# Benchmarking Vision-Language Models for Object Detection in Satellite Imagery

**Sustainability Lab** — AI and Sensing for Sustainability

**Sustainability Lab**  **IIT Gandhinagar**

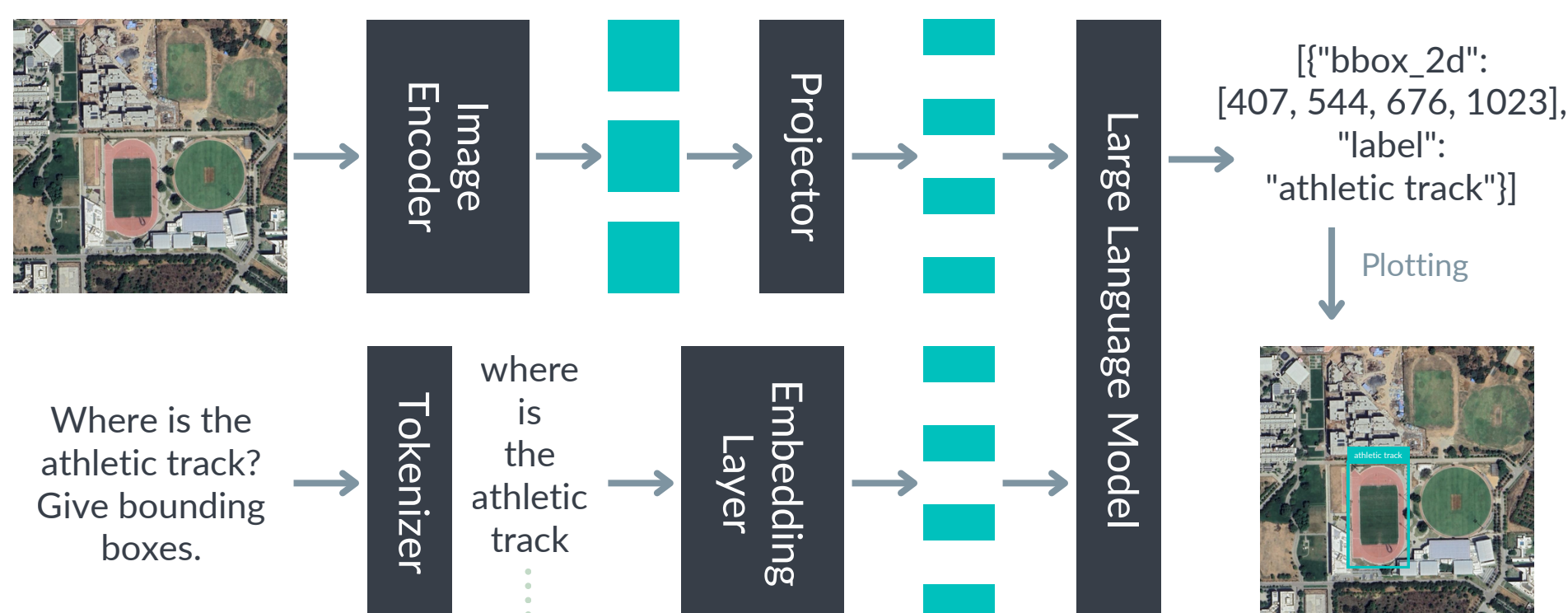| Devansh Lodha | Nupoor Assudani | Nipun Batra |
|---|---|---|
| 23110091@iitgn.ac.in | 23110224@iitgn.ac.in | nipun.batra@iitgn.ac.in |

## Motivation

**Problem:** Air pollution is a major issue in India (1.7M deaths/year)[1], with brick kilns contributing significantly (14%).[2]
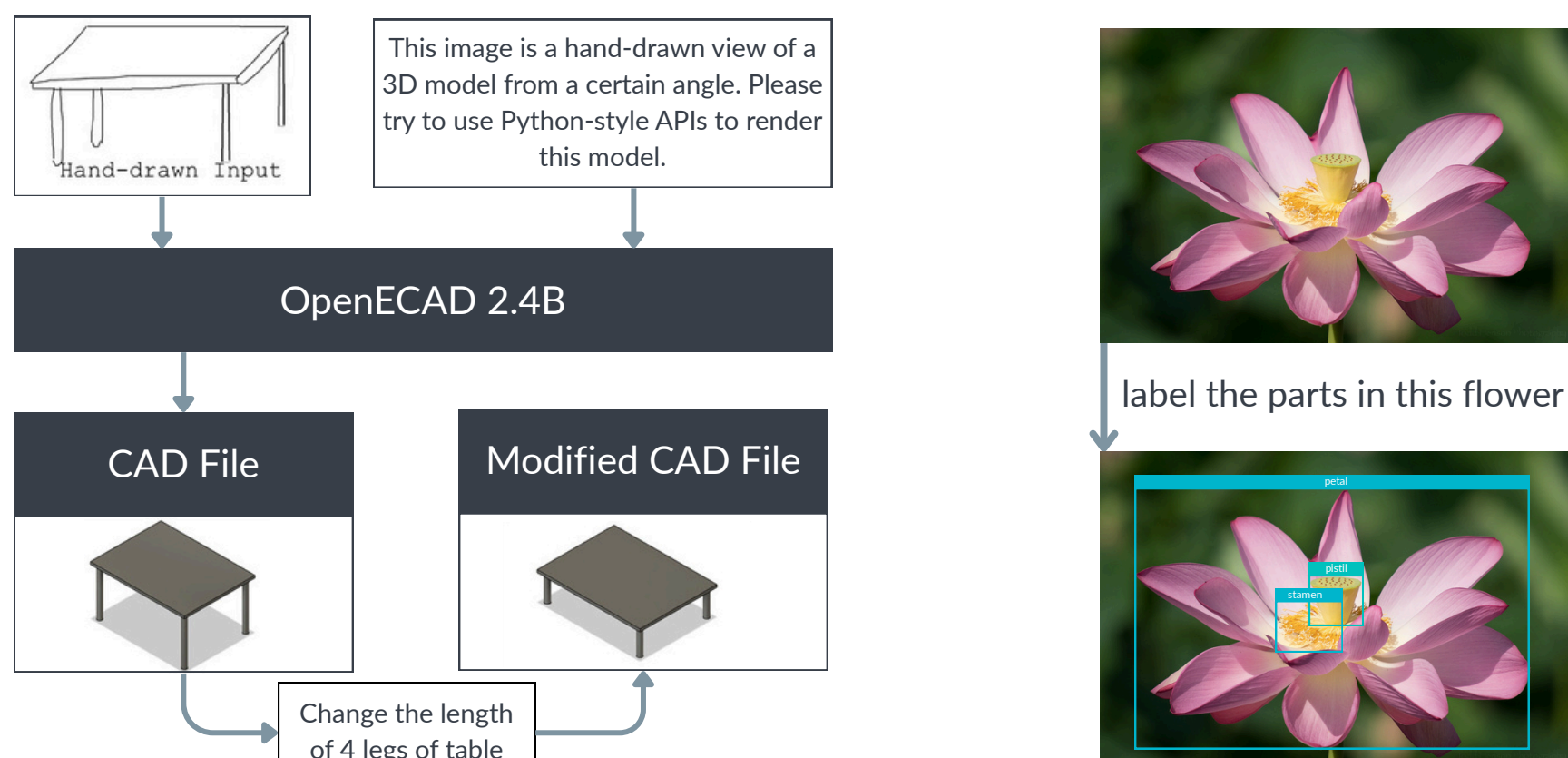
**Challenge:** Detecting these kilns is difficult as it is an unorganized sector. Manual monitoring is unscalable.

**Technical Gap:** Traditional object detection models (like YOLO) require large amounts of labelled domain-specific data, which is often scarce for satellite imagery tasks. They also struggle with domain adaptation (perform poorly when tested on regions different from their training data).

## Vision-Language Models



[{"bbox_2d": [407, 544, 676, 1023], "label": "athletic track"}]

Plotting

Where is the athletic track? Give bounding boxes.

**VLMs combine visual and text understanding**



This image is a hand-drawn view of a 3D model from a certain angle. Please try to use Python-style APIs to render this model.

OpenECAD 2.4B

Change the length of 4 legs of table

label the parts in this flower

**Models Tested**



| PaliGemma 3B | Florence2 large: 770M | Deepseek-VL2 27B | Qwen2.5-VL 72B | YOLOE 43.6M |

## Hypothesis Driven Approach



Zero-Shot VLMs

**Results:** Poor
**Hypothesis:** Image resolution is low

Super-resolution

**Results:** Still Poor
**Hypothesis:** The VLM is downscaling the image

Tiled Inference

**Results:** Better, but not good
**Hypothesis:** The VLM has not seen brick klins

Fine-tune on our data

**Results:** Better, but worse than specialist models
**Hypothesis:** Model hasn't seen enough satellite data

Fine-tune on larger satellite datasets, then on our data

Future Work

| Model Prediction | Ground Truth |

$$\frac{\text{Intersection}}{\text{Union}}$$

| IoU >= Threshold **True Positive** | IoU < Threshold **False Positive** | Unmatched Pred **False Positive** | Unmatched Truth **False Negative** |

## Key Results

**Model Comparison Across Methods**



Models: Florence2, Qwen2.5-VL, Paligemma, Deepseek-VL2, YOLOE

F1 Score @ 0.5 IoU (%)

Zero-Shot: 6.0%, 7.0%, 0.0%
Superresolution: 7.0%, 8.6%, 0.0%
Superresolution with Tiled Inference: 7.0%, 12.0%, 9.6%
Finetuning: 31.4%, 22.3%, 6.8%, 46.0%

**Model Performance vs. Training Data Size**



Models: YOLOv11, YOLOE, Florence2, Qwen2.5-VL, PaliGemma

F1 Score @ 0.5 IoU (%) vs Number of Training Images (Kiln) [Log Scale]

## Conclusions and Future Directions



N labelled Images → YOLO Training

M unlabelled Images → VLM → M labelled Images → Manual Verification → YOLO Training

Get new labels using trained YOLO

Iterative refinement with careful monitoring
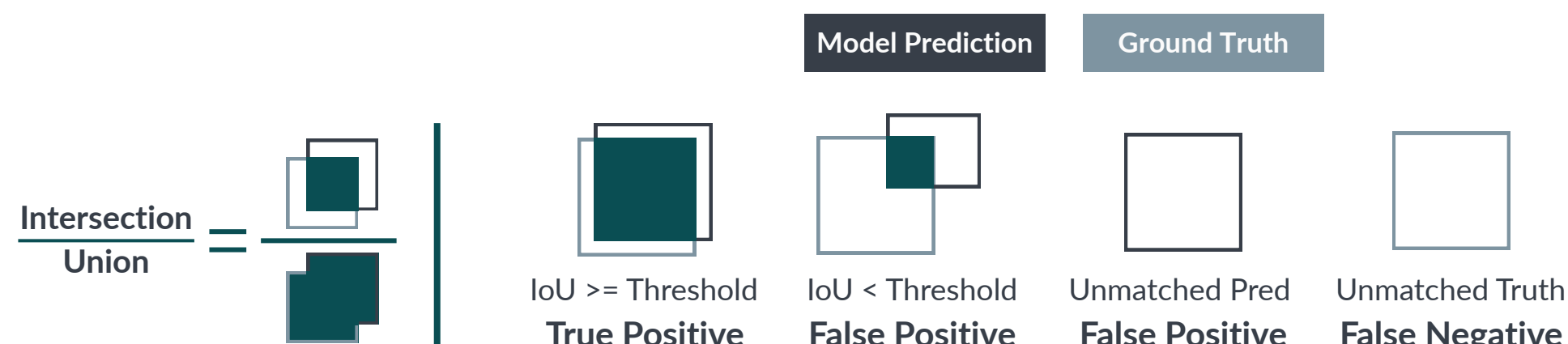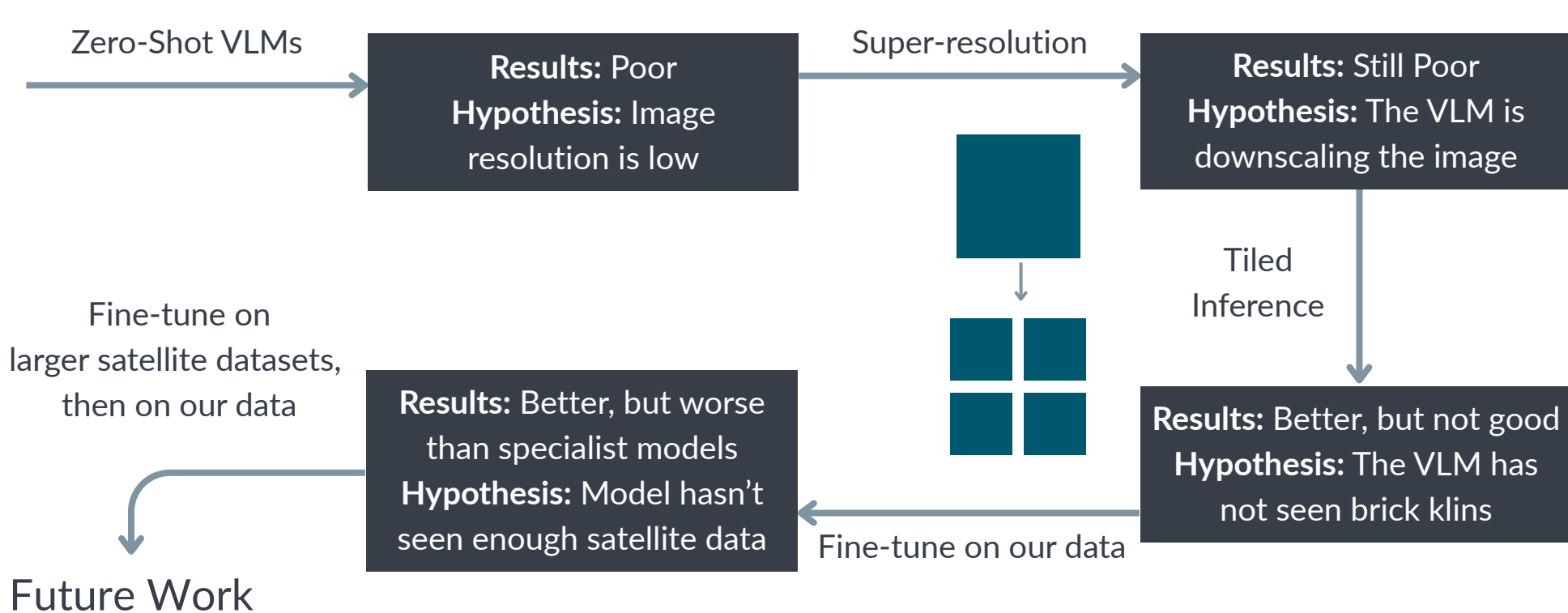
- VLMs show potential but require adaptation for satellite object detection.
- Image resolution, model architecture, and model vocabulary affect VLM performance
- Fine-tuning VLMs improves results significantly

## Acknowledgments

## References

[1] A. Pandey et al., "Health and economic impact of air pollution in the states of India: the Global Burden of Disease Study 2019," The Lancet Planetary Health, vol. 5, no. 1, Dec. 2020, doi: https://doi.org/10.1016/S2542-5196(20)30298-9.

[2] "Brick kilns embrace zigzag design to cut pollution and boost efficiency," Mongabay India, Aug. 2024. [Online]. Available: https://india.mongabay.com/2024/08/brick-kilns-embrace-zigzag-design-to-cut-pollution-and-boost-efficiency/. Accessed: Apr. 17, 2025.

[3] Planet Labs PBC, "Planet Application Program Interface: In Space for Life on Earth." [Online]. Available: https://api.planet.com. Accessed: 2024.

[4] A. Steiner et al., "PaliGemma 2: A Family of Versatile VLMs for Transfer," arXiv preprint arXiv:2412.03555, Dec. 2024. [Online]. Available: https://arxiv.org/abs/2412.03555

[5] B. Xiao et al., "Florence-2: Advancing a Unified Representation for a Variety of Vision Tasks," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2024, pp. 4818-4829.

[6] X. Chen et al., "DeepSeek-VL2: Mixture-of-Experts Vision-Language Models for Advanced Multimodal Understanding," arXiv preprint arXiv:2412.10302, Dec. 2024. [Online]. Available: https://arxiv.org/abs/2412.10302

[7] S. Bai et al., "Qwen2.5-VL Technical Report," arXiv preprint arXiv:2502.13923, Feb. 2025. [Online]. Available: https://arxiv.org/abs/2502.13923

[8] K. Kuckreja et al., "GeoChat: Grounded Large Vision-Language Model for Remote Sensing," arXiv preprint arXiv:2311.15826, Nov. 2023. [Online]. Available: https://arxiv.org/abs/2311.15826

[9] G. Jocher and J. Qiu, "Ultralytics YOLO11," ver. 11.0.0, 2024. [Online]. Available: https://github.com/ultralytics/ultralytics. Accessed: Apr. 17, 2025.

[10] A. Wang et al., "YOLOE: Real-Time Seeing Anything," arXiv preprint arXiv:2503.07465, Mar. 2025. [Online]. Available: https://arxiv.org/abs/2503.07465

[11] Z. Yuan, J. Shi, and Y. Huang, "OpenECAD: An efficient visual language model for editable 3D-CAD design," Comput. Graph., vol. 124, Art. no. 104048, 2024, doi: 10.1016/j.cag.2024.104048.