

Detecting Olives with Synthetic or Real Data? Olive the Above

Yianni Karabatis¹, Xiaomin Lin¹, Nitin J. Sanket², Michail G. Lagoudakis³, Yiannis Aloimonos¹

Abstract—Modern robotics has enabled the advancement in yield estimation for precision agriculture. However, when applied to the olive industry, the high variation of olive colors and their similarity to the background leaf canopy presents a challenge. Labeling several thousands of very dense olive grove images for segmentation is a labor-intensive task. This paper presents a novel approach to detecting olives without the need to manually label data. In this work, we present the world’s first olive detection dataset comprised of synthetic and real olive tree images. This is accomplished by generating an auto-labeled photorealistic 3D model of an olive tree. Its geometry is then simplified for lightweight rendering purposes. In addition, experiments are conducted with a mix of synthetically generated and real images, yielding an improvement of up to 66% compared to when only using a small sample of real data. When access to real, human-labeled data is limited, a combination of mostly synthetic data and a small amount of real data can enhance olive detection.

I. INTRODUCTION

The olive tree, native to the Mediterranean region, has played an essential role in the advancement of humanity for millennia. The olive harvest is a lengthy process requiring almost a year-long preparation and is essential for the global economy, with a market valued at \$13.77B and is forecasted to reach \$17.99B by 2029 [1].

Yield estimation aims to forecast the output of crops. Manual inspections help farmers better understand their crops, but patrolling a grove with thousands of trees is time-consuming. One solution is to utilize Unmanned Aerial Vehicles (UAVs) to capture and analyze aerial image data for making better decisions.

Object detection via robotics has been employed for multiple types of crops [2]–[6]. However, olives present additional problems not encountered with other crops. Specifically, one must account for moderate variation in olive color, since the olive may have multiple shades of green. Another problem is the similarity between the olive fruit and the surrounding leaves, making it challenging to differentiate between the two. This similarity can be exacerbated even further when only the end of an olive is visible, resembling a leaf. In addition, the dense canopy of an olive tree occludes many of the olives, making them partially visible and difficult to recognize as shown in Fig. 1. Finally, the average size of

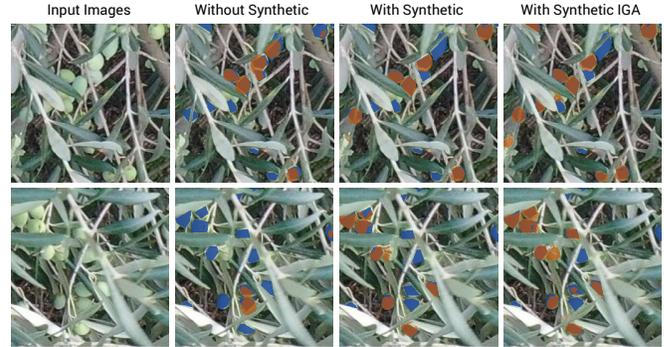


Fig. 1. Each column left to right: Input image, prediction using only real data, prediction using real and synthetic data, prediction using real and synthetic data in IGA. Predictions are shown in orange and ground truths are shown in blue. Adding synthetic data to the training set increases the number of correct predictions.

an olive is $2\text{cm} \times 1\text{cm} \times 1\text{cm}$, magnifying all the difficulties mentioned above.

However, labeling images collected from UAVs is an exhausting task, especially when there could be up to 17,500 olives [7] in a single olive tree (up to 1000 olives in a single image from our observation). An olive grove may contain thousands of olive trees, making this very labor-intensive. *Therefore, we investigate a new approach to cover the deficiency of labeled data, namely we create and incorporate synthetic data as an aid to training models for olive segmentation.* Synthetic data is automatically labeled and can be easily scaled up to higher orders of magnitude. We create synthetically generated image-mask pairs of an olive tree, including its olives, leaves, branches, and surrounding background. We propose the use of our synthetic data to create a high-performing, automated system for detecting olives in olive groves. The contributions of this paper are as follows:

- We create the world’s first olive detection dataset consisting of mostly synthetic and some real images of olive trees
 - We generate a photorealistic 3D model of an olive tree to create an auto-labeled synthetic dataset
 - We propose and simplify a geometric model for olives on an olive tree for lightweight rendering purposes
- We create a novel color input space for our synthetic data, making it more generalizable to real-world data
- We experiment with segmentation models to assist in situations where the amount of real labeled data is very limited

The remainder of this paper is structured as follows:

*This work was funded by the Fulbright foundation

¹Perception and Robotics Group, University of Maryland Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA. Emails: {yianni,xlin01,jyaloimo}@umd.edu.

²Perception and Autonomous Robotics Group, Robotics Engineering, Worcester Polytechnic Institute, MA 01609, USA. Email: nsanket@wpi.edu.

³School of ECE, Technical University of Crete, Greece. Email: lagoudakis@tuc.gr.

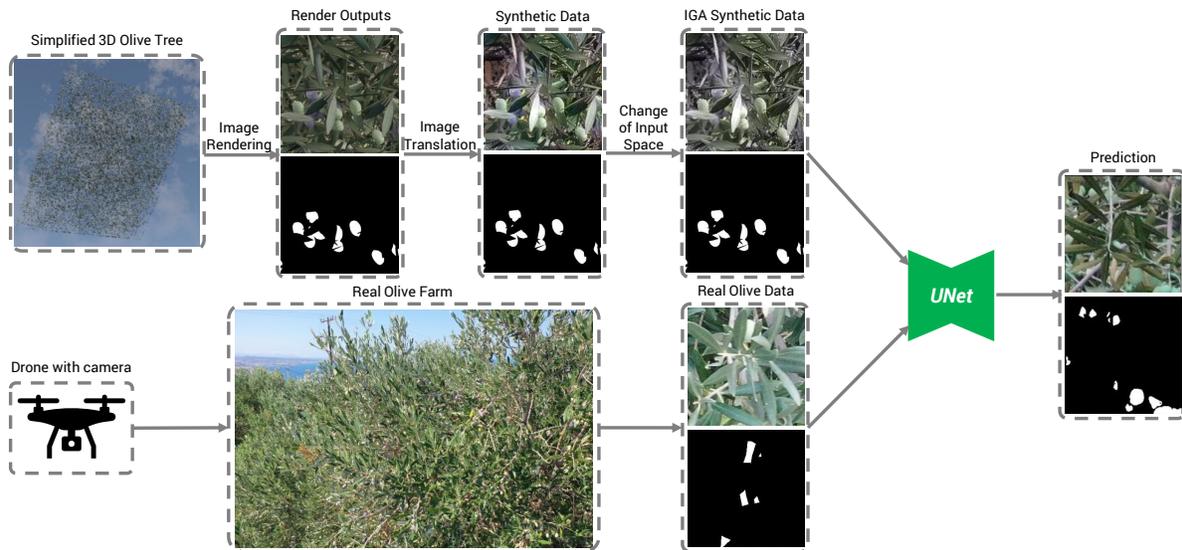


Fig. 2. An overview of our approach: The proposed geometric model is used to generate synthetic images which are further fed into a VSAIT to enable image-to-image translation. We then combine the synthetic data with real data to train an Unet for olive detection.

Section II presents related works. In Section III we describe the generation of the 3D synthetic olive tree model and discuss the image-to-image translation technique used to create realistic synthetic images. Section IV presents the experiments using a mix of synthetic and real data. Finally, Section V contains the conclusion and future work.

II. RELATED WORK

Climate change along with the rise of the human population calls for efficient and modern agricultural solutions [8]. Precision agriculture offers a partial solution by maximizing crop outputs while minimizing the environmental footprint [9].

UAVs have reduced the time required to monitor agricultural properties [10]. However, RGB cameras cannot detect phenomena invisible to the eye, such as early symptoms of plant diseases. Therefore, UAVs have been armed with thermal, multispectral, and hyperspectral cameras that see beyond the human eye [11] to calculate certain vegetation indices (NDVI, NDRE, CWSI, etc.) and diagnose diseases early [12], [13].

UAVs are also utilized to track and monitor livestock and perform search and rescue missions [14]. Alanezi *et al.* discuss tracking animals [15]. However, the number of objects involved in livestock detection is fewer than in tracking fruit production.

Robotics coupled with machine learning techniques has assisted in better segmentation of crops [16], [17]. Chen *et al.* used deep learning techniques to count apples and oranges in orchards given UAV-captured data [2]. Similar approaches have been utilized to identify apples, bananas, grapes (individually and in clusters), pears, and pineapples [18]–[23]. However, these fruits are easy to identify given their large size, and their striking color differences against their respective backgrounds (Figs. 3b, 3c, 3d). Detecting fruits of small size is even more of a challenge when the

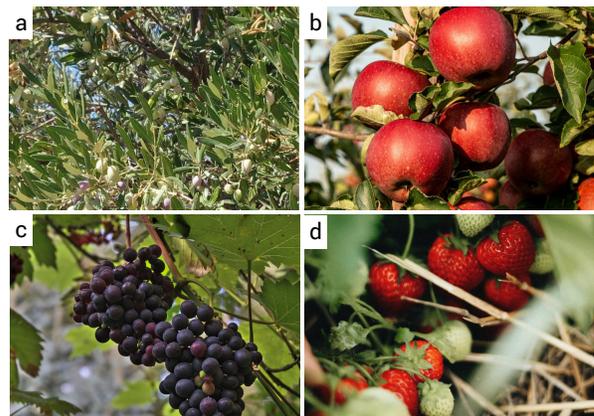


Fig. 3. (a): Olive Tree. (b): Apple Tree. (c): grape Tree. (d): Strawberry Bush. Olives are small and they have the same color as the background. Some fruits stand out among their background given their large size, and color differences.

leaves partially obstruct the fruits. Computer vision efforts to count fruits of small sizes, such as cranberries, blueberries, and cherries have been explored in [24]–[26]. Olives are of small sizes, making them more challenging to identify as we can see in Fig. 3a. Furthermore, they also blend very well into the surrounding background which has a very similar color, making detection very challenging.

The closest research to our proposed olive detection method is systems to count how many olive trees there are in a grove [27]–[29]. Ponce *et al.* proposed a direct method to count olives however only after they are harvested and sent for processing in an external image acquisition chamber [29]. Therefore, to the best of our knowledge, we propose the first autonomous on-site olive detection system.

Previous approaches utilize synthetic data to detect objects for both agricultural and non-agricultural purposes. Lin *et al.* [30] uses the 3D model of a BlueROV and creates

photo-realistic aerial image datasets with the ground truth for BlueROV detection. Sanket *et al.* [31] model the 3D geometry of a quadrotor propeller to generate a vast, automatically-labeled dataset which is then used to detect quadrotor propellers. Other works such as [32] and [33]–[35] use synthetic data to perform tomato plant and corn field reconstruction, respectively.

To the best of our knowledge, the first major breakthrough in applying synthetic data for deep learning-based fruit counting was in [36], where red circles of various sizes were scattered onto a blurred green and brown background in order to count tomatoes. This novel approach resulted in approximately 91% counting accuracy and inspired many other synthetic data approaches. Synthetic data is not just successful in fruit counting, but also in detecting oysters. Lin *et al.* [37] propose a novel open-source simulation that is used to generate photo-realistic synthetic images of oyster reefs. Using a combination of the synthetically generated oyster reef images and real images, Lin *et al.* [38] trains a semantic segmentation model resulting in a major improvement than when only trained on real data. Another synthetic data generation approach for precision agriculture is presented by Blekos *et al.* [39] where a synthetic olive tree (not containing olives) is designed and used in detecting *Verticillium Wilt*, a fungal disease that can lead to the death of an olive tree. In our paper, we create a lightweight synthetic olive tree, design a novel olive model, and train a deep semantic segmentation model for olive detection, seeing an improvement when real data is augmented by synthetically generated data.

III. SYNTHETIC OLIVE TREE

Given the arduous nature of real-world data collection and labeling for semantic segmentation, we propose a system to rapidly streamline the process of creating a rich, diverse hybrid dataset containing real and synthetic data. The need for such a swift system is highlighted further, given that each olive tree contains around 2500 to 17500 olives, on average [7]. The resulting generated synthetic data is then combined with real data to train a segmentation model to detect olives as shown in Fig. 2. In Sec. III-A.1, we illustrate the process of creating synthetic data by modeling the geometry of an olive tree utilizing Blender™, a 3D open-source graphics software to perform rendering [40]. We also ensure that our 3D model is light enough to be rendered under constrained computational resources (such as a modern-day laptop), further democratizing our synthetic data pipeline. In Sec. III-B, we explain how modern image-to-image translation methods may be applied to make the rendered synthetic images appear photorealistic.

A. Olive Tree 3D Model

1) *Synthetic Olive*: To generate synthetic data for an olive tree, we first need to create its most fundamental component: the olive. We model the geometry of an olive, by using an ellipsoid mesh as a baseline:

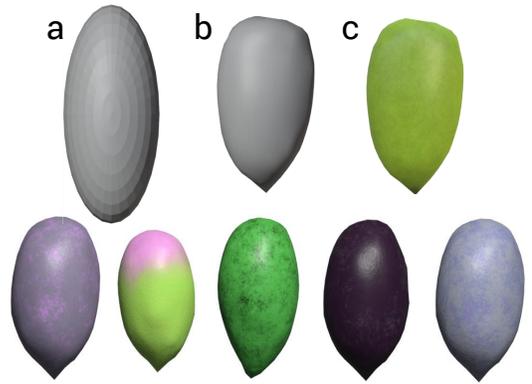


Fig. 4. (a) the 3D modeling of our ellipsoid. (b) the 3D modeling of an olive (c) an example olive rendered with texture.

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 1 \quad (1)$$

Here in Eq. 1, $a = 1.2$ cm and $b = 1.2$ cm are the width factor of the ellipsoid while $c = 2.3$ cm is the length and width of the olive as shown in Fig. 4a. We added some randomness in the scaling factor to make it more realistic. Then, the vertices near the end of the ellipsoid are extended to create a sharp point, matching the shape of a real olive.

Subsequently, a smoothing process is applied to the olive 3D model to eliminate any form of rigidness. The final dimensions of our 3D olive model are 2.5cm × 1.4cm × 1.4cm respectively as such are the approximate dimensions of olives when ready for harvest. Fig. 4b demonstrates the process of modeling our proposed olive models. We then initiate the process of adding texture to the olive model. The appearance of an olive is impacted by the lighting, shading, and albedo of the environment. All olives begin as green fruits and gradually transform into darker shades, ending in an almost black, dark purple color as the harvest season progresses (see Fig. 4c and its examples underneath). Occasionally, a single olive may have two different colors simultaneously. In addition, olives usually do not have a smooth texture but may exhibit bumps and wrinkles on their surface. Given the grove environment, it is possible for olives to accumulate dust, providing a lighter shade on certain sections of the surface. Finally, many olive groves suffer from the insect *dacus olea*, which punctures the olive causing black spot(s) to appear on its surface. We consider all the above variations when designing the texture of an olive. We model a set of eleven unique textures to be applied onto the 3D olive model, each varying in noise, roughness, distortion, bump, and color mixing. This ensures that the set of synthetic olives closely resembles real-world olives.

2) *Olive tree and leaves*: Next, we create the leaves and the branches of the olive tree. We scan a pair of leaves from a real olive tree to model the 3D geometry. The scanner automatically provides the texture of an olive leaf. To model the branches and subbranches, we utilize the bezier curve as a baseline, and is modeled as follows:



Fig. 5. Fig. 5. (a) a single bezier curve with 2 control points. (b) the bezier curve we use as a branch. (c) the bezier curve with 3D-scanned leaves scattered onto it.

$$B(t) = \sum_{i=0}^n \binom{n}{i} (1-t)^{n-i} t^i P_i \quad (2)$$

where t represents the position of the point on the curve and P_i represents the i^{th} control point. In Eq. 2 we select $n = 5$ control points per branch. We then prune the ends of the curves in Blender to make them more realistic. To match the texture of an olive tree’s branches onto our 3D model’s branches, we use a dark brown as our base color and apply dark noise to the branches. This process is illustrated in Fig. 5.

In an effort to make our 3D model more accessible to those without high-end computing devices, we decide to make a lighter version of an olive tree model without a trunk or a dense canopy. *Choosing to do so drastically reduces the average time needed for rendering synthetic image-mask pairs.* First, we eliminate the trunk. Then, we reduce the number of leaves and branches required by using an invisible 2D plane as a baseline for their scattering. Finally, to address the excessively high number of vertices on the 3D-scanned leaf, we apply a limited dissolve onto the leaf model for simplification, further democratizing (making it easily accessible for all) our 3D model.

To complete the model, we first create four separate layers of scattered leaves, two separate layers of scattered branches, and one layer of scattered olives per rendering session in the formation of 2D planes. The randomness of leaf orientation is set to no more than 9° and the size randomness is set to be no more than 10% of the original dimensions. When scattering the olives onto the leaves and branches, we place them behind one leaf layer in order to simulate the obscurity that real olive leaves provide to the olive fruit. We apply orientation randomness within 45° and size randomness within 5% of the original dimensions to all olives. In total, we scatter approximately 2,600 olives onto the leaves per rendering session. In total, we run 16 rendering sessions to produce 15,960 image-mask pairs. Example renderings are shown in Fig. 6c. Since the Mediterranean olive harvest occurs during the months of mid-Autumn, we emulate appropriate background and lighting conditions for our 3D model. We first simulate dry, sunny days with little to no clouds using Blender’s sky, Musgrave, and gradient texture nodes. We then simulate a cloudy, overcast environment using Blender’s dynamic sky. Finally, we place a plane behind the final leaf layer with ground texture projected onto it via UV mapping to simulate a UAV flying above a canopy with the ground as the background. We also UV-map an image of scattered leaves onto a background plane to simulate a UAV flying at the height of the canopy.

3) *Rendering:* At this point, the 3D olive tree model is ready to be rendered into image-mask pairs, a very computationally demanding task. Therefore, we take the following steps to shorten the rendering process, expanding the dataset’s availability to computers with lower processing power (such as modern laptops).

- We utilize Intel OpenImageDenoise [41], an open-source library of denoising filters made for images rendered with ray tracing
- We lower the maximum number of samples in Cycles render (number of paths to trace for each pixel in the final render) to 50, still producing a high-quality output
- We eliminate all subsurface scattering in all objects

B. Image-to-Image Translation

Image-to-image translation is a technique that converts an image from a source domain to a target domain (e.g. convert an image of a horse to a zebra) [42] [43]. Typically, synthetic data does not possess the realistic features of the real data domain. Models trained on untranslated synthetic data perform poorly on real-world testing data. Therefore, we conduct image-to-image translation to transform synthetically rendered images in Blender of our 3D model to the real-world olive tree domain. The objective is that training with synthetic data closely resembling the real world will result in higher accuracy models. Specifically, we employ Unpaired Image Translation via Vector Symbolic Architectures (VSAIT) [44] to translate the Blender-generated images to the real-world domain. We also choose VSAIT in order to reduce occurrences of semantic flipping (when a green olive gets translated into a leaf, corrupting the translated data) given that VSA-based methods are capable of learning high-level, abstract concepts [45].

Formally, VSAIT is trained to learn the translation between the Blender-generated synthetic olive tree images X to the real olive grove images Y . The three major components of this network are the Generator G , the Discriminator D_Y , and the Source \leftrightarrow Target Mapper F . The overall loss consists of the sum of the Hypervector Adversarial Loss and the VSA-based Cyclic Loss [44].

Hypervector Adversarial Loss: This ensures the similarity between the synthetic translated image hypervectors $v_{G(X)}$ and the real target hypervectors v_Y . This is done by applying the VSA binding operation on hypervector mapping $F(v_x)$ and source vectors v_X to yield a hypervector of synthetic features mapped to the real domain. The Hypervector Adversarial Loss is calculated as:

$$\begin{aligned} \mathcal{L}_{GAN}(G, D_Y, X, Y) = & \mathbb{E}_{y \sim p_Y(y)} [\log D_Y(v_y)] \\ & + \mathbb{E}_{x \sim p_X(x)} [\log(1 - D_Y(v_{G(x)}))] \\ & + \mathbb{E}_{x \sim p_X(x)} [\log(1 - D_Y(v_x \otimes F(v_x)))] \quad (3) \end{aligned}$$

VSA-based Cyclic Loss: This cyclic loss minimizes occurrences of semantic flipping by constraining G so that similar hypervectors may be returned when mapping translated vectors back to the synthetic domain from the

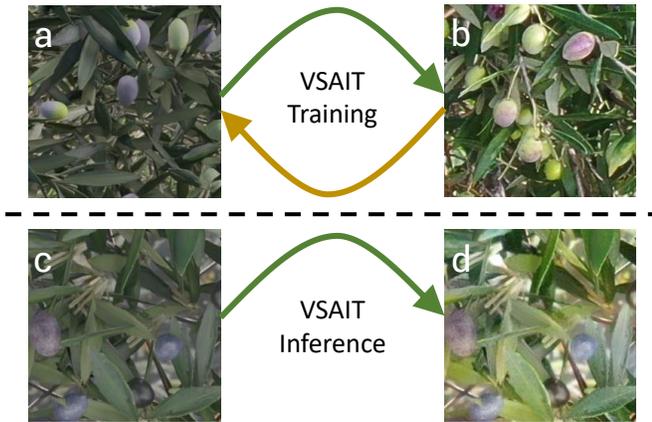


Fig. 6. Image-to-image translation (a) A single sample of the simulated olive, (b) A single sample of the real olive, (c) Simulated olive image, (d) Synthetic olive image translated into real world for photorealism.

real domain. Meaning, $v_X \approx v_G(X) \rightarrow X$. The VSA-based Cyclic Loss is calculated as:

$$\mathcal{L}_{VSA}(G, X) = \mathbb{E}_{x \sim p_X(x)} \left[\frac{1}{n} \sum_{i=1}^n \text{dist}(v_x^i, v_{G(x \rightarrow x)}^i) \right] \quad (4)$$

where the dist is the cosine distance between the source hypervectors and inverted translated hypervectors.

Therefore, the overall loss is as follows:

$$\mathcal{L}(G, D_Y, X, Y) = \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{VSA}(G, X) \quad (5)$$

To create X , the source domain, we render 12,768 images of our 3D olive tree model for training and 3,192 images as validation. To create Y , the target domain, we use 12,452 UAV-captured images for training and 3,113 for validation. Overall, we train VSAIT for 98 epochs using its default settings (see Fig. 6a, 6b). Then, we perform inference to apply VSAIT to rendered images (Fig. 6c), yielding a photorealistic synthetic output (Fig. 6d).

We ensure that the images used in the Y domain for image-to-image translation do not overlap with any of the real images that will later be used in testing segmentation models. When X gets translated to Y , it learns the features of Y . Training a segmentation model using synthetic images translated to the same Y as the real-world test set is just like training with testing data as input (and must be avoided).

C. Synthetic Data Color Input Space

After the synthetic data undergoes image-to-image translation, the olives in the synthetic data still stand out brightly, unlike in real olive groves. Furthermore, the synthetic data translated by VSAIT can easily still be distinguished from images of the real domain, if looked at closely. Therefore, in an attempt to make the olives in the VSAIT output less distinguishable (just like in the real-world), we create a new input space, with the goal of making the olives more concealable. We draw inspiration from [46] who propose a novel input space from $(RGB \rightarrow R, M=\max(G,B), I)$ where I is the overall intensity. However,

Yu *et al.* [46] apply this technique to the underwater domain where there is a deficiency of red reflections to extract a clearer image for depth estimation. We propose a color input space for the opposite purpose: to obscure the olives further into their backgrounds. This new input space may be applied to all synthetic agricultural domains with the purpose of blending any small fruit into its respective backgrounds and is as follows (called IGA):

$$RGB \rightarrow (I, G, \text{Avg}(R, B))$$

Converting the VSAIT-translated synthetic data to this input space yields a slight improvement as shown in Table I and is labeled as $O_{IGA.and.real}$.

IV. EXPERIMENTS AND RESULTS

We first describe our datasets. Then we analyze our segmentation experiments on models trained on images with and without synthetically generated data.

A. Dataset Overview

Synthetic Data: As described in Section III, we initiate multiple rendering processes of our Blender scene, each reproducing parts of an olive tree with a blue sky, leaf-only, or ground texture background. In addition, we utilize multiple lighting conditions to augment the variety of our synthetic data. Afterward, we initiate the image-to-image translation process [44]. The result of these two processes is the speedy generation of 15,960 synthetic image-mask pairs.

Real Data: We conducted several UAV flights over the island of Crete, Greece. Our UAV is the DJI Mavic 2 Pro with a Hasselblad L1D-20C camera. To make the dataset as diverse as possible, the flights were staggered to record differences in cloud cover and albedo, and to capture images of olives during different temporal stages of the olive harvest, approximately ranging from one to two months. We also chose to only set our ISO to 100 for daytime capture and 400 for evening capture. We then labeled our UAV-captured images by hand using LabelMe [47]. In total, we have 3,113 real-world image-mask pairs in patches of size $256 \times 256 \times 3$ and $256 \times 256 \times 1$ and we labeled approximately 2,600 olives. The maximum number of olives per unpatched image is around 600.

B. Experiments

We construct two training sets: one with 15,960 synthetic and 100 real image-mask pairs and another with only 100 real image-mask pairs. This was decided to reflect the reality that it is very difficult to manually create a large enough labeled olive dataset. We select UNet [48] as our model architecture and train with several backbones [49], [50] as shown in Table I. To standardize experiments across our two data sets, we only employ the Adam optimizer [51] with an initial learning rate of 0.001 and use the Jaccard index [52] as our loss function. All models are trained for a maximum of 100 epochs using a batch size of 32. We use minimal real-world data as it is arduous to obtain and label, simulating

TABLE I

IoU USING REAL ONLY, SYNTHETIC/REAL, AND SYNTHETIC IGA/REAL

Training Data	Backbone	IoU
O_{real}	EfficientNetB5	41.41%
$O_{syn.and.real}$	EfficientNetB5	53.06%
$O_{IGA.and.real}$	EfficientNetB5	54.05%
O_{real}	ResNet101	24.22%
$O_{syn.and.real}$	ResNet101	36.28%
$O_{IGA.and.real}$	ResNet101	40.22%
O_{real}	ResNet152	28.99%
$O_{syn.and.real}$	ResNet152	40.54%
$O_{IGA.and.real}$	ResNet152	45.33%

situations where one may only have limited access to real data, which are more realistic cases.

Our testing set consists of more than 3000 real-world image-mask pairs that are hand-labeled [47]. Our metric of choice is the Intersection over Union, as there is only one class of objects to segment. We opt out of using mean Average-Precision as an evaluation metric given that models trained on mostly synthetic data perform well only with low-confidence score thresholding, as is in our case.

C. Results

Based on our experiments, as shown in Table. I, we observe that the models with identical architectures and backbones trained on both synthetic and real images outperform the baseline models, which are trained without synthetic data. In our *Training Data* column, O_{real} represents the set of just real-world training data, $O_{syn.and.real}$ represents the set of real-world and VSAIT-translated synthetic data, and $O_{IGA.and.real}$ represents the VSAIT-translated synthetic data converted to the IGA input space. Our biggest marked improvement, when compared to the baseline, was observed from 24.22% to 40.22% IoU, corresponding to a percent change of 66% when using ResNet101 [49] as our backbone. We also observe that when using EfficientNetB5 [50], ResNet101, and ResNet152, the highest-performing models are trained on VSAIT-translated images converted to the IGA input space. Fig. 1 shows the difference in segmentation prediction when using models trained with and without synthetic data. When using an EfficientNetB5 backbone, we observe a 12.64% jump in IoU over the baseline when the model is supplemented with synthetic data in the IGA input space. Similarly, when we use Resnet152 as a backbone, we observe an increase of 16.34% compared to the baseline when supplemented with unedited VSAIT-generated data.

D. Discussion

Detecting olives in an olive grove is a problem in the low training data regime domain. Under this setting, it is not convenient to label vast amounts of olive tree images and hence synthetic data offers a reasonable alternative. In this paper, we simulate such a situation and have demonstrated that supplementing a low amount of real data with easy-to-generate synthetic images improves the segmentation results by 66% as shown in Table I.

V. CONCLUSION AND FUTURE WORK

In this paper, we presented a novel approach to detect olives in an olive grove without the need to manually label data. We described a technique to generate a photorealistic 3D model of an olive tree. Its geometry is then simplified for lightweight rendering purposes, thus democratizing the synthetic dataset. The goal is to supplement low levels of real data with our synthetic data to improve detection results. To evaluate our approach, we conducted experiments exhibiting promising results with a maximum observed change in IoU of 66%, when using our proposed method versus low-data regime methods. To conclude, this demonstrates that when access to real, labeled data is restricted or non-attainable, a combination of mostly synthetic data and real data can help to enhance olive detection.

ACKNOWLEDGMENTS

This work was partially sponsored by the Fulbright Foundation and USDA NIFA Award# 20206801231805. We thank all olive grove farmers, Dr. Panagiotis Partsinevelos' SenseLab, and Zisis Charokopos for data collection assistance.

REFERENCES

- [1] F. B. Insights, "Olive oil market size, trends amp; growth: Global forecast [2029]," Sep 2022. [Online]. Available: <https://www.fortunebusinessinsights.com/industry-reports/olive-oil-market-101455> 1
- [2] S. W. Chen, S. S. Shivakumar, S. Dcunha, J. Das, E. Okon, C. Qu, C. J. Taylor, and V. Kumar, "Counting apples and oranges with deep learning: A data-driven approach," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 781–788, 2017. 1, 2
- [3] M. H. Junos, A. S. Mohd Khairuddin, S. Thannirmalai, and M. Dahari, "An optimized yolo-based object detection model for crop harvesting system," *IET Image Processing*, vol. 15, no. 9, pp. 2112–2125, 2021. 1
- [4] M. T. Pratama, S. Kim, S. Ozawa, T. Ohkawa, Y. Chona, H. Tsuji, and N. Murakami, "Deep learning-based object detection for crop monitoring in soybean fields," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–7. 1
- [5] X. Chen, Z. Li, Y. Yuan, G. Yu, J. Shen, and D. Qi, "State-aware tracker for real-time video object segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9384–9393. 1
- [6] T. Liu, N. Chopra, and J. Samtani, "Information system for detecting strawberry fruit locations and ripeness conditions in a farm," in *Biology and Life Sciences Forum*, vol. 16, no. 1. MDPI, 2022, p. 22. 1
- [7] P. Dz, "How many olives are harvested on an average olive tree? (detailed answer)," Jan 2022. [Online]. Available: <https://oliveknowledge.com/how-many-olives-are-harvested-on-olive-tree/> 1, 3
- [8] R. Elijah, "Combating climate change: How precision agriculture can help," Feb 2023. [Online]. Available: <https://eos.com/blog/how-precision-farming-fights-climate-change/> 2
- [9] A. Balafoutis, B. Beck, S. Fountas, J. Vangeyte, T. Van der Wal, I. Soto, M. Gómez-Barbero, A. Barnes, and V. Eory, "Precision agriculture technologies positively contributing to ghg emissions mitigation, farm productivity and economics," *Sustainability*, vol. 9, no. 8, p. 1339, 2017. 2
- [10] G. Livanos, D. Rammalis, V. Polychronos, P. Balomenou, P. Sarigiannidis, G. Kakamoukas, T. Karamitsou, P. Angelidis, and M. Zervakis, "Extraction of reflectance maps for smart farming applications using unmanned aerial vehicles," in *2020 12th International Symposium on Communication Systems, Networks and Digital Signal Processing (CSNDSP)*. IEEE, 2020, pp. 1–6. 2

- [11] T. Adão, J. Hruška, L. Pádua, J. Bessa, E. Peres, R. Morais, and J. J. Sousa, "Hyperspectral imaging: A review on uav-based sensors, data processing and applications for agriculture and forestry," *Remote sensing*, vol. 9, no. 11, p. 1110, 2017. 2
- [12] T. Poblete, J. Navas-Cortes, C. Camino, R. Calderon, A. Hornero, V. Gonzalez-Dugo, B. Landa, and P. Zarco-Tejada, "Discriminating xylella fastidiosa from verticillium dahliae infections in olive trees using thermal- and hyperspectral-based plant traits," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 179, pp. 133–144, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0924271621001994> 2
- [13] M. Pantelidakis, A. A. Panagopoulos, K. Mykoniatis, S. Ashkan, R. C. Eravi, V. Pamula, E. C. Verduzco III, O. Babich, O. P. Panagopoulos, and G. Chalkiadakis, "Identifying sunlit leaves using convolutional neural networks: An expert system for measuring the crop water stress index of pistachio trees," *Expert Systems with Applications*, vol. 209, p. 118326, 2022. 2
- [14] D. Chatziparaschis, M. G. Lagoudakis, and P. Partsinevelos, "Aerial and ground robot collaboration for autonomous mapping in search and rescue missions," *Drones*, vol. 4, no. 4, p. 79, 2020. 2
- [15] M. A. Alanezi, M. S. Shahriar, M. B. Hasan, S. Ahmed, Y. A. Sha'aban, and H. R. Boucheqara, "Livestock management with unmanned aerial vehicles: A review," *IEEE Access*, 2022. 2
- [16] J. Champ, A. Mora-Fallas, H. Gočau, E. Mata-Montero, P. Bonnet, and A. Joly, "Instance segmentation for the fine detection of crop and weed plants by precision agricultural robots," *Applications in plant sciences*, vol. 8, no. 7, p. e11373, 2020. 2
- [17] D. Su, H. Kong, Y. Qiao, and S. Sukkarieh, "Data augmentation for deep learning based semantic segmentation and crop-weed classification in agricultural robotics," *Computers and Electronics in Agriculture*, vol. 190, p. 106418, 2021. 2
- [18] N. Häni, P. Roy, and V. Isler, "Minneapolis: a benchmark dataset for apple detection and segmentation," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 852–858, 2020. 2
- [19] A. I. B. Parico and T. Ahamed, "Real time pear fruit detection and counting using yolov4 models and deep sort," *Sensors*, vol. 21, no. 14, p. 4803, 2021. 2
- [20] R. Wan Nurazwin Syazwani, H. Muhammad Asraf, M. Megat Syahirul Amin, and K. Nur Dalila, "Automated image identification, detection and fruit counting of top-view pineapple crown using machine learning," *Alexandria Engineering Journal*, vol. 61, no. 2, pp. 1265–1276, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S111001682100418X> 2
- [21] B. Neupane, T. Horanont, and N. D. Hung, "Deep learning based banana plant detection and counting using high-resolution red-green-blue (rgb) images collected from unmanned aerial vehicle (uav)," *PloS one*, vol. 14, no. 10, p. e0223906, 2019. 2
- [22] A. K. Nellithimaru and G. A. Kantor, "Rols: Robust object-level slam for grape counting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0. 2
- [23] T. T. Santos, L. L. de Souza, A. A. dos Santos, and S. Avila, "Grape detection, segmentation, and tracking using deep neural networks and three-dimensional association," *Computers and Electronics in Agriculture*, vol. 170, p. 105247, 2020. 2
- [24] P. Akiva, K. Dana, P. Oudemans, and M. Mars, "Finding berries: Segmentation and counting of cranberries using point supervision and shape priors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 50–51. 2
- [25] S. Gonzalez, C. Arellano, and J. E. Tapia, "Deepblueberry: Quantification of blueberries in the wild using instance segmentation," *Ieee Access*, vol. 7, pp. 105 776–105 788, 2019. 2
- [26] J. F. Villacrés and F. Auat Cheein, "Detection and characterization of cherries: A deep learning usability case study in chile," *Agronomy*, vol. 10, no. 6, p. 835, 2020. 2
- [27] J. González, C. Galindo, V. Arevalo, and G. Ambrosio, "Applying image analysis and probabilistic techniques for counting olive trees in high-resolution satellite images," in *Advanced Concepts for Intelligent Vision Systems: 9th International Conference, ACIVS 2007, Delft, The Netherlands, August 28-31, 2007. Proceedings 9*. Springer, 2007, pp. 920–931. 2
- [28] A. Khan, U. Khan, M. Waleed, A. Khan, T. Kamal, S. N. K. Marwat, M. Maqsood, and F. Aadil, "Remote sensing: an automated methodology for olive tree detection and counting in satellite images," *IEEE Access*, vol. 6, pp. 77 816–77 828, 2018. 2
- [29] J. M. Ponce, A. Aquino, B. Millan, and J. M. Andujar, "Automatic counting and individual size and mass estimation of olive-fruits through computer vision techniques," *IEEE Access*, vol. 7, pp. 59 451–59 465, 2019. 2
- [30] X. Lin, C. Liu, A. Pattillo, M. Yu, and Y. Aloimonos, "Seadronesim: Simulation of aerial images for detection of objects above water," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 216–223. 2
- [31] N. J. Sanket, C. D. Singh, C. M. Parameshwara, C. Fermüller, G. C. de Croon, and Y. Aloimonos, "Evpropnet: Detecting drones by finding propellers for mid-air landing and following," *arXiv preprint arXiv:2106.15045*, 2021. 3
- [32] A. K. Burusa, E. J. van Henten, and G. Kootstra, "Attention-driven active vision for efficient reconstruction of plants and targeted plant parts," *arXiv preprint arXiv:2206.10274*, 2022. 3
- [33] A. Bacharis, H. J. Nelson, and N. Papanikolopoulos, "View planning using discrete optimization for 3d reconstruction of row crops," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 9195–9201. 3
- [34] H. Freeman, E. Schneider, C. H. Kim, M. Lee, and G. Kantor, "3d reconstruction-based seed counting of sorghum panicles for agricultural inspection," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9594–9600. 3
- [35] S. Khaki, H. Pham, Y. Han, A. Kuhl, W. Kent, and L. Wang, "Deepcorn: A semi-supervised deep learning method for high-throughput image-based corn kernel counting and yield estimation," *Knowledge-Based Systems*, vol. 218, p. 106874, 2021. 3
- [36] M. Rahnemoonfar and C. Sheppard, "Deep count: fruit counting based on deep simulated learning," *Sensors*, vol. 17, no. 4, p. 905, 2017. 3
- [37] X. Lin, N. Jha, M. Joshi, N. Karapetyan, Y. Aloimonos, and M. Yu, "Oystersim: Underwater simulation for enhancing oyster reef monitoring," in *OCEANS 2022, Hampton Roads*. IEEE, 2022, pp. 1–6. 3
- [38] X. Lin, N. J. Sanket, N. Karapetyan, and Y. Aloimonos, "Oysternet: Enhanced oyster detection using simulation," 2022. [Online]. Available: <https://arxiv.org/abs/2209.08176> 3
- [39] K. Blekos, A. Tsakas, C. Xouris, I. Evdokidis, D. Alexandropoulos, C. Alexakos, S. Katakis, A. Makedonas, C. Theoharatos, and A. Lalos, "Analysis, modeling and multi-spectral sensing for the predictive management of verticillium wilt in olive groves," *Journal of Sensor and Actuator Networks*, vol. 10, no. 1, p. 15, 2021. 3
- [40] B. O. Community, *Blender - a 3D modelling and rendering package*, Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. [Online]. Available: <http://www.blender.org> 3
- [41] Intel, "Intel® open image denoise." [Online]. Available: <https://www.openimagedenoise.org/> 4
- [42] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017. 4
- [43] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *CVPR*, 2017. 4
- [44] J. Theiss, J. Leverett, D. Kim, and A. Prakash, "Unpaired image translation via vector symbolic architectures," in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXI*. Springer, 2022, pp. 17–32. 4, 5
- [45] P. Kanerva, *Sparse Distributed Memory*. Cambridge, MA, USA: MIT Press, 1988. 4
- [46] B. Yu, J. Wu, and M. J. Islam, "Udepth: Fast monocular depth estimation for visually-guided underwater robots," *arXiv preprint arXiv:2209.12358*, 2022. 5
- [47] K. Wada, "Labelme: Image Polygonal Annotation with Python." [Online]. Available: <https://github.com/wkentaro/labelme> 5, 6
- [48] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*. Springer, 2015, pp. 234–241. 5
- [49] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 5, 6
- [50] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for

- convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114. [5](#), [6](#)
- [51] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014. [5](#)
- [52] P. Jaccard, “The distribution of the flora in the alpine zone. 1,” *New phytologist*, vol. 11, no. 2, pp. 37–50, 1912. [5](#)