

# ReLeaf: Benchmarking Leaf Segmentation across Domains and Species

## Supplementary Material

Robert Martinko<sup>1,2</sup> Daniel Steininger<sup>1</sup> Julia Simon<sup>1</sup> Andreas Trondl<sup>1</sup> Matthias Blaickner<sup>2</sup>

<sup>1</sup>AIT Austrian Institute of Technology, Center for Vision, Automation & Control

<sup>2</sup>University of Applied Sciences Technikum Wien, Computer Science & Applied Mathematics

{robert.martinko, daniel.steininger, julia.simon, andreas.trondl.fl}@ait.ac.at

matthias.blaickner@technikum-wien.at

This supplementary document complements the main paper with additional statistics regarding our benchmark dataset and more detailed quantitative and qualitative experimental results to facilitate a more thorough understanding of applied methods and insights.

### A. Benchmark dataset

Our *CropAndWeedAndLeaf* benchmark comprises a representative set of images sampled from 23 classes of the *CropAndWeed* dataset [5]. Tab. 1 gives a complete list of the exact species and the *ReLeaf* model’s performance on each of them. It shows strong variations between classes related to their appearance and leaf structure, highlighting the overall need for more variable leaf-segmentation data.

### B. Extended results

The following sections show detailed results of both our ablation and cross-dataset experiments.

#### B.1. Architecture ablation

As stated in the main paper, RF-DETR was omitted from the final ablation after initial experiments showed drastically inferior performance compared to other architectures under the applied hardware constraints. Tab. 2 provides an overview of these initial results. While inference times of RF-DETR are comparable to those of YOLO26 and even slightly lower than for Detectron2, mAP values are significantly inferior to both. At the same time, training times are higher by a factor of 3.5 compared to Detectron2 at the same image resolution and 6.3 compared to YOLO26 models trained at even higher resolutions, which are not available for RF-DETR due to its substantial VRAM requirements. Mitigating these discrepancies and making transformer-based architectures usable for our hardware and data setup could still be achieved by using specialized architecture variants (e.g., RT-DETR [3], Deformable DETR [9])

Table 1. Overview of 23 plant species included in *CropAndWeedAndLeaf* with leaf-segmentation annotations for 15 samples each, along with corresponding instance-segmentation performance (% mAP<sup>50-95</sup>). Scores are results of the selected model architecture (YOLO26 *Medium*, input size of 768<sup>2</sup> pixels), trained on the combined *ReLeaf* dataset.

Species name	Botanical name	Score
Black-bindweed	<i>Fallopia convolvulus</i>	41.6
Broad bean	<i>Vicia faba</i>	49.7
Common bean	<i>Phaseolus vulgaris</i>	66.5
Common corncockle	<i>Agrostemma githago</i>	20.2
Common sunflower	<i>Helianthus annuus</i>	73.0
Copse bindweed	<i>Fallopia dumetorum</i>	43.9
Cornflower	<i>Centaurea cyanus</i>	30.0
Creeping thistle	<i>Cirsium arvense</i>	39.9
Field sowthistle	<i>Sonchus arvensis</i>	34.7
Maize	<i>Zea mays</i>	41.2
Maple-leaf goosefoot	<i>Chenopodium hybridum</i>	31.7
Pea	<i>Pisum sativum</i>	11.2
Poppy	<i>Papaver</i>	54.5
Potato	<i>Solanum tuberosum</i>	32.5
Red-root amaranth	<i>Amaranthus retroflexus</i>	57.9
Redshank	<i>Persicaria maculosa</i>	26.7
Ribwort plantain	<i>Plantago lanceolata</i>	34.8
Soybean	<i>Glycine max</i>	27.2
Squash	<i>Cucurbita</i>	60.9
Sugar beet	<i>Beta vulgaris s. vulgaris</i>	66.0
Thornapple	<i>Datura stramonium</i>	57.6
White goosefoot	<i>Chenopodium album</i>	39.4
Zucchini	<i>Cucurbita pepo var. gir.</i>	68.2

or dedicated training strategies, such as memory-efficient attention mechanisms or staged training starting from plant-level instance-segmentation datasets. These options may be worth exploring in future extensions of this work, but were

Table 2. Initial ablation results for RF-DETR *SegPreview* with multiple input resolutions, including mask accuracy as well as training and inference times (training at an input resolution of 768<sup>2</sup> pixels was not feasible due to high VRAM requirements). Results of selected YOLO26 and Detectron2 variants (*Medium* and *ResNet-101*, respectively) are provided for reference. All models are trained for 20 epochs using the *PhenoBench* training/validation splits and evaluated on the corresponding test set.

	Size [px]	Accuracy [% mAP]	Training [h/epoch]	Inference [ms/image]
<b>RF-DETR</b>	192 <sup>2</sup>	55.9	0.85	18.0
	384 <sup>2</sup>	59.2	0.94	19.7
	576 <sup>2</sup>	59.9	1.13	23.4
<b>Detectron2</b>	576 <sup>2</sup>	80.0	0.32	33.1
<b>YOLO26</b>	768 <sup>2</sup>	78.0	0.18	21.0

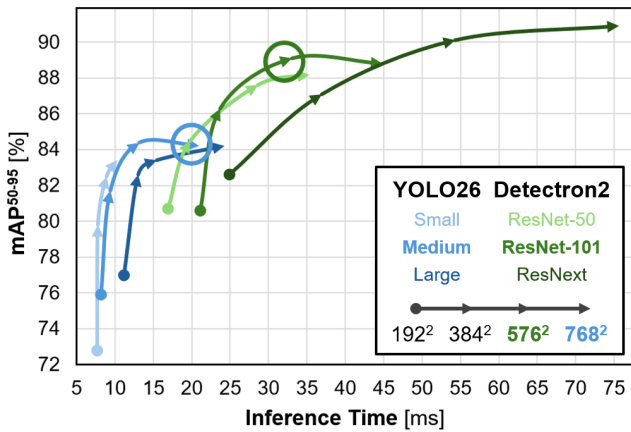


Figure 1. Comparison of bounding-box accuracy (% mAP<sup>50-95</sup>) and inference time (ms) of leaf-segmentation for multiple model-architecture variants and input resolutions on the *PhenoBench* test set. Selected variants providing a reasonable trade-off for real-time deployment are marked by circles.

not within the scope of our analysis.

In addition to the accuracy of leaf-segmentation masks provided in Fig. 4 of the main paper, we present the corresponding bounding-box accuracy in Fig. 1 to allow a differentiation between the models’ capabilities of correctly localizing individual leaves and accurately delineating their contours. As visible in the results, the discrepancy is significantly larger for Detectron2 [8] compared to YOLO26 [1], indicating that it is a superior choice for applications requiring pure detection in the form of bounding boxes without fine-grained object masks. For the purposes investigated in this work, however, YOLO26 remains a more suitable choice due to its lower processing time. Tab. 3 shows the exact values forming the basis of both charts.

Tab. 4 shows the performance of our selected configu-

Table 3. Comparison of bounding-box (**Box**) and segmentation-mask (**Mask**) accuracies (% mAP<sup>50-95</sup>) and inference times (**t**) (ms) of leaf-segmentation for multiple model-architecture variants and input resolutions on the *PhenoBench* test set. Selected variants providing a reasonable trade-off for real-time deployment are highlighted by bold text.

	Size	Box	Mask	t	
YOLO26	Small	192 <sup>2</sup>	72.8	65.0	7.7
		384 <sup>2</sup>	79.8	75.0	7.8
		576 <sup>2</sup>	82.5	79.1	8.9
		768 <sup>2</sup>	83.4	80.4	10.3
	Medium	192 <sup>2</sup>	75.9	66.8	8.2
		384 <sup>2</sup>	81.7	76.6	9.4
		576 <sup>2</sup>	84.4	80.4	13.2
		<b>768<sup>2</sup></b>	<b>84.2</b>	<b>81.4</b>	<b>21.0</b>
	Large	192 <sup>2</sup>	77.0	67.5	11.2
		384 <sup>2</sup>	82.6	76.9	13.1
		576 <sup>2</sup>	83.4	79.7	15.5
		768 <sup>2</sup>	84.2	81.3	24.3
ResNet-50	192 <sup>2</sup>	80.7	75.9	16.9	
	384 <sup>2</sup>	84.5	79.2	19.9	
	576 <sup>2</sup>	87.6	81.5	28.6	
	768 <sup>2</sup>	88.2	81.5	35.3	
Detectron2	ResNet-101	192 <sup>2</sup>	80.6	75.9	21.1
		384 <sup>2</sup>	86.3	80.4	23.7
		<b>576<sup>2</sup></b>	<b>89.1</b>	<b>82.4</b>	<b>33.1</b>
		768 <sup>2</sup>	88.8	82.1	44.9
ResNext	192 <sup>2</sup>	82.6	77.4	24.9	
	384 <sup>2</sup>	87.1	80.9	37.0	
	576 <sup>2</sup>	90.1	83.7	54.3	
	768 <sup>2</sup>	90.9	83.9	75.7	

Table 4. Comparison of bounding-box (**Box**) and segmentation-mask (**Mask**) accuracies (mAP<sup>50-95</sup>) for varying objects sizes (Small < 32<sup>2</sup> pixels < Medium < 96<sup>2</sup> pixels < Large) using the selected configurations of YOLO26 (*Medium* with an input resolution of 768<sup>2</sup> pixels) and Detectron2 (*ResNet-101* with an input resolution of 576<sup>2</sup> pixels).

	Box			Mask		
	S	M	L	S	M	L
<b>YOLO26</b>	69.9	82.2	91.3	69.5	79.3	89.1
<b>Detectron2</b>	87.4	88.0	91.7	80.9	80.9	86.3

rations of YOLO26 and Detectron2 on different leaf sizes. It confirms our prior assumption that two-stage detectors like Detectron2 excel at retrieving and accurately delineating small objects at the cost of higher parameter counts and therefore longer inference times than single-stage variants.

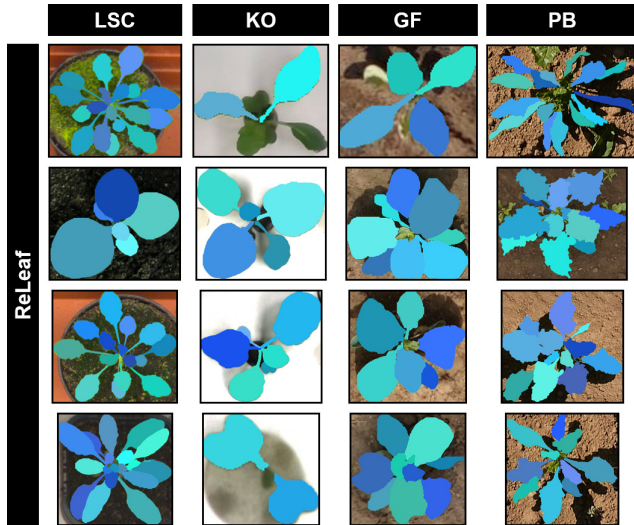


Figure 2. Representative leaf-segmentation results of selected YOLO26 configuration (*Medium* with an input resolution of  $768^2$  pixels) trained on a combination of *LSC* [4], *Komatsuna* (KO) [6], *GrowliFlower* (GF) [2] and *PhenoBench* (PB) [7] datasets and evaluated on the respective test sets.

## B.2. Cross-dataset experiments

Fig. 2 visualizes our multi-domain *ReLeaf* model’s results on a selection of representative test images from *LSC* [4], *Komatsuna* [6], *GrowliFlower* [2] and *PhenoBench* [7], confirming its stable performance across both species and environmental conditions.

Figs. 3 and 4 show the results of models trained on each dataset and evaluated on one representative image for each of the 23 plant species in the *CropAndWeedAndLeaf* benchmark. Models trained on laboratory data (*LSC*, *Komatsuna*) struggle most to generalize to novel plant species in real-world settings. While the combined *ReLeaf* model occasionally misses small inner leaves, it significantly improves adaptation to the strongly varying visual appearances and leaf configurations of different crop and weed species.

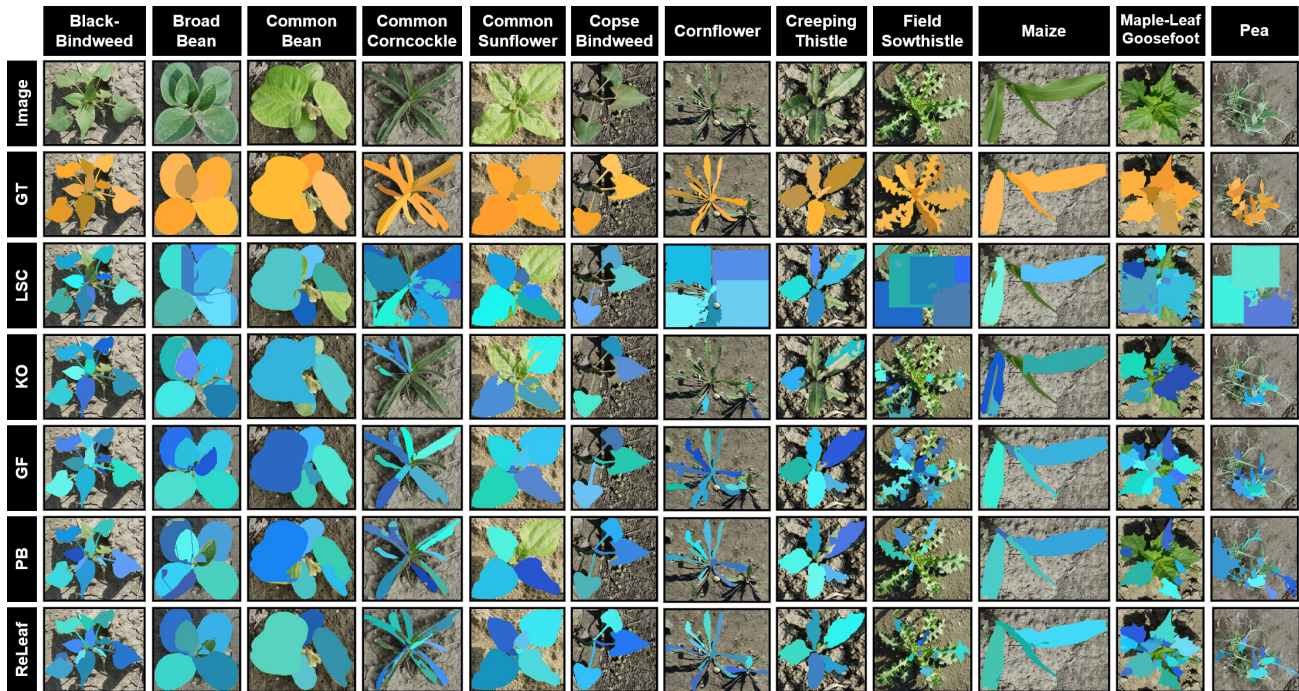


Figure 3. Representative leaf-segmentation results of YOLO26 models (*Medium*, 768<sup>2</sup> pixels input resolution) trained on *LSC* [4], *KOmatsuma* [6], *GrowliFlower* [2] and *PhenoBench* [7] datasets and their combination (*ReLeaf*), evaluated on all 23 species of the *CropAndWeedAndLeaf* benchmark, along with corresponding ground truth (GT) (Part 1).

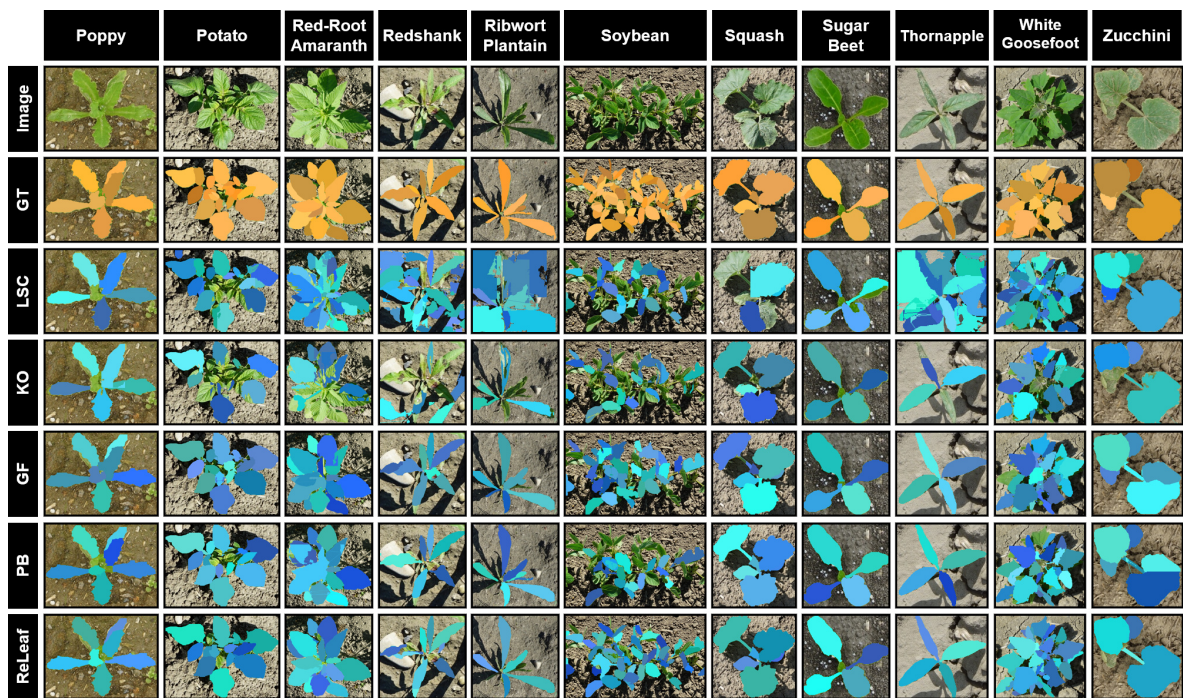


Figure 4. Representative leaf-segmentation results of YOLO26 models (*Medium*, 768<sup>2</sup> pixels input resolution) trained on *LSC* [4], *KOmatsuma* [6], *GrowliFlower* [2] and *PhenoBench* [7] datasets and their combination (*ReLeaf*), evaluated on all 23 species of the *CropAndWeedAndLeaf* benchmark, along with corresponding ground truth (GT) (Part 2).

## References

- [1] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics yolo. <https://github.com/ultralytics/ultralytics>, 2024. Accessed: 2025-10-11. 2
- [2] Jana Kierdorf, Laura Verena Junker-Frohn, Mike Delaney, Mariele Donoso Olave, Andreas Burkart, Hannah Jaenicke, Onno Muller, Uwe Rascher, and Ribana Roscher. Growliflower: An image time-series dataset for growth analysis of cauliflower. *Journal of Field Robotics*, 40(2):173–192, 2023. 3, 4
- [3] Wenyu Lv, Yian Zhao, Qinyao Chang, Kui Huang, Guanzhong Wang, and Yi Liu. Rt-detr2: Improved baseline with bag-of-freebies for real-time detection transformer. *arXiv preprint arXiv:2407.17140*, 2024. 1
- [4] Hanno Scharf, Massimo Minervini, Andreas Fischbach, and Sotirios A Tsafaris. Annotated image datasets of rosette plants. In *European conference on computer vision*, pages 6–12. Suisse Zürich, 2014. 3, 4
- [5] Daniel Steininger, Andreas Trondl, Gerardus Croonen, Julia Simon, and Verena Widhalm. The cropandweed dataset: A multi-modal learning approach for efficient crop and weed manipulation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3729–3738, 2023. 1
- [6] Hideaki Uchiyama, Shunsuke Sakurai, Masashi Mishima, Daisaku Arita, Takashi Okayasu, Atsushi Shimada, and Rin-ichiro Taniguchi. An easy-to-setup 3d phenotyping platform for komatsuna dataset. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 2038–2045, 2017. 3, 4
- [7] Jan Weyler, Federico Magistri, Elias Marks, Yue Linn Chong, Matteo Sodano, Gianmarco Roggiolani, Nived Chebrolu, Cyrill Stachniss, and Jens Behley. Phenobench: A large dataset and benchmarks for semantic image interpretation in the agricultural domain. *IEEE transactions on pattern analysis and machine intelligence*, 46(12):9583–9594, 2024. 3, 4
- [8] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 2
- [9] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2021. 1