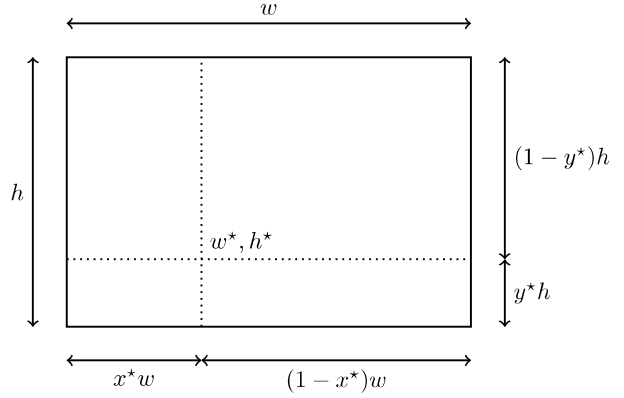


Figure I. **Hybrid Model.** The normalization of the hybrid model requires special considerations to ensure that the output (x, y) of ϕ^{aux} is such $(0, 0)$ maps to the cell’s centroid w^*, h^* , and that $[-1, 1]^2$ maps the entire cell.

For this network, a key implementation detail is the normalization of the last layer of $\phi^{relative}$. We require that for each cell a prediction of $(0, 0)$ should correspond to the centroid $h^*, w^* \in \mathcal{C}^2$ of the training set images in the cell, and that a range of prediction of $[-1, 1]^2$ covers the entire bounding box of size h, w . As illustrated in Figure I, we denote by $x^*, y^* \in [0, 1]^2$ the relative position of the centroid in the cell and by x, y the prediction of the MLP ϕ^{aux} . The output of $f^{relative}$ is defined as follows:

$$w^* + w \begin{cases} -xx^* & \text{if } x \leq 0 \\ x(1-x^*) & \text{else} \end{cases}, \quad (1)$$

$$h^* + h \begin{cases} -yy^* & \text{if } y \leq 0 \\ y(1-y^*) & \text{else} \end{cases}. \quad (2)$$

This normalization allows the network $\phi^{relative}$ to easily predict the centroid of the cell, which facilitates learning the distribution of images of that cell. This is particularly crucial for cells with an off-centered centroid, as it provides increased precision in high density areas. In practice, removing this normalization decreases the performance of the hybrid model by 59 points of geoscore, or 22% from the benefit brought by using a hybrid model over pure classification.

- *Auxiliary f^{aux} .* Finally, the auxiliary network is an MLP $d \mapsto d \mapsto 64 \mapsto A$, where A corresponds to the number of auxiliary task predictions: 11 for land cover, 31 for climate, 15 for soil type, 1 for the driving side, and 1 for the distance to the nearest sea. For all classification tasks (*i.e.* everything except the distance to the sea), we softmax the output logits.

Contrastive Learning. We use the MIL-NCE loss [45] as our contrastive objective, which extends the InfoNCE loss [51] to cases where each sample can have multiple positive

matches.

$$\sum_{i \in \mathcal{B}} \log \left(\frac{\sum_{p \in \mathcal{P}_i} e^{f^{\text{img}}(i) \top f^{\text{img}}(p)/T}}{\sum_{p \in \mathcal{P}_i} e^{f^{\text{img}}(i) \top f^{\text{img}}(p)/T} + \sum_{n \in \mathcal{B} \setminus \mathcal{P}_i} e^{f^{\text{img}}(i) \top f^{\text{img}}(n)/T}} \right), \quad (3)$$

with $\mathcal{P}_i \subset \mathcal{B}$ the set of image positively paired with i and T a temperature parameter set as 0.1. If an image has only one positive match, this equation becomes the InfoNCE loss [51].

Nearest Neighbors Retrieval To perform nearest neighbor retrieval, we create a HNSW32 index using the FAISS library [35] through the autofaiss package (<https://github.com/criteo/autofaiss>). This approach achieves fewer than 200 self-consistency errors per million with over 90% compression rate.

During retrieval, our training set is divided into five parts, each requiring 15 minutes for index computation and collectively consuming 15.6GB of storage for StreetCLIP embeddings, our most resource-intensive model. This setup enables us to predict locations for 12,000 to 32,000 test images per second, depending on the model size.

Although retrieval methods demonstrate high performance and have been made efficient with approximate methods, it is important to note that they are not a learning technique, as they rely on already geographically relevant representations that are already learned.

D. Datasheet for Dataset

D.1. Motivation

Q1 For what purpose was the dataset created? Was there a specific task in mind? Was there a particular gap that needed to be filled? Please provide a description.

- OpenStreetView-5M (OSV-5M) is the first global scale, open-access, large dataset of street view images. Its goal is to enable the training and evaluation of modern computer vision approaches for global visual geolocation, which would depend until now on proprietary or expensive APIs such as Google Street View. More broadly, OSV-5M can be used to evaluate and improve representation learning.

Q2 Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

- The dataset was created as part of the “IMAGINE Summer Hackathon”, an internal event of the LIGM/ENPC/UGE laboratory. All images of OSV-5M come from the Mapillary website, which is a platform where users upload georeferenced images.

Q3 Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.

- This work was partially supported by the ANR project READY3D ANR-19-CE23-0007 and used the HPC resources of IDRIS under the allocation 2022-AD011012096R1 made by GENCI.

Q4 Any other comments?

- All the images of OSV-5M are already openly accessible through Mapillary’s heavily moderated database. We only selected a small fraction distributed across the globe, and added metadata from public sources.

D.2. Composition

Q5 What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

- OSV-5M is composed of street view images depicting various street scenes, captured by dash-cams of different vehicles from across the world.

Q6 How many instances are there in total (of each type, if appropriate)?

- The training set contains 4,894,685 images, and the test set 210,122.

Q7 Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?

- OSV-5M is a small subset of 5.1M images from the 1.8 billion images hosted on the Mapillary website.

Q8 What data does each instance consist of?

- Each instance consists of a georeferenced street view image with a height of 512 pixels.

Q9 Is there a label or target associated with each instance?

- **Yes.** Each image is associated with the following targets: longitude and latitude, administrative division (country, region, sub-region, closest city), and labels corresponding to the local land cover, soil, and climate type at a resolution of 30 arc seconds (1km). We also add the distance to the nearest sea and the driving side of the country.

Q10 Is any information missing from individual instances?

- **Yes.** Sub-regions are not defined for all countries, about 30% of the instances do not have a value for this field.

Q11 Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network links)?

- **No.** The data is organized as a collection of images with no particular order or relations. However,

the metadata allows a user to organize them based on different geographical criteria.

Q12 Are there recommended data splits (e.g., training, development/validation, testing)?

- **Yes.** We provide an official training and test set. Our implementation also proposes a validation split.

Q13 Are there any errors, sources of noise, or redundancies in the dataset?

- **Yes.** We have heavily filtered the dataset using semi-automatic methods to discard low-quality images and wrong localization, as presented in Section A. We have estimated through the manual inspection of 4500 images that 96.1% ($\pm 0.57\%$ with a 95% confidence level) of the images in OpenStreetView-5M are perceptually localizable, *i.e.* provide a clear enough overview of their surroundings.

Q14 Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?

- **No.** OSV-5M is self-contained and will be stored and distributed on huggingface.co.

Q15 Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor–patient confidentiality, data that includes the content of individuals’ non-public communications)?

- **No.** OSV-5M relies on crowdsourced data, whose license is respected by providing usernames for each image, which is include in our metadata.

Q16 Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.

- **Highly unlikely:** OSV-5M contains 5 million images of streets that come from Mapillary, which imposes a strong crowd-sourced moderation policy.

Q17 Does the dataset relate to people?

- **Yes.** Many of the images of OSV-5M contain vehicles and some contain pedestrians, yet Mappillary performs highly accurate privacy blurring.³

Q18 Does the dataset identify any subpopulations (e.g., by age, gender)?

- **No.** The metadata contains no information about the people present in the photography beyond, who are also privacy blurred.¹

Q19 Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?

- **No.** The license plates and faces of pedestrians have been privacy blurred by Mapillary using an automatic algorithm with over 99% recall for faces and 99.9% recall for license plates.¹ Furthermore, users can signal images that violate privacy.

We also manually inspected 4500 images and observed no confidentiality leak. With a confidence of 95% we can assume that fewer than 0.067% of the dataset contains leaks.

Q20 Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)?

- **No.**

Q21 Any other comments?

- **No.**

D.3. Collection Process

Q22 How was the data associated with each instance acquired?

The images of Mapillary are taken and uploaded by users of the Mapillary platform. We downloaded the images directly from Mapillary’s API. Additional metadata was collected from the following open-access sources: (i) land cover: Global Land Cover Share Database [40] (ii) climate: Köppen-Geiger climate classification maps [10], (iii) soil type: Digital World Soil Map [59] (iv) administrative division: reverse geocoder [7].

Q23 What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)?

- We used Mapillary’s web API and a Python script running on a standard workstation.

Q24 If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

- We first defined a $100 \times 100\text{m}$ grid across the entire world and sampled one image per cell among the 1.8B images of Mapillary. We then sample the train and test sets with a weight proportional to the local image density raised to the power of -0.75 . We then filter the images based on both learned and handcrafted filters, as described in Section A.

Q25 Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

³See <https://blog.mapillary.com/update/2018/04/19/accurate-privacy-blurring-at-scale.html>

- The images are crowdsourced by Mapillary users who agree on Mapillary’s **terms of use**. To the best of our knowledge, users are not compensated.

Q26 Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?

- The images used in OSV-5M were uploaded between January 2011 and August 2023.

Q27 Were any ethical review processes conducted (e.g., by an institutional review board)?

- No.

Q28 Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?

- N.A. The images were downloaded through Mapillary’s API.

Q29 Were the individuals in question notified about the data collection?

- No. We followed the terms of use of Mapillary.

Q30 Did the individuals in question consent to the collection and use of their data?

- Yes. Following the Mapillary terms of use, a user agrees for their data to be used respecting the CC BY-SA 2.0 DEED license.

Q31 If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?

- N.A.

Q32 Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted?

- No. However, users of OSV-5M can signal potential issues with the images to the corresponding authors. Flagged images will be removed and Mapillary will be further contacted.

Q33 Any other comments?

- All the images of OSV-5M are already openly accessible through Mapillary’s heavily moderated database. We only added additional metadata from public sources.

D.4. Preprocessing, Cleaning, and/or Labeling

Q34 Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?

- Yes. We removed the images based on learned and handcrafted filters, as described in Section A. In particular, we removed images that were classified as blurry, too dark or purple, or badly exposed. We

also used a pretrained model [22] to detect and remove images with potential spurious orientation.

Q35 Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.

- Yes. The removed images are saved on a local server but are not public. Note that all these images, including the filtered ones, are still available on Mapillary’s website.

Q36 Is the software used to preprocess/clean/label the instances available?

- Yes. The script used for cleaning the dataset will be released alongside the dataset.

Q37 Any other comments?

- No.

D.5. Uses

Q38 Has the dataset been used for any tasks already?

- Yes. To train and evaluate geolocation models, the subject of the paper.

Q39 Is there a repository that links to any or all papers or systems that use the dataset?

- No. But once we release the dataset we will maintain an updated list on the project page.

Q40 What (other) tasks could the dataset be used for?

- The images of OSV-5M can be used for both self-supervised learning and generative modeling, both as a pretraining or fine-tuning dataset. The meta-data beyond geolocation can be used as targets for separate tasks.

Q41 Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?

- The density-based sampling leads to a spatial distribution that may not fit other datasets and tasks.

Q42 Are there tasks for which the dataset should not be used?

- Yes. The same limitations that apply for Mapillary data (CC BY-SA 2.0 DEED), also apply to our dataset.
- **Privacy Concerns.** Despite being heavily moderated, the dataset may contain images of individuals or private residences. Usage must avoid applications that can infringe on personal privacy or exercise surveillance and open-source intelligence (OSINT).
- **Cultural and Ethical Sensitivity.** The dataset spans a wide range of cultures and countries, each

with its own set of ethical norms and cultural sensitivities. We strongly advise against using OSV-5M in a way that might propagate stereotypes, misrepresent cultures, or otherwise harm the dignity and representation of the featured communities.

- **Manipulation and Misrepresentation.** The dataset should not be used to create misleading representations of locations or to manipulate images in a way that distorts or misrepresents the reality of the places and the depicted people.

Q43 **Any other comments?**

- No.

D.6. Distribution

Q44 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?**

- **Yes.** The dataset will be open-access and accessible to the research community.

Q45 **How will the dataset be distributed (e.g., tarball on website, API, GitHub)?**

- The data will be hosted on huggingface.co.

Q46 **When will the dataset be distributed?**

- The dataset will be distributed upon the publication of the preprint on arXiv, which should be in Q2 of 2024.

Q47 **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.**

- **Yes.** The dataset inherits from Mappillary CC-BY-SA license: free of use with attribution to the authors of the images [2].

Q48 **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**

- No.

Q49 **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**

- No.

Q50 **Any other comments?**

- No.

D.7. Maintenance

Q51 **Who will be supporting/hosting/maintaining the dataset?**

- The authors will maintain the dataset. The dataset will be hosted on huggingface.co.

Q52 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

- A dedicated email will be created.

Q53 **Is there an erratum?**

- **No.** There is no erratum for our initial release. Errata will be documented as future releases on the dataset website.

Q54 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**

- **Yes.**

Q55 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)?**

- N.A.

Q56 **Will older versions of the dataset continue to be supported/hosted/maintained?**

- **Yes.** We are dedicated to providing ongoing support for the OSV-5M dataset.

Q57 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**

- **Yes.** The data is free of use under Mappillary CC-BY-SA license. User making explicit use of our proposed split should cite our paper.

Q58 **Any other comments?**

- No.