# Weakly Supervised Panoptic Segmentation for Defect-Based Grading of Fresh Produce Supplementary Materials

## A. Dataset visualizations



Figure S.1. **Defect examples.** Example images of banana surface defects we aim to detect in this study. Bruises usually result from a dull impact, while scars are caused by a sharp impact. Old defects are usually darker than new defects as they oxidize over time.

#### **B.** Implementation details

All images were resized and padded to  $1024^2$  pixels resolution. Training images were augmented with a 50% Random Horizontal Flip and with additional random color-based augmentations (see Table S.1) to increase model robustness.

Table S.2 shows the hyperparameters used to train the Maskformer models. We evaluated the model every five epochs on the validation set and saved the model with the highest Panoptic Quality. All models were trained on a single NVIDIA RTX A5000 GPU with 24GB of memory.

Table S.1. Random color augmentations applied to training images. Values are uniformly sampled from the specified ranges.

Property	Sampling range
brightness	[0.9, 1.1]
contrast	[0.9, 1.1]
saturation	[0.9, 1.1]
hue	[-0.05, 0.05]

Batch size	2
Epochs	100
Evaluation frequency	every 5 epochs
Optimizer	Adam [12]
Learning rate	$5 \times 10^{-5}$
Learning rate schedule	constant

## C. Additional SAM evaluations



Figure S.2. **Distribution of IoU values when comparing annotated and SAM-generated defect masks.** While the ViT-B and ViT-L variants show a similar distribution, using the largest model size (ViT-H) leads to significantly higher overlap with hand-annotated masks.

### D. Results using four defect categories

In this section, we describe our results when using four different defect classes instead of a single joint category. Our results clearly show that while the detection and segmentation of bruises and scars work well, a reliable categorization into one of the four predefined classes is not possible with the current approach and dataset. We hypothesize that this limitation stems from one or more of the following factors:

- 1. Ambiguous annotation: The classification of defects across the four predefined categories may be subjective to a large degree. The distinction between "old" and "new" defects is not always clear-cut, as the transition is gradual. Additionally, distinguishing bruises from scars can be ambiguous, especially for non-experts (see Figure S.1).
- 2. Image resolution: We use 1024<sup>2</sup> pixel resolution for images. While this is generally considered high for machine learning tasks, defects can be small and, thus, only be represented by a few pixels, making them harder to categorize.
- 3. The four defect types are unevenly represented in our dataset (37/182/387/834). A more balanced distribution of categories and a larger number of defect samples are likely to improve categorization accuracy. We recommend collecting more examples from the two underrepresented classes for future work.

Table S.3. **Results using multiple defect categories.** It is evident that our models are unable to categorize defect types. Most likely due to ambiguous annotations and/or limited training data. The configuration highlighted in blue is used for the visualizations in Figure S.3.

		Defects												
		Defect masks		Old H	Bruise	New Bruise		Old Scar		New Scar		Overall		
Model	PP	train	val	AP	IoU	AP	IoU	AP	IoU	AP	IoU	mAP	mIoU	PQ
Maskformer		Anno.	Anno.	$0.031 \pm .059$	$.028\pm.043$	$.043\pm.024$	$.157 \pm .037$	$.036\pm.018$	$.225\pm.097$	$.050\pm.018$	$.316 \pm .054$	$.040\pm.019$	$.482\pm.024$	$.471\pm.023$
Maskformer	$\checkmark$	Anno.	Anno.	$0.030 \pm .059$	$.013 \pm .018$	$.088 \pm .039$	$.114\pm.035$	$.034\pm.016$	$.223\pm.083$	$.017\pm.009$	$.066\pm.021$	$.042 \pm .022$	$.436\pm.018$	$.431\pm.015$
Maskformer		SAM-L	Anno.	$.039 \pm .055$	$.060 \pm .080$	$.034 \pm .016$	$.175 \pm .074$	$.045 \pm .034$	$.179 \pm .022$	$.066\pm.018$	$.345 \pm .046$	$.046 \pm .012$	$.487 \pm .018$	$.477 \pm .013$
Maskformer	$\checkmark$	SAM-L	Anno.	$.032 \pm .049$	$.031 \pm .038$	$.074 \pm .017$	$.124 \pm .051$	$.038 \pm .021$	$.161 \pm .009$	$.032 \pm .014$	$.104 \pm .039$	$.044 \pm .015$	$.437 \pm .018$	$.435 \pm .019$
Maskformer		SAM-L	SAM-L	$.037 \pm .053$	$.063 \pm .081$	$.036 \pm .016$	$.178 \pm .075$	$.045 \pm .033$	$.182 \pm .023$	$.071 \pm .019$	$.353 \pm .049$	$.047 \pm .012$	$.489 \pm .018$	$.478 \pm .013$
Maskformer	$\checkmark$	SAM-L	SAM-L	$.032 \pm .049$	$.027\pm.034$	$.080\pm.018$	$.127\pm.051$	$.036\pm.020$	$.163\pm.010$	$.034\pm.015$	$.108\pm.041$	$.046\pm.015$	$.437 \pm .018$	$.436\pm.019$
Maskformer		SAM-H	Anno.	$0.019 \pm .026$	$.074 \pm .078$	$.043 \pm .020$	$.157 \pm .092$	$.031 \pm .021$	$.252 \pm .063$	$.056 \pm .016$	$.318 \pm .047$	$.037 \pm .014$	$.495 \pm .010$	$.474 \pm .010$
Maskformer	$\checkmark$	SAM-H	Anno.	$.032 \pm .056$	$.030 \pm .029$	$.093 \pm .022$	$.145 \pm .074$	$.025 \pm .029$	$.230 \pm .041$	$.029 \pm .011$	$.100\pm.050$	$.045 \pm .016$	$.453 \pm .013$	$.449 \pm .009$
Maskformer		SAM-H	SAM-H	$.025 \pm .035$	$.092 \pm .105$	$.049 \pm .022$	$.163 \pm .091$	$.037 \pm .024$	$.264 \pm .064$	$.062 \pm .018$	$.328 \pm .047$	$.043 \pm .015$	$.501 \pm .010$	$.482 \pm .011$
Maskformer	$\checkmark$	SAM-H	SAM-H	$0.045 \pm .0078$	$.029 \pm .027$	$.106\pm.022$	$.149 \pm .073$	$.027\pm.028$	$.239 \pm .041$	$.034\pm.010$	$.104 \pm .051$	$.053 \pm .019$	$.455\pm.013$	$.456\pm.010$



Figure S.3. Example visualizations of annotated vs. predicted masks using four defect categories. Left: Input Image, Mid: Annotation, Right: Maskformer Prediction. Segments are color-coded as follows: Foreground Banana , Background Banana , Old Bruise ,