# TREE DETECTION AND DIAMETER ESTIMATION BASED ON DEEP LEARNING

**Vincent Grondin**\*
Department of Computer Science
Laval University

**Jean-Michel Fortin**
Department of Computer Science
Laval University

**François Pomerleau**
Department of Computer Science
Laval University

**Philippe Giguère**
Department of Computer Science
Laval University

November 1, 2022

## ABSTRACT

Tree perception is an essential building block toward autonomous forestry operations. Current developments generally consider input data from lidar sensors to solve forest navigation, tree detection and diameter estimation problems. Whereas cameras paired with deep learning algorithms usually address species classification or forest anomaly detection. In either of these cases, data unavailability and forest diversity restrain deep learning developments for autonomous systems. So, we propose two densely annotated image datasets — 43 k synthetic, 100 real — for bounding box, segmentation mask and keypoint detections to assess the potential of vision-based methods. Deep neural network models trained on our datasets achieve a precision of 90.4 % for tree detection, 87.2 % for tree segmentation, and centimeter accurate keypoint estimations. We measure our models' generalizability when testing it on other forest datasets, and their scalability with different dataset sizes and architectural improvements. Overall, the experimental results offer promising avenues toward autonomous tree felling operations and other applied forestry problems. The datasets and pre-trained models in this article are publicly available on GitHub (https://github.com/norlab-ulaval/PercepTreeV1).

***Keywords*** Computer Vision · Forestry · Automation

## 1   Introduction

Heavy forest machinery can function at a capacity that is well above its operator's sustainable workload [Ringdahl, 2011, Parker et al., 2016]. Hence, automation could increase operator efficiency while reducing mental fatigue [Ortiz et al., 2014, Visser and Obi, 2021]. A technology driver behind advancements in automation is machine learning, which is an established approach to advance particular fields [Roy et al., 2021]. Therefore, applying learning methods to a problem set specifically to the field of forestry, such as harvesting (Figure 1), forwarding, or navigation, will play an important role in the future of forestry automation. However, forests are harsh, unstructured outdoor environments [Thorpe and Durrant-Whyte, 2001]. As an environment becomes less structured, its complexity increases beyond the typical well-defined boundaries and narrow operating range established in structured environments, possibly leading to operational shortcomings of autonomous systems [Roy et al., 2021].

Autonomous systems, like humans, heavily rely on the perception of their surrounding environment to make decisions. While perception systems are adequate for clearly stated tasks and well-defined objects, general object recognition is still difficult. From a forestry standpoint, the capacity to quickly and reliably detect trees in various situations remains an important challenge for the development of multiple applications. For instance, tree selection in thinning

---

\*vincent.grondin.2@ulaval.ca

Figure 1: Prediction from our vision-based model. The image is from an operator's point of view during felling operations in a Canadian forest. Each detected tree is shown in a different color and provides a bounding box, a segmentation mask and keypoint information.

operations is notably influenced by variable tree parameters such as stems and crowns density, height, diameter and species [Nurminen et al., 2006], making the tree appearance exceedingly diverse. As of now, the human operator's visual capacity is the only one capable of estimating these parameters and making informed decision during forestry operations.

However, replicating in machines the human perception system's effectiveness at distinguishing objects in different environments, and under a variety of conditions such as low illumination, color differences, and occlusions [Padilla et al., 2021], is an active area of research. In computer vision, the most successful perception methods use images combined with deep neural networks and supervised learning [Diez et al., 2021]. As a data-centric method, supervised deep learning leverage large quantities of annotated images to discover intricate structures in high-dimensional pixel arrays [LeCun et al., 2015]. An important limitation, although not unique to deep learning, is that any variation not learned within a training set can hinder performances and cause out-of-distribution generalization errors. This is particularly true for real-world applications, where the physical world is a large and rich source of information for training that it can't be seized by any dataset or even simulated worlds [Roy et al., 2021]. In this context, a central challenge for perception systems is to generalize across a range of unforeseen and changing forest environments. Parallel to generalization concerns, the lack of resources to create a large dataset limits the adoption of deep learning for many problems [Oliver et al., 2018], including tree detection.

Creating large annotated datasets typically require great human effort and financial expenditure [Oliver et al., 2018]. For instance, Bowen et al. [2021]'s analysis of human annotation time on COCO [Lin et al., 2014] measures an average annotation time of 122.4 s per instance. For a dataset with a size similar to COCO, this means thousands of annotation hours. One way to reduce these costs is to create a synthetic dataset using a simulator that can simultaneously generate and annotate images [James et al., 2019]. Provided that the perception system is pre-trained on a large synthetic dataset, one can achieve reasonably successful domain transfer (i.e. transfer learning) by fine-tuning the model on a small set of real data [Oliver et al., 2018]. This would alleviate the under generalization issue deep learning models have when facing the domain shift between synthetic and real images.

Based on the previous considerations, we propose two novel datasets to train common object detectors for tree perception, as well as felling cut, diameter and inclination estimation (see Figure 2). To do so, we create a large synthetic dataset with software-generated annotations. We also generate a smaller dataset of real images that have been hand-annotated. These datasets contain a vast variety of images and include class category, bounding box, segmentation mask and keypoint annotations. Thereby, the perception problem is framed as one of instance segmentation to leverage the Mask R-CNN [He et al., 2017] framework and its derivative architecture, Cascade Mask R-CNN [Cai and Vasconcelos, 2019]. Since these models have modular components that can be easily changed, this facilitates the integration of
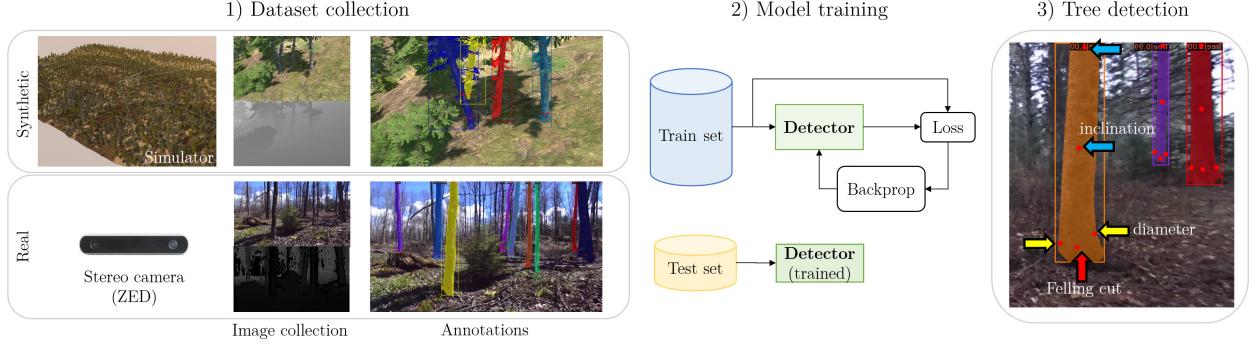
Figure 2: Workflow of the proposed tree detector. 1) We create and annotate two forest datasets – one synthetic and one real –. 2) Using these datasets, object detector models are trained and tested. 3) The model can detect, segment and estimate the felling cut, diameter and inclination.

state-of-the-art architectures as they come out. In our experiments, we compare architectures and backbones to identify the most promising model. Furthermore, limitations of our models are quantitatively evaluated by their capability to generalize on another dataset. We further analyze the benefits of synthetic images, detection reliability in presence of occlusion, the impact of train set size, and explore possible improvements for keypoint accuracy. Together, these experiments focus on discovering practical insights relevant to deep learning for autonomous tree detection.

In short, our contributions are:

- Two annotated datasets of synthetic and real RGB-D images in natural forest environments are made publicly available;
- A performance comparison between different backbones and model architectures for tree detection; and
- Modified object detectors to estimate the felling cut, diameter and inclination for each tree.

## 2  Related Work

This literature survey is an overview of important advances in forestry automation. More specifically, we present works related to robotic applications in forests at a ground level, mainly towards autonomous forwarding and tree felling. Then, we cover works related to vision-based methods and expose how individual tree detection remains a crucial open problem.

### 2.1  Robotics in Forested Environments

According to Ringdahl [2011], log transportation (i.e., forwarding) is the most likely operation to be automated, followed by tree extraction (i.e., tree felling). An early proof-of-concept for forwarding automation is proposed by Hellström et al. [2009] in 2009. They use a simulated environment where the vehicle moves along a path by using a path-search mechanism that generates feasible and obstacle-free 10 m long path segments. In real-time, obstacle detection (i.e. trees) is attempted by the vehicle's laser scanner sensors and an occupancy grid, but it is not clear if its been satisfactorily resolved. More recently, Jelavic et al. [2021] demonstrate a fully autonomous machine for tree felling operations. Using a lidar sensor as its primary sensing modality, the automated system performs mapping, localization, planning, control and tool positioning. Tree detection is based on a geometry algorithm applied on 6 m × 6 m patches from a point cloud map. During tests performed in a forest environment, the authors report successful detections in presence of vegetation, but the algorithm produces false negatives under heavy clutters. Also, the authors underscore that geometry-based detection harbor safety concerns, as the algorithm cannot distinguish between a standing person and a tree trunk.

Apart from autonomous operations, tree detection from point clouds is used in other research experiments. Zhang et al. [2019] extract the trunk's central point from clustered point clouds using a Gauss–Newton method. For their experiment they use a robot navigating in a structured, evenly spaced forest plantation. In these conditions, they obtain a tree localization error slightly above 10 cm (Root mean square (RMS), as well as a sub-cm precision for the radius (RMS = 0.66 cm), rendering the method eligible to carry out autonomous tree harvesting. In a similar way, Pierzchała et al. [2018] perform tree localization and Diameter at Breast Height (DBH) estimation in a semi-structured, little undergrowth and flat terrain forest using graph-SLAM. They achieve a mean DBH estimation error of 2 cm and a tree

position error of 4.76 cm. Tremblay et al. [2020] obtain comparable results in more challenging forests. Based on a cylinder fitting method, they report a DBH error of 2.04 cm in mature forests with well-spaced trees and visible trunks, compared to 3.45 cm error in unruly forest conditions. Liang et al. [2014] proposed a photogrammetric approach for 3D imaging, known as Structure from Motion, and achieve a 2.39 cm RMS error on DBH estimation of individual trees, which is acceptable for practical applications. The authors achieve similar results when using a terrestrial lidar [Liang et al., 2014], which was used as a benchmark.

In our work, we aim for reliable tree detection and diameter estimation, but also go one step further by extracting relevant harvesting information such as tree felling cut location and inclination. This is an important contribution, since it would provide an end-to-end perception method ready for tree felling automation.

## 2.2 Vision-based deep learning

Limited work is conducted with deep learning to perform vision-based tree detection. An image-based trunk detection system for forestry mobile robotics is presented in da Silva et al. [2021]. Using a mix of visible and thermal images, they create a dataset composed of 2895 images extracted from video sequences that contain bounding box annotations to detect the presence of tree trunks in images. They train five different one-shot detectors on their dataset and achieve 89.84 % Average Precision (AP), and 89.37 % F1 score. Their dataset only include bounding box annotation, and exclude segmentation masks to precisely delimit pixels corresponding to trunks and keypoint annotations for each tree. Meanwhile, Wang et al. [2019] use Faster R-CNN to detect tree trunks in depth images. These depth images are based on above-ground point cloud data, where the point clouds are extracted, voxelized and projected on a 2D plane to form an image. When trained on 802 and tested on 359 of these images, they achieve more than 90 % of detection accuracy. Itakura and Hosoi [2020] propose an occlusion-aware R-CNN for detecting trees in urban environments from reconstructed 3D images captured by a 360 spherical camera. Using YOLOv2 [Redmon and Farhadi, 2016], they achieved 88 % precision, 100 % recall, and can estimate diameter with less than 3 cm of error. Again, the stark difference between urban and natural forest environments render this dataset inoperative for forestry operations.

Deep learning is also used to classify trees and estimate stock volume. Liu et al. [2019] train a U-Net [Ronneberger et al., 2015] for instance segmentation on more than 3 k images, and then process the segmentation masks with a non-linear mixed effect model to estimate stock volumes, achieving 97.25 % precision, 95.68 % recall, 96.03 % classification accuracy and $30.54\,\mathrm{m^3/ha}$ stock volume RMS error. Further, Ostovar et al. [2019] tackle stump detection and classification in forests by training Faster R-CNN [Ren et al., 2017] on 800 images and then testing it on 200 images. They report precision and recall rates of 95 % and 80 %, respectively.

A common theme among the aforementioned tasks is the obligation to detect individual trees first, which allow subsequent tasks, like stock volume estimation or species classification. Insofar, precision and recall rates reported for tree detection among the precedent works are relatively high, which does not seem to truly reflect its difficulty. In this regard, a main contribution of our work is to measure the performance of deep learning on this specific task, and quantitatively demonstrate its limitations when deployed on out-of-distribution forest environments.

# 3 Method

The methodology is divided into four sections. First, we detail our two datasets, SYNTHTREE43K and CANATREE100. Then, object detection models are presented, followed by training parameters. Finally, we describe the evaluation metrics used in the experiments.

## 3.1 Datasets

One of the key issues restraining the adoption of deep learning approaches in the forestry domain is the limited availability of annotated datasets dedicated to forest applications. Although this research area has disruptive potential toward automation, no existing dataset contains enough annotation to train a tree detector with exact pixel location of the trunk, felling cut, diameter and inclination. A major advantage of simulation is the ability to fully control the data generation pipeline, which ensures reduced costs, and infinite variety and quantity [Gaidon et al., 2016]. Therefore, we propose to employ simulated virtual forests as a proxy for our tree detection tasks. As it has already been shown [Gaidon et al., 2016, Cabon et al., 2020], pre-training deep neural networks on virtual data improve object detection performances. So, we created a large novel synthetic dataset to partially fill the data gap in forestry.

Named SYNTHTREE43K, the dataset is a collection of 43 k synthetic images rendered in a simulated environment using the Unity[2] game engine. The image resolution is set to $1280 \times 720$ pixels. To build realistic forest scenes (see Figure 3),

---

[2] https://unity.com

we employ Gaia[3], a Unity asset, to efficiently terraform landscapes, texture terrain, and place objects models. Forest types described by Diez et al. [2021] — plantation, well-managed and natural forests — are taken into account by controlling trees spawn density. Then, six basic 3D tree models from Nature Manufacture[4] (from fir and beech species) are modified by changing bark and foliage textures to create 17 different trees. During scene generations, each tree is placed according to spawn rules about the terrain altitude, slope and object density in a given area. When spawned, a tree is assigned a random width and height for increased appearance diversity. Other objects and conditions commonly found in forests are also included in this simulated environment :

- grass, stumps, scrubs and branches;

- lighting conditions (i.e., morning, daytime, evening, and night light); and
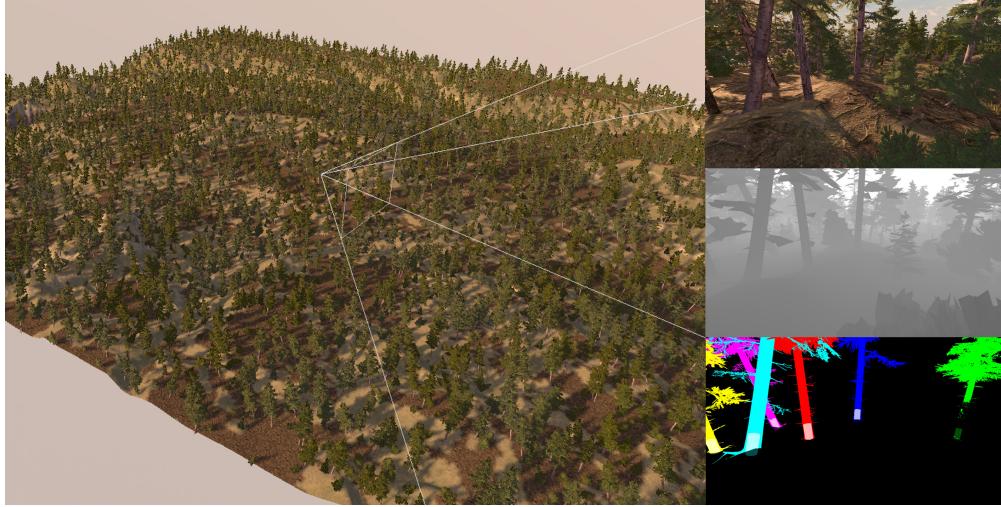
- weather effects related to fog, snow, and rain.



Figure 3: The SYNTHTREE43K dataset. **Left** a general view of a simulated scene where the camera point of view is represented by white lines. This scene has around 1500 tree instances, of which we sample 750 images. Along with a morning light effect, and the terrain texture follows an altitude rule where the soil low and the grass is high. The scene also contains scrubs, grass and branches. **Right** a sample RGB frame, its depth frame, and its segmentation mask annotations. Individual trees are in different colors, and the tree base is slightly paler.

We sample images from the scene by arbitrarily moving the camera from tree to tree with random roll, pitch and yaw angle within a certain range to mimic small changes in camera pose caused by uneven terrain navigation. For each generated scene, 200 to 1000 images are sampled depending on the number of trees in the scene. This is done in an effort to reduce the odds of the same tree appearing in multiple images. For the annotation method, a tree is only annotated if within the typical feller crane's reach, 10 m, which in our case is its distance from the camera.

Bounding box annotation is confined to the trunk to omit the tree crown foliage, and this differs from the urban tree dataset in Itakura and Hosoi [2020], where bounding box annotations contain entire trees. This decision is motivated by the fact that 1) tree density in natural forests makes them more prone to overlapping error than in urban environments, and 2) tree felling primarily requires information about the trunk. Contrary to the works of Itakura and Hosoi [2020] and da Silva et al. [2021], we add segmentation mask annotations, and provide pixel information. We equally leverage the simulator to compute the percentage of tree occlusion from the camera's point of view (see Figure 4), which is used later on to quantify the model's robustness with regard to occlusion. The percentage of occlusion takes into account occlusion from grass, other trees and the landscape. Moreover, we compute the entire tree visibility and ignore tree instances if less than 30 % of their pixels are visible. This threshold was set to restrain training on instances with too much occlusion. Such instances can be misleading, causing the detection thresholds to be lowered, which in turn increases the objects' context influence, and leads to false positive (FP) in regions with no object [Wang et al., 2020]. Lastly, we locate five keypoints on each tree to capture cues relevant to harvesting. These keypoints correspond to the tree's felling cut location, diameter and inclination. Our simulator setup can generate these annotated synthetic images

---

[3]https://assetstore.unity.com/packages/tools/terrain/gaia-2-terrain-scene-generator-42618

[4]https://naturemanufacture.com/

at a rate of approximately 20 frames/minute, for we consider RGB and depth images as one frame. Our final synthetic dataset contains over 43 k RGB-D images, and more than 162 k annotated trees.
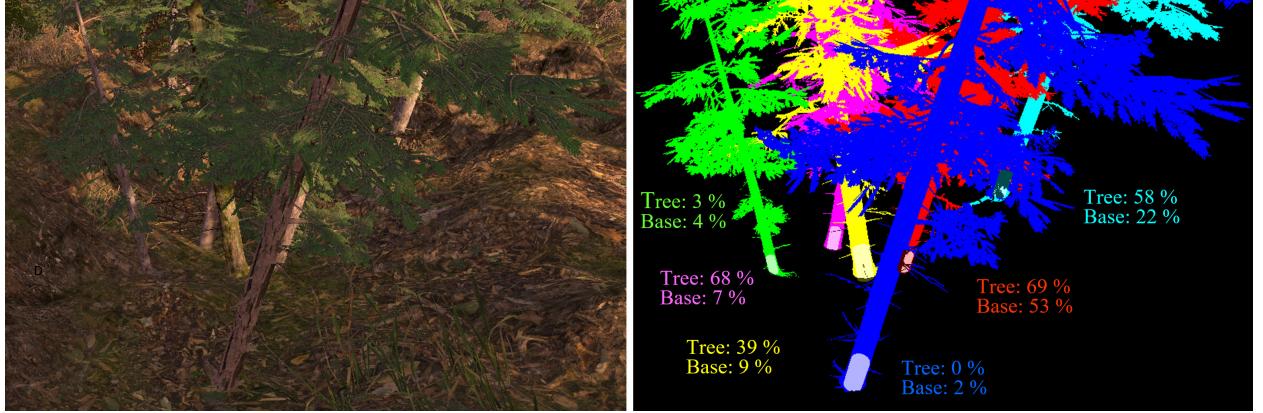


Figure 4: **Left** RGB synthetic image and **right** its segmentation mask (later on fitted to the trunk). Each tree has a segmentation mask denoted by a different color. Lighter colored areas correspond to the tree base. Occlusion is measured based on the percentage of visible tree pixels visible the image frame.

Our real dataset, referred to as the CANATREE100 dataset, was collected from public, private, and commercial forests in Quebec, Canada. The images were taken from June 2020 to April 2021, between the times of 9 am and 6 pm. Tree species commonly found in these forests include fir, spruce, pine, birch and maple. Our images are hand-annotated with the COCO-Annotator [Brooks, 2019]. During annotation, we manually draw the segmentation mask on the trunk only, fit bounding box to it and place keypoints. As opposed to the annotations in the synthetic dataset, we omit tree branches during segmentation annotation (that fit inside the bounding box), since they are complex to annotate and we have no use for it in forestry operations. The felling cut and diameter keypoints are approximately placed 10 cm above ground. For the two inclination keypoints along the tree trunk, contrary to the simulator that placed them at 90 cm intervals starting from the felling cut, we decided to place one in the middle and one on top.

Overall, CANATREE100 provides 100 RGB and 100 depth images, as well as more than 920 annotated trees. The images are garnered from 33 videos captured by a hand-held ZED stereo camera, and have HD of $1280 \times 720@60$ fps resolution. Like the synthetic dataset, the RGB and depth images are considered as one frame. We manually select approximately three images per video ($< 1\%$ of the video). We subjectively select images that are challenging and rich in information. Because of this very low sampling rate, the correlation between samples is minimized, but CANATREE100 is significantly smaller than other analogous datasets [Liu et al., 2019, da Silva et al., 2021] that employ much higher sampling rates from videos. Although using a high frame rate allows for efficient image annotation due to temporal correlation, it might result in impoverished visual diversity, which is inefficient when training deep networks.

So far, our datasets are limited to the situation where train and test sets come from forests of similar ecosystems (deciduous, mixed or boreal forests from Quebec). Since forest machinery operates in a variety of areas, there is a fundamental concern about our tree detector's ability to generalize to other regions, without fine-tuning. In an effort to address this, we use the trunk detection dataset from da Silva et al. [2021], which we refer to as the PORTUGAL dataset. This should help us measure the extent at which deep learning models trained on our datasets generalize to a different domain The domain transfer experiments in Section *Domain transfer* only measure performances for bounding box detections, as the dataset only contains bounding box annotations. This dataset is composed of 2690 RGB images collected from three different forest sites in Portugal, to which the authors applied eight unique data augmentations to every image, increasing their dataset size to a total of 24 k images. For a fair comparison, we use the same train, validation and test split as da Silva et al. [2021].

## 3.2 Model architectures

Two different model architectures are used in the experiments — Mask R-CNN and Cascade Mask R-CNN. These models have two-stages, in which candidate bounding boxes are proposed by a Region Proposal Network (RPN), and then three network heads individually predict a bounding box, segmentation mask and keypoints (see Figure 5).

In the first stage, images are first processed by the backbone to extract distinctive object features. To do so, we experiment with three backbone architectures. For instance, we use ResNeXt [Xie et al., 2017], a Convolutional Neural Network (CNN)-based backbone well suited for the grid-structure of images. For a long time, CNNs have dominated

computer vision with their intrinsic capability to impose regularizing priors for object detection like stationarity, locality and compositionality [Bronstein et al., 2017]. Image features can also be extracted using transformer-based backbones such as Swin [Liu et al., 2021]. Beyond being state-of-the-art in many domains, transformers are able to model long-range dependencies in data through a self-attention mechanism. We conjecture that this mechanism might be beneficial for forest targeted application, since short-range information is often obfuscated by overlapping trees or occlusion. The Swin transformer also employs a shifted-window strategy for self-attention, which improves efficiency at distinct object feature extraction, important for real-time operations. The third and last feature extraction method we experiment with is CBNetV2 [Liang et al., 2021]. Without introducing a new feature extraction paradigm, this method uses dual backbones (DB). Contrary to simple network deepening or widening, the CBNetV2 backbone architecture integrates high- and low-level features of multiple identical backbones that are connected through composite connections. It can be adapted to various backbones (i.e. CNN-based, Transformer-based) as well as head designs of most mainstream detectors. Table 1 presents the complete list of backbones and number of trainable parameters used in our experiments.
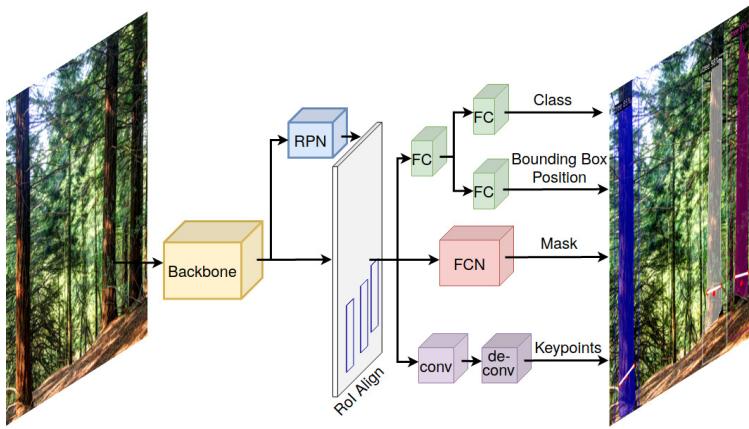


Figure 5: Model architecture of Mask R-CNN with a keypoint branch. For each input image, the backbone extracts distinctive features to predict RoIs and the three network heads. The first head (green) predicts the class and bounding box, the second head (red) predicts the segmentation mask, and the third head (purple) predicts the keypoints. This unified architecture enables end-to-end training without any post-processing.

The last step of the first stage is to pass the extracted feature to the RPN (the blue module in Figure 5), which identifies a set of regions that could contain objects. In doing so, the RPN evaluates the objectness of each pre-defined axis-aligned anchor box corresponding to three area-scales (8, 16 and 32) and three aspect ratios (0.5, 1.0, and 2.0) to generate region proposals.

Then, during the second stage, RoiAlign [He et al., 2017] extracts small feature maps for each region proposal using a bilinear interpolation to map the feature vectors extracted by the backbone, into a $7 \times 7$ input feature vector. Subsequently, feature vectors from each RoI go through the three network heads. The first head predicts the class and bounding box, the second head predicts the segmentation mask, and the third head predicts the keypoints. The segmentation head performs semantic segmentation with an output resolution of $28 \times 28$. The keypoint head accuracy is set at $56 \times 56$, because keypoint localization requires a higher resolution [He et al., 2017].

### 3.3 Training

We use the PyTorch based library MMdetection [Chen et al., 2019] for model implementations. All backbone weights are initialized from their pre-trained version on COCO [Lin et al., 2014] 2017, a large-scale dataset containing 80 object classes and 164 k images. This pre-training phase is important to help regularize the networks [Erhan et al., 2010] and facilitate transfer learning from a source to a target domain [Mahajan et al., 2018]. Hinterstoisser et al. [2018] found no significant gain to re-train the entire backbone, hence some backbone stages are frozen before training or fine-tuning the models. So, in our experiments, the first three (out of 4) stages are frozen in the Swin backbone, whereas in ResNeXt-101 is only frozen in the first (out of 5) stage. We use an initial learning rate of 1e-4 for pre-training and 5e-5 for fine-tuning, decreasing the rate by a factor of 10 at predetermined steps. The steps were obtained through a grid-search process by increments of 10. Pre-training and fine-tuning is done on two *NVIDIA A100* GPU (40 GB memory).

Pre-training on SYNTHTREE43K is conducted on three subsets — a train set, a validation set, and a test set — of 40 k, 1 k and 2 k images, respectively. At train time, we employ common data augmentation techniques (i.e. image

resize, horizontal flip, sheer, and camera distortions), and these significantly improve model generalization. The model processes images from the train set in batches of four, and it is optimized with AdamW and a weight decay of 0.05. This optimizer regularizes learning by using the weight decay. As training progresses, the validation set monitors overfitting and performs early stopping.

Table 1: Model parameters, Floating Point Operation (FLOP)s and frame per second (fps) on an NVIDIA RTX 3090.

| Model | Backbone | #Params | FLOPs | fps |
|---|---|---|---|---|
| Mask | X-101 | 63M | 238G | 12.5 |
| | Swin-T | 48M | 267G | 12.8 |
| | DB-Swin-T | 76M | 357G | 11.7 |
| Cascade | X-101 | 101M | 819G | 7.2 |
| | Swin-S | 107M | 832G | 7.0 |
| | DB-Swin-S | 156M | 1016G | 6.4 |

Although models pre-trained on synthetic images can learn the tree appearance in simulation and transfer this knowledge for detection in the real world, their performance can be improved, sometimes significantly, when fine-tuned on real images[Gaidon et al., 2016]. This fine-tuning step accounts for the reality gap between the two domains (synthetic vs. real), and the minor differences that might exist between synthetic annotations and manual annotations, for instance the bounding box height and keypoint positions. Therefore, transfer learning to real images is done by fine-tuning the model on CANATREE100. Due to its smaller size, we use a five-fold cross-validation scheme: 60 in the train set, 20 in the validation set, and 20 in the test set. Importantly, we did not segregate physical sites when doing the train, validation and test split. This means that images share similar characteristics and lighting conditions in different splits. However, the same trees are not present in more than one set. Data augmentation is used for fine-tuning, but considering the small train set size, an image batch size of two is preferred. The final results are reported on the test set without any data augmentation, as is commonly done.

## 3.4 Detection metrics

The quality of our tree detector can be measured using two criteria: 1) its ability to find all ground-truth trees, while 2) identifying only relevant trees [Padilla et al., 2021]. For the first criterion, the detector is considered to have a high recall rate when there is very few false negatives (FNs). Then, when evaluated by the second criterion, the detector is considered to have a high precision rate when there are very few false positives (FPs). In case an FP or FN happens during forestry operations, there is a risk of decreasing productivity or damaging equipment. Detection results are reported using the standard COCO metrics: Average Precision (AP) and Average Recall (AR). In our experiments, we measure the bounding box predictions ($AP^{bb}$, $AR^{bb}$) and segmentation mask predictions ($AP^{seg}$, $AR^{seg}$)

The AP takes the weighted sum of precision at each intersection over union (IoU) threshold:

$$AP = \frac{1}{N} \sum_r p_{\text{interp}}(r),$$ (1)

where IoU is defined as the ratio between the intersection and the union of the predictions and the ground truths; $p_{\text{interp}}(r)$ is the interpolated precision at each equally spaced recall level $r = \{0, 0.01, \dots, 1\}$, with $N = 101$. The AP is computed following the standard intervals of IoU = [0.5:0.05:0.95]. For clarity, we denote it $AP_{50:95}$. Similarly, we also report $AP_{50}$ results computed for IoU = 0.5, since many previous works [da Silva et al., 2021, Ostovar et al., 2019, Itakura and Hosoi, 2020, Zhang et al., 2019] only report their results using 50 % minimum IoU.

The AR measures the assertiveness of object detectors. It is computed using the average of maximum recall across several IoU thresholds $O$:

$$AR = \frac{1}{O} \sum_{o=1}^{O} \max_{k|Pr_{t(o)}(\tau(k))>0} \{Rc_{t(o)}(\tau(k)\},$$ (2)

where $Pr_{t(o)}(\tau(k))$, $Rc_{t(o)}(\tau(k))$ are the precision × recall points for $\tau(k)$ confidence. In our experiments, we use the $AR_{50:95}$, $AR_{50}$.

# 4 Results and Discussions

In order to evaluate our approach, we first test it on CANATREE100. After measuring the capability of our tree perception method on its source domain, we focus our attention on its generalization performances on another forest dataset. Finally, we conduct further studies quantifying *i)* the impact of synthetic images during pre-training, *ii)* model robustness to occlusion *iii)* performance scalability as a function of dataset size, and *iv)* keypoint architectural improvements.

## 4.1 Main results

**Tree detection:** Tree detector performances based on different backbones and model architectures are presented in Table 2. There is a consistent increase in detection results for Cascade Mask R-CNN models relative to Mask R-CNN models. Having several cascaded heads that progressively refine predictions improve results, and this is coherent with the results obtained by Cai and Vasconcelos [2019]. For backbones, transformer-based backbone (Swin) systematically outscore their CNN counterpart. Using dual backbones (DB) further improves detection and segmentation results. Overall, the model with the most (156 M) learnable parameters, Cascade Mask-RCNN DB-Swin-S, obtained a detection performance of 90.4 % $AP_{50}^{bb}$ and 94.5 % $AR_{50}^{bb}$. While these results are almost identical to the ones of Mask-RCNN DB-Swin-S — 90.3 % $AP_{50}^{bb}$ and 94.9 % $AR_{50}^{bb}$ — the Cascade model has a significant edge for all other detection metrics, notably gaining 3.8 % $AP_{50:95}^{bb}$ and 1.5 % $AR_{50:90}^{bb}$. Qualitative results can be observed in Figure 6.



Figure 6: **(Top row)** Detections are in comparison to their **(bottom row)** ground truths. The straight arrows point out the missed and false detections. Missed detections are mostly caused by tree sizes, while false detections are often due to tree distance. The curved arrows point out failed keypoint localizations commonly caused by overlapping trees; as a result, the keypoint is predicted on the wrong tree. The prediction model (Cascade Mask R-CNN DB-Swin-S) is pre-trained on SYNTHTREE43K, and fine-tuned on CANATREE100.

Table 2: Detection performance on CANATREE100. Keypoint errors relating to diameter (dia), felling cut (fc) and inclination (inc) are presented in metric units. Transformer backbones improve results, and dual-backbones consistently increase results for all tasks. Models are pre-trained on COCO and SYNTHTREE43K.

| Model | Backbone | $AP_{50}^{bb}$ | $AP_{50:95}^{bb}$ | $AP_{50}^{seg}$ | $AP_{50:95}^{seg}$ | $AR_{50}^{bb}$ | $AR_{50:95}^{bb}$ | $AR_{50}^{seg}$ | $AR_{50:95}^{seg}$ | Keypoint error | | |
| | | | | | | | | | | $dia_{(cm)}$ | $fc_{(cm)}$ | $inc_{(\circ)}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mask | X-101 | $86.3_{\pm2.2}$ | $60.4_{\pm0.9}$ | $82.8_{\pm1.3}$ | $54.2_{\pm1.5}$ | $90.7_{\pm1.3}$ | $66.9_{\pm0.3}$ | $87.5_{\pm1.3}$ | $60.2_{\pm1.0}$ | 3.1 | 6.4 | 2.5 |
| | Swin-T | $88.2_{\pm1.2}^{\uparrow1.9}$ | $59.6_{\pm2.6}^{\downarrow0.8}$ | $83.5_{\pm1.6}^{\uparrow0.7}$ | $55.8_{\pm2.2}^{\uparrow1.6}$ | $93.2_{\pm1.3}^{\uparrow2.5}$ | $67.1_{\pm1.7}^{\uparrow0.2}$ | $88.4_{\pm2.3}^{\uparrow0.9}$ | $61.8_{\pm1.9}^{\uparrow1.6}$ | $3.1^{0.0}$ | $6.4^{0.0}$ | $1.5^{\downarrow1.0}$ |
| | DB-Swin-T | $90.3_{\pm2.8}^{\uparrow4.0}$ | $61.3_{\pm1.4}^{\uparrow0.9}$ | $86.3_{\pm1.1}^{\uparrow3.5}$ | $58.8_{\pm2.0}^{\uparrow4.6}$ | $\mathbf{94.9}_{\pm1.9}^{\uparrow4.2}$ | $68.9_{\pm1.2}^{\uparrow2.0}$ | $91.1_{\pm1.3}^{\uparrow3.6}$ | $64.4_{\pm1.3}^{\uparrow4.2}$ | $\mathbf{2.8}^{\downarrow0.3}$ | $6.3^{\downarrow0.1}$ | $1.5^{\downarrow1.0}$ |
| Cascade | X-101 | $88.1_{\pm3.4}$ | $61.7_{\pm1.7}$ | $83.1_{\pm1.4}$ | $55.9_{\pm1.2}$ | $93.7_{\pm3.6}$ | $70.1_{\pm1.5}$ | $87.2_{\pm1.4}$ | $61.3_{\pm1.0}$ | 3.6 | 7.0 | 1.8 |
| | Swin-S | $88.9_{\pm2.4}^{\uparrow0.8}$ | $61.6_{\pm1.9}^{\downarrow0.1}$ | $85.9_{\pm1.5}^{\uparrow2.8}$ | $57.5_{\pm1.4}^{\uparrow0.3}$ | $94.0_{\pm2.6}^{\uparrow3.2}$ | $68.7_{\pm1.6}^{\uparrow1.4}$ | $90.4_{\pm1.5}^{\uparrow3.2}$ | $62.9_{\pm1.1}^{\uparrow1.6}$ | $3.8^{\uparrow0.2}$ | $6.4^{\uparrow0.6}$ | $1.3^{\downarrow0.5}$ |
| | DB-Swin-S | $\mathbf{90.4}_{\pm2.4}^{\uparrow2.3}$ | $\mathbf{64.1}_{\pm1.4}^{\uparrow2.4}$ | $\mathbf{87.2}_{\pm1.7}^{\uparrow4.1}$ | $\mathbf{60.0}_{\pm1.7}^{\uparrow4.1}$ | $94.5_{\pm1.9}^{\uparrow0.8}$ | $\mathbf{70.4}_{\pm1.1}^{\uparrow0.3}$ | $\mathbf{91.5}_{\pm1.4}^{\uparrow4.3}$ | $\mathbf{65.2}_{\pm1.1}^{\uparrow4.3}$ | $3.6^{0.0}$ | $\mathbf{6.1}^{\downarrow0.9}$ | $\mathbf{1.1}^{\downarrow0.7}$ |

Note: ↑,↓: gain/decrease from the baseline model; ±: standard deviation over the five folds.

**Tree segmentation:** Pixel-level segmentation allows for a rich and detailed understanding of image content, which can play an important role in precisely delimiting the boundaries of individual trees [Liu et al., 2020]. Our best model, Cascade Mask R-CNN DB-Swin-S, achieves 87.2 % $AP_{50}^{seg}$. One can also observe that models with high $AP^{seg}$ performances achieve significantly better $AP^{bb}$ and $AR^{bb}$. While these findings could simply be attributed to the model's architecture, a previous study [He et al., 2017] demonstrated that increased bounding box and keypoint detections

occur by learning features specific to segmentation. From an application standpoint, this is promising for stock volume estimation, where Liu et al. [2019] successfully did it with a lower (86.0 %) $AP_{50}^{seg}$ than us. Other applications, such as visual assistance for machine operators, can also benefit from pixel segmentation.

**Bounding box precision recall:** The quality of our tree detector can be measured by its ability to find all ground-truth trees (FN=0, high recall), while identifying only relevant trees (FP=0, high precision) [Padilla et al., 2021]. To quantitatively analyze this criterion, we show the precision $\times$ recall curve for bounding box detection in Figure 7. A large area under the curve tends to indicate both high precision and high recall. We observe that predictions with an IoU $= 0.5$ and recalls of up to 60 % obtain close to 100 % of precision. By incrementally increasing IoU, we see a significant decrease when IoU $> 0.75$, and this indicates that most bounding box predictions are within IoU $= [0.5, 0.75]$. Hence, bounding box predictions from our model can achieve high precision and recall rates at the expense of a lesser IoU. This is interesting, because our keypoint estimation method does not directly depend on bounding box localization accuracy. In other words, our method benefits from a certain margin of error on bounding box predictions without it affecting keypoint accuracy.
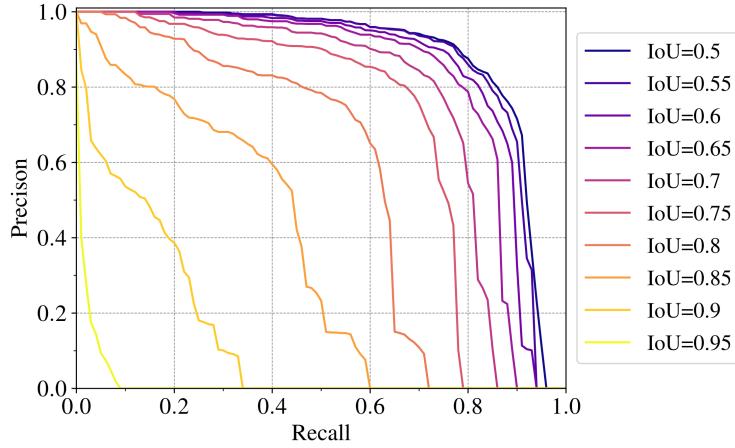


Figure 7: Precision-recall curves of predicted bounding box on trees. For each IoU threshold, we observe an expected phenomenon where the model's precision and recall diminish with each increase in IoU threshold. In our case, the precision and recall values drop steeply for IoUs above 0.75. These results are from our best performing model, Cascade Mask R-CNN DB-Swin-S.

**Felling cut position:** The qualitative keypoint estimation results can be observed in Figure 6. To obtain the error in physical units, we simply used the associated depth image gathered with the ZED stereo camera to compute the 3D spatial position of the estimated keypoint versus its ground truth. From Table 2, we can observe that using keypoints, our models can estimate the felling cut position with an average error less than 7 cm. This accuracy is, at the moment, borderline for autonomous tree felling operations, since Lindroos et al. [2015] estimate that accuracy in the range of cm to mm is essential for efficient autonomous crane movement. Fortunately, we show in Figure 8 that this error, with the image frame as referential, is mostly distributed vertically ($y$) rather than horizontally ($x$), with the $y$ error distribution twice as wide as the $x$ error. In practice, this means that the felling head would grip the tree in its center, but at a lower or higher position than its target height. Intuitively, the vertical error can be attributed to the difficulty of estimating a height from the ground. Moreover, scenes with dense understorey foliage blocking tree base visibility lead to keypoint predictions that are higher up the trunk. This phenomenon can be observed in Figure 8, where the vertical error is more densely located in the $y$ half-plane with a median offset of +1.86 cm. Notably, this vertical error $y$ can be partially mitigated by tapping the head of a feller buncher on the ground, and then pressing forward [Ireland and Kerr, 2008].

**Diameter estimation:** From Table 2, the diameter estimation error is well under 4 cm. The diameter error is almost half the size of the felling cut error (7 cm). We believe the diameter estimation is facilitated by its relative independence from the height at which it is measured. We can also see in Figure 9 that the diameter estimation error increases in relation to the distance between the camera and tree. This comes from the fact that keypoint estimation is done in pixel space, so the error grows linearly as a function of distance when converted to meters. Overall, the accuracy of our diameter estimation method is comparable to similar research using lidar, such as Tremblay et al. [2020] who report 2.04 cm DBH error for well-spaced trees and 3.45 cm in dense forests. Nevertheless, a better accuracy would lower the risks of impairing control algorithms and damaging the machine or trees [Lindroos et al., 2015]. Improving keypoint accuracy is important for autonomous tree felling operations, where accurate keypoint estimation translates to moving
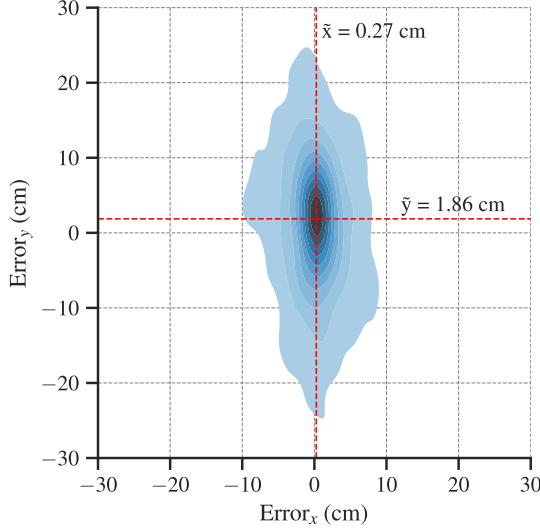
Figure 8: Distribution of felling cut error. The lateral error median ($\tilde{x}$) is mostly centered at the origin. Whereas, the vertical error median ($\tilde{y}$) has a positive offset, and this is possibly caused by the presence of foliage in the lower part of the tree base.

the harvester's head at the correct felling cut position, with the gripper opened wide enough and aligned with the tree's inclination.
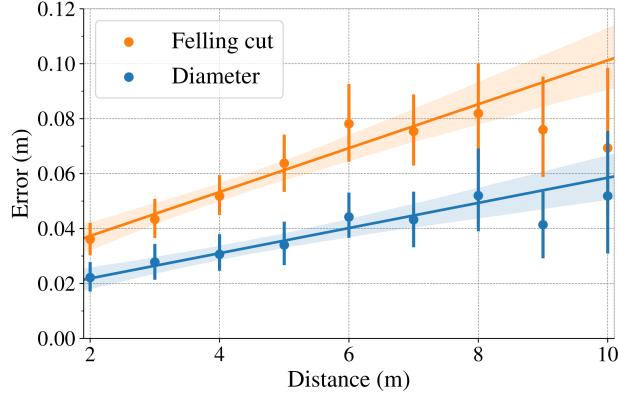


Figure 9: Diameter and felling cut errors relative to tree distance from the camera. We limit ourselves to a range of 10 m, corresponding to the typical reach of machines [Lindroos et al., 2015]. Instances further than that are ignored. Translucent bands around the regression lines correspond to a 95 % confidence interval.

## 4.2   Domain transfer

So far, the experiments suggest that deep learning approaches perform well for tree detection. Using both PORTUGAL and CANATREE100 datasets, we evaluate the domain transfer performance by training on the source domain and directly testing on the target domain.

**Tree detection generalization:** From the results in Table 3, domain transfer is underwhelming, averaging a 22 % $\mathrm{AP}_{50}^{bb}$ and 17 % $\mathrm{AR}_{50}^{bb}$ decrease between the source and target domain. For the Can→Port transfer, the $\mathrm{AP}_{50}^{bb}$ gap exceeds 40 %, which is odd since it qualitatively performs well (see Figure 10). Since a straightforward comparison of network performance between datasets produced by different people is hazardous, the reported results are probably lower than they should be. For instance, the process of image collection and annotation can vary. In our case, we observed that the bounding box annotations are based on different criteria such as minimum tree size, tree distance, and bounding box height. In CANATREE100, the bounding boxes include every visible stem pixel up to the tree top, while the

11

Table 3: Domain transfer comparison between CANATREE100 and PORTUGAL datasets. Training is done on the source domain and directly tested on the target domain, without any fine-tuning. Results are obtained with our best architecture, Cascade Mask R-CNN DB-Swin-S, pre-trained on SYNTHTREE43K.

| Transfer Source→Target | # Training images | Source Domain | | | | Target Domain | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $AP_{50}^{bb}$ | $AP_{50:95}^{bb}$ | $AP_{50}^{seg}$ | $AP_{50:95}^{seg}$ | $AR_{50}^{bb}$ | $AR_{50:95}^{bb}$ | $AR_{50}^{seg}$ | $AR_{50:95}^{seg}$ |
| Port→Can | 100 | $76.7^{\downarrow 18.8}$ | $39.2^{\downarrow 36.7}$ | $85.6^{\downarrow 12.2}$ | $51.3^{\downarrow 30.6}$ | $72.1^{\uparrow 0.1}$ | $35.6^{\downarrow 1.8}$ | $\mathbf{82.0}^{\uparrow 6.8}$ | $45.6^{\uparrow 3.0}$ |
| | 1000 | $88.9^{\downarrow 6.6}$ | $59.6^{\downarrow 16.3}$ | $94.9^{\downarrow 2.9}$ | $68.9^{\downarrow 13.0}$ | $\mathbf{73.6}^{\uparrow 1.6}$ | $\mathbf{42.3}^{\uparrow 4.9}$ | $79.9^{\uparrow 4.7}$ | $\mathbf{49.9}^{\uparrow 7.3}$ |
| | 5000 | $93.3^{\downarrow 2.2}$ | $69.5^{\downarrow 6.4}$ | $96.8^{\downarrow 1.0}$ | $69.7^{\downarrow 12.2}$ | $68.9^{\downarrow 1.1}$ | $31.8^{\downarrow 5.6}$ | $74.8^{\downarrow 0.4}$ | $37.2^{\downarrow 5.4}$ |
| | 14000 | $\mathbf{95.5}$ | $\mathbf{75.9}$ | $\mathbf{97.8}$ | $\mathbf{81.9}$ | 72.0 | 37.4 | 75.2 | 42.6 |
| Can→Port | 60 | 90.4 | 64.1 | 94.5 | 70.4 | 49.1 | 20.0 | 69.2 | 34.0 |

bounding boxes in PORTUGAL dataset are shorter (see Figure 10). Since bounding box size directly impacts IoUs, these annotation disparities can significantly decrease AP and AR results. Moreover, there seems to be missing tree annotations in some of PORTUGAL images, which further cause a drop in AP for the Can→Port transfer.
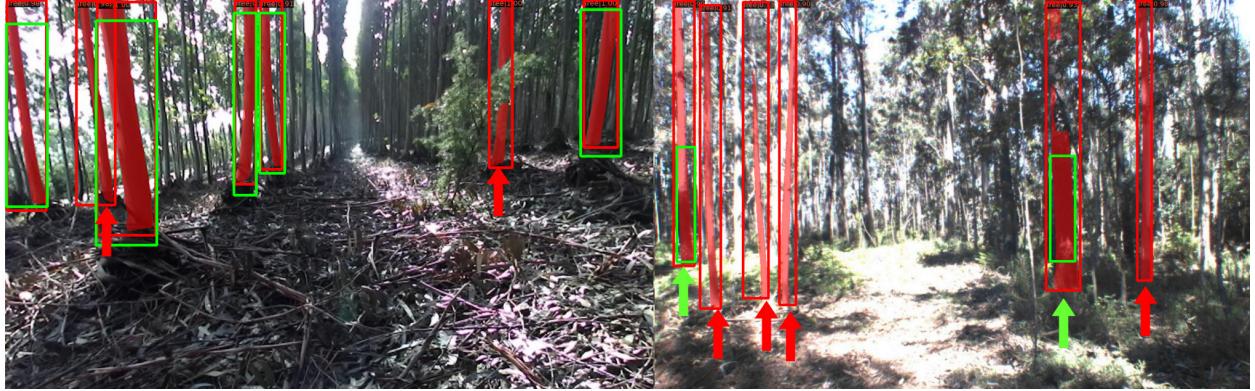


Figure 10: Domain transfer from CANATREE100 to PORTUGAL dataset and examples of annotation disparities. Ground truth annotations are in green and predictions by Cascade Mask R-CNN DB-Swin-S trained on CANATREE100 are in red, and include the predicted segmentation masks. **Left** red arrows point to trees that should be annotated, and **right** green arrows show bounding box annotations so small that they are counted as false negatives (IoU < 0.5). This leads to a decrease in the reported AP and AR results.

**Overfitting:** When training on PORTUGAL as a source domain, we notice a significant overfitting effect across domains. As can be seen in Figure 11, performance on the CANATREE100 dataset degrades beyond a handful of epochs, while it keeps improving on the PORTUGAL source domain. Moreover, we observe that the number of training images from PORTUGAL has a significant impact on source domain performances, but this is not observed on the target domain. Indeed, training on 100 or 14000 images from PORTUGAL yields no significant improvement when testing on CANATREE100 as a target domain. This might be due to the strong visual correlation between samples in the PORTUGAL dataset, as they come from a limited set of videos (we approximate less than 6). In comparison, images from our CANATREE100 are extracted from 33 videos. Note that when using CANATREE100 as a source to test on PORTUGAL, the generalization gap was more or less constant across the epochs from start to end, and it did not change by more than 3 %, in comparison to over 20 % for PORTUGAL.

**Generalization difficulties:** The underwhelming domain transfer performance can partly be attributed to annotation disparities between datasets. Yet, the current experiment still indicates a problematic decrease in perception performance when changing domain. These domain changes are caused by different camera sensor type, the camera lens, variation in forest's vegetation due to geography, etc. Drastic domain changes can lead to generalization difficulties and slow down the adoption of automation in forest operations. This is particularly relevant for international equipment manufacturers, as their machinery is expected to operate in a variety of forest sites and continents. This short experiment shows that there might be a need to sample train images from forests that are in the geographical areas where these machines are to be deployed until better domain adaptation techniques are developed.
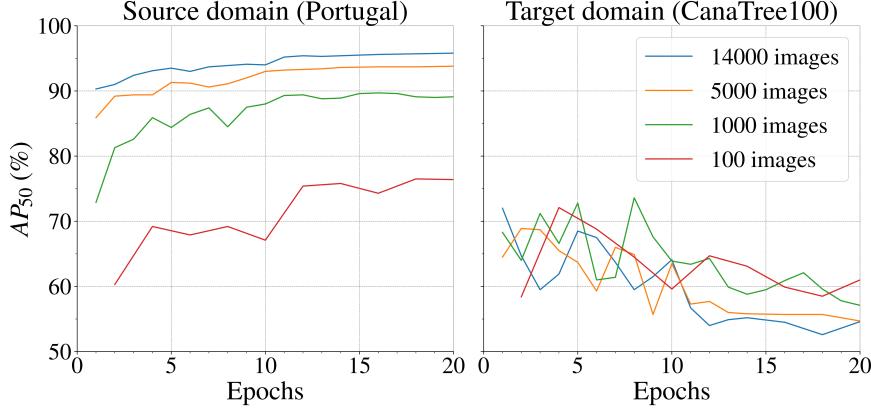
Figure 11: Impact of the number of training epochs and source dataset sizes on generalization, when going from PORTUGAL to CANATREE100. On the PORTUGAL source domain, test performance increases with the number of training images and epochs, as seen on the left. However, when evaluated on CANATREE100 (on the right), we observe the opposite: models that were trained for few epochs generalize better. Moreover, we can see that using more images in the source domain does not improve generalization results. Beyond 10 training epochs, performances degrade when passing from 100 to 14000 images.

### 4.3 Further Experiments

**Pre-training on synthetic images:** Synthetic images have the potential to improve detection results without incurring major costs like time or human resources. In this experiment, we quantitatively measure the usefulness of our synthetic dataset (SYNTHTREE43K) for tree detection and whether it benefits transfer learning. Four pre-training strategies are tested using our Cascade Mask R-CNN DB-Swin-S pre-trained on COCO. In the first case, no subsequent pre-training is performed aside from the initial COCO one. For the other three, we use either PORTUGAL, SYNTHTREE43K or both. Then, the models are fine-tuned on the downstream tasks of CANATREE100.

Table 4: Impact of the datasets used for pre-training. Cascade Mask R-CNN DB-Swin-S model is pre-trained on one or more datasets (C: COCO, S:SYNTHTREE43K, P:PORTUGAL), then fine-tuned and evaluated on CANATREE100. Pre-training on the synthetic dataset (S) improves detection performances, superseding PORTUGAL.

| Pre-training | $AP_{50}^{bb}$ | $AP_{50:95}^{bb}$ | $AP_{50}^{seg}$ | $AP_{50:95}^{seg}$ | $AR_{50}^{bb}$ | $AR_{50:95}^{bb}$ | $AR_{50}^{seg}$ | $AR_{50:95}^{seg}$ |
|---|---|---|---|---|---|---|---|---|
| C | $83.5_{\pm1.2}$ | $59.3_{\pm2.3}$ | $81.0_{\pm2.1}$ | $54.2_{\pm1.6}$ | $87.6_{\pm2.1}$ | $65.7_{\pm2.5}$ | $84.9_{\pm2.9}$ | $60.1_{\pm2.0}$ |
| C+P | $86.6_{\pm1.4}^{\uparrow3.1}$ | $63.2_{\pm1.3}^{\uparrow3.9}$ | $84.7_{\pm1.6}^{\uparrow3.7}$ | $57.5_{\pm0.8}^{\uparrow3.3}$ | $90.9_{\pm2.1}^{\uparrow3.3}$ | $69.7_{\pm1.0}^{\uparrow4.0}$ | $88.7_{\pm1.1}^{\uparrow3.8}$ | $63.0_{\pm0.7}^{\uparrow2.9}$ |
| C+S | $\mathbf{90.4}_{\pm2.4}^{\uparrow6.9}$ | $64.1_{\pm1.4}^{\uparrow4.8}$ | $\mathbf{87.2}_{\pm1.7}^{\uparrow6.1}$ | $\mathbf{60.0}_{\pm1.4}^{\uparrow5.8}$ | $\mathbf{94.5}_{\pm1.9}^{\uparrow6.9}$ | $70.4_{\pm1.1}^{\uparrow4.7}$ | $\mathbf{91.5}_{\pm1.4}^{\uparrow6.6}$ | $\mathbf{65.2}_{\pm1.1}^{\uparrow5.1}$ |
| C+S+P | $89.5_{\pm1.0}^{\uparrow6.0}$ | $\mathbf{65.2}_{\pm0.8}^{\uparrow5.9}$ | $86.3_{\pm1.7}^{\uparrow5.3}$ | $59.9_{\pm1.2}^{\uparrow5.7}$ | $93.8_{\pm1.2}^{\uparrow6.2}$ | $\mathbf{71.4}_{\pm0.5}^{\uparrow4.7}$ | $90.8_{\pm1.8}^{\uparrow5.9}$ | $\mathbf{65.2}_{\pm1.1}^{\uparrow5.1}$ |

From Table 4, we can see that the most effective pre-training strategy is based on synthetic images, outperforming COCO-only and PORTUGAL by at least 3.8 % $AP_{50}^{bb}$, and 3.6 $AR_{50}^{bb}$. It also shows that adding images from the PORTUGAL dataset does not provide significant improvements when synthetic images are already used. This clearly indicates that pre-training on a large-scale, densely annotated synthetic dataset, such as SYNTHTREE43K, is a great starting point to learn powerful and robust representation for tree detection.

**Robustness to occlusion:** In forests, computer vision algorithms are subject to many distractors such as branches, bushes or the machinery itself causing occlusion. Therefore, we conduct a pixel-level analysis to measure the robustness of our method when faced to different levels of occlusion.

From Figure 12, we observe an increasing rate of FN (missed detections) relative to the percentage of occlusion. For tree occlusion, the FN rate grows linearly, but for the tree base, it spikes when the base is over 90 % occluded. This indicates that tree base occlusion is less impactful than tree occlusion itself, which can be useful in the presence of dense understorey foliage. In fact, although at least 30 % of each tree is visible to the camera (our annotation threshold), there is only 14 % FN in cases where the tree base area is over 90 % occluded. In light of these findings, it is without surprise that the ideal conditions for tree detection are well-maintained forests with minimal occlusion. Even though
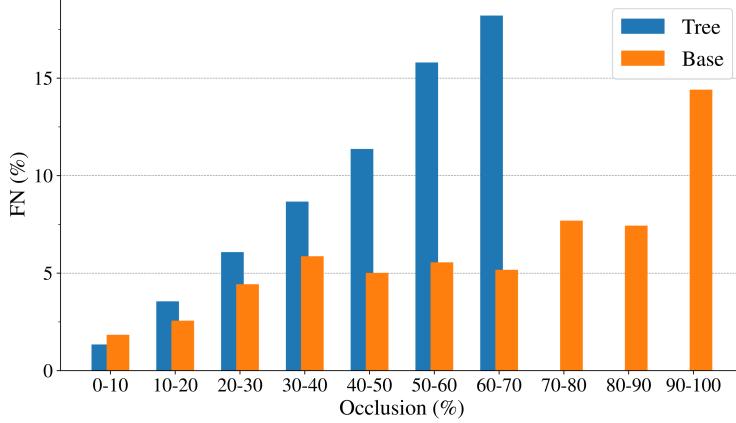
Figure 12: Impact of occlusion on tree detection, for synthetic images, in terms of FN rates. Occlusion of the tree base (0 to 30 cm above ground) has significantly less impact than for the entire trees. Without surprise, the lowest FN rate corresponds to the lowest occlusion level.

this experiment is conducted on synthetic images, due to the difficulty of measuring occlusion in real images, it provides key insights on how robust our method can be when faced with occlusion in real images.

**Dataset size:** Deep learning is well known for scaling with large amounts of data. This experiment seeks to quantitatively estimate the performance improvements if we were to increase the number of annotated samples in CANATREE100. To do so, we vary the train set size between 20 and 80 images, keeping the same 20 image test set for each fold. As shown by the logarithmic abscissa scale in Figure 13, we found that detection improvement follows a power law. Each time the number of training images doubles, we gain between 1.1 to 1.6 % $AP_{50}^{bb}$, and 1.9 % $AP_{50}^{seg}$. Moreover, by pre-training SYNTHTREE43K, we remark a near-constant gain of 4.0 % $AP_{50}^{bb}$ and 4.0 % $AP_{50}^{seg}$ compared to the model pre-trained solely on COCO. This further consolidates the fact that our synthetic forest dataset can be leveraged to achieve better performances, compared to training a limited number of real images only. Markedly, Figure 13 illustrates that pre-training on synthetic images and fine-tuning on only 20 images outperforms the model pre-trained on COCO-only and fine-tuned on 80 (four times more) real images. In a low-data regime like ours, this can allow us to partially fill the data gap and avoid hand annotation.
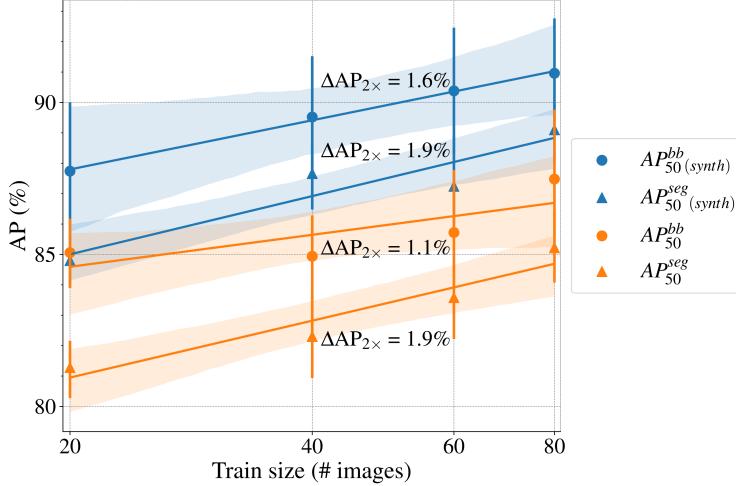


Figure 13: Performances scale with the train set size. Reported results are with and without pre-training on synthetic data, fine-tuned and tested on CANATREE100. The slope of $\Delta$AP represents the change in AP each time the train set size doubles. Translucent bands around the regression lines correspond to a 95 % confidence interval.

**Keypoint branch improvements** In order to implement autonomous movements, accurate end-effector target position is desirable. Here, we seek to reduce the felling cut error obtained with the CNN keypoint head by replacing it with

TransPose [Yang et al., 2021], a transformer-based keypoint head. The performance comparison between these two keypoint head architectures is presented in Table 5, where we observe a 20 % (1 cm) error reduction for felling cut predictions made with the TransPose approach.

Table 5: Comparison between network head architectures for keypoint predictions. TransPose achieves more accurate keypoint predictions. However, both architectures fail to scale with depth. Model is Mask R-CNN Swin-S pre-trained on the synthetic dataset and fine-tuned on CANATREE100 dataset.

| Keypoint Head | Architecture (channel×depth) | Keypoint error | | |
|---|---|---|---|---|
| | | $\text{dia}_{(cm)}$ | $\text{fc}_{(cm)}$ | $\text{inc}_{(\circ)}$) |
| CNN | $256 \times 8$ | 3.1 | 6.3 | 1.7 |
| | $512 \times 8$ | 3.3 | 6.5 | 1.9 |
| | $256 \times 16$ | 2.9 | 6.5 | **1.4** |
| TransPose | $512 \times 4$ | 3.0 | 5.6 | 1.87 |
| | $512 \times 6$ | **2.8** | 5.6 | 1.6 |
| | $1024 \times 6$ | 3.0 | 5.6 | 1.7 |
| | $512 \times 8$ | 3.0 | **5.4** | 1.5 |

However, we did not see the performance scale with increased depth, but we observe a saturation or decrease in performance. This indicates that the attention mechanism at the core of transformers is responsible for the increased accuracy and is indeed effective at capturing pertinent information from different image locations. For instance, given a query location (keypoint) and an attention map, one can explicitly reveal the dependency areas on which each keypoint prediction is based. As can be visually confirmed in Figure 14, keypoints in the tree base area have obvious dependencies on visual cues from their local area, while keypoints higher up the stem rely on long-range cues. This potentially allows the model to localize partially obstructed keypoints by looking for more significant visual cues, thus reducing its reliance on the occluded pixels.
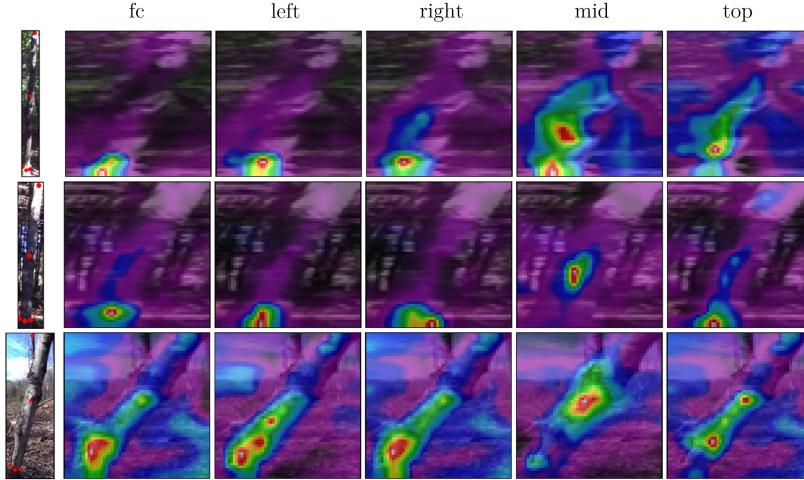


Figure 14: Heatmap of attention scores for a particular keypoint position. The felling cut, left and right diameter keypoints have obvious local dependencies, whereas mid and top keypoints exploit long range image cues. Red areas indicate higher attention scores. The rectangular images on the left are stretched into a square one by the RoIAlign layer.

## 5   Conclusion and Future Works

In this paper, we demonstrate that supervised end-to-end deep learning can successfully detect trees in forested environments. For this, we tackle the problematic lack of datasets by creating and publicly releasing two novel forest image datasets — one synthetic (SYNTHTREE43K) and one real (CANATREE100). Based on the intuitive paradigm of pre-training on a large synthetic image dataset followed by fine-tuning on a real image dataset, we achieve promising results in all three major perception tasks: 90.4 % tree detection precision, 87.2 % segmentation precision and cm

accurate keypoint estimations. Notably, we show that diameter estimation using a vision-based approach achieves comparable results to methods based on lidar, while also providing felling cut location and tree inclination.

In terms of shortcomings, the limited capacity to generalize to other forest sites suggests that training on specific forest sites is required before practical applications. Given the relative ease of training, there is optimism in the scalability of our approach to dataset size and model architecture. We show that detection performances scale with dataset size, follow a power law that increases by 1.1 to 1.6 % $\text{AP}_{50}^{bb}$, and 1.9 % $\text{AP}_{50}^{seg}$ every time the training set size doubles. In addition, changing the model backbone, architecture, or pre-training dataset can give direct performance gains. Again, these findings outline the need for much larger, rigorously annotated forest datasets. Hopefully, the paper has sufficiently highlighted the pertinence of releasing datasets so others can benchmark on it and allow the community to clearly establish a state-of-the-art method for tree detection.

Future works should investigate the creation of a larger real image dataset, but also ways to enhance model generalization to different forest sites, and conduct evaluations on independent test sets [Diez et al., 2021]. Since the reliance on annotated datasets is taking a toll on progresses made in fields like forestry, which receive far less attention from the deep learning community, semi-supervised or self-supervised learning methods are worth investigating.

# References

O. Ringdahl. *Automation in Forestry – Development of Unmanned Forwarders*. PhD thesis, Department of Computing Science Umeå University, 2011.

R. Parker, K. Bayne, and P.W Clinton. Robotics in forestry. *NZ Journal of Forestry*, 60(4):9, 2016.

M.D Ortiz, S. Westerberg, P.X La Hera, U. Mettin, L. Freidovich, and A.S Shiriaev. Increasing the level of automation in the forestry logging process with crane trajectory planning and control. *Journal of Field Robotics*, 31(3):343–363, 2014.

R. Visser and O.F Obi. Automation and robotics in forest harvesting operations: Identifying near-term opportunities. *Croatian Journal of Forest Engineering: Journal for Theory and Application of Forestry Engineering*, 42(1):13–24, 2021.

N. Roy, I. Posner, T. Barfoot, P. Beaudoin, Y. Bengio, J. Bohg, O. Brock, I. Depatie, D. Fox, D. Koditschek, and T. Lozano-Perez. From machine learning to robotics: Challenges and opportunities for embodied intelligence. *arXiv preprint arXiv:2110.15245*, 2021.

C. Thorpe and H. Durrant-Whyte. Field robots. In *Proceedings of the 10th International Symposium of Robotics Research (ISRR'01)*, 2001.

T. Nurminen, H. Korpunen, and J. Uusitalo. Time consumption analysis of the mechanized cut-to-length harvesting system. 2006.

R. Padilla, W.L Passos, T.LB Dias, S.L Netto, and E.AB da Silva. A comparative analysis of object detection metrics with a companion open-source toolkit. *Electronics*, 10(3):279, 2021.

Y. Diez, S. Kentsch, M. Fukuda, M.LL Caceres, K. Moritake, and M. Cabezas. Deep learning in forestry using uav-acquired rgb data: A practical review. *Remote Sensing*, 13(14):2837, 2021.

Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.

A. Oliver, A. Odena, C. Raffel, E.D Cubuk, and I.J Goodfellow. Realistic evaluation of semi-supervised learning algorithms. In *ICLR (Workshop)*, 2018.

C. Bowen, P. Omkar, and K. Alexander. Pointly-supervised instance segmentation. *arXiv*, 2021.

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 *IEEE European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014.

S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis. Sim-to-real via sim-to-sim: data-efficient robotic grasping via randomized-to-canonical adaptation networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12627–12637, 2019.

K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *IEEE International Conference on Computer Vision*, pages 2961–2969, 2017.

Z. Cai and N. Vasconcelos. Cascade R-CNN: high quality object detection and instance segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1483–1498, 2019.

T. Hellström, P. Lärkeryd, T. Nordfjell, and O. Ringdahl. Autonomous forest vehicles: Historic, envisioned, and state-of-the-art. *International Journal of Forest Engineering*, 20(1):31–38, 2009.

E. Jelavic, D. Jud, P. Egli, and M. Hutter. Towards autonomous robotic precision harvesting. *arXiv:2104.10110*, 2021.

C. Zhang, L. Yong, Y. Chen, S. Zhang, L. Ge, S. Wang, and W. Li. A rubber-tapping robot forest navigation and information collection system based on 2D lidar and a gyroscope. *Sensors*, 19(9):2136, 2019.

M. Pierzchała, P. Giguère, and R. Astrup. Mapping forests using an unmanned ground vehicle with 3d lidar and graph-slam. *Computers and Electronics in Agriculture*, 145:217–225, 2018.

J.-F Tremblay, M. Béland, R. Gagnon, F. Pomerleau, and P. Giguère. Automatic three-dimensional mapping for tree diameter measurements in inventory operations. *Journal of Field Robotics*, 37(8):1328–1346, 2020.

X. Liang, A. Jaakkola, Y. Wang, J. Hyyppä, E. Honkavaara, J. Liu, and H. Kaartinen. The use of a hand-held camera for individual tree 3d mapping in forest sample plots. *Remote Sensing*, 6(7):6587–6603, 2014.

D.-Q da Silva, F.-N Dos Santos, A.-J Sousa, and V Filipe. Visible and thermal image-based trunk detection with deep learning for forestry mobile robotics. *Journal of Imaging*, 7(9):176, 2021.

J. Wang, X. Chen, L. Cao, F. An, B. Chen, L. Xue, and T. Yun. Individual rubber tree segmentation based on ground-based lidar data and Faster R-CNN of deep learning. *Forests*, 10(9):793, 2019.

K. Itakura and F. Hosoi. Automatic tree detection from three-dimensional images reconstructed from 360 spherical camera using yolo v2. *Remote Sensing*, 12(6):988, 2020.

J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. arxiv 2016. *arXiv preprint arXiv:1612.08242*, 394, 2016.

J. Liu, X. Wang, and T. Wang. Classification of tree species and stock volume estimation in ground forest images using deep learning. *Computers and Electronics in Agriculture*, 166:105012, 2019.

O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention MICCAI*, pages 234–241, 2015.

A. Ostovar, B. Talbot, S. Puliti, R. Astrup, and O. Ringdahl. Detection and classification of root and butt-rot (RBR) in stumps of norway spruce using rgb images and machine learning. *Sensors*, 19(7):1579, 2019.

S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017.

A. Gaidon, Q. Wang, Y. Cabon, and E. Vig. Virtual worlds as proxy for multi-object tracking analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4340–4349, 2016.

Y. Cabon, N. Murray, and M. Humenberger. Virtual kitti 2. *arXiv preprint arXiv:2001.10773*, 2020.

A. Wang, Y. Sun, A. Kortylewski, and A.L Yuille. Robust object detection under occlusion with context-aware compositional nets. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12645–12654, 2020.

J. Brooks. COCO Annotator. `https://github.com/jsbroks/coco-annotator/`, 2019.

S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1492–1500, 2017.

M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.

Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.

T. Liang, X. Chu, Y. Liu, Y. Wang, Z. Tang, W. Chu, J. Chen, and H. Ling. Cbnetv2: A composite backbone network architecture for object detection. *arXiv preprint arXiv:2107.00420*, 2021.

K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C.C Loy, and D. Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.

D. Erhan, A. Courville, Y. Bengio, and P. Vincent. Why does unsupervised pre-training help deep learning? In *International conference on artificial intelligence and statistics*, pages 201–208. JMLR Workshop, 2010.

D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *European Conference on Computer Vision (ECCV)*, pages 181–196, 2018.

S. Hinterstoisser, V. Lepetit, P. Wohlhart, and K. Konolige. On pre-trained image features and synthetic images for deep learning. In *European Conference on Computer Vision (ECCV) Workshops*, September 2018.

L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen. Deep learning for generic object detection: A survey. *International Journal of Computer Vision*, 128(2):261–318, 2020.

O. Lindroos, O. Ringdahl, P. La Hera, P. Hohnloser, and T. Hellström. Estimating the position of the harvester head–a key step towards the precision forestry of the future? *Croatian Journal of Forest Engineering: Journal for Theory and Application of Forestry Engineering*, 36(2):147–164, 2015.

D. Ireland and G. Kerr. Ccf harvesting method development: Harvester head visibility. 2008.

S. Yang, Z. Quan, M. Nie, and W. Yang. Transpose: Keypoint localization via transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11802–11812, 2021.