# A Robust Illumination-Invariant Camera System for Agricultural Applications

Abhisesh Silwal, Tanvir Parhar, Francisco Yandun and George Kantor

*Abstract*—Object detection and semantic segmentation are two of the most widely adopted deep learning algorithms in agricultural applications. One of the major sources of variability in image quality acquired in the outdoors for such tasks is changing lighting condition that can alter the appearance of the objects or the contents of the entire image. While transfer learning and data augmentation to some extent reduce the need for large amount of data to train deep neural networks, the large variety of cultivars and the lack of shared datasets in agriculture makes wide-scale field deployments difficult. In this paper, we present a high throughput robust active lighting-based camera system that generates consistent images in all lighting conditions. We detail experiments that show the consistency in images quality leading to relatively fewer images to train deep neural networks for the task of object detection. We further present results from field experiment under extreme lighting conditions where images without active lighting significantly lack to provide consistent results. The experimental results show that on average, deep nets for object detection trained on consistent data required nearly four times less data to achieve similar level of accuracy. This proposed work could potentially provide pragmatic solutions to computer vision needs in agriculture.

## I. INTRODUCTION

Identifying fruits, vegetables, and predicting plant diseases early just from images and using that information as a layer in the decision-making process in the production cycle has the potential to revolutionize agricultural industry. Predictive tasks in agriculture such as yield estimation or detection of disease outbreak has always relied on manual workers which quickly becomes a cumbersome task as the size of the orchards and vineyards get larger. Additionally, researchers often describe the computer vision-based approach to agricultural automation as an efficient, low-cost, non-destructive, and scalable process. The combination of above-mentioned advantages and demand from the industry has generated an enormous amount of interest in machine learning-based approach to solving agricultural computer vision problems.

Detecting fruits and vegetables in images is also a fundamental requirement for any image-based automation task in agricultural fields [1]. It is often regarded as the critical factor for the success of Ag autonomous systems as it provides the ability to perceive information necessary to generate appropriate actions for tasks such as controlling a robot arm for harvesting and autonomous navigation in row crops [2] to name few. A typical computer vision pipeline in fruit detection or semantic segmentation is to identify and localize individual fruit/ fruit pixels within the image. Although, this process has received much attention, getting consistent result and generalizing across different field conditions have proven to be extremely difficult. Numerous publications in the past three

decades in vision-based agricultural robotics have regarded retrieving visual information from images as one of the major bottle necks to reach commercial maturity [3]. The quest for robust and more generalized fruit detection algorithm has received renewed focus in the research community.

With the advent of deep neural networks in Ag research communities, object detection and semantic segmentation tasks have seen great success in recent history [4]. In Ag computer vision, most of object detectors [5] [6] are trained using transfer learning. However, a significant bottleneck exists even for fine-tuning because of lack of data that can be attributed to the vast varieties of cultivars in ag industry. To name few, in the U.S. alone 100 different varieties of apples are grown commercially [7] and more than 415 variety of invasive weed species [8] infest agricultural sites. Sampling a large number of images for each cultivar and its countless sub varieties and maintaining precisely annotated datasets could overwhelm data collection and sharing efforts. Additionally, ag-specific domain expert oversights are often needed when ground truthing images for common tasks such as weed classification, disease detection, etc.

Images in agricultural sites are taken mainly in outdoor environment with changing lighting conditions which many researchers [9] often cite as a limiting factor in computer vision application in Ag. Correct image exposure is critical to preserving the perceptual quality of the images [10]. Regardless of the computer vision algorithm of choice (classical or machine learning- based), features in images are important factors. If the image is under or over-exposed, features and textures can become undetectable that could negatively impacts the performance of the computer vision algorithms. Thus a robust system to collect images with consistent quality in all lighting condition could provide robust solution.

The contributions of this paper are in the design and evaluation of a high throughput camera system that consistently produces image invariant to changes in environmental lightning. We hypothesize that the consistency in images reduces the amount of data required for training (fine-tuning) deep neural networks in agricultural applications. This paper is organized in the following way. Section II reviews relevant literature in camera systems and the use and implementation of deep learning algorithms for object detection in ag. The hardware design of the camera system is described in Section III which also provides additional details on field experiments and datasets collected and used in this paper. Then, Section IV compares the performance of various trained models using active light (AL) and natural light (NL). Finally, Section V includes concluding remarks and directions for further

research. Throughout this paper, we use the acronyms AL and NL to describe datasets and lighting environments.

## II. RELATED WORK

### A. Object detection in agriculture

One key aspect of the continuous success in the integration of agriculture and technology is the application of image processing and data analysis techniques for detecting objects within the crops canopy. Then, tasks such as yield estimation, anomaly evaluation or harvesting could be performed in a more cognizant and efficient way [11]. In the past, the images analysis required the use of significant amount of traditional vision computer techniques (e.g., erosion, dilation, contour detection) to detect and/or count objects of interest such as fruits, workers or their tools [12], [13]. For example, a radial symmetry transform and texture detection techniques were used in [14] to detect and count berries in RGB images, obtaining a $R^2$ correlation of up to 0.95.

Despite the successful results of traditional computer vision techniques in some applications, there are others that are extremely difficult to design/engineer due to the challenging characteristics of the agricultural environments (e.g., changing illumination, blurring) [9]. In these conditions, deep learning approaches have proven to be an useful tool for image classification and object detection [15]. Activities like leaf disease detection, weed/fruit counting, or plant type detection have been robustly performed in different conditions and scenarios [16]–[18]. For example, a CNN with a custom architecture was used in [19] to count oranges apples in a cluttered scene, obtaining accuracies up to 0.957 and 0.961, respectively. In contrast to using custom architectures, other authors employ known networks (or slightly modified versions) such as Faster-RCNN , Yolov3, Single Shot MultiBox Detector, among others [20]–[22].

The proliferation of large scale datasets like [23], [24], [25] and [26] has been a major factor in the recent state of the art results in object detection. But these large scale datasets are generally restricted to generic objects like cars, pets etc. While there is a body of image datasets for agricultural perception tasks, they are generally collected in well curated laboratory conditions [27], [28], [29], [30]. For robust computer vision in field settings, there is need for datasets that collected under more realistic field conditions. And, if robustness to varying lighting conditions needs to be tackled, there is a need for a dataset with that contains images representing all the varying field conditions [31].

As computer vision continues to become an important aspect of precision agriculture, to further advance research a variety of datasets have been released either through publications or independently. For example, [5] consists of 587 bounding boxes for various fruits in RGB images, [32] consists of more than 10000 images for Date fruit yield estimation, [6] consists of 49 images with pixel-wise labels for mango segmentation and [33] contains 3704 images of orchards fruits like apples and almonds. While this is a much anticipated and encouraging direction, the challenge of collecting datasets comparable to the size of [24] in agricultural cultivar has always remained illusive.

Robotics or computer vision-based autonomy in agriculture needs to be resilient to outdoor environment. Collecting data to train deep networks that represents all possible modalities of cultivar in Ag is an onerous challenge. Thus, a camera system agnostic to external lighting conditions, that can generate consistent image data could benefit Ag industry. The dataset collected with the camera system in this study will be publicly shared.

### B. Imaging sensors and systems in agriculture

A review article by [9] shows that the most popular choice of sensors for fruit detection is mostly color cameras. Nearly two dozen cited articles in [9] use color cameras in various fruit detection tasks. Other choice of sensors include hyperspectral, thermal, monochrome camera and in some occasions laser range finders. Remarkably, all research work using color camera in [9] acknowledge variation in lighting as the limiting factor. An alternative solution often used by researchers to control lighting environments is to use a mobile structure to completely block sun light to image the canopy [34] [1]. In a recent work, [35] used high resolution stereo sensor with flash to predict yield in vineyards from image-based counting for grapes. The use of flash imagery in this study generated images with uniform white balance and had minimal effects from natural illumination. Our design of the camera system in this paper is motivated by this work. To our best knowledge, the affect of consistent image quality on the size of training data for deep neural networks in agricultural has not been previously reported.

## III. METHODS

A camera essentially converts incident irradiant light energy into discrete pixel values. The conversion of the irradiant energy to brightness levels in pixel includes transforming the incident photons to electrical charges and then quantize the analog signal to digital data. Eqn.(1) from [36] shows the expression of the signal received by a pixel $P$ .

$$P = TK \frac{l^2}{N^2} \int \frac{E_\lambda R(\lambda)}{\frac{hc}{\lambda}} Q(\lambda) d\lambda \qquad (1)$$

where, $T$ is the exposure time, $K$ is a normalization constant, $l$ is the pixel pitch, $N$ is the f/number of the lens and $\lambda$ is the wavelength. The terms $Q(\lambda)$, $E_\lambda$ and $R(\lambda)$ are the spectral quantum efficiency, spectral power distribution and spectral reflectance respectively, which are dependent on the wavelength of the incident light. The remaining values $h$ and $c$ are the Plank's constant and the velocity of light. This analog value is then quantized to be represented into digital formats.

Without flash or active lighting, the terms $E_\lambda$ and $R(\lambda)$ inside the integral of Eqn.(1) are variables that depend on the intensity of natural light captured by the sensor. As natural light intensity (brightness) and wavelength (color) change drastically throughout the day, the change in the irradiant energy entering the camera sensor could alter the perceptual

TABLE I: Camera hardware components and specifications

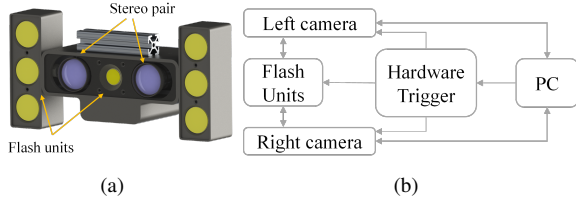| Hardware | Specifications |
| --- | --- |
| Camera | PointGrey CM3, 3.2 MP, Color, global shutter |
| Lens | 3.5mm f/2.4, 89°x73.8°x101.7° FoV, manual aperture and focus |
| Flash | 100 Watts side LEDx6, 500 Watts center LEDx1. 5600K color temp |
| Flash Trigger | 200 - 250 μs |
| Shutter Speed | as low as 11 μs |
| Acquisition Speed | 1 to 20 Hz synchronized stereo pair |



Fig. 1: Camera system (a). Image and flash synchronization flowchart (b)



Fig. 2: Image of an outdoor scene with (a) and without flash (b). The snippets next to the main image shows image formation at various times of the day.

quality of object of interest in the image. This unpredictable change in the quality of the images taken in outdoors increases the variability in the image as a data point in computer vision applications. However, with active lighting, $E_\lambda$ and $R(\lambda)$ in Eqn.(1) become constants just leaving exposure time $T$, digital gain, and aperture to control image exposure.

### A. Hardware Design

The proposed camera system shown in Figure 1 is a custom-built rugged stereo pair. This field prototype houses two industrial color cameras with spatial resolution of 2048x1536 pixels and a stereo baseline of 110 mm. The featured design is modular and enables quick changes in the orientation and quantity of flash units as required. As mentioned in Section I, the use of active lighting potentially reduces variation in image data. However, in practice the generation of such images is only possible when the flash duration is tightly synchronized with image acquisition pipeline. To achieve such synchronization, an external hardware trigger was necessary to send a control signal to all the cameras and flash units simultaneously. The synchronous signal was generated using a micro-controller on-board connected to the host PC.

The active lighting system described here uses a high-power Chip on Board (COB) LED lights (1.2 K Watts total) that floods the scene with a bright pulse of white light with color temperature of 5600K. The use of the high-power flash allowed us to set the digital gain of the camera sensor to zero dB, further reducing variables in the control of image exposure while avoiding inducing digital noise in the acquired images. In all experiments, the aperture of the lens was kept fixed at the minimum value of f/2.4 that left exposure time $T$ as the only tuneable variable. As the total number of variables in Eqn.(1) decreases, better control of image exposure is achievable. Additional hardware details are listed in Table I.

### B. Imaging in the outdoors

The aim of the first experiment under this section is to quantify the quality of images taken in an outdoor setting with and without active lighting. In this experiment, the camera system described in Section III-A imaged an outdoor scene (Figure 2) every twenty-minute interval from noon to sunset. First, an image of the scene was captured using the camera auto exposure setting and was immediately followed by the active light imaging. The camera remained in a fixed position throughout the experiment to capture the exact scene for accurate comparison. Additionally, a light meter was used to manually measure the luminance of the scene for all the intervals. For the non-active lighting images, auto exposed and HDR image sets were also collected for comparison.

The Structural Similarity (SSIM) index proposed in [37] and defined by Eqn.(2), and the Peak Signal to Noise Ratio (PSNR) are two of the metrics used to quantify the quality of the captured images in this paper. The SSIM is a quality assessment index which is a multiplicative combination of luminance, contrast, and structural terms. Where as PSNR computed the peak signal to noise ratio in images. Both SSIM and PSNR used the first image taken at noon as reference. The results of this experiment are detailed in Section IV.

$$SSIM(x,y) = f\left(l\left(\mathbf{x},\mathbf{y}\right), c\left(\mathbf{x},\mathbf{y}\right), s\left(\mathbf{x},\mathbf{y}\right)\right) \quad (2)$$

With,

$$l\left(\mathbf{x}, \mathbf{y}\right) = \frac{2\left(1+R\right)}{1 + \left(1+R\right)^2 + \frac{C_1}{\mu_x^2}}$$

$$c\left(\mathbf{x}, \mathbf{y}\right) = \frac{2\sigma_x\sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2},$$

$$s\left(\mathbf{x}, \mathbf{y}\right) = \frac{\sigma_{xy} + C_3}{\sigma_x\sigma_y + C_3}$$

where $\mu_x, \mu_y, \sigma_x, \sigma_y$, and $\sigma_{xy}$ are the local means, standard deviations, and image cross-covariance, and $C_1$, $C_2$, $C_3$ are constants.

The second experiment under this section describes the details of the camera system usage in field. Concretely, it was employed in a commercial apple orchard to generate the dataset used in this paper. Similarly to the first experiment, the apple images were collected with and without the flash. The camera system was mounted on an orchard platform that was manually driven in-between the row-space at approximately 3 mph (0.4 m/s) speed while imaging the tree canopies. The platform was driven twice in the same row to image the fruit with and without active lighting.

To demonstrate that uniformly exposed images (i.e., consistent in quality) allows a competent training requiring fewer samples, we generated two types of training and testing datasets. The first correspond to images acquired with our active lightning camera system, while the latter comprises images captured in natural lighting. Furthermore, each of the training datasets were divided in three sub-categories according to the number images: small, medium and large, which contain 20, 40 and 94 samples, respectively. Although the number of images seem relatively small, the field of view of the camera was large enough to capture significant sections of the tree canopies, thus having a larger number of fruit counts per images. The purpose of having datasets with different sizes was not only to compare the performance between NL and AL images at different scales but also to quantify the difference in training samples required to achieve similar performance (See Section IV). For testing we also included images acquired with the sun facing directly to the camera as an extreme condition for the active and natural lightning (AL Extreme and NL Extreme). Table II summarizes the details of all the datasets and also include the total number of fruits, which was counted manually.

The results were validated across four test datasets using four commonly cited deep neural networks in agricultural applications: Faster-RCNN [38], Faster-RCNN + ResNet-50, Single Shot Detector [39], and YOLO V3 [27]. The summary of the parameters used to finetune the networks is listed in Table III and Table IV summarizes the network performance for each case.

## IV. RESULTS AND DISCUSSION

The variation in light intensity from the first experiment described in Section III-B is shown in Figure 3a. As various

TABLE II: Training and test datasets description

| | Training Dataset | | | Testing Dataset | |
| Name | No. of Images | Fruit Count | Name | No. of Images | Fruit Count |
|---|---|---|---|---|---|
| AL Small | 20 | 813 | AL Test | 51 | 1974 |
| AL Medium | 40 | 1564 | NL Test | 51 | 2272 |
| AL Large | 94 | 3218 | AL Extreme | 20 | 607 |
| NL Small | 20 | 800 | NL Extreme | 20 | 697 |
| NL Medium | 40 | 1695 | | | |
| NL Large | 94 | 3315 | | | |

TABLE III: Training parameters and values.

| Network Name | Input size | Learning rate | Epochs | Augmentation |
|---|---|---|---|---|
| YoloV3 - Darknet-53 | 512x512 | 0.001 | 250 | HF, S, R |
| Faster-RCNN - VGG16 | 1024x768 | 0.001 | 300 | HF, S, R |
| Faster-RCNN - ResNet50 | 1024x768 | 0.001 | 300 | HF, S, R |
| SSD - Inception V2 | 300x300 | 0.004 | 250 | HF, S, R |

HF = Horizontal Flip; S = Scaling; R = Random Crop

natural factors affected the intensity of the light, the amount of natural light entering the sensor also changes and the camera generates image with different exposures. This change in the image luminosity and contrast captured by the SSIM index is shown in Figure 3b. For the NL and HDR images, the SSIM quality index shows small fluctuations in time intervals closer to noon (12:01 PM – 1:00 PM, Figure 3b) and relatively higher SSIM value (approximately 70). This closeness in the SSIM index resembles closeness in quality as the brightness of the sun remained relatively constant. However, beyond this point as more significant intensity changes occurred, image quality between the first image taken at noon and the newly acquired images differed drastically. The image taken at the end of the day (6:00 PM, also seen in Figure 2b) shows a high degree of quality difference compared to the images acquired at noon.

On the other hand, SSIM index of the AL images show not only consistency but higher SSIM value throughout the entire day. Additionally, the PSNR value of AL images are also significantly larger to both NL and HDR images, as depicted in Figure 3c. As shutter speed was the only variable required to tune exposure, (see Section 3.1), the high-powered flash enables us to set shutter speed to extremely low shutter durations (see Table I). Consequently, at extreme low shutter speeds the number of photons hitting the camera sensors are also very low. Without flash, the images taken under these setting would generate pitch dark images. However, a proper synchronization of flash duration and image acquisition makes the projected light from the flash as the dominating source of illumination. This reflected light that enters back into the camera sensor essentially keeps $E_\lambda$ and $R(\lambda)$ constants in Eqn.(1). Thus, providing consistent quality in all lighting conditions. The outcomes of validating trained deep neural networks on AL and NL test datasets in present in Table IV.

TABLE IV: Deep Network performance on AL and NL test datasets (AP@0.5 IoU).

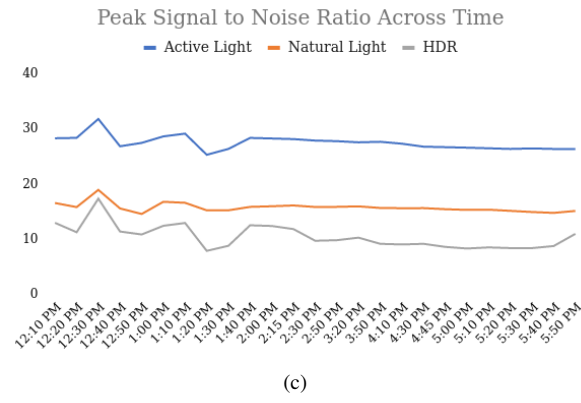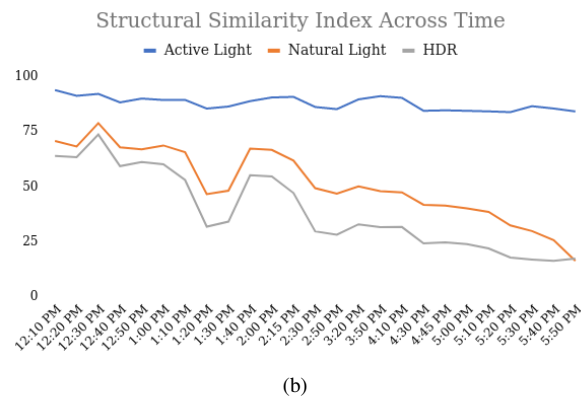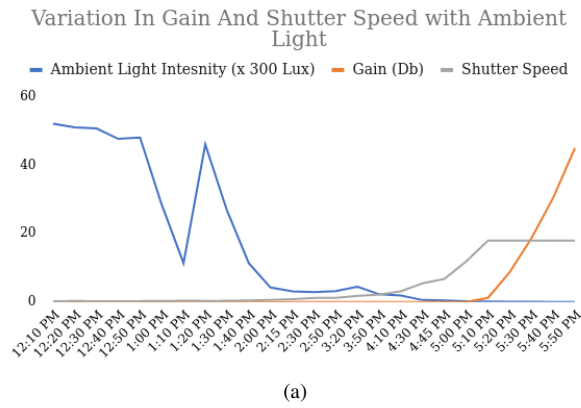| Deep Network | AL Small | NL Small | AL Medium | NL Medium | AL Large | NL Large |
|---|---|---|---|---|---|---|
| YoloV3 | **0.814** | 0.723 | **0.806** | 0.756 | **0.842** | 0.766 |
| Faster-RCNN | **0.809** | 0.630 | **0.810** | 0.714 | **0.836** | 0.714 |
| Faster-RCNN-RasNet50 | **0.801** | 0.623 | **0.813** | 0.699 | **0.829** | 0.701 |
| SSD | **0.624** | 0.515 | **0.690** | 0.524 | **0.725** | 0.542 |

Fig. 3: (a) Change in the intensity of light in a randomly selected day and its effect in the exposure variables (b) Time series plot of SSIM and (c) Time series plot of PSNR of the AL, NL, and HDR images.



Fig. 4: Precision - Recall curves for (a) YOLOv3 (b) Faster-RCNN (c) Faster-RCNN-RasNet50, and (d) SSD

For brevity, Table IV only shows the average precision (AP@0.5 IoU) as a performance measurement unit for object detection. In all instances, the object detector trained in AL datasets outperformed all their counterparts trained in the NL dataset, for each sizing category. Although the AP on both NL and AL datasets increase with more data, the trend of AP increase in AL dataset has less rate of change, showing consistency for all cases. Additionally, the AP of AL Small trained model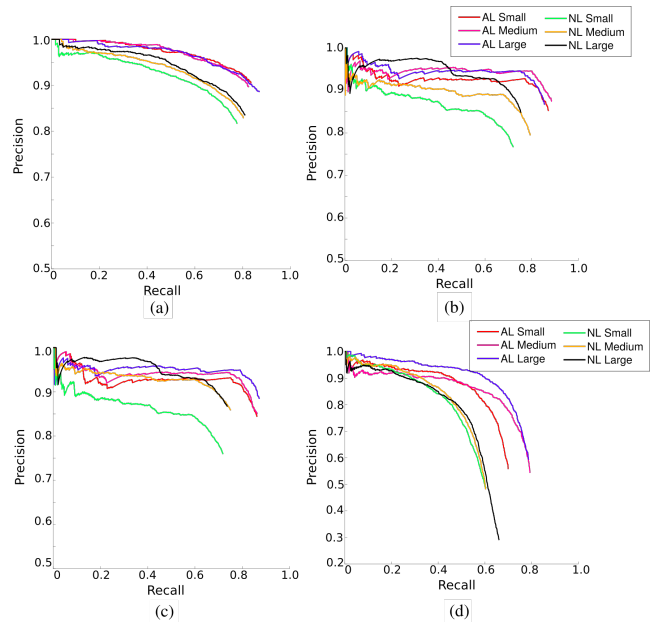 for YOLOv3, Faster-RCNN+ ResNet 50 and SSD is close to the AP of NL Large trained model. This shows that the AL Small models have learned to make better inference from relatively less data. Thus, we attribute these observations from Table IV to the quality/consistency of images acquired with active lighting.

In the AL and NL Extreme datasets (sun in the background, Figure 5), we further evaluate the robustness of the camera system in extreme situations. Although this situation is labeled as extreme test, sun flairs in images are common occurrences while imagining in the outdoor. Figure 5 shows one such situation during the data collection process. Qualitative observation of Figure 6a shows that the image with active lighting has almost no affect from sun in the background. Where as non active light image (Figure 6b) has over-saturated and most of the apples in the scene are not detectable. Faster-RCNN achieved AP of 0.71 on the AL Extreme dataset with the model in AL Large. Noticeably this detection accuracy is close to the detection accuracy of Faster-RCNN on the AL Test dataset (Table IV). On the other hand, because the images were overwhelmed with sun flairs, and perceptual information on the fruit pixels were lost, Faster-RCNN model trained on the NL Large dataset detected negligible amount of fruits, in the NL Extreme dataset counterparts.

In addition to consistent quality image and in-variance to external lighting, the proposed camera system offers additional advantages to conventional imaging systems. These advantages are summarized in the following points:

- Background subtraction: The light intensity from the flash attenuates following the principle of inverse square law [40]. This phenomenon is inherently responsible for making the background darker. This effect can also be
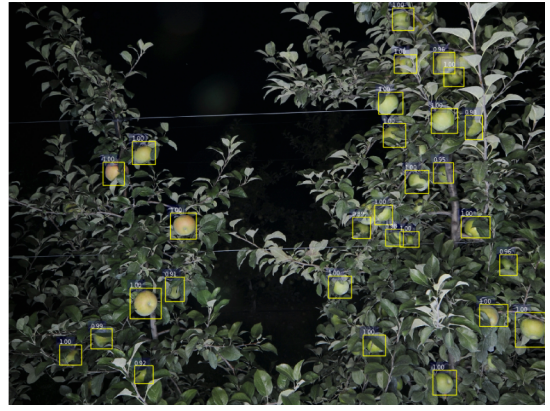
(a)



(b)

Fig. 5: Sample test images from AL Extreme (a) and NL Extreme (b) datasets with sun flairs.



(a)



(b)

Fig. 6: Sample detection results using Faster-RCNN for AL (a) and NL (b) image.

seen in Figure 6a. The tree canopy with apples and foliage in the foreground is significantly closer than the flash compared to the adjacent row in the background. The intensity of the flash that reaches the background significantly attenuates and makes it appear relatively darker. However, such affect is not possible in regular images Figure 6b. We believe that the effect of darker background could potentially facilitate and increase accuracy in image-based yield estimation of crops as fruit from the adjacent rows are inherently not visible. Further, regular computer vision algorithm for segmentation, object detection, stereo correspondence and others could also benefit.

- No motion blur: Motion blur can occur in images when the camera or the scene or both are in motion. This usually result in dis-figuration of the object and decreases the sharpness of the image. With the use of flash, time required to acquire individual image drastically decrease which consequently reduces the chances of getting motion blur. This could potentially enable robots to reliably use active lighting camera in outdoors for various tasks such as autonomous navigation and visual servoing.

- Color Consistency: With the light from the flash as the dominating source for illumination, the terms $E_\lambda$ and

$R(\lambda)$ in Eqn.(1) become constants, hence the repeatably of generating consistently exposure images get better. The combination of high color rendering index ( which is $\geqslant$ 95% for the LED flashes) and uniform exposure, render images with high color accuracy.

- High throughput: Synchronization of stereo images is a crucial requirement for stereo correspondence. Without precise image synchronization, correspondence between left and right stereo images become a difficult challenge. The stereo pairs used in our camera system are precisely triggered and synchronized with the flash. In field trials, our system could provide synchronized stereo pairs up to 20 Hz.

## V. Conclusions

Although in theory, with additional training images in the NL dataset, similar or higher detection accuracy could be achieved in fruit detection . However, in this paper we took a different approach and focused on the consistency of data rather than size. The focus of attention here was to accurately control image exposure time to generate consistent quality of images in all lighting conditions. With the use of flash as external lighting source, we demonstrated how the variables in camera exposure can be strategically reduced. Our method showed that it is possible to collect uniformly exposed images with minimum to no effects from the surrounding natural lights. This consistency in image quality allowed us to train deep neural networks with significantly less data while achieving comparable results to networks trained on larger dataset with regular images. Thus the benefit of requiring less data coupled with several advantages including background subtraction, no motion blur, color consistency, and high throughput, our proposed design could provide pragmatic solution to computer vision needs in agriculture.

Although, the use of active light generated consistently exposed images in this study, the high power of the flash and exponential attenuation of its intensity limits the work range of this camera system. In all experiments described in this paper, the camera was placed between 0.5 to 1.5 meters distance from the canopy. Beyond this range, the images could appear too dark or over exposed. By increasing shutter duration we could still get good quality images in these range, but the resilience to outdoor illuminance could get affected. However, modern vineyards and orchards grow crops in rows with uniform row separation and tree canopies trained in formal architectures to facilitate manual and mechanical operations. This trend towards almost two dimensional canopy structure provides ideal test environment for our camera system. Currently, the shutter speed during

Currently in all of our experiments, the shutter duration is intuitively selected based on qualitative observation of foreground highlight vs. background suppression. Although, this approach suffices the need for data consistency per site, temporal consistency in data could be achieved with auto shutter duration estimation. Our future work will comprise of this task of automating a calibration process to match quality between temporal datasets.

## References

[1] A. Silwal, A. Gongal, and M. Karkee, "Apple identification in field environment with over the row machine vision system," *Agricultural Engineering International: CIGR Journal*, vol. 16, no. 4, pp. 66–75, 2014.

[2] A. Silwal, J. R. Davidson, M. Karkee, C. Mo, Q. Zhang, and K. Lewis, "Design, integration, and field evaluation of a robotic apple harvester," *Journal of Field Robotics*, vol. 34, no. 6, pp. 1140–1159, 2017.

[3] K. Kapach, E. Barnea, R. Mairon, Y. Edan, and O. Ben-Shahar, "Computer vision for fruit harvesting robots–state of the art and challenges ahead," *International Journal of Computational Vision and Robotics*, vol. 3, no. 1-2, pp. 4–34, 2012.

[4] P. Chu, Z. Li, K. Lammers, R. Lu, and X. Liu, "Deepapple: Deep learning-based apple detection using a suppression mask r-cnn," *arXiv preprint arXiv:2010.09870*, 2020.

[5] I. Sa, Z. Ge, F. Dayoub, B. Upcroft, T. Perez, and C. McCool, "Deepfruits: A fruit detection system using deep neural networks," *Sensors*, vol. 16, no. 8, p. 1222, 2016.

[6] R. Kestur, A. Meduri, and O. Narasipura, "Mangonet: A deep semantic segmentation architecture for a method to detect and count mangoes in an open orchard," *Engineering Applications of Artificial Intelligence*, vol. 77, pp. 59–69, 2019.

[7] "Apple facts," https://web.extension.illinois.edu/apples/facts.cfm, accessed: 2020-11-16.

[8] USDA, "Weeds of united states and canada: Usda plants," https://plants.usda.gov/java/invasiveOne?pubID=SWSS.

[9] A. Gongal, S. Amatya, M. Karkee, Q. Zhang, and K. Lewis, "Sensors and systems for fruit detection and localization: A review," *Computers and Electronics in Agriculture*, vol. 116, pp. 8–19, 2015.

[10] D. Ilstrup and R. Manduchi, "One-shot optimal exposure control," in *European Conference on Computer Vision*. Springer, 2010, pp. 200–213.

[11] H. Tian, T. Wang, Y. Liu, X. Qiao, and Y. Li, "Computer vision technology in agricultural automation—a review," *Information Processing in Agriculture*, vol. 7, no. 1, pp. 1–19, 2020.

[12] A. Syal, D. Garg, and S. Sharma, "Apple fruit detection and counting using computer vision techniques," in *2014 IEEE International Conference on Computational Intelligence and Computing Research*. IEEE, 2014, pp. 1–6.

[13] A. B. Payne, K. B. Walsh, P. Subedi, and D. Jarvis, "Estimation of mango crop yield using image analysis–segmentation method," *Computers and electronics in agriculture*, vol. 91, pp. 57–64, 2013.

[14] S. Nuske, K. Wilshusen, S. Achar, L. Yoder, S. Narasimhan, and S. Singh, "Automated visual yield estimation in vineyards," *Journal of Field Robotics*, vol. 31, no. 5, pp. 837–860, 2014.

[15] A. Kamilaris and F. X. Prenafeta-Boldú, "Deep learning in agriculture: A survey," *Computers and electronics in agriculture*, vol. 147, pp. 70–90, 2018.

[16] J. Yu, S. M. Sharpe, A. W. Schumann, and N. S. Boyd, "Deep learning for image-based weed detection in turfgrass," *European journal of agronomy*, vol. 104, pp. 78–84, 2019.

[17] K. P. Ferentinos, "Deep learning models for plant disease detection and diagnosis," *Computers and Electronics in Agriculture*, vol. 145, pp. 311–318, 2018.

[18] J. P. Vasconez, J. Salvo, and F. Auat, "Toward semantic action recognition for avocado harvesting process based on single shot multibox detector," in *2018 IEEE International Conference on Automation/XXIII Congress of the Chilean Association of Automatic Control (ICA-ACCA)*. IEEE, 2018, pp. 1–6.

[19] 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.

[20] K. Osorio, A. Puerto, C. Pedraza, D. Jamaica, and L. Rodríguez, "A deep learning approach for weed detection in lettuce crops using multispectral images," *AgriEngineering*, vol. 2, no. 3, pp. 471–488, 2020.

[21] P. Jiang, Y. Chen, B. Liu, D. He, and C. Liang, "Real-time detection of apple leaf diseases using deep learning approach based on improved convolutional neural networks," *IEEE Access*, vol. 7, pp. 59 069–59 080, 2019.

[22] M. Rahnemoonfar and C. Sheppard, "Deep count: fruit counting based on deep simulated learning," *Sensors*, vol. 17, no. 4, p. 905, 2017.

[23] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.

[24] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[25] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.

[26] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Malloci, T. Duerig *et al.*, "The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale," *arXiv preprint arXiv:1811.00982*, 2018.

[27] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[28] S. P. Mohanty, D. P. Hughes, and M. Salathé, "Using deep learning for image-based plant disease detection," *Frontiers in plant science*, vol. 7, p. 1419, 2016.

[29] G. Lobet, "Image analysis in plant sciences: publish then perish," *Trends in plant science*, vol. 22, no. 7, pp. 559–566, 2017.

[30] G. Lobet, X. Draye, and C. Périlleux, "An online database for plant image analysis software tools," *Plant methods*, vol. 9, no. 1, p. 38, 2013.

[31] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 843–852.

[32] H. Altaheri, M. Alsulaiman, G. Muhammad, S. U. Amin, M. Bencherif, and M. Mekhtiche, "Date fruit dataset for intelligent harvesting," *Data in brief*, vol. 26, p. 104514, 2019.

[33] S. Bargoti and J. Underwood, "Deep fruit detection in orchards," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3626–3633.

[34] T. Botterill, S. Paulin, R. Green, S. Williams, J. Lin, V. Saxton, S. Mills, X. Chen, and S. Corbett-Davies, "A robot system for pruning grape vines," *Journal of Field Robotics*, vol. 34, no. 6, pp. 1100–1122, 2017.

[35] Z. Pothen and S. Nuske, "Automated assessment and mapping of grape quality through image-based color analysis," *IFAC-PapersOnLine*, vol. 49, no. 16, pp. 72–78, 2016.

[36] B. Pillman and D. Jasinski, "Camera exposure determination based on a psychometric quality model," *Journal of Signal Processing Systems*, vol. 65, no. 2, p. 147, 2011.

[37] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[38] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.

[39] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.

[40] N. Voudoukis and S. Oikonomidis, "Inverse square law for light and radiation: A unifying educational approach," *European Journal of Engineering Research and Science*, vol. 2, no. 11, pp. 23–27, 2017.