TEXTUAL SUPERVISION ENHANCES GEOSPATIAL REPRESENTATIONS IN VISION-LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Geospatial understanding is a critical yet underexplored dimension in the development of machine learning systems for tasks such as image geolocation and spatial reasoning. In this work, we analyze the geospatial representations acquired by three model families: vision-only architectures (e.g., ViT), vision-language models (e.g., CLIP), and large-scale multimodal foundation models (e.g., LLaVA, Qwen, and Gemma). By evaluating across image clusters, including people, landmarks, and everyday objects, grouped based on the degree of localizability, we reveal systematic gaps in spatial accuracy and show that textual supervision enhances fine-grained geospatial representations. Our findings suggest the role of language as an effective complementary modality for encoding spatial context and multimodal learning as a key direction for advancing geospatial AI.

1 Introduction

Vision models have undergone tremendous progress in the last decade, driven by advances in convolutional neural network (CNN) architectures (Simonyan & Zisserman, 2014; He et al., 2016) and Vision Transformers (ViT) (Dosovitskiy et al., 2021). These models are capable of capturing highlevel, transferable representations that can be utilized in zero-shot scenarios via their embeddings and be adapted through fine-tuning for specific downstream tasks. Specifically, ViTs benefit from the scalability of Transformers (Vaswani et al., 2017) and enable the development of foundation models across multiple data modalities and application domains.

Recent advances such as CLIP (Radford et al., 2021) include multimodal models that integrate text and vision to learn joint representations within a shared latent space. Another line of research focuses on vision-language models (VLMs), which integrate text and image inputs through a two-stage training pipeline: an initial phase using paired text-image data, followed by instruction tuning (Liu et al., 2024; Bai et al., 2025; Kamath et al., 2025). These models typically employ a frozen vision encoder alongside a language model, enabling multimodal understanding and generation.

Vision models increasingly demonstrate the ability to internalize diverse meta information around the world, raising the question of whether they also encode *geolocation*—even without explicit geospatial supervision. Their internal representations are shaped by architecture, pretraining, and fine-tuning, yet remain difficult to interpret (Ghiasi et al., 2022). This challenge is further amplified in emerging VLMs, where multimodal complexity obscures the mechanisms by which knowledge is encoded. Such opacity can lead to unintended outcomes, including geographic disparities (Moayeri et al., 2024). To improve fairness and transparency, we examine the capacity of vision-only and vision-language models to encode implicit geospatial information (Figure 1) by asking the following question: To what extent do these models internalize global location knowledge as an emergent property of their training and fine-tuning pipelines?

Learning geospatial representations via supervised training has already been explored by Vivanco Cepeda et al. (2023). Following on their work, we are interested in investigating what kinds of geospatial features are learned during training without additional supervision. For text-based large language models (LLMs), Gurnee & Tegmark (2024) and Godey et al. (2024) have shown that specific neurons and layers within LLMs implicitly encode latitude and longitude information and that this capacity scales with increasing model size.

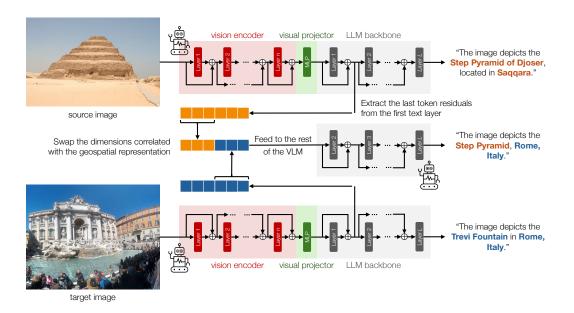


Figure 1: Schematic illustration showing that editing the geospatial representations (through dimension swapping) changes the perceived geolocation during token generation of VLMs. We demonstrated this finding on Qwen2.5-VL-3B with the methodology described in Section 4.5.

Embeddings from pretrained LLMs can also be used to create geospatial embeddings from geolocation-related prompts, as shown in LLMGeovec (He et al., 2025), where its representations improve performance on various downstream tasks requiring geospatial understanding. Additionally, Roberts et al. (2024) has shown that VLMs have spatial reasoning capabilities, being able to complete a variety of tasks through zero-shot settings. We extend these results by focusing on how ViT models separate images spatially in their learned latent space and exploring how these models learn geospatial representations.

Our main contributions are as follows.

- We investigate the emergence of geospatial representations learned by vision-only architectures, vision-language models, and large-scale multimodal foundation models, finding that the latter two groups exhibit substantially stronger geospatial structure.
- We evaluate the performance of different model representations using layer-wise probing
 for geospatial location prediction. We find that vision-only models tend to exhibit stronger
 representations on their last layer, while VLMs have better geospatial representations on
 the early layers of their language model block.
- We show that prompting VLMs allows the geospatial information to be propagated to the latter layers, in some cases leading to an increase in representation quality.

2 RELATED WORKS

2.1 VISION-BASED MODELS

ViTs emerged as a paradigm shift in computer vision, adapting Transformers (Vaswani et al., 2017) to image data by segmenting images into tokenized patches with positional embeddings (Dosovitskiy et al., 2021). The initial ViT model demonstrated that large-scale supervised pretraining on image classification tasks could yield transferable representations across domains. Researchers have also explored self-supervised approaches to vision. One such method is the masked autoencoder (MAE) (He et al., 2022), which trains models to reconstruct randomly masked image patches. This objective encourages the extraction of semantically rich features that can be effectively adapted to downstream tasks via fine-tuning.

As an alternative to self-supervised pretraining, Caron et al. (2021) proposed a self-distillation framework named DINO. This approach enables the model to learn invariant representations across multiple augmented views of the same image, resulting in linearly separable features that implicitly capture semantic structures such as object boundaries and regions. Building upon this foundation, DINOv2 (Oquab et al., 2023) extends this methodology by incorporating the IBOT loss (Zhou et al., 2021), a patch-level objective. This integration facilitates scalable pretraining, enhancing the model's capacity to learn fine-grained visual representations from large-scale unlabeled datasets.

2.2 VISION-LANGUAGE MODELS

A turning point for vision models has been the integration of language for learning shared representations. Pioneering work like CLIP jointly trained a vision encoder and a language encoder to align their representations using a large corpus of web-scraped image-caption pairs (Radford et al., 2021). Subsequent research, such as SIGLIP (Zhai et al., 2023), expanded CLIP by replacing softmax loss with a sigmoid objective, decoupling performance from batch size and allowing for improved scalability. These models are fundamental and are applied as the core vision-encoder for many of the vision-language foundation models presented in our work.

Foundation VLMs utilize specialized training pipelines. For example, LLaVA-1.5 (Liu et al., 2024) employs a two-stage training process, first aligning a frozen CLIP visual encoder with an LLM on image-text pairs, and then fine-tuning the model on a GPT-generated instruction-following dataset to enhance its conversational and reasoning abilities. Qwen2.5 (Bai et al., 2025) follows a similar process, pretraining on image-text pairs and then performing additional supervised fine-tuning (SFT) and direct preference optimization (DPO) to structure instruction-following data. Gemma 3 (Kamath et al., 2025) leverages a frozen SIGLIP vision encoder, a pretraining stage similar to previous models, and a post-training stage that includes knowledge distillation from a larger instruction-tuned model and alignment with human feedback via SFT and reinforcement learning with human feedback (RLHF).

3 METHODS

3.1 Models

To examine how geolocation capabilities emerge in vision models without explicit supervision, we curated a diverse set of architectures spanning multiple modalities and training paradigms. Our selection includes both (i) vision-only encoders, i.e., ViT (Dosovitskiy et al., 2021), ViT Masked Autoencoder (He et al., 2022), and DINOv2 (Oquab et al., 2023) and (ii) vision-language models, i.e., CLIP (Radford et al., 2021), LLaVA-1.5 (Liu et al., 2024), Qwen2.5 (Bai et al., 2025), and Gemma 3 (Kamath et al., 2025), representing supervised and self-supervised approaches. For each model family, we evaluated at least two size variants to assess the influence of scale on learned geospatial representations. Additional information about these models is given in Appendix A.

3.2 Dataset

To build our dataset, we sampled images from established benchmarks, including YFCC100M and Google Landmarks. We provide the details below.

Yahoo Flickr Creative Commons 100 Million (YFCC100M) (Thomee et al., 2016). We used a 4M-image subset from the MediaEval 2016 Placing Task competition (Choi et al., 2016), obtained via Kaggle¹. This dataset is a diverse collection of Flickr-sourced images spanning natural scenes, urban environments, and everyday objects with location data. To analyze localizability across semantic categories, we partitioned this subset via unsupervised clustering. First, we extracted image embeddings with ResNet-152 (He et al., 2016) pretrained on ImageNet (Deng et al., 2009), then we applied principal component analysis (PCA) to retain the top-100 components explaining the highest variance and performed k-means clustering (Han et al., 2022). Among 19 tested k values ($k = 10, 15, \ldots, 100$), we selected k = 40 using the elbow method (Thorndike, 1953), with manual inspection confirming semantic coherence. The resulting clusters captured meaningful categories,

https://www.kaggle.com/datasets/habedi/large-dataset-of-geotagged-images

including people, objects, cliffs, landscapes, and buildings. Complete clustering details are provided in Appendix B.

Google Landmarks (Weyand et al., 2020). This dataset contains over 5M photographs of globally recognized landmarks and tourist sites (e.g., Eiffel Tower, Mount Fuji) that are available in public data sources. We hypothesize that such iconic scenes may have been encountered during model pretraining, contributing to their high localizability. It also contains non-localizable images, such as particularly close-up shots of people or animals, and generic textures such as soil, which lack distinctive geospatial cues. In our experiments, we use a subset of 580k images², hereafter referred to as the *Landmarks* dataset, with geolocation coordinates extracted from OpenStreetMap.

3.2.1 Sampling

Across all datasets, the geospatial distribution of images was imbalanced, skewed toward major cities in Europe and North America. To mitigate this bias, we partitioned the globe into non-overlapping geocells based on global administrative area boundaries and iteratively merged regions with insufficient samples. We merged geocells hierarchically. First within regions (GID_1), then across regions of the same country (GID_0). To avoid oversampling, each geocell was defined to include at least one complete GID_2 (city-level) administrative unit, even in dense areas like Paris. These geocells were then used to balance the YFCC100M and Landmarks dataset by selecting 5,000 images per source, with at most five images per geocell. For the Landmarks dataset, we also excluded duplicate images of the same landmark to ensure diversity.

3.3 PROBING

To examine whether the evaluated models encode geospatial information, we perform linear probing (Alain & Bengio, 2017), a standard mechanistic interpretability technique for Transformer architectures (Gurnee & Tegmark, 2024; Kim et al., 2025). Transformers (Vaswani et al., 2017) consist of sequential blocks that iteratively refine token representations within the residual stream $x^{(l)} \in \mathbb{R}^{t \times d_{\text{model}}}$, where d_{model} is the hidden dimension of each evaluated model, as an input with t tokens is propagated through the l-th Transformer block (Elhage et al., 2021). This refinement is achieved through successive multi-head attention (MHA) and multi-layer perceptron (MLP) layers with residual connections. These layers are often paired with normalization, of which the formulation is omitted for brevity:

$$h_{\text{attn}}^{(l)} = x^{(l)} + \text{MHA}\left(x^{(l)}\right) \tag{1}$$

$$h_{\rm mlp}^{(l)} = \text{MLP}\left(h_{\rm attn}^{(l)}\right) \tag{2}$$

$$x^{(l+1)} = h_{\text{attn}}^{(l)} + h_{\text{mlp}}^{(l)} \tag{3}$$

Although t varies across vision models, downstream tasks typically rely on a single summary representation. Usually, the [CLS] token is used in vision models and the final token representation in VLMs. We fit ridge regression models to predict latitude and longitude (in degrees) from layer-wise token summary representations. We report the models' predictive performance using the coefficient of determination (R^2) . Formally, given hidden representations $\mathbf{A}^{(l)} \in \mathbb{R}^{n \times d_{\text{model}}}$ from layer l, number of samples n, and hidden dimension d_{model} , the regression weights $\mathbf{W} \in \mathbb{R}^{d_{\text{model}} \times 2}$, corresponding to the two-dimensional targets (latitude and longitude), are estimated as:

$$\hat{\mathbf{W}} = \arg\min_{\mathbf{W}} \|\mathbf{Y} - \mathbf{A}^{(l)}\mathbf{W}\|_F^2 + \lambda \|\mathbf{W}\|_F^2.$$
(4)

Here, $\|\cdot\|_F$ denotes the Frobenius norm, and $\lambda > 0$ is the regularization hyperparameter controlling the strength of the L_2 penalty on \mathbf{W} . In all of our experiments, λ is chosen for each probe using Leave-One-Out cross-validation (Golub et al., 1979).

²https://huggingface.co/datasets/visheratin/google_landmarks_places

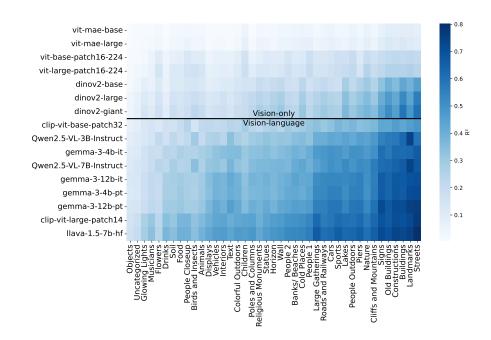


Figure 2: Model performance measured by coefficient of determination \mathbb{R}^2 across all models. The x-axis shows image clusters based on the YFCC100M dataset and the Google Landmarks dataset, and y-axis lists the models compared. Higher \mathbb{R}^2 values (darker colors in the heatmap) indicate better geolocation-prediction accuracy across clusters.

4 EXPERIMENTS

4.1 Performance on Geolocation Information Prediction

We perform layer-wise probes on all the evaluated models and report the geolocation prediction performance in Figure 2. We find that models trained jointly on text and images exhibit measurable geospatial representations: the R^2 values reach up to 0.8 for Landmarks and streets (cluster 28), while the average R^2 is above 0.4 for the larger models, suggesting some degree of geospatial representations across image types. The average R^2 for vision-only models, on the other hand, is mainly below 0.3. Among vision-only models, performance improves with model size, with the DINOv2-giant model (1B parameters) achieving the best result. This observation suggests that, when trained on a scale, geospatial representation can be learned from images alone. When compared with VLMs, DINOv2-giant is outperformed even by the much smaller CLIP-base model, suggesting the effectiveness of language pretraining in implicitly learning geospatial representations.

The cluster-wise performance presented in Figure 2 indicates that the relative difficulty of each cluster in YFCC100M remains consistent across various model architectures. For instance, the *streets* cluster, *building* cluster, and the *Landmarks* dataset consistently show higher R^2 across models, while the *objects* cluster contains little information that can serve as clues for geo-localization. Interestingly, the *signs* and *text* clusters, show a degree of localizability for VLMs not observed for vision-only models.

Figure 3 shows image samples positioned according to their \mathbb{R}^2 values for the largest model of each model family. Highly localizable images tend to be famous landmarks, e.g., pyramids, and open spaces with pieces of architecture or nature. Meanwhile, close-ups of objects and food images are the least localizable. Notice that for CLIP, the cluster of figures that contain signs achieves a notable degree of localizability.

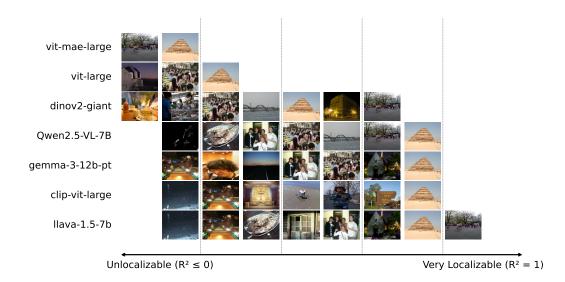


Figure 3: Image samples across the localizability spectrum, separated by R^2 . For each row, the geolocation prediction performance of a different model is shown. For example, in the first row, it is seen that ViT-MAE-large achieves a low R^2 across all images, while, in the last row, LLaVA-1.5-7B-HF achieves various R^2 levels across different clusters and datasets. For high-performing models, highly localizable landmarks usually achieve the highest performance.

4.2 Geospatial Representations Across Layers

To investigate where geospatial representation emerges within model architectures, we analyze probe performance across layers for both vision-only and multimodal models. In vision models such as ViT and DINOv2, geospatial representations tend to develop progressively with increasing layer depth, as evidenced by a consistent rise in probe R^2 values across all settings (see Figure 4(a)). For VLMs, however, an interesting observation emerges: the R^2 values increase only up to a certain point before stagnating, and in the case of Gemma, the R^2 values decrease throughout the later layers, regardless of image characteristics. This likely reflects the model's tendency to deprioritize geographic signals in the absence of a textual prompt, thus neglecting spatial information not essential for text generation. A similar, though less pronounced, effect is observed in LLaVA, where geolocation signals diminish as the image transitions into the language modeling component.

4.3 Steering the Models via Text Prompts

In prior experiments, we observed that the geospatial information in the model residuals, particularly VLMs, degrades over layers without a textual prompt. This leads to, in some extreme cases, for example, Gemma, a negative \mathbb{R}^2 , suggesting that there is no linear mapping from the model residuals to geolocation coordinates.

To check if this effect occurs because of the absence of a textual prompt related to geospatial information, we prompt the models with the query "Guess the latitude and longitude of this image. Answer only with the coordinate tuple (lat, long)". Figure 4(b) shows the per-layer R^2 of VLMs when prompted this way. We observe that, in this setting, R^2 does not decrease as drastically for Gemma and LLaVA. In fact, for LLaVA, it starts increasing over the textual layers. Interestingly, for Qwen the performance drastically increases, leading to R^2 as high as 0.88 for the *Landmarks* dataset. This effect suggests that textual and image representations of geospatial representation might be entangled in the model's latent space, especially when activated using textual prompts related to geospatial tasks.

4.4 ISOLATING GEOSPATIAL-SPECIFIC COMPONENTS FROM EMBEDDINGS

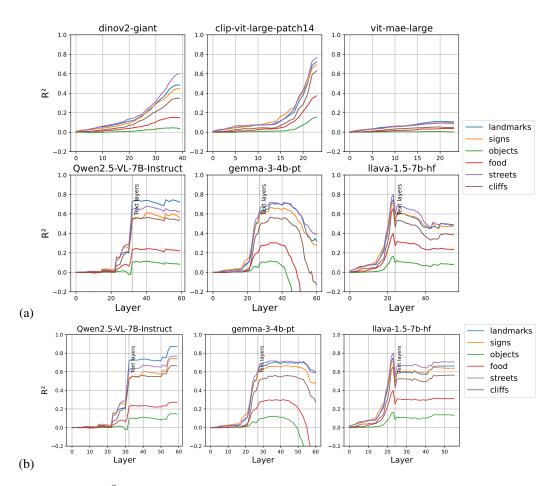


Figure 4: Probe R^2 performance by layer of the models for different clusters and datasets with varying levels of localizability. (a) R^2 performance when no textual prompt is given. (b) R^2 performance when adding a textual prompt to the input asking the model to predict the image geolocation. The R^2 is kept stable throughout the last layers when compared to the decaying performance observed in the non-prompting setup.

The linear probes operate on high-dimensional representations. In our experiments, the five selected models have $d_{\rm model}$ of 768 for CLIP-ViT-large, 1,024 for LLaVA-1.5, 2,048 for Qwen2.5-VL-3B, 3,584 for Qwen2.5-VL-7B, and 1,536 for DINOv2-giant. To explore how latent space contributes to geospatial information, we fit ridge regression probes using only a proportion of the original dimensions $p \in \{0.1, 0.2, \ldots, 1.0\}$, selecting the top p dimensions ranked by the absolute coefficients of the trained probe.

We report the predictive performance in terms of R^2 in Figure 5 as a function of the retained feature proportion p, showing that R^2 increases with p and saturates well before using the entire feature set. Across all models, we find that $p \approx 0.4$ (about 40% of dimensions) are sufficient to recover nearly the maximum R^2 , indicating that geospatial information is concentrated in a compact subset of dimensions rather than uniformly distributed throughout the em-

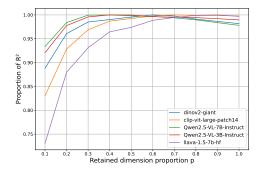


Figure 5: Probe predictive performance \mathbb{R}^2 as a function of the retained feature proportion p, illustrating the capacity of the embedding subspace needed to reach the maximum \mathbb{R}^2 for both vision-only and VLMs. Higher values of p correspond to a larger subset of the latent representation.

beddings. For Qwen2.5-VL variants, specifically, 90% of their best predictive performance is still observed using only the top 10% of the features.

4.5 STEERING THE MODEL GENERATION THROUGH REPRESENTATION SWAPPING

We examine the possibility of steering the text generated by a multimodal foundation model by swapping the top p feature dimensions related to geospatial reasoning. In this case study, the geospatial location in the predicted text should be changed, leaving other semantic information unchanged. We use Qwen2.5-VL-3B, an open-weight VLM that supports a targeted intervention in its *residual stream* or the additive hidden state passed through the transformer layers. As illustrated in Figure 1, given a *source* image and a *target* image, we replace the top geospatially relevant feature dimensions of the *source* residual-stream summary with the corresponding coordinates from the *target*, and evaluate the resulting changes in the generated text.

Let t^* be the last non-padding input token (the summary token for Qwen2.5-VL-3B), and let $x_{\text{source},t^*}^{(1)}, x_{\text{target},t^*}^{(1)} \in \mathbb{R}^{d_{\text{model}}}$ be the layer-1 residual-stream vectors for source and target images, respectively. Given $g \subseteq \{1,\ldots,d_{\text{model}}\}$ as the index set of geospatially informative dimensions (with proportion $p = |g|/d_{\text{model}}$) and its complement $g^c = \{1,\ldots,d_{\text{model}}\} \setminus g$, we intervene by replacing the source residual with the target residual on the dimensions in g as follows:

$$\tilde{x}_{\text{source},t^{\star}}^{(1)} = x_{\text{source},t^{\star}}^{(1)} \odot \mathbb{1}_{g^{C}} + x_{\text{target},t^{\star}}^{(1)} \odot \mathbb{1}_{g}, \tag{5}$$

where \odot denotes the element-wise Hadamard product and $\mathbb{1}_g$ an indicator vector with entries 1 for indices in g and 0 elsewhere. For brevity, we omit explicit indexing of the selected dimensions g. We then continue the forward pass from layer $2, \ldots, L$ using $\tilde{x}_{\text{source}}^{(1)}$ to decode the output text.

We show an example in which we can successfully steer the model in Figure 1. By swapping 50% of geospatial representations from an image of the Step Pyramid of Djoser with the geospatial representation from an image of the Trevi Fountain, the model generates the following output text: "The image depicts the Step Pyramid, Rome, Italy", altering the location of the pyramid from Saqqara to Rome. Even though our experiments show the possibility of successfully steering the model, we observed that as the text generated becomes longer, the generation becomes unstable. In some cases, the model starts to generate repetitive text or descriptions that mix the source and target locations. These results open future avenues for investigating how geospatial representations are coupled with representations related to other types of information during text generation. We discuss more details in Appendix F.

4.6 DOWNSTREAM TASK PERFORMANCE

Finally, we inspect how the quality of geospatial representations may influence downstream task performance, as these models are usually finetuned for specific downstream applications. For this analysis, we investigate a task that requires geospatial awareness: country identification.

Using the landmarks dataset, we extract country information for each picture and then subsample the dataset so that at most 100 pictures are selected for each country. Then, we finetune one large model from each studied vision-only family (ViT-MaE,ViT, and DINOv2) in addition to CLIP-ViT-large and DINOv2-giant (for the full details, see Appendix G). The models are chosen such that all take the same inputs and have similar model size, making the results comparable. We report

Table 1: Finetune performance for the country identification task.

Model	Test Acc.	Val. Loss	Train Loss
ViT-MaE-large	0.15	3.35	2.344
ViT-large	0.23	3.17	1.346
DINOv2-large	0.29	2.55	0.009
DINOv2-giant	0.32	2.78	0.001
CLIP-large	0.36	2.39	0.009

the results of each model in Table 1. We observe that the performance of the models follows the order of \mathbb{R}^2 obtained for our probe analysis, with ViT-MaE having the worst performance, while CLIP has the best performance. This corroborates the hypothesis that the presence of geospatial representations in the models is desirable for their use in downstream tasks.

5 DISCUSSION

Our findings demonstrate that the training methodology is crucial for learning geospatial representations in vision-only models and VLMs, as demonstrated in Figure 2. Models that incorporate textual supervision consistently achieve the best performance. Vision-only models improve with scale; larger models like DINOv2-giant outperform both DINOv2-large and DINOv2-base. In contrast, VLMs do not exhibit a clear correlation between model scale and probing performance. This suggests that supervision signals, particularly language, are a primary factor in learning strong geospatial representations in these models.

The optimal layer for extracting geospatial representations depends on the model family and the presence of a textual prompt. In many geolocation applications, the input image is given without a textual prompt (Figure 4(a)). For these cases, vision-only models perform best when using representations from their deepest layers. However, for VLMs, choosing which geospatial representation to use is not clear and varies between models. Across models, the layers immediately after the textual stream consistently achieve good performance. With a textual prompt, the performance in VLMs remains more stable across the post-textual layers (Figure 4(b)).

The results of our analysis have direct implications for model selection and methodology in a range of geospatial applications. For current pipelines that use traditional vision-only models and small datasets, leveraging representations from VLMs can significantly improve performance. As these representations are implicitly learned from vast datasets, they serve as strong representations for sample-efficient pipelines involving fine-tuning cases where labeled data are scarce. Moreover, our work highlights multimodal learning as a critical direction to build world models, which could be used to improve our understanding of complex social problems and empower new technologies.

Finally, the growing capability of these models poses significant privacy risks and fairness implications. Malicious actors could exploit these models to extract precise location data from images, enabling stalking and threats against individuals. The potential for mass surveillance is also a serious concern, where governments and corporations could likewise track individuals' locations and behaviors through their photos. These ethical risks underscore the need for robust regulatory policies that mandate transparency in model use and enforce explicit user consent to ensure safe deployment.

6 CONCLUSION

This study demonstrated that textual supervision significantly enhances geospatial representations in vision-language models. Through a systematic analysis of vision-only architectures, VLMs, and large-scale multimodal foundation models, we studied how geospatial understanding emerges across model families. Through layer-wise probing, we revealed that multimodal models consistently exhibit high performance for images that are localizable. Furthermore, our analysis indicated that a small subset of hidden dimensions is responsible for encoding critical geospatial features, suggesting a potential pathway for model steering and editing. In summary, our work demonstrated that multimodal learning plays an important role in improving geospatial AI. However, using these models in real-world settings should include safeguards to protect privacy and ensure fairness. As applications involving location-aware image understanding—such as environmental monitoring, urban planning, and disaster response—continue to grow, this use case is expected to become increasingly important. Future research could explore how these models handle other types of images, such as satellite data.

REPRODUCIBILITY STATEMENT

The code for reproducing our results is available through an anonymous repository for validation³, and all datasets used in this study are publicly accessible.

REFERENCES

Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *ICLR Workshop Track*, 2017.

³https://anonymous.4open.science/r/ICLR-9053

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. arXiv preprint arXiv:2502.13923,
 2025.
 - Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
 - Jaeyoung Choi, Claudia Hauff, Olivier Van Laere, and Bart Thomee. The placing task at mediaeval 2016. In *MediaEval Benchmarking Initiative for Multimedia Evaluation*, 2016. URL https://api.semanticscholar.org/CorpusID:6435875.
 - Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
 - Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
 - Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.
 - Amin Ghiasi, Hamid Kazemi, Eitan Borgnia, Steven Reich, Manli Shu, Micah Goldblum, Andrew Gordon Wilson, and Tom Goldstein. What do vision transformers learn? a visual exploration. *arXiv preprint arXiv:2212.06727*, 2022.
 - Nathan Godey, Éric de la Clergerie, and Benoît Sagot. On the scaling laws of geographical representation in language models. *arXiv preprint arXiv:2402.19406*, 2024.
 - Gene H Golub, Michael Heath, and Grace Wahba. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
 - Wes Gurnee and Max Tegmark. Language models represent space and time. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Jiawei Han, Jian Pei, and Hanghang Tong. *Data mining: concepts and techniques*. Morgan kaufmann, 2022.
 - Junlin He, Tong Nie, and Wei Ma. Geolocation representation from large language models are generic enhancers for spatio-temporal learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 17094–17104, 2025.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
 - Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
 - Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv* preprint arXiv:2503.19786, 2025.
 - Junsol Kim, James Evans, and Aaron Schein. Linear representations of political perspective emerge in large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
 - Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 26296–26306, 2024.

- Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Mazda Moayeri, Elham Tabassi, and Soheil Feizi. Worldbench: Quantifying geographic disparities in llm factual recall. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1211–1228, 2024.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
- Jonathan Roberts, Timo Lüddecke, Rehan Sheikh, Kai Han, and Samuel Albanie. Charting new territories: Exploring the geographic and geospatial capabilities of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 554–563, 2024.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.
- Robert L Thorndike. Who belongs in the family? *Psychometrika*, 18(4):267–276, 1953.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired alignment between locations and images for effective worldwide geo-localization. *Advances in Neural Information Processing Systems*, 36:8690–8701, 2023.
- Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2575–2584, 2020.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 11975–11986, 2023.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. Image bert pre-training with online tokenizer. In *International Conference on Learning Representations*, 2021.

A MODEL DETAILS

Details about the models evaluated in this work are given in Table 2.

To extract the inner representations of the model, for the probing experiments we use a single image as input and no prompt in the VLMs; then we conduct a single forward pass over the entire model, extracting all residuals $x^{(l)}$ for the last token for both the vision and the text components when applicable. A similar setup is adopted for the prompting experiments, but with added prompt tokens. Since we are not interested in generating text, we use no specific parameters (e.g., temperature) for the VLMs for the probing experiments. For the generation experiments, we use a very low temperature (0.0001) to reduce randomness.

All experiments were run on a single Nvidia A100 GPU using the HuggingFace implementation of each model. The models which only take image as input were run at full precision, while the VLMs (Gemma, Qwen, and LLaVA) were run in bfloat16 precision.

Table 2: Models evaluated in this work—ViT (Dosovitskiy et al., 2021), ViT Masked Autoencoder (He et al., 2022), DINOv2 (Oquab et al., 2023), CLIP (Radford et al., 2021), LLaVA-1.5 (Liu et al., 2024), Qwen2.5 (Bai et al., 2025), and Gemma 3 (Kamath et al., 2025)—spanning vision-only and vision-language modalities with different training paradigms. Each family is evaluated with at least two size variants to examine the effect of model scale on learned representations.

Model	Modality	Training Methodology
ViT	Vision-Only	Supervised pretraining on ImageNet-21k followed by fine-tuning on ImageNet-1k.
ViT Masked Autoencoder	Vision-Only	Self-supervised training using masked autoencoding.
DINOv2	Vision-Only	Self-supervised learning using a teacher-student framework.
CLIP	Vision-Language	Trained on 400M image-text pairs using contrastive learning.
LLaVA-1.5	Vision-Language	Combines a CLIP-based vision encoder with a language model via vision-language alignment followed by instruction tuning.
Qwen2.5	Vision-Language	CLIP-based pretraining enhanced with vision- language alignment and end-to-end instruction tuning.
Gemma 3	Vision-Language	Uses SIGLIP-based vision-text pretraining followed by alignment with instruction-tuned language models.

B CLUSTERING DETAILS

Given the massive size of the YFCC100M dataset, its content was divided into semantic clusters for a better and more comprehensive analysis. For this, we started by extracting embeddings from the final convolutional layer of a ResNet-152 pretrained on ImageNet. Then, the resulting embeddings were L_2 -normalized, reduced to 100 dimensions using PCA, and clustered with a standard k-means algorithm. In order to choose the optimal k, we tested 19 values ($k = 10, 15, \ldots, 100$), and following (Thorndike, 1953), computed the within-cluster sum of squares (WCSS) as follows:

WCSS(k) =
$$\sum_{i=1}^{N} D_i = \sum_{i=1}^{N} \min_{c \in \{1,...,k\}} ||\mathbf{x}_i - \boldsymbol{\mu}_c||^2$$
,

which corresponds to the sum of squared distances for each point to its nearest cluster center. The value of WCSS(k) decreases monotonically as k increases, meaning that the optimal k is not the one which minimizes WCSS(k), but rather the one at which the decrease plateaus. This point corresponds to the "elbow" of the curve, which in our case was k=40.

A complete overview of the resulting clusters can be seen in Figure 6. In total, eight clusters represented people, including separate clusters for sports, musicians, and children, as well as another cluster for large gatherings. Another 11 clusters were related to man-made structures and architecture, including different types of buildings, walls, streets, and decorations. Animals were represented in two groups: one for birds and insects, and the other for larger animals. Eight clusters depicted different aspects of nature and different landscapes, such as mountains, lakes, beaches, soil, etc. Finally, seven clusters showed different types of objects, mainly food, drinks, displays, and vehicles. The remaining three clusters did not appear to exhibit any clear semantic relationship.

For better visualization of the clusters obtained, we projected a small subset of 500 points per cluster into a 2D-space using Uniform Manifold Approximation and Projection (UMAP) (McInnes et al., 2018). The generated plot (Figure 7) uses different shades of the same base color to represent clusters belonging to the same macro-category (people, nature, architecture, objects, animals, and others). Clusters associated with nature, architecture, and people are tightly grouped and occupy a distinct region of the latent space, further indicating that our clustering is semantically coherent. The major-



Figure 6: Random samples for each of the 40 clusters obtained for the YFCC100M dataset.



People
Nature
Nature
Architecture
Animals
Objects
Others

Constructions
Streets
People Guitators
Constructions
Signs
Large Calherings
Musicians

Buildings
Old Buildings
Old Buildings
Plies and Columns
Religious-blonuments
People 2
People 2
People 3
People 4
People 3
People 5
People 5
People 5
People 5
People 6
People 6
People 7
People 8
People 1
People 8
People 1
People 8
People 1
People 7
People 8
People 1
People 7
People 7
People 8
People 1
People 6
People 8
People 1
People 7
People 7
People 7
People 7
People 8
People 1
People 7
People 8
People 1
People 7
People 8
People 1
People Closeup
People 8
People 1
People 8
People 9

Figure 7: 2-Dimensional visualization of the 40 clusters obtained using UMAP and colored according to macro-category.

ity of object clusters also lie in a clearly-defined area, though there is some variability depending on the image background. Notably, cars and signs—traditionally found in urban settings—appear close to architecture clusters, while general vehicles, e.g., trains and airplanes, appear in between nature and architecture. The two animal clusters are in a subregion of nature, which is consistent with their broader photographic context, as these images are typically taken in green areas. One exception is the "People Outdoors" clusters, which, though grouped together with other people clusters, extend towards both architecture and nature regions, depending on the broader context of the images.

C DATASET DESCRIPTION

C.1 YFCC100M

The YFCC100M dataset is comprised of approximately 99.2M images published on Flickr between 2004 and 2014 under a Creative Commons license (both commercial and non-commercial). It is a highly diverse set of photographs that depict natural and urban environments, people, objects, and everyday events, taken by a mix of professional photographers and casual users. In our experiments, we used a subset of 4,233,900 images, which preserves the diversity of the original set while offering reliable geolocation annotations.

The spatial distribution of the data set is highly unbalanced (Figure 8), with the majority of the samples concentrated in a few key regions. Together, the G7 countries account for more than 57% of all samples, a figure that rises to 78% with the addition of the rest of the European countries. Asia is the third most represented continent with close to 500k images (12.1% of the data), more than half of which are from East Asia. In contrast, Central Asia and the Middle East are particularly underrepresented, with no country in either region contributing more than 15k samples.

South America accounts for just over 180k images (4.2% of the data), with a sample distribution that closely matches that of the continent's population. The main outliers are Chile, overrepresented in 17% of the data versus 4.5% of the population, and Venezuela, underrepresented at less than 3% of the samples despite being 6.5% of the total population. Africa contributes with approximately 68k images (1.6% of the dataset), with only 12 countries represented by more than 1k samples. Finally, Oceania provides 136k (3.2% of the data) samples, almost entirely from Australia and New Zealand, which together account for 131k. We note that, although there is some variation, all clusters exhibit roughly the same imbalance. After sampling, the balance is marginally improved. Of the 200k total sampled images (5k per cluster), 37k (18.5%) are from Asia, almost 14k (7.0%) are from South America, 8.9k (4.5%) are from Oceania, and 7.4k (3.7%) are from Africa, while the participation of G7 countries decreases from 57% to 44%. Figure 9 shows the effects of these improvements, with our sampling methodology leading to better world coverage than a purely random strategy.

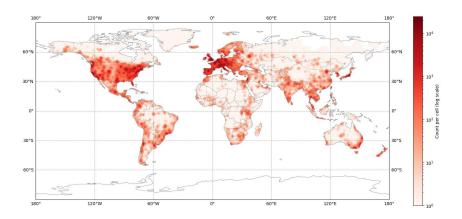


Figure 8: Spatial distribution of all samples from our subset of YFCC100M.

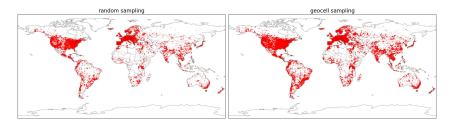


Figure 9: World Coverage of random sampling (left) and geocells-based sampling (right).

C.2 GOOGLE LANDMARKS

Google Landmarks V2 contains over 5M images of over 200k landmarks across the globe, collected mostly from Wikimedia Commons⁴. As the coordinates of images are not typically available on Google Landmarks, we used a subset of 581,215 geolocalized images from 15,453 different labeled landmarks. Though photos of relevant sites are expected to be localizable, the dataset also contains some non-localizable images, such as close-up shots of animals in a national park, as well as paintings and statues.

This data set has a pronounced spatial imbalance, as shown in Figure 10. Samples are mainly concentrated in Europe, which alone accounts for over 67% of our subset, followed by Asia (15%), North America (11%), South America (2.5%), Africa (1.1%), and Oceania (1.0%). Outside of Europe, samples are disproportionately concentrated in a few major cities, with less populated areas remaining mostly uncovered. For our experiments, we selected a single random image from each landmark and proceeded to sample 5k images as outlined in Section 3.2.1. The sampling procedure preserved the same continent-level spatial bias present in the original dataset, though the concentration of samples in major cities was reduced.

D LINEARITY OF GEOSPATIAL FEATURES

In our experiments, we focus on the linear probing for geospatial representation in vision-only and vision-language models, finding that the latter group has internal representations that can be mapped to real world locations, while the former group is much more limited in this regard. However, it could be the case that this happens only because vision-only models are representing this kind of information non-linearly. To explore this possibility, we also train non-linear probes (one hidden layer MLP regression) to check whether vision-only models may be representing geospatial features differently.

⁴https://commons.wikimedia.org/

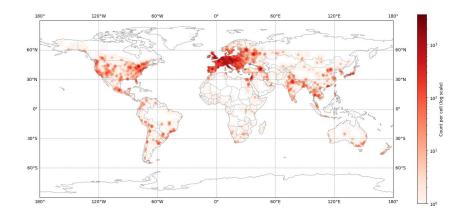


Figure 10: Spatial distribution of all samples from our subset of Google Landmarks V2.

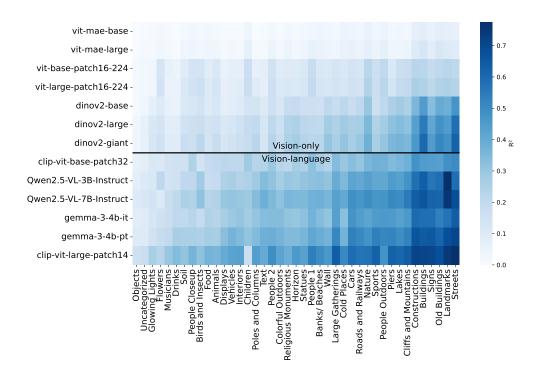


Figure 11: Performance (R^2) of each model when using a non-linear probe. The x-axis shows different clusters of the YFCC100M dataset and the Landmarks dataset, while y-axis shows the models evaluated.

Figure 11 shows the performance (R^2) for the non-linear probes for each model. Using non-linear probes did not result in a significant increase in R^2 for any model, thus strengthening the claim that vision-language models' representations encode geospatial information better.

E ABLATION STUDIES

Our probing setup utilizes the summary representation of the input image, that being either the [CLS] token or the final token representation. However, it could be the case that a geospatial representation emerges across different tokens corresponding to different image patches. To control for this, we also train the same probes using the concatenation of the min and max pooling across all tokens as the inputs for the ridge regression. Figure 12 shows the results across our datasets. It is possible to

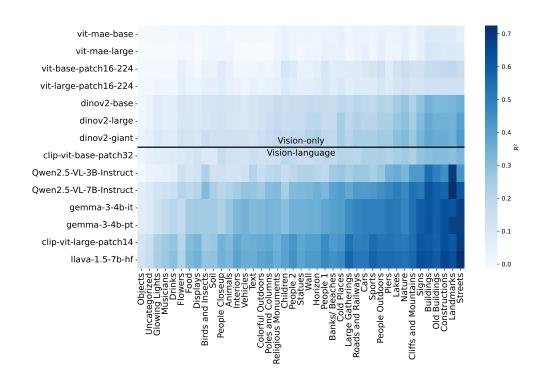


Figure 12: Performance (R^2) of each model when using a linear probe and the concatenation of max and min pooling across input tokens as features. The x-axis shows different clusters of the YFCC100M dataset and the Landmarks dataset, while y-axis shows the models evaluated.

see that, when compared to the default probing approach, the \mathbb{R}^2 is actually smaller for the visiononly models, suggesting that the summary token is adequate as a set of features for the probing setup.

F ADDITIONAL REPRESENTATION SWAPPING RESULTS

In this section, we show additional examples of location swap results obtained from the methodology described in Section 4.5. Here, we show that the intervention used the swapping methodology often leads to other changes in the text, e.g., different place names, mixing characteristics from both images, especially if the text generated has very different lengths for the source and target images.

Table 3 shows examples of generated text after swapping fraction p dimensions from the source image with the target image. We observe that some figures have strong features, while others have easily edited geospatial information. For example, with the same fraction of replaced dimensions, we can move the Cologne Cathedral to Paris, but we cannot move the Eiffel Tower to Cologne. Additionally, the St. Peter Cathedral in Vatican overwrites all other information when used as target image.

G FINETUNING DETAILS

To mitigate the spatial imbalance of our data, which could negatively affect country identification, we constructed a balanced sampling of Google Landmarks across countries. First, we selected a single image for each of the 15,453 labeled landmarks. From these, we retained only those belonging to the 51 countries with over 30 samples, and sampled up to 100 images per country, resulting in a dataset of 3,992 images.

Using this data, we finetune CLIP-large, ViT-large, ViT-MaE-large and DINOv2 (large and giant) for classification. All models were trained for 5 epochs with a 70% train, 20% validation and 10%

Table 3: Results for embedding swapping using varied landmark pictures as source and target images. The changes are very drastic for some image combinations, despite similar methodology when implementing the interventions.

Source Image	Target Image	p	Generated Text
		0.4	The image depicts the Cologne Cathedral, also known as the Cologne Cathedral, located in Cologne, France
		0.4	The image depicts the Cologne Cathedral, also known as the Cathedral of Notre-Dame de Paris, located in the of Paris, France.
		0.4	The image depicts the. Peter's Basilica, Vatican City.
		0.5	The image depicts the Taji Palace, Rome, Italy.
		0.5	The image depicts the Hagou Basil Mosque, Istanbul, Italy.

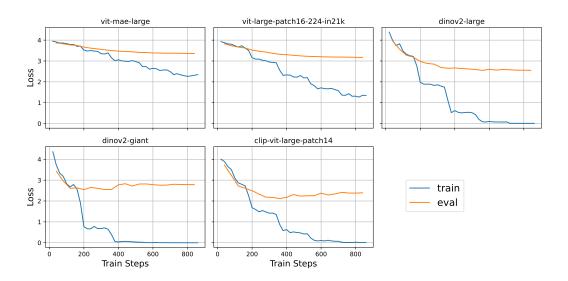


Figure 13: Train and validation loss for model finetuning.

test split using an AdamW optimizer, batch size of 16, and the following learning rates $[1 \times 10^5, 2 \times 10^5, 5 \times 10^5, 1 \times 10^4, 2 \times 10^4, 5 \times 10^4]$, with the best chosen based on validation loss. β_1 and β_2 were set as 0.9 and 0.999. Figure 13 shows the train and validation loss for the finetuned models. We can see that CLIP and DINOv2 models achieve their best validation loss within 400 training steps with CLIP's being overall lower, which also translates to better test accuracy.

H CORRELATION OF INFLUENTIAL WEIGHTS BETWEEN MODELS

Through our probing experiments we obtain a large set of regression coefficients for latitude and longitude, one pair for each layer and cluster/dataset. Using these coefficients, we investigate to what extent different types of images are represented similarly in the models' latent space. Table 4 shows the average pearson correlation between each pair of coefficients and the average R^2 for a given model on its best performing layer overall.

It is possible to see that for all models, correlations are positive, suggesting that some spatial information is common for a variety of different image subjects. However, for most of the vision-only models, the correlation is very weak ($\rho < 0.3$). Meanwhile, for the models pretrained with language data, we see that correlations are low ($0.3 \le \rho < 0.5$) to moderate ($0.5 \le \rho < 0.7$), indicating some shared neurons for geospatial representation across different image types. These correlations get much higher when we consider only the 40% most important dimensions for predicting coordinates (see Section 4.5), with CLIP and LLaVA achieving values above 0.7, as shown in Table 5.

Table 4: Average correlation (ρ) between regression coefficients for different clusters/datasets for each model. The \pm term denotes the 95% confidence interval for the mean.

Model	Avg. ρ (Longitude Coef.)	Avg. ρ (Latitude Coef.)	Avg. R2
ViT-MAE-base	0.115 ± 0.005	0.142 ± 0.007	0.035 ± 0.008
ViT-MAE-large	0.111 ± 0.004	0.143 ± 0.006	0.051 ± 0.010
ViT-base	0.153 ± 0.004	0.164 ± 0.005	0.118 ± 0.020
ViT-large	0.229 ± 0.005	0.231 ± 0.006	0.145 ± 0.020
DINOv2-base	0.294 ± 0.007	0.283 ± 0.006	0.191 ± 0.033
DINOv2-large	0.291 ± 0.006	0.312 ± 0.007	0.236 ± 0.037
DINOv2-giant	0.300 ± 0.005	0.294 ± 0.006	0.262 ± 0.040
CLIP-ViT-base	0.500 ± 0.005	0.410 ± 0.005	0.278 ± 0.034
Qwen2.5-VL-3B-Instruct	0.348 ± 0.004	0.388 ± 0.005	0.344 ± 0.047
Gemma-3-4B-IT	0.372 ± 0.005	0.426 ± 0.005	0.382 ± 0.046
Qwen2.5-VL-7B-Instruct	0.325 ± 0.004	0.373 ± 0.005	0.398 ± 0.049
Gemma-3-12B-IT	0.450 ± 0.005	0.453 ± 0.005	0.419 ± 0.050
Gemma-3-4B-PT	0.378 ± 0.005	0.421 ± 0.005	0.421 ± 0.050
Gemma-3-12B-PT	0.182 ± 0.004	0.224 ± 0.005	0.450 ± 0.055
CLIP-ViT-large	0.587 ± 0.005	0.548 ± 0.006	0.482 ± 0.048
LLaVA-1.5-7B	0.612 ± 0.005	0.573 ± 0.006	0.510 ± 0.049

Table 5: Average correlation (ρ) between the 40% most important regression coefficients for different clusters/datasets for each model. The \pm term denotes the 95% confidence interval for the mean.

Model	Avg. ρ (Longitude Coef.)	Avg. ρ (Latitude Coef.)	Avg. R2
ViT-MAE-base	0.161 ± 0.006	0.195 ± 0.009	0.035 ± 0.008
ViT-MAE-large	0.153 ± 0.005	0.198 ± 0.007	0.051 ± 0.010
ViT-base-patch16-224	0.230 ± 0.005	0.253 ± 0.007	0.118 ± 0.020
ViT-large-patch16-224	0.351 ± 0.006	0.349 ± 0.008	0.145 ± 0.020
DINOv2-base	0.421 ± 0.008	0.407 ± 0.008	0.191 ± 0.033
DINOv2-large	0.413 ± 0.007	0.440 ± 0.008	0.236 ± 0.037
DINOv2-giant	0.434 ± 0.007	0.425 ± 0.007	0.262 ± 0.040
CLIP-ViT-base-patch32	0.674 ± 0.005	0.586 ± 0.005	0.278 ± 0.034
Qwen2.5-VL-3B-Instruct	0.498 ± 0.004	0.548 ± 0.005	0.344 ± 0.047
Gemma-3-4B-IT	0.399 ± 0.005	0.455 ± 0.005	0.382 ± 0.046
Qwen2.5-VL-7B-Instruct	0.459 ± 0.004	0.521 ± 0.005	0.398 ± 0.049
Gemma-3-12B-IT	0.470 ± 0.004	0.475 ± 0.005	0.421 ± 0.049
Gemma-3-4B-PT	0.416 ± 0.005	0.463 ± 0.005	0.421 ± 0.050
Gemma-3-12B-PT	0.224 ± 0.005	0.275 ± 0.005	0.443 ± 0.055
CLIP-ViT-large-patch14	0.749 ± 0.004	0.718 ± 0.005	0.482 ± 0.048
LLaVA-1.5-7B-HF	0.771 ± 0.004	0.739 ± 0.005	0.510 ± 0.049