

PointRAFT: 3D deep learning for high-throughput prediction of potato tuber weight from partial point clouds[★]

Pieter M. Blok^{a,*}, Haozhou Wang^a, Hyun Kwon Suh^b, Peicheng Wang^a, James Burridge^a and Wei Guo^a

^aGraduate School of Agricultural and Life Sciences, The University of Tokyo, 1-1-1 Midori-cho, Nishitokyo-city, 188-0002, Tokyo, Japan

^bDepartment of Integrative Biological Sciences and Industry, Sejong University, 209 Neungdong-ro, 05006, Seoul, Republic of Korea

ARTICLE INFO

Keywords:

Potato
3D Deep Learning
RGB-D
Point Cloud
Regression

ABSTRACT

Potato yield is a key indicator for optimizing cultivation practices in precision agriculture. Potato yield can be estimated directly on a harvester using RGB-D cameras, which capture three-dimensional (3D) information of individual tubers moving along the conveyor belt. A major challenge, however, is that the 3D point clouds reconstructed from RGB-D images are incomplete due to self-occlusion, leading to systematic underestimation of tuber weight. To overcome this limitation, we introduce PointRAFT, a high-throughput point cloud regression network that directly predicts continuous 3D shape properties, such as tuber weight, from partial point clouds. Rather than reconstructing complete 3D geometry, PointRAFT infers target values directly from raw 3D data. Its key architectural novelty is an object height embedding that incorporates tuber height as an additional geometric cue, improving regression performance under practical harvesting conditions. PointRAFT was trained and evaluated on a dataset of 26,688 partial point clouds collected from 859 potato tubers across four cultivars and three growing seasons on an operational harvester in Japan. On a test set of 5,254 point clouds representing 172 unique tubers, PointRAFT achieved a mean absolute error (MAE) of 12.0 g and a root mean squared error (RMSE) of 17.2 g, substantially outperforming a linear regression baseline with an MAE of 23.0 g and an RMSE of 31.8 g. The proposed height embedding reduced RMSE by 30% compared to a standard PointNet++ regression network. With an average analysis time of 6.3 ms per point cloud, PointRAFT enables processing rates of up to 150 tubers per second, meeting the high-throughput requirements of commercial potato harvesters. Beyond potato weight estimation, PointRAFT provides a versatile regression network applicable to a wide range of 3D phenotyping and robotic perception tasks. The code, network weights, and a subset of the dataset are publicly available at <https://github.com/pieterblok/pointraft.git>.

1. Introduction

Potatoes (*Solanum tuberosum*) are an important component of the human diet, as they provide high-energy carbohydrates, vitamin C, and dietary fibers (Camire, Kubow and Donnelly, 2009). To safeguard the role of potatoes in human nutrition, further optimization of potato production is needed (Zhang, Xu, Wu, Hu and Dai, 2017). A major step toward this improvement is through precision agriculture. Precision agriculture enables site-specific application of fertilizers and crop protection products, which leads to higher yields, lower costs, and reduced environmental pressure (Bullock, Lowenberg-DeBoer and Swinton, 2002; Van Evert, Gaitán-Cremaschi, Fountas and Kempenaar, 2017). To steer precision agriculture practices, detailed information on potato yield is required. In current practice, potato yield mapping can be performed using load cells attached to the harvester's conveyor belt to measure the mass of harvested produce in real time (Zamani, Ghoşamiparashkahi, Faghavi and Ghezavati, 2014; Kabir, Myat Swe, Kim, Chung, Jeong and Lee, 2018). Although load-cell systems are easy to use and maintain, they suffer a major limitation: they measure

gross mass, including tare such as soil clods, stones, and plant residue. The inclusion of tare can lead to overestimation of tuber yield, particularly in areas where large amounts of soil or crop residue are harvested together with the potato tubers.

A more accurate alternative is the use of camera-based yield monitoring systems, which can visually distinguish potato tubers from tare. Such systems have been explored in the scientific literature since the early 2000s (Noordam, Otten, Timmermans and van Zwol, 2000; Hofstee and Molema, 2003; ElMasry, Cubero, Moltó and Blasco, 2012; Razmjoo, Mousavi and Soleymani, 2012; Lee, Kim, Lee and Shin, 2018; Long, Wang, Zhai, Wu, Li, Sun and Su, 2018; Si, Sankaran, Knowles and Pavek, 2018; Su, Kondo, Li, Sun, Al Riza and Habaragamuwa, 2018; Pandey, Kumar and Pandey, 2019; Cai, Jin, Xu and Yang, 2020; Lee and Shin, 2020; Dolata, Wróblewski, Mrzygłód and Reiner, 2021; Huynh, TonThat and Dao, 2022; Jang, Moon, Kim and Lee, 2023; Emwinghare, Al-Mallahi, Esau and Campelo, 2025). Most systems rely on RGB color cameras and estimate tuber properties such as size, volume, or weight using pixel-to-world conversions, which are calibrated under controlled laboratory conditions. Unfortunately, these externally calibrated conversions are prone to errors when the height profile of the tubers on the conveyor belt varies or when tubers are stacked or occluded. These scenarios are commonly encountered during practical harvesting.

[★]This study is funded by the Sarabetsu Village "Endowed Chair for Field Phenomics" project in Hokkaido, Japan.

*Corresponding author: pieter.blok@fieldphenomics.com (P.M. Blok).
ORCID(s): 0000-0001-9535-5354 (P.M. Blok); 0000-0001-6135-402X (H. Wang); 0000-0003-4771-9365 (H.K. Suh); 0000-0002-2194-9894 (J. Burridge); 0000-0002-3017-5464 (W. Guo)

RGB-Depth (RGB-D) cameras offer a promising solution. These cameras are inexpensive, easy to integrate, and generate both a color image and a depth image, which allows more accurate conversions. The color and depth images can also be fused into (colored) 3D point clouds from which tuber size, volume, or weight can be estimated using 3D algorithms. RGB-D-based point clouds, however, suffer from self-occlusion, meaning that the captured 3D shape is only partially observed due to the single viewpoint of the RGB-D camera (Fu, Gao, Wu, Li, Karkee and Zhang, 2020). Self-occlusion leads to underestimation of tuber size, volume, or weight, which is undesirable for yield monitoring. Multi-camera configurations can reduce this effect but cannot fully eliminate it, and their integration on a potato harvester remains challenging.

The partial observability of RGB-D point clouds has driven recent research interest in 3D shape completion algorithms (Fei, Yang, Chen, Li, Li, Ma, Hu and Ma, 2022). Most shape completion algorithms use deep learning techniques to infer the full 3D geometry of an object from a partial observation. These methods are based on voxel-based networks, point cloud-based architectures, mesh-based models, or implicit representations (Fei et al., 2022). Although many completion methods exist, including some developed for agricultural products (Ge, Xiong and From, 2020; Marangoz, Zaenker, Menon and Bennewitz, 2022; Chen, Liu, Wang, Wang, Gong, Li and Lan, 2023; Pan, Magistri, Läbe, Marks, Smitt, McCool, Behley and Stachniss, 2023; Xu, Chen and Jing, 2023; Magistri, Marcuzzi, Marks, Sodano, Behley and Stachniss, 2024; Wang, Lin, Shui, Zhang and Zhai, 2025b; Zhang, Zheng, Wang, Blok, Wang and Guo, 2025; Wang, Pan, Magistri, Kooistra, Stachniss, Wang and Valente, 2026), only two methods currently offer the millisecond-level processing speed required for use on a potato harvester: CoRe (Magistri, Marks, Nagulavancha, Vizzo, Läbe, Behley, Halstead, McCool and Stachniss, 2022) and CoRe++ (Blok, Magistri, Stachniss, Wang, Burridge and Guo, 2025).

CoRe and CoRe++ themselves face three challenges that limit their practical implementation on a potato harvester. First, training these methods requires large datasets of paired partial and complete point clouds, which are costly to obtain. The main bottleneck is the time-intensive generation of full 3D data using methods such as Gaussian Splatting or Structure-from-Motion, which typically takes about ten minutes per potato. This is followed by time-consuming object registrations using approaches such as those described in Wang, Blok, Burridge, Jiang, Miyauchi, Miyamoto, Tanaka and Guo (2025a). Second, both CoRe and CoRe++ rely on 4-channel RGB-D images, which are sensitive to camera placement and tend to generalize less robustly than raw 3D point cloud data. Third, CoRe and CoRe++ require post-processing steps to convert their output point clouds into watertight meshes for volume extraction. This post-processing adds computational overhead and depends on mesh generation parameters that generalize poorly.

Beyond 3D shape completion methods, another line of research has emerged that bypasses explicit 3D reconstruction entirely and performs direct regression on partial point clouds. With point cloud regression networks, a continuous target value, like volume or weight, can be directly estimated from 3D points using the network's dedicated regression head. Consequently, this eliminates the need for paired complete 3D shapes, improves generalizability across camera placements, removes the requirement for mesh-based post-processing, and directly estimates the weight, without relying on intermediate volume estimation. Point cloud regression networks have proven successful in estimating body weight from blanket-covered hospital patients (Bigalke, Hansen, Diesel and Heinrich, 2021) and human body fat percentage estimation (Hu, Dai, Zhao, Wang, Ma and Munteanu, 2023). In animal science, point cloud regression using PointNet (Qi, Su, Mo and Guibas, 2017a) and PointNet++ (Qi, Yi, Su and Guibas, 2017b) has been successfully applied to estimate the weight of beef cattle (Hou, Huang, Zhang and Miao, 2023), pigs (Paudel, de Sousa, Sharma and Brown-Brandl, 2024), and chickens (Zheng, Ma, Liu, Huang, Chen, Fang, Yang, Berckmans, Norton and Zhang, 2025). These studies demonstrate that complete geometry is not strictly necessary for accurate prediction of 3D shape properties.

While the above literature demonstrates that point cloud regression can replace full 3D reconstruction, existing methods are primarily developed for humans and animals, domains where objects are large, smoothly varying, and real-time processing is often not critical. In contrast, potato tubers are small, highly irregular, often piled on top of each other, and move rapidly along a harvester's conveyor belt. These conditions require a point cloud regression network designed for small, irregular objects operating at high-throughput and low latency.

To address this gap, we introduce a new algorithm called PointRAFT (Regression of Amodal Full Three-dimensional shape properties from partial point clouds). PointRAFT is a high-throughput 3D regression network that directly predicts continuous target values from a partial point cloud. PointRAFT's architectural novelty is an object height embedding that incorporates tuber height as an additional geometric cue for weight regression. We hypothesize that explicitly encoding tuber height improves weight estimation accuracy from partial point clouds under operational harvester conditions. Beyond estimating 3D shape properties such as weight, diameter or volume, PointRAFT can also be used as a generic encoder within 3D shape completion architectures, such as CoRe and CoRe++. To support the broader application of PointRAFT, we have publicly released it at <https://github.com/pieterblok/pointraft.git>. We also released a subset of our 3D dataset via Wang et al. (2025a), consisting of 9,658 partial point clouds with ground truth weights for 339 individual potato tubers across three cultivars, to facilitate further research on 3D point cloud algorithms.

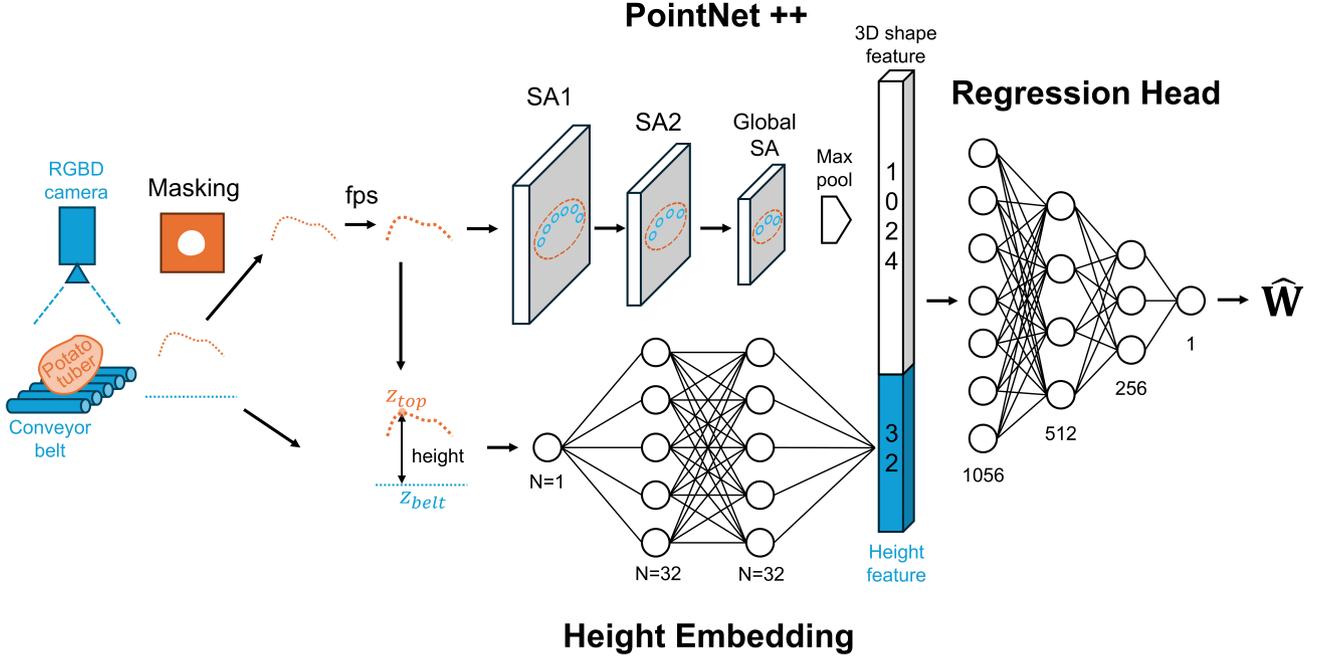


Figure 1: Schematic representation of the PointRAFT network, which directly regresses continuous 3D shape properties, such as tuber weight, from partial point clouds. PointRAFT’s architectural novelty lies in an object height embedding, which provides an additional geometric cue for weight estimation.

2. Materials and methods

2.1. PointRAFT

To estimate the weight of potato tubers from partial point clouds, we propose PointRAFT. While its overall architecture follows the hierarchical set abstraction (SA) principle of PointNet++ (Qi et al., 2017b), we added a novel object height embedding to improve the weight regression. Additionally, a PyTorch-based point cloud downsampling method was implemented to enable high-throughput processing. The entire framework was developed using the PyTorch Geometric software library (version 2.7.0), and a schematic representation of the architecture is shown in Figure 1.

2.1.1. Preprocessing of partial point clouds

Each potato tuber is represented by a 3D point cloud with a number of points that depends on the tuber’s physical size and distance to the sensor. Larger tubers typically generate point clouds with several thousand points, while smaller tubers produce sparser clouds. To allow stable training and efficient batching, these raw point clouds were downsampled and preprocessed before being passed to the network.

A common downsampling method is farthest point sampling (FPS), which produces a spatially uniform subset of points and preserves the overall geometry of the object. FPS can be computationally expensive because it repeatedly computes distances across the full point cloud, making it a bottleneck for larger point clouds. To address this limitation, we implemented a PyTorch-based FPS downsampling

method as a novel component of our pipeline, implemented using the `pytorch_fpsample` software library (Han, Wang, Xiao, Zhang, Zhang, Xu and Zhu, 2023). This software allowed efficient downsampling while retaining the geometric advantages of FPS. In our implementation, each tuber was downsampled to a fixed set \mathbf{P}_0 of 1024 points (Equation 1).

$$\mathbf{P}_0 = \{(x_i, y_i, z_i)\} \in \mathbb{R}^{N \times 3} \quad \text{for } i = 1, 2, \dots, N \quad (1)$$

In addition to downsampling, we applied PyTorch Geometric’s Center transform to each point cloud, shifting all points so that the centroid of the transformed cloud aligns with the origin (0,0,0) of the camera coordinate frame. This ensured that each point cloud was centered consistently, eliminating potential biases caused by the absolute positions of tubers in the sensor’s coordinate frame. We did not normalize the scale of the point clouds to preserve the relative sizes of the tubers. This ensured that larger tubers remained larger in the network’s representation, allowing the model to learn geometric cues related to size. The centering and PyTorch-based FPS downsampling formed the high-throughput preprocessing module of PointRAFT.

2.1.2. Neural network architecture

The backbone of PointRAFT was based on the PointNet++ single scale grouping (SSG) architecture (Qi et al., 2017b). PointNet++ SSG was chosen for its efficiency and its ability to encode hierarchical geometric features using a series of hierarchical set abstraction (SA) modules. Each

SA module follows three main steps: selecting representative points, grouping nearby points, and learning geometric features.

The first SA module downsampled the input point cloud \mathbf{P}_0 by selecting 50% of its points using farthest point sampling. The resulting subset is denoted as \mathbf{P}_1 . For each point in \mathbf{P}_1 , up to 64 neighboring points within a radius of 0.2 m were grouped for feature learning. The radius was chosen to be larger than the physical dimensions of individual potato tubers, so that each local region captured the entire partial point cloud. This was chosen because it enabled object-level feature learning in a hierarchical structure, which was deemed beneficial for self-occluded partial point clouds. Geometric features were then learned using a PointNet-style convolution implemented with shared multi-layer perceptrons (MLPs) with layer sizes [3, 64, 64, 128], producing a 128-dimensional feature vector for each sampled region in \mathbf{P}_1 .

The second SA module further reduced the point set \mathbf{P}_1 to 25% of its points using farthest point sampling, resulting in the point subset \mathbf{P}_2 . For each point in \mathbf{P}_2 , up to 64 neighboring points within a radius of 0.4 m were selected from \mathbf{P}_1 . This radius, again larger than the physical dimensions of individual tubers, allowed the network to learn object-level features from a coarser representation of the points. The input to the MLP at this stage had dimensionality 128+3, where the three additional channels corresponded to the relative 3D coordinates of each selected neighboring point in \mathbf{P}_1 , concatenated with the learned features from the previous SA layer. Shared MLPs with layer sizes [131, 128, 128, 256] then produced a 256-dimensional feature vector for each sampled region in \mathbf{P}_2 , encoding coarser geometric features of the overall shape of the tuber.

The final SA module was implemented as global set abstraction and concatenated the learned features from the second SA module with the relative 3D coordinates in \mathbf{P}_2 , forming an input of size 256 + 3. This input was then processed by a shared MLP with layer sizes [259, 256, 512, 1024], producing 1024-dimensional features for each point. A global max-pooling operation then aggregated these features into a single 1024-dimensional permutation-invariant, global 3D shape feature vector.

To enhance the global 3D shape representation, we implemented a novel object height embedding that encoded the approximate height of each potato tuber. Tuber height can serve as an additional proxy for the tuber's weight, which can improve neural network regression. Unfortunately, tuber height cannot be reliably observed from self-occluded partial point clouds, as the bottom parts of the tuber cannot be observed. Therefore, we embedded the tuber height by approximating it as the difference between the 3D point of the tuber closest to the camera and the known distance to the conveyor belt, assuming that most tubers lie flat on the conveyor belt. The approximated tuber height was encoded using a single MLP with layer sizes [1, 32, 32], producing a 32-dimensional feature vector. This vector was then concatenated with the 1024-dimensional global feature

vector learned by the SA modules, resulting in a combined 1024+32-dimensional vector that captured both 3D shape and tuber height information. The combined vector was then processed by the regression head, implemented as a single MLP with layer sizes [1056, 512, 256, 1] and dropout of 0.5, to predict the tuber weight as a continuous value.

2.2. Dataset

2.2.1. Imaging systems

Our point cloud dataset was acquired using two RGB-D imaging systems, both installed above the conveyor belt of a single-row potato harvester (Toyonoki Top-1, Figure 2a). The two systems shared the same RGB-D camera and data acquisition software, but differed in enclosure design, camera-to-belt distance, and field of view.

The first imaging system was constructed as a closed black plastic enclosure mounted above the conveyor belt, with dimensions 0.85 m (length), 0.45 m (depth), and 0.39 m (height), refer to Figure 2b. Inside the enclosure, four LED strips with a color temperature of 6000 K were installed. Reflective curtains along the sides helped to diffuse the light evenly across the conveyor belt (Figure 2c). An RGB-D camera (Intel RealSense D405) was mounted centrally inside the enclosure, capturing a top-view perspective of the conveyor belt. The distance between the camera and the conveyor belt was approximately 0.35 m, resulting in a field of view of about 0.39 m (length) by 0.64 m (width). The camera's exposure time was set to 5.0 ms to avoid motion blur caused by the fast-moving potato tubers.

The second imaging system was constructed using aluminum T-shaped profiles, forming a rigid enclosure with dimensions 0.85 m (length), 0.48 m (width), and 0.58 m (height) (Figure 3a). The sides of the enclosure were covered with plates to reduce the influence of natural sunlight on the captured images (Figure 3b). A total of 24 LED strips were mounted on the ceiling plate, providing uniform artificial illumination for measurements under both day and night conditions (Figure 3c). Each LED strip was 30 cm long (AIBOO LED) with a color temperature of 6000 K. The same Intel Realsense D405 RGB-D camera was mounted centrally between the LED strips and manually aligned with the conveyor belt. The camera-to-belt distance was approximately 0.46 m, resulting in a field of view of about 0.52 m (length) by 0.75 m (width). As with the first system, the camera's exposure time was fixed at 5.0 ms to ensure motion-blur-free image acquisition.

In both imaging systems, the camera streamed color and depth images at 30 frames per second via Robot Operating System 2 (ROS2, version: Humble Hawksbill) connected to a laptop computer (Lenovo Legion Pro 7 16IRX8H equipped with an Intel i9-13900HX CPU, 32 GB of RAM, and an NVIDIA GeForce RTX 4090 Laptop GPU with 16 GB of VRAM). The color and depth images were stored in ROS2 bag files.

2.2.2. Data collection on the harvester

RGB-D image acquisition was carried out over three growing seasons: 2023, 2024, and 2025. During the 2023



(a)



(b)



(c)

Figure 2: (a) and (b) provide overviews of the first imaging system, which was installed on a single-row potato harvester and used during the 2023 growing season. (c) Inside the imaging box, an RGB-D camera was installed, together with four LED strips that provided the necessary illumination inside the box. The sides of box were covered with a reflective curtain to diffuse the light evenly across the conveyor belt.

season, data were collected with the first imaging system (Figure 2b), while the second imaging system (Figure 3a) was used for data acquisition in the 2024 and 2025 seasons.

Image data were recorded in three potato fields located in Sarabetsu, Japan. The potato cultivars in the fields were Sayaka, and Kitahime. During the 2023 season, images from

two other cultivars, Corolle and Haruka, were collected by operating the harvester indoors, where boxes of potatoes were manually dumped onto the conveyor belt. As a result, the images collected in this study exhibited a wide range of tuber shapes and sizes, with each cultivar showing distinct

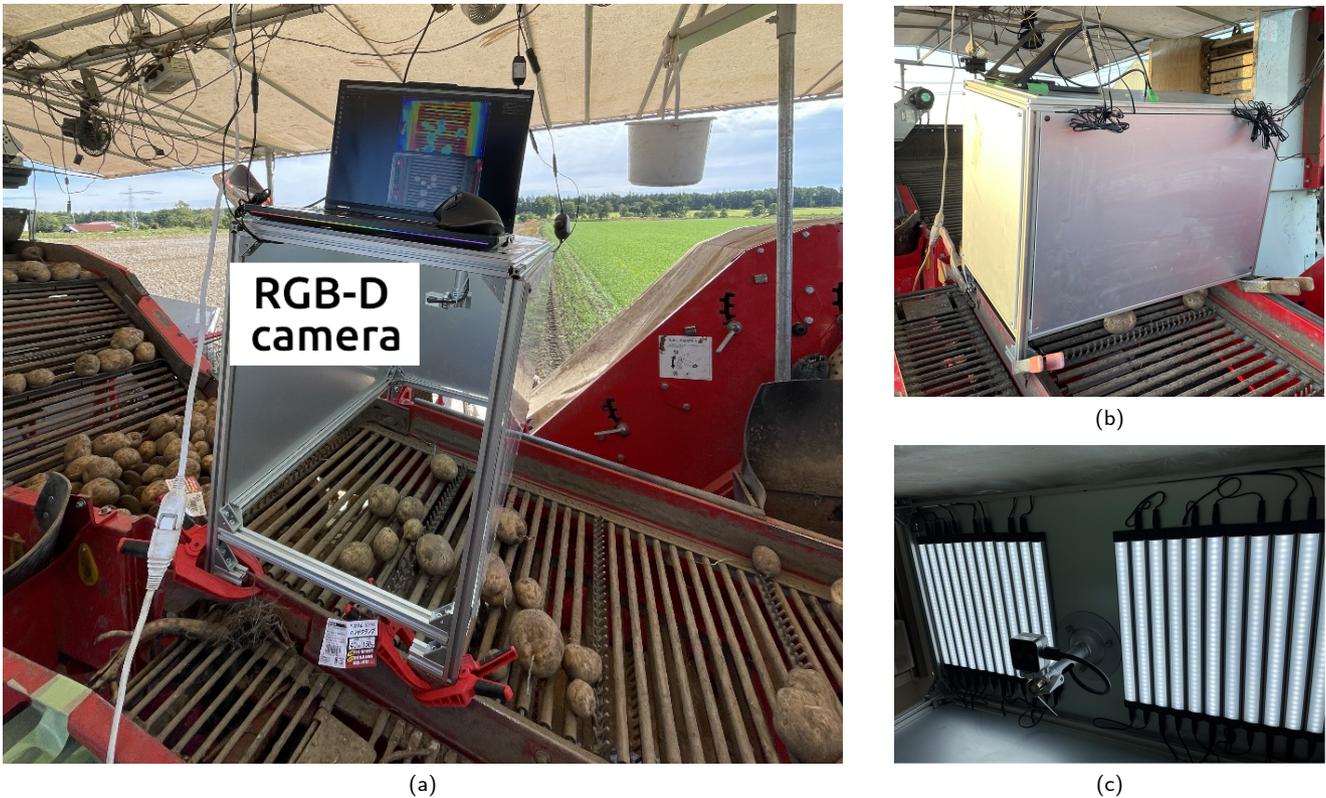


Figure 3: (a) Overview of the second imaging system, which was used during the 2024 and 2025 growing seasons. One side plate has been removed to reveal the internal components. (b) The yield monitoring system with enclosed sides during operation on the harvester. (c) On the ceiling plate of the enclosure, 24 LED strips and the RGB-D camera were installed.

visual characteristics. Haruka tubers were the most elongated, Corolle were moderately elongated, and Sayaka and Kitahime were comparatively rounder. In 2025, a fraction of Kitahime tubers were deformed due to the hot and dry growing conditions in Sarabetsu, which negatively affected tuber development and shape. Overall, the diversity of tuber shapes across the four cultivars, combined with the deformities observed in Kitahime, provided a rich and diverse dataset for analyzing the performance of the PointRAFT weight regression.

During both field harvesting and indoor operation, 859 potato tubers of various sizes and shapes were collected from the conveyor belt of the potato harvester. The tubers were collected according to the following procedure: a person positioned in front of the imaging system randomly selected a potato tuber from the conveyor belt and inserted a colored thumbtack into it. The tack was inserted so that it remained visible in the RGB images as the tuber moved through the camera's field of view. Based on the conveyor belt speed, each tagged potato appeared in approximately 20–35 consecutive image frames. After imaging recording in a ROS2 bag file, a second person positioned behind the imaging system collected the tagged tuber and placed it into an uniquely labeled paper bag. All samples were transported to the barn, where the tubers were weighed using a digital scale. Detailed information on the acquired datasets is presented in Table 1.

2.2.3. Point cloud creation

After image acquisition, the ROS2 bag files were processed to extract individual RGB and depth frames corresponding to the sampled potato tubers. The selected RGB frames were annotated using LabelMe software (version 5.9.1), with assistance from the Segment Anything 2 model (SAM2) to increase annotation accuracy and speed. Once the tubers were annotated, the corresponding masks were clipped onto the RGB and depth frames.

To generate the point clouds from the clipped RGB and depth frames, the Open3D software library (version 0.19.0) was used, along with the intrinsic camera parameters of the Intel RealSense D405. The obtained point clouds were post-processed to remove distant, low-density clusters of points that contained fewer points and were spatially separated from the main cluster representing the potato tuber. The filtering was performed using a custom method that applied the DBSCAN clustering algorithm to identify and eliminate these irrelevant clusters. In total, 26,688 point clouds were extracted from the 859 potato tubers, averaging 31 point clouds per tuber.

2.2.4. Dataset splits

The acquired dataset was split into three subsets for training, validation, and independent testing of PointRAFT. To ensure a balanced representation of potato weights across the dataset splits, we applied stratified sampling. First, the

Table 1

Overview of the three potato fields and the indoor experiment from which the data was acquired.

Year	Field	Latitude	Longitude	Camera-belt distance [m]	Potato cultivar	Sampled potatoes	Note
2023	1	42.610316	143.156753	0.35	Sayaka	280	
2023	-	-	-	0.35	Kitahime	31	Indoor
2023	-	-	-	0.35	Corolle	30	Indoor
2023	-	-	-	0.35	Haruka	30	Indoor
2024	2	42.616872	143.174312	0.46	Kitahime	258	
2025	3	42.600608	143.174741	0.46	Kitahime	230	
Total						859	

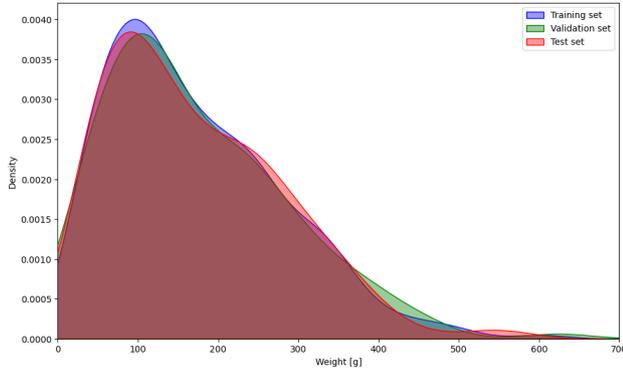


Figure 4: Kernel density estimate plot for visualizing the weight distribution in the training, validation, and test set.

weights of the 859 sampled potatoes were grouped into ten bins, each containing a similar number of samples (ranging from 83 to 88 potatoes per bin). The bin edges were set at 16 g, 51 g, 74 g, 98 g, 121 g, 150 g, 185 g, 224 g, 262 g, 325 g and 625 g. The dataset was then split with stratified labels into a training set of 515 tubers (60%), a validation set of 172 tubers (20%), and a test set of 172 tubers (20%). The total number of point clouds belonging to the training samples was 16,108, the total for validation samples was 5,326, and the total for testing samples was 5,254. Despite using stratified sampling, the dataset remained heavily imbalanced, with smaller and larger tubers underrepresented relative to moderately sized tubers (50 - 250 grams), refer to Figure 4.

2.3. Neural network training

2.3.1. Embedding of tuber height

Prior to applying the PyTorch Geometric Center transform to our point clouds (see Section 2.1.1), the height of each tuber was estimated by subtracting the z-coordinate of the 3D point in \mathbf{P}_0 closest to the camera from the fixed camera-to-belt distance (Table 1). The approximate height was included as an additional feature in PyTorch Geometric’s in-memory dataset class (`InMemoryDataset`). This class allowed efficient storage of all point clouds and their associated features in memory, enabling fast access during training and minimizing disk input/output overhead.

2.3.2. Data augmentation

The `InMemoryDataset` also incorporated the data augmentation pipeline, enabling all transformations to be applied on-the-fly during batch loading. The augmentation pipeline consisted of seven transformations, which were applied during training to enhance the generalization performance of PointRAFT.

The first transformation was a random jitter of $5 \cdot 10^{-4}$, which added small noise to the points’ 3D positions to simulate sensor inaccuracies and noise. Three random rotations were applied along the x, y, and z axes, with each axis rotated by random angles between 0° and 2° , simulating different viewing perspectives. Two random flips were applied along the x and y axes with a 50% probability, further diversifying the possible orientations. Finally, a random shear transformation with a shear factor of 0.2 was applied across all three axes to simulate geometric distortions.

2.3.3. Data sampler for imbalanced dataset

Given the significant class imbalance in our dataset (Figure 4), it was crucial to address this imbalance to prevent PointRAFT from overfitting to the more frequent classes. To ensure balanced class representation during training, we discretized the tuber weight values into 10 distinct classes: 0-50 g, 50-100 g, 100-150 g, 150-200 g, 200-250 g, 250-300 g, 300-350 g, 350-400 g, 400-450 g, and 450 g and beyond. This discretization was necessary for integration with the `ImbalancedSampler` class in PyTorch Geometric, which adjusted the sampling probabilities based on class frequency, giving higher probabilities to underrepresented classes. The approach ensured that each batch during training had a more balanced distribution of tuber weights, reducing bias and enabling the model to better generalize across the entire range of tuber sizes.

2.3.4. Training details

Instead of training PointRAFT with a fixed set of hyperparameters, the Optuna software framework (version 4.6.0) was employed to automatically identify an effective hyperparameter set based on validation performance. The optimization process explored batch sizes ranging from 16 to 64 in increments of 16, learning rates selected from $1 \cdot 10^{-2}$, $5 \cdot 10^{-3}$, $1 \cdot 10^{-3}$, $5 \cdot 10^{-4}$, $1 \cdot 10^{-4}$, and weight decay values from $1 \cdot 10^{-2}$, $1 \cdot 10^{-3}$, $1 \cdot 10^{-4}$, $1 \cdot 10^{-5}$.

In addition, Optuna evaluated multiple loss functions, including mean squared error (MSE), mean absolute error (L1), and Smooth L1 loss. All Optuna trials were trained using the Adaptive Moment Estimation (Adam) optimizer for 50 epochs, with an exponentially decaying learning rate scheduler applied using a decay factor of 0.97.

The best-performing hyperparameter set consisted of a learning rate of $1 \cdot 10^{-3}$, a weight decay of $1 \cdot 10^{-4}$, a batch size of 32, and the Smooth L1 loss function (Equation 2). The β value of the Smooth L1 loss function was set to 20.0, meaning that errors smaller than 20 were treated quadratically, like MSE, to emphasize accurate predictions for most samples, while errors larger than 20 were treated linearly, like L1, preventing extreme outliers from dominating training. Using these hyperparameters, PointRAFT was trained.

$$\mathcal{L}_{\text{SmoothL1}}(W_i, \hat{W}_i) = \begin{cases} \frac{0.5 \cdot (W_i - \hat{W}_i)^2}{\beta}, & \text{if } |W_i - \hat{W}_i| < \beta, \\ |W_i - \hat{W}_i| - 0.5 \cdot \beta, & \text{otherwise,} \end{cases} \quad (2)$$

where W_i denotes the ground truth weight, and \hat{W}_i the predicted weight.

2.4. Evaluation

2.4.1. Benchmark analysis

To assess the performance of PointRAFT in estimating tuber weight from partial point clouds, we compared it against a benchmark linear regression method, following the approaches of Su et al. (2018) and Jang et al. (2023). The linear regression model was trained using the length and width of an oriented 3D bounding box extracted from Open3D, along with the estimated height using the method described in Section 2.3.1. The linear regression model was trained and tested on the same set of point clouds as PointRAFT.

2.4.2. Performance and shape metrics

We assessed the performance of both PointRAFT and the linear regression method using the Mean Absolute Error (MAE, Equation 3), and the Root Mean Squared Error (RMSE, Equation 4). To quantify shape, we calculated an elongation factor by dividing the longest dimension of the 3D bounding box of the partial point cloud by the shortest dimension. Spherical tubers have an elongation factor close to 1, while more elongated tubers have values approaching 2. Both performance and shape metrics were computed for the entire test set, as well as for subsets grouped by the four potato cultivars and the three growing seasons. This allowed us to relate model performance to biological variation and environmental factors affecting tuber growth.

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |W_i - \hat{W}_i| \quad (3)$$

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (W_i - \hat{W}_i)^2} \quad (4)$$

where N denotes the total number of test samples, W_i the corresponding ground truth weight, and \hat{W}_i the predicted weight.

To evaluate the high-throughput processing capability, we calculated the average analysis time for PointRAFT on the point clouds of the test set. This analysis was conducted using the same laptop as described in Section 2.2.1. Since our original `InMemoryDataset` class already stored the preprocessed point clouds, we implemented a custom test script that performed point cloud loading and centering transformation during its analysis. This provided a more realistic assessment of the processing time in a practical data pipeline.

2.4.3. Ablation study

An ablation study was conducted to quantify the contribution of the individual components of PointRAFT to understand their influence on the weight estimation performance. Eleven ablations were evaluated relative to the PointRAFT configuration described earlier, using identical training and evaluation protocols.

The first two ablations focused on the input point cloud resolution, as this directly affects geometric detail and computational cost. The baseline configuration used 1024 input points. To study the trade-off between performance and throughput, the network was trained with 512 and 2048 input points. This analysis allowed us to evaluate whether fewer points are sufficient for reliable regression and whether denser point clouds provide additional benefits for self-occluded tubers.

The next four ablations examined training-related choices. For the first two ablations, PointRAFT was trained without the imbalanced dataset sampler and without data augmentation to study their influence on weight regression. The loss function was then varied by replacing the Smooth L1 loss with MSE and L1 losses, allowing us to assess the sensitivity of the regression task to less and more outlier-sensitive error formulations.

The seventh ablation studied the effect of regularization by disabling dropout in the regression head. This ablation evaluated the importance of dropout for preventing overfitting when learning from self-occluded partial point clouds.

Finally, a series of four ablations examined the role of PointRAFT's main architectural novelty, the object height embedding. First, the width of the MLP used for height encoding was varied by testing two narrower configurations, [1,8,8] and [1,16,16], and one wider configuration, [1,64,64], relative to the default PointRAFT width [1,32,32]. This allowed us to assess how sensitive the model was to the capacity of the height embedding and whether performance gains depended on a specific architectural choice. Lastly, the overall contribution of the height embedding was evaluated by completely removing it from the network, which effectively resulted in a standard PointNet++ regression model.

Table 2

Weight estimation results expressed for the four tested potato cultivars, the three growing seasons and the entire test set. LR and PR are abbreviations of the linear regression benchmark and PointRAFT, respectively.

Subset comparison	Tubers	Elongation factor	Point clouds	MAE [g] ↓		RMSE [g] ↓		
				LR	PR	LR	PR	
Potato cultivar	Sayaka	59	1.19	1670	16.0	10.1	21.0	13.7
	Kitahime	98	1.21	3144	26.5	12.0	36.3	17.4
	Corolle	8	1.30	222	18.6	14.0	24.5	21.2
	Haruka	7	1.51	218	31.2	24.2	36.0	30.2
Growing season	2023	80	1.23	2271	17.9	11.6	23.2	16.6
	2024	49	1.12	1485	22.1	11.3	30.4	15.9
	2025	43	1.30	1498	31.7	13.2	42.5	19.3
Entire test set	172	1.22	5,254	23.0	12.0	31.8	17.2	

3. Results

3.1. PointRAFT performance evaluation against benchmark

Table 2 summarizes the weight estimation performance of the linear regression benchmark (LR) and PointRAFT (PR) for the entire test set, as well as subsets by potato cultivar and growing season. PointRAFT outperformed the linear regression across all subsets, achieving an overall MAE of 12.0 g and an RMSE of 17.2 g, which were 11.0 g (47.8%) and 14.6 g (46%) lower than the linear regression, respectively. The improvement was particularly notable on the 2025 dataset, which included many deformed tubers: here, PointRAFT reduced RMSE by 23.2 g (54.5%) compared to the linear regression, demonstrating its ability to accurately estimate weights for irregularly shaped tubers. PointRAFT’s performance was consistent across the two camera placements, indicating that working directly on raw point clouds allows good generalization across different camera setups. Finally, both methods showed higher errors for the more elongated Haruka and Corolle tubers, likely reflecting the overrepresentation of the more spherical Sayaka and Kitahime tubers in the dataset.

Figure 5 shows the weight predictions for the entire test set. Both methods exhibited the largest errors for tubers heavier than 350 g, likely due to the underrepresentation of large tubers in the dataset. Overall biases were small (+1.49 g for linear regression and +0.94 g for PointRAFT), indicating that both methods had minimal systematic overestimation.

Large prediction variability was observed for different viewpoints of the same tuber, reflecting the challenge of consistently estimating weight from point clouds captured throughout the tuber’s motion along the conveyor belt. Despite this, PointRAFT produced smaller errors and a higher R^2 value (0.97) compared to the linear regression (0.91). This indicates that PointRAFT can predict weights more reliably across a wide range of tuber sizes. Visual comparisons (Figure 6) further illustrate that PointRAFT captures the weight of small, medium, deformed, and large tubers more accurately compared to the linear regression.

Figure 7 shows the cumulative ratios of absolute weight errors for both methods. PointRAFT achieved 57% of its

predictions within 10 g of the ground truth, almost double the 31% observed for linear regression. The advantage remained at larger error thresholds: 81% of PointRAFT predictions were within 20 g, compared to 58% for linear regression, and 99% were within 50 g, compared to 89%. These results highlight PointRAFT’s higher accuracy and reliability across the entire test set.

3.2. Ablation study

The ablation study (Table 3) shows that the height embedding contributed most to weight estimation performance. Removing the embedding, effectively training a standard PointNet++ regression network, increased MAE by 37.5% and RMSE by 42.4%, showing that encoding tuber height clearly improved the weight prediction. The performance of the height embedding also depended on the MLP width: reducing it from [1, 32, 32] to [1, 16, 16] increased MAE by 14.2% and RMSE by 16.9%, while a further reduction to [1, 8, 8] slightly increased RMSE (+5.8%) with no change in MAE. Increasing the width to [1, 64, 64] worsened performance, indicating that choosing an appropriate MLP dimension is important, although no clear trend emerged across different configurations.

Dropout in the regression head also had a strong effect. Removing dropout increased MAE by 9.2% and RMSE by 20.3%, highlighting its role as an effective regularization technique for training PointRAFT on a large point cloud dataset. Similarly, data augmentation also had a noticeable effect: disabling it increased MAE by 6.7% and RMSE by 12.2%, indicating improved network generalization through exposure to greater shape variability. In contrast, training without the imbalanced dataset sampler had minimal impact, with RMSE unchanged and MAE increasing only slightly (+5.8%). This is likely because the Smooth L1 loss already suppressed the influence of outliers and large errors, making the network relatively insensitive to changes in the sample distribution.

Input point density was critical: training on 512 points led to large increases in MAE (+20.8%) and RMSE (+24.4%) with only a minor speed gain (-1.3 ms). Surprisingly, using 2048 points also reduced performance (MAE +8.3%, RMSE +14.0%), suggesting that the fine-scale detail may have introduced 3D shape variation that hindered robust weight

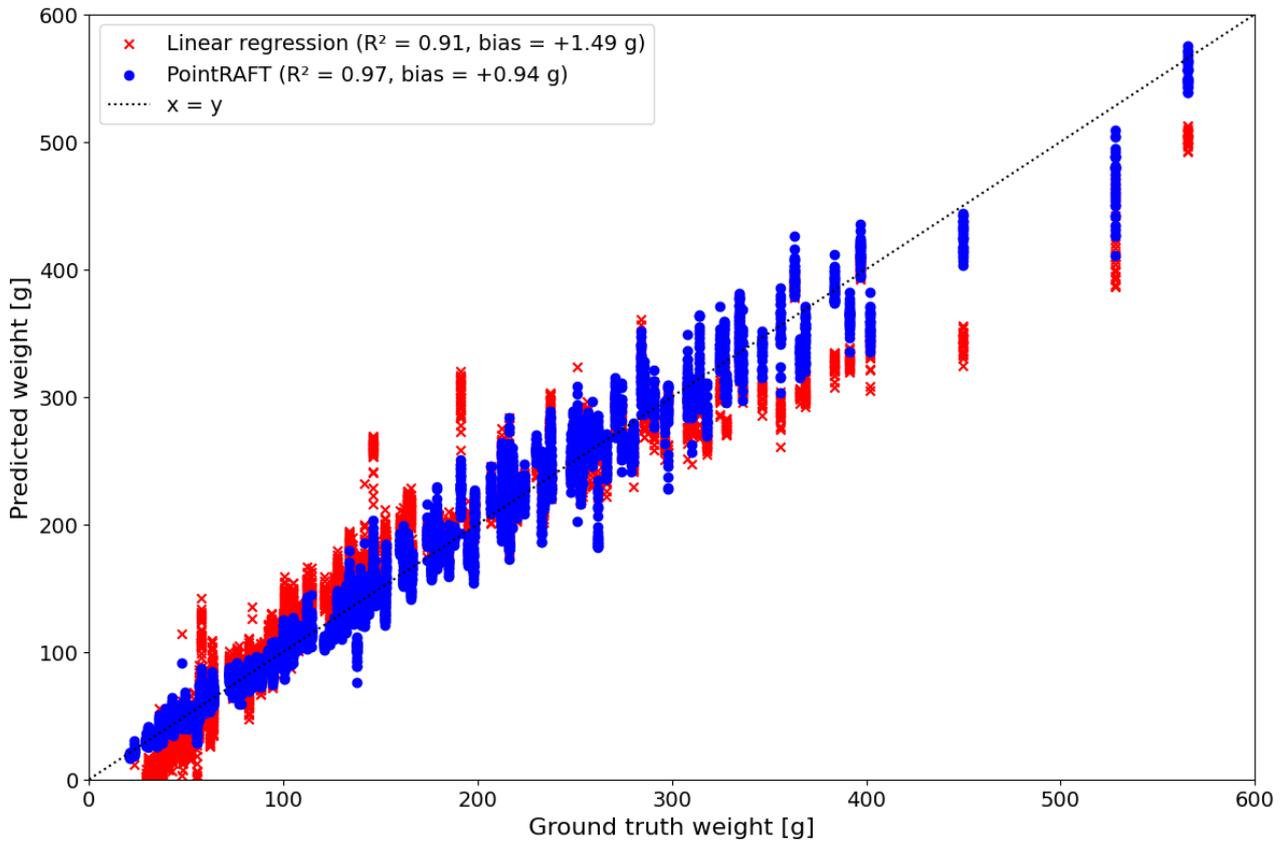


Figure 5: Predicted weights of both methods on the test samples. The blue dots are the predictions from PointRAFT and the red crosses are the predictions from the linear regression benchmark.

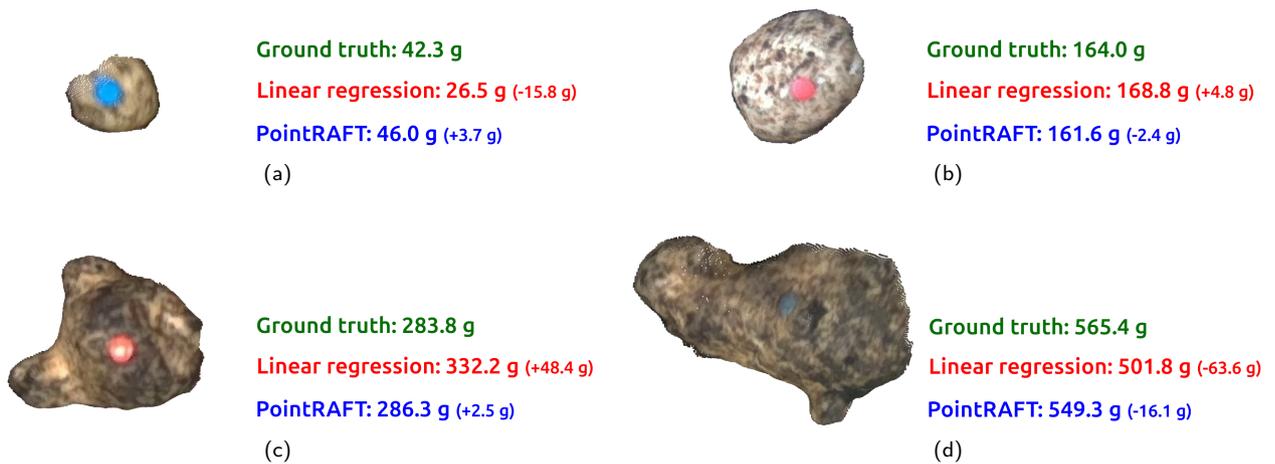


Figure 6: Visual comparison of four predictions by the linear regression benchmark (in red) and PointRAFT (in blue) on potato tubers of different shapes and sizes in our test set: small (a), medium (b), deformed (c), and large (d). The 3D point clouds are zoomed in for clarity and do not represent their actual sizes.

estimation. Training with 1024 points achieved the best balance between prediction accuracy and high-throughput processing, with an average processing time of 6.3 ms.

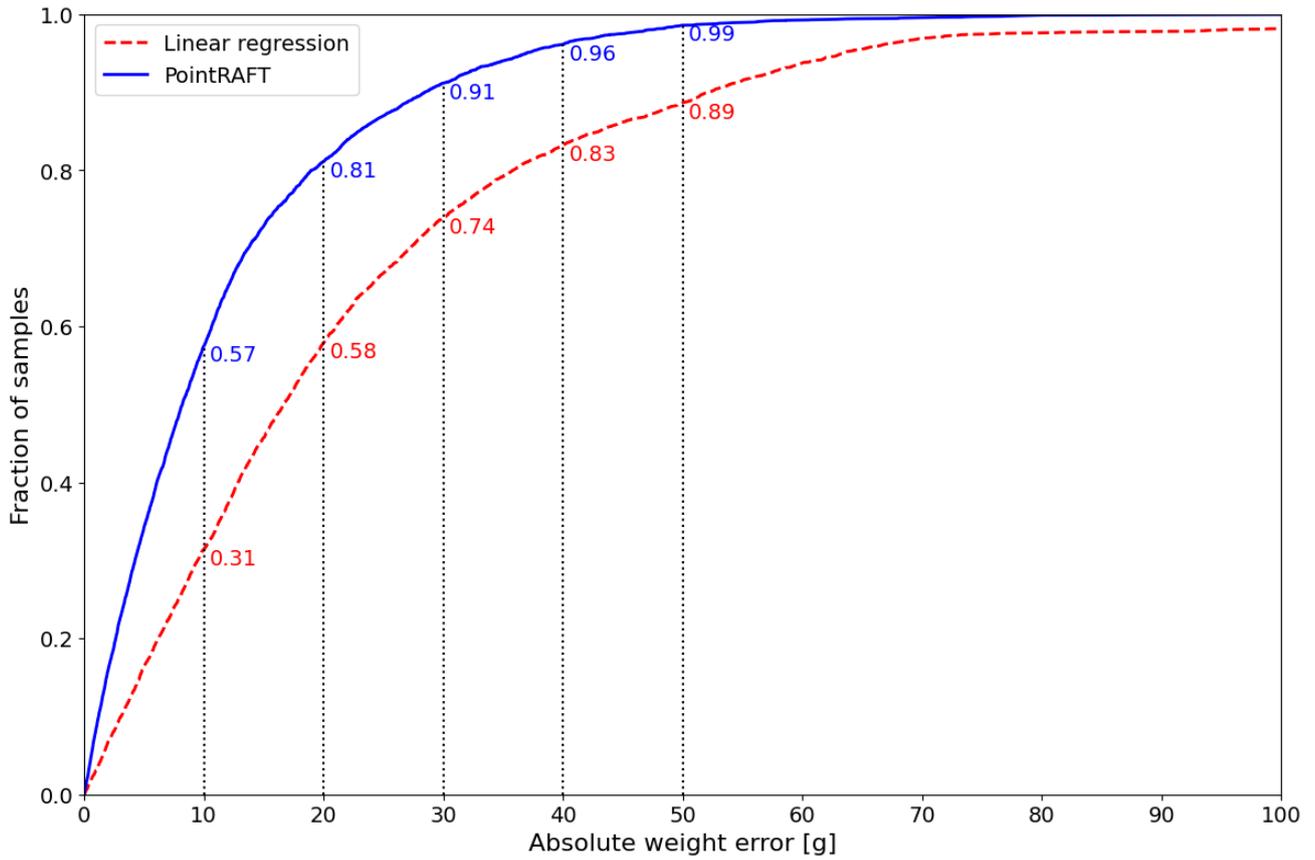


Figure 7: Cumulative ratios of the absolute weight errors of the two methods. The colored numbers summarize the ratios of the weight estimates that were within 10 g, 20 g, 30 g, 40 g and 50 g from the ground truth weight.

Table 3

Performance metrics for the eleven studied ablations relative to the baseline performance of PointRAFT.

Ablation	Category	MAE [g] ↓		RMSE [g] ↓		Time [ms] ↓	
		abs	rel	abs	rel	abs	rel
PointRAFT	Baseline	12.0	-	17.2	-	6.3	-
1024 → 512	Input points	14.5	+20.8%	21.4	+24.4%	5.0	-20.6%
1024 → 2048	Input points	13.0	+8.3%	19.6	+14.0%	9.0	+42.9%
No imbalanced sampler	Network training	12.7	+5.8%	17.2	0.0%	6.3	0.0%
No data augmentation	Network training	12.8	+6.7%	19.3	+12.2%	6.6	+4.8%
Smooth L1 → MSE	Loss function	12.7	+5.8%	19.7	+14.5%	6.8	+7.9%
Smooth L1 → L1	Loss function	12.1	+0.8%	19.1	+11.0%	6.3	0.0%
No dropout	Regularization	13.1	+9.2%	20.7	+20.3%	6.6	+4.8%
Height embed. [8,8]	Network architecture	12.0	0.0%	18.2	+5.8%	6.8	+7.9%
Height embed. [16,16]	Network architecture	13.7	+14.2%	20.1	+16.9%	6.5	+3.2%
Height embed. [64,64]	Network architecture	13.5	+12.5%	18.8	+9.3%	6.8	+7.9%
No height embedding	Network architecture	16.5	+37.5%	24.5	+42.4%	6.2	-1.6%

4. Discussion

Our results showed that PointRAFT consistently outperformed the linear regression benchmark on all test samples, including subsets per potato cultivar and growing season. This consistent improvement suggests that PointRAFT learned geometric features that generalized well across different biological and environmental conditions. One important finding was that the tuber height embedding improved weight prediction the most, while keeping processing time unchanged. This result directly confirms our hypothesis that adding tuber height as an explicit geometric cue improves weight estimation.

When compared to existing 3D methods, PointRAFT shows competitive performance. Compared to CoRe++ (Blok et al., 2025), PointRAFT achieved a 23.9% lower RMSE. This result is encouraging because CoRe++ was tested on a smaller dataset, with fewer cultivars, a fixed camera setup, and no deformed potato tubers. When compared to earlier potato studies using 2D or 3D computer vision, PointRAFT achieved lower MAE and RMSE and higher R^2 values (Long et al., 2018; Lee and Shin, 2020; Huynh et al., 2022; Jang et al., 2023). This result is promising because those studies were performed in laboratory settings with stationary potato tubers and smaller datasets. Only the work by Su et al. (2018) reported a lower MAE (7.7 g), but that study used two complementary views per tuber, capturing almost the entire 3D surface. Even when compared to studies on other crops, PointRAFT remains competitive, with an R^2 of 0.97, similar to results reported for more uniformly shaped crops such as grapes (Zhang et al., 2025).

In addition to accuracy, PointRAFT offers very high processing speed. To our knowledge, no published 3D method for crop weight or volume estimation reported a processing time faster than 6.3 ms. This processing time, however, accounted for only a fraction of the required data preprocessing. Additional steps such as RGB-D to point cloud conversion and point cloud filtering were not included and will increase the total latency in a full data pipeline. Processing time will also depend on the available GPU hardware. For practical deployment, efficient point cloud generation and filtering methods should be used to keep the total processing time low.

For full system integration, both upstream and downstream tasks must be considered. Upstream tasks include generating point clouds from instance segmentation masks. In practice, these masks may contain errors, which can lead to noisy point clouds and reduced prediction accuracy. Recent transformer-based instance segmentation models, such as RFDetr-Seg (Robinson, Robicheaux, Popov, Ramanan and Peri, 2025), are expected to provide masks of sufficient quality for this task. Downstream tasks include combining individual weight estimates into a spatial yield map, which requires integration with a global navigation satellite system (GNSS) to geo-reference predictions and visualize yield variability across the field.

Despite the promising results, this study has several limitations. First, only four potato cultivars were included,

and two of them, Haruka and Corolle, had limited sample sizes. This restricts our ability to assess how well PointRAFT generalizes to cultivars with very different shapes. In particular, the height embedding in PointRAFT may have been biased toward the dominant cultivars, Sayaka and Kitahime, which had similar length–width–height ratios. As a result, PointRAFT may generalize less effectively to tubers that are more elongated or more spherical. Second, most tubers in our dataset were not stacked, as they were manually placed on the conveyor belt. In real harvesting conditions, tuber stacking occurs more often. In such cases, using the conveyor belt as a height reference may lead to overestimation when stacked tubers are interpreted as single objects. This issue could be addressed by incorporating visual context extracted from instance segmentation, allowing tuber height to be estimated relative to neighboring tubers rather than to the conveyor belt. Finally, tuber orientation was not explicitly modeled, so tilting may have caused small height estimation errors. Including orientation information could further improve accuracy.

Several directions for future work follow from our results. First, the variation in predictions across multiple views of the same tuber suggests using a consistency loss to reduce differences between viewpoints. This could further lower MAE and RMSE. Second, the efficient inference of PointRAFT allows multiple stochastic forward passes per tuber to estimate prediction uncertainty. Methods such as Monte Carlo dropout (Gal and Ghahramani, 2016) are well suited for this and can be implemented with minimal changes. Third, if PointRAFT is used as an encoder for 3D shape completion, the search radii in the PointNet++ layers should be revisited. Smaller radii may be better suited to capture local geometric details that are important for reconstructing fine-scale shape. With the current radius settings, the set abstraction layers aggregated information over the entire tuber, so that the network already captured the global 3D shape in its early layers. This suggests that simpler architectures, such as PointNet (Qi et al., 2017a) or PointMLP (Ma, Qin, You, Ran and Fu, 2022), which learn global shape representations without hierarchical aggregation, might have sufficed for encoding the 3D shape of the tuber.

Although our study focused on potato tubers, PointRAFT is not limited to this application. Its modular design, user-friendly code base, and high processing speed make it suitable for other agricultural tasks. Examples include size estimation of occluded fruits and vegetables, as well as prediction of grasp points, 6D poses, or gripper forces for autonomous harvesting.

5. Conclusions

In this study, we introduced PointRAFT, a high-throughput point cloud regression network capable of directly predicting tuber weight from partial point clouds. Its key architectural novelty was the object height embedding, which improved weight predictions and achieved an RMSE of 17.2 g, outperforming the linear regression benchmark (31.8 g). With

an average processing time of 6.3 ms per point cloud, PointRAFT can process up to 150 tubers per second, making it well suited for high-throughput weight estimation on operational harvesters. Its accurate, real-time predictions provide farmers with actionable insights to optimize planting, fertilization, and harvesting in future seasons, helping to increase yield, reduce losses, and better meet market demand.

CRedit authorship contribution statement

Pieter M. Blok: Conceptualization, Methodology, Software, Data curation, Writing - original draft. **Haozhou Wang:** Data curation, Resources, Writing - review & editing. **Hyun Kwon Suh:** Supervision, Writing - review & editing. **Peicheng Wang:** Data curation, Resources. **James Burrige:** Conceptualization, Methodology, Data curation, Writing - review & editing. **Wei Guo:** Conceptualization, Methodology, Funding acquisition, Project administration, Supervision, Writing - review & editing.

Funding

This study is funded by the Sarabetsu Village "Endowed Chair for Field Phenomics" project in Hokkaido, Japan.

Acknowledgements

We would like to thank Okada Farm for providing the potato fields on which we sampled the potatoes and acquired the RGB-D images. We thank Minato Miyauchi, Kyosuke Miyamoto, Kunihiro Tanaka, Xinbin Zhuang, Shuai Xiang, Ting Jiang, Sylvain Grison, Yuto Imachi and Irena Drofová for their help with the image acquisition and ground truth measurements. Special thanks to Yumi Hoshino for the stimulating support and the cross-pollinating activities in Japan.

Declaration of generative AI and AI-assisted technologies in the manuscript preparation process

During the preparation of this work, the author(s) used ChatGPT to support software debugging and to assist with language refinement. All content was subsequently reviewed and edited by the author(s), who take(s) full responsibility for the final published article.

Data availability

A subset of the datasets generated and/or analyzed during this study is publicly available at: <https://huggingface.co/datasets/UTokyo-FieldPhenomics-Lab/3DPotatoTwin>

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- Bigalke, A., Hansen, L., Diesel, J., Heinrich, M.P., 2021. Seeing under the cover with a 3d u-net: point cloud-based weight estimation of covered patients. *International Journal of Computer Assisted Radiology and Surgery* 16, 2079–2087. doi:10.1007/s11548-021-02476-0.
- Blok, P.M., Magistri, F., Stachniss, C., Wang, H., Burrige, J., Guo, W., 2025. High-throughput 3d shape completion of potato tubers on a harvester. *Computers and Electronics in Agriculture* 228, 109673. URL: <https://www.sciencedirect.com/science/article/pii/S0168169924010640>, doi:10.1016/j.compag.2024.109673.
- Bullock, D.S., Lowenberg-DeBoer, J., Swinton, S.M., 2002. Adding value to spatially managed inputs by understanding site-specific yield response. *Agricultural Economics* 27, 233–245. URL: <https://www.sciencedirect.com/science/article/pii/S0169515002000786>, doi:10.1016/S0169-5150(02)00078-6.
- Cai, Z., Jin, C., Xu, J., Yang, T., 2020. Measurement of potato volume with laser triangulation and three-dimensional reconstruction. *IEEE Access* 8, 176565–176574. doi:10.1109/ACCESS.2020.3027154.
- Camire, M.E., Kubow, S., Donnelly, D.J., 2009. Potatoes and human health. *Critical Reviews in Food Science and Nutrition* 49, 823–840. URL: <https://pubmed.ncbi.nlm.nih.gov/19960391/>, doi:10.1080/10408390903041996.
- Chen, H., Liu, S., Wang, C., Wang, C., Gong, K., Li, Y., Lan, Y., 2023. Point cloud completion of plant leaves under occlusion conditions based on deep learning. *Plant Phenomics* 5. URL: <https://spj.science.org/doi/abs/10.34133/plantphenomics.0117>, doi:10.34133/plantphenomics.0117.
- Dolata, P., Wróblewski, P., Mrzygłód, M., Reiner, J., 2021. Instance segmentation of root crops and simulation-based learning to estimate their physical dimensions for on-line machine vision yield monitoring. *Computers and Electronics in Agriculture* 190, 106451. URL: <https://www.sciencedirect.com/science/article/pii/S0168169921004683>, doi:10.1016/j.compag.2021.106451.
- ElMasry, G., Cubero, S., Moltó, E., Blasco, J., 2012. In-line sorting of irregular potatoes by using automated computer-based machine vision system. *Journal of Food Engineering* 112, 60–68. URL: <https://www.sciencedirect.com/science/article/pii/S0260877412001690>, doi:10.1016/j.jfoodeng.2012.03.027.
- Emwinghare, I., Al-Mallahi, A.A., Esau, T.J., Campelo, F., 2025. Accurate estimation of tuber size in large potato throughput at potato storage using machine vision and machine learning techniques. *Smart Agricultural Technology* 10, 100860. URL: <https://www.sciencedirect.com/science/article/pii/S2772375525000930>, doi:10.1016/j.atech.2025.100860.
- Fei, B., Yang, W., Chen, W.M., Li, Z., Li, Y., Ma, T., Hu, X., Ma, L., 2022. Comprehensive review of deep learning-based 3d point cloud completion processing and analysis. *IEEE Transactions on Intelligent Transportation Systems* 23, 22862–22883. doi:10.1109/tits.2022.3195555.
- Fu, L., Gao, F., Wu, J., Li, R., Karkee, M., Zhang, Q., 2020. Application of consumer rgb-d cameras for fruit detection and localization in field: A critical review. *Computers and Electronics in Agriculture* 177, 105687. URL: <https://www.sciencedirect.com/science/article/pii/S0168169920319530>, doi:10.1016/j.compag.2020.105687.
- Gal, Y., Ghahramani, Z., 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, in: *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, JMLR.org*. p. 1050–1059.
- Ge, Y., Xiong, Y., From, P.J., 2020. Symmetry-based 3d shape completion for fruit localisation for harvesting robots. *Biosystems Engineering* 197, 188–202. URL: <https://www.sciencedirect.com/science/article/pii/S1537511020301963>, doi:10.1016/j.biosystemseng.2020.07.003.
- Han, M., Wang, L., Xiao, L., Zhang, H., Zhang, C., Xu, X., Zhu, J., 2023. Quickfps: Architecture and algorithm co-design for farthest point sampling in large-scale point clouds. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*.
- Hofstee, J., Molema, G., 2003. Volume estimation of potatoes partly covered with dirt tare, in: *Proc. of the ASAE Annual Meeting 2003*.

- Hou, Z., Huang, L., Zhang, Q., Miao, Y., 2023. Body weight estimation of beef cattle with 3d deep learning model: Pointnet++. *Computers and Electronics in Agriculture* 213, 108184. URL: <https://www.sciencedirect.com/science/article/pii/S0168169923005720>, doi:10.1016/j.compag.2023.108184.
- Hu, P., Dai, X., Zhao, R., Wang, H., Ma, Y., Munteanu, A., 2023. Point2partvolume: Human body volume estimation from a single depth image. *IEEE Transactions on Instrumentation and Measurement* 72, 1–12. doi:10.1109/TIM.2023.3284948.
- Huynh, T.T.M., TonThat, L., Dao, S.V.T., 2022. A vision-based method to estimate volume and mass of fruit/vegetable: Case study of sweet potato. *International Journal of Food Properties* 25, 717–732. doi:10.1080/10942912.2022.2057528.
- Jang, S.H., Moon, S.P., Kim, Y.J., Lee, S.H., 2023. Development of potato mass estimation system based on deep learning. *Applied Sciences* 13. URL: <https://www.mdpi.com/2076-3417/13/4/2614>, doi:10.3390/app13042614.
- Kabir, M., Myat Swe, K., Kim, Y.J., Chung, S., Jeong, D.U., Lee, S.H., 2018. Sensor comparison for yield monitoring systems of small-sized potato harvesters, in: *Proc. of the 14th International conference on precision agriculture*.
- Lee, Y.J., Kim, K.D., Lee, H.S., Shin, B.S., 2018. Vision-based potato detection and counting system for yield monitoring. *Journal of Biosystems Engineering* 43, 103–109. doi:10.5307/JBE.2018.43.2.103.
- Lee, Y.J., Shin, B.S., 2020. Development of potato yield monitoring system using machine vision. *Journal of Biosystems Engineering* 45, 282–290. doi:10.1007/s42853-020-00069-4.
- Long, Y., Wang, Y., Zhai, Z., Wu, L., Li, M., Sun, H., Su, Q., 2018. Potato volume measurement based on rgb-d camera, in: *Proc. of the 6th IFAC Conference on Bio-Robotics*. doi:10.1016/j.ifacol.2018.08.157.
- Ma, X., Qin, C., You, H., Ran, H., Fu, Y., 2022. Rethinking network design and local geometry in point cloud: A simple residual mlp framework. *arXiv preprint arXiv:2202.07123*.
- Magistri, F., Marcuzzi, R., Marks, E., Sodano, M., Behley, J., Stachniss, C., 2024. Efficient and Accurate Transformer-Based 3D Shape Completion and Reconstruction of Fruits for Agricultural Robots, in: *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*.
- Magistri, F., Marks, E., Nagulavantha, S., Vizzo, I., Labe, T., Behley, J., Halstead, M., McCool, C., Stachniss, C., 2022. Contrastive 3d shape completion and reconstruction for agricultural robots using rgb-d frames. *IEEE Robotics and Automation Letters* 7, 10120–10127. doi:10.1109/LRA.2022.3193239.
- Marangoz, S., Zaenker, T., Menon, R., Bennewitz, M., 2022. Fruit mapping with shape completion for autonomous crop monitoring, in: *Proc. of the IEEE International Conference on Automation Science and Engineering (CASE)*. doi:10.1109/CASE49997.2022.9926466.
- Noordam, J.C., Otten, G.W., Timmermans, T.J., van Zwol, B.H., 2000. High-speed potato grading and quality inspection based on a color vision system, in: *Proc. of the Machine Vision Applications in Industrial Inspection VIII*.
- Pan, Y., Magistri, F., Labe, T., Marks, E., Smitt, C., McCool, C., Behley, J., Stachniss, C., 2023. Panoptic mapping with fruit completion and pose estimation for horticultural robots, in: *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. doi:10.1109/IROS55552.2023.10342067.
- Pandey, N., Kumar, S., Pandey, R., 2019. Grading and defect detection in potatoes using deep learning, in: Verma, S., Tomar, R.S., Chaurasia, B.K., Singh, V., Abawajy, J. (Eds.), *Communication, Networks and Computing*, Springer Singapore. pp. 329–339.
- Paudel, S., de Sousa, R.V., Sharma, S.R., Brown-Brandl, T., 2024. Deep learning models to predict finishing pig weight using point clouds. *Animals* 14. URL: <https://www.mdpi.com/2076-2615/14/1/31>, doi:10.3390/ani14010031.
- Qi, C.R., Su, H., Mo, K., Guibas, L.J., 2017a. Pointnet: Deep learning on point sets for 3d classification and segmentation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 652–660.
- Qi, C.R., Yi, L., Su, H., Guibas, L.J., 2017b. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems* 30. doi:10.48550/arXiv.1706.02413.
- Razmjoooy, N., Mousavi, B.S., Soleymani, F., 2012. A real-time mathematical computer method for potato inspection using machine vision. *Computers & Mathematics with Applications* 63, 268–279. URL: <https://www.sciencedirect.com/science/article/pii/S0898122111009850>, doi:10.1016/j.camwa.2011.11.019.
- Robinson, I., Robicheaux, P., Popov, M., Ramanan, D., Peri, N., 2025. Rf-detr: Neural architecture search for real-time detection transformers. *arXiv preprint arXiv:2511.09554* arXiv:2511.09554.
- Si, Y., Sankaran, S., Knowles, N.R., Pavek, M.J., 2018. Image-based automated potato tuber shape evaluation. *Journal of Food Measurement and Characterization* 12, 702–709. doi:10.1007/s11694-017-9683-2.
- Su, Q., Kondo, N., Li, M., Sun, H., Al Riza, D.F., Habaragamuwa, H., 2018. Potato quality grading based on machine vision and 3d shape analysis. *Computers and Electronics in Agriculture* 152, 261–268. URL: <https://www.sciencedirect.com/science/article/pii/S016816991830067X>, doi:10.1016/j.compag.2018.07.012.
- Van Evert, F.K., Gaitan-Cremaschi, D., Fountas, S., Kempenaar, C., 2017. Can precision agriculture increase the profitability and sustainability of the production of potatoes and olives? *Sustainability* 9. URL: <https://www.mdpi.com/2071-1050/9/10/1863>, doi:10.3390/su9101863.
- Wang, H., Blok, P.M., Burridge, J., Jiang, T., Miyauchi, M., Miyamoto, K., Tanaka, K., Guo, W., 2025a. 3dpotatotwin: a paired potato tuber dataset for 3d multi-sensory fusion. *Plant Phenomics* 7, 100123. URL: <https://www.sciencedirect.com/science/article/pii/S2643651525001293>, doi:10.1016/j.plaphe.2025.100123.
- Wang, K., Pan, Y., Magistri, F., Kooistra, L., Stachniss, C., Wang, W., Valente, J., 2026. Uav-based monocular 3d panoptic mapping for fruit shape completion in orchard. *ISPRS Journal of Photogrammetry and Remote Sensing* 231, 608–621. URL: <https://www.sciencedirect.com/science/article/pii/S0924271625004526>, doi:10.1016/j.isprsjprs.2025.11.013.
- Wang, W., Lin, C., Shui, H., Zhang, K., Zhai, R., 2025b. Adaptive symmetry self-matching for 3d point cloud completion of occluded tomato fruits in complex canopy environments. *Plants* 14. URL: <https://www.mdpi.com/2223-7747/14/13/2080>, doi:10.3390/plants14132080.
- Xu, D., Chen, G., Jing, W., 2023. A single-tree point cloud completion approach of feature fusion for agricultural robots. *Electronics* 12. URL: <https://www.mdpi.com/2079-9292/12/6/1296>, doi:10.3390/electronics12061296.
- Zamani, D.M., Ghosamiparashkahi, M., Faghavi, A., Ghezavati, J., 2014. Design, implementation and evaluation of potato yield monitoring system. *International Journal of Technical Research and Applications* 2, 36–39.
- Zhang, H., Xu, F., Wu, Y., Hu, H., Dai, X., 2017. Progress of potato staple food research and industry development in china. *Journal of Integrative Agriculture* 16, 2924–2932. URL: <https://www.sciencedirect.com/science/article/pii/S2095311917617362>, doi:10.1016/S2095-3119(17)61736-2.
- Zhang, W., Zheng, C., Wang, C., Blok, P.M., Wang, H., Guo, W., 2025. Grapecpnet: A self-supervised point cloud completion network for 3d phenotyping of grape bunches. *Computers and Electronics in Agriculture* 237, 110595. URL: <https://www.sciencedirect.com/science/article/pii/S016816992500701X>, doi:10.1016/j.compag.2025.110595.
- Zheng, H., Ma, C., Liu, D., Huang, J., Chen, R., Fang, C., Yang, J., Berckmans, D., Norton, T., Zhang, T., 2025. Weight prediction method for individual live chickens based on single-view point cloud information. *Computers and Electronics in Agriculture* 234, 110232. URL: <https://www.sciencedirect.com/science/article/pii/S0168169925003382>, doi:10.1016/j.compag.2025.110232.