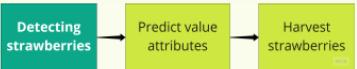


Research Question

How effective can a convolutional neural network be used to efficiently detect strawberries in the by TU Delft provided dataset?

1. Introduction

The goal of this project is to combat food waste by implementing part of the pipeline to more efficiently harvest strawberries.



Current research does tackle strawberry detection, however these papers do not compare multiple CNNs against each other. Therefore this research consists of the following sub questions:

- which CNN model does the best detection for the dataset provided by the supervisors?

3. Results & Discussion

- 100 epoch single camera
 - overfitting (figure 4)
- Doubling the data was tested.
 - On average double camera dataset 8.42% better.
 - Still overfitting

In total the following was tested:

- 1 camera trainset:
 - 10 & 100 epoch
 - learning rate from 0.002 to 0.2
- 2 camera trainset:
 - 5, 10, 25, 50 epoch
 - learning rate from 0.002 to 0.2

Model	mAP BBox	LR	Epoch	no. cams
Faster r-cnn	51.00	0.002	5	2
Mask r-cnn	51.63	0.002	5	2
RetinaNet	49.36	0.002	10	2

8. The best BBox mAP results for each CNN

Model	mAP Seg	LR	Epoch	No. cams
Mask r-cnn	73.20	0.02	10	2
Mask r-cnn	71.56	0.2	10	1
Mask r-cnn	71.27	0.02	5	1

9. The top 3 segm mAP

2. Methodology

CNNs chosen and implemented with detectron2 [1]:

- Faster R-CNN [2] (figure 1),
- Mask R-CNN [3],
- RetinaNet [4]

Training data:

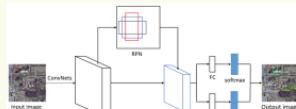
- 3 camera
 - The same one used as testing camera
- 4.370 usable images with 130.665 strawberries
- Ground truth provided (figure 2)

Test Metric:

- Precision (figure 3)
- mean Average Precision (AP)
 - From 50% to 95% confidence with steps of 5%

$$\text{Precision} = \frac{tp}{tp + fp}$$

3. How precision is calculated

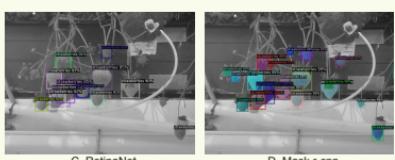
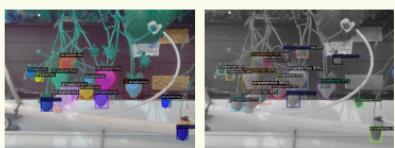


1. General layout of faster r-cnn

source: http://www.researchgate.net/profile/The_Framework_of_Faster_R-CNN_BPN_Region_Convolutional_Networks_Dual_Condition_of_Interest_FCN_dg_132046969



2. The images provided.



3. Incorrect ground truth

4. Conclusion

- Best results:
 - BBox mAP: **Mask r-cnn, lr:0.002 epochs:5, no. cams: 2**
 - Segm mAP: **Mask r-cnn, lr:0.02 epochs:10, no. cams: 2**
- Training favors lower learning rates.
- Higher learning rate sometimes made the training loss diverge.
- Adding more data improves the accuracy.

References

[1]. Wu, A., Kirillov, F., Massa, W.-Y., Lo, and R. Girshick, "Detectron2."

<https://github.com/facebookresearch/detectron2>, 2019

[2] R. Girshick, "Fast r-cnn," in Proceedings of the IEEE international conference on computer vision, pp. 1440–1448, 2015

[3] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in Proceedings of the IEEE international conference on computer vision, pp. 2961–2969, 2017.

[4] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in Proceedings

of the IEEE international conference on computer vision, pp. 2980–2988, 2017.

