






# SatCLIP: Global, General-Purpose Location Embeddings with Satellite Imagery

Konstantin Klemmer<sup>1</sup>, Esther Rolf<sup>2,3</sup>, Caleb Robinson<sup>4</sup>, Lester Mackey<sup>1</sup>,  
and Marc Rußwurm<sup>5</sup>

<sup>1</sup> Microsoft Research New England, Cambridge, MA, USA

<sup>2</sup> University of Colorado Boulder, CO, USA

<sup>3</sup> Harvard University, Cambridge, MA, USA

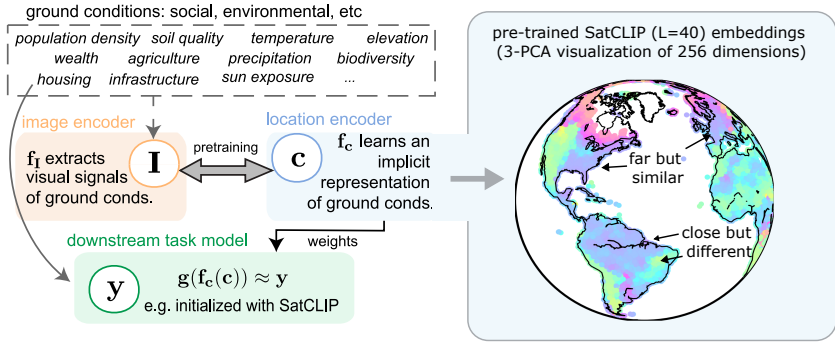
<sup>4</sup> Microsoft AI for Good Research Lab, Redmond, WA, USA

<sup>5</sup> Wageningen University, Wageningen, Netherlands

**Abstract.** Geographic information is essential for modeling tasks in fields ranging from ecology to epidemiology. However, extracting relevant location characteristics for a given task can be challenging, often requiring expensive data fusion or distillation from massive global imagery datasets. To address this challenge, we introduce Satellite Contrastive Location-Image Pretraining (SatCLIP). This *global, general-purpose geographic location encoder* learns an implicit representation of locations by matching CNN and ViT inferred visual patterns of openly available satellite imagery with their geographic coordinates. The resulting SatCLIP location encoder efficiently summarizes the characteristics of any given location for convenient use in downstream tasks. In our experiments, we use SatCLIP embeddings to improve prediction performance on nine diverse location-dependent tasks including temperature prediction, animal recognition, and population density estimation. Across tasks, SatCLIP consistently outperforms alternative location encoders and improves geographic generalization by encoding visual similarities of spatially distant environments. These results demonstrate the potential of vision-location models to learn meaningful representations of our planet from the vast, varied, and largely untapped modalities of geospatial data.

## 1 Introduction

Satellite imagery has proven to be a valuable source of input data for predictive models across a wide range of real-world applications (Rolf et al., 2021), for example, interpolating missing air pollution data (Chen et al., 2019), crop yield forecasting (Lobell et al., 2015; Tseng et al., 2022), and agro-forestry carbon stock prediction (Reiersen et al., 2022). Many geospatial modeling tasks, in fields ranging from epidemiology, the Earth system sciences, to ecology also directly leverage geographic location for improving predictions (Cole et al., 2023; Mac Aodha et al., 2019). Patterns extracted from satellite images can describe the unique characteristics of locations, by capturing their natural and built environment. These characteristics are often correlated in space: While two nearby



**Fig. 1: Motivation for SatCLIP:** Capturing ground conditions from satellite images and transferring them into a location encoder via contrastive image-location pretraining. The right globe shows a PCA representation of the pre-trained location encoder.

locations are more likely to have similar features (e.g. the same land cover), two distant locations can also share location characteristics when they share similar environmental ground conditions like climate zones (Fig. 1, right). Since the spatial patterns governing different geographic data modalities are often complex and non-linear, predictive models working with geo-data benefit from explicitly integrating intuitions for spatial and spatio-temporal dependencies (Alam et al., 2022; Cole et al., 2023; Fotheringham et al., 2009; Klemmer and Neill, 2021; Klemmer et al., 2022).

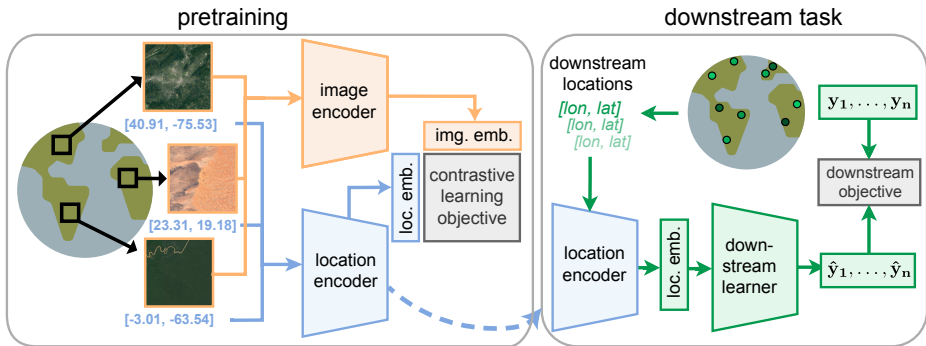
But integrating geographic information into a deep learning model is not straightforward. Even though spatial coordinates are often informative, introducing them as features can amplify geographic distribution shift problems and lead to poor evaluation accuracy. This is especially of concern for cross-regional generalization when data from evaluation areas (and their coordinates) are absent in the training data. Consequently, many location-informed models are only applicable for interpolation problems where the evaluation area overlaps with the training area. And while some settings warrant interpolation methods—e.g., species distribution modeling (Berg et al., 2014; Cole et al., 2023; Mac Aodha et al., 2019), where data is global and somewhat representative spatially—for many applications, available labeled data is patchy and sparse, and predictive models must generalize to unseen geographic areas potentially far from the training data. In these situations, utilizing location embeddings that capture ground conditions might be preferable to using raw spatial coordinates. We tackle this research problem by pretraining location encoders with globally and uniformly sampled satellite imagery and with general-purpose use in mind. Specifically, we leverage the location information indexing satellite images as an input to a contrastive pretraining objective that aims to match location-image pairs. This is analogous to the vision-language pretraining deployed in Contrastive Location Image Pretraining (CLIP) (Radford et al., 2021).



*Previous work on pretrained location encoders.* Three studies have pretrained geographic location encoders that input spatial coordinates, and return learned contextual representations. Yin et al. (2019) propose GPS2Vec, a set of UTM-zone specific location encoders using geotagged Flickr images (YFCC100M) (Thomee et al., 2016) and their corresponding semantic tags for training. Geographic generalization was out of scope for this work as embeddings are only available for UTM-zones in which training data can be found. Mai et al. (2023) introduce Contrastive Spatial Pre-Training (CSP) on the iNaturalist 2018 (iNat) (Horn et al., 2018) species imagery and the Functional Map of the World (FMoW) (Christie et al., 2018) satellite image datasets. CSP is used for unsupervised pre-training and downstream prediction on the same datasets and was not conceptualized for use on other tasks. The CSP pretraining datasets, iNaturalist and FMoW, are also unevenly distributed over space, with high image densities in North America and Europe and few images outside of Western regions. Cepeda et al. (2023) propose GeoCLIP, in which the authors pretrain image and location encoders using Flickr images from the MediaEval Placing Tasks 2016 (MP-16) dataset (Larson et al., 2017)—another dataset with strong overrepresentation of Western countries. GeoCLIP is developed for the task of geo-locating (natural) images and is not optimized or tested for general-purpose use.

*Aim and contributions.* Existing work leaves two important gaps: understanding how location encoders generalize across various downstream tasks and ensuring global coverage and approximately equal performance of location embeddings. In this work, we address both of these challenges by introducing the first *global-coverage, general-purpose pretrained geographic location encoder* based on Satellite Contrastive Location Image Pretraining – **SatCLIP**. SatCLIP distills spatially varying visual patterns from globally-distributed satellite data into an implicit neural representation in a comparatively small and efficient neural network. This location encoder projects a latitude and longitude coordinate into a higher-dimensional vector representation that is matched with a visual vector representation from a computer vision encoder (CNN or ViT), as detailed later in Sec. 2. More generally, the proposed framework represents a step towards geographically-informed “foundation models” trained with large, unlabeled datasets, that are usable for a wide range of tasks, and extrapolate to unseen geographic areas. Our contributions can be summarized as follows:

- We develop the first task-generalizable, global-coverage location encoder—**SatCLIP**—trained on Sentinel-2 multispectral satellite imagery. We release the pretrained encoder as a `PyTorch` model. We also release our new pre-training dataset, **S2-100K**.
- We compare **SatCLIP** with existing pretrained location encoders and other geographic feature generation approaches on nine diverse downstream tasks, ranging from temperature prediction to population density estimation, highlighting superior performance in prediction and geographic generalization.



**Fig. 2: The SatCLIP pretraining and deployment pipeline.** SatCLIP pretraining through image-location matching is outlined on the left. The pretrained location encoder can then be used in downstream tasks, highlighted on the right.

## 2 Satellite Contrastive Location-Image Pretraining

With **SatCLIP**, we aim to train models that (1) provide *general purpose embeddings* and (2) are *globally representative*. In this section, we first motivate SatCLIP and then outline its components and training paradigm.

Various factors influence the appearance of Earth’s surface, as captured in satellite images. In Fig. 1, we highlight that diverse environmental and socioeconomic factors are reflected in visual markers, such as the appearance of mountain ranges with their specific vegetation, the geometry and structure of agricultural fields, and the design of buildings. By training on the SatCLIP matching objective, the image encoder  $f_I$  is trained to associate an image  $\mathbf{I}$  with a location  $\mathbf{c}$  based on the various ground factors detectable in the images. At the same time, the location encoder  $f_c$  is trained to associate a given location with location-specific image characteristics. Effectively, both models learn to align the embeddings of a location and its corresponding image to maximize their similarity, as shown in the left panel of Fig. 2. Inductive biases in the location encoder model control spatial smoothness (i.e. the  $L$  hyperparameter in SatCLIP), allowing the model to interpolate to areas where no image-location pairs are present in the training data. After successful training, location features can be extracted at arbitrary locations. These features can be used to train any downstream learner  $g$  that takes locations as input (right panel of Fig. 2).

There are several advantages to training downstream models with location embeddings, as opposed to raw coordinates or images extracted at downstream locations. Models trained on raw coordinates  $\mathbf{c}$  solely rely upon spatial dependencies without taking any ground conditions into account, such as local elevation patterns or local climate zones. Models trained on full images  $\mathbf{I}$ , while able to capture ground conditions, require expensive data preprocessing (downloading images for every downstream location) and training of large vision models. This can be infeasible in many, especially low-resource settings. SatCLIP has the potential to provide the best of both worlds: location embeddings capture both

spatial effects and ground conditions while also being relatively low-dimensional (our SatCLIP location embeddings are 256-dimensional vectors) and runtime efficient due to the small size of the location encoder model. This is particularly helpful in resource-constrained settings and allows fast encoding of coordinates to embedding space without GPUs.

## 2.1 Pretraining with the SatCLIP Objective

The inputs to a geographic location encoder are latitude/longitude coordinate pairs  $\mathbf{c}_i = [\lambda_i, \phi_i]$ , where  $\lambda_i$  is the longitude,  $\phi_i$  is the latitude, and  $i$  indexes locations on the spherical surface  $\mathbb{S}^2$ . For each location  $i$ , we have a corresponding multi-spectral image  $\mathbf{I}_i \in \mathbb{R}^{m \times n \times k}$  with  $k$  channels. We now define two encoders, a *location encoder*  $f_c : \mathbb{S}^2 \rightarrow \mathbb{R}^d$  that takes in 2-dimensional coordinates  $\mathbf{c}_i$  and returns a  $d$ -dimensional latent embedding and a *image encoder*  $f_I : \mathbb{R}^{m \times n \times k} \rightarrow \mathbb{R}^d$  that takes in an image  $\mathbf{I}_i$  and also returns a  $d$ -dimensional latent embedding.

We train both encoders with the simple but highly effective CLIP (Radford et al., 2021) objective

$$\mathcal{L} = \frac{1}{2N} \left[ \sum_{i=1}^N \mathcal{L}_{\text{loc}}(\mathbf{c}_i, \mathbf{I}_1, \dots, \mathbf{I}_N) + \sum_{i=1}^N \mathcal{L}_{\text{img}}(\mathbf{I}_i, \mathbf{c}_1, \dots, \mathbf{c}_N) \right] \quad (1)$$

that matches each coordinate  $\mathbf{c}_i$  with the corresponding image  $\mathbf{I}_i$  and against all images  $\mathbf{I}_1, \dots, \mathbf{I}_N$  using

$$\mathcal{L}_{\text{loc}}(\mathbf{c}_i, \mathbf{I}_1, \dots, \mathbf{I}_N) = -\log \frac{\exp(\langle f_c(\mathbf{c}_i), f_I(\mathbf{I}_i) \rangle / \tau)}{\sum_{j=1}^N \exp(\langle f_c(\mathbf{c}_i), f_I(\mathbf{I}_j) \rangle / \tau)} \quad (2)$$

and each image with the corresponding coordinate using a complementary loss  $\mathcal{L}_{\text{img}}(\mathbf{I}_i, \mathbf{c}_1, \dots, \mathbf{c}_N)$  over a batch  $(\mathbf{c}_i, \mathbf{I}_i)_{i=1}^N$  of  $N$  coordinate-image tuples. The normalized dot-product is denoted by  $\langle \cdot, \cdot \rangle$ , and  $\tau$  is a temperature hyperparameter. This objective optimizes the weights of the location encoder  $f_c$  and image encoder  $f_I$  simultaneously to embed the feature vectors of the corresponding location  $f_c(\mathbf{c}_i) \in \mathbb{R}^d$  and image  $f_I(\mathbf{I}_i) \in \mathbb{R}^d$  nearby in a common  $d$ -dimensional feature space.

## 2.2 Encoder Architectures

Geographic **location encoders**  $f_c$  typically take the form  $f_c = \text{NN}(\text{PE}(\mathbf{c}_i))$ , (Mai et al., 2023) where  $\text{PE}(\mathbf{c}_i)$  is a nonparametric functional positional encoding and  $\text{NN}(\cdot)$  is a small neural network. The positional encodings usually have a scale hyperparameter that controls spatial smoothness of the encoding. The neural network weights encode an implicit neural representation of a signal at a specific coordinate (Cole et al., 2023). In this work, we train Siren( $\text{SH}(\mathbf{c}_i)$ ) location encoders proposed by Rußwurm et al. (2024), which use spherical harmonics basis (SH) functions as positional encoders and are particularly well-suited for coordinates on spherical surfaces. They are combined with sinusoidal representation networks (Siren) (Sitzmann et al., 2020) that are broadly used for implicit

neural representations. The spatial smoothness of the representation is controlled by the number of Legendre polynomials  $L$ . This effectively defines the resolution of the location encoding and its capacity to learn small and large-scale geospatial patterns, with larger  $L$  corresponding to finer spatial resolution.

As **image encoder**, we need a vision model that is expressive enough to learn visual patterns from satellite images. In this work, we use ResNet18, ResNet50, (He et al., 2016) and ViT16 (Dosovitskiy et al., 2020) vision encoders pretrained with momentum-contrast (MoCo) (He et al., 2020) on Sentinel-2 satellite imagery by Wang et al. (2022a). To account for the size discrepancy between the large image models and relatively smaller location encoders (for example, the image encoder of a SatCLIP ViT16 has  $\sim 22$  million parameters, while the location encoder has  $\sim 1$  million parameters) during training, we freeze the vision encoder except for the last linear projection layer.

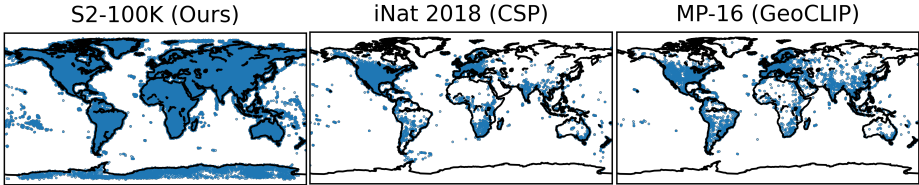
### 2.3 SatCLIP Implementation Details

We pretrain **SatCLIP** using the S2-100K dataset, which we assemble for this purpose (described in the following Sec. 3.1). We use 90% of the data points, selected uniformly at random, for pretraining and reserve the remaining 10% as a validation set to monitor overfitting. During pretraining, we found that batch sizes of  $8k$  help the model to learn more fine-grained representations, while too large batch sizes can prevent learning, as was also recently observed on CLIP models by Zhai et al. (2023). We train models for 500 epochs on an A100 GPU. More pretraining details can be found in Appendix C.

## 3 Experimental Setup

In our experiments, we focus on three research questions. From a performance perspective, we ask *how generalizable are SatCLIP embeddings* from Sentinel-2 data, both *across a diverse range of geospatial modeling tasks* (RQ 1) and *across unseen geographic areas* (RQ 2), compared to existing pretrained location encoders and location-only prediction? We design experiments to test the performance of SatCLIP embeddings for downstream tasks, both for spatial interpolation, and for geographic domain generalization, in which the training and test sets are separated geographically. Geographic generalization is an important aspect of performance, as distributional changes across geographic areas are a common challenge in environmental problems like species distribution modeling (Beery et al., 2022), land cover classification (Rußwurm et al., 2020), crop type mapping (Kondmann et al., 2021), and yield prediction (Tseng et al., 2022). While location-only prediction is generally unsuitable for geographic generalization tasks, the implicit neural representation of environmental factors within the SatCLIP location encoder may provide generalizable information supporting prediction in areas with no or few labeled data points available.

After benchmarking performance, we turn to analyses that help us to further understand why SatCLIP works and what limitations it might exhibit. Ablation studies systematically vary the image encoder architecture and the scale



**Fig. 3: Spatial distribution of the S2-100K dataset** used for training SatCLIP compared with iNaturalist 2018 (Horn et al., 2018) and MP-16 (Larson et al., 2017), which are used to pretrain CSP and GeoCLIP models. iNaturalist and MP-16 heavily overrepresent North America and Europe.

parameters of the location encoder to understand the joint effect of each design choice. Following the intuition of Fig. 1 and complementing these quantitative experiments, we ask: *Do SatCLIP embeddings capture ground conditions and incorporate similarities over space?* (RQ 3). Qualitative experiments shed light on what spatial relationships the pretrained SatCLIP embeddings harvest from multi-spectral satellite imagery. Before presenting results in Sec. 4, we detail our experimental design, including the pretraining dataset (Sec. 3.1), downstream tasks (Sec. 3.2), and comparison models (Sec. 3.3).

### 3.1 Pretraining Dataset: 100k uniformly sampled Sentinel-2 images

To construct our pretraining dataset, **S2-100K**, we sample 100,000 tiles of  $256 \times 256$  pixel, multi-spectral (12-channel) Sentinel-2 satellite imagery and their associated centroid locations. We design the S2-100K dataset with the goals of *multi-task applicability* and *geographic generalization performance* in mind. Our dataset (1) represents general location features by using multi-spectral satellite imagery (as illustrated in Fig. 1) and (2) is nearly uniformly distributed across global land mass (Fig. 3, left). More details on the S2-100K dataset and sampling procedure are provided in Appendix A. In contrast, the pretraining datasets used in comparison methods (Sec. 3.3 and Appendix B) often significantly underrepresent certain – especially non-Western – geographic areas, as they were not specifically designed to provide general-purpose embeddings. Fig. 3 illustrates the spatial coverage of S2-100K compared to the highly clustered distributions of iNaturalist, which is used as a pretraining dataset for CSP (Mai et al., 2023) and MediaEval 2016, which is used as a pretraining dataset for GeoCLIP (Cepeda et al., 2023). Similar biases are exhibited by the Yahoo-Flickr Creative Commons 100 Million (YFCC100M) dataset (Thomee et al., 2016) of image-tag-location triplets (used in GPS2Vec (Yin et al., 2019)), and the Functional Map of the World (FMoW) (Christie et al., 2018) dataset, which samples satellite imagery mostly near human-built infrastructure.

### 3.2 Downstream Tasks

To test the general applicability of SatCLIP location embeddings, we run experiments on a wide range of geospatial predictive modeling tasks. In all datasets, the inputs are raw latitude/longitude coordinates, which we transform into location embeddings. The nine downstream datasets used in this work span socioeconomic and environmental applications. We predict variables including **Air Temperature** (Hooker et al., 2018) and **Elevation** (Rolf et al., 2021) from coordinates as environmental regression objectives. To capture socioeconomic factors, we regress **Median Income** (Jia and Benson, 2020), **California Housing** prices (Pace and Barry, 2003), and logged **Population Density** (Rolf et al., 2021). We additionally classify **Biomes**, **Ecoregions** (Dinerstein et al., 2017), and compile a new country code classification task **Countries**. Lastly, we classify **iNaturalist** species (Horn et al., 2018). Here, we have additional image features extracted from an InceptionV3 model released by (Mac Aodha et al., 2019), which we concatenate with location embeddings during downstream training. Further details on downstream experiments can be found in Appendix E.

### 3.3 Comparison Methods

We compare trained SatCLIP location embeddings to GPS2Vec (Yin et al., 2019), CSP (Mai et al., 2023) and GeoCLIP (Cepeda et al., 2023) pretrained location embeddings. We refer to each comparison model by first stating the pre-training algorithm and then the pretraining dataset. For instance, CSP-FMoW represents CSP pretraining on FMoW dataset. Unless stated otherwise, we show results from SatCLIP models using a ViT-16 vision encoder. A summary of each comparison method is given in Tab. 1, with details given in Appendix B. Like SatCLIP, GeoCLIP and CSP use the CLIP loss, with CSP adding loss terms for negative location sampling and SimCSE. GPS2Vec uses a KL divergence loss on the context (image and semantic tags) and location data. To compare to an image-only embedding, we use globally precomputed MOSAIKS (Rolf et al., 2021) features, accessed via [siml.berkeley.edu](https://siml.berkeley.edu) (Carleton et al., 2022). To assess the performance improvement from the integration of contextual information over location-only prediction, we also compare to downstream learners trained on raw latitude/longitude coordinates  $g(\mathbf{c})$ . We refer to this approach as “Identity” throughout our experiments.

**Table 1: Key characteristics of SatCLIP vs. the GPS2Vec, CSP and GeoCLIP location encoders and the MOSAIKS feature extractor.**

	Contextual data	Location encoder
<b>SatCLIP</b> (ours)	S2-100K	Spherical harmonics & Siren
<b>GPS2Vec</b> (Yin et al., 2019)	YFCC images and semantic tags	Two-level soft encoding
<b>CSP</b> (Mai et al., 2023)	FMoW / iNaturalist	Sinusoidal transform & FcNet
<b>GeoCLIP</b> (Cepeda et al., 2023)	MediaEval 2016	Random Fourier Features & MLP
<b>MOSAIKS</b> (Rolf et al., 2021)	Planet Basemaps 2019 Q3	N/A (direct feature extractor)

**Table 2: Downstream task performance using SatCLIP (with ResNet50) vs. comparison location embeddings.** We report average test set  $R^2$  and accuracy  $\pm 1$  standard deviation across 10 independently initialized training runs.

Task $\downarrow$ Data $\rightarrow$	SatCLIP <sub>L=10</sub> (S2-100K)	SatCLIP <sub>L=40</sub> (S2-100K)	CSP (iNat)	GPS2Vec (tag)	MOSAICS (Planet)	GeoCLIP (MP-16)	Identity ( $y \sim g(c)$ )
<b>Regression</b>	$R^2 \uparrow$						
Air temperature	<b>0.90 <math>\pm</math> 0.13</b>	<b>0.91 <math>\pm</math> 0.01</b>	$-0.56 \pm 0.59$	$0.22 \pm 0.00$	$-0.52 \pm 2.00$	$-3.11 \pm 5.24$	$0.82 \pm 0.16$
Median income	$0.42 \pm 0.01$	<b>0.47 <math>\pm</math> 0.12</b>	$-0.01 \pm 0.02$	$0.21 \pm 0.00$	$0.02 \pm 0.05$	<b>0.50 <math>\pm</math> 0.01</b>	$-0.84 \pm 0.94$
Cali. housing	$0.35 \pm 0.04$	$0.57 \pm 0.02$	$-0.00 \pm 0.00$	$0.71 \pm 0.03$	$0.24 \pm 0.02$	<b>0.75 <math>\pm</math> 0.01</b>	$0.05 \pm 0.02$
Elevation	$0.83 \pm 0.01$	<b>0.88 <math>\pm</math> 0.00</b>	$0.11 \pm 0.05$	$0.10 \pm 0.00$	$0.21 \pm 0.01$	$0.83 \pm 0.00$	$0.25 \pm 0.08$
Population	$0.79 \pm 0.00$	<b>0.82 <math>\pm</math> 0.00</b>	$0.36 \pm 0.11$	$0.25 \pm 0.00$	$0.46 \pm 0.02$	$0.79 \pm 0.00$	$0.46 \pm 0.03$
<b>Classification</b>	% Accuracy $\uparrow$						
Countries	$94.28 \pm 0.18$	<b>96.00 <math>\pm</math> 0.14</b>	$82.11 \pm 1.72$	$70.35 \pm 0.06$	$76.16 \pm 0.50$	$90.72 \pm 0.44$	$82.94 \pm 2.23$
iNaturalist	<b>65.69 <math>\pm</math> 0.18</b>	<b>66.22 <math>\pm</math> 0.40</b>	$60.47 \pm 0.56$	$58.78 \pm 0.48$	$56.73 \pm 0.8$	$62.01 \pm 0.59$	$60.83 \pm 0.53$
Biome	$92.23 \pm 0.26$	<b>94.41 <math>\pm</math> 0.14</b>	$73.18 \pm 5.58$	$69.69 \pm 0.06$	$79.61 \pm 0.42$	$89.57 \pm 0.45$	$83.55 \pm 2.43$
Ecoregions	$89.32 \pm 0.31$	<b>91.67 <math>\pm</math> 0.15</b>	$78.43 \pm 1.71$	$68.46 \pm 0.06$	$70.48 \pm 0.21$	$84.65 \pm 0.32$	$77.07 \pm 2.54$

### 3.4 Downstream Model Training

For all downstream tasks, we train multi-layer perceptron (MLP) models  $g$  with location embeddings and raw latitude/longitude coordinates as input to predict a (continuous or discrete) outcome variable  $y$ . Regression models use a mean squared error (MSE) loss, and classification models use cross-entropy loss. Hyperparameters like learning rate, number of layers, or hidden dimensions are tuned using a random search on an independent validation set. All results are reported for an unseen test set. More details on the downstream task setups can be found in Appendix E.

## 4 Results

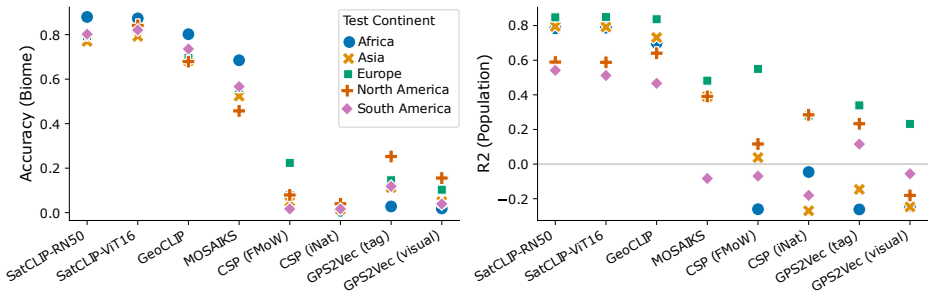
### 4.1 Downstream Task Performance

(RQ 1)

**Quantitative comparison across downstream tasks.** Tab. 2 shows performance across different downstream tasks and methods. SatCLIP embeddings (in either  $L = 10$  or  $L = 40$  configuration) achieve the best prediction scores by a large margin on seven of the nine tasks. The exceptions are the Cali. Housing dataset, which is limited to California, and the Median Income dataset, which only contains data from the continental United States (US). For both tasks, the GeoCLIP model trained on US-centric MP-16 data performs equal to or better than SatCLIP. We also observe that SatCLIP embeddings with higher spatial resolution ( $L = 40$ ) perform better than coarse-grained ( $L = 10$ ) embeddings at these regionally constrained tasks.

**Comparison across continents.** Fig. 4 shows the performance of SatCLIP and comparison methods evaluated separately by continent for the tasks of biome classification and population density estimation. SatCLIP performs well on all continents for both tasks. Prior location encoders (CSP and GPS2Vec) trained on spatially biased training data tend to perform better in Europe and North





**Fig. 4: Performance metrics aggregated by continent** highlight how location embeddings perform in different geographic areas for population density estimation and biome classification for five continents.  $L = 40$  SatCLIP models are shown.

America than in the underrepresented continents of Africa, Asia, and South America. GeoCLIP is closest in performance to SatCLIP for both tasks, performing similarly well to SatCLIP on the population density task, but worse across continents on the Biome classification task.

## 4.2 Zero/Few-Shot Geographic Adaptation

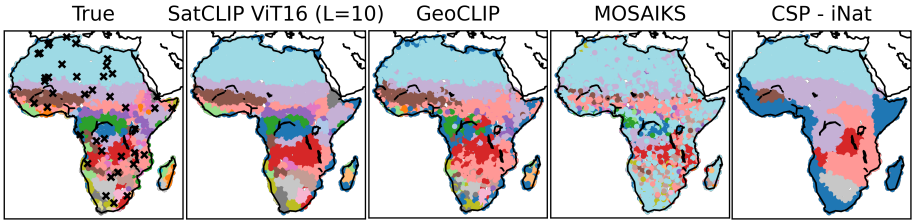
(RQ 2)

Tab. 3 shows how our embeddings perform at geographic generalization. For these experiments, we deploy a spatial train/test split strategy: We hold out entire continents, either Africa or Asia, as test sets and use the remaining data for model training and validation. Since Countries and Ecoregions are often unique to a continent, a zero-shot adaptation approach without any training points on the respective test continents is not possible. Here, we provide a small portion (1%, uniformly sampled) of test continent points in the training set to create a de-facto few-shot geographic adaptation setting. For the remaining tasks, we test zero-shot adaptation by not providing any training points from the held-out test continent. In iNaturalist 2018, for example, this means that the model will not be able to recognize any species that are endemic to the test continent (i.e., that do not live in any other continent) but only those known from training continents.

SatCLIP models are often (but not always) better than the comparison approaches across both held-out continents.  $L = 10$  SatCLIP outperforms the higher-resolution  $L = 40$  SatCLIP model (8 vs. 5 task wins). Here, smoother representations with the lower resolution  $L = 10$  allow information from spatially far points to take effect. GeoCLIP is the closest to SatCLIP, with 4 wins. While 3 wins are shared with SatCLIP, it is uniquely best in population estimation in Africa. Overall, SatCLIP embeddings perform systematically better in the few-shot geographic adaptation setting for Ecoregions and Countries tasks. Several existing location encoders even perform significantly worse than directly encoding latitude and longitude for out of sample prediction (“Identity” in Tab. 3). Fig. 5 visualizes the predictions and failure modes from different methods on Ecoregion prediction in Africa.

**Table 3: Geographic adaptation capabilities of SatCLIP (with ViT16 vision encoder) vs. baseline location embeddings to new geographic areas with no (\*) or very few (†) samples from the held-out test continent.** We report average test set  $R^2$  and accuracy in %  $\pm 1$  standard deviation across 10 independently initialized fine-tuning runs.

Test Continent	SatCLIP <sub>L=10</sub> (S2-100K)	SatCLIP <sub>L=40</sub> (S2-100K)	CSP (iNat)	GPS2Vec (tag)	MOSAIKS (Planet)	GeoCLIP (MP-16)	Identity ( $y \sim g(c)$ )
<b>Asia</b>							
Air Temp.* $R^2 \uparrow$	<b>0.75 <math>\pm</math> 0.05</b>	0.63 $\pm$ 0.04	-0.50 $\pm$ 1.32	-3.95 $\pm$ 4.89	-2.13 $\pm$ 3.50	<b>0.77 <math>\pm</math> 0.28</b>	0.20 $\pm$ 1.64
Elevation*	<b>0.46 <math>\pm</math> 0.08</b>	<b>0.48 <math>\pm</math> 0.07</b>	-0.26 $\pm$ 0.03	-0.29 $\pm$ 0.01	-0.07 $\pm$ 0.06	<b>0.50 <math>\pm</math> 0.03</b>	-0.16 $\pm$ 0.06
Pop. Density*	<b>0.42 <math>\pm</math> 0.08</b>	<b>0.45 <math>\pm</math> 0.04</b>	-1.02 $\pm$ 0.32	-0.37 $\pm$ 0.04	0.05 $\pm$ 0.12	<b>0.38 <math>\pm</math> 0.04</b>	0.03 $\pm$ 0.07
Countries† % Acc. $\uparrow$	<b>36.90 <math>\pm</math> 4.32</b>	19.17 $\pm$ 2.82	1.28 $\pm$ 0.01	1.12 $\pm$ 0.00	1.56 $\pm$ 0.47	23.12 $\pm$ 2.50	1.24 $\pm$ 0.12
iNaturalist*	19.60 $\pm$ 0.78	<b>20.91 <math>\pm</math> 0.77</b>	<b>21.49 <math>\pm</math> 0.85</b>	17.52 $\pm$ 0.38	16.14 $\pm$ 0.42	<b>20.94 <math>\pm</math> 0.38</b>	<b>21.08 <math>\pm</math> 0.69</b>
Biome*	25.89 $\pm$ 2.79	16.44 $\pm$ 1.21	3.00 $\pm$ 2.60	1.76 $\pm$ 0.04	<b>37.81 <math>\pm</math> 4.47</b>	<b>31.67 <math>\pm</math> 1.91</b>	6.24 $\pm$ 2.71
Ecoregions†	<b>21.02 <math>\pm</math> 1.09</b>	10.86 $\pm$ 1.19	1.41 $\pm$ 0.14	1.49 $\pm$ 0.03	1.36 $\pm$ 0.10	6.65 $\pm$ 1.03	1.52 $\pm$ 0.47
<b>Africa</b>							
Air Temp.* $R^2 \uparrow$	-4.71 $\pm$ 2.29	<b>-1.48 <math>\pm</math> 0.70</b>	-2.67 $\pm$ 5.80	-7.91 $\pm$ 0.04	-17.43 $\pm$ 18.37	-9.91 $\pm$ 28.82	-27.36 $\pm$ 39.46
Elevation*	-1.80 $\pm$ 1.74	<b>-0.21 <math>\pm</math> 0.09</b>	-1.20 $\pm$ 0.55	<b>-0.13 <math>\pm</math> 0.06</b>	-0.79 $\pm$ 0.43	-0.34 $\pm$ 0.10	-2.43 $\pm$ 2.67
Pop. Density*	0.17 $\pm$ 0.12	0.18 $\pm$ 0.09	-0.31 $\pm$ 0.16	-0.34 $\pm$ 0.02	0.15 $\pm$ 0.05	<b>0.32 <math>\pm</math> 0.03</b>	-0.50 $\pm$ 0.34
Countries† % Acc. $\uparrow$	<b>30.65 <math>\pm</math> 4.23</b>	10.22 $\pm$ 1.62	0.45 $\pm$ 0.04	0.47 $\pm$ 0.01	0.48 $\pm$ 0.00	10.32 $\pm$ 2.75	2.74 $\pm$ 2.52
iNaturalist*	<b>9.53 <math>\pm</math> 0.57</b>	6.23 $\pm$ 0.47	8.65 $\pm$ 0.52	7.47 $\pm$ 0.53	5.18 $\pm$ 0.38	7.69 $\pm$ 0.30	<b>9.96 <math>\pm</math> 0.33</b>
Biome*	35.72 $\pm$ 5.48	12.34 $\pm$ 1.75	1.09 $\pm$ 0.48	1.29 $\pm$ 0.04	<b>49.86 <math>\pm</math> 1.57</b>	28.28 $\pm$ 3.06	1.46 $\pm$ 0.67
Ecoregions†	<b>32.03 <math>\pm</math> 1.19</b>	12.91 $\pm$ 1.63	0.94 $\pm$ 0.04	0.88 $\pm$ 0.01	0.92 $\pm$ 0.12	12.41 $\pm$ 2.20	7.72 $\pm$ 3.93
# of wins	8	5	1	0	2	4	2

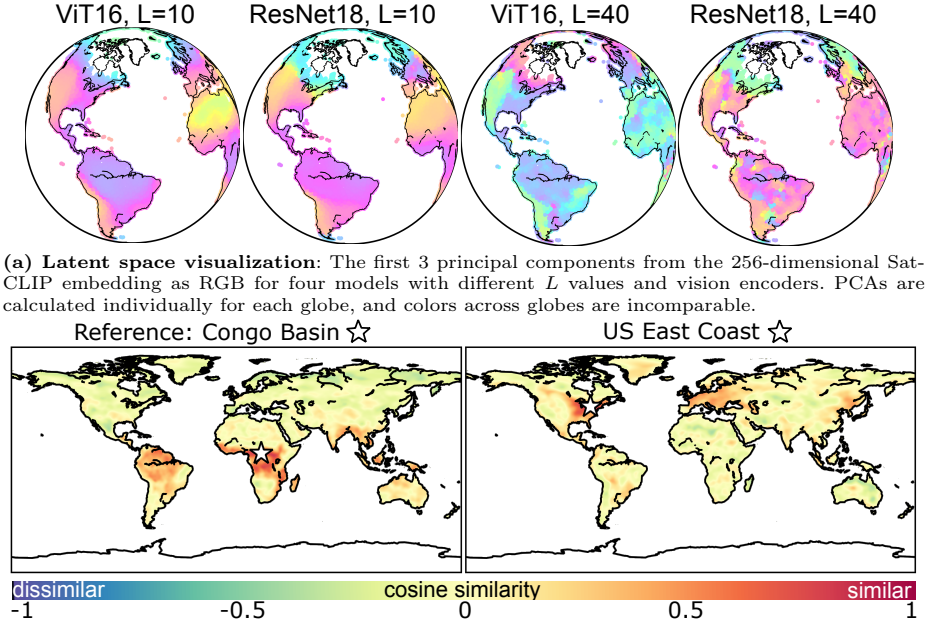


**Fig. 5: Geographic adaptation: predictions of Ecoregions for Africa.** SatCLIP with  $L = 10$  maps Ecoregions in Africa closest to the ground truth, followed by GeoCLIP. MOSAIKS provides predictions that are too fine-grained, and CSP-iNat is too coarse. “X” marked locations in the “True” panel show the sparse training locations in Africa, which are on average 480km apart from their nearest neighbor.

### 4.3 Analysis of Location Embeddings

(RQ 3)

We now investigate qualitatively to what degree the SatCLIP embeddings have learned an implicit representation of different ground conditions in the location encoder weights. We first visualize a low-dimensional projection of the latent representations learned by our location encoders. Fig. 6a shows an RGB representation of the first three principal components of SatCLIP embeddings at locations around the planet. The figure highlights how embeddings learned by SatCLIP provide fine-grained representations of different locations (expressed by different colors), capturing global patterns like climate zones. The figure also



(b) **Location similarities learned by SatCLIP** give insights into areas with similar visual features. Highlighted are similarities between location embeddings for a Congo Basin location (left) and a US East Coast location (right)—marked by a star—and the rest of the world. SatCLIP associates the Congo Basin with other Equatorial locations, e.g., in the Amazon, and associates the US East Coast with other densely populated areas, e.g., in Central Europe.

**Fig. 6: Analysis of SatCLIP location embeddings** through (a) low dimensional PCA projection as RGB visualizations and (b) cosine similarities of different locations of the full 256-dimensional embedding space.

highlights the different spatial smoothness of  $L = 10$  and  $L = 40$  SatCLIP embeddings and changes between models using different vision encoders.

Next, we examine similarities between embeddings of different locations measured in the cosine distance between the embedding of a location  $f_c(\mathbf{c})$  with respect to a reference location  $f_c(\mathbf{c}_*)$ . Fig. 6b shows the similarity of a grid of locations, i.e., the map with respect to reference locations in the Congo Basin and on the east coast of North America on the right panel. The reference locations are marked by a star  $\star$  on the map. SatCLIP location embeddings show high similarity between the Congo Basin location and other areas close to the Equator, particularly the Amazon and Indonesia (red areas on the left panel). In comparison, embeddings of the North American location are similar to areas in Europe or northern China that are similarly population-dense and industrialized.

#### 4.4 Effect of Location and Image Encoder Design on Performance

Lastly, we report what effect the vision encoder type and spatial resolution of the location encoder ( $L$  hyperparameter) have on the downstream accuracy.

**Table 4:** Comparisons of different vision encoders (% accuracy) in spatial interpolation classification tasks (comp. Tab. 2). Full results on all tasks can be found in Appendix E.4.

	ViT16,L=10	ResNet18,L=10	ResNet50,L=10	ViT16,L=40	ResNet18,L=40	ResNet50,L=40
Countries	93.97 $\pm$ 0.30	93.92 $\pm$ 0.30	94.28 $\pm$ 0.18	95.77 $\pm$ 0.14	95.92 $\pm$ 0.10	96.00 $\pm$ 0.14
iNaturalist	65.69 $\pm$ 0.50	65.56 $\pm$ 0.29	65.50 $\pm$ 0.43	65.98 $\pm$ 0.61	66.40 $\pm$ 0.49	66.03 $\pm$ 0.54
Biomes	92.07 $\pm$ 0.22	92.10 $\pm$ 0.23	92.23 $\pm$ 0.26	94.27 $\pm$ 0.15	94.33 $\pm$ 0.10	94.41 $\pm$ 0.14
Ecoregions	89.53 $\pm$ 0.28	89.57 $\pm$ 0.23	89.32 $\pm$ 0.31	91.61 $\pm$ 0.22	91.53 $\pm$ 0.15	91.67 $\pm$ 0.15

The choice of image encoder for pretraining (ViT-16, ResNet-18, ResNet-50) appears to only marginally affect results (performance differences of  $< 1\%$ ). In contrast, different location encoder resolutions with scale parameters  $L = 10$  and  $L = 40$  have a greater effect: the  $L = 40$  location encoders tend to be better for interpolation tasks, as shown in Tab. 4. The  $L = 10$  location encoders are better for the zero- and few-shot adaptation experiments (not shown in the main paper, but in Appendix E.4).

## 5 Discussion

In reference to RQ1, the results presented in Sec. 4.1 show that *SatCLIP models can provide useful information for a wide range of downstream tasks*, with the concrete benefit varying across tasks. SatCLIP outperformed other approaches on globally distributed downstream tasks, while for regional datasets such as Cali. Housing, other methods like GeoCLIP were competitive. SatCLIP embeddings perform well across continents and are less prone to geographic bias in comparison to other methods like GPS2Vec or CSP, where performance degrades outside of Europe or North America.

Regarding RQ2, the results in Sec. 4.2 indicate that the *transfer of spatial patterns in Sentinel-2 imagery into the SatCLIP location encoder enables generalization across geographic areas*. SatCLIP is the best of all models tested under conditions of geographic domain generalization and generally improves upon our satellite image-only (“MOSAICS”) and location-only (“Identity”) baseline methods. At test time, SatCLIP location encoders can be applied directly to any point on the globe, without needing to download additional imagery. Additional experiments showed that the downstream performance of our current SatCLIP is more affected by changes to the location encoder scale factor than the exact vision architecture used for the image encoder. This could relate to the relative model sizes—the vision architectures we use here are much larger than the location encoder architectures. It also suggests that when designing future iterations of SatCLIP and similar models, improvements in vision and location encoder design as well as their balancing during training should be explored. Finally, *qualitative analyses support these interpretations of our results*. Our similarity analysis of location embeddings in Sec. 4.3 showed that spatially far locations like in South America (Amazon rainforest) and Africa (Congo Basin) are embedded

nearby in the SatCLIP embedding space due to the visual similarity of satellite images from these locations. This can explain why SatCLIP models generalize to unseen geographic areas with no or little training data, as demonstrated in our geographic domain generalization experiments.

## 6 Conclusion

We presented a method to learn an implicit neural representation of visual patterns on the globe by matching satellite images and their respective coordinates using a contrastive location-image pretraining objective: SatCLIP. Experiments show that SatCLIP is effective for global prediction tasks spanning social and environmental domains, for both interpolation and out-of-sample geographic prediction, and compared to existing location encoders, image-only and location-only prediction. The effectiveness of SatCLIP is complemented by its relative simplicity in implementation (a single contrastive loss on 100,000 openly available satellite images). Two key factors contribute to the performance observed in Secs. 4.1 and 4.2. First, while our S2-100K pretraining dataset is smaller than the pretraining datasets used for other geographic location encoders, our samples are *uniformly distributed across the globe* in a way that supports globally distributed downstream tasks. Second, we use the recently proposed Siren(SH) location encoder, which has proven well-suited for the global-scale representation of data on the spherical Earth.

Our findings motivate several avenues for potential future work. First, a key limitation of the current implementation of SatCLIP is that it uses Sentinel-2 satellite imagery as the sole source of contextual grounding for our location encoder. In fact, the image encoder we utilize in SatCLIP can be seen as a special case of a more general *context* encoder that may integrate other location-specific data modalities like audio from acoustic sensors or text from geolocated social media posts for multi-source geospatial learning. Second, the current SatCLIP pre-trained weights have limited spatial scales, dictated by the  $L$  parameter of the location encoder. For extremely high resolution or local phenomena, it will be important to further study the effect of different choices of the location encoder and its use for downstream learning. Third, there is the potential to expand the SatCLIP training framework to encode locations in both time and space. Here, we effectively marginalized over all time points in the S2-100k dataset (which is sampled over two years and thus includes seasonal differences in images) but did not directly embed time in a space-time encoder of e.g., the form  $f(lat, lon, time)$ . Addressing each of these three limitations would require separate innovations (on e.g. model architecture), and would need a new set pretraining data and/or downstream tasks for evaluation. Thus, we see these extension as deserving of their own dedicated study in future work. Code for SatCLIP pretraining and downstream experiments as well as the S2-100K dataset is available at [github.com/microsoft/satclip](https://github.com/microsoft/satclip).

# Bibliography

- Md Mahbub Alam, Luis Torgo, and Albert Bifet. A survey on spatio-temporal data analytics systems. *ACM Computing Surveys*, 54(10s):1–38, 2022. 2
- Sara Beery, Guanhang Wu, Trevor Edwards, Filip Pavetic, Bo Majewski, Shreyasee Mukherjee, Stanley Chan, John Morgan, Vivek Rathod, and Jonathan Huang. The auto arborist dataset: A large-scale benchmark for multiview urban forest monitoring under domain shift. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21294–21307, 2022. 6
- Thomas Berg, Jiongxin Liu, Seung Woo Lee, Michelle L Alexander, David W Jacobs, and Peter N Belhumeur. Birdsnap: Large-scale fine-grained visual categorization of birds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2011–2018, 2014. 2
- Tamma Carleton, Trinetta Chong, Hannah Druckenmiller, Eugenio Noda, Jonathan Proctor, Esther Rolf, and Solomon Hsiang. Multi-Task Observation Using Satellite Imagery and Kitchen Sinks (MOSAICS) API. <https://siml.berkeley.edu>, 2022. 8
- Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired alignment between locations and images for effective worldwide geolocalization. *arXiv preprint arXiv:2309.16020*, 2023. 3, 7, 8
- Cen Chen, Kenli Li, Sin G. Teo, Xiaofeng Zou, Kang Wang, Jie Wang, and Zeng Zeng. Gated residual recurrent graph neural networks for traffic prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 485–492. AAAI Press, 2019. 1
- Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6172–6180, 2018. 3, 7, 1
- Elijah Cole, Grant Van Horn, Christian Lange, Alexander Shepard, Patrick Leary, Pietro Perona, Scott Loarie, and Oisín Mac Aodha. Spatial implicit neural representations for global-scale species mapping. In *Proceedings of the 40th International Conference on Machine Learning*, pages 6320–6342. PMLR, 2023. 1, 2, 5
- Eric Dinerstein, David Olson, Anup Joshi, Carly Vynne, Neil D Burgess, Eric Wikramanayake, Nathan Hahn, Suzanne Palminteri, Prashant Hedao, Reed Noss, et al. An ecoregion-based approach to protecting half the terrestrial realm. *BioScience*, 67(6):534–545, 2017. 8, 7
- 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. 6
- A Stewart Fotheringham, Chris Brunson, and ME Charlton. Geographically weighted regression. *The Sage handbook of spatial analysis*, pages 243–254, 2009. 2
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. pages 770–778, 2016. 6
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 6

- Josh Hooker, Gregory Duveiller, and Alessandro Cescatti. Data descriptor: A global dataset of air temperature derived from satellite remote sensing and weather stations. *Scientific Data*, 5:1–11, 2018. [8](#), [7](#)
- Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, Serge Belongie, Caltech 2 Google, and Cornell Tech. The iNaturalist species classification and detection dataset. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8769–8778, 2018. [3](#), [7](#), [8](#), [1](#)
- Junteng Jia and Auston R. Benson. Residual correlation in graph neural network regression. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 588–598. Association for Computing Machinery, 2020. [8](#), [7](#)
- Konstantin Klemmer and Daniel B. Neill. Auxiliary-task learning for geographic data with autoregressive embeddings. In *SIGSPATIAL: Proceedings of the ACM International Symposium on Advances in Geographic Information Systems*, 2021. [2](#)
- Konstantin Klemmer, Tianlin Xu, Beatrice Acciaio, and Daniel B. Neill. Spate-gan: Improved generative modeling of dynamic spatio-temporal patterns with an autoregressive embedding loss. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36:4523–4531, 2022. [2](#)
- Lukas Kondmann, Aysim Toker, Marc Rufwurm, Andrés Camero, Devis Peressuti, Grega Milcinski, Pierre-Philippe Mathieu, Nicolas Longépé, Timothy Davis, Giovanni Marchisio, et al. Denethor: The dynamicearthnet dataset for harmonized, inter-operable, analysis-ready, daily crop monitoring from space. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. [6](#)
- Martha Larson, Mohammad Soleymani, Guillaume Gravier, Bogdan Ionescu, and Gareth J.F. Jones. The benchmarking initiative for multimedia evaluation: Mediaeval 2016. *IEEE Multimedia*, 24:93–96, 2017. [3](#), [7](#), [2](#)
- David B Lobell, David Thau, Christopher Seifert, Eric Engle, and Bertis Little. A scalable satellite-based crop yield mapper. *Remote Sensing of Environment*, 164:324–333, 2015. [1](#)
- Oisín Mac Aodha, Elijah Cole, and Pietro Perona. Presence-only geographical priors for fine-grained image classification. In *ICCV*, 2019. [1](#), [2](#), [8](#)
- Gengchen Mai, Krzysztof Janowicz, Bo Yan, Rui Zhu, Ling Cai, and Ni Lao. Multi-scale representation learning for spatial feature distributions using grid cells. In *Proceedings in the International Conference on Learning Representations (ICLR)*, 2020. [2](#)
- Gengchen Mai, Ni Lao, Yutong He, Jiaming Song, and Stefano Ermon. CSP: Self-supervised contrastive spatial pre-training for geospatial-visual representations. *arXiv preprint arXiv:2305.01118*, 2023. [3](#), [5](#), [7](#), [8](#), [1](#)
- R. Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33:291–297, 2003. [8](#), [7](#)
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 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 *Proceedings in the International Conference on Machine Learning (ICML)*, pages 8748–8763. PMLR, 2021. [2](#), [5](#)
- Gyri Reiersen, David Dao, Björn Lütjens, Konstantin Klemmer, Kenza Amara, Attila Steinegger, Ce Zhang, and Xiaoxiang Zhu. Reforestree: A dataset for estimating tropical forest carbon stock with deep learning and aerial imagery. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36:12119–12125, 2022. [1](#)



- Esther Rolf, Jonathan Proctor, Tamma Carleton, Ian Bolliger, Vaishaal Shankar, Miyabi Ishihara, Benjamin Recht, and Solomon Hsiang. A generalizable and accessible approach to machine learning with global satellite imagery. *Nature Communications* 2021 12:1, 12:1–11, 2021. 1, 8, 7, 9
- Marc Rußwurm, Sherrie Wang, Marco Korner, and David Lobell. Meta-learning for few-shot land cover classification. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition workshops*, pages 200–201, 2020. 6
- Marc Rußwurm, Konstantin Klemmer, Esther Rolf, Robin Zbinden, and Devis Tuia. Geographic location encoding with spherical harmonics and sinusoidal representation networks. *Proceedings in the International Conference on Learning Representations (ICLR)*, 2024. 5, 3, 4
- Vincent Sitzmann, Julien N P Martel, Alexander W Bergman, David B Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020. 5
- Bart Thomee, Benjamin Elizalde, David A. Shamma, Karl Ni, Gerald Friedland, Douglas Poland, Damian Borth, and Li Jia Li. YFCC100M. *Communications of the ACM*, 59:64–73, 2016. 3, 7, 2
- Gabriel Tseng, Hannah Kerner, and David Rolnick. TIML: Task-informed meta-learning for agriculture. *arXiv preprint arXiv:2202.02124*, 2022. 1, 6
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11), 2008. 5
- Yi Wang, Nassim Ait, Ali Braham, Zhitong Xiong, Chenying Liu, Conrad M Albrecht, and Xiao Xiang Zhu. SSL4EO-S12: A large-scale multi-modal, multi-temporal dataset for self-supervised learning in Earth observation. 2022a. 6
- Yi Wang, Conrad M. Albrecht, Nassim Ait Ali Braham, Lichao Mou, and Xiao Xiang Zhu. Self-supervised learning in remote sensing: A review. *IEEE Geoscience and Remote Sensing Magazine*, 10:213–247, 2022b. 3
- Yifang Yin, Zhenguang Liu, Ying Zhang, Sheng Wang, Rajiv Ratn Shah, and Roger Zimmermann. GPS2Vec: Towards generating worldwide GPS embeddings. In *SIGSPATIAL: Proceedings of the ACM International Symposium on Advances in Geographic Information Systems*, pages 416–419. Association for Computing Machinery, 2019. 3, 7, 8, 2
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, Lucas Beyer, and Google Deepmind. Sigmoid loss for language image pre-training. In *ICCV*, 2023. 6, 3

## A S2-100K Dataset Overview

We sample the S2-100k dataset from all Sentinel 2 Level-2A scenes that meet the following criteria:

1. Are captured between January 1st, 2021 and May 17th, 2023
2. Are estimated to have less than 20% cloud cover (as reported by the European Space Agency preprocessing)
3. Are at least partially over a land mass (measured with country level boundaries from [gadm.org](https://gadm.org))

Using the Microsoft Planetary Computer, we find 2,359,972 S2-L2A scenes that meet these criteria. We sample a random scene, sample a random  $256 \times 256$  pixel patch from that scene, then keep the patch if less than 10% of the pixels include nodata values, else we reject the patch. We repeat this process until we have 100,000 patches. Each patch contains all available bands – B01, B02, B03, B04, B06, B06, B07, B08, B08A, B09, B11, and B12 – resampled to a 10m/px spatial resolution in the native UTM coordinate system and saved as a single Cloud Optimized GeoTIFF. Finally, we record the latitude and longitude (using the EPSG:4326 coordinate system) of the center point of the patch. The dataset can be downloaded at [github.com/microsoft/satclip](https://github.com/microsoft/satclip). The final dataset is 83.79 Gigabytes on disk.

## B Further Details on Other Geographic Location Encoders & Their Pretraining Datasets

### B.1 (Pre-)training Data

Here, we describe the datasets used for pretraining in previous geographic location encoders, CSP and GPS2Vec, in detail:

**iNat 2018** ([Horn et al., 2018](#)): This dataset contains crowdsourced, natural images of over 8,000 plant and animal species from around the globe. The iNat 2018 training set contains over 400,000 image-location pairs. While species distributions can be indicative of e.g., climate zones ([Mac Aodha et al., 2019](#)), they are less descriptive of socio-economic processes, as we show in our experiments. Furthermore, most iNat image locations lie in North America and Europe, with only few observations in other continents, as we show in Fig. 3. We obtain location encoders pretrained on iNat from CSP ([Mai et al., 2023](#)).

**FMoW** ([Christie et al., 2018](#)): This dataset contains satellite images annotated with bounding boxes representing 63 categories describing functional or land-use characteristics (e.g., “airport” or “gas station”). The dataset contains over 363,000 image-location pairs. FMoW is aimed to support detection of functional objects from satellite imagery and, as such, is skewed towards human-built infrastructure and includes fewer natural scenes. It is also again heavily favoring Western countries in its geographic distribution. We obtain location encoders pretrained on FMoW from CSP [Mai et al. \(2023\)](#).

**YFCC100M** (Thomee et al., 2016): This dataset contains natural images and associated semantic tags collected from the social media platform Flickr. YFCC100M includes approximately 48 million image-tag-location triplets. Images (and tags) represent diverse scenes, from street views to birthday parties. While they are more indicative of social dimensions and physical infrastructure, YFCC100M imagery is less representative of natural processes. Just like iNat and FMoW, the dataset again mostly contains locations in Western countries. We obtain models separately trained on images and semantic tags from GPS2Vec Yin et al. (2019).

**MP-16** (Larson et al., 2017): The MediaEval Placing Tasks 2016 (MP-16) dataset also consists of images from Flickr. Overall, it includes 4.27 million images and their corresponding geo-tag. Like the YFCC100M dataset used for training GPS2Vec models, MP-16 is geographically clustered in areas with higher Flickr activity and overrepresents Western countries.

## B.2 CSP

**Location and image encoders:** The authors use a **location encoder** combining sinusoidal transforms introduced by (Mai et al., 2020) with a fully-connected neural network. The **image encoder** is dataset-dependent. For the iNat dataset, the authors use a pretrained IncptionV3 network. For the FMoW dataset, they use a pretrained ResNet50 network. In both cases, neural network weights in all layers except the last linear projection layer are frozen.

**Learning objective:** Here, the authors combine (1) the standard CLIP loss, which leverages in-batch negative sampling and which is also used by us, with two other contrastive objectives: (2) Random negative location sampling (i.e. predicting the real location from a set of randomly generated locations) and (3) SimCSE sampling (i.e. matching location embeddings obtained using different dropout masks). The two extra objectives help the model to balance the location and context encoders. We find that the two additional objectives are not needed for training SatCLIP.

## B.3 GPS2Vec

**Location and context encoders:** The authors train **location encoders** specific to the UTM zones containing training data. They first encode coordinates as an exponential function of the Euclidean distance between a latitude/longitude coordinate and its respective UTM zone centroid coordinate. This is followed by a simple ReLU network with three hidden layers. GPS2Vec does not learn a **context encoder** but directly extracts text features using a vocabulary-based approach and image features using convolutional neural networks (CNN).

**Learning objective:** The authors design an objective which aims to use contextual features (images and semantic tags) as labels to train their location encoder by estimating the normalized frequency of features in the vicinity of a given location. Practically, this is achieved by minimizing the KL-divergence of context and location embedding distributions.

## B.4 GeoCLIP

**Location and image encoders:** The **location encoder** in GeoCLIP first transforms raw latitude/longitude coordinates into equal earth projection (EEP), then extracts features using Random Fourier Features (RFF). On top of this positional encoding, the authors deploy separate MLPs for each RFF which are aggregated into the final location embedding. For their **image encoder**, the authors use a CLIP pretrained ViT16 models. Just like for SatCLIP, the image encoder is frozen except the last two linear projections layers.

**Learning objective:** GeoCLIP models are trained using the CLIP objective, outlined in Eq. (1). The only difference to SatCLIP training is that the authors add additional location negatives (randomly sampled) to each batch.

## C SatCLIP (Pre)-training

### C.1 Training Details

**Batch size:** After experimenting with different batch sizes, we opt for models trained at batch sizes of  $8k$ . While traditional CLIP image-text pretraining behaves optimally at a batch size of  $32k$  (Zhai et al., 2023), we observe most effective training at  $8k$  and  $16k$ .

**Image encoder:** We train SatCLIP models with ViT16, ResNet18 and ResNet50 image encoders, all pretrained on Sentinel-2 imagery and published by Wang et al. (2022b). We keep the image encoders frozen during training, and only train a last projection layer that maps the image embeddings into the desired output space. We find this to be ideal for training at a size of 256—this is equivalent to the embedding size used by CSP Mai et al. (2023).

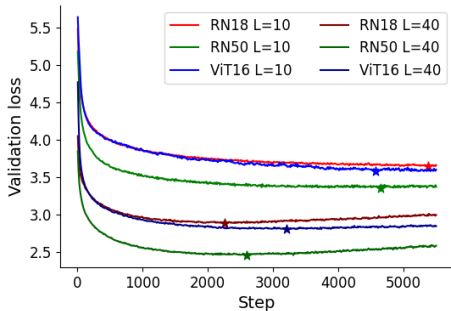
**Location encoder:** Our location encoder follows an approach recently proposed by Rufwurm et al. (2024). It first puts raw latitude/longitude coordinates through a functional transform based on orthogonal spherical harmonics. The number of spherical harmonics used corresponds to the resolution at which the model is trained. In practice this is controlled by a parameter defining the number of Legendre polynomials  $L$  used for spherical harmonics computation. We train a lower-resolution ( $L = 10$ ) and a higher-resolution ( $L = 40$ ) version of SatCLIP. The functional transform is followed by a sinusoidal representation network (Siren) consisting of two hidden layers and 512 hidden dimensions. These hyperparameters are obtained after rigorously testing different settings.

**Augmentations:** We deploy several data augmentations within our training procedure. Image augmentations include random crops, random horizontal flipping, random vertical flipping and Gaussian blurs. Point coordinates are augmented using a coordinate jitter which randomly shifts image coordinates by up to about 1km.

**Optimization:** All SatCLIP models are trained with the Adam optimizer, a learning rate of 0.0001 and a relatively high weight-decay values of 0.01 to help prevent overfitting. All final SatCLIP models are trained for 500 epochs. On our single A100 GPU, training takes around 2 days. Throughout training, we reserve

10% of the data for validation. We monitor validation loss and select SatCLIP models according to the minimum validation loss.

## C.2 Scale Sensitivity



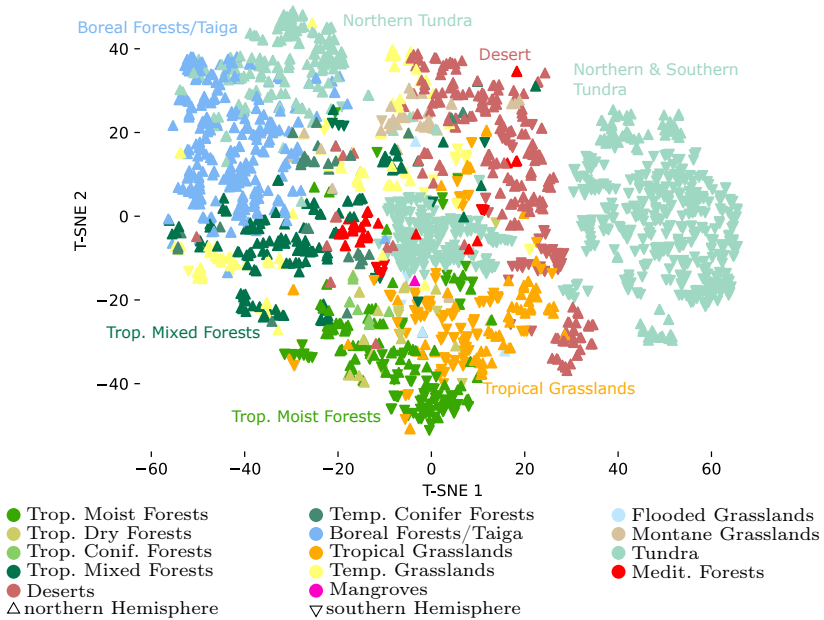
**Fig. 7:** Validation set loss of SatCLIP models trained with different vision encoder backbones (ResNet18, ResNet50 and ViT16) with  $L = 10$  and  $L = 40$ . Stars indicate the minimum validation loss, which we use for model selection.

We want to briefly comment on the sensitivity of SatCLIP training to the number of legendre polynomials  $L$  used in the model’s spherical harmonics location encoder. This hyperparameter effectively controls the spatial resolution of the obtained embeddings. Smaller values of  $L$  are computationally more efficient and ideal for representing large-scale patterns, while larger values of  $L$  are better for capturing small-scale patterns. More details on this can be found in [Rußwurm et al. \(2024\)](#). During SatCLIP training, we observe some interesting differences between smaller ( $L = 10$ ) and higher ( $L = 40$ ) resolution models. We find that higher-resolution models are more likely to exhibit overfitting, as Fig. 7 highlights. On downstream tasks, higher resolution models perform better at smaller scale, regional tasks as e.g. the Cali. Housing dataset, as Tab. 6 highlights. Lastly,  $L = 40$  models appear better for spatial interpolation (RQ1), while  $L = 10$  models seem better suited for geographic generalization (RQ2).

## D Latent Space Exploration

### D.1 SatCLIP Embeddings in Different Biomes

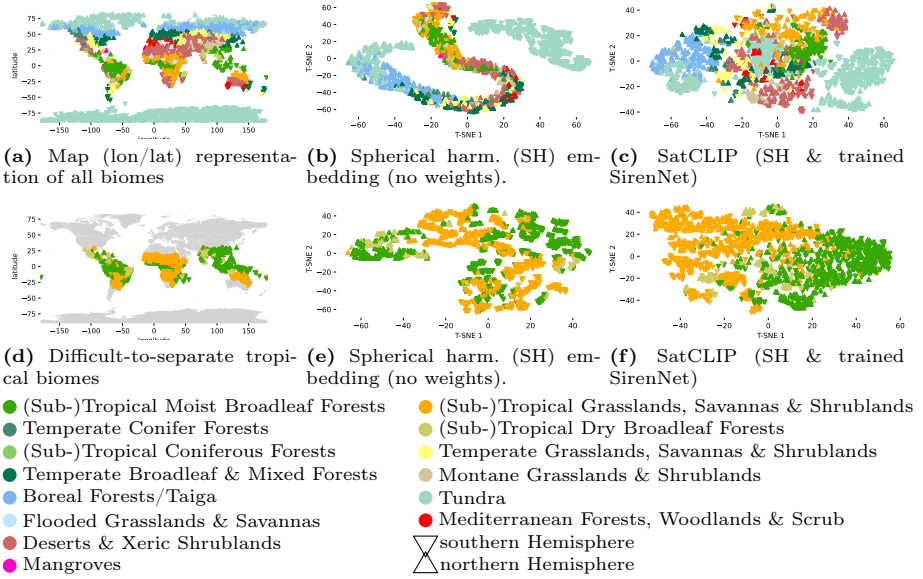
This supplementary experiment disentangles the influence of raw coordinates and learned SatCLIP embeddings in the T-SNE embedding experiment of Fig. 8 in the main paper that indicated that the SatCLIP *weights* encode environmental factors like dryness or temperature in biomes. Critically speaking, also the



**Fig. 8: Analysis of environmental factors** learned by the trained SatCLIP location encoder by a T-SNE (Van der Maaten and Hinton, 2008) dimensionality reduction at globally distributed points colored by biome type. Overall, biomes are clearly separated in this T-SNE embedding space. Note that some points of the northern  $\Delta$  and southern  $\nabla$  hemispheres are mapped nearby, which indicates that this clustering is due to the trained SatCLIP weights and not the pure coordinates, as we investigate further in the appendix Fig. 9.

direct embedding of coordinates *without any weights* could generate T-SNE clusters that look similar to Fig. 8. In Fig. 9, we investigate if this is the case by showing the embeddings of coordinates with Spherical Harmonics (SH) (no trainable weights) alone (Figs. 9b and 9e) and compare it with the trained SatCLIP location encoder, which combines Spherical Harmonics basis functions with Sinusoidal Representation Networks (SirenNets) (Sitzmann et al., 2020) (Figs. 9c and 9f).

In the top row Figs. 9a to 9c we show all biomes and can see that the T-SNE points embedded in spherical harmonic basis functions of pure coordinates (without trainable weights) (Fig. 9b) creates a clear manifold along which the different biomes cluster. However, notice that points from the northern hemisphere ( $\Delta$ ) are clearly separated from the southern hemisphere ( $\nabla$ ). Also, the arctic and antarctic tundra biomes are separated by two clusters. This is the T-SNE embedding produced from pure coordinate embeddings, that serves as the baseline to the SatCLIP embeddings that are produced by representing point coordinates as spherical harmonic basis functions *and transforming these basis*



**Fig. 9:** Supplementary experiment to Fig. 8: we disentangle the influence of spherical harmonic (SH) coordinate embedding without trainable weights in Figs b) & e) with the trained SatCLIP location encoder in Figs c) & f). Comparing the T-SNE embedding projections in these two columns reveals the additional effect of having trained SatCLIP weights in the location encoder. We perform this analysis in an easier setting with all biomes (top row) and with only tropical geographically intertwined biomes (bottom row). In the top row, also the pure spherical harmonic basis functions produce a T-SNE manifold along which biomes are located (b). In the more difficult bottom row, no coherent manifolds form in e) along biome types. These appear only in combination with the SatCLIP-trained weights in f)



functions with the SatCLIP-trained SirenNet neural network. The result is shown in Fig. 9c, which corresponds to Fig. 8 in the main paper, where the embedding clusters along dryness and temperature dimensions rather than pure geographic location. Hence, by comparing Fig. 9b (no weights) with Fig. 9c, we can see the effect of trainable SatCLIP weights on the T-SNE representation of the points.

To make sure that this analysis holds, we repeat the embedding in a more challenging setting: we show embeddings of points from only tropical biomes for (Sub-)Tropical Moist Broadleaf Forests, (Sub-)Tropical Dry Broadleaf Forests, and (Sub-)Tropical Grasslands, Savannas & Shrublands in Figs. 9d to 9f (second row). Points from these biomes are not separable using only non-parametric embeddings (Fig. 9e), since these biomes are all tropical and geographically intertwined. Here only the geographic coordinate is not sufficiently expressive to form T-SNE manifolds along biome lines. However, with the pretrained SatCLIP location encoder (Fig. 9f), we can clearly identify a dry to wet trend along T-SNE dimension 1 (left to right) where points from grassland/savanna transitions first into dry broadleaf forest and then moist broadleaf forest.

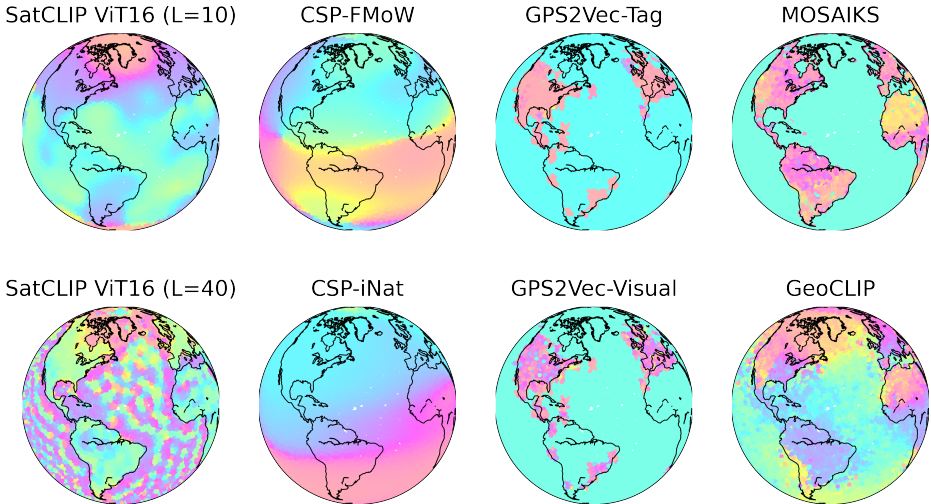
In summary, these supplementary analyses reinforce the conclusions of the main paper experiment around Fig. 8, which experimentally supports the concept that the SatCLIP embeddings represent an implicit neural representation of environmental and societal ground conditions that are visible in the Sentinel-2 images that the SatCLIP model is trained on.

## D.2 Latent Space Visualization at Different Scales

Lastly, we provide a visualization of the embeddings learned by SatCLIP models at different scales (different  $L$  values) in Fig. 10. It is clearly visible how SatCLIP  $L = 40$  resolves positional embeddings much more fine-grained. Lastly, in Fig. 11 we show the explained variance ratio by element in a respective embedding vector, extracted via PCA, including SatCLIPs with different  $L$  values and comparison methods.

**Table 5:** Description of the datasets used in downstream experiments.

Name	$n$	Spatial coverage	Locations	Outcome variable	Inputs	Task
Air Temp. (Hooker et al., 2018)	3,076	Global	Weather stations	Ann. mean temp.	lon, lat	Regr.
Med. Income (Jia and Benson, 2020)	3,107	Cont., USA	Census tract	Med. house inc.	lon, lat	Regr.
Cali. Housing (Pace and Barry, 2003)	20,640	Calif., USA	House locations	House price	lon, lat	Regr.
Elev. (Rolf et al., 2021)	99,995	Global	Reg. sampled	Elevation	lon, lat	Regr.
Pop. Dens. (Rolf et al., 2021)	74,512	Global	Reg. sampled	Pop. dens.	lon, lat	Regr.
Countries (Ours)	100,000	Global	Reg. sampled	Country code	lon, lat	Class.
Biome (Dinerstein et al., 2017)	100,000	Global	Reg. sampled	Biome type	lon, lat	Class.
Ecoregions (Dinerstein et al., 2017)	100,000	Global	Reg. sampled	Ecoregion	lon, lat	Class.
iNat 2018 (Horn et al., 2018)	460,406	Global	Image locations	Species classes	lon, lat, image	Class.



**Fig. 10:** Visualization of the top-3 principal components, plotted as RGB channels, of SatCLIP models trained with different  $L$  hyperparameters and comparison methods.

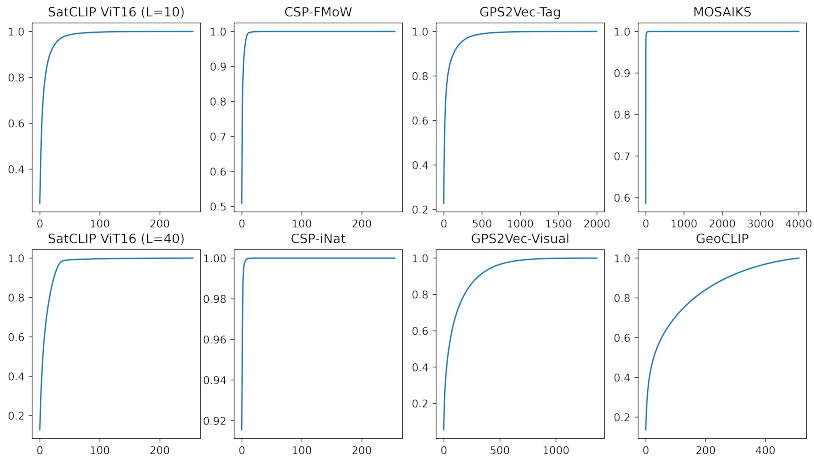
## E Details and Additional Results for Downstream Experiments

### E.1 Countries, Ecoregions, and Biome Dataset Overview

To create the Countries, Ecoregions and Biome datasets used in Section 4 we sample a dataset of 100,000 latitude and longitude points approximate uniformly at random, then record which country, ecoregion and biome each point falls within, or whether the point is over an ocean. We use country boundaries from the 4.1 release of “The Database of Global Administrative Areas” (GADM, [gadm.org](https://gadm.org)) for country codes. We rely on [Dinerstein et al. \(2017\)](#) for ecoregion and biome maps. The Countries dataset contains points over 184 countries (out of a total of 263 countries in the GADM dataset), resulting in a classification task with 185 categories (including an “ocean” class). The Ecoregions dataset contains 719 classes, also including an “ocean” class. The Biome dataset contains 15 classes and also includes an “ocean” class. Note that the latitude/longitude locations for all three datasets are the same.

### E.2 Downstream Task Overview

Tab. 5 highlights the different datasets used in downstream tasks throughout our experiments along with characteristics like their spatial coverage, outcome variables and task types. Generally, our tasks can be split into regression and (multi-)classification tasks. Our inputs are always raw longitude/latitude coordinate pairs which are processed by our pretrained location encoders to obtain location embeddings which are then used as inputs for downstream learners.



**Fig. 11:** Curves of the explained variance ratio of a principal component analysis (PCA) conducted on the different embeddings using globally distributed locations from the Population dataset.

The population density and elevation datasets are subsampled from the global datasets from [Rolf et al. \(2021\)](#). Only the iNat 2018 task contains additional data: image embeddings obtained via a pretrained InceptionV3 network.

### E.3 Training details

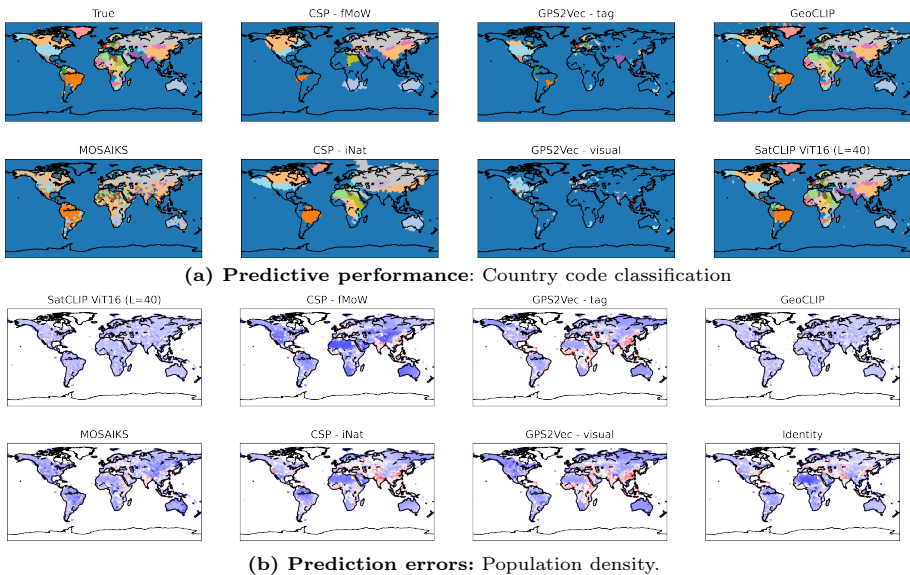
We tune all downstream task models before final training runs. All reported downstream experiments use simple MLP models. We tune the following hyperparameters by random search: number of hidden layers, hidden dimensions, learning rate, weight decay. The best setting is selected according to a validation loss. We then train 10 models until convergence (as measured on the validation loss) and report mean and standard deviations of performance metrics ( $R^2$ , MSE or accuracy). All tuning and training is conducted on a single A100 GPU.

For experiments reported in Sec. 4.1, we deploy random train/val/test splits. For the Air Temp. dataset the training size is 60%, for the Cali. Housing and Med. Income datasets it is 50% and for all other datasets it is 30%. The validation set size is 20% for Air Temp., Cali. Housing and Med. Income and 10% for all other datasets. Note that the iNat 2018 dataset comes with predefined train and test sets. Here we reserve 10% of the train set for validation. For experiments reported in Sec. 4.2 train and test set are defined by the location of points within or outside of the respective test set continent. Here, we reserve 20% of the training set for validation for Air Temp., Cali. Housing and Med. Income and 10% for all other datasets. For Countries and Ecoregions we allocated a small portion (1%) of test set points for training to allow the models to learn classes that only exist on the test set continent.

## E.4 Results of Different SatCLIP Configurations

We provide full experimental results for pretrained SatCLIPs (with different numbers of Legendre polynomials  $L$  and different vision backbones ResNet18, ResNet50 and ViT16) in Tab. 6.

## E.5 Additional Results: Predictive Performance



**Fig. 12:** Fig. 12a shows results from the Countries dataset. Shown are predicted country codes across the planet, compared to the true country codes. Fig. 12b shows test set results from the Population dataset. Shown are model residuals (errors), red values indicate positive, blue values negative errors. Errors are standardized across figures.

Fig. 12 shows results (signed errors) for air temperature, population density and country code prediction, where the improvements from SatCLIP are visually apparent. Tab. 7 shows test set results on all datasets. One interesting observation is that this use of SatCLIP embeddings are more informative for iNat classification than a location encoder pretrained on iNat (CSP-iNat). This is intuitive as SatCLIP embeddings might be able to provide auxiliary information not contained within the iNat imagery. Overall, the results confirm that SatCLIP models trained on SK-100K data provide meaningful features to help with prediction in both natural (e.g., Air Temp., Elevation) and socio-economic (e.g., Med. Income, Cali. Housing) settings.

**Table 6: SatCLIP models with different vision encoders and location encoders win varying resolutions controlled by the  $L$  parameter (maximum number of Legendre polynomial degrees). We report average test set MSE and accuracy  $\pm 1$  standard deviation across 10 independently initialized MLP training runs.**

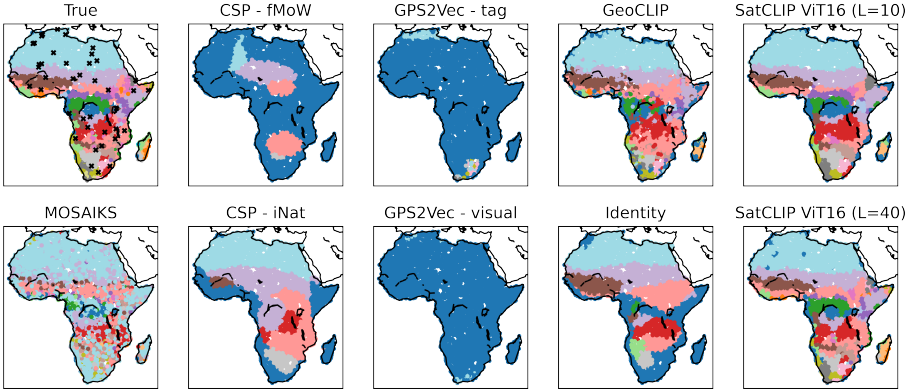
(a) Classification Accuracy in % and Regression MSE scores.

Vision Encoder	ViT16	ViT16	ResNet18	ResNet18	ResNet50	ResNet50
Spherical Harm. $L$	10	40	10	40	10	40
<b>Regression MSE <math>\downarrow</math></b>						
Air Temp.	<b>0.20 <math>\pm</math> 0.31</b>	0.25 $\pm$ 0.02	0.59 $\pm$ 0.75	0.81 $\pm$ 0.77	0.29 $\pm$ 0.38	0.27 $\pm$ 0.03
Median Income	0.92 $\pm$ 0.24	<b>0.67 <math>\pm</math> 0.01</b>	0.87 $\pm$ 0.22	<b>0.67 <math>\pm</math> 0.01</b>	0.78 $\pm$ 0.02	0.71 $\pm$ 0.16
Cali Housing	3.55 $\pm$ 0.06	2.62 $\pm$ 0.28	3.66 $\pm$ 0.13	<b>2.46 <math>\pm</math> 0.08</b>	3.68 $\pm$ 0.23	<b>2.42 <math>\pm</math> 0.12</b>
Elevation	0.22 $\pm$ 0.01	<b>0.15 <math>\pm</math> 0.01</b>	0.20 $\pm$ 0.01	<b>0.15 <math>\pm</math> 0.00</b>	0.21 $\pm$ 0.01	<b>0.15 <math>\pm</math> 0.00</b>
Population Density	0.56 $\pm$ 0.01	0.50 $\pm$ 0.02	0.55 $\pm$ 0.01	<b>0.48 <math>\pm</math> 0.01</b>	0.55 $\pm$ 0.01	<b>0.48 <math>\pm</math> 0.01</b>
<b>Classification % Acc. <math>\uparrow</math></b>						
Countries	93.97 $\pm$ 0.30	95.77 $\pm$ 0.14	93.92 $\pm$ 0.30	95.92 $\pm$ 0.10	94.28 $\pm$ 0.18	<b>96.00 <math>\pm</math> 0.14</b>
iNaturalist	65.69 $\pm$ 0.50	65.98 $\pm$ 0.61	65.56 $\pm$ 0.29	<b>66.40 <math>\pm</math> 0.49</b>	65.50 $\pm$ 0.43	66.03 $\pm$ 0.54
Biomes	92.07 $\pm$ 0.22	<b>94.27 <math>\pm</math> 0.15</b>	92.10 $\pm$ 0.23	<b>94.33 <math>\pm</math> 0.10</b>	92.23 $\pm$ 0.26	<b>94.41 <math>\pm</math> 0.14</b>
Ecoregions	89.53 $\pm$ 0.28	<b>91.61 <math>\pm</math> 0.22</b>	89.57 $\pm$ 0.23	<b>91.53 <math>\pm</math> 0.15</b>	89.32 $\pm$ 0.31	<b>91.67 <math>\pm</math> 0.15</b>

(b) Held-out countries regression and classification tasks.

Vision Encoder	ViT16	ViT16	ResNet18	ResNet18	ResNet50	ResNet50
Spherical Harm. $L$	10	40	10	40	10	40
<b>Asia</b>						
Air Temp.* MSE $\downarrow$	<b>0.85 <math>\pm</math> 0.17</b>	1.26 $\pm$ 0.15	1.22 $\pm$ 0.56	1.56 $\pm$ 0.14	1.17 $\pm$ 0.52	1.50 $\pm$ 0.10
Elevation*	2.13 $\pm$ 0.37	<b>2.06 <math>\pm</math> 0.28</b>	2.36 $\pm$ 0.21	2.58 $\pm$ 0.26	2.22 $\pm$ 0.35	3.28 $\pm$ 0.09
Population Density*	2.03 $\pm$ 0.26	<b>1.94 <math>\pm</math> 0.15</b>	2.18 $\pm$ 0.21	<b>1.92 <math>\pm</math> 0.23</b>	2.06 $\pm$ 0.21	2.82 $\pm$ 0.22
Countries $^\dagger$ % Acc. $\uparrow$	36.90 $\pm$ 4.32	19.17 $\pm$ 2.82	37.15 $\pm$ 3.28	13.71 $\pm$ 3.26	<b>40.36 <math>\pm</math> 3.21</b>	14.29 $\pm$ 1.62
iNaturalist*	19.60 $\pm$ 0.78	<b>20.91 <math>\pm</math> 0.77</b>	<b>20.75 <math>\pm</math> 1.00</b>	19.67 $\pm$ 0.97	19.75 $\pm$ 0.49	17.67 $\pm$ 0.32
Biome*	25.89 $\pm$ 2.79	16.44 $\pm$ 1.21	21.97 $\pm$ 3.50	15.02 $\pm$ 0.97	21.83 $\pm$ 3.44	<b>30.26 <math>\pm</math> 3.00</b>
Ecoregions $^\dagger$	<b>21.02 <math>\pm</math> 1.09</b>	10.86 $\pm$ 1.19	<b>20.54 <math>\pm</math> 1.55</b>	8.87 $\pm$ 1.70	<b>20.32 <math>\pm</math> 1.41</b>	8.46 $\pm$ 0.79
<b>Africa</b>						
Air Temp.* MSE $\downarrow$	4.13 $\pm$ 1.66	<b>1.79 <math>\pm</math> 0.50</b>	3.36 $\pm$ 0.91	17.60 $\pm$ 25.27	6.56 $\pm$ 1.53	2.17 $\pm$ 0.33
Elevation*	1.34 $\pm$ 0.83	<b>0.57 <math>\pm</math> 0.04</b>	1.25 $\pm$ 0.36	0.73 $\pm$ 0.05	1.41 $\pm$ 0.51	0.81 $\pm$ 0.06
Pop. D.*	<b>1.97 <math>\pm</math> 0.29</b>	<b>1.96 <math>\pm</math> 0.22</b>	2.39 $\pm$ 0.46	2.14 $\pm$ 0.25	2.63 $\pm$ 0.50	2.99 $\pm$ 0.23
Countries $^\dagger$ % Acc. $\uparrow$	30.65 $\pm$ 4.23	10.22 $\pm$ 1.62	<b>34.41 <math>\pm</math> 3.18</b>	9.96 $\pm$ 1.56	<b>35.30 <math>\pm</math> 2.00</b>	8.95 $\pm$ 1.04
iNaturalist*	9.53 $\pm$ 0.57	6.23 $\pm$ 0.47	8.96 $\pm$ 0.88	6.68 $\pm$ 0.44	<b>9.90 <math>\pm</math> 0.52</b>	5.22 $\pm$ 0.26
Biome*	35.72 $\pm$ 5.48	12.34 $\pm$ 1.75	34.58 $\pm$ 5.14	20.60 $\pm$ 1.40	<b>39.93 <math>\pm</math> 4.41</b>	33.77 $\pm$ 2.69
Ecoregions $^\dagger$	<b>32.03 <math>\pm</math> 1.19</b>	12.91 $\pm$ 1.63	<b>31.66 <math>\pm</math> 2.29</b>	12.87 $\pm$ 2.47	<b>32.69 <math>\pm</math> 1.79</b>	13.54 $\pm$ 2.06

\* Denotes zero-shot,  $^\dagger$  denotes few-shot domain generalization tasks.



**Fig. 13: Few-shot ecoregion prediction.** Predictions on the test continent Africa are shown for the different embeddings. Models are trained on only 1% (marked by black crosses) of the data points on the test continent to evaluate their capacity for generalizing to an unseen environment.

## E.6 Additional Results: Geographic Generalization

Fig. 13 shows results for the few-shot domain generalization setting with the Ecoregions dataset. Here, we highlight true values and predictions on the test set continent Africa. We can see that SatCLIP embeddings perform best, followed by GeoCLIP and MOSAICS embeddings. This experiment helps us to evaluate the embeddings capacity for overcoming geographic distribution shift. The full results table for this experiment can be found in Tab. 8.

## E.7 Additional Results: Combining Embeddings

We test whether combinations of embeddings obtained from location encoders trained on different datasets further improves performance. We do this by con-

**Table 7: Downstream task performance using SatCLIP (with ViT16 vision encoder) vs. baseline location embeddings.** We report average test set  $R^2$  and accuracy  $\pm 1$  standard deviation across 10 independently initialized MLP training runs.

Task ↓ Data →	SatCLIP <sub>L=10</sub> (S2-100K)	SatCLIP <sub>L=40</sub> (S2-100K)	CSP (fMoW)	CSP (iNat)	GPS2Vec (tag)	GPS2Vec (visual)	MOSAICS (Planet)	GeoCLIP (MP-16)	Identity (y ~ g(c))
<b>Regression</b>	$R^2 \uparrow$								
Air temperature	<b>0.90 ± 0.13</b>	<b>0.91 ± 0.01</b>	0.07 ± 0.37	-0.56 ± 0.59	0.22 ± 0.00	0.03 ± 0.00	-0.52 ± 2.00	-3.11 ± 5.24	0.82 ± 0.16
Median income	0.42 ± 0.01	<b>0.47 ± 0.12</b>	-0.03 ± 0.05	-0.01 ± 0.02	0.21 ± 0.00	0.03 ± 0.00	0.02 ± 0.05	<b>0.50 ± 0.01</b>	-0.84 ± 0.94
Cali. housing	0.35 ± 0.04	0.57 ± 0.02	0.00 ± 0.00	-0.00 ± 0.00	0.71 ± 0.03	0.61 ± 0.03	0.24 ± 0.02	<b>0.75 ± 0.01</b>	0.05 ± 0.02
Elevation	0.83 ± 0.01	<b>0.88 ± 0.00</b>	0.36 ± 0.04	0.11 ± 0.05	0.10 ± 0.00	0.06 ± 0.00	0.21 ± 0.01	0.83 ± 0.00	0.25 ± 0.08
Population	<b>0.79 ± 0.00</b>	<b>0.82 ± 0.00</b>	0.37 ± 0.06	0.36 ± 0.11	0.25 ± 0.00	0.15 ± 0.00	0.46 ± 0.02	0.79 ± 0.00	0.46 ± 0.03
<b>Classification</b>	% Accuracy ↑								
Countries	94.28 ± 0.18	<b>96.00 ± 0.14</b>	77.78 ± 1.66	82.11 ± 1.72	70.35 ± 0.06	67.80 ± 0.03	76.16 ± 0.50	90.72 ± 0.44	82.94 ± 2.23
iNaturalist	<b>65.69 ± 0.18</b>	66.22 ± 0.40	56.73 ± 0.83	60.47 ± 0.56	58.78 ± 0.48	53.39 ± 0.67	56.73 ± 0.8	62.01 ± 0.59	60.83 ± 0.53
Biome	92.23 ± 0.26	<b>94.41 ± 0.14</b>	75.81 ± 1.53	73.18 ± 5.58	69.69 ± 0.06	68.29 ± 0.11	79.61 ± 0.42	89.57 ± 0.45	83.55 ± 2.43
Ecoregions	89.32 ± 0.31	<b>91.67 ± 0.15</b>	76.87 ± 1.27	78.43 ± 1.71	68.46 ± 0.06	67.26 ± 0.02	70.48 ± 0.21	84.65 ± 0.32	77.07 ± 2.54

**Table 8: Geographic adaptation capabilities of SatCLIP (with ViT16 vision encoder) vs. baseline location embeddings to new geographic areas with no (\*) or very few (†) samples from the held-out test continent.** We report average test set  $R^2$  and accuracy in %  $\pm 1$  standard deviation across 10 independently initialized fine-tuning runs.

Test Continent	SatCLIP <sub>L=10</sub> (S2-100K)	SatCLIP <sub>L=40</sub> (S2-100K)	CSP (FMoW)	CSP (iNat)	GPS2Vec (tag)	GPS2Vec (visual)	MOSAICS (Planet)	GeoCLIP (MP-16)	Identity ( $y \sim g(c)$ )
<b>Asia</b>									
Air Temp.* $R^2 \uparrow$	<b>0.75 <math>\pm</math> 0.05</b>	0.63 $\pm$ 0.04	0.09 $\pm$ 0.37	-0.50 $\pm$ 1.32	-3.95 $\pm$ 4.89	-0.23 $\pm$ 0.17	-2.13 $\pm$ 3.50	<b>0.77 <math>\pm</math> 0.28</b>	0.20 $\pm$ 1.64
Elevation*	<b>0.46 <math>\pm</math> 0.09</b>	<b>0.48 <math>\pm</math> 0.07</b>	-0.20 $\pm$ 0.07	-0.26 $\pm$ 0.03	-0.29 $\pm$ 0.01	-0.25 $\pm$ 0.00	-0.07 $\pm$ 0.06	<b>0.50 <math>\pm</math> 0.03</b>	-0.16 $\pm$ 0.06
Pop. Density*	<b>0.42 <math>\pm</math> 0.08</b>	<b>0.45 <math>\pm</math> 0.04</b>	-0.29 $\pm$ 0.11	-1.02 $\pm$ 0.32	-0.37 $\pm$ 0.04	-0.57 $\pm$ 0.01	0.05 $\pm$ 0.12	0.38 $\pm$ 0.04	0.03 $\pm$ 0.07
Countries <sup>†</sup> % Acc. $\uparrow$	<b>36.90 <math>\pm</math> 4.32</b>	19.17 $\pm$ 2.82	1.22 $\pm$ 0.05	1.28 $\pm$ 0.01	1.12 $\pm$ 0.00	0.92 $\pm$ 0.02	1.56 $\pm$ 0.47	23.12 $\pm$ 2.50	1.24 $\pm$ 0.12
iNaturalist*	19.60 $\pm$ 0.78	<b>20.91 <math>\pm</math> 0.77</b>	19.85 $\pm$ 0.55	<b>21.49 <math>\pm</math> 0.85</b>	17.52 $\pm$ 0.38	18.11 $\pm$ 0.34	16.14 $\pm$ 0.42	<b>20.94 <math>\pm</math> 0.38</b>	<b>21.08 <math>\pm</math> 0.69</b>
Biome*	25.89 $\pm$ 2.79	16.44 $\pm$ 1.21	1.98 $\pm$ 0.62	3.00 $\pm$ 2.60	1.76 $\pm$ 0.04	2.79 $\pm$ 0.19	<b>37.81 <math>\pm</math> 4.47</b>	31.67 $\pm$ 1.91	6.24 $\pm$ 2.71
Ecoregions <sup>†</sup>	<b>21.02 <math>\pm</math> 1.09</b>	10.86 $\pm$ 1.19	1.55 $\pm$ 0.17	1.41 $\pm$ 0.14	1.49 $\pm$ 0.03	1.48 $\pm$ 0.00	1.36 $\pm$ 0.10	6.65 $\pm$ 1.03	1.52 $\pm$ 0.47
<b>Africa</b>									
Air Temp.* $R^2 \uparrow$	-4.71 $\pm$ 2.29	<b>-1.48 <math>\pm</math> 0.70</b>	-3.64 $\pm$ 3.18	-2.67 $\pm$ 5.80	-7.91 $\pm$ 0.04	-7.29 $\pm$ 0.27	-17.43 $\pm$ 18.37	-9.91 $\pm$ 28.82	-27.36 $\pm$ 39.46
Elevation*	-1.80 $\pm$ 1.74	<b>-0.21 <math>\pm</math> 0.09</b>	-0.39 $\pm$ 0.09	-1.20 $\pm$ 0.55	-0.13 $\pm$ 0.06	-0.15 $\pm$ 0.02	-0.79 $\pm$ 0.43	-0.34 $\pm$ 0.10	-2.43 $\pm$ 2.67
Pop. Density*	0.17 $\pm$ 0.12	0.18 $\pm$ 0.09	-0.53 $\pm$ 0.21	-0.31 $\pm$ 0.16	-0.34 $\pm$ 0.02	-0.37 $\pm$ 0.01	0.15 $\pm$ 0.05	<b>0.32 <math>\pm</math> 0.03</b>	-0.50 $\pm$ 0.34
Countries <sup>†</sup> % Acc. $\uparrow$	<b>30.65 <math>\pm</math> 4.23</b>	10.22 $\pm$ 1.62	0.47 $\pm$ 0.01	0.45 $\pm$ 0.04	0.47 $\pm$ 0.01	0.45 $\pm$ 0.00	0.48 $\pm$ 0.00	10.32 $\pm$ 2.75	2.74 $\pm$ 2.52
iNaturalist*	<b>9.53 <math>\pm</math> 0.57</b>	6.23 $\pm$ 0.47	6.63 $\pm$ 0.57	8.65 $\pm$ 0.52	7.47 $\pm$ 0.53	6.85 $\pm$ 0.39	5.18 $\pm$ 0.38	7.69 $\pm$ 0.30	<b>9.96 <math>\pm</math> 0.33</b>
Biome*	35.72 $\pm$ 5.48	12.34 $\pm$ 1.75	0.94 $\pm$ 0.00	1.09 $\pm$ 0.48	1.29 $\pm$ 0.04	1.17 $\pm$ 0.21	<b>49.86 <math>\pm</math> 1.57</b>	28.28 $\pm$ 3.06	1.46 $\pm$ 0.67
Ecoregions <sup>†</sup>	<b>32.03 <math>\pm</math> 1.19</b>	12.91 $\pm$ 1.63	0.90 $\pm$ 0.00	0.94 $\pm$ 0.04	0.88 $\pm$ 0.01	0.90 $\pm$ 0.00	0.92 $\pm$ 0.12	12.41 $\pm$ 2.20	7.72 $\pm$ 3.93
# of wins	8	5	0	1	0	0	2	4	2

catenating the embeddings before feeding them to the downstream learner. We show in our results in Tab. 9 that combinations do not help and performance is not improved over the respective best single embedding.

**Table 9: Downstream task performance using combinations of different embeddings.** We report average test set MSE ( $\downarrow$ )  $\pm 1$  standard deviation across 10 independently initialized MLP training runs on all regression tasks. Here, SatCLIP corresponds to a model with ResNet-50 vision encoder.

Embedding 1	Embedding 2	Task					
		Air Temp.	Med. Income	Cali. Housing	Elevation	Population	
SatCLIP <sub>L</sub> = 10	GeoCLIP	0.23 $\pm$ (0.01)	0.67 $\pm$ (0.01)	1.81 $\pm$ (0.11)	0.14 $\pm$ (0.01)	0.48 $\pm$ (0.01)	
SatCLIP <sub>L</sub> = 10	CSP (iNat)	0.23 $\pm$ (0.02)	1.08 $\pm$ (0.03)	4.04 $\pm$ (0.06)	0.26 $\pm$ (0.01)	0.69 $\pm$ (0.02)	
SatCLIP <sub>L</sub> = 10	GPS2Vec (tag)	0.26 $\pm$ (0.02)	0.95 $\pm$ (0.02)	3.76 $\pm$ (0.06)	0.22 $\pm$ (0.01)	0.6 $\pm$ (0.01)	
SatCLIP <sub>L</sub> = 10	GPS2Vec (visual)	0.31 $\pm$ (0.02)	0.97 $\pm$ (0.02)	3.82 $\pm$ (0.1)	0.24 $\pm$ (0.02)	0.61 $\pm$ (0.01)	
CSP (iNat)	GPS2Vec (tag)	1.56 $\pm$ (0.08)	1.29 $\pm$ (0.02)	5.66 $\pm$ (0.09)	0.88 $\pm$ (0.04)	2.06 $\pm$ (0.16)	
CSP (iNat)	GPS2Vec (visual)	1.85 $\pm$ (0.22)	1.3 $\pm$ (0.01)	5.72 $\pm$ (0.07)	0.91 $\pm$ (0.02)	1.78 $\pm$ (0.04)	
GPS2Vec (visual)	GPS2Vec (tag)	2.99 $\pm$ (0.18)	1.19 $\pm$ (0.02)	2.53 $\pm$ (0.13)	0.9 $\pm$ (0.0)	2.03 $\pm$ (0.01)	