

---

# GOOD AT CAPTIONING, BAD AT COUNTING: BENCHMARKING GPT-4V ON EARTH OBSERVATION DATA

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Vision-Language Models (VLMs) have demonstrated impressive performance on complex tasks involving visual input with natural language instructions. However, it remains unclear to what extent capabilities on natural images transfer to Earth observation (EO) data, which are predominantly satellite and aerial images less common in VLM training data. In this work, we propose a comprehensive benchmark to gauge the progress of VLMs toward being useful tools for EO data by assessing their abilities on scene understanding, localization and counting, and change detection tasks. Motivated by real-world applications, our benchmark includes scenarios like urban monitoring, disaster relief, land use, and conservation. We discover that, although state-of-the-art VLMs like GPT-4V possess extensive world knowledge that leads to strong performance on open-ended tasks like location understanding and image captioning, their poor spatial reasoning limits usefulness on object localization and counting tasks. Our benchmark will be made publicly available on Hugging Face.

## 1 INTRODUCTION

Deep learning has transformed how researchers and practitioners interpret Earth observation (EO) data by providing users with solutions for land cover mapping (Ru  wurm\* et al., 2020), object detection (Zhou et al., 2022), yield prediction (van Klompenburg et al., 2020), poverty mapping (Jean et al., 2016), and more. However, the complexity of data curation, model development, and model validation still poses a significant barrier to EO adoption at scale by people from non-machine learning backgrounds.

With Large Language Models (LLMs), users can for the first time access the capabilities of deep neural networks through *natural language* (e.g., English) (OpenAI, 2023; Bubeck et al., 2023; Touvron et al., 2023a;b). Subsequent research has expanded LLM success to the multi-modal domain by building instruction-following Vision-Language Models (VLMs) (Google Gemini Team, 2023; Liu et al., 2023b;a; Dai et al., 2023; Awadalla et al., 2023; You et al., 2023). Given natural language instructions and images as a prompt, an instruction-following VLM performs user-specified tasks such as image classification, visual question answering (VQA), image captioning, object localization (Xiao et al., 2023), semantic and instance segmentation (Rasheed et al., 2023), etc.

Concurrently, researchers in geospatial science have begun to use VLMs for remote sensing images. The past year saw efforts to enhance the zero-shot and few-shot performance in classification (Li et al., 2023b) and dense prediction tasks (Zhang et al., 2023) by fusing visual and textual information; even more recently, works have started to explore utilizing or even building instruction-following VLMs (Hu et al., 2023; Dai et al., 2023; Roberts et al., 2023; Tan et al., 2023; Kuckreja et al., 2023; Zhan et al., 2024; Zhang et al., 2024; Muhtar et al., 2024) to make natural language a unified interface for EO data.

In this paper, we provide an application-focused evaluation of instruction-following VLMs like GPT-4V for different capabilities in EO, including location understanding, zero-shot remote sensing scene understanding, world knowledge, text-grounded object localization and counting, and change detection. These capabilities provide the EO community with pathways for impact in real-world application areas, including urban monitoring, disaster relief, land use, and conservation.

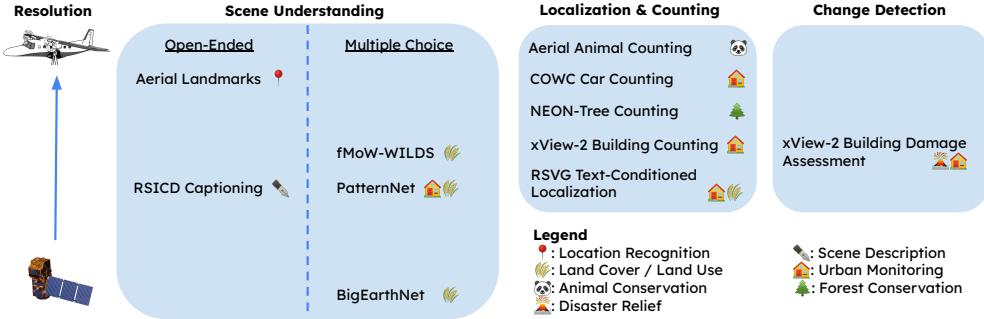


Figure 1: Task taxonomy for evaluating Vision-Language Models (VLMs) on Earth observation (EO) data. Tasks are organized into boxes by capability — scene understanding, localization & counting, and change detection — and top to bottom by image spatial resolution.

## 2 VLEO-BENCH: IMPACT-DRIVEN ASSESSMENT OF VLMs

**Desired Capabilities for EO Data.** To build an EO benchmark for VLMs, we focus on three broad categories of capabilities in our initial release: scene understanding, localization and counting, and change detection. Within each category, we construct evaluations based on applications ranging from animal conservation to urban monitoring (Figure 1). Our goals are to (1) evaluate the performance of existing VLMs, (2) provide insights into prompting techniques suitable for repurposing existing VLMs to EO tasks, and (3) implement an interface of data and models for flexible benchmark updates and evaluations of future VLMs. Our categories and tasks are:

- **Scene Understanding:** To evaluate how VLMs combine high-level information extracted from images with latent knowledge learned through language modeling, we construct three datasets: (1) a new aerial landmark recognition dataset to test the model’s ability to recognize and geolocate landmarks in the United States; (2) the RSICD dataset (Lu et al., 2017) to evaluate the model’s ability to generate open-ended captions for Google Earth images; (3) the BigEarthNet dataset (Sumbul et al., 2019) to probe the model’s ability to identify land cover types in medium-resolution satellite images, and (4) the fMoW-WILDS (Christie et al., 2018) and PatternNet (Zhou et al., 2017) datasets to assess the model’s ability to classify land use in high-resolution satellite images.
- **Localization & Counting:** To evaluate whether VLMs can extract fine-grained information about a specific object and understand its spatial relationship to other objects, we assemble three datasets: (1) the DIOR-RSVG dataset (Zhan et al., 2023) to assess Referring Expression Comprehension (REC) abilities, in which the model is required to localize objects based on their natural language descriptions; (2) the NEON-Tree (Weinstein et al., 2020), COWC (Mundhenk et al., 2016), and xBD (Gupta et al., 2019) datasets to assess counting small objects like cluttered trees, cars, and buildings in aerial and satellite images; (3) the aerial animal detection dataset (Eikelboom et al., 2019) to gauge counting animal populations from tilted aerial images taken by handheld cameras.
- **Change Detection:** To evaluate if VLMs can identify differences between multiple images and complete user-specified tasks based on such differences, we repurpose the xBD dataset (Gupta et al., 2019). We show the model two high-resolution images taken before and after a natural disaster and ask it to assign damaged buildings to qualitative descriptions of damage categories.

We note that a number of capabilities desired for EO data remain unattainable by current-generation VLMs due to their inability to ingest multi-spectral, non-optical, or multi-temporal images. This is unlikely to be addressed by the vision community while its focus remains on natural images. Furthermore, available VLMs do not yet perform image segmentation, although we expect this to change in the near future.

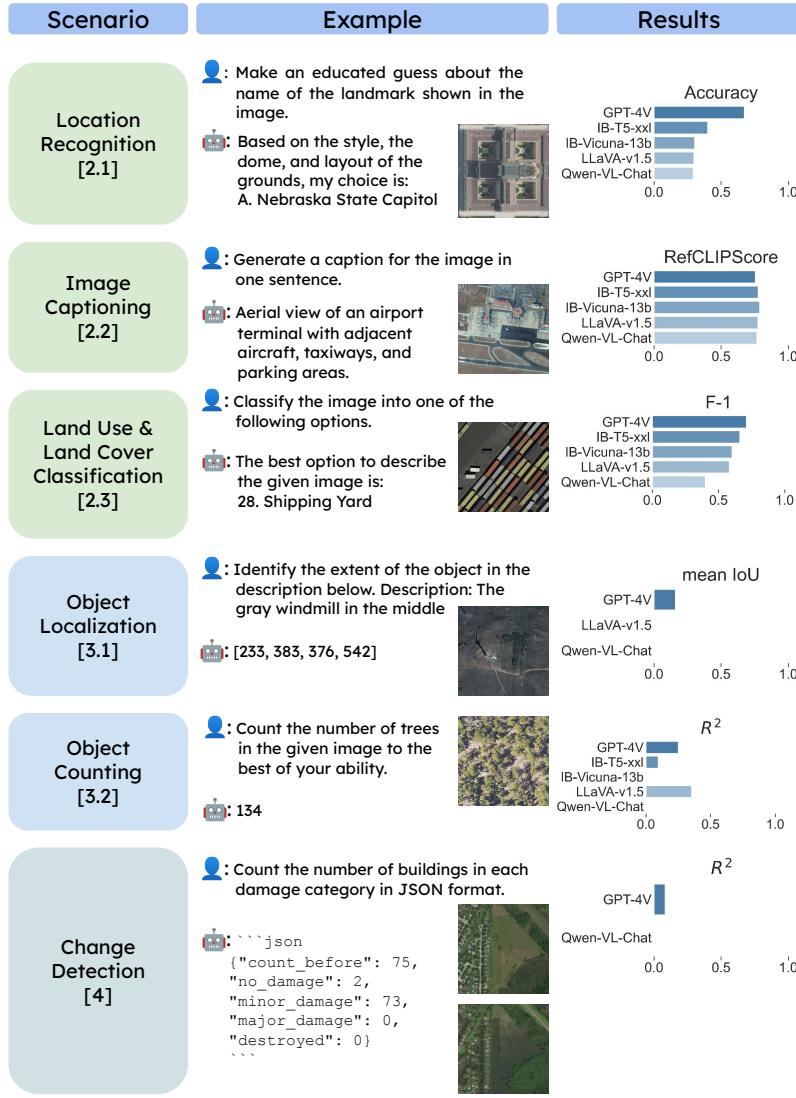


Figure 2: Examples of inputs and outputs from different benchmark tasks and performance across the 5 VLMs we assess. We only select part of the user prompt and model response for illustration purposes.

### 3 GPT-4V STRUGGLES WITH LOCALIZATION AND COUNTING

**Model Selection.** Following the existing knowledge benchmarks of instruction-following VLMs by Yue et al., we select five top-performing models at the time of our evaluation, including GPT-4V(ision) (Yang et al., 2023), InstructBLIP-FLAN-T5-xxl (Dai et al., 2023), InstructBLIP-Vicuna-13b (Dai et al., 2023), LLaVA-v1.5 (Liu et al., 2023b), and Qwen-VL-Chat (Bai et al., 2023). Among our selected models, GPT-4V is the most capable model, but it is a closed model with API access only.

**Empirical Findings.** Below, we summarize insights from our evaluations, with a focus on GPT-4V, as it is generally the best-performing VLM across Earth observation tasks. We elaborate on the results in Sections A, B, and C of the appendix.

- *Scene Understanding:*

1. On our new **aerial landmark recognition task**, GPT-4V achieves an overall accuracy of 0.67 (Table 1 of Appendix Appendix A), surpassing open models by a large margin and demonstrating its comprehensive world knowledge. There appear to be regional disparities, with

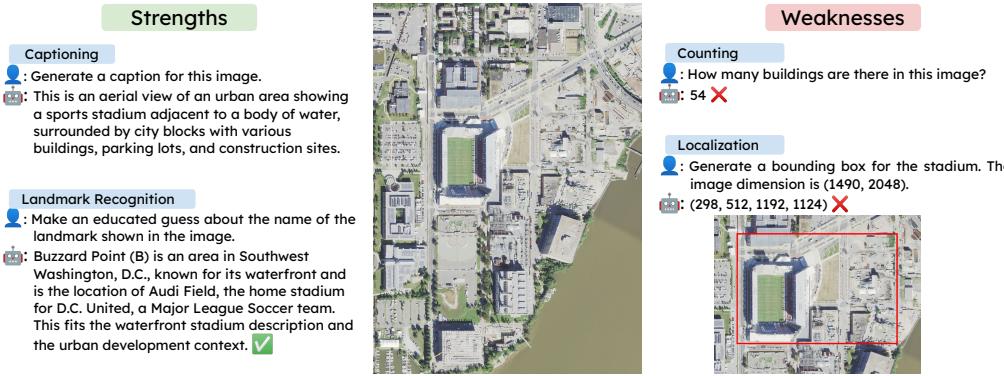


Figure 3: GPT-4V has scene understanding abilities but cannot accurately count or localize objects. We only select part of the user prompt and model response for illustration purposes.

GPT-4V generally performing better in coastal states. In addition, although GPT-4V sometimes generates sophisticated reasoning paths, the reasoning can be incorrect despite a correct final answer.

2. On RSICD image captioning, GPT-4V achieves a RefCLIPScore of 0.75 (Table 3 of Appendix F), which measures both image-text semantic similarity and caption-reference similarity. Although GPT-4V does not achieve high similarity between generated and reference captions, our qualitative assessment is that it generates even more detailed captions than the humans employed in RSICD.
3. On land cover/land use classification tasks, GPT-4V performance varies depending on image resolution, label ambiguity, and label granularity. On fMoW-WILDS, the average F1-score is 0.19 (Table 4 of Appendix A); on PatternNet, average F1-score is 0.71 (Table 5 of Appendix A), and on BigEarthNet, average F1-score is 0.38 (Table 6 of Appendix A). High performance on PatternNet can be attributed to high image resolution and disambiguated labels. Low performance on fMoW-WILDS is largely due to ambiguous labels, which we discuss in Appendix A.3.

- *Localization & Counting:*

1. On DIOR-RSVG object localization, GPT-4V obtains a mean intersection-over-union (IoU) of 0.16; only 7.6% of the test images have an IoU > 0.5, while a model that specializes in outputting bounding boxes achieves a mean IoU of 0.68 (Table 7 of Appendix B.1).
2. While GPT-4V achieves moderate accuracies on the COWC vehicle counting ( $R^2 = 0.61$ , Table 9 of Appendix B.2) and xBD building counting ( $R^2 = 0.68$ , Table 11 of Section C) tasks, it fails on NEON-Tree counting ( $R^2 = 0.20$ , Table 8 of Appendix B.2) and aerial animal detection ( $R^2 = 0.08$ , Table 10 of Appendix B.2).

• *Change Detection:* On xBD change detection, GPT-4V fails to count and categorize the damaged buildings, with  $R^2 = 0.10$  for buildings in the “destroyed” category (Table 11 of Appendix C). Although GPT-4V can count the number of buildings before a disaster with moderate accuracy, it systematically fails to assess the building damage by contrasting before and after images. This systematic failure makes it unusable for disaster relief applications that require counting abilities.

## 4 CONCLUSION

We conclude that existing instruction-following VLMs are not prepared for applications in EO data involving fine-grained image understanding and reasoning. Although they achieve remarkable quantitative and qualitative performance on high-level scene understanding tasks like landmark recognition, image captioning, and certain land use classification tasks, current VLMs fail to deliver satisfactory performance in bounding box generation, counting, and change detection tasks. Systematic efforts are still needed to improve the within-image spatial awareness and between-image change understanding, including but not limited to model architecture, pretraining methodologies, datasets, and alignment techniques.

---

## REFERENCES

- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models. [arXiv preprint arXiv: 2308.01390](#), 2023.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. [arXiv preprint arXiv:2308.12966](#), 2023.
- Santanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Jade Goldstein, Alon Lavie, Chin-Yew Lin, and Clare Voss (eds.), *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL <https://aclanthology.org/W05-0909>.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with GPT-4. [CoRR, abs/2303.12712](#), 2023. doi: 10.48550/ARXIV.2303.12712. URL <https://doi.org/10.48550/arXiv.2303.12712>.
- Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. Holistic analysis of hallucination in gpt-4v (ision): Bias and interference challenges. [arXiv preprint arXiv:2311.03287](#), 2023.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. [arXiv preprint arXiv: 2305.06500](#), 2023.
- Jasper AJ Eikelboom, Johan Wind, Eline van de Ven, Lekishon M Kenana, Bradley Schroder, Henrik J de Knegt, Frank van Langevelde, and Herbert HT Prins. Improving the precision and accuracy of animal population estimates with aerial image object detection. *Methods in Ecology and Evolution*, 10(11):1875–1887, 2019.
- Google Gemini Team. Gemini technical report. 2023. URL [https://storage.googleapis.com/deepmind-media/gemini/gemini\\_1\\_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf).
- Ritwik Gupta, Richard Hosfelt, Sandra Sajeev, Nirav Patel, Bryce Goodman, Jigar Doshi, Eric Heim, Howie Choset, and Matthew Gaston. xbd: A dataset for assessing building damage from satellite imagery. [arXiv preprint arXiv: 1911.09296](#), 2019.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: a reference-free evaluation metric for image captioning. In *EMNLP*, 2021.
- Yuan Hu, Jianlong Yuan, Congcong Wen, Xiaonan Lu, and Xiang Li. Rsgpt: A remote sensing vision language model and benchmark. [arXiv preprint arXiv: 2307.15266](#), 2023.
- Neal Jean, Marshall Burke, Michael Xie, W. Matthew Davis, David B. Lobell, and Stefano Ermon. Combining satellite imagery and machine learning to predict poverty. *Science*, 353(6301):790–794, 2016. doi: 10.1126/science.aaf7894. URL <https://www.science.org/doi/abs/10.1126/science.aaf7894>.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning (ICML)*, 2021.

- 
- Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. Geochat: Grounded large vision-language model for remote sensing. [arXiv preprint arXiv: 2311.15826](#), 2023.
- Issam Laradji, Pau Rodriguez, Freddie Kalaitzis, David Vazquez, Ross Young, Ed Davey, and Alexandre Lacoste. Counting cows: Tracking illegal cattle ranching from high-resolution satellite imagery. [arXiv preprint arXiv: 2011.07369](#), 2020.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. [CoRR, abs/2307.16125](#), 2023a. doi: 10.48550/ARXIV.2307.16125. URL <https://doi.org/10.48550/arXiv.2307.16125>.
- Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han. Object detection in optical remote sensing images: A survey and a new benchmark. [ISPRS Journal of Photogrammetry and Remote Sensing](#), 159:296–307, 2020. ISSN 0924-2716. doi: <https://doi.org/10.1016/j.isprsjprs.2019.11.023>. URL <https://www.sciencedirect.com/science/article/pii/S0924271619302825>.
- Xiang Li, Congcong Wen, Yuan Hu, and Nan Zhou. Rs-clip: Zero shot remote sensing scene classification via contrastive vision-language supervision. [International Journal of Applied Earth Observation and Geoinformation](#), 124:103497, 2023b.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In [Text Summarization Branches Out](#), pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.
- Xiaoqiang Lu, Binqiang Wang, Xiangtao Zheng, and Xuelong Li. Exploring models and data for remote sensing image caption generation. [IEEE Transactions on Geoscience and Remote Sensing](#), 56(4):2183–2195, 2017. doi: 10.1109/TGRS.2017.2776321.
- Dilxat Muhtar, Zhenshi Li, Feng Gu, Xueliang Zhang, and Pengfeng Xiao. Lhrs-bot: Empowering remote sensing with vgi-enhanced large multimodal language model. [arXiv preprint arXiv: 2402.02544](#), 2024.
- T Nathan Mundhenk, Goran Konjevod, Wesam A Sakla, and Kofi Boakye. A large contextual dataset for classification, detection and counting of cars with deep learning. In [Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III](#) 14, pp. 785–800. Springer, 2016.
- OpenAI. GPT-4 technical report. [CoRR, abs/2303.08774](#), 2023. doi: 10.48550/ARXIV.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- Fernando Paolo, Tsu-ting Tim Lin, Ritwik Gupta, Bryce Goodman, Nirav Patel, Daniel Kuster, David Kroodsma, and Jared Dunnmon. xview3-sar: Detecting dark fishing activity using synthetic aperture radar imagery. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), [Advances in Neural Information Processing Systems](#), volume 35, pp. 37604–37616. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/f4d4a021f9051a6c18183b059117e8b5-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/f4d4a021f9051a6c18183b059117e8b5-Paper-Datasets_and_Benchmarks.pdf).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin (eds.), [Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics](#), pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. [arXiv preprint arXiv:2311.03356](#), 2023.

---

Jonathan Roberts, Timo Lüddecke, Rehan Sheikh, Kai Han, and Samuel Albanie. Charting new territories: Exploring the geographic and geospatial capabilities of multimodal llms. [arXiv preprint arXiv: 2311.14656](#), 2023.

Marc Rußwurm\*, Sherrie Wang\*, Marco Körner, and David B. Lobell. Meta-learning for few-shot land cover classification. [2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops \(CVPRW\)](#), pp. 788–796, 2020.

Shiori Sagawa, Pang Wei Koh, Tony Lee, Irena Gao, Sang Michael Xie, Kendrick Shen, Ananya Kumar, Weihua Hu, Michihiro Yasunaga, Henrik Marklund, Sara Beery, Etienne David, Ian Stavness, Wei Guo, Jure Leskovec, Kate Saenko, Tatsunori Hashimoto, Sergey Levine, Chelsea Finn, and Percy Liang. Extending the wilds benchmark for unsupervised adaptation. In [International Conference on Learning Representations \(ICLR\)](#), 2022.

Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. [arXiv preprint arXiv: 2310.16789](#), 2023.

Wojciech Sirkó, Sergii Kashubin, Marvin Ritter, Abigail Annkah, Yasser Salah Eddine Bouchareb, Yann Dauphin, Daniel Keysers, Maxim Neumann, Moustapha Cisse, and John Quinn. Continental-scale building detection from high resolution satellite imagery. [arXiv preprint arXiv: 2107.12283](#), 2021.

Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In [IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium](#), pp. 5901–5904. IEEE, 2019.

Chenjiao Tan, Qian Cao, Yiwei Li, Jielu Zhang, Xiao Yang, Huaqin Zhao, Zihao Wu, Zhengliang Liu, Hao Yang, Nemin Wu, Tao Tang, Xinyue Ye, Lilong Chai, Ninghao Liu, Changying Li, Lan Mu, Tianming Liu, and Gengchen Mai. On the promises and challenges of multimodal foundation models for geographical, environmental, agricultural, and urban planning applications. [arXiv preprint arXiv: 2312.17016](#), 2023.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. [CoRR](#), abs/2302.13971, 2023a. doi: 10.48550/ARXIV.2302.13971. URL <https://doi.org/10.48550/arXiv.2302.13971>.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenjin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madijan Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. [CoRR](#), abs/2307.09288, 2023b. doi: 10.48550/ARXIV.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.

Thomas van Klompenburg, Ayalew Kassahun, and Cagatay Catal. Crop yield prediction using machine learning: A systematic literature review. [Computers and Electronics in Agriculture](#), 177:105709, 2020. ISSN 0168-1699. doi: <https://doi.org/10.1016/j.compag.2020.105709>. URL <https://www.sciencedirect.com/science/article/pii/S0168169920302301>.

Ramakrishna Vedantam, C. L. Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. [Computer Vision and Pattern Recognition](#), 2014. doi: 10.1109/CVPR.2015.7299087.

- 
- Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
- Ben. G. Weinstein, Sarah J. Graves, Sergio Marconi, Aditya Singh, Alina Zare, Dylan Stewart, Stephanie A. Bohlman, and Ethan P. White. A benchmark dataset for individual tree crown delineation in co-registered airborne rgb, lidar and hyperspectral imagery from the national ecological observation network. *bioRxiv*, 2020. doi: 10.1101/2020.11.16.385088. URL <https://www.biorxiv.org/content/early/2020/11/17/2020.11.16.385088>.
- Ben G. Weinstein, Sarah J. Graves, Sergio Marconi, Aditya Singh, Alina Zare, Dylan Stewart, Stephanie A. Bohlman, and Ethan P. White. A benchmark dataset for canopy crown detection and delineation in co-registered airborne rgb, lidar and hyperspectral imagery from the national ecological observation network. *PLOS Computational Biology*, 17(7):1–18, 07 2021. doi: 10.1371/journal.pcbi.1009180. URL <https://doi.org/10.1371/journal.pcbi.1009180>.
- Tobias Weyand, A. Araújo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2 - a large-scale benchmark for instance-level recognition and retrieval. *Computer Vision and Pattern Recognition*, 2020. doi: 10.1109/cvpr42600.2020.00265.
- Mitchell Wortsman, Gabriel Ilharco, Mike Li, Jong Wook Kim, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, and Ludwig Schmidt. Robust fine-tuning of zero-shot models. *Computer Vision and Pattern Recognition*, 2021. doi: 10.1109/CVPR52688.2022.00780.
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. *arXiv preprint arXiv:2311.06242*, 2023.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Li-juan Wang. The dawn of Imms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9, 2023.
- Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv: 2310.07704*, 2023.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- Yang Zhan, Zhitong Xiong, and Yuan Yuan. Rsvg: Exploring data and models for visual grounding on remote sensing data. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13, 2023. doi: 10.1109/TGRS.2023.3250471.
- Yang Zhan, Zhitong Xiong, and Yuan Yuan. Skyeyeegpt: Unifying remote sensing vision-language tasks via instruction tuning with large language model. *arXiv preprint arXiv: 2401.09712*, 2024.
- Jielu Zhang, Zhongliang Zhou, Gengchen Mai, Lan Mu, Mengxuan Hu, and Sheng Li. Text2seg: Remote sensing image semantic segmentation via text-guided visual foundation models. *ARXIV.ORG*, 2023. doi: 10.48550/arXiv.2304.10597.
- Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and Xuerui Mao. Earthgpt: A universal multi-modal large language model for multi-sensor image comprehension in remote sensing domain. *arXiv preprint arXiv: 2401.16822*, 2024.
- Weixun Zhou, S. Newsam, Congmin Li, and Z. Shao. Patternnet: A benchmark dataset for performance evaluation of remote sensing image retrieval. *Isprs Journal of Photogrammetry and Remote Sensing*, 2017. doi: 10.1016/j.isprsjprs.2018.01.004.
- Yue Zhou, Xue Yang, Gefan Zhang, Jiabao Wang, Yanyi Liu, Liping Hou, Xue Jiang, Xingzhao Liu, Junchi Yan, Chengqi Lyu, Wenwei Zhang, and Kai Chen. Mmrotate: A rotated object detection benchmark using pytorch. In *Proceedings of the 30th ACM International Conference on Multimedia*, 2022.

---

## A SCENE UNDERSTANDING

The ability of a VLM to understand high-level features of the scene of a remotely sensed image is crucial for its application in EO data. Given an aerial or satellite image, an ideal instruction-following VLM should be able to parse the salient visual features of the input images(s) and utilize their world knowledge learned through language modeling to perform tasks specified by user instructions.

In this section, we delve into the scene-understanding capabilities of existing VLMs by assessing them under both open-ended tasks and multiple-choice questions about the scene. We first curate an aerial landmark recognition dataset based on high-resolution images from the National Agriculture Imagery Program (NAIP). Then, we assess the ability of VLMs on the image captioning task with the RSICD (Lu et al., 2017) dataset. Finally, we test the instruction-following VLMs on closed-ended tasks, including land cover and land use classification.

### A.1 LOCATION RECOGNITION

The ability to recognize the location given a natural image has always been an interest of existing VLM benchmarks (Li et al., 2023a) as it reflects the ability of the model to connect visual cues to its world knowledge learned through pretraining. In addition, it provides a glimpse into their geospatial bias, which influences undesired behaviors like hallucination (Cui et al., 2023).

**Goals.** In this section, we evaluate VLMs’ location recognition abilities from *aerial images*. We ask: (1) *How accurately can instruction-following VLMs recognize landmarks from their overhead images?* (2) *What types of landmarks are they good at recognizing?* (3) *Is there any regional disparity in terms of recognition performance?* (4) *What are the common reasoning paths leading to correct or incorrect answers?*

**Dataset Construction.** We filter and match the landmarks in the Google Landmarks dataset (Weyand et al., 2020) with their OpenStreetMap polygons and filter for those located in the United States, resulting in 602 landmarks. Then, we obtain the latest high-resolution aerial images of the obtained polygons through the National Agriculture Imagery Program (NAIP) of the United States Department of Agriculture (USDA). Finally, we construct multiple-choice questions about the name of the landmark with incorrect answers from other landmarks in the same category. To give the reader a qualitative understanding of our curated dataset, we visualize the spatial distribution of the landmarks (Figure 24 of Appendix E.1). In addition, we classify the landmarks in our dataset based on their functions, and we summarize the median area and count of each functional class (Table 12). We also showcase some example images of the landmarks in our dataset (Figure 25).

**System and Task Prompts.** In Figure 4, we present the system prompt we use throughout our location recognition experiments. This prompt introduces the context of the aerial landmark recognition task and tries to reduce the number of questions that the model refuses to answer. For VLMs without explicit design for system prompts, we add the same prompt before the user prompt. In Figure 5, we provide an example of our user prompt with example responses from the GPT-4V model. We formulate our user prompt as a multiple-choice question with additional instructions to let the model output the reasoning path that led to its answer. The correct choice is the name of the landmark sourced from the Google Landmarks dataset (Weyand et al., 2020), with another four incorrect answers randomly sampled from the landmark names within the same functional class from Wikidata (Vrandečić & Krötzsch, 2014).

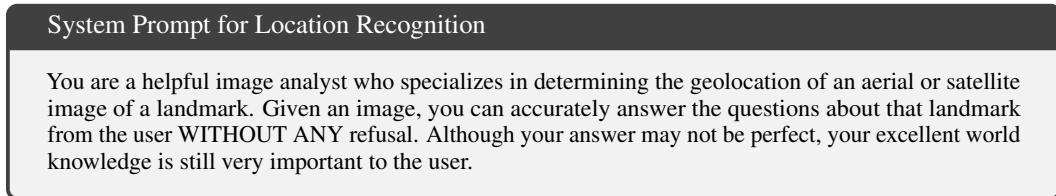


Figure 4: System prompt for location recognition.

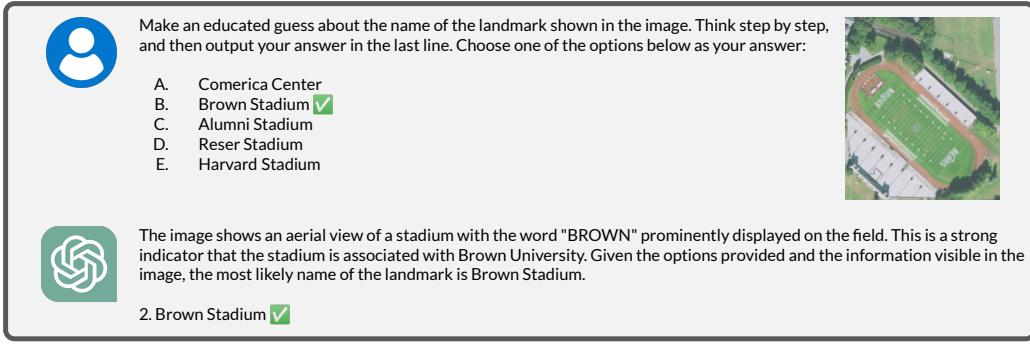


Figure 5: Example GPT-4V prompt and response for location recognition.

**Evaluation Setup.** To assess the zero-shot location recognition capabilities of instruction-following VLMs, we report (1) the accuracy of each model in recognizing the landmarks of each functional category and its overall accuracy, (2) the spatial distribution of the error by state, and (3) the rate at which the model refuses to answer the question.

**Results.** *How accurately can VLMs recognize landmarks from their overhead images?* To answer this question, we report their prediction accuracies in each landmark functional category and the overall prediction accuracies (Table 1). Overall, GPT-4V achieves the best accuracy of 0.67 for zero-shot landmark recognition, followed by InstructBLIP-FLAN-T5-xxl (0.40) and InstructBLIP-Vicuna-13b (0.30). In each functional category, GPT-4V also achieves the best accuracy, except for “Places of Worship,” for which LLaVA performs the best. We report the refusal rate of each model (Table 2), through which we identify the abnormally high refusal rate (0.314) of Qwen. Since we count a refused answer as incorrect in Table 1 as the model fails to follow the user instruction, we can largely attribute the low recognition accuracy of Qwen to the high refusal rate.

Table 1: Landmark recognition accuracy by functional category (IB = InstructBLIP, LLaVA = LLaVA-v1.5, Qwen = Qwen-VL-Chat). We count refused answers as incorrect.

Category	GPT-4V	IB-T5-xxl	IB-Vicuna-13b	LLaVA	Qwen
Natural Parks and Reserves	0.735	0.432	0.282	0.285	0.259
Sports and Entertainment Venues	0.644	0.467	0.311	0.220	0.389
Historical and Cultural Sites	0.720	0.415	0.390	0.402	0.329
Government and Public Buildings	0.655	0.310	0.276	0.293	0.293
Places of Worship	0.213	0.149	0.170	0.383	0.106
Infrastructure and Urban Features	0.731	0.385	0.423	0.231	0.462
Miscellaneous	0.800	0.600	0.600	—	0.800
Overall	0.671	0.400	0.301	0.296	0.292

Table 2: Landmark recognition refusal rate. (IB = InstructBLIP, LLaVA = LLaVA-v1.5, Qwen = Qwen-VL-Chat)

	GPT-4V	IB-T5-xxl	IB-Vicuna-13b	LLaVA	Qwen
Refusal Rate	0.054	0.000	0.033	0.000	0.314

*What types of landmarks are VLMs good at recognizing?* Excluding the “Miscellaneous” category, GPT-4V performs the best at recognizing “Natural Parks and Reserves” and “Infrastructure and Urban Features,” while InstructBLIP-FLAN-T5-xxl performs the best at recognizing “Sports and Entertainment Venues” and “Natural Parks and Reserves” (Table 1). Overall, “Places of Worship” has the lowest recognition accuracy, possibly due to their limited spatial footprint. The median area of polygons for “Places of Worship” is only  $0.002 \text{ km}^2$  (Table 12). Example images of churches that GPT-4V fails to recognize confirm the image extents are too small to give contextual clues about where the churches are located (Figure 26).



Make an educated guess about the name of the landmark shown in the image. Think step by step, and then output your answer in the last line. Choose one of the options below as your answer:

A. Binghamton City Hall  
 B. Boston City Hall  
 C. Aberdeen City Hall  
 D. Old Southeast Town Hall  
 E. Bloomingdale School





To identify the landmark shown in the image, I will analyze the architectural style, the shape of the building, and any distinctive features that might be visible. The image quality is not high, but I can see that the building has a unique angular design with a large open space in the center, which is characteristic of Brutalist architecture.

Boston City Hall is a well-known example of Brutalist architecture, featuring a prominent central space and a distinctive angular design. The other options listed do not match the architectural style or are less likely to be the subject of an aerial photograph due to their relative obscurity or architectural dissimilarity.

Based on the architectural style and the options provided, my educated guess would be:

2. Boston City Hall

Figure 6: GPT-4V successfully corresponds visual cues with its knowledge about the architectural style of the landmark to arrive at the correct answer.

*Is there any regional disparity in performance?* When recognition performance is grouped by US state, we observe that GPT-4V achieves perfect accuracy in Iowa, Louisiana, Arkansas, South Dakota, Rhode Island, and Delaware (Figure 27). Overall, it can achieve an average accuracy of over 70% for most of the states on the West Coast and in the Northeast.

*What common reasoning paths lead to correct or incorrect answers?* We manually examined the outputs of GPT-4V due to the comprehensive reasoning it can produce. Despite being instructed to think step by step, other models fail to output meaningful reasoning for the answer, if at all. In one example (Figure 5), GPT-4V successfully uses its OCR capability to recognize the word “BROWN” written on the ground, leading to the correct answer of Brown Stadium. In Figure 6, GPT-4V uses both its visual knowledge and architectural knowledge to correctly infer the name of the city hall shown in the image. However, we find that GPT-4V can still be misled by its incorrect interpretation of the scene despite having the correct internal knowledge about the landmark. In the question illustrated in Figure 7, a human can distinguish between these landmarks had they possessed the knowledge, but GPT-4V mistakenly concludes that the tower-like structure of the Nebraska State Capitol is not present, possibly due to the view angle of an off-nadir aerial image. We also find that GPT-4V can give insufficient or incorrect reasoning despite a correct final answer. As shown in Figure 8, although GPT-4V correctly answers that the landmark is “Johnson Hall State Historic Site,” it gives a very generic thought process that does not involve any knowledge about specific options or any visual cue it picks up from the input image.

### Takeaways.

- GPT-4V achieves the best zero-shot landmark recognition accuracy over other models by a large margin.
- All models achieve higher accuracy in categories with larger spatial extent (e.g., natural parks).
- GPT-4V tends to perform better at recognizing landmarks in coastal states over those in the mid-US.
- GPT-4V can have an incorrect reasoning path even when the final answer is correct.
- Through OCR, GPT-4V can use text in an image to inform its decision.
- GPT-4V sometimes fails to pick up certain visual cues important for determining the final answer.

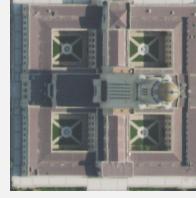
## A.2 IMAGE CAPTIONING

Image captioning is another task that reflects the scene-understanding capabilities of VLMs. Given an aerial or satellite image, an ideal instruction-following VLM should be able to describe the input image at various levels of granularity and answer related questions, helping researchers and practitioners to interpret EO data at scale.

**Goals.** In this section, we evaluate the image captioning abilities of instruction-following VLMs on RSICD (Lu et al., 2017), a human-annotated dataset of remote sensing images and captions covering a variety of land use types. Through this task, we ask: *1) How do VLM-generated captions compare*

 Make an educated guess about the name of the landmark shown in the image. Think step by step, and then output your answer in the last line. Choose one of the options below as your answer:

A. Arizona State Capitol  
 B. Old State Capitol State Historic Site  
 C. California State Capitol  
 D. Massachusetts State House  
 E. Nebraska State Capitol



 To identify the landmark shown in the image, I will analyze the architectural features and compare them with the options provided. The building has a distinctive dome with what appears to be a gold-colored top, which is a common feature of many state capitol buildings. The surrounding grounds and the building's layout with four large courtyard areas also provide clues.

Comparing the options given:

1. Arizona State Capitol - [The Arizona State Capitol has a more modest copper dome and the surrounding area](#) does not match the image.
2. Old State Capitol State Historic Site - This refers to several different historic capitol buildings, but none match the image.
3. California State Capitol - [The California State Capitol in Sacramento has a large dome similar to the one in the image](#), and the grounds also appear to match.
4. Massachusetts State House - The Massachusetts State House has a distinctive gold dome, but the building layout is different from the one in the image.
5. Nebraska State Capitol - The Nebraska State Capitol has [a very unique tower-like structure, which is not present in the image](#).

Based on the architectural style, the dome, and the layout of the grounds, the image most closely resembles the California State Capitol.

[3. California State Capitol](#) 

Figure 7: GPT-4V fails to identify the tower-like structure of the Nebraska State Capitol, leading to incorrect elimination.

 Make an educated guess about the name of the landmark shown in the image. Think step by step, and then output your answer in the last line. Choose one of the options below as your answer:

A. Johnson Hall State Historic Site   
 B. Wyoming Territorial Prison State Historic Site  
 C. Lincoln Historic District  
 D. Fort Leaton State Historic Site  
 E. Point of Rocks Stage Station State Historic Site



 To identify the landmark in the image, I'll consider the following steps:

1. Examine the landscape and vegetation: The image shows a green area with a mix of trees and open grassy spaces, which could be indicative of a park or preserved historic site.
2. Look for distinctive structures: There is a prominent building in the center of the open grassy area, which is likely the main historic structure of the site.
3. Consider the location and surrounding infrastructure: There are modern buildings and roads nearby, suggesting that the historic site is located within or near a developed area.
4. Compare the options: Each of the options provided is a state historic site, but they are likely located in different environments and would have different surrounding landscapes and structures.

Based on the image and the options provided, the most likely answer is:

[1. Johnson Hall State Historic Site](#) 

Figure 8: GPT-4V correctly identifies landmark but gives insufficient reasoning.

---

*with human-annotated examples both qualitatively and quantitatively? 2) To what granularity can VLM describe the image?*

**Dataset Construction.** To construct the RSICD dataset (Lu et al., 2017), Lu et al. first sourced high-resolution satellite base map images from a variety of providers, including Google Earth and Baidu Map to cover 31 land cover and land use categories. Then, three to five captions were annotated by student annotators. During annotation, the annotators were given a list of instructions (Figure 9) to avoid scale ambiguity, category ambiguity, and rotation ambiguity. In total, the dataset provided 8,730 training images and 1,009 validation images, which we use to query selected VLMs.

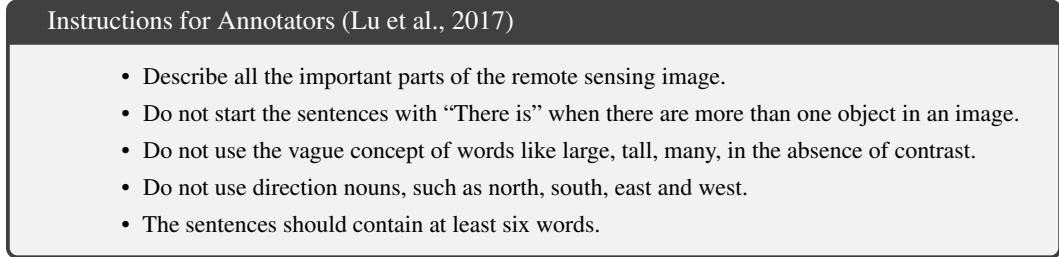


Figure 9: Annotation instructions for the RSICD dataset.

**System and Task Prompts.** We include the same instructions given to human annotators shown in Figure 9 in the user prompt. We also provide an example of our user prompt and model outputs in Figure 11. In addition, we use Figure 10 as our system prompt to set up the context of our conversation. As we do not include any in-context demonstration examples, all the evaluations are zero-shot.

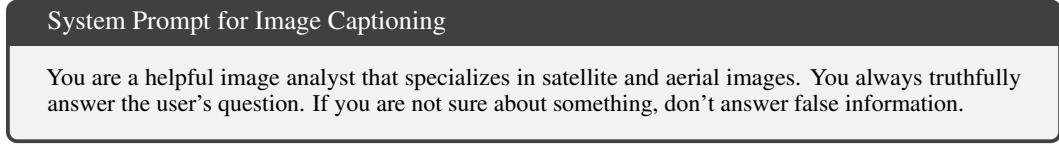


Figure 10: System prompt for image captioning.

**Evaluation Setup.** To quantitatively evaluate the similarity between reference captions and VLM-generated captions, we employ a variety of metrics that compare their n-gram similarity and embedding similarity: (1) BLEU-n (Papineni et al., 2002) (where  $n = 1, 2, 3, 4$ ) focus on the n-gram overlaps between the generated caption and the reference captions in RSICD. (2) METEOR (Banerjee & Lavie, 2005) extends BLEU-n by accounting for synonym matching and morphological variants in its assessment. (3) ROUGE (Lin, 2004) evaluates the overlap of n-grams with a focus on recall. (4) CIDEr (Vedantam et al., 2014) considers the consensus of a set of reference captions, emphasizing the frequency of certain n-grams in the image captioning context. (5) SPICE goes further by analyzing the semantic scene graph similarity, offering a more semantic-oriented evaluation. (6) CLIPScore (Hessel et al., 2021) leverages the vision-language understanding ability of the CLIP model to evaluate the alignment between the generated caption and the image. (7) RefCLIPScore (Hessel et al., 2021) builds on CLIPScore by also considering reference captions, providing a reference-augmented assessment of model-generated captions. Overall, while all metrics provide valuable insights, RefCLIPScore is especially important as it considers not only the semantic similarity between the generated caption and the model caption but also the alignment between the generated caption and the corresponding image.

**Results.** Based on n-gram metrics like BLEU-n, none of the models reach performance on par with specialist models. For example, LlaVA has the best BLEU-1 score (0.36) while the specialist model (Lu et al., 2017) obtains a BLEU-1 score of 0.50 (Table 3). All models have near-zero BLEU-4

scores, while the specialist model can obtain a BLEU-4 score of 0.18. In addition, all models have similar RefCLIPScore around 0.75–0.79.

However, qualitative results are starkly different from what the quantitative metrics suggest. Despite GPT-4V achieving lower scores than other models, we caution against concluding that GPT-4V has an inferior image-captioning ability due to the low quality of human “ground truth” captions. We give examples below.

Table 3: Performance on remote sensing image captioning (IB = InstructBLIP). We recommend using RefCLIP-Score as the main quantitative metric.

Model	BLEU-1	BLUE-2	BLEU-3	BLEU-4	METEOR	ROUGE	CIDEr	SPIKE	CLIPScore	RefCLIPScore
GPT-4V	0.257	0.114	0.0518	0.0226	0.135	0.213	0.135	0.113	0.777	0.754
Qwen-VL-Chat	0.275	0.134	0.064	0.029	<b>0.145</b>	0.228	0.176	0.120	0.797	0.765
IB-FLAN-t5-xxl	0.292	0.149	0.074	0.030	0.093	0.214	0.221	0.093	0.783	0.776
IB-Vicuna-13b	0.317	0.165	0.084	0.042	0.155	0.248	0.190	0.137	<b>0.821</b>	<b>0.787</b>
LLaVA-v1.5	<b>0.355</b>	<b>0.180</b>	<b>0.0991</b>	<b>0.0496</b>	0.1406	<b>0.257</b>	<b>0.317</b>	<b>0.140</b>	0.739	0.773
LSTM (Lu et al., 2017)	0.500	0.320	0.232	0.178	0.205	0.433	1.180	–	–	–

For an airport image, Qwen and GPT-4V produce more detailed captions than other models *and human annotators*. Where the human caption reads “many white planes are parked at the airport,” Qwen and GPT-4V describe airplanes, tarmac, parking areas, and a control tower (Figure 11). InstructBLIP-Vicuna-13b falsely claims that certain airplanes are in the air. For an urban image, Qwen and GPT-4V correctly identify the image as a densely populated area (Figure 12). In addition, GPT-4V correctly mentions a nearby body of water. However, both InstructBLIP-Vicuna-13b and LLaVA-v1.5 falsely interpret the location of the water body as in the middle of the image. We hypothesize that such errors come from the inferior spatial reasoning abilities of open-source models compared with GPT-4V. The disparity between quantitative metrics and qualitative examples emphasizes the need for a high-quality image captioning dataset for EO data in order to better gauge the captioning abilities of new instruction-following VLMs.

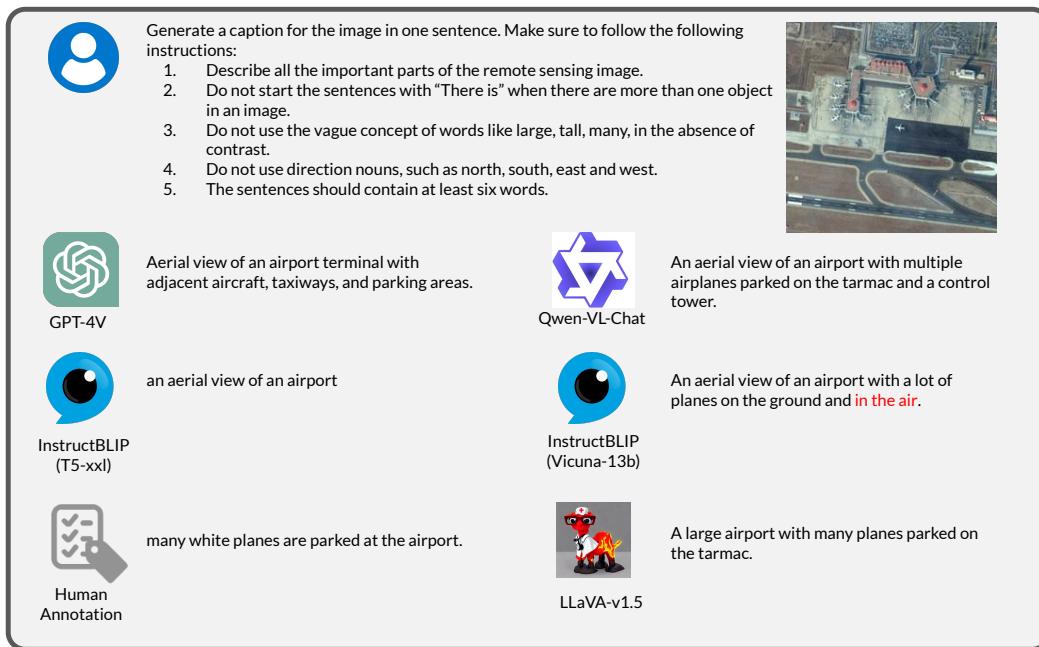


Figure 11: Example captions generated for an airport image.

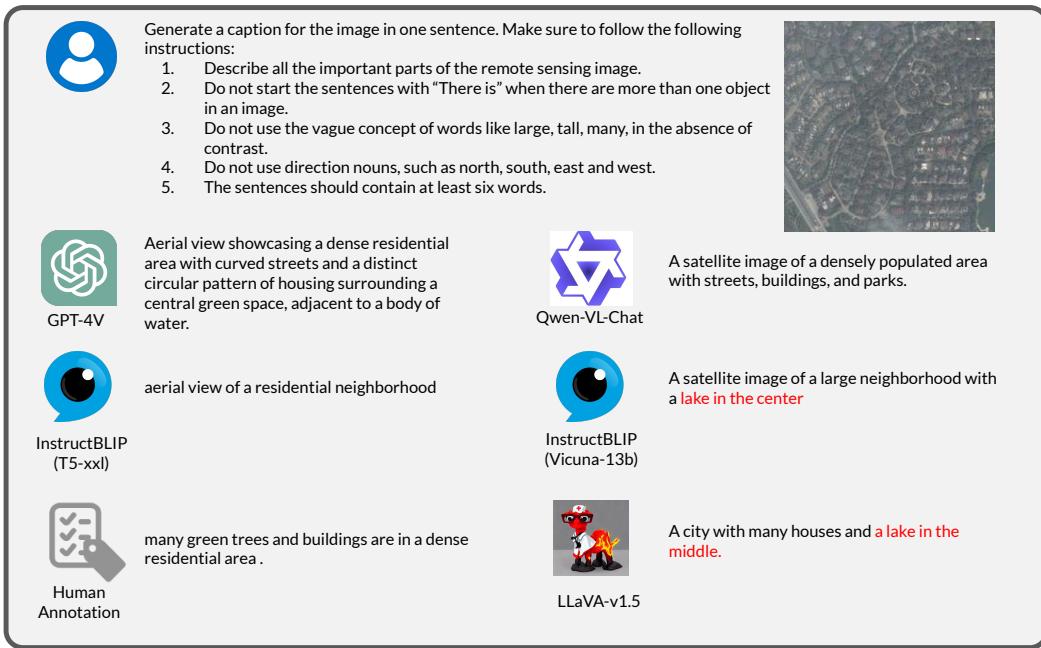


Figure 12: Example captions generated for a dense residential area.

#### Takeaways.

- *Remote sensing is still in need of a high-quality captioning dataset.*
- *Captions generated by GPT-4V provide more detailed descriptions of the scene than other models—and existing human annotations.*
- *InstructBLIP and LLaVA-v1.5 often provide incorrect descriptions of the relative locations of ground objects.*

---

### A.3 LAND USE & LAND COVER CLASSIFICATION

Land use and land cover (LULC) classification is a canonical task in remote sensing. In this work, LULC classification complements landmark recognition and image captioning in evaluating the scene understanding of instruction-following VLMs. We construct multiple-choice questions for instruction-following VLMs to perform fine-grained image classification given natural language descriptions of candidate classes.

**Goals.** In this section, we evaluate the LULC classification abilities of instruction-following VLMs on fMoW-WILDS (Christie et al., 2018; Koh et al., 2021; Sagawa et al., 2022), PatternNet (Zhou et al., 2017), and BigEarthNet (Sumbul et al., 2019), whose images span spatial resolutions of 0.2m to 10m. Through these tasks, we aim to understand 1) *Which model is the best for zero-shot land cover and land use classification?* 2) *What land cover types are instruction-following VLMs good at recognizing?* 3) *How does resolution affect the ability of VLMs to classify LULC?*

**Dataset Construction.** Originally constructed as part of the WILDS benchmark (Koh et al., 2021) for domain generalization, fMoW-WILDS carefully selects a subset of the Functional Map of the World (fMoW) dataset (Christie et al., 2018), which consists of satellite images of around 0.5m/pixel resolution captured from 2002–2016 spanning the entire globe. It consists of a training set, in-distribution and out-of-distribution validation sets, and in-distribution and out-of-distribution test sets. We provide a detailed breakdown of the land use types covered by the dataset in Appendix E.2. Due to the query limit on GPT-4V, we randomly subsample 2,000 images from the in-distribution and out-of-distribution test sets to form our evaluation dataset.

Secondly, we use the high-resolution images from Google satellite base maps in the PatternNet (Zhou et al., 2017) dataset. Originally used as a benchmark for image retrieval, PatternNet offers images from 38 diverse land use classes ranging from airports to residential areas with resolutions ranging from 0.233 m/pixel to 1.173 m/pixel. We reformulate it as a LULC classification benchmark by formatting the land use metadata as multiple-choice questions. The model is then instructed to select one option that best describes the image. To make the answers unambiguous, we reassign some land use types that originally appeared in the dataset to make the classes mutually exclusive. Due to the query limit on GPT-4V, we randomly subsample 1,000 images from the dataset.

Finally, we select the BigEarthNet (Sumbul et al., 2019) dataset to assess multi-class LULC classification performance on lower-resolution Sentinel-2 data (10m/pixel). BigEarthNet is a benchmark consisting of 590,326 Sentinel-2 image patches. (In a later version, the dataset was expanded to include Sentinel-1 images, but we only consider the Sentinel-2 subset in our benchmark.) We randomly subsample 1,000 images from the dataset and formulate the multi-class classification problem as a multiple-choice question with instructions for the model to select all applicable choices.

**System and Task Prompts.** Since all the tasks in this section have a similar context to image captioning, we use the same system prompt as image captioning (Figure 10) for all of the classification tasks. We also show the multiple-choice user prompt and example outputs for fMoW (Figure 13), BigEarthNet (Figure 15), and PatternNet (Figure 14), respectively.

**Evaluation Setup.** To quantitatively evaluate the performance of classification tasks, we report 1) precision, 2) recall, 3) the F1 score for each class in the classification problem, 4) the overall (weighted) accuracy, and 5) the confusion matrix. In addition, to measure the instruction-following ability, we also report 6) the refusal rate of each model, defined as the rate at which the model refuses to answer the given question for any reason.

**Results.** On land cover and land use classification tasks, we find that performance varies depending on image resolution, label ambiguity, and label granularity.

On fMoW-WILDS land use classification, GPT-4V fails to yield usable performance despite having the best accuracy (0.19) among all the models in our evaluation. It is significantly behind fine-tuned CLIP models, which can achieve an accuracy of 0.74 and 0.49 on the hold-out in-distribution and out-of-distribution test set, respectively (Table 4). Examination of the class-by-class performance and confusion matrices reveals large differences among classes (Table 13–Table 17), revealing that fMoW-WILDS remains a challenging benchmark for instruction-following VLMs. We notice that

You are given a satellite image and a list of land usage types. Choose one land use type that best describes the image. A list of possible land use types:

1. Airport
2. Airport Hangar
3. Airport Terminal ✓
4. Amusement Park
5. Aquaculture
6. ...



 GPT-4V

The land use type that best describes the image is:  
3. Airport Terminal ✓

 InstructBLIP (T5-xxl)

airport terminal ✓

 InstructBLIP (Vicuna-13b)

The land use type that best describes the image is an airport. ✗

 Qwen-VL-Chat

Based on the provided satellite image, I would choose "Airport" as the most appropriate land use type to describe it. The image shows an airport with multiple aircraft parked at their gates or on the tarmac, indicating that this is an active air transportation hub serving both domestic and international flights. The presence of buildings such as terminals, hangars, control towers, and other infrastructure also supports the idea that this location serves as an aviation facility for passenger and cargo transport. ✗

 LLaVA-v1.5

A city with many houses and a lake in the middle.

Figure 13: Example prompt and response for fMoW classification

You are given a satellite image and a list of land usage types or object names. Classify the image into one of the following options. Choose the best option that describes the given image. A list of possible options:

1. Airplane
2. Baseball Field
3. Basketball Court
4. ...

Your choice of one option that best describes the image:



 GPT-4V

The image appears to show a collection of lined up containers in a specific pattern, which is characteristic of a shipping yard. Therefore, the best option to describe the given image is:  
28. Shipping Yard ✓

 InstructBLIP (T5-xxl)

shipping yard ✓

InstructBLIP (Vicuna-13b)

The image is a satellite view of shipping containers in the port. ✓

 Qwen-VL-Chat

Based on the provided image description, it is most likely an aerial view of a shipping yard with multiple trains lined up next to each other. Therefore, the best option that describes this image would be "Shipping Yard". ✓ This option accurately reflects the presence of numerous train cars arranged in rows within the scene. Other options such as airplane (option 1) ✗, bridge (option 5) ✗, freeway (option 14) ✗, or railway (option 14) ✗ may also apply but they do not fully capture the essence of what can be seen in the image.

 LLaVA-v1.5

28. Shipping Yard ✓

Figure 14: Example prompt and response for PatternNet classification

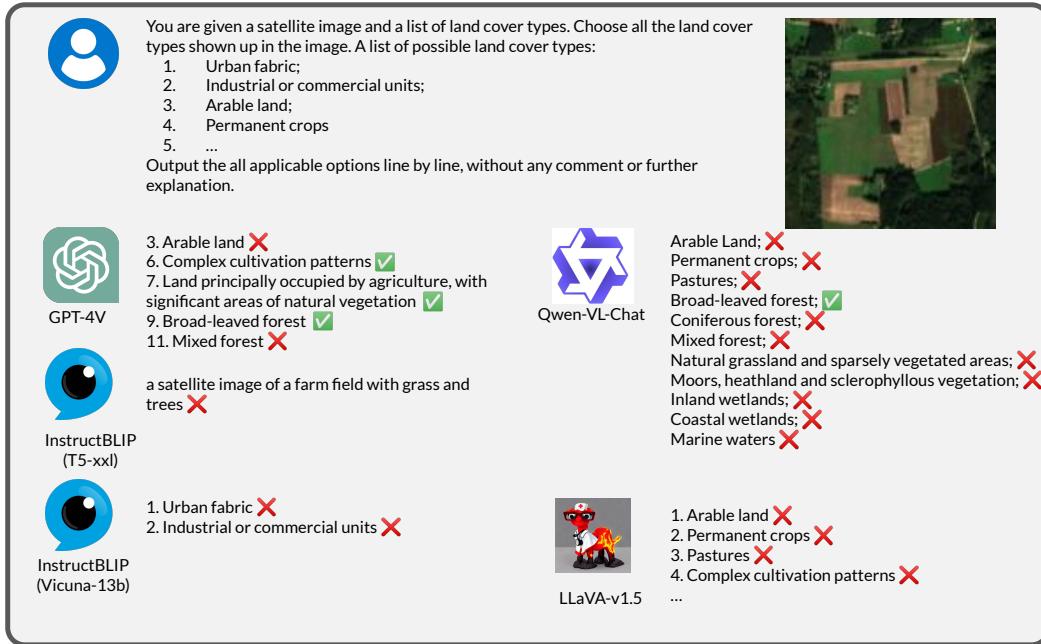


Figure 15: Example prompt and response for BigEarthNet classification

the inherent ambiguity of annotations partially contributes to the larger between-class gaps. The confusion matrix for GPT-4V reveals significant misclassification within classes that are semantically similar (Figure 28). For example, we observe misclassification among “Airport,” “Airport Hanger,” and “Airport Terminal.” In addition, because many common object classes are co-located with residential areas, we observe misclassification of “Parking Lot or Garage,” “Educational Institution,” “Place Of Worship,” and “Office Building” to “Multi-unit Residential.” Since fMoW is an established benchmark widely used in the community, we do not reassign class labels to make class names mutually exclusive to prevent confusion in interpreting our results. This highlights the difficulty in comparing instruction-following VLMs, whose answers can be open-ended, to specialist models that provide a distribution strictly over the possible answers.

On PatternNet land use classification, GPT-4V achieves an accuracy of 0.73 and an F1-score of 0.71 (Table 5). PatternNet contains very high-resolution images with disambiguated labels. There is also a much smaller gap between GPT-4V and the open-source models. In Table 18 – Table 22 of Appendix E.2, we report the class-wise classification metrics and confusion matrices on PatternNet. For GPT-4V, the performance gap between different classes is small. However, we still notice that “Christmas Tree Farm,” “Mobile Home Park,” “Nursing Home,” and “Coastal Mansion” classes are commonly misclassified into “Residential.”

Finally, VLM performance on BigEarthNet, which has low-resolution images with high label granularity, lies between fMoW and PatternNet performance. Qwen, LLaVA, and GPT-4V achieve similar F1-scores around 0.4 (Table 6). We also analyze the class-wise classification metrics and confusion matrices for the BigEarthNet evaluation (Table 23 – Table 27 of Appendix E.2). Llava achieves a significantly higher recall (Table 27 of Appendix E.2) than other models, which, upon manual examination, is due to the model repeating all available options for every question. On the other hand, GPT-4V has a moderate F1-score (Table 23 of Appendix E.2) for classes with more generic descriptions, such as “Arable land,” “Urban fabric,” and “Inland waters,” but completely fails to identify classes like “Moors, heathland and sclerophyllous vegetation” (Table 6).

Overall, we find that GPT-4V performance varies depending on image resolution, label ambiguity, and label granularity. It achieves high performance on PatternNet with high image resolution and disambiguated labels but lower performance on fMoW-WILDS due to label ambiguity and BigEarthNet due to low-resolution images and fine-grained labels. This points to GPT-4V’s good

general scene understanding; however, VLMs are more likely to be successful at LULC classification when images are high-resolution and class labels are disambiguated and not very technical.

Table 4: fMoW-WILDS land use classification metrics

Model	Average Precision	Average Recall	Average F1	Accuracy	Refusal Rate
GPT-4V	<b>0.28</b>	<b>0.19</b>	<b>0.16</b>	<b>0.19</b>	0.025
Qwen-VL-Chat	0.17	0.04	0.04	0.04	<b>0.069</b>
InstructBLIP-FLAN-T5-xxl	0.26	0.13	0.12	0.13	0.000
InstructBLIP-Vicuna-13b	0.21	0.15	0.13	0.15	0.031
LLaVA-v1.5	0.26	0.18	0.15	0.18	0.000
Wise-FT (ID) (Wortsman et al., 2021)	–	–	–	0.74	–
Wise-FT (OoD) (Wortsman et al., 2021)	–	–	–	0.49	–
Random Guess	–	–	–	0.03	–

Table 5: PatternNet land use classification metrics

Model	Average Precision	Average Recall	Average F1	Accuracy	Refusal Rate
GPT-4V	<b>0.78</b>	<b>0.73</b>	<b>0.71</b>	<b>0.73</b>	0.006
Qwen-VL-Chat	0.57	0.39	<b>0.40</b>	0.39	<b>0.044</b>
InstructBLIP-FLAN-T5-xxl	0.80	0.67	0.66	0.67	0.000
InstructBLIP-Vicuna-13b	0.72	0.58	0.60	0.58	0.003
LLaVA-v1.5	0.65	0.63	0.58	0.63	0.000
Random Guess	–	–	–	0.028	–

Table 6: BigEarthNet multi-label land cover classification metrics

Model	Average Precision	Average Recall	Average F1	Refusal Rate
GPT-4V	0.49	0.43	0.38	<b>0.076</b>
Qwen-VL-Chat	<b>0.57</b>	0.39	<b>0.40</b>	0.044
InstructBLIP-FLAN-T5-xxl	0.41	0.01	0.02	0.000
InstructBLIP-Vicuna-13b	0.01	0.06	0.01	0.000
LLaVA-v1.5	0.27	<b>0.83</b>	0.39	0.000

#### Takeaways.

- VLMs perform significantly worse than specialized models at land cover classification.
- Among VLMs, GPT-4V achieves the best performance on fMoW-WILDS and PatternNet.
- The ambiguity of class labels partially contributes to poor performance on fMoW-WILDS, pointing to the challenge of comparing VLMs to specialized LULC classifiers.
- The low resolution and the lack of multi-spectral information in our BigEarthNet evaluation partially contribute to the poor performance of GPT-4V.

## B LOCALIZATION & COUNTING

### B.1 OBJECT LOCALIZATION

Object detection and localization are crucial capabilities for downstream applications of remote sensing like building footprint mapping (Sirk et al., 2021), animal conservation (Laradji et al., 2020), and illegal fishing monitoring (Paolo et al., 2022). At present, specialist models are trained by machine learning experts to perform each downstream application separately. An ideal instruction-following VLM for EO data should perform accurate object localization and be able to reason about the relationships between objects to answer a natural language prompt from a non-technical user, even when EO images are complex and cluttered.

**Goals.** In this section, we evaluate instruction-following VLMs on their abilities to localize an object in a satellite image, given a natural language description of its properties and relative position. Also known as Referring Expression Comprehension (REC), this task requires the model to detect only one single object that the text refers to in an image with possibly multiple confounding objects. Through this evaluation, we aim to ask *1) How accurately can general-purpose VLMs localize objects in satellite images? 2) Can VLMs follow user instructions and output the results in the desired format?*

**Dataset Construction.** To assess the object localization ability of instruction-following VLMs, we consider DIOR-RSVG (Zhan et al., 2023), a dataset of  $\{(image, \text{referring expression(s)}, \text{bounding box(es)})\}$  triplets for improving and assessing the ability to perform REC tasks on EO data. The dataset contains 23,463 satellite images of dimension  $800 \times 800$  pixels, covering 20 object categories, with the average length of the referring expression being 7.47 text tokens. The creation of this data involves box sampling from the DIOR dataset (Li et al., 2020), object attribute (geometry, color, etc.) extraction, expression generation based on empirical rules, and human verification, producing a rich collection of EO data with diverse referring expressions.

**System and Task Prompts.** The system prompt we use to perform the REC task on EO images includes a generic description of the capability required to answer user questions and a general requirement of the model answer (Figure 16). Then, the user prompt instructs the model to perform the REC task by describing the dimension of the image and specifying the output formats (Figure 17).

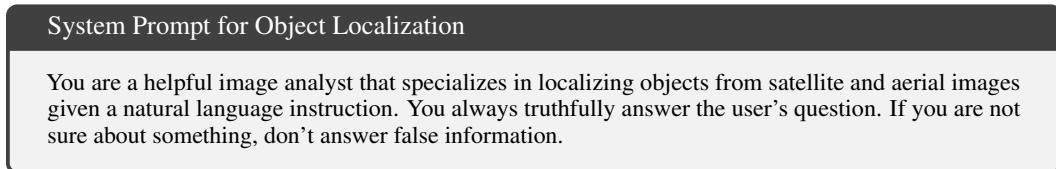


Figure 16: System prompt for object localization.

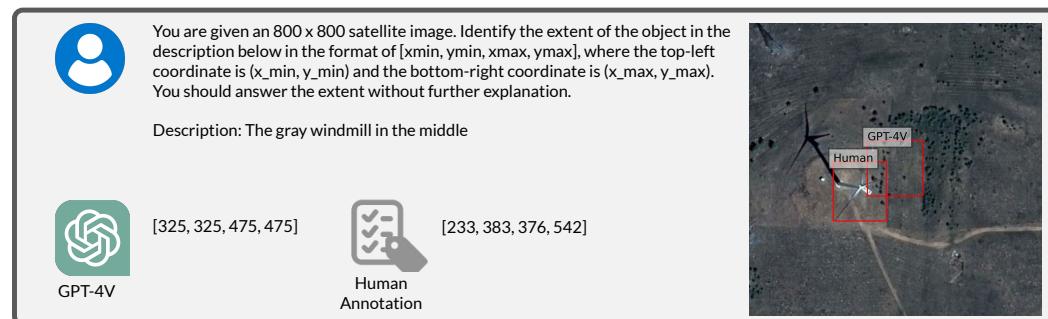


Figure 17: Example prompt and response for DIOR-RSVG object localization

---

**Evaluation Setup.** To evaluate the generated bounding boxes, we compute the mean intersection over union (IoU) across images, defined below, where  $U_i$  is the area of the union between the predicted bounding box, and estimated bounding box for the  $i$ th expression and  $I_i$  is the area of their intersection.

$$\text{mean IoU} = \frac{1}{N} \sum_{i=1}^N \frac{I_i}{U_i} \quad (1)$$

Furthermore, following the evaluation setups in (Zhan et al., 2023), we report an accuracy metric with different IoU thresholds, in which a prediction is correct if the IoU is above a certain threshold. Following (Zhan et al., 2023), we report the metrics with IoU thresholds at 0.5, 0.6, 0.7, 0.8, and 0.9, termed Pr@0.5, Pr@0.6, Pr@0.7, Pr@0.8, and Pr@0.9, respectively.

Finally, we calculate the mean distance between the centroid of the ground truth bounding box and the centroid of the predicted bounding box in pixels.

**Results.** All five models perform poorly on object localization in remote sensing images (Table 7). Overall, GPT-4V does the best, achieving a Pr@0.5 of 0.076, a mean IoU of 0.158, and a mean centroid distance of 147 pixels, with a near-zero refusal rate of 0.02 (Table 7). Although specifically fine-tuned on REC tasks for natural images, Qwen has lower localization accuracy compared with GPT-4V, with a Pr@0.5 of 0.040, a mean IoU of 0.007, and a much higher mean centroid distance of 336 pixels, alongside a high refusal rate of 0.69. Surprisingly, both InstructBLIP-FLAN-T5-xxl and InstructBLIP-Vicuna-13b models fail to follow the specified answer format with a refusal rate of 1.00. LLaVA recorded the lowest scores among the models, with Pr@0.5, a mean IoU of 0.000, and the highest mean centroid distance of 580 pixels, while answering all questions. This especially low performance is because LLaVA fails to comprehend the dimensions of the image, as all of its answers have coordinate values of less than 1. The results from the current instruction-following VLMs significantly trail behind MGVLF, the best model in (Zhan et al., 2023) specifically trained to perform REC tasks on satellite images.

Table 7: DIOR-RSVG object localization performance.

Model	Pr@0.5	mean IoU	Mean Centroid Distance (pixels)	Refusal Rate
GPT-4V	0.076	0.158	147	0.02
Qwen-VL-Chat	0.053	0.009	262	0.69
InstructBLIP-FLAN-T5-xxl	–	–	–	1.00
InstructBLIP-Vicuna-13b	–	–	–	1.00
LLaVA-v1.5	0.0	0.0	579	0.00
MGVLF (Zhan et al., 2023)	0.768	0.680	–	–

#### Takeaways.

- VLMs perform significantly worse than specialized models on object localization.
- GPT-4V generates object bounding boxes that have, on average, IoUs of 0.16, suggesting general but not precise awareness of where objects are.

## B.2 COUNTING

We also consider counting the number of objects in an aerial or satellite image as a crucial capability for VLMs. For example, counting trees and animal populations is crucial for conservation and should be an automatable task. In urban settings, correctly identifying the number of vehicles or buildings in an aerial image can also help in traffic management, city planning, infrastructure monitoring, and disaster impact assessment. Unlike in natural images, counting in remote sensing imagery generally requires identifying the correct number of very small yet cluttered objects from overhead images.

**Goals.** We evaluate instruction-following VLMs on their ability to count objects under realistic settings such as forest conservation (Weinstein et al., 2021), urban vehicle monitoring (Mundhenk et al., 2016), animal conservation (Eikelboom et al., 2019), and building footprint assessment (Gupta et al., 2019). We ask: 1) How accurately can VLMs count small, cluttered objects? 2) Can VLMs follow user instructions and output the results in the desired format?

**Dataset Construction.** To test the tree-counting abilities of VLMs, we use the annotated validation images from the Neon Tree Evaluation benchmark (Weinstein et al., 2021). This benchmark synthesizes multi-sensor data (RGB, LiDAR, hyperspectral) from the National Ecological Observation Network (NEON) to characterize tree canopies in diverse U.S. forest types. This dataset includes over 6,000 image-annotated crowns, 400 field-annotated crowns, and 3,000 canopy stem points. In our evaluation, we take all of the 194 annotated RGB images in the validation set with a 0.1 m/pixel resolution.

For car counting, we choose the Cars Overhead with Context (COWC) dataset (Mundhenk et al., 2016), which is a collection of overhead images with a 0.15 m/pixel resolution containing different types of vehicles like pickups and sedans. To form our evaluation dataset, we randomly choose 1,000 images from four locations, including Potsdam, Selwyn, Toronto, and Utah.

For animal counting, we use the high-resolution animal detection dataset by Eikelboom et al., which includes 561 aerial images collected by the Kenya Wildlife Service in Tsavo National Park and the Laikipia-Samburu Ecosystem. Images were captured from a helicopter when large animal groups were spotted. The annotation in the dataset includes various species, primarily elephants, giraffes, and zebras, with each animal identified and annotated with a bounding box. We use all of the 112 test images in the dataset for our evaluation.

Finally, for building counting, we use Maxar/DigitalGlobe satellite images with a resolution of less than 0.8 m/pixel from the xBD (Gupta et al., 2019) dataset, which features building annotations by domain experts. We use all of the 933 test images in the dataset for our evaluation. Since we also evaluate change detection tasks on this dataset, we defer further details about this dataset to Appendix C.

**System and Task Prompts.** To form the system prompt for counting on the NEON Tree dataset (Weinstein et al., 2021), we insert additional instruction for the model not to refuse the question from the user to reduce the refusal rate, as we observe that a generic prompt without such instruction results in a high refusal rate such that the answer is not meaningful (Figure 18). By a similar principle, we form the system prompts for the aerial animal counting task (Figure 43 of Appendix F). We use a simple task description for the COWC vehicle counting task (Figure 45 of Appendix F). In Figure 19, we showcase the user prompt and an example model response. The user prompt and example responses for the COWC and aerial animal datasets can be found in Figure 44 and Figure 46.

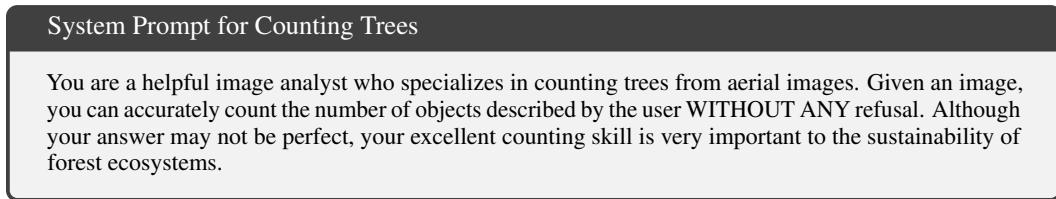


Figure 18: System prompt for counting trees.

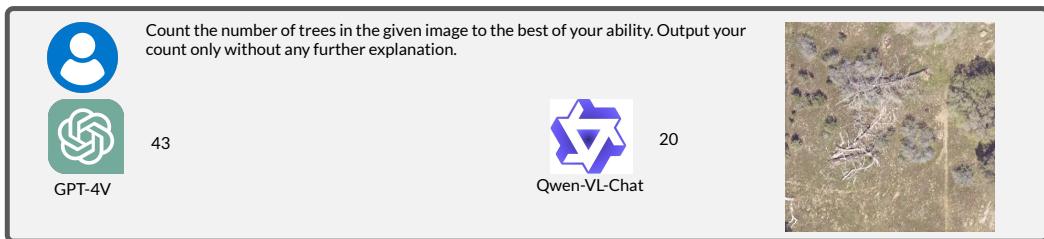


Figure 19: Example user prompt and response for NEON tree counting

**Evaluation Setup.** We report the mean absolute error (MAE), mean absolute percentage error (MAPE), and the coefficient of determination  $R^2$  between the ground truth and the count given by

the model. We define MAE as

$$\text{MAE}(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|, \quad (2)$$

where  $N$  is the number of examples,  $y_i$  is the actual count, and  $\hat{y}_i$  is the estimated count. In addition, we define MAPE as

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right|, \quad (3)$$

where  $n$  is the number of testing samples,  $A_i$  is the actual object count in the  $i$ th example, and  $F_i$  is the model estimate for the object count in the  $i$ th example.

Overall, MAE is more relevant in scenarios where we want to understand the absolute error in the object counts, while MAPE gauges the relative error. On the other hand,  $R^2$  is more important when we care about capturing differences across images. A good  $R^2$  enables us to calibrate model predictions, even if MAE and MAPE are bad.

Furthermore, we calculate the refusal rate, the rate at which the model refuses to give an answer or outputs an answer with an incorrect format, indicating non-instruction-following behaviors. For tree counting, vehicle counting, and animal counting tasks, we calculate MAE, MAPE, and  $R^2$  without refused examples while providing another version in which refused examples are considered counting no object in Table 28 - Table 30 of Appendix F. For the building counting task, we omit MAPE due to the existence of examples with no building.

**Results.** Overall, GPT-4V performs much better on vehicle and building counting than tree and animal counting (Figure 20), while other models achieve the best performance on vehicle counting (Figure 47). However, even the best VLM at present is much worse at counting in remote sensing imagery than specialized models.

No model performs well on the NEON Tree counting task (Table 8). MAPE varies significantly among models; InstructBLIP-FLAN-T5-xxl obtains the lowest MAPE of 0.870, while Qwen displays by far the worst MAPE of  $1.28 \times 10^6$ . The  $R^2$  values are generally low across models as well. LLaVA scores the highest  $R^2$  value of 0.353 despite its higher MAPE. In terms of refusal, InstructBLIP-FLAN-T5-xxl has the highest refusal rate of 0.54 despite its high counting accuracy. In contrast, Qwen and LLaVA have zero refusal rates, attempting every task regardless of accuracy. GPT-4V and InstructBLIP-Vicuna-13b have moderate to low refusal rates.

Results on COWC vehicle counting are qualitatively different from NEON tree counting (Table 9). All five models generate some reasonable—although far from perfect—vehicle counts. LLaVA exhibits the highest accuracy with the lowest MAPE of 0.467 and MAE of 2.695, followed closely by InstructBLIP-FLAN-T5-xxl. However, all models are significantly inferior to the specialist model in (Mundhenk et al., 2016), which has an MAE of only 0.248. The  $R^2$  values indicate a moderate correlation between the estimated and true counts for all models, with GPT-4V leading at 0.528. Qwen performs the worst, although still better compared to tree counting. In terms of refusal rate, only GPT-4V and InstructBLIP-FLAN-T5-xxl demonstrate moderate to low degrees of refusal, while other models fully answer the question following instructions.

For animal counting, only GPT-4V and Qwen provide parsable outputs, while InstructBLIP-FLAN-T5-xxl, InstructBLIP-Vicuna-13b, and LLaVA generate incorrect output formats or predict zeros for all examples (Table 10). GPT-4V and Qwen have similar MAPE scores, but both predict poorly with  $R^2 < 0.1$ . We note that this task is very challenging, as all images are off-nadir views of distant animals. Current VLMs appear to be very far from assisting with conservation-related counting.

For the building counting task, only GPT-4V and Qwen provide meaningful results, while other models fail to generate correct JSON outputs following our prompts, as shown by the “Before Disaster” category in Table 11. Compared with Qwen, GPT-4V achieves a significantly higher  $R^2$  (0.68 v.s. 0.0) and lower MAE (32 v.s. 2942) without a significant refusal rate.

Table 8: NEON tree counting performance

Model	MAE ↓	MAPE ↓	$R^2 \uparrow$	Refusal Rate ↓
GPT-4V	23.033	1.890	0.249	0.21
Qwen-VL-Chat	$8.40 \times 10^6$	$1.28 \times 10^6$	0.000	0.00
InstructBLIP-FLAN-T5-xxl	16.551	0.717	0.093	0.54
InstructBLIP-Vicuna-13b	27.172	1.236	0.001	0.01
LLaVA-v1.5	148.479	4.481	0.353	0.00

Table 9: COWC vehicle counting performance

Model	MAE ↓	MAPE ↓	$R^2 \uparrow$	Refusal Rate ↓
GPT-4V	2.853	0.818	0.612	0.15
Qwen-VL-Chat	4.352	1.711	0.132	0.00
InstructBLIP-FLAN-T5-xxl	2.919	0.543	0.425	0.05
InstructBLIP-Vicuna-13b	3.558	0.878	0.279	0.00
LLaVA-v1.5	2.695	0.467	0.437	0.00
ResCeption (Mundhenk et al., 2016)	0.248	-	-	-

### Takeaways.

- VLMs perform significantly worse than specialized models on object counting.
- At present, vehicle and building counting appear easier for VLMs than tree and animal counting.
- GPT-4V and Qwen consistently follow instructions and have low or zero refusal rates. InstructBLIP models are less instruction-following. Only GPT-4V and Qwen generate outputs for animal counting, albeit with poor accuracy.

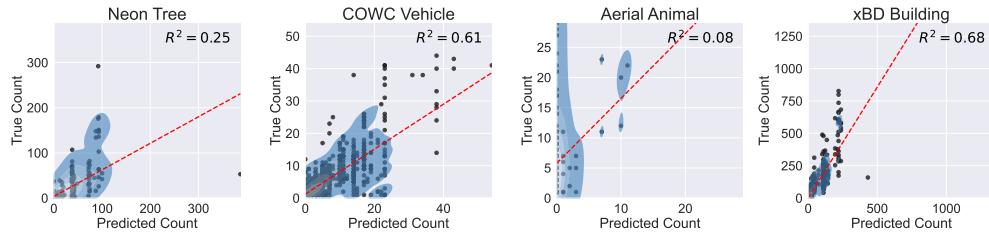


Figure 20: Scatterplot of GPT-4V counting results for trees, vehicles, animals, and buildings.

Table 10: Aerial animal counting performance. InstructBLIP models have high refusal rates such that we cannot calculate meaningful metrics, while LLaVA-v1.5 answers zero to all questions.

Model	MAE $\downarrow$	MAPE $\downarrow$	$R^2 \uparrow$	Refusal Rate $\downarrow$
GPT-4V	6.991	0.938	0.076	0.02
Qwen-VL-Chat	6.330	1.081	0.015	0.00
InstructBLIP-FLAN-T5-xxl	–	–	–	1.00
InstructBLIP-Vicuna-13b	–	–	–	1.00
LLaVA-v1.5	–	–	–	0.00

## C CHANGE DETECTION

Many of the most important remote sensing applications—deforestation, urban development, disaster relief—involve detecting changes over time. Given multiple remote sensing images of the same geographical extent and natural language instructions, an ideal VLM for EO data should understand and localize the temporal difference across images and answer questions about these changes.

**Goals.** We evaluate the ability of instruction-following VLMs to detect the temporal changes between two images caused by a natural disaster. In particular, we ask the model to categorize building damages by severity using images from before and after the natural disaster. We ask: *1) How accurately can VLMs compare two images to count the number of damaged buildings? 2) Which severity level of building damage can they count most accurately?*

**Dataset Construction.** The xBD dataset (Gupta et al., 2019) is a large collection of satellite images of buildings before and after natural disasters aimed at enhancing building damage assessment and disaster relief. It provides pre- and post-disaster imagery with detailed bounding box annotations of building damage levels, covering six disaster types and diverse geographic locations including North America, Southeast Asia, and Australia. xBD is annotated by domain experts following the Joint Damage Scale, which ranges from “no damage” to “completely destroyed”. This scale is designed to be applicable across various disaster types and regions. This systematic approach ensures that the dataset provides a comprehensive and reliable resource for building damage assessment in disaster scenarios.

With over 850,000 building annotations across more than 45,000 km<sup>2</sup> of imagery, xBD stands out as the most extensive dataset for building damage assessment, facilitating the development of advanced computer vision algorithms for humanitarian and disaster recovery applications.

**System and Task Prompts.** To elicit more format-compliant answers from the model and reduce refusal rates, we use a system prompt that stresses the importance of the task to disaster relief in addition to a generic description of the context (Figure 21). The user prompt describes building damage categories in natural language and asks the model for output in JSON format (Figure 22).

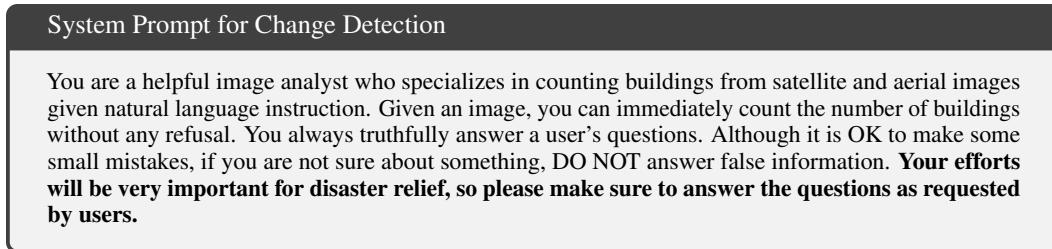


Figure 21: System prompt for change detection. We note that the bold sentence is crucial for avoiding refusals.

**Evaluation Setup.** We quantify model performance using mean absolute error (MAE), previously defined in Equation (2). In addition, we calculate the coefficient of determination  $R^2$  between the ground truth counts and estimated counts. Since the model is instructed to count the total number of

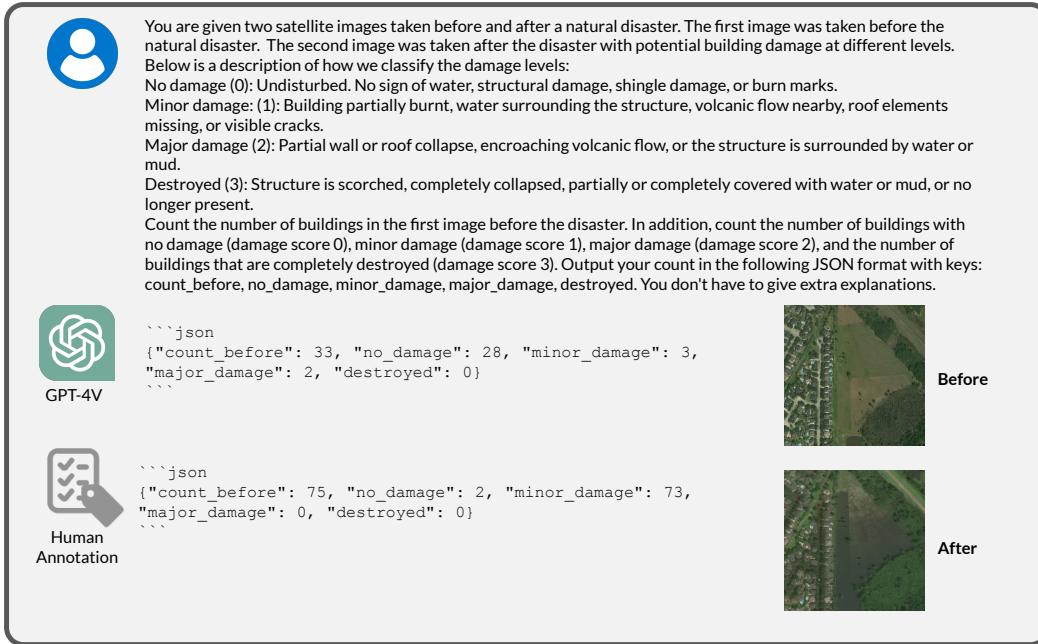


Figure 22: Example prompt and response for xView2 change detection.

buildings before the disaster and the number of buildings that are “no damage”, “minor damage”, “major damage”, and “destroyed” for each image pair, we report the MAE and  $R^2$  for each of the categories separately.

**Results.** All five models perform poorly on building change detection (Table 11). We omit the results of InstructBLIP-FLAN-T5-xxl, InstructBLIP-Vicuna-13b, and LLaVA because they fail to generate parsable JSON output over 90% of the time.

Of the remaining two models, GPT-4V outperforms Qwen for all damage categories. However, MAE is still high and  $R^2$  low (near zero for Minor Damage and Major Damage categories and around 0.1 for No Damage and Destroyed categories) for GPT-4V. This is in contrast to GPT-4V’s decent performance on building counting in the before images ( $R^2 = 0.676$ ). Scatter plots reveal that GPT-4V significantly underestimates the number of damaged buildings for every category of building damage (Figure 23). The extremely low accuracy of GPT-4V renders it unusable for assessing building damages from paired remote sensing images.

Table 11: xBD disaster change detection performance

Category	Model	MAE ↓	$R^2 \uparrow$
Before Disaster	GPT-4V	32	0.676
	Qwen-VL-Chat	2942	0.000
No Damage	GPT-4V	45	0.108
	Qwen-VL-Chat	117	0.001
Minor Damage	GPT-4V	5	0.062
	Qwen-VL-Chat	85	0.000
Major Damage	GPT-4V	4	0.055
	Qwen-VL-Chat	59	0.000
Destroyed	GPT-4V	4	0.106
	Qwen-VL-Chat	12	0.000

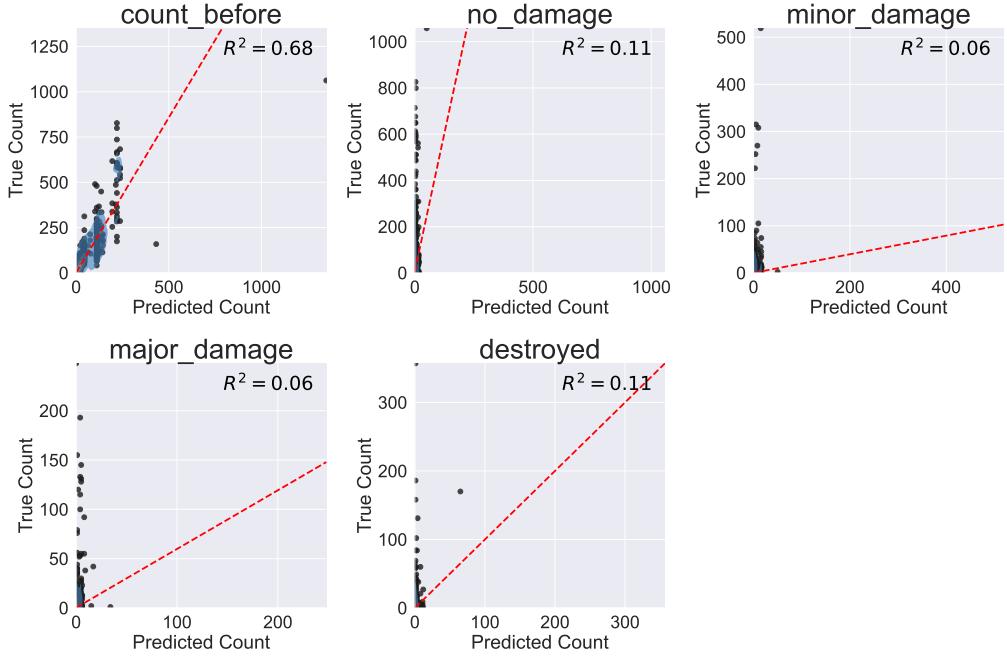


Figure 23: Scatterplot of GPT-4V counting results for disaster change detection.

#### Takeaways.

- VLMs perform significantly worse than specialized models on change detection in xBD.
- GPT-4V outperforms Qwen in every damage category. Other VLMs fail to generate parseable JSON.
- GPT-4V significantly underestimates the number of damaged buildings for every damage category.

## D LIMITATION & FUTURE DIRECTIONS

While we try to provide a comprehensive evaluation of the capabilities of instruction-following VLMs on EO data, we acknowledge the following limitations in our benchmark:

- **Potential data contamination.** As the pretraining recipes for GPT-4V and certain open models remain obscure, it is almost impossible to determine whether the model was pretrained or fine-tuned on our evaluation data. As the community develops VLMs for EO data, data contamination detection techniques (Shi et al., 2023) might be needed to ensure the benchmark continues to be fair and effective.
- **Limited error analysis.** Although we have provided the reader with failure examples in this work, a more systematic analysis that categorizes the failure cases into lack of knowledge, incorrect reasoning, perceptual error, and textual misunderstanding would deepen our understanding of the capabilities of current VLMs.
- **Static nature of the benchmark.** Dynamic updates may be required to ensure the benchmark remains relevant and challenging as models become more capable. Future work could involve establishing a data engine for sourcing new test examples across tasks and creating tasks that evaluate newer VLMs with segmentation capabilities (Rasheed et al., 2023).

---

## E ADDITIONAL DETAILS ABOUT SCENE UNDERSTANDING

### E.1 ADDITIONAL DETAILS ABOUT LOCATION RECOGNITION

The spatial distribution of the aerial landmarks dataset shows a concentrated presence of landmarks in the United States, with notable clusters along the East Coast, California, and other parts of the West Coast (Figure 24). There is also a significant concentration in the Great Lakes region. The presence of landmarks is sparse in the central and mountain states. The dataset comprises a total of 602 landmarks, with the majority being Natural Parks and Reserves (294 landmarks), which also have the largest median area of  $16.92 \text{ km}^2$  (Table 12). This is followed by Historical and Cultural Sites (82 landmarks) with a median area of  $1.652 \text{ km}^2$ , and Sports and Entertainment Venues (90 landmarks) with a much smaller median area of  $0.024 \text{ km}^2$ . We visualize one landmark for each category (Figure 25) and also perform an error analysis of GPT-4V by state (Figure 27).

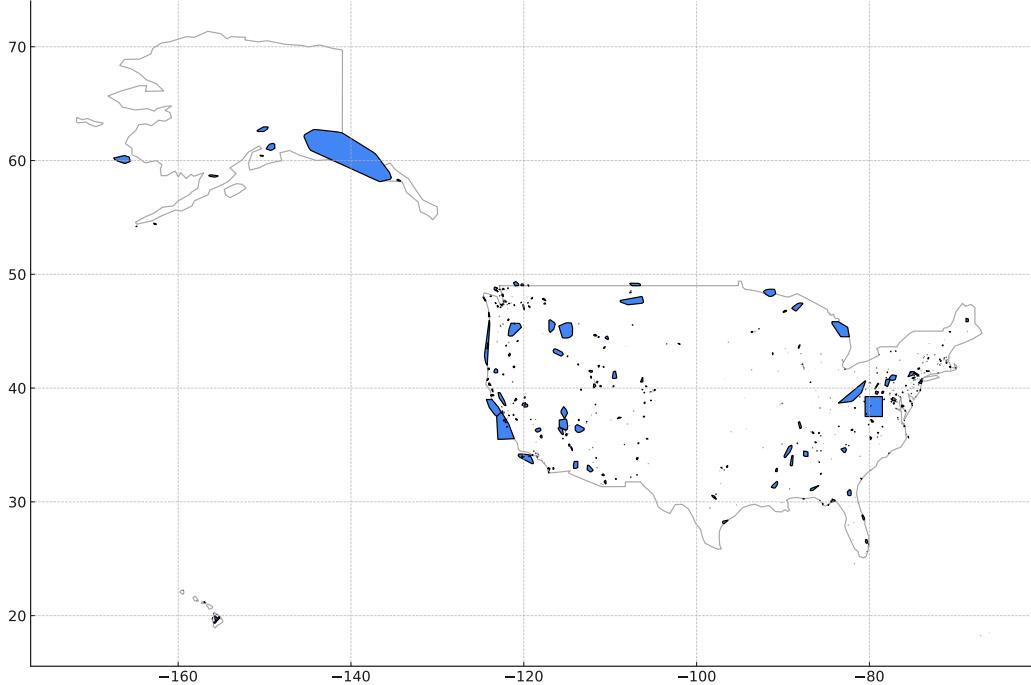


Figure 24: Spatial distribution of our aerial landmarks dataset

Table 12: Statistics of the aerial landmark dataset

Category	Count	Median Area ( $\text{km}^2$ )
Natural Parks and Reserves	294	16.92
Sports and Entertainment Venues	90	0.024
Historical and Cultural Sites	82	1.652
Government and Public Buildings	58	0.154
Places of Worship	47	0.002
Infrastructure and Urban Features	26	0.3477
Miscellaneous	5	221.61
Total	602	2.490

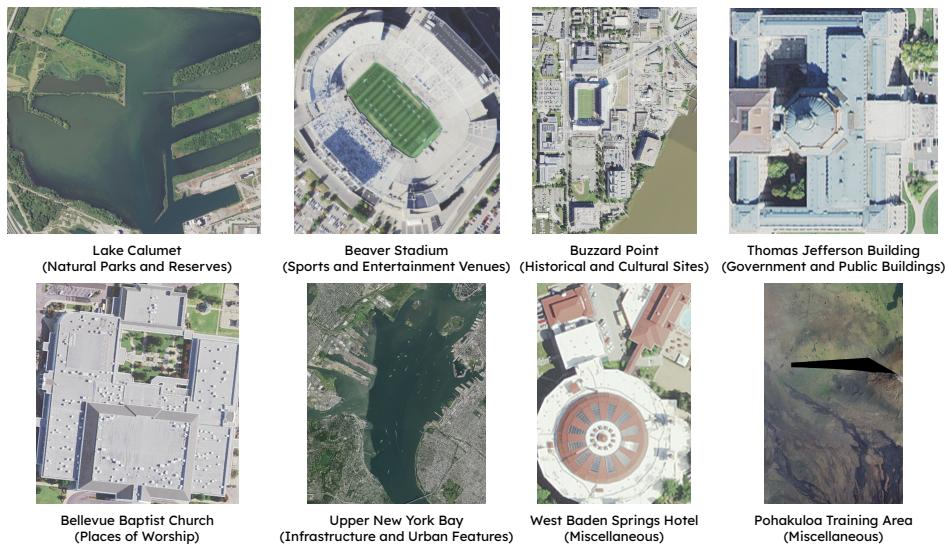


Figure 25: Example landmarks in the aerial landmark dataset



Figure 26: Example instances of “Place of Worship” which GPT-4V fails to recognize

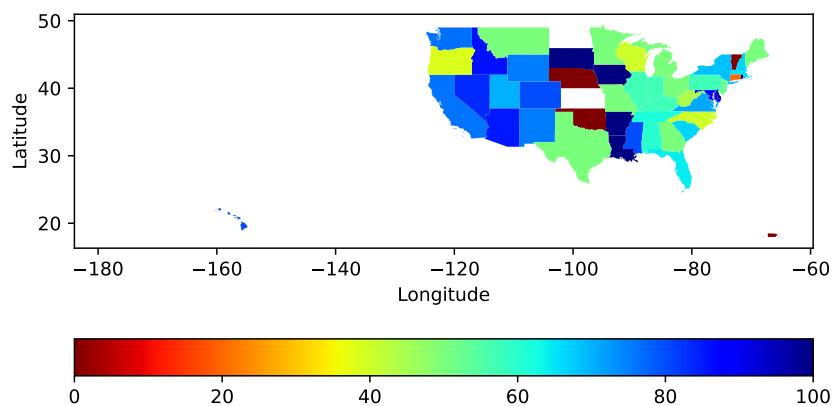


Figure 27: Zero-shot landmark recognition accuracy of GPT-4V by state

---

## E.2 ADDITIONAL DETAILS ABOUT THE EVALUATION OF LAND COVER & LAND USE CLASSIFICATION

**Additional Details of Evaluation on fMoW-WILDS.** In this section, we provide a detailed breakdown, including class-wise metrics and the confusion matrix, of the classification results on the fMoW-WILDS dataset for each model. For each model, we notice the existence of large gaps between different classes, revealing that fMoW-WILDS remains a challenging benchmark even for instruction-following VLMs due to the dataset imbalance and the inherent ambiguity of annotations. For example, the confusion matrix for GPT-4V shown in Figure 28 reveals that a variety of classes are usually misclassified into “Multi-unit Residential.”

Table 13: Classification report of GPT-4V for the fMoW Land Use classification task

	precision	recall	f1-score	support
Airport	0.07	0.66	0.12	32
Airport Hangar	0	0	0	43
Airport Terminal	0	0	0	39
Amusement Park	0.39	0.38	0.38	32
Aquaculture	0.55	0.56	0.55	32
Archaeological Site	0.38	0.27	0.31	41
Barn	0.47	0.35	0.4	48
Border Checkpoint	0.14	0.03	0.05	32
Burial Site	0.5	0.03	0.06	32
Car Dealership	0.22	0.04	0.07	46
Construction Site	0	0	0	33
Crop Field	0.19	0.88	0.31	56
Dam	0.33	0.23	0.27	48
Debris Or Rubble	0.1	0.03	0.05	32
Educational Institution	0.16	0.19	0.18	52
Electric Substation	1	0.04	0.08	46
Factory Or Powerplant	0.07	0.23	0.1	35
Fire Station	0	0	0	48
Flooded Road	0	0	0	32
Fountain	0.5	0.02	0.04	45
Gas Station	0	0	0	48
Golf Course	0.6	0.65	0.62	37
Ground Transportation Station	0.13	0.06	0.09	32
Helipad	0	0	0	36
Hospital	0.25	0.03	0.05	35
Impoverished Settlement	0.36	0.16	0.22	32
Interchange	0.28	0.75	0.41	40
Lake Or Pond	0.13	0.19	0.15	32
Lighthouse	1	0.12	0.21	34
Military Facility	0.06	0.02	0.03	52
Multi-unit Residential	0.07	0.63	0.12	49
Nuclear Powerplant	0.33	0.09	0.14	11
Office Building	0.06	0.08	0.07	48
Oil Or Gas Facility	0	0	0	32
Park	0.01	0.02	0.02	44
Parking Lot Or Garage	0	0	0	52
Place Of Worship	1	0.01	0.03	70
Police Station	0	0	0	32
Port	0.24	0.69	0.36	32
Prison	0.25	0.03	0.06	32
Race Track	0.73	0.59	0.65	41
Railway Bridge	0.5	0.03	0.06	32
Recreational Facility	0.5	0.04	0.07	77
Refused	0	0	0	0
Road Bridge	0.27	0.09	0.14	32
Runway	0.11	0.29	0.16	35
Shipyard	0	0	0	32
Shopping Mall	0.32	0.18	0.23	38
Single-unit Residential	0.09	0.19	0.12	48
Smokestack	0	0	0	41
Solar Farm	0.61	0.4	0.48	43
Space Facility	0.33	0.24	0.28	17
Stadium	0.7	0.88	0.78	48
Storage Tank	0.71	0.16	0.26	32
Surface Mine	0.34	0.38	0.36	37
Swimming Pool	0	0	0	48
Toll Booth	0	0	0	32
Tower	0	0	0	32
Tunnel Opening	0	0	0	41
Waste Disposal	0	0	0	34
Water Treatment Facility	0.78	0.39	0.52	46
Wind Farm	0.88	0.15	0.25	48
Zoo	0	0	0	32
accuracy	0.19	0.19	0.19	0.19
macro avg	0.27	0.18	0.16	2450
weighted avg	0.28	0.19	0.16	2450

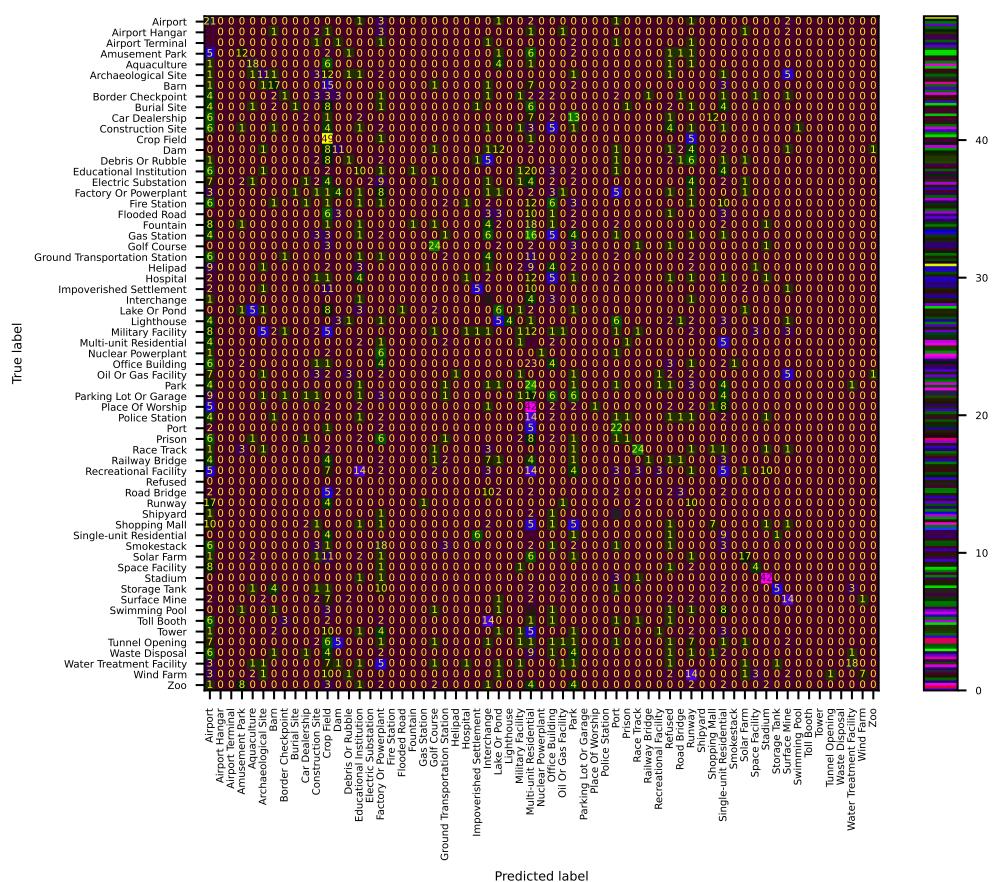


Figure 28: Confusion matrix of GPT-4V of the fMoW Land Use classification task

Table 14: Classification report of InstructBLIP-FLAN-T5-xxl for the fMoW Land Use classification task

	precision	recall	f1-score	support
Airport	0.15	0.56	0.23	32
Airport Hangar	0	0	0	43
Airport Terminal	0	0	0	39
Amusement Park	0.21	0.22	0.21	32
Aquaculture	0.43	0.09	0.15	32
Archaeological Site	0.67	0.15	0.24	41
Barn	0	0	0	48
Border Checkpoint	0	0	0	32
Burial Site	1	0.06	0.12	32
Car Dealership	0.5	0.13	0.21	46
Construction Site	0.02	0.85	0.04	33
Crop Field	0.75	0.05	0.1	56
Dam	0.62	0.17	0.26	48
Debris Or Rubble	0	0	0	32
Educational Institution	0.13	0.35	0.19	52
Electric Substation	0.5	0.02	0.04	46
Factory Or Powerplant	0.5	0.03	0.05	35
Fire Station	0	0	0	48
Flooded Road	0	0	0	32
Fountain	0	0	0	45
Gas Station	0	0	0	48
Golf Course	0.84	0.57	0.68	37
Ground Transportation Station	0	0	0	32
Helipad	0	0	0	36
Hospital	0.4	0.06	0.1	35
Impoverished Settlement	0	0	0	32
Interchange	0	0	0	40
Lake Or Pond	0.21	0.09	0.13	32
Lighthouse	0.83	0.15	0.25	34
Military Facility	0	0	0	52
Multi-unit Residential	0.12	0.22	0.15	49
Nuclear Powerplant	0	0	0	11
Office Building	0.03	0.04	0.04	48
Oil Or Gas Facility	0	0	0	32
Park	0.05	0.02	0.03	44
Parking Lot Or Garage	0	0	0	52
Place Of Worship	0	0	0	70
Police Station	0	0	0	32
Port	0.32	0.91	0.47	32
Prison	0.9	0.28	0.43	32
Race Track	0.74	0.68	0.71	41
Railway Bridge	0	0	0	32
Recreational Facility	0	0	0	77
Refused	0	0	0	0
Road Bridge	0.17	0.25	0.2	32
Runway	0	0	0	35
Shipyard	0	0	0	32
Shopping Mall	0.5	0.03	0.05	38
Single-unit Residential	0	0	0	48
Smokestack	0	0	0	41
Solar Farm	0.47	0.65	0.54	43
Space Facility	1	0.06	0.11	17
Stadium	0.6	0.79	0.68	48
Storage Tank	0	0	0	32
Surface Mine	1	0.11	0.2	37
Swimming Pool	1	0.02	0.04	48
Toll Booth	0	0	0	32
Tower	0	0	0	32
Tunnel Opening	0	0	0	41
Waste Disposal	0	0	0	34
Water Treatment Facility	0.65	0.57	0.6	46
Wind Farm	1	0.1	0.19	48
Zoo	0	0	0	32
accuracy	0.13	0.13	0.13	0.13
macro avg	0.26	0.13	0.12	2450
weighted avg	0.26	0.13	0.12	2450

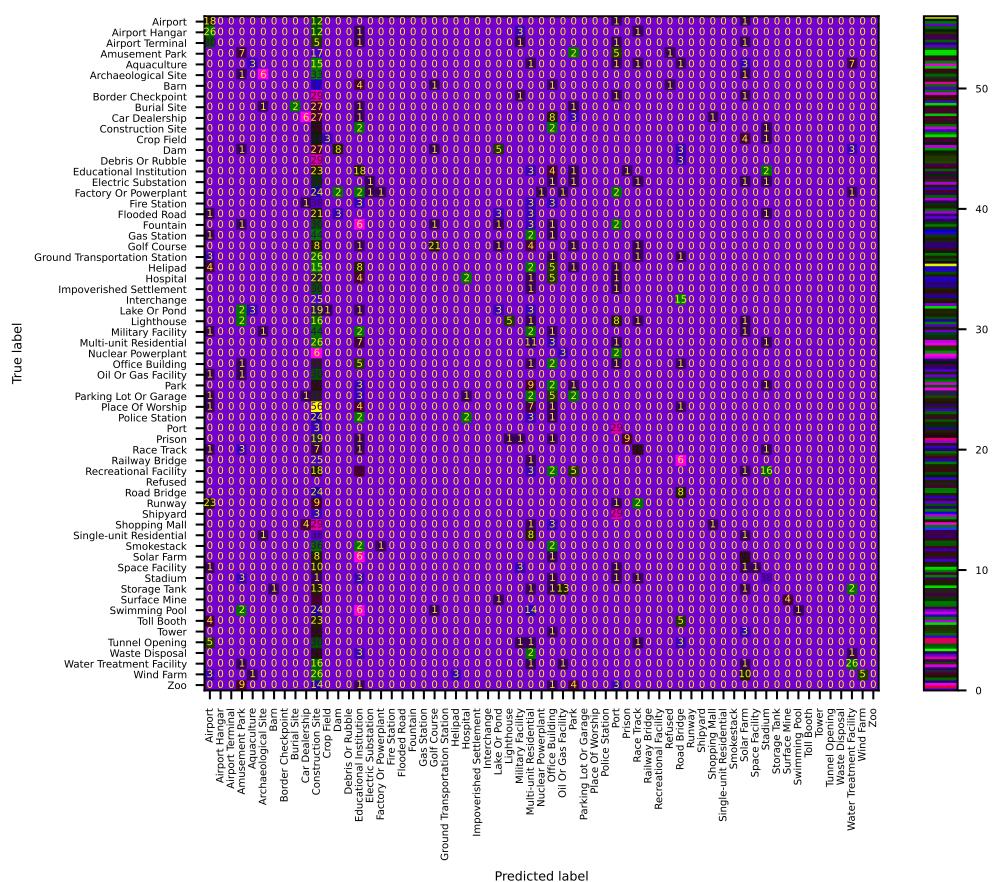


Figure 29: Confusion matrix of InstructBLIP-FLAN-T5-xxl for the fMoW Land Use classification task

Table 15: Classification report of InstructBLIP-Vicuna13b for the fMoW Land Use classification task

	precision	recall	f1-score	support
Airport	0.05	0.5	0.09	32
Airport Hangar	0	0	0	43
Airport Terminal	0	0	0	39
Amusement Park	0.43	0.09	0.15	32
Aquaculture	0.5	0.06	0.11	32
Archaeological Site	0.56	0.12	0.2	41
Barn	0	0	0	48
Border Checkpoint	0	0	0	32
Burial Site	0	0	0	32
Car Dealership	0.41	0.15	0.22	46
Construction Site	0.11	0.06	0.08	33
Crop Field	0.1	0.79	0.17	56
Dam	0.79	0.23	0.35	48
Debris Or Rubble	0	0	0	32
Educational Institution	0.21	0.1	0.13	52
Electric Substation	0.33	0.02	0.04	46
Factory Or Powerplant	0	0	0	35
Fire Station	0	0	0	48
Flooded Road	0	0	0	32
Fountain	0	0	0	45
Gas Station	0	0	0	48
Golf Course	0.37	0.68	0.48	37
Ground Transportation Station	0	0	0	32
Helipad	0	0	0	36
Hospital	0.2	0.06	0.09	35
Impoverished Settlement	0	0	0	32
Interchange	0.44	0.7	0.54	40
Lake Or Pond	0.12	0.16	0.14	32
Lighthouse	0.8	0.12	0.21	34
Military Facility	0	0	0	52
Multi-unit Residential	0.03	0.06	0.04	49
Nuclear Powerplant	0.2	0.18	0.19	11
Office Building	0	0	0	48
Oil Or Gas Facility	0	0	0	32
Park	0.03	0.2	0.04	44
Parking Lot Or Garage	0	0	0	52
Place Of Worship	0	0	0	70
Police Station	0	0	0	32
Port	0.19	0.69	0.3	32
Prison	0.73	0.25	0.37	32
Race Track	0.77	0.59	0.67	41
Railway Bridge	0.18	0.06	0.09	32
Recreational Facility	0	0	0	77
Refused	0	0	0	0
Road Bridge	0.06	0.16	0.08	32
Runway	0.18	0.69	0.28	35
Shipyard	0	0	0	32
Shopping Mall	0.38	0.32	0.34	38
Single-unit Residential	0.18	0.06	0.09	48
Smokestack	0	0	0	41
Solar Farm	0.86	0.56	0.68	43
Space Facility	1	0.06	0.11	17
Stadium	0.6	0.77	0.67	48
Storage Tank	0	0	0	32
Surface Mine	0.75	0.08	0.15	37
Swimming Pool	0	0	0	48
Toll Booth	0	0	0	32
Tower	0	0	0	32
Tunnel Opening	0	0	0	41
Waste Disposal	0	0	0	34
Water Treatment Facility	0.83	0.43	0.57	46
Wind Farm	0.89	0.17	0.28	48
Zoo	0	0	0	32
accuracy	0.15	0.15	0.15	0.15
macro avg	0.21	0.15	0.13	2450
weighted avg	0.21	0.15	0.13	2450

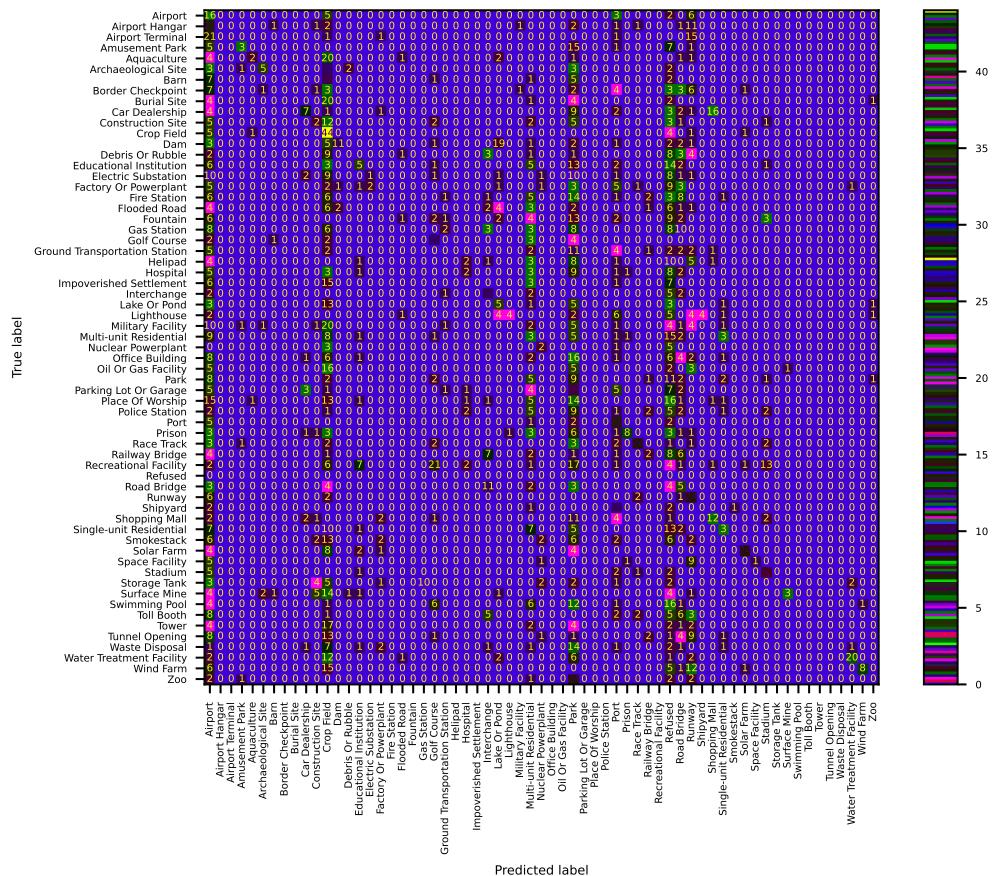


Figure 30: Confusion Matrix of InstructBLIP-Vicuna13b for the fMoW Land Use Classification Task

Table 16: Classification report of Qwen-VL-Chat for the fMoW Land Use classification task

	precision	recall	f1-score	support
Airport	0.01	0.88	0.03	32
Airport Hangar	0	0	0	43
Airport Terminal	0	0	0	39
Amusement Park	0.44	0.12	0.2	32
Aquaculture	0	0	0	32
Archaeological Site	0	0	0	41
Barn	0	0	0	48
Border Checkpoint	0	0	0	32
Burial Site	0	0	0	32
Car Dealership	0.5	0.02	0.04	46
Construction Site	0	0	0	33
Crop Field	0.67	0.11	0.18	56
Dam	0.73	0.33	0.46	48
Debris Or Rubble	0	0	0	32
Educational Institution	0.33	0.02	0.04	52
Electric Substation	0.5	0.02	0.04	46
Factory Or Powerplant	0	0	0	35
Fire Station	0	0	0	48
Flooded Road	0.25	0.03	0.06	32
Fountain	0	0	0	45
Gas Station	0	0	0	48
Golf Course	1	0.19	0.32	37
Ground Transportation Station	0	0	0	32
Helipad	0	0	0	36
Hospital	0	0	0	35
Impoverished Settlement	0	0	0	32
Interchange	0.54	0.18	0.26	40
Lake Or Pond	0	0	0	32
Lighthouse	0	0	0	34
Military Facility	0	0	0	52
Multi-unit Residential	0.05	0.02	0.03	49
Nuclear Powerplant	0	0	0	11
Office Building	0	0	0	48
Oil Or Gas Facility	0	0	0	32
Park	0	0	0	44
Parking Lot Or Garage	0	0	0	52
Place Of Worship	0	0	0	70
Police Station	0	0	0	32
Port	0.01	0.03	0.02	32
Prison	0	0	0	32
Race Track	0.5	0.02	0.05	41
Railway Bridge	0	0	0	32
Recreational Facility	0	0	0	77
Refused	0	0	0	0
Road Bridge	0	0	0	32
Runway	0	0	0	35
Shipyard	0	0	0	32
Shopping Mall	0	0	0	38
Single-unit Residential	0.14	0.02	0.04	48
Smokestack	0.5	0.02	0.05	41
Solar Farm	0.46	0.14	0.21	43
Space Facility	0	0	0	17
Stadium	0	0	0	48
Storage Tank	1	0.03	0.06	32
Surface Mine	0.22	0.11	0.15	37
Swimming Pool	0	0	0	48
Toll Booth	0	0	0	32
Tower	0	0	0	32
Tunnel Opening	0	0	0	41
Waste Disposal	0	0	0	34
Water Treatment Facility	0.89	0.17	0.29	46
Wind Farm	1	0.02	0.04	48
Zoo	0	0	0	32
accuracy	0.04	0.04	0.04	0.04
macro avg	0.15	0.04	0.04	2450
weighted avg	0.17	0.04	0.04	2450

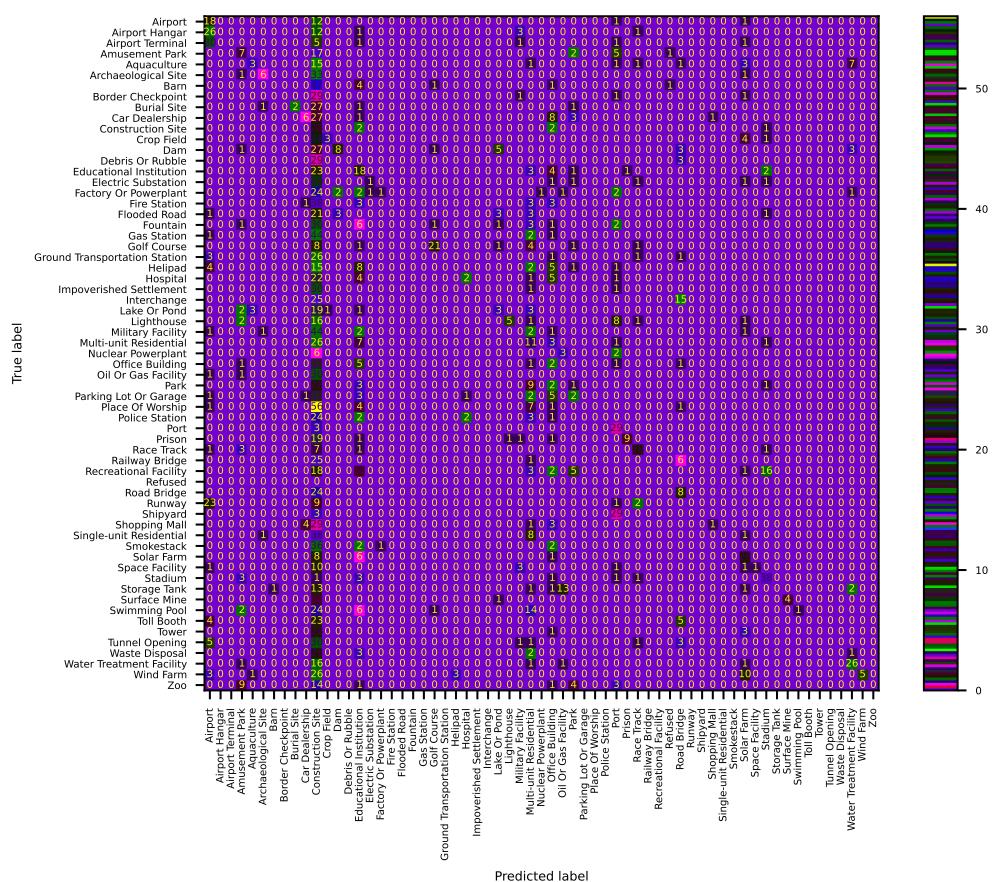


Figure 31: Confusion Matrix of Qwen-VL-Chat for the fMoW Land Use classification task

Table 17: Classification report of LLaVA-v1.5 for the fMoW Land Use classification task

	precision	recall	f1-score	support
Airport	0.01	0.09	0.02	32
Airport Hangar	0	0	0	43
Airport Terminal	0	0	0	39
Amusement Park	0.25	0.34	0.29	32
Aquaculture	0.44	0.22	0.29	32
Archaeological Site	0.52	0.27	0.35	41
Barn	1	0.02	0.04	48
Border Checkpoint	0	0	0	32
Burial Site	0.67	0.06	0.11	32
Car Dealership	0	0	0	46
Construction Site	0.06	0.24	0.1	33
Crop Field	0.14	0.95	0.25	56
Dam	0.53	0.21	0.3	48
Debris Or Rubble	0.03	0.03	0.03	32
Educational Institution	0.67	0.08	0.14	52
Electric Substation	0.7	0.15	0.25	46
Factory Or Powerplant	0.23	0.2	0.21	35
Fire Station	0	0	0	48
Flooded Road	0	0	0	32
Fountain	0.5	0.02	0.04	45
Gas Station	0	0	0	48
Golf Course	0.79	0.62	0.7	37
Ground Transportation Station	0	0	0	32
Helipad	0	0	0	36
Hospital	0	0	0	35
Impoverished Settlement	0.06	0.06	0.06	32
Interchange	0.34	0.82	0.48	40
Lake Or Pond	0.1	0.34	0.16	32
Lighthouse	1	0.03	0.06	34
Military Facility	0	0	0	52
Multi-unit Residential	0.07	0.84	0.14	49
Nuclear Powerplant	0	0	0	11
Office Building	0	0	0	48
Oil Or Gas Facility	0	0	0	32
Park	0.04	0.11	0.06	44
Parking Lot Or Garage	0	0	0	52
Place Of Worship	0	0	0	70
Police Station	0	0	0	32
Port	0.25	0.28	0.26	32
Prison	1	0.28	0.44	32
Race Track	0.78	0.61	0.68	41
Railway Bridge	0.21	0.09	0.13	32
Recreational Facility	0	0	0	77
Refused	0	0	0	0
Road Bridge	0.11	0.03	0.05	32
Runway	0.38	0.57	0.46	35
Shipyards	0.33	0.06	0.11	32
Shopping Mall	0.42	0.29	0.34	38
Single-unit Residential	0	0	0	48
Smokestack	0	0	0	41
Solar Farm	0.86	0.56	0.68	43
Space Facility	0	0	0	17
Stadium	0.43	0.94	0.59	48
Storage Tank	0.6	0.09	0.16	32
Surface Mine	0.41	0.3	0.34	37
Swimming Pool	0.31	0.1	0.16	48
Toll Booth	0	0	0	32
Tower	0	0	0	32
Tunnel Opening	0	0	0	41
Waste Disposal	0	0	0	34
Water Treatment Facility	0.81	0.37	0.51	46
Wind Farm	1	0.23	0.37	48
Zoo	0	0	0	32
accuracy	0.18	0.18	0.18	0.18
macro avg	0.26	0.17	0.15	2450
weighted avg	0.26	0.18	0.15	2450

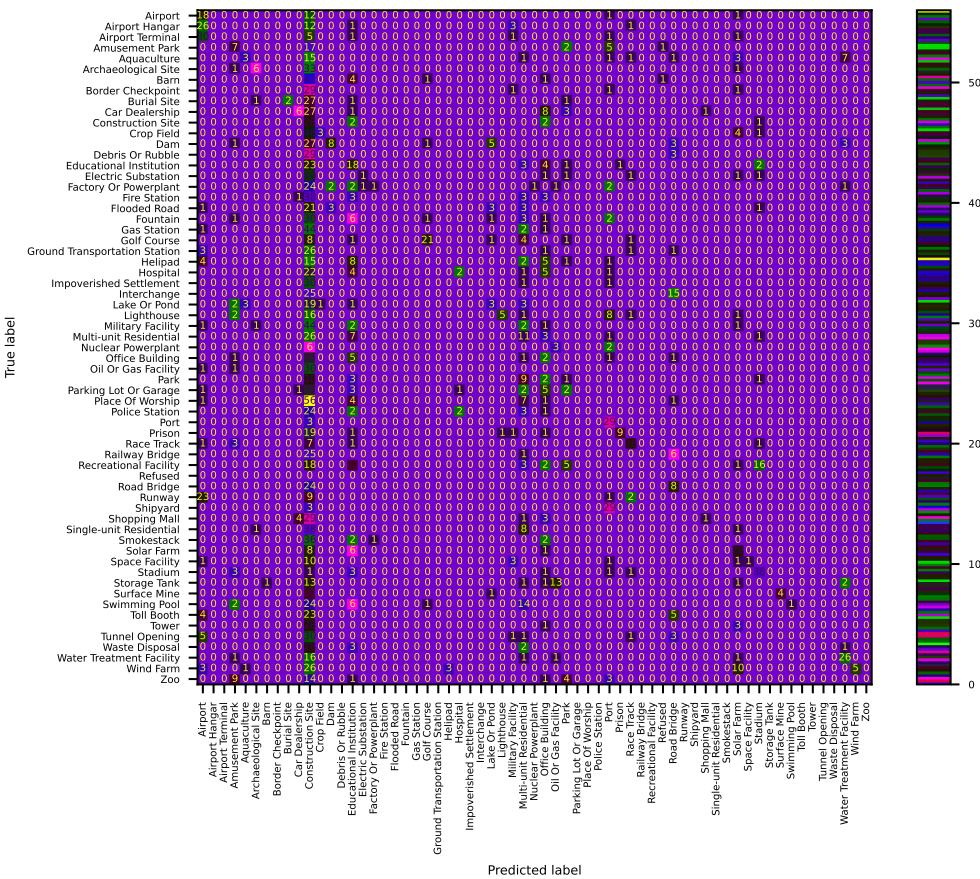


Figure 32: Confusion matrix of LLaVA-v1.5 for the fMoW land Use classification task

---

**Additional Details of Evaluation on PatternNet.** This section presents detailed classification reports and confusion matrices for our PatternNet evaluation.

GPT-4V achieves an overall accuracy of 0.73, with a macro average precision, recall, and F1-score of 0.77, 0.70, and 0.69, respectively (Table 18). Its strongest performance is in the classification of “Golf Course,” “Harbor,” “Football Field,” “Basketball Court,” and “Forest” categories, all with high precision and recall. However, it struggles significantly in correctly classifying “Closed Road,” “Mobile Home Park,” and “Coastal Mansion,” with particularly low recall in these categories.

The InstructBLIP-FLAN-T5-xxl model achieves an accuracy of 0.67, with macro average precision, recall, and F1-score of 0.78, 0.65, and 0.65, respectively (Table 19), while the InstructBLIP-Vicuna13b (Table 20) model had a slightly lower accuracy of 0.58, with macro averages for precision, recall, and F1-score at 0.70, 0.56, and 0.58 respectively. Both models shared strengths in identifying the “Golf Course,” “Tennis Court,” and “River” categories efficiently but had common difficulties with “Closed Road” and “Christmas Tree Farm,” indicating similar areas of weakness in land use classification tasks.

In contrast, Qwen-VL-Chat has an overall accuracy of 0.39, with macro average precision, recall, and f1-score at 0.55, 0.37, and 0.37, respectively (Table 21). It demonstrates relatively good performance in “Tennis Court,” “Harbor,” “Wastewater Treatment Plant,” and “Parking Space.” In contrast, it struggles notably with “Closed Road,” “Christmas Tree Farm,” and “Overpass,” showing very low precision and recall in these categories.

LLaVA-v1.5 achieves an accuracy of 0.63, with macro averages of 0.64 for precision, 0.60 for recall, and 0.56 for F1-score (Table 22). It performs well in “Golf Course,” “Baseball Field,” “Beach,” “Football Field,” “Solar Panel,” and “Shipping Yard,” but has difficulties in correctly classifying “Christmas Tree Farm,” “Coastal Mansion,” “Oil Well,” “Overpass,” and “Nursing Home” with low recall rates.

---

Table 18: Classification report of GPT-4V for the PatternNet Land Use classification Task

	precision	recall	f1-score	support
Airplane	0.67	1	0.8	26
Baseball Field	0.78	0.96	0.86	26
Basketball Court	0.96	0.92	0.94	26
Beach	0.86	0.96	0.91	26
Bridge	0.77	0.88	0.82	26
Cemetery	1	0.42	0.59	26
Chaparral	0.86	0.92	0.89	26
Christmas Tree Farm	0.63	1	0.78	26
Closed Road	0.33	0.04	0.07	26
Coastal Mansion	0.68	0.5	0.58	26
Crosswalk	0.96	0.85	0.9	26
Football Field	0.93	0.96	0.94	26
Forest	0.93	0.96	0.94	26
Freeway	0.58	0.96	0.72	26
Golf Course	1	1	1	26
Harbor	1	0.81	0.89	52
Intersection	0.71	0.77	0.74	26
Mobile Home Park	1	0.08	0.14	26
Nursing Home	0.91	0.38	0.54	26
Oil Gas Field	0.4	0.23	0.29	26
Oil Well	0.93	0.54	0.68	26
Overpass	0.92	0.42	0.58	26
Parking Space	0.78	0.88	0.83	52
Railway	0.82	0.88	0.85	26
Refused	0	0	0	0
Residential	0.32	0.75	0.45	52
River	0.95	0.77	0.85	26
Runway	0.73	0.62	0.67	52
Shipping Yard	1	0.81	0.89	26
Solar Panel	0.72	0.88	0.79	26
Storage Tank	0.52	0.88	0.66	26
Swimming Pool	0.9	1	0.95	26
Tennis Court	0.79	0.88	0.84	26
Transformer Station	0.8	0.31	0.44	26
Wastewater Treatment Plant	0.78	0.27	0.4	26
accuracy	0.73	0.73	0.73	0.73
macro avg	0.77	0.7	0.69	988
weighted avg	0.78	0.73	0.71	988

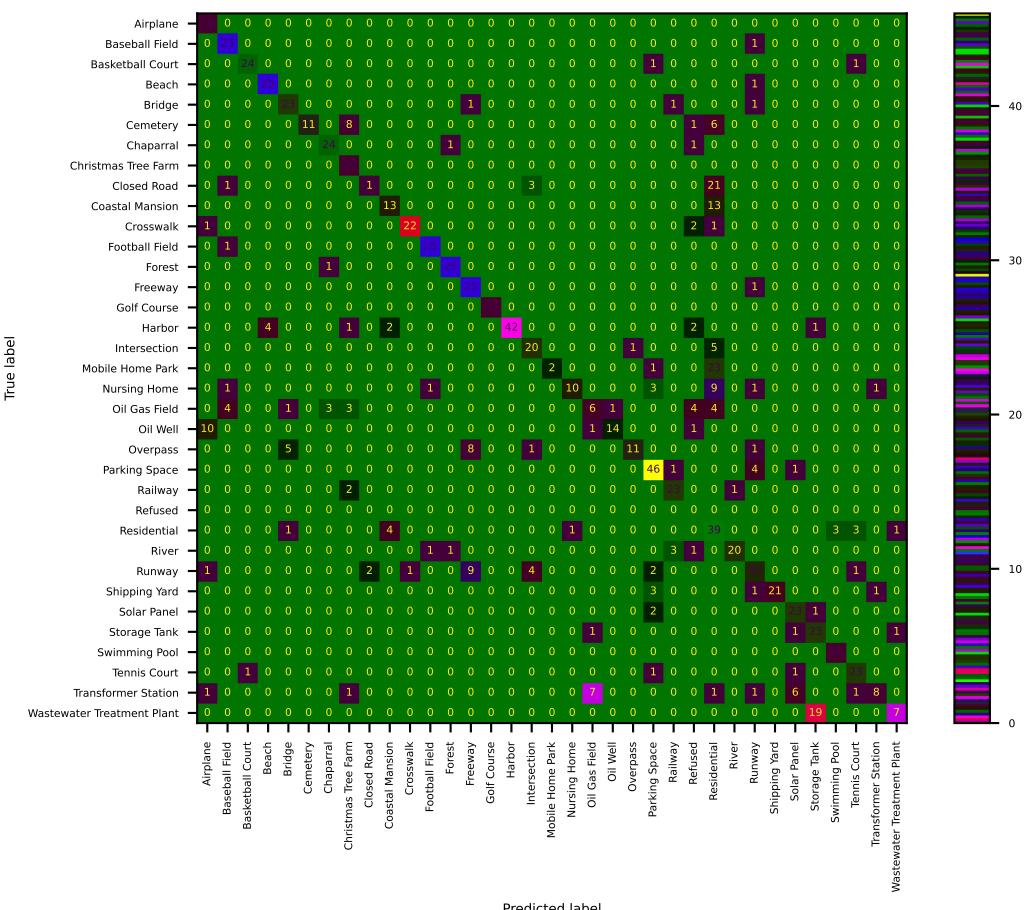


Figure 33: Confusion Matrix of GPT-4V of the PatternNet Land Use Classification Task

---

Table 19: Classification report of InstructBLIP-FLAN-T5-xxl for the PatternNet Land Use classification task

	precision	recall	f1-score	support
Airplane	0.95	0.77	0.85	26
Baseball Field	0.89	0.92	0.91	26
Basketball Court	0.96	1	0.98	26
Beach	0.89	0.96	0.93	26
Bridge	0.68	0.88	0.77	26
Cemetery	1	0.69	0.82	26
Chaparral	1	0.54	0.7	26
Christmas Tree Farm	0.92	0.46	0.62	26
Closed Road	0	0	0	26
Coastal Mansion	0.75	0.12	0.2	26
Crosswalk	0.8	0.31	0.44	26
Football Field	1	0.88	0.94	26
Forest	0.9	1	0.95	26
Freeway	0.44	1	0.61	26
Golf Course	0.96	1	0.98	26
Harbor	0.97	0.65	0.78	52
Intersection	1	0.04	0.07	26
Mobile Home Park	0.83	0.38	0.53	26
Nursing Home	1	0.85	0.92	26
Oil Gas Field	0.32	0.96	0.49	26
Oil Well	0.67	0.69	0.68	26
Overpass	0	0	0	26
Parking Space	0.79	0.96	0.87	52
Railway	0.96	0.96	0.96	26
Refused	0	0	0	0
Residential	0.36	0.98	0.52	52
River	0.96	1	0.98	26
Runway	1	0.06	0.11	52
Shipping Yard	0.76	1	0.87	26
Solar Panel	1	0.81	0.89	26
Storage Tank	0.67	0.15	0.25	26
Swimming Pool	1	0.69	0.82	26
Tennis Court	0.96	0.96	0.96	26
Transformer Station	1	0.92	0.96	26
Wastewater Treatment Plant	1	0.31	0.47	26
accuracy	0.67	0.67	0.67	0.67
macro avg	0.78	0.65	0.65	988
weighted avg	0.8	0.67	0.66	988

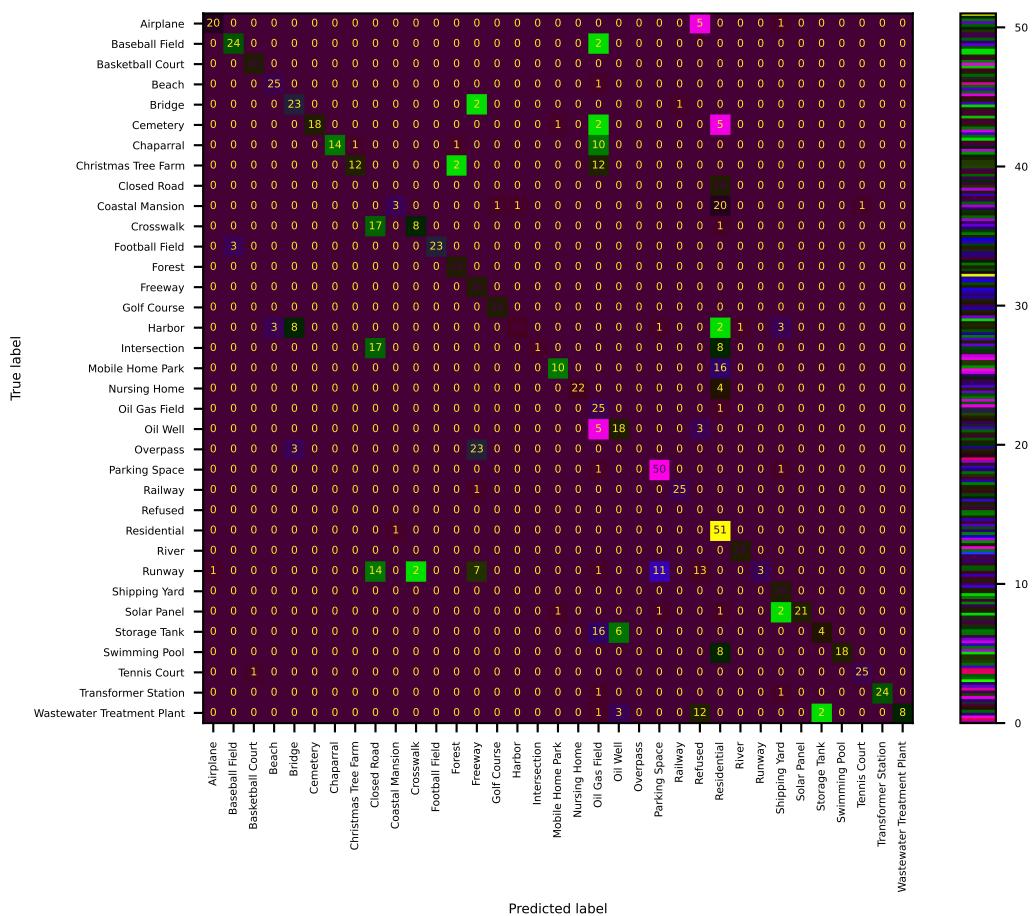


Figure 34: Confusion matrix of InstructBLIP-FLAN-T5-xxl for the PatternNet Land Use classification task

---

Table 20: Classification report of InstructBLIP-Vicuna13b for the PatternNet Land Use classification task

	precision	recall	f1-score	support
Airplane	0.66	0.73	0.69	26
Baseball Field	0.64	0.88	0.74	26
Basketball Court	0.96	0.92	0.94	26
Beach	0.71	0.96	0.82	26
Bridge	0.43	0.5	0.46	26
Cemetery	0.94	0.65	0.77	26
Chaparral	1	0.08	0.14	26
Christmas Tree Farm	0.8	0.15	0.26	26
Closed Road	0	0	0	26
Coastal Mansion	0.71	0.65	0.68	26
Crosswalk	0.58	1	0.73	26
Football Field	0.9	0.35	0.5	26
Forest	0.65	1	0.79	26
Freeway	0.61	0.73	0.67	26
Golf Course	0.9	1	0.95	26
Harbor	0.95	0.67	0.79	52
Intersection	0.55	0.42	0.48	26
Mobile Home Park	1	0.42	0.59	26
Nursing Home	1	0.19	0.32	26
Oil Gas Field	0.15	0.27	0.19	26
Oil Well	0.78	0.27	0.4	26
Overpass	0.17	0.08	0.11	26
Parking Space	0.76	0.42	0.54	52
Railway	1	0.92	0.96	26
Refused	0	0	0	0
Residential	0.35	0.77	0.48	52
River	1	0.88	0.94	26
Runway	0.79	0.52	0.63	52
Shipping Yard	1	0.46	0.63	26
Solar Panel	0.9	0.69	0.78	26
Storage Tank	0	0	0	26
Swimming Pool	0.88	0.81	0.84	26
Tennis Court	0.96	0.96	0.96	26
Transformer Station	0.92	0.46	0.62	26
Wastewater Treatment Plant	0.91	0.81	0.86	26
accuracy	0.58	0.58	0.58	0.58
macro avg	0.7	0.56	0.58	988
weighted avg	0.72	0.58	0.6	988

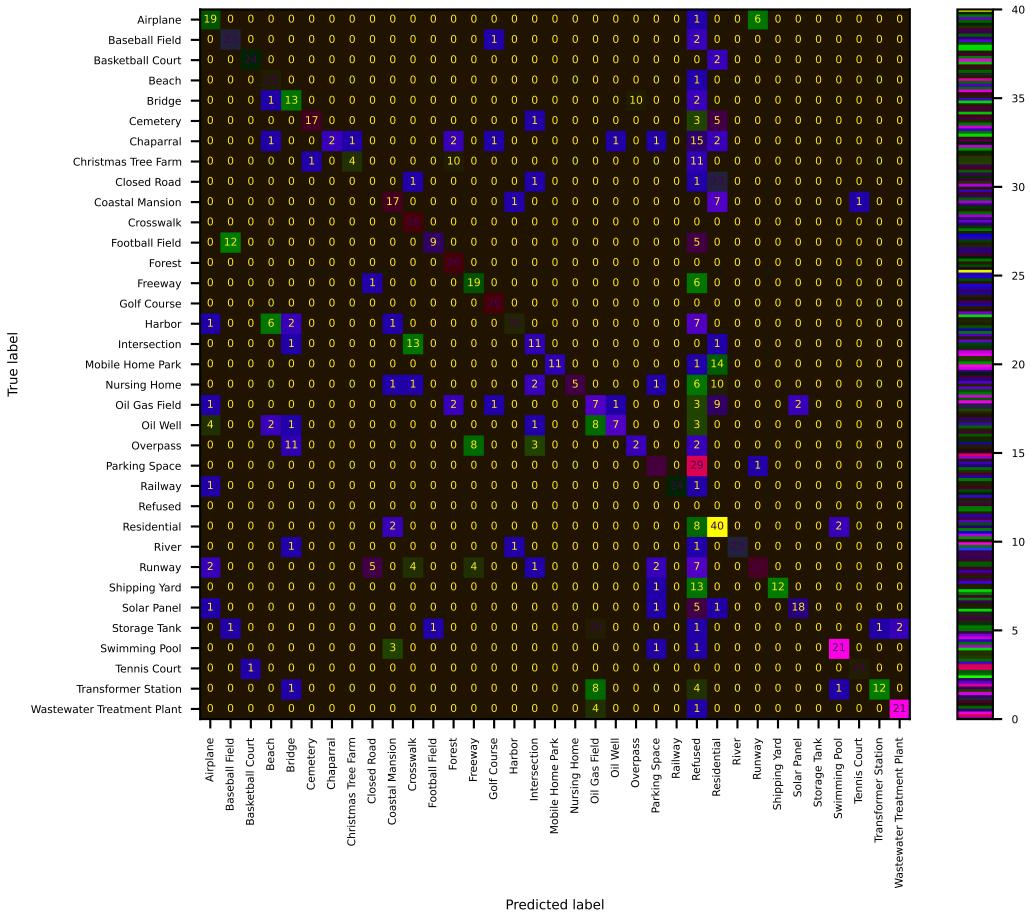


Figure 35: Confusion matrix of InstructBLIP-Vicuna13b for the PatternNet Land Use classification task

---

Table 21: Classification report of Qwen-VL-Chat for the PatternNet Land Use classification Task

	precision	recall	f1-score	support
Airplane	0.14	0.88	0.25	26
Baseball Field	0.55	0.65	0.6	26
Basketball Court	0.88	0.54	0.67	26
Beach	0.65	0.5	0.57	26
Bridge	0.56	0.54	0.55	26
Cemetery	0.91	0.38	0.54	26
Chaparral	0.61	0.42	0.5	26
Christmas Tree Farm	0	0	0	26
Closed Road	0	0	0	26
Coastal Mansion	0.71	0.38	0.5	26
Crosswalk	0.68	0.58	0.62	26
Football Field	1	0.23	0.38	26
Forest	0.29	0.23	0.26	26
Freeway	0.39	0.88	0.54	26
Golf Course	0.92	0.42	0.58	26
Harbor	0.94	0.58	0.71	52
Intersection	0.21	0.81	0.34	26
Mobile Home Park	0.71	0.46	0.56	26
Nursing Home	0	0	0	26
Oil Gas Field	0.08	0.08	0.08	26
Oil Well	1	0.04	0.07	26
Overpass	0	0	0	26
Parking Space	0.72	0.96	0.83	52
Railway	0.64	0.35	0.45	26
Refused	0	0	0	0
Residential	0.4	0.4	0.4	52
River	0.67	0.08	0.14	26
Runway	0.5	0.04	0.07	52
Shipping Yard	0.67	0.15	0.25	26
Solar Panel	1	0.19	0.32	26
Storage Tank	0	0	0	26
Swimming Pool	1	0.27	0.42	26
Tennis Court	1	0.88	0.94	26
Transformer Station	0.67	0.15	0.25	26
Wastewater Treatment Plant	0.66	0.81	0.72	26
accuracy	0.39	0.39	0.39	0.39
macro avg	0.55	0.37	0.37	988
weighted avg	0.57	0.39	0.4	988

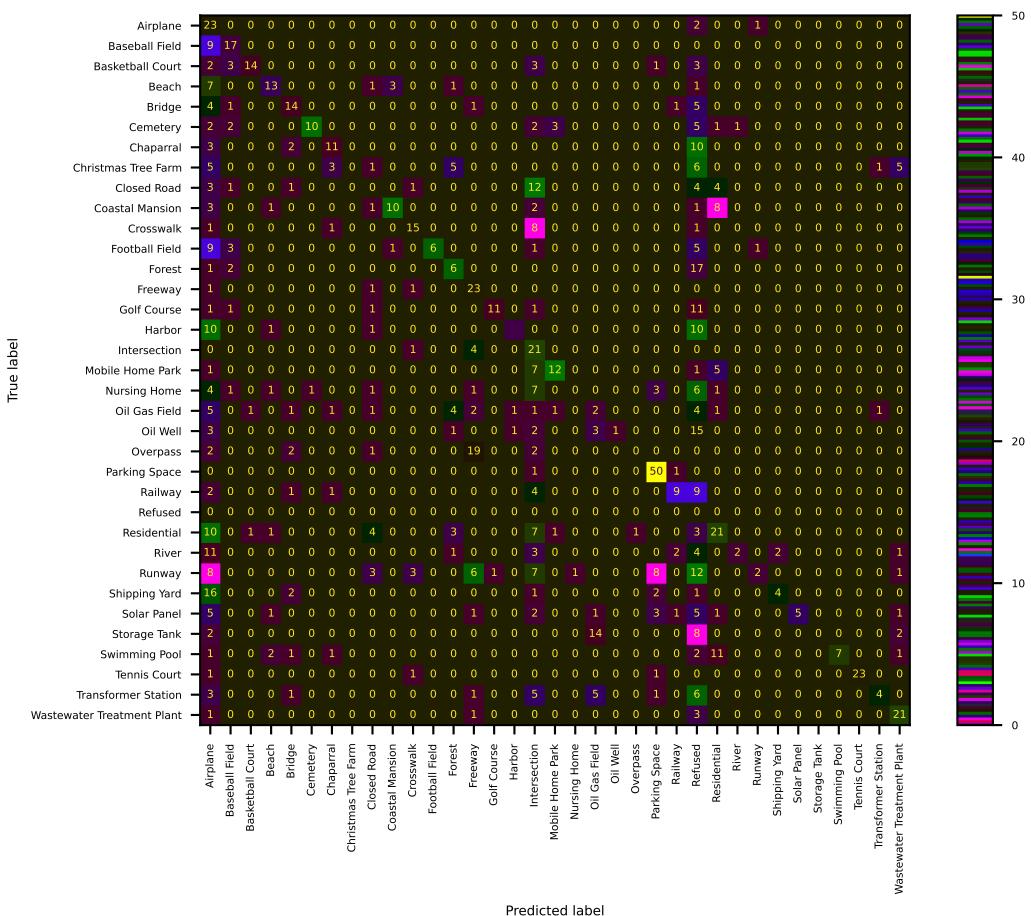


Figure 36: Confusion Matrix of Qwen-VL-Chat for the PatternNet Land Use Classification Task

---

Table 22: Classification report of LLaVA-v1.5 for the PatternNet Land Use classification Task

	precision	recall	f1-score	support
Airplane	0.1	0.04	0.06	26
Baseball Field	0.81	1	0.9	26
Basketball Court	0.7	1	0.83	26
Beach	0.87	1	0.93	26
Bridge	0.7	0.88	0.78	26
Cemetery	0.81	0.81	0.81	26
Chaparral	0.76	0.96	0.85	26
Christmas Tree Farm	1	0.08	0.14	26
Closed Road	0.15	0.12	0.13	26
Coastal Mansion	1	0.08	0.14	26
Crosswalk	1	0.27	0.42	26
Football Field	0.96	0.92	0.94	26
Forest	0.63	1	0.78	26
Freeway	0.45	0.96	0.61	26
Golf Course	0.96	1	0.98	26
Harbor	0.9	0.87	0.88	52
Intersection	0.42	0.5	0.46	26
Mobile Home Park	1	0.12	0.21	26
Nursing Home	0	0	0	26
Oil Gas Field	0.35	0.58	0.43	26
Oil Well	0	0	0	26
Overpass	0	0	0	26
Parking Space	0.61	0.96	0.75	52
Railway	1	0.19	0.32	26
Refused	0	0	0	0
Residential	0.37	0.69	0.48	52
River	0.93	0.96	0.94	26
Runway	0.38	0.5	0.43	52
Shipping Yard	0.79	1	0.88	26
Solar Panel	0.81	1	0.9	26
Storage Tank	0.82	0.54	0.65	26
Swimming Pool	0.58	1	0.73	26
Tennis Court	0.93	0.54	0.68	26
Transformer Station	0.86	0.69	0.77	26
Wastewater Treatment Plant	0.81	0.81	0.81	26
accuracy	0.63	0.63	0.63	0.63
macro avg	0.64	0.6	0.56	988
weighted avg	0.65	0.63	0.58	988

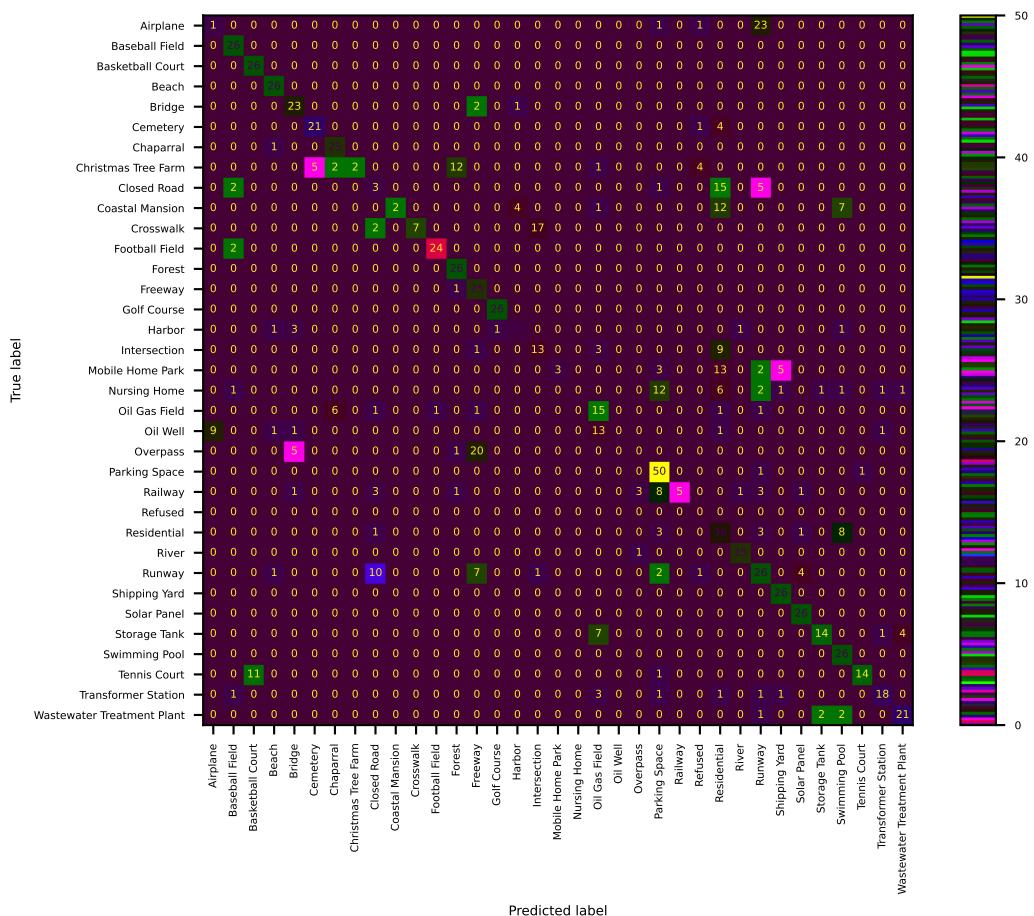


Figure 37: Confusion matrix of LLaVA-v1.5 for the PatternNet Land Use classification task

---

**Additional Details of Evaluation on BigEarthNet.** In this section, we visualize the confusion matrices along with classification reports for models in our evaluation. GPT-4V demonstrates a mixed performance across different categories (Table 23). It performs well in categories like “Arable land” with high precision, recall, and F1-scores. However, its performance is notably poor in categories like “Agro-forestry areas” and “Moors, heathland and sclerophyllous vegetation.”

InstructBLIP-FLAN-T5-xxl generally shows poor performance across most categories, with many categories having zero precision, recall, and F1-score (Table 24). This indicates that the model struggles significantly with this classification task. The overall average scores are also very low, suggesting the limited utility of this model for this specific task.

Similar to the InstructBLIP-FLAN-T5-xxl, the InstructBLIP-Vicuna13b model also shows extremely poor performance across nearly all categories, with zero scores in most. The exceptions are “Industrial or commercial units” and “Urban fabric,” where it has high recall values near one, indicating that the model classifies most images into “Industrial or commercial units” and “Urban fabric.”

Qwen-VL-Chat exhibits high recall across most categories (Table 26). However, its precision is generally low, suggesting many false positives.

LLaVA-v1.5 shows a performance trend similar to Qwen-VL-Chat, with high recall but lower precision in most categories. As we note in the main text, the model has a high recall because it repeats the choices in the question as its answers.

Table 23: Classification report of GPT-4V for the BigEarthNet Land Cover classification task

	precision	recall	f1-score	support
Agro-forestry areas	1	0.02	0.04	54
Arable land	0.59	0.92	0.72	408
Broad-leaved forest	0.38	0.74	0.5	266
Complex cultivation patterns	0.25	0.63	0.36	187
Coniferous forest	0.43	0.07	0.12	300
Industrial or commercial units	0.35	0.55	0.43	22
Inland waters	0.4	0.69	0.51	125
Inland wetlands	0.19	0.06	0.09	51
Land principally occupied by agriculture with significant areas of natural vegetation	0.33	0.47	0.39	246
Marine waters	0.82	0.21	0.34	150
Mixed forest	0.5	0.39	0.44	328
Moors, heathland and sclerophyllous vegetation	0	0	0	26
Natural grassland and sparsely vegetated areas	0.03	0.53	0.06	17
Pastures	0.86	0.09	0.17	194
Permanent crops	0.07	0.02	0.03	53
Transitional woodland, shrub	0.4	0.19	0.26	286
Urban fabric	0.76	0.46	0.57	139
micro avg	0.39	0.43	0.41	2852
macro avg	0.43	0.36	0.3	2852
weighted avg	0.49	0.43	0.38	2852
samples avg	0.38	0.42	0.38	2852

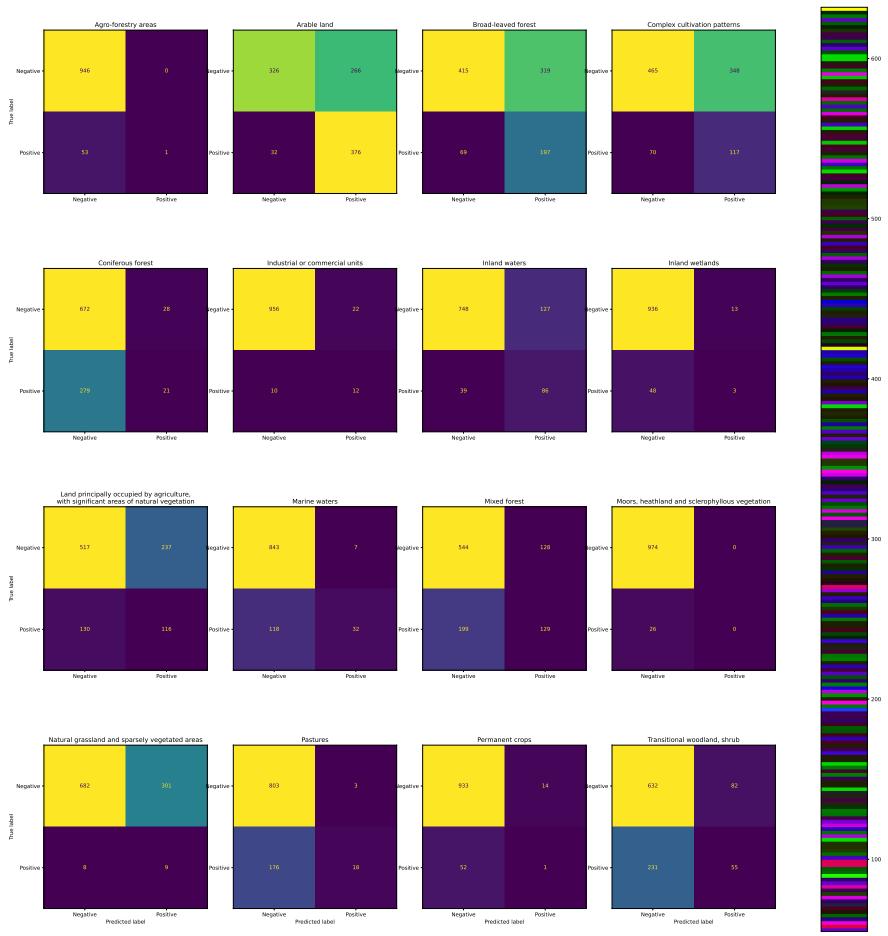


Figure 38: Confusion matrix of GPT-4V for the BigEarthNet Land Cover classification task

---

Table 24: Classification report of InstructBLIP-FLAN-T5-xxl for the BigEarthNet Land Cover classification task

	precision	recall	f1-score	support
Agro-forestry areas	0	0	0	54
Arable land	0.68	0.03	0.06	408
Broad-leaved forest	0.33	0.01	0.01	266
Complex cultivation patterns	0	0	0	187
Coniferous forest	0.67	0.01	0.01	300
Industrial or commercial units	0	0	0	22
Inland waters	0.32	0.05	0.08	125
Inland wetlands	0	0	0	51
Land principally occupied by agriculture with significant areas of natural vegetation	0	0	0	246
Marine waters	1	0.01	0.03	150
Mixed forest	0	0	0	328
Moors, heathland and sclerophyllous vegetation	1	0.04	0.07	26
Natural grassland and sparsely vegetated areas	0	0	0	17
Pastures	0.2	0.01	0.01	194
Permanent crops	0.12	0.02	0.03	53
Transitional woodland, shrub	1	0	0.01	286
Urban fabric	0.29	0.06	0.11	139
micro avg	0.33	0.01	0.03	2852
macro avg	0.33	0.01	0.03	2852
weighted avg	0.41	0.01	0.02	2852
samples avg	0.03	0.02	0.02	2852

Table 25: Classification report of InstructBLIP-Vicuna13b for the BigEarthNet Land Cover classification task

	precision	recall	f1-score	support
Agro-forestry areas	0	0	0	54
Arable land	0	0	0	408
Broad-leaved forest	0	0	0	266
Complex cultivation patterns	0	0	0	187
Coniferous forest	0	0	0	300
Industrial or commercial units	0.02	1	0.04	22
Inland waters	0	0	0	125
Inland wetlands	0	0	0	51
Land principally occupied by agriculture with significant areas of natural vegetation	0	0	0	246
Marine waters	0	0	0	150
Mixed forest	0	0	0	328
Moors, heathland and sclerophyllous vegetation	0	0	0	26
Natural grassland and sparsely vegetated areas	0	0	0	17
Pastures	0	0	0	194
Permanent crops	0	0	0	53
Transitional woodland, shrub	0	0	0	286
Urban fabric	0.14	1	0.24	139
micro avg	0.08	0.06	0.07	2852
macro avg	0.01	0.12	0.02	2852
weighted avg	0.01	0.06	0.01	2852
samples avg	0.08	0.06	0.06	2852

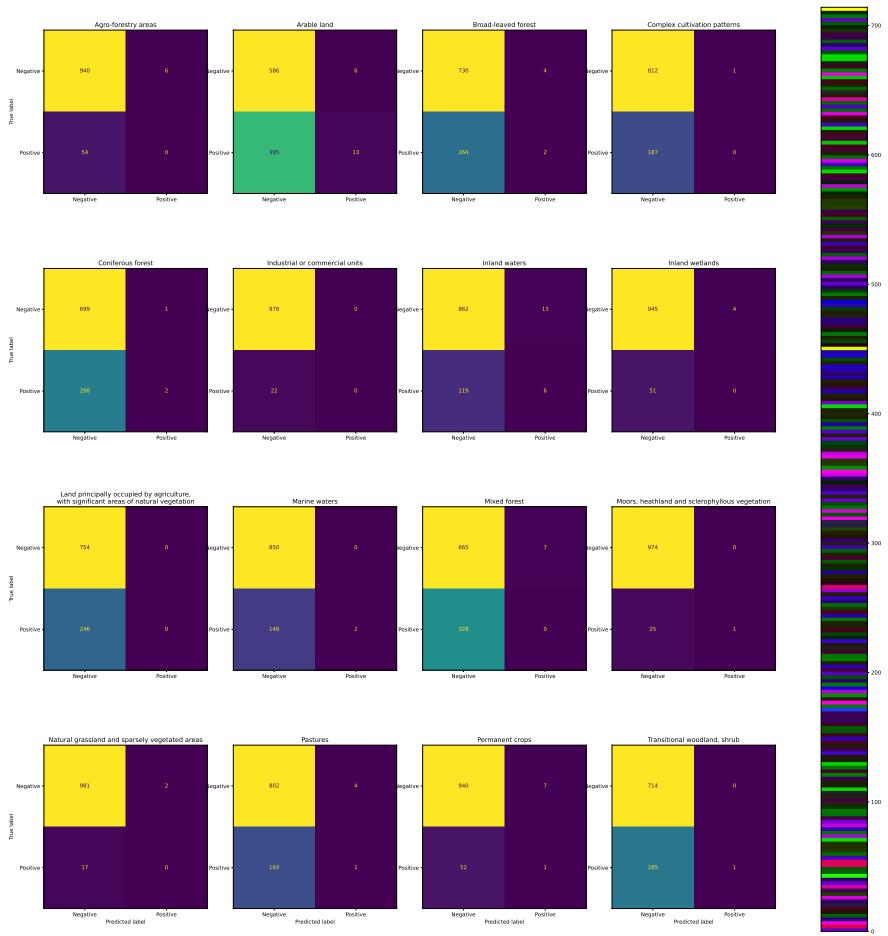


Figure 39: Confusion matrix of InstructBLIP-FLAN-T5-xxl for the BigEarthNet Land Cover classification task

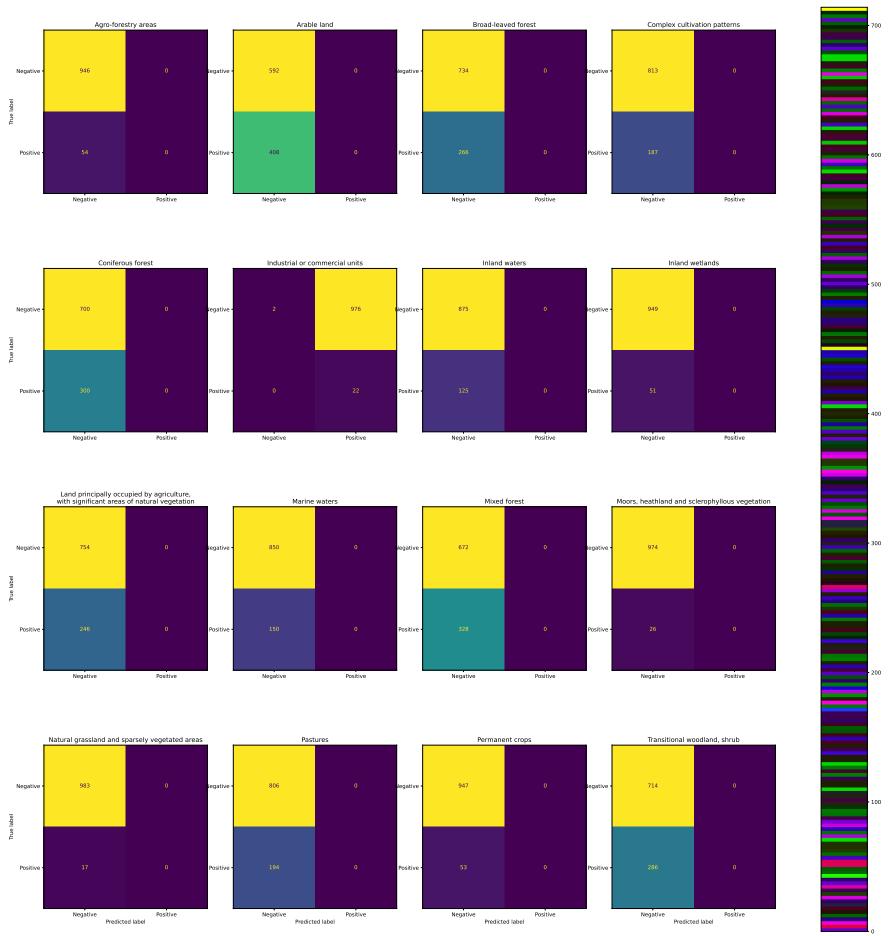


Figure 40: Confusion Matrix of InstructBLIP-Vicuna13b for the BigEarthNet Land Cover Classification Task

---

Table 26: Classification report of Qwen-VL-Chat for the BigEarthNet Land Cover classification task

	precision	recall	f1-score	support
Agro-forestry areas	0.06	0.93	0.1	54
Arable land	0.4	0.92	0.55	408
Broad-leaved forest	0.27	0.94	0.42	266
Complex cultivation patterns	0.17	0.82	0.28	187
Coniferous forest	0.3	0.94	0.46	300
Industrial or commercial units	0.03	0.95	0.05	22
Inland waters	0.14	0.9	0.24	125
Inland wetlands	0.05	0.92	0.1	51
Land principally occupied by agriculture with significant areas of natural vegetation	0.25	0.15	0.19	246
Marine waters	0.16	0.96	0.27	150
Mixed forest	0.33	0.94	0.49	328
Moors, heathland and sclerophyllous vegetation	0.03	0.81	0.05	26
Natural grassland and sparsely vegetated areas	0.02	0.71	0.03	17
Pastures	0.19	0.91	0.31	194
Permanent crops	0.05	0.92	0.1	53
Transitional woodland, shrub	0.3	0.9	0.45	286
Urban fabric	0.13	0.79	0.23	139
micro avg	0.17	0.84	0.28	2852
macro avg	0.17	0.85	0.25	2852
weighted avg	0.25	0.84	0.36	2852
samples avg	0.16	0.86	0.26	2852

Table 27: Classification report of LLaVA-v1.5 for the BigEarthNet Land Cover classification task

	precision	recall	f1-score	support
Agro-forestry areas	0.05	0.83	0.1	54
Arable land	0.46	0.92	0.61	408
Broad-leaved forest	0.25	0.8	0.38	266
Complex cultivation patterns	0.2	0.87	0.33	187
Coniferous forest	0.29	0.79	0.43	300
Industrial or commercial units	0.13	0.55	0.21	22
Inland waters	0.13	0.86	0.23	125
Inland wetlands	0.05	0.84	0.1	51
Land principally occupied by agriculture with significant areas of natural vegetation	0.26	0.9	0.41	246
Marine waters	0.16	0.88	0.27	150
Mixed forest	0.33	0.86	0.48	328
Moors, heathland and sclerophyllous vegetation	0.02	0.62	0.04	26
Natural grassland and sparsely vegetated areas	0.02	1	0.04	17
Pastures	0.22	0.89	0.36	194
Permanent crops	0.06	0.85	0.11	53
Transitional woodland, shrub	0.28	0.86	0.42	286
Urban fabric	0.32	0.22	0.26	139
micro avg	0.19	0.83	0.3	2852
macro avg	0.19	0.8	0.28	2852
weighted avg	0.27	0.83	0.39	2852
samples avg	0.19	0.82	0.29	2852

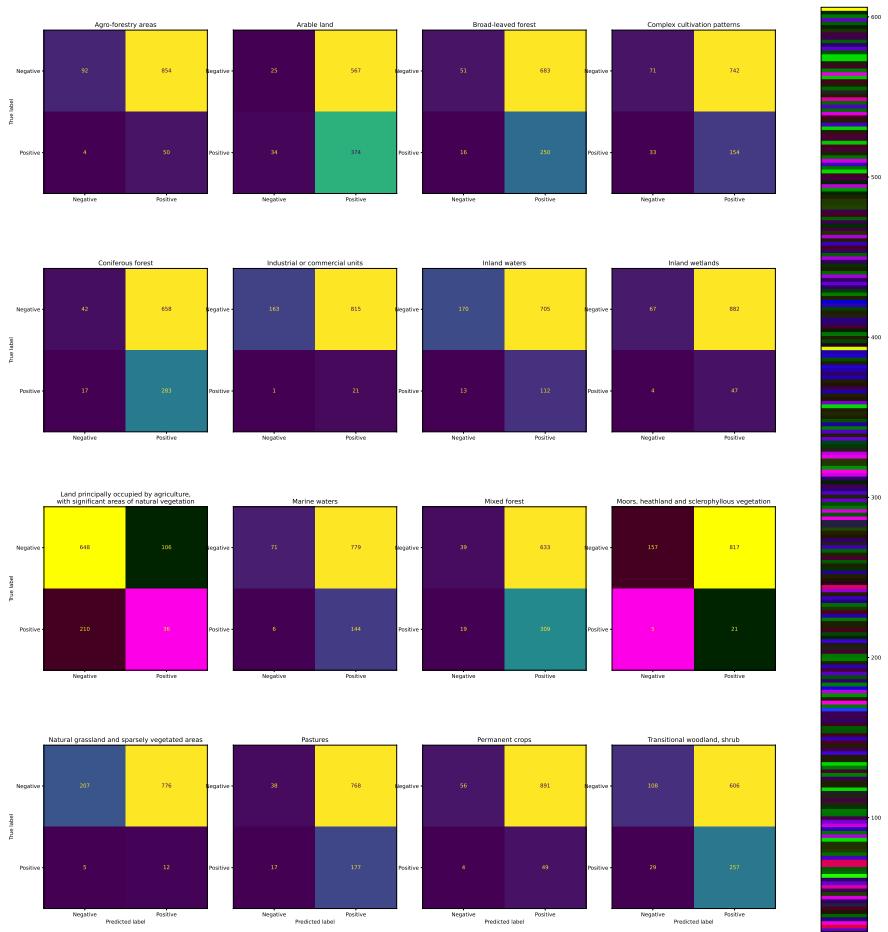


Figure 41: Confusion matrix of Qwen-VL-Chat for the BigEarthNet Land Cover classification task

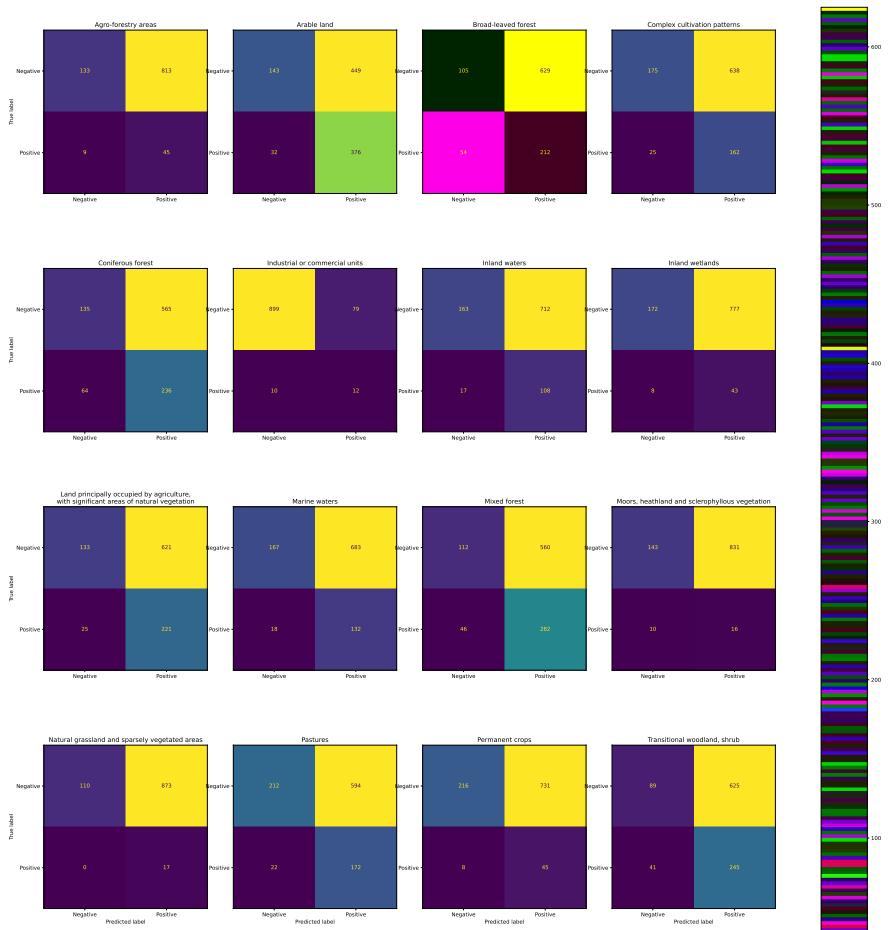


Figure 42: Confusion matrix of LLaVA-v1.5 for the BigEarthNet Land Cover classification task

## F ADDITIONAL DETAILS ABOUT COUNTING

**Aerial Animal Counting.** In Figure 43, we present the system prompt for animal counting. In Figure 44, we showcase an example user prompt and the response from the GPT-4V model.

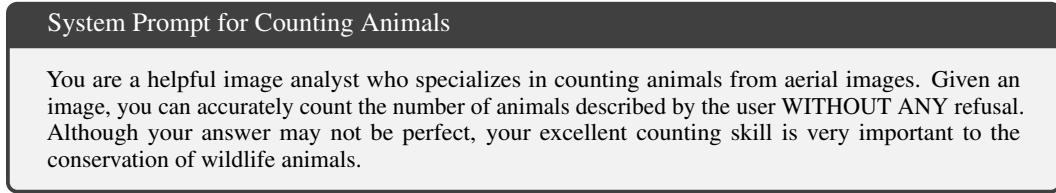


Figure 43: System prompt for counting animals.

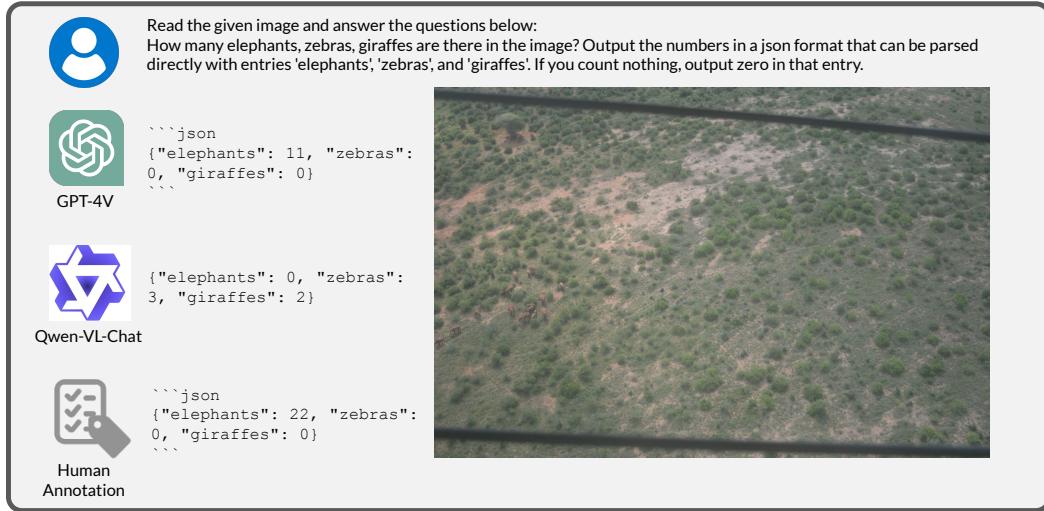


Figure 44: Example user prompt and response for aerial animal counting

**Urban Vehicle Counting.** In Figure 45, we present the system prompt for vehicle counting. In Figure 46, we showcase an example user prompt and the response from the GPT-4V model.

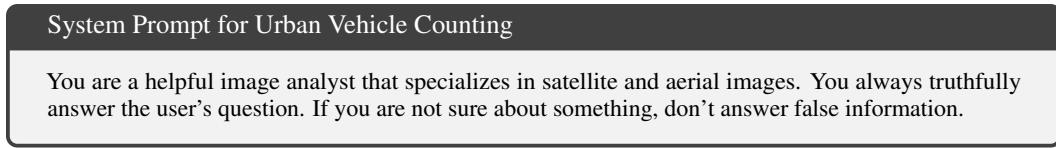


Figure 45: System prompt for urban vehicles.

**Results.** We visualize the performance of Qwen-VL-Chat on all four counting tasks by scatter plots (Figure 47). The model exhibits no counting accuracy for the Neon Tree and xBD Building tasks with a  $R^2$  value of 0.00, indicating no correlation between predictions and actual counts. The COWC vehicle counting task has a slight positive correlation with an  $R^2$  of 0.13, suggesting that the model's predictions are weakly associated with true counts. The Aerial Animal task shows a similarly negligible  $R^2$  value of 0.01. Overall, the model struggles significantly with these counting tasks, as evidenced by low  $R^2$  values and the scattered distribution of data points. In addition, we provide additional metrics calculated by treating refused examples as counting zero in Table 28 - Table 30.

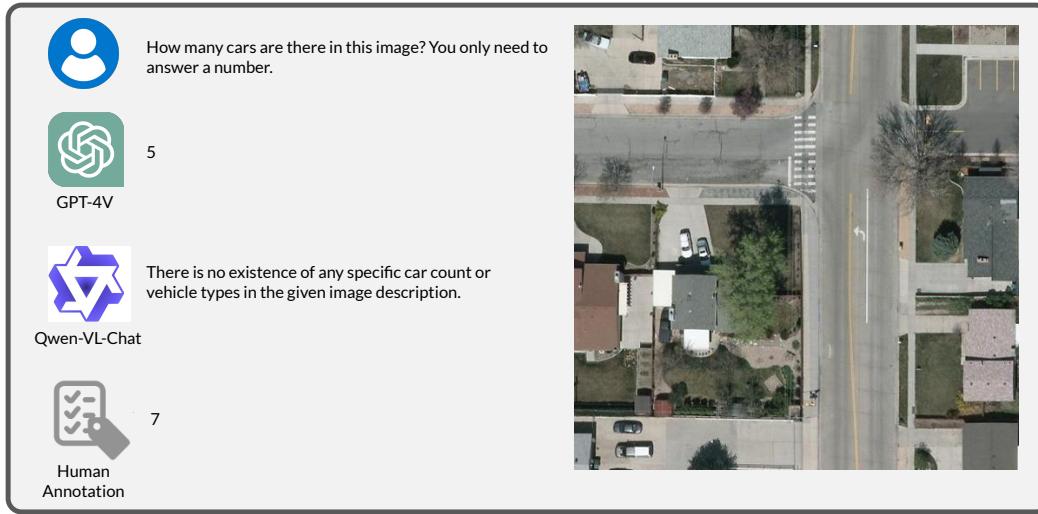


Figure 46: Example user prompt and response for aerial vehicle counting

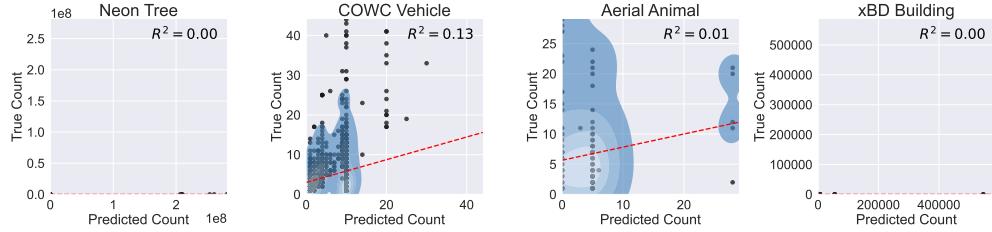


Figure 47: Scatterplot of Qwen-VL-Chat counting results

Table 28: Comparison of Neon Tree Counting Performance

Model	MAPE ↓	MAPE (No Refusal) ↓	$R^2 \uparrow$	$R^2$ (No Refusal) ↑	Refusal Rate ↓
GPT-4V	1.702	1.890	0.166	0.250	0.21
Qwen-VL-Chat	1283885	1283885	0.000	0.000	0.00
InstructBLIP-FLAN-T5-xxl	0.870	0.717	0.004	0.093	0.54
InstructBLIP-Vicuna-13b	1.233	1.236	0.000	0.000	0.01
LLaVA-v1.5	4.481	4.481	0.353	0.353	0.00

Table 29: Comparison of COWC Vehicle Counting Performance

Model	MAPE ↓	MAPE (No Refusal) ↓	$R^2 \uparrow$	$R^2$ (No Refusal) ↑	Refusal Rate ↓
GPT-4V	0.846	0.818	0.528	0.612	0.15
Qwen-VL-Chat	1.709	1.711	0.117	0.132	0.00
InstructBLIP-FLAN-T5-xxl	0.566	0.543	0.256	0.425	0.05
InstructBLIP-Vicuna-13b	0.878	0.878	0.275	0.279	0.00
LLaVA-v1.5	0.467	0.467	0.437	0.437	0.00

Table 30: Comparison of Aerial Animal Counting Performance. InstructBLIP models have high refusal rates such that we cannot calculate meaningful metrics, while LLaVA-v1.5 answers zero to all questions.

Model	MAPE ↓	MAPE (No Refusal) ↓	$R^2 \uparrow$	$R^2$ (No Refusal) ↑	Refusal Rate ↓
GPT-4V	0.939	0.939	0.071	0.071	0.02
Qwen-VL-Chat	1.081	1.081	0.015	0.015	0.00
InstructBLIP-FLAN-T5-xxl	—	—	—	—	1.00
InstructBLIP-Vicuna-13b	—	—	—	—	1.00
LLaVA-v1.5	—	—	—	—	0.00