

SUMMARY

WILD-OV: A Unified Open-Vocabulary Detection Framework for Ecology and Agriculture

- Supports few- or zero-shot object detection using state-of-the-art vision–language models in both text- and image-conditioning modes.
- Effective in detecting objects across real-world ecological and agricultural scenarios.
- Enhances small-object detection with Slicing Aided Hyper Inference (SAHI).
- Provides a modular, user-friendly pipeline that lets domain experts choose between full-frame or sliced inference and swap between models with ease.

RESEARCH MOTIVATION

- New vantage points** (e.g., UAVs, UGVs) enable scalable, remote monitoring of complex agro-ecological scenes.
- High-resolution imagery** reveals fine-grained object interactions—critical for plant phenology, species monitoring, and yield estimation.
- Manual annotation bottleneck:** Systematic labeling of object/pixel-level data is labor-intensive, error-prone, and unscalable.
- Low-supervision approaches** (few-/zero-shot) and **open-vocabulary detectors** (OWLv2, Grounding DINO, Florence-2) offer promising alternatives.

KEY CHALLENGES

- Tiny and dense objects** in high-res imagery yield weak objectness and low detection recall.
- Domain shift** from natural images to aerial, under-canopy, and biological scenes breaks generalization of pretrained models.
- Image-guided inference** requires resolution-matched slicing and feature extraction for robust patch-to-patch comparison.
- Need for systematic evaluation:** Real-world deployment depends on holistic benchmarking across detectors, slicing strategies, and support-image conditioning.

SCHEMATIC OVERVIEW OF WILD-OV

Text-based (top): User queries are passed to OWLv2 or Grounding DINO to extract text embeddings and compare them with visual features for object localization.

Image-based (bottom): Exemplar bounding boxes in source images are used by OWLv2 to extract region-level features and match similar proposals in the target image. When slicing-aided hyperinference is enabled, detection operates on overlapping patches to improve small-object recall.

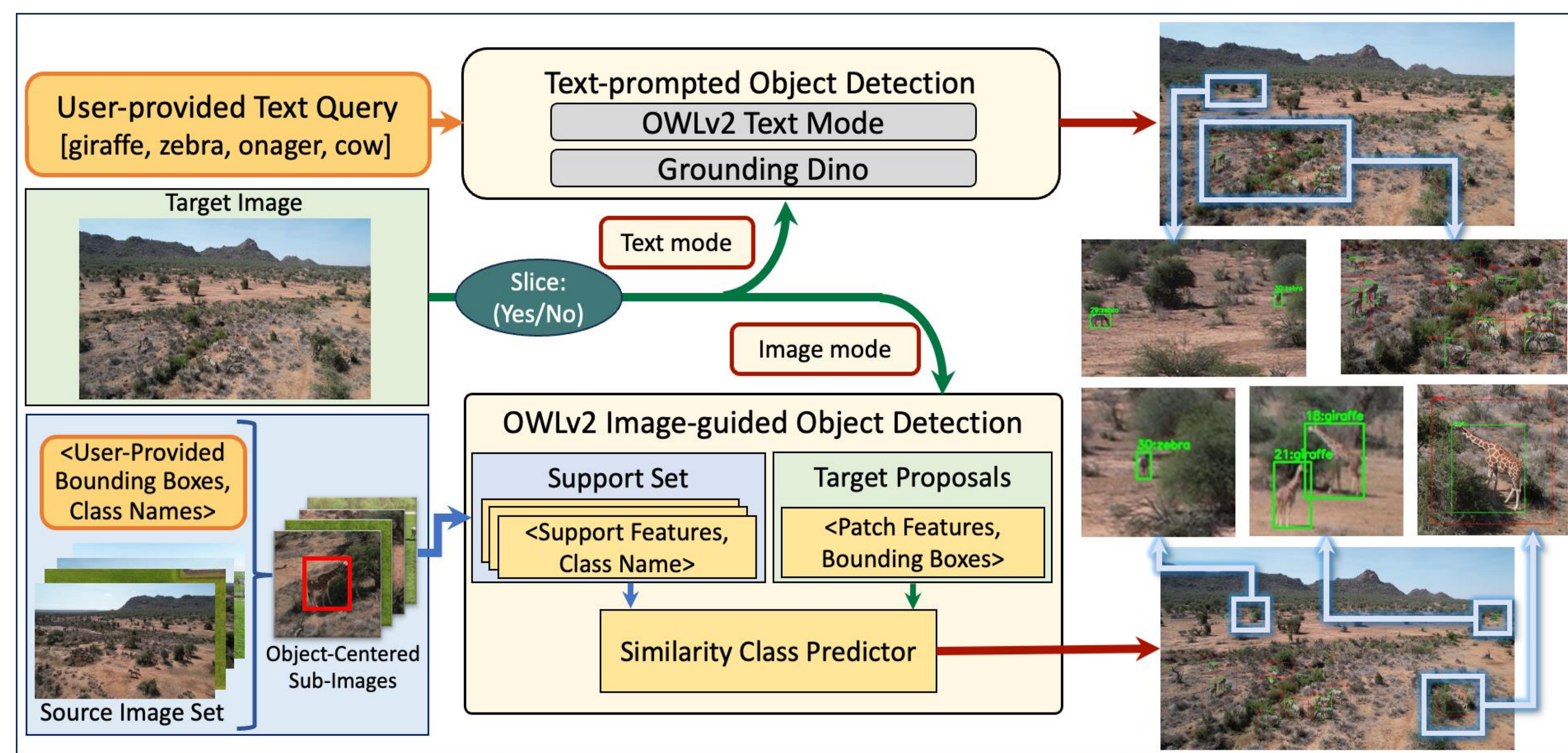


Figure 1: Schematic overview of WILD-OV: Text-based: Queries are embedded and matched to image features; Image-based: Sample boxes provide region features matched to proposals.

In this example, both the text-based and image-based variants of OWLv2 were applied independently to the same scene. It is observed that each variant was able to detect several additional animals in the background that human annotators had overlooked, highlighting the potential of open-vocabulary models to complement manual labeling. Red boxes denote the human-provided ground truth, while green boxes indicate the detections produced by the models.

PERFORMANCE ANALYSIS

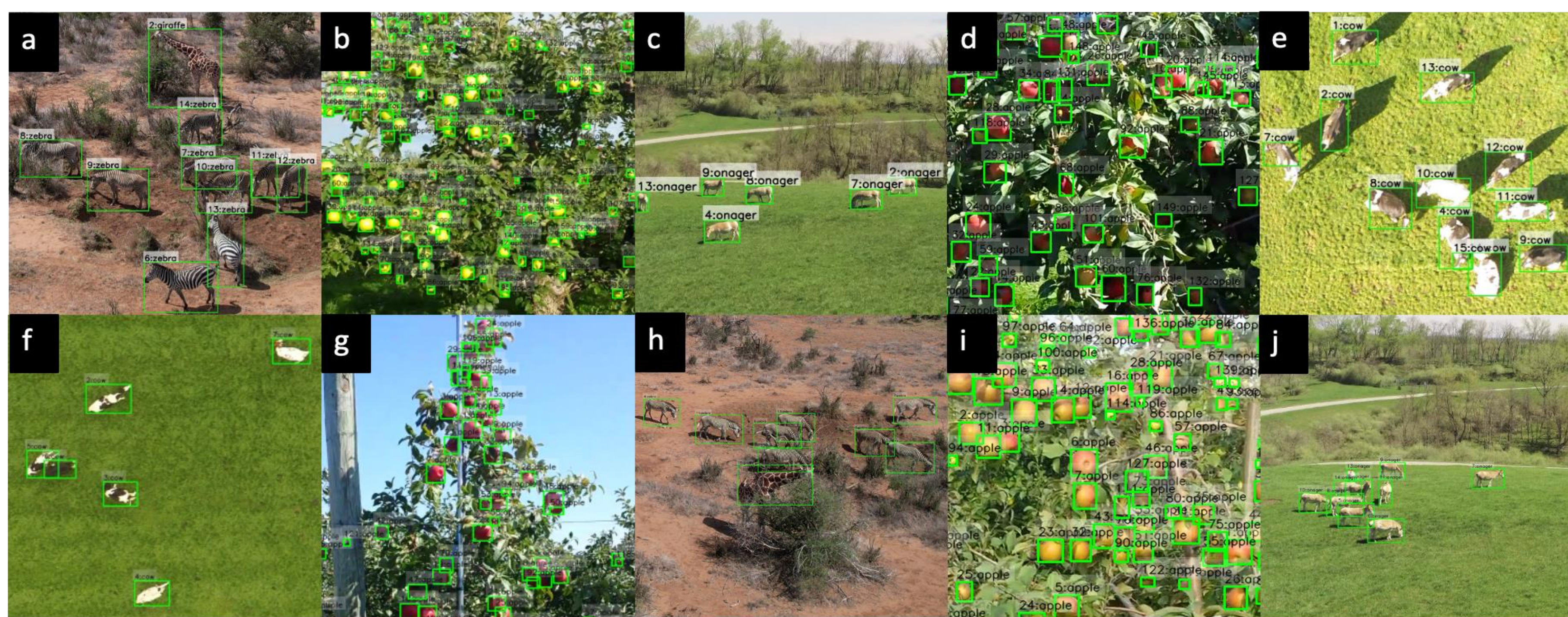


Figure 2: The image depicts the performance of the framework on our target datasets. The model accurately detects and labels different object categories even in complex backgrounds and high occlusion. The results demonstrate the model’s ability to generalize well across diverse settings without task-specific training or fine-tuning.

MAJOR CONTRIBUTIONS

- Introduced **WILD-OV**, a unified pipeline combining text- and image-guided open-vocabulary detectors with slicing-aided inference.
- Achieved **strong performance without extra training**, uncovering objects often missed by human annotators.
- OWLv2’s dual modes** provide flexibility: text prompts generalize to unseen classes, while image prompts detect small or occluded instances.
- Validated on **ecological (aerial animals)** and **agricultural (orchard apples)** datasets, showing broad applicability.
- Provides domain experts with a **modular, practical toolkit** to accelerate annotation and analysis.

ONGOING WORK

- Expand to **new crop types** and **diverse wildlife datasets**.
- Test detection of **rare/unseen species** with text-guided prompts.
- Benchmark additional VLMs (**Florence-2**, **Paligemma**, **BioCLIP**) against OWLv2 and Grounding DINO.
- Explore **lightweight fine-tuning** (LoRA, adapters) for specialized tasks.
- Optimize slicing for **faster, real-time edge deployment**.
- Develop **user-friendly interfaces** to simplify adoption by non-ML experts.

References

- Matthias Minderer, Alexey A. Gritsenko, and Neil Houlsby. "Scaling Open-Vocabulary Object Detection." In Thirty-seventh Conference on Neural Information Processing Systems. 2023.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyu Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. "Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection." (2024).
- Akyon, F., Onur Altinuc, S., & Temizel, A. (2022). Slicing Aided Hyper Inference and Fine-Tuning for Small Object Detection. In 2022 IEEE International Conference on Image Processing (ICIP) (pp. 966–970). IEEE.